

Ontology Based Data Access and Data Independence

(an alternative look)

David Toman

D.R. Cheriton School of Computer Science

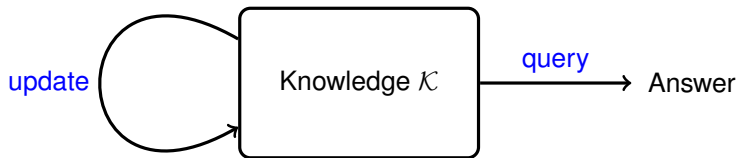
University of

Waterloo



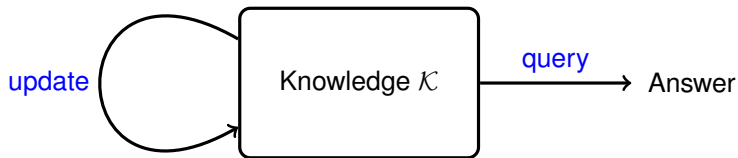
Joint work with Alexander Hudek and Grant Weddell

Knowledge Representation: a Big Picture



What is "Knowledge" (how is it represented, and does the user care?)
⇒ not really as long as the updates and queries "play nicely together"

Knowledge Representation: a Big Picture

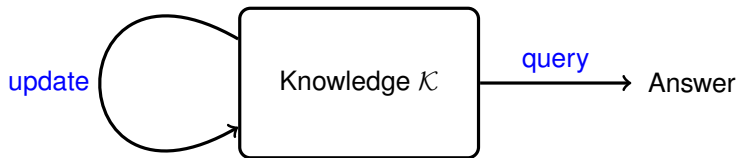


What is “Knowledge” (how is it represented, and does the user care?)
⇒ not really as long as the **updates** and **queries** “play nicely together”

Structured World:

- \mathcal{K} is a (first order) theory,
- queries are (FO) formulæ with answers defined by entailment, and
- updates are (variations on) belief revision.

Knowledge Representation: a Big Picture



What is “Knowledge” (how is it represented, and does the user care?)
⇒ not really as long as the **updates** and **queries** “play nicely together”

Probabilistic World:

- \mathcal{K} is a ML model (e.g., neural net),
- queries are inputs (e.g., photos) and answers are labels
- updates are pairs of, e.g., photos with their labels.

Ontology-based Data Access (OBDA) [Calvanese et al.: Mastro, 2011]

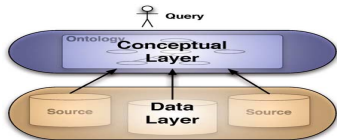
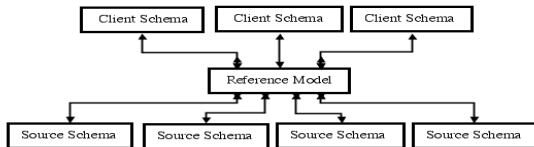


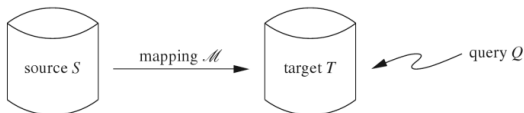
Fig. 1. Ontology-based data access.

Information Integration [Genesereth: Data Integration, 2010]

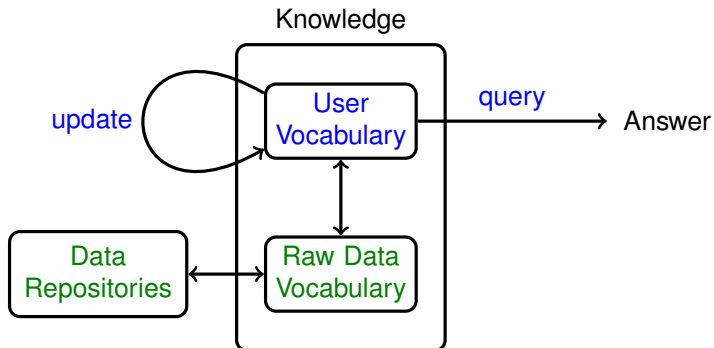


Data Exchange [Arenas et al.: Data Exchange, 2014]

The general setting of data exchange is this:

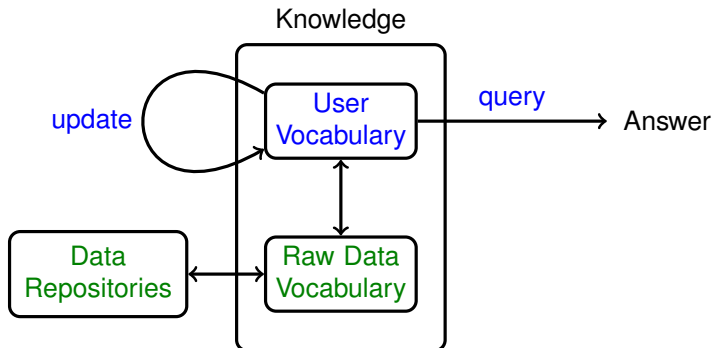


Data vs. Metadata



- Metadata: constraints formulated in FOL (static)
- Data: ground tuples (can be "modified")
 - ⇒ user queries and updates only about data.

Data vs. Metadata



- 1 Metadata: constraints formulated in FOL (static)
- 2 Data: ground tuples (can be “modified”)
⇒ user **queries** and **updates** only about data.

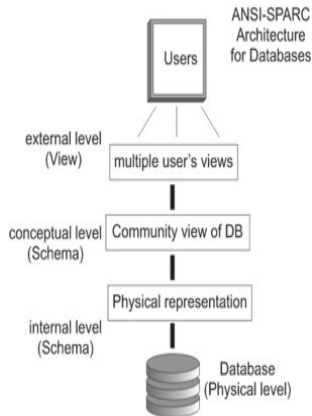
(Physical) Data Independence

IDEA:

Separate the users' view(s) of the data from the way it is physically represented.

Originally just two levels: conceptual/logical and physical [Codd 1970]

- data independence [Codd, 1970] and [Bachman, 1968, Date and Hopewell, 1971]
- ADTs [Liskov and Zilles, 1974]



[ANSI/X3/SPARC Standards Planning and Requirements Committee, Bachman, 1975]

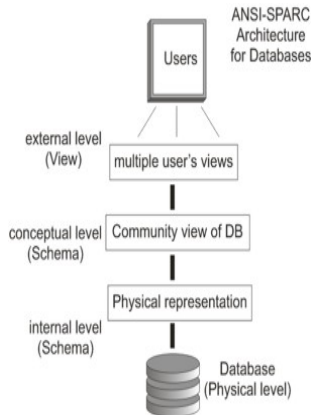
(Physical) Data Independence

IDEA:

Separate the users' view(s) of the data from the way it is physically represented.

Originally just two levels: **conceptual/logical** and **physical** [Codd1970].

- data independence [Codd, 1970] and [Bachman, 1968, Date and Hopewell, 1971]
- ADTs [Liskov and Zilles, 1974]



[ANSI/X3/SPARC Standards Planning and Requirements Committee, Bachman, 1975]

(Physical) Data Independence

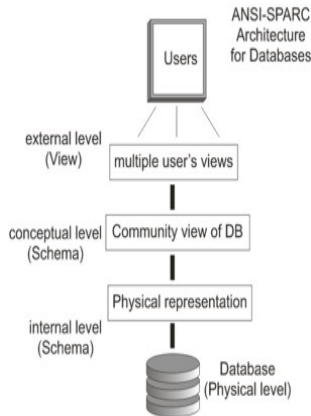
IDEA:

Separate the users' view(s) of the data from the way it is physically represented.

Originally just two levels: **conceptual/logical** and **physical** [Codd1970].

Physical Data Independence and ADTs

- data independence [Codd, 1970] and [Bachman, 1969, Date and Hopewell, 1971]
- ADTs [Liskov and Zilles, 1974]



[ANSI/X3/SPARC Standards Planning and Requirements Committee, Bachman, 1975]

Common Threads and Issues

- In general two *schemas*: **Conceptual/Logical** and **Physical**
 - ⇒ both endowed with *metadata* (vocabulary, constraints, ...)
 - ⇒ mappings that connect the two schemas
- Rules: (1) queries (updates) over the *conceptual/logical* schema only,
(2) raw data as instance of the *physical* schema only,
(3) *conceptual/logical* and *physical* schemata are *disjoint*.

■ Issues to be formalized/fixed:

- Formal description of the two schemas (same formalism for both?)
- Language(s) for metadata, mappings, etc.
- User interface:
 - Data Model
 - Query (and Update) Language (e.g., when is an answer really an answer?)

■ Algorithms/Execution model for user queries and updates

Common Threads and Issues

- In general two *schemas*: **Conceptual/Logical** and **Physical**
 - ⇒ both endowed with *metadata* (vocabulary, constraints, ...)
 - ⇒ mappings that connect the two schemas
- Rules: (1) queries (updates) over the *conceptual/logical* schema only,
(2) raw data as instance of the *physical* schema only,
(3) *conceptual/logical* and *physical* schemata are *disjoint*.
- Issues to be formalized/fixed:
 - 1 Formal description of the two schemas (same formalism for both?)
 - 2 Language(s) for metadata, mappings, etc.
 - 3 **User Interface:**
 - **Data Model**
 - **Query (and Update) Language (e.g., when is an answer really an answer?)**

■ Algorithms/Execution model for user queries and updates

Common Threads and Issues

- In general two *schemas*: **Conceptual/Logical** and **Physical**
 - ⇒ both endowed with *metadata* (vocabulary, constraints, ...)
 - ⇒ mappings that connect the two schemas
- Rules: (1) queries (updates) over the *conceptual/logical* schema only,
(2) raw data as instance of the *physical* schema only,
(3) *conceptual/logical* and *physical* schemata are *disjoint*.
- Issues to be formalized/fixed:
 - 1 Formal description of the two schemas (same formalism for both?)
 - 2 Language(s) for metadata, mappings, etc.
 - 3 User Interface:
 - Data Model
 - Query (and Update) Language (e.g., when is an answer really an answer?)
- Algorithms/Execution model for user queries and updates

Outline

1 Queries

- Review of Standard Approaches to OBDA
- Alternative (DB) Approach

2 Updates

3 How does it Work and (Performance) Bonus

4 Future Research/Open Issues

QUERIES AND QUERY ANSWERING

The Structured/Logical Way (via an OBDA example)

Queries and Ontologies

Queries are answered not only w.r.t. *explicit data* (\mathcal{A})

but also w.r.t. *background knowledge* (\mathcal{T})

⇒ Ontology-based Data Access (OBDA)

Example

■ Socrates is a MAN (explicit data)

■ Every MAN is MORTAL (ontology)

List all MORTALS ⇒ {Socrates} (query)

Using *logical implication* (to define certain answers):

$$\text{Ans}(\varphi, \mathcal{A}, \mathcal{T}) := \{\varphi(a_1, \dots, a_k) \mid \mathcal{T} \cup \mathcal{A} \models \varphi(a_1, \dots, a_k)\}$$

⇒ answers are *ground φ -atoms* logically implied by $\mathcal{A} \cup \mathcal{T}$.

The Structured/Logical Way (via an OBDA example)

Queries and Ontologies

Queries are answered not only w.r.t. *explicit data* (\mathcal{A})
but also w.r.t. *background knowledge* (\mathcal{T})
 \Rightarrow Ontology-based Data Access (OBDA)

Example

- Socrates is a MAN (explicit data)
 - Every MAN is MORTAL (ontology)
- List all MORTALS* \Rightarrow {Socrates} (query)

How do we answer queries?

Using *logical implication* (to define *certain answers*):

$$\text{Ans}(\varphi, \mathcal{A}, \mathcal{T}) := \{\varphi(\mathbf{a}_1, \dots, \mathbf{a}_k) \mid \mathcal{T} \cup \mathcal{A} \models \varphi(\mathbf{a}_1, \dots, \mathbf{a}_k)\}$$

\Rightarrow answers are *ground φ -atoms* logically implied by $\mathcal{A} \cup \mathcal{T}$.

The Logical Way: Complexity

How do we answer queries?

Using *logical implication* (to define *certain answers*):

$$\text{Ans}(\varphi, \mathcal{A}, \mathcal{T}) := \{\varphi(\mathbf{a}_1, \dots, \mathbf{a}_k) \mid \mathcal{T} \cup \mathcal{A} \models \varphi(\mathbf{a}_1, \dots, \mathbf{a}_k)\}$$

\Rightarrow answers are *ground φ -atoms* logically implied by $\mathcal{A} \cup \mathcal{T}$.

- undecidable in general
- coNP-complete (data complexity) from \mathcal{ECI} up to \mathcal{SROIQ} and (U)CQs

LOGSPACE/PETIME (data complexity) for query answering:

- DL-Lite/ \mathcal{EL}_1 / $\mathcal{CFD}_{\text{tree}}^*$ /"rules"-lite (Horn), s-t dependencies, ...
- and (still only) (U)CQ

The Logical Way: Complexity

How do we answer queries?

Using *logical implication* (to define *certain answers*):

$$\text{Ans}(\varphi, \mathcal{A}, \mathcal{T}) := \{\varphi(\mathbf{a}_1, \dots, \mathbf{a}_k) \mid \mathcal{T} \cup \mathcal{A} \models \varphi(\mathbf{a}_1, \dots, \mathbf{a}_k)\}$$

\Rightarrow answers are *ground φ -atoms* logically implied by $\mathcal{A} \cup \mathcal{T}$.

The Bad News

- undecidable in general
- coNP-complete (data complexity) from \mathcal{ELI} up to \mathcal{SROIQ} and (U)CQs

LOGSPACE/P-TIME (data complexity) for query answering:

- DL-Lite/ \mathcal{EL}_1 / $\mathcal{CFD}_{\text{core}}^*$ /"rules"-lite (Horn), s-t dependencies,
- and (still only) (U)CQ

The Logical Way: Complexity

How do we answer queries?

Using *logical implication* (to define *certain answers*):

$$\text{Ans}(\varphi, \mathcal{A}, \mathcal{T}) := \{\varphi(\mathbf{a}_1, \dots, \mathbf{a}_k) \mid \mathcal{T} \cup \mathcal{A} \models \varphi(\mathbf{a}_1, \dots, \mathbf{a}_k)\}$$

\Rightarrow answers are *ground φ -atoms* logically implied by $\mathcal{A} \cup \mathcal{T}$.

The Bad News

- undecidable in general
- coNP-complete (data complexity) from \mathcal{ELI} up to \mathcal{SROIQ} and (U)CQs

The Good News

LOGSPACE/PETIME (data complexity) for query answering:

- DL-Lite/ \mathcal{EL}_{\perp} / $\mathcal{CFD}_{nc}^{\forall}$ /"rules"-lite (Horn), s-t dependencies,...
- and (still only) (U)CQ

Approaches to Ontology-based Data Access

Main Task

INPUT: $\underbrace{\text{Ontology } (\mathcal{T}), \text{ Data } (\mathcal{A})}_{\text{Knowledge Base } (\mathcal{K})}$, and a Query (Q)

OUTPUT: $\{a \mid \mathcal{K} \models Q[a]\}$

1 Reduction to *standard reasoning* (e.g., satisfiability)

2 Reduction to *querying a relational database*

\Rightarrow very good at $\{a \mid \mathcal{A} \models Q[a]\}$ for range restricted Q

\Rightarrow what to do with \mathcal{T} ??

- 1 incorporate into Q (perfect rewriting for DL-Lite et al. (AC^0 logics)); or
- 2 incorporate into \mathcal{A} (combined approach for \mathcal{EL} (PTIME-complete logics));
or sometimes both ($CFDI$ or Horn- ALC^* logics).

Combined Combined Approach to OMQ

Theorem (Combined Combined Approach

[Eiter et al., 2012] [Toman and Weddell, 2013])

Let $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ be a Horn-SHIQ (or Horn-DLFD) knowledge base and φ a conjunctive query. Then there is

- 1 a FO (typically UCQ) query $\varphi_{\mathcal{T}}$, called a **query rewriting**, and
- 2 a Datalog program $\Pi_{\mathcal{T}}$, called an **ABox completion** program,

such that

$$\mathcal{K} \models \varphi[\vec{x} \mapsto \vec{a}] \iff \Pi_{\mathcal{T}}(\mathcal{A}) \models \varphi_{\mathcal{T}}[\vec{x} \mapsto \vec{a}]$$

where $\Pi_{\mathcal{T}}(\mathcal{A})$ is the minimal model of $\Pi_{\mathcal{T}}$ over \mathcal{A} , the completion of \mathcal{A} w.r.t. \mathcal{T} .

Prefix Reformulation (Cavarese et al., 2007)

$\Pi_{\mathcal{T}}$ is an identity on \mathcal{A} ;

Combined Approach (Lutz et al., 2009; Kartchuk et al., 2010)

$\varphi_{\mathcal{T}}$ does not depend on \mathcal{T} .

Combined Combined Approach to OMQ

Theorem (Combined Combined Approach

[Eiter et al., 2012] [Toman and Weddell, 2013])

Let $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ be a Horn-SHIQ (or Horn-DLFD) knowledge base and φ a conjunctive query. Then there is

- 1 a FO (typically UCQ) query $\varphi_{\mathcal{T}}$, called a *query rewriting*, and
- 2 a Datalog program $\Pi_{\mathcal{T}}$, called an *ABox completion* program,

such that

$$\mathcal{K} \models \varphi[\vec{x} \mapsto \vec{a}] \iff \Pi_{\mathcal{T}}(\mathcal{A}) \models \varphi_{\mathcal{T}}[\vec{x} \mapsto \vec{a}]$$

where $\Pi_{\mathcal{T}}(\mathcal{A})$ is the minimal model of $\Pi_{\mathcal{T}}$ over \mathcal{A} , the completion of \mathcal{A} w.r.t. \mathcal{T} .

Perfect Reformulation [Calvanese et al., 2007]:

$\Pi_{\mathcal{T}}$ is an identity on \mathcal{A} ;

Combined Approach [Lutz et al., 2009, Kontchakov et al., 2010]:

$\varphi_{\mathcal{T}}$ does not depend on \mathcal{T} .

The *really* BAD news for OBDA

- no negative queries/sub-queries (SQL??)
- no negations in ABox
- **no closed-world assumption**
- **counter-intuitive query answers**

⇒ the same goes for *information integration, data exchange, etc.*

The *really* BAD news for OBDA

- no negative queries/sub-queries (SQL??)
- no negations in ABox
- **no closed-world assumption**
- **counter-intuitive query answers**

⇒ the same goes for *information integration*, *data exchange*, etc.

Difficulties: (Data) Complexity

Example (Why can't we have CWA?)

- relations: “ $ColNode(x, y)$ ” and “ $Edge(x, y)$ ”;
- the data: graph $(Node^I, Edge^I)$, and $Colour^I = \{r, g, b\}$.

■ ontology: $\forall x. Node(x) \rightarrow \exists y. ColNode(x, y)$

$\forall x, y. ColNode(x, y) \rightarrow Colour(y)$;

■ query: $\exists x, y, z. Edge(x, y) \wedge ColNode(x, z) \wedge ColNode(y, z)$

PROBLEM: having exactly 3 colours is not in DL-Lite (or even HORN).

⇒ same examples for other non-Horn features, e.g., disjunction

Difficulties: (Data) Complexity

Example (Why can't we have CWA?)

- relations: “ $ColNode(x, y)$ ” and “ $Edge(x, y)$ ”;
- the data: graph $(Node^I, Edge^I)$, and $Colour^I = \{r, g, b\}$.
- ontology: $\forall x. Node(x) \rightarrow \exists y. ColNode(x, y)$,
 $\forall x, y. ColNode(x, y) \rightarrow Colour(y)$;

■ query: $\exists x, y, z. Edge(x, y) \wedge ColNode(x, z) \wedge ColNode(y, z)$

PROBLEM: having exactly 3 colours is not in DL-Lite (or even HORN).
⇒ same examples for other non-Horn features, e.g., disjunction

Difficulties: (Data) Complexity

Example (Why can't we have CWA?)

- relations: “ $ColNode(x, y)$ ” and “ $Edge(x, y)$ ”;
- the data: graph $(Node^I, Edge^I)$, and $Colour^I = \{r, g, b\}$.
- ontology: $\forall x. Node(x) \rightarrow \exists y. ColNode(x, y)$,
 $\forall x, y. ColNode(x, y) \rightarrow Colour(y)$;
- query: $\exists x, y, z. Edge(x, y) \wedge ColNode(x, z) \wedge ColNode(y, z)$
... says “the graph $(Node^I, Edge^I)$ is NOT 3-colourable”

PROBLEM: having exactly 3 colours is not in DL-Lite (or even HORN).
⇒ same examples for other non-Horn features, e.g., disjunction

Difficulties: (Data) Complexity

Example (Why can't we have CWA?)

- relations: “ $ColNode(x, y)$ ” and “ $Edge(x, y)$ ”;
- the data: graph $(Node^I, Edge^I)$, and $Colour^I = \{r, g, b\}$.
- ontology: $\forall x. Node(x) \rightarrow \exists y. ColNode(x, y)$,
 $\forall x, y. ColNode(x, y) \rightarrow Colour(y)$;
- query: $\exists x, y, z. Edge(x, y) \wedge ColNode(x, z) \wedge ColNode(y, z)$
... says “the graph $(Node^I, Edge^I)$ is NOT 3-colourable”

PROBLEM: having *exactly 3 colours* is not in DL-Lite (or even HORN).

\Rightarrow same examples for other non-Horn features, e.g., disjunction

(First-order) Query Rewritability

Rewritability (Decision Problem)

Given:

- 1 a TBox \mathcal{T} and
- 2 a Query φ .

Decide whether there is a FO query ψ such that

$$\text{Ans}(\varphi, \mathcal{A}, \mathcal{T}) = \text{Ans}(\psi, \mathcal{A}, \emptyset)$$

for every ABox \mathcal{A} (optionally where ψ is over a sub-vocabulary of \mathcal{T}).

[Bienvenu, Lutz, Wolter: First-Order Rewritability of Atomic Queries in Horn Description Logics. IJCAI 2013. (and many papers followed...)]

Difficulties: Unintuitive Answers

Example

- $EMP(Sue)$
- $EMP \sqsubseteq \exists PHONENUM$ (or $\forall x.EMP(x) \rightarrow \exists y.PHONENUM(x, y)$)

Difficulties: Unintuitive Answers

Example

- $EMP(Sue)$
- $EMP \sqsubseteq \exists PHONENUM$ (or $\forall x.EMP(x) \rightarrow \exists y.PHONENUM(x, y)$)

User: *Does Sue have a phone number?*

Information System: **YES**

Difficulties: Unintuitive Answers

Example

- $EMP(Sue)$
- $EMP \sqsubseteq \exists PHONENUM$ (or $\forall x.EMP(x) \rightarrow \exists y.PHONENUM(x, y)$)

User: *Does Sue have a phone number?*

Information System: **YES**

User: *OK, tell me Sue's phone number!*

Information System: **(no answer)**

Difficulties: Unintuitive Answers

Example

- $EMP(Sue)$
- $EMP \sqsubseteq \exists PHONENUM$ (or $\forall x.EMP(x) \rightarrow \exists y.PHONENUM(x, y)$)

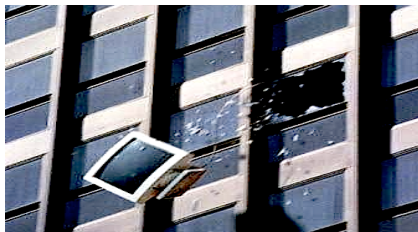
User: *Does Sue have a phone number?*

Information System: **YES**

User: *OK, tell me Sue's phone number!*

Information System: **(no answer)**

User:



The problem: Users (essentially) **EXPECT** CWA

What does $\mathcal{A} = \{EMP(Bob), EMP(Sue)\}$ mean?

OWA: $Bob^I \in EMP^I, Sue^I \in EMP^I$ (KR folks)

CWA: $\{Bob^I, Sue^I\} = EMP^I$ (DB folks and **users**)

... at least for *their* relations (i.e., in the conceptual schema).

The problem: Users (essentially) **EXPECT** CWA

What does $\mathcal{A} = \{EMP(Bob), EMP(Sue)\}$ mean?

OWA: $Bob^I \in EMP^I, Sue^I \in EMP^I$ (KR folks)

CWA: $\{Bob^I, Sue^I\} = EMP^I$ (DB folks and **users**)

... at least for *their* relations (i.e., in the conceptual schema).

Simulations:

CWA in OWA: **closure axioms**: $\forall x. EMP(x) \rightarrow (x = Bob) \vee (x = Sue)$;

OWA in CWA: **auxiliary symbols**: $ExpEMP(Bob), ExpEMP(Sue)$
and **constraints**: $\forall x. ExpEMP(x) \rightarrow EMP(x)$

What does a User Want? ... but is afraid to ask

- 1 what I know and what I don't is just a **single model** (CWA);
- 2 queries are **model-checked** against this model;
- 3 updates change the model into another **single model**.

What does a User Want? ... but is afraid to ask

- 1 what I know and what I don't is just a **single model** (CWA);
- 2 queries are **model-checked** against this model;
- 3 updates change the model into another **single model**.

YES, BUT:

- it better run fast!!
 - ⇒ preferably without having to code algorithms/data structures by hand
- algorithms/performance/data storage-representation/. . .
 - can be changed *without changes to user queries/updates*

User Queries and Updates – for TODAY

Queries: First-order (open) formulae over the user vocabulary

⇒ only *range-restricted* formulae
(i.e., with appropriate *binding pattern restrictions*)

Updates: Instances of *delta-relations* (tuples to be inserted/deleted)
for ALL relations in the user vocabulary

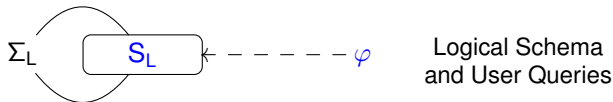
⇒ only *consistency-preserving transactions* allowed

... a.k.a. the Relational Model and Relational Calculus [Codd, 1972].

Rewritability and Definability

User and System Expectations

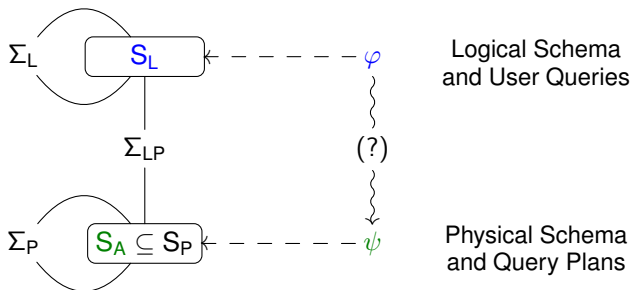
Queries	range-restricted FOL (a.k.a. SQL)
Ontology/Schema	range-restricted FOL $\Sigma := \Sigma_L \cup \Sigma_{LP} \cup \Sigma_P$
Data	CWA (complete information)



Rewritability and Definability

User and System Expectations

Queries	range-restricted FOL over S_L <i>definable w.r.t. Σ and S_A</i>
Ontology/Schema	range-restricted FOL $\Sigma := \Sigma_L \cup \Sigma_{LP} \cup \Sigma_P$
Data	CWA (complete information for S_A symbols)



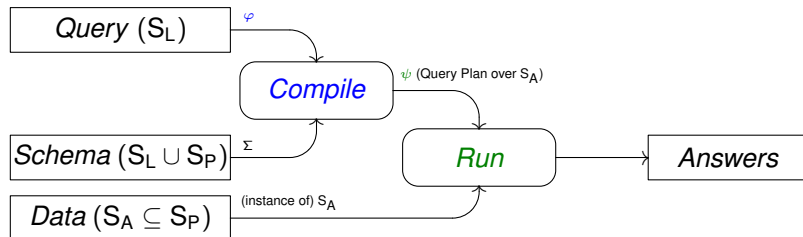
[Borgida, de Bruijn, Franconi, Seylan, Straccia, Toman, Weddell: On Finding Query Rewritings under Expressive Constraints. SEBD 2010: 426-437]

Rewritability and Definability

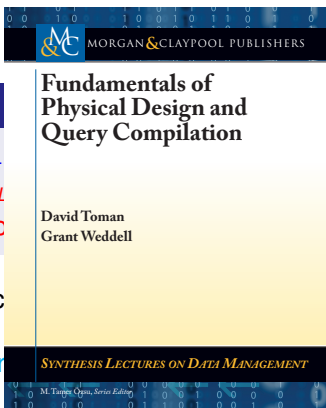
User and System Expectations

Queries	range-restricted FOL over S_L <i>definable</i> w.r.t. Σ and S_A
Ontology/Schema	range-restricted FOL $\Sigma := \Sigma_L \cup \Sigma_{LP} \cup \Sigma_P$
Data	CWA (complete information for S_A symbols)

- to users it looks like a *single model* (of the logical schema)
- implementation can pick from many models
but *definable* queries answer the same in each of them



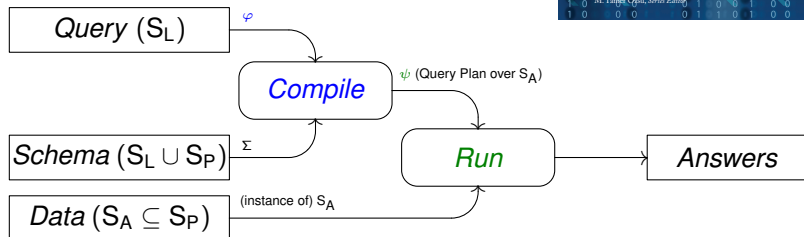
Rewritability and Definability



User and System Expectations

Queries	range-restricted FOL over S_L
Ontology/Schema	range-restricted FOL $\Sigma := \Sigma_L$
Data	CWA (complete information for)

- to users it looks like a *single model* (of the logic)
- implementation can pick from many models
but *definable queries answer*



©2011

This is NOT OMQ/OBDA (by example)

Example (Horizontal Partition)

$$S_L = \{\text{emp}/1, \text{wkr}/1, \text{mgr}/1\} \text{ and } \Sigma_L = \left\{ \begin{array}{l} \text{mgr}(\mathbf{x}) \vee \text{wkr}(\mathbf{x}) \leftrightarrow \text{emp}(\mathbf{x}) \\ \text{mgr}(\mathbf{x}) \wedge \text{wkr}(\mathbf{x}) \rightarrow \perp \end{array} \right\}$$

This is NOT OMQ/OBDA (by example)

Example (Horizontal Partition)

$$S_L = \{\text{emp}/1, \text{wkr}/1, \text{mgr}/1\} \text{ and } \Sigma_L = \left\{ \begin{array}{l} \text{mgr}(\mathbf{x}) \vee \text{wkr}(\mathbf{x}) \leftrightarrow \text{emp}(\mathbf{x}) \\ \text{mgr}(\mathbf{x}) \wedge \text{wkr}(\mathbf{x}) \rightarrow \perp \end{array} \right\}$$

$S_A = \{\text{emp}/1, \text{mgr}/1\}$, note that $\text{wkr}/1$ is NOT in S_A ;

This is NOT OMQ/OBDA (by example)

Example (Horizontal Partition)

$$S_L = \{\text{emp}/1, \text{wkr}/1, \text{mgr}/1\} \text{ and } \Sigma_L = \left\{ \begin{array}{l} \text{mgr}(x) \vee \text{wkr}(x) \leftrightarrow \text{emp}(x) \\ \text{mgr}(x) \wedge \text{wkr}(x) \rightarrow \perp \end{array} \right\}$$

$S_A = \{\text{emp}/1, \text{mgr}/1\}$, note that $\text{wkr}/1$ is NOT in S_A ;

Query $\{x \mid \text{wkr}(x)\}$ over an instance $\text{mgr} = \{\text{Fred}\}$, and $\text{emp} = \{\text{Fred}, \text{Wilma}\}$

This is NOT OMQ/OBDA (by example)

Example (Horizontal Partition)

$$S_L = \{\text{emp}/1, \text{wkr}/1, \text{mgr}/1\} \text{ and } \Sigma_L = \left\{ \begin{array}{l} \text{mgr}(x) \vee \text{wkr}(x) \leftrightarrow \text{emp}(x) \\ \text{mgr}(x) \wedge \text{wkr}(x) \rightarrow \perp \end{array} \right\}$$

$S_A = \{\text{emp}/1, \text{mgr}/1\}$, note that $\text{wkr}/1$ is NOT in S_A ;

Query $\{x \mid \text{wkr}(x)\}$ over an instance $\text{mgr} = \{\text{Fred}\}$, and $\text{emp} = \{\text{Fred}, \text{Wilma}\}$

Certain Answer under OWA: $\{\}$

Answer under CWA: $\{\text{Wilma}\}$

(obtained by executing the *plan* $\{x \mid \text{emp}(x) \wedge \neg \text{mgr}(x)\}$).

What can we do with this?

Goal #2

Generate query plans *that compete with hand-written programs in C*

- 1 standard RDBMS physical designs (and more),
 - access to search structures (index access and selection),
 - horizontal partitioning/sharding,
 - column store/index-only plans,
- 2 pointer-based data structures (including main mamory),
- 3 hash-based access to data (including hash-joins),
- 4 multi-level storage (aka disk/remote/distributed files), ...
- 5 materialized views,
- 6 updates through logical schema
- 7 ...

... all **without** having to code (too much) in C/C++ !

Standard Physical Designs

- 1 scanning (flat) files
- 2 primary and secondary indices (via record ids/addresses)
- 3 horizontal partitioning/sharding
- 4 column store/index-only plans
- 5 (disjoint) generalizations

Pointers in Main Memory-Logical Schema

```
CREATE TABLE employee (  
  num      INTEGER NOT NULL,  
  name     CHAR(20),  
  worksin  INTEGER NOT NULL  
  PRIMARY KEY (num),  
  FOREIGN KEY (worksin)  
           REFERENCES department  
)  
  
CREATE TABLE department (  
  num      INTEGER NOT NULL,  
  name     CHAR(50),  
  manager  INTEGER NOT NULL,  
  PRIMARY KEY (num),  
  FOREIGN KEY (manager)  
           REFERENCES employee  
)
```

this corresponds to

- $S_L = \{\text{employee}/3, \text{department}/3\}$ and
- $\Sigma_L = \{\text{employee}(x, y_1, z_1) \wedge \text{employee}(x, y_2, z_2) \rightarrow y_1 = y_2 \wedge z_1 = z_2, \\ \text{employee}(x, y, z) \rightarrow \exists u, v. \text{department}(z, u, v), \dots \text{and many more}\}.$

Pointers in Main Memory-Logical Schema

```
CREATE TABLE employee (  
  num      INTEGER NOT NULL,  
  name     CHAR(20),  
  worksin  INTEGER NOT NULL  
  PRIMARY KEY (num),  
  FOREIGN KEY (worksin)  
           REFERENCES department  
)
```

```
CREATE TABLE department (  
  num      INTEGER NOT NULL,  
  name     CHAR(50),  
  manager  INTEGER NOT NULL,  
  PRIMARY KEY (num),  
  FOREIGN KEY (manager)  
           REFERENCES employee  
)
```

this corresponds to

- $S_L = \{\text{employee}/3, \text{department}/3\}$ and
- $\Sigma_L = \{\text{employee}(x, y_1, z_1) \wedge \text{employee}(x, y_2, z_2) \rightarrow y_1 = y_2 \wedge z_1 = z_2,$
 $\text{employee}(x, y, z) \rightarrow \exists u, v. \text{department}(z, u, v), \dots \text{and many more}\}.$

additional logical constraints (for example):

- managers are employees that manage a department (a view)
- managers work in their own departments (business rule)
- workers and managers partition employees (partition), etc.

Pointers in Main Memory-Physical Design

1 Records:

```
struct emp {                struct dept {
    int      num;           int      num;
    char[20] name;         char[50] name;
    dept*    dept; };     mgr*    emp;  };
```

2 a linked list of emp records.

that corresponds to

■ Access paths (S_A):

- $empfile/1/0$: set (list) of *addresses* of emp records;
- $emp-num/2/1$: pairs emp record address-emp number (pointer navigation)
same for $emp-name/2/1$ and $emp-dept/2/1$;
- $dept-num/2/1$: pairs dept record address-dept number
same for $dept-name/2/1$ and $dept-mgr/2/1$.

■ Integrity constraints ($\Sigma_P \cup \Sigma_{LP}$):

$$\forall x, y, z. employee(x, y, z) \rightarrow \exists w. empfile(w) \wedge emp-num(w, x),$$
$$\forall a, x. empfile(a) \wedge emp-num(a, x) \rightarrow \exists y, z. employee(x, y, z), \dots$$

Query Plans that Navigate Pointers

1 List employee numbers, names, and departments (`employee(x, y, z)`):

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \text{dept-num}(d, z)$$

or, in C-like syntax: for `e` in `empfile` do

`x := e->num;`

`y := e->name;`

`d := e->dept;`

`z := d->num;`

Query Plans that Navigate Pointers

- 1 List employee numbers, names, and departments ($\text{employee}(x, y, z)$):

$$\begin{aligned} \exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \text{dept-num}(d, z) \end{aligned}$$

- 2 List worker numbers and names ($\exists z. \text{worker}(x, y, z)$):

$$\begin{aligned} \exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \neg \text{dept-mgr}(d, e) \end{aligned}$$

Query Plans that Navigate Pointers

- 1 List employee numbers, names, and departments ($\text{employee}(x, y, z)$):

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \text{dept-num}(d, z)$$

- 2 List worker numbers and names ($\exists z. \text{worker}(x, y, z)$):

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \neg \text{dept-mgr}(d, e)$$

- 3 List all department numbers and their names ($\exists z. \text{department}(x, y, z)$):

➤ Caveat: we do NOT have a (direct) way to “scan” departments! ⚡

Query Plans that Navigate Pointers

- 1 List employee numbers, names, and departments ($\text{employee}(x, y, z)$):

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \text{dept-num}(d, z)$$

- 2 List worker numbers and names ($\exists z. \text{worker}(x, y, z)$):

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \neg \text{dept-mgr}(d, e)$$

- 3 List all department numbers and their names ($\exists z. \text{department}(x, y, z)$):

$$\exists d, e. \text{empfile}(e) \wedge \text{emp-dept}(e, d) \\ \wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y)$$

\Rightarrow needs “departments have at least one employee”.

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-dept}(e, d) \\ \wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y) \wedge \text{dept-mgr}(d, e)$$

\Rightarrow needs “managers work in their own departments”.

Query Plans that Navigate Pointers

- 1 List employee numbers, names, and departments ($\text{employee}(x, y, z)$):

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \text{dept-num}(d, z)$$

- 2 List worker numbers and names ($\exists z. \text{worker}(x, y, z)$):

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-num}(e, x) \wedge \text{emp-name}(e, y) \\ \wedge \text{emp-dept}(e, d) \wedge \neg \text{dept-mgr}(d, e)$$

- 3 List all department numbers and their names ($\exists z. \text{department}(x, y, z)$):

$$\exists d, e. \text{empfile}(e) \wedge \text{emp-dept}(e, d) \\ \wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y)$$

\Rightarrow needs “departments have at least one employee”.

... needs *duplicate elimination* during projection.

$$\exists e, d. \text{empfile}(e) \wedge \text{emp-dept}(e, d) \\ \wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y) \wedge \text{dept-mgr}(d, e)$$

\Rightarrow needs “managers work in their own departments”.

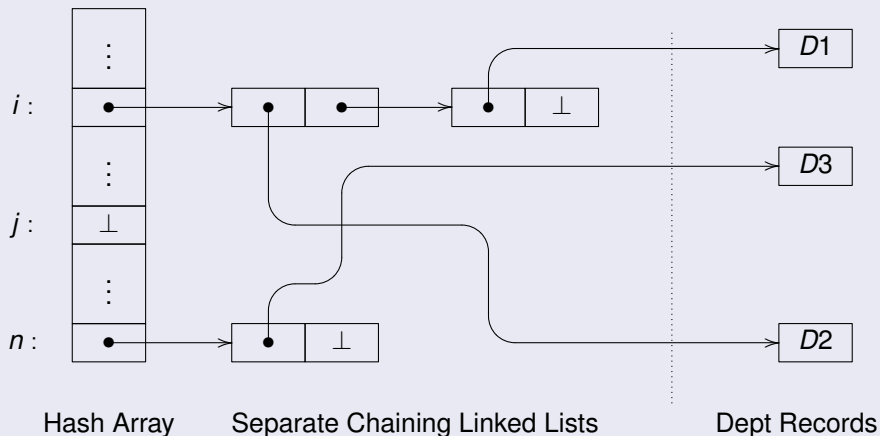
... *NO duplicate elimination* during projection.

... and we can actually synthesize this!

```
david$ compile tests/new_fe/book-em-v4-new-query.fol
query(dept2,2,0,[var(0,0,1,int),var(0,0,2,int)]) <->
  ex(var(0,76,4),
    ex(var(0,81,5),
      and (
        and (
          empfile(var(0,76,4))
          emp_dept(var(0,76,4),var(0,81,5))
        )
        and (
          and (
            dept_num(var(0,81,5),var(0,0,1))
            dept_name(var(0,81,5),var(0,0,2))
          )
          dept_mgr(var(0,81,5),var(0,76,4))
        )
      )
    )
  )
```

What can it do: Hashing, Lists, et al.

Hash Index with (list-based) Separate Chaining



What can it do: Hashing, Linked lists, et al.

Hash Index on department's name:

Access paths:

$S_A \supseteq \{\text{hash}/2/1, \text{hasharraylookup}/2/1, \text{listscan}/2/1\}.$

Physical Constraints:

$$\begin{aligned} \Sigma_{LP} \supseteq \{ & \forall x, y. ((\text{deptfile}(x) \wedge \text{dept-name}(x, y)) \rightarrow \exists z, w. (\text{hash}(y, z) \\ & \wedge \text{hasharraylookup}(z, w) \wedge \text{listscan}(w, x))), \\ & \forall x, y. (\text{hash}(x, y) \rightarrow \exists z. \text{hasharraylookup}(y, z)), \\ & \forall x, y. (\text{listscan}(x, y) \rightarrow \text{deptfile}(y)) \quad \} \end{aligned}$$

What can it do: Hashing, Linked lists, et al.

Hash Index on department's name:

Access paths:

$$S_A \supseteq \{\text{hash}/2/1, \text{hasharraylookup}/2/1, \text{listscan}/2/1\}.$$

Physical Constraints:

$$\begin{aligned} \Sigma_{LP} \supseteq \{ & \forall x, y. ((\text{deptfile}(x) \wedge \text{dept-name}(x, y)) \rightarrow \exists z, w. (\text{hash}(y, z) \\ & \wedge \text{hasharraylookup}(z, w) \wedge \text{listscan}(w, x))), \\ & \forall x, y. (\text{hash}(x, y) \rightarrow \exists z. \text{hasharraylookup}(y, z)), \\ & \forall x, y. (\text{listscan}(x, y) \rightarrow \text{deptfile}(y)) \quad \} \end{aligned}$$

4 List departments and their managers given dept name p

$$(\exists y. \text{department}(x_1, p, y) \wedge \text{employee}(y, x_2)\{p\}):$$

$$\begin{aligned} \exists h, l, d, e. & \text{hash}(p, h) \wedge \text{hasharraylookup}(h, l) \wedge \\ & \text{listscan}(l, d) \wedge \text{dept-name}(d, p) \wedge \\ & \text{dept-num}(d, x_1) \wedge \text{dept-mgr}(d, e) \wedge \text{emp-name}(e, x_2) \end{aligned}$$

What can this do: two-level store

The access path `empfile` is refined by `emppages/1/0` and `emprecords/2/1`:

`emppages` returns (sequentially) disk pages containing `emp` records, and `emprecords` given a disc page, returns `emp` records in that page.

- 5 List all employees with the same name

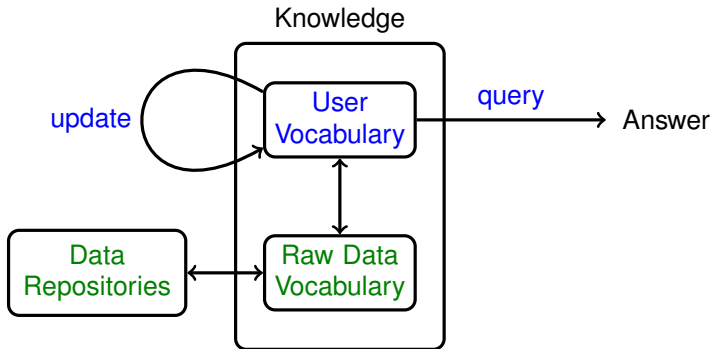
$(\exists z. \text{employee}(x_1, z) \wedge \text{employee}(x_2, z))$:

$\exists y, z, w, v, p, q. \text{emppages}(p) \wedge \text{emppages}(q)$
 $\wedge \text{emprecords}(p, y) \wedge \text{emp-num}(y, x_1) \wedge \text{emp-name}(y, w)$
 $\wedge \text{emprecords}(q, z) \wedge \text{emp-num}(z, x_2) \wedge \text{emp-name}(z, v)$
 $\wedge \text{compare}(w, v).$

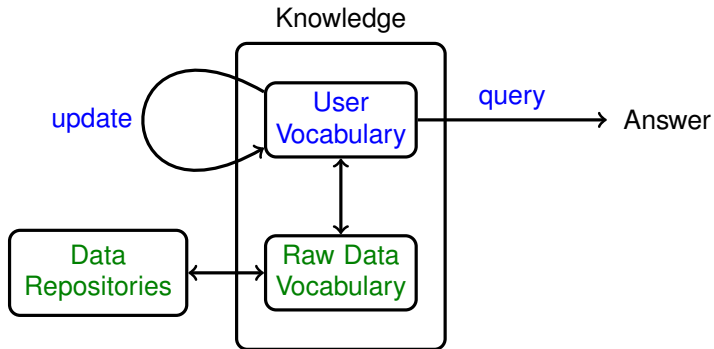
\Rightarrow this plan implements the *block nested loops join* algorithm.

UPDATES

Updates

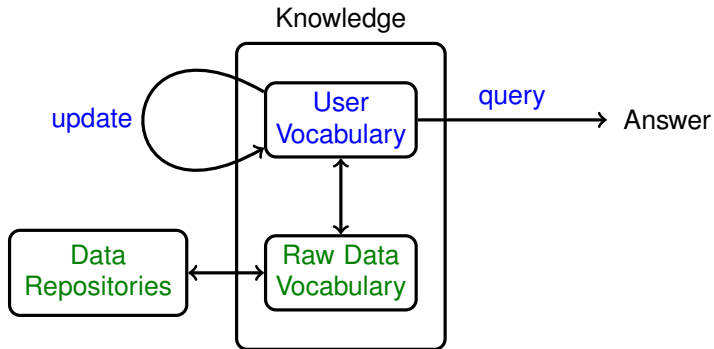


Updates



- 1 Katsuno, Mendelzon: On the Difference between Updating a Knowledge Base and Revising It. KR 1991.
- 2 De Giacomo, Lenzerini, Poggi, Rosati: On Instance-level Update and Erasure in Description Logic Ontologies. J. Log. Comput. 19(5) 2009.

Updates



- 1 Katsuno, Mendelzon: On the Difference between Updating a Knowledge Base and Revising It. KR 1991.
- 2 De Giacomo, Lenzerini, Poggi, Rosati: On Instance-level Update and Erasure in Description Logic Ontologies. J. Log. Comput. 19(5) 2009.

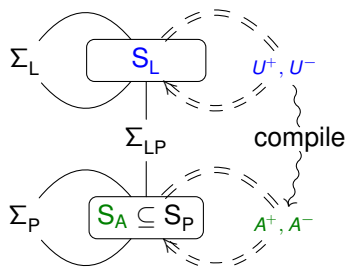
... we use *definable updates* approach instead...

Updates and Definability

User updates *through logical schema ONLY*:

⇒ supplying “delta” relations (sets of tuples)

- Delta relations: R^+ (insertions) and R^- (deletions);

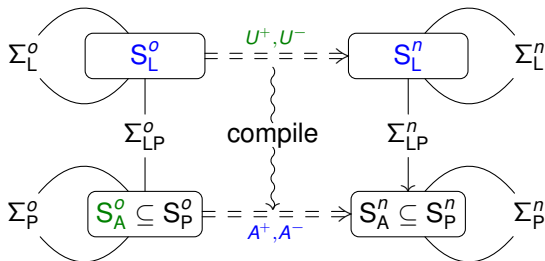


Updates and Definability

User updates *through logical schema ONLY*:

⇒ supplying “delta” relations (sets of tuples)

- Delta relations: R^+ (insertions) and R^- (deletions);



Update turned into *definability* question

Is A^n (or A^+, A^-) definable in terms of $A_j^o \in S_A^o$ (old access paths)
and U_j^+, U_j^- (user updates) for every access path $A \in S_A$?

Update Example

Example (Update Schema)

$$S_L = \{\text{emp}^o/1, \text{wkr}^o/1, \text{mgr}^o/1\} \cup \{\text{emp}^n/1, \text{wkr}^n/1, \text{mgr}^n/1\} \cup$$

$$\Sigma = \left\{ \begin{array}{l} \text{mgr}^o(x) \vee \text{wkr}^o(x) \leftrightarrow \text{emp}^o(x) \\ \text{mgr}^o(x) \wedge \text{wkr}^o(x) \rightarrow \perp \end{array} \right\} \cup \left\{ \begin{array}{l} \text{mgr}^n(x) \vee \text{wkr}^n(x) \leftrightarrow \text{emp}^n(x) \\ \text{mgr}^n(x) \wedge \text{wkr}^n(x) \rightarrow \perp \end{array} \right\} \cup$$

Update Example

Example (Update Schema)

$$S_L = \{\text{emp}^o/1, \text{wkr}^o/1, \text{mgr}^o/1\} \cup \{\text{emp}^n/1, \text{wkr}^n/1, \text{mgr}^n/1\} \cup \\ \{\text{emp}^+/1, \text{wkr}^+/1, \text{mgr}^+/1\} \cup \{\text{emp}^-/1, \text{wkr}^-/1, \text{mgr}^-/1\}$$

$$\Sigma = \left\{ \begin{array}{l} \text{mgr}^o(x) \vee \text{wkr}^o(x) \leftrightarrow \text{emp}^o(x) \\ \text{mgr}^o(x) \wedge \text{wkr}^o(x) \rightarrow \perp \end{array} \right\} \cup \left\{ \begin{array}{l} \text{mgr}^n(x) \vee \text{wkr}^n(x) \leftrightarrow \text{emp}^n(x) \\ \text{mgr}^n(x) \wedge \text{wkr}^n(x) \rightarrow \perp \end{array} \right\} \cup \\ \left\{ \begin{array}{l} \text{emp}^o(x) \vee \text{emp}^+(x) \leftrightarrow \text{emp}^n(x) \vee \text{emp}^-(x), \text{emp}^o(x) \wedge \text{emp}^+(x) \rightarrow \perp, \dots \\ \text{mgr}^o(x) \vee \text{mgr}^+(x) \leftrightarrow \text{mgr}^n(x) \vee \text{mgr}^-(x), \text{mgr}^o(x) \wedge \text{mgr}^+(x) \rightarrow \perp, \dots \\ \text{wkr}^o(x) \vee \text{wkr}^+(x) \leftrightarrow \text{wkr}^n(x) \vee \text{wkr}^-(x), \text{wkr}^o(x) \wedge \text{wkr}^+(x) \rightarrow \perp, \dots \end{array} \right\}$$

Update Example

Example (Update Schema)

$$\mathcal{S}_L = \{\text{emp}^o/1, \text{wkr}^o/1, \text{mgr}^o/1\} \cup \{\text{emp}^n/1, \text{wkr}^n/1, \text{mgr}^n/1\} \cup \\ \{\text{emp}^+/1, \text{wkr}^+/1, \text{mgr}^+/1\} \cup \{\text{emp}^-/1, \text{wkr}^-/1, \text{mgr}^-/1\}$$

$$\Sigma = \left\{ \begin{array}{l} \text{mgr}^o(\mathbf{x}) \vee \text{wkr}^o(\mathbf{x}) \leftrightarrow \text{emp}^o(\mathbf{x}) \\ \text{mgr}^o(\mathbf{x}) \wedge \text{wkr}^o(\mathbf{x}) \rightarrow \perp \end{array} \right\} \cup \left\{ \begin{array}{l} \text{mgr}^n(\mathbf{x}) \vee \text{wkr}^n(\mathbf{x}) \leftrightarrow \text{emp}^n(\mathbf{x}) \\ \text{mgr}^n(\mathbf{x}) \wedge \text{wkr}^n(\mathbf{x}) \rightarrow \perp \end{array} \right\} \cup \\ \left\{ \begin{array}{l} \text{emp}^o(\mathbf{x}) \vee \text{emp}^+(\mathbf{x}) \leftrightarrow \text{emp}^n(\mathbf{x}) \vee \text{emp}^-(\mathbf{x}), \text{emp}^o(\mathbf{x}) \wedge \text{emp}^+(\mathbf{x}) \rightarrow \perp, \dots \\ \text{mgr}^o(\mathbf{x}) \vee \text{mgr}^+(\mathbf{x}) \leftrightarrow \text{mgr}^n(\mathbf{x}) \vee \text{mgr}^-(\mathbf{x}), \text{mgr}^o(\mathbf{x}) \wedge \text{mgr}^+(\mathbf{x}) \rightarrow \perp, \dots \\ \text{wkr}^o(\mathbf{x}) \vee \text{wkr}^+(\mathbf{x}) \leftrightarrow \text{wkr}^n(\mathbf{x}) \vee \text{wkr}^-(\mathbf{x}), \text{wkr}^o(\mathbf{x}) \wedge \text{wkr}^+(\mathbf{x}) \rightarrow \perp, \dots \end{array} \right\}$$

Result(s)

$$\text{emp}^+(\mathbf{x}) := (\text{wkr}^+(\mathbf{x}) \vee \text{mgr}^+(\mathbf{x})) \wedge \neg \text{emp}^o(\mathbf{x})$$

$$\text{emp}^-(\mathbf{x}) := (\text{wkr}^-(\mathbf{x}) \vee \text{mgr}^-(\mathbf{x})) \wedge \neg(\text{wkr}^+(\mathbf{x}) \vee \text{mgr}^+(\mathbf{x}))$$

... and we can again synthesize that

```
david$ compile tests/new_fe/mgrwkr-upd.fol
query(empi,1,0,[var(0,0,1,int)]) <->
  and ( or (
    wkri(var(0,0,1))
    mgri(var(0,0,1))
  ) not (
    empo(var(0,0,1))
  ) )
```

```
query(empd,1,0,[var(0,0,1,int)]) <->
  and ( or (
    wkrd(var(0,0,1))
    mgrd(var(0,0,1))
  ) not ( or (
    mgri(var(0,0,1))
    wkri(var(0,0,1))
  ) ) )
```


Unknown/Anonymous Values?

Example (Add a new Worker record (needs an address))

```
INSERT into worker values (1234);
```

⇒ the request then needs to be translated to

```
INSERT into worker-physical values (0xFE1234, 1234);
```

⇒ but where did 0xFE1234 came from? (definability issue!)

Unknown/Anonymous Values?

Example (Add a new Worker record (needs an address))

```
INSERT into worker values (1234);
```

⇒ the request then needs to be translated to

```
INSERT into worker-physical values (0xFE1234, 1234);
```

⇒ but where did 0xFE1234 came from? (definability issue!)

Constant Complement: [Bancilhon, Spyrtos: Update semantics of relational views. ACM Trans. Database Syst. 6(4), 1981.]

additional access paths that *provide* such values:

⇒ in our case `worker-addr(id, address)`

⇒ and where $worker^+ = \{(1234)\}$

$worker-physical^+(x_2, x_1) = worker^+(x_1) \wedge worker-addr(x_1, x_2)$

Unknown/Anonymous Values?

Example (Add a new Worker record (needs an address))

```
INSERT into worker values (1234);
```

⇒ the request then needs to be translated to

```
INSERT into worker-physical values (0xFE1234, 1234);
```

⇒ but where did 0xFE1234 came from? (definability issue!)

Constant Complement: [Bancilhon, Spyrtos: Update semantics of relational views. ACM Trans. Database Syst. 6(4), 1981.]

additional access paths that *provide* such values:

⇒ in our case `worker-addr(id, address)`

⇒ and where $worker^+ = \{(1234)\}$

$worker-physical^+(x_2, x_1) = worker^+(x_1) \wedge worker-addr(x_1, x_2)$

The additional access path(s) correspond to *space allocation*

... and cyclic dependencies are broken via *reification*.

... more details and examples in

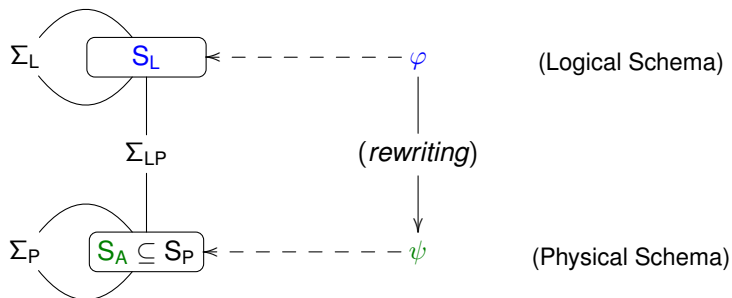


HOW DOES IT ALL WORK?

The Plan

Definability and Rewriting

Queries	range-restricted FOL over S_L <i>definable</i> w.r.t. Σ and S_A
Ontology/Schema	range-restricted FOL
Data	CWA (complete information for S_A symbols)



Query Plans via Interpolation

IDEA #1: Plans as Formulas

Represent *query plans* as (annotated) range-restricted formulas ψ over S_A :

atomic formula	\mapsto	access path (<code>get-first-get-next iterator</code>)
conjunction	\mapsto	nested loops join
existential quantifier	\mapsto	projection (annotated w/duplicate info)
disjunction	\mapsto	concatenation
negation	\mapsto	simple complement

Query Plans via Interpolation

IDEA #1: Plans as Formulas

Represent *query plans* as (annotated) range-restricted formulas ψ over S_A :

atomic formula	\mapsto	access path (<code>get-first-get-next iterator</code>)
conjunction	\mapsto	nested loops join
existential quantifier	\mapsto	projection (annotated w/duplicate info)
disjunction	\mapsto	concatenation
negation	\mapsto	simple complement

\Rightarrow reduces correctness of ψ to logical implication $\Sigma \models \varphi \leftrightarrow \psi$

Query Plans via Interpolation

IDEA #1: Plans as Formulas

Represent *query plans* as (annotated) range-restricted formulas ψ over S_A :

atomic formula	\mapsto	access path (<code>get-first-get-next iterator</code>)
conjunction	\mapsto	nested loops join
existential quantifier	\mapsto	projection (annotated w/duplicate info)
disjunction	\mapsto	concatenation
negation	\mapsto	simple complement

\Rightarrow reduces correctness of ψ to logical implication $\Sigma \models \varphi \leftrightarrow \psi$

Non-logical (but necessary) Add-ons

1 Non-logical properties/operators

- binding patterns
- duplication of data and duplicate-preserving/eliminating projections
- sortedness of data (with respect to the *iterator semantics*) and sorting

2 Cost model

Beth Definability and Craig Interpolation

IDEA #2: What Queries do we allow?

We only allow queries that have *the same answer* in every model of Σ for a fixed interpretation of the signature S_A (i.e., where the actual data is).

Beth Definability and Craig Interpolation

IDEA #2: What Queries do we allow?

We only allow queries that have *the same answer* in every model of Σ for a fixed interpretation of the signature S_A (i.e., where the actual data is).

How do we test for this?

φ is *Beth definable* [Beth'56] if

$$\Sigma \cup \Sigma' \models \varphi \rightarrow \varphi'$$

where Σ' (φ') is Σ (φ) in which symbols *NOT in S_A* are *primed*, respectively.

Beth Definability and Craig Interpolation

IDEA #2: What Queries do we allow?

We only allow queries that have *the same answer* in every model of Σ for a fixed interpretation of the signature S_A (i.e., where the actual data is).

How do we test for this?

φ is *Beth definable* [Beth'56] if

$$\Sigma \cup \Sigma' \models \varphi \rightarrow \varphi'$$

where Σ' (φ') is Σ (φ) in which symbols *NOT in S_A* are *primed*, respectively.

How do we find ψ ?

If $\Sigma \cup \Sigma' \models \varphi \rightarrow \varphi'$ then there is ψ s.t. $\Sigma \cup \Sigma' \models \varphi \rightarrow \psi \rightarrow \varphi'$ with $\mathcal{L}(\psi) \subseteq \mathcal{L}(S_A)$.
... ψ is called the *Craig Interpolant* [Craig'57].

... and we can extract ψ from a (TABLEAU) proof of $\Sigma \cup \Sigma' \models \varphi \rightarrow \varphi'$

TABLEAU Proofs

1 Make entailment into formula to be proven:

$$\begin{aligned} \Sigma \cup \Sigma' \models \varphi \rightarrow \varphi' \quad \text{iff} \quad & \models ((\bigwedge \Sigma) \wedge (\bigwedge \Sigma')) \rightarrow (\varphi \rightarrow \varphi') \quad \text{iff} \\ & \models (\bigwedge \Sigma) \rightarrow ((\bigwedge \Sigma') \rightarrow (\varphi \rightarrow \varphi')) \quad \text{iff} \\ & \models (\bigwedge \Sigma) \rightarrow (\varphi \rightarrow ((\bigwedge \Sigma') \rightarrow \varphi')) \quad \text{iff} \\ & \models ((\bigwedge \Sigma) \wedge (\varphi) \rightarrow ((\bigwedge \Sigma') \rightarrow \varphi')) \end{aligned}$$

or, equivalently, $(\bigwedge \Sigma) \wedge (\bigwedge \Sigma') \wedge \alpha \wedge \neg \alpha'$ is inconsistent;

2 Apply TABLEAU expansion rules to (*) until all branches are closed;

$$\frac{S \cup \{\alpha, \beta\}}{S} \alpha \wedge \beta \in S, \quad \frac{S \cup \{\alpha\} \quad S \cup \{\beta\}}{S} \alpha \vee \beta \in S, \quad \dots$$

where a branch S is *closed* if $\{\alpha, \neg \alpha\} \subseteq S$ (or $\perp \in S$).

TABLEAU Proofs

- 1 Make entailment into formula to be proven:

$$\begin{aligned} \Sigma \cup \Sigma' \models \varphi \rightarrow \varphi' \quad \text{iff} \quad & \models ((\bigwedge \Sigma) \wedge (\bigwedge \Sigma')) \rightarrow (\varphi \rightarrow \varphi') \quad \text{iff} \\ & \models (\bigwedge \Sigma) \rightarrow ((\bigwedge \Sigma') \rightarrow (\varphi \rightarrow \varphi')) \quad \text{iff} \\ & \models (\bigwedge \Sigma) \rightarrow (\varphi \rightarrow ((\bigwedge \Sigma') \rightarrow \varphi')) \quad \text{iff} \\ & \models ((\bigwedge \Sigma) \wedge (\varphi \rightarrow ((\bigwedge \Sigma') \rightarrow \varphi'))) \end{aligned}$$

or, equivalently, $(\bigwedge \Sigma)^L \wedge (\bigwedge \Sigma')^R \wedge \alpha^L \wedge \neg \alpha'^R$ is inconsistent;

- 2 Apply TABLEAU expansion rules to (*) until all branches are closed;

$$\frac{S \cup \{\alpha, \beta\}}{S} \alpha \wedge \beta \in S, \quad \frac{S \cup \{\alpha\} \quad S \cup \{\beta\}}{S} \alpha \vee \beta \in S, \quad \dots$$

where a branch S is *closed* if $\{\alpha, \neg \alpha\} \subseteq S$ (or $\perp \in S$).

- 3 Extract an interpolant from (2); this needs extra $.^L$ and $.^R$ annotations.

TABLEAU and Interpolant Extraction (sketch)

- an interpolant $S \xrightarrow{\text{int}} \psi$; invariant $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$
where S^L and S^R are the $.^L$ and $.^R$ -tagged subsets of S ;

TABLEAU and Interpolant Extraction (sketch)

- an interpolant $S \xrightarrow{\text{int}} \psi$; invariant $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$
where S^L and S^R are the $.^L$ and $.^R$ -tagged subsets of S ;
- tableau rules (sample):
 - LR clash $S \cup \{R^L, \neg R^R\} \xrightarrow{\text{int}} R$, where $R \in S_A$
as $(\bigwedge S^L \wedge R^L) \rightarrow R$ and $R \rightarrow (R^R \vee \neg \bigwedge S^R)$;

TABLEAU and Interpolant Extraction (sketch)

- an interpolant $S \xrightarrow{\text{int}} \psi$; invariant $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$
where S^L and S^R are the $.^L$ and $.^R$ -tagged subsets of S ;

- tableau rules (sample):

- LR clash $S \cup \{R^L, \neg R^R\} \xrightarrow{\text{int}} R$, where $R \in S_A$

as $(\bigwedge S^L \wedge R^L) \rightarrow R$ and $R \rightarrow (R^R \vee \neg \bigwedge S^R)$;

- L-conjunction
$$\frac{S \cup \{\alpha^L, \beta^L\} \xrightarrow{\text{int}} \delta}{S \cup \{(\alpha \wedge \beta)^L\} \xrightarrow{\text{int}} \delta}$$

as $(\bigwedge S^L \wedge \alpha^L \wedge \beta^L) \rightarrow \delta$ implies $(\bigwedge S^L \wedge (\alpha \wedge \beta)^L) \rightarrow \delta$;

TABLEAU and Interpolant Extraction (sketch)

- an interpolant $S \xrightarrow{\text{int}} \psi$; invariant $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$
where S^L and S^R are the $.^L$ and $.^R$ -tagged subsets of S ;
- tableau rules (sample):

- LR clash $S \cup \{R^L, \neg R^R\} \xrightarrow{\text{int}} R$, where $R \in S_A$

as $(\bigwedge S^L \wedge R^L) \rightarrow R$ and $R \rightarrow (R^R \vee \neg \bigwedge S^R)$;

- L-conjunction
$$\frac{S \cup \{\alpha^L, \beta^L\} \xrightarrow{\text{int}} \delta}{S \cup \{(\alpha \wedge \beta)^L\} \xrightarrow{\text{int}} \delta}$$

as $(\bigwedge S^L \wedge \alpha^L \wedge \beta^L) \rightarrow \delta$ implies $(\bigwedge S^L \wedge (\alpha \wedge \beta)^L) \rightarrow \delta$;

- R-Disjunction
$$\frac{S \cup \{\alpha^R\} \xrightarrow{\text{int}} \delta_\alpha \quad S \cup \{\beta^R\} \xrightarrow{\text{int}} \delta_\beta}{S \cup \{(\alpha \vee \beta)^R\} \xrightarrow{\text{int}} \delta_\alpha \wedge \delta_\beta}$$

as $\bigwedge S^L \rightarrow \delta_\alpha, \delta_\alpha \rightarrow \neg(\alpha^R \wedge \bigwedge S^R)$ and $\bigwedge S^L \rightarrow \delta_\beta, \delta_\beta \rightarrow \neg(\beta^R \wedge \bigwedge S^R)$
implies $(\bigwedge S^L) \rightarrow \delta_\alpha \wedge \delta_\beta, \delta_\alpha \wedge \delta_\beta \rightarrow \neg((\alpha \vee \beta)^R \wedge \bigwedge S^R)$.

- etc. (see [Fitting] for details)

Issues with TABLEAU

Dealing with the *subformula property* of Tableau

- ⇒ analytic tableau *explores* formulas *structurally*
- ⇒ (to large degree) the structure of interpolant depends on where access paths are present in queries/constraints.

Factoring *logical reasoning* from *plan enumeration*

- ⇒ backtracking tableau to get alternative plans: too slow, too few plans

Issues with TABLEAU

Dealing with the *subformula property* of Tableau

- ⇒ analytic tableau *explores* formulas *structurally*
- ⇒ (to large degree) the structure of interpolant depends on where access paths are present in queries/constraints.

IDEA #3:

Separate *general constraints* from *physical rules* in the formulation of the definability question (and the subsequent interpolant extraction):

$$\Sigma^L \cup \Sigma^R \cup \Sigma^{LR} \models \varphi^L \rightarrow \varphi^R \text{ where } \Sigma^{LR} = \{\forall \bar{x}. P^L \leftrightarrow P \leftrightarrow P^R \mid P \in S_A\}$$

Factoring *logical reasoning* from *plan enumeration*

- ⇒ backtracking tableau to get alternative plans: too slow, too few plans

IDEA #4:

Define *conditional tableau* exploration (using general constraints) and separate it from plan generation (using physical rules)

Conditional Formulæ and Tableau

Conditional Formulæ

$\varphi[C]\{B\}$ where C is a set of (ground) atoms over S_A and B branch descriptors
 $\Rightarrow \varphi$ only exists in T if all atoms in C are “used” in a plan tableau

Absorbed Range-restricted Formulæ: ANF

$$Q ::= R(\bar{x}) \mid \perp \mid Q \wedge Q \mid Q \vee Q \mid \forall \bar{x}. R(\bar{x}) \rightarrow Q,$$

... and all \exists 's are Skolemized.

Conditional Tableau Rules for ANF

$$\frac{S \cup \{\varphi[C], \psi[C]\}}{(\varphi \wedge \psi)[C] \in S} \text{ (conj)}$$

$$\frac{S \cup \{\varphi[C]\} \quad S \cup \{\psi[C]\}}{(\varphi \vee \psi)[C] \in S} \text{ (disj)}$$

$$\frac{S \cup \{(\varphi[\bar{t}/\bar{x}])[C \cup D]\}}{\{R(\bar{t})[C], (\forall \bar{x}. R(\bar{x}) \rightarrow \varphi)[D]\} \subseteq S} \text{ (abs)}$$

$$\frac{S \cup \{R(\bar{t})[R(\bar{t})]\}}{S} R(\bar{x}) \in S_A \text{ (phys)}$$

Conditional Tableau and Interpolation

Conditional Tableau for (Q, Σ, S_A)

Proof trees (T^L, T^R) : T^L for $\Sigma^L \cup \{Q^L(\bar{a})\}$ over $\{P^L \mid P \in S_A\}$
 T^R for $\Sigma^R \cup \{Q^R(\bar{a}) \rightarrow \perp\}$ over $\{P^R \mid P \in S_A\}$

Closing Set(s)

We call a set C of literals over S_A a *closing set* for T if, for every branch

- 1 there is an atom $R(\bar{t})[D]$ such that $D \cup \{\neg R(\bar{t})\} \subseteq C$.
- 2 there is $\perp[D]$ such that $D \subseteq C$.

\Rightarrow there are many different *minimal* closing sets for T .

Observation

For an arbitrary closing set C , the interpolant for $T^L(T^R)$ is $\perp(T)$.

Conditional Tableau through an Example

Example (Horizontal Partition)

$$S_L = \{\text{emp}/1, \text{wkr}/1, \text{mgr}/1\} \text{ and } \Sigma_L = \left\{ \begin{array}{l} \text{mgr}(x) \vee \text{wkr}(x) \leftrightarrow \text{emp}(x) \\ \text{mgr}(x) \wedge \text{wkr}(x) \rightarrow \perp \end{array} \right\}$$

$S_A = \{\text{emp}/1, \text{mgr}/1\}$, note that $\text{wkr}/1$ is NOT in S_A ;

Conditional Tableau through an Example

Example (Horizontal Partition)

$$S_L = \{\text{emp}/1, \text{wkr}/1, \text{mgr}/1\} \text{ and } \Sigma_L = \left\{ \begin{array}{l} \text{mgr}(\mathbf{x}) \vee \text{wkr}(\mathbf{x}) \leftrightarrow \text{emp}(\mathbf{x}) \\ \text{mgr}(\mathbf{x}) \wedge \text{wkr}(\mathbf{x}) \rightarrow \perp \end{array} \right\}$$

$S_A = \{\text{emp}/1, \text{mgr}/1\}$, note that $\text{wkr}/1$ is NOT in S_A ;

ANF of Σ

$$\left\{ \begin{array}{ll} \text{emp}(\mathbf{x}) \rightarrow \text{wkr}(\mathbf{x}) \vee \text{mgr}(\mathbf{x}), & \text{mgr}(\mathbf{x}) \rightarrow \text{emp}(\mathbf{x}), \\ \text{mgr}(\mathbf{x}) \wedge \text{wkr}(\mathbf{x}) \rightarrow \perp, & \text{wkr}(\mathbf{x}) \rightarrow \text{emp}(\mathbf{x}) \end{array} \right\}$$

Conditional Tableau through an Example

ANF of Σ

$$\left\{ \begin{array}{l} \text{emp}(\mathbf{x}) \rightarrow \text{wkr}(\mathbf{x}) \vee \text{mgr}(\mathbf{x}), \quad \text{mgr}(\mathbf{x}) \rightarrow \text{emp}(\mathbf{x}), \\ \text{mgr}(\mathbf{x}) \wedge \text{wkr}(\mathbf{x}) \rightarrow \perp, \quad \text{wkr}(\mathbf{x}) \rightarrow \text{emp}(\mathbf{x}) \end{array} \right\}$$

ANF of the Query

$$\begin{array}{l} \text{query}(\mathbf{x}) \rightarrow \text{wkr}(\mathbf{x}), \quad \text{wrk}(\mathbf{x}) \rightarrow \text{query}(\mathbf{x}) \quad \text{and} \\ \text{query}^L(\mathbf{0}), \quad \text{query}^R(\mathbf{0}) \rightarrow \perp \end{array}$$

Conditional Tableau through an Example

ANF of Σ

$$\left\{ \begin{array}{l} \text{emp}(x) \rightarrow \text{wkr}(x) \vee \text{mgr}(x), \quad \text{mgr}(x) \rightarrow \text{emp}(x), \\ \text{mgr}(x) \wedge \text{wkr}(x) \rightarrow \perp, \quad \text{wkr}(x) \rightarrow \text{emp}(x) \end{array} \right\}$$

ANF of the Query

$$\begin{array}{l} \text{query}(x) \rightarrow \text{wkr}(x), \quad \text{wrk}(x) \rightarrow \text{query}(x) \quad \text{and} \\ \text{query}^L(0), \quad \text{query}^R(0) \rightarrow \perp \end{array}$$

$$\begin{array}{c} \mathcal{T}^L \qquad \qquad \qquad \mathcal{T}^R \\ \hline \begin{array}{l} \text{query}(0)[]\{\} \\ \text{wkr}(0)[]\{\} \\ \text{emp}(0)[]\{\} \end{array} \quad \rightarrow \quad \begin{array}{l} \text{emp}(0)[\text{emp}(0)]\{\} \\ \text{wkr}(0)[\text{emp}(0)]\{1.0\} \\ \text{query}(0)[\text{emp}(0)]\{1.0\} \\ \perp[\text{emp}(0)]\{1.0\} \end{array} \\ \begin{array}{l} \text{mgr}(0)[\text{mgr}(0)]\{\} \\ \perp[\text{mgr}(0)]\{\} \end{array} \quad \leftarrow \quad \begin{array}{l} \text{mgr}(0)[\text{emp}(0)]\{1.1\} \end{array} \end{array}$$

Conditional Tableau through an Example

T^L	T^R
query(0)[]{} wkr(0)[]{} emp(0)[]{} mgr(0)[mgr(0)]{} \perp [mgr(0)]{} 	\rightarrow emp(0)[emp(0)]{} wkr(0)[emp(0)]{1.0} query(0)[emp(0)]{1.0} \perp [emp(0)]{1.0} \leftarrow mgr(0)[emp(0)]{1.1}

Closing sets:

$\{\neg\text{emp}(0)\}, \{\text{mgr}(0)\}$

$\{\neg\text{mgr}(0), \text{emp}(0)\}$

Conditional Tableau through an Example

T^L	T^R
$query(0)[]\{\}$	
$wkr(0)[]\{\}$	
$emp(0)[]\{\}$	$\rightarrow emp(0)[emp(0)]\{\}$
	$wkr(0)[emp(0)]\{1.0\}$
	$query(0)[emp(0)]\{1.0\}$
	$\perp[emp(0)]\{1.0\}$
$mgr(0)[mgr(0)]\{\}$	$\leftarrow mgr(0)[emp(0)]\{1.1\}$
$\perp[mgr(0)]\{\}$	

Closing sets:

$\{\neg emp(0)\}, \{mgr(0)\}$

$\{\neg mgr(0), emp(0)\}$

Query Plan:

$query(x) \leftrightarrow emp(x) \wedge \neg mgr(x)$

Plan Enumeration

Physical Tableau T^P for a Plan P

$P : L_P$	R_P
$R(\bar{t}) : \{\{\neg R^L(\bar{t})\}\}$	$\{\{R^R(\bar{t})\}\}$
$P_1 \wedge P_2 : L_{P_1} \cup L_{P_2}$	$\{S_1 \cup S_2 \mid S_1 \in R_{P_1}, S_2 \in R_{P_2}\}$
$P_1 \vee P_2 : \{S_1 \cup S_2 \mid S_1 \in L_{P_1}, S_2 \in L_{P_2}\}$	$R_{P_1} \cup R_{P_2}$
$\neg P_1 : \{\{L^L(\bar{t}) \mid L^R(\bar{t}) \in S\} \mid S \in R_{P_1}\}$	$\{\{L^R(\bar{t}) \mid L^L(\bar{t}) \in S\} \mid S \in L_{P_1}\}$
$\exists x.P_1 : L_{P_1[t/x]}$	$R_{P_1[t/x]}$

Observation

For a range-restricted formula P over S_A there is an analytic tableau tree T^P that uses only formulæ in $\Sigma^{LR} \cup \{\forall x.\text{true}^R(x)\}$ such that:

- 1 Open branches of T^P correspond to *sets of literals* $C \in L_P$ (left branch) or $C \in R_P$ (right branch); and
- 2 The interpolant extracted from the closed tableau $T^P[T^L, T^R]$, the *closure of (T^L, T^R)* by (the *branches* of) T^P , is logically equivalent to P .

Logical&Physical Combined, Controlling the Search

Basic Strategy

- 1 build (T^L, T^R) for (Q, Σ, S_A) to a *certain depth*,
- 2 build T^P and test if each element in $L_P(R_P)$ closes $T^L(T^R)$.
if so, $T^P[T^L, T^R]$ is closed tableau yielding an interpolant equivalent to P ;
(... otherwise extend depth in step 1 and repeat.)

NOTE: in step 2 we can “test” many P s (plan enumeration), but
how do we know which ones to try? while building these bottom-up?

Controlling the Search

- only use the (phys) rule in $T^L(T^R)$ for $R(\bar{t})$ that appears in $T^R(T^L)$,
- only consider *fragments* that help closing (T^L, T^R)
⇒ this is determined using the minimal closing sets for (T^L, T^R) .

... combine with A^* search (among P s) with respect to a *cost model*.

Closing Sets (more complex example)

1 Byte code generation for $q/2$

```
q(x, y) <-> ex(z, table(x, x, z) and table(z, y, y)
              and not table(x, x, x))
```

2 Conditional Tableau Construction

```
L { -p0basetable(s119:7, s114:3, s10:2, s10:2) }
L { -p0basetable(s119:5, s10:1, s10:1, s114:3) }
L { +p0basetable(sr19:8, s10:1, s10:1, s10:1) }
R { -p0basetable(sr19:8, s10:1, s10:1, s10:1),
    +p0basetable(s119:7, s114:3, s10:2, s10:2),
    +p0basetable(s119:5, s10:1, s10:1, s114:3) }
```

3 Cost-based Optimization (A*)

4 C code Generation (+ compilation/linking w/runtime library)

[Hudek, Toman, Weddell: On Enumerating Query Plans Using Analytic Tableau. TABLEAUX 2015.]

[Toman, Weddell: An Interpolation-based Compiler and Optimizer for Relational Queries (System design Report). IWIL-LPAR 2017.]

CONDITIONAL TABLEAU: Result

```
query(q, 2, 0, [var(0, 0, 1, int), var(0, 0, 2, int)]) <->
  ex(var(0, 14, 3),
    ex(var(0, 19, 5),
      ex(var(0, 19, 7),
        and (
          and (
            p0basetable(var(0, 19, 7), var(0, 14, 3),
              var(0, 0, 2), var(0, 0, 2))
            p0basetable(var(0, 19, 5), var(0, 0, 1),
              var(0, 0, 1), var(0, 14, 3))
          )
        )
      )
    )
  )
  not (
    ex(var(1, 19, 8),
      p0basetable(var(1, 19, 8), var(0, 0, 1),
        var(0, 0, 1), var(0, 0, 1))
    )
  )
) ) ) ) )
```

Postprocessing: Duplicate Elimination Elimination

IDEA:

Separate the projection operation ($\exists \bar{x}.$) to

- a duplicate preserving projection (\exists) and
- an explicit (idempotent) duplicate elimination operator ($\{\cdot\}$).

Postprocessing: Duplicate Elimination Elimination

IDEA:

Separate the projection operation ($\exists \bar{x}$.) to

- a duplicate preserving projection (\exists) and
- an explicit (idempotent) duplicate elimination operator ($\{\cdot\}$).

Use the following rewrites to eliminate/minimize the use of $\{\cdot\}$:

$$\psi[\{R(x_1, \dots, x_k)\}] \leftrightarrow \psi[R(x_1, \dots, x_k)]$$

$$\psi[\{\psi_1 \wedge \psi_2\}] \leftrightarrow \psi[\{\psi_1\} \wedge \{\psi_2\}]$$

$$\psi[\{\neg\psi_1\}] \leftrightarrow \psi[\neg\psi_1]$$

$$\psi[\neg\{\psi_1\}] \leftrightarrow \psi[\neg\psi_1]$$

$$\psi[\{\psi_1 \vee \psi_2\}] \leftrightarrow \psi[\{\psi_1\} \vee \{\psi_2\}] \quad \text{if } \Sigma \cup \{\psi[]\} \models \psi_1 \wedge \psi_2 \rightarrow \perp$$

$$\psi[\{\exists x.\psi_1\}] \leftrightarrow \psi[\exists x.\{\psi_1\}]$$

if

$$\Sigma \cup \{\psi[] \wedge \psi_1[y_1/x] \wedge \psi_1[y_2/x]\} \models y_1 \approx y_2$$

Postprocessing: Duplicate Elimination Elimination

IDEA:

Separate the projection operation ($\exists \bar{x}.$) to

- a duplicate preserving projection (\exists) and
- an explicit (idempotent) duplicate elimination operator ($\{\cdot\}$).

Use the following rewrites to eliminate/minimize the use of $\{\cdot\}$:

$$\psi[\{R(x_1, \dots, x_k)\}] \leftrightarrow \psi[R(x_1, \dots, x_k)]$$

$$\psi[\{\psi_1 \wedge \psi_2\}] \leftrightarrow \psi[\{\psi_1\} \wedge \{\psi_2\}]$$

$$\psi[\{\neg\psi_1\}] \leftrightarrow \psi[\neg\psi_1]$$

$$\psi[\neg\{\psi_1\}] \leftrightarrow \psi[\neg\psi_1]$$

$$\psi[\{\psi_1 \vee \psi_2\}] \leftrightarrow \psi[\{\psi_1\} \vee \{\psi_2\}] \quad \text{if } \Sigma \cup \{\psi[]\} \models \psi_1 \wedge \psi_2 \rightarrow \perp$$

$$\psi[\{\exists x.\psi_1\}] \leftrightarrow \psi[\exists x.\{\psi_1\}] \quad \text{if } \Sigma \cup \{\psi[] \wedge \psi_1[y_1/x] \wedge \psi_1[y_2/x]\} \models y_1 \approx y_2$$

... reasoning abstracted in FunDL-Lite (a PTIME fragment)

[Toman, Weddell: Using Feature-Based Description Logics to avoid Duplicate Elimination in Object-Relational Query Languages. *Künstliche Intell.* 34(3): 2020]

Postprocessing: Duplicate Elimination Elimination

IDEA:

Separate the projection operation ($\exists \bar{x}.$) to

- a duplicate preserving projection (\exists) and
- an explicit (idempotent) duplicate elimination operator ($\{\cdot\}$).

$\exists d, e. \text{empfile}(e) \wedge \text{emp-dept}(e, d)$
 $\wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y)$
 \Rightarrow dept number (x) and name (y)
do NOT functionally determine employee (e).

$\exists e, d. \text{empfile}(e) \wedge \text{emp-dept}(e, d)$
 $\wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y) \wedge \text{dept-mgr}(d, e)$
 \Rightarrow dept number (x) and name (y)
do functionally determine employee (e) and dept (d).

Postprocessing: Duplicate Elimination

IDEA:

Separate the projection operation ($\exists \bar{x}.$) to

- a duplicate preserving projection (\exists) and
- an explicit (idempotent) duplicate elimination operator ($\{\cdot\}$).

$\{\exists d, e. \text{empfile}(e) \wedge \text{emp-dept}(e, d)$
 $\wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y)\}$

\Rightarrow dept number (x) and name (y)

do NOT functionally determine employee (e).

... needs *duplicate elimination* during projection.

$\exists e, d. \text{empfile}(e) \wedge \text{emp-dept}(e, d)$
 $\wedge \text{dept-num}(d, x) \wedge \text{dept-name}(d, y) \wedge \text{dept-mgr}(d, e)$

\Rightarrow dept number (x) and name (y)

do functionally determine employee (e) and dept (d).

... does NOT need *duplicate elimination* during projection.

Interpolation (summary)

(A) solution: a *conditional tableau*:

1 reformulate the interpolation problem to

$$\Sigma^L \cup \Sigma^R \cup \Sigma^{LR} \models \varphi^L \rightarrow \varphi^R \text{ where } \Sigma^{LR} = \{\forall \bar{x}. P^L \leftrightarrow P \leftrightarrow P^R \mid P \in S_A\}$$

2 use conditional (ground) atoms to generate closing sets: sets of S_A literals that (fully) close the tableau

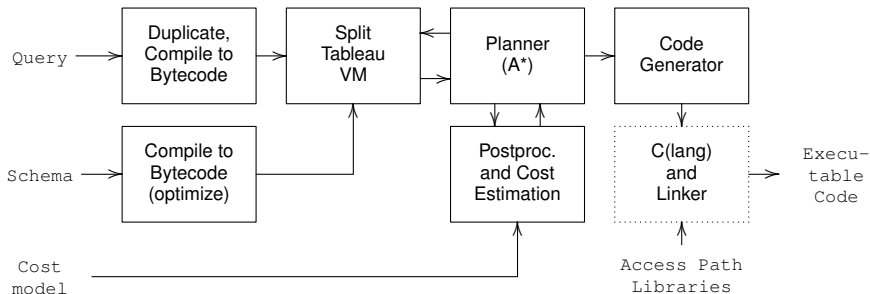
3 separate general reasoning from interpolant enumeration

- 1) VM-driven conditional tableau for $\Sigma^L \cup \{\varphi^L\}$ and for $\Sigma^R \cup \{\varphi^R \rightarrow \perp\}$
- 2) A*-based interpolant generator w.r.t. closing sets and Σ^{LR}

4 FunDL-Lite *postprocessing* to deal w/duplicates *et al.*

Details: [Hudek et al., 2015, Toman and Weddell, 2017]

Compiler Architecture



Summary

Take Home

While in theory *interpolation* essentially solves the *query rewriting over FO schemas/views* problem, **the devil is (as usual) in the details.**

[Borgida, de Bruijn, Franconi, Seylan, Straccia, Toman, Weddell: On Finding Query Rewritings under Expressive Constraints. SEBD 2010: 426-437
... **but an (almost) working system only this year.**

1 FO tableau based interpolation algorithm

- ⇒ enumeration of plans factored from of tableau reasoning
- ⇒ extra-logical binding patterns and cost model

2 Post processing (using $CFDI_{NC}$ approximation)

- ⇒ duplicate elimination
- ⇒ cut insertion

3 Run time

- ⇒ library of common data/legacy structures+schema constraints
- ⇒ finger data structures to simulate merge joins et al.

Research Directions and Open Issues

- 1 Dealing with ordered data? (merge-joins etc.: we have a partial solution)
- 2 Decidable schema languages (decidable interpolation problem)?
- 3 More powerful schema languages (inductive types, etc.)?
- 4 Beyond FO Queries/Views (e.g., count/sum aggregates)?
- 5 Coding extra-logical bits (e.g., **binding patterns**, postprocessing, etc.)
in the schema itself?
- 6 Standard Designs (a plan can always be found as in SQL)?
- 7 Explanation(s) of non-definability?
- 8 Fine(r)-grained updates?
- 9 ...

... and, as always, performance, performance, performance!



Bachman, C. W. (1969).

CODASYL data base task group: October 1969 report.



Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., and Rosati, R. (2007).

Tractable reasoning and efficient query answering in description logics: The DL-Lite family.
J. Autom. Reasoning, 39(3):385–429.



Codd, E. F. (1970).

A relational model of data for large shared data banks.
Commun. ACM, 13(6):377–387.



Codd, E. F. (1972).

Relational completeness of data base sub-languages.
In Rustin, R., editor, *Data Base Systems*, pages 33–64. Prentice-Hall.



Date, C. J. and Hopewell, P. (1971).

Storage structure and physical data independence.
In *Proceedings of the 1971 ACM SIGFIDET (now SIGMOD) Workshop on Data Description, Access and Control*, SIGFIDET '71, pages 139–168. ACM.



Eiter, T., Ortiz, M., Simkus, M., Tran, T., and Xiao, G. (2012).

Query rewriting for horn-shiq plus rules.
In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence, July 22-26, 2012, Toronto, Ontario, Canada*.



Hudek, A. K., Toman, D., and Weddell, G. E. (2015).

On enumerating query plans using analytic tableau.
In *Automated Reasoning with Analytic Tableaux and Related Methods - 24th International Conference, TABLEAUX 2015, Wrocław, Poland, September 21-24, 2015. Proceedings*, pages 339–354.



Kontchakov, R., Lutz, C., Toman, D., Wolter, F., and Zakharyashev, M. (2010).

The combined approach to query answering in DL-Lite.

In *Principles of Knowledge Representation and Reasoning: Proceedings of the Twelfth International Conference, KR 2010, Toronto, Ontario, Canada, May 9-13, 2010*.



Liskov, B. H. and Zilles, S. N. (1974).
Programming with abstract data types.
SIGPLAN Notices, 9(4):50–59.



Lutz, C., Toman, D., and Wolter, F. (2009).
Conjunctive query answering in the description logic EL using a relational database system.
In *Proc. Int. Joint Conf. on Artificial Intelligence (IJCAI)*, pages 2070–2075.



Toman, D. and Weddell, G. E. (2013).
Conjunctive Query Answering in \mathcal{CFD}_{nc} : A PTIME Description Logic with Functional Constraints and Disjointness.
In *Australasian Conference on Artificial Intelligence*, pages 350–361.



Toman, D. and Weddell, G. E. (2017).
An interpolation-based compiler and optimizer for relational queries (system design report).
In *IWIL@LPAR 2017 Workshop and LPAR-21 Short Presentations, Maun, Botswana, May 7-12, 2017*.