

# Özgür L. Özçep

# Query Reformulation

Lecture 13

Foundations of Ontologies and Databases for Information Systems CS5130 (Winter 16/17)



### Recap

- Talked about higher-level declarative stream processing using STARQL as example
  - Declarative: streams have assertional status
  - ► High-level: Have to incorporate (reason over) a background KB

#### This Lecture

- Back to general topic of query rewriting
  - Here we call it "Query reformulation"
  - ► Examples from DB theory: Reformulating a query w.r.t. views
- Uses the "last significant property of FOL that has come to light" namely interpolation (Benthem 2008)
   Lit: J. van Benthem. The many faces of interpolation. Synthese,

164(3):451?460, 2008.

- Sources used for this lecture
  - Mainly: Slides from an invited talk of M. Benedikt at DL 2014 with my annotations
  - B. ten Cate: "Craig Interpolation Theorems and Database Applications", talk given at UC Berkeley - Logic Colloquium, November 7, 2014
  - M. Benedikt, J. Leblay, B. ten Cate, and E. Tsamoura. Generating Plans from Proofs: The interpolation-based Approach to Query Reformulation. Synthesis Lectures on Data Management, 2016.

#### Example for Rewriting

#### Following Road Example from (tenCate 2014)

#### First Example: View-Based Query Reformulation

- Road network database: Road(x,y)
- Views:
  - $V_2(x,y) = "\exists$  path of length 2 from x to  $y'' = \exists u \operatorname{Road}(x,u) \land \operatorname{Road}(u,y)$
  - $V_3(x,y) = "\exists$  path of length 3 from x to  $y'' = \exists u, v \operatorname{Road}(x,u) \land \operatorname{Road}(u,v) \land \operatorname{Road}(v,y)$
  - ...
- Observation: V<sub>4</sub> can be expressed in terms of V<sub>2</sub>.
- Puzzle (Afrati'07): can V<sub>5</sub> be expressed (in FO logic) in terms of V<sub>3</sub> and V<sub>4</sub>?

## Solution to the puzzle

 $V_5(x,y) \Leftrightarrow \exists u \ ( \ V_4(x,u) \land \forall v \ ( \ V_3(v,u) \rightarrow V_4(v,y) \ ) \ ) \ )$ 

Proof:



### Why this Example is Important

- A conjunctive query (CQ) is a FO formula built up using only A, J.
  - Conjunctive queries are the most common type of database queries.
  - Every positive-existential FO formula is equivalent to a union of CQs.
- Remarkable fact:
  - V<sub>3</sub>, V<sub>4</sub> and V<sub>5</sub> are all defined by CQs over the base relation (Road).
  - V<sub>5</sub> is definable in terms of V<sub>3</sub> and V<sub>4</sub> but not by means of a CQ.

### **Classic Results**

Querying using views has been around since the 1980s. E.g.,

- **Theorem** (Levy Mendelzon Sagiv Srivastava '95): there is an effective procedure to decide whether a conjunctive query is rewritable as a conjunctive query over a given set of conjunctive views.
- **Open problem** (Nash, Segoufin, Vianu '10): is there an effective procedure to decide if a conjunctive query is answerable on the basis of a set of conjunctive views (a.k.a., is determined by the views)? if so, in what language can we express the rewriting?

**NB**: The Beth definability theorem (1953) tell us that, if a FO query is answerable on the basis of a set of FO views, then, it has a FO rewriting.

## New Perspectives on Query Reformulation

Michael Benedikt, based on joint work with: Julien Leblay, Efi Tsamoura, Michael Vanden Boom (Oxford) Balder ten Cate (LogicBlox & UC-Santa Cruz)

- Overview the theory behind proof-driven querying, which takes:
  - a conjunctive query query Q =  $\exists x_1 \dots x_m A_1(x_1 \dots) \land \dots \land A_m(x_m \dots)$
  - metadata consisting of relation descriptions, access methods, first order constraints, and costs of access

and generates (if possible) a low cost plan that gives the same answer as Q on every instance satisfying the constraints

- · Give connections to:
  - proof theory
  - · preservation theorems in model theory
  - data integration
  - · work of other DL '14 invited speakers

## William Craig's Approach to Query Reformulation Meta-Algorithm



In two papers from 1957, W. Craig developed a general methodology for reformulating a query **Q** in some restricted target language with respect to constraints:

- Pick out a property that **Q** should possess in order to have the desired reformulation (E.g.: Determinedness)
- Write out the property as an implication/proof goal in logic
- Look for a proof fulfilling the proof goal
- Generate a reformulation from the proof using an **interpolation** algorithm

## Craig in Action

**Goal:** given boolean conjunctive query **Q** find a **first-order logic reformulation** over some restricted signature **T** with respect to constraints  $\Sigma$  (constraints also in FO)

Show that **Q** is **determined** over **T** with respect to  $\varSigma$ : for any two instances I and I' satisfying the constraints  $\varSigma$ if I and I' have the same **T** facts, then **Q**(I)=**Q**(I')

Do this by showing that the following implication, using two copies of relations, is valid:  $\mathbf{Q} \wedge \Sigma \wedge \Sigma' \wedge [\Lambda_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_1 \dots \mathbf{x}_n \mathbf{R}(\mathbf{x}_1 \dots \mathbf{x}_n) \leftrightarrow \mathbf{R}'(\mathbf{x}_1 \dots \mathbf{x}_n)] \rightarrow \mathbf{Q}'$   $\Sigma' = \text{copy of constraints on primed relations}$  $\mathbf{Q}' = \text{copy of query on } \mathbf{Q} \text{ primed relations}$ 

To witness the validity, search for a proof that:

 $\mathbf{Q} \wedge \boldsymbol{\Sigma} \wedge \boldsymbol{\Sigma}' \wedge [\boldsymbol{\wedge}_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_1 \dots \mathbf{x}_n \ \mathbf{R}(\mathbf{x}_1 \dots \mathbf{x}_n) \leftrightarrow \mathbf{R}'(\mathbf{x}_1 \dots \mathbf{x}_n)] \vdash \mathbf{Q}'$ Re-arranging:  $\mathbf{Q} \wedge \boldsymbol{\Sigma} \vdash [\boldsymbol{\Sigma}' \wedge \boldsymbol{\wedge}_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_1 \dots \mathbf{x}_n \ \mathbf{R}(\mathbf{x}_1 \dots \mathbf{x}_n)] \leftrightarrow \mathbf{R}'(\mathbf{x}_1 \dots \mathbf{x}_n)] \rightarrow \mathbf{Q}'$ 

Craig 1957 paper: from any interpolant for re-arranged proof goal, can obtain a first-order reformulation of Q over T with respect to  $\Sigma$ 

**Interpolant**: If  $\rho_1 \vdash \rho_2$ , an interpolant is a formula  $\rho$  such that  $\rho_1 \vdash \rho \vdash \rho_2$ , and  $\rho$  uses only relations common to  $\rho_1$  and  $\rho_2$ 

## Craig's Reformulation Recipe

To find a **first-order reformulation** of query **Q** over some restricted signature **T** with respect to constraints  $\varSigma$  .

Show that **Q** is **determined** over **T** with respect to  $\varSigma$ : for any two instance I and I' satisfying the constraints  $\varSigma$  if I and I' have the same **T** facts, then **Q**(I)=**Q**(I')

This condition holds iff there is a **proof** that:

 $\mathbf{Q} \land \varSigma \land \varSigma \land [\land_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_{1} ... \mathbf{x}_{n} \mathsf{R}(\mathbf{x}_{1} ... \mathbf{x}_{n}) \leftrightarrow \mathsf{R}'(\mathbf{x}_{1} ... \mathbf{x}_{n}) ] \vdash \mathbf{Q}'$ 

Craig '57:

- Gave a proof system complete for finding these implications
- Gave an interpolation algorithm: from a proof witnessing  $\rho_1 \vdash \rho_2$ , efficiently extracts a formula  $\rho$  such that  $\rho_1 \vdash \rho \vdash \rho_2$ , and  $\rho$  uses only relations common to  $\rho_1$  and  $\rho_2$
- Showed that any interpolant gives a reformulation, and if there is **any** first-order reformulation, this process gives one

# First Order Tableau proofs

To Prove:  $\exists x A(x) \land \neg B(x) \land C(x) \vdash \neg \forall y [(\neg A(y) \land E(y)) \lor B(y)]$ 

Derive a contradiction from:  $\exists x A(x) \land \neg B(x) \land C(x), \forall y [(\neg A(y) \land E(y)) \lor B(y)]$  $A(c) \land \neg B(c) \land C(c), \forall v [(\neg A(v) \land E(v)) \lor B(v)]$  $A(c) \land \neg B(c) \land C(c)$ ,  $[(\neg A(c) \land E(c)) \lor B(c)]$ A(c),  $\neg B(c)$ , C(c),  $[(\neg A(c) \land E(c)) \lor B(c)]$ A(c),  $\neg B(c)$ , C(c),  $\neg A(c) \land E(c)$ A(c),  $\neg B(c)$ , C(c), B(c)A(c),  $\neg B(c)$ , C(c),  $\neg A(c)$ , E(c)

# Tableau proofs

To Prove:  $\exists x A(x) \land \neg B(x) \land C(x) \vdash \neg \forall y [(\neg A(y) \land E(y)) \lor B(y)]$ 

Derive a contradiction from:  

$$\exists x \ A(x) \land \neg B(x) \land C(x), \ \forall y [(\neg A(y) \land E(y)) \lor B(y)]$$

$$A(c) \land \neg B(c) \land C(c), \ \forall y [(\neg A(y) \land E(y)) \lor B(y)]$$

$$A(c) \land \neg B(c) \land C(c), \ [(\neg A(c) \land E(c)) \lor B(c)]$$

$$A(c), \ \neg B(c), C(c), \ [(\neg A(c) \land E(c)) \lor B(c)]$$

$$A(c), \ \neg B(c), C(c), \ \neg A(c) \land E(c)$$

$$A(c), \ \neg B(c), C(c), \ \neg A(c), \ E(c)$$

# Tableaux and Craig's Interpolation Algorithm

 $\exists x \ A(x) \land \neg B(x) \land C(x) \vdash \neg \forall y [(\neg A(y) \land E(y)) \lor B(y)]$ 

 $\exists x \ A(x) \land \neg B(x) \land C(x) \ , \ \forall y \ [ \ (\neg A(y) \land E(y) \ ) \lor B(y) ]$ 

 $A(c) \land \neg B(c) \land C(c) \ , \ \forall y [ (\neg A(y) \land E(y) ) \lor B(y)]$ 

 $A(c) \land \neg B(c) \land C(c)$ ,  $[(\neg A(c) \land E(c)) \lor B(c)]$ 

A(c),  $\neg B(c)$ , C(c),  $[(\neg A(c) \land E(c)) \lor B(c)]$ 

 $A(c) , \ \neg B(c) , C(c), \ \neg A(c) \land E(c)$   $A(c) , \ \neg B(c) , C(c), \ B(c)$ 

A(c),  $\neg B(c)$ , C(c),  $\neg A(c)$ , E(c)

# Tableaux and Craig's Interpolation Algorithm

 $\exists x A(x) \land \neg B(x)$  $\exists x A(x) \land \neg B(x) \land C(x) \vdash \neg \forall y [(\neg A(y) \land E(y)) \lor B(y)]$ 

 $\exists x \ A(x) \land \neg B(x) \land C(x) , \ \forall y [(\neg A(y) \land E(y)) \lor B(y)]$ 

 $\begin{array}{c} \mathsf{A}(\mathsf{c}) \land \neg \mathsf{B}(\mathsf{c}) \land \mathsf{C}(\mathsf{c}) , \ \forall \ \mathbf{y} \left[ \left( \neg \mathsf{A}(\mathsf{y}) \land \mathsf{E}(\mathsf{y}) \right) \lor \mathsf{B}(\mathsf{y}) \right] \quad \mathsf{Exists} \times \mathsf{A}(\mathsf{x}) \And_{\neg} \mathsf{B}(\mathsf{x}) \\ & \\ \mathsf{A}(\mathsf{c}) \land \neg \mathsf{B}(\mathsf{c}) \land \mathsf{C}(\mathsf{c}) , \ \left[ \left( \neg \mathsf{A}(\mathsf{c}) \land \mathsf{E}(\mathsf{c}) \right) \lor \mathsf{B}(\mathsf{c}) \right] \quad \mathsf{A}(\mathsf{c}) \And_{\neg} \mathsf{B}(\mathsf{c}) \\ & \\ \mathsf{A}(\mathsf{c}) , \ \neg \mathsf{B}(\mathsf{c}) , \mathsf{C}(\mathsf{c}) , \ \left[ \left( \neg \mathsf{A}(\mathsf{c}) \land \mathsf{E}(\mathsf{c}) \right) \lor \mathsf{B}(\mathsf{c}) \right] \quad \mathsf{A}(\mathsf{c}) \And_{\neg} \mathsf{B}(\mathsf{c}) \\ & \\ \mathsf{a}(\mathsf{c}) , \ \neg \mathsf{B}(\mathsf{c}) , \mathsf{C}(\mathsf{c}) , \ \mathsf{C}(\mathsf{c}) , \ \mathsf{E}(\mathsf{c}) \land \mathsf{E}(\mathsf{c}) \right) \lor \mathsf{B}(\mathsf{c}) \right] \quad \mathsf{A}(\mathsf{c}) \And_{\neg} \mathsf{B}(\mathsf{c}) \\ & \\ \mathsf{A}(\mathsf{c}) , \ \neg \mathsf{B}(\mathsf{c}) , \mathsf{C}(\mathsf{c}) , \ \neg \mathsf{A}(\mathsf{c}) \land \mathsf{E}(\mathsf{c}) \land \mathsf{A}(\mathsf{c}) \\ & \\ \mathsf{A}(\mathsf{c}) , \ \neg \mathsf{B}(\mathsf{c}) , \mathsf{C}(\mathsf{c}) , \ \mathsf{B}(\mathsf{c}) \end{array}$ 

A(c),  $\neg B(c)$ , C(c),  $\neg A(c)$ , E(c)

Interpolant (in red) construction by bottom up strategy

## Summary: Craig's Methodology Episode 1

To find a first-order reformulation of query  ${\bf Q}$  over some restricted signature  ${\bf T}$  with respect to constraints  $\varSigma$  .

Show that **Q** is **determined** over **T** with respect to  $\varSigma$ : for any two instance I and I' satisfying the constraints  $\varSigma$  if I and I' have the same **T** facts, then **Q**(I)=**Q**(I')

This condition holds iff:

 $\mathbf{Q} \land \varSigma \land \varSigma \land [\land_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_{1} ... \mathbf{x}_{n} \mathsf{R}(\mathbf{x}_{1} ... \mathbf{x}_{n}) \leftrightarrow \mathsf{R}'(\mathbf{x}_{1} ... \mathbf{x}_{n})] \vdash \mathbf{Q}'$ 

Apply Craig's **interpolation algorithm** to any tableau proof witnessing the re-arranged version of this: it will give a reformulation.

#### But wait...

- Craig's Lemma works for FOL where we have natural domain semantics and allow for infinite structures
- In DB one normally has active domain semantics and is interested in finite structures (but in the finite interpolation does not hold for FOL)
- Idea
  - Consider fragments of FOL where natural domain semantics (alias classical FOL semantics) = active domain semantics (alias DB semantics) = RQFO
  - Consider further restrictions on constraints (Guarded fragment logics)
- RQFO = FOL with relativized quantifiers
  - $\blacktriangleright \exists \vec{x} R(\vec{y}, \vec{x}) \land \phi(\vec{y}, \vec{x}, \vec{t}))$
  - $\forall \vec{x} R(\vec{y}, \vec{x}) \rightarrow \phi(\vec{y}, \vec{x}, \vec{t}))$ where  $\phi$  again a RQFO
  - Note that RQFO are not necessarily safe (A(x) ∨ B(y) still allowed).

Interpolation does not hold in the Finite for FOL

- $\mathcal{L}$  a logic;  $K = \text{class of } \tau \text{ structures}$
- Call K good iff for class

 $\mathcal{K}_{<} := \{(\mathfrak{A}, <) \mid \mathfrak{A} \in \mathcal{K}, < \text{ an ordering on } \mathfrak{A}\}$ 

there is sentence  $\phi$  in vocabulary  $\tau \cup <$  such that

$$K_{<} = Mod(\phi)$$

L is closed under order-invariant sentences in the finite iff every good structure K is axiomatizable by a τ sentence ψ

#### Observation

Logics with interpolation property are closed under order-invariant sentences in the finite.

Proof: Consider φ = φ(<) with Mod(K) = φ. Then φ(<) ⊨<sub>fin</sub> (" <′ is an ordering" → φ(<′)). If ψ is interpolant, then Mod(K) = ψ. Interpolation does not hold in the Finite for FOL

- $\mathcal{L}$  a logic;  $K = \text{class of } \tau \text{ structures}$
- Call K good iff for class

 $\mathit{K}_{<} := \{(\mathfrak{A}, <) \mid \mathfrak{A} \in \mathit{K}, < \mathsf{an ordering on } \mathfrak{A}\}$ 

there is sentence  $\phi$  in vocabulary  $\tau \cup <$  s.t.  $K_{<} = Mod(\phi)$ 

L is closed under order-invariant sentences in the finite iff every good structure K is axiomatizable by a τ sentence ψ

#### Proposition

FOL is not closed under order-invariant sentences in the finite
 FOL does not have interpolation property in the finite

- Proof of 1.:
  - K = boolean algebras with even number of atoms.
  - ► K not axiomatizable in FOL (use Fraissee game)
  - But  $K_{<}$  is axiomatizable by  $\phi$ :
  - $\phi$  = boolean algebra axioms + order axioms + sentence A
  - ► A = there is an element covering exactly atoms in even position and the last atom.

## Craig works for database-style queries + constraints



To find a **relational algebra** reformulation of conjunctive query **Q** over subschema **T** with respect to **relational algebra constraints**  $\Sigma$  equivalently, constraints and reformulation built up via **relative quantifiers**:  $\exists x_1...x_m R(x_1...x_m, y_1...y_n) \land \phi$  "Safe quantification for domain independence"  $\forall x_1...x_m R(x_1...x_m, y_1...y_n) \rightarrow \phi$ 

Show that **Q** is determined over **T** with respect to  $\Sigma$ : for any two instance I and I' satisfying the constraints  $\Sigma$ if I and I' have the same **T** facts, then **Q**(I)=**Q**(I')

To do this, need to show:

 $\mathsf{Q} \land \varSigma \land \varSigma \land [ \land_{\mathcal{R} \, \in \, \mathcal{T}} \forall \, \mathsf{x}_1 ... \, \mathsf{x}_n \, \mathsf{R}(\mathsf{x}_1 ... \mathsf{x}_n) \leftrightarrow \mathsf{R'}(\mathsf{x}_1 \, ... \, \mathsf{x}_n)] \vdash \mathsf{Q'}$ 

Apply Craig's **interpolation algorithm** to any tableau proof of re-arranged proof goal to extract a reformulation. **Observe** that if all formulas involved are relativized, the resulting interpolant can be put in relativized form as well.

## Based on a true story

This is a reconstruction of Craig's work.

The exposition of interpolation comes from Fitting's 1996 textbook.

The connection between interpolation and query reformulation was discovered by Segoufin and Vianu [PODS 2005] and explored more fully in Nash, Segoufin, Vianu [TODS 2010]. They make use of interpolation results of Otto [BSL 2000].

A. Nash, L. Segoufin, and V. Vianu. Views and queries: Determinacy and rewriting. ACM Trans. Database Syst., 35(3):21:1–21:41, 2010.

M. Otto. An interpolation theorem. Bulletin Symbolic Logic, 6(4):447-462, 12 2000.

The fact that interpolation can be used to generate reformulations constructively has been observed several places:

Toman and Weddell, Fundamentals of Physical Design and Query Reformulation (Book)

• Franconi, Kerhet and Ngo, JAIR 2013



L. Segoufin and V. Vianu. Views and queries: determinacy and rewriting. PODS 2005, pages 49–60. ACM, 2005.

E. Franconi, V. Kerhet, and N. Ngo.Exact query reformulation over databases with first-order and description logics ontologies. J. Artif. Intell. Res. (JAIR), 48:885–922, 2013.













## Craig works for finding USPJ Reformulations

To find a **positive existential** (i.e. **UCQ/USPJ** with inequalities) reformulation of conjunctive query **Q** over subschema **T** with respect to (relativized FO) constraints  $\Sigma$ .

Show that **Q** is **monotone** over **T** with respect to  $\varSigma$ :

for any two instances I and I' satisfying the constraints  $\varSigma$ if T facts of I are contained in the T facts of I', then  $Q(I) \subseteq Q(I')$ 

To do this need to show:

 $\mathbf{Q}\wedge\varSigma\wedge\varSigma^{'}\wedge\llbracket\wedge_{\mathcal{R}\,\in\,\mathcal{T}}\forall\ \mathbf{x}_{1}...\ \mathbf{x}_{n}\ \mathbf{R}(\mathbf{x}_{1}...\mathbf{x}_{n})\rightarrow\mathbf{R}'(\mathbf{x}_{1}\,...\,\mathbf{x}_{n})]\vdash\ \mathbf{Q}'$ 

Apply Craig's **interpolation algorithm** to a tableau proof of the re-arranged proof goal and extract a reformulation. *Argue that the resulting reformulation will be positive existential.* (Lyndon: There is positive (negative) occurrence of relation symbol in interpolant iff is positive (negative) occurrence in premise and conclusion)

Interpolants produced by Craig's algorithm "respect positivity of relations": Lyndon Interpolation Theorem

## Craig works for Existential Reformulations

To find an **existential** (UCQ with atomic negation) reformulation of conjunctive query **Q** over some restricted signature **T** with respect to constraints  $\varSigma$ .

Show that **Q** is **induced-subinstance monotone** with respect to  $\varSigma$  :

for any two instances I and I' satisfying the constraints  $\Sigma$  if T facts of I contained in the T facts of I', and I' has no T fact using elements from I, then  $Q(I) \subseteq Q(I')$ 

To do this, need to show:  $Q \land \Sigma \land \Sigma \land [\land_{\mathcal{R} \in \mathcal{T}} \forall x_1...x_n R(x_1...x_n) \rightarrow R'(x_1...x_n) \land ...] \vdash Q'$ 

Rewrite and apply Craig's interpolation algorithm to tableau proof of re-arranged goal. *Argue that this gives the desired existential reformulation.* 

#### Summary: Craig's Recipe for Restricted Vocabulary Reformulations

Reformulation Goal	Semantic Property	Search for proof of this entailment, apply interpolation	Parallel with Theorem in Model Theory
<b>Q</b> is FO-rewritable over <b>T</b> w.r.t. constraints	<b>Q</b> is determined by <b>T</b> w.r.t constraints	$\mathbf{Q} \land \varSigma \land \land_{T} \mathbf{R}_{i} = \mathbf{R'}_{i}$ $\vdash \varSigma' \to \mathbf{Q'}$	Projective Beth Definability
<b>Q</b> is ∃*-rewritable over <b>T</b> w.r.t constraints	<b>Q</b> is monotone in <b>T</b> w.r.t. constraints	$\mathbf{Q} \land \varSigma \land \land_{T} \mathbf{R} \subseteq \mathbf{R'}_{i} \\ \vdash \varSigma' \to \mathbf{Q'}$	Rough analogy: Lyndon Preservation Theorem/ Hom. Pres. Thm
Q ∃-rewritable over T w.r.t. constraints	<b>Q</b> induced-subinst. monotone in <b>T</b> w.r.t constraints	$\mathbf{Q} \land \boldsymbol{\varSigma} \land \land_{T} \mathbf{R}_{i} \subseteq \mathbf{R}'_{i} \\ \land \dots \vdash \boldsymbol{\varSigma}' \to \mathbf{Q}'$	(Projective) Los-Tarski Preservation Theorem

## Craig works for Access Patterns



#### To find a **relational algebra plan using a fixed set of access methods** equivalent to query **Q** with respect to constraints $\Sigma$ .

SELECT a1.name FROM YahooPlaces / JOIN YahooBelo JOIN YahooPlac WHERE a0.name='As AND a1.placeTypeNi	NS a0 ongsTo AS a2 ON a0. es AS a1 sia' ame='Country'	woeid=a2.	target		
Search type	Chasing	type	Cost model	Cost	
OPTIMIZED	KTERMINATIO	N	BLACKBOX_CARD	2.5000501001276E13	
Search type Chasing type OPTIMI2ED Timeout Blocking intervals Infinity N/A Max. iterations Match intervals		T1 <= Ya T2 <= Ya T3 := T4 <= Ya T5 := T :=	$\begin{array}{l} T1 \in YahooCountries \notin \varnothing\\ T2 \in YahooBelongsTo/yh_geo_belongs \in m[c18](\sigma[\#2='Countr\\T3 := \sigma[\#2='Country'](p[c16, c20, Country, c19](T1)) \models [\#0='T4 \in YahooPlaces/h_geo_woeid \notin m[c1](T3)\\ T5 := T3 \models [\#5=\#9] \sigma (\#1='Asia'](p[c1, Asia, c2, c3, c4, c5, c6)\\T := m[c19](T5) \end{array}$		
Infinity Cost model BLACKBOX ~ Executor type DEFAULT ~	1 Plan				

### Access Methods

- Access method: a pair (R,X) where R is an n-ary relation and X⊆{1, ..., n} is a set of "input positions"
  - Relation R can be accessed if specific values are provided for the positions in X.
- Examples:
  - (Yellowpages(name,city,address,phone#), {1,2})
  - (R,Ø) means free (unrestricted) access to R.
  - (R,{1, ..., n}) means only **membership tests** for specific tuples.
- There may be any number of access methods for a given relation. The allowed access methods for a relation can be assumed to be an upwards closed set.

## Craig works for Access Patterns

#### To find a **relational algebra plan using a fixed set of access methods** equivalent to guery **Q** with respect to constraints $\Sigma$ .

Show that **Q** is **access-determined** with respect to access methods and  $\varSigma$ : for any two instances I and I' satisfying the constraints  $\varSigma$  if I and I' have the same accessible data, then **Q**(I) = **Q**(I')

Need to prove:  $\mathbf{Q} \land \boldsymbol{\Sigma} \land \boldsymbol{\Sigma}' \land$ ["Accessibility Axioms"]  $\vdash \mathbf{Q}'$ 

#### Example accessibility axiom:

Suppose  $R(x_1, x_2)$  has an access method on the first position.

Instead of previous "transfer axiom"  $\forall x_1 x_2 R(x_1, x_2) \rightarrow R'(x_1, x_2)$ 

We have, for any n-ary relation F in the signature, for every position in F, an axiom:

$$\forall \mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_n \mathbf{F'}(\dots \mathbf{x}_1 \dots) \land \mathbf{R}(\mathbf{x}_1, \mathbf{x}_2) \rightarrow \mathbf{R'}(\mathbf{x}_1, \mathbf{x}_2)$$

## Craig works for Access Patterns

To find a **relational algebra plan using a fixed set of access methods** equivalent to query **Q** with respect to constraints  $\Sigma$ .

Show that **Q** is **access-determined** with respect to access methods and  $\Sigma$ : for any two instances I and I' satisfying the constraints  $\Sigma$  if I and I' have the same accessible data, then **Q**(I) = **Q**(I')

Need to prove:  $\mathbf{Q} \land \Sigma \land \Sigma' \land$ ["Accessibility Axioms"]  $\vdash \mathbf{Q}'$ 

Re-arrange and apply Craig's interpolation algorithm to the proof to extract a reformulation. Argue that this can be converted to a plan that fits the access methods.

Interpolants produced by Craig's algorithm "reflect the quantification patterns on the left and right of the proof symbol": **Access Interpolation Theorem** 

## Craig works for Access-Pattern based Restrictions

To find a negation-free (USPJ) plan using a fixed set of access methods equivalent to query **Q** with respect to constraints  $\Sigma$ .

Show that **Q** is **access-monotone** with respect to methods and  $\Sigma$ : for any two instances I and I' satisfying the constraints  $\Sigma$  the accessible data in I is contained in the accessible data of I' then **Q**(I)  $\subseteq$  **Q**(I')

Need to show:  $\mathbf{Q} \land \Sigma \land \Sigma' \land$  ["Uni-directional Accessibility Axioms"]  $\vdash \mathbf{Q}'$ 

Re-arrange and apply Craig's interpolation algorithm to any proof to extract a reformulation. Argue using Access Interpolation Theorem that this gives a plan of the desired shape.

## **Binding Patterns**

Consider first order logic built up from equalities and true/false via relative quantifiers:

$$\forall \mathbf{y}_1 \dots \mathbf{y}_n \ \mathbf{R}(\mathbf{x}_1 \dots \mathbf{x}_m, \mathbf{y}_1 \dots \mathbf{y}_n) \rightarrow \phi$$
  
$$\exists \mathbf{y}_1 \dots \mathbf{y}_n \ \mathbf{R}(\mathbf{x}_1 \dots \mathbf{x}_m, \mathbf{y}_1 \dots \mathbf{y}_n) \land \phi(\mathbf{x}_1 \dots \mathbf{x}_m, \mathbf{y}_1 \dots \mathbf{y}_n)$$

The **binding pattern** of a formula tells where in each relation we have free variables in relative quantifications.

 $bindpatt(\forall \ y_1 \ y_2 \ R(x_1, \ y_1, \ y_2) \rightarrow S(x_1, y_2)) = \{(R, \ \{1\}), \ (S, \ \{1, 2\})\}$ 

#### Access Methods "Used" by a Formula

BindPatt( $\phi$ ) is the set of access methods "used" by  $\phi$ .

- For example BindPatt( $\forall y(Rxy \rightarrow Sxy)$ ) = { (R,{1}), (S,{1,2}) }
- A FO formula  $\phi$  is executable if BindPatt( $\phi$ ) consists of allowed access methods.
- Fact: Each executable FO formula admits a query plan, and, conversely, every formula that admits a query plan is equivalent to an executable FO formula.
  - Query plan = sequence of allowed accesses and/or relational algebra operations.

## Craig gives decision procedures



To find a positive existential reformulation of query **Q** over restricted signature **T** with respect to **Tuple-generating Dependencies (TGDs)**  $\Sigma$ 

$$\forall \mathbf{x}_1 \dots \mathbf{x}_m \ [\mathbf{R}_1(\mathbf{x}_1 \dots) \land \dots \land \mathbf{R}_m(\mathbf{x}_1 \dots) \rightarrow \exists \mathbf{y}_1 \dots \mathbf{y}_m \ \mathbf{S}_1(\mathbf{x}_1 \dots \mathbf{y}_1 \dots) \land \dots]$$

Very common in databases: e.g. inclusion dependencies (referential constraints), data integration mappings

## Craig gives decision procedures



the only point is that we can

To find a positive existential reformulation of query Q over restricted signature T with respect to Tuple-generating Dependencies  $\Sigma$ .

 $\forall \mathbf{x}_1 \dots \mathbf{x}_m \ [\mathbf{R}_1(\mathbf{x}_1 \dots) \land \dots \land \mathbf{R}_m(\mathbf{x}_1 \dots) \rightarrow \exists \mathbf{y}_1 \dots \mathbf{y}_m \ \mathbf{S}_1(\mathbf{x}_1 \dots \mathbf{y}_1 \dots) \land \dots]$ 

Show that **Q** is monotone with respect to  $\Sigma$ : for any two instance I and I' satisfying  $\Sigma$  if I and I' have the same T facts. Note: This is the same as then Q(I)=Q(I') above:

Search for proof that.

Q

earch for **proof** that:  

$$\land \Sigma \land \Sigma' \land [\land_{\mathcal{R} \in \mathcal{T}} \forall x_1 \dots x_n R(x_1 \dots x_n) \rightarrow R'(x_1 \dots x_n)] \vdash Q'$$

Q contained in Q' with respect to an augmented set of TGDs

Specialized proof procedures have been developed for proving implications of this form: the chase

## Craig gives decision procedures



To find a positive existential reformulation of query **Q** over restricted signature **T** with respect to **Tuple-generating Dependencies**  $\Sigma$ .

Show that **Q** is monotone with respect to  $\boldsymbol{\varSigma}$  :

for any two instance I and I' satisfying  $\varSigma$  if I and I' have the same T facts, then  $\mathbf{Q}(\mathbf{I}){=}\mathbf{Q}(\mathbf{I}')$ 

Search for proof that:

$$\begin{array}{c} \mathbf{Q} \land \varSigma \land \varSigma \land & [\land_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_{1} \dots \mathbf{x}_{n} \ \mathbf{R}(\mathbf{x}_{1} \dots \mathbf{x}_{n}) \rightarrow \mathbf{R}'(\mathbf{x}_{1} \dots \mathbf{x}_{n})] \vdash \mathbf{Q}' \\ \uparrow & \uparrow \\ \end{array}$$

**Q** contained in **Q'** with respect to an augmented set of TGDs

Apply Craig's algorithm to a chase proof (which is a tableau proof!) to extract a positive existential reformulation

### Chase again

- Remember the chase construction from our lecture on data exchange
  - Use TGDs as firing rules.
  - Chase procedure complete and correct; termination under some constraints f
    ür TGDs
- Here for testing  $Q \land \Sigma \models Q^*$ , where Q a CQ
- Canonical database DB(Q) for input Q:

instance with elements all constants in Q and additional distinct constants for variables; fact in DB(Q) iff it is contained in Q (modulo substituting new constants with variables)

- Chase procedure
  - 1. start with DB(Q)
  - 2. iteratively apply TGDs
  - 3. Stop if some fact contained (via a homomorphism) in  ${\it Q}^{\ast}$
- Chase derivation can be transformed into a tableau derivation

## End-to-end reformulation algorithms via Craig's technique

To find a positive existential reformulation of query **Q** over restricted signature **T** with respect to some set of TGD constraints  $\Sigma$  where the chase terminates (i.e. there is a maximal tableau).

Show that **Q** is **monotone** over **T** with respect to  $\Sigma$ : for any two instances I and I' satisfying the constraints  $\Sigma$ if **T** facts of I are contained in the **T** facts of I', then **Q**(I)  $\subseteq$  **Q**(I')

Perform the chase effectively to determine if there is a proof that:

 $\mathbf{Q} \land \varSigma \land \varSigma \land [\land_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_1 ... \mathbf{x}_n \ \mathbf{R}(\mathbf{x}_1 ... \mathbf{x}_n) \rightarrow \mathbf{R}'(\mathbf{x}_1 ... \mathbf{x}_n)] \vdash \mathbf{Q}'$ 

Apply Craig's interpolation algorithm to the proof to extract a reformulation.

This yields many prior results in the database literature.

## End-to-end reformulation algorithms via Craig's technique

#### Theorem [Levy Mendelzon Sagiv Srivistava PODS 95]:

If the constraints only say that each **R** in **T** is defined by a conjunctive view definition over another set of relations **B**: Compare the road example from the very beginning

 $\forall x_1 \dots x_n \mathbf{R}(x_1 \dots x_n) \leftrightarrow \mathbf{Q}_{\mathbf{R}}(\mathbf{x}_1 \dots \mathbf{x}_n),$ 

**Q**<sub>R</sub> a conjunctive query using a set of base relations **B** disjoint from **T**.

One can decide whether a query **Q** written over **B** can be rewritten as a conjunctive query using only the view relations.

Perform the chase effectively to determine if there is a proof that:  $\mathbf{Q} \land \Sigma \land \Sigma' \land [\land_{\mathcal{R} \in \mathcal{T}} \forall \mathbf{x}_1 \dots \mathbf{x}_n \mathbf{R}(\mathbf{x}_1 \dots \mathbf{x}_n) \rightarrow \mathbf{R}'(\mathbf{x}_1 \dots \mathbf{x}_n)] \vdash \mathbf{Q}'$ 

Apply Craig's interpolation algorithm to the proof to extract a reformulation.

## New end-to-end reformulation algorithms via Craig's technique

To find a FO reformulation (resp. positive existential, existential reformulation, RA-plan ... ) of query **Q** with respect to some set of TGD constraints  $\Sigma$  where the chase is well-behaved. e.g. Inclusion dependencies, LAV schema mappings, Guarded TGDs.

Show that **Q** is determined over **T** (resp. monotone overt **T**, induced-subinst. monotone, access-determined over the access methods. ..) with respect to  $\Sigma$ :

Effectively determine whether there is a chase proof that:  $Q \land \Sigma \land \Sigma \land \Lambda$  [......]  $\vdash Q'$ 

Apply an interpolation algorithm to the proof to extract a reformulation.

Gives first effective reformulation algorithms for very expressive constraint languages (Guarded TGDs, Guarded Fragment, ...)

# Craig gives completeness over finite instances

To find a first-order reformulation (resp. positive existential reformulation, RA-plan...) that is **equivalent to Q over every finite instance**, with respect to some "well-behaved" class of constraints  $\Sigma$ 

Show that **Q** is determined (resp. **monotone....**) with respect to  $\Sigma$ .:

Search for a proof effectively to determine if:  $\mathbf{Q} \land \Sigma \land \Sigma' \land [...] \vdash \mathbf{Q'}$ 

Apply Craig's interpolation algorithm to the proof to extract a reformulation

Use **finite controllability results** to argue that this process is complete for finding plans that work over all finite instances.

E.g. for Guarded TGDs, Guarded Fragment constraints, terminating chase classes....

## Craig can find interesting plans

To search for a **USPJ plan** using a fixed set of access methods equivalent to query **Q** with respect to constraints  $\varSigma$  .

Search for a proof that **Q** is **access-monotone** with respect to  $\Sigma$ : for any two instances I and I' satisfying the constraints  $\Sigma$  the accessible data in I is contained in the accessible data of I' then **Q**(I)  $\subseteq$  **Q**(I')

Do this by searching the space of proofs showing that:  $Q \land \Sigma \land \Sigma' \land$  ["Uni-directional Accessibility Axioms"]  $\vdash Q'$ 

The plans produced by applying Craig's interpolation algorithm to the proofs will include many natural optimizations used in plan search.

## Craig can find low-cost plans

To search for a **low-cost USPJ plan** using a fixed set of access methods equivalent to query **Q** with respect to constraints  $\varSigma$ .

Search for a proof that **Q** is **access-monotone** with respect to  $\Sigma$ : for any two instances I and I' satisfying the constraints  $\Sigma$  the accessible data in I is contained in the accessible data of I' then **Q**(I)  $\subseteq$  **Q**(I')

Do this by searching the space of proofs showing that:  $Q \land \Sigma \land \Sigma' \land$  ["Uni-directional Accessibility Axioms"]  $\vdash Q'$ 

While searching, apply Craig's interpolation algorithm to the proofs, and measure the cost of the resulting plan on-the-fly.



## Craig can find low-cost plans

See PDQ: Proof-Driven Querying over Web-based Datasources (B., Leblay, Tsamoura) in VLDB 2014.

API/examples; http://www.cs.ox.ac.uk/pdq/