# From Data Independence to Ontology Based Data Access (and back)

David Toman

D.R. Cheriton School of Computer Science
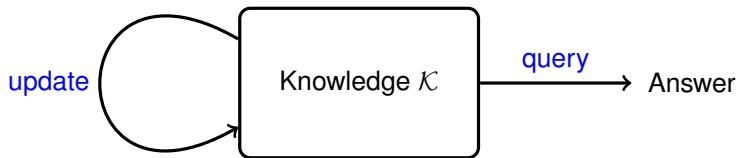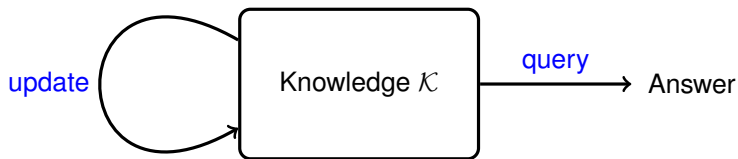
University of

## 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?)
$\Rightarrow$ 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.

Waterloo

# Knowledge Representation: a Big Picture
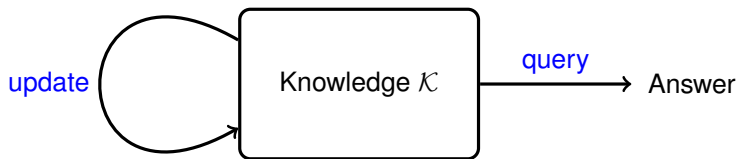


What is "Knowledge" (how is it represented, and does the user care?)

⇒ 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.

Waterloo

David Toman (et al.)    Physical Data Independence    Motivation    2/41

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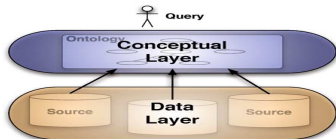


Fig. 1. Ontology-based data access.

# Information Integration  [Genesereth: Data Integration, 2010]



# Data Exchange  [Arenas et el.: Data Exchange, 2014]

The general setting of data exchange is this:

University of Waterloo

# 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: physical
and conceptual/logical [Codd1970]

- data Independence [Bachman, 1969,
  Date and Hopewell, 1971] and
  [Codd, 1970]
- 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: <span style="color:red">physical</span> and <span style="color:blue">conceptual/logical</span> [Codd1970].

data independence [Bachman, 1969,
Date and Hopewell, 1971] and
[Codd, 1970]

ADTs [Liskov and Zilles, 1974]



[ANSI/X3/SPARC Standards Planning and Requirements Committee, Bachman, 1975]
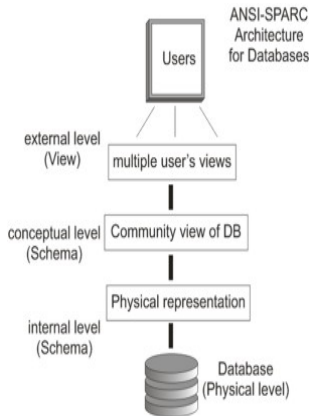
# (Physical) Data Independence

**IDEA:**

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

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

## Physical Data Independence and ADTs

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

ANSI-SPARC Architecture for Databases

Users

external level (View) — multiple user's views

conceptual level (Schema) — Community view of DB

internal level (Schema) — Physical representation

Database (Physical level)

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

# Outline

# QUERIES AND QUERY COMPILATION

# 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)

# 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*):

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

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

# The Logical Way: Complexity

## The Good News

LOGSPACE/PTIME (data complexity) for query answering:

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

# The Logical Way: Complexity

## The Good News

LOGSPACE/PTIME (data complexity) for query answering:

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

## The Bad News

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

⇒ the same goes for Information integration, data exchange, etc.

University of Waterloo

# The Logical Way: Complexity

## The Good News

LOGSPACE/PTIME (data complexity) for query answering:

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

## The Bad News

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

$\Rightarrow$ the same goes for *information integration*, *data exchange*, etc.

# Difficulties: Unintuitive Answers

## Example

- *EMP*(*Sue*)
- *EMP* ⊑ ∃*PHONENUM*    (or ∀*x.EMP*(*x*) → ∃*y.PHONENUM*(*x, y*))

Waterloo

David Toman  (et al.)                    Physical Data Independence                    OBDA Basics          10/41

# 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:

# 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!! (and preferably without having to code algorithms/data structures by hand)
- and performance/data storage-representation/...can all be improved/changed without changing the user queries/updates

# User Queries and Updates – for TODAY

Queries: First-order (open) formulae over the user vocabulary
  ⇒ only *range-restricted* formulae

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) |



$\Sigma_L$  $S_L$  $\leftarrow - - - - - \varphi$    Logical Schema and User Queries

# Rewritability and Definability

## User and System Expectations

| | |
|---|---|
| Queries | range-restricted FOL  over $S_L$ *definable* w.r.t. $\Sigma$ 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) |



Logical Schema
and User Queries

Physical Schema
and Query Plans

[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 $S_L$ *definable* w.r.t. $\Sigma$ 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

*Query* ($S_L$)  $\varphi$

*Compile*  $\psi$ (Query Plan over $S_A$)

*Schema* ($S_L \cup S_P$)  $\Sigma$

*Run*  →  *Answers*

*Data* ($S_A \subseteq S_P$)  (instance of) $S_A$

# 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



```
Query (S_L) ──φ──┐
                 ↓
              Compile ──ψ (Query Plan over S_A)──┐
                 ↑                                ↓
Schema (S_L ∪ S_P) ──Σ──┘                      Run ──→ Answers
                                                ↑
Data (S_A ⊆ S_P) ──(instance of) S_A───────────┘
```

# This is NOT OMQ/OBDA (by example)

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

$S_A = \{\mathrm{emp}/1/0, \mathrm{mgr}/1/0\}$

Waterloo

David Toman (et al.)                Physical Data Independence                Definability/Interpolation        14/41

# This is NOT OMQ/OBDA (by example)

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

$S_A = \{\text{emp}/1/0, \text{mgr}/1/0\}$

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

Certain Answer under OWA: { }
      Answer under CWA: {Wilma}
               (obtained by executing the plan $\{x \mid \text{emp}(x) \wedge \neg\text{mgr}(x)\}$).

Waterloo

David Toman (et al.)      Physical Data Independence      Definability/Interpolation   14/41

# What can we do with this?

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

Waterloo

David Toman (et al.)                    Physical Data Independence                    What can it do?    16/41

# Pointers in Main Memory-Logical Schema

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

this corresponds to

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

Waterloo

David Toman  (et al.)              Physical Data Independence              What can it do?        17/41

# Pointers in Main Memory-Logical Schema

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

this corresponds to

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

additional logical constraints (for example):

- managers are employees that manage a department (a view)
- managers work in their own departnemts (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 ($\Sigma_P \cup \Sigma_{LP}$):

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

# Query Plans that Navigate Pointers

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

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

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 (`employee`($x, y, z$)):

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

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

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

Waterloo

David Toman (et al.)          Physical Data Independence          What can it do?          19/41

# Query Plans that Navigate Pointers

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

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

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

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

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

   ≫ Caveat: we do NOT have a (direct) way to "scan" $\texttt{depatment}$s! ≪

# Query Plans that Navigate Pointers

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

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

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

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

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

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

$\Rightarrow$ needs "departments have at least one employee".

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

$\Rightarrow$ needs "managers work in their own departments".

Waterloo

David Toman (et al.)     Physical Data Independence     What can it do?     19/41

# Query Plans that Navigate Pointers

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

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

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

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

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

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

$\Rightarrow$ needs "departments have at least one employee".

. . . needs *duplicate elimination* during projection.

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

$\Rightarrow$ needs "managers work in their own departments".

. . . NO *duplicate elimination* during projection.

Waterloo

# ...and we can actually synthesize this!

```
david$ compile tests/new_fe/book-em-v4-new-query.fol
query(q0dept2,2,0,[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



Hash Array     Separate Chaining Linked Lists     Dept Records

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

## Hash Index on `department`'s `name`:

Access paths:

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

Physical Constraints:

$\Sigma_{LP} \supseteq \{ \forall x, y.((\texttt{deptfile}(x) \wedge \texttt{dept-name}(x, y)) \rightarrow \exists z, w.(\texttt{hash}(y, z)$
$\wedge \texttt{hasharraylookup}(z, w) \wedge \texttt{listscan}(w, x))),$
$\forall x, y.(\texttt{hash}(x, y) \rightarrow \exists z.\texttt{hasharraylookup}(y, z)),$
$\forall x, y.(\texttt{listscan}(x, y) \rightarrow \texttt{deptfile}(y))$ $\}$

University of Waterloo

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

**Hash Index on `department`'s `name`:**

Access paths:

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

Physical Constraints:

$\Sigma_{LP} \supseteq \{\forall x, y.((\texttt{deptfile}(x) \wedge \texttt{dept-name}(x, y)) \rightarrow \exists z, w.(\texttt{hash}(y, z)$
$\wedge \texttt{hasharraylookup}(z, w) \wedge \texttt{listscan}(w, x))),$
$\forall x, y.(\texttt{hash}(x, y) \rightarrow \exists z.\texttt{hasharraylookup}(y, z)),$
$\forall x, y.(\texttt{listscan}(x, y) \rightarrow \texttt{deptfile}(y)) \qquad \}$

Query:

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

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

# 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
($\exists z$.`employee`$(x_1, z) \wedge$ `employee`$(x_2, z)$):

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

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

...more examples in

# UPDATES

University of Waterloo

David Toman (et al.)                     Physical Data Independence                     What can it do?         24 / 41

# Updates

# Updates



Knowledge

update

User Vocabulary

query → Answer

Data Repositories ↔ Raw Data Vocabulary

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.

University of Waterloo

David Toman (et al.)    Physical Data Independence    What can it do?    25/41

# Updates

Knowledge



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 follow a *definable updates* approach here instead...

# Updates and Definability

$\Rightarrow$ supplying "delta" relations (sets of tuples)

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

# Updates and Definability

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



## Update turned into *definability* question

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

# Unknown/Anonymous Values?

## Example (Add a new Undergraduate student)

```
INSERT into undergrad values (1234, 'Wilma');
```

$\Rightarrow$ the request then needs to be translated to

```
    INSERT into student values (0xFE1234, 1234, 'Wilma');
```

$\Rightarrow$ but where did `0xFE1234` came from? (definability issue!)

# Unknown/Anonymous Values?

## Example (Add a new Undergraduate student)

```
INSERT into undergrad values (1234, 'Wilma');
```
$\Rightarrow$ the request then needs to be translated to
```
     INSERT into student values (0xFE1234, 1234, 'Wilma');
```
$\Rightarrow$ but where did `0xFE1234` came from? (definability issue!)

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

additional access paths that *provide* such values:

$\Rightarrow$ in our case `student-addr(id,adress)`
$\Rightarrow$ and where `undergrad`$^+$ = $\{(1234, \text{Vilma})\}$

$\texttt{student}^+(x_1, x_2, x_3) = \texttt{undergrad}^+(x_1, x_3) \land \texttt{student-addr}(x_2, x_1)$

# Unknown/Anonymous Values?

## Example (Add a new Undergraduate student)

```
INSERT into undergrad values (1234, 'Wilma');
```
$\Rightarrow$ the request then needs to be translated to
```
    INSERT into student values (0xFE1234, 1234, 'Wilma');
```
$\Rightarrow$ but where did `0xFE1234` came from? (definability issue!)

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

additional access paths that *provide* such values:

$\Rightarrow$ in our case `student-addr(id,adress)`
$\Rightarrow$ and where $\texttt{undergrad}^+ = \{(1234, \textit{Vilma})\}$

$$\texttt{student}^+(x_1, x_2, x_3) = \texttt{undergrad}^+(x_1, x_3) \wedge \texttt{student-addr}(x_2, x_1)$$

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

| Queries | range-restricted FOL over $S_L$ *definable* w.r.t. $\Sigma$ and $S_A$ |
|---|---|
| Ontology/Schema | range-restricted FOL |
| Data | CWA (complete information for $S_A$ symbols) |



$\Sigma_L$    $S_L$ $\longleftarrow$ — — — — $\varphi$    (Logical Schema)

$\Sigma_{LP}$     (*rewriting*)

$\Sigma_P$    $S_A \subseteq S_P$ $\longleftarrow$ — — — — $\psi$    (Physical Schema)

# Query Plans via Interpolation

## IDEA #1: Plans as Formulas

Represent *query plans* as (annotated) range-restricted formulas $\psi$ over $S_A$:

| | | |
|---|---|---|
| atomic formula | $\mapsto$ | access path (`get-first`–`get-next` iterator) |
| 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 $\psi$ over $S_A$:

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

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

University of Waterloo

David Toman (et al.)    Physical Data Independence    How does it work?    30/41

# Query Plans via Interpolation

## IDEA #1: Plans as Formulas

Represent *query plans* as (annotated) range-restricted formulas $\psi$ over $S_A$:

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

$\Rightarrow$ reduces correctness of $\psi$ 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

University of Waterloo

David Toman (et al.)     Physical Data Independence     How does it work?     30/41

# 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 $\Sigma$
for a fixed interpretation of the signature $S_A$ (i.e., where the actual data is).

Waterloo
University of

David Toman (et al.)                     Physical Data Independence                     How does it work?          31/41

# 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 $\Sigma$
for a fixed interpretation of the signature $S_A$ (i.e., where the actual data is).

## How do we test for this?

$\varphi$ is *Beth definable* [Beth'56] if

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

where $\Sigma'$ ($\varphi'$) is $\Sigma$ ($\varphi$) in which symbols *NOT in* $S_A$ are *primed*, respectively.

University of Waterloo

# 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 $\Sigma$
  for a fixed interpretation of the signature $S_A$ (i.e., where the actual data is).

## How do we test for this?

$\varphi$ is *Beth definable* [Beth'56] if

$$\Sigma \cup \Sigma' \models \varphi \to \varphi'$$

where $\Sigma'$ ($\varphi'$) is $\Sigma$ ($\varphi$) in which symbols *NOT in* $S_A$ are *primed*, respectively.

## How do we find $\psi$?

If $\Sigma \cup \Sigma' \models \varphi \to \varphi'$ then there is $\psi$ s.t. $\Sigma \cup \Sigma' \models \varphi \to \psi \to \varphi'$ with $\mathcal{L}(\psi) \subseteq \mathcal{L}(S_A)$.

  ... $\psi$ is called the *Craig Interpolant* [Craig'57].

  ... we extract an *interpolant* $\psi$ from a (TABLEAU) proof of $\Sigma \cup \Sigma' \models \varphi \to \varphi'$

# 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

Waterloo

David Toman (et al.)                Physical Data Independence                How does it work?        32/41

# 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 Tableau through (a simple) Example

## Example Schema

Rules: $q(x) \rightarrow a(x)$, $q(x) \rightarrow c(x)$, $a(x) \wedge c(x) \rightarrow q(x)$, $c(x) \rightarrow \exists y.b(x,y)$

# Conditional Tableau through (a simple) Example

## Example Schema

Rules: $q(x) \rightarrow a(x)$, $q(x) \rightarrow c(x)$, $a(x) \wedge c(x) \rightarrow q(x)$, $c(x) \rightarrow \exists y.b(x,y)$

## Conditional Tableau

| $\Sigma^L$ | | $\Sigma^R$ |
|:---:|:---:|:---:|
| $q(0)[]\{\}$ | | |
| $a(0)[]\{\}$ | $\rightarrow$ | $a(0)[a(0)]\{\}$ |
| $c(0)[]\{\}$ | $\rightarrow$ | $c(0)[c(0)]\{\}$ |
| $b(0,1)[]\{\}$ | | $q(0)[a(0),c(0)]\{\}$ |
| | | $\perp[a(0),c(0)]\{\}$ |

# Conditional Tableau through (a simple) Example

## Example Schema

Rules: $q(x) \rightarrow a(x)$, $q(x) \rightarrow c(x)$, $a(x) \wedge c(x) \rightarrow q(x)$, $c(x) \rightarrow \exists y.b(x,y)$

## Conditional Tableau

| $\Sigma^L$ | | $\Sigma^R$ |
|---|---|---|
| $q(0)[]\{\}$ | | |
| $a(0)[]\{\}$ | $\rightarrow$ | $a(0)[a(0)]\{\}$ |
| $c(0)[]\{\}$ | $\rightarrow$ | $c(0)[c(0)]\{\}$ |
| $b(0,1)[]\{\}$ | | $q(0)[a(0),c(0)]\{\}$ |
| | | $\bot[a(0),c(0)]\{\}$ |

$$a^L(\bar{x}) \rightarrow a(\bar{x}) \rightarrow a^R(\bar{x})$$

# Conditional Tableau through (a simple) Example

## Example Schema

Rules: $q(x) \rightarrow a(x)$, $q(x) \rightarrow c(x)$, $a(x) \wedge c(x) \rightarrow q(x)$, $c(x) \rightarrow \exists y.b(x,y)$

## Conditional Tableau

| $\Sigma^L$ | | $\Sigma^R$ |
|---|---|---|
| $q(0)[]\{\}$ | | |
| $a(0)[]\{\}$ | $\rightarrow$ | $a(0)[a(0)]\{\}$ |
| $c(0)[]\{\}$ | $\rightarrow$ | $c(0)[c(0)]\{\}$ |
| $b(0,1)[]\{\}$ | | $q(0)[a(0),c(0)]\{\}$ |
| | | $\perp[a(0),c(0)]\{\}$ |

Closing Sets for $\Sigma^L$: $\{\neg a(0)\}$, $\{\neg c(0)\}$, $\{\neg b(0,1)\}$     Only atoms that are in

Closing Sets for $\Sigma^R$: $\{a(0),c(0)\}$     closing sets appear in query plans!

University of Waterloo

David Toman (et al.)     Physical Data Independence     How does it work?     33/41

# Conditional Tableau through (a simple) Example

## Example Schema

Rules: $q(x) \to a(x)$, $q(x) \to c(x)$, $a(x) \land c(x) \to q(x)$, $c(x) \to \exists y.b(x,y)$

## Conditional Tableau

| $\Sigma^L$ | | $\Sigma^R$ |
|---|---|---|
| $q(0)[]\{\}$ | | |
| $a(0)[]\{\}$ | $\to$ | $a(0)[a(0)]\{\}$ |
| $c(0)[]\{\}$ | $\to$ | $c(0)[c(0)]\{\}$ |
| $b(0,1)[]\{\}$ | | $q(0)[a(0),c(0)]\{\}$ |
| | | $\perp[a(0),c(0)]\{\}$ |

Closing Sets for $\Sigma^L$: $\{\neg a(0)\}, \{\neg c(0)\}, \{\neg b(0,1)\}$    Only atoms that are in
Closing Sets for $\Sigma^R$: $\{a(0), c(0)\}$    closing sets appear in query plans!

## Plans

Query plan 1: $a(x) \land c(x)$
Query plan 2: $a(x) \land (\exists y.b(x,y)) \land c(x)$

# CONDITIONAL TABLEAU AND CLOSING SETS

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(sl19:7,sl14:3,sl0:2,sl0:2) }
   L { -p0basetable(sl19:5,sl0:1,sl0:1,sl14:3) }
   L { +p0basetable(sr19:8,sl0:1,sl0:1,sl0:1) }
   R { -p0basetable(sr19:8,sl0:1,sl0:1,sl0:1),
       +p0basetable(sl19:7,sl14:3,sl0:2,sl0:2),
       +p0basetable(sl19:5,sl0:1,sl0:1,sl14:3) }
   ```

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.]

University of Waterloo

David Toman  (et al.)          Physical Data Independence          How does it work?          34/41

# 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 ($\{\cdot\}$).

Waterloo

David Toman  (et al.)      Physical Data Independence      How does it work?      36 / 41

# 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 $\{\cdot\}$:

$$Q[\{R(x_1, \ldots, x_k)\}] \leftrightarrow Q[R(x_1, \ldots, x_k)]$$
$$Q[\{Q_1 \wedge Q_2\}] \leftrightarrow Q[\{Q_1\} \wedge \{Q_2\}]$$
$$Q[\{\neg Q_1\}] \leftrightarrow Q[\neg Q_1]$$
$$Q[\neg \{Q_1\}] \leftrightarrow Q[\neg Q_1]$$
$$Q[\{Q_1 \vee Q_2\}] \leftrightarrow Q[\{Q_1\} \vee \{Q_2\}] \quad \text{if } \Sigma \cup \{Q[]\} \models Q_1 \wedge Q_2 \rightarrow \bot$$
$$Q[\{\exists x.Q_1\}] \leftrightarrow Q[\exists x.\{Q_1\}] \quad \text{if}$$
$$\Sigma \cup \{Q[] \wedge (Q_1)[y_1/x] \wedge (Q_1)[y_2/x] \models y_1 \approx y_2$$

# 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 $\{\cdot\}$:

$$Q[\{R(x_1, \ldots, x_k)\}] \leftrightarrow Q[R(x_1, \ldots, x_k)]$$
$$Q[\{Q_1 \wedge Q_2\}] \leftrightarrow Q[\{Q_1\} \wedge \{Q_2\}]$$
$$Q[\{\neg Q_1\}] \leftrightarrow Q[\neg Q_1]$$
$$Q[\neg\{Q_1\}] \leftrightarrow Q[\neg Q_1]$$
$$Q[\{Q_1 \vee Q_2\}] \leftrightarrow Q[\{Q_1\} \vee \{Q_2\}] \quad \text{if } \Sigma \cup \{Q[]\} \models Q_1 \wedge Q_2 \rightarrow \bot$$
$$Q[\{\exists x.Q_1\}] \leftrightarrow Q[\exists x.\{Q_1\}] \quad \text{if}$$
$$\Sigma \cup \{Q[] \wedge (Q_1)[y_1/x] \wedge (Q_1)[y_2/x] \models y_1 \approx y_2$$

… reasoning abstracted: a DL $\mathcal{CFD}_{nc}^{\forall-}$ (a PTIME fragment)

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

# Interpolation (in practice)

## Difficulties with naive implementation/Obstacles

1. Structural properties of proofs (e.g., subformula property)
2. Alternative interpolants/plans (can we just *backtrack* the proof system?)

Waterloo

David Toman (et al.)    Physical Data Independence    How does it work?    37/41

# Interpolation (in practice)
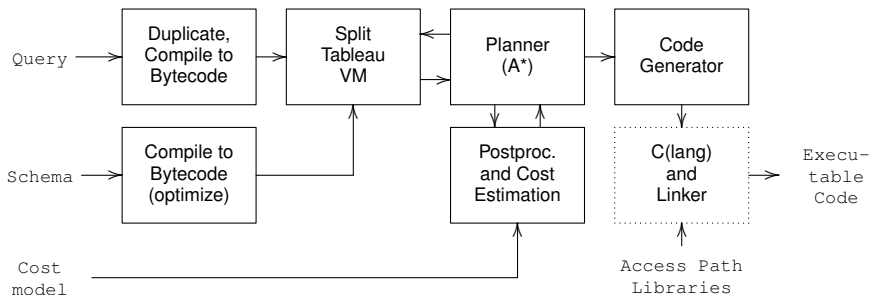
## Difficulties with naive implementation/Obstacles

1. Structural properties of proofs (e.g., subformula property)
2. Alternative interpolants/plans (can we just *backtrack* the proof system?)

(A) solution: a *conditional tableau*:

1. reformulate the interpolation problem to
   $\Sigma^L \cup \Sigma^R \cup \Sigma^{LR} \models \varphi^L \to \varphi^R$ where $\Sigma^{LR} = \{\forall \bar{x}.P^L \leftrightarrow P \leftrightarrow P^R \mid P \in S_A\}$
2. use conditional (ground) atoms and
   closing sets: sets of $S_A$ literals that (fully) close a tableaau
3. separate general reasoning from interpolant enumeration
   VM-driven conditional tableau for $\Sigma^L \cup \{\varphi^L\}$ and for $\Sigma^R \cup \{\varphi^R \to \bot\}$
   A$^*$-based interpolant generator w.r.t. closing sets and $\Sigma^{LR}$

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

# Compiler Architecture

Waterloo

David Toman (et al.)    Physical Data Independence    How does it work?    38/41

# 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 $\mathcal{CFDI}_{nc}$ approximation)
   - ⇒ duplicate elimination 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.

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.

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.

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

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*.

University of Waterloo

David Toman (et al.)  Physical Data Independence  Summary  40 / 41