

Logical Approach to Physical Data Independence and Query Compilation

Query Rewriting

David Toman

D.R. Cheriton School of Computer Science

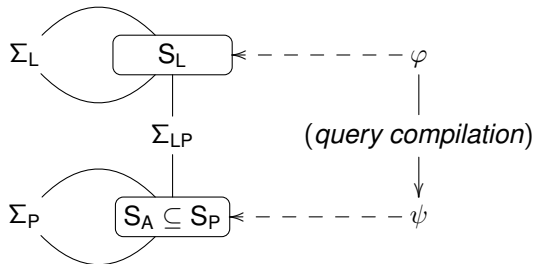
University of

Waterloo

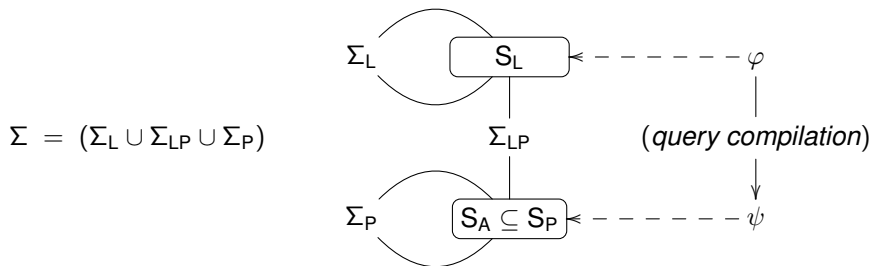


The Story So Far...

$$\Sigma = (\Sigma_L \cup \Sigma_{LP} \cup \Sigma_P)$$



The Story So Far...

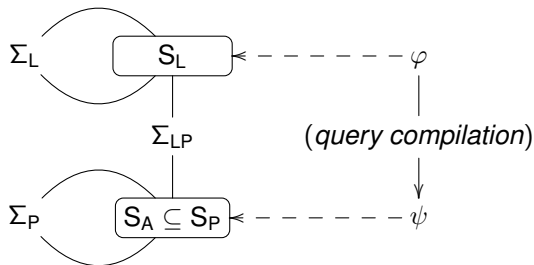


Features:

- Flexible *physical design*: constraints $\Sigma_P \cup \Sigma_{LP}$ and code for S_A
⇒ main-memory operations, disk access, external sources of data, ... ;
- Query plans are efficient
⇒ all combination of *access paths* and simple operators;
⇒ often comparable to hand-written programs.

The Story So Far...

$$\Sigma = (\Sigma_L \cup \Sigma_{LP} \cup \Sigma_P)$$



- 1 How do we find ψ such that $\psi \in \mathcal{L}(S_A)$ and $\Sigma \models \varphi \leftrightarrow \psi$?
- 2 How do we deal with non-logical issues (e.g., duplicates)?

Goal and Steps

- 1 Find ψ such that $\psi \in \mathcal{L}(S_A)$ and $\Sigma \models \varphi \leftrightarrow \psi$?
 - search for *optimal* ψ (according to a *cost model*)
 - in general many *candidates* (even for CQ: join-order optimization)

Goal and Steps

- 1 Find ψ such that $\psi \in \mathcal{L}(S_A)$ and $\Sigma \models \varphi \leftrightarrow \psi$?
 - search for *optimal* ψ (according to a *cost model*)
 - in general many *candidates* (even for CQ: join-order optimization)
- 2 How do we deal with non-logical issues?
 - elimination of unnecessary *duplicate elimination* operations
 - *cut* insertion (when one solution suffices)

QUERY REWRITING

Chase and Backchase

- Input: φ a *CQ*, Σ a set of *dependencies*, and S_A .
 \Rightarrow a *dependency* is a formula $\forall \bar{x}. \alpha \rightarrow \beta$ where α and β are CQs.

Chase and Backchase

- Input: φ a **CQ**, Σ a set of **dependencies**, and S_A .
 \Rightarrow a **dependency** is a formula $\forall \bar{x}. \alpha \rightarrow \beta$ where α and β are CQs.
- Algorithm:
 - 1 chase φ with Σ producing a CQ $\mathbf{chase}_\Sigma(\varphi)$;
 - $\mathbf{chase}_\Sigma^0 = \varphi$
 - $\mathbf{chase}_\Sigma^{i+1} = \mathbf{chase}_\Sigma^i \wedge (\beta\theta)$ for $\forall \bar{x}. \alpha \rightarrow \beta \in \Sigma$ and $\theta : \alpha \mapsto \mathbf{chase}_\Sigma^i$;
 - $\mathbf{chase}_\Sigma = \lim_{i \rightarrow \infty} \mathbf{chase}_\Sigma^i$.
 - 2 select $\psi \in \mathcal{L}(S_A)$ such that $\text{atoms}(\psi) \subseteq \text{atoms}(\mathbf{chase}_\Sigma(\varphi))$;
 - 3 chase ψ with Σ producing $\mathbf{chase}_\Sigma(\psi)$;
 - 4 test whether $\mathbf{chase}_\Sigma(\psi)$ implies φ
 \Rightarrow essentially $\text{atoms}(\varphi) \subseteq \text{atoms}(\mathbf{chase}_\Sigma(\psi))$.

Chase and Backchase

- Input: φ a **CQ**, Σ a set of **dependencies**, and S_A .
 \Rightarrow a **dependency** is a formula $\forall \bar{x}. \alpha \rightarrow \beta$ where α and β are CQs.
- Algorithm:
 - 1 chase φ with Σ producing a CQ **chase** $_{\Sigma}(\varphi)$;
 - 2 select $\psi \in \mathcal{L}(S_A)$ such that $\text{atoms}(\psi) \subseteq \text{atoms}(\mathbf{chase}_{\Sigma}(\varphi))$;
 - 3 chase ψ with Σ producing **chase** $_{\Sigma}(\psi)$;
 - 4 test whether **chase** $_{\Sigma}(\psi)$ implies φ
 \Rightarrow essentially $\text{atoms}(\varphi) \subseteq \text{atoms}(\mathbf{chase}_{\Sigma}(\psi))$.
- Problems:
 - **chase** $_{\Sigma}(\varphi)$ may be infinite (non-termination);
 \Rightarrow in theory restrict Σ to constraints w/terminating chase;
 \Rightarrow in practice fair interleaving of the steps of the algorithm
 - it only works well for CQs.

Won't work in General

- Chase extensions
 - disjunctions in heads of dependencies: UCQ plans
 - denial dependencies: pruning of disjuncts in such UCQ

Won't work in General

- Chase extensions
 - disjunctions in heads of dependencies: UCQ plans
 - denial dependencies: pruning of disjuncts in such UCQ
- does the algorithm find a plan if one exists?

Won't work in General

- Chase extensions
 - disjunctions in heads of dependencies: UCQ plans
 - denial dependencies: pruning of disjuncts in such UCQ
- does the algorithm find a plan if one exists?

Example

- $S_L = \{R/2\}$, $S_P = S_A = \{V_1/2/0, V_2/2/0, V_3/2/0\}$,
- $\Sigma = \left\{ \begin{array}{l} \forall x, y. V_1(x, y) \equiv \exists u, w. (R(u, x) \wedge R(u, w) \wedge R(w, y)), \\ \forall x, y. V_2(x, y) \equiv \exists u, w. (R(x, u) \wedge R(u, w) \wedge R(w, y)), \\ \forall x, y. V_3(x, y) \equiv \exists u. (R(x, u) \wedge R(u, y)) \end{array} \right\}$,
- $\varphi = \exists u, v, w. (R(u, x) \wedge R(u, w) \wedge R(w, v) \wedge R(v, y))$,

Won't work in General

- Chase extensions
 - disjunctions in heads of dependencies: UCQ plans
 - denial dependencies: pruning of disjuncts in such UCQ
- does the algorithm find a plan if one exists?

Example

- $S_L = \{R/2\}$, $S_P = S_A = \{V_1/2/0, V_2/2/0, V_3/2/0\}$,
- $\Sigma = \left\{ \begin{array}{l} \forall x, y. V_1(x, y) \equiv \exists u, w. (R(u, x) \wedge R(u, w) \wedge R(w, y)), \\ \forall x, y. V_2(x, y) \equiv \exists u, w. (R(x, u) \wedge R(u, w) \wedge R(w, y)), \\ \forall x, y. V_3(x, y) \equiv \exists u. (R(x, u) \wedge R(u, y)) \end{array} \right\}$
- $\varphi = \exists u, v, w. (R(u, x) \wedge R(u, w) \wedge R(w, v) \wedge R(v, y))$,
- $\psi = \exists u. (V_1(x, u) \wedge \forall w. (V_3(w, u) \rightarrow V_2(w, y)))$.

... but there is not a CQ rewriting.

Won't work in General

- Chase extensions
 - disjunctions in heads of dependencies: UCQ plans
 - denial dependencies: pruning of disjuncts in such UCQ
- does the algorithm find a plan if one exists?

Example

- $S_L = \{R/2\}$, $S_P = S_A = \{V_1/2/0, V_2/2/0, V_3/2/0\}$,
- $\Sigma = \left\{ \begin{array}{l} \forall x, y. V_1(x, y) \equiv \exists u, w. (R(u, x) \wedge R(u, w) \wedge R(w, y)), \\ \forall x, y. V_2(x, y) \equiv \exists u, w. (R(x, u) \wedge R(u, w) \wedge R(w, y)), \\ \forall x, y. V_3(x, y) \equiv \exists u. (R(x, u) \wedge R(u, y)) \end{array} \right\}$
- $\varphi = \exists u, v, w. (R(u, x) \wedge R(u, w) \wedge R(w, v) \wedge R(v, y))$,
- $\psi = \exists u. (V_1(x, u) \wedge \forall w. (V_3(w, u) \rightarrow V_2(w, y)))$.

... but there is not a CQ rewriting.

\Rightarrow cannot be found by chase-backchase

INTERPOLATION

Definition (Beth Definability)

A formula φ is *definable w.r.t. Σ and S_A* if $\varphi^{M_1} = \varphi^{M_2}$
for every pair M_1, M_2 of models of Σ such that $R^{M_1} = R^{M_2}$ for all $R \in S_A$.

\Rightarrow sometimes called *parametric definability* (due to S_A).

Definability and Interpolation

Definition (Beth Definability)

A formula φ is *definable w.r.t. Σ and S_A* if $\varphi^{M_1} = \varphi^{M_2}$
for every pair M_1, M_2 of models of Σ such that $R^{M_1} = R^{M_2}$ for all $R \in S_A$.

\Rightarrow sometimes called *parametric definability* (due to S_A).

Theorem (Craig'57)

Let α and β be FO formulæ such that $\models \alpha \rightarrow \beta$. Then there is a FO formula $\gamma \in \mathcal{L}(\alpha) \cap \mathcal{L}(\beta)$, called *an interpolant*, such that $\models \alpha \rightarrow \gamma$ and $\models \gamma \rightarrow \beta$.

How do we Use it?

IDEA:

Only allow queries that are Beth definable w.r.t. Σ and S_A

⇒ provides users with an illusion of a *single model*

How do we Use it?

IDEA:

Only allow queries that are Beth definable w.r.t. Σ and S_A

\Rightarrow provides users with an illusion of a *single model*

Definability Test:

φ is definable w.r.t. Σ and S_A if and only if $\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^*$
for Σ^* and φ^* having all $R \notin S_A$ replaced by R^* (Beth).

How do we Use it?

IDEA:

Only allow queries that are Beth definable w.r.t. Σ and S_A

\Rightarrow provides users with an illusion of a *single model*

Definability Test:

φ is definable w.r.t. Σ and S_A if and only if $\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^*$
for Σ^* and φ^* having all $R \notin S_A$ replaced by R^* (Beth).

Interpolant Existence:

If φ is definable w.r.t. Σ and S_A then
there is a FO $\psi \in \mathcal{L}(S_A)$ such that $\Sigma \models \varphi \leftrightarrow \psi$ (Craig).

How do we Use it?

IDEA:

Only allow queries that are Beth definable w.r.t. Σ and S_A

\Rightarrow provides users with an illusion of a *single model*

Definability Test:

φ is definable w.r.t. Σ and S_A if and only if $\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^*$
for Σ^* and φ^* having all $R \notin S_A$ replaced by R^* (Beth).

Interpolant Existence:

If φ is definable w.r.t. Σ and S_A then
there is a FO $\psi \in \mathcal{L}(S_A)$ such that $\Sigma \models \varphi \leftrightarrow \psi$ (Craig).

NOTE: this does NOT account for *binding patterns*.

Derivation

INPUT: finite Σ and φ .; output: ψ

$$\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^* \quad \Rightarrow$$

$$\models (\bigwedge \Sigma) \rightarrow ((\bigwedge \Sigma^*) \rightarrow (\varphi \rightarrow \varphi^*)) \quad \Rightarrow$$

$$\models (\bigwedge \Sigma) \rightarrow (\varphi \rightarrow ((\bigwedge \Sigma^*) \rightarrow \varphi^*)) \quad \Rightarrow$$

$$\models ((\bigwedge \Sigma) \wedge \varphi) \rightarrow ((\bigwedge \Sigma^*) \rightarrow \varphi^*) \quad \Rightarrow$$

$$\models ((\bigwedge \Sigma) \wedge \varphi) \rightarrow \psi \text{ and } \models \psi \rightarrow ((\bigwedge \Sigma^*) \rightarrow \varphi^*) \quad \Rightarrow$$

$$\models (\bigwedge \Sigma) \rightarrow (\varphi \rightarrow \psi) \text{ and } \models (\bigwedge \Sigma^*) \rightarrow (\psi \rightarrow \varphi^*) \quad \Rightarrow$$

$$\Sigma \models \varphi \rightarrow \psi \text{ and } \Sigma^* \models \psi \rightarrow \varphi^* \quad \Rightarrow$$

$$\Sigma \cup \Sigma^* \models \varphi \rightarrow \psi \text{ and } \Sigma \cup \Sigma^* \models \psi \rightarrow \varphi^*$$

$$\Sigma \cup \Sigma^* \models \varphi^* \rightarrow \varphi \quad \Rightarrow$$

$$\vdots \quad \Rightarrow$$

$$\Sigma \cup \Sigma^* \models \varphi^* \rightarrow \psi \text{ and } \Sigma \cup \Sigma^* \models \psi \rightarrow \varphi$$

Constructive Interpolation via Tableau

IDEA:

We try to *prove* $\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^*$ producing a proof (in a form of closed tableau) from which *extract the interpolant*.

Constructive Interpolation via Tableau

IDEA:

We try to *prove* $\Sigma \cup \Sigma^* \models \varphi \rightarrow \varphi^*$ producing a proof (in a form of closed tableau) from which *extract the interpolant*.

(Biased) Analytic Tableau

A *refutation* proof system for FOL:

- instead of $\vdash \alpha \rightarrow \beta$ we show $S = \{\alpha^L, \neg\beta^R\}$ is inconsistent
formulæ in S are *adorned* by L and R (needed for interpolant extraction);
- we use *inference rules* to generate *successors* of S in a proof tree;
- a proof is complete if all leaves contain a clash, a pair $\delta, \neg\delta$
otherwise the tableau saturates and we can extract a counterexample.

Interpolant Extraction (by example)

- an *invariant* for interpolation $S \xrightarrow{int} \psi$ is $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$ where S^L and S^R are subsets of S derived from adornments of formulas.

Interpolant Extraction (by example)

- an *invariant* for interpolation $S \xrightarrow{int} \psi$ is $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$ where S^L and S^R are subsets of S derived from adornments of formulas.
- tableau rules (sample):
 - LR clash $S \cup \{R^L, \neg R^R\} \xrightarrow{int} R$, $R \in S_A$ because $(\bigwedge S^L \wedge R^L) \rightarrow R$ and $R \rightarrow (R^R \vee \neg \bigwedge S^R)$

Interpolant Extraction (by example)

- an *invariant* for interpolation $S \xrightarrow{int} \psi$ is $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$ where S^L and S^R are subsets of S derived from adornments of formulas.

- tableau rules (sample):

- LR clash $\boxed{S \cup \{R^L, \neg R^R\} \xrightarrow{int} R}$, $R \in S_A$ because
 $(\bigwedge S^L \wedge R^L) \rightarrow R$ and $R \rightarrow (R^R \vee \neg \bigwedge S^R)$

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

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

Interpolant Extraction (by example)

- an *invariant* for interpolation $S \xrightarrow{int} \psi$ is $(\bigwedge S^L) \rightarrow \psi$ and $\psi \rightarrow (\neg \bigwedge S^R)$ where S^L and S^R are subsets of S derived from adornments of formulas.

- tableau rules (sample):

- LR clash $\boxed{S \cup \{R^L, \neg R^R\} \xrightarrow{int} R}$, $R \in S_A$ because
 $(\bigwedge S^L \wedge R^L) \rightarrow R$ and $R \rightarrow (R^R \vee \neg \bigwedge S^R)$

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

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

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

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

- etc. (see [Fitting] for details)

Implementation “details”

1 Plan enumeration:

⇒ enumeration of proofs \sim enumeration all equivalent rewritings? (NO)

Implementation “details”

① Plan enumeration:

- ⇒ enumeration of proofs \sim enumeration all equivalent rewritings? (NO)
- ⇒ do we want to enumerate all equivalent rewritings? (NO)

Implementation “details”

1 Plan enumeration:

- ⇒ enumeration of proofs \sim enumeration all equivalent rewritings? (NO)
- ⇒ do we want to enumerate all equivalent rewritings? (NO, **why?**)

Implementation “details”

1 Plan enumeration:

- ⇒ enumeration of proofs \sim enumeration all equivalent rewritings? (NO)
- ⇒ do we want to enumerate all equivalent rewritings? (NO)
- ⇒ do we get “enough”? (NO)

Implementation “details”

1 Plan enumeration:

- ⇒ enumeration of proofs \sim enumeration all equivalent rewritings? (NO)
- ⇒ do we want to enumerate all equivalent rewritings? (NO)
- ⇒ do we get “enough”? (NO, **needs tableau modifications**)

Implementation “details”

1 Plan enumeration:

- ⇒ enumeration of proofs \sim enumeration all equivalent rewritings? (NO)
- ⇒ do we want to enumerate all equivalent rewritings? (NO)
- ⇒ do we get “enough”? (NO)

2 Is backtracking of the tableau proofs feasible approach? (NO)

⇒ in Σ we separate

- “logical” (lots, complex) and
- “physical” (few, simple) constraints

... limits backtracking during plan search to physical constraints;

Implementation “details”

- 1 Plan enumeration:
 - ⇒ enumeration of proofs \sim enumeration all equivalent rewritings? (NO)
 - ⇒ do we want to enumerate all equivalent rewritings? (NO)
 - ⇒ do we get “enough”? (NO)
- 2 Is backtracking of the tableau proofs feasible approach? (NO)
 - ⇒ in Σ we separate
 - “logical” (lots, complex) and
 - “physical” (few, simple) constraints
 - ... limits backtracking during plan search to physical constraints;
- 3 Still needs to check for satisfaction of *binding patterns*.

POST-PROCESSING

Duplicate Elimination Elimination

In general $\exists x.\psi$ has to *eliminate duplicates* in the result (expensive)

\Rightarrow we want to detect when duplicate elimination can be safely omitted.

Duplicate Elimination Elimination

In general $\exists x.\psi$ has to *eliminate duplicates* in the result (expensive)

\Rightarrow we want to detect when duplicate elimination can be safely omitted.

IDEA:

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

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

Duplicate Elimination Elimination

In general $\exists x.\psi$ has to *eliminate duplicates* in the result (expensive)

\Rightarrow we want to detect when duplicate elimination can be safely omitted.

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 $\{\cdot\}$:

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

$$Q[\{Q_1 \wedge Q_2\}] \leftrightarrow Q[\{Q_1\} \wedge \{Q_2\}]$$

$$Q[\{\neg Q_1\}] \leftrightarrow Q[\neg Q_1]$$

$$Q[\{\neg\{Q_1\}\}] \leftrightarrow Q[\neg Q_1]$$

$$Q[\{Q_1 \vee Q_2\}] \leftrightarrow Q[\{Q_1\} \vee \{Q_2\}] \quad \text{if } \Sigma \cup \{Q[\cdot]\} \models Q_1 \wedge Q_2 \rightarrow \perp$$

$$Q[\{\exists x.Q_1\}] \leftrightarrow Q[\exists x.\{Q_1\}] \quad \text{if } \Sigma \cup \{Q[\cdot] \wedge (Q_1)[y_1/x] \wedge (Q_1)[y_2/x]\} \models y_1 \approx y_2$$

where y_1 and y_2 are fresh variable names not occurring in Q , Q_1 , and Q_2 .

Summary

- 1 interpolation provides a powerful tool for query optimization, but
 - efficiency of reasoning is an issue (single proof is not sufficient)
 - generating enough candidate plans (at odds with structural proofs)
 - but needs to avoid useless plans (e.g., co-joining tautologies, etc.)

Summary

- 1 interpolation provides a powerful tool for query optimization, but
 - efficiency of reasoning is an issue (single proof is not sufficient)
 - generating enough candidate plans (at odds with structural proofs)
 - but needs to avoid useless plans (e.g., co-joining tautologies, etc.)
- 2 postprocessing needed to deal with non-FO features
 - duplicate semantics (hard to even define query equivalence!)
 - cuts (see textbook for details)