

# Equivalence of Queries with Nested Aggregation

by

David E. DeHaan

A thesis  
presented to the University of Waterloo  
in fulfillment of the  
thesis requirement for the degree of  
Doctor of Philosophy  
in  
Computer Science

Waterloo, Ontario, Canada, 2009

© David E. DeHaan 2009



I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.



## Abstract

Query equivalence is a fundamental problem within database theory. The correctness of all forms of logical query rewriting—join minimization, view flattening, rewriting over materialized views, various semantic optimizations that exploit schema dependencies, federated query processing and other forms of data integration—requires proving that the final executed query is equivalent to the original user query. Hence, advances in the theory of query equivalence enable advances in query processing and optimization.

In this thesis we address the problem of deciding query equivalence between conjunctive SQL queries containing aggregation operators that may be nested. Our focus is on understanding the interaction between nested aggregation operators and the other parts of the query body, and so we model aggregation functions simply as abstract collection constructors. Hence, the precise language that we study is a conjunctive algebraic language that constructs complex objects from databases of flat relations. Using an encoding of complex objects as flat relations, we reduce the query equivalence problem for this algebraic language to deciding equivalence between relational encodings output by traditional conjunctive queries (not containing aggregation). This *encoding-equivalence* cleanly unifies and generalizes previous results for deciding equivalence of conjunctive queries evaluated under various processing semantics. As part of our study of aggregation operators that can construct empty sub-collections—so-called “scalar” aggregation—we consider query equivalence for conjunctive queries extended with a left outer join operator, a very practical class of queries for which the general equivalence problem has never before been analyzed. Although we do not completely solve the equivalence problem for queries with outer joins or with scalar aggregation, we do propose useful sufficient conditions that generalize previously known results for restricted classes of queries. Overall, this thesis offers new insight into the fundamental principles governing the behaviour of nested aggregation.



## Acknowledgements

I would like to thank my supervisor, Frank Wm. Tompa, for his continued encouragement and support throughout my tenure as a graduate student at Waterloo. I greatly appreciate the freedom that Frank provided me to pursue research problems of my own interest, especially when charting the course of this dissertation. In his patient approach to graduate supervision, Frank prioritized giving me the opportunity to grow as an academic. I am extremely grateful for that. I could not ask for a better mentor.

Thank you to my external examiner, Dirk Van Gucht from Indiana University, and to the other members of my examining committee, Ashraf Aboulnaga, David Toman, and Ross Willard, for their valuable comments and questions. My oral defence of this dissertation was a highlight of my time as a graduate student; I am grateful for their part in making that experience an enjoyable yet meaningful culmination of my work as a graduate student.

I did not initially set out to write a theoretical study of query equivalence. In late 2003, Glenn Paulley from Sybase iAnywhere suggested that there were many unsolved problems regarding the use of materialized views. I am very grateful to Glenn for that suggestion. My preliminary investigations led me to an ongoing research project by Paul Larson at Microsoft Research which modelled the use of distributed caches as a materialized view rewriting problem. Paul invited me to work with him during the summer of 2004, and I am very thankful to him and to MSR for furnishing that opportunity. Working within the Microsoft SQL Server code base was a wonderful experience, affording many insights into the design of commercial-grade query optimizers. Reflecting on my work at the completion of my internship, I had some lingering discontent over my inability to determine when there existed a query rewriting over a materialized view that our algorithm was failing to find. The seeds of the thesis topic had been sown.

Returning to graduate study with an interesting problem in hand, I found myself woefully unprepared to tackle it. I am indebted to both David Toman and Grant Weddell for their patient handling of the countless questions of a database-theory neophyte. Their combined help in pointing me toward useful concepts and/or literature was instrumental in the formation and writing of this dissertation.

Many friends and colleagues within the UW database group have encouraged and supported me. I owe a special debt of gratitude to Khuzaima Daudjee, a close friend and former office-mate. Khuzaima's patient listening and sage counsel over many years and during uncountable coffee runs got me past many unforeseen obstacles.

Finally, and most of all, I thank my wife Cara for her confidence in my abilities and firm support of my goal of completing this dissertation. As draft proofs coalesced only to disintegrate like sand castles on a beach, she journeyed beside me and helped me to keep perspective and focus.

*Soli Deo gloria.*





# Contents

<b>List of Figures</b>	<b>xiii</b>
<b>List of Symbols</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivating Examples . . . . .	2
1.2 Related Problems and Technologies . . . . .	7
1.3 Road Map . . . . .	10
<b>2 Equivalence of Conjunctive Queries that Construct Objects</b>	<b>11</b>
2.1 Complex Objects . . . . .	11
2.1.1 Objects and Sorts . . . . .	11
2.1.2 Object Linearization . . . . .	17
2.1.3 Object-Constructing Queries . . . . .	24
2.2 Encoding Objects as Relations . . . . .	34
2.2.1 Encoding and Decoding Chain Objects . . . . .	34
2.2.2 Encoding and Decoding Linearized Objects . . . . .	38
2.2.3 Encoding-Equality . . . . .	41
2.3 Encoding Queries . . . . .	44
2.3.1 Hierarchical Conjunctive Encoding Queries . . . . .	46
2.3.2 Reducing Equivalence to Encoding-Equivalence . . . . .	52
2.4 Relevant Literature . . . . .	65
2.4.1 Query Languages for Complex Objects . . . . .	65
2.4.2 Relational Encodings of Complex Objects . . . . .	67
2.4.3 Algebraic Transformations . . . . .	68
2.4.4 Query Equivalence . . . . .	69

<b>3</b>	<b>Encoding-Equivalence of Conjunctive Encoding Queries</b>	<b>73</b>
3.1	A Normal Form for CEQs . . . . .	74
3.1.1	Dependencies Over Query Results . . . . .	74
3.1.2	Converting CEQs to Normal Form . . . . .	79
3.2	CEQ Encoding-Equivalence . . . . .	84
3.3	Relevant Literature . . . . .	86
3.3.1	Conjunctive Queries under Set Semantics . . . . .	87
3.3.2	Other Processing Semantics . . . . .	87
3.3.3	Aggregation Queries . . . . .	88
<b>4</b>	<b>Equivalence of Conjunctive Queries with Outer Joins</b>	<b>91</b>
4.1	Conjunctive Queries with Outer Joins . . . . .	91
4.1.1	The SPCL Algebra . . . . .	92
4.1.2	Directed Acyclic Conjunctive Queries . . . . .	94
4.2	A Canonical Form for DACQs . . . . .	97
4.2.1	Functionally-Determined Nullability . . . . .	98
4.2.2	Converting DACQs to Canonical Form . . . . .	102
4.3	Equivalence of Constant-Headed Queries . . . . .	107
4.3.1	Head Graphs and Canonical Tuples . . . . .	107
4.3.2	Tuple-Specific Equivalence . . . . .	112
4.4	Queries with Variable Heads . . . . .	118
4.4.1	Canonical Freezings and Canonical Tuples . . . . .	119
4.4.2	Avoiding Tuple-Specific Equivalence Tests . . . . .	123
4.5	Soundness and Completeness . . . . .	126
4.5.1	Early Termination . . . . .	126
4.5.2	Constraint Relaxation . . . . .	127
4.6	Relevant Literature . . . . .	127
4.6.1	Algebraic Optimization of Outer Joins . . . . .	128
4.6.2	Queries with Negation . . . . .	129
4.6.3	Query Containment with Null Values . . . . .	130
4.6.4	Query Containment under Dependencies . . . . .	130
4.6.5	Dependencies over Relations with Nulls . . . . .	132

<b>5</b>	<b>Encoding-Equivalence of Hierarchical CEQ's</b>	<b>135</b>
5.1	A Normal Form for HCEQs . . . . .	136
5.2	Encoding-Equivalence of Derived DACEQs . . . . .	142
5.3	Food for Thought . . . . .	146
<b>6</b>	<b>Conclusions</b>	<b>151</b>
6.1	Summary of Results . . . . .	151
6.2	Comprehensive Examples . . . . .	152
6.3	Future Work . . . . .	159
	<b>Appendices</b>	<b>162</b>
<b>A</b>	<b>Proof of Theorem 3.2.4</b>	<b>163</b>
A.1	Pseudo-Trivial Certificate Nodes . . . . .	164
A.1.1	Trivial Nodes . . . . .	165
A.1.2	Pseudo-Trivial Nodes . . . . .	165
A.2	Bag Equivalence: A Counting Argument . . . . .	166
A.2.1	Colour-Diversification . . . . .	168
A.2.2	$k$ -Distinguishing Diversification . . . . .	171
A.2.3	Guaranteeing Pseudo-Trivial Certificate Nodes . . . . .	172
A.3	Set Equivalence: A Symmetry Argument . . . . .	174
A.3.1	Labels, Prefixes and Canonical Database $\mathbb{D}_Q$ . . . . .	177
A.3.2	Canonical Objects . . . . .	180
A.3.3	Pseudo-Trivial Certificate Nodes . . . . .	181
A.4	Normalized Bag Equivalence: Symmetry + Counting . . . . .	183
A.4.1	Labels, Prefixes, and Canonical Database $\mathbb{D}_Q$ . . . . .	185
A.4.2	Canonical Objects . . . . .	187
A.4.3	Pseudo-Trivial Certificate Nodes . . . . .	187
A.5	Arbitrary Signatures . . . . .	191
	<b>References</b>	<b>192</b>



# List of Figures

1.1	A simple relational database schema for tracking company sales . . .	3
1.2	Quarterly average order value, per agent . . . . .	4
1.3	Maximal residential customer order counts, per quarter . . . . .	6
2.1	Graphical representation of sorts . . . . .	14
2.2	Completion of objects $o_7$ , $o_8$ , and $o_9$ . . . . .	18
2.3	Linearization of sorts $\tau_7$ , $\tau_8$ , and $\tau_9$ . . . . .	21
2.4	Linearization of objects $o_7$ , $o_8$ , and $o_9$ . . . . .	22
2.5	COCQL <sup><math>\mathcal{F}</math></sup> algebra operators . . . . .	27
2.6	Encoding of complete chain object $o_4$ . . . . .	34
2.7	Encoding relation $R_3$ and select sub-relations . . . . .	36
2.8	Encoding relations $R_4$ , $R_5$ , and $R_6$ . . . . .	39
2.9	Encoding relation $R_7$ . . . . .	41
2.10	ns-certificate proving $R_1 \doteq_{\text{ns}} R_3$ . . . . .	45
2.11	Hierarchical conjunctive query $Q_{17}$ . . . . .	48
2.12	The six trees in $\text{trees}(Q_{17})$ . . . . .	50
2.13	Embeddings of $\text{trees}(Q_{17})$ into database $\mathbb{D}_2$ . . . . .	51
2.14	Terminal embeddings and their partial applications to $\overline{\mathcal{W}}$ . . . . .	51
2.15	COCQL <sup><math>\{f_b, f_s, f_n\}</math></sup> queries $Q_{18}$ and $Q_{19}$ . . . . .	53
2.16	Output sort $\tau$ and its linearization $\text{CHAIN}(\tau) = (\text{snsb}, 5)$ . . . . .	54
2.17	Query result $o_{12} = (Q_{18})^{\mathbb{D}_3} = (Q_{19})^{\mathbb{D}_3}$ and its linearization w.r.t. $\tau$ . .	55
2.18	Encoding query $Q_{20} = \text{ENCODE}(Q_{18})$ . . . . .	57
2.19	Encoding query $Q_{21} = \text{ENCODE}(Q_{19})$ . . . . .	59
2.20	Encoding query result $R_8 = (Q_{20})^{\mathbb{D}_3}$ satisfying $\text{DECODE}(R_8, \text{snsb}) =$ $o_{13}$ . . . . .	61

2.21	Encoding query result $R_9 = (Q_{21})^{\mathbb{D}_3}$ satisfying $\text{DECODE}(R_9, \text{snsb}) = o_{13}$ . . . . .	62
2.22	Query result $o_{13} = \text{DECODE}((Q_{20})^{\mathbb{D}_3}, \text{snsb}) = \text{DECODE}((Q_{21})^{\mathbb{D}_3}, \text{snsb})$	63
2.23	COQL queries $Q_{22}$ , $Q_{23}$ , and $Q_{24}$ . . . . .	70
2.24	Database instance $\mathbb{D}_4$ . . . . .	71
2.25	Evaluating $Q'_{22}$ , $Q'_{23}$ , and $Q'_{24}$ over $\mathbb{D}_4$ . . . . .	71
3.1	Database $\mathbb{D}_5 = \{\text{Prof}, \text{Student}\}$ . . . . .	76
3.2	Query result $R_{10} = (Q_{25})^{\mathbb{D}_5} = (Q_{26})^{\mathbb{D}_5}$ . . . . .	77
3.3	Two hypergraphs . . . . .	79
3.4	Query results $R_{11} = (Q_{29})^{\mathbb{D}_5}$ , $R_{12} = (Q_{31})^{\mathbb{D}_5}$ , and $R_{13} = (Q_{30})^{\mathbb{D}_5} = (Q_{32})^{\mathbb{D}_5}$ . . . . .	82
4.1	SPCL query $Q_{33}$ equivalent to HCQ $Q_{17}$ . . . . .	93
4.2	SPCL query results $R_{14} = (Q_{33})^{\mathbb{D}_2}$ and $R_{15} = (Q_{34})^{\mathbb{D}_2}$ . . . . .	94
4.3	Directed acyclic conjunctive query $Q_{35}$ . . . . .	95
4.4	Terminal embeddings of $Q_{35}$ into $\mathbb{D}_2$ . . . . .	95
4.5	DACQ $Q_{36}$ equivalent to SPCL query $Q_{34}$ . . . . .	96
4.6	Query $Q_{37}$ and its terminal embeddings into $\mathbb{D}_6$ . . . . .	98
4.7	HCQ $Q_{40}$ prior to conversion to canonical form . . . . .	106
4.8	Query $Q_{41} = \text{CANONICAL}(Q_{40})$ . . . . .	106
4.9	A constant-headed DACQ with its head graph and canonical tuples	108
4.10	Four CH-DACQs with dual-node head graphs . . . . .	114
4.11	Three DACQs with their canonical tuples . . . . .	121
5.1	Encoding queries $Q_{55} = \text{CANONICAL}(Q_{20})$ and $Q_{56} = \text{CANONICAL}(Q_{21})$ 137	
5.2	The <b>snsb</b> -normal forms of $Q_{55}$ and $Q_{56}$ and their results over $\mathbb{D}_3$ . .	143
5.3	Encoding queries $Q_{59}$ , $Q_{60}$ , and $Q_{61}$ evaluated over database $\mathbb{D}_9$ . . .	148
5.4	HCQs $Q_{62}$ and $Q_{63}$ . . . . .	148
6.1	$Q_{68}$ , the $\text{COCQL}^{\{f_b, f_s, f_n\}}$ translation of SQL query $Q_1$ . . . . .	154
6.2	$Q_{69}$ , the $\text{COCQL}^{\{f_b, f_s, f_n\}}$ translation of SQL query $Q_3$ . . . . .	155
6.3	Output sort $\tau_{14}$ of $Q_{68}$ and $Q_{69}$ , and its transformation to a chain .	155
6.4	Queries $Q_{70} := \text{ENCODE}(Q_1)$ , $Q_{71} := \text{ENCODE}(Q_3)$ , and $Q_{72} := \text{chase}_\Sigma(Q_{70})$ 156	

6.5	$Q_{73}$ and $Q_{74}$ , the COCQL $\{f_b, f_s, f_n\}$ translations of SQL queries $Q_4$ and $Q_6$	158
6.6	Encoding queries derived from $Q_{73}$ and $Q_{74}$	158
A.1	Databases $\mathbb{D}_{10}$ and $\mathbb{D}_{11}$ for Example 65	167
A.2	Query results $R_{21} = (Q_{78})^{\mathbb{D}_{10}}$ , $R_{22} = (Q_{79})^{\mathbb{D}_{10}}$ , and $R_{23} = (Q_{80})^{\mathbb{D}_{10}}$	167
A.3	Query results $R_{24} = (Q_{78})^{\mathbb{D}_{11}}$ , $R_{25} = (Q_{79})^{\mathbb{D}_{11}}$ , and $R_{26} = (Q_{80})^{\mathbb{D}_{11}}$	167
A.4	Database $\Delta^{[2,1,3,2]}(\mathbb{D}_{10})$	169
A.5	Database $\mathbb{D}_{12} = \{\text{Wrote}\}$ , query result $R_{27} = (Q_{81})^{\mathbb{D}_{12}}$ , and object $o_{16} = \text{DECODE}(R_{27}, \mathbf{ss})$	175
A.6	Database $\mathbb{D}_{14} = \Delta^{[2,3,2]}(\mathbb{D}_{13})$ , query results $R_{28} = (Q_{86})^{\mathbb{D}_{14}}$ and $R_{29} = (Q_{87})^{\mathbb{D}_{14}}$ , and object $o_{17} = \text{DECODE}(R_{28}, \mathbf{n}) = \text{DECODE}(R_{29}, \mathbf{n})$	184
A.7	Database $\mathbb{D}_{16} = \Delta^{[2,3,2,1,1,1]}(\mathbb{D}_{15})$ , query results $R_{30} = (Q_{86})^{\mathbb{D}_{16}}$ and $R_{31} = (Q_{87})^{\mathbb{D}_{16}}$ , and objects $o_{18} = \text{DECODE}(R_{30}, \mathbf{n})$ and $o_{19} = \text{DECODE}(R_{31}, \mathbf{n})$	184





# List of Symbols

*Page numbers indicate the location of definitions.*

<b>A</b>	
aggregation functions	$\mathcal{F}$ (set of aggregation functions) ..... 26
	$\mathcal{F}$ (set of aggregation functions) ..... 28
	$f_b$ (bag constructor) ..... 28
	$f_n$ (normalized bag constructor) ..... 28
	$f_s$ (set constructor) ..... 28
algebraic operators	$\times$ (Cartesian product) ..... 27
	$\bowtie_p$ (inner join) ..... 27
	$\overrightarrow{\bowtie}_p$ (left outer join (LOJ)) ..... 92
	$\oplus$ (minimal union) ..... 92
	$\uplus$ (outer union) ..... 92
	$\Pi_{\overline{A}}$ (projection, duplicate-eliminating) 27
	$\Pi_{\overline{A}}^{dup}$ (projection, duplicate-preserving) 27
	$\Pi_{\overline{A}}^{X=f(Y)}$ (projection, generalized) ..... 27
	$\downarrow$ (removal of subsumed tuples) ..... 92
	$\lambda_{X=f(T)}$ (scalar subquery) ..... 27
	$\sigma_p$ (selection) ..... 27
$\text{atoms}_v^+$ (atoms within $v$ or its ancestors) .. 47	
$\text{atoms}_v^{\text{anc}}$ (atoms within ancestors of $v$ ) .... 47	
$\text{atoms}_Q$ (atoms within query $Q$ ) ..... 46	
$\text{atoms}_T$ (atoms within tree $T$ ) ..... 49	
$\text{atoms}_v$ (atoms within vertex $v$ ) ..... 47	
<b>B</b>	
$\text{body}_v^+$ (body of CQ formed from $v$ and its ancestors) ..... 47	
$\text{body}_T^+$ (body of CQ formed from all vertices in $T$ ) ..... 49	
$\text{body}_G^+$ (body of CQ formed from all vertices of $G$ ) ..... 95	
$\text{body}_v^{\text{anc}}$ (body of CQ formed from ancestors of $v$ ) ..... 47	
$\text{body}_v$ (body of CQ within vertex $v$ ) ..... 47	
$\text{body}_Q$ (body of query $Q$ ) ..... 46	
<b>C</b>	
$\mathbb{D}_Q$ (canonical database for query $Q$ ) 173, 179, 186	
	$\mathbf{CF}_Q$ (canonical freezings of $\text{CQ}^= Q$ ) ..... 120
	$\mathbf{CF}_Q^+$ (valid freezings of $\text{CQ}^= Q$ ) ..... 120
	$o_p$ (canonical object for $p \in \text{LGP}$ ) .... 180, 187
	$\mathbf{CT}_Q$ (canonical tuples of $\text{DACQ } Q$ ) .. 111, 120
	$\overline{\S}$ -certificate ..... 42–44
	$n_{R,R'}^b$ (bag node) ..... 43
	$n_{R,R'}^n$ (normalized bag node) ..... 43
	$n_{R,R'}^s$ (set node) ..... 42
	$n_{R,R'}^t$ (tuple node) ..... 44
	$(\overline{\S}, k)$ (chain sort) ..... 15
collection sorts	$\{\cdot\}$ (bag) ..... 12
	$\{\!\!\! \{ \cdot \}\!\!\!$ (normalized bag) ..... 12
	$\{\cdot\}$ (set) ..... 12
	$\text{CF-LGC}_i$ (conflict-free subset of $\text{LGC}_i$ ) . 177, 185
	$\text{CF-LGP}$ (conflict-free subset of $\text{LGP}$ ) 178, 185
	$\text{CF-LGP}_m$ (conflict-free subset of $\text{LGP}_m$ ) 178, 185
	$\text{CF-LGS}$ (conflict-free subset of $\text{LGS}$ ) 178, 185
	$\text{CEQ}$ (conjunctive encoding query) ..... 50
	$\text{COCQL}^{\mathcal{F}}$ (conjunctive object-constructing query language) ..... 26
	$\text{COCQL}^{\mathcal{F}-}$ ( $\text{COCQL}^{\mathcal{F}}$ without scalar aggregation) 30
	$\text{CQ}$ (conjunctive query) ..... 46
	$\text{CQ}^=$ (conjunctive query with equality) .... 46
	$\text{CSQL}^{\mathcal{F}-}$ ( $\text{CSQL}^{\mathcal{F}}$ without scalar aggregation) 31
	$\text{CSQL}^{\mathcal{F}}$ (conjunctive $\text{SQL}^{\mathcal{F}}$ ) ..... 31
	$\text{CH-DACQ}$ (constant-headed $\text{DACQ}$ ) ..... 107
	$\mathcal{C}_Q$ (constants occurring in query $Q$ ) ... 46, 47
	$\mathcal{C}$ (constants, infinite set of) ..... 46, 47
	$\mathcal{C}_S$ (constants, infinite set of symbolic) ... 119
<b>D</b>	
	$\text{dags}(G)$ (DAGs form-able by deletion from $G$ , set of) ..... 95
	$\text{dags}(\text{body}_Q, \mathfrak{ts})$ (DAGs in $\text{dags}(\text{body}_Q)$ that output tuples of shape $\mathfrak{ts}$ ) ..... 109
	$\mathbb{D}$ (database) ..... 25

$\mathbb{S}(\mathbb{D})$ (database schema) .....	25	$\mathbf{E}(T)$ (environment, query evaluation) .....	26
DACEQ (directed acyclic conjunctive encoding query) .....	137	$\equiv$ (equivalence) .....	33
DACCQ (directed acyclic conjunctive query) .....	94	$\stackrel{t}{\equiv}$ (equivalence specific to tuple $t$ ) .....	112
$\Delta^{\bar{r}}(\mathbb{D})$ ( $\bar{r}$ -diversification of database $\mathbb{D}$ ) ..	169	<b>F</b>	
$\Delta^{\bar{r}}(t)$ ( $\bar{r}$ -diversification of tuple $t$ ) .....	168	$Q _t^{\dagger}$ (freezing of DACCQ $Q$ with tuple $t$ ) ...	122
$\bar{r}$ (diversification point for $\widehat{\text{dom}}$ ) .....	168	<b>G</b>	
domains		$G _{U'}$ (graph induced from graph $G$ by vertex set $U'$ ) .....	78
$\text{adom}(A_i, T)$ (active domain of attribute $A_i$ in table $T$ ) .....	25	<b>H</b>	
$\text{adom}(\mathbb{D})$ (active domain of database $\mathbb{D}$ ) ..	25	$W_i^{\square}$ (head attribute labelled by vertex $v$ ) ..	48
$\text{dom}$ (atomic constants, infinite set of) ..	12	$\mathcal{W}$ (head attributes of a query, set of) ..	46, 48
$\widehat{\text{dom}}$ (atomic constants, totally-ordered finite set of) .....	168	$\overline{\mathcal{W}}$ (head attributes of a query, tuple of) ...	46
$\text{aname}$ (attribute names, infinite set of) ..	24	$\theta$ (head graph isomorphism) .....	108
$\text{dom}^{\infty}$ (complex objects, infinite set of) ..	16	$G_Q^{\text{head}}$ (head graph of DACCQ $Q$ ) .....	108
$\text{dom}(A_i)$ (domain of attribute $A_i$ ) .....	25	$L$ (head labels of a query, set of) .....	48
$\text{dom}^i$ (finite objects with sort depth $\leq i$ , infinite set of) .....	16	$\overline{L}$ (head labels of a query, tuple of) .....	48
$\widehat{\text{dom}}_i$ (isomorphic “painting” of $\widehat{\text{dom}}$ ) ..	168	$\overline{\mathcal{W}} _{\mathfrak{ts}}$ (head tuple of shape- $\mathfrak{ts}$ query $\hat{Q} _{\mathfrak{ts}}$ ) ...	109
$\text{aname}$ (table names, infinite set of) ...	24	HCEQ (hierarchical conjunctive encoding query) .....	50
<b>E</b>		HCQ (hierarchical conjunctive query) .....	47
embedded dependencies		$h: Q' \rightarrow Q$ (homomorphism) .....	74, 84
$C^{\xi}$ (conclusion query for $\xi$ ) .....	113, 118	$H _{U'}$ (hypergraph induced from hypergraph $H$ by vertex set $U'$ ) .....	78
$\xi^{(G,t)}$ (dependency for $t \in \text{CT}_Q$ and $G \in \text{dags}(\text{body}_Q, \mathfrak{ts}_i^Q)$ ) .....	113–115, 118	<b>I</b>	
$\Sigma^{(Q,t)}$ (dependency set for canonical tuple $t \in \text{CT}_Q$ ) .....	113–115, 118	$\hat{\mathcal{I}}_i$ (index attributes, set of labelled) .....	52
$\xi$ (embedded dependency) .....	113, 118	$\mathcal{I}_i$ (index attributes, set of unlabelled) ..	35, 52
$\equiv_{\text{fin}}$ (finite equivalence) .....	113	$\mathcal{I}_i^{\bar{\mathfrak{s}}}$ ( <i>core</i> index attributes, set of unlabelled) ..	79, 138
$\models_{\text{fin}}$ (finite implication) .....	113	$\overline{\mathcal{I}}_i$ (index attributes, tuple of) .....	34, 50
$P^{\xi}$ (premise query for $\xi$ ) .....	113, 118	$\overline{\mathcal{I}}_i^{\bar{\mathfrak{s}}}$ ( <i>core</i> index attributes, tuple of) ...	80, 138
$\Sigma$ (set of dependencies) .....	113	$[[\tau]]$ (interpretation of sort $\tau$ ) .....	15
$\equiv_{\text{unr}}$ (unrestricted equivalence) .....	113	<b>L</b>	
$\models_{\text{unr}}$ (unrestricted implication) .....	113	$LGC_i$ (label-generating components at level $i$ ) ..	177, 185
embeddings		$\beta_p$ (label-generating function for $p \in LGP$ ) ..	179, 186
$\gamma: G \rightarrow \mathbb{D}$ (DAG embedding) .....	95	$\beta_s$ (label-generating function for $s \in LGS$ ) ..	179, 185
$\Gamma[\bar{a}]$ (embeddings in $\Gamma$ indexed by $\bar{a}$ ) ..	163	$LGP$ (label-generating prefixes) .....	178, 185
$\Gamma[\bar{a}]_o$ (embeddings in $\Gamma[\bar{a}]$ encoding sub-object $o$ ) .....	172	$LGP_m$ (label-generating prefixes of length $m$ ) ..	178, 185
$\Phi$ (embeddings yielding relation $(Q')^{\mathbb{D}}$ ) ..	163	$LGS$ (label-generating sequences) .....	178, 185
$\Gamma$ (embeddings yielding relation $(Q)^{\mathbb{D}}$ ) ..	163	$\mathcal{B}^l$ (labelled set of constants isomorphic to $\mathcal{B}$ ) ..	179, 185
$\hat{\gamma}(W_i)$ (partial application of $\gamma$ to $W_i$ ) ..	49	$\alpha^l$ (labelling function) .....	179, 185
$\gamma: Q \rightarrow \mathbb{D}$ (query embedding) .....	47	$\alpha^{-1}$ (labelling function inverse) .....	179, 185
$\gamma: \text{body}_Q \rightarrow \mathbb{D}$ (query embedding) .....	47	$\mathcal{L}$ (local variables of a query) .....	46, 47
$\gamma: G \rightsquigarrow \mathbb{D}$ (terminal DAG embedding) ..	95	$\mathcal{L}_T$ (local variables within tree $T$ ) .....	49
$\gamma: Q \rightsquigarrow \mathbb{D}$ (terminal query embedding) ..	49	$\mathcal{L}_v$ (local variables within vertex $v$ ) .....	47
$\gamma: T \rightsquigarrow \mathbb{D}$ (terminal tree embedding) ..	49		
$\gamma: T \rightarrow \mathbb{D}$ (tree embedding) .....	49		
$\stackrel{\bar{\mathfrak{s}}}{\equiv}$ (encoding-equality w.r.t. $\bar{\mathfrak{s}}$ ) .....	42		
$\stackrel{\bar{\mathfrak{s}}}{\equiv}$ (encoding-equivalence w.r.t. $\bar{\mathfrak{s}}$ ) .....	64		

<b>M</b>	
$G_{\mathfrak{s}}^{\min}$ (minimal DAG in $\mathbf{dags}(\mathbf{body}_Q, \mathfrak{s})$ ) . . . . .	109
$\rightarrow$ (multivalued dependency (MVD)) . . . . .	75
<b>N</b>	
$\bar{\mathbf{a}}$ (null constant) . . . . .	38
$\bar{\rightarrow}$ (nullability dependency) . . . . .	98
$\mathcal{S}_Q$ (nullability set of DACQ $Q$ ) . . . . .	102
$\#(t, \mathbf{c}_i)$ (number of occurrences of $\mathbf{c}_i$ within $t$ )	169
<b>O</b>	
$\hat{\mathcal{V}}$ (output attributes, set of labelled) . . . . .	50
$\mathcal{V}$ (output attributes, set of unlabelled) . . . . .	35, 50
$\bar{\mathcal{V}}$ (output attributes, tuple of) . . . . .	35, 50
<b>P</b>	
$\delta_i$ (painting function) . . . . .	168
$\delta^{-1}$ (painting function inverse (whitewash)) . . . . .	168
$\mathcal{P}_v^{\text{anc}}$ (parameter variables (inter-query) from root to vertex $v$ , exclusive) . . . . .	48
$\mathcal{P}_v^+$ (parameter variables (inter-query) from root to vertex $v$ , inclusive) . . . . .	48
$\mathcal{P}_Q^+$ (parameter variables (inter-query) within $\mathbf{body}_Q$ ) . . . . .	48
$\mathcal{P}_v$ (parameter variables (inter-query) within vertex $v$ ) . . . . .	47
$\mathcal{P}_{\mathfrak{s}}$ (parameter variables for shape $\mathfrak{s}$ ) . . . . .	110
$\mathcal{P}$ (parameter variables of a query) . . . . .	46, 47
$\Pi_2^p$ (polynomial hierarchy, second tier) . . . . .	87
$\Sigma_2^p$ (polynomial hierarchy, second tier) . . . . .	87
$\text{pred}_v^+$ (predicates within $v$ or its ancestors) . . . . .	47
$\text{pred}_v^{\text{anc}}$ (predicates within ancestors of $v$ ) . . . . .	47
$\text{pred}_Q$ (predicates within query $Q$ ) . . . . .	46
$\text{pred}_T$ (predicates within tree $T$ ) . . . . .	49
$\text{pred}_v$ (predicates within vertex $v$ ) . . . . .	47
<b>Q</b>	
$H^Q$ (query hypergraph) . . . . .	78
$Q^{\mathbb{D}}$ (query result over database $\mathbb{D}$ ) . . . . .	47
<b>S</b>	
$\S$ (semantic indicator $\mathfrak{s}, \mathfrak{z}$ or $\mathfrak{n}$ ) . . . . .	15
<b>T</b>	
$\hat{Q}$ (shape query for DACQ $Q$ ) . . . . .	109
$Q _{\mathfrak{s}}$ (shape- $\mathfrak{s}$ query for DACQ $Q$ ) . . . . .	109
$\bar{\S}$ (signature) . . . . .	15
$\leq_d$ (simulation to depth $d$ ) . . . . .	69
$\tau$ (sort) . . . . .	13
$\text{sort}(A_i)$ (sort of attribute $A_i$ ) . . . . .	25
$\mathcal{T}$ (sorts, set of all) . . . . .	13
$\mathcal{T}^i$ (sorts, set of all with depth $\leq i$ ) . . . . .	13
$\ll_d$ (strong simulation to depth $d$ ) . . . . .	70
SQL (structured query language) . . . . .	30
SQL $^{\mathcal{F}}$ (SQL with aggregation functions $\mathcal{F}$ ) . . . . .	31
$R[\bar{\mathbf{a}}]$ (sub-relation indexed by $\bar{\mathbf{a}}$ ) . . . . .	35
symbolic constants	
$\bar{\subseteq}$ ( $\mathcal{C}_S$ -contained (tuple sets)) . . . . .	119
$\bar{\approx}$ ( $\mathcal{C}_S$ -equivalent (tuple sets)) . . . . .	119
$\bar{\succ}$ ( $\mathcal{C}_S$ -homomorphic (tuples)) . . . . .	119
$\bar{\approx}$ ( $\mathcal{C}_S$ -isomorphic (tuples)) . . . . .	119
$\mathcal{C}_S$ (infinite set of symbolic constants) . . . . .	119
$\succ$ (more general than (tuples)) . . . . .	119
<b>T</b>	
$\$(T), \mathbb{T}(T)$ (table) . . . . .	24
$\mathbb{T}(T)$ (table instance) . . . . .	24
$\$(T)$ (table schema) . . . . .	24
$\mathbf{trees}(T)$ (trees form-able by deletion from $T$ , set of) . . . . .	48
$\langle \cdot \rangle$ (tuple) . . . . .	13
$\mathfrak{s}$ (tuple shape) . . . . .	108
$\mathfrak{s}_t^Q$ (tuple shape in $\mathbb{TS}_Q$ corresponding to tuple $t \in \mathbf{CT}_Q$ ) . . . . .	111, 120
$\mathbb{TS}_Q$ (tuple shapes of DACQ $Q$ , set of) . . . . .	108
<b>V</b>	
$\mathcal{B}$ (variables occurring in a query) . . . . .	46, 47
$\mathcal{B}_v^{\text{anc}}$ (variables within $\mathbf{body}_v^{\text{anc}}$ ) . . . . .	47
$\mathcal{B}_v^+$ (variables within $\mathbf{body}_v^+$ ) . . . . .	47
$\mathcal{B}_v$ (variables within $\mathbf{body}_v$ ) . . . . .	47
$\mathcal{B}_T$ (variables within tree $T$ ) . . . . .	49
<b>W</b>	
$\delta^{-1}$ (whitewash function) . . . . .	168



# Chapter 1

## Introduction

In this thesis we consider the problem of deciding when a pair of queries posed to a database system are *equivalent*—that is, guaranteed to return identical results. The query equivalence problem (and the closely related query *containment* problem) is a fundamental problem that appears within a variety of contexts in the design of database management systems.

The most obvious application of query equivalence is to the “global” query optimization problem: given a query and a cost model, find within the space of equivalent query formulations the query with minimal cost [6, Ch. 6]. For the class of simple conjunctive queries, query equivalence is well understood and global optimization techniques are widely implemented [16, 19]; however, query equivalence (and hence global optimization) is in general undecidable for arbitrary relational algebra expressions [6, Ch. 6]. For this reason, most database systems rely upon “local” optimization strategies for complex queries—i.e., the generation of alternative query formulations via localized (often algebraic) transformation rules, often called “query rewriting” or “logical rewriting” [17]. Within this thesis we study equivalence for queries defined by conjunctive relational algebra expressions extended with an aggregation operator *that may be arbitrarily nested within the algebra expression*. Our work thus contrasts with much of the literature on aggregation queries which assumes a single top-level aggregation operation. Our results are applicable to both forms of query optimization: by characterizing equivalence over more powerful classes of queries, we both enable global optimization techniques for those classes and enable development of more powerful rewriting rules for complex queries.

Beyond its application to query optimization, the query equivalence problem is directly applicable to various problems within information integration. The *view rewriting* problem—rewriting a query to use a set of materialized views—may arise either in the context of cost-based query optimization or in the context of information integration where the views model the underlying data sources [57]. The *data exchange* problem—transforming a database instance conforming to a source schema into an instance conforming to a target schema—is also based upon

query equivalence [33]. Finally, the recent advent of *object-relational mappings* promises convenient integration of relational data into object-oriented programming languages, but translating operations over objects into operations over relations requires reasoning over declarative mappings that makes heavy use of query equivalence tests (in the guise of materialized-view answering and maintenance problems) [83].

## 1.1 Motivating Examples

Our interest in query equivalence for conjunctive queries with nested aggregation stems from the problem of optimizing complex SQL queries. In particular, we are interested in how existing query optimizers might be extended with robust algorithms for cost-based view rewriting in a context where complex materialized view definitions are formed by “stacking” (i.e., joining and re-aggregating) simpler views. This class of so-called “stacked” materialized views is interesting in practice both because it allows efficient incremental maintenance algorithms, and because stacked views are useful for answering the complex aggregation queries that appear within decision-support workloads [29]. Unfortunately, as illustrated in the following examples, view rewriting algorithms implemented within contemporary relational database management systems (RDBMSs) are not robust when rewriting queries over complex view definitions.

**Example 1** Consider the database schema shown in Figure 1.1, which stores information about customer orders solicited by a company’s agents. The database schema contains six base tables with primary and foreign key constraints (primary keys denoted by underlined attributes), along with six views defined by SQL queries,<sup>1</sup> four of which are materialized.

The attribute `Customer.ctype` is code that classifies customers as either *Residential* or *Corporate*, and sales from the two sectors are often reported separately. Suppose that an end user desires a report that lists for each agent the quarterly average order value, with the Residential and Corporate metrics shown in separate columns. Equipped with a reporting tool capable of generating single-block conjunctive SQL queries (with top-level aggregation) only, the user could accomplish the desired report by using the logical view `AgentSales`, generating the query  $Q_1$  shown in Figure 1.2.

---

<sup>1</sup>For conciseness, we have abbreviated the SQL syntax by shortening relation names to their capital letters (using subscripts to distinguish repeated relations), and by using an equi-join operator to denote the join predicates. For example, the following query  $Q$  (in proper SQL syntax) is abbreviated as query  $Q'$ .

<pre> Q: SELECT *   FROM Customer, Order o1, Order o2  WHERE Customer.cid = o1.cid        AND Customer.cid = o2.cid </pre>	<pre> Q': SELECT *      FROM C ⋈<sub>cid</sub> O<sub>1</sub> ⋈<sub>cid</sub> O<sub>2</sub> </pre>
--	---

Abbr	Table Schema	Foreign Keys
C	Customer( <u>cid</u> , cname, ctype)	
O	Order( <u>oid</u> , cid, date)	date:Date(date)
LI	LineItem( <u>oid</u> , <u>lineno</u> , price, qty)	oid:Order(oid)
A	Agent( <u>aid</u> , aname)	
OA	OrderAgent( <u>oid</u> , <u>aid</u> )	oid:Order(oid), aid:Agent(aid)
D	Date( <u>date</u> , qtr)	

Logical (non-Materialized) Views:

```
AS    AgentSales(aid, aname, date, ctype, oid, oval)
      SELECT aid, aname, date, ctype, oid, SUM(price * qty)
      FROM C ⋈cid O ⋈oid LI ⋈oid OA ⋈aid A
      GROUP BY aid, aname, date, ctype, oid
```

```
QC    QuarterlyCustomer(qtr, cid, cname, ctype)
      SELECT qtr, cid, cname, ctype
      FROM C ⋈cid QCO
```

Materialized Views:

```
OV    OrderValues(oid, oval)
      SELECT oid, SUM(price * qty)
      FROM LI
      GROUP BY oid
```

```
QAS   QuarterlyAgentSales(aid, qtr, ctype, avgOval)
      SELECT aid, qtr, ctype, AVG(oval)
      FROM C ⋈cid O ⋈oid OV ⋈oid OA ⋈date D
      GROUP BY aid, qtr, ctype
```

```
QCO   QuarterlyCustomerOrders(qtr, cid, ordCount)
      SELECT qtr, cid, COUNT(*)
      FROM O ⋈date D
      GROUP BY qtr, cid
```

```
QMCO  QuarterlyMaxCustomerOrders(qtr, ctype, maxOrdCount)
      SELECT qtr, ctype, MAX(ordCount)
      FROM C ⋈cid QCO
      GROUP BY qtr, ctype
```

Figure 1.1: A simple relational database schema for tracking company sales

```

Q1:  SELECT AS1.aname, qtr,
        AVG(AS1.oval) AS avgRsale,
        AVG(AS2.oval) AS avgCsale
FROM (AS1 ⋈date D1) ⋈{aid, qtr} (AS2 ⋈date D2)
WHERE AS1.ctype = 'R'
      AND AS2.ctype = 'C'
GROUP BY aid, AS1.aname, qtr

Q2:  SELECT A1.aname, qtr,
        AVG(OV1.oval) AS avgRsale,
        AVG(OV2.oval) AS avgCsale
FROM (C1 ⋈cid O1 ⋈oid OV1 ⋈oid OA1 ⋈aid A1 ⋈date D1) ⋈{aid, qtr}
      (C2 ⋈cid O2 ⋈oid OV2 ⋈oid OA2 ⋈aid A2 ⋈date D2)
WHERE C1.ctype = 'R'
      AND C2.ctype = 'C'
GROUP BY aid, A1.aname, qtr

Q3:  SELECT aname, qtr,
        QAS1.avgOval AS avgRsale,
        QAS2.avgOval AS avgCsale
FROM A ⋈aid QAS1 ⋈{aid, qtr} QAS2
WHERE QAS1.ctype = 'R'
      AND QAS2.ctype = 'C'

```

Figure 1.2: Quarterly average order value, per agent



We implemented query  $Q_1$  on the latest versions of several commercial RDBMSs. The best rewriting of  $Q_1$  found by any system that we tested uses schema information to push down the SUM aggregate within AgentSales in order to rewrite over two occurrences of the materialized view OrderValues—resulting in query  $Q_2$  (cf. Figure 1.2). Observe that both  $Q_1$  and  $Q_2$  perform a Cartesian product between each agent’s quarterly Residential and Corporate orders prior to the aggregation.<sup>2</sup>

In contrast, within query  $Q_3$  (cf. Figure 1.2) the Residential and Corporate orders are aggregated independently. This allows the query to be rewritten over two instances of the materialized view QuarterlyAgentSales, yielding a drastically more efficient execution plan. (Even if view QuarterlyAgentSales were not materialized,  $Q_3$  would still allow much more efficient execution than either  $Q_1$  or  $Q_2$  because it avoids computing the Cartesian product as an intermediate result.) In Examples 62 and 63 (Chapter 6) we will show that  $Q_3$  is equivalent to  $Q_1$  over all database instances that conform to the primary and foreign key constraints in the database schema, but not equivalent over arbitrary instances.

Complex SQL queries containing aggregation are not restricted to decision-support workloads. The popularity of object-relational mapping tools means that the SQL queries within contemporary database workloads are increasingly program-generated, rather than hand-crafted. The trade-off in providing this layer of abstraction to the application programmers is that queries generated via compositions of mappings or logical views can be highly nested, with significant redundancy that is challenging for the optimizer to simplify. We illustrate this with an example.

**Example 2** Consider again the database schema in Figure 1.1. The sales division of the company commonly refers to its “Quarterly Customers,” meaning the set of customers who have placed an order during a given quarter. To support this concept, the database administrator created the materialized view QuarterlyCustomerOrders which totals the number of orders by quarter and customer id. She then created the logical view QuarterlyCustomer for the convenience of users who only need to know the identity of the quarterly customers, not the count of their orders.

Application programmers at the company use an object-oriented query framework (on top of an object-relational mapping) to interface with the database. Defined within the entity schema of this framework are entities corresponding to each table and view within the database schema, along with entities *QuarterlyResidentialCustomer*

---

<sup>2</sup>Some of the RDBMSs were restricted to query optimization over “single-block” materialized-view definitions. In this case, we tested their hypothetical rewriting ability over the “stacked” materialized views in Figure 1.1 by first testing their ability to rewrite over the leaf level views, and then—if successful—we converted the leaf level views to base tables and manually rewrote the queries over them in order to test rewriting at the next level. Some of the RDBMSs did not support the MAX or AVG aggregation functions within materialized-view definitions. Because an average can be computed as a sum divided by a count, for systems that did not support AVG in materialized views we added the appropriate SUM and COUNT expressions to the view and then tested the query both with the original AVG expressions and with manual decompositions into SUM and COUNT.

```

Q4:  SELECT qtr, MAX(countOrders)
      FROM (SELECT qtr, cid, cname,
                  (SELECT COUNT(*)
                   FROM D ⋈date O
                   WHERE D.qtr = QC.qtr
                        AND O.cid = QC.cid
                  ) AS countOrders
            FROM (SELECT qtr, cid, cname
                  FROM QC
                  WHERE ctype = 'R'
                  ) AS QuarterlyResidentialCustomer
            ) AS QuarterlyResCustWithOrderCounts
      GROUP BY qtr

Q5:  SELECT qtr, MAX(countOrders)
      FROM (SELECT qtr, cid, cname,
                  (SELECT COUNT(*)
                   FROM QCO2
                   WHERE QCO2.qtr = QCO1.qtr
                        AND QCO2.cid = QCO1.cid
                  ) AS countOrders
            FROM C ⋈cid QCO1
            WHERE ctype = 'R'
            ) AS QuarterlyResCustWithOrderCounts
      GROUP BY qtr

Q6:  SELECT qtr, maxOrdCount
      FROM QMCO
      WHERE ctype = 'R'

```

Figure 1.3: Maximal residential customer order counts, per quarter

and *QuarterlyCorporateCustomer* that inherit from entity *QuarterlyCustomer*. Assume that an application programmer (who does not know SQL) requires a list of the maximum number of orders placed by quarterly residential customers, grouped by quarter. The programmer is not familiar with either the *QuarterlyCustomerOrders* or *QuarterlyMaxCustomerOrders* entities. Instead, within his object-oriented query framework he first creates an entity *QuarterlyResCustWithOrderCounts*, formed by supplementing *QuarterlyResidentialCustomer* with a computed attribute *countOrders* containing the count of the corresponding order values. Next, he selects the maximum value of *countOrders* across all instances of *QuarterlyResCustWithOrderCounts*, grouped by quarter. Query  $Q_4$  in Figure 1.3 shows the SQL query generated by the programming framework. Observe that derived attribute *countOrders* is rendered as a SELECT-clause sub-query. From the perspective of entity-to-SQL translation this is the correct translation of a computed attribute, because *QuarterlyResidentialCustomer* entities without corresponding orders should yield instances of *QuarterlyResCustWithOrderCounts* entities with *countOrders*=0, and these zeros could potentially affect the calculation of the maximum. The SQL generator within the entity framework does not know that each instance of *QuarterlyResidentialCustomer* is guaranteed to have at least one corresponding order, because that information is encapsulated in the logical view *QuarterlyCustomer* inside the RDBMS.

Consider queries  $Q_5$  and  $Q_6$  in Figure 1.3. Query  $Q_5$  is the best rewriting of  $Q_4$  found by any RDBMS that we tested. None of the RDBMSs were able to rewrite  $Q_4$  over the materialized view *QuarterlyMaxCustomerOrders*; however, in Example 64 (Chapter 6) we will show that queries  $Q_4$  and  $Q_6$  are equivalent.

## 1.2 Related Problems and Technologies

**Materialized Views and Data Warehousing** Halevy has written a thorough survey of research on the use of materialized views to answer queries [57]. He draws an important distinction separating works whose purpose is *query optimization* (including maintenance of physical data independence) from works whose purpose is *data integration*. Data integration typically utilizes rewritings that are *contained* within the original query, whereas query optimization requires rewritings to be *equivalent* to the original query. Because our motivation stems from query optimization, we restrict our attention to query equivalence and do not consider query containment in this thesis. Within the context of query optimization, Halevy distinguishes works treating view rewriting as a *logical optimization* (i.e., as a query language-level transformation) from works treating view rewriting as a *physical optimization* (i.e., during generation of query execution plans). Although this distinction is pertinent to how our results might be integrated into existing RDBMSs (see Section 6.3), it is orthogonal to the issue of guaranteeing equivalence between the original query and the rewritten query. Hence, for the purposes of motivating our study of query equivalence it suffices to treat view rewriting as a purely logical optimization (which is why in Examples 1 and 2 we show the rewritings performed

by RDBMS optimizers as SQL queries, even though the output we obtained from the optimizers were actually physical execution plans).

For the restricted case of conjunctive rewritings of conjunctive queries over conjunctive views, the query equivalence problem is solved [16, 19], as is the view rewriting problem both as a logical optimization [74] and as a physical optimization integrated into contemporary optimizer architectures [18, 44, 45]. From the perspective of query optimization, pre-materialization of views offers the greatest potential savings when the size of the query output is vastly smaller than the size of the raw data that must be examined in order to evaluate the query over a database instance. For this reason, materialized views are especially important within data warehouse environments, due to the heavy use of aggregation within decision-support queries (e.g., see the queries within the TPC-H [105] or TPC-DS [104] benchmark workloads). One restricted form of data warehousing depends upon a *star schema* (i.e., a (very large) central *fact table* connected by foreign key joins to (small) *dimension tables*), and both queries and materialized views are limited to defining portions of a multidimensional *data cube* [48]. Within this context, testing query equivalence is trivial, and view rewriting can be efficiently effected by analyzing the *data cube lattice* [59, 90]. A less restricted form of data warehousing assumes that queries and materialized views are limited to expressions with a single top-level aggregation operator. Within this context, the study of query equivalence has focused on numeric or algebraic properties of specific aggregation functions [25, 51], while the study of view rewriting ranges from formal aspects of modelling properties of aggregation functions [24] and determining the complexity of the rewriting language needed to fully utilize the views [53] to more pragmatic approaches that utilize various algebraic transformations for re-shaping the query to obtain sound (but incomplete) rewriting algorithms [54, 101].

There is only limited previous research that considers view rewriting in contexts where the query and view definitions contain nested aggregation. In one work, a bottom-up matching algorithm is proposed that applies templates based upon (previously known) algebraic transformations [115]. Another work focuses on increasing the robustness and efficiency with which transformational optimizers apply sets of (previously known) transformations to enact view rewriting [29]. All of this work presumes a set of algebraic transformations for re-shaping queries, and there is no consideration of whether the set of transformations is complete for traversing the space of equivalent queries.

**Query Simplification and Logical Rewriting** Work on simplifying complex SQL queries composed of arbitrarily nested queries was pioneered by Kim, who characterizes different forms of nested sub-queries and focuses on heuristics for flattening or de-correlating them to increase the search space of the query optimizer [69]. (Traditional query optimizer architectures only optimize join orders within the individual blocks of a complex query; hence, flattening nested blocks increases the space of possible join orders, while de-correlating nested blocks allows

the use of physical join operators other than nested loops.) Many of Kim’s heuristics are implemented as logical rewriting rules within modern query optimizers, along with various other algebraic transformations that re-shape the nesting blocks of a query by commuting join and aggregation operators [54, 113, 114]. Query simplification via logical rewriting can vastly reduce the execution cost of certain queries; however, the pattern-based rule languages of contemporary optimizers are quite restrictive [46, 56]. Performing the degree of query simplification we illustrate in Example 2 requires a formal analysis of the input expression that cannot be conveniently specified as an algebraic template.

**Complex Objects and Object-Relational Mappings** Extensions to the relational model for handling complex objects (a.k.a. complex values or nested relations) were studied extensively during the 1980s; the reader is referred to Abiteboul et al.’s textbook for an introduction to the topic [6, Ch. 20–21], and to Section 2.4 of this thesis where we include a brief survey of some of the pertinent literature in this area. Database systems based upon native nested-relational and/or object-oriented data models remain niche products, however. Instead, the current mainstream approach to integrating object-oriented programming languages with persistent database storage is to utilize *object-relational mapping (ORM)* tools that present the programmer with an abstraction of the relational schema in terms of abstract entities. A fixed language of query operations over these entities is then translated by the ORM framework into SQL queries over the underlying relations. Popular examples include Hibernate, which provides the object-oriented query language HQL [1]; ActiveRecord, the ORM component of the Ruby on Rails programming framework [3]; and the ADO.NET Entity Framework, the ORM component of Microsoft’s .NET framework which supports its Language Integrated Query (LINQ) facilities [2, 82].

The results in this thesis pertain to ORM technologies in two different ways. First, as mentioned in the previous section, ORM tools are quickly emerging as a source of very complex SQL queries, which motivates the study of query equivalence as a means towards deriving more powerful query simplification techniques for relational query optimizers. Second, the data models that ORM tools expose to the programmer are based upon nested collections, and many of the SQL queries issued by the ORM tools are used to construct in-memory complex objects out of data stored in flat relations. Within our study of nested aggregation we abstract aggregation functions as collection constructors, and so the query language we study constructs complex objects out of databases of flat relations. Hence, it is foreseeable that our results could be applied not only to optimization of SQL queries at the RDBMS level, but also to improve the quality of SQL translations generated by ORM packages.

## 1.3 Road Map

The query language for which we propose to study equivalence is a conjunctive algebra extended with aggregation functions that construct abstract collections which—when nested—form complex objects. We consider objects constructed by arbitrary nesting of three different types of unordered collections—sets, bags, and a collection type we call *normalized bags*—which abstract common aggregation functions by modelling their sensitivity to the tuple cardinalities within their inputs. We have chosen this abstraction to highlight the impact of *nested* aggregation on the query equivalence problem, without getting bogged down in complexities specific to particular aggregation functions.

In Chapter 2 we lay the necessary groundwork for considering the query equivalence problem for queries with nested aggregation. We start by formally defining the data model, algebraic query language, and a procedure for encoding arbitrary complex objects within a single flat relation (i.e., containing no aggregated values). The main contribution of this chapter is a reduction of the query equivalence problem to deciding a relationship we call *encoding equivalence* between conjunctive queries that do not contain any aggregation operators, but may contain an operator that introduces null values (a restricted form of the *left outer-join* operator).

Chapters 3–5 address the problem raised in Chapter 2 of deciding encoding equivalence between conjunctive queries containing null-introducing operators. In Chapter 3 we focus on understanding encoding equivalence (which is a generalization of standard query equivalence) in the restricted case where the null-introducing operator is not present. In Chapter 4 we focus on standard query equivalence in the case where null-introducing operators are present. We combine these orthogonal results in Chapter 5, where we consider encoding equivalence in the presence of null-introducing operators.

In Chapter 6 we summarize our results. We also revisit the motivating examples from Section 1.1 and use them to illustrate both how our theoretical results can be applied in practice and where shortcomings in our theory exist that require further attention.

# Chapter 2

## Equivalence of Conjunctive Queries that Construct Objects

Our primary purpose in this chapter is to lay a theoretic foundation for discussing queries—over relational databases—which return complex objects as their result. In future chapters we will build upon this foundation to reason about the equivalence of such queries. In Section 2.1 we define complex objects and related notation, including an algebraic language for constructing complex objects. In Section 2.2 we define a representation of complex objects as relations of flat tuples, and show how our complex object algebra queries can be mapped to queries in a (slightly extended) relational query language. Finally, relevant literature is discussed in Section 2.4.

### 2.1 Complex Objects

Informally, a complex object is a recursive structure built up from *atomic values*, *tuples* of complex objects, and *collections* of complex objects. In Section 2.1.1 we give a more formal definition of complex objects. In Section 2.1.2 we describe a transformation that we call *linearization* which converts objects to a canonical form. Finally, we introduce a query language for generating complex objects from a database of flat relations in Section 2.1.3.

#### 2.1.1 Objects and Sorts

Our investigation of complex objects is motivated by asking about equivalence of SQL queries containing aggregation functions, and so our choice of which collection types to consider is determined by the semantics of commonly-supported aggregation functions. All SQL aggregation functions ignore the order of elements within their input, and so we restrict our attention to unordered collections. Aggregation functions such as `MIN`, `MAX`, and any function using the `DISTINCT` keyword (e.g.



COUNT(DISTINCT *column-name*)) also ignore the element cardinalities within their input; thus, we model their input as a *set*. Aggregate functions such as COUNT and SUM are sensitive to the element cardinalities and so we model their input as a *bag* (a.k.a. *multi-set*). Finally, many statistical functions such as AVERAGE, MEDIAN, or STDEV are sensitive only to *relative* element cardinalities (e.g., doubling the cardinality of each value does not change the aggregated result), and so we model their input using a collection type we call a *normalized bag*.

Sets and bags are well-known collection types and we denote them using the delimiters  $\{\cdot\}$  and  $\{\!\!\{ \cdot \}\!\!\}$ , respectively. Given an arbitrary finite collection  $c$ , the *bag projection of  $c$* —denoted  $\text{Bag}(c)$ —is the bag of elements that occur in  $c$ , while the *set projection of  $c$* — $\text{Set}(c)$ —is the set of elements in  $c$ . A normalized bag is a special case of a bag in which the greatest common divisor of the element frequencies is one; we denote a normalized bag with the delimiters  $\{\!\!\{ \cdot \}\!\!\}$ . Given an arbitrary collection  $c$ , the *normalized bag projection of  $c$* — $\text{NormBag}(c)$ —is the normalized bag formed by dividing each element frequency by the GCD of all element frequencies.

**Example 3** The following four bags are all distinct.

$$\begin{aligned} B_1 &:= \{x, x, y, y, y\} & B_3 &:= \{\!\!\{x, x, y\}\!\!\} \\ B_2 &:= \{x, x, x, x, y, y, y, y, y\} & B_4 &:= \{\!\!\{x, x, x, x, y, y\}\!\!\} \end{aligned}$$

The first and second have the same normalized bag projection

$$N_1 := \{\!\!\{x, x, y, y, y\}\!\!\} = \text{NormBag}(B_1) = \text{NormBag}(B_2)$$

as do the third and fourth.

$$N_2 := \{\!\!\{x, x, y\}\!\!\} = \text{NormBag}(B_3) = \text{NormBag}(B_4)$$

Finally, the four bags and two normalized bags all have the same set projection.

$$\text{Set}(B_1) = \text{Set}(B_2) = \text{Set}(B_3) = \text{Set}(B_4) = \text{Set}(N_1) = \text{Set}(N_2) = \{x, y\}$$

We now formalize our definition of complex objects. Definitions 2.1.1 and 2.1.4 below borrow heavily from Abiteboul et al.’s definition of *complex value databases* [6, Ch. 20], the primary difference being our extension to multiple collection types. Another minor difference is that we do not assign attribute names to elements of tuples; the reason for this will be explained when we present our query algebra.

## Sorts

We start by fixing a countably infinite set of atomic values  $\text{dom}$ , called the underlying domain. We desire complex objects to be typed, and so we next define a set of complex *sorts*, which are meta-entities that each define a class of objects of a restricted structure.



**Definition 2.1.1 (Sort)** A sort (or type)  $\tau$  is syntactically defined by the following recursive grammar:

$$\tau := \text{dom} \mid \{ \tau \} \mid \{ \{ \tau \} \} \mid \{ \{ \{ \tau \} \} \} \mid \langle \tau, \dots, \tau \rangle$$

We call sort  $\text{dom}$  an *atomic* sort; sorts  $\{ \tau \}$ ,  $\{ \{ \tau \} \}$ , and  $\{ \{ \{ \tau \} \} \}$  *collection* sorts; and sort  $\langle \tau, \dots, \tau \rangle$  a *tuple* sort (which may have any finite arity  $k \geq 0$ ). We call a tuple sort *flat* if it is composed of only atomic sorts. We use  $\mathcal{T}$  to denote the set of all sorts.

**Definition 2.1.2 (Sort Depth)** The depth of a sort  $\tau$  is defined recursively as follows.

$$\begin{aligned} \text{depth}(\text{dom}) &= 0 \\ \text{depth}(\{ \tau \}) &= 1 + \text{depth}(\tau) \\ \text{depth}(\{ \{ \tau \} \}) &= 1 + \text{depth}(\tau) \\ \text{depth}(\{ \{ \{ \tau \} \} \}) &= 1 + \text{depth}(\tau) \\ \text{depth}(\langle \tau_1, \dots, \tau_k \rangle) &= \max(0, \text{depth}(\tau_1), \dots, \text{depth}(\tau_k)) \end{aligned}$$

*Intuitively, the depth of  $\tau$  is the maximum number of collection sorts occurring along any root-to-leaf path in the hierarchical definition of  $\tau$ .*

Although not enforced by the grammar in Definition 2.1.1, without loss of generality we adopt a convention that tuple sorts are always structured with the child sorts of zero depth occurring to the left of child sorts with non-zero depth.

**Example 4** Consider the following sorts.

$$\begin{aligned} \tau_1 &:= \langle \text{dom}, \text{dom}, \text{dom} \rangle & \tau_2 &:= \{ \langle \text{dom}, \text{dom} \rangle \} \\ \tau_3 &:= \{ \langle \langle \text{dom} \rangle, \langle \text{dom}, \text{dom} \rangle \rangle \} & \tau_4 &:= \{ \{ \langle \text{dom} \rangle \} \} \\ \tau_5 &:= \{ \{ \langle \langle \text{dom} \rangle \rangle \} \} & \tau_6 &:= \{ \langle \text{dom}, \{ \langle \text{dom} \rangle \} \rangle \} \\ \tau_7 &:= \{ \langle \langle \text{dom}, \{ \langle \langle \text{dom} \rangle \rangle \rangle \rangle \} \} & \tau_8 &:= \{ \langle \{ \langle \langle \langle \langle \text{dom} \rangle \rangle \rangle \rangle \rangle \}, \{ \langle \langle \text{dom} \rangle \rangle \} \rangle \} \\ \tau_9 &:= \{ \langle \langle \{ \langle \{ \langle \langle \text{dom} \rangle \rangle \} \rangle \rangle \rangle \}, \{ \langle \{ \langle \langle \text{dom} \rangle \rangle \} \rangle \}, \{ \langle \langle \text{dom} \rangle \rangle \} \rangle \} \end{aligned}$$

An equivalent tree representation is shown in Figure 2.1. Sort  $\tau_1$  has depth zero,  $\tau_2$  and  $\tau_3$  have depth one,  $\tau_4$ – $\tau_7$  have depth two, and  $\tau_8$  and  $\tau_9$  have depth three.

For each integer  $i \geq 0$ , we use  $\mathcal{T}^i \subset \mathcal{T}$  to denote the (infinite) set of all sorts with depth at most  $i$ .

$$\mathcal{T}^i := \{ \tau \mid \text{depth}(\tau) \leq i \} \quad (2.1)$$

**Branching, Non-Branching and Chain Sorts** We classify sorts as either *branching* or *non-branching* as follows.

- Atomic sort  $\text{dom}$  is non-branching.

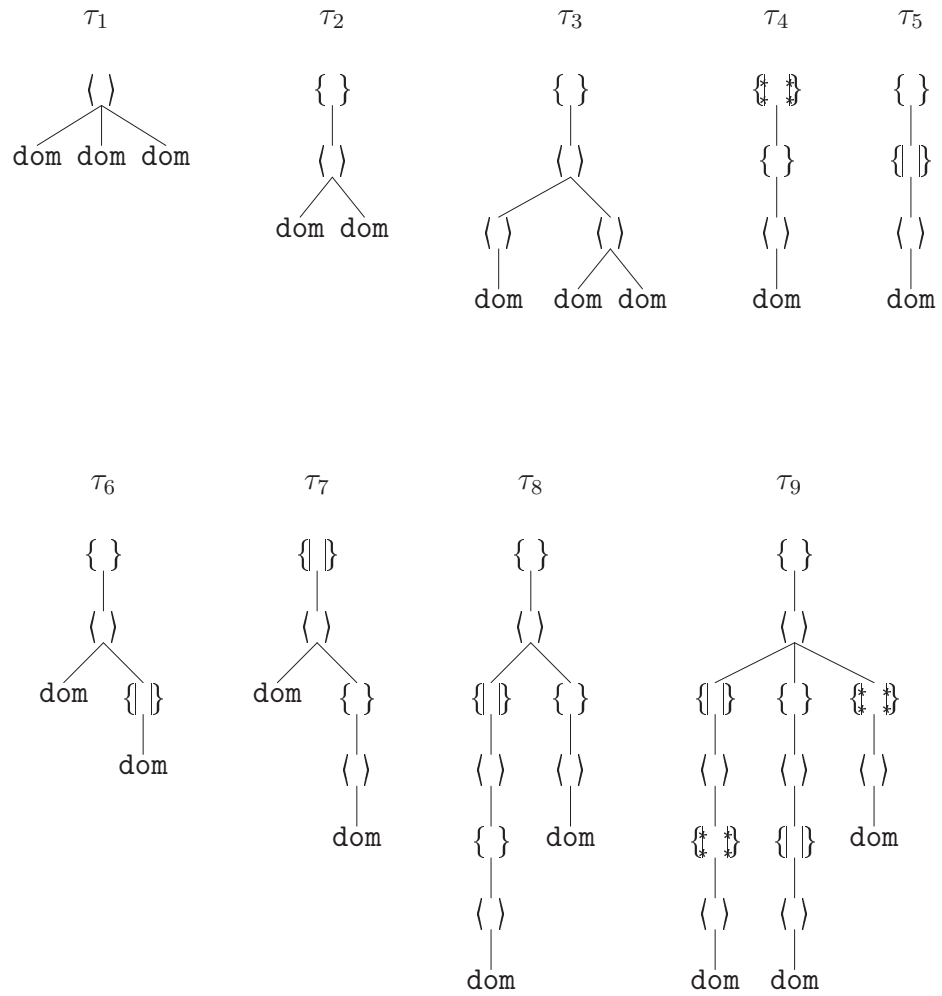


Figure 2.1: Graphical representation of sorts

- Collection sorts  $\{\tau\}$ ,  $\{\!\{\tau\}\!\}$ , and  $\{\!\!\{\tau\}\!\!\}$  are branching if  $\tau$  is a branching sort; otherwise they are non-branching.
- Tuple sort  $\langle \tau_1, \dots, \tau_k \rangle$  is branching if either
  1. one of its child sorts  $\tau_i$  is branching; or
  2. at least two of its child sorts  $\tau_i$  and  $\tau_j$  have depth greater than zero.

In other words, branching sorts include at least two collection types where neither type is nested within the other. For example,  $\tau_8$  and  $\tau_9$  in Figure 2.1 are branching sorts, while  $\tau_1$ – $\tau_7$  are non-branching sorts. By our convention on the structuring of tuple sorts, within a non-branching tuple sort  $\langle \tau_1, \dots, \tau_k \rangle$  only  $\tau_k$  can have depth greater than zero. Hence, non-branching sorts contain only “right deep” nesting of collections, as illustrated by sorts  $\tau_6$  and  $\tau_7$  in Figure 2.1.

We say that a sort is a *chain sort* if it contains precisely one descendant tuple sort, and that tuple sort is *flat*. Hence, a chain sort is a special case of a non-branching sort. For example, sorts  $\tau_1$ ,  $\tau_2$ ,  $\tau_4$ , and  $\tau_5$  in Figure 2.1 are chain sorts. The object encoding method that we propose later in this chapter is based upon chain sorts. Corresponding to each chain sort we define a sequence  $\bar{\xi}$  called its *signature*.

**Definition 2.1.3 (Signature)** A semantic indicator  $\xi$  is one of the characters **s**, **b**, or **n** and denotes that a collection is of type set, bag, or normalized bag, respectively. A signature  $\bar{\xi}$  of length  $d$  is a finite sequence of semantic indicators  $\xi_1 \xi_2 \dots \xi_d$ .

Any chain sort of depth  $d$  can be represented by a pair  $(\bar{\xi}, k)$ , where  $\bar{\xi}$  is a signature of length  $d$  that indicates from left-to-right the semantics of the successive descendant collection sorts, and  $k$  is the arity of the tuple at the leaf of the sort. For example, chain sorts  $\tau_1$ ,  $\tau_2$ ,  $\tau_4$ , and  $\tau_5$  from Figure 2.1 can be represented as  $(\emptyset, 3)$ ,  $(\mathbf{s}, 2)$ ,  $(\mathbf{ns}, 1)$ , and  $(\mathbf{sb}, 1)$ , respectively.

## Complex Objects

Corresponding to each sort  $\tau \in \mathcal{T}$  is a space of possible values that conform to it, called its *interpretation*.

**Definition 2.1.4 (Sort Interpretation)** The interpretation of sort  $\tau$ , denoted  $[[\tau]]$ , is defined recursively as follows.

$$\begin{aligned}
[[\text{dom}]] &= \text{dom} \\
[[\{\tau\}]] &= \{\{v_1, \dots, v_j\} \mid j \geq 0, v_i \in [[\tau]], i \in [1, j]\} \\
[[\{\!\{\tau\}\!\}]] &= \{\{\!\{v_1, \dots, v_j\}\!\} \mid j \geq 0, v_i \in [[\tau]], i \in [1, j]\} \\
[[\{\!\!\{\tau\}\!\!\}]] &= \{\{\!\!\{v_1, \dots, v_j\}\!\!\} \mid j \geq 0, v_i \in [[\tau]], i \in [1, j]\} \\
[[\langle \tau_1, \dots, \tau_k \rangle]] &= \{\langle v_1, \dots, v_k \rangle \mid v_j \in [[\tau_j]], j \in [1, k]\}
\end{aligned}$$

For each integer  $i \geq 0$ , we use  $\text{dom}^i$  to denote the (infinite) set of all finite elements belonging to interpretations of sorts with depth at most  $i$ .

$$\text{dom}^i := \bigcup_{\tau \in \mathcal{T}^i} [[\tau]] \quad (2.2)$$

We use  $\text{dom}^\infty$  to denote the (infinite) set of all finite elements that conform to any sort.

$$\text{dom}^\infty := \bigcup_{\tau \in \mathcal{T}} [[\tau]] \quad (2.3)$$

**Definition 2.1.5 (Complex Object)** *A complex object is a member of  $\text{dom}^\infty$ .*

Any complex object  $o$  necessarily belongs to some sort interpretation  $[[\tau]]$ , in which case we say that  $o$  *conforms* to  $\tau$ . (If  $o$  contains one or more empty collections, it is possible that  $o$  conforms to infinitely many sorts.) Object  $o$  is a tuple if  $\tau$  is a tuple sort, and a collection if  $\tau$  is a collection sort.

**Example 5** Consider the following objects. Our convention throughout this thesis will be to use both integers and lowercase letters near the beginning of the alphabet to denote atomic elements (i.e., constants).

$$\begin{aligned} o_1 &:= \langle a, b, 0 \rangle \\ o_2 &:= \{ \langle a, 1 \rangle, \langle a, 2 \rangle, \langle b, 2 \rangle, \langle b, 3 \rangle \} \\ o_3 &:= \{ \langle \langle a \rangle, \langle 2, 5 \rangle \rangle, \langle \langle a \rangle, \langle 3, 4 \rangle \rangle, \langle \langle b \rangle, \langle 3, 4 \rangle \rangle \} \\ o_4 &:= \{ \{ \langle 1 \rangle, \langle 2 \rangle \}, \{ \langle 1 \rangle, \langle 2 \rangle \}, \{ \langle 1 \rangle, \langle 3 \rangle \} \} \\ o_5 &:= \{ \{ \langle 1 \rangle, \langle 1 \rangle, \langle 2 \rangle \}, \{ \langle 1 \rangle, \langle 2 \rangle \}, \{ \langle 1 \rangle, \langle 3 \rangle, \langle 3 \rangle \} \} \\ o_6 &:= \{ \langle 1, \{ a, a, c \} \rangle, \langle 2, \{ a, c \} \rangle \} \\ o_7 &:= \{ \langle a, \{ \langle 1 \rangle, \langle 2 \rangle \} \rangle, \langle b, \{ \langle 1 \rangle, \langle 3 \rangle \} \rangle, \langle c, \{ \} \rangle \} \\ o_8 &:= \{ \{ \{ \} \}, \{ \} \} \\ o_9 &:= \{ \{ \{ \{ \langle a \rangle, \langle c \rangle, \langle c \rangle \} \}, \\ &\quad \{ \langle a \rangle, \langle c \rangle, \langle c \rangle \}, \\ &\quad \{ \langle \} \} \}, \\ &\quad \{ \}, \\ &\quad \{ \langle 5 \rangle \} \} \\ o_{10} &:= \{ \} \\ o_{11} &:= \{ a, \{ a, a, c \}, c, \{ a, c \} \} \end{aligned}$$

Objects  $o_1$ – $o_9$  are complex objects because they conform to sorts  $\tau_1$ – $\tau_9$ , respectively, from Example 4. Object  $o_{10}$  conforms to any sort rooted by a set, including  $\tau_2$ ,  $\tau_3$ ,  $\tau_5$ ,  $\tau_6$ ,  $\tau_8$ , and  $\tau_9$ ; therefore it is also a complex object. Finally, although  $o_{11}$  satisfies our informal characterization from the start of Section 2.1, it is not formally a complex object (as per Definition 2.1.5) because the elements of the set are not uniformly typed and so  $o_{11}$  does not conform to a sort.

We will not consider further untyped objects such as  $o_{11}$  in the above example. Henceforth, when we refer to an “object” the reader should interpret this to mean a “complex object” as defined above. We also represent objects as trees, analogous to the representation of sorts in Figure 2.1. At times we may refer to the “nodes” or “leaves” of an object or a sort, by which we mean the nodes or leaves of the tree representation.

**Branching, Non-Branching and Chain Objects** We classify complex objects as either *branching* or *non-branching*. An object is non-branching if it belongs to the interpretation of a non-branching sort; otherwise it is branching. For example, consider the objects  $o_1$ – $o_{10}$  from Example 5 and the sorts  $\tau_1$ – $\tau_9$  from Example 4. Objects  $o_1$ – $o_7$  are all non-branching because they conform to  $\tau_1$ – $\tau_7$ , respectively, which are all non-branching sorts. Object  $o_8$  is a branching object because any sort to which it conforms is necessarily branching, including  $\tau_8$ . Similarly,  $o_9$  is a branching object because any sort to which it conforms is necessarily branching, including  $\tau_9$ . Object  $o_{10}$  is a *non-branching* object even though it conforms to the *branching* sort  $\tau_9$ —it also conforms to the *non-branching* sort  $\tau_2$  (among others).

We say that a complex object is a *chain* if there exists a chain sort to which it conforms. Objects  $o_1$ ,  $o_2$ ,  $o_4$ ,  $o_5$ , and  $o_{10}$  are chain objects because they conform to the chain sorts  $\tau_1$ ,  $\tau_2$ ,  $\tau_4$ ,  $\tau_5$ , and  $\tau_2$ , respectively.

**Complete, Incomplete, and Trivial Objects** We say that a complex object  $o$  is *complete* if it does not contain an empty collection. We say that a complex object  $o$  is *trivial* if it is either an empty collection or a tuple of trivial objects. For example, consider complex objects  $o_1$ – $o_{10}$  from Example 5 (recall that  $o_{11}$  is not a complex object). Objects  $o_1$ – $o_6$  are complete because they do not contain empty collections, while objects  $o_7$ – $o_{10}$  are incomplete. Object  $o_{10}$  is trivial.

### 2.1.2 Object Linearization

In this section we describe a transformation from an object of arbitrary sort to a chain object that is either trivial or complete.<sup>1</sup> This transformation—which we call *linearization*—takes place in two steps:

1. If the object is non-trivial, *complete it* to remove empty collections.
2. *Remove branching* via a Cartesian product of siblings.

In the case where the original object was a *collection* object that does not contain any 0-ary tuples, this transformation is injective (i.e., one-to-one). We will use object linearization at the end of this chapter to reduce the problem of testing equivalence of queries returning arbitrary complex objects to the equivalence of queries returning trivial or complete chain objects.

---

<sup>1</sup>Our reason for handling trivial objects as a special case will become clear in Section 2.2 when we encode objects as flat relations. Trivial objects correspond to an empty encoding relation.

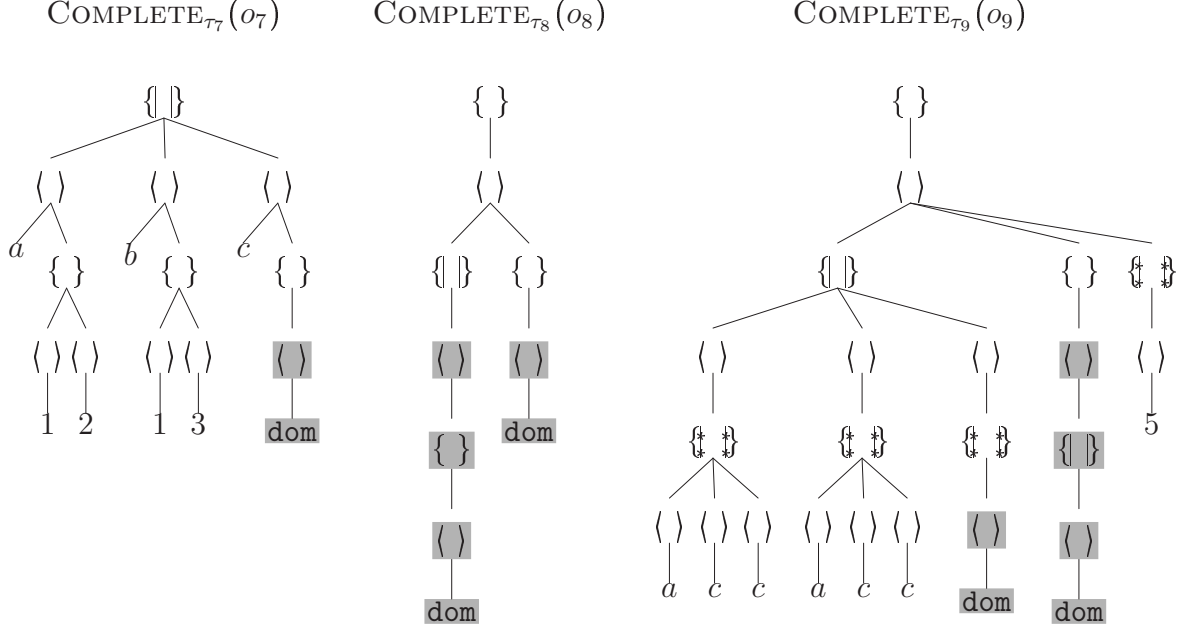


Figure 2.2: Completion of objects  $o_7$ ,  $o_8$ , and  $o_9$

## Object Completion

The first step in linearizing an object is to make it *complete* (cf. Section 2.1.1) by padding with placeholder nodes. Object completion is necessary to ensure that the second step remains lossless.

Given an object  $o$  that conforms to sort  $\tau$ , we define the *completion of  $o$  with respect to  $\tau$* , denoted  $\text{COMPLETE}_\tau(o)$ , as the new object created from  $o$  by the following process. If  $o$  is either trivial or already complete, then  $\text{COMPLETE}_\tau(o) = o$ . Otherwise, let  $o'$  be any descendant empty collection within  $o$ , and let  $\tau'$  be the descendant sort within  $\tau$  corresponding to  $o'$ . We know that  $\tau'$  is a collection sort, and so let  $\tau''$  be the child sort of  $\tau'$  (i.e.  $\tau'$  is a collection of  $\tau''$ 's). Then, we insert into the empty collection  $o'$  a “shaded copy” of  $\tau''$ . Repeating this process until  $o$  does not contain any more empty collections, we obtain  $\text{COMPLETE}_\tau(o)$ .

**Example 6** Consider the (non-trivial) incomplete objects  $o_7$ ,  $o_8$ , and  $o_9$  from Example 5 that conform to sorts  $\tau_7$ ,  $\tau_8$ , and  $\tau_9$ , respectively, from Example 4. Figure 2.2 illustrates the completion of each object with respect to its associated sort.

Object  $\text{COMPLETE}_\tau(o)$  may contain “shaded” delimiter pairs  $\{\{\}\}$ ,  $\{\}\}$ ,  $\{\}\}$ , or  $\{\}$ ; it also may contain the symbol **dom** which is used as an atomic value. If  $\text{COMPLETE}_\tau(o)$  contains any of these symbols, it is technically not a complex object (as per Definition 2.1.5). Hence, we want to extend our definition of complex objects slightly to allow them to contain both shaded delimiter pairs and the atomic value **dom**. Our intention is that shaded delimiter pairs retain the same semantics as their unshaded counterparts, but affect object equality (that is, two objects with

identical structure but different shading of nodes are not considered equal). To accomplish this, we modify our interpretation of sorts (Definition 2.1.4) as follows.

$$\begin{aligned}
[[ \mathbf{dom} ]] &= \mathbf{dom} \cup \{ \mathbf{dom} \} \\
[[ \{ \tau \} ]] &= \{ \{ v_1, \dots, v_j \} \mid j \geq 0, v_i \in [[ \tau ]], i \in [1, j] \} \\
&\quad \cup \{ \{ \mathbf{dom}, v_1, \dots, v_j \} \mid j \geq 0, v_i \in [[ \tau ]], i \in [1, j] \} \\
[[ \langle \tau_1, \dots, \tau_k \rangle ]] &= \{ \langle v_1, \dots, v_k \rangle \mid v_j \in [[ \tau_j ]], j \in [1, k] \} \\
&\quad \cup \{ \langle \mathbf{dom}, v_1, \dots, v_k \rangle \mid v_j \in [[ \tau_j ]], j \in [1, k] \}
\end{aligned}$$

The interpretations  $[[ \{ \tau \} ]]$  and  $[[ \langle \tau \rangle ]]$  are modified analogously to  $[[ \{ \tau \} ]]$ . Under this modified interpretation of sorts, completion with respect to sort  $\tau$  is a function

$$\text{COMPLETE}_\tau : [[ \tau ]]$$

and so, by Definition (2.1.5),  $\text{COMPLETE}_\tau(o)$  is a complex object.

**Observation 2.1.6** *Given an initial object  $o$  that does not contain any shaded nodes and a sort  $\tau$  to which  $o$  conforms, object  $\text{COMPLETE}_\tau(o)$  obeys the following properties:*

1. *The root node is always unshaded.*
2. *The only shaded atomic value is  $\mathbf{dom}$ .*
3. *No unshaded node occurs beneath a shaded node.*
4. *Any unshaded node that has a shaded child must be a collection type.*
5. *Any collection node*
  - (a) *has at least one child node;*
  - (b) *has at most one shaded child node; and*
  - (c) *cannot have both shaded and unshaded child nodes.*

Finally, we note that completion with respect to a sort  $\tau$  is both idempotent and invertible (by deleting shaded nodes).

## Removing Branching

Once an object  $o$  is completed, the next step in linearizing it is to transform it into a chain object. A recursive procedure for performing this transformation is shown in Algorithm 1. The critical step in calculating  $\text{CHAIN}(o)$  is the invocation of  $\text{DISTRIBUTE}$  on line 10 to remove tuple sorts with arity of two or more. The  $\text{DISTRIBUTE}$  function combines two chain objects  $o_a$  and  $o_b$  into a single chain object by distributing  $o_b$  across each leaf tuple  $t$  in  $o_a$  and then distributing the atomic values from  $t$  across the leaves of its copy of  $o_b$ .

Algorithm 1 can be applied to any complete complex object. By temporarily treating the symbol  $\mathbf{dom}$  as an atomic value, any sort  $\tau$  can be viewed as a complete complex object ( $\tau$  conforms to itself, and  $\text{COMPLETE}_\tau(\tau) = \tau$ ), allowing Algorithm 1 to be applied to sorts. This leads us to the following observation.

---

**Algorithm 1** Transforming objects into chains

---

CHAIN( $o$ )▷ **Input:** *complete* or *trivial* object  $o$ ▷ **Output:** chain object formed from  $o$ 

```
1  if  $o$  is atomic
2    then return  $\langle o \rangle$ 
3  elseif  $o = [o_1, \dots, o_n]$  ▷ Where  $[\cdot]$  is any (possibly shaded) collection type
4    then return  $[ \text{CHAIN}(o_1), \dots, \text{CHAIN}(o_n) ]$  ▷ preserving collection type  $[\cdot]$ 
5  elseif  $o = \langle \rangle$  or  $o = \langle \rangle$ 
6    then return  $\langle \rangle$ 
7  elseif  $o = \langle o_1 \rangle$  or  $o = \langle \langle o_1 \rangle \rangle$ 
8    then return CHAIN( $o_1$ )
9  elseif ( $o = \langle o_1, \dots, o_n \rangle$  or  $o = \langle \langle o_1, \dots, o_n \rangle \rangle$ ) and  $n > 1$ 
10   then return DISTRIBUTE(CHAIN( $o_1$ ), CHAIN( $\langle o_2, \dots, o_n \rangle$ ))
```

DISTRIBUTE( $o_a, o_b$ )▷ **Input:** *chain* object  $o_a$  of sort  $(\bar{\xi}^a, k)$ Assume that  $o_a$  is a tree whose  $m$  leaves are the  $k$ -ary tuples  $\langle a_1^1, \dots, a_k^1 \rangle, \dots, \langle a_1^m, \dots, a_k^m \rangle$ ▷ **Input:** *chain* object  $o_b$  of sort  $(\bar{\xi}^b, l)$ Assume that  $o_b$  is a tree whose  $n$  leaves are the  $l$ -ary tuples  $\langle b_1^1, \dots, b_l^1 \rangle, \dots, \langle b_1^n, \dots, b_l^n \rangle$ ▷ **Output:** chain object of sort  $(\bar{\xi}^a \cdot \bar{\xi}^b, k + l)$ formed by distributing  $o_b$  over each leaf of  $o_a$  and pushing down atomic values

```
1   $o \leftarrow$  copy of  $o_a$ 
2  foreach  $i \in [1, m]$ 
3    do  $o^i \leftarrow$  copy of  $o_b$ 
4    foreach  $j \in [1, n]$ 
5      do substitute tuple  $\langle a_1^i, \dots, a_k^i, b_1^j, \dots, b_l^j \rangle$ 
           for tuple  $\langle b_1^j, \dots, b_l^j \rangle$  within  $o^i$ 
6    substitute  $o^i$  for tuple  $\langle a_1^i, \dots, a_k^i \rangle$  within  $o$ 
7  return  $o$ 
```

---



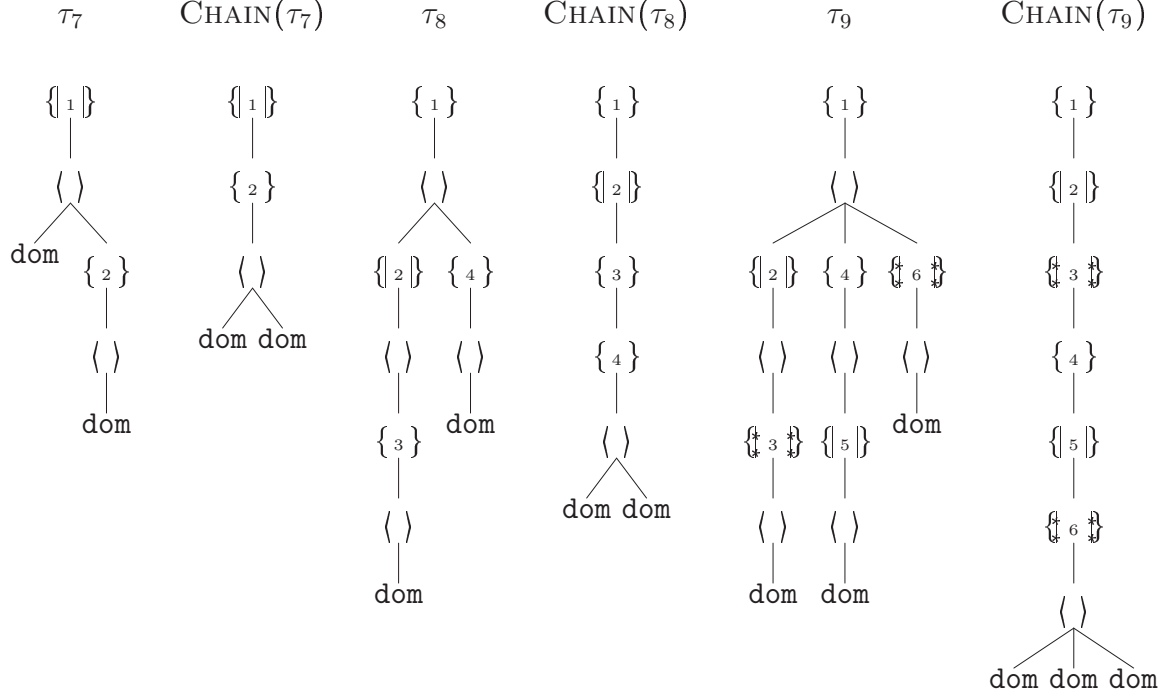


Figure 2.3: Linearization of sorts  $\tau_7$ ,  $\tau_8$ , and  $\tau_9$

**Observation 2.1.7** *If  $o$  is a complete object conforming to sort  $\tau$ , then  $\text{CHAIN}(o)$  is a complete chain object conforming to sort  $\text{CHAIN}(\tau)$ .*

**Definition 2.1.8 (Object Linearization)** *Given an object  $o$  conforming to sort  $\tau$ , the linearization of  $o$  with respect to  $\tau$  is the complete chain object given by the following function  $\text{LINEARIZE}_\tau : \tau \rightarrow \text{CHAIN}(\tau)$ .*

$$\text{LINEARIZE}_\tau(o) := \text{CHAIN} \circ \text{COMPLETE}_\tau(o) \quad (2.4)$$

Given an arbitrary sort  $\tau$ , we may at times refer to the chain sort  $\text{CHAIN}(\tau)$  as the *linearization of  $\tau$* . This is consistent with Definition 2.1.8, because if we treat  $\tau$  as an object of sort  $\tau$  to be linearized, we deduce the following.

$$\text{LINEARIZE}_\tau(\tau) = \text{CHAIN} \circ \text{COMPLETE}_\tau(\tau) = \text{CHAIN}(\tau)$$

**Example 7** Sorts  $\tau_7$ ,  $\tau_8$ , and  $\tau_9$  are shown in Figure 2.3 along with their linearizations  $\text{CHAIN}(\tau_7)$ ,  $\text{CHAIN}(\tau_8)$ , and  $\text{CHAIN}(\tau_9)$ . Collection types have been numbered to help illustrate the effect of Algorithm 1. Observe that in each case the ordering of collections in the chain corresponds to a pre-order traversal of the original sort.

**Example 8** Figure 2.4 shows the linearization of objects  $o_7$ ,  $o_8$ , and  $o_9$  with respect to sorts  $\tau_7$ ,  $\tau_8$ , and  $\tau_9$ . The reader is encouraged to step through the application of Algorithm 1 to the completed objects in Figure 2.2 to yield the objects shown here.

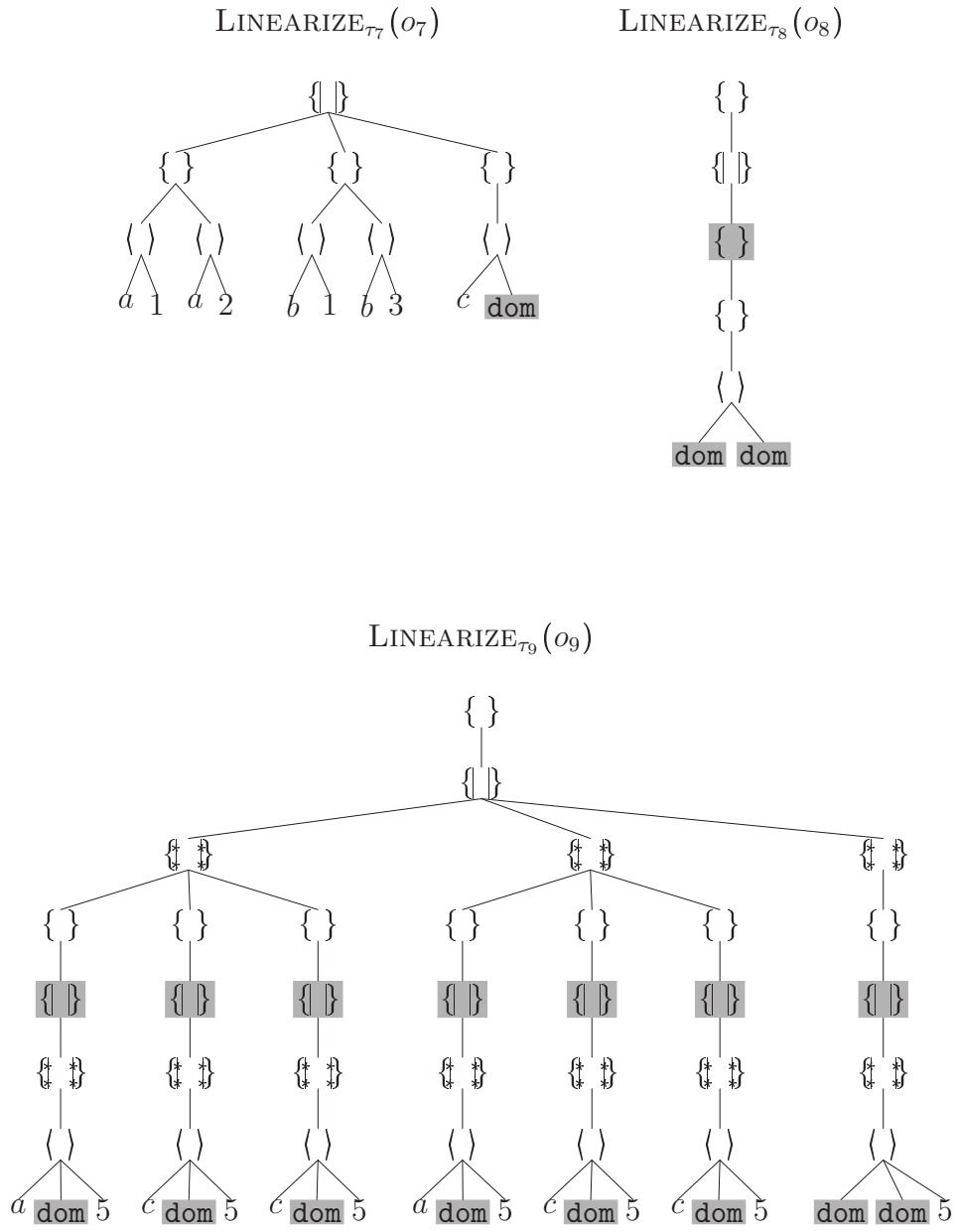


Figure 2.4: Linearization of objects  $o_7$ ,  $o_8$ , and  $o_9$

With certain restrictions (to be explained shortly), Algorithm 1 has a retraction (i.e., it is reversible). The retraction operation depends upon the sort of the original object; we don't go into detail because we never have need to perform the retraction operation. It suffices to point out that from the result of  $\text{DISTRIBUTE}(o_a, o_b)$  we can reconstruct both  $o_a$  and  $o_b$  because either both are trivial chain objects (in which case  $\text{DISTRIBUTE}(o_a, o_b)$  is also a trivial chain object), or neither contains an empty collection and so the output has redundancy akin to a relational multi-value dependency. There is a minor difficulty when dealing with shaded tuples because lines 6, 8, and 10 of procedure  $\text{CHAIN}$  discard the delimiters of the original tuples (and hence the bit of information denoting whether a tuple is shaded). For example, in the transformation of  $\text{COMPLETE}_{\tau_9}(o_9)$  (cf. Figure 2.2) to  $\text{LINEARIZE}_{\tau_9}(o_9)$  (cf. Figure 2.4) thirteen pairs of tuple delimiters were discarded, three of which were shaded. If the original object contained arbitrarily-shaded tuples, then information has been lost. However, if we assume that all shaded nodes in  $\text{COMPLETE}_{\tau}(o)$  were inserted by the completion of  $o$  with respect to some sort  $\tau$ , and  $\tau$  does not contain any zero-ary tuple sorts, then the properties in Observation 2.1.6 are sufficient for deducing which tuple delimiters require shading.

**Theorem 2.1.9** *Given sort  $\tau$  not containing any zero-ary tuple sorts,  $\text{LINEARIZE}_{\tau}$  has a retraction.*

**Corollary 2.1.10** *Given sort  $\tau$  not containing any zero-ary tuple sorts, any two objects  $o, o' \in [[\tau]]$  satisfy the following.*

$$o = o' \iff \text{LINEARIZE}_{\tau}(o) = \text{LINEARIZE}_{\tau}(o')$$

Finally, we note that the properties in Observation 2.1.6 may be violated by application of Algorithm 1. For example, linearized objects  $\text{LINEARIZE}_{\tau_8}(o_8)$  and  $\text{LINEARIZE}_{\tau_9}(o_9)$  in Figure 2.4 both violate the first property by having unshaded nodes occurring beneath shaded nodes.

**Observation 2.1.11** *Given an initial object  $o$  that does not contain any shaded nodes and a sort  $\tau$  to which it conforms, linearized object  $\text{LINEARIZE}_{\tau}(o)$  obeys the following properties:*

1. *The root node is always unshaded.*
2. *The only shaded atomic value is `dom`.*
3. *Any collection node*
  - (a) *has at least one child node;*
  - (b) *has at most one shaded child node; and*
  - (c) *cannot have both shaded and unshaded child nodes.*
4. *Tuple nodes are never shaded.*

The first three properties in Observation 2.1.11 hold because they hold in the input object  $\text{COMPLETE}_{\tau}(o)$  (by Observation 2.1.6) and are preserved by Algorithm 1 (although preservation of the third property is not obvious and relies upon

$\text{COMPLETE}_\tau(o)$  satisfying all of the properties in Observation 2.1.6). The fourth property holds because procedure `CHAIN` deletes all tuple nodes from its input and inserts unshaded tuples of atomic values at the leaves.

### 2.1.3 Object-Constructing Queries

In this subsection we propose the “Conjunctive Object-Constructing Query Language” ( $\text{COCQL}^{\mathcal{F}}$ )<sup>2</sup> which is intended to model a specific subclass of `SQL-92` [63] queries containing aggregation. As such,  $\text{COCQL}^{\mathcal{F}}$  uses a data model based on “tables,” which differ from traditional relations both by allowing duplicate tuples and by allowing attribute values to range over complex objects. After first defining the data model we then describe  $\text{COCQL}^{\mathcal{F}}$ , which is an algebraic language for specifying functions (queries) that map a database of flat tables to a complex object. We also define a subclass of  $\text{COCQL}^{\mathcal{F}}$  that is incapable of constructing non-trivial incomplete objects (i.e., objects containing empty nested collections), which we will denote as  $\text{COCQL}^{\mathcal{F}-}$ .

#### Data Model

Let `aname` and `tname` denote countably infinite sets of attribute names and table names, respectively, disjoint from each other and from the underlying domain `dom`. (Our convention will be to use uppercase letters to denote elements of `aname` and `tname`.) Our data model is defined as follows.

**Definition 2.1.12 (Table Schema)** *A table schema of arity  $k$  has the form*

$$T(A_1 : \tau_1, \dots, A_k : \tau_k) \quad (2.5)$$

for some  $k \geq 0$ , where  $T$  is a table name from `tname`,  $A_1, \dots, A_k$  are distinct attribute names from `aname`, and  $\tau_1, \dots, \tau_k$  are sorts. If  $\tau_1, \dots, \tau_k$  are all `dom`, then we say that the table schema is flat.

**Definition 2.1.13 (Relation Schema)** *A relation schema is a flat table schema.*

**Definition 2.1.14 (Table)** *A table  $T$  is a pair  $(\mathcal{S}(T), \mathbb{T}(T))$  where  $\mathcal{S}(T)$  is a table schema  $T(A_1 : \tau_1, \dots, A_k : \tau_k)$  and  $\mathbb{T}(T)$  is a finite bag-instance of  $\mathcal{S}(T)$ .*

$$\mathbb{T}(T) \in [\{ \{ \tau_1, \dots, \tau_k \} \}]$$

We call  $T$  flat if  $\mathcal{S}(T)$  is a flat schema.

**Definition 2.1.15 (Relation)** *A relation  $R$  is a table with a flat schema and whose instance  $\mathcal{S}(R)$  does not contain duplicate tuples. In other words, if  $\mathcal{S}(R)$  has arity  $k$ , then*

$$\mathbb{T}(R) \in [\{ \{ \overbrace{\text{dom}, \dots, \text{dom}}^k \} \}]$$

---

<sup>2</sup>Pronounced “cockle”.

We often omit the sorts and write a table schema as  $T(A_1, \dots, A_k)$ . If necessary, we use  $\mathbf{sort}(A_i)$  to denote sort  $\tau_i$  corresponding to attribute  $A_i$  (within the schema of table  $T$ , which will be obvious from the context). We use  $\mathbf{dom}(A_i)$  to denote the domain of values for attribute  $A_i$ ,

$$\mathbf{dom}(A_i) := [[ \mathbf{sort}(A_i) ]] = [[ \tau_i ]] \quad (2.6)$$

and we use  $\mathbf{adom}(A_i, T)$  to denote the *active domain* of attribute  $A_i$  within  $T$ , which is the finite subset of  $\mathbf{dom}(A_i)$  that appears within the projection of attribute  $A_i$  from table instance  $\mathbb{T}(T)$  (when table instance  $\mathbb{T}(T)$  is obvious from the context, we simply write  $\mathbf{adom}(A_i)$ ). Similarly, given a sequence  $\bar{A} = A_{i_1} \dots A_{i_j}$  of  $j$  (not necessarily distinct) attributes of  $T$ , we use  $\mathbf{sort}(\bar{A})$  to denote the sort of arity- $j$  tuples that can be projected from  $T$ , with  $\mathbf{dom}(\bar{A})$  denoting the corresponding domain of tuple values,

$$\mathbf{sort}(\bar{A}) := \langle \tau_{i_1}, \dots, \tau_{i_j} \rangle \quad (2.7)$$

$$\begin{aligned} \mathbf{dom}(\bar{A}) := & \{ \langle v_{i_1}, \dots, v_{i_j} \rangle \mid \langle v_{i_1}, \dots, v_{i_j} \rangle \in [[ \mathbf{sort}(\bar{A}) ]] \\ & \wedge \forall l, m \in [1, j] (A_{i_l} = A_{i_m} \implies v_{i_l} = v_{i_m}) \} \end{aligned} \quad (2.8)$$

and we use  $\mathbf{adom}(\bar{A}, T)$  (or  $\mathbf{adom}(\bar{A})$ ) to denote the active domain of  $\bar{A}$  within table instance  $\mathbb{T}(T)$ . Finally, we use  $\mathbf{sort}(\mathbb{S}(T))$  to denote the sort of tuples within  $T$

$$\mathbf{sort}(\mathbb{S}(T)) := \mathbf{sort}(A_1 \dots A_k) = \langle \tau_1, \dots, \tau_k \rangle \quad (2.9)$$

with  $\mathbf{dom}(\mathbb{S}(T))$  denoting the domain of possible tuples.

$$\mathbf{dom}(\mathbb{S}(T)) := \mathbf{dom}(A_1 \dots A_k) = [[ \langle \tau_1, \dots, \tau_k \rangle ]] = [[ \mathbf{sort}(\mathbb{S}(T)) ]] \quad (2.10)$$

Equations 2.8 and 2.10 are consistent because  $A_1, \dots, A_k$  are distinct names.

Observe from Definition 2.1.14 that each table instance  $\mathbb{T}(T)$  is a bag of tuples, which is a complex object. However, our data model does not define “nested tables” per se, because table schemas are not recursive and so attribute names are only assigned to the elements of the outermost tuple—a distinction that becomes important for the design of our query language.

**Definition 2.1.16 (Relational Database)** *A relational database  $\mathbb{D}$  is a finite set of table names*

$$\{T_1, \dots, T_m\} \subset \mathbf{tname}$$

*that have associated with each table name  $T_i$  a relation  $(\mathbb{S}(T_i), \mathbb{T}(T_i))$ . We call the tables in  $\mathbb{D}$  base tables (or base relations); the set of base table schemas are collectively called the database schema, denoted  $\mathbb{S}(\mathbb{D})$ .*

We define  $\mathbf{adom}(\mathbb{D}) \subset \mathbf{dom}$  to be the active domain of database  $\mathbb{D}$ , which is the set of values that occur within any attribute of any tuple within  $\mathbb{D}$ .

## COCQL <sup>$\mathcal{F}$</sup> and COCQL <sup>$\mathcal{F}$</sup> - Specification

Traditional conjunctive queries can be specified by relational algebra expressions combining base relations with the operators *selection*, *projection*, *Cartesian product* (or *join*), and *attribute renaming* (sometimes merged into the projection operator) [6, Ch. 4]. Standard approaches for extending relational algebra to bags of tuples include introducing two variations of projection (namely, *duplicate-preserving* and *duplicate-eliminating*), or defining projection as duplicate-preserving and introducing a new *duplicate-elimination* operator [42]. Well-formed algebra queries do not contain variable names; atomic values are only referenced either explicitly (as constants) or via attribute names occurring within the schema of the operator’s input. As its name implies, COCQL <sup>$\mathcal{F}$</sup>  is intended as an extension of conjunctive queries that allows for the construction of complex objects.

**Definition 2.1.17 (COCQL <sup>$\mathcal{F}$</sup>  Query)** *A COCQL <sup>$\mathcal{F}$</sup>  query  $Q$  is an expression of the form*

$$\text{EVAL}(T, f, \mathbf{E}(T))$$

where

- $T$  is an algebraic expression composed of the operators in Figure 2.5;
- $f \in \mathcal{F}$  is an aggregation function that takes a table as input and returns a complex object as output;
- $\mathbf{E}(T)$  is a query evaluation environment; and
- $\text{EVAL}()$  is the query evaluation function which evaluates algebra expression  $T$  within the context of  $\mathbf{E}(T)$  to yield the table to which  $f$  is applied.

Given a relational database  $\mathbb{D}$ , the result of evaluating  $Q$  over  $\mathbb{D}$ —denoted  $(Q)^{\mathbb{D}}$ —is a complex object of sort  $\text{sort}(Q) := \text{sort}(f(T))$ .

If for every database instance  $\mathbb{D}$  the object  $(Q)^{\mathbb{D}}$  is trivial, then  $Q$  is called unsatisfiable; otherwise  $Q$  is called satisfiable. If the environment  $\mathbf{E}(T)$  is empty, then  $Q$  is called unparameterized; otherwise  $Q$  is called parameterized.

Although a COCQL <sup>$\mathcal{F}$</sup>  query contains an algebraic expression, it also contains some additional non-algebraic syntax. These non-algebraic constructs lead to several features that distinguish COCQL <sup>$\mathcal{F}$</sup>  from the usual language of relational algebra expressions over bags.

1. **Environment Variables:** Associated with the evaluation of table-generating expression  $T$  is an environment  $\mathbf{E}(T)$ , which is a set of assignments of the form

$$\mathbf{E}(T) = \{V_1 \leftarrow v_1, \dots, V_n \leftarrow v_n\}$$

for some  $n \geq 0$ , where for each  $i \in [1, n]$ ,  $V_i$  is a distinct name from **aname** with which a sort  $\text{sort}(V_i)$  is implicitly associated, and  $v_i \in \text{dom}(V_i) = \llbracket \text{sort}(V_i) \rrbracket$ . Environment  $\mathbf{E}(T)$  is inherited down the algebraic operator tree so that any reference to  $V_i$  anywhere within expression  $T$  is interpreted as a reference to the value  $v_i$ .

		Output Table $T$	
Operator Name	Query Expression	Schema $S(T)$	Instance $\mathbb{T}(T)$
Base Relation	$T(A_1, \dots, A_k)$ where $A_1, \dots, A_k \in \text{aname}$ are “fresh” attribute names	$T(A_1, \dots, A_k)$	$\mathbb{T}(T) \in \mathbb{D}$
Cartesian Product	$T_1 \times T_2$	$T(A_1, \dots, A_{k_1}, B_1, \dots, B_{k_2})$	$\{ \langle t_1[A_1], \dots, t_1[A_{k_1}], t_2[B_1], \dots, t_2[B_{k_2}] \rangle \mid t_1 \in T_1, t_2 \in T_2 \}$
Selection	$\sigma_{A=B}(T_1)$ where $\text{sort}(A) = \text{sort}(B) = \text{dom}$	$T(A_1, \dots, A_{k_1})$	$\{ \langle t_1[A_1], \dots, t_1[A_{k_1}] \rangle \mid t_1 \in T_1, t_1[A] = t_1[B] \}$
Duplicate-Preserving Projection	$\Pi_{A_{i_1}, \dots, A_{i_k}}^{dup}(T_1)$ where $\text{sort}(A_{i_1}) = \dots = \text{sort}(A_{i_k}) = \text{dom}$	$T(A_{i_1}, \dots, A_{i_k})$	$\{ \langle t_1[A_{i_1}], \dots, t_1[A_{i_k}] \rangle \mid t_1 \in T_1 \}$
Duplicate-Eliminating Projection	$\Pi_{A_{i_1}, \dots, A_{i_k}}(T_1)$ where $\text{sort}(A_{i_1}) = \dots = \text{sort}(A_{i_k}) = \text{dom}$	$T(A_{i_1}, \dots, A_{i_k})$	$\{ \langle t_1[A_{i_1}], \dots, t_1[A_{i_k}] \rangle \mid t_1 \in T_1 \}$
Scalar Subquery	$\lambda X = f(T_2)(T_1)$ where $X \in \text{aname}$ is a “fresh” attribute name and $f \in \mathcal{F}$	$T(A_1, \dots, A_{k_1}, X)$ where $\text{sort}(X) = \text{sort}(f(T_2))$	$\{ \langle t_1[A_1], \dots, t_1[A_{k_1}], \text{EVAL}(T_2, f, \mathbf{E}(T)) \rangle \mid t_1 \in T_1 \}$
Conjunctive Selection	$\sigma_{p_1 \wedge p_2}(T_1)$		$\sigma_{p_2}(\sigma_{p_1}(T_1))$
Join	$T_1 \bowtie_p T_2$		$\sigma_p(T_1 \times T_2)$
Generalized Projection	$\Pi_{A_{i_1}, \dots, A_{i_k}}^{X=f(A_{j_1}, \dots, A_{j_l})}(T_1)$ where $X \in \text{aname}$ is a “fresh” attribute name and $\text{sort}(A_{i_1}) = \dots = \text{sort}(A_{i_k}) = \text{dom}$	$\lambda X = f(\Pi_{A_{j_1}, \dots, A_{j_l}}^{dup}(\sigma_{A_{i_1} \wedge \dots \wedge A_{i_k}'}(T_1')))(\Pi_{A_{i_1}, \dots, A_{i_k}}(T_1))$ where $T_1'$ is a copy of $T_1$ with schema $T_1'(A_{i_1}, \dots, A_{i_k})$	

Figure 2.5: COCQL $\mathcal{F}$  algebra operators

2. **Aggregation Functions:** We assume a set  $\mathcal{F}$  of functions that each aggregate a table into a complex object. Given function  $f \in \mathcal{F}$  and a table (or table-generating expression)  $T$ , we use  $\text{sort}(f(T))$  to denote the sort of object  $f(T)$ . For this thesis we will restrict our attention to the case where  $\mathcal{F}$  only contains three functions

$$\mathcal{F} = \{f_b, f_n, f_s\} \quad (2.11)$$

that return the bag, normalized bag, and set projections, respectively, of the collection of tuples in the instance of table  $T$ . That is:

$$\begin{aligned} f_b(T) &:= \text{Bag}(\mathbb{T}(T)) & \text{sort}(f_b(T)) &= \{\!\! \{ \text{sort}(\mathbb{S}(T)) \}\!\! \} \\ f_n(T) &:= \text{NormBag}(\mathbb{T}(T)) & \text{sort}(f_n(T)) &= \{\!\! \{ \text{sort}(\mathbb{S}(T)) \}\!\! \} \\ f_s(T) &:= \text{Set}(\mathbb{T}(T)) & \text{sort}(f_s(T)) &= \{ \text{sort}(\mathbb{S}(T)) \} \end{aligned}$$

3. **Query Composability:** A  $\text{COCQL}^{\mathcal{F}}$  query  $Q = \text{EVAL}(T, f, \mathbf{E}(T))$  takes as inputs a set of tables and outputs a complex object of sort  $\text{sort}(Q) = \text{sort}(f(T))$ . Hence,  $\text{COCQL}^{\mathcal{F}}$  is not a fully-composable query language, per se. However, the sub-language of algebraic expressions is fully-composable, and arbitrary nesting of queries is made possible by the  $\lambda$  algebraic operator described below.

The algebra operators for specifying expression  $T$  are defined in Figure 2.5, using the following notational conventions:

1. We use  $T_1$  and  $T_2$  to denote algebra expressions yielding tables with schema  $T_1(A_1, \dots, A_{k_1})$  and  $T_2(B_1, \dots, B_{k_2})$ , respectively.
2. We use  $\{\!\! \{ t_1 \mid t_1 \in T_1 \}\!\! \}$  to denote a bag constructed by letting tuple variable  $t_1$  range over tuple *occurrences* in  $\mathbb{T}(T_1)$ .<sup>3</sup> This requires an inherently procedural (i.e. variable-binding) interpretation of  $t_1 \in T_1$  that differs from the usual interpretation of  $t_1 \in T_1$  as a predicate.
3. Given some  $A \in \text{aname} \cup \text{dom}$  and a binding of tuple variable  $t_1 \in T_1$  to value  $t_1 = \langle a_1, \dots, a_{k_1} \rangle$ , we use  $t_1[A]$  to denote the following complex object:<sup>3</sup>
  - If  $A = A_i$  for some  $i \in [1, k_1]$ , then  $t_1[A] := a_i$ .
  - If  $A \in \text{aname}$  but is not an attribute in  $\mathbb{S}(T_1)$ , then  $\mathbf{E}(T)$  must contain some assignment  $A \leftarrow v$  and  $t_1[A] := v$ .
  - If  $A \in \text{dom}$ , then  $t_1[A] := A$  and  $\text{dom}(A) := \text{dom}$ .

The definitions in Figure 2.5 for the core algebraic operators *base relation*, *Cartesian product*, *selection*, and *projection* (duplicate-preserving and duplicate-eliminating) are all fairly standard, as are the syntactic extensions *conjunctive selection* and *join*. In a minor break from tradition, we completely ignore the attribute names in the database schema and instead embed attribute renaming into the base relation operator—this avoids the need for attribute renaming at intermediate stages of the query. With one caveat, all of these algebraic operators are

---

<sup>3</sup>Analogous for  $t_2 \in T_2$ .



blind to the sort of their attributes and so treat  $\text{dom}^\infty$  as an abstract atomic domain. The exception is a syntactic restriction on the selection, duplicate-eliminating projection, and generalized projection operators to prevent the (explicit or implicit) comparison of higher-order objects.

**Example 9** Consider the following simple student enrolment database  $\mathbb{D}_1$ .

Prof	pid	Course	cid pid	Student	sid	Enrol	cid sid
	a		c <sub>1</sub> a		1		c <sub>1</sub> 1
	b		c <sub>2</sub> a		2		c <sub>1</sub> 2
	c		c <sub>3</sub> b		3		c <sub>2</sub> 2
			c <sub>4</sub> b		4		c <sub>3</sub> 1
			c <sub>5</sub> b				c <sub>4</sub> 3

The following COCQL $_{\{f_b, f_s, f_n\}}$  query  $Q_7$  returns the set of pairs  $(p, s)$  such that professor  $p$  teaches a course in which student  $s$  is enrolled.

$$Q_7: \quad \text{EVAL}(\Pi_{B,D}^{\text{dup}}(\text{Course}(A, B) \bowtie_{A=C} \text{Enrol}(C, D)), f_s, \emptyset)$$

The reader can verify that  $Q_7$  returns the following sort

$$\text{sort}(Q_7) = \{ \langle \text{dom}(B), \text{dom}(D) \rangle \} = \{ \langle \text{dom}, \text{dom} \rangle \} = \tau_2 \quad \text{from Example 4}$$

and that  $Q_7^{\mathbb{D}_1}$ —the result of evaluating  $Q_7$  over database  $\mathbb{D}_1$ —yields the following complex object.

$$Q_7^{\mathbb{D}_1} = \{ \langle a, 1 \rangle, \langle a, 2 \rangle, \langle b, 2 \rangle, \langle b, 3 \rangle \} = o_2 \quad \text{from Example 5}$$

The most common method for extending relational algebra to handle aggregation is to explicitly define a *generalized projection* operator [42, 54], which is a unary operator that first partitions the input table and then applies an aggregation function to each partition. (If grouped tables are explicitly supported by the data model, generalized projection can be split into separate *partition* and *aggregation* steps [63, 92].) In our algebra, aggregation is instead handled by the *scalar subquery* operator, which is a *pseudo-binary* operator. A true binary relational algebra operator specifies a logical operation that can be performed on two independent inputs. In contrast, our scalar subquery operator requires a nested loops semantics in which algebra expression  $T_2$ —the “(correlated) subquery”—is re-evaluated for each tuple in  $T_1$  to yield a new complex object  $f(T_2)$ ; we use the symbol  $\lambda$  to emphasize the obvious similarity to a lambda-calculus expression. Our  $\lambda$  operator is similar to the *Apply* operator that Galindo-Legaria and Joshi introduce to handle correlated subqueries [38, 39], but whereas *Apply* returns a table for the subquery which is then crossed with  $T_1$ ,  $\lambda$  always returns a single complex object from the subquery. We define the generalized projection operator in Figure 2.5 as a syntactic extension that is specified in terms of  $\lambda$ . The following example illustrates the scalar subquery operator; the reader is referred to Example 13 for an example of generalized projection.

**Example 10** Assuming the same database  $\mathbb{D}_1$  as in Example 9, consider the following  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  query  $Q_8$  which returns the sets of students who have been taught by each professor.

$$Q_8: \quad \text{EVAL}(\lambda_{F=f_s(\Pi_D^{\text{dup}}(\sigma_{B=E}(\text{Course}(A,B)\bowtie_{A=C}\text{Enrol}(C,D))))}(\text{Prof}(E)), f_b, \emptyset)$$

The reader can verify that  $Q_8$  returns the following sort

$$\text{sort}(Q_8) = \{\{\text{dom}(E), \text{dom}(F)\}\} = \{\{\text{dom}, \{\{\text{dom}\}\}\}\} = \tau_7 \quad \text{from Example 4}$$

and that  $Q_8^{\mathbb{D}_1}$  yields the following complex object.

$$Q_8^{\mathbb{D}_1} = \{\{a, \{\{1\}, \{2\}\}\}, \{b, \{\{1\}, \{3\}\}\}, \{c, \{\}\}\} = o_7 \quad \text{from Example 5}$$

Query  $Q_8$  is able to generate *non-trivial incomplete* object  $o_7$  because it is explicitly written using a scalar subquery. If the query were written to use generalized projection instead of a scalar subquery, professor  $c$  would not appear in the output.

As mentioned in the above example, the scalar subquery operator is required in order to construct non-trivial incomplete objects. A very practical subclass of  $\text{COCQL}^{\mathcal{F}}$  queries is those that are known to always output complete objects.

**Definition 2.1.18** ( $\text{COCQL}^{\mathcal{F}-}$  Query) *The language  $\text{COCQL}^{\mathcal{F}-}$  is the subclass of  $\text{COCQL}^{\mathcal{F}}$  queries that do not contain an explicit scalar subquery operator  $\lambda$ . Queries in  $\text{COCQL}^{\mathcal{F}-}$  may contain generalized projection.*

We emphasize that none of the algebraic operators have the ability to look inside a complex object. All knowledge of the internal structure of objects is encapsulated within the aggregation functions in  $\mathcal{F}$ , while the query algebra simply treats  $\text{dom}^\infty$  as an abstract domain. This is why we christen  $\text{COCQL}^{\mathcal{F}}$  as an “object-constructing” query language, rather than an “object query language.”

## Relating $\text{COCQL}^{\mathcal{F}}$ and SQL

There are three main differences between the data model of the SQL standard [63] and the data model we have defined for the query language  $\text{COCQL}^{\mathcal{F}}$ . First, the SQL data model allows the base tables stored by the RDBMS to contain duplicates, whereas  $\text{COCQL}^{\mathcal{F}}$  queries are expressed over databases containing relations. Second, the SQL data model designates a particular atomic value—NULL—to represent “unknown” data, and requires comparisons between atomic elements to be evaluated using *three-valued* logic. In contrast, we have assumed a standard *boolean* interpretation of the equality operator over elements of  $\text{dom}$  (the underlying domain for the  $\text{COCQL}^{\mathcal{F}}$  data model). Third, SQL tables only allow atomic-valued attributes, whereas  $\text{COCQL}^{\mathcal{F}}$  queries can construct arbitrary complex-valued attributes. As

previously mentioned, this third difference is relatively minor since the  $\text{COCQL}^{\mathcal{F}}$  algebraic operators cannot look inside a complex-valued attribute, and so  $\text{COCQL}^{\mathcal{F}}$  essentially treats  $\text{dom}^{\infty}$  as an atomic domain.

SQL and  $\text{COCQL}^{\mathcal{F}}$  also differ in the aggregation functions that they support. For comparison purposes, let  $\text{SQL}^{\mathcal{F}}$  denote an “SQL-like” query language operating over the data model we have defined for  $\text{COCQL}^{\mathcal{F}}$  (i.e. complex-valued table attributes; no NULL values; and databases containing relations) and in which the usual SQL aggregation functions (SUM, COUNT, AVG, etc.) have been replaced with the set of functions  $\mathcal{F}$ . Furthermore, let  $\text{COCQL}^{\mathcal{F},f_b}$  denote the slight extension of the language  $\text{COCQL}^{\mathcal{F}}$  in which the aggregate function  $f_b$  may be used in the root-level invocation of EVAL, even if  $f_b \notin \mathcal{F}$  (clearly  $\text{COCQL}^{\mathcal{F},f_b} = \text{COCQL}^{\mathcal{F}}$  if  $f_b \in \mathcal{F}$ ).

**Proposition 2.1.19**  *$\text{COCQL}^{\mathcal{F},f_b}$  corresponds to a restricted class of queries that can be expressed in  $\text{SQL}^{\mathcal{F}}$  using only the following constructs:*

1. *SELECT* clauses composed of attribute names, constants, scalar subqueries, and aggregation expressions of the form  $f(A_1, \dots, A_k)$ ,  $f \in \mathcal{F}$ ;
2. *FROM* clauses composed of a list of base relation names and nested  $\text{SQL}^{\mathcal{F}}$  table expressions;
3. *WHERE* clauses composed of conjunctions of equality comparisons, with each comparison containing only constants and/or references to attributes originating from base relations; and
4. *GROUP BY* clauses composed only of constants and attributes originating from base relations.

We denote this subset of  $\text{SQL}^{\mathcal{F}}$  expressions as  $\text{CSQL}^{\mathcal{F}}$  (for “Conjunctive  $\text{SQL}^{\mathcal{F}}$ ”).

**Proposition 2.1.20**  *$\text{COCQL}^{\mathcal{F},f_b^-}$  corresponds to the restricted class of  $\text{CSQL}^{\mathcal{F}}$  queries in which*

- every nested *SELECT-FROM-WHERE-GROUPBY* block that contains an aggregation function in the *SELECT* clause also has as *GROUPBY* clause, and
- the *SELECT* clause does not contain scalar subqueries

(i.e., the subset of  $\text{CSQL}^{\mathcal{F}}$  that permits scalar aggregation only at the root block). We denote this subset of  $\text{CSQL}^{\mathcal{F}}$  expressions as  $\text{CSQL}^{\mathcal{F}-}$ .

A  $\text{COCQL}^{\mathcal{F},f_b}$  query always returns a single complex object (treated as a scalar value from abstract domain  $\text{dom}^{\infty}$ ). In contrast, a  $\text{CSQL}^{\mathcal{F}}$  query always returns a table which—per Definition 2.1.14—is a table schema paired with a bag-object. Corresponding to any  $\text{CSQL}^{\mathcal{F}}$  query is a  $\text{COCQL}^{\mathcal{F},f_b}$  query whose root invocation of EVAL applies the function  $f_b$  to construct the same bag-object as the table instance returned by the  $\text{CSQL}^{\mathcal{F}}$  query.

**Example 11** Assume that a database schema contains relation  $\text{Person}(\text{id}, \text{fname}, \text{lname}, \text{phone})$ .  $\text{CSQL}^{\mathcal{F}}$  query  $Q_9$  returns a table of (fname, phone) pairs for persons with the last name DeHaan.

$Q_9$ :     SELECT fname, phone FROM Person WHERE lname = ‘DeHaan’

Query  $Q_{10}$  shows a corresponding  $\text{COCQL}^{\mathcal{F}, f_b}$  query which returns a bag of binary tuples.

$Q_{10}$ :     EVAL( $\Pi_{B,D}^{dup}(\sigma_{C='DeHaan'}(\text{Person}(A, B, C, D)))$ ),  $f_b, \emptyset$ )

Given an arbitrary  $\text{COCQL}^{\mathcal{F}}$  (or  $\text{COCQL}^{\mathcal{F}, f_b}$ ) query whose root invocation of EVAL applies the function  $f$ , the translation of the query into  $\text{CSQL}^{\mathcal{F}}$  depends upon  $f$ . If  $f = f_b$ , then the corresponding  $\text{CSQL}^{\mathcal{F}}$  query returns a table whose instance is the same bag-object as the output of the  $\text{COCQL}^{\mathcal{F}}$  query. If  $f \neq f_b$ , then the corresponding  $\text{CSQL}^{\mathcal{F}}$  query uses scalar aggregation in the root block to apply the function  $f$  and the resulting object is wrapped within a single-row, single-column table.

**Example 12** Assuming the same Person relation as the previous example,  $\text{COCQL}^{\mathcal{F}}$  query  $Q_{11}$  returns the result of applying an arbitrary function  $f \in \mathcal{F}$  to a collection of (fname, phone) pairs for persons with the last name DeHaan.

$Q_{11}$ :     EVAL( $\Pi_{B,D}^{dup}(\sigma_{C='DeHaan'}(\text{Person}(A, B, C, D)))$ ),  $f, \emptyset$ )

If  $f = f_b$ , then the corresponding  $\text{CSQL}^{\mathcal{F}}$  query is  $Q_9$  from Example 11. Otherwise, the corresponding  $\text{CSQL}^{\mathcal{F}}$  query is the scalar aggregation query  $Q_{12}$  below.

$Q_{12}$ :     SELECT  $f(\text{fname}, \text{phone})$  FROM Person WHERE lname = ‘DeHaan’

We have defined the  $\text{COCQL}^{\mathcal{F}}$  generalized projection operator as a syntactic extension for a correlated scalar subquery. In contrast, the semantics of the GROUP BY construct in SQL are usually described in terms of partitioning—and for good reason, because physical implementations of grouping/partitioning are typically much more efficient than nested-loop-based evaluation of correlated subqueries. Nevertheless, a logical rewriting of GROUP BY in terms of correlated subqueries is also possible within  $\text{CSQL}^{\mathcal{F}}$  as illustrated in the next example.

**Example 13** Using the same Person(id, fname, lname, phone) relation as the previous example, consider  $\text{CSQL}^{\mathcal{F}}$  query  $Q_{13}$ , which returns the result of applying arbitrary aggregation function  $f \in \mathcal{F}$  to the collection of first names for each group of persons with the same last name. This corresponds to  $\text{COCQL}^{\mathcal{F}, f_b}$  query  $Q_{14}$ .

$Q_{13}$ :     SELECT  $f(\text{fname})$  FROM Person GROUP BY lname

$Q_{14}$ :     EVAL( $\Pi_E^{dup}(\Pi_C^{E=f(B)}(\text{Person}(A, B, C, D)))$ ),  $f_b, \emptyset$ )

Expanding the generalized projection in  $Q_{14}$  yields  $\text{COCQL}^{\mathcal{F}, f_b}$  query  $Q_{15}$ , which corresponds to  $\text{CSQL}^{\mathcal{F}}$  query  $Q_{16}$ .

$Q_{15}$ :  $\text{EVAL}(\Pi_E^{dup}(\lambda_{E=f(\Pi_{B'}^{dup}(\sigma_{C'=C}(\text{Person}(A, B, C, D))))}(\Pi_C(\text{Person}(A, B, C, D))))$ ,  
 $f_b, \emptyset$ )  
 $Q_{16}$ :  $\text{SELECT}(\text{SELECT } f(\text{fname})$   
 $\text{FROM Person}$   
 $\text{WHERE lname} = \text{t.lname})$   
 $\text{FROM}(\text{SELECT lname FROM Person GROUP BY lname}) \text{ AS t}$

CSQL <sup>$\mathcal{F}$</sup>  queries  $Q_{13}$  and  $Q_{16}$  are equivalent; however, this equivalence relies upon the fact that attribute `lname` does not contain `NULL` values (as per our definition of language SQL <sup>$\mathcal{F}$</sup> ). If the data model were extended with `NULL` values, the `WHERE` clause of the subquery in  $Q_{16}$  would need minor adjustment to account for SQL's three-valued logic.

Although we do not prove it formally here, the correspondence between COCQL <sup>$\mathcal{F}, f_b$</sup>  and CSQL <sup>$\mathcal{F}$</sup>  posited in Proposition 2.1.19 can be verified by

1. defining a canonical form for COCQL <sup>$\mathcal{F}, f_b$</sup>  algebra expressions,
2. demonstrating the correspondence between CSQL <sup>$\mathcal{F}$</sup>  expressions and COCQL <sup>$\mathcal{F}, f_b$</sup>  expressions in canonical form, and
3. showing that every COCQL <sup>$\mathcal{F}, f_b$</sup>  can be converted to the canonical form by a sequence of equivalence-preserving algebraic transformations.

The correspondence between COCQL <sup>$\mathcal{F}, f_b^-$</sup>  and CSQL <sup>$\mathcal{F}^-$</sup>  posited in Proposition 2.1.20 can be verified in a similar fashion.

## Deciding Equivalence of COCQL <sup>$\mathcal{F}$</sup> Queries

In the introduction to this thesis we purport to study the equivalence of SQL queries containing nested aggregation functions. Having just introduced the algebraic query language COCQL <sup>$\mathcal{F}$</sup>  and shown its relationship to a specific class of SQL queries, we briefly define formally what it means for two COCQL <sup>$\mathcal{F}$</sup>  queries to be equivalent.

**Definition 2.1.21 (COCQL <sup>$\mathcal{F}$</sup>  Query Equivalence)** *Two COCQL <sup>$\mathcal{F}$</sup>  queries  $Q, Q'$  over database schema  $\mathbb{S}(\mathbb{D})$  with the same output sort  $\tau$  are equivalent—denoted  $Q \equiv Q'$ —if for each relational database  $\mathbb{D}$  with schema  $\mathbb{S}(\mathbb{D})$  and finite instance  $\mathbb{T}(\mathbb{D})$  the complex objects  $(Q)^{\mathbb{D}}$  and  $(Q')^{\mathbb{D}}$  are equal.*

**Proposition 2.1.22** *The query equivalence problems for COCQL <sup>$\mathcal{F}, f_b$</sup>  and CSQL <sup>$\mathcal{F}$</sup>  are polynomial-time reducible to each other.*

**Proposition 2.1.23** *The query equivalence problems for COCQL <sup>$\mathcal{F}, f_b^-$</sup>  and CSQL <sup>$\mathcal{F}^-$</sup>  are polynomial-time reducible to each other.*

Much of the remainder of this thesis concerns how to test COCQL <sup>$\mathcal{F}$</sup>  query equivalence. The main goal of the next two sections is to reduce COCQL <sup>$\mathcal{F}$</sup>  query equivalence to testing a special form of equivalence between queries that return flat relations.

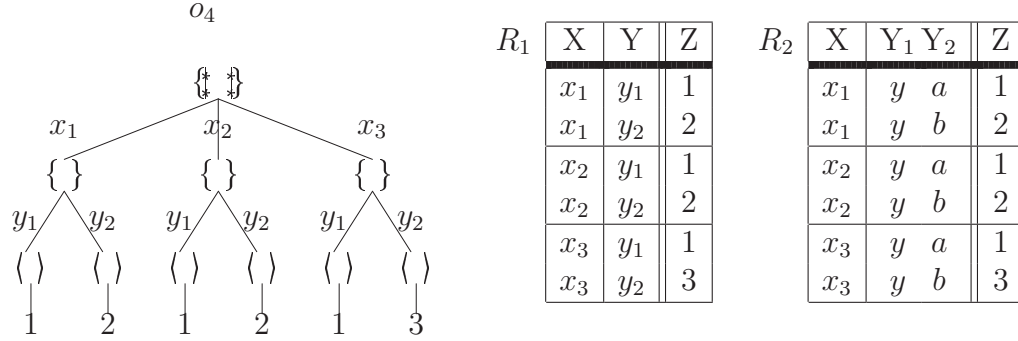


Figure 2.6: Encoding of complete chain object  $o_4$

## 2.2 Encoding Objects as Relations

Section 2.1.2 presented a process, called *linearization*, for transforming an arbitrary complex object into a complete chain object that possibly contains “shaded” nodes. In this section we describe a method for converting between a complete chain object and its representation as a set of flat tuples. In Section 2.2.1 we only consider complete chain objects that do not have shaded nodes—a restriction that is lifted in Section 2.2.2. Combined with Section 2.1.2, this provides a method for encoding arbitrary objects as a single relation. (This encoding technique is similar to one used previously by Suciu [102]; see Section 2.4.2 for a brief comparison.) Although each object has a unique linearization, the same cannot be said for relational encodings—there are an infinite number of different encodings of the same object. Therefore, Section 2.2.3 defines a binary relationship between encoding relations that characterizes when two different relations encode the same object.

### 2.2.1 Encoding and Decoding Chain Objects

Figure 2.6 illustrates the basic technique for encoding a complete chain object as a set of tuples, showing both the complete chain object  $o_4$  from Example 5 and its encoding as relation  $R_1$ . Starting with the tree representation of the object, we label each outgoing edge from a collection node with a locally-unique value (shown) and then create a tuple for each root-to-leaf path down the tree, incorporating the edge labels as attribute values. Note that edge labels could be drawn from a composite domain that corresponds to multiple relational attributes. For example, relation  $R_2$  shows an alternative encoding that could be formed by replacing labels  $y_1$  and  $y_2$  with  $y \cdot a$  and  $y \cdot b$ , respectively. We now give a more formal definition of encoding relations.

**Definition 2.2.1 (Encoding Schema)** *An encoding schema of depth  $d \geq 0$  is a relation schema  $R(A_1, \dots, A_k)$  annotated with the following additional information:*

- For each  $i \in [1, d]$ , the “tuple of index attributes at level  $i$ ”  $\bar{I}_i$  is a sequence

$$\bar{I}_i := I_{i,1} \cdots I_{i,k_i}$$

where  $k_i \geq 0$  and for  $j \in [1, k_i]$  each  $I_{i,j} \in \{A_1, \dots, A_k\}$ . Because sequences allow repeated elements, we use  $\mathcal{I}_i$  to denote the set of attributes occurring in  $\bar{\mathcal{I}}_i$ . For convenience, we use  $\bar{\mathcal{I}}_{[i,j]}$  to denote the sequence  $\bar{\mathcal{I}}_i \bar{\mathcal{I}}_{i+1} \dots \bar{\mathcal{I}}_j$ , and  $\mathcal{I}_{[i,j]}$  to denote the corresponding set  $\bigcup_{l \in [i,j]} \mathcal{I}_l$ .

- The “tuple of output attributes”  $\bar{\mathcal{V}}$  is a sequence

$$\bar{\mathcal{V}} := V_1 \dots V_{k_v}$$

where  $k_v \geq 0$  and for  $j \in [1, k_v]$  each  $V_j \in \{A_1, \dots, A_k\}$ . We use  $\mathcal{V}$  to denote the set of output attributes occurring in  $\bar{\mathcal{V}}$ .

Furthermore, the sets of index and output attributes must together cover all the attributes of relation  $R$ .

$$\mathcal{I}_1 \cup \dots \cup \mathcal{I}_d \cup \mathcal{V} = \{A_1, \dots, A_k\} \quad (2.12)$$

**Definition 2.2.2 (Encoding Relation)** An encoding relation  $R$  of depth  $d$  is a relation whose schema is an encoding schema of depth  $d$  and whose instance satisfies the following constraints.

1.  $\mathbb{T}(R)$  satisfies the functional dependency  $\bar{\mathcal{I}}_{[1,d]} \rightarrow \bar{\mathcal{V}}$ .
2. If  $d = 0$ , then  $\mathbb{T}(R)$  contains precisely one tuple.

For convenience, we use the notation

$$R(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$$

to represent the schema of encoding relation  $R$ . When depicting an encoding relation graphically, we separate index levels with a single vertical line and separate output variables from index variables with a double vertical line. For example, Figure 2.6 depicts encoding relation  $R_1$  with schema  $R_1(X; Y; Z)$ , while Figure 2.7 depicts an encoding relation  $R_3$  with schema<sup>4</sup>  $R_3(A, B; C, D; C)$ .

Given an encoding relation  $R = (\mathbb{S}(R), \mathbb{T}(R))$  and a value  $\bar{a} \in \text{adom}(\bar{\mathcal{I}}_{[1,l-1]}, R)$  we use  $R[\bar{a}]$  to denote the sub-relation indexed by  $\bar{a}$ , which is an encoding relation  $(\mathbb{S}(R[\bar{a}]), \mathbb{T}(R[\bar{a}]))$  defined as

$$\mathbb{S}(R[\bar{a}]) := R(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}}) \quad (2.13)$$

$$\mathbb{T}(R[\bar{a}]) := \Pi_{\bar{\mathcal{I}}_{[l,d]} \bar{\mathcal{V}}}(\sigma_{\bar{\mathcal{I}}_{[1,l-1]} = \bar{a}}(R)) \quad (2.14)$$

where  $\Pi$  and  $\sigma$  are relational projection and selection, respectively. Intuitively, if  $R$  encodes an object of depth  $d$ , then  $R[\bar{a}]$  encodes a sub-object of depth  $d - l + 1$  that

---

<sup>4</sup>At times we may refer to encoding relation  $R$  as having “relation schema  $R(\bar{\mathcal{I}}_{[1,d]} \bar{\mathcal{V}})$ ” or containing “tuples over  $\text{dom}(\bar{\mathcal{I}}_{[1,d]} \bar{\mathcal{V}})$ .” This is technically incorrect because tables (and hence relations) are required to have distinct attribute names, while the sequence  $\bar{\mathcal{I}}_{[1,d]} \bar{\mathcal{V}}$  allows repetitions. For example,  $R_3$  in Figure 2.7 appears to be a *quinary* relation with schema  $R_3(ABCDC)$ , but it is properly a *quaternary* relation over the attributes  $\{A, B, C, D\}$ . This distinction is moot, however, because we can treat  $R_3$  as a set of quinary tuples over  $\text{dom}(\bar{\mathcal{I}}_{[1,d]} \bar{\mathcal{V}}) = \text{dom}(ABCDC)$  since agreement of values for repeated attributes is guaranteed by equation 2.8.



A	B	C	D	C
$a_1$	$b_1$	1	$d_1$	1
$a_1$	$b_1$	2	$d_1$	2
$a_1$	$b_2$	1	$d_1$	1
$a_1$	$b_2$	2	$d_1$	2
$a_2$	$b_1$	1	$d_3$	1
$a_2$	$b_1$	1	$d_4$	1
$a_2$	$b_1$	2	$d_3$	2
$a_2$	$b_2$	1	$d_3$	1
$a_2$	$b_2$	2	$d_3$	2
$a_3$	$b_1$	1	$d_5$	1
$a_3$	$b_1$	3	$d_5$	3
$a_3$	$b_1$	3	$d_6$	3
$a_3$	$b_2$	1	$d_5$	1
$a_3$	$b_2$	3	$d_5$	3
$a_3$	$b_2$	3	$d_6$	3

C	D	C
1	$d_3$	1
1	$d_4$	1
2	$d_3$	2

C
1

Figure 2.7: Encoding relation  $R_3$  and select sub-relations

occurs as a descendant within the object encoded by  $R$ . For example, Figure 2.7 also shows the encoding relations  $R_3[a_2b_1]$  and  $R_3[a_2b_1d_3]$  derived from  $R_3$ .

We say that an encoding relation and a signature are *compatible* if the length of the signature matches the depth of the encoding schema. Each combination of an encoding relation and a compatible signature decodes to a unique complete chain object. Given an encoding relation  $R$  of depth  $d$  and a signature  $\bar{\xi}$  of length  $d$ , the function  $\text{DECODESIMPLE}(R, \bar{\xi})$  shown in Algorithm 2 computes a complete chain object of sort  $(\bar{\xi}, |\bar{\mathcal{V}}|)$ . If  $R$  is non-empty, then  $\text{DECODESIMPLE}(R, \bar{\xi})$  returns a hierarchy of  $d$  (possibly zero) levels of nested collections with tuples of arity  $|\bar{\mathcal{V}}|$  at the leaves. If  $R$  is empty, then Definition 2.2.2 requires that  $d > 0$ , and so  $\text{DECODESIMPLE}(R, \bar{\xi})$  returns a single empty collection of type  $\xi_1$ . The structuring of Algorithm 2 into two procedures  $\text{DECODESIMPLE}$  and  $\text{INNERDECODESIMPLE}$  is not necessary, but it facilitates our extending of the algorithm to handle shading in the next subsection.

**Example 14** Decoding relation  $R_1$  (cf. Figure 2.6) with signature  $\text{ns}$  yields object  $o_4$ .

$$\begin{aligned}
\text{DECODESIMPLE}(R_1, \text{ns}) &= \text{INNERDECODESIMPLE}(R_1, \text{ns}) \\
&= \{ \text{INNERDECODESIMPLE}(R_1[x_1], \mathbf{s}), \\
&\quad \text{INNERDECODESIMPLE}(R_1[x_2], \mathbf{s}), \\
&\quad \text{INNERDECODESIMPLE}(R_1[x_3], \mathbf{s}) \} \\
&= \{ \{ \text{INNERDECODESIMPLE}(R_1[x_1y_1], \emptyset), \text{INNERDECODESIMPLE}(R_1[x_1y_2], \emptyset) \}, \\
&\quad \{ \text{INNERDECODESIMPLE}(R_1[x_2y_1], \emptyset), \text{INNERDECODESIMPLE}(R_1[x_2y_2], \emptyset) \}, \\
&\quad \{ \text{INNERDECODESIMPLE}(R_1[x_3y_1], \emptyset), \text{INNERDECODESIMPLE}(R_1[x_3y_2], \emptyset) \} \}
\end{aligned}$$



---

**Algorithm 2** Reconstructing an encoded chain object
 

---

DECODESIMPLE( $R, \bar{\xi}$ )

▷ **Input:** encoding relation  $R(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$   
 ▷ **Input:** compatible signature  $\bar{\xi}$   
 ▷ **Output:** complete or trivial chain object of sort  $(\bar{\xi}, |\bar{\mathcal{V}}|)$

- 1 **if**  $R$  is non-empty
- 2     **then return** INNERDECODESIMPLE( $R, \bar{\xi}$ )
- 3     **else**                     ▷  $d > 0$  and  $R$  encodes a trivial empty collection
- 4         **if**  $\bar{\xi} = s\bar{Y}$
- 5             **then return**  $\{ \}$
- 6         **elseif**  $\bar{\xi} = b\bar{Y}$
- 7             **then return**  $\{ \}$
- 8         **else**     ▷  $\bar{\xi} = n\bar{Y}$
- 9             **return**  $\{ \}$

INNERDECODESIMPLE( $R, \bar{\xi}$ )

▷ **Input:** non-empty encoding relation  $R(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$   
 ▷ **Input:** compatible signature  $\bar{\xi}$   
 ▷ **Output:** complete or trivial chain object of sort  $(\bar{\xi}, |\bar{\mathcal{V}}|)$

- 1 **if**  $d = 0$                      ▷  $R$  contains singleton tuple  $\langle a_1, \dots, a_{|\bar{\mathcal{V}}}| \rangle$
- 2     **then return**  $\langle a_1, \dots, a_{|\bar{\mathcal{V}}}| \rangle$
- 3     **elseif**  $\bar{\xi} = s\bar{Y}$
- 4         **then return**  $\{ \text{INNERDECODESIMPLE}(R[\bar{a}], \bar{Y}) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$
- 5     **elseif**  $\bar{\xi} = b\bar{Y}$
- 6         **then return**  $\{ \text{INNERDECODESIMPLE}(R[\bar{a}], \bar{Y}) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$
- 7     **else**     ▷  $\bar{\xi} = n\bar{Y}$
- 8         **return**  $\{ \text{INNERDECODESIMPLE}(R[\bar{a}], \bar{Y}) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$

---

$$\begin{aligned}
&= \{\!\!\{ \langle 1 \rangle, \langle 2 \rangle \}, \langle 1 \rangle, \langle 2 \rangle \}\!\!\} \\
&= o_4
\end{aligned}$$

We started this subsection by informally describing how to generate an encoding relation from a complete chain object (not containing shaded nodes), using  $o_4$  and  $R_1$  as an example. Clearly, our method could generate an infinite number of unique encoding relations for  $o_4$  simply by choosing different edge labels; however all of these relations would be isomorphic to  $R_1$  modulo renaming of index values. Because there is a one-to-one mapping between tuples in  $R_1$  and the leaves of  $o_4$ , there cannot exist another relation that contains fewer tuples than  $R_1$  and still decodes to  $o_4$ ; hence,  $R_1$  could be considered a “minimal encoding” of  $o_4$ . There are also infinitely many encodings of  $o_4$  that are not isomorphic to  $R_1$  and hence contain varying amounts of redundancy with respect to sort  $\tau_4$ .

**Example 15** Relation  $R_3$  in Figure 2.7 also yields object  $o_4$  when decoded with signature  $\text{ns}$ .

$$\begin{aligned}
\text{DECODESIMPLE}(R_3, \text{ns}) &= \{\!\!\{ \langle 1 \rangle, \langle 2 \rangle \}, \langle 1 \rangle, \langle 2 \rangle \}, \{\!\!\{ \langle 1 \rangle, \langle 1 \rangle, \langle 2 \rangle \}, \\
&\quad \langle 1 \rangle, \langle 2 \rangle \}, \{\!\!\{ \langle 1 \rangle, \langle 3 \rangle, \langle 3 \rangle \}, \langle 1 \rangle, \langle 3 \rangle, \langle 3 \rangle \}\!\!\} \\
&= \{\!\!\{ \langle 1 \rangle, \langle 2 \rangle \}, \langle 1 \rangle, \langle 2 \rangle \}, \{\!\!\{ \langle 1 \rangle, \langle 2 \rangle \}, \\
&\quad \langle 1 \rangle, \langle 2 \rangle \}, \{\!\!\{ \langle 1 \rangle, \langle 3 \rangle \}, \langle 1 \rangle, \langle 3 \rangle \}\!\!\} \\
&= \{\!\!\{ \langle 1 \rangle, \langle 2 \rangle \}, \langle 1 \rangle, \langle 2 \rangle \}, \{\!\!\{ \langle 1 \rangle, \langle 3 \rangle \}\!\!\} \\
&= o_4
\end{aligned}$$

The reader can verify that when relation  $R_3$  is instead decoded with signature  $\text{sb}$  it yields object  $o_5$ .

## 2.2.2 Encoding and Decoding Linearized Objects

The *linearization* process of Section 2.1.2 outputs complete chain objects that may contain special “shaded” nodes (cf. Figure 2.4). The encoding and decoding procedures described in the previous subsection must be extended if we want to reconstruct linearized objects. From Observation 2.1.11 we know that the root node of a linearized object is never shaded; therefore, we can store information about a node’s shading within the label we assign to its *incoming* edge.

One approach would be to use an extra boolean attribute for each index level of an encoding schema to indicate shading. An alternative would be to assume the existence of a new domain of “shaded labels” distinguishable from unshaded labels. Because labels are only required to be locally-unique (i.e. unique among siblings), we can improve upon this second approach—from Observation 2.1.11 we know that a node can have at most one shaded child node, and so we need only one shaded label.

We extend the definition of an encoding relation as follows. Let the symbol  $\blacksquare$  denote a designated atomic element not occurring in the atomic domain  $\text{dom}$ .

A	B	C	D
$a_1$	$b_1$	$a$	1
$a_1$	$b_2$	$a$	2
$a_2$	$b_1$	$b$	1
$a_2$	$b_2$	$b$	3
$a_3$	$b_1$	$c$	■

A	B	C	D	E	F
$a_1$	■	$c_1$	$d_1$	■	■

A	B	C	D	E	F	G	H	I
$a_1$	$b_1$	$c_1$	■	$e_1$	$f_1$	$a$	■	5
$a_1$	$b_1$	$c_2$	■	$e_1$	$f_1$	$c$	■	5
$a_1$	$b_1$	$c_3$	■	$e_1$	$f_1$	$c$	■	5
$a_1$	$b_2$	$c_1$	■	$e_1$	$f_1$	$a$	■	5
$a_1$	$b_2$	$c_2$	■	$e_1$	$f_1$	$c$	■	5
$a_1$	$b_2$	$c_3$	■	$e_1$	$f_1$	$c$	■	5
$a_1$	$b_3$	■	■	$e_1$	$f_1$	■	■	5

Figure 2.8: Encoding relations  $R_4$ ,  $R_5$ , and  $R_6$

Similar to how we extended the underlying domain of complex objects to allow for shaded nodes, we slightly extend the underlying domain of encoding relations<sup>5</sup> to include the atomic value ■. The encoding process described in Section 2.2.1 now requires two minor adjustments:

1. The label for any given edge should be composed solely of (one or more occurrences of) the ■ symbol if and only if the edge leads into a shaded node.
2. Each occurrence of the atomic value `dom` should be replaced by ■ within the output attributes of the encoding relation.

**Example 16** Applying the modified encoding algorithm to objects  $\text{LINEARIZE}_{\tau_7}(o_7)$ ,  $\text{LINEARIZE}_{\tau_8}(o_8)$ , and  $\text{LINEARIZE}_{\tau_9}(o_9)$  from Figure 2.4 yields the encoding relations  $R_4$ ,  $R_5$ , and  $R_6$ , respectively, shown in Figure 2.8.

Algorithm 3 shows the modified decoding algorithm `DECODE`. The key modification is that in the recursive invocation to decode sub-relation  $R[\bar{a}]$ , we test the predicate  $\bar{a} = \text{■}^{|\bar{a}|}$  which returns `TRUE` when  $\bar{a}$  is a tuple composed *entirely* of the ■ symbol. The result of this test is passed to the decoding of  $R[\bar{a}]$  to indicate whether shading needs to be performed on the constructed sub-object. At the leaf level ( $d = 0$ ) this boolean parameter is ignored because tuple nodes of linearized objects are never shaded (see Observation 2.1.11); however, at this level the ■ symbols within the output are translated into the atomic value `dom` that occurs within the leaves of linearized objects.

<sup>5</sup>Technically, an encoding relation is defined as a special case of a flat table, and so we actually need to extend the underlying domain of tables. The important point is that ■ may not occur within the base relations of a database, or as an explicit constant within an object-constructing  $\text{COCQL}^{\mathcal{F}}$  query.

---

**Algorithm 3** Reconstructing an encoded *linearized* object
 

---

DECODE( $R, \bar{\xi}$ )

▷ **Input:** encoding relation  $R(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$   
 ▷ **Input:** compatible signature  $\bar{\xi}$   
 ▷ **Output:** complete or trivial chain object of sort  $(\bar{\xi}, |\bar{\mathcal{V}}|)$ ,  
 possibly containing shaded nodes

- 1 **if**  $R$  is non-empty
- 2   **then return** INNERDECODE( $R, \bar{\xi}, \text{FALSE}$ )
- 3   **else**                   ▷  $d > 0$  and  $R$  encodes a trivial empty collection
- 4       **if**  $\bar{\xi} = \text{s}\bar{Y}$
- 5           **then return**  $\{ \}$
- 6       **elseif**  $\bar{\xi} = \text{b}\bar{Y}$
- 7           **then return**  $\{ \}$
- 8       **else**   ▷  $\bar{\xi} = \text{n}\bar{Y}$
- 9           **return**  $\{ \}$

INNERDECODE( $R, \bar{\xi}, \text{shade}$ )

▷ **Input:** non-empty encoding relation  $R(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$   
 ▷ **Input:** compatible signature  $\bar{\xi}$   
 ▷ **Input:** boolean *shade* indicating if shading is required  
 ▷ **Output:** complete or trivial chain object of sort  $(\bar{\xi}, |\bar{\mathcal{V}}|)$ ,  
 possibly containing shaded nodes

- 1 **if**  $d = 0$                    ▷  $R$  contains singleton tuple  $\langle a_1, \dots, a_{|\bar{\mathcal{V}}}| \rangle$
- 2   **then**  $t \leftarrow \langle a_1, \dots, a_{|\bar{\mathcal{V}}}| \rangle$
- 3       Replace each occurrence of  $\blacksquare$  within  $t$  with **dom**
- 4       **return**  $t$
- 5 **elseif**  $\bar{\xi} = \text{s}\bar{Y}$
- 6    **then if**  $\text{shade} = \text{TRUE}$
- 7       **then return**  $\{ \text{INNERDECODE}(R[\bar{a}], \bar{Y}, (\bar{a} = \blacksquare^{|\bar{a}})) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$
- 8       **else return**  $\{ \text{INNERDECODE}(R[\bar{a}], \bar{Y}, (\bar{a} = \blacksquare^{|\bar{a}})) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$
- 9 **elseif**  $\bar{\xi} = \text{b}\bar{Y}$
- 10 **then if**  $\text{shade} = \text{TRUE}$
- 11    **then return**  $\{ \text{INNERDECODE}(R[\bar{a}], \bar{Y}, (\bar{a} = \blacksquare^{|\bar{a}})) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$
- 12    **else return**  $\{ \text{INNERDECODE}(R[\bar{a}], \bar{Y}, (\bar{a} = \blacksquare^{|\bar{a}})) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$
- 13 **else**   ▷  $\bar{\xi} = \text{n}\bar{Y}$
- 14    **if**  $\text{shade} = \text{TRUE}$
- 15       **then return**  $\{ \text{INNERDECODE}(R[\bar{a}], \bar{Y}, (\bar{a} = \blacksquare^{|\bar{a}})) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$
- 16       **else return**  $\{ \text{INNERDECODE}(R[\bar{a}], \bar{Y}, (\bar{a} = \blacksquare^{|\bar{a}})) \mid \bar{a} \in \Pi_{\bar{\mathcal{I}}_1}(R) \}$

---

	S	T	U	V	W	X	Y	Z
$s_1$	■	■	■	$w_1$	$x_1$		■	■
$s_2$	■	■	$v_1$	$w_1$	$x_1$		■	■
$s_2$	■	■	$v_1$	■	$x_1$		■	■
$s_2$	■	■	$v_1$	■	■		■	■

Figure 2.9: Encoding relation  $R_7$

**Example 17** The reader can verify the following for the relations shown in Figure 2.8.

$$\begin{aligned} \text{DECODE}(R_4, \mathbf{bs}) &= \text{LINEARIZE}_{\tau_7}(o_7) \\ \text{DECODE}(R_5, \mathbf{sbss}) &= \text{LINEARIZE}_{\tau_8}(o_8) \\ \text{DECODE}(R_6, \mathbf{sbnsbn}) &= \text{LINEARIZE}_{\tau_8}(o_8) \end{aligned}$$

Decoding relation  $R_7$  in Figure 2.9 with signature  $\mathbf{sbss}$  also yields object  $\text{LINEARIZE}_{\tau_8}(o_8)$ .

$$\begin{aligned} \text{DECODE}(R_7, \mathbf{sbss}) &= \text{INNERDECODE}(R_7, \mathbf{sbss}, \text{FALSE}) \\ &= \{ \text{INNERDECODE}(R_7[s_1], \mathbf{bss}, \text{FALSE}), \\ &\quad \text{INNERDECODE}(R_7[s_2], \mathbf{bss}, \text{FALSE}) \} \\ &= \{ \{ \text{INNERDECODE}(R_7[s_1--], \mathbf{ss}, \text{TRUE}) \}, \\ &\quad \{ \text{INNERDECODE}(R_7[s_2--], \mathbf{ss}, \text{TRUE}) \} \} \\ &= \{ \{ \{ \text{INNERDECODE}(R_7[s_1---w_1], \mathbf{s}, \text{FALSE}) \} \}, \\ &\quad \{ \{ \text{INNERDECODE}(R_7[s_2--v_1w_1], \mathbf{s}, \text{FALSE}), \\ &\quad \text{INNERDECODE}(R_7[s_2--v_1], \mathbf{s}, \text{FALSE}) \} \} \} \\ &= \{ \{ \{ \{ \text{INNERDECODE}(R_7[s_1---w_1x_1], \emptyset, \text{FALSE}) \} \} \}, \\ &\quad \{ \{ \{ \text{INNERDECODE}(R_7[s_2--v_1w_1x_1], \emptyset, \text{FALSE}) \}, \\ &\quad \text{INNERDECODE}(R_7[s_2--v_1x_1], \emptyset, \text{FALSE}), \\ &\quad \text{INNERDECODE}(R_7[s_2--v_1--], \emptyset, \text{TRUE}) \} \} \} \} \\ &= \{ \{ \{ \{ \{ \mathbf{dom}, \mathbf{dom} \} \} \} \}, \\ &\quad \{ \{ \{ \{ \mathbf{dom}, \mathbf{dom} \}, \{ \mathbf{dom}, \mathbf{dom} \}, \{ \mathbf{dom}, \mathbf{dom} \} \} \} \} \} \\ &= \{ \{ \{ \{ \{ \mathbf{dom}, \mathbf{dom} \} \} \} \} \} \\ &= \text{LINEARIZE}_{\tau_8}(o_8) \end{aligned}$$

Algorithm 3 is clearly backwards-compatible with Algorithm 2. That is, for any encoding relation  $R$  and compatible signature  $\bar{\xi}$  such that  $R$  does not contain the symbol ■,  $\text{DECODE}(R, \bar{\xi}) = \text{DECODESIMPLE}(R, \bar{\xi})$ . Therefore, we will not discuss Algorithm 2 and any future mention of a decoding algorithm is implicitly a reference to Algorithm 3.

### 2.2.3 Encoding-Equality

In the previous section we defined a decoding algorithm that takes an encoding relation and a compatible signature as inputs and returns a complete chain object.

Examples 14 and 15 illustrate two encoding relations  $R_1$  and  $R_3$  that, when decoded with signature  $\mathbf{ns}$ , each yield object  $o_4$ . We would like to recognize that  $R_1$  and  $R_3$  are in a sense “equal” because they encode the same object; however, this notion of equality is dependent upon the signature  $\mathbf{ns}$  used in the decoding. When decoded with signature  $\mathbf{sb}$ ,  $R_3$  yields  $o_5$  but  $R_1$  does not.

**Definition 2.2.3 (Encoding-Equality)** *Given two encoding relations  $R, R'$  and a compatible signature  $\bar{\xi}$ , if*

$$\text{DECODE}(R, \bar{\xi}) = \text{DECODE}(R', \bar{\xi})$$

*then we say that  $R$  and  $R'$  are encoding-equal w.r.t. signature  $\bar{\xi}$  (abbreviated  $\bar{\xi}$ -equal, and denoted  $R \doteq_{\bar{\xi}} R'$ ).*

We use  $R \doteq_{\bar{\xi}} R'$  to denote  $\bar{\xi}$ -equality, as distinct from  $R = R'$  which—as an extension of traditional relational equality—implies that  $\mathcal{S}(R)$  and  $\mathcal{S}(R')$  have identical arities of index attributes at each index level and of output attributes, and that  $\mathbb{T}(R) = \mathbb{T}(R')$ . Observe that  $R \doteq_{\bar{\xi}} R'$  implies that the encoding schemas  $\mathcal{S}(R)$  and  $\mathcal{S}(R')$  both have depth  $|\bar{\xi}|$  and the same arity of output variables; however, it does not require attribute names to be the same, nor does it require the arity of corresponding index levels to be the same. For example,  $\mathbf{s}$ -equality of encoding relations  $R(\bar{\mathcal{I}}_1; \bar{\mathcal{V}})$  and  $R'(\bar{\mathcal{I}}'_1, \bar{\mathcal{V}})$  corresponds to relational equality of *projections onto the output attributes*—that is,  $\Pi_{\bar{\mathcal{V}}}(R) = \Pi_{\bar{\mathcal{V}}}(R')$ , but it does not imply  $\mathbb{T}(R) = \mathbb{T}(R')$  because index values could differ.

While Definition 2.2.3 captures the desired semantics of  $\bar{\xi}$ -equality, reasoning with it is awkward due to the reliance on the `DECODE` procedure. We next introduce a mechanism called a  $\bar{\xi}$ -certificate that allows us to characterize  $\bar{\xi}$ -equality in a more declarative fashion. A  $\bar{\xi}$ -certificate is essentially a proof that the decoding of two relations with signature  $\bar{\xi}$  yields the same object.

**Definition 2.2.4 ( $\bar{\xi}$ -Certificate)** *Given encoding relations  $R$  and  $R'$  and signature  $\bar{\xi}$ , a  $\bar{\xi}$ -certificate is a tree rooted by a node  $n_{(R,R')}$  that proves  $R \doteq_{\bar{\xi}} R'$ . When  $\xi_1 = \mathbf{s}$  then  $n_{(R,R')}$  is a set node; when  $\xi_1 = \mathbf{b}$  then  $n_{(R,R')}$  is a bag node; when  $\xi_1 = \mathbf{n}$  then  $n_{(R,R')}$  is a normalized bag node; and when  $\bar{\xi} = \emptyset$  then  $n_{(R,R')}$  is a tuple node.*

**Definition 2.2.5 (Set Node)** *A set node  $n_{(R,R')}^{\mathbf{s}}$  proves  $R \doteq_{\mathbf{s}\bar{\mathcal{Y}}} R'$  for some  $R, R'$ , and  $\bar{\mathcal{Y}}$ . It contains a function  $f : \text{adom}(\bar{\mathcal{I}}'_1, R') \rightarrow \text{adom}(\bar{\mathcal{I}}_1, R)$  satisfying*

$$\forall \bar{x}' \in \text{adom}(\bar{\mathcal{I}}'_1) \left[ \begin{array}{l} (R[f(\bar{x}')] \doteq_{\bar{\mathcal{Y}}} R'[\bar{x}']) \wedge \\ (\bar{x}' = \blacksquare^{|\bar{x}'|} \iff f(\bar{x}') = \blacksquare^{|f(\bar{x}')|}) \end{array} \right] \quad (2.15)$$

*and a function  $f' : \text{adom}(\bar{\mathcal{I}}_1) \rightarrow \text{adom}(\bar{\mathcal{I}}'_1)$  satisfying*

$$\forall \bar{x} \in \text{adom}(\bar{\mathcal{I}}_1) \left[ \begin{array}{l} (R[\bar{x}] \doteq_{\bar{\mathcal{Y}}} R'[f'(\bar{x})]) \wedge \\ (\bar{x} = \blacksquare^{|\bar{x}|} \iff f'(\bar{x}) = \blacksquare^{|f'(\bar{x})|}) \end{array} \right] \quad (2.16)$$

For each pair  $(\bar{x}, \bar{x}')$  such that either  $\bar{x}' = f(\bar{x})$  or  $\bar{x} = f(\bar{x}')$ , node  $n_{(R,R')}^s$  has a child  $\bar{Y}$ -certificate that proves  $R[\bar{x}] \doteq_{\bar{Y}} R'[\bar{x}']$ .

A set node needs to guarantee that each child object in one of the decodings matches a child object in the other decoding; it does this as follows. For each (non-empty) sub-relation  $R'[\bar{x}']$  of  $R'$ , the first conjunct in formula 2.15 enforces that

$$\begin{aligned} \text{DECODE}(R'[\bar{x}'], \bar{Y}) &= \text{INNERDECODE}(R'[\bar{x}'], \bar{Y}, \text{FALSE}) \\ &= \text{INNERDECODE}(R[f(\bar{x}')], \bar{Y}, \text{FALSE}) \\ &= \text{DECODE}(R[f(\bar{x}')], \bar{Y}) \end{aligned}$$

for some (non-empty) sub-relation  $R[f(\bar{x}')]$  of  $R$ . Observe that in the decoding of  $R$  and  $R'$  by Algorithm 3, function  $\text{DECODE}$  is never invoked on either  $R'[\bar{x}']$  or  $R[f(\bar{x}')]$ ; the actual invocations would be  $\text{INNERDECODE}(R'[\bar{x}'], \bar{Y}, (\bar{x}' = \blacksquare^{|\bar{x}'|}))$  and  $\text{INNERDECODE}(R[f(\bar{x}')], \bar{Y}, (f(\bar{x}') = \blacksquare^{|f(\bar{x}')|}))$ , respectively. Therefore, the second conjunct in formula 2.15 is needed to ensure that the predicates  $\bar{x}' = \blacksquare^{|\bar{x}'|}$  and  $f(\bar{x}') = \blacksquare^{|f(\bar{x}')|}$  agree, so that the roots of the decoded sub-objects are either both or neither shaded. Formula 2.15 thus guarantees  $\text{DECODE}(R', \mathbf{s}\bar{Y}) \subseteq \text{DECODE}(R, \mathbf{s}\bar{Y})$ . Formula 2.16 likewise guarantees  $\text{DECODE}(R, \mathbf{s}\bar{Y}) \subseteq \text{DECODE}(R', \mathbf{s}\bar{Y})$ .

**Definition 2.2.6 (Bag Node)** A bag node  $n_{(R,R')}^b$  proves  $R \doteq_{\mathbf{b}\bar{Y}} R'$  for some  $R$ ,  $R'$ , and  $\bar{Y}$ . It contains a bijective function  $f : \text{adom}(\bar{\mathcal{I}}'_1, R') \rightarrow \text{adom}(\bar{\mathcal{I}}_1, R)$  satisfying

$$\forall \bar{x}' \in \text{adom}(\bar{\mathcal{I}}'_1) \left[ \begin{array}{l} (R[f(\bar{x}')] \doteq_{\bar{Y}} R'[\bar{x}']) \wedge \\ (\bar{x}' = \blacksquare^{|\bar{x}'|} \iff f(\bar{x}') = \blacksquare^{|f(\bar{x}')|}) \end{array} \right] \quad (2.17)$$

For each pair  $(\bar{x}, \bar{x}')$  such that  $\bar{x} = f(\bar{x}')$ , node  $n_{(R,R')}^b$  has a child  $\bar{Y}$ -certificate that proves  $R[\bar{x}] \doteq_{\bar{Y}} R'[\bar{x}']$ .

A bag node needs to guarantee that equal sub-objects occur with the same *absolute frequency* in both decodings. The bag node functions similarly to a set node; the main difference is that by requiring  $f$  to be a *bijection* between  $\text{adom}(\bar{\mathcal{I}}'_1, R')$  and  $\text{adom}(\bar{\mathcal{I}}_1, R)$ , we enforce that equal sub-objects of  $\text{DECODE}(R, \mathbf{s}\bar{Y})$  and  $\text{DECODE}(R', \mathbf{s}\bar{Y})$  can be paired up, guaranteeing identical absolute frequencies.

**Definition 2.2.7 (Normalized Bag Node)** A normalized bag node  $n_{(R,R')}^n$  proves  $R \doteq_{\mathbf{n}\bar{Y}} R'$  for some  $R$ ,  $R'$ , and  $\bar{Y}$ . It contains two finite domains  $D_1$  and  $D_2$ , and two surjective functions  $\rho : \text{adom}(\bar{\mathcal{I}}_1, R) \rightarrow D_1$  and  $\varrho : \text{adom}(\bar{\mathcal{I}}'_1, R') \rightarrow D_2$  that satisfy

$$\forall p \in D_1. \forall q \in D_2 [(R)^p \doteq_{\mathbf{b}\bar{Y}} (R')^q]$$

where  $(R)^p$  is the (non-empty) encoding relation defined by  $\mathbb{S}((R)^p) := \mathbb{S}(R)$  and  $\mathbb{T}((R)^p) := \sigma_{\rho(\bar{\mathcal{I}}_1)=p}(R)$ , and  $(R')^q$  is defined analogously using  $\varrho$ . For each pair  $(p, q)$  such that  $p \in D_1$  and  $q \in D_2$ , node  $n_{(R,R')}^n$  has a child  $\mathbf{b}\bar{Y}$ -certificate that proves  $(R)^p \doteq_{\mathbf{b}\bar{Y}} (R')^q$ .

A normalized bag node needs to guarantee that equal sub-objects occur with the same *relative frequency* in both decodings. To understand how this is enforced, let  $B_R$  and  $B_{R'}$  denote the *bags* (of sub-objects) obtained by decoding relations  $R$  and  $R'$ , respectively, with signature  $\mathbf{b}\bar{Y}$ . Let  $B \otimes n$  denote the *n-expansion of bag*  $B$ , which is the bag obtained from  $B$  by multiplying the cardinality of each element by  $n$  [23]. Definition 2.2.7 ensures that

$$B_R \otimes |D_2| = B_{R'} \otimes |D_1| \quad (2.18)$$

which implies that  $B_R$  and  $B_{R'}$  have the same normalized bag projection.

**Definition 2.2.8 (Tuple Node)** *A tuple node  $n_{(R,R')}^{\dagger}$  proves  $R \doteq_{\emptyset} R'$  for some  $R$  and  $R'$  with depth zero. By Definition 2.2.2, an encoding relation of depth zero contains precisely one tuple (of only output values). Therefore, node  $n_{(R,R')}^{\dagger}$  contains a single comparison showing that  $R$  and  $R'$  contain the same tuple.*

The main result of this subsection is summarized by the next theorem, which is a result of how Definitions 2.2.3–2.2.8 have been constructed to mirror the semantics of Algorithm 3.

**Theorem 2.2.9** *Given two encoding relations  $R$  and  $R'$  and a compatible signature  $\bar{\xi}$ ,  $R \doteq_{\bar{\xi}} R'$  if and only if there exists a  $\bar{\xi}$ -certificate between  $R$  and  $R'$ .*

**Example 18** Consider encoding relations  $R_1$  (cf. Figure 2.6) and  $R_3$  (cf. Figure 2.7). Figure 2.10 depicts a (partial) ns-certificate proving that  $R_1 \doteq_{\mathbf{ns}} R_3$ .

Observe that the nodes within a  $\bar{\xi}$ -certificate are organized into levels of identical node type. For any index level  $i$ , if  $\bar{\xi}_i = \mathbf{b}$  then “the nodes in (bag) level  $i$ ” are the equi-depth bag nodes that each prove  $\mathbf{b}\bar{\xi}_{[i+1,d]}$ -equality of some pair of sub-relations. Similarly, if  $\bar{\xi}_i = \mathbf{s}$  then “the nodes in (set) level  $i$ ” are the equi-depth set nodes that prove  $\mathbf{s}\bar{\xi}_{[i+1,d]}$ -equality. If  $\bar{\xi}_i = \mathbf{n}$  then “the nodes in (normalized bag) level  $i$ ” are the equi-depth normalized bag nodes that prove  $\mathbf{n}\bar{\xi}_{[i+1,d]}$ -equality—note, however, that immediately underneath a normalized bag level is an extra layer of bag nodes proving  $\mathbf{b}\bar{\xi}_{[i+1,d]}$ -equality, and so when we refer to “nodes at level  $i$ ” it is *not* the same as tree depth  $i$ . For example, in Figure 2.10 the “nodes at level 1” refers to the normalized bag node at the root (tree depth 0), while the “nodes at level 2” refers to the set nodes at tree depth 2. Finally, the leaves of a certificate are all equi-depth and guaranteed to be tuple nodes.

## 2.3 Encoding Queries

An “encoding query” is a query that outputs a (flat) encoding relation. The goal of this section is to reduce the problem of deciding  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  query equivalence (see Definition 2.1.21) to deciding “encoding-equivalence” between encoding queries. In



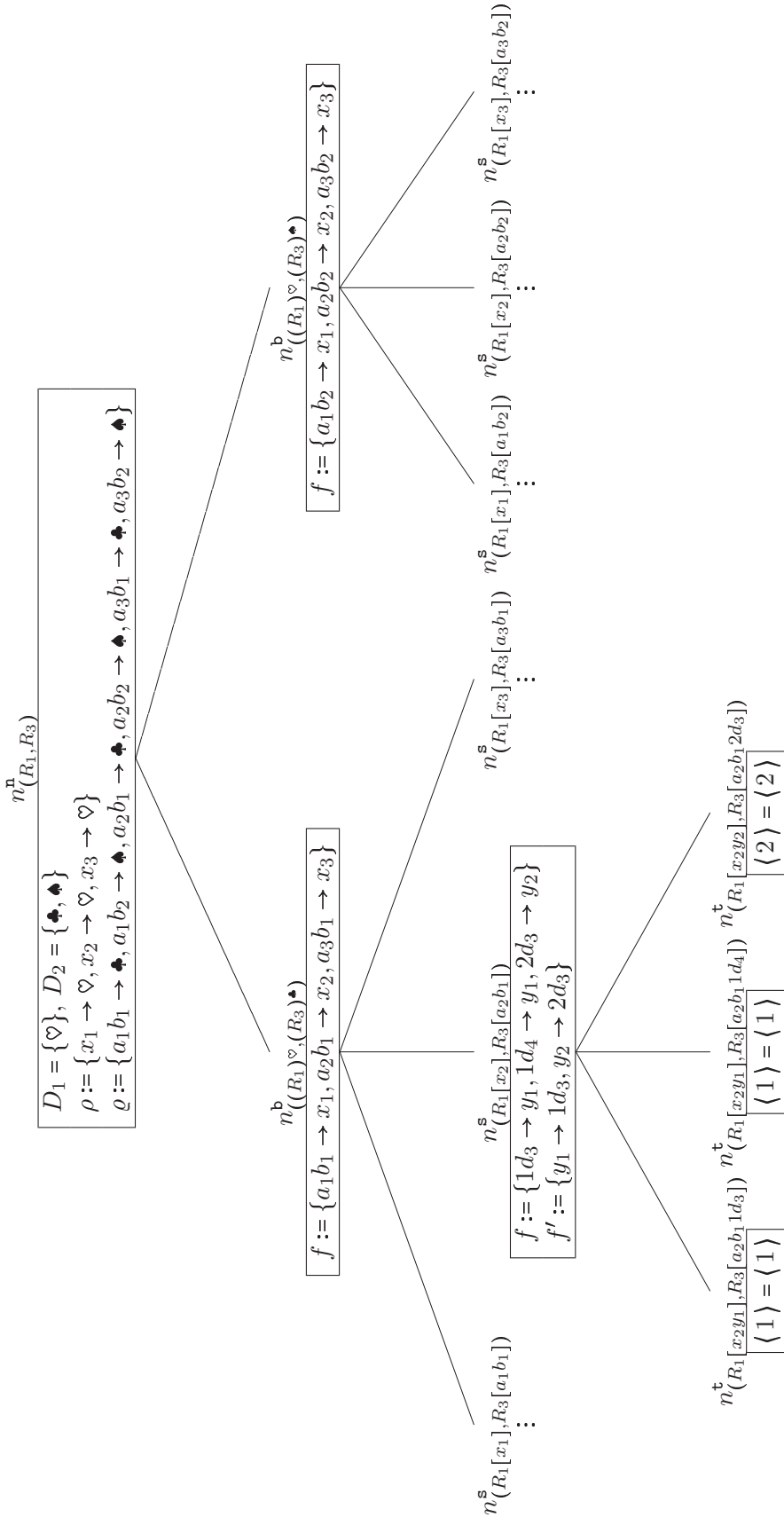


Figure 2.10: ns-certificate proving  $R_1 \stackrel{\text{ns}}{=} R_3$

Section 2.3.1 we define a language for specifying encoding queries. In Section 2.3.2 we define a mapping from each  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  query  $Q$  to an encoding query  $Q'$  such that over any relational database  $D$ , the output of  $Q'$  is an encoding of the linearization of the object output by  $Q$ .

### 2.3.1 Hierarchical Conjunctive Encoding Queries

Our language for specifying encoding queries is based on the traditional rule-based syntax of conjunctive queries. Our primary extension is to allow the right-hand side of the rule to be a hierarchy of conjunctive query bodies. Although there is a connection between our hierarchical syntax and the well-known *left outer join* algebra operator, we defer that discussion until Chapter 4 and restrict this section to defining our language for encoding queries. In the following definitions, let  $\mathbb{D}$  be a relational database with schema  $\mathbb{S}(\mathbb{D})$ .

**Definition 2.3.1 (Conjunctive Query with Equality (CQ<sup>=</sup>))** A conjunctive query with equality  $Q$  is a rule of the form:

$$Q(\overline{\mathcal{W}}) :- R_1(\overline{X}_1), \dots, R_m(\overline{X}_m), Y_1 = Z_1, \dots, Y_n = Z_n \quad (2.19)$$

where  $\overline{\mathcal{W}}$  is a tuple composed of constants and variables, each  $R_i$  is a relation name in  $\mathbb{S}(\mathbb{D})$ , each  $\overline{X}_i$  is a tuple composed of constants and variables, and each  $Y_i$  and  $Z_i$  is either a constant or a variable. When  $n = 0$  query  $Q$  is simply called a conjunctive query (CQ).

Expression  $Q(\overline{\mathcal{W}})$  is termed the query *head*, each expression  $R_i(\overline{X}_i)$  is called an *atom*, and each expression  $Y_j = Z_j$  is called a *predicate*. The *body* of  $Q$ —denoted  $\text{body}_Q$ —is the set of atoms  $\text{atoms}_Q$  and predicates  $\text{pred}_Q$  in  $Q$ . We define the following sets of constants and variables:

- $\mathcal{B}$  denotes the set of variables occurring anywhere in  $Q$ .
- $\mathcal{W}$  denotes the set of variables and constants occurring within head tuple  $\overline{\mathcal{W}}$ .
- $\mathcal{C}$  denotes an infinite set of constants, and  $\mathcal{C}_Q \subset \mathcal{C}$  is the set of constants occurring anywhere in  $Q$ .
- $\mathcal{P} \subseteq \mathcal{B}$  is a distinguished set of variables called the *parameters*.
- $\mathcal{L} := \mathcal{B} \setminus \mathcal{P}$  denotes the *local* variables. To ensure query safety, we require that each local variable must occur within at least one atom.

Query  $Q$  is called *parameterized* when  $|\mathcal{P}| > 0$ . Query  $Q$  is called *boolean* if  $|\overline{\mathcal{W}}| = 0$ . Query  $Q$  is called *unsatisfiable* if  $\text{pred}_Q$  transitively imply a predicate of the form  $c_i = c_j$ , where  $c_i$  and  $c_j$  are distinct constants in  $\mathcal{C}$ , can be derived from  $\text{pred}_Q$  using transitivity. It is well-known that any parameter-free, satisfiable CQ with equality can be converted to an equivalent CQ (without equality) by a process of variable substitution. We therefore assume without loss of generality that that predicates within  $\text{pred}_Q$  only reference constants and parameter variables.

An *embedding of  $\text{body}_Q$  into  $\mathbb{D}$* —denoted  $\gamma : \text{body}_Q \rightarrow \mathbb{D}$ , but for conciseness we often write  $\gamma : Q \rightarrow \mathbb{D}$ —is a mapping  $\gamma : \mathcal{B} \rightarrow \text{dom}$  (extended with identity on the constants in  $\mathcal{C}$ ) such that  $\gamma$  satisfies every atom in  $\text{atoms}_Q$  and every predicate in  $\text{pred}_Q$ . Atom  $R_i(\overline{X}_i)$  is satisfied by  $\gamma$  if relation  $R_i \in \mathbb{D}$  and tuple  $\gamma(\overline{X}_i) \in \mathbb{T}(R_i)$ . Predicate  $Y_j = Z_j$  is satisfied by  $\gamma$  if  $\gamma$  maps  $Y_j$  and  $Z_j$  to the same element of  $\text{dom}$ .

The result of evaluating unparameterized conjunctive query  $Q$  over database  $\mathbb{D}$ —denoted  $Q^{\mathbb{D}}$ —is a flat relation  $(\mathbb{S}(Q^{\mathbb{D}}), \mathbb{T}(Q^{\mathbb{D}}))$ . Schema  $\mathbb{S}(Q^{\mathbb{D}})$  is deduced from the query head by associating with each element<sup>6</sup> of  $\overline{\mathcal{W}}$  a unique attribute name from  $\text{aname}$ . Instance  $\mathbb{T}(Q^{\mathbb{D}})$  is the set of tuples for which there exists a supporting embedding of  $\text{body}_Q$  into  $\mathbb{D}$ .

$$\mathbb{T}(Q^{\mathbb{D}}) := \{\gamma(\overline{\mathcal{W}}) \mid \gamma : Q \rightarrow \mathbb{D}\} \quad (2.20)$$

If  $Q$  is parameterized, then it can only be evaluated within the context of a specific assignment of values to the parameter variables. Let  $v : \mathcal{P} \rightarrow \text{dom}$  be any such assignment; then, the parameterized evaluation of  $Q$  over  $\mathbb{D}$  relative to  $v$  restricts the supporting embeddings to those that are extensions of  $v$ .

$$\mathbb{T}(Q^{\mathbb{D}}|_v) := \{\gamma(\overline{\mathcal{W}}) \mid \gamma : Q \rightarrow \mathbb{D} \wedge \gamma(\overline{\mathcal{P}}) = v(\overline{\mathcal{P}})\} \quad (2.21)$$

(Within equation 2.21,  $\overline{\mathcal{P}}$  denotes an arbitrary sequencing of set  $\mathcal{P}$  into a tuple.)

**Definition 2.3.2 (Hierarchical Conjunctive Query (HCQ))** *A hierarchical conjunctive query  $Q$  is a generalization of a parameter-free CQ, defined as follows.*

1.  $\text{body}_Q$  is a rooted (unordered) tree in which each vertex has a unique label, the root vertex contains an unparameterized CQ body, and the non-root vertices each contain a parameterized  $CQ^{\bar{}}$  body. For each vertex  $v \in \text{body}_Q$ , we use  $\text{body}_v$  to denote the  $CQ^{\bar{}}$  body occurring within vertex  $v$ ,  $\text{body}_v^{\text{anc}}$  to denote the  $CQ^{\bar{}}$  body formed by conjoining the  $CQ^{\bar{}}$  bodies of all of the ancestors of  $v$ , and  $\text{body}_v^+$  to denote the  $CQ^{\bar{}}$  body formed by conjoining  $\text{body}_v$  with  $\text{body}_v^{\text{anc}}$ ; the sets  $\text{atoms}_v$ ,  $\text{atoms}_v^{\text{anc}}$ ,  $\text{atoms}_v^+$ ,  $\text{pred}_v$ ,  $\text{pred}_v^{\text{anc}}$ , and  $\text{pred}_v^+$ , are defined analogously.
2. Sets  $\mathcal{B}$ ,  $\mathcal{C}_Q \subset \mathcal{C}$ ,  $\mathcal{P}$ , and  $\mathcal{L}$  are defined the same as for  $CQ^{\bar{}}$ s. Because the overall query  $Q$  is parameter-free, set  $\mathcal{P} = \emptyset$  and therefore  $\mathcal{L} = \mathcal{B}$ .
3. We use  $\mathcal{B}_v$  to denote the set of variables occurring anywhere in  $\text{body}_v$ , with  $\mathcal{B}_v^{\text{anc}}$  and  $\mathcal{B}_v^+$  likewise corresponding to  $\text{body}_v^{\text{anc}}$  and  $\text{body}_v^+$ . We use  $\mathcal{P}_v$  to denote the set of designated parameter variables in  $\text{body}_v$ , and we require that  $\mathcal{P}_v := \mathcal{B}_v \cap \mathcal{B}_v^{\text{anc}}$ . Similarly, we use  $\mathcal{L}_v := \mathcal{B}_v \setminus \mathcal{B}_v^{\text{anc}}$  to denote the set of local variables of  $\text{body}_v$ , and for any two distinct vertices  $u, v$  in  $\text{body}_Q$ , we require that local variables  $\mathcal{L}_u$  and  $\mathcal{L}_v$  be disjoint. As a consequence, for each variable

---

<sup>6</sup>Technically, the attribute names in a relation schema must be unique (per Definitions 2.1.12 and 2.1.13) and so if  $\overline{\mathcal{W}}$  contains repeated variables or constants, then each occurrence must be assigned a different attribute name. When discussing containment or equivalence of queries we will generally ignore attribute names, relying instead upon positional alignment of columns in the depiction of the schema or instance.

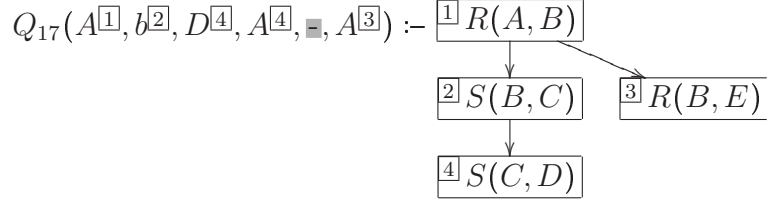


Figure 2.11: Hierarchical conjunctive query  $Q_{17}$

$X \in \mathcal{B}$  there exists precisely one node  $u$  such that  $X \in \mathcal{L}_u$ , and  $X$  only occurs within the subtree rooted at  $u$ . We use  $\mathcal{P}_v^{\text{anc}}$  and  $\mathcal{P}_v^+$  to denote the sets of “inter-query parameters” along the path from the root to  $v$ , with  $\mathcal{P}_Q^+$  denoting the set of all inter-query parameters within  $\text{body}_Q$ .

$$\mathcal{P}_v^{\text{anc}} := \bigcup_{u \in \text{ancestors}(v)} \mathcal{P}_u \quad \mathcal{P}_v^+ := \mathcal{P}_v^{\text{anc}} \cup \mathcal{P}_v \quad \mathcal{P}_Q^+ := \bigcup_{u \in \text{body}_Q} \mathcal{P}_u$$

(The analogous sets  $\mathcal{L}_v^{\text{anc}}$ ,  $\mathcal{L}_v^+$ , and  $\mathcal{L}_Q^+$  would be redundant, because the above definitions imply  $\mathcal{L}_v^{\text{anc}} = \mathcal{B}_v^{\text{anc}}$ ,  $\mathcal{L}_v^+ = \mathcal{B}_v^+$ , and  $\mathcal{L}_Q^+ = \mathcal{B}$ .)

4. The head tuple  $\overline{W}$  is composed from elements of  $\mathcal{B} \cup \mathcal{C} \cup \{\blacksquare\}$  (that is, variables, constants, and the special  $\blacksquare$  constant). We use  $\mathcal{W}$  to denote the set of elements occurring within  $\overline{W}$ .
5. Each component of head tuple  $\overline{W}$  that is not the  $\blacksquare$  constant is labelled with a vertex  $v \in \text{body}_Q$ . We write  $W_i^{\boxed{v}}$  to denote that the  $i^{\text{th}}$  component of  $\overline{W}$  is labelled with vertex  $v$ . If  $W_i \in \mathcal{C}$ , then  $v$  is allowed to be any vertex in  $\text{body}_Q$ . If  $W_i \in \mathcal{B}$ , then  $v$  must satisfy  $W_i \in \mathcal{B}_v^+$ .
6. We define the head label tuple  $\overline{L}$  as the tuple with the same arity as head tuple  $\overline{W}$  in which each position  $L_i$  contains the label of the corresponding head component  $W_i$  (or is empty if  $W_i = \blacksquare$ ). We define the head label set  $L$  as the set of labels occurring in sequence  $\overline{L}$ .

We again adopt the (non-restrictive) syntactic convention that for each vertex  $v$  in  $\text{body}_Q$  the equality predicates in  $\text{body}_v$  only reference  $\mathcal{C} \cup \mathcal{P}_v$  (that is, constants and variables that are parameters within  $v$ ).

**Example 19** Hierarchical conjunctive query  $Q_{17}$  is shown in Figure 2.11.  $Q_{17}$  has head label tuple  $\overline{L} = \langle \boxed{1}, \boxed{2}, \boxed{4}, \boxed{4}, \boxed{3} \rangle$  and head label set  $L = \{\boxed{1}, \boxed{2}, \boxed{3}, \boxed{4}\}$ .

The result of evaluating hierarchical conjunctive query  $Q$  over database  $\mathbb{D}$  is a relation  $Q^{\mathbb{D}} = (\mathcal{S}(Q^{\mathbb{D}}), \mathbb{T}(Q^{\mathbb{D}}))$  whose instance may contain the special atomic value  $\blacksquare$  not occurring in the domain of  $\mathbb{D}$  (see Section 2.2.2). Schema  $\mathcal{S}(Q^{\mathbb{D}})$  is defined the same as for conjunctive queries. Similar to conjunctive queries, instance  $\mathbb{T}(Q^{\mathbb{D}})$  is defined in terms of *terminal embeddings* of  $\text{body}_Q$  into  $\mathbb{D}$  which generalize the definition of  $\text{CQ}^-$  embeddings given above.

Let  $\text{trees}(T)$  denote the set of all non-empty trees that can be formed from some non-empty tree  $T$  by deleting zero or more of its strict subtrees; note that

every tree in  $\mathbf{trees}(T)$  contains at least the root of  $T$ . We now restrict our attention to trees in the set  $\mathbf{trees}(\mathbf{body}_Q)$ . For each such tree  $T$ , define  $\mathcal{B}_T := \cup_{v \in T} \mathcal{B}_v$  and define  $\mathcal{L}_T$ ,  $\mathbf{atoms}_T$ , and  $\mathbf{pred}_T$  analogously; it follows from Definition 2.3.2 that  $T$  is a parameter-free query body (i.e.,  $\mathcal{L}_T = \mathcal{B}_T$ ).

Given any  $T \in \mathbf{trees}(\mathbf{body}_Q)$ , let  $\mathbf{body}_T^+$  denote the  $\text{CQ}^-$  body formed by conjoining the  $\text{CQ}^-$  bodies of all vertices within  $T$ . An *embedding* of  $T \in \mathbf{trees}(\mathbf{body}_Q)$  into  $\mathbb{D}$ —denoted  $\gamma : T \rightarrow \mathbb{D}$ —is a mapping  $\gamma : \mathcal{B}_T \rightarrow \mathbf{dom}$  (extended with identity on the constants in  $\mathcal{C}$ ) that is an embedding of  $\mathbf{body}_T^+$  into  $\mathbb{D}$ . Observe that in the degenerate case where  $\mathbf{body}_Q$  is a single vertex, this definition of an embedding corresponds to the definition given earlier for  $\text{CQ}^-$ s. We use  $\mathbf{body}_Q^+$  to denote the single  $\text{CQ}^-$  formed by conjoining the  $\text{CQ}^-$ s from all of the vertices in  $\mathbf{body}_Q$ , which corresponds to the maximal case where  $T$  contains all of the vertices within  $\mathbf{body}_Q$ .

A *terminal embedding* of  $T \in \mathbf{trees}(\mathbf{body}_Q)$  into  $\mathbb{D}$ —denoted  $\gamma : T \rightsquigarrow \mathbb{D}$ —is an embedding  $\gamma : T \rightarrow \mathbb{D}$  for which there does not exist some other tree  $T' \in \mathbf{trees}(\mathbf{body}_Q)$  and some other embedding  $\gamma' : T' \rightarrow \mathbb{D}$  such that  $\gamma'$  is an extension of  $\gamma$ . Note that for  $\gamma'$  to be an extension of  $\gamma$ , tree  $T'$  must be an extension of tree  $T$ .

Given any embedding  $\gamma : T \rightarrow \mathbb{D}$  of some  $T \in \mathbf{trees}(\mathbf{body}_Q)$  and any component  $W_i$  of the head tuple  $\overline{W}$ , the *partial application of  $\gamma$  to  $W_i$* —denoted  $\hat{\gamma}(W_i)$ —is defined as follows.

$$\hat{\gamma}(W_i) := \begin{cases} \gamma(W_i) & \text{if } \exists v \in T \text{ such that } W_i \text{ is labelled with } v \\ \blacksquare & \text{otherwise (includes case } W_i = \blacksquare) \end{cases} \quad (2.22)$$

Instance  $\mathbb{T}(Q^{\mathbb{D}})$  is then defined as the set of tuples for which there exists a supporting partial application of some terminal embedding.

$$\mathbb{T}(Q^{\mathbb{D}}) := \{\hat{\gamma}(\overline{W}) \mid \exists T \in \mathbf{trees}(\mathbf{body}_Q), \gamma : T \rightsquigarrow \mathbb{D}\} \quad (2.23)$$

As a shorthand notation, we use  $\gamma : Q \rightsquigarrow \mathbb{D}$  to denote that  $\gamma$  is a terminal embeddings of  $\mathbf{body}_Q$ , and so Equation 2.23 can be rewritten as follows.

$$\mathbb{T}(Q^{\mathbb{D}}) := \{\hat{\gamma}(\overline{W}) \mid \gamma : Q \rightsquigarrow \mathbb{D}\} \quad (2.24)$$

**Example 20** Consider evaluating HCQ  $Q_{17}$  from Figure 2.11 over a database instance  $\mathbb{D}_2$  containing the following two relations  $R$  and  $S$ .

$R$	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th style="border-bottom: 1px solid black; padding: 2px 5px;"><math>col_1</math></th> <th style="border-bottom: 1px solid black; padding: 2px 5px;"><math>col_2</math></th> </tr> </thead> <tbody> <tr><td style="padding: 2px 5px;">1</td><td style="padding: 2px 5px;">2</td></tr> <tr><td style="padding: 2px 5px;">2</td><td style="padding: 2px 5px;">3</td></tr> <tr><td style="padding: 2px 5px;">2</td><td style="padding: 2px 5px;">4</td></tr> <tr><td style="padding: 2px 5px;">3</td><td style="padding: 2px 5px;">4</td></tr> <tr><td style="padding: 2px 5px;">3</td><td style="padding: 2px 5px;">5</td></tr> <tr><td style="padding: 2px 5px;">5</td><td style="padding: 2px 5px;">7</td></tr> <tr><td style="padding: 2px 5px;">6</td><td style="padding: 2px 5px;">8</td></tr> </tbody> </table>	$col_1$	$col_2$	1	2	2	3	2	4	3	4	3	5	5	7	6	8	$S$	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th style="border-bottom: 1px solid black; padding: 2px 5px;"><math>col_1</math></th> <th style="border-bottom: 1px solid black; padding: 2px 5px;"><math>col_2</math></th> </tr> </thead> <tbody> <tr><td style="padding: 2px 5px;">2</td><td style="padding: 2px 5px;">4</td></tr> <tr><td style="padding: 2px 5px;">4</td><td style="padding: 2px 5px;">5</td></tr> <tr><td style="padding: 2px 5px;">5</td><td style="padding: 2px 5px;">6</td></tr> <tr><td style="padding: 2px 5px;">7</td><td style="padding: 2px 5px;">8</td></tr> </tbody> </table>	$col_1$	$col_2$	2	4	4	5	5	6	7	8
$col_1$	$col_2$																												
1	2																												
2	3																												
2	4																												
3	4																												
3	5																												
5	7																												
6	8																												
$col_1$	$col_2$																												
2	4																												
4	5																												
5	6																												
7	8																												

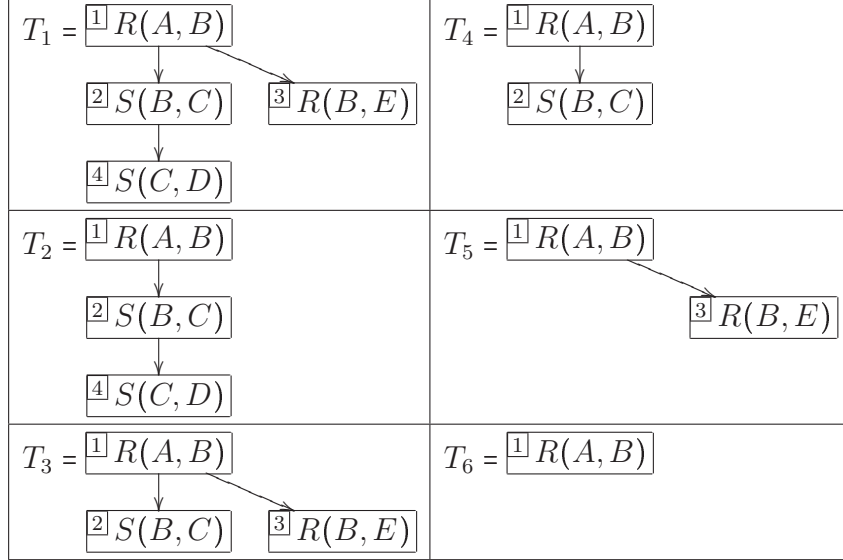


Figure 2.12: The six trees in  $\text{trees}(Q_{17})$

The set  $\text{trees}(Q_{17})$  contains the six trees  $T_1, \dots, T_6$  shown in Figure 2.12. For each  $i \in [1, 6]$ , let  $\Gamma^i$  denote the set of embeddings of  $T_i$  into  $\mathbb{D}_2$ . In Figure 2.13 we tabularly depict each  $\Gamma^i$  and we have shaded rows corresponding to terminal embeddings. Finally, let  $\Gamma$  denote the set of terminal embeddings of members of  $\text{trees}(\text{body}_Q)$  into  $\mathbb{D}_2$ . Relation  $(Q_{17})^{\mathbb{D}_2}$  is obtained by duplicate-eliminating projection of columns  $A^{\boxed{1}}$ ,  $b^{\boxed{2}}$ ,  $D^{\boxed{4}}$ ,  $A^{\boxed{4}}$ ,  $\blacksquare$ , and  $A^{\boxed{3}}$  from “relation”  $\Gamma$  shown in Figure 2.14.

**Definition 2.3.3 (Encoding Query (CEQ or HCEQ))** A hierarchical conjunctive encoding query  $Q$  of depth  $d \geq 1$  is an HCQ in which the head tuple  $\overline{W}$  is partitioned into  $d + 1$  sub-tuples as follows:

$$Q(\overline{\mathcal{I}}_1; \dots; \overline{\mathcal{I}}_d; \overline{\mathcal{V}})$$

In the degenerate case where  $\text{body}_Q$  is a single CQ body and  $\overline{W}$  does not contain the special constant  $\blacksquare$ , we call  $Q$  a conjunctive encoding query.

The result of evaluating encoding query  $Q$  over database  $\mathbb{D}$  is a relation  $Q^{\mathbb{D}} = (\mathcal{S}(Q^{\mathbb{D}}), \mathbb{T}(Q^{\mathbb{D}}))$  whose schema  $\mathcal{S}(Q^{\mathbb{D}})$  is an *encoding schema* (see Definition 2.2.1) deduced by associating a unique attribute name from  $\text{aname}$  with each occurrence of a variable or constant in the query head. Instance  $\mathbb{T}(Q^{\mathbb{D}})$  is defined the same as for HCQs.

We call  $\overline{\mathcal{V}}$  the “output tuple” because it corresponds to the output attributes in encoding schema  $\mathcal{S}(Q^{\mathbb{D}})$ ; we use  $\hat{\mathcal{V}}$  to denote the set of *labelled* constants/variables composing  $\overline{\mathcal{V}}$ , and we use  $\mathcal{V}$  to denote the set of *unlabelled* constants/variables composing  $\overline{\mathcal{V}}$  (i.e., for each  $X \in \mathcal{V}$  there exists some  $v$  such that  $X^{\boxed{v}}$  occurs in both  $\hat{\mathcal{V}}$  and  $\overline{\mathcal{V}}$ ). Likewise, for each  $i \in [1, d]$  we call  $\overline{\mathcal{I}}_i$  the “tuple of indexes at level  $i$ ” because it corresponds to the tuple of index attributes at level  $i$  in  $\mathcal{S}(Q^{\mathbb{D}})$ ; we use

$\Gamma^1$	A	B	C	D	E
$\gamma_1^1$	1	2	4	5	3
$\gamma_2^1$	1	2	4	5	4

$\Gamma^2$	A	B	C	D
$\gamma_1^2$	1	2	4	5
$\gamma_2^2$	2	4	5	6
$\gamma_3^2$	3	4	5	6

$\Gamma^3$	A	B	C	E
$\gamma_1^3$	1	2	4	3
$\gamma_2^3$	1	2	4	4
$\gamma_3^3$	3	5	6	7

$\Gamma^4$	A	B	C
$\gamma_1^4$	1	2	4
$\gamma_2^4$	2	4	5
$\gamma_3^4$	3	4	5
$\gamma_4^4$	3	5	6
$\gamma_5^4$	5	7	8

$\Gamma^5$	A	B	E
$\gamma_1^5$	1	2	3
$\gamma_2^5$	1	2	4
$\gamma_3^5$	2	3	4
$\gamma_4^5$	2	3	5
$\gamma_5^5$	3	5	7

$\Gamma^6$	A	B
$\gamma_1^6$	1	2
$\gamma_2^6$	2	3
$\gamma_3^6$	2	4
$\gamma_4^6$	3	4
$\gamma_5^6$	3	5
$\gamma_6^6$	5	7
$\gamma_7^6$	6	8

Figure 2.13: Embeddings of  $\text{trees}(Q_{17})$  into database  $\mathbb{D}_2$

$\Gamma$	A	B	C	D	E	b	$A^{[1]}$	$b^{[2]}$	$D^{[4]}$	$A^{[4]}$	-	$A^{[3]}$
$\gamma_1^1$	1	2	4	5	3	b	1	b	5	1	-	1
$\gamma_2^1$	1	2	4	5	4	b	1	b	5	1	-	1
$\gamma_3^5$	2	3			4	b	2	-	-	-	-	2
$\gamma_4^5$	2	3			5	b	2	-	-	-	-	2
$\gamma_2^2$	2	4	5	6		b	2	b	6	2	-	-
$\gamma_3^2$	3	4	5	6		b	3	b	6	3	-	-
$\gamma_3^3$	3	5	6		7	b	3	b	-	-	-	3
$\gamma_5^4$	5	7	8			b	5	b	-	-	-	-
$\gamma_7^6$	6	8				b	6	-	-	-	-	-

Figure 2.14: Terminal embeddings and their partial applications to  $\overline{\mathcal{W}}$

$\hat{\mathcal{I}}_i$  to denote the set of *labelled* constants/variables composing  $\bar{\mathcal{I}}_i$ , and we use  $\mathcal{I}_i$  to denote the set of *unlabelled* constants/variables composing  $\bar{\mathcal{I}}_i$ . We further define tuple  $\bar{\mathcal{I}}_{[i,j]}$  and sets  $\hat{\mathcal{I}}_{[i,j]}$  and  $\mathcal{I}_{[i,j]}$  as in Definition 2.2.1.

**Definition 2.3.4 (Validity)** *We call encoding query  $Q$  valid if for any database  $\mathbb{D}$  the relation  $Q^{\mathbb{D}}$  is a valid encoding relation (i.e., satisfies Definition 2.2.2).*

We are only interested in valid encoding queries because our definition of “encoding-equivalence” in the next sub-section presupposes that encoding queries always output encoding relations. The following lemma yields a sufficient condition for query validity that meets our needs for this chapter; we re-visit query validity in Chapters 3 and 5 when we address the problem of encoding-equivalence.

**Lemma 2.3.5** *Encoding query  $Q$  is valid if it satisfies the following conditions.*

1. *For every component  $V_i^{\square}$  of  $\bar{\mathcal{V}}$  where  $V_i \in \mathcal{B}$  (i.e.,  $V_i$  is a variable), index set  $\hat{\mathcal{I}}_{[1,d]}$  contains  $V_i^{\square}$ .*
2. *For every component  $V_i^{\square}$  of  $\bar{\mathcal{V}}$  where  $V_i \in \mathcal{C}$  (i.e.,  $V_i$  is a constant) and  $v$  is a non-root vertex of  $\mathbf{body}_Q$ , index set  $\hat{\mathcal{I}}_{[1,d]}$  contains an index labelled with vertex  $v$  (i.e.,  $\exists X.(X^{\square} \in \hat{\mathcal{I}}_{[1,d]})$ ).*

PROOF. Suppose that  $Q$  satisfies the two conditions above. Since Definition 2.3.3 requires that  $d \geq 1$ , by Definition 2.2.2 we must show that for any database  $\mathbb{D}$  the relation  $Q^{\mathbb{D}}$  satisfies the functional dependency  $\bar{\mathcal{I}}_{[1,d]} \rightarrow \bar{\mathcal{V}}$ .

Consider any two terminal embeddings  $\gamma_1, \gamma_2 : Q \rightsquigarrow \mathbb{D}$  that satisfy  $\hat{\gamma}_1(\bar{\mathcal{I}}_{[1,d]}) = \hat{\gamma}_2(\bar{\mathcal{I}}_{[1,d]})$ . Consider in turn each element  $V_i^{\square}$  of tuple  $\bar{\mathcal{V}}$ . If  $V_i$  is a variable then  $\hat{\gamma}_1(V_i^{\square}) = \hat{\gamma}_2(V_i^{\square})$  follows trivially from  $V_i^{\square} \in \hat{\mathcal{I}}_{[1,d]}$ . Otherwise,  $V_i$  must be a constant. If  $v$  is the root vertex of  $\mathbf{body}_Q$ , it follows from the definition of partial application that  $\hat{\gamma}_1(V_i^{\square}) = \hat{\gamma}_2(V_i^{\square}) = V_i$ . Otherwise  $v$  is a non-root vertex of  $\mathbf{body}_Q$ , and so  $\hat{\mathcal{I}}_{[1,d]}$  must contain some index  $X^{\square}$ . It again follows from the definition of partial application that if  $\hat{\gamma}(X^{\square}) \in \text{dom}$  then  $\hat{\gamma}_1(V_i^{\square}) = \hat{\gamma}_2(V_i^{\square}) = V_i$ , whereas if  $\hat{\gamma}(X^{\square}) = \blacksquare$  then  $\hat{\gamma}_1(V_i^{\square}) = \hat{\gamma}_2(V_i^{\square}) = \blacksquare$ .  $\square$

## 2.3.2 Reducing Equivalence to Encoding-Equivalence

In this section we reduce the problem of deciding equivalence of satisfiable unparameterized  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  queries to deciding “encoding-equivalence” of HCEQs. We will illustrate this process with a running example using the following queries and database instance.

**Example 21** Consider the two  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  queries shown in Figure 2.15. Following our convention, uppercase letters denote variables while integers and lowercase letters denote atomic constants; for reasons that will become clear shortly, we have also numbered the query nodes in which an aggregation function is applied. Over any database,



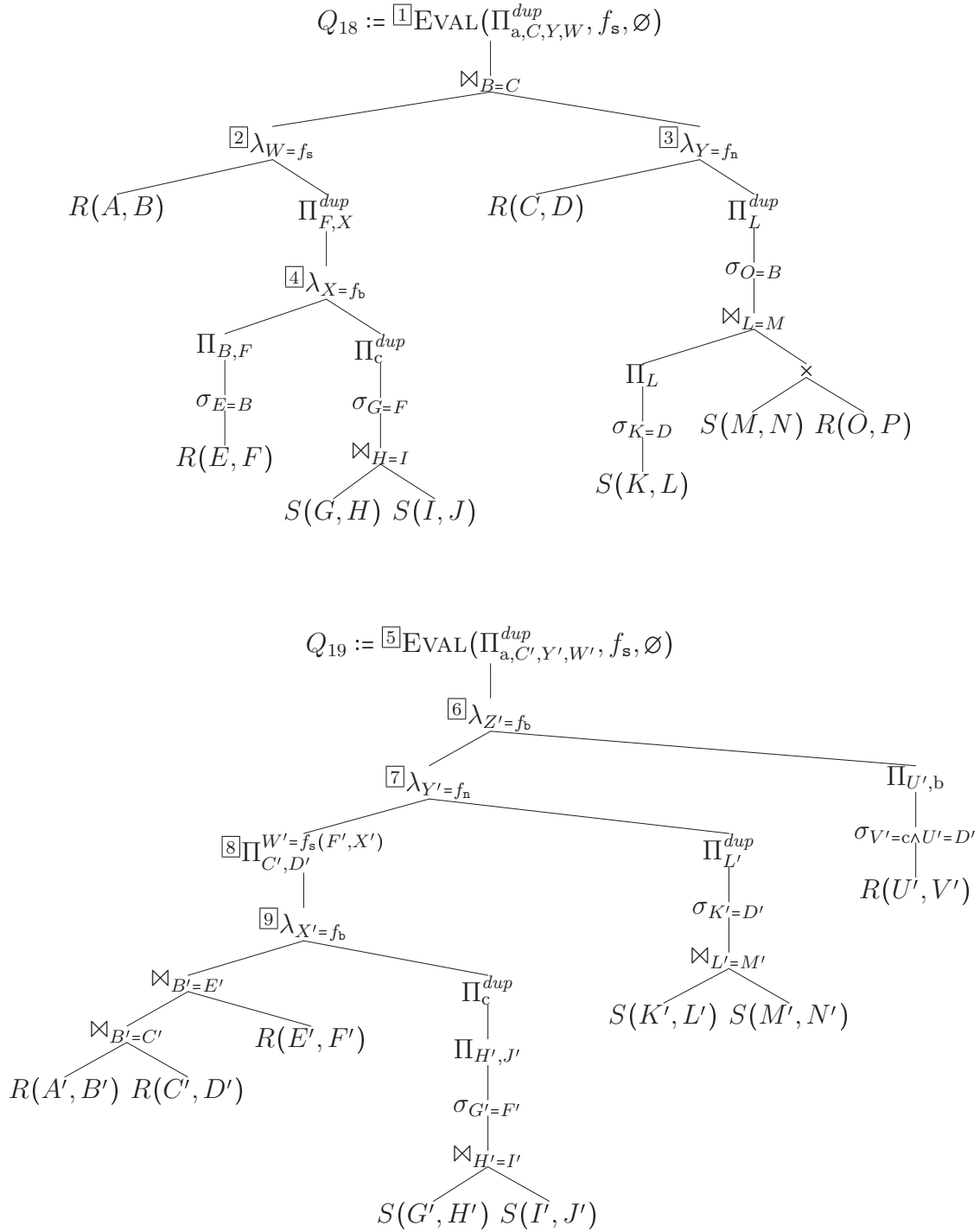


Figure 2.15: COCQL $\{f_b, f_s, f_n\}$  queries  $Q_{18}$  and  $Q_{19}$

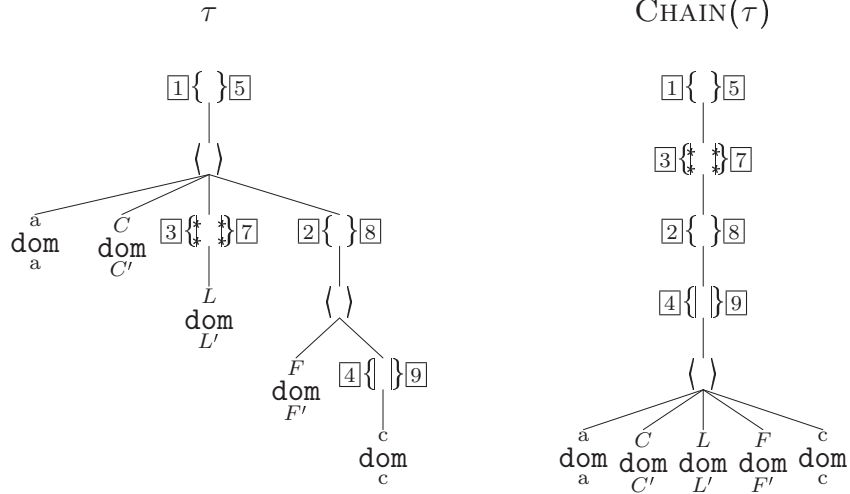
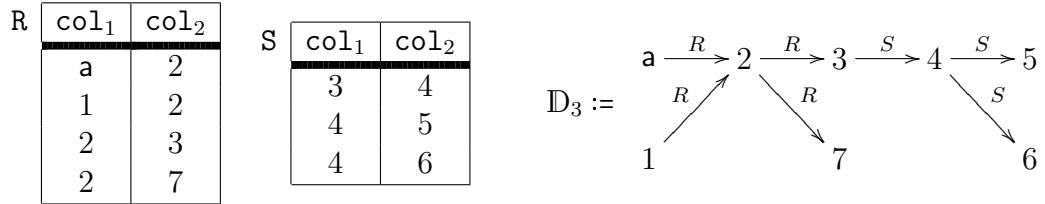


Figure 2.16: Output sort  $\tau$  and its linearization  $\text{CHAIN}(\tau) = (\text{snsb}, 5)$

both  $Q_{18}$  and  $Q_{19}$  output objects conforming to the sort  $\tau$  shown in Figure 2.16. We have labelled the atomic sorts in  $\tau$  with the corresponding output variables/constants from  $Q_{18}$  (above) and  $Q_{19}$  (below). We have labelled the collection sorts in  $\tau$  with the corresponding node numbers from  $Q_{18}$  (left) and  $Q_{19}$  (right). Now consider evaluating  $Q_{18}$  and  $Q_{19}$  over database  $\mathbb{D}_3$  containing the relations  $R$  and  $S$  below. To aid in visualizing the database instance, we also provide a graphical depiction of  $\mathbb{D}_3$ .



Both  $Q_{18}$  and  $Q_{19}$  return the object  $o_{12}$  depicted in Figure 2.17. The figure also shows the linearization of  $o_{12}$  with respect to  $\tau$ ; observe that it conforms to the chain sort  $\text{CHAIN}(\tau) = (\text{snsb}, 5)$  shown in Figure 2.16.

We will now describe a function `ENCODE` that maps satisfiable unparameterized  $\text{COCQL}_{\{f_b, f_s, f_n\}}$  queries to hierarchical conjunctive encoding queries. Given an arbitrary satisfiable  $\text{COCQL}_{\{f_b, f_s, f_n\}}$  query  $Q$ , the encoding query  $\text{ENCODE}(Q)$  is constructed by the following steps. For now, assume that  $Q$  does not contain the *generalized projection* operator (this assumption is not restrictive since generalized projection can always be rewritten using scalar aggregation as per Figure 2.5).

1. Number each of the nodes in  $Q$  that applies an aggregation function (i.e., the root `EVAL` as well as each  $\lambda$ ), as we have done to query  $Q_{18}$  in Figure 2.15.
2. For each node numbered  $\boxed{i}$ , create an algebraic expression  $E_i$  as follows:

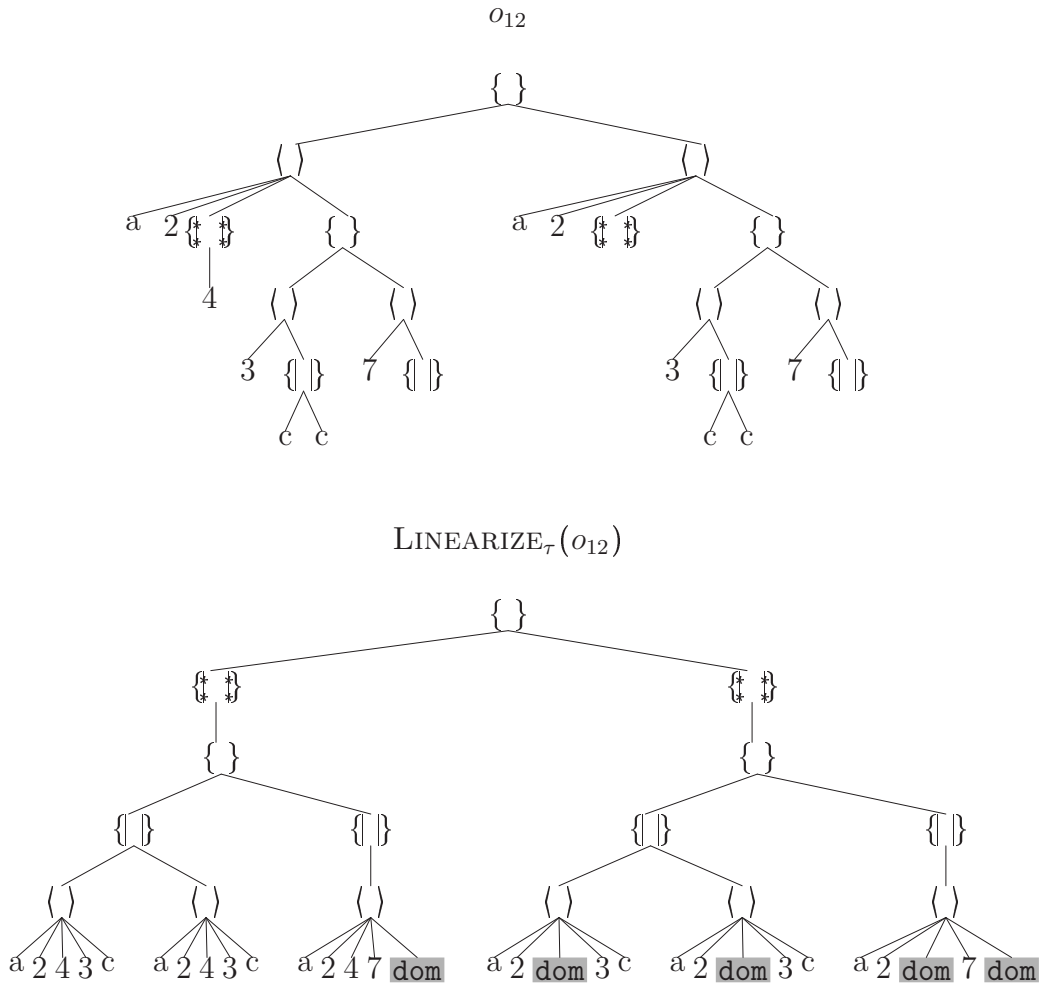


Figure 2.17: Query result  $o_{12} = (Q_{18})^{\text{D}_3} = (Q_{19})^{\text{D}_3}$  and its linearization w.r.t.  $\tau$

- (a) Initialize  $E_i$  with the algebraic expression whose result is the input to the aggregation function at node  $\boxed{i}$  (i.e., the right input if node  $\boxed{i}$  is an  $\lambda$  operator).
- (b) Traverse  $E_i$  depth-first starting from its root.
  - Whenever a scalar subquery operator ( $\lambda$ ) is encountered, delete that operator and its entire right input expression, and continue traversing the left input expression.
  - Whenever a duplicate-preserving projection operator ( $\Pi^{dup}$ ) is encountered
    - delete from its projection list all attributes of non-atomic sort,
    - add to its projection list any constants or attributes of atomic sort that are currently projected away,
    - change it to duplicate-eliminating projection, and
    - continue traversing its input expression.
- (c)  $E_i$  now contains only base relations, selection, Cartesian product, join, and duplicate-eliminating projection. Convert  $E_i$  to a canonical form having a single duplicate-eliminating projection operator above a conjunctive selection operator above the Cartesian product of all of the relations.

For example, query  $Q_{18}$  yields the following expressions.

$$\begin{aligned}
 E_1 &= \Pi_{a,A,B,C,D}(\sigma_{B=C}(R(A, B) \times R(C, D))) \\
 E_2 &= \Pi_{B,F}(\sigma_{E=B}(R(E, F))) \\
 E_3 &= \Pi_{L,M,N,O,P}(\sigma_{K=D \wedge O=B \wedge L=M}(S(K, L) \times S(M, N) \times R(O, P))) \\
 E_4 &= \Pi_{c,G,H,I,J}(\sigma_{G=F \wedge H=I}(S(G, H) \times S(I, J)))
 \end{aligned}$$

3. To construct the body of  $\text{ENCODE}(Q)$ :
  - (a) From each expression  $E_i$  above, create a vertex labelled  $\boxed{i}$  containing a  $\text{CQ}^\#$  body constructed from the set of relations and predicates occurring in  $E_i$ .
  - (b) For each pair of expressions  $E_i$  and  $E_j$ , if the traversal of  $E_i$  visited (and deleted) the  $\lambda$  node labelled  $\boxed{j}$ , then draw an edge from vertex  $\boxed{i}$  to vertex  $\boxed{j}$ .
4. To construct the head of  $\text{ENCODE}(Q)$ :
  - (a) Let  $\tau$  be the output sort of  $Q$ . Label the atomic sorts in  $\tau$  with the corresponding output variables/constants and the collection types with the corresponding node numbers (see Example 21).
  - (b) The depth  $d$  of the encoding query head is the depth of the chain sort  $\text{CHAIN}(\tau)$  (e.g.,  $d = 4$  for  $Q_{18}$ ). We calculate the index variables at each level  $i \in [1, d]$  as follows.
    - i. Enumerate the collection types in  $\text{CHAIN}(\text{sort})$  starting at the root.
    - ii. If the  $i^{\text{th}}$  collection type enumerated has label  $\boxed{j}$ , form  $\bar{\mathcal{I}}_i$  from the output list of  $E_j$  by labelling each constant/variable with  $\boxed{j}$ .

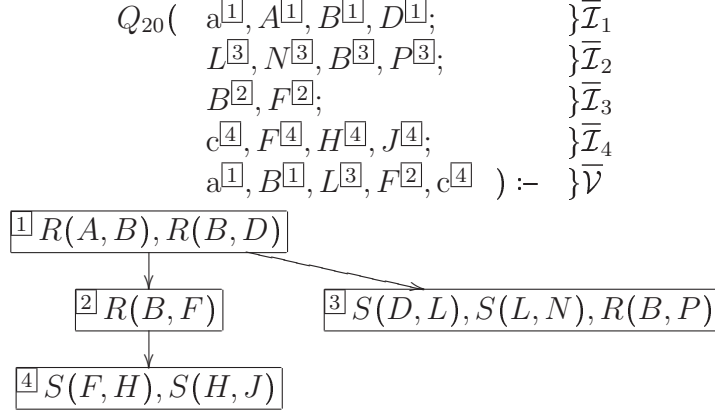


Figure 2.18: Encoding query  $Q_{20} = \text{ENCODE}(Q_{18})$

For example, the collection types in  $\text{CHAIN}(\tau)$  from query  $Q_{18}$  are enumerated in the order  $\boxed{1}$ ,  $\boxed{3}$ ,  $\boxed{2}$ ,  $\boxed{4}$ , yielding the following indexes.

$$\begin{array}{ll}
 \bar{\mathcal{I}}_1 = a^{\boxed{1}}, A^{\boxed{1}}, B^{\boxed{1}}, C^{\boxed{1}}, D^{\boxed{1}} & \bar{\mathcal{I}}_2 = L^{\boxed{3}}, M^{\boxed{3}}, N^{\boxed{3}}, O^{\boxed{3}}, P^{\boxed{3}} \\
 \bar{\mathcal{I}}_3 = B^{\boxed{2}}, F^{\boxed{2}} & \bar{\mathcal{I}}_4 = c^{\boxed{4}}, G^{\boxed{4}}, H^{\boxed{4}}, I^{\boxed{4}}, J^{\boxed{4}}
 \end{array}$$

- (c) The output list is constructed from  $\tau$  by a pre-order traversal that emits at each atomic domain the associated variable/constant, labelled with the label of the closest containing collection type. For example, for query  $Q_{18}$  we obtain  $\bar{\mathcal{V}} = a^{\boxed{1}}, C^{\boxed{1}}, L^{\boxed{3}}, F^{\boxed{2}}, c^{\boxed{4}}$ .
5. Finally, we merge variables to simplify the query representation and to conform to our syntactic convention for HCQs.
    - (a) For every vertex  $v$  in the HCQ body and every variable  $X_i \in \mathcal{L}_v$ , if  $\text{pred}_v \models X_i = Y$  then replace all occurrences of  $X_i$  in  $Q$  with  $Y$ . This causes the following substitutions for our running example:  $C/B$ ,  $E/B$ ,  $G/F$ ,  $I/H$ ,  $K/D$ ,  $O/B$ , and  $M/L$ .
    - (b) Delete any predicates that are trivial.
    - (c) For every component  $X_i^{\boxed{v}} \in \hat{\mathcal{I}}_{[1,d]}$ , delete all but the *left-most* occurrence of  $X_i^{\boxed{v}}$  from the index sequence  $\bar{\mathcal{I}}_{[1,d]}$ .

Figure 2.18 shows the encoding query  $Q_{20}$  which is the final result of our running example  $\text{ENCODE}(Q_{18})$ .

We now extend the definition of  $\text{ENCODE}(Q)$  to allow  $Q$  to contain generalized projection operators. Handling generalized projection operators directly is desirable because it simplifies both queries  $Q$  and  $\text{ENCODE}(Q)$ . Generalized projection occurs much more frequently than scalar aggregation in real workloads, and so restricting our attention to  $\text{COCQL}^{\mathcal{F}}$ -expressions will allow for a simpler decision procedure to test equivalence for this common class of queries. (Specifically, reasoning tools presented in Chapter 3 for CEQs will suffice, rather than the more

general methods for HCEQs required by arbitrary  $\text{COCQL}^{\mathcal{F}}$  queries.) We make the following modifications to  $\text{ENCODE}(Q)$ :

1. When numbering nodes that apply aggregation functions, generalized projection nodes also need to be numbered. For example, see query  $Q_{19}$  in Figure 2.15.
2. During the traversal of each expression  $E_i$ , whenever a generalized projection operator is encountered, delete from the attribute list the aggregation expression  $X = f(A_{j_1}, \dots, A_{j_l})$  (effectively converting the operator into duplicate-eliminating projection) and continue traversing its input expression. For example, query  $Q_{19}$  yields the following expressions.

$$\begin{aligned}
E_5 &= \Pi_{a,C',D'}(\sigma_{B'=C' \wedge B'=E'}(R(A', B') \times R(C', D') \times R(E', F'))) \\
E_6 &= \Pi_{U',b}(\sigma_{V'=c \wedge U'=D'}(R(U', V'))) \\
E_7 &= \Pi_{K',L',M',N'}(\sigma_{K'=D' \wedge L'=M'}(S(K', L') \times S(M', N'))) \\
E_8 &= \Pi_{A',B',C',D',E',F'}(\sigma_{B'=C' \wedge B'=E'}(R(A', B') \times R(C', D') \times R(E', F'))) \\
E_9 &= \Pi_{c,H',J'}(\sigma_{G'=F' \wedge H'=I'}(S(G', H') \times S(I', J')))
\end{aligned}$$

3. When creating the body of  $\text{ENCODE}(Q)$ , we ignore the expressions corresponding to generalized projection nodes (i.e, don't create  $\text{CQ}^{\bar{=}}$  bodies for them and don't consider them when constructing edges); the rest of the algorithm remains unchanged. For example, the traversal yielding expression  $E_8$  above was a sub-traversal of the one yielding expression  $E_5$ . Because of this, all of the base relations (with their local variables) and predicates occurring in  $E_8$  already occur in  $E_5$ .
4. The head of  $\text{ENCODE}(Q)$  is calculated from  $\tau$  and each expression  $E_j$  almost as before, except that when constructing index tuple  $\bar{\mathcal{I}}_i$  for collection type  $\boxed{j}$ , we label the components of  $\bar{\mathcal{I}}_i$  with the label of the *nearest enclosing collection type within  $\tau$  that is constructed by an  $\lambda$  operator*. For example, the collection types in query  $Q_{19}$  are enumerated in the order  $\boxed{5}$ ,  $\boxed{7}$ ,  $\boxed{8}$ ,  $\boxed{9}$ , but type  $\boxed{8}$  is constructed by a generalized projection operator whose nearest enclosing type within  $\tau$  that is constructed by an  $\lambda$  operator is type  $\boxed{5}$ .

$$\begin{aligned}
\bar{\mathcal{I}}_1 &= a^{\boxed{5}}, C'^{\boxed{5}}, D'^{\boxed{5}} & \bar{\mathcal{I}}_2 &= K'^{\boxed{7}}, L'^{\boxed{7}}, M'^{\boxed{7}}, N'^{\boxed{7}} \\
\bar{\mathcal{I}}_3 &= A'^{\boxed{5}}, B'^{\boxed{5}}, C'^{\boxed{5}}, D'^{\boxed{5}}, E'^{\boxed{5}}, F'^{\boxed{5}} & \bar{\mathcal{I}}_4 &= c^{\boxed{9}}, H'^{\boxed{9}}, J'^{\boxed{9}}
\end{aligned}$$

Output attributes are labelled similarly, yielding  $\bar{\mathcal{V}} = a^{\boxed{5}}, C'^{\boxed{5}}, L'^{\boxed{7}}F'^{\boxed{5}}, c^{\boxed{9}}$ .

5. A final syntax simplification is performed as before. This causes the following variable substitutions for our running example:  $C'/B'$ ,  $E'/B'$ ,  $G'/F'$ ,  $I'/H'$ ,  $K'/D'$ ,  $M'/L'$ ,  $U'/D'$ , and  $V'/c$ .

Figure 2.19 shows the encoding query  $Q_{21}$  which is the final result of our running example  $\text{ENCODE}(Q_{19})$ . Although expression  $E_8$  was not used when creating the body of  $Q_{21}$ , it was needed to determine the index tuple  $\bar{\mathcal{I}}_3$  in the query head.

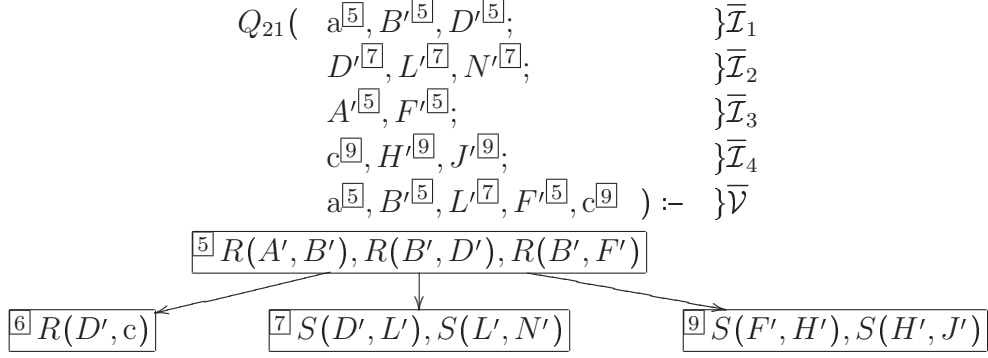


Figure 2.19: Encoding query  $Q_{21} = \text{ENCODE}(Q_{19})$

**Theorem 2.3.6** *Given any satisfiable unparameterized  $\text{CQCQL}^{\{f_b, f_s, f_n\}}$  query  $Q$ , the query  $\text{ENCODE}(Q)$  is a valid encoding query.*

PROOF. When  $\text{body}_{\text{ENCODE}(Q)}$  is constructed by the algorithm  $\text{ENCODE}(Q)$ , the  $\lambda$  operators effectively partition the traversals of the original algebra tree (ignoring the traversals corresponding to generalized projection operators, since they don't introduce new nodes into  $\text{body}_{\text{ENCODE}(Q)}$ ). This means that each labelled node in  $\text{body}_Q$  is only ever visited once, and hence every non-root node of  $\text{body}_{\text{ENCODE}(Q)}$  has exactly one incoming edge. This partitioning also guarantees that the sets of local variables in each CQ body are pairwise disjoint. Syntactic validity of query  $Q$  requires that every attribute referenced in a selection predicate must either be visible in the schema of the current algebra expression or must be bound within the evaluation environment inherited by the nearest  $\lambda$  operator. The fact that  $Q$  is unparameterized then guarantees that any non-local variable used in a CQ body occurs as a local variable in one of its ancestor CQ bodies. Hence,  $\text{body}_{\text{ENCODE}(Q)}$  is a proper HCQ body (as per Definition 2.3.2).

When  $\text{head}_{\text{ENCODE}(Q)}$  is constructed by algorithm  $\text{ENCODE}(Q)$ , the variables in the head originate from attribute references in  $Q$ . Because  $Q$  is unparameterized, every variable in  $\text{head}_{\text{ENCODE}(Q)}$  will be local to  $\text{body}_{\text{ENCODE}(Q)}$ . In order for  $\text{head}_{\text{ENCODE}(Q)}$  to contain some variable  $W_i^{[u]}$ , it follows from the algorithm that expression  $E_v$  must contain some variable  $W_j$  such that variables  $W_i$  and  $W_j$  were merged during the final step. Syntactic validity of query  $Q$  requires that either  $W_j$  is local to expression  $E_v$ , or  $W_j$  was inherited from the parent evaluation environment in  $\text{body}_Q$ . The construction of  $\text{body}_{\text{ENCODE}(Q)}$  entails that either  $W_i \in \mathcal{L}_v$  or  $W_i \in \mathcal{L}_u$  for some vertex  $u$  that is an ancestor of  $v$  in  $\text{body}_{\text{ENCODE}(Q)}$ . Hence,  $\text{head}_{\text{ENCODE}(Q)}$  satisfies Definitions 2.3.2 and 2.3.3 and so  $\text{ENCODE}(Q)$  is a proper HCEQ.

Finally, we need to establish that  $\text{ENCODE}(Q)$  is *valid* in that it satisfies Definition 2.3.4. Algorithm  $\text{ENCODE}(Q)$  causes each component of the output tuple to also appear in one of the index tuples (because it necessarily occurs within the projection list of one of the expressions  $E_i$ ). Therefore, by Lemma 2.3.5,  $\text{ENCODE}(Q)$  is valid.  $\square$

**Theorem 2.3.7** *Given any satisfiable unparameterized  $\text{COCQL}^{\{f_b, f_s, f_n\}}$ -query  $Q$ , the query  $\text{ENCODE}(Q)$  is a valid conjunctive encoding query.*

PROOF. Similar to the proof of Theorem 2.3.6, but relies on the fact that the extended definition of  $\text{ENCODE}(Q)$  avoids introducing hierarchical edges into the query body for generalized projection operators.  $\square$

**Example 22** Consider the encoding queries  $Q_{20} = \text{ENCODE}(Q_{18})$  from Figure 2.18 and  $Q_{21} = \text{ENCODE}(Q_{19})$  from Figure 2.19. Figures 2.20 and 2.21 shows the encoding relations  $R_8$  and  $R_9$  that result from evaluating  $Q_{20}$  and  $Q_{21}$ , respectively, over database  $\mathbb{D}_3$  from Example 21. The reader can verify that decoding either relation with the signature  $\text{snsb}$  (cf. Algorithm 3) yields the object  $o_{13}$  shown in Figure 2.22.

Observe that object  $o_{13}$  in Figure 2.22 differs from object  $\text{LINEARIZE}_\tau(o_{12})$  in Figure 2.17 only by the additional shading of one of the collection nodes. This “spurious shading” is a side-effect of the decoding process, but it is entirely predictable. The spuriously shaded set within chain object  $o_{13}$  corresponds to an (unshaded) set within object  $o_{12}$  that was linearized underneath an instance of an *empty* normalized bag. By adding additional bookkeeping to Algorithm 1, it is possible to create a modified function  $\widehat{\text{CHAIN}}$  that adds spurious shadings whenever a collection object is linearized beneath an instance of an empty collection. Analogous to Definition 2.1.8, we can then define a modified linearization function  $\widehat{\text{LINEARIZE}}_\tau$  as follows

$$\widehat{\text{LINEARIZE}}_\tau(o) := \widehat{\text{CHAIN}} \circ \text{COMPLETE}_\tau(o) \quad (2.25)$$

which then satisfies  $o_{13} = \widehat{\text{LINEARIZE}}_\tau(o_{12})$ . It follows from Observation 2.1.6 that for each tuple sub-object within  $\text{COMPLETE}_\tau(o)$ , the roots of the tuple components are either all shaded or all unshaded. Given sort  $\tau$  and object  $\widehat{\text{LINEARIZE}}_\tau(o)$ , we can therefore identify precisely the shadings that are spurious. The following theorem and corollary are suitably modified versions of Theorem 2.1.9 and Corollary 2.1.10.

**Theorem 2.3.8** *Given sort  $\tau$  not containing any zero-ary tuple sorts,  $\widehat{\text{LINEARIZE}}_\tau$  has a retraction.*

**Corollary 2.3.9** *Given sort  $\tau$  not containing any zero-ary tuple sorts, any two objects  $o, o' \in [[\tau]]$  satisfy the following.*

$$o = o' \iff \widehat{\text{LINEARIZE}}_\tau(o) = \widehat{\text{LINEARIZE}}_\tau(o')$$

**Theorem 2.3.10** *Given any finite database  $\mathbb{D}$  and any  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  query  $Q$  with output sort  $\tau$  not containing any zero-ary tuples, let  $\bar{\xi}$  be the signature corresponding to the linearization of  $\tau$ . Then, the complex object  $(Q)^\mathbb{D}$  and the encoding relation  $(\text{ENCODE}(Q))^\mathbb{D}$  are related as follows.*

$$\text{DECODE}((\text{ENCODE}(Q))^\mathbb{D}, \bar{\xi}) = \widehat{\text{LINEARIZE}}_\tau((Q)^\mathbb{D}) \quad (2.26)$$



$R_8$	$a^{[1]}$	$A^{[1]}$	$B^{[1]}$	$D^{[1]}$	$L^{[3]}$	$N^{[3]}$	$B^{[3]}$	$P^{[3]}$	$B^{[2]}$	$F^{[2]}$	$c^{[4]}$	$F^{[4]}$	$H^{[4]}$	$J^{[4]}$	$a^{[1]}$	$B^{[1]}$	$L^{[3]}$	$F^{[2]}$	$c^{[4]}$
	a	a	2	3	4	5	2	3	2	3	c	3	4	5	a	2	4	3	c
	a	a	2	3	4	5	2	3	2	3	c	3	4	6	a	2	4	3	c
	a	a	2	3	4	5	2	3	2	7	-	-	-	-	a	2	4	7	-
	a	a	2	3	4	5	2	7	2	3	c	3	4	5	a	2	4	3	c
	a	a	2	3	4	5	2	7	2	3	c	3	4	6	a	2	4	3	c
	a	a	2	3	4	5	2	7	2	7	-	-	-	-	a	2	4	7	-
	a	a	2	3	4	6	2	3	2	3	c	3	4	5	a	2	4	3	c
	a	a	2	3	4	6	2	3	2	3	c	3	4	6	a	2	4	3	c
	a	a	2	3	4	6	2	3	2	7	-	-	-	-	a	2	4	7	-
	a	a	2	7	-	-	-	-	2	3	c	3	4	5	a	2	-	3	c
	a	a	2	7	-	-	-	-	2	3	c	3	4	6	a	2	-	3	c
	a	a	2	7	-	-	-	-	2	7	-	-	-	-	a	2	-	7	-
	a	1	2	3	4	5	2	3	2	3	c	3	4	5	a	2	4	3	c
	a	1	2	3	4	5	2	3	2	3	c	3	4	6	a	2	4	3	c
	a	1	2	3	4	5	2	3	2	7	-	-	-	-	a	2	4	7	-
	a	1	2	3	4	5	2	7	2	3	c	3	4	5	a	2	4	3	c
	a	1	2	3	4	5	2	7	2	3	c	3	4	6	a	2	4	3	c
	a	1	2	3	4	5	2	7	2	7	-	-	-	-	a	2	4	7	-
	a	1	2	3	4	6	2	3	2	3	c	3	4	5	a	2	4	3	c
	a	1	2	3	4	6	2	3	2	3	c	3	4	6	a	2	4	3	c
	a	1	2	3	4	6	2	3	2	7	-	-	-	-	a	2	4	7	-
	a	1	2	3	4	6	2	7	2	3	c	3	4	5	a	2	4	3	c
	a	1	2	3	4	6	2	7	2	3	c	3	4	6	a	2	4	3	c
	a	1	2	3	4	6	2	7	2	7	-	-	-	-	a	2	4	7	-
	a	1	2	7	-	-	-	-	2	3	c	3	4	5	a	2	-	3	c
	a	1	2	7	-	-	-	-	2	3	c	3	4	6	a	2	-	3	c
	a	1	2	7	-	-	-	-	2	7	-	-	-	-	a	2	-	7	-

Figure 2.20: Encoding query result  $R_8 = (Q_{20})^{D_3}$  satisfying  $\text{DECODE}(R_8, \text{snsb}) = o_{13}$

$R_9$															
$a^{[5]}$	$B^{[5]}$	$D^{[5]}$	$D^{[7]}$	$L^{[7]}$	$N^{[7]}$	$A^{[5]}$	$F^{[5]}$	$c^{[9]}$	$H^{[9]}$	$J^{[9]}$	$a^{[5]}$	$B^{[5]}$	$L^{[7]}$	$F^{[5]}$	$c^{[9]}$
a	2	3	3	4	5	a	3	c	4	5	a	2	4	3	c
a	2	3	3	4	5	a	3	c	4	6	a	2	4	3	c
a	2	3	3	4	5	a	7	-	-	-	a	2	4	7	-
a	2	3	3	4	5	1	3	c	4	5	a	2	4	3	c
a	2	3	3	4	5	1	3	c	4	6	a	2	4	3	c
a	2	3	3	4	5	1	7	-	-	-	a	2	4	7	-
a	2	3	3	4	6	a	3	c	4	5	a	2	4	3	c
a	2	3	3	4	6	a	3	c	4	6	a	2	4	3	c
a	2	3	3	4	6	1	3	c	4	6	a	2	4	3	c
a	2	3	3	4	6	1	7	-	-	-	a	2	4	7	-
a	2	7	7	-	-	a	3	c	4	5	a	2	-	3	c
a	2	7	7	-	-	a	3	c	4	6	a	2	-	3	c
a	2	7	7	-	-	a	7	-	-	-	a	2	-	7	-
a	2	7	7	-	-	1	3	c	4	5	a	2	-	3	c
a	2	7	7	-	-	1	3	c	4	6	a	2	-	3	c
a	2	7	7	-	-	1	7	-	-	-	a	2	-	7	-

Figure 2.21: Encoding query result  $R_9 = (Q_{21})^{D_3}$  satisfying  $\text{DECODE}(R_9, \text{snsh}) = O_{13}$

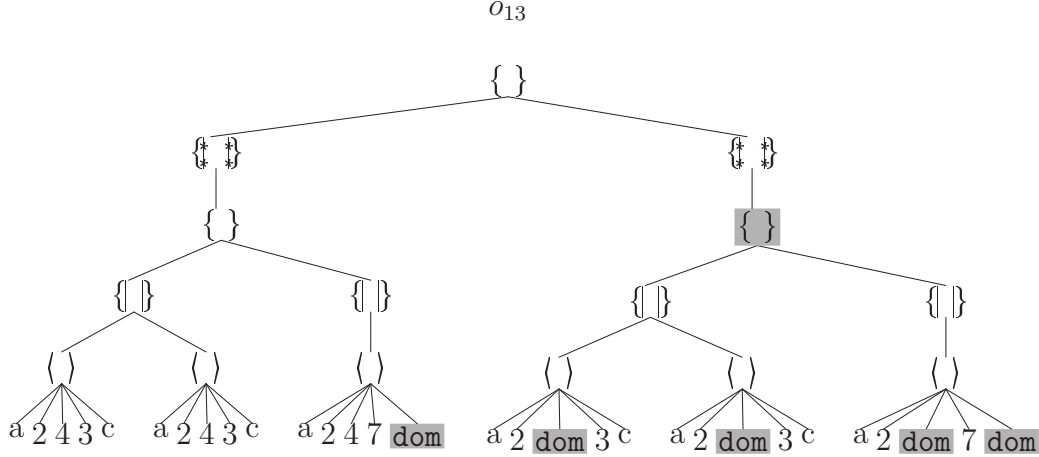


Figure 2.22: Query result  $o_{13} = \text{DECODE}((Q_{20})^{\mathbb{D}^3}, \text{snsb}) = \text{DECODE}((Q_{21})^{\mathbb{D}^3}, \text{snsb})$

**PROOF (SKETCH).** For conciseness of presentation, let  $o$  denote the object  $(Q)^{\mathbb{D}}$ ,  $Q'$  denote the encoding query  $\text{ENCODE}(Q)$ , and  $R$  denote the encoding relation  $(\text{ENCODE}(Q))^{\mathbb{D}}$ .

If  $o$  is a trivial object, then  $\text{LINEARIZE}_{\tau}(o)$  is one of  $\{\{\}, \{\}, \{\!\!\{ \!\!\}\}$  (depending upon the first character of  $\bar{\xi}$ ). In this case it is easy to show that  $R$  will be an empty encoding relation, and so the result follows from the definition of  $\text{DECODE}$  (Algorithm 3). Therefore, assume that  $o$  is not trivial.

Consider the simple case where  $o$  is a chain object that does not contain empty collections, and so object linearization (cf. Section 2.1.2) is idempotent.

$$\widehat{\text{LINEARIZE}}_{\tau}(o) = \widehat{\text{CHAIN}} \circ \text{COMPLETE}_{\tau}(o) = o \quad (2.27)$$

The procedure  $\text{ENCODE}$  essentially partitions the algebraic tree of  $Q$  based upon the  $\lambda$  operators and reshapes the hierarchy of these partitions into a hierarchy of CQ bodies. It is straightforward to verify that because  $o$  does not contain any empty collections,  $R$  does not contain the  $\blacksquare$  symbol. In this simple case, procedure  $\text{DECODE}$  (Algorithm 3) is equivalent to procedure  $\text{DECODESIMPLE}$  (Algorithm 2). Because  $\text{DECODESIMPLE}$  essentially performs the same collection-construction operations as would be required to evaluate  $Q$  over  $\mathbb{D}$ , it is tedious but straightforward to verify the following.

$$\text{DECODE}(R, \bar{\xi}) = \text{DECODESIMPLE}(R, \bar{\xi}) = o \quad (2.28)$$

The result then follows from equations 2.26, 2.27, and 2.28.

Now suppose that  $o$  is a chain object but contains an empty collection generated by some  $\lambda$  operator.

$$\widehat{\text{LINEARIZE}}_{\tau}(o) = \widehat{\text{CHAIN}} \circ \text{COMPLETE}_{\tau}(o) = \text{COMPLETE}_{\tau}(o) \quad (2.29)$$

In this case, function  $\text{COMPLETE}_\tau$  introduces shaded nodes to complete the object. Appealing to the semantics of HCQs, we can show that there exist a partial application of some terminal embedding of  $Q'$  into  $\mathbb{D}$  that assigns the value  $\blacksquare$  to all of the variables in the index level corresponding to the shaded subobject. These embeddings yield index tuples satisfying the predicate  $\bar{a} = \blacksquare^{\bar{a}}$  in Algorithm 3, causing the constructed nodes to be shaded. The result then follows from equations 2.26 and 2.29.

Finally, suppose that  $o$  is not a chain object;

$$\widehat{\text{LINEARIZE}}_\tau(o) = \widehat{\text{CHAIN}} \circ \text{COMPLETE}_\tau(o) \quad (2.30)$$

in this case, the effect of the  $\text{COMPLETE}_\tau$  function is mimicked by  $\blacksquare$  values in  $R$  and in the decoding process as described above, but the effect of the  $\widehat{\text{CHAIN}}$  function is non-trivial. The net effect of the  $\widehat{\text{CHAIN}}$  function is three-fold: first, to linearize sibling collection types consistent with a pre-order traversal of  $\tau$  (see Example 7); second, to push atomic values down by distributing them across all of the leaves; and third, to add spurious shadings to certain collection nodes. The  $\text{ENCODE}$  and  $\text{DECODE}$  functions together mirror all three of these behaviours. The type linearization occurs both when the indexes within encoding head of  $Q'$  are chosen based upon a pre-order traversal of  $\tau$  and when the signature  $\bar{\xi}$  used in the decoding process is calculated via a pre-order traversal of  $\tau$ . The push-down of atomic values is implicit in the semantics of HCQs—specifically, that the value for a particular attribute originating from a parent block is repeated within all of the tuples originating from different terminal embeddings that agree on the variable bindings of the parent block but differ on the bindings of the child block. Lastly, consider any two sibling collection types within sort  $\tau$  that correspond to collection types in  $\text{CHAIN}(\tau)$  at depths  $i, j$  with  $i < j$ . Given any tuple that satisfies  $\bar{\mathcal{I}}_{i-1} = \blacksquare^{\bar{\mathcal{I}}_{i-1}}$  but  $\bar{\mathcal{I}}_j \neq \blacksquare^{\bar{\mathcal{I}}_j}$ , the predicate test  $\bar{\mathcal{I}}_{i-1} = \blacksquare^{\bar{\mathcal{I}}_{i-1}}$  in Algorithm 3 will cause a spurious shading of the child collection object constructed at depth  $i$ .  $\square$

**Definition 2.3.11 (Encoding-Equivalence)** *Given two encoding queries  $Q, Q'$  over database schema  $\mathbb{S}(\mathbb{D})$  and a signature  $\bar{\xi}$ , if for each relational database  $\mathbb{D}$  with schema  $\mathbb{S}(\mathbb{D})$  and finite instance  $\mathbb{T}(\mathbb{D})$  the encoding-equality  $(Q)^{\mathbb{D}} \stackrel{\bar{\xi}}{\doteq} (Q')^{\mathbb{D}}$  holds, then we say that  $Q$  and  $Q'$  are encoding-equivalent w.r.t. signature  $\bar{\xi}$  (abbreviated  $\bar{\xi}$ -equivalent, and denoted  $Q \stackrel{\bar{\xi}}{\equiv} Q'$ ).*

**Theorem 2.3.12** *Given two satisfiable, unparameterized  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  queries  $Q, Q'$  with the same output sort  $\tau$  not containing zero-ary tuples,*

$$Q \equiv Q' \iff \text{ENCODE}(Q) \stackrel{\bar{\xi}}{\doteq} \text{ENCODE}(Q')$$

where  $\bar{\xi}$  is the signature corresponding to the linearization of  $\tau$ .

**PROOF.** Follows from the relationship between the output of queries  $Q$  and  $\widehat{\text{ENCODE}}(Q)$  (Theorem 2.3.10) and the existence of a retraction for transformation  $\widehat{\text{LINEARIZE}}_\tau$  (Theorem 2.3.8 and Corollary 2.3.9).  $\square$

Note that the ban on zero-ary tuples in Theorem 2.3.12 is easily overcome. If the output sort  $\tau$  of  $Q$  and  $Q'$  contains a collection of zero-ary tuples, we can convert this to a collection of unary tuples by modifying both  $Q$  and  $Q'$  so that the corresponding query operators project a constant into the collection constructor. This modification does not change the equivalence relation between  $Q$  and  $Q'$ .

**Corollary 2.3.13** *Deciding equivalence of satisfiable  $CSQL^{\{f_b, f_s, f_n\}}$  queries is polynomial-time reducible to  $\bar{\xi}$ -equivalence of hierarchical conjunctive encoding queries.*

PROOF. Follows from Proposition 2.1.22. □

**Corollary 2.3.14** *Deciding equivalence of satisfiable  $CSQL^{\{f_b, f_s, f_n\}-}$  queries is polynomial-time reducible to  $\bar{\xi}$ -equivalence of conjunctive encoding queries.*

PROOF. Follows from Proposition 2.1.23 and Theorem 2.3.7. □

**Example 23** The equivalence of  $COCQL^{\{f_b, f_s, f_n\}}$  queries  $Q_{18}$  and  $Q_{19}$  can be decided by testing the condition  $Q_{20} \stackrel{\text{sbsb}}{=} Q_{21}$  (which we show to be true in Example 59, Chapter 5).

Now suppose that  $Q_{18}$  and  $Q_{19}$  were each modified by replacing the aggregate function  $f_n$  with  $f_b$ . The modified queries would not be equivalent, because  $Q_{20}$  and  $Q_{21}$  are not sbsb-equivalent. (The reader can easily verify from Figures 2.20 and 2.21 that  $R_8 \not\stackrel{\text{sbsb}}{=} R_9$  because object  $\text{DECODE}(R_8, \text{sbsb})$  contains a sub-bag with cardinality four, while object  $\text{DECODE}(R_9, \text{sbsb})$  does not.)

## 2.4 Relevant Literature

Here we briefly survey literature relevant to deciding equivalence between queries that output complex objects.

### 2.4.1 Query Languages for Complex Objects

Much of the previous work on complex objects has been performed within the context of *nested relations*. Because a relation is by definition a *set* of tuples, a nested relation corresponds to a complex object whose sort contains only *set* collections. The literature surrounding nested relations and (set-based) complex objects is vast; Abiteboul, Hull, and Vianu provide a good summary in their textbook [6, Ch. 20].

Some of the earliest work on nested structures includes the *quotient relations* of Furtado and Kerschberg, in which tuples are organized into blocks via the *partition* and *departition* operators [35]. The partitioning of quotient relations is not recursive, however, and so the nested relational model along with the restructuring operators *nest* and *unnest* of Jaeschke and Schek [65] can be viewed as a recursive generalization of quotient relations. Jaeschke and Schek's nest/unnest operators

are limited to nesting/unnesting of a single attribute, however; the generalization of nest/unnest to handle nested relations of arbitrary arity is due to subsequent work by Thomas and Fischer [103].

Gyssens, Paradaens, and Van Gucht analyze the space of possible nested relations with the goal of identifying useful subclasses, since they contend that not all nested relations correspond to real world situations [55]. They organize subclasses of nested relations (NR) into the following hierarchy.

$$\text{HNR} < \text{NFR} < \text{NL} < \text{NR}$$

The class of *normalization-lossless relations* (NL) consists of all relations obtainable from flat relations using arbitrary sequences of nest and unnest operators. Note that relations containing empty nested sub-relations are members of NR but not of NL. The class of *nested flat relations* (NFR) consists of relations obtainable from flat relations using only consecutive nest operations. Finally, the class of *hierarchical nested relations* (HNR) is the subset of NFR in which at each level of nesting the atomic attributes form a superkey. Relations within HNR are also known as *Verso-relations* [5, 13], and have the useful property that the size of an instance is polynomially-bounded by the number of atomic elements in its active domain [6, Ch. 20].

Whereas Gyssens et al.’s taxonomy only considers restructuring performed by sequences of nest and unnest operators, the query language usually called the *nested relational algebra* (NRA) is formed by adding Thomas and Fischer’s nest and unnest operators to the relational algebra. The expressive power of NRA is very limited compared to many complex object query languages; not only is it unable to create empty sub-relations, Paradaens and Van Gucht showed that its expressive power does not exceed the standard relational algebra for queries that both input and output flat relations [89]. A flurry of research in late 1980s and early 1990s yielded a variety of query languages for complex objects with expressive power exceeding NRA (in some cases, *far* exceeding NRA by including various combinations of operators such as empty collection constructors, transitive closure, fixed points, powersets, etc.). For example, pairs of algebraic and calculus-based languages for complex objects along with a proof of equivalent expressive power are proposed by Ozsoyoglu et al. [88], Buneman et al. [14], and Abiteboul and Beeri [4]. Grumbach considers a calculus for *bags* [49], while Libkin and Wong consider a variety of calculus-based languages that include nested sets, bags, and various aggregation functions [78, 79, 112]. A common theme throughout this body of literature is the focus on the *expressive power* of the query languages. Since all of these languages meet or exceed the expressive power of the traditional relational algebra/calculus, the query equivalence problem is undecidable for any of these languages.

To situate our work within the context of this body of literature, our query language  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  can be described informally as follows. The restricted subclass  $\text{COCQL}^{\{f_b, f_s, f_n\}-}$  that we study in Chapter 3—queries that do not contain scalar sub-queries but may contain generalized projection—is essentially an extension of the

*bag semantic conjunctive relational algebra* with three variants of the *nest* operator (for constructing different collection types), but with no *unnest* operator and with certain restrictions on the operations permitted on non-atomic attributes. If the three variants of *nest* are replaced with analogous variants of the *outer-nest* operator described by Levy and Suciu [77], then the full language  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  is obtained.

Even though language  $\text{COCQL}^{\mathcal{F}}$  does not contain an *unnest* operator, it is straightforward to show that any object in Gyssens et al.’s class NR can be constructed by language  $\text{COCQL}^{\{f_s\}}$  out of some database of flat relations. Similarly, language  $\text{COCQL}^{\{f_s\}-}$  is capable of constructing any object in class NL. This apparent discrepancy is due to the simple fact that  $\text{COCQL}^{\{f_s\}-}$  allows duplicate-eliminating projection to appear above generalized projection, whereas Gyssens et al. restrict their attention to operator sequences containing only *nest*/*unnest*. While on one hand this seems to call into question the usefulness of Gyssens et al.’s taxonomy when considering more powerful query languages, there is an interesting connection between their taxonomy and the equivalence problem that we study in this thesis. Because encoding relations rely on atomic-valued index attributes to distinguish the encoded objects, encoding relations resemble “flattened” Verso-relations. Our main result in Chapter 3 can then be construed as saying that equivalence of  $\text{COCQL}^{\{f_b, f_s, f_n\}-}$  expressions can be reduced to equivalence of queries that return objects in the simplest class HNR. This does not, however, prove the decidability of the open problem of equivalence between two *nest*/*unnest* sequences, because such a sequence may not be equivalent to any  $\text{COCQL}^{\{f_b, f_s, f_n\}-}$  query.

## 2.4.2 Relational Encodings of Complex Objects

Our approach of encoding arbitrary complex objects into flat relations and subsequently translating queries over objects into queries over flat relations is not novel. Suciu proposes a query language that extends NRA with a bounded fixed point operator, and his proof that this language is a conservative extension of non-nested relational algebra plus bounded fixed points—that is, that the languages have the same expressive power over flat inputs and outputs—uses a similar technique [102]. Suciu’s encoding method uses a separate relation for each collection type in the object’s sort (he only considers nested sets), and uses index values to distinguish the elements of each collection. Suciu’s encoding handles arbitrary typed objects, whereas our encoding method only handles chain objects, and thus necessitates prior application of the linearization transformation described in Section 2.1.2. The two encoding methods are fundamentally related: our relational representation obtained by object linearization followed by subsequent encoding of the chain object is equivalent to the *left outer join* of the various relations within Suciu’s encoding. Our encoding requires a special  $\blacksquare$  symbol because it corresponds to the NULL values introduced by an outer join operator.

### 2.4.3 Algebraic Transformations

Algebraic query optimization is a well-studied problem. The literature pertaining to algebraic optimization of SQL queries can be divided into a two different veins.

One line of research focuses on algebraic transformations that change the “shape” of the query, where a query’s “shape” is loosely defined as the nesting structure of the `SELECT-FROM-WHERE-GROUPBY` blocks (which corresponds visually to the shape of the associated Query Graph Model representation [56]). Early work in this vein includes work by Kim on flattening of nested query blocks [69], and subsequent work by Dayal on both merging and decorrelating nested blocks [28]. The advent of rule-based optimization paradigms pioneered by systems like `STARBURST` [56] and `VOLCANO` [47] caused a focus on local manipulations of algebra trees. Various algebraic transformation rules for commuting aggregation blocks, partitioning or merging aggregation operations, avoiding duplicate elimination, or commuting aggregation with join were proposed by Yan and Larson [113, 114], Paulley and Larson [91], Gupta, Harinarayan, and Quass [54], and Khizder and Weddell [68].

The second vein of research generalizes the well-known heuristic of pushing down predicates. Starting with the idea of “magic sets”—an optimization technique proposed for bottom-up evaluation of recursive Datalog programs [9]—Mumick and various co-authors apply the idea to nested SQL queries. As a result, they derive conditions for introducing or eliminating redundant (semi)joins into existing query blocks, or moving selection or join predicates around the algebra tree [76, 85, 86, 87].

Algebraic transformation rules have also been proposed for nested relational algebra queries. Scholl proposes a set of algebraic equivalences that hold over NRA expressions, with the stated goal of eliminating redundant join operations [100]. More recently, Liu and Yu propose an optimization algorithm for NRA expressions based upon heuristic application of algebraic transformations that also attempt to reduce the execution cost of certain types of nested join operations [80].

While algebraic transformations have proved to be a powerful and flexible optimization technique, many of the transformation rules appearing in the literature are given with only sound conditions characterizing their applicability. Complete characterizations for a given transformation rule are difficult to prove because applicability could easily depend upon global features of the query plan, whereas algebraic transformations by their nature perform very localized query rewriting. The query transformation literature as a whole fails to provide a systematic understanding of underlying principles governing global interactions between different nested query blocks. As such, this body of work does not furnish a coherent characterization of query equivalence that would enable the development of a provably complete decision procedure.



## 2.4.4 Query Equivalence

The equivalence of conjunctive queries that input and output flat relational structures is a well-studied problem which we will survey in Chapter 3 when we consider encoding-equivalence of CEQs. We restrict our attention here to the small amount of literature that considers the query equivalence problem for queries that output nested structures.

When a query language allows the ability to both create and explicitly test for empty collections, it can simulate the full relational algebra and query equivalence becomes undecidable. Van den Bussche et al. recently proved a more interesting undecidability result. They show that the query equivalence problem is undecidable for the Positive-Existential fragment of the Nested Relational Calculus [107]. NRC is a conservative extension of relational algebra, and it includes the ability to both create and test for empty sets. PENRC is the fragment of NRC formed by removing relational difference and the emptiness test. Although PENRC lacks the ability to explicitly test set-emptiness, Van den Bussche et al.’s proof of undecidability relies on the ability to both construct objects containing empty subsets and to union sets, but does not appear to require the ability to unnest sets. As such, their undecidability result would likely transfer to any extension of our language  $\text{COCQL}^{\mathcal{F}}$  that includes the ability to union collections, but not necessarily to a similar extension to  $\text{COCQL}^{\mathcal{F}^-}$  since that language does not have the ability to create empty sets.

Containment and equivalence of queries returning set-based complex objects is studied by Levy and Suciu [77], who consider “conjunctive OQL” (COQL) queries (a language with expressive power equivalent to extending our language  $\text{COCQL}^{\{f_s\}}$  with an unnest operator). Whereas containment of flat relations indisputably corresponds to set inclusion, there is no single definition for containment of nested sets. Levy and Suciu use an inductive definition of containment previously proposed for Verso relations [13], and they reduce containment (under this definition) of COQL queries constructing objects with nesting depth  $d$  to testing a relationship between CQs that they call “simulation to depth  $d$ ,” defined as follows. Let  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  be a CQ whose head has been annotated to distinguish  $d$  sets of *index variables*, and define  $\bar{\mathcal{I}} := (\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d)$ . Given two such queries,  $Q$  *simulates*  $Q'$  to depth  $d$ —denoted  $Q \leq_d Q'$ —iff over every database instance the following equation holds:

$$\forall \bar{\mathcal{I}}_1. \exists \bar{\mathcal{I}}'_1 \dots \forall \bar{\mathcal{I}}_d. \exists \bar{\mathcal{I}}'_d. \forall \bar{\mathcal{V}} \left[ Q(\bar{\mathcal{I}}; \bar{\mathcal{V}}) \implies Q'(\bar{\mathcal{I}}'; \bar{\mathcal{V}}) \right] \quad (2.31)$$

a condition characterized by the existence of a *simulation mapping*, and hence NP-complete to decide [77]. For COQL queries that cannot construct empty sets (i.e., our language  $\text{COCQL}^{\{f_s\}^-}$  plus unnest), containment reduces to a single simulation test. Levy and Suciu claim that arbitrary COQL containment reduces to testing an exponential number of simulation conditions; however, Dong et al. point out that this is insufficient for implying containment [32]. Dong et al. consider containment of a restricted class of COQL queries (corresponding to XQuery), showing it to be

$$\begin{aligned}
Q_{22}: \quad & \left\{ \begin{array}{l} \text{SELECT } \{u.C\} \text{ FROM E AS } x, \\ \quad (\text{SELECT } z.P, \{z.C\} \text{ AS } C \text{ FROM E AS } z \\ \quad \text{GROUP BY } z.P) \text{ AS } u \\ \text{WHERE } x.C = u.P \text{ GROUP BY } x.C \end{array} \right\} \\
Q_{23}: \quad & \left\{ \begin{array}{l} \text{SELECT } \{u.C\} \text{ FROM E AS } x, \text{ E AS } y, \\ \quad (\text{SELECT } z.P, \{z.C\} \text{ AS } C \text{ FROM E AS } z \\ \quad \text{GROUP BY } z.P) \text{ AS } u \\ \text{WHERE } x.C = u.P \text{ AND } y.C = u.P \text{ GROUP BY } x.P, y.P \end{array} \right\} \\
Q_{24}: \quad & \left\{ \begin{array}{l} \text{SELECT } \{C\} \text{ FROM E AS } x, \\ \quad (\text{SELECT } z.P, \{z.C\} \text{ AS } C \\ \quad \text{FROM E AS } y, \text{ E AS } z \text{ WHERE } y.C = z.P \\ \quad \text{GROUP BY } y.P, z.P) \text{ AS } u \\ \text{WHERE } x.C = u.P \text{ GROUP BY } x.P \end{array} \right\}
\end{aligned}$$

Figure 2.23: COQL queries  $Q_{22}$ ,  $Q_{23}$ , and  $Q_{24}$

in co-NEXPTIME, but NP-complete or co-NP-complete for a variety of further restrictions. To the best of our knowledge, the complexity of the general COQL containment problem remains open.

The containment relationship used by Levy and Suciu is not antisymmetric (mutual containment does not imply equivalence) and so they define a separate “strong simulation” relationship between CQs for testing COQL equivalence. Query  $Q$  *strongly simulates*  $Q'$  to depth  $d$ —denoted  $Q \ll_d Q'$ —iff:

$$\forall \bar{\mathcal{I}}_1. \exists \bar{\mathcal{I}}'_1 \dots \forall \bar{\mathcal{I}}_d. \exists \bar{\mathcal{I}}'_d. \forall \bar{\mathcal{V}} \left[ Q(\bar{\mathcal{I}}; \bar{\mathcal{V}}) \iff Q'(\bar{\mathcal{I}}'; \bar{\mathcal{V}}) \right] \quad (2.32)$$

a condition which they claim is characterized by the existence of a *strong simulation mapping* [77], and hence still NP-complete to decide (although they define this mapping only for  $d \leq 1$ ). While equivalence of general COQL queries is left open, they claim that equivalence for COQL queries that cannot construct empty sets (i.e., our COCQL<sup>{fs}</sup>- plus unnest) reduces to testing a single strong simulation condition in each direction (Proposition 6.3 [77]). We end this section with an example that demonstrates that this reduction of query equivalence to strong simulation is incorrect.

**Example 24** Consider a database containing a relation  $E(P, C)$  that denotes parent-child relationships, along with the three queries shown in Figure 2.23 (written in an SQL-like syntax that corresponds to empty-set-free COQL). Query  $Q_{22}$  returns sets of related grandchildren, grouped first into sets with a common parent, and then into sets with a common grandparent. Query  $Q_{23}$  is similar to  $Q_{22}$ , but the outer aggregation groups by pairs of grandparents. Query  $Q_{24}$  is also similar to  $Q_{22}$ , but the inner aggregation groups by both parent and grandparent. Levy and Suciu’s technique associates  $Q_{22}$ ,

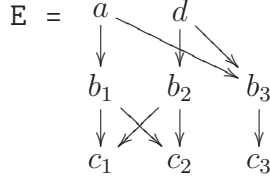


Figure 2.24: Database instance  $\mathbb{D}_4$

$A$	$B$	$C$
$a$	$b_1$	$c_1$
$a$	$b_1$	$c_2$
$a$	$b_3$	$c_3$
$d$	$b_2$	$c_1$
$d$	$b_2$	$c_2$
$d$	$b_3$	$c_3$

$AD$	$B$	$C$
$a a$	$b_1$	$c_1$
$a a$	$b_1$	$c_2$
$a a$	$b_3$	$c_3$
$a d$	$b_3$	$c_3$
$d a$	$b_3$	$c_3$
$d d$	$b_2$	$c_1$
$d d$	$b_2$	$c_2$
$d d$	$b_3$	$c_3$

$A$	$DB$	$C$
$a$	$a b_1$	$c_1$
$a$	$a b_1$	$c_2$
$a$	$a b_3$	$c_3$
$a$	$d b_3$	$c_3$
$d$	$a b_3$	$c_3$
$d$	$d b_2$	$c_1$
$d$	$d b_2$	$c_2$
$d$	$d b_3$	$c_3$

Figure 2.25: Evaluating  $Q'_{22}$ ,  $Q'_{23}$ , and  $Q'_{24}$  over  $\mathbb{D}_4$

$Q_{23}$ , and  $Q_{24}$  with the following indexed CQs.

$$\begin{aligned}
 Q'_{22}(\overline{A} ; \overline{B} ; \overline{C}) &: \neg E(A, B), E(B, C) \\
 Q'_{23}(A, D ; \overline{B} ; \overline{C}) &: \neg E(A, B), E(B, C), E(D, B) \\
 Q'_{24}(\overline{A} ; D, B ; \overline{C}) &: \neg E(A, B), E(B, C), E(D, B)
 \end{aligned}$$

Consider the database  $\mathbb{D}_4$  in Figure 2.24 and the corresponding query results in Figure 2.25 (index groups have been visually separated for clarity). The reader can verify that over database  $\mathbb{D}_4$  all six strong simulation conditions  $Q'_{22} \ll_2 Q'_{23}$ ,  $Q'_{23} \ll_2 Q'_{22}$ ,  $Q'_{22} \ll_2 Q'_{24}$ ,  $Q'_{24} \ll_2 Q'_{22}$ ,  $Q'_{23} \ll_2 Q'_{24}$ , and  $Q'_{24} \ll_2 Q'_{23}$  are satisfied (cf. equation 2.32); in fact, we can show that they are all satisfied over any database. However, the queries are not all equivalent since over  $\mathbb{D}_4$  queries  $Q_{22}$  and  $Q_{24}$  output the object  $\{\{\{c_1, c_2\}, \{c_3\}\}\}$  while  $Q_{23}$  outputs  $\{\{\{c_1, c_2\}, \{c_3\}\}, \{\{c_3\}\}\}$ . Using the results from the next chapter, we can show that queries  $Q_{22}$  and  $Q_{24}$  are equivalent.



## Chapter 3

# Encoding-Equivalence of Conjunctive Encoding Queries

The purpose of this chapter is to present a necessary and sufficient condition for encoding-equivalence in the case where both queries are conjunctive encoding queries (see Section 2.3). Combined with Theorems 2.3.12 and 2.3.7 from Chapter 2, the results in this chapter yield a procedure for deciding query equivalence for the class of  $\text{COCQL}^{\{f_b, f_s, f_n\}\text{-}}$  expressions. This subclass is of practical interest because  $\text{COCQL}^{\{f_b, f_s, f_n\}\text{-}}$  includes queries that can be created by “stacking” views defined by SPJG algebra expressions [29] (see Proposition 2.1.20).

In Theorem 2.3.7 we showed that when the ENCODE algorithm is applied to a  $\text{COCQL}^{\{f_b, f_s, f_n\}\text{-}}$  query, it yields a valid CEQ (cf. Definition 2.3.4). In this chapter, however, we consider *all* valid (and satisfiable) CEQs, independent of their origin. As per Definition 2.3.3, the body of any CEQ is a single CQ body; therefore, without loss of generality we dispense with the vertex labels  $\boxed{v}$  in the query head, since all of the components of the head must have the same label. We restrict our attention to satisfiable CEQs, and so—by our syntactic convention for CQs—we can assume without loss of generality that the query bodies do not contain explicit equality predicates. The following lemma then characterizes CEQ validity.

**Lemma 3.0.1** *Assuming that the set of possible databases is not restricted by schema constraints, conjunctive encoding query  $Q$  is valid if and only if every output variable also occurs as an index—that is,  $(\mathcal{V} \cap \mathcal{B}) \subseteq \mathcal{I}_{[1,d]}$ .*

PROOF. Sufficiency holds by Lemma 2.3.5. Conversely, if  $(\mathcal{V} \cap \mathcal{B}) \not\subseteq \mathcal{I}_{[1,d]}$  then construct database  $\mathbb{D}$  by unioning two copies of  $\text{body}_Q$  in which the variables in  $(\mathcal{V} \cap \mathcal{B}) \setminus \mathcal{I}_{[1,d]}$  have been isomorphically renamed within one of the copies. (Because we have assumed that  $Q$  does not contain explicit equality predicates,  $\text{body}_Q$  is simply a set of atoms and can therefore be treated as a set of canonical tuples.) The query result  $(Q)^\mathbb{D}$  then violates the functional dependency  $\overline{\mathcal{I}}_{[1,d]} \rightarrow \overline{\mathcal{V}}$ , and so is not a valid encoding relation.  $\square$

The remainder of this chapter will proceed as follows. We propose a normal form for CEQs in Section 3.1 and prove that the transformation to the normal form preserves encoding-equivalence. In Section 3.2 we characterize equivalence between CEQs in normal form in terms of the existence of a special type of homomorphism. We survey related literature in Section 3.3.

## 3.1 A Normal Form for CEQs

The main contribution of this section is a normal form for conjunctive encoding queries that preserves encoding-equivalence relative to a given signature  $\bar{\Sigma}$ . Because calculating the normal form requires reasoning about dependencies—specifically, multivalued dependencies—that hold over the result of a conjunctive query, we start by discussing how to reason about multivalued dependencies that are implied by a query definition.

### 3.1.1 Dependencies Over Query Results

We say that a query definition implies a particular dependency if that dependency is guaranteed to hold over the result of that query when evaluated over any possible database instance. Before we consider query-implied dependencies, we quickly review the concept of a *query homomorphism*. The reduction of the query equivalence problem to the existence of query homomorphisms is a fundamental result within database theory, and standard textbook material [6].

**Definition 3.1.1 (CQ Homomorphism)** *Given two conjunctive queries  $Q(\bar{W})$  and  $Q'(\bar{W}')$ , a homomorphism  $h$  from  $Q'$  to  $Q$  is a mapping  $h : \mathcal{B}' \rightarrow \mathcal{B} \cup \mathcal{C}$  (suitably extended to tuples, sets, and atoms, and with identity on constants) satisfying*

1.  $h(\text{body}_{Q'}) \subseteq \text{body}_Q$ ; and
2.  $h(\bar{W}') = \bar{W}$ .

We write  $h : Q' \rightarrow Q$  to denote that  $h$  is a homomorphism from  $Q'$  to  $Q$ . A satisfiable CQ  $Q(\bar{W})$  is *contained* under CQ  $Q'(\bar{W}')$  if and only if there exists a homomorphism  $h : Q' \rightarrow Q$ . Two CQs are *equivalent* if and only if they are mutually contained. A conjunctive query  $Q$  is *minimal* if there is no equivalent query  $Q'$  such that  $|\text{body}_{Q'}| < |\text{body}_Q|$ . Each query  $Q$  has a unique minimal form (up to isomorphic renaming of variables), but computing it is NP-complete [6, Sect. 6.2].<sup>1</sup> The following theorem proves a property of minimal CQs that will be useful for later proofs.

---

<sup>1</sup>Technically, the NP-complete decision problem is testing query equivalence/containment, while minimizing a query requires a linear number of such tests.

**Theorem 3.1.2** *Given any minimal conjunctive query  $Q(\overline{\mathcal{W}})$ , any mapping  $h : \mathcal{B} \rightarrow \mathcal{B} \cup \mathcal{C}$  (extended with identity on constants) that satisfies  $h(\mathcal{W}) = \mathcal{W}$  and  $h(\text{body}_Q) \subseteq \text{body}_Q$  forms a bijection between  $\mathcal{B}$  and itself.*

PROOF. Assume w.l.o.g. that  $\overline{\mathcal{W}}$  does not contain repeated variables or constants.

Suppose that mapping  $h : \mathcal{B} \rightarrow \mathcal{B} \cup \mathcal{C}$  satisfies  $h(\mathcal{W}) = \mathcal{W}$  and  $h(\text{body}_Q) \subseteq \text{body}_Q$  but does not form a bijection between  $\mathcal{B}$  and itself. Then,  $h$  must satisfy  $|h(\text{body}_Q)| < |\text{body}_Q|$ . Furthermore,  $h(\overline{\mathcal{W}})$  must be a permutation of  $\overline{\mathcal{W}}$ . Because every permutation has a finite period, there exists some positive integer  $j$  such that the composite mapping  $h^j := \underbrace{h \circ \dots \circ h}_j$  satisfies  $h^j(\overline{\mathcal{W}}) = \overline{\mathcal{W}}$ ; however,

then  $h^j$  is a homomorphism from  $Q$  to itself that contradicts the minimality of  $Q$ .  $\square$

## Implied Multivalued Dependencies

A *multivalued dependency* (MVD) over some attribute set  $\mathcal{W}$  is an expression of the form  $X \twoheadrightarrow Y$ , where  $X, Y$  are disjoint subsets of  $\mathcal{W}$ . Let  $Z$  denote the attribute set  $\mathcal{W} \setminus (X \cup Y)$ . Then, a relation  $R$  over  $\mathcal{W}$  satisfies  $X \twoheadrightarrow Y$  if the following equation holds [6, Sect. 8.3].

$$R = \Pi_{XY}(R) \bowtie_X \Pi_{XZ}(R) \quad (3.1)$$

Given an CQ  $Q$  that yields a relation over  $\mathcal{W}$ , we say that  $Q$  implies  $X \twoheadrightarrow Y$ —denoted  $Q(\mathcal{W}) \models X \twoheadrightarrow Y$ —if for every database  $\mathbb{D}$  the relation  $(Q)^\mathbb{D}$  satisfies  $X \twoheadrightarrow Y$ .

$$(Q)^\mathbb{D} = \Pi_{XY}((Q)^\mathbb{D}) \bowtie_X \Pi_{XZ}((Q)^\mathbb{D}) \quad (3.2)$$

Projection and natural join are conjunctive operations, and so the entire right-hand side of equation 3.2 can be perceived as a single conjunctive query evaluated over  $\mathbb{D}$ . The following theorem formalizes this.

**Theorem 3.1.3**  $Q(\overline{\mathcal{W}}) \models X \twoheadrightarrow Y \iff Q \equiv \Pi_{XY}(Q) \bowtie_X \Pi_{XZ}(Q)$

**Corollary 3.1.4**  $Q(\overline{\mathcal{W}}) \models X \twoheadrightarrow Y$  iff there exists a homomorphism  $h : Q \rightarrow Q$  (thereby satisfying  $h(\overline{\mathcal{W}}) = \overline{\mathcal{W}}$ ) such that  $h(\text{atoms}_Q)$  can be partitioned into two sets  $S_Y$  and  $S_Z$  satisfying

1. variables  $Y$  only occur in  $S_Y$ ,
2. variables  $Z$  only occur in  $S_Z$ , and
3. all variables in common between  $S_Y$  and  $S_Z$  occur in  $X$ .

**Corollary 3.1.5** If (satisfiable) conjunctive query  $Q(\overline{\mathcal{W}})$  is equivalent to a constant selection operation over conjunctive query  $Q'(\overline{\mathcal{W}})$

$$Q \equiv \sigma_{\overline{A}=\overline{a}}(Q')$$

then  $Q \models X \twoheadrightarrow Y$  iff  $Q' \models X \cup A \twoheadrightarrow Y \setminus A$ .

Prof	school	dept	name	Student	school	dept	name
	UW	CS	Frank		UW	CS	Dave
	UW	CS	David		UW	CS	Xin
	UW	CS	Ashraf		UW	ENG	Cara
	UW	ENG	Linda		UW	ENG	Julia
	UW	ENG	John				

Figure 3.1: Database  $\mathbb{D}_5 = \{\text{Prof}, \text{Student}\}$

**Corollary 3.1.6 (Proposition 10.2.4 in [6])** *If conjunctive query  $Q(\overline{\mathcal{W}})$  is equivalent to a projection operation over conjunctive query  $Q'(\overline{\mathcal{W}}')$*

$$Q(\overline{\mathcal{W}}) \equiv \Pi_{\overline{\mathcal{W}}} (Q'(\overline{\mathcal{W}}'))$$

then  $Q \models X \twoheadrightarrow Y$  iff  $Q' \models X \twoheadrightarrow Y \cup Y'$  for some  $Y' \subseteq \mathcal{B}' \setminus \mathcal{W}$ , where  $\mathcal{B}'$  is the set of variables in  $\text{body}_{Q'}$ .

**Theorem 3.1.7** *Deciding whether a given conjunctive query implies a given multivalued dependency is NP-complete.*

PROOF. Theorem 3.1.3 establishes that the problem is in NP. We establish NP-hardness using a reduction from testing containment of boolean CQs. Let  $Q_a$  and  $Q_b$  be two boolean CQs whose bodies contain the disjoint sets of variables  $\mathcal{B}_a$  and  $\mathcal{B}_b$ , respectively. Let  $A, Z$  be two fresh variables. Let  $Q(\overline{\mathcal{W}})$  be a new conjunctive query whose output tuple satisfies  $\mathcal{W} = \mathcal{B}_a \cup \{A, Z\}$ , and whose body is defined as follows.

$$\text{body}_Q = \text{body}_{Q_a} \cup \text{body}_{Q_b} \cup \bigcup_{x \in \mathcal{B}_a \cup \mathcal{B}_b} \{R(A, x), R(x, Z)\}$$

Then,  $Q_a \subseteq Q_b$  iff there exists a homomorphism  $h : \mathcal{B}_b \rightarrow \mathcal{B}_a$  such that  $h(\text{body}_{Q_b}) \subseteq \text{body}_{Q_a}$  iff  $Q$  implies  $\mathcal{B}_a \twoheadrightarrow A$  (and  $\mathcal{B}_a \twoheadrightarrow Z$ ). The last step follows from Corollary 3.1.6.  $\square$

**Example 25** Assume a database schema containing relations  $\text{Student}(\text{school}, \text{dept}, \text{name})$  and  $\text{Prof}(\text{school}, \text{dept}, \text{name})$ , and consider the following two conjunctive queries.

$$\begin{aligned} Q_{25}(A, B, C) &:- \text{Prof}('UW', A, B), \text{Student}('UW', A, C) \\ Q_{26}(A, B, C) &:- \text{Prof}('UW', A, B), \text{Student}('UW', A, C), \\ &\quad \text{Prof}(F, D, B), \text{Student}(F, E, C) \end{aligned}$$

Figure 3.1 shows a sample database instance  $\mathbb{D}_5$ , while Figure 3.2 shows the relation  $R_{10}$  resulting from evaluating either query  $Q_{25}$  or  $Q_{26}$  over  $\mathbb{D}_5$ . The reader can verify that relation  $R_{10}$  satisfies  $A \twoheadrightarrow C$ .



A	B	C
CS	Ashraf	Dave
CS	Ashraf	Xin
CS	David	Dave
CS	David	Xin
CS	Frank	Dave
CS	Frank	Xin
ENG	Linda	Cara
ENG	Linda	Julia
ENG	John	Cara
ENG	John	Julia

Figure 3.2: Query result  $R_{10} = (Q_{25})^{\mathbb{D}_5} = (Q_{26})^{\mathbb{D}_5}$

Query  $Q_{27}$  embodies the appropriate query for testing directly—per Theorem 3.1.3—whether  $Q_{25}$  satisfies the MVD  $A \twoheadrightarrow C$ .

$$Q_{27}(A, B, C) :- \text{Prof}('UW', A, B), \text{Student}('UW', A, C'), \\ \text{Prof}('UW', A, B''), \text{Student}('UW', A, C)$$

The following homomorphisms  $h : Q_{25} \rightarrow Q_{27}$  and  $h' : Q_{27} \rightarrow Q_{25}$  prove that  $Q_{25}$  and  $Q_{27}$  are equivalent, and so  $Q_{25}$  satisfies  $A \twoheadrightarrow C$ .

$$h := \{A/A, B/B, C/C\} \quad h' := \{A/A, B/B, C/C, C'/C, B''/B\}$$

Similarly, query  $Q_{28}$  embodies the appropriate query for testing directly—per Theorem 3.1.3—whether  $Q_{26}$  satisfies the MVD  $A \twoheadrightarrow C$ .

$$Q_{28}(A, B, C) :- \text{Prof}('UW', A, B), \text{Student}('UW', A, C'), \\ \text{Prof}(F', D', B), \text{Student}(F', E', C') \\ \text{Prof}('UW', A, B''), \text{Student}('UW', A, C), \\ \text{Prof}(F'', D'', B''), \text{Student}(F'', E'', C)$$

The following homomorphisms  $h'' : Q_{26} \rightarrow Q_{28}$  and  $h''' : Q_{28} \rightarrow Q_{26}$  prove that  $Q_{26}$  and  $Q_{28}$  are equivalent, and so  $Q_{26}$  also satisfies  $A \twoheadrightarrow C$ .

$$h'' := \{A/A, B/B, C/C, D/A, E/A, F/'UW'\} \\ h''' := \{A/A, B/B, B''/B, C/C, C'/C, D'/D, D''/D, E'/E, E''/E, F'/F, F''/F\}$$

Alternatively, we could have deduced that  $Q_{26}$  satisfies  $A \twoheadrightarrow C$  by instead recognizing that  $Q_{25}$  and  $Q_{26}$  are equivalent ( $Q_{25}$  is a minimal form of  $Q_{26}$ ).

## A Hypergraph-based Characterization of Query-Implied MVDs

Corollary 3.1.4 provides a straightforward homomorphic condition for testing query-implied MVDs. However, homomorphisms are sometimes difficult to visualize for

large queries. We now provide an alternate characterization of query-implied dependencies based upon the topology of a query’s hypergraph; this characterization is both easier to visualize and provides more intuition into for the role that MVDs play in the normal form we define later in this section. It also simplifies many of our proofs. As a simple example, we can see that  $Q_{25}$  in Example 25 implies  $A \twoheadrightarrow C$  because variables  $C$  and  $B$  are “separated” by variable  $A$  within  $\text{body}_{Q_{25}}$ . We now formalize this idea of “separation” using graph-theoretic concepts.

Given an undirected graph  $G = (U, E)$  and some subset of its vertices<sup>2</sup>  $U' \subseteq U$ , the subgraph of  $G$  induced by  $U'$ —denoted  $G|_{U'}$ —is formed from  $G$  by deleting the vertices not in  $U'$  as well as any edges incident to them. Given vertices  $s, t \in U$ , an  $(s, t)$ -path in  $G$  is a sequence of  $k \geq 1$  adjacent edges starting at  $s$  and ending at  $t$ . Given  $X \subset U$  and  $s, t \in (U \setminus X)$ , we call  $X$  an  $(s, t)$ -articulation set of  $G$  if the subgraph  $G|_{U \setminus X}$  does not contain any  $(s, t)$ -paths.<sup>3</sup> Given disjoint sets  $X, S, T \subset U$ , we call  $X$  a *weak*  $(S, T)$ -articulation set if there exists some pair  $s \in S$  and  $t \in T$  for which  $X$  is an  $(s, t)$ -articulation set, and we call  $X$  a *strong*  $(S, T)$ -articulation set if it is an  $(s, t)$ -articulation set for every pair  $s \in S$  and  $t \in T$ .

A (undirected) *hypergraph* is a pair  $H = (U, F)$  where  $U$  is a set of vertices and  $F$  is a set of non-empty subsets of  $U$ , called *hyperedges*. Given some set  $U' \subseteq U$ , the induced subgraph  $H|_{U'}$  is formed from  $H$  by deleting the vertices not in  $U'$  and then deleting any hyperedge that has become empty. Given vertices  $s, t \in U$ , an  $(s, t)$ -path in  $H$  is a sequence of  $k \geq 1$  hyperedges  $f_1, \dots, f_k$  such that  $s \in f_1$ ,  $t \in f_k$ , and  $f_j \cap f_{j+1} \neq \emptyset$  for  $j \in [1, k - 1]$  [6, Sect. 6.4]. We define the  $(s, t)$ - and  $(S, T)$ -articulation sets for hypergraphs in terms of  $(s, t)$ -paths exactly as we did for standard graphs.

**Example 26** Consider the two hypergraphs depicted in Figure 3.3. The set  $\{A\}$  is a  $(B, C)$ -articulation set in  $H_{25}$  but not in  $H_{26}$ . Within  $H_{26}$ , the set  $\{B, F\}$  is a weak (but not strong)  $(E, AD)$ -articulation set, while  $\{C, F\}$  is a strong  $(E, AD)$ -articulation set.

Associated with each conjunctive query  $Q$  we define the *query hypergraph*  $H^Q$  as follows. The query variables  $\mathcal{B}$  from  $Q$  form the vertex set of  $H^Q$ . The hyperedges correspond to atoms of  $\text{body}_Q$  such that for each atom  $R_i(\overline{X}_i)$  there is a hyperedge  $f$  containing precisely the set of variables in  $\overline{X}_i$ .

**Example 27** The query hypergraphs  $H^{Q_{25}}$  and  $H^{Q_{26}}$  are equal to the hypergraphs  $H_{25}$  and  $H_{26}$ , respectively, from Figure 3.3. Observe that the query constant ‘UW’ does not appear in the hypergraphs.

<sup>2</sup>To avoid confusion, we use  $U$  rather than  $V$  to denote the set of graph vertices since we already use  $\mathcal{V}$  for the output attributes of CEQs.

<sup>3</sup>An articulation vertex (or set) is typically defined as vertex (or set of vertices) whose deletion increases the number of connected components in the graph [11, 58]. In network flow analysis, an  $(s, t)$ -cut is defined as a set of *edges* whose deletion disconnects  $s$  from  $t$  [84, 93]. Although our definition of an  $(s, t)$ -articulation set is similar in spirit to each of these concepts, it differs in the case where the initial graph is disconnected—specifically, we do not assume that  $s$  and  $t$  were connected in the original graph  $G$ .

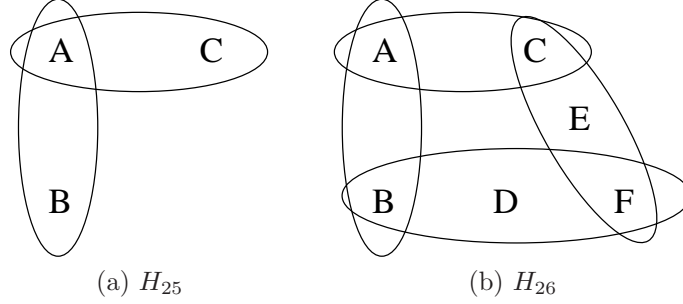


Figure 3.3: Two hypergraphs

Given a hypergraph  $H = (U, F)$  and a database  $\mathbb{D}$ , an *embedding of  $H$  into  $\mathbb{D}$*  is a mapping  $\gamma : U \rightarrow \text{dom}$  such that for each hyperedge  $f \in F$ , there exists some tuple  $t$  in some relation in  $\mathbb{D}$  such that all of the constants in  $\gamma(f)$  occur in  $t$ .

**Lemma 3.1.8** *Given a database  $\mathbb{D}$  and conjunctive query  $Q$ , any query embedding  $\gamma : \text{body}_Q \rightarrow \mathbb{D}$ , is also an embedding of hypergraph  $H^Q$  into  $\mathbb{D}$ .*

**Lemma 3.1.9** *Given a database  $\mathbb{D}$ , a hypergraph  $H = (U, F)$ , and a set  $U' \subseteq U$ , any hypergraph embedding  $\gamma : U \rightarrow \mathbb{D}$  is also an embedding of the induced hypergraph  $H|_{U'}$  into  $\mathbb{D}$ .*

**Theorem 3.1.10** *Given a minimal conjunctive query  $Q(\overline{\mathcal{W}})$  and disjoint attribute sets  $X, Y \subset \mathcal{W}$  query  $Q$  implies MVD  $X \twoheadrightarrow Y$  iff  $X$  is a strong  $(Y, Z)$ -articulation set of the hypergraph  $H^Q$ , where  $Z = \mathcal{W} \setminus (X \cup Y)$ ,*

PROOF. Follows from Corollary 3.1.4, because any homomorphism  $h : Q \rightarrow Q$  is necessarily isomorphic over  $\mathcal{B}$  due to the minimality of  $Q$ .  $\square$

It is trivial to extend Theorem 3.1.10 (and Corollary 3.1.4) to allow constants in sets  $X$  and  $Y$ ; we assume this slight extension for the MVDs in the next section.

**Example 28** Queries  $Q_{25}$  and  $Q_{26}$  both imply  $A \twoheadrightarrow C$  because  $A$  is a  $(B, C)$ -articulation set of  $H^{Q_{25}}$ . In contrast,  $A$  is *not* a  $(B, C)$ -articulation set of  $H^{Q_{26}}$ .

### 3.1.2 Converting CEQs to Normal Form

We now define a normal form for CEQs based upon query-implied MVDs. Our normal form is defined by recursively identifying the *core* indexes  $\mathcal{I}_i^{\overline{\S}} \subseteq \mathcal{I}_i$  at each level  $i$ .

**Definition 3.1.11 (Core Indexes)** *Given a CEQ  $Q(\overline{\mathcal{L}}_1; \dots; \overline{\mathcal{L}}_d; \overline{\mathcal{V}})$ , a length- $d$  signature  $\overline{\S}$ , and an integer  $i \in [1, d]$ , define the core indexes at level  $i$  relative to  $\overline{\S}$ —denoted  $\mathcal{I}_i^{\overline{\S}}$ —as follows. Let  $Q_{(i)}$  be the following CQ.*

$$Q_{(i)}(\mathcal{I}_{[1,i]}^{\overline{\S}} \mathcal{I}_{[i+1,d]}^{\overline{\S}}) :- \text{body}_Q \quad (3.3)$$

Set  $\mathcal{I}_i^{\bar{s}}$  is a minimal subset of  $\mathcal{I}_i$  satisfying both  $\mathcal{V} \subseteq \mathcal{I}_{[1,i-1]} \cup \mathcal{I}_{[i,d]}^{\bar{s}} \cup \mathcal{C}$  and the following signature-specific conditions.

$\xi_i$	Condition
<b>b</b>	$\mathcal{I}_i \subseteq \mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{s}} \cup \mathcal{C}$
<b>s</b>	$Q_{(i)} \models (\mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{s}}) \twoheadrightarrow \mathcal{I}_{[i+1,d]}^{\bar{s}}$
<b>n</b>	$Q_{(i)} \models \mathcal{I}_{[1,i-1]} \twoheadrightarrow \mathcal{I}_{[i,d]}^{\bar{s}}$

We denote by  $\bar{\mathcal{I}}_i^{\bar{s}}$  the tuple formed by deleting from  $\bar{\mathcal{I}}_i$  all attributes not in  $\mathcal{I}_i^{\bar{s}}$ .

**Observation 3.1.12** *Constant columns are never core; that is,  $\bar{\mathcal{I}}_{[1,d]}^{\bar{s}} \subseteq \mathcal{B}$ .*

**Observation 3.1.13** *The sets of core indexes at different levels are disjoint.*

**Observation 3.1.14** *If  $\xi_d = \mathbf{s}$ , then  $\mathcal{I}_d^{\bar{s}} = (\mathcal{I}_d \cap \mathcal{V} \cap \mathcal{B}) \setminus \mathcal{I}_{[1,i-1]}$ .*

**Lemma 3.1.15** *Definition 3.1.11 determines a unique minimal set  $\mathcal{I}_i^{\bar{s}}$ .*

PROOF. First, let  $Q'_{(i)}(\mathcal{I}_{[1,i]} \bar{\mathcal{I}}_{[i+1,d]}^{\bar{s}})$  be a minimal CQ equivalent to  $Q_{(i)}$  from equation 3.3. Next, let a “candidate for  $\mathcal{I}_i^{\bar{s}}$ ” denote any set  $X \subseteq \mathcal{I}_i$  such that if we choose  $\mathcal{I}_i^{\bar{s}} := X$  then all conditions in Definition 3.1.11 except for minimality are satisfied. We now show that if sets  $X_1, X_2$  are both candidates for  $\mathcal{I}_i^{\bar{s}}$ , then  $X_1 \cap X_2$  is also a candidate for  $\mathcal{I}_i^{\bar{s}}$ . If  $\mathcal{V} \subseteq \mathcal{I}_{[1,i-1]} \cup X_1 \cup \mathcal{C}$  and  $\mathcal{V} \subseteq \mathcal{I}_{[1,i-1]} \cup X_2 \cup \mathcal{C}$  then  $\mathcal{V} \subseteq \mathcal{I}_{[1,i-1]} \cup (X_1 \cap X_2) \cup \mathcal{C}$  follows directly. Now consider the signature-specific conditions.

**Case  $\xi_i = \mathbf{b}$ :**

$\mathcal{I}_i \subseteq \mathcal{I}_{[1,i-1]} \cup X_1$  and  $\mathcal{I}_i \subseteq \mathcal{I}_{[1,i-1]} \cup X_2$  implies  $\mathcal{I}_i \subseteq \mathcal{I}_{[1,i-1]} \cup (X_1 \cap X_2)$ .

**Case  $\xi_i = \mathbf{s}$ :**

The MVD  $Q_{(i)} \models \mathcal{I}_{[1,i-1]} \cup (X_1 \cap X_2) \twoheadrightarrow \mathcal{I}_{[i+1,d]}^{\bar{s}}$  cannot be derived from axioms for MVDs [6, Sect. 8.3], but can be reasoned from the query structure as follows.

1. By candidacy of  $X_1$ ,  $Q_{(i)} \models (\mathcal{I}_{[1,i-1]} \cup X_1) \twoheadrightarrow \mathcal{I}_{[i+1,d]}^{\bar{s}}$ . By Theorem 3.1.10,  $\mathcal{I}_{[1,i-1]} \cup X_1$  is a strong  $(\mathcal{I}_{[i+1,d]}^{\bar{s}}, (\mathcal{I}_i \setminus X_1))$ -articulation set of  $H^{Q'_{(i)}}$ .
2. By candidacy of  $X_2$ ,  $Q_{(i)} \models \mathcal{I}_{[1,i-1]} \cup X_2 \twoheadrightarrow \mathcal{I}_{[i+1,d]}^{\bar{s}}$ . By Theorem 3.1.10,  $\mathcal{I}_{[1,i-1]} \cup X_2$  is a strong  $(\mathcal{I}_{[i+1,d]}^{\bar{s}}, (\mathcal{I}_i \setminus X_2))$ -articulation set of  $H^{Q'_{(i)}}$ .
3. The two articulation sets together imply that deleting  $\mathcal{I}_{[1,i-1]} \cup (X_1 \cap X_2)$  from  $H^{Q'_{(i)}}$  causes the two sets  $\mathcal{I}_{[i+1,d]}^{\bar{s}}$  and  $\mathcal{I}_i \setminus (X_1 \cap X_2)$  to occur in separate partitions of the remaining hypergraph. The MVD then follows from Theorem 3.1.10.

**Case  $\xi_i = \mathbf{n}$ :**

Similar to case  $\xi_i = \mathbf{s}$ . The candidacies of  $X_1$  and  $X_2$  imply two MVDs which imply two strong articulation set of  $H^{Q^{(i)}}$ . This time, the two articulation sets together imply that deleting  $\mathcal{I}_{[1,i-1]}$  from  $H^{Q^{(i)}}$  causes the four sets  $X_1 \setminus X_2$ ,  $X_2 \setminus X_1$ ,  $\mathcal{I}_i \setminus (X_1 \cup X_2)$ , and  $(X_1 \cap X_2) \cup \mathcal{I}_{[i+1,d]}^{\bar{\xi}}$  to occur in separate partitions of the remaining hypergraph. The MVD  $Q^{(i)} \models \mathcal{I}_{[1,i-1]} \rightarrow (X_1 \cap X_2) \cup \mathcal{I}_{[i+1,d]}^{\bar{\xi}}$  then follows from Theorem 3.1.10. □

We say that CEQ  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  is in  $\bar{\xi}$ -normal form ( $\bar{\xi}$ -NF) if  $\mathcal{I}_{[1,d]} = \mathcal{I}_{[1,d]}^{\bar{\xi}}$ , and so an arbitrary CEQ is converted to  $\bar{\xi}$ -NF by deleting all non-core indexes from the query head. Assuming that the original query was valid (cf. Lemma 3.0.1), Definition 3.1.11 implies  $(\mathcal{V} \cap \mathcal{B}) \subseteq \mathcal{I}_{[1,d]}^{\bar{\xi}}$  and so the new query is guaranteed to be valid.

**Theorem 3.1.16** *Performing  $\bar{\xi}$ -normalization of CEQs is NP-complete.*

PROOF. NP-hardness follows directly from Theorem 3.1.7. Identifying the core indexes at each level can be done in NP time using an algorithm that traverses query hypergraphs (due to Observation 3.1.12, it is not a problem that query hypergraphs do not contain constants).

**Case  $\xi_i = \mathbf{b}$ :** Trivial.

**Case  $\xi_i = \mathbf{s}$ :** Minimize the body of  $Q^{(i)}$ , then construct hypergraph  $H^{Q^{(i)}}$ . Delete from  $H^{Q^{(i)}}$  all nodes corresponding to variables in the set  $\mathcal{I}_{[1,i-1]} \cup (\mathcal{I}_i \cap \mathcal{V})$ .

Identify any non-output members of  $\mathcal{I}_i^{\bar{\xi}}$  incrementally by traversing the connected components containing  $\mathcal{I}_{[i+1,d]}^{\bar{\xi}}$  and deleting the “nearest” member of  $\mathcal{I}_i$ .

**Case  $\xi_i = \mathbf{n}$ :** Minimize the body of  $Q^{(i)}$ , then construct hypergraph  $H^{Q^{(i)}}$ . Delete from  $H^{Q^{(i)}}$  all nodes corresponding to variables in the set  $\mathcal{I}_{[1,i-1]}$ . Identify  $\mathcal{I}_i^{\bar{\xi}}$  by traversing the connected components containing any variable in  $(\mathcal{I}_i \cap \mathcal{V}) \cup \mathcal{I}_{[i+1,d]}^{\bar{\xi}}$ . □

**Example 29** Consider the following four CEQs.

$$\begin{aligned}
Q_{29}(\overbrace{A, B}^{\bar{\mathcal{I}}_1}; \overbrace{C}^{\bar{\mathcal{I}}_2}; \overbrace{C}^{\bar{\mathcal{V}}}) &:- \text{Prof}('UW', A, B), \text{Student}('UW', A, C) \\
Q_{30}(A; C; C) &:- \text{Prof}('UW', A, B), \text{Student}('UW', A, C) \\
Q_{31}(A; B, C; C) &:- \text{Prof}('UW', A, B), \text{Student}('UW', A, C), \\
&\quad \text{Prof}(F, D, B), \text{Student}(F, E, C) \\
Q_{32}(A; C; C) &:- \text{Prof}('UW', A, B), \text{Student}('UW', A, C), \\
&\quad \text{Prof}(F, D, B), \text{Student}(F, E, C)
\end{aligned}$$

A	B	C	C
CS	Ashraf	Dave	Dave
CS	Ashraf	Xin	Xin
CS	David	Dave	Dave
CS	David	Xin	Xin
CS	Frank	Dave	Dave
CS	Frank	Xin	Xin
ENG	Linda	Cara	Cara
ENG	Linda	Julia	Julia
ENG	John	Cara	Cara
ENG	John	Julia	Julia

A	B	C	C
CS	Ashraf	Dave	Dave
CS	Ashraf	Xin	Xin
CS	David	Dave	Dave
CS	David	Xin	Xin
CS	Frank	Dave	Dave
CS	Frank	Xin	Xin
ENG	Linda	Cara	Cara
ENG	Linda	Julia	Julia
ENG	John	Cara	Cara
ENG	John	Julia	Julia

A	C	C
CS	Dave	Dave
CS	Xin	Xin
ENG	Cara	Cara
ENG	Julia	Julia

Figure 3.4: Query results  $R_{11} = (Q_{29})^{\mathbb{D}_5}$ ,  $R_{12} = (Q_{31})^{\mathbb{D}_5}$ , and  $R_{13} = (Q_{30})^{\mathbb{D}_5} = (Q_{32})^{\mathbb{D}_5}$ .

Observe that  $Q_{29}$  and  $Q_{30}$  have the same body as  $Q_{25}$  from Example 28, while  $Q_{31}$  and  $Q_{32}$  have the same body as  $Q_{26}$ . Therefore, both  $Q_{29}$  and  $Q_{31}$  imply  $A \rightarrow C$ .

Consider signature  $sb$ . The dependency  $A \rightarrow C$  implies that  $B$  is not a core index of  $Q_{29}$  relative to  $sb$ . In contrast, queries  $Q_{30}$ ,  $Q_{31}$ , and  $Q_{32}$  are all in  $sb$ -NF.

Now consider signature  $nn$ . This time queries  $Q_{29}$ ,  $Q_{30}$ , and  $Q_{32}$  are already in  $nn$ -NF, but  $Q_{31}$  is not because  $A \rightarrow C$  implies that  $B$  is not a core index.

The next theorem establishes that  $\bar{\xi}$ -normalization does not change the complex object obtained after decoding the CEQ result with signature  $\bar{\xi}$ . Before we prove this formally, however, consider the following example which we will use to illustrate the core ideas of the proof.

**Example 30** Figure 3.4 shows the encoding relations that result when the four queries from Example 29 are evaluated over the database  $\mathbb{D}_5$  from Figure 3.1. Query  $Q_{30}$  is the  $sb$ -normal form of  $Q_{29}$ , and the reader can verify that both encoding relations  $R_{13}$  and  $R_{11}$  yield the object  $\{ \{ \{ \text{Dave}, \text{Xin} \} \}, \{ \{ \text{Cara}, \text{Julia} \} \} \}$  when decoded with signature  $sb$ . Similarly,  $Q_{32}$  is the  $nn$ -normal form of  $Q_{31}$ , and both encoding relations  $R_{13}$  and  $R_{12}$  yield the object  $\{ \{ \{ \text{Dave}, \text{Xin} \} \}, \{ \{ \text{Cara}, \text{Julia} \} \} \}$  when decoded with signature  $nn$ .

**Theorem 3.1.17**  $\bar{\xi}$ -Normalization preserves  $\bar{\xi}$ -equivalence.

PROOF. Let  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  CEQ, and let  $Q'(\bar{\mathcal{I}}_1^{\bar{\xi}}; \dots; \bar{\mathcal{I}}_d^{\bar{\xi}}; \bar{\mathcal{V}})$  be its  $\bar{\xi}$ -normal form. For every  $i \in [0, d]$ , let  $Q_i$  denote the following CEQ.

$$Q_i(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_i; \bar{\mathcal{I}}_{i+1}^{\bar{\xi}}; \dots; \bar{\mathcal{I}}_d^{\bar{\xi}}; \bar{\mathcal{V}}) :- \text{body}_Q \quad (3.4)$$

Because  $Q' = Q_0$ , we need to show that  $Q_0 \doteq_{\bar{\xi}} Q$ , which we do by induction.

**Base Case:** Equation 3.4 implies  $Q_d = Q$ , and so  $Q_d \doteq_{\bar{\xi}} Q$  is trivial.

**Induction Hypothesis:** Suppose there exists some  $i \in [1, d]$  such that  $Q_i \doteq_{\bar{\xi}} Q$ .

**Inductive Step:** We need to show that  $Q_{i-1} \doteq_{\bar{\xi}} Q_i$ . For any database  $\mathbb{D}$ , let  $\Gamma$  and  $\Phi$  denote the sets of embeddings of  $\text{body}_{Q_{i-1}}$  and  $\text{body}_{Q_i}$  into  $\mathbb{D}$ ; because  $Q_{i-1}$  and  $Q_i$  have the same body,  $\Gamma$  and  $\Phi$  contain the same embeddings. W.l.o.g., assume that  $\bar{\mathcal{I}}_i = \bar{\mathcal{I}}_i^{\bar{\xi}} \cdot \bar{\mathcal{J}}_i$ , where  $\bar{\mathcal{J}}_i$  is the set of non-core indexes in  $\bar{\mathcal{I}}_i$ . The encoding relations  $(Q_{i-1})^{\mathbb{D}}$  and  $(Q_i)^{\mathbb{D}}$  can now be characterized as follows.

$$\mathbb{S}((Q_{i-1})^{\mathbb{D}}) = Q_{i-1}(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_{i-1}; \bar{\mathcal{I}}_i^{\bar{\xi}}; \bar{\mathcal{I}}_{i+1}^{\bar{\xi}}; \dots; \bar{\mathcal{I}}_d^{\bar{\xi}}; \bar{\mathcal{V}}) \quad (3.5)$$

$$\mathbb{T}((Q_{i-1})^{\mathbb{D}}) = \{\gamma[\bar{\mathcal{I}}_{[1, i-1]} \bar{\mathcal{I}}_{[i, d]}^{\bar{\xi}} \bar{\mathcal{V}}] \mid \gamma \in \Gamma\} \quad (3.6)$$

$$\mathbb{S}((Q_i)^{\mathbb{D}}) = Q_i(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_{i-1}; \bar{\mathcal{I}}_i^{\bar{\xi}} \bar{\mathcal{J}}_i; \bar{\mathcal{I}}_{i+1}^{\bar{\xi}}; \dots; \bar{\mathcal{I}}_d^{\bar{\xi}}; \bar{\mathcal{V}}) \quad (3.7)$$

$$\mathbb{T}((Q_i)^{\mathbb{D}}) = \{\phi[\bar{\mathcal{I}}_{[1, i-1]} \bar{\mathcal{I}}_i^{\bar{\xi}} \bar{\mathcal{J}}_i \bar{\mathcal{I}}_{[i+1, d]}^{\bar{\xi}} \bar{\mathcal{V}}] \mid \phi \in \Phi\} \quad (3.8)$$

Abusing notation, we will now use  $\Gamma$  and  $\Phi$  in place of the encoding relations  $(Q_{i-1})^{\mathbb{D}}$  and  $(Q_i)^{\mathbb{D}}$  by implicitly associating schema  $\mathbb{S}((Q_{i-1})^{\mathbb{D}})$  with  $\Gamma$  and  $\mathbb{S}((Q_i)^{\mathbb{D}})$  with  $\Phi$ . Then, we need to show that  $\Gamma \doteq_{\bar{\xi}} \Phi$ .

**Case  $\xi_i = \mathbf{b}$ :**

By Definition 3.1.11,  $\bar{\mathcal{J}}_i \subseteq \bar{\mathcal{I}}_{[1, i-1]} \cup \mathcal{C}$ . Because  $\bar{\mathcal{I}}_{[1, i-1]}$  functionally determines  $\bar{\mathcal{J}}_i$ , removing  $\bar{\mathcal{J}}_i$  from the head does not affect the cardinality or contents of any encoding sub-relations; hence, there is a trivial isomorphism between the tuples of  $\Gamma$  and  $\Phi$  that proves  $\Gamma \doteq_{\bar{\xi}} \Phi$ .

**Case  $\xi_i = \mathbf{s}$ :**

For each value  $\bar{a} \in \text{adom}(\bar{\mathcal{I}}_{[1, i-1]}, \Gamma) = \text{adom}(\bar{\mathcal{I}}_{[1, i-1]}, \Phi)$  let  $C_{\bar{a}}$  be a  $\bar{\xi}_{[i, d]}$ -certificate rooted by an initially empty set node (see Definition 2.2.5). We will incrementally construct  $C_{\bar{a}}$  until it proves that  $\Gamma[\bar{a}] \doteq_{\bar{\xi}_{[i, d]}} \Phi[\bar{a}]$ . Because  $\text{adom}(\bar{\mathcal{I}}_{[1, i-1]}, \Gamma) = \text{adom}(\bar{\mathcal{I}}_{[1, i-1]}, \Phi)$ , it is then trivial to construct the upper levels of a  $\bar{\xi}$ -certificate proving  $\Gamma \doteq_{\bar{\xi}} \Phi$ .

Let  $\{\bar{b}_1, \dots, \bar{b}_k\}$  be the set of values in  $\text{adom}(\bar{\mathcal{I}}_i^{\bar{\xi}}, \Gamma[\bar{a}]) = \text{adom}(\bar{\mathcal{I}}_i^{\bar{\xi}}, \Phi[\bar{a}])$ , and let  $\{\bar{b}_1 \bar{c}_1^1, \dots, \bar{b}_1 \bar{c}_1^{l_1}, \dots, \bar{b}_k \bar{c}_k^1, \dots, \bar{b}_k \bar{c}_k^{l_k}\}$  be the set of values in  $\text{adom}(\bar{\mathcal{I}}_i^{\bar{\xi}} \bar{\mathcal{J}}_i, \Phi[\bar{a}])$ . For each  $\bar{b}_j$ , the sub-relation  $\Gamma[\bar{a} \bar{b}_j]$  encodes an object of sort  $(\bar{\xi}_{[i+1, d]}, d-i)$  occurring in the set at level  $i$ . By Definition 3.1.11, relation  $\Phi$  satisfies  $\bar{\mathcal{I}}_{[1, i-1]} \cup \bar{\mathcal{I}}_i^{\bar{\xi}} \twoheadrightarrow \bar{\mathcal{I}}_{[i+1, d]}^{\bar{\xi}}$ , which—by Corollaries 3.1.5 and 3.1.6—implies that

$\Phi[\bar{a}]$  satisfies  $\mathcal{I}_i^{\bar{s}} \rightarrow \mathcal{I}_{[i+1,d]}^{\bar{s}}$  (and  $\mathcal{I}_i^{\bar{s}} \rightarrow \mathcal{J}_i$ ). It follows from that all of the sub-relations  $\Phi[\bar{a}\bar{b}_j\bar{c}_j^1], \dots, \Phi[\bar{a}\bar{b}_j\bar{c}_j^{l_j}]$  are identical to each other and to the sub-relation  $\Gamma[\bar{a}\bar{b}_j]$ . For each  $\bar{c}_j^{l_j}$  add to  $C_{\bar{a}}$  the mapping  $f(\bar{b}_j\bar{c}_j^{l_j}) := \bar{b}_j$ , as well as a child  $\bar{s}_{[i+1,d]}$ -certificate proving that  $\Phi[\bar{a}\bar{b}_j\bar{c}_j^{l_j}] \doteq_{\bar{s}_{[i+1,d]}} \Gamma[\bar{a}\bar{b}_j]$  (which is trivial, because relations  $\Phi[\bar{a}\bar{b}_j\bar{c}_j^{l_j}] = \Gamma[\bar{a}\bar{b}_j]$ ). Then, add the mapping  $f'(\bar{b}_j) := \bar{b}_j \cdot \bar{c}_j^1$  (we could choose any  $\bar{c}_j^{l_j}$ ). Certificate  $C_{\bar{a}}$  is complete when this has been performed for all values of  $\bar{b}_j$ .

**Case  $\bar{s}_i = \mathbf{n}$ :**

The proof is almost identical to case  $\bar{s}_i = \mathbf{s}$ . Because  $\mathcal{I}_{[1,i-1]} \rightarrow \mathcal{I}_{[i,d]}^{\bar{s}}$  implies  $\mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{s}} \rightarrow \mathcal{I}_{[i+1,d]}^{\bar{s}}$ , all of the relations  $\Phi[\bar{a}\bar{b}_j\bar{c}_j^1], \dots, \Phi[\bar{a}\bar{b}_j\bar{c}_j^{l_j}]$  are identical to each other and to the sub-relation  $\Gamma[\bar{a}\bar{b}_j]$ . Furthermore,  $\mathcal{I}_{[1,i-1]} \rightarrow \mathcal{I}_{[i,d]}^{\bar{s}}$  also implies that across all values of  $j$ , the sets  $\{\bar{c}_j^1, \dots, \bar{c}_j^{l_j}\}$  are identical. Therefore, the degree of replication remains constant across all values of  $\bar{b}_j \in \text{adom}(\bar{\mathcal{I}}_i^{\bar{s}}, \Gamma[\bar{a}])$ .

□

## 3.2 CEQ Encoding-Equivalence

In this section we fully characterize  $\bar{s}$ -equivalence of CEQs by tailoring the traditional homomorphism test for CQs to encoding queries.

**Definition 3.2.1 (Index-Covering Homomorphism)** *Given two conjunctive encoding queries  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  and  $Q'(\bar{\mathcal{I}}'_1; \dots; \bar{\mathcal{I}}'_d; \bar{\mathcal{V}}')$  of equal depth, an index-covering homomorphism from  $Q'$  to  $Q$  is a mapping  $h : \mathcal{B}' \rightarrow \mathcal{B} \cup \mathcal{C}$  (suitably extended to tuples, sets, and atoms, and with identity on constants) satisfying*

1.  $h(\text{body}_{Q'}) \subseteq \text{body}_Q$ ,
2.  $h(\bar{\mathcal{V}}') = \bar{\mathcal{V}}$ , and
3.  $\forall i \in [1, d]: \mathcal{I}_i \subseteq h(\mathcal{I}'_i)$ .

Comparing Definitions 3.2.1 and 3.1.1, we note that if  $h$  is an index-covering homomorphism from CEQ  $Q'$  to CEQ  $Q$ , then  $h$  is also a CQ homomorphism  $h : Q'(\bar{\mathcal{V}}') \rightarrow Q(\bar{\mathcal{V}})$ , where  $Q'(\bar{\mathcal{V}}')$  and  $Q(\bar{\mathcal{V}})$  denote the CQs formed from  $Q'$  and  $Q$  by restricting the head tuples to  $\bar{\mathcal{V}}'$  and  $\bar{\mathcal{V}}$ , respectively. Conversely, if  $h'$  is a CQ homomorphism from CQ  $Q'''(\bar{\mathcal{W}}''')$  to CQ  $Q''(\bar{\mathcal{W}}'')$ , then  $h'$  is also an index-covering homomorphism from the depth-one CEQ  $Q'''(\bar{\mathcal{W}}'''; \bar{\mathcal{W}}''')$  to the depth-one CEQ  $Q''(\bar{\mathcal{W}}''; \bar{\mathcal{W}}'')$ . For this reason, we use the same notation  $h : Q' \rightarrow Q$  to denote index-covering homomorphisms as for CQ homomorphisms, and we will clearly indicate within the context which type of homomorphism (and queries) we are



discussing. The next theorem and its corollary follow from well-known results for conjunctive queries.

**Theorem 3.2.2** *Given conjunctive encoding queries  $Q$  and  $Q'$ , each index-covering homomorphism from  $Q'$  to  $Q$  corresponds to a tuple in the result of evaluating  $Q'$  over a database isomorphic to the body of  $Q$ .*

**Corollary 3.2.3** *Determining the existence of an index-covering homomorphism is NP-complete.*

We now combine our tailored homomorphism with the normal form from the previous section to obtain a necessary and sufficient test for  $\bar{\xi}$ -equivalence. The necessity proof for the following theorem is intricate—as we will discuss further in Section 3.3, just considering nested sets is already complex [77], and handling arbitrary signatures requires integrating three different proof techniques (for sets, bags, and normalized bags). For readability, we give only a short sketch here and defer the full proof of necessity to Appendix A.

**Theorem 3.2.4** *Two conjunctive encoding queries are  $\bar{\xi}$ -equivalent iff there exists index-covering homomorphisms in both directions between their  $\bar{\xi}$ -normal forms.*

PROOF (SKETCH). Assume without loss of generality that  $Q$  and  $Q'$  are already in  $\bar{\xi}$ -normal form (justified by Theorem 3.1.17).

If index-covering homomorphisms exist in both directions, then for any database  $\mathbb{D}$  the encoding relations  $(Q)^{\mathbb{D}}$  and  $(Q')^{\mathbb{D}}$  differ at most by permutation of attributes within each index level. A  $\bar{\xi}$ -certificate proving  $(Q)^{\mathbb{D}} \doteq_{\bar{\xi}} (Q')^{\mathbb{D}}$  is therefore straightforward to construct, since each node is simply an isomorphism between sub-relations;  $Q \doteq_{\bar{\xi}} Q'$  follows immediately.

We now sketch the argument for the existence of index-covering homomorphisms (see Appendix A for the full proof). We start by defining a degenerate class of “pseudo-trivial”  $\bar{\xi}$ -certificate nodes that imply certain regularities between the encoding relations that they relate. The overall goal is then to construct a canonical database  $\mathbb{D}$  out of the bodies of both  $Q$  and  $Q'$  such that if we can find a  $\bar{\xi}$ -certificate between  $(Q)^{\mathbb{D}}$  and  $(Q')^{\mathbb{D}}$  that contains these pseudo-trivial nodes, we can use it to generate the index-covering homomorphisms in each direction. We then argue for the existence of such a certificate by proving that *any*  $\bar{\xi}$ -certificate between  $(Q)^{\mathbb{D}}$  and  $(Q')^{\mathbb{D}}$  must contain these trivial nodes. The proof uses induction level-by-level from inside-out to prove triviality of certificate nodes, and the argument differs depending upon the semantics of the current level. Bag equivalence is based upon counting, and so for bag levels we rely upon certain algebraic properties of polynomials combined with a procedure to inflate the cardinalities of constants in any database in a regular way such that non-pseudo-trivial certificate nodes would imply differing element cardinalities, violating bag equivalence. Set equivalence is based upon existence, and so for set levels we use the query bodies to inject large amounts of symmetry into the construction of  $\mathbb{D}$ , with the constants corresponding to the index

variables in  $Q$  and  $Q'$  forming the “axes” of the symmetry; we then show that non-pseudo-trivial certificate nodes would imply the existence (non-existence) of extra (missing) “canonical sub-objects,” violating set equivalence. Finally, normalized bag equivalence exhibits both count- and existence-based features, and so for these levels the argument synthesizes both the polynomial and the symmetry techniques.

□

**Corollary 3.2.5** *Testing  $\bar{\xi}$ -equivalence of CEQs is NP-complete.*

**Corollary 3.2.6** *Testing equivalence of  $CDCQL^{\{f_b, f_s, f_n\}}$ - queries is NP-complete.*

Corollary 3.2.6 is the culmination of this chapter, and one of the major results within this thesis. It is important because it demonstrates that as long as queries do not compare aggregated values, non-scalar forms of nesting do not increase the hardness of the query equivalence problem.

**Example 31** Consider again queries  $Q_{29}$ – $Q_{32}$  from Example 29. Queries  $Q_{30}$  and  $Q_{32}$  are each in both sb-NF and nn-NF, and so the following index-covering homomorphisms  $h : Q_{30} \rightarrow Q_{32}$  and  $h' : Q_{32} \rightarrow Q_{30}$  prove  $Q_{30} \stackrel{\cdot}{\equiv}_{sb} Q_{32}$  and  $Q_{30} \stackrel{\cdot}{\equiv}_{nn} Q_{32}$ .

$$\begin{aligned} h &= \{A/A, B/B, C/C\} \\ h' &= \{A/A, B/B, C/C, F/'UW', D/A, E/A\} \end{aligned}$$

Performing sb-normalization on  $Q_{29}$  yields  $Q_{30}$ , and so  $Q_{29} \stackrel{\cdot}{\equiv}_{sb} Q_{32}$ . In contrast,  $Q_{31}$  is in sb-NF and there clearly cannot exist an index-covering homomorphism from  $Q_{32}$  to  $Q_{31}$ , so  $Q_{31} \not\stackrel{\cdot}{\equiv}_{sb} Q_{32}$  follows by Theorem 3.2.4. By similar analysis we conclude  $Q_{29} \not\stackrel{\cdot}{\equiv}_{nn} Q_{30} \stackrel{\cdot}{\equiv}_{nn} Q_{31} \stackrel{\cdot}{\equiv}_{nn} Q_{32}$ . The encoding relations in Figure 3.4 illustrate a counter-example proving both  $Q_{31} \not\stackrel{\cdot}{\equiv}_{sb} Q_{32}$  and  $Q_{29} \not\stackrel{\cdot}{\equiv}_{nn} Q_{30}$ .

### 3.3 Relevant Literature

In this section we survey literature relevant to deciding various forms of equivalence between conjunctive queries operating within a data model containing only flat relational structures. The reader is referred to the previous chapter for literature pertaining to query languages that construct nested objects. In particular, in Section 2.4.4 we discuss Levy and Suciu’s reduction of COQL equivalence to testing *strong simulation* between CQs [77], and so we do not discuss that work further here, even though it is in many ways the closest research to the CEQ encoding-equivalence problem which is our focus in this chapter.

### 3.3.1 Conjunctive Queries under Set Semantics

Equivalence and containment of conjunctive queries under set semantics was first studied by Chandra and Merlin [16] who show that the two problems are mutually reducible, and that they are each mutually reducible to the NP-complete problem of finding homomorphisms between queries. They also show that every query has a unique minimal form, and that the minimal forms of equivalent queries are isomorphic (hence minimization is also NP-complete). Chekuri and Rajaraman [20] take a previously-known result that homomorphisms are computable in PTIME when the query bodies correspond to acyclic hypergraphs [6], and they generalize the notion of acyclicity to *query width*. Their work furnishes algorithms for the containment, equivalence, and minimization problems that run in time polynomial in  $n^k$ , for queries of size  $n$  and width  $k$ . These efficient algorithms are directly applicable to the problems of performing  $\bar{\xi}$ -normalization and finding index-covering homomorphisms that we discuss in this chapter.

### 3.3.2 Other Processing Semantics

Chaudhuri and Vardi [19] and Ioannidis and Ramakrishnan [64] independently propose bag/bag-set semantics as a way to model the input to cardinality-sensitive aggregation functions and proceed to study the query equivalence and containment problems under the new query processing semantics. Under bag semantics, conjunctive queries are equivalent precisely when they are isomorphic; the obvious corollaries are that every conjunctive query is already minimal under bag semantics, and that testing bag equivalence is GraphIsomorphism-complete. Under *bag-set* semantics (bag semantics where database relations are known to be sets), queries are equivalent if and only if they are *variable*-isomorphic; that is, query minimization simply entails removing duplicate relational atoms (linear time), and two queries are equivalent under bag-set semantics precisely when their minimized forms are equivalent under bag semantics. Chaudhuri and Vardi also show a simple reduction from bag equivalence to bag-set equivalence, which justifies our design of  $\text{COCQL}^{\mathcal{F}}$  in Section 2.1.3 as a bag semantic algebra operating over a database of relations (sets of tuples).

Although both Chaudhuri and Vardi as well as Ioannidis and Ramakrishnan characterize the CQ equivalence problem under bag/bag-set semantics, the CQ containment problem under bag/bag-set semantics is still open, although known to be at least  $\Pi_2^p$ -hard [19].<sup>4</sup> Ioannidis and Ramakrishnan show bag containment to be linear time when queries do not contain repeated predicate names, and undecidable when extended to unions of CQs. More recently, Jayram, Kolaitis and Vee show that bag containment is undecidable for CQs extended with inequality predicates (i.e.,  $\neq$ ) [66].

---

<sup>4</sup> $\Pi_2^p$  is a class within the second tier of the polynomial hierarchy. A decision problem is in  $\Pi_2^p$  if its complement problem is in  $\Sigma_2^p$ , the class of problems that permit an NP-time reduction to a problem that is in NP [43, Sect. 7.2].

Variations of bag/bag-set semantics have also been proposed in the literature. Grumbach et al. relax bag-set semantics in order model the input to certain normalizing aggregation functions such as **AVG**. Under their relaxed semantics, two queries are equivalent if over every database they return two bags containing the same elements with the same *relative* cardinalities. They find that CQs are equivalent under these semantics precisely when the query bodies are isomorphic “modulo a product” which roughly means that they differ at most by the addition of certain Cartesian product operations [50]. Cohen also relaxes bag-set semantics in order to study equivalence of SQL queries that contain a nested **SELECT DISTINCT** block. She proposes “combined semantics” as a generalization of bag-set semantics in which cardinality depends only on a specified subset of the query variables, and deduces that CQs are equivalent under these semantics precisely when there exist homomorphisms in each direction that are isomorphic with respect to the designated “multi-set variables” [22].

Our results in this chapter for characterizing encoding-equivalence between conjunctive encoding queries generalizes CQ equivalence under all of these previously-proposed processing semantics. More specifically, given two CQs  $Q(\bar{\mathcal{V}})$  and  $Q'(\bar{\mathcal{V}}')$ , testing  $Q \equiv Q'$  under

- *set semantics* reduces to  $Q(\bar{\mathcal{V}}; \bar{\mathcal{V}}) \dot{\equiv}_s Q'(\bar{\mathcal{V}}'; \bar{\mathcal{V}}')$ ;
- *bag-set semantics* reduces to  $Q(\mathcal{B}; \bar{\mathcal{V}}) \dot{\equiv}_b Q'(\mathcal{B}'; \bar{\mathcal{V}}')$  where  $\mathcal{B}$  and  $\mathcal{B}'$  are the query body variables;
- *bag-set semantics modulo a product* reduces to  $Q(\mathcal{B}; \bar{\mathcal{V}}) \dot{\equiv}_n Q'(\mathcal{B}'; \bar{\mathcal{V}}')$ ; and
- *combined semantics* reduces to  $Q(\mathcal{V} \cup \mathcal{M}; \bar{\mathcal{V}}) \dot{\equiv}_b Q'(\mathcal{V}' \cup \mathcal{M}'; \bar{\mathcal{V}}')$  where  $\mathcal{M}$  and  $\mathcal{M}'$  are the specified multi-set variables.

### 3.3.3 Aggregation Queries

The equivalence and containment problems for queries containing non-nested aggregation have been studied before. Cohen, Nutt, and Sagiv address the equivalence problem for conjunctive queries followed by a single aggregation operation applying one of the functions **SUM**, **COUNT**, **COUNT-DISTINCT**, **MAX**, or **MIN** [25]. When the query bodies do not contain comparison predicates, the equivalence problem for **SUM** and **COUNT** queries reduces to equivalence under bag semantics, while the equivalence problem for **MAX**, **MIN**, and **COUNT-DISTINCT** reduces to equivalence under set semantics. When the query bodies contain comparison predicates, they show the problem to be  $\Pi_2^p$ -complete for **MAX** queries, **MIN** queries, and **COUNT-DISTINCT** queries under a slight syntactic restriction that limits interaction between the aggregated variables and the variables involved in comparison predicates. For **SUM** and **COUNT** queries with comparisons they show the problem to be in PSPACE. These authors also extend their detailed study of aggregation functions to query bodies containing disjunction [26] and safe atomic negation [27]. In Cohen’s doctoral dissertation she uses abstract algebra concepts such as commutative semigroups, idempotency,

and abelian groups to define five properties of aggregation functions—*shiftable*, *order-decidable*, *decomposable*, *singleton-determining*, and *expandable*—with which she taxonomizes commonly-used aggregation functions [21]. She then generalizes previous results for equivalence of queries with specific aggregation functions to black-box functions characterized by this taxonomy, and uses these equivalence results to characterize the problem of rewriting non-nested aggregation queries over non-nested aggregation views, an extension that is fleshed out further within a subsequent journal publication [24].

Grumbach, Rafanelli, and Tininini also consider the problem of rewriting non-nested aggregation queries over non-nested aggregation views, and as part of this work characterize the equivalence problem for queries that apply the **AVG** or **PERCENT** aggregation functions to the result of conjunctive queries [51]. Their characterization reduces this problem to equivalence of CQs under the relaxation of bag-set semantics discussed in Section 3.3.2. In related work, Grumbach and Tininini consider the view rewriting problem for the restricted case where the query and view languages are aggregation queries over a single relation (i.e., no joins) [53]; this case has very practical application to OLAP-style queries data warehouses with star schemas. Their main contribution is a polynomial-time rewriting algorithm which they prove to be complete for finding single-view rewritings of queries containing the **SUM**, **COUNT**, and **PRODUCT** aggregation functions.

The various works by both Cohen et al. and Grumbach et al. surveyed above only considered *non-nested* aggregation queries and views. More specifically, their various query and view languages all correspond to relational algebra expressions of the form

$$\Pi_x^{z=f(y)}(E) \tag{3.9}$$

where  $E$  is a relational algebra query not containing aggregation functions. Observe that when generalized projection only occurs as the top-most query operator, *every grouping variable is also an output variable*. This syntactic restriction is significant, and highlights the orthogonality between our work and the work surveyed in this section. Whereas previous work focuses on understanding and characterizing specific aggregation operations, our work focuses on the interaction between aggregation and query nesting. It is precisely the ability to project away attributes that have been used for grouping at lower levels that makes the equivalence problem for nested queries a non-trivial extension of work on non-nested aggregation queries. Our abstraction of aggregation functions as collection constructors is comparatively primitive, but arguably necessary both as a simplifying assumption and in order to preserve decidability for the query equivalence problem. Other authors have shown that too much domain-specific knowledge about the functioning of particular aggregation functions quickly causes undecidable implication of aggregation constraints [75, 97]. Within our context, the obvious corollary is that the query equivalence problem becomes undecidable when queries contain aggregation functions whose results can be both arbitrarily nested and compared with explicit (and possibly, only implicit) higher-order comparisons.



# Chapter 4

## Equivalence of Conjunctive Queries with Outer Joins

Our goal in this thesis is to characterize equivalence of conjunctive object-constructing queries. In Chapter 2 we reduced this problem to deciding encoding-equivalence between *hierarchical* conjunctive encoding queries. In Chapter 3 we showed that encoding-equivalence between CEQs generalizes query equivalence between CQs. Therefore, before we address encoding-equivalence between *hierarchical CEQs (HCEQs)*, it is helpful for us to first address query equivalence between *hierarchical CQs (HCQs)*.

In Section 2.3.1 we defined the semantics of HCQs in terms of *hierarchical applications of terminal embeddings* of the query body into the database instance. There is a strong relationship between HCQs and relational algebra expressions containing the *left outer-join* (LOJ) operator. In Section 4.1 we show that extending the conjunctive algebra with LOJ yields a language whose expressive power corresponds to a generalization of HCQs that we call *directed acyclic CQs (DACQs)*. In Sections 4.3–4.5 we address the DACQ equivalence problem (of which the HCQ equivalence problem is a special case). We end the chapter by surveying relevant literature in Section 4.6.

### 4.1 Conjunctive Queries with Outer Joins

In this section we propose two languages for conjunctive queries with outer joins. In Section 4.1.1 we formally define the SPCL algebra, which extends the conjunctive algebra with a left outer-join operator. Next, in Section 4.1.2 we define *directed acyclic CQs*, a straightforward generalization of hierarchical CQs, and we prove that DACQs have the same expressive power as SPCL expressions.

### 4.1.1 The SPCL Algebra

The conjunctive relational algebra—i.e., the operators *base relation*, *selection*, *duplicate-eliminating projection*, and *Cartesian product* along with the syntactic extensions *conjunctive selection* and *join* from Figure 2.5—is sometimes called the *SPC algebra*. We define the *SPCL algebra* to be language formed by extending the SPC algebra with the *left outer-join* operator [42].

The left outer join (LOJ) of relation  $R$  with  $S$ —denoted  $R \vec{\bowtie}_p S$ —is equivalent to the *outer union* of the inner- and anti-joins of  $R$  and  $S$ ,

$$R \vec{\bowtie}_p S \equiv (R \bowtie_p S) \uplus (R \triangleright_p S) \equiv (R \bowtie_p S) \uplus (R - (R \ltimes_p S)) \quad (4.1)$$

where the outer union operator  $\uplus$  pads the tuples originating from the anti-join with a special null constant<sup>1</sup>  $\blacksquare$  so that they have the same arity as the tuples originating from the inner join. Alternatively, the LOJ operator can be defined using the notion of *tuple subsumption*. Given two tuples  $t_1$  and  $t_2$  of the same arity,  $t_1$  subsumes  $t_2$  if  $t_2$  contains strictly more  $\blacksquare$  values than  $t_1$ , and  $t_1$  and  $t_2$  coincide on all non- $\blacksquare$  attributes of  $t_1$ . The *removal of subsumed tuples* from  $R$ —denoted  $R \downarrow$ —returns only the tuples in  $R$  that are not subsumed by another tuple. The LOJ expression  $R \vec{\bowtie}_p S$  can then be characterized in terms of inner join, outer union, and removal of subsumed tuples, which can be combined into the *minimal union* operator  $\oplus$  [37].

$$R \vec{\bowtie}_p S \equiv ((R \bowtie_p S) \uplus R) \downarrow \equiv (R \bowtie_p S) \oplus R \quad (4.2)$$

An advantage of equation 4.2 is that all of the non-monotonic behaviour is encapsulated within the  $\downarrow$  or  $\oplus$  operators, which can then be commuted above the inner join operators (although not above projection). This allows a tree of inner and outer joins to be abstracted as a tree of conjunctive and disjunctive operations, which can then be transformed into a canonical form; this is the conceptual basis for Galindo-Legaria’s *join-disjunctive normal form (JDNF)* [37], which we discuss further in Section 4.6.1.

In order to combine LOJ with the other conjunctive algebra operators, we must extend the definitions of the selection and duplicate-eliminating projection operators to handle input relations that contain nulls (introduced by lower LOJ operators). Following the three-valued logic of the SQL standard [63], we define the selection predicate  $A = B$  to evaluate to **false** if either attribute  $A$  or  $B$  is bound to the null constant; in other words, equality predicates are *null-rejecting*. Also following the SQL standard, we extend the definition of duplicate-eliminating projection to treat the null constant the same as any other database constant; in other words, identical tuples are always considered duplicate copies of the same tuple. Due to the interaction between LOJ and projection operators that introduce constants, we also need to slightly modify the syntax of projection attribute lists so

---

<sup>1</sup>We intentionally use the symbol  $\blacksquare$  rather than the SQL constant NULL to highlight our assumption that this special null constant does not occur within the input to an SPCL query, and therefore all occurrences of  $\blacksquare$  within the query output are introduced by LOJ operators.



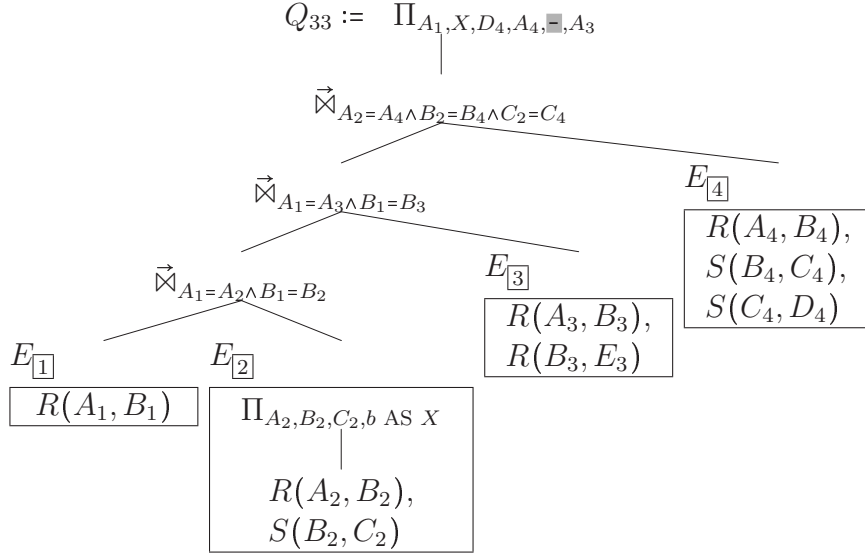


Figure 4.1: SPCL query  $Q_{33}$  equivalent to HCQ  $Q_{17}$

that constants are assigned fresh attribute names; for this we adopt the common “ $c$  AS  $A_{ij}$ ” syntax.

The SPCL algebra is capable of expressing any hierarchical CQ. The following example illustrates the basic technique behind the formal proof.

**Example 32** Consider the SPCL query  $Q_{33}$  shown in Figure 4.1.<sup>2</sup> Query  $Q_{33}$  is equivalent to the hierarchical CQ  $Q_{17}$  from Figure 2.11. For example,  $E_{\boxed{1}}$  corresponds to the tree  $T_6$  in Figure 2.12; similarly,  $E_{\boxed{2}}$  corresponds to  $T_4$ ,  $E_{\boxed{3}}$  to  $T_5$ , and  $E_{\boxed{4}}$  to  $T_2$ . The reader can verify that when  $Q_{33}$  is evaluated over database instance  $\mathbb{D}_2$  from Example 20, the result—relation  $R_{14}$  in Figure 4.2—is obtained from the relations shown in Figure 2.13 by projecting the appropriate attributes from the result of expression  $((\Gamma^6 \bowtie \Gamma^4) \bowtie \Gamma^5) \bowtie \Gamma^2$ .

**Theorem 4.1.1** *Every HCQ can be expressed as an equivalent SPCL expression.*

PROOF. Given any HCQ  $Q$ , create an equivalent SPCL expression as follows.

1. For each node  $v$  in  $\text{body}_Q$ , create an SPC expression  $E_{\boxed{v}}$  from the CQ bodies of  $v$  and all of  $v$ 's ancestors, with “fresh” copies of each attribute name. If the head of  $Q$  contains a constant that is labelled with node  $v$ , then add a projection operator to the top of  $E_{\boxed{v}}$  to introduce an appropriate constant attribute.
2. Arrange these conjunctive queries into a left-deep tree of LOJ operators whose left-to-right ordering corresponds to a breadth-first traversal of the

<sup>2</sup>For convenience we have used a hybrid notation rather than pure SPCL; the subexpressions  $E_{\boxed{1}}$ ,  $E_{\boxed{2}}$ ,  $E_{\boxed{3}}$ , and  $E_{\boxed{4}}$  represent CQ bodies that are easily translated into pure SPC expressions.

$A_1$	$X$	$D_4$	$A_4$	$\text{--}$	$A_3$
1	$b$	5	1	$\text{--}$	1
2	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$	2
2	$b$	6	2	$\text{--}$	$\text{--}$
3	$b$	6	3	$\text{--}$	$\text{--}$
3	$b$	$\text{--}$	$\text{--}$	$\text{--}$	3
5	$b$	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$
6	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$

$A_1$	$X$	$D_4$	$A_4$	$\text{--}$	$A_3$
1	$b$	$\text{--}$	$\text{--}$	$\text{--}$	1
1	$b$	5	1	$\text{--}$	1
2	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$	2
2	$b$	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$
3	$b$	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$
3	$b$	$\text{--}$	$\text{--}$	$\text{--}$	3
5	$b$	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$
6	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$	$\text{--}$

Figure 4.2: SPCL query results  $R_{14} = (Q_{33})^{\mathbb{D}_2}$  and  $R_{15} = (Q_{34})^{\mathbb{D}_2}$

nodes of  $\text{body}_Q$ . For each non-root node  $v$  whose parent node in  $\text{body}_Q$  is  $u$ , the LOJ operator immediately above  $E_{\square v}$  has a join predicate that equates all of the attributes from  $E_{\square u}$  with their corresponding attributes in  $E_{\square v}$ .

3. Create a top-most projection operator, and for each component  $W_i$  labelled with  $v$  in the head  $Q(\overline{W})$ , output the associated attribute from  $E_{\square v}$ .

Consider any tree  $T_i \in \text{trees}(\text{body}_Q)$ , and let  $u_1, \dots, u_k$  denote the nodes of  $T_i$  in breadth-first order. Given any database instance  $\mathbb{D}$ , there is a one-to-one correspondence between embeddings of  $T_i$  into  $\mathbb{D}$  and tuples in the *inner* join  $(E_{\square u_1} \bowtie \dots \bowtie E_{\square u_k})^{\mathbb{D}}$ . An embedding  $\gamma : T_i \rightarrow \mathbb{D}$  is terminal precisely when the corresponding tuple in  $(E_{\square u_1} \bowtie \dots \bowtie E_{\square u_k})^{\mathbb{D}}$  is not subsumed by any other tuples in the result of the join tree (prior to the application of the top-most projection operator).  $\square$

### 4.1.2 Directed Acyclic Conjunctive Queries

The following example demonstrates that SPCL is strictly more expressive than the language of hierarchical CQs.

**Example 33** Let  $Q_{34}$  be the SPCL query formed from  $Q_{33}$  in Figure 4.1 by adding the predicate  $C_2 = E_3$  to the join predicate of the LOJ operator immediately above  $E_{\square 4}$ . Evaluating  $Q_{34}$  over database instance  $\mathbb{D}_2$  from Example 20 yields the relation  $R_{15}$  shown in Figure 4.2. There is no HCQ equivalent to  $Q_{34}$ .

The reason that SPCL query  $Q_{34}$  cannot be expressed as an HCQ is that the variables  $C_2$  and  $E_3$  are inherited from separate blocks that are not themselves in an ancestor-descendant relationship. Syntactic well-formed-ness of HCQs requires that node  $u$  must be an ancestor of  $v$  if  $\mathcal{P}_v \cap \mathcal{L}_u \neq \emptyset$ ; however, satisfying this requirement for this query requires a body that is not a tree. We now generalize the language of HCQs to encompass this case.

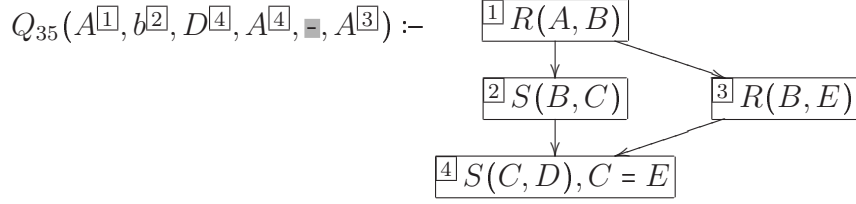


Figure 4.3: Directed acyclic conjunctive query  $Q_{35}$

$\Gamma$	$A$	$B$	$C$	$D$	$E$	$b$	$A^{[1]}$	$b^{[2]}$	$D^{[4]}$	$A^{[4]}$	$\blacksquare$	$A^{[3]}$
	1	2	4		3	$b$	1	$b$	$\blacksquare$	$\blacksquare$	$\blacksquare$	1
	1	2	4	5	4	$b$	1	$b$	5	1	$\blacksquare$	1
	2	3			4	$b$	2	$\blacksquare$	$\blacksquare$	$\blacksquare$	$\blacksquare$	2
	2	3			5	$b$	2	$\blacksquare$	$\blacksquare$	$\blacksquare$	$\blacksquare$	2
	2	4	5			$b$	2	$b$	$\blacksquare$	$\blacksquare$	$\blacksquare$	$\blacksquare$
	3	4	5			$b$	3	$b$	$\blacksquare$	$\blacksquare$	$\blacksquare$	$\blacksquare$
	3	5	6		7	$b$	3	$b$	$\blacksquare$	$\blacksquare$	$\blacksquare$	3
	5	7	8			$b$	5	$b$	$\blacksquare$	$\blacksquare$	$\blacksquare$	$\blacksquare$
	6	8				$b$	6	$\blacksquare$	$\blacksquare$	$\blacksquare$	$\blacksquare$	$\blacksquare$

Figure 4.4: Terminal embeddings of  $Q_{35}$  into  $\mathbb{D}_2$

**Definition 4.1.2 (Directed Acyclic Conjunctive Query (DACQ))** A directed acyclic conjunctive query is a generalization of an HCQ whose body is a rooted, connected DAG.

Let  $\text{dags}(G)$  denote the set of all non-empty DAGs that can be formed from some non-empty DAG  $G$  by zero or more iterations of selecting and deleting a leaf node. The evaluation of DACQs is defined analogously to the evaluation of HCQs (see Section 2.3.1), in terms of embeddings  $\gamma : G \rightarrow \mathbb{D}$  and terminal embeddings  $\gamma : G \rightsquigarrow \mathbb{D}$  for  $G \in \text{dags}(\text{body}_Q)$ . The CQ  $\text{body}_G^+$  is defined analogously to  $\text{body}_T^+$ .

**Example 34** Figure 4.3 depicts DACQ  $Q_{35}$ , and Figure 4.4 shows the terminal embeddings of  $Q_{35}$  into database instance  $\mathbb{D}_2$ .

We now show that every DACQ can be expressed as an SPCL expression and vice-versa. The basic technique behind the formal proof is illustrated by the following example.

**Example 35** Consider the SPCL query  $Q_{34}$  from Example 33. The DACQ  $Q_{36}$  shown in Figure 4.5 is equivalent to  $Q_{34}$ . Observe the close correspondence between the vertices of  $\text{body}_{Q_{36}}$  and the SPC blocks within  $Q_{34}$ .

**Theorem 4.1.3** *DACQs and SPCL expressions have the same expressive power.*

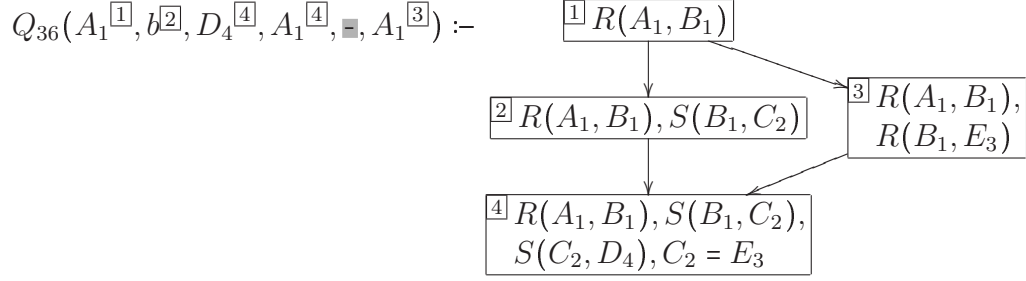


Figure 4.5: DACQ  $Q_{36}$  equivalent to SPCL query  $Q_{34}$

PROOF. The construction in the proof of Theorem 4.1.1 applies identically to DACQs, which proves that SPCL expressions are at least as expressive as DACQs.

To prove the other direction, it suffices to restrict our attention to a particular subclass of SPCL expressions. It follows from the work of previous authors on algebraic transformations for outer-join queries [12, 41] that given an arbitrary SPCL query, there exists an equivalent SPCL query that is rooted by a projection operator whose input is a left-deep tree of LOJ operators with SPC expressions occurring at its leaves. Let  $Q$  be any SPCL query with this restricted form, and let  $E_{[u_1]}, \dots, E_{[u_k]}$  denote the SPC subqueries at the leaves of the left-deep tree, as they appear in left-to-right order. An equivalent DACQ can be constructed as follows.

1. Create a root node labelled  $[u_1]$  containing all of the base relations and join/selection predicates occurring within SPC subquery  $E_{[u_1]}$ .
2. Enumerate SPC queries  $E_{[u_2]}, \dots, E_{[u_k]}$ . For each query  $E_{[u_i]}$ , let  $p_i$  denote the join predicate of the LOJ operator that has  $E_{[u_i]}$  as its right input.
  - (a) Create a node labelled  $[u_i]$  containing all of the base relations and join/selection predicates occurring within SPC subquery  $E_{[u_i]}$ , as well as the predicates in  $p_i$ .
  - (b) For each  $j \neq i$  such that  $p_i$  references an attribute originating from  $E_{[u_j]}$ , add an edge from node  $u_j$  to node  $u_i$ . Note that well-formed-ness of SPCL expressions guarantees  $j < i$ , and so node  $u_j$  already exists.
3. Create the head of the DACQ from the projection list of the top projection operator of  $Q$ . For each attribute, label it with the node label corresponding to the SPC expression from which it originates.
4. Use variable/constant substitution to eliminate any equality predicate that refers to a variable that is local to the node containing the predicate. Stop when the equality predicates in each node reference only constants and parameter variables.

□

**Corollary 4.1.4** *The query equivalence problems for DACQs and SPCL expressions are polynomial-time reducible to each other.*

## 4.2 A Canonical Form for DACQs

In order to standardize the input to the query equivalence problem which we consider in the next section, in this section we define a canonical form for DACQs. The canonical form is defined in terms of several declarative properties. In Section 4.2.1 we introduce a new type of dependency that we call a *nullability dependency*. In Section 4.2.2 we use nullability dependencies to characterize whether a query adheres to the properties of the canonical form, and we give an algorithm for converting an arbitrary DACQ into an equivalent canonical form.

**Definition 4.2.1 (DACQ Canonical Form)** *Given any satisfiable DACQ  $Q$ , we say that  $Q$  is in canonical form if it satisfies all of the following properties.*

- [P1] *For each vertex  $u \in \text{body}_Q$ , CQ body  $\text{body}_u^+$  is satisfiable.*
- [P2] *For each non-root leaf vertex  $u \in \text{body}_Q$ , the head tuple  $\overline{W}$  of  $Q$  contains a component  $W_i$  that is labelled with vertex  $u$ .*
- [P3] *For each vertex  $u \in \text{body}_Q$  and each parameter variable  $X \in \mathcal{P}_u$ :*
  - (a) *if  $\text{pred}_u^+$  implies that  $X$  is equal to a constant, then  $X$  does not occur within  $\text{atoms}_u$ ; and*
  - (b) *if  $\text{pred}_u^+$  implies that  $X$  is equal to some other parameter  $Y \in \mathcal{P}_u$ , then  $\text{atoms}_u$  does not contain both  $X$  and  $Y$  (it may contain one).**(This presumes that predicates in  $\text{pred}_u$  do not reference variables from  $\mathcal{L}_u$ , which is the (non-restrictive) syntactic convention we previously adopted for HCQs—see Section 2.3.1.)*
- [P4] *The body of  $Q$  does not contain any redundant atoms, predicates, or edges. An atom within  $\text{atoms}_u$  is redundant if  $\text{pred}_u^{\text{anc}}$  implies that it is identical to an atom within  $\text{atoms}_u^{\text{anc}}$ . A predicate within  $\text{pred}_u$  is redundant if it is implied by  $\text{pred}_u^{\text{anc}}$ . An edge  $(u, v)$  is redundant if there remains a path from  $u$  to  $v$  after deletion of edge  $(u, v)$ .*
- [P5] *For each pair of distinct vertices  $u, v \in \text{body}_Q$ , if  $v$  is not an ancestor of  $u$ , then there exists some terminal embedding of  $Q$  into some database  $\mathbb{D}$  that embeds  $\text{body}_u$  into  $\mathbb{D}$  but not  $\text{body}_v$  into  $\mathbb{D}$ .*

The purpose of properties [P1] and [P2] is to avoid having vertices in  $\text{body}_Q$  that do not affect the query output. The purpose of properties [P3] and [P4] is to standardize the syntax of the CQ bodies within the vertices by removing redundancy and minimizing the number of parameter variables. The purpose of [P5] is to standardize the shape of the DAG formed by the vertices of  $\text{body}_Q$ —if property [P5] is violated, then an edge can be added from vertex  $v$  to vertex  $u$  without affecting the output of  $Q$  over any database instance.

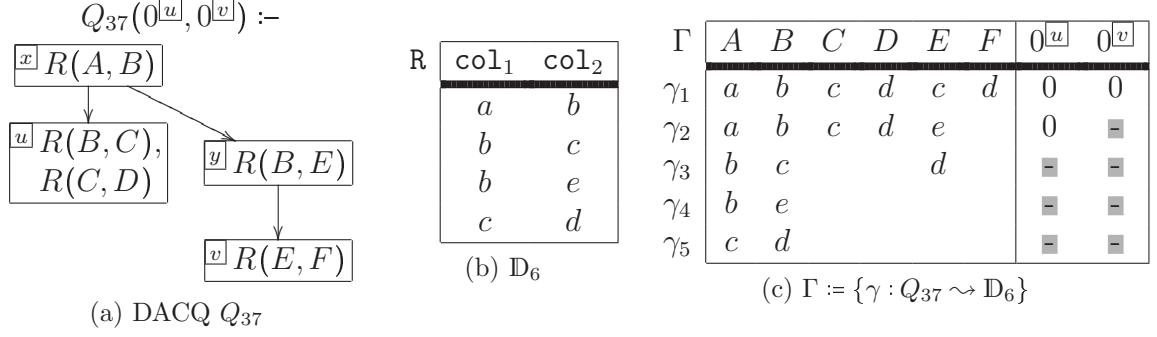


Figure 4.6: Query  $Q_{37}$  and its terminal embeddings into  $\mathbb{D}_6$

Given a query  $Q$ , properties [P1]–[P4] can all be tested in polynomial time by examining the syntax of  $Q$  and using transitivity to reason about implication of equality predicates. However, it is not obvious how to test property [P5]. On first glance, it may appear that this is simply a CQ containment problem—specifically, that  $\text{body}_u^+$  is not contained under  $\text{body}_v^+$ . However, the following example demonstrates that the problem is more nuanced than that.

**Example 36** Consider the DACQ  $Q_{37}$  shown in Figure 4.6a. The variable mapping  $h := \{A/A, B/B, E/C, F/D\}$  is a homomorphism  $h : \text{body}_v^+ \rightarrow \text{body}_u^+$  proving that any embedding of vertex  $x$  that extends to vertex  $u$  will also extend to vertex  $v$ . Nevertheless, database instance  $\mathbb{D}_6$  in Figure 4.6b proves that we cannot safely add an edge from  $v$  to  $u$ . Doing so would eliminate the terminal embedding  $\gamma_2$  shown in Figure 4.6c, thereby removing the output tuple  $\langle 0, \blacksquare \rangle$  from  $(Q_{37})^{\mathbb{D}_6}$ .

### 4.2.1 Functionally-Determined Nullability

In this section we define *nullability dependencies*, which identify the variables in a DACQ body that determine whether an embedding extends to some node  $v$ .

**Definition 4.2.2 (Nullability Dependency (ND))** *Given a DACQ  $Q$ , a non-root node  $v \in \text{body}_Q$ , and a set of variables  $\mathcal{S} \subseteq \mathcal{B}_v^{\text{anc}}$ , we say that  $\mathcal{S}$  null-determines  $v$  within the embeddings of  $Q$ —denoted  $Q \models \mathcal{S} \Rightarrow v$ —if for every database instance  $\mathbb{D}$  and every pair of embeddings  $\gamma, \phi : \text{body}_v^{\text{anc}} \rightarrow \mathbb{D}$  that coincide over  $\mathcal{S}$ , embedding  $\gamma$  extends to an embedding of  $v$  iff embedding  $\phi$  extends to an embedding of  $v$ . (Because all embeddings can be assumed to be extended with identity over constants in  $\mathcal{C}$ , it is a trivial extension to allow  $\mathcal{S}$  to contain constants.)*

**Example 37** Continuing Example 36, terminal embeddings  $\gamma_1$  and  $\gamma_2$  in Figure 4.6c prove that  $Q_{37} \not\models \{A, B\} \Rightarrow v$ . However,  $Q_{37} \models \{E\} \Rightarrow v$  holds (by the next lemma).

**Lemma 4.2.3** *For any DACQ  $Q$  and non-root node  $v \in \text{body}_Q$  containing parameter variables  $\mathcal{P}_v$ ,  $Q \models \mathcal{P}_v \Rightarrow v$ .*

PROOF. For any embedding  $\gamma : \text{body}_v^{\text{anc}} \rightarrow \mathbb{D}$ , whether or not  $\gamma$  extends to  $v$  is determined only by the tuple of parameter values  $\gamma(\overline{\mathcal{P}}_v)$ .  $\square$

We now present a homomorphic condition that characterizes the existence of a nullability dependency. For simplicity, we start by considering the very restricted case where the query body contains only two nodes, does not contain any equality predicates, and obeys properties **[P1]**–**[P4]** of Definition 4.2.1.

**Theorem 4.2.4** *Given*

- DACQ  $Q$  obeying properties **[P1]**–**[P4]** and composed of a single edge from vertex  $u$  to vertex  $v$  in which  $\text{pred}_u = \text{pred}_v = \emptyset$ , and
- a set  $\mathcal{S} \subseteq \mathcal{B}_u$ ;

$Q \models \mathcal{S} \Rightarrow v$  iff there exists a mapping  $h : \mathcal{B}_v^+ \rightarrow \mathcal{B}_v^+ \cup \mathcal{C}$ , extended with identity over constants  $\mathcal{C}$ , such that

1.  $h$  is the identity mapping over  $\text{atoms}_u$ , and
2. for any atom  $R(\overline{X}) \in \text{atoms}_v$ , if  $h(R(\overline{X})) \in \text{atoms}_v$  then tuple  $h(\overline{X})$  does not contain any variables from  $\mathcal{B}_u \setminus \mathcal{S}$ .

PROOF. Suppose that  $h$  exists. Choose any database instance  $\mathbb{D}$ , and any two embeddings  $\gamma, \phi : \text{body}_u \rightarrow \mathbb{D}$  that coincide over  $\mathcal{S}$  and for which there exists an embedding  $\phi' : \text{body}_v^+ \rightarrow \mathbb{D}$  that extends  $\phi$ . Define the variable mapping  $\gamma' : \mathcal{B}_v^+ \rightarrow \mathbb{D}$  as follows.

$$\gamma'(X) := \begin{cases} \gamma \circ h(X) & \text{if } h(X) \in \mathcal{B}_u \\ \phi' \circ h(X) & \text{otherwise} \end{cases}$$

Because  $h$  is the identity mapping over  $\text{atoms}_u$ ,  $\gamma'$  is consistent with  $\gamma$  over  $\text{body}_u$ . For each atom  $R(\overline{X}) \in \text{atoms}_v$ , if  $h(R(\overline{X})) \in \text{atoms}_u$  then  $\gamma'(R(\overline{X})) = \gamma \circ h(R(\overline{X})) \in \mathbb{D}$ . If  $h(R(\overline{X})) \in \text{atoms}_v$ , then because  $h(\overline{X})$  does not contain any variables from  $\mathcal{B}_u \setminus \mathcal{S}$  and because  $\gamma, \phi$  are consistent over  $\mathcal{S}$ ,  $\gamma'(R(\overline{X})) = \phi'(R(\overline{X})) \in \mathbb{D}$ . Therefore  $\gamma'$  is an embedding of  $\text{body}_v^+$  into  $\mathbb{D}$  that extends  $\gamma$ , which implies  $Q \models \mathcal{S} \Rightarrow v$ .

Now suppose  $Q \models \mathcal{S} \Rightarrow v$ . Let  $f_1$  be a function that isomorphically freezes  $\mathcal{B}_u$  to constants in  $\mathcal{C} \setminus \mathcal{C}_Q$ , extended with identity over  $\mathcal{C}_Q$ . Let  $f_2$  be a function that isomorphically freezes  $\mathcal{B}_v^+$  to constants in  $\mathcal{C} \setminus \mathcal{C}_Q$ , extended with identity over  $\mathcal{C}_Q$ . Choose  $f_1, f_2$  so that they are consistent over  $\mathcal{S}$  and satisfy  $f_1(\mathcal{B}_u \setminus \mathcal{S}) \cap f_2(\mathcal{B}_v^+ \setminus \mathcal{S}) = \emptyset$ . Define database instances  $\mathbb{D}_1 := f_1(\text{atoms}_u)$ ,  $\mathbb{D}_2 := f_2(\text{atoms}_v^+)$ , and  $\mathbb{D} := \mathbb{D}_1 \cup \mathbb{D}_2$ . Because  $f_1^{-1}$  and  $f_2^{-1}$  are consistent over their common domain, we can define a single inverse mapping  $f^{-1} = f_1^{-1} \cup f_2^{-1}$  that satisfies  $f^{-1}(\mathbb{D}) = \text{atoms}_v^+$ . Choose embedding  $\gamma : \text{body}_u \rightarrow \mathbb{D}$  as  $\gamma := f_1$ . Choose embedding  $\phi : \text{body}_v^+ \rightarrow \mathbb{D}$  as  $\phi := f_2$ . Because  $Q \models \mathcal{S} \Rightarrow v$  and  $\gamma, \phi$  are consistent over  $\mathcal{S}$ , there must exist some embedding  $\gamma' : \text{body}_v^+ \rightarrow \mathbb{D}$  that extends  $\gamma$ .

Define  $h_1 := f^{-1} \circ \gamma'$ , and for each  $i > 1$  define  $h_{i+1} := h_1 \circ h_i$ . Observe that each  $h_i$  is a homomorphism from  $\text{body}_v^+$  to  $\text{body}_v^+$  that is the identity over  $\mathcal{B}_u$  (and therefore over  $\text{body}_u$ ). For  $h_1$  and each variable  $X \in \mathcal{L}_v$ , either

1.  $\gamma'(X) \in \mathcal{C}_Q$  and so  $h_1(X) \in \mathcal{C}_Q$ ;
2.  $\gamma'(X) \in f_1(\mathcal{S}) = f_2(\mathcal{S})$  and so  $h_1(X) \in \mathcal{S}$ ;
3.  $\gamma'(X) \in f_1(\mathcal{B}_u \setminus \mathcal{S})$  and so  $h_1(X) \in \mathcal{B}_u \setminus \mathcal{S}$ ; or
4.  $\gamma'(X) \in f_2(\mathcal{B}_v^+ \setminus \mathcal{S})$  and so  $h_1(X) \in \mathcal{B}_v^+ \setminus \mathcal{S}$ .

Consider the infinite sequence  $h_1(X), h_2(X), \dots$ ; after  $k = |\mathcal{L}_v|$  steps, either

1.  $\exists i \in [1, k]$  such that  $h_i(X) \in \mathcal{C}_Q \cup \mathcal{B}_u$ , which implies that  $\forall j \geq i. [h_j(X) = h_i(X)]$ , and so  $\gamma' \circ h_k(X) \in f_1(\mathcal{C}_Q \cup \mathcal{B}_u)$ ; or
2.  $\forall i \in [1, k], h_i(X) \in \mathcal{L}_v$ , which implies that the mapping composition is periodic with period at most  $k$ , and so  $\forall j. [h_j(X) \in \mathcal{L}_v]$ . In this case,  $\gamma' \circ h_k(X) \in f_2(\mathcal{L}_v)$ .

Consider any atom  $R(\overline{X}) \in \mathbf{atoms}_v$  such that  $h_{k+1}(R(\overline{X})) \in \mathbf{atoms}_v$ . By property **[P4]**,  $h_{k+1}(R(\overline{X})) \notin \mathbf{atoms}_u$  and therefore  $\gamma' \circ h_k(R(\overline{X})) \notin \mathbb{D}_1$ . Suppose that there exists some variable  $Y \in \overline{X}$  such that  $h_{k+1}(Y) \in \mathcal{B}_u \setminus \mathcal{S}$ ; then,  $h_{k+1}(Y)$  must belong to the first case above, implying that  $h_k(Y) = h_{k+1}(Y)$  and  $\gamma' \circ h_k(Y) = f_1(h_k(Y)) \notin \mathbf{adom}(\mathbb{D}_2)$ , which is a contradiction. Therefore, variable  $Y$  cannot exist, and so tuple  $h_{k+1}(\overline{X})$  does not contain variables from  $\mathcal{B}_u \setminus \mathcal{S}$ .  $\square$

**Corollary 4.2.5** *There exists a unique minimal set  $\mathcal{S}_v \subseteq \mathcal{B}_u$  satisfying  $Q \models \mathcal{S}_v \Rightarrow v$ , and  $\mathcal{S}_v$  can be computed in NP time.*

PROOF. It suffices to show that for any pair of sets  $\mathcal{S}_1, \mathcal{S}_2 \subseteq \mathcal{B}_u$ , if  $Q \models \mathcal{S}_1 \Rightarrow v$  and  $Q \models \mathcal{S}_2 \Rightarrow v$  then  $Q \models (\mathcal{S}_1 \cap \mathcal{S}_2) \Rightarrow v$ . If this is true, then

$$\mathcal{S}_v = \mathcal{P}_v \setminus \{X \in \mathcal{P}_v \mid Q \models (\mathcal{P}_v \setminus \{X\}) \Rightarrow v\}$$

and so computing  $\mathcal{S}_v$  requires testing  $|\mathcal{P}_v|$  nullability-dependencies (each of which can be tested in NP time á la Theorem 4.2.4).

Suppose that  $Q \models \mathcal{S}_1 \Rightarrow v$ , and let  $h_1$  be the corresponding homomorphism that satisfies Theorem 4.2.4. Suppose that  $Q \models \mathcal{S}_2 \Rightarrow v$ , with  $h_2$  the corresponding homomorphism. Consider the composite function  $h_2 \circ h_1$ . For any atom  $R(\overline{X}) \in \mathbf{atoms}_v$ , if  $h_1(R(\overline{X})) \in \mathbf{atoms}_u$  then  $h_2 \circ h_1(R(\overline{X})) \in \mathbf{atoms}_u$ , while if  $h_1(R(\overline{X})) \in \mathbf{atoms}_v$  and  $h_1(\overline{X})$  contains a variable from  $\mathcal{B}_u \setminus \mathcal{S}_2$  then  $h_2 \circ h_1(R(\overline{X})) \in \mathbf{atoms}_u$  as well (because  $h_2$  is the identity over  $\mathcal{B}_u$  and any atom that  $h_2$  maps into  $\mathbf{atoms}_v$  can not contain  $\mathcal{B}_u \setminus \mathcal{S}_2$ ). The symmetric argument can then be applied to the role of  $h_2$  within the composite function  $h_1 \circ h_2 \circ h_1$ , and then to the second  $h_1$  within  $h_2 \circ h_1 \circ h_2 \circ h_1$ , etc. Let  $h$  be a mapping formed by composing at least  $|\mathbf{atoms}_v|$  alternations of  $h_1$  and  $h_2$ ;  $h$  is the identity over  $\mathbf{atoms}_u$ , and for any atom  $R(\overline{X}) \in \mathbf{atoms}_v$  that  $h$  maps into  $\mathbf{atoms}_v$ , tuple  $h(\overline{X})$  does not contain any variable from  $\mathcal{B}_u \setminus (\mathcal{S}_1 \cap \mathcal{S}_2)$ . Therefore,  $h$  proves  $Q \models (\mathcal{S}_1 \cap \mathcal{S}_2) \Rightarrow v$ .  $\square$

**Example 38** Consider the following two DACQs.

$$Q_{38}(A^{\overline{b}}) :- \boxed{a} R(A, B), R(X, Y) \longrightarrow \boxed{b} R(A, A), R(X, Z)$$

$$Q_{39}(A^{\overline{d}}) :- \boxed{c} R(A, B), R(X, Y) \longrightarrow \boxed{d} R(A, Z), R(X, Z)$$



The homomorphism  $h_b = \{A/A, B/B, X/X, Y/Y, Z/Y\}$  proves  $Q_{38} \models \{A\} \Rightarrow b$ , and  $\mathcal{S}_b = \{A\}$  is minimal because no homomorphism from  $\text{body}_b^+$  to itself can map atom  $R(A, A)$  into  $\text{body}_a$ . For  $Q_{39}$ , the only homomorphism from  $\text{body}_d^+$  to itself that is the identity over  $\text{body}_c$  is  $h_d = \{A/A, B/B, X/X, Y/Y, Z/Z\}$ ; therefore,  $\mathcal{S}_d = \{A, X\} = \mathcal{P}_d$  is the minimal set satisfying  $Q_{39} \models \mathcal{S}_d \Rightarrow d$ .

We now extend Theorem 4.2.4 to the case where the query body contains equality predicates.

**Theorem 4.2.6** *Given*

- *DACQ*  $Q$  obeying properties [P1]–[P4] and composed of a single edge from vertex  $u$  to vertex  $v$ , and
- a set  $\mathcal{S} \subseteq \mathcal{B}_u$ ;

*assume a total ordering over  $\mathcal{B} \cup \mathcal{C}$  that places all constants before any variables. Let  $\kappa : \mathcal{B} \rightarrow \mathcal{B} \cup \mathcal{C}$  be the mapping of each variable to the minimal element  $Y \in \mathcal{B} \cup \mathcal{C}$  such that  $\text{pred}_v^+ \models X = Y$ . Then,  $Q \models \mathcal{S} \Rightarrow v$  iff*

1. *for each variable  $X$  in  $\text{pred}_v$  there exists a variable  $Y \in \mathcal{S}$  such that  $\text{pred}_u \models X = Y$ ; and*
2.  *$\kappa(Q) \models \kappa(\mathcal{S}) \Rightarrow v$ .*

PROOF. Suppose that  $Q \models \mathcal{S} \Rightarrow v$ .

1. Suppose that  $\text{pred}_v$  contains some predicate  $X = Z$  such that there does not exist a variable  $Y \in \mathcal{S}$  satisfying  $\text{pred}_u \models X = Y$ . Let  $\gamma$  be any embedding of  $\text{body}_v^+$  into any database  $\mathbb{D}$ . Let  $c$  denote some fresh constant not in  $\text{adom}(D)$ . Initialize variable mapping  $\phi$  as identical to  $\gamma$ ; then, for each variable  $Y$  such that  $\text{pred}_u \models X = Y$ , set  $\phi(Y) := c$ . Define database  $\mathbb{D}' := \mathbb{D} \cup \phi(\text{atoms}_u)$ . Both  $\gamma$  and  $\phi$  embed  $\text{body}_u$  into  $\mathbb{D}'$  and they coincide on  $\mathcal{S}$ , yet  $\gamma$  extends to  $v$  while  $\phi$  does not because  $\phi(X) = c \neq \phi(Z)$  (property [P4] requires that  $\text{pred}_u \models X = Z$ ). Hence,  $Q \not\models \mathcal{S} \Rightarrow v$ , which is a contradiction, and so  $\text{pred}_v$  cannot contain predicate  $X = Z$ .
2. By construction, mapping  $\kappa$  enacts the equality predicates within  $\text{pred}_v^+$ ; therefore, each embedding of  $\kappa(\text{body}_u)$  is also an embedding of  $\text{body}_u$ , and each embedding of  $\kappa(\text{body}_v^+)$  is also an embedding of  $\text{body}_v^+$ . Therefore, any two embeddings of  $\kappa(\text{body}_u)$  proving  $\kappa(Q) \not\models \kappa(\mathcal{S}) \Rightarrow v$  would also prove  $Q \not\models \mathcal{S} \Rightarrow v$ , which is a contradiction. Therefore,  $\kappa(Q) \models \kappa(\mathcal{S}) \Rightarrow v$ .

Now suppose that for each variable  $X$  in  $\text{pred}_v$  there exists a variable  $Y \in \mathcal{S}$  satisfying  $\text{pred}_u \models X = Y$ , and suppose that  $\kappa(Q) \models \kappa(\mathcal{S}) \Rightarrow v$ . Given any database  $\mathbb{D}$  and any two embeddings  $\gamma, \phi : \text{body}_u \rightarrow \mathbb{D}$  that coincide over  $\mathcal{S}$ , embeddings  $\gamma, \phi$  are also embeddings of  $\kappa(\text{body}_u) \rightarrow \mathbb{D}$  that coincide over  $\kappa(\mathcal{S})$ . Furthermore,  $\gamma$  and  $\phi$  coincide over  $\mathcal{P}_v$ , and so  $\gamma$  satisfies  $\text{pred}_v$  iff  $\phi$  satisfies  $\text{pred}_v$ . Suppose that  $\gamma'$  extends  $\gamma$  to an embedding of  $\text{body}_v^+$ ; then  $\gamma'$  is also an embedding of  $\kappa(\text{body}_v^+)$ . By  $\kappa(Q) \models \kappa(\mathcal{S}) \Rightarrow v$ , there exists an embedding  $\phi'$  that extends  $\phi$  to  $\kappa(\text{body}_v^+)$ , and  $\phi'$  is also an embedding of  $\text{body}_v^+$ . Therefore,  $Q \models \mathcal{S} \Rightarrow v$ .  $\square$

Each predicate in  $\kappa(\text{pred}_Q)$  is a tautology and can be deleted. Hence,  $\kappa(Q)$  does not contain equality predicates and so  $\kappa(Q) \models \kappa(\mathcal{S}) \Rightarrow v$  can be tested via Theorem 4.2.4.

Theorem 4.2.6 can be used to test whether the embeddings of an arbitrary DACQ (satisfying properties [P1]–[P4]) satisfy a given nullability dependency. Given DACQ  $Q$  (obeying properties [P1]–[P4]), vertex  $v \in \text{body}_Q$ , and some set  $\mathcal{S} \subseteq \mathcal{B}_v^{\text{anc}}$ , create query  $Q'$  from  $Q$  by first deleting all nodes not in  $\{v\} \cup \text{ancestors}(v)$  and then merging together all of the nodes in  $\text{ancestors}(v)$ . It follows from Definition 4.2.2 that  $Q \models \mathcal{S} \Rightarrow v$  if and only if  $Q' \models \mathcal{S} \Rightarrow v$ . This further yields the following two corollaries of Theorem 4.2.6.

**Corollary 4.2.7** *Given any DACQ  $Q$ , there exists a unique minimal set  $\mathcal{S}_Q \subseteq \mathcal{B}$  such that for every non-root node  $v \in \text{body}_Q$  and every set  $\mathcal{S} \subseteq \mathcal{B}_v^{\text{anc}}$ , if  $Q \models \mathcal{S} \Rightarrow v$  then  $Q \models (\mathcal{S}_Q \cap \mathcal{S}) \Rightarrow v$ . Furthermore,  $\mathcal{S}_Q$  contains all variables occurring within  $\text{pred}_Q$ .*

**Corollary 4.2.8** *Set  $\mathcal{S}_Q$  can be computed in NP time.*

We call  $\mathcal{S}_Q$  the *nullability set* of  $Q$ . We will not actually use set  $\mathcal{S}_Q$  within this chapter; however, it plays a prominent role within the next chapter when we define a normal form for HCEQs.

## 4.2.2 Converting DACQs to Canonical Form

We now show how property [P5] can be characterized in terms of nullability dependencies. Using this test, we can then define an algorithm for converting an arbitrary DACQ to an equivalent canonical form.

**Theorem 4.2.9** *Given*

1. *any DACQ  $Q$  conforming to properties [P1]–[P4] of Definition 4.2.1, and*
2. *any two distinct vertices  $u, v \in \text{body}_Q$  such that  $v$  is not an ancestor of  $u$ ;*

*create the query  $Q_u$  from  $Q$  by merging  $u$  with all of its ancestors into the root vertex of the query body. Let  $\vartheta$  be a mapping from the vertex labels of  $Q$  to the vertex labels of  $Q_u$ . Then, the following property holds*

$$\forall D. \forall G \in \text{dags}(Q). [(u \in G \wedge v \notin G) \implies (\{\gamma : G \rightsquigarrow D\} = \emptyset)] \quad (4.3)$$

*if and only if*

1. *there exists a homomorphism  $h : \text{body}_v^+ \rightarrow \text{body}_u^+$  such that  $h$  is the identity over  $\mathcal{B}_u^+ \cap \mathcal{B}_v^+$ , and*
2. *for all  $t \in \text{ancestors}(v) \cup \{v\}$ ,  $Q_u \models \mathcal{B}_u^+ \cap \mathcal{B}_v^+ \Rightarrow \vartheta(t)$ .*

PROOF. Suppose that property 4.3 holds.

1. Property 4.3 implies that any embedding  $\gamma : \text{body}_u^+ \rightarrow \mathbb{D}$  can be extended into an embedding  $\gamma' : \text{body}_u^+ \cup \text{body}_v^+ \rightarrow \mathbb{D}$  (note that  $\text{body}_u^+ \cup \text{body}_v^+$  is a CQ body since  $\text{body}_u^+$  and  $\text{body}_v^+$  are each *sets* of atoms and predicates). This implies that the following two CQs  $Q'$  and  $Q''$  satisfy  $Q' \sqsubseteq Q''$

$$\begin{aligned} Q'(\overline{\mathcal{B}_u^+}) &:- \text{body}_u^+ \\ Q''(\overline{\mathcal{B}_u^+}) &:- \text{body}_u^+ \cup \text{body}_v^+ \end{aligned}$$

which, by the Homomorphism Theorem for CQ containment [6], implies that there exists a homomorphism  $h' : Q'' \rightarrow Q'$  satisfying  $h'(\overline{\mathcal{B}_u^+}) = \overline{\mathcal{B}_u^+}$ . This mapping  $h'$  therefore satisfies the criteria for mapping  $h$  given above.

2. For any  $\mathbb{D}$ , let  $\gamma, \phi$  be any two embeddings of  $Q_u$  that are consistent over  $\mathcal{B}_u^+ \cap \mathcal{B}_v^+$ . For any  $t \in \text{ancestors}(v) \cup v$ , property 4.3 guarantees that both  $\gamma$  and  $\phi$  can be extended into embeddings of  $\vartheta(t)$ , and so  $Q_u \models \mathcal{B}_u^+ \cap \mathcal{B}_v^+ \Rightarrow \vartheta(t)$  follows by Definition 4.2.2.

We prove the opposite direction by contradiction. Suppose that  $h$  exists and that  $Q_u$  satisfies the nullability dependency for each  $t \in \text{ancestors}(v) \cup \{v\}$ . Suppose further than property 4.3 is violated; that is, there exists some database  $\mathbb{D}$  and  $G \in \text{dags}(Q)$  with  $u \in G$  and  $v \notin G$  that allows some terminal embedding  $\gamma : G \rightsquigarrow \mathbb{D}$ . Choose any vertex  $t \in \text{ancestors}(v) \cup \{v\}$  that does not appear in  $G$  but whose ancestors all appear in  $G$  (there must exist a choice for  $t$ , since  $v \notin G$  and the root node is both an ancestor of  $v$  and in  $G$ ). Then,

1.  $\gamma$  is an embedding of  $\text{body}_{\vartheta(t)}^{\text{anc}}$  into  $\mathbb{D}$ ;
2.  $\gamma \circ h$  is an embedding of  $\text{body}_{\vartheta(t)}^{\text{anc}}$  into  $\mathbb{D}$ ;
3.  $\gamma$  and  $\gamma \circ h$  are consistent over  $\mathcal{B}_u^+ \cap \mathcal{B}_v^+$  (by definition of  $h$  and  $\mathcal{B}_u^+ \cap \mathcal{B}_v^+ \subseteq \mathcal{B}_t^{\text{anc}}$ );
4.  $\gamma \circ h$  extends to  $\vartheta(v)$  and therefore also to  $\vartheta(t)$ ; and
5.  $Q_u \models \mathcal{B}_u^+ \cap \mathcal{B}_v^+ \Rightarrow \vartheta(t)$  implies that  $\gamma$  must extend to  $\vartheta(t)$  within  $\text{body}_{Q_u}$ , which further implies that  $\gamma$  extends to  $t$  within  $\text{body}_Q$ .

The above reasoning contradicts the choice of  $\gamma$  as a *terminal* embedding of  $G$ . Therefore, property 4.3 must hold.  $\square$

**Corollary 4.2.10** *For a given  $Q$  and pair of vertices  $u, v$ , testing property 4.3 is NP-complete.*

**Corollary 4.2.11** *Testing whether  $Q$  satisfies property [P5] is co-NP-complete.*

PROOF. To prove that  $Q$  does *not* satisfy [P5], it suffices to guess the pair of vertices  $u, v \in \text{body}_Q$  that cause the violation (i.e., that satisfy property 4.3), along with the homomorphism  $h$  from Theorem 4.2.4 and, for each  $t \in \text{ancestors}(v) \cup \{v\}$  the homomorphisms needed to verify  $Q_u \models \mathcal{B}_u^+ \cap \mathcal{B}_v^+ \Rightarrow \vartheta(t)$  (cf. Theorems 4.2.4 and 4.2.6).  $\square$

**Corollary 4.2.12** *For a given  $Q$  and pair of vertices  $u, v$ , property 4.3 holds if*

1. *there exists a homomorphism  $h : \text{body}_v^+ \rightarrow \text{body}_u^+$  such that  $h$  is the identity over  $\mathcal{B}_u^+ \cap \mathcal{B}_v^+$ , and*

---

**Algorithm 4** Converting DACQs into canonical form
 

---

CANONICAL( $Q$ )

▷ **Input:** satisfiable DACQ  $Q(\overline{W})$   
 ▷ **Output:** an equivalent DACQ in canonical form  
 1 **while**  $\exists u \in \text{body}_Q$  such that  $\text{body}_u^+$  is unsatisfiable  
 2     **do** delete vertex  $u$  and all its descendants from  $\text{body}_Q$   
 3         substitute  $\blacksquare$  for each component of  $\overline{W}$  labelled with a deleted vertex  
 4 **while**  $\exists u \in \text{body}_Q$  such that  $u$  is a leaf,  $u$  is not the root,  
    and  $u$  does not occur as a label in  $\overline{W}$   
 5     **do** delete vertex  $u$  from  $\text{body}_Q$   
 6 **foreach**  $u \in \text{body}_Q, X \in \mathcal{P}_u$   
 7     **do if**  $\text{pred}_u^+$  equates  $X$  to  $Y \in \mathcal{P}_u \cup \mathcal{C}$   
 8         **then** replace any occurrence of  $X$  with  $Y$  within  $\text{atoms}_u$   
 9         replace any occurrence of  $X^{\overline{u}}$  with  $Y^{\overline{u}}$  within  $\overline{W}$   
 10 delete redundant atoms, predicates, and edges from  $\text{body}_Q$   
 11 **foreach** pair  $u, v \in \text{body}_Q$  such that  $u \neq v$  and  $v \notin \text{ancestors}(u)$   
 12     **do if**  $\text{TESTPROPERTY4.3}(u, v) = \text{true}$   
 13         **then** add an edge from  $v$  to  $u$   
 14             **while**  $\text{body}_Q$  contains a cycle  
 15                 **do** merge all participating vertices  
 16             delete redundant atoms and predicates from  $\text{body}_Q$   
 17 **return**  $Q$

---

2. for all  $t \in \text{ancestors}(v) \cup \{v\}$ ,  $\mathcal{P}_t \subseteq (\mathcal{B}_u^+ \cap \mathcal{B}_v^+)$ .

PROOF. By Lemma 4.2.3,  $\mathcal{P}_t \subseteq (\mathcal{B}_u^+ \cap \mathcal{B}_v^+)$  implies that  $Q \models \mathcal{B}_u^+ \cap \mathcal{B}_v^+ \Rightarrow t$ . By construction of  $Q_u$ , this further implies that  $Q_u \models \mathcal{B}_u^+ \cap \mathcal{B}_v^+ \Rightarrow \vartheta(t)$ .  $\square$

We now consider how to transform a query into an equivalent canonical form. We start with an example.

**Example 39** Reconsider query  $Q_{36}$  from Figure 4.5. It is straightforward to verify that  $Q_{36}$  satisfies properties **[P1]–[P3]**. To verify property **[P5]**, it suffices to consider each of the following pairs as bindings for vertices  $(u, v)$ .

$$\{(\overline{1}, \overline{2}), (\overline{1}, \overline{3}), (\overline{1}, \overline{4}), (\overline{2}, \overline{3}), (\overline{2}, \overline{4}), (\overline{3}, \overline{2}), (\overline{3}, \overline{4})\}$$

For each pair, the reader can verify that no homomorphism exists from  $\text{body}_v^+$  to  $\text{body}_u^+$ ; therefore,  $Q_{36}$  satisfies **[P5]**. Hence, to convert  $Q_{36}$  to canonical form we simply need to delete the redundant atoms so that it satisfies **[P4]**. The resulting query is variable-isomorphic to query  $Q_{35}$ .

The function CANONICAL( $Q$ ) shown in Algorithm 4 transforms an arbitrary satisfiable DACQ into an equivalent query in canonical form. It relies upon invoking

the function  $\text{TESTPROPERTY4.3}(u, v)$ , which is not shown but directly implements the homomorphic conditions of Theorem 4.2.9. We now briefly argue the correctness of this algorithm.

- **Lines 1–3:** Because  $\text{body}_u^+$  is unsatisfiable, neither  $u$  nor any of its descendants can participate in any embedding. Deleting  $u$  and its descendants from  $\text{body}_Q$  does not change the set of terminal embeddings, and the partial application of any terminal embedding to any component of  $\overline{W}$  labelled with  $u$  or any of its descendants always returns  $\blacksquare$ . Hence, when the loop terminates the modified query remains equivalent and obeys property **[P1]**.
- **Lines 4–5:** Let  $G$  be any DAG in  $\text{dags}(\text{body}_Q)$  such that  $u \in G$ , and let  $G' \in \text{dags}(\text{body}_Q)$  be the DAG formed by deleting  $u$  from  $G$ . Let  $\gamma : G \rightsquigarrow \mathbb{D}$  be any terminal embedding of  $G$  into some database  $\mathbb{D}$ , and let  $\gamma' : G' \rightarrow \mathbb{D}$  be the restriction of  $\gamma$  to  $G'$ . Deleting  $u$  from  $\text{body}_Q$  deletes terminal embedding  $\gamma$  but causes  $\gamma'$  to become a terminal embedding; furthermore, the partial applications of  $\gamma$  and  $\gamma'$  to  $\overline{W}$  are identical. Hence, when the loop terminates the modified query remains equivalent and obeys properties **[P1]** and **[P2]**.
- **Lines 6–9:** For any  $G \in \text{dags}(\text{body}_Q)$  containing vertex  $u$ , any embedding  $\gamma$  of  $G$  already satisfies  $\gamma(X) = \gamma Y$ . Hence, when the loop terminates the modified query remains equivalent and obeys properties **[P1]**–**[P3]**.
- **Line 10:** For any vertex  $u$  and any  $G \in \text{dags}(\text{body}_Q)$  containing  $u$ ,  $G$  necessarily contains all of  $u$ 's ancestors and so deleting redundant atoms and predicate from  $u$  does not change the embeddings for  $G$  nor the set of terminal embeddings. Deleting a redundant edge does not change the set  $\text{dags}(\text{body}_Q)$ , and so does not affect the set of terminal embeddings. Hence, after this step the modified query remains equivalent and obeys properties **[P1]**–**[P4]**.
- **Lines 11–16:** If  $\text{TESTPROPERTY4.3}(u, v)$  evaluates to **true**, then any *terminal* embedding that includes vertex  $u$  must also include vertex  $v$ . Therefore, adding edge  $(v, u)$  does not change the set of terminal embeddings, although it can both introduce cycles into  $\text{body}_Q$  and cause some atoms or predicates to become redundant. Any terminal embedding must include either all or none of the vertices in a cycle, and so merging cycles does not affect the set of terminal embeddings. Hence, after all cycles have been merged and redundant atoms and predicates deleted, the modified query is again a proper DACQ that remains equivalent and obeys properties **[P1]**–**[P4]**. Once this process has been repeated for all pairs of vertices in  $\text{body}_Q$ , the query also satisfies property **[P5]**.

**Proposition 4.2.13** *Given any satisfiable DACQ  $Q$ ,  $\text{CANONICAL}(Q)$  is a DACQ in canonical form satisfying  $Q \equiv \text{CANONICAL}(Q)$ .*

**Example 40** Consider query  $Q_{40}$  shown in Figure 4.7. The equivalent canonical form  $Q_{41} = \text{CANONICAL}(Q_{40})$ —shown in Figure 4.8—is calculated as follows.

1.  $\text{body}_5^+$  is unsatisfiable, so delete nodes  $\boxed{5}$  and  $\boxed{8}$ .

$$\text{body}_5^+ := R(A, B, C), R(A, D, 2), S(G, H), A = 2, A = 3$$

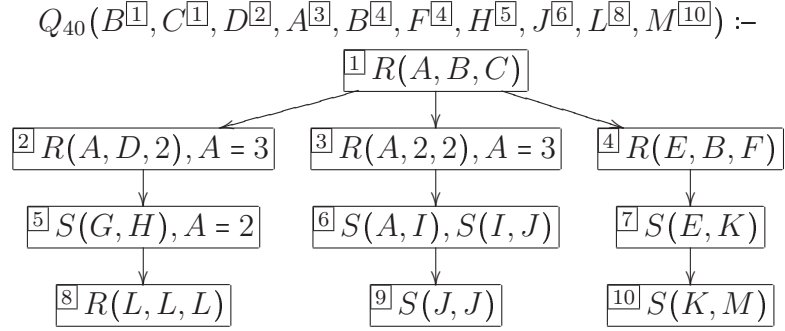


Figure 4.7: HCQ  $Q_{40}$  prior to conversion to canonical form

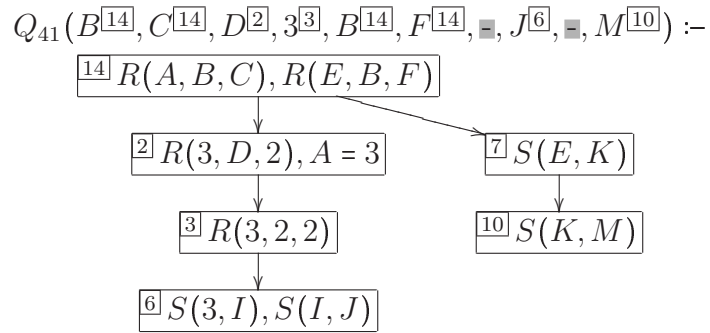


Figure 4.8: Query  $Q_{41} = \text{CANONICAL}(Q_{40})$

2. Delete non-root leaf node  $\boxed{9}$  because label  $\boxed{9}$  does not occur in the query head.
3. Replace variable  $A$  with the constant 3 within the atoms of nodes  $\boxed{2}$ ,  $\boxed{3}$ , and  $\boxed{6}$ , and replace  $A^{\boxed{3}}$  with  $3^{\boxed{3}}$  within the query head.
4. Add an edge from  $\boxed{2}$  to  $\boxed{3}$  because  $\text{TESTPROPERTY4.3}(\boxed{3}, \boxed{2})$  returns true, justified by Corollary 4.2.12 as follows. Sets  $\mathcal{B}_3^+ = \{A, B, C\}$  and  $\mathcal{B}_2^+ = \{A, B, C, D\}$ , and so homomorphism  $h : \text{body}_2^+ \rightarrow \text{body}_3^+$  is the identity over  $\mathcal{B}_3^+ \cap \mathcal{B}_2^+$ .

$$\begin{aligned} \text{body}_3^+ &:= R(A, B, C), R(A, 2, 2), A = 3 \\ \text{body}_2^+ &:= R(A, B, C), R(A, D, 2), A = 3 \\ h &:= \{A/A, B/B, C/C, D/2\} \end{aligned}$$

Next,  $\text{ancestors}(\boxed{2}) \cup \{\boxed{2}\} = \{\boxed{1}, \boxed{2}\}$ , and so  $\mathcal{P}_1 = \emptyset$  and  $\mathcal{P}_2 = \{A\}$  are both subsets of  $\mathcal{B}_3^+ \cap \mathcal{B}_2^+$ . Adding this edge makes the predicate  $A = 3$  within node  $\boxed{3}$  redundant, and makes the edge  $(\boxed{1}, \boxed{3})$  redundant, so delete them.

5. Add an edge from  $\boxed{4}$  to  $\boxed{1}$  because  $\text{TESTPROPERTY4.3}(\boxed{1}, \boxed{4})$  returns true, justified by Corollary 4.2.12 as follows. Sets  $\mathcal{B}_1^+ = \{A, B, C\}$  and  $\mathcal{B}_4^+ = \{A, B, C, E, F\}$ , and so homomorphism  $h : \text{body}_4^+ \rightarrow \text{body}_1^+$  is the identity over  $\mathcal{B}_1^+ \cap \mathcal{B}_4^+$ .

$$\begin{aligned} \text{body}_1^+ &:= R(A, B, C) \\ \text{body}_4^+ &:= R(A, B, C), R(E, B, F) \\ h &:= \{A/A, B/B, C/C, E/A, F/C\} \end{aligned}$$

Next,  $\text{ancestors}(\boxed{4}) \cup \{\boxed{4}\} = \{\boxed{1}, \boxed{4}\}$ , and so  $\mathcal{P}_1 = \emptyset$  and  $\mathcal{P}_4 = \{B\}$  are both subsets of  $\mathcal{B}_1^+ \cap \mathcal{B}_4^+$ . Adding this edge causes a cycle, so merge nodes  $\boxed{1}$  and  $\boxed{4}$  into a single node  $\boxed{14}$ .

## 4.3 Equivalence of Constant-Headed Queries

Our goal in this chapter is to characterize the query equivalence problem for HCQs and, more generally, for DACQs. In this section we study the equivalence problem for a subclass of DACQs that we call *constant-headed*, which generalize boolean CQs. We say that a DACQ is *constant-headed* if its head tuple does not contain any variables, and we use CH-DACQ to denote the class of constant-headed DACQs.

We describe a method for testing CH-DACQ equivalence via reduction to a set of logical implication problems. Unfortunately, although this algorithm is sound, it is not complete for concluding equivalence because it is possible that some of the generated implication problems are only semi-decidable. As a result, our results in this section are inconclusive in that we have neither proved nor disproved whether CH-DACQ equivalence is decidable.

### 4.3.1 Head Graphs and Canonical Tuples

We will address CH-DACQ equivalence by considering separate equivalence problems for each *canonical tuple*. Defining a query's canonical tuples requires understanding a query's *tuple shapes*, which are characterized by the query *head graph*.

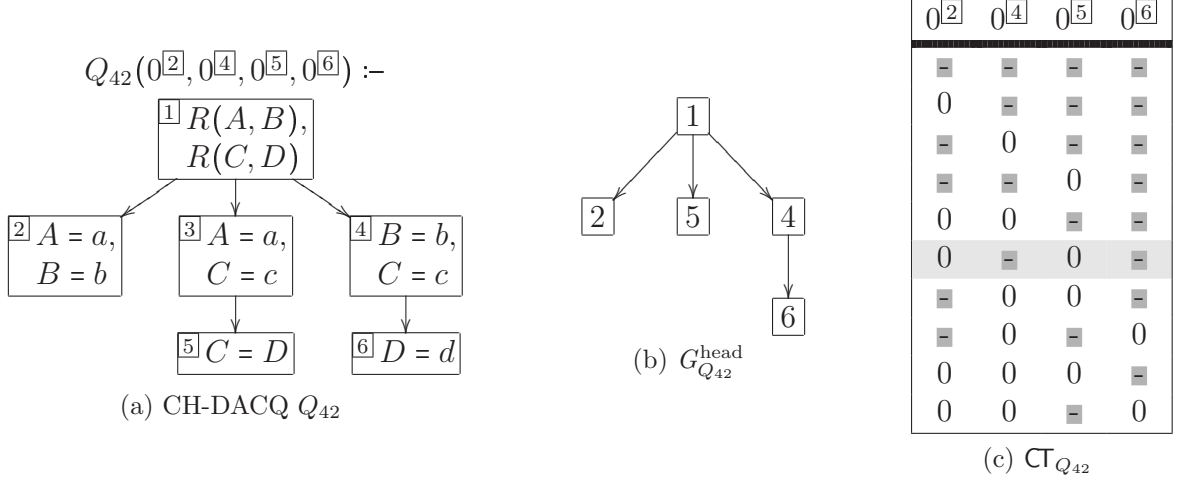


Figure 4.9: A constant-headed DACQ with its head graph and canonical tuples

(Although our focus here is on CH-DACQ equivalence, we will make certain definitions and theorems general to arbitrary DACQs so that we can reuse them in Section 4.4 without modification.)

**Definition 4.3.1 (Head Graph)** *Given DACQ  $Q(\overline{\mathcal{W}})$  with head label set  $L$ , the head graph of  $Q$  is a rooted DAG  $G_Q^{\text{head}} = (U, E)$  defined as follows.*

$$\begin{aligned}
 U &:= L \cup \{\boxed{u} \mid u \text{ is the root node of } \text{body}_Q\} \\
 E &:= \{(\boxed{v_i}, \boxed{v_k}) \mid \text{body}_Q \text{ contains path } v_i \rightsquigarrow v_k, \text{ and} \\
 &\quad \nexists \boxed{v_j} \in L. [\text{body}_Q \text{ contains path } v_i \rightsquigarrow v_j \rightsquigarrow v_k]\}
 \end{aligned}$$

**Definition 4.3.2 (Tuple Shapes)** *Given DACQ  $Q(\overline{\mathcal{W}})$  with head label set  $L$ , we define a tuple shape  $\mathfrak{ts}$  of  $Q$  to be a subset of  $L$  corresponding to the vertices of some DAG within  $\text{dags}(G_Q^{\text{head}})$ . We use  $\text{TS}_Q$  to denote the set of all tuple shapes of  $Q$ .*

$$\text{TS}_Q := \{\mathfrak{ts} \mid \exists (U, E) \in \text{dags}(G_Q^{\text{head}}) \wedge \mathfrak{ts} = L \cap U\} \quad (4.4)$$

We call shape  $\mathfrak{ts} = L$  the terminal shape for  $Q$ .

**Example 41** Figures 4.9a and 4.9b show constant-headed query  $Q_{42}$  and its head graph  $G_{Q_{42}}^{\text{head}}$ , respectively. The tuple shapes of  $Q_{42}$  are as follows,

$$\text{TS}_{42} = \{\mathfrak{ts}_{\emptyset}, \mathfrak{ts}_2, \mathfrak{ts}_4, \mathfrak{ts}_5, \mathfrak{ts}_{24}, \mathfrak{ts}_{25}, \mathfrak{ts}_{45}, \mathfrak{ts}_{46}, \mathfrak{ts}_{245}, \mathfrak{ts}_{246}, \mathfrak{ts}_{2456}\}$$

where the subscripts denote the set of labels in each shape. Observe that the relation in Figure 4.9c (which we will describe further in Example 43) contains one tuple corresponding to each tuple shape in  $\text{TS}_{42}$  *except* for the terminal shape  $\mathfrak{ts}_{2456}$ .

Next we show that one necessary condition for equivalence between queries in canonical form is that they have isomorphic head graphs.



**Definition 4.3.3 (Head Graph Isomorphism)** Given DACQs  $Q(\overline{W})$  and  $Q'(\overline{W}')$  with head graphs  $G_Q^{\text{head}} = (U, E)$  and  $G_{Q'}^{\text{head}} = (U', E')$ , a head graph isomorphism  $\theta$  is a bijective mapping  $\theta : U \rightarrow U'$  satisfying the two properties

1.  $\theta(\overline{L}) = \overline{L}'$ , and
2.  $\theta(G_Q^{\text{head}}) = G_{Q'}^{\text{head}}$ .

We say that  $Q$  and  $Q'$  are head graph isomorphic if  $\theta$  exists.

The existence of a head graph isomorphism can be tested in linear time because the property  $\theta(\overline{L}) = \overline{L}'$  dictates the precise bijection between the vertices of the head graphs. Observe that if  $\theta$  is a head graph isomorphism between  $Q$  and  $Q'$ , then  $\theta(\text{TS}_Q) = \text{TS}_{Q'}$  follows from Definitions 4.3.2 and 4.3.3.

**Theorem 4.3.4** Let  $Q(\overline{W})$  and  $Q'(\overline{W}')$  be DACQs in canonical form. If  $Q \equiv Q'$  then  $Q$  and  $Q'$  are head graph isomorphic.

PROOF (SKETCH). We use properties [P1] and [P5] of Definition 4.2.1 and an analysis of cases under which  $\theta$  does not exist to establish that if  $\theta$  does not exist, then there exists a tuple that only one of the queries is able to output.  $\square$

The shapes in  $\text{TS}_Q$  can be used to partition the set  $\text{dags}(\text{body}_Q)$  so that associated with each shape  $\mathfrak{ts} \in \text{TS}_Q$  is a subset of  $\text{dags}(\text{body}_Q)$  comprised of the DAGs whose terminal embeddings yield the tuples corresponding to shape  $\mathfrak{ts}$ .

$$\text{dags}(\text{body}_Q, \mathfrak{ts}) := \{G \mid G = (U, E) \in \text{dags}(\text{body}_Q) \quad (4.5)$$

and  $U$  contains all of the labels in  $\mathfrak{ts}$

and  $U$  contains none of the labels in  $L \setminus \mathfrak{ts}$

$$\text{dags}(\text{body}_Q) = \bigcup_{\mathfrak{ts} \in \text{TS}_Q} \text{dags}(\text{body}_Q, \mathfrak{ts}) \quad (4.6)$$

Observe that equations 4.4 and 4.5 together imply that set  $\text{dags}(\text{body}_Q, \mathfrak{ts})$  contains a unique minimal DAG—which we denote as  $G_{\mathfrak{ts}}^{\text{min}}$ —such that all other DAGs in  $\text{dags}(\text{body}_Q, \mathfrak{ts})$  are extensions of  $G_{\mathfrak{ts}}^{\text{min}}$ .

$$G_{\mathfrak{ts}}^{\text{min}} := G \in \text{dags}(\text{body}_Q, \mathfrak{ts}). \text{ all leaves of } G \text{ are in } \mathfrak{ts} \quad (4.7)$$

We use  $G_{\mathfrak{ts}}^{\text{min}}$  to define a particular  $\text{CQ}^-$  we call the *shape- $\mathfrak{ts}$  query*.

**Definition 4.3.5 (Shape Queries)** Given DACQ  $Q(\overline{W})$ , for each shape  $\mathfrak{ts} \in \text{TS}_Q$  let  $\overline{W}|_{\mathfrak{ts}}$  denote the head tuple formed by changing to  $\blacksquare$  any component of  $\overline{W}$  whose label is not in  $\mathfrak{ts}$ . We define the shape- $\mathfrak{ts}$  query  $\hat{Q}|_{\mathfrak{ts}}$  as the  $\text{CQ}^-$  with head tuple  $\overline{W}|_{\mathfrak{ts}}$  and body  $\text{body}_{G_{\mathfrak{ts}}^{\text{min}}}^+$ , and we define shape query  $\hat{Q}$  as the union of all shape- $\mathfrak{ts}$  queries.

$$\hat{Q}|_{\mathfrak{ts}}(\overline{W}|_{\mathfrak{ts}}) := \text{body}_{G_{\mathfrak{ts}}^{\text{min}}}^+ \quad (4.8)$$

$$\hat{Q} := \bigcup_{\mathfrak{ts} \in \text{TS}_Q} \hat{Q}|_{\mathfrak{ts}} \quad (4.9)$$

**Example 42** Continuing Example 41, consider shape  $\mathfrak{ts}_\emptyset \in \mathbf{TS}_{Q_{42}}$ , which corresponds to the following subset of  $\text{dags}(\text{body}_{Q_{42}})$  (subscripts denote labels within each DAG).

$$\text{dags}(\text{body}_{Q_{42}}, \mathfrak{ts}_\emptyset) = \{G_1, G_{13}\}$$

Shape  $\mathfrak{ts}_\emptyset$  has minimal dag  $G_{\mathfrak{ts}_\emptyset}^{\min} = G_1$ , yielding shape- $\mathfrak{ts}_\emptyset$  query  $\hat{Q}_{42}|_{\mathfrak{ts}_\emptyset}$  as follows.

$$\hat{Q}_{42}|_{\mathfrak{ts}_\emptyset}(\blacksquare, \blacksquare, \blacksquare, \blacksquare) :- R(A, B), R(C, D)$$

**Theorem 4.3.6** *Given DACQ  $Q$  and any database  $\mathbb{D}$ , the output instance of shape query  $\hat{Q}$  corresponds precisely to the partial applications of all (not just terminal) embeddings of  $\text{body}_Q$  into  $\mathbb{D}$ .*

$$\mathbb{T}((\hat{Q})^{\mathbb{D}}) = \{\hat{\gamma}(\overline{\mathcal{W}}) \mid \exists G \in \text{dags}(\text{body}_Q), \gamma : G \rightarrow \mathbb{D}\} \quad (4.10)$$

PROOF. By equations 4.5 and 4.6, every  $G \in \text{dags}(\text{body}_Q)$  corresponds to precisely one shape  $\mathfrak{ts} \in \mathbf{TS}_Q$ , and so  $G$  is an extension of the minimal DAG  $G_{\mathfrak{ts}}^{\min}$ . Any embedding  $\gamma : G \rightarrow \mathbb{D}$  is therefore also an embedding of  $G_{\mathfrak{ts}}^{\min}$  into  $\mathbb{D}$ , and head tuple  $\overline{\mathcal{W}}|_{\mathfrak{ts}}$  (Definition 4.3.5) is designed so that  $\hat{\gamma}(\overline{\mathcal{W}}) = \gamma(\overline{\mathcal{W}}|_{\mathfrak{ts}})$  (see equation 2.22).  $\square$

**Corollary 4.3.7**  $Q \equiv \hat{Q}$

**Corollary 4.3.8** *Given DACQs  $Q$  and  $Q'$  in canonical form, if  $Q \equiv Q'$  then*

$$\forall \mathfrak{ts} \in \mathbf{TS}_Q. [\hat{Q}|_{\mathfrak{ts}} \equiv \hat{Q}'|_{\theta(\mathfrak{ts})}]$$

where  $\theta$  is the head graph isomorphism guaranteed by Theorem 4.3.4.

Corollary 4.3.8 is important because it provides another necessary condition for DACQ equivalence; it can be tested in NP-time per shape  $\mathfrak{ts} \in \mathbf{TS}_Q$ . Roughly speaking, the condition is analogous to explicitly verifying equivalence between each pair of terms in the JDNF of two join queries (see Section 4.6.1). In our context this is only a necessary condition for equivalence, whereas for join queries it is necessary and sufficient for equivalence (due to lack of a projection operator). The complicating effect of projection is that it may project away variables used in outer-join predicates (i.e., any variable used as a parameter variable in a child  $\text{CQ}^-$ ). We can obtain a sufficient condition for query equivalence by dictating a particular agreement between the attributes used in the outer-join predicates. To do this, first define for each shape  $\mathfrak{ts} \in \mathbf{TS}_Q$  the set of parameter variables  $\mathcal{P}_{\mathfrak{ts}} \subseteq \mathcal{P}_Q^+$ , as follows.

$$\mathcal{P}_{\mathfrak{ts}} := \bigcup_{v \in G_{\mathfrak{ts}}^{\min}} \mathcal{P}_v$$

**Theorem 4.3.9** *Given DACQs  $Q$  and  $Q'$  in canonical form, if*

- *there exists a head graph isomorphism  $\theta$  between  $Q$  and  $Q'$ ,*
- *there exists a bijection  $\nu : \mathcal{P}_Q^+ \rightarrow \mathcal{P}_{Q'}^+$ , satisfying  $\forall \mathfrak{ts} \in \mathbf{TS}_Q. [\nu(\mathcal{P}_{\mathfrak{ts}}) = \mathcal{P}_{\theta(\mathfrak{ts})}]$ , and*

- for each shape  $\mathfrak{ts} \in \mathbf{TS}_Q$ , the CQ equivalence  $\hat{Q}|_{\mathfrak{ts}} \equiv \hat{Q}'|_{\theta(\mathfrak{ts})}$  can be proved by homomorphisms  $h_{\mathfrak{ts}} : \hat{Q}'|_{\theta(\mathfrak{ts})} \rightarrow \hat{Q}|_{\mathfrak{ts}}$  and  $h'_{\mathfrak{ts}} : \hat{Q}|_{\mathfrak{ts}} \rightarrow \hat{Q}'|_{\theta(\mathfrak{ts})}$  that are consistent with  $\nu$ ,

then  $Q \equiv Q'$ .

PROOF. Suppose that the conditions of the theorem are met. For any database  $\mathbb{D}$ , and any tuple shape  $\mathfrak{ts} \in \mathbf{TS}_Q$ , let  $\gamma$  be any terminal embedding of  $Q$  into  $\mathbb{D}$  yielding a tuple  $t$  with shape  $\mathfrak{ts}$ . Then,  $\gamma$  is also an embedding of  $\hat{Q}|_{\mathfrak{ts}}$  into  $\mathbb{D}$ , and so  $\gamma \circ h_{\mathfrak{ts}}$  is an embedding of  $\hat{Q}'|_{\theta(\mathfrak{ts})}$  into  $\mathbb{D}$ , and therefore  $\gamma \circ h_{\mathfrak{ts}}$  is an embedding of  $Q'$  into  $\mathbb{D}$  whose partial application yields  $t$ .

If  $\gamma \circ h_{\mathfrak{ts}}$  is a terminal embedding then we are done; hence, suppose that there exists some embedding  $\phi : Q' \rightarrow \mathbb{D}$  that extends  $\gamma \circ h_{\mathfrak{ts}}$  and generates a tuple  $t'$  with shape  $\mathfrak{ts}'$  that subsumes shape  $\mathfrak{ts}$ . In this case,  $\phi \circ h'_{\mathfrak{ts}'}$  is an embedding of  $Q$  into  $\mathbb{D}$  yielding  $t'$ . Although  $\phi \circ h'_{\mathfrak{ts}'}$  is not necessarily an extension of  $\gamma$ , it is consistent with  $\gamma$  over any variable in  $\mathcal{P}_Q^+$  over which  $\gamma$  is defined. It follows that  $\gamma$  must have an extension that generates tuple  $t'$ , which contradicts the choice of  $\gamma$  as a terminal embedding.  $\square$

The sufficient condition in Theorem 4.3.9 can actually be tested in NP-time, because agreement of all of the homomorphisms on the parameter variables means that there must exist a single homomorphism in each direction that can prove equivalence for each pair of shape queries. We defer an example of Theorem 4.3.9 until the end of Example 47.

Having defined both head graphs and shape queries, we can now formally define the canonical tuples for a given CQ-DACQ.

**Definition 4.3.10 (Canonical Tuples)** *Given CH-DACQ  $Q$  in canonical form, the set of canonical tuples for  $Q$ —denoted  $\mathbf{CT}_Q$ —corresponds to the head tuples of satisfiable shape queries.*

$$\mathbf{CT}_Q := \{\overline{\mathcal{W}}|_{\mathfrak{ts}} \mid \mathfrak{ts} \in \mathbf{TS}_Q \wedge \hat{Q}|_{\mathfrak{ts}} \text{ is satisfiable}\} \quad (4.11)$$

For a given canonical tuple  $t \in \mathbf{CT}_Q$ , we use  $\mathfrak{ts}_t^Q$  to denote the corresponding tuple shape from  $\mathbf{TS}_Q$ .

**Lemma 4.3.11** *Given CH-DACQ  $Q$  and any database  $\mathbb{D}$ ,  $\mathbf{T}((Q)^{\mathbb{D}}) \subseteq \mathbf{CT}_Q$ .*

PROOF. By Definition 4.3.10 and Corollary 4.3.7, each  $t \in \mathbf{T}((Q)^{\mathbb{D}})$  corresponds to some satisfiable query  $\hat{Q}|_{\mathfrak{ts}}$  with  $t = \overline{\mathcal{W}}|_{\mathfrak{ts}}$ .  $\square$

**Lemma 4.3.12** *Given CH-DACQs  $Q$  and  $Q'$  in canonical form, if  $Q \equiv Q'$  then  $\mathbf{CT}_Q = \mathbf{CT}_{Q'}$ .*

PROOF. Follows from Definition 4.3.10 because  $\theta(\mathbf{TS}_Q) = \mathbf{TS}_{Q'}$  (Theorem 4.3.4) and  $\hat{Q}|_{\mathfrak{ts}} \equiv \hat{Q}'|_{\theta(\mathfrak{ts})}$  (Corollary 4.3.8).  $\square$

Even though the set of canonical tuples  $\text{CT}_Q$  is defined in terms of satisfiable shape queries,  $\text{CT}_Q$  can still contain spurious tuples that  $Q$  cannot actually output. We illustrate this in the following example.

**Example 43** Continuing Examples 41 and 42, Figure 4.9c shows the canonical tuples of  $Q_{42}$  (depicted as a relation). Consider the following two shape queries  $\hat{Q}_{42|\mathfrak{ts}_{25}}$  and  $\hat{Q}_{42|\mathfrak{ts}_{2456}}$ , whose bodies correspond to DAGs  $G_{125}, G_{12456} \in \text{dags}(\text{body}_{Q_{42}})$ .

$$\begin{aligned} \hat{Q}_{42|\mathfrak{ts}_{25}}(0, \blacksquare, 0, \blacksquare) &:- R(A, B), R(C, D), A = a, B = b, C = c, C = D \\ \hat{Q}_{42|\mathfrak{ts}_{2456}}(0, 0, 0, 0) &:- R(A, B), R(C, D), A = a, B = b, C = c, C = D, D = d \end{aligned}$$

The shaded row in Figure 4.9c is canonical tuple  $t_{25} = \text{head}_{\hat{Q}_{42|\mathfrak{ts}_{25}}}$ . Shape query  $\hat{Q}_{42|\mathfrak{ts}_{25}}$  is satisfiable, yet  $Q_{42}$  will never output  $t_{25}$  because any embedding  $\gamma$  of  $\hat{Q}_{42|\mathfrak{ts}_{25}}$  is an embedding of  $G_{125}$  that can be extended to an embedding of  $G_{1245}$ , and therefore  $\gamma$  is not terminal. The tuple  $t_{2456} = \langle 0, 0, 0, 0 \rangle = \text{head}_{\hat{Q}_{42|\mathfrak{ts}_{2456}}}$  is not a canonical tuple—even though  $\mathfrak{ts}_{2456} \in \text{TS}_{Q_{42}}$ —because shape query  $\hat{Q}_{42|\mathfrak{ts}_{2456}}$  is unsatisfiable.

### 4.3.2 Tuple-Specific Equivalence

Having defined a query’s canonical tuples, we now reduce CH-DACQ equivalence to a set of *tuple-specific* equivalence tests, one per canonical tuple.

**Definition 4.3.13 (Tuple-Specific Equivalence)** *Given any tuple  $t$  and any two DACQs  $Q$  and  $Q'$ , we write  $Q \stackrel{t}{\equiv} Q'$  to mean  $\forall \mathbb{D}. [t \in (Q)^{\mathbb{D}} \iff t \in (Q')^{\mathbb{D}}]$ .*

**Theorem 4.3.14** *Given CH-DACQs  $Q$  and  $Q'$  in canonical form,  $Q \equiv Q'$  iff*

1.  $\text{CT}_Q = \text{CT}_{Q'}$ , and
2.  $\forall t \in \text{CT}_Q, Q \stackrel{t}{\equiv} Q'$ .

PROOF. “If” follows from Lemma 4.3.11 and Definition 4.3.13. “Only if” follows from Lemma 4.3.12 and Definition 4.3.13.  $\square$

In the remainder of this section we consider how to decide the tuple-specific equivalence condition that arises from Theorem 4.3.14. Our presentation of the problem is organized by the complexity of the head graph.

#### Tuple-Specific Equivalence: Single-Node Head Graphs

Let  $Q$  denote a CH-DACQ in canonical form whose head graph contains only one node  $\underline{u}$ . This implies  $\text{dags}(G_Q^{\text{head}}) = \{G_u\}$ , and so  $\text{CT}_Q = \{t\}$  for some (possibly empty) tuple  $t$ . By Definition 4.3.1, vertex  $u$  must be the root vertex of  $\text{body}_Q$ . Because  $Q$  is in canonical form, property [P2] of Definition 4.2.1 implies that  $\text{body}_Q$  is also only a single vertex—hence,  $Q$  is simply a boolean CQ. Given query  $Q'$  with

$\text{CT}_{Q'} = \text{CT}_Q$ , testing tuple-specific equivalence  $Q \stackrel{t}{\equiv} Q'$  therefore corresponds to boolean CQ equivalence, which corresponds to the NP-complete problem of finding homomorphisms in each direction between  $Q$  and  $Q'$ .

Before moving on to queries with more complex head graphs, we first present an alternative characterization of equivalence in terms of finite implication of dependencies. While this may seem like overkilling the problem for this case, in later sections we will extend this alternative characterization to handle more complex head graphs.

An embedded dependency  $\xi$  can be expressed as a rule of the form

$$P^\xi \rightarrow C^\xi$$

where  $P^\xi$ —called the *premise*—and  $C^\xi$ —called the *conclusion*—are both CQ<sup>+</sup> bodies [6, 30]. Any variable occurring in both the premise and the conclusion is called *free*, and is implicitly universally quantified; all other variables are implicitly existentially quantified. Dependency  $\xi$  is satisfied by database instance  $\mathbb{D}$ —denoted  $\mathbb{D} \models \xi$ —if each embedding of the premise into  $\mathbb{D}$  can be extended to an embedding of the conclusion. Without loss of generality we can assume that  $P^\xi \cap C^\xi = \emptyset$ , because any atom or predicate occurring on both sides of the dependency is trivially implied and can therefore be deleted from  $C^\xi$ . Given a set of dependencies  $\Sigma$  and a database instance  $\mathbb{D}$ , we write  $\mathbb{D} \models \Sigma$  to mean  $\forall \xi \in \Sigma. (\mathbb{D} \models \xi)$ . Given two dependency sets  $\Sigma_1$  and  $\Sigma_2$ , we write  $\Sigma_1 \models_{\text{fin}} \Sigma_2$  and  $\Sigma_1 \equiv_{\text{fin}} \Sigma_2$  to mean finite implication  $\forall \mathbb{D}. (\mathbb{D} \models \Sigma_1 \implies \mathbb{D} \models \Sigma_2)$  and finite equivalence  $\forall \mathbb{D}. (\mathbb{D} \models \Sigma_1 \iff \mathbb{D} \models \Sigma_2)$ , respectively ( $\models_{\text{unr}}$  and  $\equiv_{\text{unr}}$  are analogous for when the quantification over  $\mathbb{D}$  allows infinite database instances). Deciding  $\models_{\text{fin}}$  and  $\equiv_{\text{fin}}$  are mutually-reducible problems,

$$\Sigma_1 \models_{\text{fin}} \Sigma_2 \iff \Sigma_1 \equiv_{\text{fin}} (\Sigma_1 \cup \Sigma_2) \tag{4.12}$$

$$\Sigma_1 \equiv_{\text{fin}} \Sigma_2 \iff (\Sigma_1 \models_{\text{fin}} \Sigma_2) \wedge (\Sigma_2 \models_{\text{fin}} \Sigma_1) \tag{4.13}$$

and so although we use finite equivalence for conceptual simplicity, it corresponds to the widely-studied problem of finite implication.

Corresponding to node  $u$  of query  $Q$  and canonical tuple  $t$ , we define the following dependency  $\xi^{(u,t)}$  and (singleton) dependency set  $\Sigma^{(Q,t)}$ .

$$\xi^{(u,t)} := \text{body}_u^+ \rightarrow \perp \tag{4.14}$$

$$\Sigma^{(Q,t)} := \{\xi^{(u,t)}\} \tag{4.15}$$

**Lemma 4.3.15** *For any database  $\mathbb{D}$ ,  $t \in (Q)^{\mathbb{D}}$  iff  $\mathbb{D} \not\models \Sigma^{(Q,t)}$ .*

PROOF. By equation 4.14,  $\mathbb{D} \models \xi^{(u,t)}$  iff  $\text{body}_u^+$  does not embed into  $\mathbb{D}$ . □

**Corollary 4.3.16**  *$Q \stackrel{t}{\equiv} Q'$  iff  $\Sigma^{(Q,t)} \equiv_{\text{fin}} \Sigma^{(Q',t)}$ .*

Deciding finite equivalence between sets  $\Sigma^{(Q,t)}$  and  $\Sigma^{(Q',t)}$  (as defined by equation 4.15) is NP-complete due to the constrained forms of the dependencies.

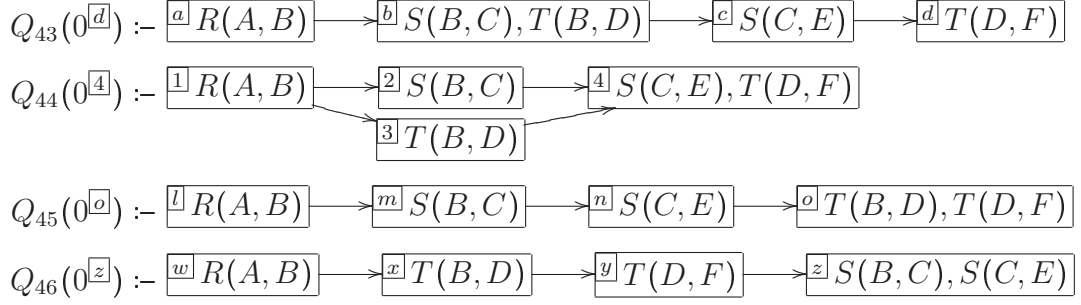


Figure 4.10: Four CH-DACQs with dual-node head graphs

### Tuple-Specific Equivalence: Dual-Node Head Graphs

Let  $Q$  denote a CH-DACQ in canonical form whose head graph contains a single edge from parent  $\boxed{u}$  to child  $\boxed{v}$ . This implies  $\text{dags}(G_Q^{\text{head}}) = \{G_{uv}, G_u\}$ , and so  $\text{CT}_Q = \{t_1, t_2\}$  for some tuple  $t_1$  whose corresponding shape  $\mathfrak{ts}_{t_1}^Q \in \text{TS}_Q$  is terminal, and a second tuple  $t_2$  that is subsumed by  $t_1$  (and whose corresponding shape  $\mathfrak{ts}_{t_2}^Q \in \text{TS}_Q$  does not contain  $\boxed{v}$ ). Because  $Q$  is in canonical form, property **[P2]** of Definition 4.2.1 implies that  $\text{body}_Q$  is a rooted DAG with root  $u$  and the single leaf  $v$ , but with an arbitrary number of intermediate vertices. Given a second query  $Q'$  with  $\text{CT}_{Q'} = \text{CT}_Q$ , we can characterize the relationships  $\stackrel{t_1}{\equiv}$  and  $\stackrel{t_2}{\equiv}$  by reducing them to deciding finite equivalence between sets of dependencies.

**Example 44** Figure 4.10 shows four queries that each have a dual-node head graph and canonical tuples  $\{\langle 0 \rangle, \langle \blacksquare \rangle\}$ . All four queries are in DACQ canonical form.

Testing equivalence for tuple  $t_1$  is analogous to the case of single-node head graphs. Corresponding to node  $v$  of query  $Q$  and tuple  $t_1$  we define the dependency  $\xi^{(v, t_1)}$  and the dependency set  $\Sigma^{(Q, t_1)}$  as follows.

$$\xi^{(v, t_1)} := \text{body}_v^+ \rightarrow \perp \quad (4.16)$$

$$\Sigma^{(Q, t_1)} := \{\xi^{(v, t_1)}\} \quad (4.17)$$

As before, deciding finite equivalence between  $\Sigma^{(Q, t_1)}$  and  $\Sigma^{(Q', t_1)}$  is NP-complete.

**Lemma 4.3.17** *For any database  $\mathbb{D}$ ,  $t_1 \in (Q)^{\mathbb{D}}$  iff  $\mathbb{D} \neq \Sigma^{(Q, t_1)}$ .*

**Corollary 4.3.18**  $Q \stackrel{t_1}{\equiv} Q'$  iff  $\Sigma^{(Q, t_1)} \equiv_{\text{fin}} \Sigma^{(Q', t_1)}$ .

**Example 45** Continuing Example 44, for canonical tuple  $\langle 0 \rangle$  all four queries in Figure 4.10 yield the same  $\langle 0 \rangle$ -dependency set

$$\begin{aligned}
\Sigma^{(Q_{43}, \langle 0 \rangle)} &= \Sigma^{(Q_{44}, \langle 0 \rangle)} = \Sigma^{(Q_{45}, \langle 0 \rangle)} = \Sigma^{(Q_{46}, \langle 0 \rangle)} \\
&= \{[R(A, B), S(B, C), T(B, D), S(C, E), T(D, F)] \rightarrow \perp\}
\end{aligned}$$

and so  $Q_{43} \stackrel{(0)}{\equiv} Q_{44} \stackrel{(0)}{\equiv} Q_{45} \stackrel{(0)}{\equiv} Q_{46}$  follows trivially.

Deciding equivalence for tuple  $t_2$  is more challenging, because the definition of DACQ evaluation semantics in terms of terminal embeddings makes query  $Q$  non-monotonic with respect to tuple  $t_2$ . As a result, the dependencies we derive no longer have the constrained form of equations 4.14 and 4.16, but instead have conclusions with non-empty  $\text{CQ}^-$  bodies. In particular, it is the existentially-quantified variables within the conclusions that are problematic. Implication of arbitrary sets of embedded dependencies is a widely-studied problem, and in Section 4.6.4 we survey some of the literature pertaining to it. For now, the salient features are that although it is only semi-decidable in general, there are useful sub-classes of the problem which allow decidability—in particular, classes for which the well-studied *chase* algorithm is guaranteed to terminate [6, 30].

To decide  $\Sigma \models_{\text{fin}} \xi$  using the chase we construct a new dependency  $\hat{\xi}$  by isomorphically “freezing” both the free variables and the non-free premise variables of  $\xi$  into fresh constants. Because premise  $P^{\hat{\xi}}$  contains only constants, it can be treated as a database instance to be chased. Let database  $\mathbb{D} := \text{chase}_{\Sigma}(P^{\hat{\xi}})$  denote the (possibly infinite) database resulting from any sequence of chase steps applying dependencies from  $\Sigma$  to the initial database  $P^{\hat{\xi}}$ . If  $\mathbb{D}$  is finite then  $\mathbb{D} \models \Sigma$ ; if  $\mathbb{D}$  is infinite then the truth of  $\mathbb{D} \models \Sigma$  is independent of the order of chase steps<sup>3</sup> [6, Ch. 10]. In either case, the truth of  $\mathbb{D} \models \hat{\xi}$  is independent of the order of chase steps.<sup>3</sup>

The chase procedure can return the conclusion  $\Sigma \models_{\text{fin}} \xi$  in finite time by demonstrating any finite database  $\hat{\mathbb{D}} \subseteq \mathbb{D}$  corresponding to a prefix of the chase sequence yielding  $\mathbb{D}$ , such that either

1. there exists an embedding  $\gamma : C^{\hat{\xi}} \rightarrow \hat{\mathbb{D}}$ , or
2.  $\hat{\mathbb{D}}$  contains a contradiction.

This proves  $\hat{\mathbb{D}} \models \hat{\xi}$  and hence  $\mathbb{D} \models \hat{\xi}$  (because premise  $P^{\hat{\xi}}$  does not contain variables) which—by the isomorphic freezing—implies  $\Sigma \models_{\text{unr}} \xi$  and hence  $\Sigma \models_{\text{fin}} \xi$ .

The chase procedure only returns the conclusion  $\Sigma \not\models_{\text{fin}} \xi$  when  $\mathbb{D}$  is finite and  $\mathbb{D} \not\models \hat{\xi}$  (i.e.,  $\mathbb{D}$  does not contain a contradiction and there does not exist an embedding of  $C^{\hat{\xi}}$  into  $\mathbb{D}$ ). Hence,  $\Sigma \models_{\text{fin}} \xi$  can only be decided by the chase algorithm in cases where there exists a finite chase sequence for  $\text{chase}_{\Sigma}(P^{\hat{\xi}})$ —a condition which is itself undecidable to identify [30]. In Section 4.5 we discuss how to deal with cases where chase termination cannot be guaranteed.

Recall from equation 4.5 that  $\text{dags}(\text{body}_Q, \mathfrak{t}_{t_2}^Q)$  contains all of the DAGs in  $\text{dags}(\text{body}_Q)$  containing node  $u$  but not node  $v$ . For each  $G \in \text{dags}(\text{body}_Q, \mathfrak{t}_{t_2}^Q)$  we define the following embedded dependency  $\xi^{(G, t_2)}$ , and we collect all of these dependencies into the dependency set  $\Sigma^{(Q, t_2)}$ .

$$\xi^{(G, t_2)} := \text{body}_G^+ \rightarrow \text{body}_v^+ \quad (4.18)$$

$$\Sigma^{(Q, t_2)} := \{\xi^{(G, t_2)} \mid G \in \text{dags}(\text{body}_Q, \mathfrak{t}_{t_2}^Q)\} \quad (4.19)$$

---

<sup>3</sup>As long as the infinite chase sequence obeys certain basic properties—primarily, that no dependency is “starved.”



(When constructing  $\xi^{(G,t_2)}$  from equation 4.18, any trivially-implied atoms are deleted from the conclusion.)

**Example 46** Continuing Examples 44 and 45, for canonical tuple  $\langle \blacksquare \rangle$  the queries in Figure 4.10 yield the following  $\langle \blacksquare \rangle$ -dependency sets.

$$\begin{aligned}
\Sigma^{(Q_{43}, \langle \blacksquare \rangle)} &= \{ \xi^{(a, \langle \blacksquare \rangle)} : [R(A, B)] \rightarrow [S(B, C), T(B, D), S(C, E), T(D, F)], \\
&\quad \xi^{(ab, \langle \blacksquare \rangle)} : [R(A, B), S(B, C), T(B, D)] \rightarrow [S(C, E), T(D, F)], \\
&\quad \xi^{(abc, \langle \blacksquare \rangle)} : [R(A, B), S(B, C), T(B, D), S(C, E)] \rightarrow [T(D, F)] \} \\
\Sigma^{(Q_{44}, \langle \blacksquare \rangle)} &= \{ \xi^{(1, \langle \blacksquare \rangle)} : [R(A, B)] \rightarrow [S(B, C), T(B, D), S(C, E), T(D, F)], \\
&\quad \xi^{(12, \langle \blacksquare \rangle)} : [R(A, B), S(B, C)] \rightarrow [T(B, D), S(C, E), T(D, F)], \\
&\quad \xi^{(13, \langle \blacksquare \rangle)} : [R(A, B), T(B, D)] \rightarrow [S(B, C), S(C, E), T(D, F)], \\
&\quad \xi^{(123, \langle \blacksquare \rangle)} : [R(A, B), S(B, C), T(B, D)] \rightarrow [S(C, E), T(D, F)] \} \\
\Sigma^{(Q_{45}, \langle \blacksquare \rangle)} &= \{ \xi^{(l, \langle \blacksquare \rangle)} : [R(A, B)] \rightarrow [S(B, C), S(C, E), T(B, D), T(D, F)], \\
&\quad \xi^{(lm, \langle \blacksquare \rangle)} : [R(A, B), S(B, C)] \rightarrow [S(C, E), T(B, D), T(D, F)], \\
&\quad \xi^{(lmn, \langle \blacksquare \rangle)} : [R(A, B), S(B, C), S(C, E)] \rightarrow [T(B, D), T(D, F)] \} \\
\Sigma^{(Q_{46}, \langle \blacksquare \rangle)} &= \{ \xi^{(w, \langle \blacksquare \rangle)} : [R(A, B)] \rightarrow [T(B, D), T(D, F), S(B, C), S(C, E)], \\
&\quad \xi^{(wx, \langle \blacksquare \rangle)} : [R(A, B), T(B, D)] \rightarrow [T(D, F), S(B, C), S(C, E)], \\
&\quad \xi^{(wxy, \langle \blacksquare \rangle)} : [R(A, B), T(B, D), T(D, F)] \rightarrow [S(B, C), S(C, E)] \}
\end{aligned}$$

**Lemma 4.3.19** For any database  $\mathbb{D}$ ,  $t_2 \in (Q)^{\mathbb{D}}$  iff  $\mathbb{D} \neq \Sigma^{(Q, t_2)}$ .

PROOF. Suppose that  $t_2 \in (Q)^{\mathbb{D}}$ . Then, there exists some terminal embedding  $\gamma : G \rightarrow \mathbb{D}$  of some DAG  $G \in \text{dags}(\text{body}_Q, \mathfrak{ts}_{t_2}^Q)$  which doesn't contain  $v$ ; however,  $\gamma$  is also an embedding of  $P^{\xi^{(G, t_2)}}$  that does not extend to  $C^{\xi^{(G, t_2)}}$  and so  $\mathbb{D} \neq \xi^{(G, t_2)}$ . Conversely, let  $G$  be a maximal DAG in  $\text{dags}(\text{body}_Q, \mathfrak{ts}_{t_2}^Q)$  such that  $\mathbb{D} \neq \xi^{(G, t_2)}$ . Then, there exists embedding  $\gamma : P^{\xi^{(G, t_2)}} \rightarrow \mathbb{D}$  that does not extend to  $C^{\xi^{(G, t_2)}}$ ; however,  $\gamma$  is a terminal embedding of  $G$  (by choice of  $G$ ) and so generates tuple  $t_2 \in (Q)^{\mathbb{D}}$ .  $\square$

**Corollary 4.3.20**  $Q \stackrel{t_2}{\equiv} Q'$  iff  $\Sigma^{(Q, t_2)} \equiv_{\text{fin}} \Sigma^{(Q', t_2)}$ .

**Example 47** Continuing Examples 44–46, we now prove that  $Q_{43} \equiv Q_{44}$ , but for all other pairs of queries, equivalence does not hold.

$$\begin{aligned}
Q_{43} \stackrel{\langle \blacksquare \rangle}{\equiv} Q_{44}: \quad &\Sigma^{(Q_{43}, \langle \blacksquare \rangle)} \models_{\text{fin}} \xi^{(1, \langle \blacksquare \rangle)}: \text{Trivial} \\
&\Sigma^{(Q_{43}, \langle \blacksquare \rangle)} \models_{\text{fin}} \xi^{(12, \langle \blacksquare \rangle)}: \\
&\quad \mathbb{D} := \text{chase}_{\Sigma^{(Q_{43}, \langle \blacksquare \rangle)}}([R(a, b), S(b, c)]) \\
&\quad = \{R(a, b), S(b, c), S(b, c_1), T(b, d_1), S(c_1, e_1), T(d_1, f_1), S(c, e_2)\} \\
&\quad \models [T(b, D), S(c, E), T(D, F)] \\
&\Sigma^{(Q_{43}, \langle \blacksquare \rangle)} \models_{\text{fin}} \xi^{(13, \langle \blacksquare \rangle)}: \\
&\quad \mathbb{D} := \text{chase}_{\Sigma^{(Q_{43}, \langle \blacksquare \rangle)}}([R(a, b), T(b, d)]) \\
&\quad = \{R(a, b), T(b, d), S(b, c_1), T(b, d_1), S(c_1, e_1), T(d_1, f_1), T(d, f_2)\} \\
&\quad \models [S(b, C), S(C, E), T(d, F)] \\
&\Sigma^{(Q_{43}, \langle \blacksquare \rangle)} \models_{\text{fin}} \xi^{(123, \langle \blacksquare \rangle)}: \text{Trivial}
\end{aligned}$$



$$\begin{aligned}
\Sigma(Q_{44}, \langle \blacksquare \rangle) &\models_{\text{fin}} \xi(a, \langle \blacksquare \rangle): \text{ Trivial} \\
\Sigma(Q_{44}, \langle \blacksquare \rangle) &\models_{\text{fin}} \xi(ab, \langle \blacksquare \rangle): \text{ Trivial} \\
\Sigma(Q_{44}, \langle \blacksquare \rangle) &\models_{\text{fin}} \xi(abc, \langle \blacksquare \rangle): \\
\mathbb{D} &:= \text{chase}_{\Sigma(Q_{44}, \langle \blacksquare \rangle)} ([R(a, b), S(b, c), T(b, d), S(c, e)]) \\
&= \{R(a, b), S(b, c), T(b, d), S(c, e), S(c, e_1), T(d, f_1)\} \\
&\models [T(d, F)]
\end{aligned}$$

$$\begin{aligned}
Q_{43} \stackrel{\langle \blacksquare \rangle}{\not\equiv} Q_{45}: \quad &\Sigma(Q_{45}, \langle \blacksquare \rangle) \not\models_{\text{fin}} \xi(ab, \langle \blacksquare \rangle): \\
\mathbb{D} &:= \text{chase}_{\Sigma(Q_{45}, \langle \blacksquare \rangle)} ([R(a, b), S(b, c), T(b, d)]) \\
&= \{R(a, b), S(b, c), T(b, d), S(c, e_1), T(b, d_1), T(d_1, f_1)\} \\
&\not\models [S(c, E), T(d, F)]
\end{aligned}$$

$$\begin{aligned}
Q_{43} \stackrel{\langle \blacksquare \rangle}{\not\equiv} Q_{46}: \quad &\Sigma(Q_{46}, \langle \blacksquare \rangle) \not\models_{\text{fin}} \xi(ab, \langle \blacksquare \rangle): \\
\mathbb{D} &:= \text{chase}_{\Sigma(Q_{46}, \langle \blacksquare \rangle)} ([R(a, b), S(b, c), T(b, d)]) \\
&= \{R(a, b), S(b, c), T(b, d), T(d, f_1), S(b, c_1)S(c_1, e_1)\} \\
&\not\models [S(c, E), T(d, F)]
\end{aligned}$$

$$\begin{aligned}
Q_{45} \stackrel{\langle \blacksquare \rangle}{\not\equiv} Q_{46}: \quad &\Sigma(Q_{46}, \langle \blacksquare \rangle) \not\models_{\text{fin}} \xi(lm, \langle \blacksquare \rangle): \\
\mathbb{D} &:= \text{chase}_{\Sigma(Q_{46}, \langle \blacksquare \rangle)} ([R(a, b), S(b, c)]) \\
&= \{R(a, b), S(b, c), T(b, d_1), T(d_1, f_1), S(b, c_1)S(c_1, e_1)\} \\
&\not\models [S(c, E), T(b, D), T(D, F)]
\end{aligned}$$

The equivalence  $Q_{43} \equiv Q_{44}$  also follows from Theorem 4.3.9. Queries  $Q_{43}$  and  $Q_{44}$  have parameter sets  $\mathcal{P}_{Q_{43}}^+ = \mathcal{P}_{Q_{44}}^+ = \{B, C, D\}$  with bijection  $\nu$  being the obvious identity mapping. The obvious identity mapping also forms the homomorphism between the shape queries for shapes  $\mathfrak{ts}_{\langle \blacksquare \rangle}^{Q_{43}} = \emptyset = \mathfrak{ts}_{\langle \blacksquare \rangle}^{Q_{44}}$

$$\begin{aligned}
\hat{Q}_{43}|_{\emptyset}(\blacksquare) &:- R(A, B) \\
\hat{Q}_{44}|_{\emptyset}(\blacksquare) &:- R(A, B)
\end{aligned}$$

and shape queries for shapes  $\mathfrak{ts}_{(0)}^{Q_{43}} = \{\underline{d}\} \in \mathfrak{TS}_{Q_{43}}$  and  $\mathfrak{ts}_{(0)}^{Q_{44}} = \{\underline{4}\} \in \mathfrak{TS}_{Q_{46}}$ .

$$\begin{aligned}
\hat{Q}_{43}|_d(0) &:- R(A, B), S(B, C), T(B, D), S(C, E), T(D, F) \\
\hat{Q}_{44}|_4(0) &:- R(A, B), S(B, C), T(B, D), S(C, E), T(D, F)
\end{aligned}$$

In contrast, queries  $Q_{45}$  and  $Q_{46}$  have parameter sets  $\mathcal{P}_{Q_{45}}^+ = \{B, C\}$  and  $\mathcal{P}_{Q_{46}}^+ = \{B, D\}$ , respectively, and so Theorem 4.3.9 does not apply to any other pairings of the four queries. (Even though  $\mathcal{P}_{Q_{45}}^+$  and  $\mathcal{P}_{Q_{46}}^+$  are bijective, no homomorphisms exist between  $\hat{Q}_{45}|_o$  and  $\hat{Q}_{46}|_z$  that are consistent with either possible bijection.)

## Tuple-Specific Equivalence: Arbitrary Head Graphs

We now consider the problem of deciding tuple-specific equivalence for constant-headed DACQs with arbitrary head graphs. Our approach is a straightforward

extension of the technique we used for queries with single-edge head graphs; the additional complication is that we require reasoning about implication of *disjunctive* embedded dependencies. A disjunctive embedded dependency  $\xi$  can be expressed as a rule of the form

$$P^\xi \rightarrow C_1^\xi \mid \dots \mid C_n^\xi \quad (4.20)$$

where the premise  $P^\xi$  and each conclusion  $C_i^\xi$  are CQ<sup>=</sup> bodies;  $\xi$  is satisfied by database instance  $\mathbb{D}$  if each embedding of the premise into  $\mathbb{D}$  can be extended to an embedding of at least one of the conclusions. Deciding implication between sets of disjunctive embedded dependencies is discussed in Section 4.6.4.

Given any constant-headed DACQ  $Q$  and any canonical tuple  $t \in \text{CT}_Q$  corresponding to shape  $\mathfrak{ts}_t^Q \in \text{TS}_Q$ , let  $\{\boxed{v_1}, \dots, \boxed{v_n}\} \subseteq L$  be the set of vertices in head graph  $G_Q^{\text{head}}$  that do not occur in shape  $\mathfrak{ts}_t^Q$  but are children of vertices in  $\mathfrak{ts}_t^Q$ . For each  $G \in \text{dags}(\text{body}_Q, \mathfrak{ts}_t^Q)$  we define a single embedded disjunctive dependency  $\xi^{(G,t)}$

$$\xi^{(G,t)} := \begin{cases} \text{body}_G^+ \rightarrow \perp & \text{if } n = 0 \\ \text{body}_G^+ \rightarrow \text{body}_{v_1}^+ \mid \dots \mid \text{body}_{v_n}^+ & \text{otherwise} \end{cases} \quad (4.21)$$

and we define dependency set  $\Sigma^{(Q,t)}$  as follows.

$$\Sigma^{(Q,t)} := \{\xi^{(G,t)} \mid G \in \text{dags}(\text{body}_Q, \mathfrak{ts}_t^Q)\} \quad (4.22)$$

Note that  $n = 0$  if and only if  $\mathfrak{ts}_t^Q$  is terminal if and only if  $\text{dags}(\text{body}_Q, \mathfrak{ts}_t^Q) = \{\text{body}_Q\}$  (by property [P2] of Definition 4.2.1).

**Lemma 4.3.21** *For any database  $\mathbb{D}$  and canonical tuple  $t, t \in (Q)^\mathbb{D}$  iff  $\mathbb{D} \not\models \Sigma^{(Q,t)}$ .*

PROOF. Suppose that  $t \in (Q)^\mathbb{D}$ . Then, there exists some terminal embedding  $\gamma : G \rightsquigarrow \mathbb{D}$  of some DAG  $G \in \text{dags}(\text{body}_Q, \mathfrak{ts}_t^Q)$ . However,  $\gamma$  is also an embedding of premise query  $P^{\xi^{(G,t)}}$ . If  $\mathfrak{ts}_t^Q$  is not the terminal shape then  $\gamma$  does not extend to any conclusion query  $C_i^{\xi^{(G,t)}}$ , whereas if  $\mathfrak{ts}_t^Q$  is terminal then the conclusion of  $\xi^{(G,t)}$  is a contradiction to which no embedding can extend; in either case,  $\mathbb{D} \not\models \xi^{(G,t)}$ .

Conversely, let  $G$  be a maximal DAG in  $\text{dags}(\text{body}_Q, \mathfrak{ts}_t^Q)$  such that  $\mathbb{D} \not\models \xi^{(G,t)}$ . There exists embedding  $\gamma : P^{\xi^{(G,t)}} \rightarrow \mathbb{D}$  that does not extend to any conclusion of  $\xi^{(G,t)}$ ; however,  $\gamma$  is a terminal embedding of  $G$  that generates tuple  $t \in (Q)^\mathbb{D}$ .  $\square$

**Corollary 4.3.22** *For any canonical tuple  $t \in \text{CT}_Q$ ,  $Q \stackrel{t}{\equiv} Q'$  iff  $\Sigma^{(Q,t)} \equiv_{\text{fin}} \Sigma^{(Q',t)}$ .*

## 4.4 Queries with Variable Heads

In this section consider the equivalence problem for arbitrary DACQs. Given a DACQ  $Q$  whose head tuple  $\overline{W}$  may contain variables, we start by observing that

the definitions of head graphs, tuple shapes, and shape queries given in Section 4.3.1 are not specific to CH-DACQs. As a result, both necessary conditions for query equivalence that we derived in that section—*isomorphism between head graphs* (Theorem 4.3.4) and *equivalence of shape queries* (Corollary 4.3.8)—apply to arbitrary DACQs.

The definition of canonical tuples given in Section 4.3.1 is specific to CH-DACQs. In Section 4.4.1 we extend the definition of canonical tuples to handle queries with variables in the head tuple, using a standard approach of “freezing” the head variables into constants [6, 34]. This allows us to reduce DACQ equivalence to tuple-specific equivalence between CH-DACQs, analogous to the similar reduction for CH-DACQs (Theorem 4.3.14). In the worst case our decision algorithm enumerates an exponential number of canonical freezings per tuple shape; however, in Section 4.4.2 we discuss an optimization that allows us to avoid considering certain freezings. For DACQs that can be expressed without explicit equality—including all CQs and most practical queries expressible in SPCL—this reduces the problem to considering a single canonical freezing per tuple shape.

#### 4.4.1 Canonical Freezings and Canonical Tuples

The canonical freezings of a query are defined in terms of a set of *symbolic constants*  $\mathcal{C}_S$ . Choose  $\mathcal{C}_S$  as any strict subset of  $\mathcal{C}$  such that both  $\mathcal{C}_S$  and  $\mathcal{C} \setminus \mathcal{C}_S$  are infinite. Without loss of generality, we will assume that  $\mathcal{C}_S$  is chosen so that for any given set of queries, all constants appearing in the queries are from the set  $\mathcal{C} \setminus \mathcal{C}_S$ . Because the constants in  $\mathcal{C}_S$  do not appear in query definitions, their precise identity is unimportant (hence the name “symbolic”).

**Proposition 4.4.1** *Given any two DACQs  $Q$  and  $Q'$  (satisfying  $(\mathcal{C}_Q \cup \mathcal{C}_{Q'}) \cap \mathcal{C}_S = \emptyset$ ), if  $Q \not\equiv Q'$  then there exists a database instance  $\mathbb{D}$  with  $\text{adom}(\mathbb{D}) \subset \mathcal{C}_S \cup \mathcal{C}_Q \cup \mathcal{C}_{Q'}$  such that  $(Q)^\mathbb{D} \neq (Q')^\mathbb{D}$ .*

Given any two tuples  $t$  and  $t'$  over domain  $\mathcal{C} \cup \{\blacksquare\}$ , we say that  $t$  is  $\mathcal{C}_S$ -*homomorphic* to  $t'$ —denoted  $t \succsim t'$ —if there exists a mapping  $h : \mathcal{C}_S \rightarrow \mathcal{C}$ , extended with identity over  $(\mathcal{C} \setminus \mathcal{C}_S) \cup \{\blacksquare\}$ , such that  $h(t) = t'$ . We say that  $t$  is  $\mathcal{C}_S$ -*isomorphic* to  $t'$ —denoted  $t \approx t'$ —if  $t \succsim t'$  and  $t' \succsim t$ . We say that  $t$  is *more general* than  $t'$ —denoted  $t > t'$ —if  $t \succsim t'$  and  $t' \not\approx t$ .

Given two sets of tuples  $S_1, S_2$  over  $\mathcal{C}$ , we say that  $S_1$  is  $\mathcal{C}_S$ -*contained* in  $S_2$ —denoted  $S_1 \preceq S_2$ —if  $\forall t_1 \in S_1 \exists t_2 \in S_2. (t_1 \approx t_2)$ . We say that  $S_1$  is  $\mathcal{C}_S$ -*equivalent* to  $S_2$ —denoted  $S_1 \approx S_2$ —if  $S_1 \preceq S_2$  and  $S_2 \preceq S_1$ .

In order to define the canonical tuples for a DACQ, we first define the *canonical freezings* of a satisfiable CQ $^\bar{\cdot}$ . Any tuple that can be output by the CQ $^\bar{\cdot}$  can be obtained by an isomorphic substitution for the symbolic constants within one of the canonical freezings.

**Definition 4.4.2 (Canonical  $\text{CQ}^\pm$  Freezings)** *Given any satisfiable  $\text{CQ}^\pm Q(\overline{\mathcal{W}})$  with  $\mathcal{C}_B \cap \mathcal{C}_S = \emptyset$ , let  $\mathcal{H}_Q$  denote the set of all homomorphisms  $h : \mathcal{W} \cap \mathcal{B} \rightarrow \mathcal{C}_S \cup \mathcal{C}_Q$  that have been extended with identity over all symbols not in  $\mathcal{W} \cap \mathcal{B}$ . The (infinite) set of valid freezings of  $Q$ —denoted  $\text{CF}_Q^+$ —contains the homomorphic images of head tuple  $\overline{\mathcal{W}}$  under mappings from  $\mathcal{H}_Q$  that are consistent with the equality predicates in the query body.*

$$\text{CF}_Q^+ := \{h(\overline{\mathcal{W}}) \mid h \in \mathcal{H}_Q \wedge h(\text{pred}_Q) \neq \perp\} \quad (4.23)$$

*The set of canonical freezings of  $Q$ —denoted  $\text{CF}_Q$ —is a minimal subset of  $\text{CF}_Q^+$  satisfying  $\text{CF}_Q \approx \text{CF}_Q^+$  (if many such sets exist, we arbitrarily choose one).*

Set  $\text{CF}_Q$  is finite, with size bounded by a function exponential in both  $|\mathcal{C}_Q|$  and  $|\mathcal{W} \cap \mathcal{B}|$ .

**Lemma 4.4.3** *Given  $\text{CQ}^\pm Q$  and any database instance  $\mathbb{D}$  over  $\mathcal{C}_S \cup \mathcal{C}_Q$ ,  $\mathbb{T}((Q)^\mathbb{D}) \approx \text{CF}_Q$ .*

PROOF. Any embedding  $\gamma : Q \rightarrow \mathbb{D}$  yields  $\gamma(\overline{\mathcal{W}}) \in \text{CF}_Q^+$  (by Definition 4.4.2), and so  $\mathbb{T}((Q)^\mathbb{D}) \subseteq \text{CF}_Q^+ \approx \text{CF}_Q$ .  $\square$

**Definition 4.4.4 (Canonical Tuples)** *Given DACQ  $Q$  in canonical form, the set  $\text{CT}_Q$  of canonical tuples for  $Q$  corresponds to the canonical freezings of all satisfiable shape queries.*

$$\text{CT}_Q := \{t \mid \exists \mathfrak{ts} \in \text{TS}_Q \wedge \hat{Q}|_{\mathfrak{ts}} \text{ is satisfiable} \wedge t \in \text{CF}_{\hat{Q}|_{\mathfrak{ts}}}\}$$

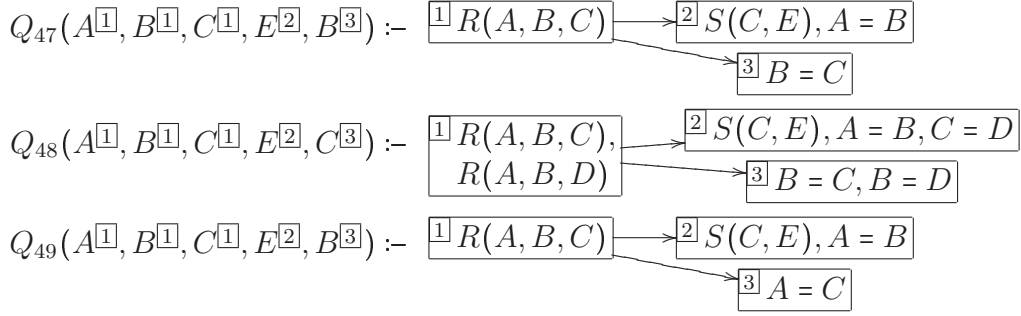
*As defined previously for CH-DACQs, we use  $\mathfrak{ts}_t^Q$  to denote the tuple shape from  $\text{TS}_Q$  corresponding to a given canonical tuple  $t \in \text{CT}_Q$ .*

**Example 48** Consider the three DACQs  $Q_{47}$ ,  $Q_{48}$ , and  $Q_{49}$  shown in Figure 4.11a. Observe that all three queries have the following set of tuple shapes,

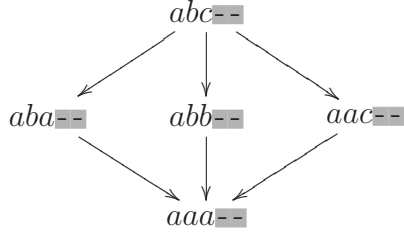
$$\text{TS}_{Q_{47}} = \text{TS}_{Q_{48}} = \text{TS}_{Q_{49}} = \{\mathfrak{ts}_1, \mathfrak{ts}_{12}, \mathfrak{ts}_{13}, \mathfrak{ts}_{123}\}$$

where the subscripts denote the set of labels in each shape. Figures 4.11b–4.11f show the canonical tuples for  $Q_{47}$ ,  $Q_{48}$ , and  $Q_{49}$ , partitioned into sets of canonical freezings for the different tuples shapes. For each tuple shape we have organized the canonical freezings into a lattice in which ancestors are more general than their descendants. Observe that  $\text{CT}_{Q_{47}} = \text{CT}_{Q_{48}}$  even though the head tuples of  $Q_{47}$  and  $Q_{48}$  are not isomorphic. In contrast,  $\text{CT}_{Q_{47}} \neq \text{CT}_{Q_{49}}$  (due to shape  $\mathfrak{ts}_{13}$ ) even though the head tuples of  $Q_{47}$  and  $Q_{49}$  are isomorphic.

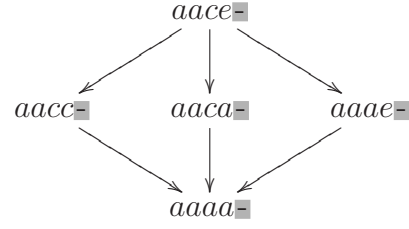
**Lemma 4.4.5** *Given DACQ  $Q$  in canonical form and any database instance  $\mathbb{D}$  over  $\mathcal{C}_S \cup \mathcal{C}_Q$ ,  $\mathbb{T}((Q)^\mathbb{D}) \approx \text{CT}_Q$ .*



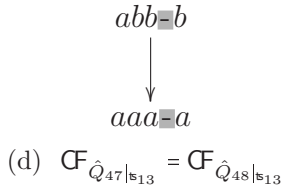
(a) DACQs  $Q_{47}$ ,  $Q_{48}$ , and  $Q_{49}$



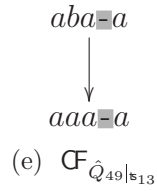
(b)  $\mathbf{CF}_{\hat{Q}_{47}|s_1} = \mathbf{CF}_{\hat{Q}_{48}|s_1} = \mathbf{CF}_{\hat{Q}_{49}|s_1}$



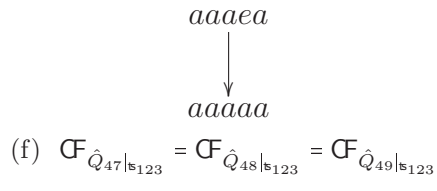
(c)  $\mathbf{CF}_{\hat{Q}_{47}|s_{12}} = \mathbf{CF}_{\hat{Q}_{48}|s_{12}} = \mathbf{CF}_{\hat{Q}_{49}|s_{12}}$



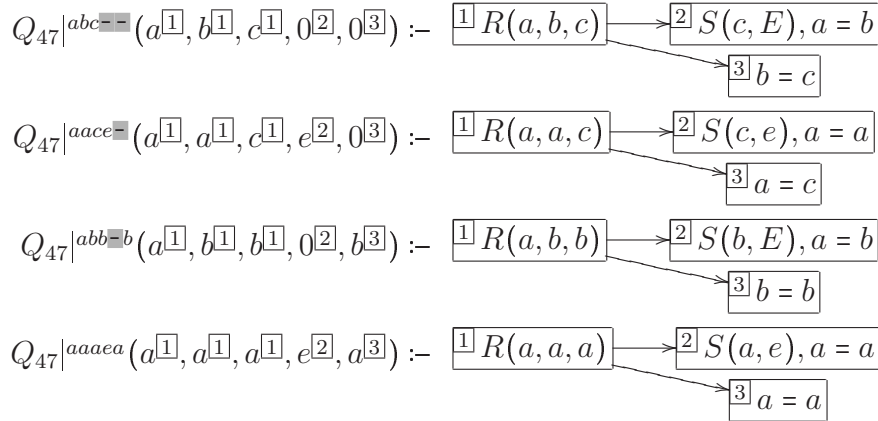
(d)  $\mathbf{CF}_{\hat{Q}_{47}|s_{13}} = \mathbf{CF}_{\hat{Q}_{48}|s_{13}}$



(e)  $\mathbf{CF}_{\hat{Q}_{49}|s_{13}}$



(f)  $\mathbf{CF}_{\hat{Q}_{47}|s_{123}} = \mathbf{CF}_{\hat{Q}_{48}|s_{123}} = \mathbf{CF}_{\hat{Q}_{49}|s_{123}}$



(g) Frozen DACQs obtained from  $Q_{47}$

Figure 4.11: Three DACQs with their canonical tuples

PROOF. Follows from Definition 4.4.4 and Lemma 4.4.3.  $\square$

**Lemma 4.4.6** *Given two DACQs  $Q$  and  $Q'$  in canonical form, if  $Q \equiv Q'$  then  $\text{CT}_Q \approx \text{CT}_{Q'}$ .*

PROOF. Follows from Definition 4.4.4 because  $\theta(\text{TS}_Q) = \text{TS}_{Q'}$  (Theorem 4.3.4), and  $\hat{Q}|_{\mathfrak{s}} \equiv \hat{Q}'|_{\theta(\mathfrak{s})}$  (Corollary 4.3.8) implies  $\text{CF}_{\hat{Q}|_{\mathfrak{s}}} \approx \text{CF}_{\hat{Q}'|_{\theta(\mathfrak{s})}}$  (by Definition 4.4.2).  $\square$

**Corollary 4.4.7** *If  $Q \equiv Q'$  then  $\text{CT}_Q$  and  $\text{CT}_{Q'}$  can be chosen so that  $\text{CT}_Q = \text{CT}_{Q'}$ .*

**Definition 4.4.8 (Frozen DACQ)** *Given DACQ  $Q(\overline{W})$  in canonical form and any canonical tuple  $t \in \text{CT}_Q$  with corresponding shape  $\mathfrak{ts}_t^Q \in \text{TS}_Q$ ,*

- recall that  $t$  is a canonical freezing of shape query  $\hat{Q}|_{\mathfrak{ts}_t^Q}$  (cf. Definition 4.4.4),
- recall that  $\hat{Q}|_{\mathfrak{ts}_t^Q}$  is a  $CQ^\exists$  with head tuple  $\overline{W}|_{\mathfrak{ts}_t^Q}$  (cf. Definition 4.3.5), and
- recall that there exists some homomorphism  $h_t : \mathcal{W}|_{\mathfrak{ts}_t^Q} \cap \mathcal{B} \rightarrow \mathcal{C}_S \cup \mathcal{C}_Q$  extended with identity over all symbols not in  $\mathcal{W}|_{\mathfrak{ts}_t^Q} \cap \mathcal{B}$  such that  $t = h_t(\overline{W}|_{\mathfrak{ts}_t^Q})$  (cf. Definition 4.4.2).

We define the freezing of  $Q$  with  $t$ —denoted  $Q|_t$ —as the following CH-DACQ.

- $\text{head}_{Q|_t} := \overline{W}|_t$ , where  $\overline{W}|_t$  is constructed from  $\overline{W}$  by locating each component  $W_i^{\boxed{v}}$  in which  $W_i$  is a variable, and if  $v \in \mathfrak{ts}_t^Q$  replacing  $W_i$  with  $h_t(W_i)$ , whereas if  $v \notin \mathfrak{ts}_t^Q$  replacing  $W_i$  with an arbitrary constant from  $\mathcal{C} \setminus \mathcal{C}_S$  (for simplicity, we always use the constant 0).
- $\text{body}_{Q|_t} := h_t(\text{body}_Q)$

**Example 49** Continuing Example 48, Figure 4.11g shows four frozen DACQs obtained from  $Q_{47}$  by freezing with the most general canonical freezing for each of the four tuple shapes. Note that all of the DACQs in Figure 4.11g are constant-headed, but none of them are in canonical form.

**Lemma 4.4.9** *Given any DACQ  $Q$  in canonical form,  $\forall t \in \text{CT}_Q. [Q \stackrel{t}{\equiv} Q|_t]$ .*

PROOF. Consider any database  $\mathbb{D}$ . Suppose that  $t \in (Q)^{\mathbb{D}}$ ; then there exists some terminal embedding  $\gamma : Q \rightsquigarrow \mathbb{D}$  satisfying  $\hat{\gamma}(\overline{W}) = t$ . This implies that there exists some  $G \in \text{dags}(\text{body}_Q, \mathfrak{ts}_t^Q)$  such that  $\gamma$  is a terminal embedding of  $G$  (see equation 4.5). By Theorem 4.3.6,  $\gamma$  is also an embedding of  $\hat{Q}|_{\mathfrak{ts}_t^Q}$  into  $\mathbb{D}$  satisfying  $\hat{\gamma}(\overline{W}|_{\mathfrak{ts}_t^Q}) = t$ , and so, by definition of  $\overline{w}|_t$ ,  $\hat{\gamma}(\overline{W}|_t) = \hat{\gamma}(\overline{W}|_{\mathfrak{ts}_t^Q}) = t = h_t(\overline{W}|_{\mathfrak{ts}_t^Q})$ . Because mappings  $\gamma$  and  $h_t$  agree over  $\overline{W}|_{\mathfrak{ts}_t^Q}$  and  $h_t$  is the identity mapping over all other symbols, we can decompose embedding  $\gamma$  as  $\gamma = \phi \circ h_t$  for some mapping  $\phi : \mathcal{B} \rightarrow \text{adom}(\mathbb{D})$  extended with identity over constants  $\mathcal{C}$ . Because  $Q|_t$  has body  $h_t(\text{body}_Q)$ , mapping  $\phi$  is an embedding of  $h_t(G)$  into  $\mathbb{D}$  satisfying  $\phi(\overline{W}|_t) = t$ . This implies  $t \in (Q|_t)^{\mathbb{D}}$ .

Now suppose that  $t \in (Q|t)^{\mathbb{D}}$ ; then there exists some  $G' \in \text{dags}(h_t(\text{body}_Q))$  and some terminal embedding  $\phi : G' \rightsquigarrow \mathbb{D}$  satisfying  $\hat{\phi}(\overline{W}|t) = t$ , implying<sup>4</sup>  $G' \in \text{dags}(h_t(\text{body}_Q), \mathfrak{ts}_t^Q)$ . Then  $\gamma := \phi \circ h_t$  is a terminal embedding of the DAG  $G \in \text{dags}(\text{body}_Q)$  satisfying  $G' = h_t(G)$ , and  $\hat{\gamma}(\overline{W}) = t$  because  $G \in \text{dags}(\text{body}_Q, \mathfrak{ts}_t^Q)$ .  $\square$

The following theorem is the main result of this section. It reduces DACQ equivalence to a set of tuple-specific equivalence tests between CH-DACQs, which was the topic of Section 4.3.2.

**Theorem 4.4.10** *Given two DACQs  $Q$  and  $Q'$  in canonical form,  $Q \equiv Q'$  iff*

1.  $\text{CT}_Q \approx \text{CT}_{Q'}$ , and
2.  $\forall t \in \text{CT}_Q, Q|t \stackrel{t}{\equiv} Q'|t$ .

PROOF. “If” follows from Proposition 4.4.1, Lemma 4.4.5, and Corollary 4.4.7. “Only if” follows from Lemmas 4.4.6 and 4.4.9.  $\square$

## 4.4.2 Avoiding Tuple-Specific Equivalence Tests

In Section 4.3 we reduced CH-DACQ equivalence to testing one tuple-specific CH-DACQ equivalence problem per canonical tuple (Theorem 4.3.14). While Theorem 4.4.10 proves a similar result for arbitrary DACQs, the difference is that the number of canonical tuples per tuple shape depends upon the number of canonical freezings, which can be exponential in the number of variables in the head tuple. It is well-known that to test CQ equivalence it suffices to test boolean CQ equivalence for *only* the most general freezing (a consequence of the Homomorphism Theorem), and so the algorithm yielded by Theorem 4.4.10 is clearly not optimal for CQs (i.e., DACQs with one vertex in the body). However, the following example illustrates that for arbitrary DACQs it is not sufficient to only consider the most general freezings.

**Example 50** Continuing Example 49, consider the two queries obtained by freezing  $Q_{48}$  with canonical tuples  $abc\text{--}$  and  $aace$ .

$$\begin{array}{l}
 Q_{48}|^{abc\text{--}}(a^{\boxed{1}}, b^{\boxed{1}}, c^{\boxed{1}}, 0^{\boxed{2}}, 0^{\boxed{3}}) :- \begin{array}{l} \boxed{1} R(a, b, c), \\ R(a, b, D) \end{array} \begin{array}{l} \xrightarrow{\quad} \boxed{2} S(c, E), a = b, c = D \\ \xrightarrow{\quad} \boxed{3} b = c, b = D \end{array} \\
 \\
 Q_{48}|^{aace}(a^{\boxed{1}}, a^{\boxed{1}}, c^{\boxed{1}}, e^{\boxed{2}}, 0^{\boxed{3}}) :- \begin{array}{l} \boxed{1} R(a, a, c), \\ R(a, a, D) \end{array} \begin{array}{l} \xrightarrow{\quad} \boxed{2} S(c, e), a = a, c = D \\ \xrightarrow{\quad} \boxed{3} a = c, a = D \end{array}
 \end{array}$$

<sup>4</sup>We needed to insert arbitrary constants into head  $\overline{W}|t$  of query  $Q|t$  so that  $\overline{W}|t$  retains the same labels and tuple shapes as  $\overline{W}$ . If we instead defined  $\overline{W}|t := t$ —which effectively removes any output labels not in  $\mathfrak{ts}_t^Q$ —then  $Q|t$  would output the tuple  $t$  too often (in particular,  $Q|t$  would incorrectly output  $t$  for any terminal embedding  $\phi' : G' \rightsquigarrow \mathbb{D}$  of any  $G' \in \text{dags}(h_t(\text{body}_Q), \mathfrak{ts}')$  for which  $\mathfrak{ts}_t^Q \subset \mathfrak{ts}'$ ).

From the queries in Figure 4.11g it is straightforward to verify the tuple-specific equivalences.  $Q_{47}|_{abc\text{--}} \stackrel{\langle abc\text{--} \rangle}{\equiv} Q_{48}|_{abc\text{--}}$  and  $Q_{47}|_{aace\text{--}} \stackrel{\langle aace\text{--} \rangle}{\equiv} Q_{48}|_{aace\text{--}}$ . In fact, for any canonical tuple  $t$  with shape  $\mathfrak{t}_{12}$ ,  $\mathfrak{t}_{13}$ , or  $\mathfrak{t}_{123}$  we can show that the  $t$ -specific equivalence  $Q_{47}|^t \stackrel{t}{\equiv} Q_{48}|^t$  holds (even though the generic equivalence  $Q_{47}| \equiv Q_{48}|$  may not hold). In contrast, for shape  $\mathfrak{t}_1$ , tuple-specific equivalence holds for the most general freezing  $abc\text{--}$ , but the database  $\mathbb{D}_7 = \{R(a, a, c), R(a, a, d), S(c, e)\}$  proves  $Q_{47}|_{aac\text{--}} \stackrel{\langle aac\text{--} \rangle}{\not\equiv} Q_{48}|_{aac\text{--}}$  and therefore  $Q_{47} \stackrel{\langle aac\text{--} \rangle}{\not\equiv} Q_{48}$ , as demonstrated by the following relations  $(Q_{47})^{\mathbb{D}_7}$  and  $(Q_{48})^{\mathbb{D}_7}$ .

$(Q_{47})^{\mathbb{D}_7}$	<table style="border-collapse: collapse; text-align: center;"> <tr> <th style="border: 1px solid black; padding: 2px 5px;"><math>A^{[1]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>B^{[1]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>C^{[1]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>E^{[2]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>B^{[3]}</math></th> </tr> <tr> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>c</math></td> <td style="padding: 2px 5px;"><math>e</math></td> <td style="padding: 2px 5px;">-</td> </tr> <tr> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>d</math></td> <td style="padding: 2px 5px;">-</td> <td style="padding: 2px 5px;">-</td> </tr> </table>	$A^{[1]}$	$B^{[1]}$	$C^{[1]}$	$E^{[2]}$	$B^{[3]}$	$a$	$a$	$c$	$e$	-	$a$	$a$	$d$	-	-
$A^{[1]}$	$B^{[1]}$	$C^{[1]}$	$E^{[2]}$	$B^{[3]}$												
$a$	$a$	$c$	$e$	-												
$a$	$a$	$d$	-	-												

$(Q_{48})^{\mathbb{D}_7}$	<table style="border-collapse: collapse; text-align: center;"> <tr> <th style="border: 1px solid black; padding: 2px 5px;"><math>A^{[1]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>B^{[1]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>C^{[1]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>E^{[2]}</math></th> <th style="border: 1px solid black; padding: 2px 5px;"><math>C^{[3]}</math></th> </tr> <tr> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>c</math></td> <td style="padding: 2px 5px;"><math>e</math></td> <td style="padding: 2px 5px;">-</td> </tr> <tr> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>c</math></td> <td style="padding: 2px 5px;">-</td> <td style="padding: 2px 5px;">-</td> </tr> <tr> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>a</math></td> <td style="padding: 2px 5px;"><math>d</math></td> <td style="padding: 2px 5px;">-</td> <td style="padding: 2px 5px;">-</td> </tr> </table>	$A^{[1]}$	$B^{[1]}$	$C^{[1]}$	$E^{[2]}$	$C^{[3]}$	$a$	$a$	$c$	$e$	-	$a$	$a$	$c$	-	-	$a$	$a$	$d$	-	-
$A^{[1]}$	$B^{[1]}$	$C^{[1]}$	$E^{[2]}$	$C^{[3]}$																	
$a$	$a$	$c$	$e$	-																	
$a$	$a$	$c$	-	-																	
$a$	$a$	$d$	-	-																	

The “problem” illustrated by Example 50 is that the tuple-specific equivalence  $Q_{47}|_{abc\text{--}} \stackrel{\langle abc\text{--} \rangle}{\equiv} Q_{48}|_{abc\text{--}}$  does not imply  $Q_{47}|_{aac\text{--}} \stackrel{\langle aac\text{--} \rangle}{\equiv} Q_{48}|_{aac\text{--}}$ , even though canonical tuple  $\langle abc\text{--} \rangle$  is more general than  $\langle aac\text{--} \rangle$ . In order to do better than the naive algorithm that tests tuple-specific equivalence for every canonical tuple, we require an admissible heuristic for identifying when these “implication failures” are guaranteed not to occur. The next theorem proposes a sufficient condition for identifying when the reasoning that  $Q|^{t_1} \stackrel{t_1}{\equiv} Q|^{t_1}$  also implies  $Q|^{t_2} \stackrel{t_2}{\equiv} Q|^{t_2}$ , for two canonical tuples satisfying  $t_1 > t_2$ . First, however, we need to revisit how the chase procedure is used to decide implication of disjunctive embedded dependencies. (The reader is referred to page 115 for discussion of how the chase procedure is used to decide implication of embedded dependencies without disjunction.)

Given disjunctive embedded dependency  $\xi$  (see equation 4.20) and a set of disjunctive embedded dependencies  $\Sigma$ , the chase  $\text{chase}_\Sigma(P^\xi)$  is a (possibly infinite) *chase tree* in which each chase step has one branch for each conclusion in the chased dependency. Each path  $p$  in the tree corresponds to a unique chase sequence of (non-disjunctive) embedded dependencies that defines some (possibly infinite) database  $\mathbb{D}_p$ . The chase procedure returns the conclusion  $\Sigma \models_{\text{fin}} \xi$  if and only if for each path  $p$  it can prove  $\mathbb{D}_p \models \hat{\xi}$ , which it does by exhibiting a finite database  $\hat{\mathbb{D}}_p \subseteq \mathbb{D}_p$ —corresponding to some finite prefix of path  $p$ —that satisfies  $\hat{\mathbb{D}}_p \models \hat{\xi}$ . Given any mapping  $h : \mathcal{C} \rightarrow \mathcal{C}$ , and a specific chase invocation that yielded the conclusion  $\Sigma \models_{\text{fin}} \xi$ , we say that the conclusion is *sound w.r.t. h* if for each path  $p$ , whenever the database  $\hat{\mathbb{D}}_p$  contains a contradiction, so does the database  $h(\hat{\mathbb{D}}_p)$ .

**Theorem 4.4.11** *Given DACQs  $Q$  and  $Q'$  in canonical form satisfying  $\text{CT}_Q = \text{CT}_{Q'}$ , and given two canonical tuples  $t_1, t_2 \in \text{CT}_Q$  satisfying  $t_1 > t_2$ , let  $h$  be the  $\mathcal{C}_S$ -homomorphism (extended with identity) satisfying  $h(t_1) = t_2$ . If  $Q|^{t_1} \stackrel{t_1}{\equiv} Q'|^{t_1}$  (i.e.,  $\Sigma(Q|^{t_1}, t_1) \equiv_{\text{fin}} \Sigma(Q'|^{t_1}, t_1)$ ) can be proved using only chase invocations that are sound w.r.t.  $h$ , then  $Q|^{t_2} \stackrel{t_2}{\equiv} Q'|^{t_2}$  (i.e.,  $\Sigma(Q|^{t_2}, t_2) \equiv_{\text{fin}} \Sigma(Q'|^{t_2}, t_2)$ ).*



PROOF. Suppose that  $Q|^{t_1} \stackrel{t_1}{\equiv} Q|^{t_1}$ . Due to Definition 4.4.8 and equations 4.21 and 4.22, we observe that  $\Sigma^{(Q|^{t_2}, t_2)} = h(\Sigma^{(Q|^{t_1}, t_1)})$  and  $\Sigma^{(Q'|^{t_2}, t_2)} = h(\Sigma^{(Q'|^{t_1}, t_1)})$ .

Consider any chase invocation that proves some implication  $\Sigma \models_{\text{fin}} \xi$  as part of the proof that  $\Sigma^{(Q|^{t_1}, t_1)} \equiv_{\text{fin}} \Sigma^{(Q'|^{t_1}, t_1)}$ . Let  $T$  denote a (possibly infinite) chase tree for  $\text{chase}_{\Sigma}(P^{\hat{\xi}})$ , with  $\mathbb{D}_p$  denoting the (possibly infinite) database defined by path  $p$ . Let  $\hat{T}$  denote the finite prefix of  $T$  that was actually explored by the chase invocation, with  $\hat{\mathbb{D}}_p \subseteq \mathbb{D}_p$  the finite database that terminated the search of the chase sequence along path  $p$ .

Define  $\mathbb{D}_1 := P^{\hat{\xi}}$ , and let  $\mathbb{D}_1 \xrightarrow{(\xi_i, \gamma)} \{\mathbb{D}_{1.1}, \dots, \mathbb{D}_{1.k}\}$  be the first chase step in chase tree  $T$  enacted by some variable mapping  $\gamma : X \rightarrow \text{adom}(\mathbb{D}_1)$ , where  $X$  is the set of premise and free variables of some dependency  $\xi_i \in \Sigma$  having  $k$  conclusions, and  $\gamma$  embeds premise  $P^{\xi_i}$  into  $\mathbb{D}_1$ . Let  $\mathbb{D}_{1.j}$  denote the database obtained from  $\mathbb{D}_1$  by chasing with the  $j^{\text{th}}$  conclusion of  $\xi_i$ .

Now consider the dependency set  $h(\Sigma)$  and the dependency  $h(\xi)$  which has premise  $P^{h(\xi)} = h(P^{\xi}) = h(\mathbb{D}_1)$ . Because  $\xi_i$  and  $h(\xi_i)$  have the same variables ( $h$  only remaps certain constants), it follows that  $h \circ \gamma : X \rightarrow \text{adom}(h(\mathbb{D}_1))$  is an embedding of  $h(P^{\xi_i}) = P^{h(\xi_i)}$  into  $h(\mathbb{D}_1)$ . Furthermore,  $h(\mathbb{D}_{1.j})$  is the database obtained from  $h(\mathbb{D}_1)$  by chasing with the  $j^{\text{th}}$  conclusion of  $h(\xi_i)$ . It follows that  $h(\mathbb{D}_1) \xrightarrow{(h(\xi_i), h \circ \gamma)} \{h(\mathbb{D}_{1.1}), \dots, h(\mathbb{D}_{1.k})\}$  is a valid—although possibly spurious<sup>5</sup>—first step for some chase invocation proving  $h(\Sigma) \models_{\text{fin}} h(\xi)$ . By induction, it follows that  $h(\hat{T})$  is the prefix of a valid chase tree—possibly containing spurious chase steps—for  $\text{chase}_{h(\Sigma)}(P^{h(\hat{\xi})})$ .

Now consider any database  $\hat{\mathbb{D}}_p$  at the leaves of  $\hat{T}$ . We know that  $\hat{\mathbb{D}}_p \models \hat{\xi}$ , which implies that either there exists an embedding  $\gamma : C_j^{\xi} \rightarrow \hat{\mathbb{D}}_p$  of one of the conclusions of  $\xi$ , or  $\hat{\mathbb{D}}_p$  contains a contradiction (see page 115). If  $\gamma$  exists, then, by the argument above,  $h \circ \gamma$  is an embedding of  $h(C_j^{\xi}) = C_j^{h(\xi)}$  into  $h(\hat{\mathbb{D}}_p)$ . Alternatively, the final chase step leading into  $\hat{\mathbb{D}}_p$  may have derived a contradiction, either explicitly (i.e., a dependency of the form  $P^{\xi_i} \rightarrow \perp$ ) or implicitly by equating two distinct constants  $a, b \in \mathcal{C}$ ; as long as the chase invocation is sound w.r.t.  $h$ ,  $h(\hat{\mathbb{D}}_p)$  will also contain a contradiction. In either case,  $h(\hat{\mathbb{D}}_p) \models h(\xi)$ . Hence, the chase tree  $h(\hat{T})$  proves  $h(\Sigma) \models_{\text{fin}} h(\xi)$ .  $\square$

**Corollary 4.4.12** *For any shape  $\mathfrak{ts} \in \text{TS}_Q$  and corresponding shape  $\mathfrak{ts}' \in \text{TS}_{Q'}$ , if*

- *all of the equality predicates in  $\text{body}_Q$  occur within  $G_{\mathfrak{ts}}^{\min} \in \text{dags}(\text{body}_Q)$ , and*
- *all of the equality predicates in  $\text{body}_{Q'}$  occur within  $G_{\mathfrak{ts}'}^{\min} \in \text{dags}(\text{body}_{Q'})$ ,*

*then only the most general canonical tuple of shape  $\mathfrak{ts}$  needs to be considered when deciding  $Q \equiv Q'$ .*

---

<sup>5</sup>The chase step  $\mathbb{D}_1 \xrightarrow{(\xi_i, \gamma)} \{\mathbb{D}_{1.1}, \dots, \mathbb{D}_{1.k}\}$  is spurious if  $\mathbb{D}_1 \models \gamma(\xi_i)$ . The introduction of a spurious chase step can cause an otherwise finite chase sequence to become infinite, but it does not change the soundness of the conclusion  $\Sigma \models_{\text{fin}} \xi$ .

Corollary 4.4.12 can yield an exponential reduction in the number of canonical tuples that need to be considered. In particular, for queries expressible without explicit equality predicates, Corollary 4.4.12 implies that the number of tuple-specific equality tests required is independent of the number of variables in the query head. Explicit equality is only required to equate ancestor variables to constants or to each other, which corresponds to a class of LOJ predicates that occurs only rarely in real world workloads. The vast majority of LOJ predicates in real-world SPCL workloads equate one attribute from each side of the LOJ, which can be expressed as a DACQ without explicit equality by using shared variables.

**Example 51** Consider the following queries  $Q_{50}$  and  $Q_{51}$  which are written over the sample company database schema from Chapter 1 (cf. Figure 1.1).

$$\begin{aligned}
 Q_{50}(M^{\boxed{a}}, O^{\boxed{b}}, P^{\boxed{c}}, Y^{\boxed{c}}) &:- \boxed{a} \mathbf{C}(C, M, \text{'C'}) \longrightarrow \boxed{b} \mathbf{O}(O, C, D) \longrightarrow \boxed{c} \mathbf{LI}(O, L, P, Y) \\
 Q_{51}(M^{\boxed{x}}, O^{\boxed{y}}, P^{\boxed{z}}, Y^{\boxed{z}}) &:- \boxed{x} \mathbf{C}(C, M, T) \longrightarrow \boxed{y} \mathbf{O}(O, C, D), \\
 &\quad T = \text{'C'} \longrightarrow \boxed{z} \mathbf{LI}(O, L, P, Y)
 \end{aligned}$$

Query  $Q_{50}$  returns the names of all Corporate customers, along with the order ids of any associated orders (if they exist), and the price and quantity of any associated lineitems (if they exist); this corresponds to a typical usage of the LOJ operator, and does not require explicit equality when expressed as a DACQ. Query  $Q_{51}$  is more esoteric in that it returns the names of *all* customers, but only returns the associated orders and lineitems for the *Corporate* customers; this type of LOJ predicate necessitates explicit equality when expressed as a DACQ.

## 4.5 Soundness and Completeness

We have reduced DACQ query equivalence to deciding tuple-specific equivalence between CH-DACQs (Theorems 4.3.14 and 4.4.10), and tuple-specific CH-DACQ equivalence to finite equivalence of sets of disjunctive embedded dependencies (Corollary 4.3.22). Certainly it is not ideal to reduce a problem we only know to be NP-hard to one that is undecidable in general. Given a DACQ equivalence problem, our algorithm might construct an instance of an implication problem for which we cannot guarantee chase termination (which does not prove that DACQ equivalence is undecidable). In this section we briefly discuss two strategies for dealing with difficult implication problems, each of which preserves soundness (but not completeness) of the equivalence test.

### 4.5.1 Early Termination

As described in Section 4.3.2, even when database  $\mathbb{D} := \text{chase}_{\Sigma}(P^{\hat{\xi}})$  is infinite, the chase algorithm may be able to conclude  $\Sigma \models_{\text{fin}} \xi$  based upon a finite sequence of

chase steps yielding some  $\hat{D} \subset D$ . Early termination is the obvious pragmatic solution: run the chase algorithm for a bounded amount of time, and return “I don’t know” if by the end of the time limit neither  $\Sigma \models_{\text{fin}} \xi$  nor  $\Sigma \not\models_{\text{fin}} \xi$  has been derived. In the simplest form, the number of chase steps can simply be fixed; unfortunately, this yields a non-deterministic result depending upon the order of chase steps. Alternatively, bounding either the *dependency depth*—the maximum length of any chain of order-dependencies between chase steps—or the *Skolem depth*—the nesting depth of the Skolem functions (constants) introduced by chase steps—yields a deterministic result with specific properties that can be used to characterize subclasses of database instances over which the implication does hold. Recently, more sophisticated methods have been proposed that dynamically monitor the chase execution and terminate based upon violation of some boundary condition [81].

## 4.5.2 Constraint Relaxation

The alternative to early termination is to invoke the chase only when it is guaranteed to terminate. Given a dependency set  $\Sigma$  that does not guarantee a terminating chase procedure, the main idea behind constraint relaxation is to replace  $\Sigma$  with a set  $\hat{\Sigma}$  that both satisfies  $\Sigma \models_{\text{fin}} \hat{\Sigma}$  and guarantees a terminating chase procedure.

**Proposition 4.5.1** *Given dependency  $\xi$  and set  $\Sigma$ , let  $\hat{\Sigma}$  be any dependency set satisfying  $\Sigma \models_{\text{fin}} \hat{\Sigma}$ . If  $\hat{\Sigma} \models_{\text{fin}} \xi$  then  $\Sigma \models_{\text{fin}} \xi$ .*

The trick is in constructing a useful set  $\hat{\Sigma}$  without requiring a chase procedure to explicitly prove  $\Sigma \models_{\text{fin}} \hat{\Sigma}$ . One technique is to construct  $\hat{\Sigma}$  from  $\Sigma$  by identifying dependencies that cause violations of the termination guarantee, and for each such dependency either delete it from  $\Sigma$  or relax it by adding atoms/predicates to the premise and/or deleting atoms/predicates from the conclusion.

**Example 52** Consider the following dependency  $\xi$  and dependency set  $\Sigma = \{\xi_1, \xi_2\}$ .

$$\begin{aligned} \xi &: [R(X, Y)] \rightarrow [S(Y, Z)] \\ \xi_1 &: [R(A, B), R(A, C)] \rightarrow [S(B, D), S(C, D)] \\ \xi_2 &: [R(A, B), R(A, C)] \rightarrow [R(B, E)] \end{aligned}$$

Set  $\Sigma$  does not guarantee chase termination, and indeed the database  $\text{chase}_{\Sigma}([R(x, y)])$  is not finite. If we choose  $\hat{\Sigma} = \{\xi_1\}$  then  $\Sigma \models_{\text{fin}} \hat{\Sigma}$  is trivial, and  $\hat{\Sigma}$  guarantees chase termination (due to *weak acyclicity*—see Section 4.6.4). Then,  $\text{chase}_{\hat{\Sigma}}([R(x, y)]) = \{R(x, y), S(y, d_1)\}$  proves  $\hat{\Sigma} \models_{\text{fin}} \xi$  and therefore  $\Sigma \models_{\text{fin}} \xi$ .

## 4.6 Relevant Literature

To the best of our knowledge, there is no previous literature that directly addresses the problem of query equivalence for queries containing outer joins. There is, however, a wide variety of literature relevant to various aspects of the problem and/or

our solution to it, and so in this section we briefly survey literature surrounding five problems that seem the most pertinent—algebraic optimization of outer joins, containment for query languages with negation, query containment in the presence of nulls, finite implication of dependencies (equivalently, query containment under a set of dependencies), and implication of dependencies over relations containing nulls.

### 4.6.1 Algebraic Optimization of Outer Joins

Research into algebraic optimization of queries containing outer joins was pioneered by Galindo-Legaria in his doctoral research with his supervisor, Rosenthal [96, 40, 36]. They propose and study the problem of how to generate all equivalent join re-orderings of a given “outer-join tree,” which is an algebraic expression composed of inner, left (right) outer, and full outer joins, together with selection. While this problem is straightforward for inner joins because both commutativity and associativity hold, outer joins are not associative in general. The culmination of this work is a set of transformation rules that is complete for enumerating equivalent outer-join trees [41]. While a complete set of algebraic transformations could serve as the basis for a decision procedure for query equivalence, it is important to note the lack of a projection operator—which means that relation names can be assumed to be unique (an assumption that simplifies the equivalence problem for CQs from NP-complete to linear time). Instead, the main contribution of this work is the study of the conditions under which outer joins are associative.

Subsequent work by Galindo-Legaria proposes *join-distinct normal form (JDNF)* as a canonical representation for outer-join trees [37]. JDNF is formed by rewriting the outer joins in terms of inner join and *minimal union*—see equation 4.2—and then commuting all inner joins below all minimal unions. Given a DACQ  $Q$  containing unique relation names and whose head tuple does not project away any attributes, JDNF corresponds to replacing union with minimal union within “shape query”  $\hat{Q}$  (see Definition 4.3.5); hence, the JDNF is exponential in the size of  $Q$ . JDNF has been used as the basis for a view rewriting algorithm that handles outer joins [71], but arguably its main contribution is conceptual—it explicitly exposes the interactions between the predicates of the inner and outer-join operators in the join tree.

Later authors have proposed alternative canonical representations for outer-join trees. Bhargava, Goel, and Iyer define a canonical representation that encodes the join predicates within a hypergraph—offering a visual representation of the interactions between join predicates that is much more concise than JDNF [12]. They then show how hypergraph connectivity can be used to guide the algebraic manipulations in order to enumerate all equivalent outer-join trees. Rao, Pirahesh, and Zuzarte propose a canonical representation for outer-join trees that generalizes the standard normal form for SPC queries [6, Ch. 4], by replacing Cartesian product with *outer Cartesian product* and selection with *nullification* [94]. Whereas the techniques of

both Galindo-Legaria and Bhargava et al. are geared towards transformation-based enumeration of equivalent join trees, the primary contribution of Rao et al.’s canonical form is that it facilitates integration into bottom-up dynamic-programming join enumeration.

Recently, Hill and Ross have proposed an algorithm for enumerating equivalent join trees that converts outer-join trees into inner join trees over “derived relations” [60]. Their focus is on encapsulating all of the logic pertaining to the non-associativity of outer joins into a pre-computation step, allowing the join enumeration phase to assume associativity. By removing the outer-join-specific logic from the join enumeration phase, they are able to extend existing methods for expanding an optimizer’s search space—in particular, semijoin reduction, which is an important technique within distributed query processing environments—to queries containing outer joins.

## 4.6.2 Queries with Negation

The LOJ operator can be defined in terms of set difference (see equation 4.1), which accounts for its non-monotonicity with respect to null-padded tuples. It is well known that adding unrestricted set difference to the conjunctive algebra yields the full relational algebra, for which query containment is undecidable [6]. However, there are interesting restrictions on the negation operator that yield decidable query containment.

The most fundamental result is due to Sagiv and Yannakakis, who show that the containment problem is only NP-complete for unions of “elementary differences” (an elementary difference is a difference operations between two CQs) [99]. Their complexity result hinges upon the fact that within this class of expressions, projection never occurs after a union or difference operator. As a consequence, query containment is  $\Pi_2^p$ -complete for arbitrary relational algebra expressions as long as the projection operator never occurs above difference—which implies that SPCL equivalence is decidable as long as projection never occurs above LOJ. Note the close correspondence between this class of SPCL expressions and the class of outer-join trees studied by Galindo-Legaria, described above. The analogous syntactic restriction for DACQs is that all local variables of all non-leaf nodes appear in the query head, labelled by the node to which they are local. For this special case, the necessary condition in Corollary 4.3.8 and the sufficient condition in Theorem 4.3.9 converge to form a complete characterization of DACQ equivalence.

Other restricted forms of negation have also been explored:

- Ullman shows that extending CQs with safe atomic negation (i.e., negation of individual atoms whose variables all occur within positive atoms) causes containment to be  $\Pi_2^p$ -complete [106]. Whereas Ullman’s algorithm naively generates canonical databases, for this same class Wei and Lausen propose a recursive decision procedure that essentially performs a chase algorithm over

sets of disjunctive *full* dependencies. They then extend their algorithm—without any increase in complexity—to decide containment for *unions* of CQs with safe atomic negation [111].

- Klug shows that extending CQs with the predicates  $\neq$ ,  $\leq$ , and  $<$  keeps the CQ containment (and equivalence) problem in  $\Pi_2^P$  [70], assuming an abstract totally-ordered underlying domain. The matching lower bound is due to van der Meyden, who proves that just  $\neq$  is enough to make the CQ containment problem  $\Pi_2^P$ -hard [108]. Ibarra and Su generalize this result further to consider CQs with linear arithmetic constraints, and show the containment and equivalence problems to be decidable in double exponential time but not non-deterministic single exponential time when queries are evaluated over the integers, and decidable in double exponential time but not non-deterministic polynomial time when constant-free queries are evaluated over the real numbers [61].

Unfortunately, there does not appear to be an obvious way to reduce equivalence of arbitrary DACQs either to or from the equivalence or containment problems for any of these extensions to CQs.

### 4.6.3 Query Containment with Null Values

Only recently has the interaction between null values and query containment been considered. Farré, Nutt, Teniente, and Urpí study the containment problem for CQs over database instances containing null values, under the assumption that equality and join predicates are evaluated using 3-valued logic [34]. They show that the containment problem remains NP-complete for boolean queries, and corresponds to the existence of homomorphisms that map join variables onto other join variables. Interestingly, the problem becomes  $\Pi_2^P$ -complete as soon as the language allows either output variables or an explicit *isNull* predicate (or both).

### 4.6.4 Query Containment under Dependencies

Because embedded dependencies can be expressed as a pair of CQs—the premise and the conclusion—finite implication of embedded dependencies corresponds precisely to containment of CQs under a set of dependencies. It is standard practise to transform sets of embedded dependencies into an equivalent form in which every dependency is either *tuple-generating* or *equation-generating* [6, Ch. 10], and most of the literature that we mention in this section assume this distinction. We did not make this distinction when we introduced embedded dependencies in Section 4.3 because it is orthogonal to our reduction of query containment/equivalence to dependency implication; furthermore, sets of *disjunctive* embedded dependencies can not be decomposed in this way.

The CQ containment problem is decidable for any class of dependencies that allows a terminating chase procedure. Well-known examples include the class of



*full* dependencies—i.e., dependencies that do not contain existentially-quantified variables—and the class FDs + JDs + acyclic INDs [6, Ch. 9–10]. Johnson and Klug show that containment of CQs under FDs and INDs is decidable—even though the chase does not terminate—as long as the INDs are *key-based*; however, their proof allows for infinite counter-example databases [67]. Rosati extends Johnson and Klug’s work and shows that CQ containment under INDs or key-based FDs + INDs is *finitely controllable*, i.e., that the finite-model and infinite-model problems coincide [95]. To do this, he uses a complex process for generating chase constants that is periodic (thereby bounding the total number of constants) but whose period is provably long enough so that any infinite chase sequence that serves as a counter-example to containment can be mapped to a “periodic” chase sequence that is also a counter-example to containment.

For arbitrary embedded dependencies the problem is known to be undecidable; however, there has been significant research into trying to identify when CQ containment is decidable for a particular set of embedded dependencies. Calì, Gottlob, and Kifer generalize Johnson and Klug’s work by considering CQ containment under arbitrary tuple-generating and equality-generating dependencies [15]. They define classes of *guarded* dependencies for which CQ containment is EXPTIME-complete with bounded schema size, or 2-EXPTIME-complete in general. With guarded dependencies the chase is still unbounded, but the chase sequence can be modelled by what they call a *squid decomposition*, which restricts the possible interactions between constants generated by different chase paths. Like Johnson and Klug, their decidability results assume a data model allowing infinite databases; it is not known whether Rosati’s technique can be extended to show finite controllability for this more general problem.

Other authors have focused on identifying when the chase is guaranteed to terminate. Deutsch and Popa propose a termination condition called *weak acyclicity*<sup>6</sup> [31, 33]; if set  $\Sigma$  is weakly acyclic, then any chase over set  $\Sigma$  is guaranteed to terminate. Weak acyclicity is directly generalized to a termination condition called *stratification* by Deutsch, Nash, and Rimmel [30]. Those authors also show that it is undecidable to identify when a given  $\Sigma$  allows a terminating chase procedure. Meier, Schmidt, and Lausen generalize the stratification condition further, defining an infinite hierarchy of sufficient termination conditions, where each level  $k$  in the hierarchy corresponds to a proof that any chase over set  $\Sigma$  will not generate fresh constants with Skolem depth greater than  $k$ .

Most of the research on query containment under dependencies focuses on embedded dependencies for which the chase is a linear search of a sequence of chase steps. Beeri and Vardi propose extending the chase procedure to a tree search to handle dependencies containing disjunction in the conclusion [10]. Although not widely used, this extension does add expressive power. Consider a set  $\Sigma$  containing disjunctive dependencies. Let  $\hat{\Sigma}$  be the corresponding set of “flattened”

---

<sup>6</sup>The condition was derived collaboratively by Deutsch and Popa, but published independently (with other authors) under the names “weak acyclicity” [33] and “stratified-witness” [31].

dependencies, formed by decomposing each  $\xi \in \Sigma$  into separate (non-disjunctive) dependencies for each conclusion. Any path within a chase tree over  $\Sigma$  corresponds to a valid prefix of a chase sequence over  $\hat{\Sigma}$ . This implies that decidability under  $\hat{\Sigma}$  is a sufficient (although not necessary) condition for decidability under  $\Sigma$ , making all of the research on decidability for non-disjunctive dependencies directly applicable to disjunctive dependencies.

#### 4.6.5 Dependencies over Relations with Nulls

As mentioned above, the problem of conjunctive query containment under embedded dependencies corresponds precisely to (finite) implication of embedded dependencies. Much of the historic work on dependency analysis is motivated by problems related to database schema design rather than query equivalence, and often considers much simpler dependency classes than arbitrary embedded dependencies. The most fundamental result is the axiomatization of inference over sets of functional dependencies proposed by Armstrong [7], which has been extended to inference over sets of FDs and MVDs [6, Sect. 8.2]. A key assumption in this work (and in the work surveyed on query equivalence under dependencies) is that the data model does not include NULL values. Literature that extends dependency analysis to relations containing NULL values is pertinent to Section 4.2—specifically, to our definition and use of *nullability dependencies* to characterize property [P5] of DACQ canonical form.

Vassiliou pioneered the analysis of dependencies over relations containing NULL values [109, 110]. He points out that a NULL value can be interpreted in two different ways—either to mean “unknown value,” or to mean “no information” [109]. The “unknown value” interpretation assumes that there exists some substitution of domain values for NULLS yielding tuples that occur within some unknown “complete” version of the relation; hence, relations with NULLS under this interpretation are often referred to as “incomplete relations”—a concept further developed by Imielinski and Lipski [62]—and reasoning about them requires considering the set of possible worlds. In contrast, the “no information” interpretation is more primitive, making no presumption that there exists any underlying completed version of the relation.

Vassiliou extends the concept of functional dependencies to relations containing NULLS, assuming the unknown value interpretation [110]. He classifies FDs as either *strong* or *weak*. Vassiliou’s definitions of strong and weak FDs are in terms of lists of syntactic conditions, but more intuitive characterizations used by Levene and Loizou [73] are that an FD is strong if under all possible worlds the completed relation obeys the FD, whereas an FD is weak if there exists at least one possible world under which the completed relation obeys the FD. Vassiliou shows that Armstrong’s axioms are sound and complete for inference over strong FDs. Also assuming the unknown value interpretation, Levene and Loizou extend Armstrong’s axioms to obtain a set of axioms that is sound and complete



for inference over sets of both strong and weak FDs [73]. Motivated by schema design for databases containing *nested relations*, Levene and Loizou also define an extension to the nested relational model to handle NULL values, and they consider “null extended data dependencies”—specifically, null extended FDs, MVDs, and join dependencies (JDs)—which extend traditional FDs/MVDs/JDs to both nested relations and NULL values (interpreted as unknown values) [72].

Atzeni and Morfuni study implication of FDs under the “no information” interpretation [8]. They define FD  $X \rightarrow Y$  to hold over a relation  $R$  when for any two tuples  $t, t' \in R$ , if  $t[X] = t'[X]$  and  $t[X], t'[X]$  do not contain NULL, then  $t[Y] = t'[Y]$ . Under this definition of FDs, Armstrong’s axiom of transitivity fails to hold, and so they introduce a concept they call “null transitivity,” which they use to devise a modified set of axioms that are sound and complete for implication of FDs over relations containing NULLS. Our definition of a *nullability dependency* (Definition 4.2.2) is convenient for our purposes in this chapter, but it can be equivalently expressed in terms of Atzeni and Morfuni’s definition of an FD, as follows.

**Proposition 4.6.1** *Given a query  $Q$  and a nullability dependency  $ABC \Rightarrow v$  with  $v \in \text{body}_Q$  and  $\{A, B, C\} \subseteq \mathcal{B}_v^{\text{anc}}$ , let  $c$  be any constant, let  $\{u_1, \dots, u_k\}$  be the set of ancestor nodes of  $v$ , and let  $u_A$  denote the ancestor node of  $v$  satisfying  $A \in \mathcal{L}_{u_A}$ , with  $u_B$  and  $u_C$  defined analogously. Then,  $Q \models ABC \Rightarrow v$  if and only if over every database  $\mathbb{D}$ , the relation  $(\text{body}_Q)^{\mathbb{D}}$  (i.e., the result of  $Q$  prior to projection onto the head attributes) satisfies the FD  $A^{\boxed{u_A}} B^{\boxed{u_B}} C^{\boxed{u_C}} c^{\boxed{u_1}} \dots c^{\boxed{u_k}} \rightarrow c^{\boxed{v}}$  (under Atzeni and Morfuni’s definition).*

Whereas all of the literature discussed above addresses dependencies in the context of database schema design, Paulley in his doctoral dissertation analyzes how FDs are introduced by and propagated through algebraic queries [92]. Paulley defines a “strict” FD  $X \rightarrow Y$  to hold over a relation with NULLS if the relation satisfies the FD under the classical definition, treating NULL as a normal domain value. He defines a “lax” FD  $X \mapsto Y$  to hold if the subset of the relation not containing a NULL value in either  $X$  or  $Y$  satisfies the FD  $X \rightarrow Y$ . Of particular note, Paulley works out the interaction between strict/lax FDs and the left outer-join operator [92, Ch. 3]. In the introduction to the next chapter we raise but leave open the problem of testing query validity for encoding queries with DACQ query bodies and encoding heads containing arbitrary assignments of node labels to index attributes. This problem requires testing that the query body implies that the query result—which is a relation containing NULL (i.e.,  $\blacksquare$ ) values—always satisfies the strict FD  $\bar{\mathcal{I}}_{[1,d]} \rightarrow \bar{\mathcal{V}}$ . We conjecture that this problem could be solved using Paulley’s results to characterize the effect of arbitrary assignments of node labels within the index attributes.



# Chapter 5

## Encoding-Equivalence of Hierarchical CEQ's

In Chapter 3 we considered the encoding-equivalence problem for CEQs—that is, encoding queries whose bodies are CQs—and we demonstrated an equivalence procedure that hinges upon using the encoding signature to convert the query head to a normal form. In this chapter we extend that approach to encoding queries whose bodies contain hierarchical edges (HCEQs).

We started Chapter 3 by characterizing CEQ validity, after which we characterized encoding-equivalence between any pair of valid CEQs. CEQ validity corresponds to a simple syntactic condition (see Definition 2.3.4 and Lemma 3.0.1); unfortunately, hierarchical edges interact with functional dependencies, making the validity problem non-trivial for arbitrary HCEQs. Because the HCEQ validity problem is orthogonal to our original problem of deciding  $\text{COCQL}^{\mathcal{F}}$  equivalence, when considering HCEQ encoding-equivalence we will restrict our attention to the subclass of HCEQs that can be output by the `ENCODE()` procedure from Section 2.3.2, which are known to be valid (see Theorem 2.3.6).

Our main result in this chapter is a homomorphism-based sufficient condition for HCEQ encoding-equivalence akin to a synthesis of Theorem 3.2.4 (complete characterization of CEQ encoding-equivalence) with Theorem 4.3.9 (sufficient condition for HCQ/DACQ equivalence). In Section 5.1 we enhance the encoding normal form from Chapter 3 to handle HCEQs. In Section 5.2 we prove a homomorphic condition between HCEQs in encoding normal form to be sufficient for implying encoding-equivalence, while in Section 5.3 we give two examples for which the homomorphic condition fails, but for which we can show that encoding-equivalence holds. We are not aware of any literature specific to the equivalence of queries outputting complex objects (or relational encodings of complex objects) containing empty collections, beyond that already discussed in previous chapters, and so this chapter does not include any further survey of literature.

## 5.1 A Normal Form for HCEQs

In Chapter 3 we defined  $\bar{\xi}$ -normal form for CEQs by identifying a set of *core* indexes. This definition of core indexes (Definition 3.1.11) explicitly uses MVDs, and implicitly relies upon the FD  $\bar{\mathcal{I}}_{[1,d]} \rightarrow \bar{\mathcal{V}}$ . However, hierarchical edges (and the subsequent introduction of node labels into the query head) interact with both MVDs and FDs. The following example illustrates this interaction, and relates it to the *nullability dependencies* that we defined in Chapter 4.

**Example 53** Consider the following HCQ  $Q_{52}$ , along with the two CQs  $Q_{53}$  and  $Q_{54}$  corresponding to the two elements of  $\text{dags}(\text{body}_{Q_{52}})$ .

$$\begin{aligned}
 Q_{52}(A^{\boxed{x}}, B^{\boxed{x}}, 0^{\boxed{y}}) &:- \boxed{x} R(A, B) \longrightarrow \boxed{y} R(B, C) \\
 Q_{53}(A, B, 0) &:- R(A, B) \\
 Q_{54}(A, B, 0) &:- R(A, B), R(B, C)
 \end{aligned}$$

Over attribute set  $\mathcal{W} = \{A, B, 0\}$ , the FD  $A \rightarrow 0$  and the MVD  $A \twoheadrightarrow 0$  are both trivial (because 0 is a constant), and so  $Q_{53} \models A \rightarrow 0$ ,  $Q_{53} \models A \twoheadrightarrow 0$ ,  $Q_{54} \models A \rightarrow 0$ , and  $Q_{54} \models A \twoheadrightarrow 0$ . In contrast, over labelled attribute set  $\mathcal{W} = \{A^{\boxed{x}}, B^{\boxed{x}}, 0^{\boxed{y}}\}$ , neither the FD  $A^{\boxed{x}} \rightarrow 0^{\boxed{y}}$  nor the MVD  $A^{\boxed{x}} \twoheadrightarrow 0^{\boxed{y}}$  is trivial, because the database instance  $\mathbb{D}_8 := \{R(a, a), R(a, b)\}$  proves  $Q_{52} \not\models A^{\boxed{x}} \rightarrow 0^{\boxed{y}}$  and  $Q_{52} \not\models A^{\boxed{x}} \twoheadrightarrow 0^{\boxed{y}}$ .

$$\begin{array}{ccc}
 (Q_{52})^{\mathbb{D}_8} & \begin{array}{|c|c|c|} \hline A^{\boxed{x}} & B^{\boxed{x}} & 0^{\boxed{y}} \\ \hline a & a & 0 \\ \hline a & b & \blacksquare \\ \hline \end{array} &
 (Q_{53})^{\mathbb{D}_8} & \begin{array}{|c|c|c|} \hline A & B & 0 \\ \hline a & a & 0 \\ \hline a & b & 0 \\ \hline \end{array} &
 (Q_{54})^{\mathbb{D}_8} & \begin{array}{|c|c|c|} \hline A & B & 0 \\ \hline a & a & 0 \\ \hline \end{array}
 \end{array}$$

The violation of the FD  $A^{\boxed{x}} \rightarrow 0^{\boxed{y}}$  within relation  $(Q_{52})^{\mathbb{D}_8}$  is caused by two terminal embeddings of  $Q_{52}$  that agree on variable  $A$ , but one embedding extends to vertex  $y$  while the other does not. Restated in the terms of nullability dependencies (see Definition 4.2.2), the source of this FD violation is the nullability dependency violation  $Q_{52} \not\models A \twoheadrightarrow y$ . Similarly, the violation of MVD  $A^{\boxed{x}} \twoheadrightarrow 0^{\boxed{y}}$  is related to the fact that the *nullability set*  $\mathcal{S}_{Q_{52}} = \{B\}$  (see Corollary 4.2.7) is not subsumed by the left-hand-side of the MVD.

We now extend the definition of  $\bar{\xi}$ -NF from Chapter 3 to handle HCEQs originating from the `ENCODE()` procedure. Our definition of core indexes relies upon testing whether a CQ $^{\bar{\cdot}}$  satisfies a particular MVD. In Section 3.1.1 we considered the comparable problem for CQs, and those results can be extended to CQ $^{\bar{\cdot}}$ s by accounting for the equality predicates in a straightforward (albeit tedious) manner; hence, we do not discuss that problem further here.

Our definition of the core indexes also relies on nullability set  $\mathcal{S}_Q$  (see Corollary 4.2.7). To facilitate calculating  $\mathcal{S}_Q$  via the theorems in Section 4.2.1, we first convert the query body to DACQ canonical form (assuming  $\bar{\mathcal{W}} := \bar{\mathcal{I}}_{[1,d]} \bar{\mathcal{V}}$  for the



Figure 5.1: Encoding queries  $Q_{55} = \text{CANONICAL}(Q_{20})$  and  $Q_{56} = \text{CANONICAL}(Q_{21})$

purposes of Definition 4.2.1). Conversion to DACQ canonical form may cause the tree-shaped query body to become a DAG, and so the resulting query is no longer a hierarchical CEQ. We also wish to restrict our attention to encoding queries originating from the reduction in Chapter 2; for that reason, we define the following class of encoding queries which we call *derived DACEQs*.

**Definition 5.1.1 (Derived DACEQs)** A directed acyclic conjunctive encoding query (DACEQ) is a generalization of an HCEQ (see Definition 2.3.3) to have a DACQ body (see Definition 4.1.2). We call DACEQ  $Q$  derived if there exists some  $\text{COCQL}^{\{f_b, f_s, f_a\}}$  query  $Q'$  such that  $Q = \text{CANONICAL}(\text{ENCODE}(Q'))$ .

**Example 54** Queries  $Q_{55}$  and  $Q_{56}$  shown in Figure 5.1 are both derived DACEQs, because they are the result of converting HCEQs  $Q_{20}$  and  $Q_{21}$  (cf. Figures 2.18 and 2.19) to DACQ canonical form. (To improve readability, we have merged the labels for each index tuple.)

The following lemma lists properties of derived DACEQs that follow directly from the definitions of procedures  $\text{ENCODE}()$  (see Section 2.3.2) and  $\text{CANONICAL}()$  (see Algorithm 4 in Section 4.2.2). Recall from the definition of encoding queries (Definition 2.3.3) that  $\bar{\mathcal{I}}_i$  is a sequence of labelled indexes, and that we use  $\hat{\mathcal{I}}_i$  and  $\mathcal{I}_i$  to denote the corresponding sets of *labelled* and *unlabelled* indexes, respectively.

**Observation 5.1.2** Given any derived DACEQ  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$ , query  $Q$  satisfies the following properties:

- [P6]  $Q$  is in DACQ canonical form (see Definition 4.2.1).
- [P7] For any  $i \in [1, d]$ , index set  $\hat{\mathcal{I}}_i$  is either
  - (a) a subset of  $\{\blacksquare\}$  (i.e.,  $\{\}$  or  $\{\blacksquare\}$ ), or
  - (b) a non-empty set of variables and constants all labelled with the same node  $v_i \in \text{body}_Q$ . It is possible that  $v_i = v_j$  for  $i \neq j$ .

- [P8] For any variable  $X \in \mathcal{I}_{[1,d]}$ , let  $i \in [1,d]$  be the smallest integer such that  $X \in \mathcal{I}_i$ . Then,  $X \in \mathcal{L}_{v_i}$ .
- [P9] For any  $i \in [1,d]$  such that  $\hat{\mathcal{I}}_i \not\subseteq \{\blacksquare\}$ , the parameters of  $\text{body}_{v_i}$  satisfy  $\mathcal{P}_{v_i} \subseteq \mathcal{I}_{[1,i-1]}$ , implying  $Q \models \mathcal{I}_{[1,i-1]} \Rightarrow v_i$  by Lemma 4.2.3.
- [P10] For every node  $v \in \text{body}_Q$ , there exists at least one level  $i \in [1,d]$  such that  $v = v_i$ .

**Example 55** Continuing Example 54, the reader can verify that queries  $Q_{55}$  and  $Q_{56}$  each satisfy all of the properties [P6]–[P10]. Notably, property [P7] guarantees that any valuation of index tuple  $\bar{\mathcal{I}}_i$  is either entirely  $\blacksquare$  values or does not include any  $\blacksquare$  value. Property [P9] guarantees that for any valuation of index tuples  $\bar{\mathcal{I}}_{[1,i-1]}$ , the possible valuations of index tuple  $\bar{\mathcal{I}}_i$  cannot include both the entirely  $\blacksquare$  tuple and a tuple not containing  $\blacksquare$ . For a concrete example, the reader is invited to refer back to encoding relations  $R_8 = (Q_{55})^{\mathbb{D}^3}$  and  $R_9 = (Q_{56})^{\mathbb{D}^3}$  in Figures 2.20 and 2.21, respectively. (Proposition 4.2.13 guarantees  $Q_{55} \equiv Q_{20}$  and  $Q_{56} \equiv Q_{21}$ .)

**Definition 5.1.3 (Core Indexes)** Given any derived DACEQ  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  with nullability set  $\mathcal{S}_Q$ , a length- $d$  signature  $\bar{\mathcal{S}}$ , and an integer  $i \in [1,d]$ , for each  $G \in \text{dags}(\text{body}_Q)$  satisfying  $v_i \in G$ , define set  $\mathcal{W}_{(G,i)} := \mathcal{B}_G^+ \cap (\mathcal{I}_{[1,i]} \cup \mathcal{I}_{[i+1,d]}^{\bar{\mathcal{S}}})$ , and let  $Q_{(G,i)}$  be the following CQ $^\pm$ .

$$Q_{(G,i)}(\mathcal{W}_{(G,i)}) := \text{body}_G^+ \quad (5.1)$$

Define the set of core indexes at level  $i$  relative to  $\bar{\mathcal{S}}$ —denoted  $\hat{\mathcal{I}}_i^{\bar{\mathcal{S}}}$ —as a minimal subset of  $\hat{\mathcal{I}}_i$  such that  $\mathcal{I}_i^{\bar{\mathcal{S}}}$  (the unlabelled version the indexes in  $\hat{\mathcal{I}}_i^{\bar{\mathcal{S}}}$ —see Section 2.3.1) satisfies both  $(\mathcal{V} \cup \mathcal{S}_Q) \subseteq \mathcal{I}_{[1,i-1]} \cup \mathcal{I}_{[i,d]}^{\bar{\mathcal{S}}} \cup \mathcal{C}$  and the following signature-specific conditions.

$\mathcal{S}_i$	Condition
<b>b</b>	$\mathcal{I}_i \subseteq \mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{\mathcal{S}}} \cup \mathcal{C}$
<b>s</b>	$\forall G \text{ containing } v_i : \left[ Q_{(G,i)} \models \mathcal{B}_G^+ \cap (\mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{\mathcal{S}}}) \Rightarrow \mathcal{B}_G^+ \cap \mathcal{I}_{[i+1,d]}^{\bar{\mathcal{S}}} \right]$
<b>n</b>	$\forall G \text{ containing } v_i : \left[ Q_{(G,i)} \models \mathcal{B}_G^+ \cap \mathcal{I}_{[1,i-1]} \Rightarrow \mathcal{B}_G^+ \cap \mathcal{I}_{[i,d]}^{\bar{\mathcal{S}}} \right]$

We denote by  $\bar{\mathcal{I}}_i^{\bar{\mathcal{S}}}$  the tuple formed by deleting from  $\bar{\mathcal{I}}_i$  all attributes not in  $\hat{\mathcal{I}}_i^{\bar{\mathcal{S}}}$ .

**Lemma 5.1.4** Definition 5.1.3 determines a unique minimal set  $\hat{\mathcal{I}}_i^{\bar{\mathcal{S}}}$ .

**PROOF (SKETCH).** Identical to the proof of Lemma 3.1.15, except that the MVD arguments for cases  $\mathcal{S}_i = \mathbf{s}$  and  $\mathcal{S}_i = \mathbf{n}$  need to be repeated for each  $Q_{(G,i)}$  (equation 5.1).  $\square$

**Example 56** Continuing Examples 54 and 55, Query  $Q_{55}$  satisfies  $\{F\} \Rightarrow 4$  and  $\{D\} \Rightarrow 3$ , and so  $\mathcal{S}_{Q_{55}} = \{D, F\}$ . The core indexes of  $Q_{55}$  with respect to encoding signature  $\bar{\mathcal{S}} = \text{snsb}$  are calculated as follows.

- $i = 4$ :  $\mathfrak{I}_4 = \mathfrak{b}$ ,  $\hat{\mathcal{I}}_4^{\bar{\mathfrak{s}}} = \{H^{\bar{4}}, J^{\bar{4}}\}$   
 $c^{\bar{4}} \notin \hat{\mathcal{I}}_4^{\bar{\mathfrak{s}}}$  because  $c \in \mathcal{C}$ ;  $F^{\bar{4}} \notin \hat{\mathcal{I}}_4^{\bar{\mathfrak{s}}}$  because  $F \in \mathcal{I}_3$ ;  $\{H, J\} \subseteq \mathcal{I}_4^{\bar{\mathfrak{s}}}$  in order to satisfy condition  $\mathcal{I}_4 = \{c, F, H, J\} \subseteq \mathcal{I}_{[1,3]} \cup \mathcal{I}_4^{\bar{\mathfrak{s}}} \cup \mathcal{C}$ .
- $i = 3$ :  $\mathfrak{I}_3 = \mathfrak{s}$ ,  $\hat{\mathcal{I}}_3^{\bar{\mathfrak{s}}} = \{F^{\bar{12}}\}$   
 $B^{\bar{12}} \notin \hat{\mathcal{I}}_3^{\bar{\mathfrak{s}}}$  because  $B \in \mathcal{I}_2$ ;  $F^{\bar{12}} \in \hat{\mathcal{I}}_3^{\bar{\mathfrak{s}}}$  because  $F \in \mathcal{S}_{Q_{55}} \setminus (\mathcal{I}_{[1,2]} \cup \mathcal{I}_4^{\bar{\mathfrak{s}}} \cup \mathcal{C})$ .
- $i = 2$ :  $\mathfrak{I}_2 = \mathfrak{n}$ ,  $\hat{\mathcal{I}}_2^{\bar{\mathfrak{s}}} = \{L^{\bar{3}}, N^{\bar{3}}\}$   
 $L^{\bar{3}} \in \hat{\mathcal{I}}_2^{\bar{\mathfrak{s}}}$  because  $L \in \mathcal{V} \setminus (\mathcal{I}_1 \cup \mathcal{I}_{[3,4]}^{\bar{\mathfrak{s}}} \cup \mathcal{C})$ ;  $B^{\bar{3}} \notin \hat{\mathcal{I}}_2^{\bar{\mathfrak{s}}}$  because  $B \in \mathcal{I}_1$ ;  $P^{\bar{3}} \notin \hat{\mathcal{I}}_2^{\bar{\mathfrak{s}}}$  due to the MVDs  $Q_{(123,2)} \models \mathfrak{a}ABD \twoheadrightarrow LNF$  and  $Q_{(1234,2)} \models \mathfrak{a}ABD \twoheadrightarrow LNFHJ$ , while  $Q_{(123,2)}$  requires  $N^{\bar{3}} \in \hat{\mathcal{I}}_2^{\bar{\mathfrak{s}}}$  in order to satisfy the necessary MVD.

$$Q_{(123,2)}(\mathfrak{a}, A, B, D, L, N, P, F) :- R(A, B), R(B, D), R(B, F), \\ S(D, L), S(L, N), R(B, P)$$

$$Q_{(1234,2)}(\mathfrak{a}, A, B, D, L, N, P, F, H, J) :- R(A, B), R(B, D), R(B, F), S(D, L), \\ S(L, N), R(B, P), S(F, H), S(H, J)$$

- $i = 1$ :  $\mathfrak{I}_1 = \mathfrak{s}$ ,  $\hat{\mathcal{I}}_1^{\bar{\mathfrak{s}}} = \{B^{\bar{12}}, D^{\bar{12}}\}$   
 $a^{\bar{12}} \notin \hat{\mathcal{I}}_1^{\bar{\mathfrak{s}}}$  because  $a \in \mathcal{C}$ ;  $B^{\bar{12}} \in \hat{\mathcal{I}}_1^{\bar{\mathfrak{s}}}$  because  $B \in \mathcal{V} \setminus (\mathcal{I}_{[2,4]}^{\bar{\mathfrak{s}}} \cup \mathcal{C})$ ;  $D^{\bar{12}} \in \hat{\mathcal{I}}_1^{\bar{\mathfrak{s}}}$  because  $D \in \mathcal{S}_{Q_{55}} \setminus (\mathcal{I}_{[2,4]}^{\bar{\mathfrak{s}}} \cup \mathcal{C})$ ;  $A^{\bar{12}} \notin \hat{\mathcal{I}}_1^{\bar{\mathfrak{s}}}$  due to the MVDs  $Q_{(12,1)} \models BD \twoheadrightarrow F$ ,  $Q_{(123,1)} \models BD \twoheadrightarrow LNF$ ,  $Q_{(124,1)} \models BD \twoheadrightarrow FHJ$ , and  $Q_{(1234,1)} \models BD \twoheadrightarrow LNFHJ$ .

$$Q_{(12,1)}(\mathfrak{a}, A, B, D, F) :- R(A, B), R(B, D), R(B, F)$$

$$Q_{(123,1)}(\mathfrak{a}, A, B, D, L, N, F) :- R(A, B), R(B, D), R(B, F), \\ S(D, L), S(L, N), R(B, P)$$

$$Q_{(124,1)}(\mathfrak{a}, A, B, D, F, H, J) :- R(A, B), R(B, D), R(B, F), \\ S(F, H), S(H, J)$$

$$Q_{(1234,1)}(\mathfrak{a}, A, B, D, L, N, F, H, J) :- R(A, B), R(B, D), R(B, F), S(D, L), \\ S(L, N), R(B, P), S(F, H), S(H, J)$$

Given derived DACEQ  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$ , we convert  $Q$  to  $\bar{\mathfrak{s}}$ -normal form ( $\bar{\mathfrak{s}}$ -NF) by deleting all non-core indexes from the query head. Performing  $\bar{\mathfrak{s}}$ -normalization of derived DACEQs is clearly NP-hard, as Definition 5.1.3 directly generalizes the CEQ case. We can also derive a (possibly loose) upper bound directly from Definition 5.1.3.

**Theorem 5.1.5** *Performing  $\bar{\mathfrak{s}}$ -normalization of derived DACEQs is in  $\Pi_2^P$ .*

PROOF. For each level  $i$  and each dag  $G \in \text{dags}(\text{body}_Q)$ , the MVD associated with query  $Q_{(G,i)}$  (equation 5.1) dictates that a certain subset of  $\mathcal{I}_i$  needs to be core, and this subset can be identified in NP time (per Theorem 3.1.16). Therefore, verifying each core  $I \in \hat{\mathcal{I}}_{[1,d]}^{\bar{\mathfrak{s}}}$  is in  $\Sigma_2^P$  (by verifying that either  $I \in \mathcal{V} \cup \mathcal{S}_Q$  or by guessing the appropriate query  $Q_{(G,i)}$  and verifying in NP time that its MVD

requires  $I$  to be core). Hence, verifying the non-core attributes  $\mathcal{I}_i \setminus \hat{\mathcal{I}}_{[1,d]}^{\bar{\S}}$  is in  $\text{co-}\Sigma_2^p = \Pi_2^p$ .  $\square$

**Conjecture 5.1.6** *Performing  $\bar{\S}$ -normalization of derived DACEQs is in NP.*

COMMENT. It would suffice to show

$$\left\{ \forall v \in \text{body}_Q. \left[ Q_{(v,i)} \models \mathcal{B}_v^+ \cap (\mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{\S}}) \rightarrow \mathcal{B}_v^+ \cap \mathcal{I}_{[i+1,d]}^{\bar{\S}} \right] \right\} \implies \left\{ \forall G \text{ containing } v_i : \left[ Q_{(G,i)} \models \mathcal{B}_G^+ \cap (\mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{\S}}) \rightarrow \mathcal{B}_G^+ \cap \mathcal{I}_{[i+1,d]}^{\bar{\S}} \right] \right\} \quad (5.2)$$

and

$$\left\{ \forall v \in \text{body}_Q. \left[ Q_{(v,i)} \models \mathcal{B}_v^+ \cap \mathcal{I}_{[1,i-1]} \rightarrow \mathcal{B}_v^+ \cap \mathcal{I}_{[i,d]}^{\bar{\S}} \right] \right\} \implies \left\{ \forall G \text{ containing } v_i : \left[ Q_{(G,i)} \models \mathcal{B}_G^+ \cap \mathcal{I}_{[1,i-1]} \rightarrow \mathcal{B}_G^+ \cap \mathcal{I}_{[i,d]}^{\bar{\S}} \right] \right\} \quad (5.3)$$

where  $Q_{(v,i)}$  with body  $\text{body}_v^+$  is defined analogously to  $Q_{(G,i)}$  from equation 5.1. Then, for each level  $i$  it is only necessary to examine each query  $Q_{(v,i)}$ , and for each such query the hypergraph-based algorithm from Theorem 3.1.16 would identify a subset of  $\hat{\mathcal{I}}_i$  that must be core.  $\hat{\mathcal{I}}_i^{\bar{\S}}$  would then be the union of these sets.

Unfortunately, due to the intricacy of the various mappings involved, we have been unable to definitively prove (or disprove) equations 5.2 and 5.3.  $\square$

**Lemma 5.1.7** *Given any derived DACEQ  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$ , any index level  $i \in [1, d]$ , and any node  $v \in \text{body}_Q$ : if  $Q \models \mathcal{I}_{[1,i]} \Rightarrow v$  then  $Q \models \mathcal{I}_{[1,i]}^{\bar{\S}} \Rightarrow v$ .*

PROOF. Definition 5.1.3 implies that  $(\mathcal{S}_Q \cap \mathcal{I}_{[1,i]}) \subseteq \mathcal{I}_{[1,i]}^{\bar{\S}}$ . By Corollary 4.2.7, if  $Q \models \mathcal{I}_{[1,i]} \Rightarrow v$  then  $Q \models (\mathcal{S}_Q \cap \mathcal{I}_{[1,i]}) \Rightarrow v$ .  $\square$

**Theorem 5.1.8**  *$\bar{\S}$ -Normalization preserves  $\bar{\S}$ -equivalence*

PROOF. For every  $i \in [0, d]$ , let  $Q_i$  denote the following HCEQ,

$$Q_i(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_i; \bar{\mathcal{I}}_{i+1}^{\bar{\S}}; \dots; \bar{\mathcal{I}}_d^{\bar{\S}}; \bar{\mathcal{V}}) :- \text{body}_Q \quad (5.4)$$

and assume that  $\bar{\mathcal{I}}_i = \bar{\mathcal{I}}_i^{\bar{\S}} \cdot \bar{\mathcal{J}}_i$ , where  $\hat{\mathcal{J}}_i$  is the set of non-core indexes in  $\hat{\mathcal{I}}_i$ . The induction proceeds exactly as in the proof of Theorem 3.1.17. The only complication is that the conditions in Definition 5.1.3 are stated in terms of the query  $Q_{(G,i)}$  which is a  $\text{CQ}^-$ , whereas query  $Q_i$  has a hierarchical body. (In contrast, in the proof of Theorem 3.1.17, queries  $\hat{Q}_i$  and  $Q_i$  had the same body.) We prove here for each case within the inductive step that the conditions on the output of  $\text{CQ}^-$   $Q_{(G,i)}$  in Definition 5.1.3 imply corresponding conditions for the output of HCEQ  $Q_i$ .

**Case  $\bar{\S}_i = \text{b}$ :**

It suffices to show that for any database  $\mathbb{D}$ , the relation  $\Phi := (Q_i)^{\mathbb{D}}$  satisfies the FD  $\hat{\mathcal{I}}_{[1,i-1]} \rightarrow \hat{\mathcal{J}}_i$ ; the remainder then follows identically to the proof of Theorem 3.1.17. By Definition 3.1.11, either



1.  $\mathcal{J}_i \subseteq \{\blacksquare\}$ , or
2.  $\mathcal{J}_i \subseteq \mathcal{I}_{[1,i-1]} \cup \mathcal{C}$ .

In the second case, property **[P7]** (Observation 5.1.2) implies that all elements of  $\hat{\mathcal{J}}_i$  are labelled with  $v_i$  satisfying  $Q \models \mathcal{I}_{[1,i-1]} \Rightarrow v_i$ ; therefore, in either case the relation  $\Phi$  satisfies the FD  $\hat{\mathcal{I}}_{[1,i-1]} \rightarrow \hat{\mathcal{J}}_i$ .

**Case  $\S_i = \mathbf{s}$ :**

It suffices to show that relation  $\Phi$  satisfies the MVD  $\hat{\mathcal{I}}_{[1,i-1]} \cup \hat{\mathcal{I}}_i^{\bar{\bar{\mathbb{S}}}} \twoheadrightarrow \hat{\mathcal{I}}_{[i+1,d]}^{\bar{\bar{\mathbb{S}}}} \cup \hat{\mathcal{V}}$ ; the remainder then follows identically to the proof of Theorem 3.1.17.

By property **[P7]**, either  $\hat{\mathcal{I}}_i \subseteq \{\blacksquare\}$  or all elements of  $\hat{\mathcal{I}}_i$  are labelled with  $v_i$ . If  $\hat{\mathcal{I}}_i \subseteq \{\blacksquare\}$  then attributes  $\bar{\mathcal{I}}_i$  are constant within relation  $\Phi$  and the MVD holds trivially; therefore, assume that some  $v_i$  labels all elements of  $\hat{\mathcal{I}}_i$ .

Let  $\phi_1 : G_1 \rightsquigarrow \mathbb{D}$  and  $\phi_2 : G_2 \rightsquigarrow \mathbb{D}$  be any two terminal embeddings in  $\Phi$ —of  $G_1, G_2 \in \text{dags}(\text{body}_Q)$ , respectively—that satisfy  $\hat{\phi}_1(\bar{\mathcal{I}}_{[1,i-1]} \bar{\mathcal{I}}_i^{\bar{\bar{\mathbb{S}}}}) = \hat{\phi}_2(\bar{\mathcal{I}}_{[1,i-1]} \bar{\mathcal{I}}_i^{\bar{\bar{\mathbb{S}}}})$ . (Recall from equation 2.22 that  $\hat{\phi}$  denotes the partial application of  $\phi$ .) For each  $j \in [1, i]$  such that  $\hat{\mathcal{I}}_j \not\subseteq \{\blacksquare\}$ ,  $v_j \in G_1$  iff  $v_j \in G_2$ . If  $v_i \notin G_1$  then  $\hat{\phi}_1(\bar{\mathcal{J}}_i) = \hat{\phi}_2(\bar{\mathcal{J}}_i)$  is a constant tuple composed only of the  $\blacksquare$  symbol, and so the MVD holds trivially; therefore, assume that  $v_i \in G_1$  and  $v_i \in G_2$ .

Define  $G_3 = G_1 \cap G_2$ .<sup>1</sup> Because

$$Q_{(G_3,i)} \models \mathcal{B}_{G_3}^+ \cap (\mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{\bar{\mathbb{S}}}}) \twoheadrightarrow \mathcal{B}_{G_3}^+ \cap \mathcal{I}_{[i+1,d]}^{\bar{\bar{\mathbb{S}}}}$$

and  $\phi_1, \phi_2$  are embeddings of  $\text{body}_{G_3}^+$  into  $\mathbb{D}$  that agree on  $\mathcal{B}_{G_3}^+ \cap (\mathcal{I}_{[1,i-1]} \cap \mathcal{I}_i^{\bar{\bar{\mathbb{S}}}})$ , there exists an embedding  $\phi_3 : G_3 \rightarrow \mathbb{D}$  satisfying

- $\hat{\phi}_3(\bar{\mathcal{I}}_{[1,i-1]} \bar{\mathcal{I}}_i^{\bar{\bar{\mathbb{S}}}}) = \hat{\phi}_2(\bar{\mathcal{I}}_{[1,i-1]} \bar{\mathcal{I}}_i^{\bar{\bar{\mathbb{S}}}}) = \hat{\phi}_1(\bar{\mathcal{I}}_{[1,i-1]} \bar{\mathcal{I}}_i^{\bar{\bar{\mathbb{S}}}})$ ,
- $\hat{\phi}_3(\bar{\mathcal{J}}_i) = \hat{\phi}_2(\bar{\mathcal{J}}_i)$  which, by  $\mathcal{I}_{[1,i-1]} \Rightarrow v_i$ , implies  $\hat{\phi}_3(\hat{\mathcal{J}}_i) = \hat{\phi}_2(\hat{\mathcal{J}}_i)$ , and
- $\forall j \in [i+1, d]$ , if  $v_j \notin (G_1 \setminus G_3)$  then  $\hat{\phi}_3(\bar{\mathcal{I}}_j^{\bar{\bar{\mathbb{S}}}}) = \hat{\phi}_1(\bar{\mathcal{I}}_j^{\bar{\bar{\mathbb{S}}}})$ .

Now consider the smallest  $j \in [i+1, d]$  such that  $v_j \in (G_1 \setminus G_3)$ , and define  $G_4 := G_3 \cup \{v_j\}$ . By property **[P9]** and Lemma 5.1.7,  $Q \models \mathcal{I}_{[1,i-1]} \cup \mathcal{I}_{[i,j-1]}^{\bar{\bar{\mathbb{S}}}} \Rightarrow v_j$ . By  $\hat{\phi}_3(\bar{\mathcal{I}}_{[1,i-1]} \bar{\mathcal{I}}_{[i,j-1]}^{\bar{\bar{\mathbb{S}}}}) = \hat{\phi}_1(\bar{\mathcal{I}}_{[1,i-1]} \bar{\mathcal{I}}_{[i,j-1]}^{\bar{\bar{\mathbb{S}}}})$  and the definition of nullability dependencies (Definition 4.2.2),  $\phi_3$  must extend to  $G_4$ . We now iterate the above argument, appealing to the following MVD

$$Q_{(G_4,i)} \models \mathcal{B}_{G_4}^+ \cap (\mathcal{I}_{[1,i-1]} \cup \mathcal{I}_i^{\bar{\bar{\mathbb{S}}}}) \twoheadrightarrow \mathcal{B}_{G_4}^+ \cap \mathcal{I}_{[i+1,d]}^{\bar{\bar{\mathbb{S}}}}$$

and to  $\phi_1, \phi_3$  being embeddings of  $\text{body}_{G_4}^+$  into  $\mathbb{D}$  agreeing on  $\mathcal{B}_{G_4}^+ \cap (\mathcal{I}_{[1,i-1]} \cap \mathcal{I}_i^{\bar{\bar{\mathbb{S}}}})$  to deduce the existence of some embedding  $\phi_4 : G_4 \rightarrow \mathbb{D}$ . By induction, we eventually arrive at some terminal embedding  $\phi_k : G_1 \rightsquigarrow \mathbb{D}$  satisfying

---

<sup>1</sup>We are abusing notation by using  $G$  to refer to both a DAG within  $\text{dags}(\text{body}_Q)$  and the set of node labels occurring within that DAG.

- $\hat{\phi}_k(\overline{\mathcal{I}}_{[1,i-1]}\overline{\mathcal{I}}_i^{\overline{\mathfrak{s}}}) = \hat{\phi}_2(\overline{\mathcal{I}}_{[1,i-1]}\overline{\mathcal{I}}_i^{\overline{\mathfrak{s}}}) = \hat{\phi}_1(\overline{\mathcal{I}}_{[1,i-1]}\overline{\mathcal{I}}_i^{\overline{\mathfrak{s}}})$ ,
- $\hat{\phi}_k(\hat{\mathcal{J}}_i) = \hat{\phi}_2(\hat{\mathcal{J}}_i)$ , and
- $\forall j \in [i+1, d], \hat{\phi}_k(\overline{\mathcal{I}}_{[i+1,d]}^{\overline{\mathfrak{s}}}) = \hat{\phi}_1(\overline{\mathcal{I}}_{[i+1,d]}^{\overline{\mathfrak{s}}})$ .

Embedding  $\phi_k$  proves that that relation  $\Phi$  satisfies the desired MVD.

**Case  $\mathfrak{s}_i = \mathbf{n}$ :**

It suffices to show that relation  $\Phi$  satisfies the MVD  $\hat{\mathcal{I}}_{[1,i-1]} \rightarrow \hat{\mathcal{I}}_{[i,d]}^{\overline{\mathfrak{s}}} \cup \hat{\mathcal{V}}$ . The proof is identical to case  $\mathfrak{s}_i = \mathbf{s}$ , with appropriate modifications to the variables in the MVD. The remainder then follows identically to the proof of Theorem 3.1.17. □

**Example 57** Continuing Example 56, query  $Q_{57}$  in Figure 5.2a is the snsb-normal form of  $Q_{55}$  and therefore  $Q_{57} \dot{\equiv}_{\text{snsb}} Q_{55} \dot{\equiv}_{\text{snsb}} Q_{20}$ . Query  $Q_{58}$  in Figure 5.2b is the snsb-normal form of  $Q_{56}$  and therefore  $Q_{58} \dot{\equiv}_{\text{snsb}} Q_{56} \dot{\equiv}_{\text{snsb}} Q_{21}$ . Figure 5.2c shows the encoding relation  $R_{16} = (Q_{57})^{\mathbb{D}_3}$ , which can be obtained by projecting away the non-core index attributes from encoding relation  $R_8 = (Q_{20})^{\mathbb{D}_3}$  (cf. Figure 2.20). Figure 5.2d shows the encoding relation  $R_{17} = (Q_{58})^{\mathbb{D}_3}$ , which can be obtained by projecting away the non-core index attributes from encoding relation  $R_9 = (Q_{21})^{\mathbb{D}_3}$  (cf. Figure 2.21); observe that  $R_9$  has the identical instance as  $R_8$ . In the next section (Examples 58 and 59) we will prove that over *any* database, the pair of encoding relations yielded by queries  $Q_{57}$  and  $Q_{58}$  will have identical instances.

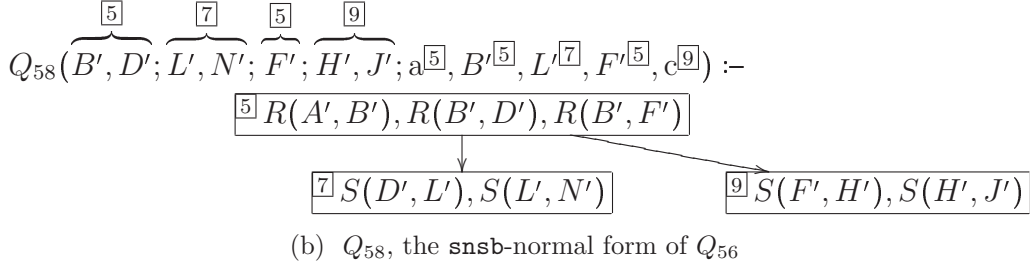
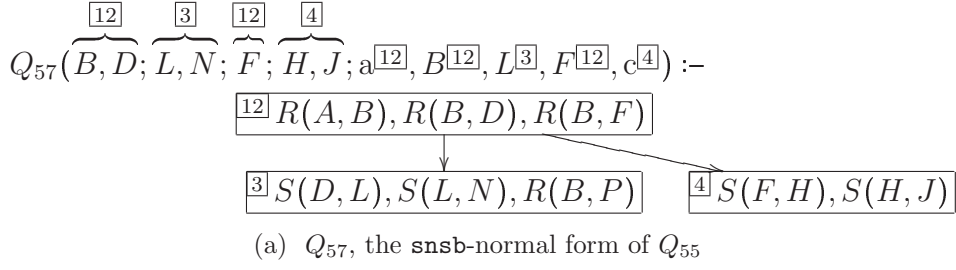
## 5.2 Encoding-Equivalence of Derived DACEQs

In this section we partially characterize  $\overline{\mathfrak{s}}$ -equivalence between derived DACEQs. We start by deriving a necessary condition that relates  $\overline{\mathfrak{s}}$ -equivalence to DACQ equivalence, whose corollary is a necessary isomorphic condition between DACEQ query bodies analogous to the head graph isomorphisms of Section 4.3.1. We then extend the definition of index-covering homomorphisms (see Chapter 3) to derived DACEQs, and show that the existence of index-covering homomorphisms between the  $\overline{\mathfrak{s}}$ -normal forms constitutes a sufficient condition for  $\overline{\mathfrak{s}}$ -equivalence.

**Definition 5.2.1 (Structural DACQ)** *Given derived DACEQ  $Q(\overline{\mathcal{I}}_1; \dots; \overline{\mathcal{I}}_d; \overline{\mathcal{V}})$ , for each  $i \in [1, d]$  define the attribute  $X_i$  as follows.*

$$X_i := \begin{cases} \blacksquare & \text{if } \hat{\mathcal{I}}_i \subseteq \{\blacksquare\} \\ 0^{v_i} & \text{otherwise } (v_i \text{ defined as in property [P7] of Observation 5.1.2}) \end{cases}$$

We define the structural query  $\text{STRUCTURAL}(Q)$  as the DACQ with the same body as  $Q$  and with the head tuple  $\overline{\mathcal{W}} := \overline{X}_{[1,d]}\overline{\mathcal{V}}$ .


 $R_{16}$ 

$B^{[12]}$	$D^{[12]}$	$L^{[3]}$	$N^{[3]}$	$F^{[12]}$	$H^{[4]}$	$J^{[4]}$	$a^{[12]}$	$B^{[12]}$	$L^{[3]}$	$F^{[12]}$	$c^{[4]}$
2	3	4	5	3	4	5	a	2	4	3	c
2	3	4	5	3	4	6	a	2	4	3	c
2	3	4	5	7	■	■	a	2	4	7	■
2	3	4	6	3	4	5	a	2	4	3	c
2	3	4	6	3	4	6	a	2	4	3	c
2	3	4	6	7	■	■	a	2	4	7	■
2	7	■	■	3	4	5	a	2	■	3	c
2	7	■	■	3	4	6	a	2	■	3	c
2	7	■	■	7	■	■	a	2	■	7	■

(c) Query result  $R_{16} = (Q_{57})^{\mathbb{D}_3}$  satisfying  $\text{DECODE}(R_{16}, \text{snsb}) = o_{13}$

 $R_{17}$ 

$B'^{[5]}$	$D'^{[5]}$	$L'^{[7]}$	$N'^{[37]}$	$F'^{[5]}$	$H'^{[9]}$	$J'^{[9]}$	$a^{[5]}$	$B'^{[5]}$	$L'^{[7]}$	$F'^{[5]}$	$c^{[9]}$
2	3	4	5	3	4	5	a	2	4	3	c
2	3	4	5	3	4	6	a	2	4	3	c
2	3	4	5	7	■	■	a	2	4	7	■
2	3	4	6	3	4	5	a	2	4	3	c
2	3	4	6	3	4	6	a	2	4	3	c
2	3	4	6	7	■	■	a	2	4	7	■
2	7	■	■	3	4	5	a	2	■	3	c
2	7	■	■	3	4	6	a	2	■	3	c
2	7	■	■	7	■	■	a	2	■	7	■

(d) Query result  $R_{17} = (Q_{58})^{\mathbb{D}_3}$  satisfying  $\text{DECODE}(R_{17}, \text{snsb}) = o_{13}$

Figure 5.2: The **snsb**-normal forms of  $Q_{55}$  and  $Q_{56}$  and their results over  $\mathbb{D}_3$

**Theorem 5.2.2** *Given any two derived DACEQs  $Q$  and  $Q'$  and any encoding signature  $\bar{\xi}$ , if  $Q \dot{\equiv}_{\bar{\xi}} Q'$  then  $\text{STRUCTURAL}(Q) \equiv \text{STRUCTURAL}(Q')$ .*

PROOF. Given any database  $\mathbb{D}$ , let  $\Gamma$  denote the set of terminal embeddings of  $\text{body}_Q$  into  $\mathbb{D}$ , and let  $\Phi$  be the analogous set for  $\text{body}_{Q'}$ . Suppose that  $Q \dot{\equiv}_{\bar{\xi}} Q'$ ; then, there must exist some  $\bar{\xi}$ -certificate  $C$  proving  $(Q)^{\mathbb{D}} \dot{\equiv}_{\bar{\xi}} (Q')^{\mathbb{D}}$ . Consider any  $\gamma \in \Gamma$ —then, there must exist some  $\phi \in \Phi$  such that certificate  $C$  contains a tuple node  $n^\dagger$  that proves  $\hat{\gamma}(\bar{\mathcal{V}}) = \hat{\phi}(\bar{\mathcal{V}}')$ . By Definitions 2.2.5–2.2.7, the certificate node at any level  $i \in [1, d]$  along the path to  $n^\dagger$  must satisfy  $\hat{\gamma}(\bar{\mathcal{I}}_i) = \#^{|\bar{\mathcal{I}}_i|} \iff \hat{\phi}(\bar{\mathcal{I}}'_i) = \#^{|\bar{\mathcal{I}}'_i|}$ . It follows that  $\hat{\gamma}(\bar{X}_{[1,d]}\bar{\mathcal{V}}) = \hat{\phi}(\bar{X}'_{[1,d]}\bar{\mathcal{V}}')$ , which proves  $\text{STRUCTURAL}(Q) \equiv \text{STRUCTURAL}(Q')$ . Repeating the argument in the other direction proves  $\text{STRUCTURAL}(Q') \equiv \text{STRUCTURAL}(Q)$ .  $\square$

**Corollary 5.2.3** *If  $Q \dot{\equiv}_{\bar{\xi}} Q'$  then  $\text{STRUCTURAL}(Q)$  and  $\text{STRUCTURAL}(Q')$  are head graph isomorphic.*

PROOF. Theorem 4.3.4.  $\square$

**Corollary 5.2.4** *Given any derived DACEQ  $Q$  and any signature  $\bar{\xi}$ , the  $\bar{\xi}$ -normal form of  $Q$  contains within its head the label of every node in  $\text{body}_Q$ .*

PROOF. Property [P10] of Observation 5.1.2 implies that the label of every node in  $\text{body}_Q$  appears within the head graph of  $\text{STRUCTURAL}(Q)$ , and so the claim follows by Theorem 5.1.8 and Corollary 5.2.3.  $\square$

Our previous definition of index-covering homomorphisms between CEQs (Definition 3.2.1) was a simple generalization of homomorphisms between CQs (Definition 3.1.1). We now extend the definition of index-covering homomorphisms to handle DACEQs. Unlike CEQ bodies, DACEQ bodies may contain equality predicates; therefore, it is helpful to first review how the definition of CQ homomorphisms can be extended to to handle CQ $^\bar{\cdot}$ s.

**Definition 5.2.5 (CQ $^\bar{\cdot}$  Homomorphism)** *Given two satisfiable CQ $^\bar{\cdot}$ s  $Q(\bar{\mathcal{W}})$  and  $Q'(\bar{\mathcal{W}}')$ , a homomorphism  $h$  from  $Q'$  to  $Q$  is a mapping  $h : \mathcal{B}' \rightarrow \mathcal{B} \cup \mathcal{C}$  satisfying*

1.  $\text{body}_Q \equiv h(\text{body}_{Q'})$ ; in other words,
  - (a)  $\text{pred}_Q \equiv h(\text{pred}_{Q'})$ , and
  - (b)  $\forall R(\bar{X}') \in \text{atoms}_{Q'}. \exists R(\bar{X}) \in \text{atoms}_Q. [\text{pred}_Q \equiv h(\bar{X}') = \bar{X}]$ ; and
2. for each attribute  $W'_i \in \bar{\mathcal{W}}'$ :  $\text{pred}_Q \equiv h(W'_i) = W_i$ , where  $W_i$  is the corresponding attribute in  $\bar{\mathcal{W}}$ .

The reader is invited to verify that in the restricted case where the queries do not contain equality predicates, Definition 5.2.5 is equivalent to Definition 3.1.1.

**Definition 5.2.6 (Index-Covering Homomorphism)** *Given any two DACEQs  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  and  $Q'(\bar{\mathcal{I}}'_1; \dots; \bar{\mathcal{I}}'_d; \bar{\mathcal{V}}')$  such that  $\text{STRUCTURAL}(Q')$  and  $\text{STRUCTURAL}(Q)$*

are head graph isomorphic, for each node  $u \in \text{body}_Q$  let  $u'$  denote the corresponding node in  $\text{body}_{Q'}$ , and vice-versa. Then, an index-covering homomorphism from  $Q'$  to  $Q$  is a mapping  $h : \mathcal{B}' \rightarrow \mathcal{B} \cup \mathcal{C}$  satisfying

1. for each node  $u' \in \text{body}_{Q'}$ ,  $\text{body}_u^+ \models \text{body}_{u'}$ ; in other words,
  - (a)  $\text{pred}_u^+ \models h(\text{pred}_{u'})$ , and
  - (b)  $\forall R(\overline{X}') \in \text{atoms}_{u'} \cdot \exists R(\overline{X}) \in \text{atoms}_u^+ \cdot [\text{pred}_u^+ \models h(\overline{X}') = \overline{X}]$ ;
2. for each attribute  $V_i^{\square u'} \in \overline{\mathcal{V}}' : \text{pred}_u^+ \models h(V_i^{\square u'}) = V_i$ , where  $V_i^{\square u'}$  is the corresponding attribute in  $\overline{\mathcal{V}}'$ ; and
3.  $\forall i \in [1, d] : \hat{\mathcal{I}}_i \subseteq h(\hat{\mathcal{I}}'_i)$ .

Again, the reader is invited to verify that in the restricted case where the query bodies contain only a single node (i.e., both  $Q$  and  $Q'$  are CEQs), Definition 5.2.6 is equivalent to Definition 3.2.1.

**Example 58** The following mappings  $h : Q_{57} \rightarrow Q_{58}$  and  $h' : Q_{58} \rightarrow Q_{57}$  between the queries in Figure 5.2 are both index-covering homomorphisms.

$$\begin{aligned} h &= \{A/A', B/B', D/D', F/F', L/L', N/N', P/F', H/H', J/J'\} \\ h' &= \{A'/A, B'/B, D'/D, F'/F, L'/L, N'/N, H'/H, J'/J\} \end{aligned}$$

**Theorem 5.2.7** *Two derived DACEQs are  $\bar{\xi}$ -equivalent if there exists index-covering homomorphisms in both directions between their  $\bar{\xi}$ -normal forms.*

PROOF. Let  $Q(\overline{\mathcal{I}}_1; \dots; \overline{\mathcal{I}}_d; \overline{\mathcal{V}})$  and  $Q'(\overline{\mathcal{I}}'_1; \dots; \overline{\mathcal{I}}'_d; \overline{\mathcal{V}}')$  be the  $\bar{\xi}$ -normal forms of the given queries, and let  $h : Q' \rightarrow Q$  and  $h' : Q \rightarrow Q'$  be index-covering homomorphisms between them. By Theorem 5.1.8, the given queries are  $\bar{\xi}$ -equivalent iff  $Q \stackrel{\bar{\xi}}{\equiv} Q'$ . The existence of  $h$  and  $h'$  implies by definition that there exists a head graph isomorphism  $\theta$  from  $Q$  to  $Q'$ , and by Corollary 5.2.4  $\theta$  must be a bijection between the nodes in  $\text{body}_Q$  and  $\text{body}_{Q'}$ .

By Definition 5.2.6, composite mapping  $h \circ h'$  satisfies  $\forall i \in [1, d]. [h \circ h'(\hat{\mathcal{I}}'_i) = \hat{\mathcal{I}}_i]$ . Therefore, we can assume w.l.o.g. that  $h \circ h'$  satisfies the following property

$$\forall i \in [1, d]. [h \circ h'(\overline{\mathcal{I}}_i) = \overline{\mathcal{I}}_i] \quad (5.5)$$

(by the finite period of any permutation of tuple  $\overline{\mathcal{I}}_{[1,d]}$ ).

Consider any database  $\mathbb{D}$  and any terminal embedding  $\gamma : G \rightsquigarrow \mathbb{D}$  of any  $G \in \text{dags}(\text{body}_Q)$ ; then,  $\gamma \circ h$  is an embedding of  $\theta(G) \in \text{dags}(\text{body}_{Q'})$  into  $\mathbb{D}$ . Definition 5.1.3 implies  $\mathcal{S}_Q \subseteq \mathcal{I}_{[1,d]}$ , and equation 5.5 implies  $\gamma \circ h \circ h'(\overline{\mathcal{I}}_{[1,d]}) = \gamma(\overline{\mathcal{I}}_{[1,d]})$ ; therefore, the definition of  $\mathcal{S}_Q$  (see Corollary 4.2.7) mandates that  $\gamma \circ h \circ h'$  is a terminal embedding. Hence, embedding  $\gamma \circ h$  must also be a terminal embedding of  $\theta(G)$ .

Repeating this argument in the other direction proves that encoding relations  $(Q)^\mathbb{D}$  and  $(Q')^\mathbb{D}$  differ only by permutation of attributes within each index level, and so it is trivial to construct a  $\bar{\xi}$ -certificate proving  $(Q)^\mathbb{D} \stackrel{\bar{\xi}}{\equiv} (Q')^\mathbb{D}$ . By our unrestricted choice of  $\mathbb{D}$ ,  $Q \stackrel{\bar{\xi}}{\equiv} Q'$  follows.  $\square$

**Example 59** The homomorphisms in Example 58 prove that over any database, the encoding relations yielded by  $Q_{57}$  and  $Q_{58}$  will have identical<sup>2</sup> instances, and hence  $Q_{57} \stackrel{\text{snsb}}{\equiv} Q_{58}$ . We conclude that  $Q_{20} \stackrel{\text{snsb}}{\equiv} Q_{21}$  (see Example 57) and therefore  $Q_{18} \equiv Q_{19}$  (see Example 23).

### 5.3 Food for Thought

In this section we show two examples that provide further insight into the equivalence results of the last section. In particular, the examples illustrate that the condition in Theorem 5.2.7—the existence of index-covering homomorphisms in both directions between the  $\bar{\xi}$ -normal forms—is not necessary for  $\bar{\xi}$ -equivalence between derived DACEQs.

The first example illustrates a shortcoming in our characterization of core indexes. Ideally, the set of core indexes should be minimal in the sense that no strict subset would preserve  $\bar{\xi}$ -equivalence. Our definition of core indexes for CEQs in Chapter 3 satisfies this minimality property, as evidenced by the proof of Theorem 3.2.4 (see Appendix A). This example demonstrates that our definition of core indexes for derived DACEQs (Definition 5.1.3) does not share this property.

**Example 60** Consider the three HCEQs shown in Figure 5.3a, all sharing a query body that is already in DACQ canonical form. Consider the sss-normalization of  $Q_{59}$  (using Definition 5.1.3).

$$\begin{aligned}
i = 3: \quad & \mathfrak{I}_3 = \mathbf{s}, \hat{\mathcal{I}}_3^{\bar{\xi}} = \{\} \\
& C^{\square} \notin \hat{\mathcal{I}}_3^{\bar{\xi}} \text{ because } C \in \mathcal{I}_2. \\
i = 2: \quad & \mathfrak{I}_2 = \mathbf{s}, \hat{\mathcal{I}}_2^{\bar{\xi}} = \{A^{\square}, C^{\square}\} \\
& C^{\square} \in \hat{\mathcal{I}}_2^{\bar{\xi}} \text{ because } C \in \mathcal{V} \setminus (\mathcal{I}_1 \cup \mathcal{I}_3^{\bar{\xi}} \cup \mathcal{C}); \text{ meanwhile, } A^{\square} \in \hat{\mathcal{I}}_2^{\bar{\xi}} \text{ because } Q_{(u,2)} \models \\
& BA \rightarrow C \text{ and } Q_{(uv,2)} \models BA \rightarrow C \text{ (both trivially), but } Q_{(u,2)} \not\models B \rightarrow C.
\end{aligned}$$

$$\begin{aligned}
Q_{(u,2)}(B, A, C) & :- R(A, B), S(B, C), R(A, D), S(D, C), T(D) \\
Q_{(uv,2)}(B, A, C) & :- R(A, B), S(B, C), R(A, D), S(D, C), T(D), T(B)
\end{aligned}$$

$$\begin{aligned}
i = 1: \quad & \mathfrak{I}_1 = \mathbf{s}, \hat{\mathcal{I}}_1^{\bar{\xi}} = \{B^{\square}\} \\
& B^{\square} \in \hat{\mathcal{I}}_1^{\bar{\xi}} \text{ because } \mathcal{S}_{Q_{59}} = \{B\}.
\end{aligned}$$

Therefore,  $Q_{60}$  is the sss-NF of  $Q_{59}$ . However, we can show that  $Q_{60} \not\stackrel{\text{sss}}{\equiv} Q_{61}$ , implying that  $A$  should not need to be a core index of  $Q_{59}$  to preserve ss-equivalence.

<sup>2</sup>In general, the existence of index-covering homomorphisms in both directions only implies only that query results are identical *modulo permutation of attributes within each index level*. However, the index-covering homomorphisms  $h$  and  $h'$  from Example 58 have the additional property that for each index level  $i$ ,  $h$  and  $h'$  form an isomorphism between the pair-wise components of the two index tuples. This justifies the stronger claim that the query results are identical, without any permutation of attributes.

Suppose that we replace the attribute  $C^{\boxed{v}}$  with the constant  $0^{\boxed{v}}$  within the head of  $Q_{59}$ . Within this new query, variable  $C$  would not be an output variable, and so the trivial MVDs  $Q_{(u,2)} \models B \twoheadrightarrow \emptyset$  and  $Q_{(uv,2)} \models B \twoheadrightarrow \emptyset$  would justify  $\hat{\mathcal{I}}_2^{\S} = \{\}$  (i.e., neither  $A^{\boxed{u}}$  nor  $C^{\boxed{u}}$  would be core indexes). Therefore, the assessment that both  $A^{\boxed{u}}$  and  $C^{\boxed{u}}$  are core within  $Q_{59}$  (as per Definition 5.1.3) is directly attributable to head attribute  $C^{\boxed{v}}$ .

Now reconsider CQ  $Q_{(uv,2)}$  shown above;  $Q_{(uv,2)}$  is not a minimal CQ because it allows the automorphism  $h = \{A/A, B/B, C/C, D/B\}$ . After using  $h$  to minimize the query body, it becomes clear that  $Q_{(uv,2)}$  satisfies the MVD  $B \twoheadrightarrow C$ . In contrast,  $h$  is not an automorphism for  $Q_{(u,2)}$ , and  $Q_{(u,2)} \not\models B \twoheadrightarrow C$ . As a concrete example, consider the database instance  $\mathbb{D}_9$  shown in Figure 5.3b and the set of terminal embeddings  $\Gamma = \{\gamma : Q_{59} \rightsquigarrow \mathbb{D}_9\}$  shown in Figure 5.3c. If we project “relation”  $\Gamma$  onto attributes  $ABC$ , it violates the MVD  $B \twoheadrightarrow C$ ; however, if we first restrict  $\Gamma$  to only embeddings that extend to  $\boxed{v}$ , then MVD  $B \twoheadrightarrow C$  is satisfied.

In summary, index  $C^{\boxed{u}}$  needs to be core in order to functionally determine the value of the output attribute  $C^{\boxed{v}}$ , but *only for embeddings that extend to node  $\boxed{v}$* . Index  $A^{\boxed{u}}$  was deemed core because in general it is not separated from core index  $C^{\boxed{u}}$  by an MVD, but the subset of embeddings that extend to node  $\boxed{v}$ —i.e., the embeddings that require  $C^{\boxed{u}}$  to be core—satisfy an MVD separating  $A^{\boxed{u}}$  from  $C^{\boxed{u}}$  (or more precisely,  $A$  from  $C$ ).

It is likely that the shortcomings of Definition 5.1.3 illustrated by Example 60 could be overcome by a more sophisticated characterization of core indexes. One possibility would be to devise a system of bookkeeping that associates with each core index a set of reasons justifying why that index needs to be core. The reasons associated with known core indexes would then need to factor into later decisions identifying further core indexes.

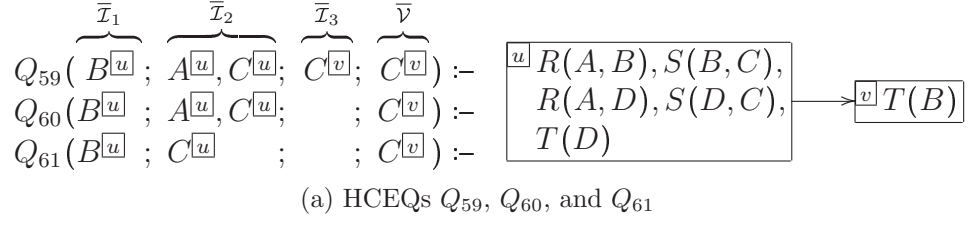
In the next example we use HCEQ encoding-equivalence to decide equivalence between a pair of HCQs composed of a single edge between two nodes. The two HCQs in this example have been carefully crafted to illustrate why index-covering homomorphisms are not sufficient for characterizing HCEQ encoding-equivalence.

**Example 61** Consider the two HCQs  $Q_{62}$  and  $Q_{63}$  shown in Figure 5.4. Using the results of Section 4.3, we can show  $Q_{62} \equiv Q_{63}$  as follows.

Both queries have canonical tuples  $\{\{0\}, \{\blacksquare\}\}$ . The tuple-specific equivalence  $Q_{62} \stackrel{\{0\}}{\equiv} Q_{63}$  follows as a result of the following two homomorphisms.

$$\begin{aligned} h &= \{A/W, B/X, C/Y, D/Z, E/W, F/X, B_1/X_1, \dots, B_n/X_n\} \\ h' &= \{W/A, X/B, Y/A, Z/B, X_1/B_1, \dots, X_n/B_n, Z_1/B_1, \dots, Z_n/B_n\} \end{aligned}$$

To prove  $Q_{62} \stackrel{\{\blacksquare\}}{\equiv} Q_{63}$  requires proving  $\Sigma^{(Q_{62}, \{\blacksquare\})} \equiv_{\text{fin}} \Sigma^{(Q_{63}, \{\blacksquare\})}$ , where  $\Sigma^{(Q_{62}, \{\blacksquare\})}$  and



R	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr style="border-bottom: 1px solid black;"> <th style="padding: 2px 10px;">col<sub>1</sub></th> <th style="padding: 2px 10px;">col<sub>2</sub></th> </tr> </thead> <tbody> <tr><td style="padding: 2px 10px;">a<sub>1</sub></td><td style="padding: 2px 10px;">b</td></tr> <tr><td style="padding: 2px 10px;">a<sub>1</sub></td><td style="padding: 2px 10px;">d<sub>1</sub></td></tr> <tr><td style="padding: 2px 10px;">a<sub>2</sub></td><td style="padding: 2px 10px;">b</td></tr> <tr><td style="padding: 2px 10px;">a<sub>2</sub></td><td style="padding: 2px 10px;">d<sub>2</sub></td></tr> </tbody> </table>	col <sub>1</sub>	col <sub>2</sub>	a <sub>1</sub>	b	a <sub>1</sub>	d <sub>1</sub>	a <sub>2</sub>	b	a <sub>2</sub>	d <sub>2</sub>	S	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr style="border-bottom: 1px solid black;"> <th style="padding: 2px 10px;">col<sub>1</sub></th> <th style="padding: 2px 10px;">col<sub>2</sub></th> </tr> </thead> <tbody> <tr><td style="padding: 2px 10px;">b</td><td style="padding: 2px 10px;">c<sub>1</sub></td></tr> <tr><td style="padding: 2px 10px;">d<sub>1</sub></td><td style="padding: 2px 10px;">c<sub>1</sub></td></tr> <tr><td style="padding: 2px 10px;">b</td><td style="padding: 2px 10px;">c<sub>2</sub></td></tr> <tr><td style="padding: 2px 10px;">d<sub>2</sub></td><td style="padding: 2px 10px;">c<sub>2</sub></td></tr> </tbody> </table>	col <sub>1</sub>	col <sub>2</sub>	b	c <sub>1</sub>	d <sub>1</sub>	c <sub>1</sub>	b	c <sub>2</sub>	d <sub>2</sub>	c <sub>2</sub>	T	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr style="border-bottom: 1px solid black;"> <th style="padding: 2px 10px;">col<sub>1</sub></th> </tr> </thead> <tbody> <tr><td style="padding: 2px 10px;">d<sub>1</sub></td></tr> <tr><td style="padding: 2px 10px;">d<sub>2</sub></td></tr> </tbody> </table>	col <sub>1</sub>	d <sub>1</sub>	d <sub>2</sub>
col <sub>1</sub>	col <sub>2</sub>																											
a <sub>1</sub>	b																											
a <sub>1</sub>	d <sub>1</sub>																											
a <sub>2</sub>	b																											
a <sub>2</sub>	d <sub>2</sub>																											
col <sub>1</sub>	col <sub>2</sub>																											
b	c <sub>1</sub>																											
d <sub>1</sub>	c <sub>1</sub>																											
b	c <sub>2</sub>																											
d <sub>2</sub>	c <sub>2</sub>																											
col <sub>1</sub>																												
d <sub>1</sub>																												
d <sub>2</sub>																												

(b) Database  $D_9$

Γ	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr style="border-bottom: 1px solid black;"> <th style="padding: 2px 10px;">A</th> <th style="padding: 2px 10px;">B</th> <th style="padding: 2px 10px;">C</th> <th style="padding: 2px 10px;">D</th> <th style="padding: 2px 10px;"><math>B^{\overline{u}}</math></th> <th style="padding: 2px 10px;"><math>A^{\overline{u}}</math></th> <th style="padding: 2px 10px;"><math>C^{\overline{u}}</math></th> <th style="padding: 2px 10px;"><math>C^{\overline{v}}</math></th> </tr> </thead> <tbody> <tr><td style="padding: 2px 10px;">γ<sub>1</sub></td><td style="padding: 2px 10px;">a<sub>1</sub></td><td style="padding: 2px 10px;">b</td><td style="padding: 2px 10px;">c<sub>1</sub></td><td style="padding: 2px 10px;">d<sub>1</sub></td><td style="padding: 2px 10px;">b</td><td style="padding: 2px 10px;">a<sub>1</sub></td><td style="padding: 2px 10px;">c<sub>1</sub></td><td style="padding: 2px 10px;">■</td></tr> <tr><td style="padding: 2px 10px;">γ<sub>2</sub></td><td style="padding: 2px 10px;">a<sub>1</sub></td><td style="padding: 2px 10px;">d<sub>1</sub></td><td style="padding: 2px 10px;">c<sub>1</sub></td><td style="padding: 2px 10px;">d<sub>1</sub></td><td style="padding: 2px 10px;">d<sub>1</sub></td><td style="padding: 2px 10px;">a<sub>1</sub></td><td style="padding: 2px 10px;">c<sub>1</sub></td><td style="padding: 2px 10px;">c<sub>1</sub></td></tr> <tr><td style="padding: 2px 10px;">γ<sub>3</sub></td><td style="padding: 2px 10px;">a<sub>2</sub></td><td style="padding: 2px 10px;">b</td><td style="padding: 2px 10px;">c<sub>2</sub></td><td style="padding: 2px 10px;">d<sub>2</sub></td><td style="padding: 2px 10px;">b</td><td style="padding: 2px 10px;">a<sub>2</sub></td><td style="padding: 2px 10px;">c<sub>2</sub></td><td style="padding: 2px 10px;">■</td></tr> <tr><td style="padding: 2px 10px;">γ<sub>4</sub></td><td style="padding: 2px 10px;">a<sub>2</sub></td><td style="padding: 2px 10px;">d<sub>2</sub></td><td style="padding: 2px 10px;">c<sub>2</sub></td><td style="padding: 2px 10px;">d<sub>2</sub></td><td style="padding: 2px 10px;">d<sub>2</sub></td><td style="padding: 2px 10px;">a<sub>2</sub></td><td style="padding: 2px 10px;">c<sub>2</sub></td><td style="padding: 2px 10px;">c<sub>2</sub></td></tr> </tbody> </table>	A	B	C	D	$B^{\overline{u}}$	$A^{\overline{u}}$	$C^{\overline{u}}$	$C^{\overline{v}}$	γ <sub>1</sub>	a <sub>1</sub>	b	c <sub>1</sub>	d <sub>1</sub>	b	a <sub>1</sub>	c <sub>1</sub>	■	γ <sub>2</sub>	a <sub>1</sub>	d <sub>1</sub>	c <sub>1</sub>	d <sub>1</sub>	d <sub>1</sub>	a <sub>1</sub>	c <sub>1</sub>	c <sub>1</sub>	γ <sub>3</sub>	a <sub>2</sub>	b	c <sub>2</sub>	d <sub>2</sub>	b	a <sub>2</sub>	c <sub>2</sub>	■	γ <sub>4</sub>	a <sub>2</sub>	d <sub>2</sub>	c <sub>2</sub>	d <sub>2</sub>	d <sub>2</sub>	a <sub>2</sub>	c <sub>2</sub>	c <sub>2</sub>
A	B	C	D	$B^{\overline{u}}$	$A^{\overline{u}}$	$C^{\overline{u}}$	$C^{\overline{v}}$																																						
γ <sub>1</sub>	a <sub>1</sub>	b	c <sub>1</sub>	d <sub>1</sub>	b	a <sub>1</sub>	c <sub>1</sub>	■																																					
γ <sub>2</sub>	a <sub>1</sub>	d <sub>1</sub>	c <sub>1</sub>	d <sub>1</sub>	d <sub>1</sub>	a <sub>1</sub>	c <sub>1</sub>	c <sub>1</sub>																																					
γ <sub>3</sub>	a <sub>2</sub>	b	c <sub>2</sub>	d <sub>2</sub>	b	a <sub>2</sub>	c <sub>2</sub>	■																																					
γ <sub>4</sub>	a <sub>2</sub>	d <sub>2</sub>	c <sub>2</sub>	d <sub>2</sub>	d <sub>2</sub>	a <sub>2</sub>	c <sub>2</sub>	c <sub>2</sub>																																					

(c) Terminal embeddings of  $Q_{59}$  into  $D_9$

Figure 5.3: Encoding queries  $Q_{59}$ ,  $Q_{60}$ , and  $Q_{61}$  evaluated over database  $D_9$ .

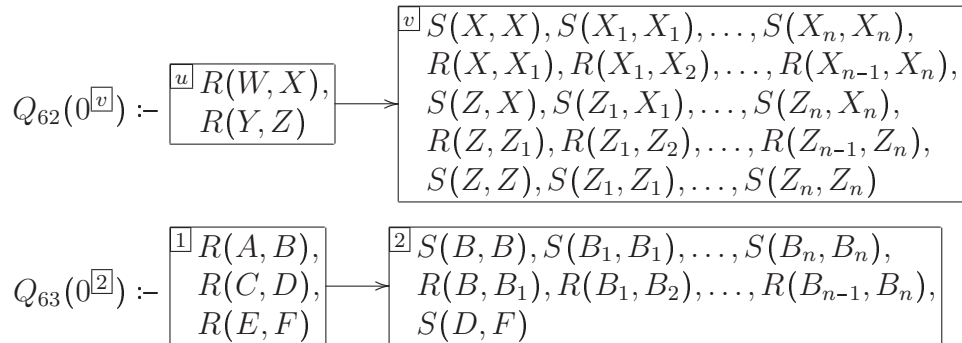


Figure 5.4: HCQs  $Q_{62}$  and  $Q_{63}$



$\Sigma(Q_{62}, \langle \blacksquare \rangle)$  are defined as follows.

$$\begin{aligned}\Sigma(Q_{62}, \langle \blacksquare \rangle) &= \{\xi^{(u, \langle \blacksquare \rangle)} : \text{body}_u \rightarrow \text{body}_v\} \\ \Sigma(Q_{63}, \langle \blacksquare \rangle) &= \{\xi^{(a, \langle \blacksquare \rangle)} : \text{body}_1 \rightarrow \text{body}_2\}\end{aligned}$$

The implication  $\Sigma(Q_{62}, \langle \blacksquare \rangle) \models_{\text{fin}} \Sigma(Q_{63}, \langle \blacksquare \rangle)$  can be proven by the chase procedure using two chase steps:

1. chase  $\xi^{(u, \langle \blacksquare \rangle)}$  with binding  $X/b, Z/*$  to generate the chains  $S(b, b), S(b_1, b_1), \dots, S(b_n, b_n)$  and  $R(b, b_1), \dots, R(b_{n-1}, b_n)$ ;
2. chase  $\xi^{(a, \langle \blacksquare \rangle)}$  with binding  $X/f, Z/d$  to generate the atom  $S(d, f)$ .

The implication  $\Sigma(Q_{63}, \langle \blacksquare \rangle) \models_{\text{fin}} \Sigma(Q_{62}, \langle \blacksquare \rangle)$  can be proven by the chase procedure using  $n + 2$  chase steps:

1. chase  $\xi^{(1, \langle \blacksquare \rangle)}$  with binding  $B/x, D/z, F/x$  to generate the atom  $S(z, x)$  and the chains  $S(x, x), S(x_1, x_1), \dots, S(x_n, x_n)$  and  $R(x, x_1), \dots, R(x_{n-1}, x_n)$ ;
2. chase  $\xi^{(1, \langle \blacksquare \rangle)}$  with binding  $B/z, D/*, F/*$  to generate the chains  $S(z, z), S(z_1, z_1), \dots, S(z_n, z_n)$  and  $R(z, z_1), \dots, R(z_{n-1}, z_n)$ ;
3. for each  $i \in [1, n]$ , chase  $\xi^{(1, \langle \blacksquare \rangle)}$  with binding  $B/*, D/z_i, F/x_i$  to generate the atom  $S(z_i, x_i)$ .

Now consider the following HCEQ  $Q_{64}$ .

$$Q_{64}(\overbrace{W^{\blacksquare}, X^{\blacksquare}, Y^{\blacksquare}, Z^{\blacksquare}}^{\bar{x}_1}; \overbrace{X_1^{\blacksquare}, \dots, X_n^{\blacksquare}, Z_1^{\blacksquare}, \dots, Z_n^{\blacksquare}}^{\bar{x}_2}; \overbrace{0^{\blacksquare}}^{\bar{y}}) :- \text{body}_{Q_{62}}$$

Theorem 4.2.4 implies that  $\mathcal{S}_{Q_{64}} = \{X, Z\}$  is the minimal set satisfying  $Q_{64} \models \mathcal{S}_{Q_{64}} \rightarrow v$ , and so ss-normalization of  $Q_{64}$  yields the following query  $Q_{65}$ .

$$Q_{65}(X^{\blacksquare}, Z^{\blacksquare}; ; 0^{\blacksquare}) :- \text{body}_{Q_{62}}$$

It straightforward to verify that over any database  $\mathbb{D}$ , the ss-decoding of  $(Q_{65})^{\mathbb{D}}$  yields one of four possible values:

1.  $\{\{0\}, \{\blacksquare\}\}$  if and only if  $(Q_{62})^{\mathbb{D}} = \{\{0\}, \{\blacksquare\}\}$ ,
2.  $\{\{0\}\}$  if and only if  $(Q_{62})^{\mathbb{D}} = \{\{0\}\}$ ,
3.  $\{\{\blacksquare\}\}$  if and only if  $(Q_{62})^{\mathbb{D}} = \{\{\blacksquare\}\}$ , or
4.  $\{\}$  if and only if  $(Q_{62})^{\mathbb{D}} = \{\}$ .

Repeating the above argument for  $Q_{63}$ , we define  $Q_{66}$

$$Q_{66}(\overbrace{A^{\blacksquare}, B^{\blacksquare}, C^{\blacksquare}, D^{\blacksquare}, E^{\blacksquare}, F^{\blacksquare}}^{\bar{x}_1}; \overbrace{B_1^{\blacksquare}, \dots, B_n^{\blacksquare}}^{\bar{x}_2}; \overbrace{0^{\blacksquare}}^{\bar{y}}) :- \text{body}_{Q_{63}}$$

which has  $\mathcal{S}_{Q_{66}} = \{B, D, F\}$  and ss-normal form  $Q_{67}$ .

$$Q_{67}(B^{\blacksquare}, D^{\blacksquare}, F^{\blacksquare}; ; 0^{\blacksquare}) :- \text{body}_{Q_{63}}$$

Due to correspondence between the outputs of  $Q_{67}$  and  $Q_{63}$ , we conclude the following.

$$Q_{65} \stackrel{\dot{=}}{\equiv}_{\text{ss}} Q_{67} \iff Q_{62} \equiv Q_{63}$$

We have already shown that  $Q_{62} \equiv Q_{63}$ , and therefore  $Q_{65} \stackrel{\dot{=}}{\equiv}_{\text{ss}} Q_{67}$ . Clearly index-covering homomorphisms can not exist in both directions between  $Q_{65}$  and  $Q_{67}$ , since they have different numbers of core indexes. Furthermore, this phenomenon is not due to a failure to minimize the core indexes—as in Example 60—because for both  $Q_{65}$  and  $Q_{67}$  the core indexes are already minimal (in the sense that no strict subset preserves ss-equality).

The reduction in Example 61 can be repeated for any pair of two-node HCQs. Consequently, fully solving the HCEQ encoding-equivalence problem necessitates first solving the HCQ equivalence problem (at least for the restricted case of two-node HCQs). Unfortunately, our results from Chapter 4 where we considered HCQ/DACQ equivalence do not guarantee a decidable equivalence condition even for the simple class of two-node HCQs.

# Chapter 6

## Conclusions

In Section 6.1 we summarize the main results of this thesis. In Section 6.2 we revisit the motivating examples introduced in Chapter 1 and use them to illustrate both concrete applications of our results and shortcomings in the theory which motivate some of the avenues of future work we discuss in Section 6.3.

### 6.1 Summary of Results

To the best of our knowledge, this thesis is the only systematic study of the query equivalence problem for conjunctive queries extended with nested aggregation. Our study yielded the following results:

- 1. A novel encoding of complex objects within a single flat relation.**  
Most previous encodings of complex objects as flat relations utilize multiple relations, with the existence of a collection stored separately from the values of its elements. Our definition of object “linearization” and subsequent use of null values to encode empty sub-collections highlights an interesting connection between scalar aggregation (i.e., aggregation operators capable of aggregating over empty inputs) and the left outer-join operator (see Chapter 2).
- 2. A normal form for encoding queries that illuminates the interactions between nested aggregation operations.**  
In particular, our characterization of the *core* indexes of an encoding query in terms of query-implied multi-valued dependencies reveals how the joins within the query body interact with nested aggregation (see Chapters 3 and 5).
- 3. A complete characterization of equivalence between conjunctive queries that construct objects without empty sub-collections.**  
For the special case of object-constructing queries that do not construct empty sub-collections, deciding query equivalence reduces to deciding encoding-equivalence

between standard conjunctive queries. We prove that deciding this encoding-equivalence condition is NP-complete and corresponds to a homomorphic condition between the normal forms of the encoding queries (see Chapter 3 and Appendix A).

4. **Necessary conditions and sufficient conditions for equivalence between conjunctive queries with outer joins.**

To the best of our knowledge, this thesis is the first work addressing the query equivalence problem for queries containing the left outer-join operator. We do not completely solve this problem; however, we do introduce a helpful high-level representation of these queries as DAGs whose nodes are (rule-based) conjunctive queries. This variable-based (rather than algebraic) syntax facilitates converting the query bodies to a canonical form, over which we are able to state several necessary conditions and sufficient conditions for query equivalence in terms of homomorphisms between query bodies (see Chapter 4).

5. **Necessary conditions and sufficient conditions for equivalence between conjunctive queries that construct arbitrary objects.**

These conditions are stated in terms of homomorphisms between normal forms of the associated encoding queries, and could be used to develop more powerful logical rewriting rules for relational query optimizers. We also show that query equivalence in the presence of nested aggregation is at least as hard as query equivalence in the presence of left outer joins (see Chapter 5).

## 6.2 Comprehensive Examples

Recall the motivating examples from Chapter 1. We now show that, without the primary and foreign key constraints shown in Figure 1.1, the logical rewriting we suggest within Example 1 does not preserve equivalence.

**Example 62** Consider queries  $Q_1$  and  $Q_3$  from Example 1. If we model the outputs of the SQL aggregation functions SUM and AVG as bags and normalized bags, respectively, then  $Q_1$  and  $Q_3$  translate into the COCQL $_{\{f_b, f_s, f_n\}}$  queries  $Q_{68}$  and  $Q_{69}$ , shown in Figures 6.1 and 6.2. The translation of  $Q_1$  makes use of a simple transformation from an aggregation block with  $k$  aggregation expressions into a join of  $k$  such blocks, each with a single aggregation expression. (For space reasons, we use a hybrid notation in Figures 6.1 and 6.2 that replaces blocks of joins between base tables with an equivalent rule-based notation. We have used frames to indicate the logical expansion of view definitions, with thin frames denoting logical views—which are logically expanded during query evaluation—and thick frames denoting materialized views—which are not logically expanded during query evaluation.)

Figure 6.3 depicts the output sort  $\tau_{14}$  of  $Q_{68}$  and  $Q_{69}$ , as well as its transformation to the chain sort (b**nb**nb, 6). Figure 6.4 illustrates the CEQs  $Q_{70} := \text{ENCODE}(Q_{68})$  and

$Q_{71} := \text{ENCODE}(Q_{69})$ . (Components of the queries have been labelled for the sake of clarity.) Converting query  $Q_{70}$  to bnb-NF removes the shaded indexes from  $\bar{I}_4$  and  $\bar{I}_2$ . Query  $Q_{71}$  is already in bnb-NF. Clearly, index-covering homomorphisms cannot exist in both directions between the normalized versions of  $Q_{70}$  and  $Q_{71}$ , because the sizes of the index sets differ. Therefore, we conclude  $Q_{70} \not\equiv_{\text{bnb}} Q_{71}$ , which implies  $Q_{68} \not\equiv Q_{69}$ , which further implies that the SQL queries  $Q_1$  and  $Q_3$  are not equivalent over unrestricted database instances.<sup>1</sup>

Our characterization of CEQ encoding-equivalence in terms of homomorphisms is a direct generalization of the traditional characterization of CQ equivalence. Because of this, it is straightforward to adapt techniques for testing CQ equivalence with respect to a set  $\Sigma$  of schema constraints (denoted  $Q \equiv^\Sigma Q'$ ) to CEQ encoding-equivalence. Assuming that  $\Sigma$  contains only functional dependencies (FDs) and acyclic inclusion dependencies (INDs)—which generalize primary keys and acyclic foreign keys—we can decide *encoding-equivalence w.r.t.  $\Sigma$*  as follows:

1. Prior to the conversion to  $\bar{\xi}$ -normal form, we pre-process each CEQ by first applying the standard  $\Sigma$ -chase, as defined over CQs [6, Ch. 8–9]. (Chasing a query with an FD causes two variables to be merged, while chasing with an INDs introduces a new atom into the query body.)
2. After the queries have been chased, we use only the FDs in  $\Sigma$  to expand out the index sets in the query head. For each  $i \in [1, d]$  starting with  $i = 1$ , we expand  $\bar{I}_i$  with all variables functionally determined by  $\bar{I}_{[1,i]}$ , deleting variables from inner index sets whenever they are added to outer index sets.
3. The conversion to  $\bar{\xi}$ -NF is unchanged, but the testing of the query-implied MVDs when identifying core indexes (Definition 3.1.11) may need to account for the FDs in  $\Sigma$ . (More specifically, if query-implied MVDs are tested via the hypergraph condition of Theorem 3.1.10, then no changes are necessary; however, if query-implied MVDs are tested via the query equivalence condition of Theorem 3.1.3, then that condition needs to use  $\equiv^\Sigma$  instead.)

We now repeat Example 62, showing that the logical rewriting we suggest within Example 1 preserves equivalence with respect to database instances that conform to the schema constraints.

**Example 63** Reconsider queries  $Q_{70}$  and  $Q_{71}$  in Figure 6.4. Let  $\Sigma$  denote the set of primary and foreign key constraints shown in Figure 1.1. Chasing the body of  $Q_{70}$  with  $\Sigma$  does not introduce any new atoms, but it does merge the variables  $N, N_2, N_4$ . After expanding the index sets in the head of  $Q_{70}$ , we obtain the new query  $Q_{72} := \text{chase}_\Sigma(Q_{70})$  shown in Figure 6.4. Shaded attributes again indicate redundant index columns that get removed by bnb-normalization. Query  $Q_{71}$  is unchanged after chasing with  $\Sigma$ .

---

<sup>1</sup>Cohen et al. prove that queries performing top-level application of SUM are equivalent if and only if the relational inputs to the SUM function are always equivalent under bag semantics [25]. Grumbach et al. prove a similar result relating AVG to normalized bags [50].

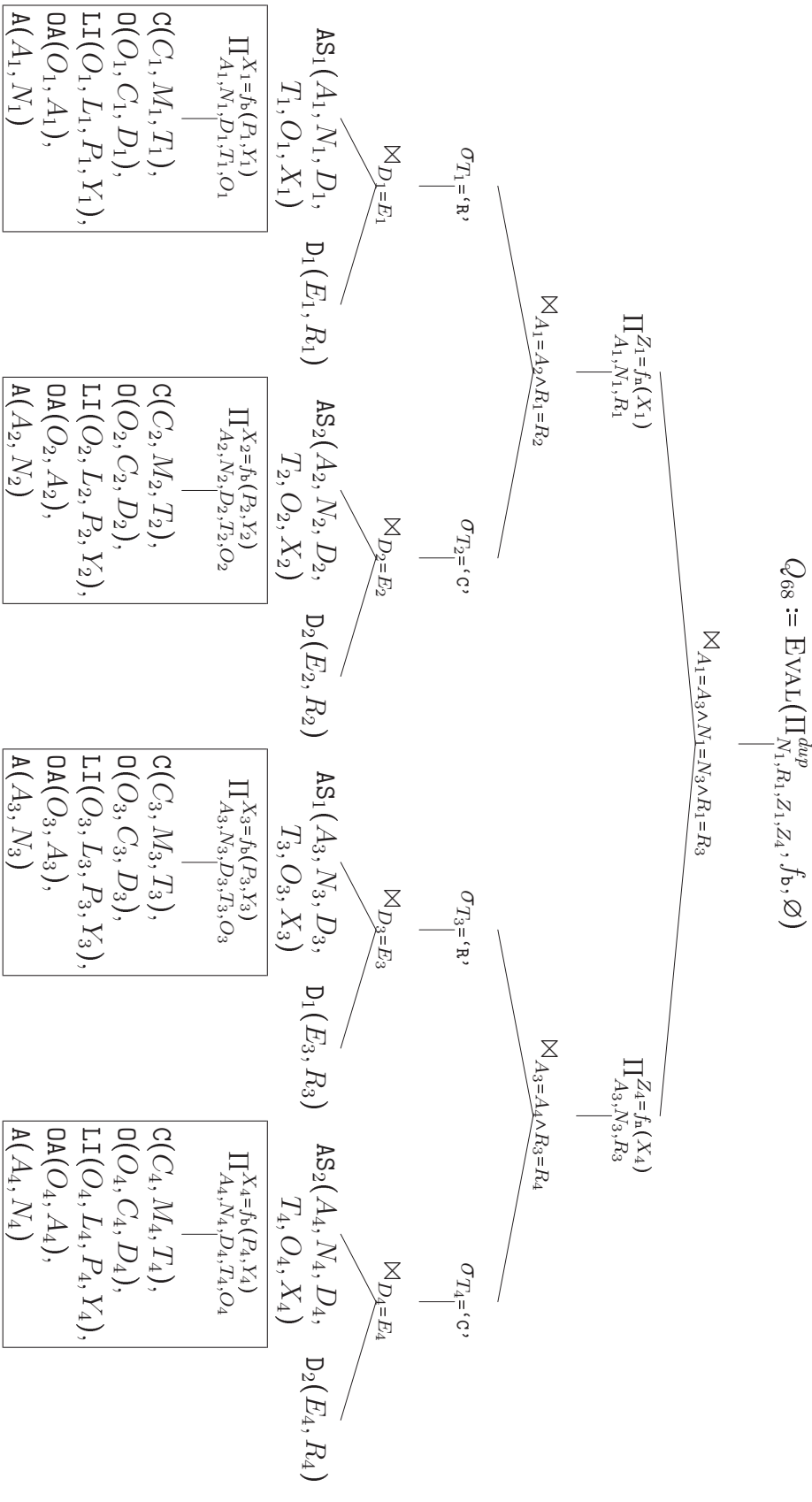


Figure 6.1:  $Q_{68}$ , the COCQL  $\{f_b, f_s, f_n\}$  translation of SQL query  $Q_1$

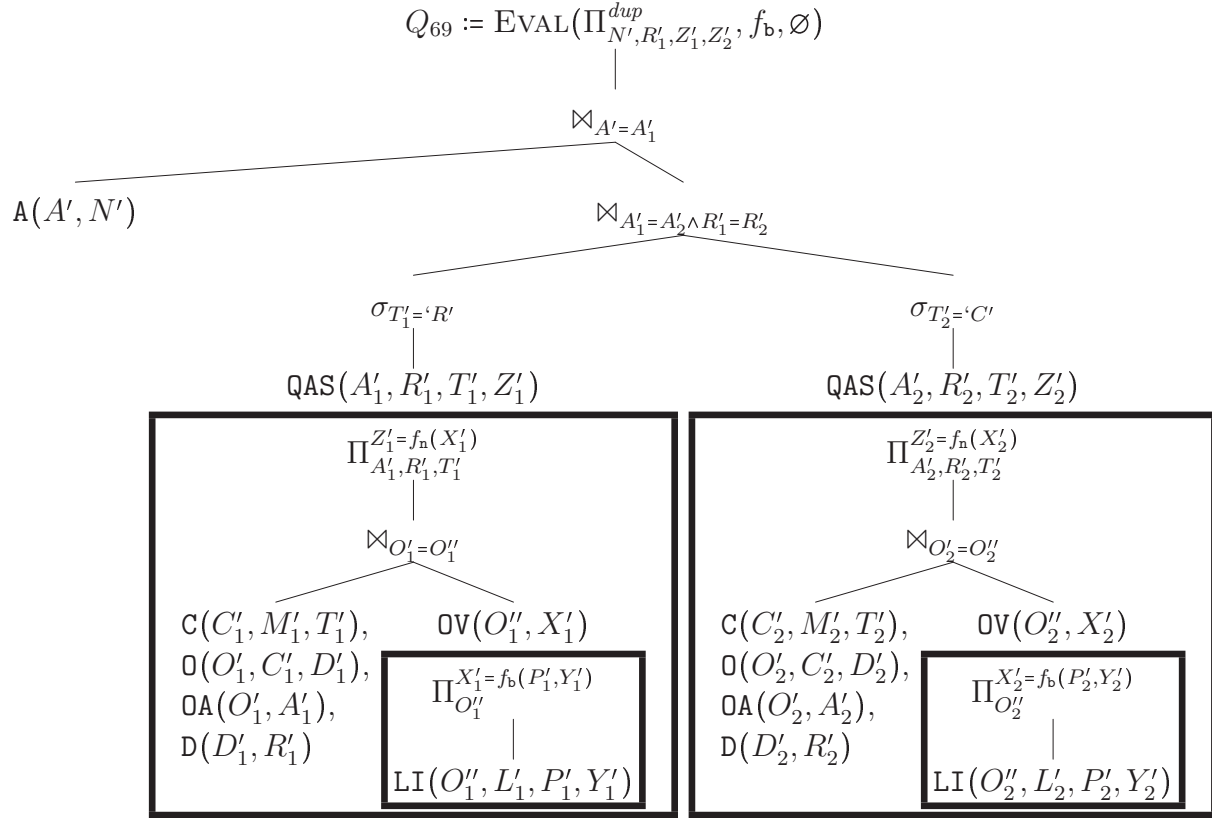


Figure 6.2:  $Q_{69}$ , the  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  translation of SQL query  $Q_3$

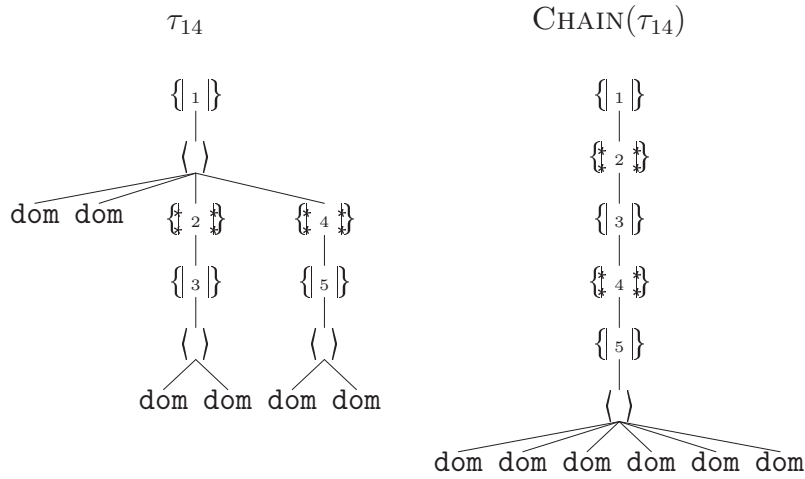


Figure 6.3: Output sort  $\tau_{14}$  of  $Q_{68}$  and  $Q_{69}$ , and its transformation to a chain

$$\begin{aligned}
Q_{70} & \underbrace{(A, N, R; D_1, O_1, N_2, D_2, O_2; C_1, M_1, L_1, P_1, Y_1; D_3, O_3, M_4, D_4, O_4; C_4, M_4, L_4, P_4, Y_4; N, R, P_1, Y_1, P_4, Y_4)}_{\bar{I}_1 \bar{I}_2 \bar{I}_3 \bar{I}_4 \bar{I}_5 \bar{Y}} \text{ :-} \\
& \underbrace{c(C_1, M_1, 'R^2'), 0(O_1, C_1, D_1), \text{LI}(O_1, L_1, P_1, Y_1), 0A(O_1, A), A(A, N), D(D_1, R), \}_{(AS_1 \bowtie D_1)}_{\bar{I}_1} \text{ avgRrsale} \\
& \underbrace{c(C_2, M_2, 'C^1'), 0(O_2, C_2, D_2), \text{LI}(O_2, L_2, P_2, Y_2), 0A(O_2, A), A(A, N_2), D(D_2, R), \}_{(AS_2 \bowtie D_2)}_{\bar{I}_2} \text{ avgRrsale} \\
& \underbrace{c(C_3, M_3, 'R^2'), 0(O_3, C_3, D_3), \text{LI}(O_3, L_3, P_3, Y_3), 0A(O_3, A), A(A, N), D(D_3, R), \}_{(AS_1 \bowtie D_1)}_{\bar{I}_3} \text{ avgCsale} \\
& \underbrace{c(C_4, M_4, 'C^1'), 0(O_4, C_4, D_4), \text{LI}(O_4, L_4, P_4, Y_4), 0A(O_4, A), A(A, N_4), D(D_4, R) \}_{(AS_2 \bowtie D_2)}_{\bar{I}_4 \bar{I}_5} \text{ avgCsale} \\
\\
Q_{71} & \underbrace{(A, N', R; D'_1, O'_1, C'_1, M'_1; L'_1, P'_1, Y'_1; D_2, O_2, C'_2, M'_2; L_2, P_2, Y'_2; N', R, P'_1, Y'_1, P'_2, Y'_2)}_{\bar{I}'_1 \bar{I}'_2 \bar{I}'_3 \bar{I}'_4 \bar{I}'_5 \bar{Y}} \text{ :-} \\
& \underbrace{c(C'_1, M'_1, 'R^2'), 0(O'_1, C'_1, D'_1), \text{LI}(O'_1, L'_1, P'_1, Y'_1), 0A(O'_1, A), D(D'_1, R), A(A', N')}_{\bar{I}'_1} \}_{QAS_1 \bowtie A} \\
& \underbrace{c(C'_2, M'_2, 'C^1'), 0(O'_2, C'_2, D'_2), \text{LI}(O'_2, L'_2, P'_2, Y'_2), 0A(O'_2, A), D(D'_2, R) \}_{QAS_2}}_{\bar{I}'_2} \\
\\
Q_{72} & \underbrace{(A, N, R; D_1, O_1, C_1, M_1, D_2, O_2, C_2, M_2; L_1, P_1, Y_1; D_3, O_3, C_3, M_3, D_4, O_4, C_4, M_4; L_4, P_4, Y_4; N, R, P_1, Y_1, P_4, Y_4)}_{\bar{I}_1 \bar{I}_2 \bar{I}_3 \bar{I}_4 \bar{I}_5 \bar{Y}} \text{ :-} \\
& \underbrace{c(C_1, M_1, 'R^2'), 0(O_1, C_1, D_1), \text{LI}(O_1, L_1, P_1, Y_1), 0A(O_1, A), A(A, N), D(D_1, R), \}_{(AS_1 \bowtie D_1)}_{\bar{I}_1} \text{ avgRrsale} \\
& \underbrace{c(C_2, M_2, 'C^1'), 0(O_2, C_2, D_2), \text{LI}(O_2, L_2, P_2, Y_2), 0A(O_2, A), A(A, N), D(D_2, R), \}_{(AS_2 \bowtie D_2)}_{\bar{I}_2} \text{ avgRrsale} \\
& \underbrace{c(C_3, M_3, 'R^2'), 0(O_3, C_3, D_3), \text{LI}(O_3, L_3, P_3, Y_3), 0A(O_3, A), A(A, N), D(D_3, R), \}_{(AS_1 \bowtie D_1)}_{\bar{I}_3} \text{ avgCsale} \\
& \underbrace{c(C_4, M_4, 'C^1'), 0(O_4, C_4, D_4), \text{LI}(O_4, L_4, P_4, Y_4), 0A(O_4, A), A(A, N), D(D_4, R) \}_{(AS_2 \bowtie D_2)}_{\bar{I}_4 \bar{I}_5} \text{ avgCsale}
\end{aligned}$$

Figure 6.4: Queries  $Q_{70} := \text{ENCODE}(Q_1)$ ,  $Q_{71} := \text{ENCODE}(Q_3)$ , and  $Q_{72} := \text{chases}_{\Sigma}(Q_{70})$



The following two mappings  $h : Q_{72} \rightarrow Q_{71}$  and  $h' : Q_{71} \rightarrow Q_{72}$  are index-covering homomorphisms between the bnb-nb-normal forms of  $Q_{72}$  and  $Q_{71}$ .

$$\begin{aligned}
h &:= \{A/A', N/N', R/R', \\
&\quad C_1/C'_1, M_1/M'_1, O_1/O'_1, D_1/D'_1, L_1/L'_1, P_1/P'_1, Y_1/Y'_1, \\
&\quad C_2/C'_2, M_2/M'_2, O_2/O'_2, D_2/D'_2, L_2/L'_2, P_2/P'_2, Y_2/Y'_2, \\
&\quad C_3/C'_1, M_3/M'_1, O_3/O'_1, D_3/D'_1, L_3/L'_1, P_3/P'_1, Y_3/Y'_1, \\
&\quad C_4/C'_2, M_4/M'_2, O_4/O'_2, D_4/D'_2, L_4/L'_2, P_4/P'_2, Y_4/Y'_2 \} \\
h' &:= \{A'/A, N'/N, R'/R, \\
&\quad C'_1/C_1, M'_1/M_1, O'_1/O_1, D'_1/D_1, L'_1/L_1, P'_1/P_1, Y'_1/Y_1, \\
&\quad C'_2/C_4, M'_2/M_4, O'_2/O_4, D'_2/D_4, L'_2/L_4, P'_2/P_4, Y'_2/Y_4 \}
\end{aligned}$$

Homomorphisms  $h$  and  $h'$  together prove  $Q_{72} \stackrel{\Sigma}{\equiv}_{\text{bnb-nb}} Q_{71}$ . Because  $Q_{72} := \text{chase}_{\Sigma}(Q_{70})$  implies  $Q_{72} \stackrel{\Sigma}{\equiv}_{\xi} Q_{70}$  for any compatible signature  $\xi$ , this entails  $Q_{70} \stackrel{\Sigma}{\equiv}_{\text{bnb-nb}} Q_{71}$  and therefore  $Q_{68} \stackrel{\Sigma}{\equiv} Q_{69}$ . It follows that the SQL queries  $Q_1$  and  $Q_3$  are equivalent over all database instances conforming to the schema constraints in  $\Sigma$ .

We now turn our attention to Example 2 and show how the techniques in our thesis allow for simplification of the nested scalar aggregation in query  $Q_4$ .

**Example 64** Consider queries  $Q_4$  and  $Q_6$  from Example 2. If we model the output of aggregation function  $\text{MAX}$  as a set, and the output of  $\text{COUNT}(\ast)$  as a bag of zero-ary tuples, then  $Q_4$  and  $Q_6$  translate into the  $\text{COCQL}^{\{f_b, f_s, f_n\}}$  queries  $Q_{73}$  and  $Q_{74}$  shown in Figure 6.5 (thin and thick frames denote the expansions of logical and materialized views, respectively). Queries  $Q_{73}$  and  $Q_{74}$  both have output sort  $\{\{\text{dom}, \{\{\text{dom}\}\}\}\}$ , which we linearize into the chain sort  $(\text{bsb}, 1)$ .

The encoding query  $Q_{75} := \text{ENCODE}(Q_{73})$  is shown in Figure 6.6; observe that  $Q_{75}$  is a *hierarchical* CEQ due to the scalar aggregation operator in  $Q_4$ . As described in Chapter 5, before converting  $Q_{75}$  into bsb-normal form we must first convert its body into DACQ canonical form (see Chapter 4). Because there exists a homomorphism  $h_b : \text{body}_b^+ \rightarrow \text{body}_a^+$  that is the identity over  $\mathcal{B}_a^+$ ,

$$h_b = \{C/C, R/R, O_1/O_1, D_1/D_1, M_1/M_1, O_2/O_1, D_2/D_1\}$$

the conversion to DACQ canonical form merges node  $\boxed{b}$  into  $\boxed{a}$ . This yields the CEQ  $Q_{76}$  shown in Figure 6.6; note that because the body does not contain any hierarchical edges, we can drop the node labels from the query head. We have shaded the attributes that are removed from the head when  $Q_{76}$  is converted to bsb-NF.

Figure 6.6 also shows encoding query  $Q_{77} := \text{ENCODE}(Q_{74})$ ; the absence of scalar aggregation in  $Q_{74}$  means that  $Q_{77}$  does not contain any hierarchical edges, and so the attributes in the head do not require node labels. We have shaded the attributes that are removed from the head when  $Q_{77}$  is converted to bsb-NF.

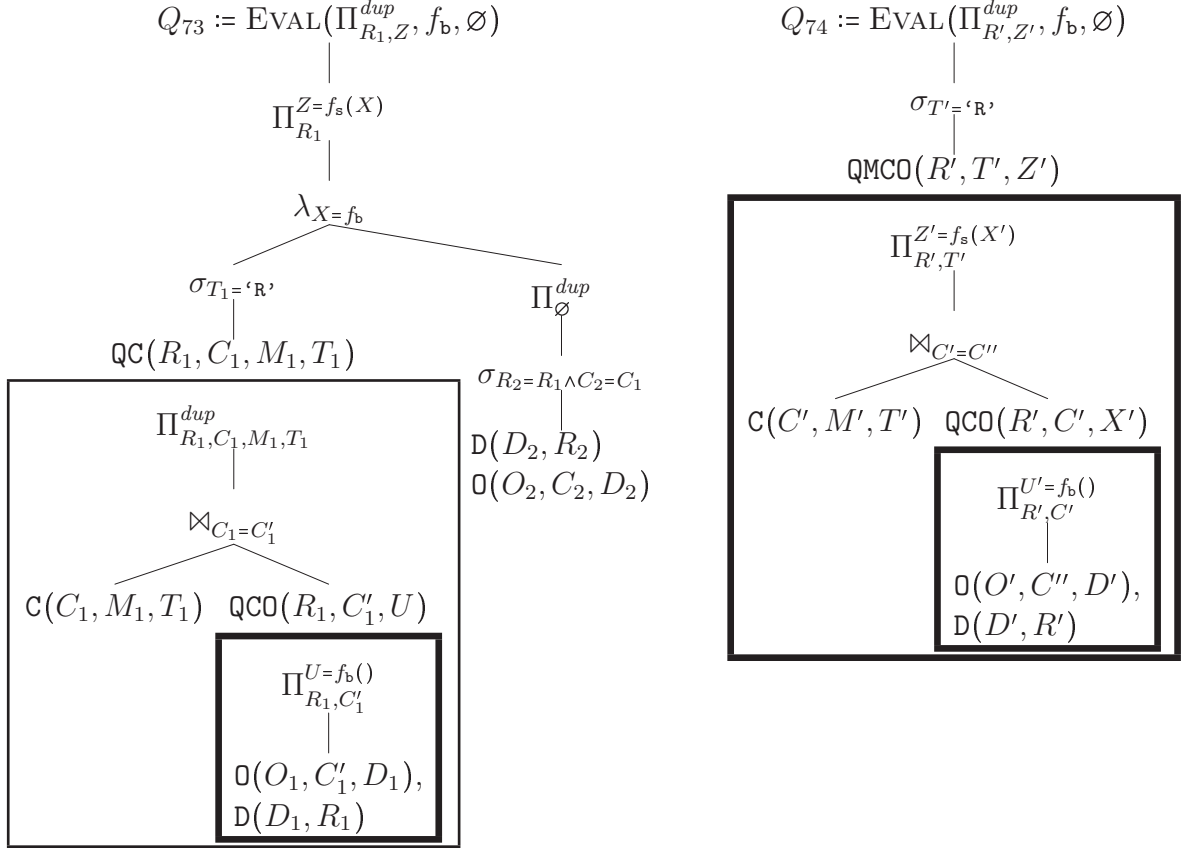


Figure 6.5:  $Q_{73}$  and  $Q_{74}$ , the COCQL $\{f_b, f_s, f_n\}$  translations of SQL queries  $Q_4$  and  $Q_6$

$$\begin{aligned}
Q_{75}(\overline{R^a}; \overline{C^a}, \overline{M_1^a}, \overline{R'}^a; \overline{D_2^b}, \overline{R^b}, \overline{O_2^b}, \overline{C^b}; \overline{R^a}) &:- \boxed{\overline{a} \text{O}(O_1, C, D_1), \text{D}(D_1, R), \text{C}(C, M_1, \overline{R'})} \\
&\quad \downarrow \\
&\quad \boxed{\overline{b} \text{O}(O_2, C, D_2), \text{D}(D_2, R)} \\
Q_{76}(\overline{R}; \overline{C}, \overline{M_1}, \overline{R'}; \overline{D_2}, \overline{R}, \overline{O_2}, \overline{C}; \overline{R}) &:- \text{O}(O_1, C, D_1), \text{D}(D_1, R), \text{C}(C, M_1, \overline{R'}), \\
&\quad \text{O}(O_2, C, D_2), \text{D}(D_2, R) \\
Q_{77}(\overline{R'}, \overline{R'}; \overline{C'}, \overline{M'}; \overline{D'}, \overline{O'}; \overline{R'}) &:- \text{O}(O', C', D'), \text{D}(D', R'), \text{C}(C', M', \overline{R'})
\end{aligned}$$

Figure 6.6: Encoding queries derived from  $Q_{73}$  and  $Q_{74}$

The following two mappings  $h : Q_{76} \rightarrow Q_{77}$  and  $h' : Q_{77} \rightarrow Q_{76}$  are index-covering homomorphisms between the bsb-normal forms of  $Q_{76}$  and  $Q_{77}$ .

$$\begin{aligned} h &:= \{O_1/O', C/C', D_1/D', R/R', M_1/M', O_2/O', D_2/D'\} \\ h' &:= \{O'/O_2, C'/C, D'/D_2, R'/R\} \end{aligned}$$

Homomorphisms  $h$  and  $h'$  together prove  $Q_{76} \dot{\equiv}_{\text{bsb}} Q_{77}$ . Because  $Q_{76}$  is the canonical form of  $Q_{75}$  and therefore equivalent to it, this entails  $Q_{75} \dot{\equiv}_{\text{bsb}} Q_{77}$  and therefore  $Q_{73} \equiv Q_{74}$ . It follows that the SQL queries  $Q_4$  and  $Q_6$  are equivalent.

## 6.3 Future Work

There are many directions in which the work in this thesis can be extended. We briefly mention some of them here.

**Schema Dependencies** Throughout Chapters 2–5 we made the assumption that the set of possible database instances was not constrained by any schema dependencies. Clearly, taking schema constraints into consideration is desirable for query optimization (cf. Examples 62 and 63). In Section 6.2 we briefly described how sets of FDs and acyclic INDs can be utilized when testing CEQ encoding-equivalence (if we restrict our attention to sets of primary and foreign keys, then we can drop the acyclicity requirement [67]). This technique can be extended to HCEQ encoding-equivalence in a straightforward manner: prior to the conversion to DACQ canonical form, the nodes in the query body are visited in any order that places ancestors before their descendants, and when each node  $v$  is visited the chase is invoked upon the CQ  $\text{body}_v^+$ .

From a practical perspective, the above technique is sufficient to handle most schema constraints permitted by commercial RDBMSs. An interesting question is whether our equivalence results can be extended to other types of schema constraints that guarantee terminating chases sequences. In particular, our characterization of *core indexes* in terms of MVDs assumes that all MVDs holding over the output must have originated from the query structure. If MVDs can also originate from the database schema, then the definition of core indexes breaks down—specifically, there is no longer a guarantee that there is one unique set of core indexes.

**Distributive Aggregation Functions** The most significant drawback of our modelling of aggregation functions as collection constructors is the implicit assumption that queries yielding objects of different sorts are not equivalent. For example, our techniques would fail to recognize that the following two queries are equivalent, because  $Q$  and  $Q'$  translate to COCQL $\{f_b, f_s, f_n\}$  queries with output sorts  $\{\{\text{dom}\}\}$  and  $\{\{\{\text{dom}\}\}\}$ , respectively.

```

Q:  SELECT SUM(qty)      Q': SELECT SUM(X)
     FROM LI              FROM (SELECT oid, SUM(qty) AS X
                              FROM LI
                              GROUP BY oid)

```

Many common SQL aggregation functions are *distributive* [48], and the ability to “re-aggregate” their output values is crucial for view rewriting within data warehousing environments. There are well-known algebraic transformations for manipulating operators that apply distributive aggregation [114]. A very practical extension to our work would be to define a normalization algorithm that performs the maximal flattening of distributive aggregation functions prior to applying the techniques from our thesis. In the case of query  $Q'$  above, a normalization algorithm would merge the two SUM aggregations, yielding query  $Q$ .

**Inequalities and Higher-Order Comparisons** Decision support workloads commonly include comparison predicates. While we have assumed an abstract domain of database constants, it would be very useful to extend our results to permit non-equality atomic comparisons over ordered concrete domains (e.g., integers or real numbers). We expect that this should be a relatively straightforward generalization of techniques known for extending CQ equivalence to atomic comparisons [70]. A more significant extension along the same vein would be to allow higher-order comparisons between aggregated values that occur either implicitly (within GROUP BY clauses) or explicitly (within WHERE or HAVING clauses). This would be particularly useful for optimizing decision support workloads, whose complex queries often contain predicates over the result of scalar aggregated subqueries (e.g., see the TPC-H [105] or TPC-DS [104] workloads). Explicit higher-order comparisons are known to make the equivalence problem undecidable with respect to the semantics of specific aggregation functions (both for numeric functions such as COUNT, SUM, etc. and for collection constructors such as the functions  $f_s$ ,  $f_b$ ,  $f_s$  we used in this thesis) [14, 52, 97]; however, it is possible that the equivalence problem with higher-order comparisons would remain decidable if we modelled aggregation functions as black boxes known only to satisfy certain arithmetic properties (e.g., see Cohen’s work on queries with non-nested aggregation [24]).

**Outer Joins and Null Values** In Chapter 4 we considered queries with explicit outer-join operators but no aggregation, whereas in the rest of the thesis we considered queries with aggregation but no explicit outer join. We expect that without too much difficulty our results could be extended to handle queries mixing both aggregation and explicit outer join, likely by using two types of hierarchical edges outputting distinct types of null values ( $\blacksquare$  vs. NULL). A more significant extension along the same vein would be to consider database instances that contain NULL values within the base relations/tables, along with queries that permit an explicit test for NULL. This extension was only recently studied for the CQ equivalence

problem by Farré et al. [34], and it would be interesting to determine whether their results can be extended to our encoding-equivalence problem.

**View Rewriting and Cost-based Optimization** As described in Chapter 1, an equivalence test is only one component of a view-rewriting algorithm. Given a query, our results in this thesis help to characterize the (infinite) space of equivalent logical rewritings. From a practical perspective, the most important next step is to use these results to design algorithms that generate logical query rewritings (either for the purposes of view rewriting or for query simplification). The design of any such algorithm and how it would fit into cost-based query optimization would be dependent upon the architecture of the target query optimizer, and thus this research needs to be conducted for each architecture. The results would be a set of cost-based optimizers that can efficiently process arbitrary queries with nested aggregation and outer joins in the presence of materialized views that also include nested aggregation and outer joins.



# Appendix A

## Proof of Theorem 3.2.4

In this section we formally prove that  $\bar{\xi}$ -equivalence of CEQs implies the existence of index-covering homomorphisms between  $\bar{\xi}$ -normal forms. Throughout this section we will make the following assumptions:

- [A1]  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  and  $Q'(\bar{\mathcal{I}}'_1; \dots; \bar{\mathcal{I}}'_d; \bar{\mathcal{V}}')$  are two CEQs satisfying  $Q \dot{\equiv}_{\bar{\xi}} Q'$ .
- [A2]  $Q$  and  $Q'$  are already in  $\bar{\xi}$ -NF (justified by Theorem 3.1.17).
- [A3] The body of  $Q$  is minimal relative to treating  $Q$  as a CQ with head  $Q(\bar{\mathcal{I}}_{[1,d]})$  (justified by  $(\mathcal{V} \cap \mathcal{B}) \subseteq \mathcal{I}_{[1,d]}$  and because minimization of CQs preserves relational equivalence). Similarly, the body of  $Q'$  is minimal relative to CQ head  $Q'(\bar{\mathcal{I}}'_{[1,d]})$ .
- [A4] Within the context of some database  $\mathbb{D}$ , we use  $\Gamma$  to denote both the embeddings of  $\text{body}_Q$  into  $\mathbb{D}$  and—by abuse of notation—the encoding relation  $(Q)^{\mathbb{D}}$  with schema  $Q(\bar{\mathcal{I}}_1; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$ . Then, given some value  $\bar{a} \in \text{dom}(\bar{\mathcal{I}}_{[1,l-1]})$ , we use  $\Gamma[\bar{a}]$  to denote the sub-relation indexed by  $\bar{a}$ , which has schema  $Q(\bar{\mathcal{I}}_l; \dots; \bar{\mathcal{I}}_d; \bar{\mathcal{V}})$  and instance  $\Pi_{\bar{\mathcal{I}}_{[l,d]}\bar{\mathcal{V}}}(\sigma_{\bar{\mathcal{I}}_{[1,l-1]}=\bar{a}}(\Gamma))$  (see Section 2.2.1). We define  $\Phi$  analogously for  $Q'$ .

The traditional proof of the Homomorphism Theorem for equivalence/containment of CQs [6, Theorem 6.2.3] can be broken down into the following four steps:

1. Use  $\text{body}_Q$  to construct a “canonical database”  $\mathbb{D}_Q$ .
2. Choose a particular embedding  $\gamma : \text{body}_Q \rightarrow \mathbb{D}_Q$  that yields a “canonical tuple”  $\gamma(\bar{\mathcal{V}})$  within  $(Q)^{\mathbb{D}_Q}$ .
3. Use the definition of query equivalence/containment to argue the existence of an embedding  $\phi : \text{body}_{Q'} \rightarrow \mathbb{D}_Q$  that yields the same canonical tuple  $\phi(\bar{\mathcal{V}}') = \gamma(\bar{\mathcal{V}})$ .
4. Define a mapping  $h$  in terms of  $\phi$ , and use both the definition of  $\mathbb{D}_Q$  and the properties of the chosen canonical tuple to prove that  $h$  is a homomorphism from  $Q'$  to  $Q$ .

At a high level, our proof of Theorem 3.2.4 uses the same four steps. Unfortunately, encoding-equality between encoding relations is characterized by the existence of a  $\bar{\xi}$ -certificate (see Theorem 2.2.9), which is a significantly more complicated

condition than equality/containment of two flat relations. This extra complexity arises for two distinct reasons—first, the recursive nature of comparing nested (encoded) objects rather than flat tuples; and second, the differing semantics required at each level of the comparison depending upon the encoding signature  $\bar{\xi}$ .

The complexity due to nested (encoded) objects forces two modifications. First, we need to generalize the idea of “equal canonical tuples”  $\gamma(\bar{\mathcal{V}}) = \phi(\bar{\mathcal{V}}')$  to “ $\bar{\xi}_{[i+1,d]}$ -equal canonical sub-relations”  $\Gamma[\bar{a}_{[1,i]}] \doteq_{\bar{\xi}_{[i+1,d]}} \Phi[\bar{a}'_{[1,i]}]$  (that is, corresponding sub-relations that encode the same “canonical sub-object”). Second, to ensure that  $h$  is index-covering requires demonstrating a relationship between corresponding index assignments  $\bar{a}_i$  and  $\bar{a}'_i$ , which requires proving properties about specific certificate nodes. To this end, Section A.1 defines certain degenerate classes of “pseudo-trivial” certificate nodes whose existence allows us to conclude the necessary properties for  $h$ .

To argue for the existence of a  $\bar{\xi}$ -certificate containing the necessary pseudo-trivial nodes, we require a carefully designed canonical database  $\mathbb{D}_Q$  and corresponding definitions of canonical sub-relations/sub-objects. This is where the complexity due to varying semantics manifests itself—each semantic type places different requirements on the design of  $\mathbb{D}_Q$ . Bag equivalence is based upon counting, and so for bag levels we rely upon certain algebraic properties of polynomials combined with a procedure to inflate cardinalities of database constants in a regular way such that non-pseudo-trivial certificate nodes would imply differing element cardinalities. In Section A.2 we describe this process, including the complete proof for the restricted case where every level requires bag semantics (i.e.,  $\bar{\xi} = \mathbf{b}^d$ ). Set equivalence is based upon existence, and so for set levels we inject large amounts of symmetry into the construction of  $\mathbb{D}_Q$ , with the constants corresponding to the index variables in  $Q$  forming the “axes” of the symmetry. We then show that non-pseudo-trivial certificate nodes would imply the existence/non-existence of extra/missing canonical sub-objects. This process is described in Section A.3, including the complete proof for the case  $\bar{\xi} = \mathbf{s}^d$ . Normalized bag equivalence exhibits both count- and existence-based features, and so for these levels our argument synthesizes both the polynomial and the symmetry techniques; Section A.4 details normalized bag equivalence, including the complete proof for the case  $\bar{\xi} = \mathbf{n}^d$ . Finally, Section A.5 shows how the three arguments can be interleaved to handle arbitrary signatures.

## A.1 Pseudo-Trivial Certificate Nodes

Given two encoding relations  $\Gamma$  and  $\Phi$ , consider any  $\bar{\xi}$ -certificate between them. A  $\bar{\xi}$ -certificate is a tree built out of set, bag, normalized bag, and tuple nodes (cf. Definitions 2.2.4–2.2.8). Choose any node  $n$  at any level  $i$  within the certificate; then, node  $n$  is itself the root of a  $\bar{\xi}_{[i,d]}$  certificate that proves  $\bar{\xi}_{[i,d]}$ -equality between two sub-relations  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$ , for some  $\bar{a} \in \text{adom}(\bar{\mathcal{I}}_{[1,i-1]}, \Gamma)$  and  $\bar{a}' \in \text{adom}(\bar{\mathcal{I}}'_{[1,i-1]}, \Phi)$ .



### A.1.1 Trivial Nodes

Regardless of the node type of  $n$  (i.e., set, bag, or normalized bag), the internal mechanisms of  $n$  effect one or more mappings from  $\text{adom}(\overline{\mathcal{I}}_i, \Gamma[\overline{a}])$  to  $\text{adom}(\overline{\mathcal{I}}'_i, \Phi[\overline{a}'])$  and from  $\text{adom}(\overline{\mathcal{I}}'_i, \Phi[\overline{a}'])$  to  $\text{adom}(\overline{\mathcal{I}}_i, \Gamma[\overline{a}])$ , and these mappings govern which  $\overline{\mathfrak{S}}_{[i+1,d]}$ -certificates node  $n$  must have as children. We say that  $n$  is *trivial w.r.t. index value*  $\overline{x}$  if any child  $\overline{\mathfrak{S}}_{[i+1,d]}$ -certificate of  $n$  that involves sub-relation  $\Gamma[\overline{a}\overline{x}]$  also involves sub-relation  $\Phi[\overline{a}'\overline{x}]$ , and vice versa (i.e., all of the mappings within  $n$  act as the identity for the value  $\overline{x}$ ). We say that  $n$  is *trivial* if it is trivial w.r.t. every value within  $\text{adom}(\overline{\mathcal{I}}_i, \Gamma[\overline{a}]) \cup \text{adom}(\overline{\mathcal{I}}'_i, \Phi[\overline{a}'])$ .

We say that relations  $\Gamma[\overline{a}]$  and  $\Phi[\overline{a}']$  *permit a proof that is trivial w.r.t. index value*  $\overline{x}$  if

1.  $\overline{x} \in \text{adom}(\overline{\mathcal{I}}_i, \Gamma[\overline{a}])$ ,
2.  $\overline{x} \in \text{adom}(\overline{\mathcal{I}}'_i, \Phi[\overline{a}'])$ , and
3.  $\Gamma[\overline{a}\overline{x}] \doteq_{\overline{\mathfrak{S}}_{[i+1,d]}} \Phi[\overline{a}'\overline{x}]$ .

If  $\Gamma[\overline{a}]$  and  $\Phi[\overline{a}']$  permit a proof that is trivial w.r.t. index value  $\overline{x}$ , we call the node  $n$  *potentially trivial w.r.t. index value*  $\overline{x}$ . It is straightforward to show that if node  $n$  is potentially trivial w.r.t.  $\overline{x}$  then we can re-arrange the mappings within  $n$  until we obtain a node that is trivial w.r.t. value  $\overline{x}$  (while still forming a valid  $\overline{\mathfrak{S}}_{[i,d]}$ -certificate between  $\Gamma[\overline{a}]$  and  $\Phi[\overline{a}']$ ).

We say that  $\Gamma[\overline{a}]$  and  $\Phi[\overline{a}']$  *permit a proof that is trivial* if for every value  $\overline{x}$  in  $\text{adom}(\overline{\mathcal{I}}_i, \Gamma[\overline{a}]) \cup \text{adom}(\overline{\mathcal{I}}'_i, \Phi[\overline{a}'])$  they permit a proof that is trivial w.r.t.  $\overline{x}$ . In this case, we call node  $n$  *potentially trivial*; again, it is straightforward to show for any potentially trivial node that we can re-arrange the mappings to obtain a trivial node.

### A.1.2 Pseudo-Trivial Nodes

Observe that  $\Gamma[\overline{a}]$  and  $\Phi[\overline{a}']$  only permit a trivial proof if  $\Gamma[\overline{a}] \doteq_{\mathfrak{b}\overline{\mathfrak{S}}_{[i+1,d]}} \Phi[\overline{a}']$ ; this follows because triviality implies equality (and hence bijectivity) between  $\text{adom}(\overline{\mathcal{I}}_i, \Gamma)$  and  $\text{adom}(\overline{\mathcal{I}}'_i, \Phi)$ . If  $\mathfrak{S}_i \neq \mathfrak{b}$ , proving that  $\Gamma[\overline{a}]$  and  $\Phi[\overline{a}']$  permit a trivial proof (and hence, that  $n$  can be reorganized into a trivial node) is challenging since it requires demonstrating bag-equality, which is an encoding relationship that is potentially stronger than the one enforced by  $n$ . For this reason, we want to relax the notion of a trivial proof/node.

We say that  $n$  is *pseudo-trivial w.r.t. value*  $\overline{x}$  if  $n$  contains a child  $\overline{\mathfrak{S}}_{[i+1,d]}$ -certificate that proves  $\Gamma[\overline{a}\overline{x}_1] \doteq_{\overline{\mathfrak{S}}_{[i+1,d]}} \Phi[\overline{a}'\overline{x}_2]$ , where  $\overline{x}_1 \in \text{adom}(\overline{\mathcal{I}}_i, \Gamma[\overline{a}])$  and  $\overline{x}_2 \in \text{adom}(\overline{\mathcal{I}}'_i, \Phi[\overline{a}'])$  are two (not necessarily distinct) permutations of  $\overline{x}$ . We say that  $n$  is *pseudo-trivial* if it is pseudo-trivial w.r.t. every value within  $\text{adom}(\overline{\mathcal{I}}_i, \Gamma[\overline{a}]) \cup \text{adom}(\overline{\mathcal{I}}'_i, \Phi[\overline{a}'])$ .

We say that  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$  *permit a proof that is pseudo-trivial w.r.t. value  $\bar{x}$*

1.  $\exists \bar{x}_1 \in \text{adom}(\bar{\mathcal{I}}_i, \Gamma[\bar{a}])$  that is a permutation of  $\bar{x}$ ,
2.  $\exists \bar{x}_2 \in \text{adom}(\bar{\mathcal{I}}'_i, \Phi[\bar{a}'])$  that is a permutation of  $\bar{x}$ , and
3.  $\Gamma[\bar{a}\bar{x}_1] \doteq_{\bar{\mathfrak{s}}_{[i+1,d]}} \Phi[\bar{a}'\bar{x}_2]$ .

If  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$  permit a proof that is pseudo-trivial w.r.t. value  $\bar{x}$ , we say that node  $n$  is *potentially pseudo-trivial w.r.t.  $\bar{x}$*  because we can re-arrange the mappings within  $n$  to obtain a node that is pseudo-trivial w.r.t.  $\bar{x}$ .

Similarly, we say that  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$  *permit a proof that is pseudo-trivial* if for every value  $\bar{x}$  in  $\text{adom}(\bar{\mathcal{I}}_i, \Gamma[\bar{a}]) \cup \text{adom}(\bar{\mathcal{I}}'_i, \Phi[\bar{a}'])$  they permit a proof that is pseudo-trivial w.r.t.  $\bar{x}$ . In this case, we say that  $n$  is *potentially pseudo-trivial* because we can re-arrange the mappings within  $n$  to obtain a pseudo-trivial node.

Observe that pseudo-triviality allows the permutation orderings to vary between different values of  $\bar{x}$ . Because of this, the existence of a pseudo-trivial proof only implies *set-equality* between the collections encoded by  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$ . Set-equality is the weakest form of encoding-equality, and so proving that  $n$  can be re-organized into a pseudo-trivial node does not require demonstrating an encoding relationship stronger than the one already enforced by  $n$ .

**Lemma A.1.1** *Let  $I \in (\mathcal{V} \cap \mathcal{B})$  and  $I' \in (\mathcal{V}' \cap \mathcal{B}')$  be output variables that occur in corresponding positions within  $\bar{\mathcal{V}}$  and  $\bar{\mathcal{V}}'$ . Then,  $I$  and  $I'$  must be indexes at the same level.*

PROOF. By Lemma 3.0.1, we know  $\exists i, j \in [1, d]$  such that  $I \in \mathcal{I}_i$  and  $I' \in \mathcal{I}'_j$ .

Suppose that  $i > j$ . Then, it is easy to construct a database  $\mathbb{D}$  such that for some  $\bar{a} \in \text{adom}(\bar{\mathcal{I}}_{[1,j]}, \Gamma)$ , sub-relation  $\Gamma[\bar{a}]$  contains two tuples with differing values for output column  $I$ . In contrast, for every  $\bar{a}' \in \text{adom}(\bar{\mathcal{I}}'_{[1,j]}, \Phi)$  the output column  $I'$  is constant across all tuples in  $\Phi[\bar{a}']$ , implying  $\Gamma[\bar{a}] \not\#_{\bar{\mathfrak{s}}_{[j+1,d]}} \Phi[\bar{a}']$  (a contradiction). Hence, conclude  $i \not> j$ . Conclude  $i \not< j$  similarly.  $\square$

**Corollary A.1.2** *If for some index level  $i \in [1, d]$  both  $\mathcal{I}_i \subseteq \mathcal{V}$  and  $\mathcal{I}'_i \subseteq \mathcal{V}'$ , then for any database  $\mathbb{D}$  and any  $\bar{\mathfrak{s}}$ -certificate between  $\Gamma$  and  $\Phi$ , all of the certificate nodes at level  $i$  are pseudo-trivial.*

## A.2 Bag Equivalence: A Counting Argument

Although bag nodes are defined in terms of a bijective function between index values, their effect can be equivalently stated in terms of counting—bag nodes effect agreement of element cardinalities. In order to argue that within a given  $\bar{\mathfrak{s}}$ -certificate a particular bag node is trivial, we use a counting argument based upon uniqueness of polynomials. The following example illustrates the core idea.

Student			

Student			

(a)  $\mathbb{D}_{10} = \{\text{Student}\}$ (b)  $\mathbb{D}_{11} = \{\text{Student}\}$ Figure A.1: Databases  $\mathbb{D}_{10}$  and  $\mathbb{D}_{11}$  for Example 65

$R_{21}$			

$R_{22}$				

$R_{23}$				

Figure A.2: Query results  $R_{21} = (Q_{78})^{\mathbb{D}_{10}}$ ,  $R_{22} = (Q_{79})^{\mathbb{D}_{10}}$ , and  $R_{23} = (Q_{80})^{\mathbb{D}_{10}}$ .

**Example 65** Consider the following three conjunctive encoding queries

$$Q_{78}(A, C; C) \text{ :- Student}(A, B, C)$$

$$Q_{79}(A, B, C; C) \text{ :- Student}(A, B, C)$$

$$Q_{80}(B, C, D; C) \text{ :- Student}(A, B, C), \text{Student}(A, D, C)$$

and the database  $\mathbb{D}_{10}$  shown in Figure A.1a. The results of evaluating the queries over  $\mathbb{D}_{10}$  are shown in Figure A.2. Clearly, the three encoding relations  $R_{21}$ ,  $R_{22}$ , and  $R_{23}$  are b-equal (and hence n-equal). Now consider database  $\mathbb{D}_{11}$  shown in Figure A.1b which was created from  $\mathbb{D}_{10}$  by “diversifying” the constant CS into two distinct constants (which we show as CS and **CS**). Under  $\mathbb{D}_{11}$ , no two of the resulting encoding relations (shown in Figure A.3) are b-equal (or n-equal).

In the above example, *diversifying* a single constant by a factor of two effectively breaks any nodes that prove the b-equality of the original encoding relations. (Because normalized bag nodes utilize a layer of bag sub-nodes, in this example the same process also breaks any normalized bag nodes that prove n-equality.) We formalize this process for bag equality next. Formalizing it for normalized bag equality requires additional mechanisms which we defer until Section A.4.

$R_{24}$			

$R_{25}$				

$R_{26}$				

Figure A.3: Query results  $R_{24} = (Q_{78})^{\mathbb{D}_{11}}$ ,  $R_{25} = (Q_{79})^{\mathbb{D}_{11}}$ , and  $R_{26} = (Q_{80})^{\mathbb{D}_{11}}$ .

## A.2.1 Colour-Diversification

Suppose that we are given an infinite palette of colours, each indexed by a positive integer:

$$\text{colour}_1 \text{ colour}_2 \text{ colour}_3 \text{ colour}_4 \text{ colour}_5 \dots$$

(It is intentional and significant that  $\text{colour}_1$  appears transparent.)

Let  $\widehat{\text{dom}} \subset \text{dom}$  be any set of  $n$  constants adhering to some total ordering  $<$ .

$$\widehat{\text{dom}} = \{c_1, \dots, c_n\} \quad \forall 1 \leq i < j \leq n : c_i < c_j$$

For each integer  $i \geq 2$ , let  $\widehat{\text{dom}}_i$  be a fresh set of constants isomorphic to  $\widehat{\text{dom}}$ , and let  $\delta_i : \widehat{\text{dom}} \rightarrow \widehat{\text{dom}}_i$  be a function that “paints” each  $c_j \in \widehat{\text{dom}}$  with colour  $\text{colour}_i$  to yield the constant  $c_j \in \widehat{\text{dom}}_i$ . As implied by the transparency of  $\text{colour}_1$ , we define  $\widehat{\text{dom}}_1 = \widehat{\text{dom}}$  and the painting function  $\delta_1 : \widehat{\text{dom}} \rightarrow \widehat{\text{dom}}_1$  as the identity function. Finally, because the different paintings of  $\widehat{\text{dom}}$  are mutually disjoint, we can define a single whitewash<sup>1</sup> function  $\delta^{-1} : \bigcup_i \widehat{\text{dom}}_i \rightarrow \widehat{\text{dom}}$  as follows.

$$\forall i \forall c_j \in \widehat{\text{dom}}_i : \delta^{-1}(c_j) := c_j \tag{A.1}$$

We now describe a procedure we call “colour-diversification” which we define for individual tuples, but generalize to arbitrary sets of tuples (and hence relations and databases).

**Definition A.2.1 (Diversification Point)** *Given totally-ordered domain  $\widehat{\text{dom}}$  of  $n$  constants, a diversification point for  $\widehat{\text{dom}}$  is any  $n$ -dimensional point  $\bar{r} \in \mathbb{N}^n$  whose coordinates are positive integers.*

**Definition A.2.2 ( $\bar{r}$ -Colour-Diversification)** *Given domain  $\widehat{\text{dom}}$ , diversification point  $\bar{r}$  for  $\widehat{\text{dom}}$ , and any tuple  $t$  over  $\widehat{\text{dom}}$ ,*

$$t = \langle c_{i_1}, c_{i_2}, \dots, c_{i_m} \rangle$$

*the colour-diversification of  $t$  w.r.t.  $\bar{r}$ —denoted  $\Delta^{\bar{r}}(t)$ —is the set of all possible “paintings” of  $t$  generated by independently choosing for each tuple component  $c_{i_j}$  one of the first  $r_{i_j}$  colours in the colour palette. We also refer to  $\Delta^{\bar{r}}(t)$  the  $\bar{r}$ -colour-diversification of  $t$  or the  $\bar{r}$ -diversification of  $t$ .*

We can colour-diversify an entire database  $\mathbb{D}_i$  by colour-diversifying each tuple  $t \in \mathbb{D}_i$  with the same point  $\bar{r}$  (where  $\bar{r}$  is a diversification point for active domain  $\text{adom}(\mathbb{D}_i)$ ). We then obtain a new database  $\mathbb{D}_j$

$$\mathbb{D}_j = \Delta^{\bar{r}}(\mathbb{D}_i) = \bigcup_{t \in \mathbb{D}_i} \Delta^{\bar{r}}(t)$$

---

<sup>1</sup>The inverse of painting

Student	school	dept	name
	UW	CS	Dave
	UW	CS	Dave
	UW	CS	Dave
	UW	Art	Art
	UW	Art	Art
	UW	Art	Art
	UW	Art	Art
	UW	Art	Art
	UW	CS	Dave
	UW	CS	Dave
	UW	CS	Dave
	UW	Art	Art
	UW	Art	Art
	UW	Art	Art
	UW	Art	Art

Figure A.4: Database  $\Delta^{[2,1,3,2]}(\mathbb{D}_{10})$

which we call the  $\bar{r}$ -(*colour-*)*diversification* of  $\mathbb{D}_i$ . We also say that  $\mathbb{D}_j$  is *colour-diverse w.r.t. point  $\bar{r}$* —a.k.a.  *$\bar{r}$ -colour-diverse* or  *$\bar{r}$ -diverse*—if  $\mathbb{D}_j$  is the result of performing  $\bar{r}$ -colour-diversification on some other database  $\mathbb{D}_i$ . Because colour-diversification is defined in terms of painting, the whitewashing function  $\delta^{-1}$  from equation A.1 is also the inverse of  $\bar{r}$ -diversification (when extended to sets of tuples). Hence, given a database  $\mathbb{D}_j$  that is  $\bar{r}$ -colour-diverse, we can always re-obtain the original database  $\mathbb{D}_i$  by whitewashing  $\mathbb{D}_j$ .

$$\mathbb{D}_i = \delta^{-1}(\mathbb{D}_j) = \delta^{-1} \circ \Delta^{\bar{r}}(\mathbb{D}_i)$$

Note that when we say  $\mathbb{D}_j$  is  $\bar{r}$ -colour-diverse, it implies that  $\bar{r}$  is a diversification point for active domain  $\text{adom}(\delta^{-1}(\mathbb{D}_j)) = \text{adom}(\mathbb{D}_i)$ , *not* active domain  $\text{adom}(\mathbb{D}_j)$  (which is already colour-diverse).

**Example 66** Consider database  $\mathbb{D}_{10}$  in Figure A.1a which contains four distinct constants in its active domain. Assume the following ordering for  $\text{adom}(\mathbb{D}_{10})$ .

$$\text{adom}(\mathbb{D}_{10}) = \{\text{UW}, \text{CS}, \text{Dave}, \text{Art}\}$$

Then, database  $\mathbb{D}_{11}$  in Figure A.1b is the  $[1, 2, 1, 1]$ -diversification of  $\mathbb{D}_{10}$  because constant  $c_2 = \text{CS}$  is painted with two colours. As another example, the  $[2, 1, 3, 2]$ -diversification of  $\mathbb{D}_{10}$  is shown in Figure A.4.

The cardinality of a tuple's colour-diversification depends upon both  $\bar{r}$  and the number of occurrences of each constant  $c_i \in \widehat{\text{dom}}$  within the tuple. For a given tuple  $t$ , let  $\#(t, c_i)$  denote the number of occurrences of constant  $c_i$  within  $t$ . Then, the set  $\Delta^{\bar{r}}(t)$  has the following cardinality.

$$|\Delta^{\bar{r}}(t)| = \prod_{c_i \in \widehat{\text{dom}}} (r_i)^{\#(t, c_i)} \quad (\text{A.2})$$

Observe that if tuple  $t$  is fixed, then each value  $\#(t, \mathbf{c}_i)$  is also fixed; therefore, equation A.2 defines a monomial of variables  $r_1, \dots, r_n$  with coefficient one and degree equal to the arity of  $t$ . Given an arbitrary set  $S$  of tuples over  $\widehat{\text{dom}}$ , the cardinality of the colour-diversification of  $S$  is then given by the following multivariate polynomial  $f_S$  over variables  $\bar{r}$ .

$$f_S(\bar{r}) = |\Delta^{\bar{r}}(S)| = \left| \bigcup_{t \in S} \Delta^{\bar{r}}(t) \right| = \sum_{t \in S} |\Delta^{\bar{r}}(t)| = \sum_{t \in S} \prod_{\mathbf{c}_i \in \widehat{\text{dom}}} (r_i)^{\#(t, \mathbf{c}_i)} \quad (\text{A.3})$$

Furthermore, given two arbitrary sets  $S_1$  and  $S_2$  of tuples over  $\widehat{\text{dom}}$ , equation A.3 implies the following relationships between polynomials  $f_{S_1}$  and  $f_{S_2}$ .

$$f_{S_1}(\bar{r}) + f_{S_2}(\bar{r}) = f_{(S_1 \cup S_2)}(\bar{r}) \quad (\text{A.4})$$

$$f_{S_1}(\bar{r}) \times f_{S_2}(\bar{r}) = f_{(S_1 \times S_2)}(\bar{r}) \quad (\text{A.5})$$

**Example 67** Consider again Example 66 and observe that database  $\mathbb{D}_{10}$  contains two tuples  $t_1 := \text{Student}(\text{UW}, \text{CS}, \text{Dave})$  and  $t_2 := \text{Student}(\text{UW}, \text{Art}, \text{Art})$ . The colour-diversification of  $t_1$  and  $t_2$  yields the following monomials

$$\begin{aligned} f_{t_1}(\bar{r}) &= |\Delta^{\bar{r}}(\langle \text{UW}, \text{CS}, \text{Dave} \rangle)| = (r_1)^1 (r_2)^1 (r_3)^1 (r_4)^0 = (r_1)(r_2)(r_3) \\ f_{t_2}(\bar{r}) &= |\Delta^{\bar{r}}(\langle \text{UW}, \text{Art}, \text{Art} \rangle)| = (r_1)^1 (r_2)^0 (r_3)^0 (r_4)^2 = (r_1)(r_4)^2 \end{aligned}$$

and so the  $[2, 1, 3, 2]$ -diversification of  $\mathbb{D}_{10}$  (shown in Figure A.4) has the following cardinality.

$$\begin{aligned} f_{\mathbb{D}_{10}}([2, 1, 3, 2]) &= f_{t_1}([2, 1, 3, 2]) + f_{t_2}([2, 1, 3, 2]) \\ &= 2 \times 1 \times 3 + 2 \times 2^2 \\ &= 14 \end{aligned}$$

The symmetry that colour-diversification injects into the base relations of a database also propagates to the result of any queries. More specifically, let  $\mathbb{D}$  be any database,  $\bar{r}$  any diversification point for  $\text{adom}(\mathbb{D})$ , and  $Q''$  be any conjunctive query over  $\mathbb{D}$  that does not contain constants in the query head (i.e.,  $Q''$  does not output constant columns). Then, the output of  $Q''$  exhibits the following characteristic.

$$(Q'')^{\Delta^{\bar{r}}(\mathbb{D})} = \Delta^{\bar{r}}((Q'')^{\mathbb{D}}) \quad (\text{A.6})$$

The reader can verify that  $[1, 2, 1, 1]$ -diversifying the query results in Figure A.2 yields the relations in Figure A.3.

Similarly, by whitewashing the result of any query over a colour-diverse database, we obtain the result of that query over the original database.

$$\delta^{-1} \left( (Q'')^{\Delta^{\bar{r}}(\mathbb{D})} \right) = \delta^{-1} \circ \Delta^{\bar{r}} \left( (Q'')^{\mathbb{D}} \right) = (Q'')^{\mathbb{D}} \quad (\text{A.7})$$

The reader can verify that whitewashing the query results in Figure A.3 yields the relations in Figure A.2.

## A.2.2 $k$ -Distinguishing Diversification

Examples 65 and 66 illustrate that by colour-diversifying a database we can differentiate queries that are not equivalent; however, the choice of the diversification point is important. By choosing diversification point  $[1, 2, 1, 1]$  for  $\mathbb{D}_{10}$  we distinguished between the results of  $Q_{78}$ ,  $Q_{79}$ , and  $Q_{80}$  under  $\mathbf{b}$ -equality (and  $\mathbf{n}$ -equality). However, the reader can verify that point  $[1, 2, 2, 2]$  fails to distinguish between any of the queries under  $\mathbf{b}$ -equality (or  $\mathbf{n}$ -equality). Ideally, we would like to show that we can always use colour-diversification to distinguish between non-equivalent queries.

**Definition A.2.3 ( $k$ -Distinguishing Diversification Point)** *Given any positive integer  $k$  and any diversification point  $\bar{r}$  for domain  $\widehat{\mathbf{dom}}$ , we say that  $\bar{r}$  is  $k$ -distinguishing if for any two sets  $S_1, S_2$  composed of tuples over  $\widehat{\mathbf{dom}}$  with maximum arity  $k$ , either  $f_{S_1} = f_{S_2}$  or  $f_{S_1}(\bar{r}) \neq f_{S_2}(\bar{r})$ .*

**Definition A.2.4 ( $k$ -Distinguishing Database)** *We say that database  $\mathbb{D}$  is  $k$ -distinguishing if there exists some  $k$ -distinguishing diversification point  $\bar{r}$  for domain  $\delta^{-1}(\mathbf{adom}(\mathbb{D}))$  such that  $\mathbb{D}$  is colour-diverse w.r.t.  $\bar{r}$ .*

In light of the polynomials within Definition A.2.3, we turn our attention to known results about equivalence of polynomials. The Fundamental Theorem of Algebra (proved by Gauss) states that a univariate polynomial of degree  $d$  has  $d$  roots (some of which might be degenerate). The Schwartz-Zippel Theorem is essentially a generalization of this result to multivariate polynomials. The following formulation of the theorem is due to Rudich [98].

**Theorem A.2.5 (Schwartz-Zippel)** *Let  $f(x_1, \dots, x_n) \neq 0$  be a degree  $k$  multivariate polynomial over some field  $F$ . For any finite set  $S \subset F$ , if you pick random (not necessarily distinct)  $r_1, r_2, \dots, r_n \in S$  then*

$$\Pr[f(r_1, r_2, \dots, r_n) = 0] \leq \frac{k}{|S|}.$$

**Corollary A.2.6** *Let  $F = \{f_1, f_2, \dots, f_m\}$  be any finite set of distinct multivariate polynomials over the real numbers with maximum degree  $k$  and  $m \geq 2$ . Let  $S \subset \mathbb{N}$  be any finite set of positive integers. If you pick a random point  $\bar{r} \in S^n$ , then*

$$\Pr[\exists i, j \in [1, m]. (i < j \wedge f_i(\bar{r}) = f_j(\bar{r}))] < \frac{km^2}{2|S|}$$

**Corollary A.2.7** *For any finite totally-ordered domain  $\widehat{\mathbf{dom}}$  and non-negative integer  $k$ , there exists an infinite number of  $k$ -distinguishing diversification points.*

**Corollary A.2.8** *Given any database  $\mathbb{D}$  and any non-negative integer  $k$ , there exists a colour-diversification of  $\mathbb{D}$  that is  $k$ -distinguishing.*

### A.2.3 Guaranteeing Pseudo-Trivial Certificate Nodes

We now show how colour-diversification can be used to guarantee that any  $\bar{\xi}$ -certificate between the results of queries  $Q$  and  $Q'$  contains pseudo-trivial nodes. Consider any colour-diverse database  $\mathbb{D}$  and any  $\bar{\xi}$ -certificate proving  $\Gamma \doteq_{\bar{\xi}} \Phi$ . For any  $i \in [1, d]$ , let  $n$  be any node in level  $i$ ; then, node  $n$  proves  $\Gamma[\bar{a}] \doteq_{\bar{\xi}_{[i, d]}} \Phi[\bar{a}']$  for some  $\bar{a} \in \text{adom}(\bar{\mathcal{I}}_{[1, i-1]}, \Gamma)$  and  $\bar{a}' \in \text{adom}(\bar{\mathcal{I}}'_{[1, i-1]}, \Phi)$ . Regardless of the value of semantic indicator  $\bar{\xi}_i$ , node  $n$  must at the very least enforce that  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$  encode equivalent sets of sub-objects; therefore, let  $\text{subobjects}(n)$  denote the set of sub-objects encoded within the sub-relations of  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$ .

$$\text{subobjects}(n) := \text{DECODE}(\Gamma[\bar{a}], \mathbf{s}\bar{\xi}_{[i+1, d]}) = \text{DECODE}(\Phi[\bar{a}'], \mathbf{s}\bar{\xi}_{[i+1, d]}) \quad (\text{A.8})$$

For each sub-object  $o \in \text{subobjects}(n)$ , let  $\Gamma[\bar{a}]|_o$  denote the subset of  $\Gamma[\bar{a}]$  containing only the embeddings that correspond to an encoding of sub-object  $o$ .

$$\Gamma[\bar{a}]|_o := \{\gamma \mid \exists \bar{b} \in \text{adom}(\bar{\mathcal{I}}_i, \Gamma[\bar{a}]). (\gamma \in \Gamma[\bar{a}\bar{b}] \wedge \text{DECODE}(\Gamma[\bar{a}\bar{b}], \bar{\xi}_{[i+1, d]}) = o)\} \quad (\text{A.9})$$

When viewed as an encoding relation,  $\Gamma[\bar{a}]|_o$  has the same schema as  $\Gamma[\bar{a}]$ . Define relation  $\Phi[\bar{a}']|_o$  analogously.

In order to invoke count-based reasoning, we need to model sub-object cardinalities as polynomials. Consider the relation  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}])$ —because  $\mathbb{D}$  is colour-diverse w.r.t. some point  $\bar{r}$ , it follows from equation A.6 that  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}])$  is  $\bar{r}$ -diverse. Now, for each sub-object  $o \in \text{subobjects}(n)$ , consider the relation  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)$ —unfortunately, this relation *is not* necessarily  $\bar{r}$ -diverse, because all of the values in attributes  $\mathcal{I}_i \cap \mathcal{V}$  will be constant (they are functionally determined by  $o$ ). However, if we define  $\bar{\mathcal{J}}_i$  as the non-output columns in  $\bar{\mathcal{I}}_i$ , then the relation  $\Pi_{\bar{\mathcal{J}}_i}(\Gamma[\bar{a}]|_o)$  *is*  $\bar{r}$ -diverse. Next, define set  $S_o$  as the whitewashing of the tuples in this relation instance

$$S_o := \delta^{-1}(\mathbb{T}(\Pi_{\bar{\mathcal{J}}_i}(\Gamma[\bar{a}]|_o))) \quad (\text{A.10})$$

and define polynomial  $f_{S_o}$  as in equation A.3. Finally, define  $\bar{\mathcal{J}}'_i$ ,  $S'_o$ , and  $f_{S'_o}$  analogously for  $\Phi[\bar{a}']$ . If we let  $\#(\Gamma[\bar{a}], o)$  and  $\#(\Phi[\bar{a}'], o)$  denote the frequency of sub-object  $o$  within the encoding relations  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$ , then the following equations hold by construction and the fact that  $\mathbb{D}$  is  $\bar{r}$ -diverse.

$$\#(\Gamma[\bar{a}], o) = |\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)| = |\Pi_{\bar{\mathcal{J}}_i}(\Gamma[\bar{a}]|_o)| = |\Delta^{\bar{r}}(S_o)| = f_{S_o}(\bar{r}) \quad (\text{A.11})$$

$$\#(\Phi[\bar{a}'], o) = |\Pi_{\bar{\mathcal{I}}'_i}(\Phi[\bar{a}']|_o)| = |\Pi_{\bar{\mathcal{J}}'_i}(\Phi[\bar{a}']|_o)| = |\Delta^{\bar{r}}(S'_o)| = f_{S'_o}(\bar{r}) \quad (\text{A.12})$$

Using the definitions above, we can now identify situations where certificates contain pseudo-trivial nodes. This is enough to prove Theorem 3.2.4 in the restricted case where every level requires bag semantics (i.e.,  $\bar{\xi} = \mathbf{b}^d$ ).

**Lemma A.2.9** *Given*

- *some database  $\mathbb{D}$  such that  $\Gamma$  and  $\Phi$  are non-empty,*



- some  $\bar{\xi}$ -certificate that proves  $\Gamma \doteq_{\bar{\xi}} \Phi$ , and
- some level  $i \in [1, d]$  such that  $\xi_i = \mathbf{b}$ ;

if  $\mathbb{D}$  is  $(|\mathcal{I}_{[1,d]}| + |\mathcal{I}'_{[1,d]}|)$ -distinguishing then every node  $n^{\mathbf{b}}$  in level  $i$  of the certificate is potentially pseudo-trivial.

PROOF. Suppose that  $\mathbb{D}$  is  $k$ -distinguishing for  $k = |\mathcal{I}_{[1,d]}| + |\mathcal{I}'_{[1,d]}|$ . If  $|\mathcal{I}_i| = |\mathcal{I}'_i| = 0$  then any node at level  $i$  is trivial; therefore assume  $|\mathcal{I}_i| + |\mathcal{I}'_i| > 0$ . Choose any  $n^{\mathbf{b}}$  in level  $i$  of the certificate. Then,  $n^{\mathbf{b}}$  proves  $\Gamma[\bar{a}] \doteq_{\mathbf{b}\bar{\xi}_{[i+1,d]}} \Phi[\bar{a}']$  for some  $\bar{a} \in \text{adom}(\bar{\mathcal{I}}_{[1,i-1]})$  and  $\bar{a}' \in \text{adom}(\bar{\mathcal{I}}'_{[1,i-1]})$ ;

For each  $o \in \text{subobjects}(n^{\mathbf{b}})$ , bag equivalence requires that  $\#(\Gamma[\bar{a}], o) = \#(\Phi[\bar{a}'], o)$ . Because  $\mathbb{D}$  is colour-diverse w.r.t. some  $k$ -distinguishing point  $\bar{r}$ , by equations A.11 and A.12, this implies  $f_{S_o}(\bar{r}) = f_{S'_o}(\bar{r})$ . Observe that sets  $S_o$  and  $S'_o$  both contain tuples whose arity does not exceed the value  $k$ ; hence,  $f_{S_o} = f_{S'_o}$  follows by Definitions A.2.3 and A.2.4.

By virtue of equations A.3 and A.2,  $f_{S_o} = f_{S'_o}$  implies that each tuple  $t \in S_o$  can be mapped to a permutation of itself in  $S'_o$ , and vice versa. (It additionally implies that there is a bijection between  $S_o$  and  $S'_o$  that satisfies this property, but we don't require that to prove pseudo-triviality). This further implies that a similar “permuting mapping” exists in each direction between the tuples of relations  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)$  and  $\Pi_{\bar{\mathcal{I}}'_i}(\Phi[\bar{a}']|_o)$ .

Lemma A.1.1 implies that the same set of constant columns both extends relation  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)$  into  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)$  and extends relation  $\Pi_{\bar{\mathcal{I}}'_i}(\Phi[\bar{a}']|_o)$  into  $\Pi_{\bar{\mathcal{I}}'_i}(\Phi[\bar{a}']|_o)$ . Hence, the “permuting mappings” between  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)$  and  $\Pi_{\bar{\mathcal{I}}'_i}(\Phi[\bar{a}']|_o)$  are easily extended into a “permuting mappings” between  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)$  and  $\Pi_{\bar{\mathcal{I}}'_i}(\Phi[\bar{a}']|_o)$ —in other words,  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$  permit pseudo-trivial proofs for all of the index values in  $\Pi_{\bar{\mathcal{I}}_i}(\Gamma[\bar{a}]|_o)$  and  $\Pi_{\bar{\mathcal{I}}'_i}(\Phi[\bar{a}']|_o)$ . Repeating this argument for each  $o \in \text{subobjects}(n^{\mathbf{b}})$ , we conclude that  $\Gamma[\bar{a}]$  and  $\Phi[\bar{a}']$  permit a pseudo-trivial proof (which means  $n^{\mathbf{b}}$  is potentially pseudo-trivial).  $\square$

**Lemma A.2.10** *If  $\bar{\xi} = \mathbf{b}^d$  then there exists index-covering homomorphisms in both directions between  $Q$  and  $Q'$ .*

PROOF. To construct an index-covering homomorphism  $h$  from  $Q'$  to  $Q$ , we start by defining the canonical database  $\mathbb{D}_Q$  as follows

$$\mathbb{D}_Q := \Delta^{\bar{r}}(\text{body}_Q) \tag{A.13}$$

where  $\bar{r}$  is any  $(|\mathcal{I}_{[1,d]}| + |\mathcal{I}'_{[1,d]}|)$ -distinguishing diversification point for (any total ordering of) the set  $\mathcal{B} \cup \mathcal{C}_Q$  (i.e., the variables and constants occurring in  $\text{body}_Q$ ).

Consider any  $\mathbf{b}^d$ -certificate proving  $\Gamma \doteq_{\mathbf{b}^d} \Phi$ . Choose any  $\gamma \in \Gamma$  that satisfies  $\delta^{-1} \circ \gamma(\bar{\mathcal{I}}_{[1,d]}) = \bar{\mathcal{I}}_{[1,d]}$  (by construction of  $\mathbb{D}_Q$ , there exist many such  $\gamma$ ). Let  $\bar{a}_{[1,d]}$  denote the tuple of index values  $\gamma(\bar{\mathcal{I}}_{[1,d]})$ . By Lemma A.2.9, the root node of the certificate is *potentially* pseudo-trivial w.r.t.  $\bar{a}_1$ ; therefore, re-arrange

the mappings within the node until it *is* pseudo-trivial w.r.t.  $\bar{a}_1$ . By recursing into any child  $\mathbf{b}^{d-1}$ -certificate involving sub-relation  $\Gamma[\bar{a}_1]$  and repeating this procedure, we eventually arrive at a tuple node proving  $\Gamma[\bar{a}_{[1,d]}] \doteq_{\emptyset} \Phi[\bar{a}'_{[1,d]}]$  for some  $\bar{a}'_{[1,d]} \in \mathbf{adom}(\bar{\mathcal{I}}'_{[1,d]}, \Phi)$ . Choose any  $\phi \in \Phi[\bar{a}'_{[1,d]}]$ .

Define  $h := \delta^{-1} \circ \phi$ . Then,  $h$  is a mapping from  $\mathcal{B}'$  to  $\delta^{-1}(\mathbf{adom}(\mathbb{D}_Q)) = \mathcal{B} \cup \mathcal{C}_Q$  that satisfies the following properties.

1.  $h(\mathbf{body}_{Q'}) = \delta^{-1} \circ \phi(\mathbf{body}_{Q'}) \subseteq \delta^{-1}(\mathbb{D}_Q) = \mathbf{body}_Q$
2.  $h(\bar{\mathcal{V}}') = \delta^{-1} \circ \phi(\bar{\mathcal{V}}') = \delta^{-1} \circ \gamma(\bar{\mathcal{V}}) = \bar{\mathcal{V}}$
3.  $\forall i \in [1, d], h(\mathcal{I}'_i) = \delta^{-1} \circ \phi(\mathcal{I}'_i) = \delta^{-1} \circ \gamma(\mathcal{I}_i) = \mathcal{I}_i$   
 which holds because the node proving  $\Gamma[\bar{a}_{[1,i-1]}] \doteq_{\bar{\mathcal{S}}_{[i,d]}} \Phi[\bar{a}'_{[1,i-1]}]$  was rearranged to be pseudo-trivial w.r.t.  $\bar{a}_i$  and so  $\phi(\bar{\mathcal{I}}'_i) = \bar{a}'_i$  is a permutation of  $\bar{a}_i = \gamma(\bar{\mathcal{I}}_i)$ .

Hence,  $h$  satisfies Definition 3.2.1. The proof is analogous for  $h'$  from  $Q$  to  $Q'$ .  $\square$

### A.3 Set Equivalence: A Symmetry Argument

The proof for bag equivalence in the previous section focuses on object frequencies without paying much attention to the objects themselves. In particular, within the proof of Lemma A.2.10 where we argue that any bag node involving sub-relation  $\Gamma[\bar{a}_{[1,i-1]}]$  must be pseudo-trivial w.r.t.  $\bar{a}_i$ , the argument complete ignores the actual sub-object encoded by sub-relation  $\Gamma[\bar{a}_{[1,i]}]$ . Instead, the focus is entirely upon the effect that index tuple  $\bar{a}_i$  has on the sub-object's frequency.

In contrast, the argument for set equivalence depends upon examining the contents of the encoded sub-objects. To argue that any set node involving some sub-relation  $\Gamma[\bar{a}_{[1,i-1]}]$  must be pseudo-trivial w.r.t. some  $\bar{a}_i$ , we use a carefully crafted “canonical sub-object” that can be uniquely identified only by the precise collection of constants within index tuple  $\bar{a}_{[1,i-1]}$ . We illustrate the basic idea using the following (rather lengthy) example. For the duration of this section we assume that the encoding signature contains only set semantics ( $\bar{\mathcal{S}} = \mathbf{s}^d$ ).

**Example 68** Consider a relation schema `Wrote(Name, Title, Year)` that models the fact that an author published a particular title in a particular year. Assume that each publication is uniquely identified by the combination of its title and year, so multiple tuples with the same title and year represent a co-authored publication. Figure A.5 shows a sample database instance  $\mathbb{D}_{12}$  containing relation `Wrote`.

Now consider the following CEQ  $Q_{81}$ . Figure A.5 shows the encoding relation  $R_{27}$  that results from evaluating  $Q_{81}$  over  $\mathbb{D}_{12}$ , as well as the object  $o_{16}$  that results from decoding  $R_{27}$  with signature `ss`.

$$Q_{81}(N, M, Y; T; T) :- \text{Wrote}(N, T, Y), \text{Wrote}(M, T, Y)$$

Wrote	Name	Title	Year
	Amy	A1	2007
	Amy	AB1	2007
	Amy	A2	2008
	Amy	AB2	2008
	Amy	AB3	2008
	Amy	AC1	2008
	Amy	AC2	2008
	Bob	B1	2007
	Bob	AB1	2007
	Bob	B2	2008
	Bob	AB2	2008
	Bob	AB3	2008
	Chad	C1	2008
	Chad	AC1	2008
	Chad	AC2	2008

$R_{27}$	N	M	Y	T	T
	Amy	Amy	2007	A1	A1
	Amy	Amy	2007	AB1	AB1
	Amy	Amy	2008	A2	A2
	Amy	Amy	2008	AB2	AB2
	Amy	Amy	2008	AB3	AB3
	Amy	Amy	2008	AC1	AC1
	Amy	Amy	2008	AC2	AC2
	Amy	Bob	2007	AB1	AB1
	Amy	Bob	2008	AB2	AB2
	Amy	Bob	2008	AB3	AB3
	Amy	Chad	2008	AC1	AC1
	Amy	Chad	2008	AC2	AC2
	Bob	Amy	2007	AB1	AB1
	Bob	Amy	2008	AB2	AB2
	Bob	Amy	2008	AB3	AB3
	Bob	Bob	2007	B1	B1
	Bob	Bob	2007	AB1	AB1
	Bob	Bob	2008	B2	B2
	Bob	Bob	2008	AB2	AB2
	Bob	Bob	2008	AB3	AB3
	Chad	Amy	2008	AC1	AC1
	Chad	Amy	2008	AC2	AC2
	Chad	Chad	2008	C1	C1
	Chad	Chad	2008	AC1	AC1
	Chad	Chad	2008	AC2	AC2

$$o_{16} = \{ \{ A1, AB1 \}, \{ A2, AB2, AB3, AC1, AC2 \}, \{ AB1 \}, \{ AB2, AB3 \}, \{ AC1, AC2 \}, \{ B1, AB1 \}, \{ B2, AB2, AB3 \}, \{ C1, AC1, AC2 \} \}$$

Figure A.5: Database  $D_{12} = \{\text{Wrote}\}$ , query result  $R_{27} = (Q_{81})^{D_{12}}$ , and object  $o_{16} = \text{DECODE}(R_{27}, ss)$ .

Intuitively, encoding query  $Q_{81}$  groups together publication titles related to pairs of (not necessarily distinct) author names and a particular year; hence, object  $o_{16}$  could be called the “set of annual co-authorships,” where each “annual co-authorship” is simply a set of related titles. The crucial point is that  $o_{16}$  does not actually contain the author names and year that relate the titles within a particular annual co-authorship.

Now consider three more CEQs that all contain fewer index variables than  $Q_{81}$ .

$$\begin{aligned} Q_{82}(N, Y; T; T) &:- \text{Wrote}(N, T, Y) \\ Q_{83}(T; ; T) &:- \text{Wrote}(N, T, Y) \\ Q_{84}(N, M; T; T) &:- \text{Wrote}(N, T, Y), \text{Wrote}(M, T, Y) \end{aligned}$$

Ideally, we would like  $o_{16}$  to contain a sub-object that can be uniquely identified only via the author pair and year that index it in  $R_{27}$ . Such a sub-object could be used to

distinguish  $o_{16}$  from the result of any CEQ whose resulting encoding relation does not contain an index tuple with the same combination of authors and year (including the three queries above).

Consider sub-object  $\{ A1, AB1 \}$  corresponding to index tuple Amy-Amy-2007 in  $R_{27}$ . The reader can easily verify that this sub-object also occurs in the result of  $Q_{82}$  corresponding to index tuple Amy-2007. The same argument applies to any sub-object corresponding to an index tuple in  $R_{27}$  with a repeated author name, including  $\{A2, AB2, AB3, AC1, AC2\}$ ,  $\{B1, AB1\}$ ,  $\{B2, AB2, AB3\}$ , and  $\{C1, AC1, AC2\}$ .

Consider sub-object  $\{ B1 \}$  corresponding to index tuples Amy-Bob-2007 and Bob-Amy-2007 in  $R_{27}$ . Clearly, this sub-object also occurs in the result of  $Q_{83}$  and the same argument would apply to any other singleton set.

Consider sub-object  $\{ AC1, AC2 \}$  corresponding to index tuples Amy-Chad-2008 and Chad-Amy-2008 in  $R_{27}$ . This sub-object also occurs in the result of  $Q_{84}$  corresponding to index tuples Amy-Chad and Chad-Amy.

Finally, consider sub-object  $\{ AB2, AB3 \}$  corresponding to index tuples Amy-Bob-2008 and Bob-Amy-2008 in  $R_{27}$ . Because the sub-object is not singleton it does not occur in the result of  $Q_{83}$ , and it also does not occur in the results of  $Q_{82}$  or  $Q_{84}$  because each of those encoding relations only contains index tuples with two out of the three necessary values needed to isolate the set  $\{ AB2, AB3 \}$ .

Unfortunately, sub-object  $\{ AB2, AB3 \}$  is not sufficient to distinguish object  $o_{16}$  from the result of *any* CEQ. Consider the following query  $Q_{85}$ , which is also in ss-NF.

$$Q_{85}(T', T''; T; T) := \text{Wrote}(N, T, Y), \text{Wrote}(M, T, Y), \\ \text{Wrote}(N, T', Y), \text{Wrote}(M, T'', Y)$$

The reader can verify that the relation  $(Q_{85})^{\mathbb{D}_{12}}$  contains an encoding of sub-object  $\{ AB2, AB3 \}$  (corresponding to index tuples A2-B2 and B2-A2). Query  $Q_{85}$  exploits the fact that within  $\mathbb{D}_{12}$ , occurrences of the values A2 and B2 are correlated with occurrences of other constants. Specifically, any embedding of  $Q_{85}$  into  $\mathbb{D}_{12}$  that maps attribute  $T'$  to A2 must also map  $N$  to Amy and  $Y$  to 2008; similarly, mapping  $T''$  to B2 forces mappings of  $M$  to Bob and  $Y$  to 2008. By “triangulating” index values A2 and B2,  $Q_{85}$  effectively isolates the values Amy, Bob, and 2008 without explicitly naming them. Indiscriminately adding a few tuples with titles A2 or B2 to the database instance will not necessarily prevent this, either. For example, adding the tuples  $\text{Wrote}(\text{Don}, \text{A2}, 2008)$  and  $\text{Wrote}(\text{Chad}, \text{A2}, 2007)$  to  $\mathbb{D}_{12}$  does not change the sub-object corresponding to A2-B2 or B2-A2 in  $(Q_{85})^{\mathbb{D}_{12}}$ , because these tuples do not increase the possible choices for authors who co-authored with Bob in 2008. However, by adding tuple  $\text{Wrote}(\text{Chad}, \text{B2}, 2008)$  to  $\mathbb{D}_{12}$  we allow variable  $M$  to map to either Bob or Chad, which changes the sub-object A2-B2 (and B2-A2) in  $(Q_{85})^{\mathbb{D}_{12}}$  to  $\{ AB2, AB3, AC1, AC2 \}$  without affecting the value of sub-object Amy-Bob-2008/Bob-Amy-2008 in  $(Q_{27})^{\mathbb{D}_{12}}$ .

Example 68 only considers encoding queries/relations of depth two, which is actually the simplest interesting case for nested sets (because a depth one query in

s-NF is essentially just a CQ). Nevertheless, from the example we can still draw several lessons about the necessary features of any canonical sub-object and the implications for the structure of the database instance that yields it; these lessons generalize to considering canonical sub-objects of varying depth. First of all, the canonical object needs to be a strict subset of the objects formed when one or more of the index columns is removed from the query head. Second, the canonical object cannot be a singleton set, or else it could be constructed by grouping on values from inner index columns. Finally, the database needs to contain enough symmetry to guarantee that isolating an index value via data triangulation requires more triangulation points than there are index columns in either  $Q$  or  $Q'$ .

### A.3.1 Labels, Prefixes and Canonical Database $\mathbb{D}_Q$

We now describe a method for constructing canonical database  $\mathbb{D}_Q$  from the body of query  $Q$  that contains large amounts of symmetry. The symmetry is designed so that each constant is related to a large number of other constants in a symmetric manner, which forms the basis of the argument against isolating particular index values via triangulation. Whereas in Section A.2 we “painted” with different colours to generate multiple copies of database constants, in this section we will embellish database constants with labels which we denote as superscripts.

Choose  $N$  to be any integer satisfying the following constraint.

$$N > \max(|\mathcal{I}_{[1,d]}|, |\mathcal{I}'_{[1,d]}|) + 1 \quad (\text{A.14})$$

The symmetry will be specified using mechanisms we will call *label-generating components*, *sequences*, and *prefixes*. For each level  $i \in [1, d]$ , let  $LGC_i$  denote the set of *label-generating components* at level  $i$ , defined as follows.

$$LGC_i := \{(\bar{y}_i, \bar{z}_i) \mid \bar{y}_i, \bar{z}_i \in [1, N]^{|\mathcal{I}_i|}\} \quad (\text{A.15})$$

We say that component  $c = (\bar{y}_i, \bar{z}_i)$  contains a *conflict at position  $i.j$*  if  $y_{i.j} = z_{i.j}$ , and we use  $CF-LGC_i$  to denote the conflict-free subset of  $LGC_i$ .

$$CF-LGC_i := LGC_i \setminus \{(\bar{y}_i, \bar{z}_i) \mid \exists j. (y_{i.j} = z_{i.j})\} \quad (\text{A.16})$$

Observe that set  $LGC_i$  has size  $N^{2|\mathcal{I}_i|}$ , while set  $CF-LGC_i$  has size  $N^{|\mathcal{I}_i|}(N-1)^{|\mathcal{I}_i|}$ .

**Example 69** Suppose  $Q = Q_{81}$  and  $N = 6$ . Because every integer in  $[1, N]$  is a single digit, we will represent each  $\bar{x} \in [1, N]^{|\mathcal{I}_i|}$  simply as a  $|\mathcal{I}_i|$ -digit number. Then, sets  $LGC_2$  and  $CF-LGC_2$  are as follows

$$\begin{aligned} LGC_2 &= \{(1, 1), \dots, (1, 6), (2, 1), \dots, (6, 6)\} \\ CF-LGC_2 &= LGC_2 \setminus \{(1, 1), (2, 2), \dots, (6, 6)\} \end{aligned}$$

with sizes  $N^{2|I_2|} = 36$  and  $N^{|I_2|}(N-1)^{|I_2|} = 30$ , respectively. Similarly,  $LGC_1$  and  $CF-LGC_1$  are as follows

$$\begin{aligned} LGC_1 &= \{(111, 111), \dots, (111, 666), (112, 111), \dots, (666, 666)\} \\ CF-LGC_1 &= LGC_1 \setminus \{(111, 111), \dots, (111, 661), (112, 111), \dots, (666, 666)\} \end{aligned}$$

with sizes  $N^{2|I_1|} = 46656$  and  $N^{|I_1|}(N-1)^{|I_1|} = 27000$ , respectively.

Let  $LGS$  denote the set of *label-generating sequences*, composed out of label-generating components as follows.

$$LGS := \{c_1.c_2 \dots c_d \mid \forall i \in [1, d] : c_i \in LGC_i\} \quad (\text{A.17})$$

We say that a sequence  $s \in LGS$  is *conflict-free* if it is composed entirely of conflict-free components, and we use  $CF-LGS$  to denote the conflict-free subset of  $LGS$ .

For each integer  $m \in [1, d]$ , let  $LGP_m$  denote the set of *label-generating prefixes of length  $m$* , defined as follows.

$$\begin{aligned} LGP_m &:= \{c_1.c_2 \dots c_{(m-1)}.\bar{y}_m \mid \forall i \in [1, m-1]. (c_i \in LGC_i) \\ &\quad \wedge \exists \bar{z}_m. ((\bar{y}_m, \bar{z}_m) \in LGC_m)\} \end{aligned} \quad (\text{A.18})$$

Let  $LGP$  denote the set of all label-generating prefixes. Every sequence  $s \in LGS$  corresponds to a unique prefix in each of  $LGP_1, \dots, LGP_d$ . Conversely, every prefix  $p \in LGP_m$  with  $m < d$  can be extended in  $N^{|I_m|+|I_{(m+1)}|}$  different ways to yield a prefix in  $LGP_{m+1}$ , while every prefix  $p \in LGP_d$  can be extended in  $N^{|I_d|}$  different ways to yield a complete sequence in  $LGS$ . We say that a prefix  $p \in LGP$  is *conflict-free* if it can be iteratively extended into a conflict-free sequence, and we use  $CF-LGP_m$  ( $CF-LGP$ ) to denote the conflict-free subset of  $LGP_m$  ( $LGP$ ). Every conflict-free prefix  $p \in CF-LGP_m$  with  $m < d$  can be extended in  $(N-1)^{|I_m|}N^{|I_{(m+1)}|}$  different ways to yield a conflict-free prefix in  $CF-LGP_{m+1}$ , while every conflict-free prefix  $p \in CF-LGP_d$  can be extended in  $(N-1)^{|I_d|}$  different ways to yield a conflict-free sequence in  $CF-LGS$ .

**Example 70** Continuing Example 69, consider the following label-generating sequences  $s_1, s_2, s_3 \in LGS$  and prefixes  $p_{12}, p_3, p_{123} \in LGP$ .

$$\begin{aligned} s_1 &= (354, 133).(6, 5) & p_{12} &= (354, 133).6 \\ s_2 &= (354, 133).(6, 6) & p_3 &= (354, 134).1 \\ s_3 &= (354, 134).(1, 2) & p_{123} &= 354 \end{aligned}$$

Sequence  $s_1$  is conflict-free, while  $s_2$  and  $s_3$  contain conflicts at positions 2.1 and 1.3, respectively. Prefix  $p_{12}$  is the length 2 prefix of both  $s_1$  and  $s_2$ ;  $p_3$  is the length 2 prefix of  $s_3$ ; and  $p_{123}$  is the length 1 prefix of  $s_1, s_2$ , and  $s_3$ . Finally,  $p_{12}$  and  $p_{123}$  are conflict-free, but  $p_3$  is not.

We now define a complex system of labels in terms of the label-generating components/sequences/prefixes defined above. The easiest way to describe the labels

in the labelling system is to show how they are applied; hence, for the moment, let  $l$  denote an abstract label. Corresponding to each such label  $l$ , let  $\mathcal{B}^l$  be a fresh set of constants isomorphic to  $\mathcal{B}$  (the variables in query  $Q$ ), and let  $\alpha^l : \mathcal{B} \rightarrow \mathcal{B}^l$  be a “labelling function” from each  $x \in \mathcal{B}$  to the corresponding constant in  $\mathcal{B}^l$  (which we will denote as  $x^l$ ). Because all such sets  $\mathcal{B}^l$  are mutually disjoint, we use  $\alpha^{-1}$  to denote the obvious “de-labelling function”  $\alpha^{-1} : (\bigcup_l \mathcal{B}^l) \rightarrow \mathcal{B}$  which we extend with identity over constants in  $\mathcal{C}_Q$  (the constants in query  $Q$ ).

For each label-generating sequence  $s \in LGS$ , we now define the following label-generating function  $\beta_s : \mathcal{B} \rightarrow (\bigcup_l \mathcal{B}^l)$  which performs a particular labelling of the variables of  $Q$ .

$$\beta_{(\bar{y}_1, \bar{z}_1) \dots (\bar{y}_d, \bar{z}_d)}(x) := \begin{cases} \alpha^{y_{1,j}}(x) & \text{if } \exists j \in [1, |\mathcal{I}_1|] \text{ such that } x = I_{1,j} \\ \alpha^{\bar{z}_1 \cdot y_{2,j}}(x) & \text{if } \exists j \in [1, |\mathcal{I}_2|] \text{ such that } x = I_{2,j} \\ \alpha^{\bar{z}_1 \cdot \bar{z}_2 \cdot y_{3,j}}(x) & \text{if } \exists j \in [1, |\mathcal{I}_3|] \text{ such that } x = I_{3,j} \\ \vdots & \vdots \\ \alpha^{\bar{z}_1 \dots \bar{z}_{(d-1)} \cdot y_{d,j}}(x) & \text{if } \exists j \in [1, |\mathcal{I}_d|] \text{ such that } x = I_{d,j} \\ \alpha^{\bar{z}_1 \dots \bar{z}_d}(x) & \text{otherwise} \end{cases} \quad (\text{A.19})$$

We immediately extend function  $\beta_s$  to tuples, sets, and atoms and with identity on constants in  $\mathcal{C}$ . We make the following observations about function  $\beta_s$ .

1. Because  $\beta_s$  is defined purely in terms of labelling functions, the de-labelling function  $\alpha^{-1}$  acts as an inverse of  $\beta_s$ .
2. For any level  $m \in [1, d]$ , the labels that  $\beta_s$  assigns to index variables in  $\mathcal{I}_m$  depend only on the length  $m$  prefix of  $s$ . Therefore, for each  $p \in LGS_m$  we can define the function  $\beta_p : \mathcal{I}_{[1,m]} \rightarrow (\bigcup_l \mathcal{B}^l)$  such that  $\beta_p(\mathcal{I}_{[1,m]}) = \beta_s(\mathcal{I}_{[1,m]})$  for every extension  $s$  of  $p$ .
3. Suppose  $s \in LGS$  contains a conflict at position  $i.j$ . Given any tuple  $t$  containing both variable  $I_{i,j}$  and some variable  $I \in \mathcal{B} \setminus \mathcal{I}_{[1,i]}$  (i.e. a non-index variable or a member of  $\mathcal{I}_{[i+1,d]}$ ), the labelled tuple  $\beta_s(t)$  *evidences the conflict at  $i.j$* . That is, because  $\beta_s(I_{i,j}) = I_{i,j}^{\bar{z}_1 \dots \bar{z}_{(i-1)} \cdot y_{i,j}}$  and  $\beta_s(I)$  is assigned a label that starts with  $\bar{z}_1 \dots \bar{z}_i$ , from tuple  $\beta_s(t)$  we can infer that  $y_{i,j} = z_{i,j}$  and so conclude that sequence  $s$  has a conflict at position  $i.j$ .

We now define the canonical database  $\mathbb{D}_Q$  as follows.

$$\mathbb{D}_Q := \bigcup_{s \in CF-LGS} \beta_s(\text{body}_Q) \quad (\text{A.20})$$

Because variable  $s$  ranges over only conflict-free label-generating sequences, database  $\mathbb{D}_Q$  does not contain any tuple that evidences any conflicts. Furthermore, observe that by de-labelling  $\mathbb{D}_Q$  we re-obtain  $\text{body}_Q$ .

$$\alpha^{-1}(\mathbb{D}_Q) = \text{body}_Q \quad (\text{A.21})$$



**Example 71** Continuing Example 70, functions  $\beta_{p_{123}}$ ,  $\beta_{p_{12}}$ , and  $\beta_{p_3}$  work as follows (recall that we assumed  $Q = Q_{81}$ ).

$$\begin{aligned}\beta_{p_{123}}(\mathcal{I}_1) &= \beta_{s_1}(\mathcal{I}_1) = \beta_{s_2}(\mathcal{I}_1) = \beta_{s_3}(\mathcal{I}_1) = \{N^3, M^5, Y^4\} \\ \beta_{p_{12}}(\mathcal{I}_{[1,2]}) &= \beta_{s_1}(\mathcal{I}_{[1,2]}) = \beta_{s_2}(\mathcal{I}_{[1,2]}) = \{N^3, M^5, Y^4, T^{133.6}\} \\ \beta_{p_3}(\mathcal{I}_{[1,2]}) &= \beta_{s_3}(\mathcal{I}_{[1,2]}) = \{N^3, M^5, Y^4, T^{134.1}\}\end{aligned}$$

Thus, applying functions  $\beta_{s_1}$ ,  $\beta_{s_2}$ , and  $\beta_{s_3}$  to  $\text{body}_{Q_{81}}$  yields the following sets of tuples.

$$\begin{aligned}\beta_{s_1}(\text{body}_{Q_{81}}) &= \{\text{Wrote}(N^3, T^{133.6}, Y^4), \text{Wrote}(M^5, T^{133.6}, Y^4)\} \\ \beta_{s_2}(\text{body}_{Q_{81}}) &= \{\text{Wrote}(N^3, T^{133.6}, Y^4), \text{Wrote}(M^5, T^{133.6}, Y^4)\} \\ \beta_{s_3}(\text{body}_{Q_{81}}) &= \{\text{Wrote}(N^3, T^{134.1}, Y^4), \text{Wrote}(M^5, T^{134.1}, Y^4)\}\end{aligned}$$

Database  $\mathbb{D}_{Q_{81}}$  is constructed out of  $|\text{CF-LGS}| = |\text{CF-LGC}_1| \times |\text{CF-LGC}_2| = 162000$  labelled copies of  $\text{body}_{Q_{81}}$ , one for each conflict-free label-generating sequence.<sup>2</sup> Note that in this particular example,  $\text{body}_{Q_{81}}$  does not contain any non-index variables which implies

$$\beta_{s_2}(\text{body}_{Q_{81}}) = \beta_{p_{12}}(\text{body}_{Q_{81}}) = \beta_{s_1}(\text{body}_{Q_{81}}) \subset \mathbb{D}_{Q_{81}}$$

even though  $s_2 \notin \text{CF-LGS}$ ; however, this is a special case. Indeed, any tuple containing both  $T^{134.1}$  and  $Y^4$  evidences a conflict at position 1.3, and so  $\beta_{s_3}(\text{body}_{Q_{81}}) \not\subset \mathbb{D}_{Q_{81}}$ .

### A.3.2 Canonical Objects

Having just defined the canonical database  $\mathbb{D}_Q$ , we now define the *canonical objects* encoded within  $\Gamma$ .

**Definition A.3.1 (Canonical Object)** *Given any integer  $m \in [1, d]$  and any label-generating prefix  $p \in \text{LGP}_m$ , we define the canonical object for prefix  $p$ —denoted  $o_p$ —as follows.*

$$o_p := \text{DECODE}(\Gamma[\beta_p(\bar{\mathcal{I}}_{[1,m]})], \bar{\mathfrak{S}}_{[m+1,d]})$$

Observe that if  $m = d$ , the canonical object for prefix  $p \in \text{LGP}_d$  is simply the tuple  $\beta_p(\bar{\mathcal{V}})$ . The following two lemmas reveal how the design of database  $\mathbb{D}_Q$  endows canonical objects with specific properties.

**Lemma A.3.2** *Given any prefix  $p \in \text{LGP}_d$ ,*

1. *if  $p$  is conflict-free then  $\beta_p(\bar{\mathcal{I}}_{[1,d]}) \in \text{adom}(\bar{\mathcal{I}}_{[1,d]}, R)$ ; and*
2. *if  $\beta_p(\bar{\mathcal{I}}_{[1,d]}) \in \text{adom}(\bar{\mathcal{I}}_{[1,d]}, R)$  then for any position  $i.j$  for which  $p$  contains a conflict, that conflict is evidenced within the tuple  $\beta_p(\bar{\mathcal{I}}_{[i+1,d]}\bar{\mathcal{V}})$ .*

---

<sup>2</sup>That database  $\mathbb{D}_Q$  has exponential size compared to  $Q$  is not important because we never materialize  $\mathbb{D}_Q$  nor evaluate any query against it. For the purposes of the proof, it suffices that  $\mathbb{D}_Q$  is finite.



PROOF. If  $p$  is conflict-free, then there exists some  $s \in CF-LGS$  that extends  $p$ , and so  $\beta_p(\bar{\mathcal{I}}_{[1,d]}) = \beta_s(\bar{\mathcal{I}}_{[1,d]}) \in \text{adom}(\bar{\mathcal{I}}_{[1,d]}, R)$  follows by equation A.20. Conversely, suppose both that  $\beta_p(\bar{\mathcal{I}}_{[1,d]}) \in \text{adom}(\bar{\mathcal{I}}_{[1,d]}, R)$  and that  $p$  contains a conflict at position  $i.j$ ; this entails  $i \in [1, d-1]$ ,  $j \in [1, |\mathcal{I}_i|]$ , and  $|\mathcal{I}_{[i+1,d]}| > 0$ . If attribute  $I_{i,j} \in \mathcal{V}$ , then tuple  $\beta_p(\bar{\mathcal{I}}_{[i+1,d]}\bar{\mathcal{V}})$  is guaranteed to evidence the conflict (see Section A.3.1). If  $I_{i,j} \notin \mathcal{V}$ , let  $s \in LGS$  be any label-generating sequence that extends  $p$ . By the design of  $\mathbb{D}_Q$ , embedding  $\gamma := \beta_s$  proves  $Q \models (\mathcal{I}_{[1,i]} \setminus \{I_{i,j}\}) \{ \rightarrow I_{i,j} \}$  which, by Definition 3.1.11, implies that  $Q$  is not in  $\mathbf{s}^d\text{-NF}$  (a contradiction).  $\square$

**Lemma A.3.3** *Given any integer  $m \in [1, d]$  and any prefix  $p \in CF-LGP_m$ , canonical object  $o_p$  contains enough information to reconstruct the value  $\beta_p(\bar{\mathcal{I}}_m)$ .*

PROOF. By induction on  $m$ . Our assumption that  $Q$  is in  $\mathbf{s}^d\text{-NF}$  implies that  $\mathcal{I}_d \subseteq \mathcal{V}$  (see Definition 3.1.11); therefore, the base case  $m = d$  is trivial because object  $o_p$  is a single tuple whose output values contain all of the values in  $\beta_p(\bar{\mathcal{I}}_m)$ . Now suppose that the lemma holds for all prefixes of length in the range  $[m+1, d]$ . Recall from equations A.18 and A.19 that  $p$  has the form  $p = (\bar{y}_1, \bar{z}_1) \dots (\bar{y}_{(m-1)}, \bar{z}_{(m-1)}) \cdot \bar{y}_m$ , but computing  $\beta_p(\bar{\mathcal{I}}_m)$  only requires deducing the sequence  $\bar{z}_1 \dots \bar{z}_{(m-1)} \cdot \bar{y}_m$ . Consider each prefix  $q \in LGP_{(m+1)}$  that extends  $p$ . If  $q$  is conflict-free, then by the induction hypothesis  $o_q$  contains enough information to identify the value  $\beta_q(\bar{\mathcal{I}}_{(m+1)})$ , from which we can deduce the values of sequence  $\bar{z}_1 \dots \bar{z}_{(m-1)} \cdot \bar{z}_m^q \cdot \bar{y}_{(m+1)}^q$ , where the values  $\bar{z}_m^q$  and  $\bar{y}_{(m+1)}^q$  vary for different choices of  $q$ . If  $q$  contains a conflict, then Lemma A.3.2 guarantees that the existence of the conflict can be detected by examining  $o_q$ . By examining all of the sub-objects of  $o_p$  we can determine which values of  $\bar{z}_m^q$  corresponds to a conflict-free extension of  $p$ , from which we can deduce the value of  $\bar{y}_m$ .  $\square$

### A.3.3 Pseudo-Trivial Certificate Nodes

We now show that sub-relations of  $\Gamma$  that encode canonical objects can only be equated to sub-relations of  $\Phi$  by pseudo-trivial certificate nodes.

**Lemma A.3.4** *Given*

- any integer  $m \in [1, d]$ ,
- any prefix  $p \in CF-LGP_m$ , and
- any sub-relation  $\Phi[\bar{a}'_{[1,m]}]$  that encodes canonical object  $o_p$ ;

*index tuple  $\bar{a}'_{[1,m]}$  must contain all of the values in  $\beta_p(\mathcal{I}_m)$ .*

PROOF. For each  $j \in [1, |\bar{\mathcal{I}}_m|]$ , consider the index attribute  $I_{m,j} \in \mathcal{I}_m$  and the corresponding value  $\beta_p(I_{m,j}) = I_{m,j}^{\bar{z}_1 \dots \bar{z}_{(m-1)} \cdot \bar{y}_{m,j}}$ . By Lemma A.3.3, index tuple  $\bar{a}'_{[1,m]}$  must contain enough information to compute the value  $\beta_p(I_{m,j})$ . In order to do this, either  $\bar{a}'_{[1,m]}$  must contain  $\beta_p(I_{m,j})$ , or  $\bar{a}'_{[1,m]}$  must contain enough values to isolate  $\beta_p(I_{m,j})$  via data triangulation. By the symmetry in the construction

of  $\mathbb{D}_Q$ , every database constant that co-occurs with  $\beta_p(I_{m,j})$  in a tuple also co-occurs symmetrically with at least  $N - 2$  other constants of the form  $I_{m,j}^{\bar{z}_1 \dots \bar{z}_{(m-1)} \cdot n}$ . Therefore, isolating  $\beta_p(I_{m,j})$  via data triangulation requires at least  $N - 1 > |\mathcal{I}'_{[1,d]}|$  triangulation points, which exceeds the size of  $\bar{a}'_{[1,m]}$ . Hence,  $\bar{a}'_{[1,m]}$  must contain the value  $\beta_p(I_{m,j})$ .  $\square$

**Lemma A.3.5** *Given*

- any conflict-free label-generating sequence  $s \in CF-LGS$ ,
- any  $\bar{\xi}$ -certificate that proves  $\Gamma \doteq_{\bar{\xi}} \Phi$ , and
- any integer  $m \in [1, d]$ ;

every certificate node involving sub-relation  $\Gamma[\beta_s(\bar{\mathcal{I}}_{[1,m-1]})]$  is pseudo-trivial w.r.t.  $\beta_s(\bar{\mathcal{I}}_m)$ .

PROOF. Suppose that  $Q \doteq_{s^d} Q'$ . We have not constrained which query was chosen as  $Q$  and which was chosen as  $Q'$ , and so assume without loss of generality that if  $Q$  and  $Q'$  do not have the same number of index variables at every level, then  $Q$  was chosen so that  $|\bar{\mathcal{I}}_m| > |\bar{\mathcal{I}}'_m|$  at the first level  $m$  upon which they disagree, starting at level 1 and counting up to  $d$ .

Choose any conflict-free label-generating sequence  $s \in CF-LGS$ ; within any  $s^d$ -certificate proving  $R \doteq_{s^d} R'$ , the certificate must contain a tuple node proving  $\Gamma[\beta_s(\bar{\mathcal{I}}_{[1,d]})] \doteq_{\emptyset} \Phi[\bar{a}'_{[1,d]}]$  for some  $\bar{a}'_{[1,d]} \in \text{adom}(\bar{\mathcal{I}}'_{[1,d]}, \Phi)$ . Lemma A.3.4 implies that  $\bar{a}'_1$  contains all values in  $\beta_s(\bar{\mathcal{I}}_1)$ , and so by choice of  $Q$  above conclude that  $|\bar{\mathcal{I}}_1| = |\bar{\mathcal{I}}'_1|$  and therefore  $\bar{a}'_1$  is a permutation of  $\beta_s(\bar{\mathcal{I}}_1)$ . A simple induction proves that  $\bar{a}'_m$  is a permutation of  $\beta_s(\bar{\mathcal{I}}_m)$  for every  $m \in [1, d]$ .  $\square$

**Lemma A.3.6** *If  $\bar{\xi} = s^d$  then there exists index-covering homomorphisms in both directions between  $Q$  and  $Q'$ .*

PROOF. Choose any conflict-free label-generating sequence  $s \in CF-LGS$ . Consider any  $s^d$ -certificate proving  $\Gamma \doteq_{s^d} \Phi$ . The certificate must contain a tuple node proving  $\Gamma[\beta_s(\bar{\mathcal{I}}_{[1,d]})] \doteq_{\emptyset} \Phi[\bar{a}'_{[1,d]}]$  for some  $\bar{a}'_{[1,d]} \in \text{adom}(\bar{\mathcal{I}}'_{[1,d]}, \Phi)$ . Choose any  $\phi \in \Phi[\bar{a}'_{[1,d]}]$ .

Define  $h := \alpha^{-1} \circ \phi$ . Then,  $h$  is a mapping from  $\mathcal{B}'$  to  $\alpha^{-1}(\text{adom}(\mathbb{D}_Q)) = \mathcal{B} \cup \mathcal{C}_Q$  that satisfies the following properties.

1.  $h(\text{body}_{Q'}) = \alpha^{-1} \circ \phi(\text{body}_{Q'}) \subseteq \alpha^{-1}(\mathbb{D}_Q) = \text{body}_Q$
2.  $h(\bar{\mathcal{V}}') = \alpha^{-1} \circ \phi(\bar{\mathcal{V}}') = \alpha^{-1} \circ \beta_s(\bar{\mathcal{V}}) = \bar{\mathcal{V}}$
3.  $\forall i \in [1, d], h(\mathcal{I}'_i) = \alpha^{-1} \circ \phi(\mathcal{I}'_i) = \alpha^{-1} \circ \beta_s(\mathcal{I}_i) = \mathcal{I}_i$

By Lemma A.3.5 the node proving  $\Gamma[\beta_s(\bar{\mathcal{I}}_{[1,i-1]})] \doteq_{\bar{\xi}_{[i,d]}} \Phi[\bar{a}'_{[1,i-1]}]$  is pseudo-trivial w.r.t.  $\beta_s(\bar{\mathcal{I}}_i)$ , and so  $\phi(\bar{\mathcal{I}}'_i) = \bar{a}'_i$  is a permutation of  $\beta_s(\bar{\mathcal{I}}_i)$ .

Hence,  $h$  satisfies Definition 3.2.1. The proof is analogous for  $h'$  from  $Q$  to  $Q'$ .  $\square$

## A.4 Normalized Bag Equivalence: Symmetry + Counting

Whereas bag equivalence requires agreement of absolute frequencies between corresponding sub-objects, normalized bag equivalence only requires agreement of relative frequencies. We introduced *colour-diversification* in Section A.2 as a technique for controlling the frequency of sub-objects, and in Lemma A.2.9 we showed that by applying colour-diversification to *any* database we can guarantee the existence of pseudo-trivial nodes within the  $\bar{\S}$ -certificate (assuming  $Q \doteq_{\bar{\S}} Q'$ ). While colour-diversification is still a key technique in our proof for normalized bag equivalence, the following example shows that we require more care when selecting the database to diversify. For the duration of this section we assume that the encoding signature contains only normalized bag semantics ( $\bar{\S} = \mathbf{n}^d$ ).

**Example 72** Consider the following two CEQs, which assume the same relation schema `Wrote(Name, Title, Year)` as used in Example 68.

$$\begin{aligned} Q_{86}(T, Y; Y) &:- \text{Wrote}(N, T, Y) \\ Q_{87}(N, T, Y; Y) &:- \text{Wrote}(N, T, Y) \end{aligned}$$

Intuitively, both queries characterize the relative annual publication rate. They differ in that  $Q_{86}$  counts titles (per year), whereas  $Q_{87}$  counts (author name,title) pairs (per year).

Consider our approach for testing bag equivalence (Section A.2), using the canonical database defined by equation A.13. Start with database  $\mathbb{D}_{13}$  formed from  $\text{body}_{Q_{86}}$ ,

$$\mathbb{D}_{13} = \{\text{Wrote}(N, T, Y)\}$$

assume  $\text{adom}(\mathbb{D}_{13})$  has order  $\{N, T, Y\}$ , and choose any 3-distinguishing diversification point  $\bar{r}$ . For conciseness of presentation we will use  $\bar{r} = [2, 3, 2]$  even though it is not 3-distinguishing. Encoding relations  $R_{28}$  and  $R_{29}$  in Figure A.6 shows the result of evaluating  $Q_{86}$  and  $Q_{87}$ , respectively, over database  $\mathbb{D}_{14} = \Delta^{[2,3,2]}(\mathbb{D}_{13})$ . Whereas  $\mathbb{D}_{14}$  discriminates between  $R_{28}$  and  $R_{29}$  under  $\mathbf{b}$ -equality, it does not discriminate between them under  $\mathbf{n}$ -equality because the effect of multiplicative factors  $r_1 = 2$  and  $r_2 = 3$  (corresponding to number of colours for  $N$  and  $T$ , respectively) is voided by the normalization.

We can create a database that discriminates between  $Q_{86}$  and  $Q_{87}$  under  $\mathbf{n}$ -equality by a slight modification of the above approach. Start with database  $\mathbb{D}_{15}$  formed from two distinct copies of  $\text{body}_{Q_{86}}$ ,

$$\mathbb{D}_{15} = \{\text{Wrote}(N, T, Y), \text{Wrote}(N', T', Y')\}$$

assume  $\text{adom}(\mathbb{D}_{15})$  has order  $\{N, T, Y, N', T', Y'\}$ , and choose  $\bar{r} = [2, 3, 2, 1, 1, 1]$ . Figure A.7 shows the result of evaluating  $Q_{86}$  and  $Q_{87}$  over database  $\Delta^{[2,3,2,1,1,1]}(\mathbb{D}_{15})$ , along with the differing normalized bags obtained after decoding.



The key lesson from the above example is that in order for colour diversification to distinguish the results of two encoding queries under n-equality, the encoded normalized bag must contain two sub-objects whose bag cardinalities are described by polynomials (equations A.11 and A.12 in Section A.2) without any common multiplicative factors. In this case, the choice of a  $(|\mathcal{I}_{[1,d]}| + |\mathcal{I}'_{[1,d]}|)$ -distinguishing diversification point guarantees that the encoded bag cannot be normalized further. Extending this to nested normalized bags requires guaranteeing that there exists a recursive descent of the encoding relation in which at each level the encoded normalized bag contains two sufficiently-different encoded sub-objects. To ensure this, we next devise a canonical database containing canonical objects defined in terms of a system of label-generating prefixes, similar to the one we devised for nested sets.

#### A.4.1 Labels, Prefixes, and Canonical Database $\mathbb{D}_Q$

The requirements of our labelling system in this section are significantly simpler than the requirements in Section A.3; however, it is expedient to adopt the same terminology of label-generating components/sequences/prefixes and simply modify the definitions.

For each level  $i \in [1, d]$ , let  $LGC_i$  denote the set of *label-generating components* at level  $i$ , defined as follows.

$$LGC_i := \{1, 2\} \tag{A.22}$$

Let  $LGS$  denote the set of *label-generating sequences*, composed out of label-generating components as follows.

$$LGS := \{c_1.c_2 \dots c_d \mid \forall i \in [1, d] : c_i \in LGC_i\} \tag{A.23}$$

For each integer  $m \in [1, d]$ , let  $LGP_m$  denote the set of *label-generating prefixes of length  $m$* , defined as follows.

$$LGP_m := \{c_1.c_2 \dots c_m \mid \forall i \in [1, m]. (c_i \in LGC_i)\} \tag{A.24}$$

Let  $LGP$  denote the set of all label-generating prefixes. We do not require the notion of a *conflict*, and so for uniformity of terminology we simply define  $CF-LGC_i := LGC_i$  which then entails  $CF-LGP_m = LGP_m$ ,  $CF-LGP = LGP$ , and  $CF-LGS = LGS$ .

As in Section A.3, we define a system of labels in terms of the label-generating components/sequences/prefixes defined above. For each label  $l$ , we define the set  $\mathcal{B}^l$  and the labelling function  $\alpha^l : \mathcal{B} \rightarrow \mathcal{B}^l$  and the de-labelling function  $\alpha^{-1} : (\bigcup_l \mathcal{B}^l) \rightarrow \mathcal{B}$  exactly as in Section A.3. For each label-generating sequence  $s \in LGS$ , we define the following label-generating function  $\beta_s : \mathcal{B} \rightarrow (\bigcup_l \mathcal{B}^l)$  which performs a particular

labelling of the variables of  $Q$ .

$$\beta_{c_1.c_2\dots c_d}(x) := \begin{cases} \alpha^{c_1}(x) & \text{if } x \in \mathcal{I}_1 \\ \alpha^{c_1.c_2}(x) & \text{if } x \in \mathcal{I}_2 \\ \vdots & \vdots \\ \alpha^{c_1\dots c_{(d-1)}}(x) & \text{if } x \in I_{(d-1)} \\ \alpha^{c_1\dots c_d}(x) & \text{if } x \in I_d \\ \alpha^{c_1\dots c_d}(x) & \text{otherwise} \end{cases} \quad (\text{A.25})$$

We immediately extend function  $\beta_s$  to tuples, sets, and atoms and with identity on constants in  $\mathcal{C}$ . Equation A.25 is much simpler than equation A.19, and in this case it suffices that  $\mathbb{L} := LGP$  (which was not true for nested sets). We observe that  $\alpha^{-1}$  again acts as an inverse of  $\beta_s$ , and for any level  $m \in [1, d]$  and each  $p \in LGP_m$  we again define the function  $\beta_p : \mathcal{I}_{[1,m]} \rightarrow (\bigcup_l \mathcal{B}^l)$  as in Section A.3.

We now define the canonical database  $\mathbb{D}_Q$  in two stages. First, we form the database  $\mathbb{D}_Q^{\text{pre}}$  in terms of the label-generating sequences just defined, analogous to the canonical database used for nested sets (equation A.20).

$$\mathbb{D}_Q^{\text{pre}} := \bigcup_{s \in CF\text{-}LGS} \beta_s(\text{body}_Q) \quad (\text{A.26})$$

Next, we let  $\bar{r}$  be any  $(|\mathcal{I}_{[1,d]}| + |\mathcal{I}'_{[1,d]}|)$ -distinguishing diversification point for  $\text{adom}(\mathbb{D}_Q^{\text{pre}})$  and we use colour-diversification to define canonical database  $\mathbb{D}_Q$ , as we did with nested bags (equation A.13).

$$\mathbb{D}_Q := \Delta^{\bar{r}}(\mathbb{D}_Q^{\text{pre}}) \quad (\text{A.27})$$

Observe that by first whitewashing and then de-labelling  $\mathbb{D}_Q$  we re-obtain  $\text{body}_Q$ .

$$\alpha^{-1} \circ \delta^{-1}(\mathbb{D}_Q) = \text{body}_Q \quad (\text{A.28})$$

**Lemma A.4.1** *Given*

- any integer  $i \in [1, d]$ ,
- any embedding  $\gamma \in \Gamma$ , and
- any label-generating sequence  $s \in LGS$ ;

if  $\gamma(\mathcal{I}_{[1,d]} \setminus (\mathcal{I}_i \setminus \mathcal{V})) = \beta_s(\mathcal{I}_{[1,d]} \setminus (\mathcal{I}_i \setminus \mathcal{V}))$  and  $\delta^{-1} \circ \gamma(\mathcal{I}_i \setminus \mathcal{V}) \neq \beta_s(\mathcal{I}_i \setminus \mathcal{V})$  then  $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\mathcal{I}_i \setminus \mathcal{V}) \neq (\mathcal{I}_i \setminus \mathcal{V})$ .

PROOF. Assume w.l.o.g. that  $s = c_1 \dots c_d$ . Assume to the contrary that

- $\gamma(\mathcal{I}_{[1,d]} \setminus (\mathcal{I}_i \setminus \mathcal{V})) = \beta_s(\mathcal{I}_{[1,d]} \setminus (\mathcal{I}_i \setminus \mathcal{V}))$ ,
- $\delta^{-1} \circ \gamma(\mathcal{I}_i \setminus \mathcal{V}) \neq \beta_s(\mathcal{I}_i \setminus \mathcal{V})$ , and
- $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\mathcal{I}_i \setminus \mathcal{V}) = (\mathcal{I}_i \setminus \mathcal{V})$ .

Let  $j$  be any integer such that  $I_{i,j} \notin \mathcal{V}$ . Variable  $I_{i,j}$  is a core index (assumption [A2]), and so by Definition 3.1.11 (with the assumption  $\bar{\mathfrak{s}} = \mathfrak{n}^d$ ) and Theorem 3.1.10 it follows that

1.  $|(\mathcal{I}_i \cap \mathcal{V}) \cup \mathcal{I}_{[i+1,d]}| > 0$ , and
2. set  $\mathcal{I}_{[1,i-1]}$  is *not* a strong  $(I_{i,j}, (\mathcal{I}_i \cap \mathcal{V}) \cup \mathcal{I}_{[i+1,d]})$ -articulation set of hypergraph  $H^Q$ . In other words, if we delete vertices  $\mathcal{I}_{[1,i-1]}$  from hypergraph  $H^Q$ , the remaining hypergraph must contain an  $(I_{i,j}, I)$ -path for some variable  $I \in (\mathcal{I}_i \cap \mathcal{V}) \cup \mathcal{I}_{[i+1,d]}$ .

Now consider the function  $\alpha^{-1} \circ \delta^{-1} \circ \gamma$ . Because whitewashing function  $\delta^{-1}$  acts as the identity over  $\text{colour}_1$  and because  $\alpha^{-1}$  is the inverse of  $\beta_s$ , it follows that  $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\mathcal{I}_{[1,d]} \setminus (\mathcal{I}_i \setminus \mathcal{V})) = (\mathcal{I}_{[1,d]} \setminus (\mathcal{I}_i \setminus \mathcal{V}))$ . Combined with the above assumption and with equation A.28, this implies that  $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\overline{\mathcal{I}}_{[1,d]}) = \overline{\mathcal{I}}_{[1,d]}$  and  $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\text{body}_Q) \subseteq \text{body}_Q$ . Because  $\text{body}_Q$  is minimal (assumption [A3]), it follows from Theorem 3.1.2 that  $\alpha^{-1} \circ \delta^{-1} \circ \gamma$  is a bijection from  $\mathcal{B}$  onto  $\mathcal{B}$ ; hence, hypergraphs  $H^Q$  and  $\gamma(H^Q)$  must be isomorphic.

Now consider deleting vertices  $\gamma(\mathcal{I}_{[1,i-1]})$  from hypergraph  $\gamma(H^Q)$ . By isomorphism of  $H^Q$  and  $\gamma(H^Q)$  the remaining hypergraph contains a path between vertices  $\gamma(I_{i,j})$  and  $\gamma(I)$ . By the bijectivity of  $\alpha^{-1} \circ \delta^{-1} \circ \gamma$ , all of the remaining vertices are labelled copies of either non-index variables or variables in  $\mathcal{I}_{[i,d]}$ —and so by equations A.27 and A.25 all of the other vertices in the same connected component as  $\gamma(I)$  (including  $\gamma(I_{i,j})$ ) must have labels that begin with  $c_1 \dots c_i$ . By  $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\mathcal{I}_i \setminus \mathcal{V}) = (\mathcal{I}_i \setminus \mathcal{V})$ , conclude that  $\delta^{-1} \circ \gamma(I_{i,j}) \in \beta_s(\mathcal{I}_i \setminus \mathcal{V})$ .

Repeating this argument for each value of  $j \in [1, |\mathcal{I}_i|]$  such that  $I_{i,j} \notin \mathcal{V}$ , and combining it with  $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\mathcal{I}_i \setminus \mathcal{V}) = (\mathcal{I}_i \setminus \mathcal{V})$ , conclude the contradiction  $\delta^{-1} \circ \gamma(\mathcal{I}_i \setminus \mathcal{V}) = \beta_s(\mathcal{I}_i \setminus \mathcal{V})$ .  $\square$

## A.4.2 Canonical Objects

Having just defined the canonical database  $D_Q$ , we now define the *canonical objects* encoded within  $\Gamma$  exactly the same as we defined them for nested sets (Definition A.3.1). That is, the canonical object  $o_p$  for prefix  $p \in LGP_m$  is defined as follows.

$$o_p := \text{DECODE}(\Gamma[\beta_p(\overline{\mathcal{I}}_{[1,m]})], \overline{\mathfrak{s}}_{[m+1,d]})$$

Observe that given any integer  $m \in [1, d-1]$  and any prefix  $p \in LGP_m$ , canonical object  $o_p$  always contains canonical sub-objects  $o_{p.1}$  and  $o_{p.2}$ .

## A.4.3 Pseudo-Trivial Certificate Nodes

We now show that sub-relations of  $\Gamma$  that encode canonical objects can only be equated to sub-relations of  $\Phi$  by pseudo-trivial certificate nodes.

**Lemma A.4.2** *Given*

- any label-generating sequence  $s \in LGS$ ,
- any  $\overline{\mathfrak{s}}$ -certificate that proves  $\Gamma \stackrel{\overline{\mathfrak{s}}}{\doteq} \Phi$ , and
- any level  $m \in [1, d]$ ;

every certificate node involving sub-relation  $\Gamma[\beta_s(\bar{\mathcal{I}}_{[1,m-1]})]$  is potentially pseudo-trivial w.r.t.  $\beta_s(\bar{\mathcal{I}}_m)$ .

PROOF. The proof is by induction on  $m$ . First, however, note that we have never constrained which query was chosen as  $Q$  and which was chosen as  $Q'$  (i.e, which query should be used to create  $\mathbb{D}_Q$ ). W.l.o.g. assume that if  $Q$  and  $Q'$  do not have the same number of index variables at every level, then  $Q$  was chosen so that  $|\mathcal{I}_l| > |\mathcal{I}'_l|$  at the first level  $l$  upon which  $Q$  and  $Q'$  disagree, starting at level  $d$  and counting down. Also, assume  $s = c_1 \dots c_d$  and without loss of generality assume that  $c_m = 1$  (the proof is symmetric if  $c_m = 2$ ).

**Base Case:** Suppose  $m = d$ . Define prefix  $p := c_1 \dots c_{(d-1)}$ . If  $|\mathcal{I}_d| = 0$  then every certificate node at level  $d$  is trivially pseudo-trivial; therefore, assume that  $|\mathcal{I}_d| > 0$ . The proof is lengthy, and so we will break it into numbered steps for readability.

1. Let  $n^a$  be any certificate node proving  $\Gamma[\bar{a}] \doteq_{\S_d} \Phi[\bar{a}']$  for  $\bar{a} := \beta_s(\bar{\mathcal{I}}_{[1,d-1]})$  and some  $\bar{a}' \in \text{adom}(\bar{\mathcal{I}}'_{[1,d-1]}, \Phi)$ . Relation  $\Gamma[\bar{a}]$  encodes canonical object  $o_p$ , which contains canonical sub-objects  $o_{p.1}$  and  $o_{p.2}$ . Let sets  $S_{o_{p.1}}$  and  $S_{o_{p.2}}$  be defined as in equation A.10, containing the tuples of non-output values that index sub-objects  $o_{p.1}$  and  $o_{p.2}$ , respectively, within  $\Gamma[\bar{a}]$ . Define  $S'_{o_{p.1}}$  and  $S'_{o_{p.2}}$  analogously for  $\Phi[\bar{a}']$ . Define polynomials  $f_{S_{o_{p.1}}}$ ,  $f_{S_{o_{p.2}}}$ ,  $f_{S'_{o_{p.1}}}$ , and  $f_{S'_{o_{p.2}}}$  as in equations A.11 and A.12.
2. By definition of  $\Gamma[\bar{a}]$ , every assignment  $\gamma \in \Gamma[\bar{a}]$  satisfies

$$\gamma(\mathcal{I}_{[1,d-1]}) = \beta_{p.1}(\mathcal{I}_{[1,d-1]}) = \beta_s(\mathcal{I}_{[1,d-1]}) \quad (\text{A.29})$$

while by definition of canonical object  $o_{p.1}$  every  $\gamma \in \Gamma[\bar{a}]|_{o_{p.1}}$  satisfies the following.

$$\gamma(\mathcal{I}_d \cap \mathcal{V}) = \beta_{p.1}(\mathcal{I}_d \cap \mathcal{V}) = \beta_s(\mathcal{I}_d \cap \mathcal{V}) \quad (\text{A.30})$$

Tuples in  $S_{o_{p.1}}$  are formed by projecting  $\bar{\mathcal{J}}_d$  from the assignments in  $\Gamma[\bar{a}]|_{o_{p.1}}$  (see equation A.10, and recall  $\mathcal{J}_d = \mathcal{I}_d \setminus \mathcal{V}$ ). Therefore, there cannot exist a tuple  $t \in S_{o_{p.1}}$  such that

- (a)  $\alpha^{-1} \circ \delta^{-1}(t)$  contains all of the variables in  $\mathcal{J}_d$ , and
- (b) the label  $p.2$  occurs within  $t$ .

If  $t$  did exist, then it would correspond to an assignment  $\gamma$  satisfying

$$\alpha^{-1} \circ \delta^{-1} \circ \gamma(\mathcal{J}_d) = \mathcal{J}_d \quad (\text{A.31})$$

$$\delta^{-1} \circ \gamma(\mathcal{J}_d) \neq \beta_{p.1}(\mathcal{J}_d) = \beta_s(\mathcal{J}_d) \quad (\text{A.32})$$

which contradicts Lemma A.4.1. In contrast, the design of  $\mathbb{D}_Q$  implies that  $\beta_{p.1}(\bar{\mathcal{I}}_d) \in \text{adom}(\bar{\mathcal{I}}_d, \Gamma[\bar{a}])$ , and the definition of canonical object  $o_{p.1}$  further implies  $\beta_{p.1}(\bar{\mathcal{I}}_d) \in \text{adom}(\bar{\mathcal{I}}_d, \Gamma[\bar{a}]|_{o_{p.1}})$ ; therefore,  $S_{o_{p.1}}$  does contain the tuple  $\beta_s(\bar{\mathcal{J}}_d) = \beta_{p.1}(\bar{\mathcal{J}}_d)$ , for some ordering  $\bar{\mathcal{J}}_d$  of  $\mathcal{J}_d$ .



3. Each monomial in  $f_{S_{o_{p.1}}}$  (resp.,  $f_{S_{o_{p.2}}}$ ) corresponds to a collection of values occurring within a tuple in  $S_{o_{p.1}}$  (resp.,  $S_{o_{p.2}}$ ). Define polynomial  $f_{\text{GCD}} := \text{GCD}(f_{S_{o_{p.1}}}, f_{S_{o_{p.2}}})$ . By definition of greatest common divisor, there exists two polynomials  $f_1, f_2$  with  $\text{GCD}(f_1, f_2) = 1$  that satisfy  $f_{S_{o_{p.1}}} = f_1 \times f_{\text{GCD}}$  and  $f_{S_{o_{p.2}}} = f_2 \times f_{\text{GCD}}$ . We will now prove by contradiction that  $\text{degree}(f_{\text{GCD}}) = 0$ . Suppose that  $\text{degree}(f_{\text{GCD}}) = k > 0$ . Because  $S_{o_{p.1}}$  contains tuple  $\beta_{p.1}(\overline{\mathcal{I}}_d)$ , polynomial  $f_{S_{o_{p.1}}}$  contains a monomial corresponding to the values  $\beta_{p.1}(\mathcal{I}_d)$ . Therefore, there must exist some subset  $J'_d \subseteq J_d$  with  $|J'_d| = k$  such that polynomial  $f_{\text{GCD}}$  contains a monomial corresponding to values  $\beta_{p.1}(J'_d)$ . By symmetry of  $\mathbb{D}_Q$ , polynomial  $f_{\text{GCD}}$  must be symmetric between labels  $p.1$  and  $p.2$  which means that  $f_{\text{GCD}}$  also contains a monomial corresponding to values  $\beta_{p.2}(J'_d)$ . However, this implies that  $f_{S_{o_{p.1}}}$  contains a monomial corresponding to values  $\beta_{p.1}(\mathcal{I}_d \setminus J'_d) \cup \beta_{p.2}(J'_d)$ , which further implies the existence of tuple  $t \in S_{o_{p.1}}$  with  $\alpha^{-1} \circ \delta^{-1} \circ \gamma(\mathcal{I}_d) = \mathcal{I}_d$  and  $\delta^{-1} \circ \gamma(\mathcal{I}_d) \neq \beta_{p.1}(\mathcal{I}_d) = \beta_s(\mathcal{I}_d)$ —which contradicts the previous paragraph. Hence,  $\text{GCD}(f_{S_{o_{p.1}}}, f_{S_{o_{p.2}}})$  is a constant.
4. Because  $\Gamma[\overline{a}]$  and  $\Phi[\overline{a}']$  encode the same normalized bag, we know

$$\frac{f_{S_{o_{p.1}}}(\overline{r})}{f_{S_{o_{p.2}}}(\overline{r})} = \frac{f_{S'_{o_{p.1}}}(\overline{r})}{f_{S'_{o_{p.2}}}(\overline{r})}$$

and so

$$f_{S_{o_{p.1}}}(\overline{r}) \times f_{S'_{o_{p.2}}}(\overline{r}) = f_{S_{o_{p.2}}}(\overline{r}) \times f_{S'_{o_{p.1}}}(\overline{r}) \quad (\text{A.33})$$

follows. During construction of database  $\mathbb{D}_Q$  point  $\overline{r}$  was chosen to be  $(|\mathcal{I}_{[1,d]}| + |\mathcal{I}'_{[1,d]}|)$ -distinguishing, while each side of equation A.33 is a polynomial of degree  $(|\mathcal{I}_d| + |\mathcal{I}'_d|)$ ; therefore the two sides must be the same polynomial.

$$f_{S_{o_{p.1}}} \times f_{S'_{o_{p.2}}} = f_{S_{o_{p.2}}} \times f_{S'_{o_{p.1}}} \quad (\text{A.34})$$

5. From Lemma A.1.1 we know that  $|\mathcal{I}_d \cap \mathcal{V}| = |\mathcal{I}'_d \cap \mathcal{V}'|$ . From equation A.30 it follows that every  $\phi \in \Phi[\overline{a}']|_{o_{p.1}}$  must satisfy  $\phi(\mathcal{I}'_d \cap \mathcal{V}') = \beta_{p.1}(\mathcal{I}_d \cap \mathcal{V})$ .
6. We know  $|\mathcal{I}_d| \geq |\mathcal{I}'_d|$  by choice of  $Q$ , and so by  $|\mathcal{I}_d \cap \mathcal{V}| = |\mathcal{I}'_d \cap \mathcal{V}'|$  conclude the following.

$$\text{degree}(f_{S_{o_{p.1}}}) = |\mathcal{I}_d \setminus \mathcal{V}| \geq |\mathcal{I}'_d \setminus \mathcal{V}'| = \text{degree}(f_{S'_{o_{p.1}}}) \quad (\text{A.35})$$

Because  $\text{GCD}(f_{S_{o_{p.1}}}, f_{S_{o_{p.2}}})$  is a constant, from equation A.34 conclude that  $f_{S_{o_{p.1}}}$  must be a constant multiple of  $f_{S'_{o_{p.1}}}$ . Therefore, for each tuple  $t \in S_{o_{p.1}}$  there exists a permutation of  $t$  within  $S'_{o_{p.1}}$ , and vice versa.

7. By the definitions of  $\overline{a}$  and canonical object  $o_{p.1}$  By  $\beta_s(\overline{\mathcal{I}}_{[1,d-1]}) = \overline{a}$  and  $\beta_s(\overline{\mathcal{I}}_d) = \beta_{p.1}(\overline{\mathcal{I}}_d)$ , we know that  $\Gamma[\overline{a}\beta_s(\overline{\mathcal{I}}_d)]$  is a sub-relation of  $\Gamma[\overline{a}]$  that encodes  $o_{p.1}$ . From the previous two steps it follows that there must be some tuple  $\overline{x}' \in \text{adom}(\overline{\mathcal{I}}'_d, \Phi[\overline{a}'])$  such that  $\Phi[\overline{a}'\overline{x}']$  encodes  $o_{p.1}$  and  $\overline{x}'$  is

a permutation of  $\beta_s(\overline{\mathcal{I}}_d)$ . Therefore, node  $n^n$  is potentially pseudo-trivial w.r.t.  $\beta_s(\overline{\mathcal{I}}_d)$ .

**Induction Hypothesis:** Suppose that there exists  $i \in [1, d-1]$  such that the lemma is true whenever  $m > i$ .

**Inductive Step:** Suppose  $m = i$ . Define prefix  $p := c_1 \dots c_{(i-1)}$ . If  $|\mathcal{I}_i| = 0$  then every certificate node at level  $i$  is trivially pseudo-trivial; therefore, assume that  $|\mathcal{I}_i| > 0$ . The proof is identical to the base case, with the following additions.

- By the induction hypothesis, for every level  $j \in [i+1, d]$  every certificate node involving sub-relation  $\Gamma[\beta_s(\overline{\mathcal{I}}_{[1, j-1]})]$  is *potentially* pseudo-trivial w.r.t.  $\beta_s(\overline{\mathcal{I}}_j)$ . Because potentially pseudo-trivial nodes can always be re-arranged until they *are* pseudo-trivial, assume that that within the given certificate every certificate node involving sub-relation  $\Gamma[\beta_s(\overline{\mathcal{I}}_{[1, j-1]})]$  *is* pseudo-trivial w.r.t.  $\beta_s(\overline{\mathcal{I}}_j)$ .
- Within the second step we need to note that there exists an assignment  $\gamma \in \Gamma[\overline{a}]|_{o_{p,1}}$  that satisfies  $\gamma(\overline{\mathcal{I}}_{[i, d]}) = \beta_s(\overline{\mathcal{I}}_{[i, d]})$ . Every tuple  $t \in S_{o_{p,1}}$  corresponds to a sub-relation of  $\Gamma[\overline{ax}]$  that is  $\overline{\mathfrak{s}}_{\mathfrak{s}[i+1, d]}$ -equal to  $\Gamma[\overline{a}\beta_s(\overline{\mathcal{I}}_i)]$ , and so by the induction hypothesis there must exist an assignment  $\gamma' \in \Gamma[\overline{ax}]$  satisfying  $\gamma'(\mathcal{I}_j) = \beta_s(\mathcal{I}_j)$  for all  $j \in [i+1, d]$ . This implies  $\alpha^{-1} \circ \delta^{-1} \circ \gamma'(\mathcal{I}_{[i+1, d]}) = \beta_s(\mathcal{I}_{[i+1, d]})$ , which is needed to contradict Lemma A.4.1.
- Within the sixth step,  $|\mathcal{I}_i| \geq |\mathcal{I}'_i|$  follows from the choice of  $Q$  because the induction hypothesis implies that  $|\mathcal{I}_j| = |\mathcal{I}'_j|$  for every  $j \in [i+1, d]$ .

□

**Lemma A.4.3** *If  $\overline{\mathfrak{s}} = \mathbf{n}^d$  then there exists index-covering homomorphisms in both directions between  $Q$  and  $Q'$ .*

PROOF. Choose any label-generating sequence  $s \in LGS$ . Consider any  $\mathbf{n}^d$ -certificate proving  $\Gamma \doteq_{\mathbf{n}^d} \Phi$ . By Lemma A.4.2, the root node of the certificate is *potentially* pseudo-trivial w.r.t.  $\beta_s(\overline{\mathcal{I}}_1)$ ; therefore, re-arrange the mappings within the node until it *is* pseudo-trivial w.r.t.  $\beta_s(\overline{\mathcal{I}}_1)$ . By recursing into any child  $\mathbf{n}^{d-1}$ -certificate involving sub-relation  $\Gamma[\beta_s(\overline{\mathcal{I}}_1)]$  and repeating this procedure, we eventually arrive at a tuple node proving  $\Gamma[\beta_s(\overline{\mathcal{I}}_{[1, d]})] \doteq_{\emptyset} \Phi[\overline{a}'_{[1, d]}]$  for some  $\overline{a}'_{[1, d]} \in \text{adom}(\overline{\mathcal{I}}'_{[1, d]}, \Phi)$  such that for each  $i \in [1, d]$ ,  $\overline{a}'_i$  is a permutation of  $\beta_s(\overline{\mathcal{I}}_i)$ . Choose any  $\phi \in \Phi[\overline{a}'_{[1, d]}]$ .

Define  $h := \alpha^{-1} \circ \delta^{-1} \circ \phi$ . Then,  $h$  is a mapping from  $\mathcal{B}'$  to  $\alpha^{-1} \circ \delta^{-1}(\text{adom}(\mathbb{D}_Q)) = \mathcal{B} \cup \mathcal{C}_Q$  that satisfies the following properties.

1.  $h(\text{body}_{Q'}) = \alpha^{-1} \circ \delta^{-1} \circ \phi(\text{body}_{Q'}) \subseteq \alpha^{-1} \circ \delta^{-1}(\mathbb{D}_Q) = \text{body}_Q$
2.  $h(\overline{\mathcal{V}}') = \alpha^{-1} \circ \delta^{-1} \circ \phi(\overline{\mathcal{V}}') = \alpha^{-1} \circ \delta^{-1} \circ \beta_s(\overline{\mathcal{V}}) = \overline{\mathcal{V}}$
3.  $\forall i \in [1, d]$ ,  $h(\mathcal{I}'_i) = \alpha^{-1} \circ \delta^{-1} \circ \phi(\mathcal{I}'_i) = \alpha^{-1} \circ \delta^{-1} \circ \beta_s(\mathcal{I}_i) = \mathcal{I}_i$   
which follows because  $\phi(\overline{\mathcal{I}}'_i) = \overline{a}'_i$  is a permutation of  $\beta_s(\overline{\mathcal{I}}_i)$ .

Hence,  $h$  satisfies Definition 3.2.1. The proof is analogous for  $h'$  from  $Q$  to  $Q'$ .  $\square$

## A.5 Arbitrary Signatures

We can combine the proof from the previous three sections to obtain a general proof of Theorem 3.2.4 that handles signatures denoting arbitrary combinations of nested sets, bags, and normalized bags. Like the proof for normalized bags, the general proof relies on both colour-diversification and a labelling system in respect to which certain objects are deemed to be canonical.

We use the same vocabulary of label-generating components/prefixes/sequences as we used in both Sections A.3 and A.4. For any level  $m$  in which  $\xi_m = \mathbf{s}$ , we define  $LGC_m$ ,  $CF-LGC_m$ ,  $LGP_m$ , and  $CF-LGP_m$  as we defined them in Section A.3. When  $\xi_m = \mathbf{n}$  we instead use the corresponding definitions from Section A.4. When  $\xi_m = \mathbf{b}$  we could use either labelling system, but we prefer the one from Section A.4 because it is simpler. The set of label-generating sequences  $LGS$  is defined relative to the interleaved definitions of  $LGC_m$ , and so the label-generating function  $\beta_s : \mathcal{B} \rightarrow (\cup_l \mathcal{B}^l)$  is formed from equations A.19 and A.25 by selecting the labelling system appropriate to the semantics of each level. The canonical database  $\mathbb{D}_Q$  is then constructed using both colour-diversification and labelling (exactly as in Section A.4), and the same definition of canonical objects is used as we used for both sets and normalized bags.

We can prove the existence of (potentially) pseudo-trivial certificate nodes inductively, with the inductive argument varying depending upon the semantics of the level. When  $\xi_m = \mathbf{b}$ , we use Lemma A.2.9 to argue that nodes involving  $\Gamma[\beta_s(\overline{\mathcal{I}}_{[1,m-1]})]$  are potentially pseudo-trivial w.r.t.  $\beta_s(\overline{\mathcal{I}}_m)$ . When  $\xi_m = \mathbf{s}$ , the argument uses Lemmas A.3.4 and A.3.5, and when  $\xi_m = \mathbf{n}$  the argument uses Lemma A.4.2. The construction of index-covering homomorphisms is then identical to Lemma A.4.3.



# References

- [1] Hibernate: Relational persistence for Java and .NET. Available at <http://www.hibernate.org>. Retrieved 17 September 2009.
- [2] Microsoft SQL Server 2008: ADO.NET entity framework. Available at <http://www.microsoft.com/sqlserver/2008/en/us/ado-net-entity.aspx>. Retrieved 17 September 2009.
- [3] Ruby on Rails. Available at <http://www.rubyonrails.org>. Retrieved 17 September 2009.
- [4] Serge Abiteboul and Catriel Beeri. The power of languages for the manipulation of complex values. *VLDB Journal*, 4(4):727–794, 1995.
- [5] Serge Abiteboul and Nicole Bidoit. Non first normal form relations to represent hierarchically organized data. In *Proceedings of the 3rd ACM SIGACT-SIGMOD Symposium on Principles of Database Systems (PODS'84)*, pages 191–200. ACM Press, 1984.
- [6] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [7] William Ward Armstrong. Dependency structures of data base relationships. In *Proceedings of the IFIP Congress*, pages 580–583, August 1974.
- [8] Paolo Atzeni and Nicola M. Morfuni. Functional dependencies in relations with null values. *Information Processing Letters*, 18(4):233–238, May 1984.
- [9] Francois Bancilhon, David Maier, Yehoshua Sagiv, and Jeffrey D Ullman. Magic sets and other strange ways to implement logic programs (extended abstract). In *Proceedings of the 5th ACM SIGACT-SIGMOD Symposium on Principles of Database Systems (PODS'86)*, pages 1–15, New York, NY, USA, 1986. ACM.
- [10] Catriel Beeri and Moshe Y. Vardi. A proof procedure for data dependencies. *Journal of the ACM*, 31(4):718–741, 1984.
- [11] Claude Berge. *Graphs and Hypergraphs*. North-Holland Publishing Company, 1973.

- [12] Gautam Bhargava, Piyush Goel, and Bala Iyer. Hypergraph based reorderings of outer join queries with complex predicates. In *Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data*, pages 304–315, New York, NY, USA, 1995. ACM.
- [13] Nicole Bidoit. The Verso algebra or how to answer queries with fewer joins. *Journal of Computer and System Sciences*, 35(3):321–364, 1987.
- [14] Peter Buneman, Shamim A. Naqvi, Val Tannen, and Limsoon Wong. Principles of programming with complex objects and collection types. *Theoretical Computer Science*, 149(1):3–48, 1995.
- [15] Andrea Cali, Georg Gottlob, and Michael Kifer. Taming the infinite chase: Query answering under expressive relational constraints. In *Proceedings of the 11th International Conference on Principles of Knowledge Representation and Reasoning (KR 2008)*, pages 70–80, 2008.
- [16] Ashok K. Chandra and Philip M. Merlin. Optimal implementation of conjunctive queries in relational data bases. In *Proceedings of the 9th annual ACM Symposium on Theory of Computing*, pages 77–90. ACM Press, 1977.
- [17] Surajit Chaudhuri. An overview of query optimization in relational systems. In *Proceedings of the 17th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS 1998)*, pages 34–43. ACM Press, 1998.
- [18] Surajit Chaudhuri, Ravi Krishnamurthy, Spyros Potamianos, and Kyuseok Shim. Optimizing queries with materialized views. In *Proceedings of the 11th International Conference on Data Engineering (ICDE'95)*, pages 190–200, 1995.
- [19] Surajit Chaudhuri and Moshe Y. Vardi. Optimization of real conjunctive queries. In *Proceedings of the 12th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS 1993)*, pages 59–70. ACM Press, 1993.
- [20] Chandra Chekuri and Anand Rajaraman. Conjunctive query containment revisited. In *Proceedings of the 6th International Conference on Database Theory (ICDT '97)*, volume 1186 of *Lecture Notes in Computer Science*, pages 56–70. Springer, 1997.
- [21] Sara Cohen. *Equivalence, Containment and Rewriting of Aggregate Queries*. PhD thesis, Hebrew University, May 2004.
- [22] Sara Cohen. Equivalence of queries combining set and bag-set semantics. In *Proceedings of the 25th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS'06)*, pages 70–79. ACM Press, 2006.

- [23] Sara Cohen, Werner Nutt, and Yehoshua Sagiv. Containment of aggregate queries. In *Proceedings of the 9th International Conference on Database Theory (ICDT 2003)*, volume 2572 of *Lecture Notes in Computer Science*, pages 111–125. Springer, 2003.
- [24] Sara Cohen, Werner Nutt, and Yehoshua Sagiv. Rewriting queries with arbitrary aggregation functions using views. *ACM Transactions on Database Systems*, 31(2):672–715, 2006.
- [25] Sara Cohen, Werner Nutt, and Yehoshua Sagiv. Deciding equivalences among conjunctive aggregate queries. *Journal of the ACM*, 54(2):5, 2007.
- [26] Sara Cohen, Werner Nutt, and Alexander Serebrenik. Rewriting aggregate queries using views. In *Proceedings of the 18th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS 1999)*, pages 155–166. ACM Press, 1999.
- [27] Sara Cohen, Yehoshua Sagiv, and Werner Nutt. Equivalences among aggregate queries with negation. *ACM Transactions on Computational Logic*, 6(2):328–360, 2005.
- [28] Umeshwar Dayal. Of nests and trees: A unified approach to processing queries that contain nested subqueries, aggregates, and quantifiers. In *Proceedings of 13th International Conference on Very Large Data Bases (VLDB'87)*, pages 197–208, 1987.
- [29] David DeHaan, Per-Åke Larson, and Jingren Zhou. Stacked indexed views in Microsoft SQL Server. In *Proceedings of the 2005 ACM SIGMOD International Conference on Management of Data*, pages 179–190, New York, NY, USA, 2005. ACM.
- [30] Alin Deutsch, Alan Nash, and Jeff Remmel. The chase revisited. In *Proceedings of the 27th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS'08)*, pages 149–158, New York, NY, USA, 2008. ACM.
- [31] Alin Deutsch and Val Tannen. Reformulation of XML queries and constraints. In *Proceedings of the 9th International Conference on Database Theory (ICDT 2003)*, pages 225–241, 2003.
- [32] Xin Dong, Alon Y. Halevy, and Igor Tatarinov. Containment of nested XML queries. In *Proceedings of the 30th International Conference on Very Large Data Bases (VLDB 2004)*, pages 132–143, 2004.
- [33] Ronald Fagin, Phokion G. Kolaitis, Renée J. Miller, and Lucian Popa. Data exchange: semantics and query answering. *Theoretical Computer Science*, 336(1):89–124, 2005.

- [34] Carles Farré, Werner Nutt, Ernest Teniente, and Toni Urpí. Containment of conjunctive queries over databases with null values. In *Proceedings of the 11th International Conference on Database Theory (ICDT 2007)*, number 4353 in Lecture Notes in Computer Science, pages 389–403. Springer, 2007.
- [35] Antonio L. Furtado and Larry Kerschberg. An algebra of quotient relations. In *Proceedings of the 1977 ACM SIGMOD International Conference on Management of Data*, pages 1–8, New York, NY, USA, 1977. ACM Press.
- [36] César Galindo-Legaria. *Algebraic optimization of outerjoin queries*. PhD thesis, Harvard University, June 1992.
- [37] César Galindo-Legaria. Outerjoins as disjunctions. In *Proceedings of the 1994 ACM SIGMOD International Conference on Management of Data*, pages 348–358. ACM Press, 1994.
- [38] César Galindo-Legaria. Parameterized queries and nesting equivalencies. Technical Report MSR-TR-2000-31, Microsoft Research, 2000.
- [39] César Galindo-Legaria and Milind Joshi. Orthogonal optimization of subqueries and aggregation. In *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data*, pages 571–581, New York, NY, USA, 2001. ACM.
- [40] César Galindo-Legaria and Arnon Rosenthal. How to extend a conventional optimizer to handle one- and two-sided outerjoin. In Forouzan Golshani, editor, *Proceedings of the 8th International Conference on Data Engineering (ICDE'92)*, pages 402–409. IEEE Computer Society, 1992.
- [41] César Galindo-Legaria and Arnon Rosenthal. Outerjoin simplification and reordering for query optimization. *ACM Transactions on Database Systems*, 22(1):43–73, 1997.
- [42] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. *Database Systems: The Complete Book*. Prentice Hall, 2002.
- [43] Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W.H. Freeman, 1979.
- [44] Jonathan Goldstein and Per-Åke Larson. Optimizing queries using materialized views: a practical, scalable solution. In *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data*, pages 331–342. ACM Press, 2001.
- [45] Gang Gou, Maxim Kormilitsin, and Rada Chirkova. Query evaluation using overlapping views: Completeness and efficiency. In *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data*, pages 37–48. ACM Press, 2006.



- [46] Goetz Graefe and David J. DeWitt. The EXODUS optimizer generator. In *Proceedings of the 1987 ACM SIGMOD International Conference on Management of Data*, pages 160–172, New York, NY, USA, 1987. ACM.
- [47] Goetz Graefe and William J. McKenna. The Volcano optimizer generator: Extensibility and efficient search. In *Proceedings of the 9th International Conference on Data Engineering (ICDE'93)*, pages 209–218. IEEE Computer Society, 1993.
- [48] Jim Gray, Adam Bosworth, Andrew Layman, and Hamid Pirahesh. Data cube: A relational aggregation operator generalizing Group-By, Cross-Tab, and Sub-Total. In Stanley Y. W. Su, editor, *Proceedings of the 12th International Conference on Data Engineering (ICDE'96)*, pages 152–159. IEEE Computer Society, 1996.
- [49] Stéphane Grumbach, Tova Milo, and Yoram Kornatzky. Calculi for bags and their complexity. In *Proceedings of the 4th International Workshop on Database Programming Languages (DBPL-4)*, Workshops in Computing, pages 65–79. Springer, 1993.
- [50] Stéphane Grumbach, Maurizio Rafanelli, and Leonardo Tininini. Querying aggregate data. In *Proceedings of the 18th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS'99)*, pages 174–184. ACM Press, 1999.
- [51] Stéphane Grumbach, Maurizio Rafanelli, and Leonardo Tininini. On the equivalence and rewriting of aggregate queries. *Acta Informatica*, 40(8):529–584, July 2004.
- [52] Stéphane Grumbach and Leonardo Tininini. Automatic aggregation using explicit metadata. In *Proceedings of the 12th International Conference on Scientific and Statistical Database Management (SSDBM 2000)*, pages 85–94, 2000.
- [53] Stéphane Grumbach and Leonardo Tininini. On the content of materialized aggregate views. *Journal of Computer and System Sciences*, 66(1):133–168, February 2003.
- [54] Ashish Gupta, Venky Harinarayan, and Dallen Quass. Aggregate-query processing in data warehousing environments. In *Proceedings of 21st International Conference on Very Large Data Bases (VLDB'95)*, pages 358–369. Morgan Kaufmann, 1995.
- [55] Marc Gyssens, Jan Paredaens, and Dirk Van Gucht. On a hierarchy of classes for nested databases. *Information Processing Letters*, 36(5):259–266, 1990.

- [56] Laura M. Haas, J. C. Freytag, G. M. Lohman, and H. Pirahesh. Extensible query processing in Starburst. In *Proceedings of the 1989 ACM SIGMOD International Conference on Management of Data*, pages 377–388. ACM Press, 1989.
- [57] Alon Y. Halevy. Answering queries using views: A survey. *VLDB Journal*, 10(4):270–294, 2001.
- [58] Frank Harary. *Graph Theory*. Addison-Wesley, 1969.
- [59] Venky Harinarayan, Anand Rajaraman, and Jeffrey D. Ullman. Implementing data cubes efficiently. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data*, pages 205–216. ACM Press, 1996.
- [60] Gerhard Hill and Andrew Ross. Reducing outer joins. *The VLDB Journal*, 18(3):599–610, June 2009.
- [61] Oscar H. Ibarra and Jianwen Su. On the containment and equivalence of database queries with linear constraints (extended abstract). In *Proceedings of the 16th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'97)*, pages 32–43. ACM Press, 1997.
- [62] Tomasz Imielinski and Jr. Witold Lipski. Incomplete information in relational databases. *Journal of the ACM*, 31(4):761–791, 1984.
- [63] International Standards Organization. Information technology—database language SQL (proposed revised text of DIS 9075). ISO/IEC 9075:1992, July 1992. Available at <http://www.contrib.andrew.cmu.edu/%7Eshadow/sql/sql1992.txt>. Retrieved 17 September 2009.
- [64] Yannis E. Ioannidis and Raghu Ramakrishnan. Containment of conjunctive queries: Beyond relations as sets. *ACM Transactions on Database Systems*, 20(3):288–324, 1995.
- [65] G. Jaeschke and H. J. Schek. Remarks on the algebra of non first normal form relations. In *Proceedings of the 1st ACM SIGACT-SIGMOD Symposium on Principles of Database Systems (PODS'82)*, pages 124–138, New York, NY, USA, 1982. ACM.
- [66] T.S. Jayram, Phokion G. Kolaitis, and Erik Vee. The containment problem for real conjunctive queries. In *Proceedings of the 25th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS'06)*, pages 80–89. ACM Press, 2006.
- [67] D. S. Johnson and A. Klug. Testing containment of conjunctive queries under functional and inclusion dependencies. *Journal of Computer and System Sciences*, 28(1):167–189, 1984.

- [68] Vitaliy L. Khizder and Grant E. Weddell. Reasoning about uniqueness constraints in object relational databases. *IEEE Transactions on Knowledge and Data Engineering*, 15(5):1295–1306, 2003.
- [69] Won Kim. On optimizing an SQL-like nested query. *ACM Transactions on Database Systems*, 7(3):443–469, 1982.
- [70] Anthony Klug. On conjunctive queries containing inequalities. *Journal of the ACM*, 35(1):146–160, 1988.
- [71] Per-Åke Larson and Jingren Zhou. View matching for outer-join views. In *Proceedings of the 31st International Conference on Very Large Data Bases (VLDB 2005)*, 2005.
- [72] Mark Levene and George Loizou. Semantics for null extended nested relations. *ACM Transactions on Database Systems*, 18(3):414–459, 1993.
- [73] Mark Levene and George Loizou. Axiomatisation of functional dependencies in incomplete relations. *Theoretical Computer Science*, 206(1-2):283–300, October 1998.
- [74] Alon Y. Levy, Alberto O. Mendelzon, Yehoshua Sagiv, and Divesh Srivastava. Answering queries using views (extended abstract). In *Proceedings of the 14th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS 1995)*, pages 95–104. ACM Press, 1995.
- [75] Alon Y. Levy and Inderpal Singh Mumick. Reasoning with aggregation constraints. In *Advances in Database Technology - Proceedings of the 5th International Conference on Extending Database Technology (EDBT'96)*, volume 1057 of *Lecture Notes in Computer Science*, pages 514–534. Springer, 1996.
- [76] Alon Y. Levy, Inderpal Singh Mumick, and Yehoshua Sagiv. Query optimization by predicate move-around. In *Proceedings of the 20th International Conference on Very Large Data Bases (VLDB'94)*, pages 96–107. Morgan Kaufmann, 1994.
- [77] Alon Y. Levy and Dan Suciu. Deciding containment for queries with complex objects (extended abstract). In *Proceedings of the 16th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'97)*, pages 20–31. ACM Press, 1997.
- [78] Leonid Libkin and Limsoon Wong. Aggregate functions, conservative extensions, and linear orders. In *Proceedings of the 4th International Workshop on Database Programming Languages (DBPL-4)*, Workshops in Computing, pages 282–294. Springer, 1993.
- [79] Leonid Libkin and Limsoon Wong. New techniques for studying set languages, bag languages and aggregate functions. In *Proceedings of the 13th ACM*

- SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS 1994)*, pages 155–166. ACM Press, 1994.
- [80] Hong-Cheu Liu and Jeffery X. Yu. Algebraic equivalences of nested relational operators. *Information Systems*, 30(3):167–204, 2005.
- [81] Michael Meier, Michael Schmidt, and Georg Lausen. On chase termination beyond stratification. In *Proceedings of the 35th International Conference on Very Large Data Bases (VLDB'09)*, 2009.
- [82] Erik Meijer, Brian Beckman, and Gavin Bierman. LINQ: Reconciling objects, relations and XML in the .NET framework. In *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data*, page 706, New York, NY, USA, 2006. ACM.
- [83] Sergey Melnik, Atul Adya, and Philip A. Bernstein. Compiling mappings to bridge applications and databases. *ACM Transactions on Database Systems*, 33(4):1–50, 2008.
- [84] Edward F. Moore and Claude E. Shannon. Reliable circuits using less reliable relays. *Journal of the Franklin Institute*, 262(3,4):191–208, 281–297, 1956.
- [85] I. S. Mumick, S. J. Finkelstein, Hamid Pirahesh, and Raghu Ramakrishnan. Magic is relevant. In *Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data*, pages 247–258, New York, NY, USA, 1990. ACM.
- [86] Inderpal Singh Mumick, Sheldon J. Finkelstein, Hamid Pirahesh, and Raghu Ramakrishnan. Magic conditions. In *Proceedings of the 9th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'90)*, pages 314–330, New York, NY, USA, 1990. ACM.
- [87] Inderpal Singh Mumick, Hamid Pirahesh, and Raghu Ramakrishnan. The magic of duplicates and aggregates. In *Proceedings of the 16th International Conference on Very Large Data Bases (VLDB'90)*, pages 264–277. Morgan Kaufmann, 1990.
- [88] G. Özsoyoglu, Z. M. Özsoyoglu, and V. Matos. Extending relational algebra and relational calculus with set-valued attributes and aggregate functions. *ACM Transactions on Database Systems*, 12(4):566–592, December 1987.
- [89] Jan Paredaens and Dirk Van Gucht. Converting nested algebra expressions into flat algebra expressions. *ACM Transactions on Database Systems*, 17(1):65–93, March 1992.
- [90] Chang-Sup Park, Myoung Ho Kim, and Yoon-Joon Lee. Finding an efficient rewriting of OLAP queries using materialized views in data warehouses. *Decision Support Systems*, 32(4):379–399, March 2002.

- [91] Glenn N. Paulley and Per-Åke Larson. Exploiting uniqueness in query optimization. In *Proceedings of the 10th International Conference on Data Engineering (ICDE 1994)*, pages 68–79. IEEE Computer Society, 1994.
- [92] Glenn Norman Paulley. *Exploiting Functional Dependence in Query Optimization*. PhD thesis, University of Waterloo, May 2000.
- [93] J. S. Provan and M. O. Ball. Computing network reliability in time polyomial in the number of cuts. *Operations Research*, 32(3):516–526, May-June 1984.
- [94] Jun Rao, Hamid Pirahesh, and Calisto Zuzarte. Canonical abstraction for outerjoin optimization. In *Proceedings of the 2004 ACM SIGMOD International Conference on Management of Data*, pages 671–682, New York, NY, USA, 2004. ACM Press.
- [95] Riccardo Rosati. On the decidability and finite controllability of query processing in databases with incomplete information. In *Proceedings of the 25th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS '06)*, pages 356–365, New York, NY, USA, 2006. ACM.
- [96] Arnon Rosenthal and César Galindo-Legaria. Query graphs, implementing trees, and freely-reorderable outerjoins. In *Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data*, pages 291–299, New York, NY, USA, 1990. ACM.
- [97] Kenneth A. Ross, Divesh Srivastava, Peter J. Stuckey, and S. Sudarshan. Foundations of aggregation constraints. *Theoretical Computer Science*, 193(1-2):149–179, February 1998.
- [98] Steven Rudich and Avi Wigderson, editors. *Computational Complexity Theory*, volume 10 of *IAS/Park City Mathematics Series*, chapter Complexity Theory: from Godel to Feynman, pages 5–87. American Math Society, 2004.
- [99] Yehoshua Sagiv and Mihalis Yannakakis. Equivalences among relational expressions with the union and difference operators. *Journal of the ACM*, 27(4):633–655, 1980.
- [100] Marc H. Scholl. Theoretical foundation of algebraic optimization utilizing unnormalized relations. In *Proceedings of the International Conference on Database Theory (ICDT'86)*, volume 243 of *Lecture Notes in Computer Science*, pages 380–396. Springer, 1986.
- [101] Divesh Srivastava, Shaul Dar, H. V. Jagadish, and Alon Y. Levy. Answering queries with aggregation using views. In T. M. Vijayaraman, Alejandro P. Buchmann, C. Mohan, and Nandlal L. Sarda, editors, *Proceedings of 22nd International Conference on Very Large Data Bases (VLDB'96)*, pages 318–329. Morgan Kaufmann, 1996.

- [102] Dan Suciu. Bounded fixpoints for complex objects. *Theoretical Computer Science*, 176(1-2):283–328, April 1997.
- [103] Stan J. Thomas and Patrick C. Fischer. Nested relational structures. In Paris C. Kanellakis, editor, *The Theory of Databases*, volume 3 of *Advances in Computing Research*, pages 269–307. JAI Press, 1986.
- [104] Transaction Processing Performance Council. TPC benchmark DS (decision support): Standard specification, version 1.0.0d, December 2007. Available at <http://www.tpc.org/tpcds/spec/tpcds1.0.0.d.pdf>. Retrieved 17 September 2009.
- [105] Transaction Processing Performance Council. TPC benchmark H (decision support): Standard specification, revision 2.8.0, September 2008. Available at <http://www.tpc.org/tpch/spec/tpch2.8.0.pdf>. Retrieved 17 September 2009.
- [106] Jeffrey D. Ullman. Information integration using logical views. In *Proceedings of the 6th International Conference on Database Theory (ICDT '97)*, volume 1186 of *Lecture Notes in Computer Science*, pages 19–40. Springer, 1997.
- [107] Jan Van den Bussche, Dirk Van Gucht, and Stijn Vansummeren. Well-definedness and semantic type-checking for the nested relational calculus. *Theoretical Computer Science*, 371(3):183–199, March 2007.
- [108] Ron van der Meyden. The complexity of querying infinite data about linearly ordered domains. *Journal of Computer and System Sciences*, 52(1):113–135, Feb. 1997.
- [109] Yannis Vassiliou. Null values in data base management: a denotational semantics approach. In *Proceedings of the 1979 ACM SIGMOD International Conference on Management of Data*, pages 162–169, New York, NY, USA, 1979. ACM.
- [110] Yannis Vassiliou. Functional dependencies and incomplete information. In *Proceedings of the 6th International Conference on Very Large Data Bases (VLDB'80)*, pages 260–269. VLDB Endowment, 1980.
- [111] Fang Wei and Georg Lausen. Containment of conjunctive queries with safe negation. In *Proceedings of the 9th International Conference on Database Theory (ICDT 2003)*, volume 2572 of *Lecture Notes in Computer Science*, pages 346–360. Springer, 2003.
- [112] Limsoon Wong. Normal forms and conservative properties for query languages over collection types. In *Proceedings of the 12th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'93)*, pages 26–36. ACM Press, 1993.

- [113] Weipeng P. Yan and Per-Åke Larson. Performing group-by before join. In *Proceedings of the 10th International Conference on Data Engineering (ICDE'94)*, pages 89–100. IEEE Computer Society, 1994.
- [114] Weipeng P. Yan and Per-Åke Larson. Eager aggregation and lazy aggregation. In *Proceedings of 21st International Conference on Very Large Data Bases (VLDB'95)*, pages 345–357. Morgan Kaufmann, 1995.
- [115] Markos Zaharioudakis, Roberta Cochrane, George Lapis, Hamid Pirahesh, and Monica Urata. Answering complex SQL queries using automatic summary tables. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, pages 105–116. ACM Press, 2000.