# KR in Database Systems Implementation
## (or Life beyond Lite Logics and CQ/UCQ)

David Toman

D.R. Cheriton School of Computer Science
University of
**Waterloo**

Joint work with Alexander Hudek and Grant Weddell

# Data and Constraints: the Database Recap

## The Textbook View

| | |
|---|---|
| Data | represented as an instance of a *relational structure* |
| Queries | access to data via *open formulæ* (in an appropriate logic) |
| Constraints | data integrity enforces by *sentences* (in the same logic) |

$\Rightarrow$ the instance is a *model* of the constraints

## What about CREATE VIEW Statements?

View declaration ∼ a sentence $\forall \mathbf{x}. V(\mathbf{x}) \leftrightarrow \varphi$ (in our logic)
where $V$ is a (new) relational symbol and $\varphi$ is a *query*.

## Much Bigger Deal: Physical Data Independence

| | |
|---|---|
| Logical Symbols | user (visible) relations/tables |
| Mapping | |
| Physical Symbols | data structures (indices) |

# Data and Constraints: the Database Recap

## The Textbook View

| | |
|---|---|
| Data | represented as an instance of a *relational structure* |
| Queries | access to data via *open formulæ* (in an appropriate logic) |
| Constraints | data integrity enforces by *sentences* (in the same logic) |

$\Rightarrow$ the instance is a *model* of the constraints

## What about `CREATE VIEW` Statements?

View declaration $\sim$ a sentence $\forall \mathbf{x}.V(\mathbf{x}) \leftrightarrow \varphi$ (in our logic)
where *V* is a (new) relational symbol and $\varphi$ is a *query*.

## Much Bigger Deal: Physical Data Independence

| | |
|---|---|
| Logical Symbols | user (visible) relations/tables |
| Mapping | |
| Physical Symbols | data structures (indices) |

# Data and Constraints: the Database Recap

## The Textbook View

| Data | represented as an instance of a *relational structure* |
|------|-------------------------------------------------------|
| Queries | access to data via *open formulæ* (in an appropriate logic) |
| Constraints | data integrity enforces by *sentences* (in the same logic) |

$\Rightarrow$ the instance is a *model* of the constraints

## What about `CREATE VIEW` Statements?

View declaration $\sim$ a sentence $\forall \mathbf{x}. V(\mathbf{x}) \leftrightarrow \varphi$ (in our logic)
where *V* is a (new) relational symbol and $\varphi$ is a *query*.

## Much Bigger Deal: Physical Data Independence

| Logical Symbols | user (visible) relations/tables |
|-----------------|--------------------------------|
| Mapping | |
| Physical Symbols | data structures (indices) |

University of Waterloo

David Toman (et al.)  KR in DBMS Implementation  Motivation  2 / 18

# Data and Constraints: the Database Recap

## The Textbook View

| | |
|---|---|
| Data | represented as an instance of a *relational structure* |
| Queries | access to data via *open formulæ* (in an appropriate logic) |
| Constraints | data integrity enforces by *sentences* (in the same logic) |

$\Rightarrow$ the instance is a *model* of the constraints

## What about `CREATE VIEW` Statements?

View declaration $\sim$ a sentence $\forall \mathbf{x}. V(\mathbf{x}) \leftrightarrow \varphi$ (in our logic)
where $V$ is a (new) relational symbol and $\varphi$ is a *query*.

## Much Bigger Deal: Physical Data Independence

| | |
|---|---|
| Logical Symbols | user (visible) relations/tables |
| Mapping | C/C++ goo |
| Physical Symbols | data structures (indices) |

# Data and Constraints: the Database Recap

## The Textbook View

Data      represented as an instance of a *relational structure*
Queries      access to data via *open formulæ* (in an appropriate logic)
Constraints      data integrity enforces by *sentences* (in the same logic)

$\Rightarrow$ the instance is a *model* of the constraints

## What about `CREATE VIEW` Statements?

View declaration $\sim$ a sentence $\forall \mathbf{x}. V(\mathbf{x}) \leftrightarrow \varphi$ (in our logic)
where $V$ is a (new) relational symbol and $\varphi$ is a *query*.

## Much Bigger Deal: Physical Data Independence

Logical Symbols      user (visible) relations/tables
Mapping      constraints (+ a minimal runtime)
Physical Symbols      data structures (indices)

# The KR Way

## Queries and Ontologies

Queries are answered not only w.r.t. *explicit data* ($\mathcal{A}$)
but also w.r.t. *background knowledge* ($\mathcal{T}$) under OWA

$\Rightarrow$ Ontology-based Data Access (OBDA)

## Example

- Socrates is a MAN                                          (explicit data)
- Every MAN is MORTAL                                        (ontology)

*List all MORTALs* $\Rightarrow$ {Socrtes}                  (query)

How do we answer queries?

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

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

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

University of Waterloo

# The KR Way

## Queries and Ontologies

Queries are answered not only w.r.t. *explicit data* ($\mathcal{A}$)
but also w.r.t. *background knowledge* ($\mathcal{T}$) under OWA

$\Rightarrow$ Ontology-based Data Access (OBDA)

## Example

- Socrates is a MAN (explicit data)
- Every MAN is MORTAL (ontology)

*List all MORTALs* $\Rightarrow$ {Socrtes} (query)

## How do we answer queries?

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

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

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

# Complexity

## Good/Standard News

LOGSPACE/PTIME (data complexity) for query answering:

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

## Bad News

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

# Complexity

## Good/Standard News

LOGSPACE/PTIME (data complexity) for query answering:

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

## Bad News

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

University of Waterloo

# 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 $\forall x.EMP(x) \rightarrow \exists y.PHONENUM(x, y)$)
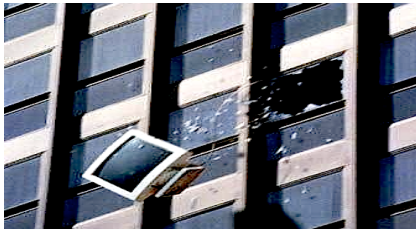
User: *Does Sue have a phone number?*

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:

# What to do?

## Definability and Rewriting

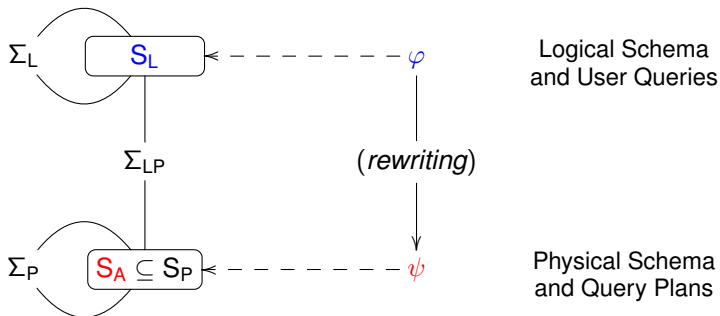| | |
|---|---|
| Queries | range-restricted FOL (a.k.a. SQL) |
| Ontology/Schema | range-restricted FOL $\Sigma := \Sigma^L \cup \Sigma^{LP} \cup \Sigma^P$ |
| Data | CWA (complete information) |

Waterloo

David Toman (et al.)          KR in DBMS Implementation          Definability/Interpolation    6 / 18

# What to do?

## Definability and Rewriting

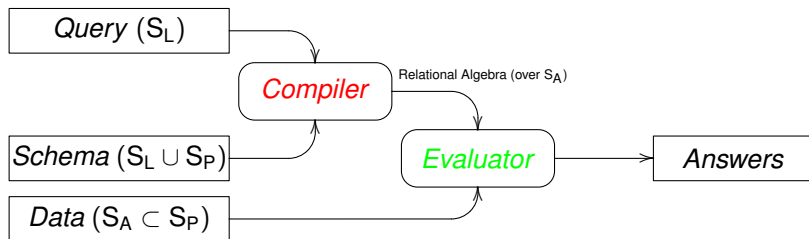| | |
|---|---|
| Queries | range-restricted FOL over $S_L$ *definable* w.r.t. $\Sigma$ and $S_A$ |
| Ontology/Schema | range-restricted FOL $\Sigma := \Sigma^L \cup \Sigma^{LP} \cup \Sigma^P$ |
| Data | CWA (complete information for $S_A$ symbols) |



$\Sigma_L$ — $S_L$ ⟵ — — — — — $\varphi$ — Logical Schema and User Queries

$\Sigma_{LP}$ — (*rewriting*)

$\Sigma_P$ — $S_A \subseteq S_P$ ⟵ — — — — — $\psi$ — Physical Schema and Query Plans

# What to do?

## Definability and Rewriting

| | |
|---|---|
| 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) |

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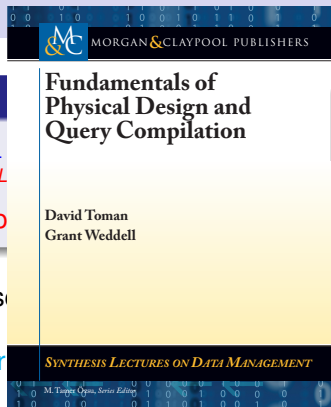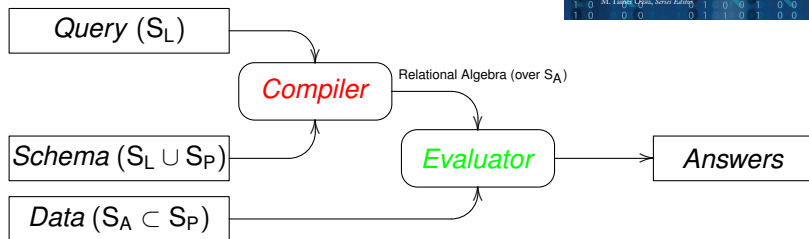  but *definable* queries answer the same in each of them

Waterloo

# What to do?

## Definability and Rewriting

| Queries | range-restricted FOL over $S_L$ |
|---|---|
| Ontology/Schema | range-restricted FOL $\Sigma := \Sigma^L$ |
| Data | CWA (complete information fo... |

- users: looks like a *single model* (of the logical s...
- implementation: many models

  but *definable* queries answer...



©2011

Fundamentals of
Physical Design and
Query Compilation

David Toman
Grant Weddell

MORGAN & CLAYPOOL PUBLISHERS

# GRAND UNIFIED APPROACH TO QUERY COMPILATION

## PART I: WHAT CAN IT DO?

# What can this do?

## GOAL

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

1. linked data structures, pointers, . . .
2. access to search structures (index access and selection),
3. hash-based access to data (including hash-joins),
4. multi-level storage (aka disk/remote/distributed files), . . .
5. materialized views (FO-definable),
6. updates through logical schema (needs *id invention*!), . . .

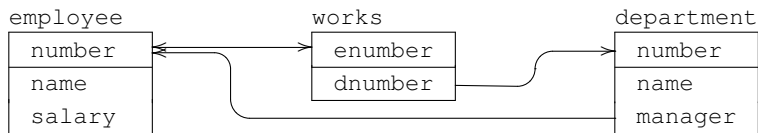. . . all without having to code (too much) in C/C++ !

# What can this do?

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

1. linked data structures, pointers, . . .
2. access to search structures (index access and selection),
3. hash-based access to data (including hash-joins),
4. multi-level storage (aka disk/remote/distributed files), . . .
5. materialized views (FO-definable),
6. updates through logical schema (needs *id invention*!), . . .

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

# Lists and Pointers

1. Logical Schema



2. Physical Design: a *linked list of* `emp` *records pointing to* `dept` *records*.

```
record emp of                    record dept of
        integer    num                   integer    num
        string     name                  string     name
        integer    salary                reference  manager
        reference  dept
```

3. Access Paths: `empfile`/1/0, `emp-num`/2/1, . . . (but no `deptfile`)

4. Integrity Constraints (many), e.g.,

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

# What can this do: navigating pointers

Example queries:

1. List all employee numbers and names ($\exists z, w.\mathtt{employee}(x, y, z, w)$):

$$\exists a.\mathtt{empfile}(a) \wedge \mathtt{emp\text{-}num}(a, x) \wedge \mathtt{emp\text{-}name}(a, y)$$

Waterloo

David Toman (et al.)      KR in DBMS Implementation      What can it do?    10 / 18

# What can this do: navigating pointers

Example queries:

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

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

2. List all department numbers with their manager names
   ($\exists z, u, v, w.\texttt{department}(x, z, u) \land \texttt{employee}(u, y, v, w)$):

Waterloo

David Toman (et al.)                KR in DBMS Implementation                What can it do?      10 / 18

# What can this do: navigating pointers

Example queries:

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

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

2. List all department numbers with their manager names
   ($\exists z, u, v, w.\texttt{department}(x, z, u) \wedge \texttt{employee}(u, y, v, w)$):

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

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

Waterloo

David Toman (et al.)        KR in DBMS Implementation                What can it do?    10 / 18

# What can this do: navigating pointers

Example queries:

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

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

2. List all department numbers with their manager names
   ($\exists z, u, v, w.\texttt{department}(x, z, u) \wedge \texttt{employee}(u, y, v, w)$):

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

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

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

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

Waterloo

David Toman (et al.)          KR in DBMS Implementation          What can it do?     10 / 18

# What can this do: navigating pointers

Example queries:

1. List all employee numbers and names ($\exists z, w.\text{employee}(x, y, z, w)$):

$$\exists a.\text{empfile}(a) \land \text{emp-num}(a, x) \land \text{emp-name}(a, y)$$

2. List all department numbers with their manager names
   ($\exists z, u, v, w.\text{department}(x, z, u) \land \text{employee}(u, y, v, w)$):

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

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

. . . needs *duplicate elimination* during projection.

$$\exists a, b, d.\text{empfile}(a) \land \text{emp-name}(a, y) \land \text{emp-dept}(a, d)$$
$$\land \text{dept-num}(d, x) \land \text{dept-mgr}(d, b) \land \text{compare}(a, b)$$

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

. . . NO *duplicate elimination* during projection.

Waterloo

David Toman (et al.)          KR in DBMS Implementation          What can it do?     10 / 18

# 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
$(\exists z, u, v, w, t.\texttt{employee}(x_1, z, u, v) \wedge \texttt{employee}(x_2, z, w, t))$:

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

⇒ this plan implements the *block nested loops join* algorithm.
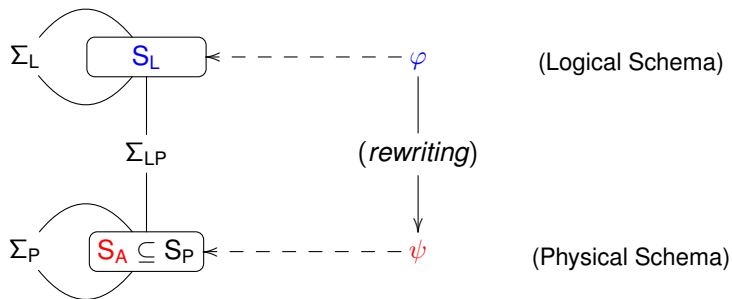
. . . many more examples in

# GRAND UNIFIED APPROACH TO QUERY COMPILATION

## PART II: HOW DOES IT WORK?

Waterloo

David Toman (et al.)                    KR in DBMS Implementation                    12 / 18

# The Plan

## Definability and Rewriting

| | |
|---|---|
| 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) |



$\Sigma_L$ — $S_L$ ←------- $\varphi$ (Logical Schema)

$\Sigma_{LP}$

(*rewriting*)

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

# Query Plans via Interpolation

## IDEA #1:

Represent *physical design* as *access paths* ($S_A$) and constraints ($\Sigma$).
Represent *query plans* as (annotated) range-restricted formulas $\psi$ over $S_A$.

| | | |
|---|---|---|
| atomic formula | $\mapsto$ | access path |
| 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:

Represent *physical design* as *access paths* ($S_A$) and constraints ($\Sigma$).
Represent *query plans* as (annotated) range-restricted formulas $\psi$ over $S_A$.

$\Rightarrow$ reduces correctness of $\psi$ w.r.t. the user query $\varphi$ to $\Sigma \models \varphi \leftrightarrow \psi$

# Query Plans via Interpolation

## IDEA #1:

Represent *physical design* as *access paths* ($S_A$) and constraints ($\Sigma$).
Represent *query plans* as (annotated) range-restricted formulas $\psi$ over $S_A$.

$\Rightarrow$ reduces correctness of $\psi$ w.r.t. the user query $\varphi$ to $\Sigma \models \varphi \leftrightarrow \psi$

## IDEA #2:

Use *interpolation* to search for $\psi$:

extract an *interpolant* $\psi$ from a (TABLEAU) proof of $\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^*$

# Query Plans via Interpolation

### IDEA #1:

Represent *physical design* as *access paths* ($S_A$) and constraints ($\Sigma$).
Represent *query plans* as (annotated) range-restricted formulas $\psi$ over $S_A$.

$\Rightarrow$ reduces correctness of $\psi$ w.r.t. the user query $\varphi$ to $\Sigma \models \varphi \leftrightarrow \psi$

### IDEA #2:

Use *interpolation* to search for $\psi$:
   extract an *interpolant* $\psi$ from a (TABLEAU) proof of $\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^*$

$\Rightarrow$ *Beth Definability* of $\varphi$ over $\Sigma$ and $S_A$ resolves the existence of $\psi$
(except for *binding patterns*)

# Engineering Issues

Subformula (structural) Property: not enough rewritings (plans)

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

... for details see

# Engineering Issues

Subformula (structural) Property: not enough rewritings (plans)

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

Alternative Proofs/Plans: backtracking is too slow

- *conditional formulæ*: $\varphi[C]$ where $C$ is a set of (ground) literals over $S_A$
- logical (non-backtrackable) *conditional tableau* ($T^L, T^R$)
- cost-based plan enumeration based on *closing sets* in ($T^L, T^R$) and $\Sigma^{LR}$

... for details see

# Engineering Issues

**Subformula (structural) Property:** not enough rewritings (plans)

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

**Alternative Proofs/Plans:** backtracking is too slow

- *conditional formulæ*: $\varphi[C]$ where $C$ is a set of (ground) literals over $S_A$
- logical (non-backtrackable) *conditional tableau* ($T^L$, $T^R$)
- cost-based plan enumeration based on *closing sets* in ($T^L$, $T^R$) and $\Sigma^{LR}$

**Non-logical Features:** dealing with duplicates et al.

- $Q[\exists x.Q_1] \mapsto Q[\exists x.Q_1]$ if $\Sigma \cup \{Q[] \wedge Q_1[y_1/x] \wedge Q_1[y_2/x]\} \models y_1 \approx y_2$
- $Q[Q_1 \vee Q_2] \mapsto Q[Q_1 \vee Q_2]$ if $\Sigma \cup \{Q[]\} \models Q_1 \wedge Q_2 \to \bot$

  $\Rightarrow \mathcal{CFDI}_{\text{nc}}$ description logic approximation of $\Sigma$ (PTIME reasoning).

... for details see .

# Summary of the Approach

1. FO ($\mathcal{DLFDE}$) tableau based interpolation algorithm
   - $\Rightarrow$ enumeration of plans factored from reasoning
   - $\Rightarrow$ range-restricted queries and constraints $\rightarrow$ ground terms only
   - $\Rightarrow$ extra-logical binding patterns and cost model

2. Post processing (using $\mathcal{CFDI}_{nc}$ approximation)
   - $\Rightarrow$ duplicate elimination elimination
   - $\Rightarrow$ cut insertion

3. Run time
   - $\Rightarrow$ library of common data structures+schema constraints

     or an interface to a legacy system
   - $\Rightarrow$ finger data structures to simulate merge joins et al.

Waterloo

David Toman (et al.)     KR in DBMS Implementation     Summary     16 / 18

# 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!

# Message from our Sponsors

## Database Group at the University of Waterloo

- 7 professors, affiliated faculty, postdocs, 30+ graduate students, . . .
- wide range of research interests
  - Advanced query processing/Knowledge representation
  - System aspects of database systems and Distributed data management
  - Data quality/Managing uncertain data/Data mining
  - New(-ish) domains (text, streaming, graph data/RDF, OLAP)
- research sponsored by governments, and local/global companies

  NSERC/CFI/OIT and Google, IBM, SAP, OpenText, . . .
- part of a School of CS with 75+ professors, 300+ grad students, etc.

  AI&ML, Algorithms&Data Structures, PL, Theory, Systems, . . .

  Cheriton School of Computer Science has been ranked #18 in CS by
  the world by *US News and World Report* (#1 in Canada).

. . . and we are always looking for good graduate students (MMath/PhD)

⇒ comes with full support over multiple years