

# Efficient Fixpoint Methods for Approximate Query Answering in Locally Complete Databases

Alvaro Cortés<sup>◇</sup> Marc Denecker<sup>◇</sup>  
Ofer Arieli\* Maurice Bruynooghe<sup>◇</sup>

◇: Katholieke Universiteit Leuven, Belgium

\*: The Academic College of Tel-Aviv/Yaffo, Israel

## Reiter's Closed World Assumption (1978)

In an arbitrary relational database, the Closed World Assumption (*CWA*) rules as *false* all those tuples that are not in the database.

**Train Time Table**

<i>Destination</i>	<i>Time</i>
Brussels Centraal	8:03
Antwerpen Centraal	8:05
Ghent St. Pieters	8:13
Brugge	8:22

Reiter's *CWA* allows us to conclude that there is no train to Hasselt at 8:04, for instance.

## Incomplete databases

The CWA applies to stand-alone databases storing *complete* (and correct) information about the world

What happens when the database is not complete?

<b>Telephone</b>		<b>Department</b>	
<i>Name</i>	<i>Telephone</i>	<i>Name</i>	<i>Department</i>
Leen Desmet	6531421	Bart Delvaux	Computer Sci.
Leen Desmet	09-23314	Leen Desmet	Philosophy
Bart Delvaux	5985625	David Finner	Computer Sci.

This database does not store *complete knowledge* about collaborators from the Philosophy department.

⇒ The *CWA* is not the correct approach in this context

## An alternative view: Open-World Assumption

The other extreme: The **Open-World Assumption (OWA)**

- The world can be in any state in which all database atoms are true
- A common approach in data integration systems
- The **OWA** is often too incomplete and underestimates the knowledge in a database.

How to identify those parts of the database that *are* complete?

Different approaches were presented to specify that the database is partially complete.

- First approach : [Motro 88]
- We follow the approach of Local Closed-World Assumption of [Levy96] and [CDAB05]

A specification of the “areas” in the real world in which the tables of the database contain *all* (true) tuples.

## Expressing Local Closed-World Assumptions

A *Local Closed-World Assumption* ( $\mathcal{LCWA}$ ) is an expression [CDAB05] (extended from [Levy96]):

$$\mathcal{LCWA}(P(\bar{x}), \Psi[\bar{x}])$$

Where:

- $\Psi[\bar{x}]$  : The window of expertise of the  $\mathcal{LCWA}$
- $P(\bar{x})$ : A database predicate, called the object of the  $\mathcal{LCWA}$
- $\Psi[\bar{x}]$ : The free variables in  $\Psi$  are a subset of  $\bar{x}$

”For all  $\bar{x}$  such that  $\Psi$  holds in the *real world*, if  $P(\bar{x})$  is true in the real world, then  $P(\bar{x})$  appears in the database”

## Example of the LCWA

Database knows about all the Telephone numbers of people in the CS (computer Science) department:

$$\mathcal{LCWA}(\text{Telephone}(x, y), \text{Dept}(x, \text{CS}))$$

- Windows of expertise:  $\text{Dept}(x, \text{CS})$
- Object of the LCWA:  $\text{Telephone}(x, y)$

”For all persons  $x$  of the department of computer science, all true facts of the form  $\text{Telephone}(x, y)$  appear in the database.”

## Semantics of the LCWA

Let  $D$  a set of ground atoms (a database) and the LCWA expression

$$\theta = \text{LCWA}(P(\bar{x}), \Psi[\bar{x}]),$$

the meaning of  $\theta$  given  $D$  is

$$\mathcal{M}_D(\theta) = \forall \bar{x} \left( \Psi[\bar{x}] \supset (P(\bar{x}) \equiv (P(\bar{x}) \in P^D)) \right)$$

Where “ $P(\bar{x}) \in P^D$ ” is a shorthand for

$$\bigvee_{\bar{a} \in P^D} (\bar{x} = \bar{a})$$

## Semantics of a Locally Complete Database

A locally closed database  $\mathcal{D}$  over  $\Sigma$  is pair  $(D, \mathcal{L})$ , where  $D$  is a database and  $\mathcal{L} = \{\theta_1, \dots, \theta_m\}$  is a finite set of LCWAs.

The *semantics* of  $\mathcal{D}$  is:

$$\mathcal{M}(\mathcal{D}) = \mathcal{A} \wedge \bigwedge_{j=1}^m \mathcal{M}_D(\theta_j) \wedge \text{UNA}(\Sigma) \wedge \text{DCA}(\Sigma).$$

Where

- $\mathcal{A}$ : The conjunction of atoms in  $D$ .
- $\text{UNA}(\Sigma)$ : Unique names axioms.
- $\text{DCA}(\Sigma)$ : Domain closure axioms.

## Query Answering

We are interested in evaluating queries  $Q$  with respect to  $\mathcal{M}(\mathfrak{D})$ :

- $\bar{t}$  is a *certain answer* for  $Q[\bar{x}]$  in  $\mathcal{M}(\mathfrak{D})$ , if

$$\mathcal{M}(\mathfrak{D}) \models Q[\bar{t}/\bar{x}].$$

Set of certain answers:  $Cert_{\mathfrak{D}}(Q[\bar{x}])$ .

- $\bar{t}$  is a *possible answer* for  $Q[\bar{x}]$  in  $\mathcal{M}(\mathfrak{D})$ , if

$$\mathcal{M}(\mathfrak{D}) \cup Q[\bar{t}/\bar{x}] \text{ is satisfiable.}$$

Set of possible answers:  $Poss_{\mathfrak{D}}(Q[\bar{x}])$ .

Negative complexity result [CDAB05]:

- Checking whether  $\bar{t}$  in  $Cert_{\mathfrak{D}}(Q[\bar{x}])$ : *coNP – complete*.
- $Poss_{\mathfrak{D}}(Q[\bar{x}])$  is *NP – complete*.

## Query Answering II

We present a tractable method for *approximate query answering*

- Under approximation of **certain answers**.
- Over approximation of **possible answers**.

The approach:

- Posing fixpoint queries that symbolically describe the construction of a 3-valued structure that approximates  $\mathcal{D}$ .
  - **Efficient** and **sound**.
  - In important cases, also **complete**.

## Constructing a 3-valued approximation of $\mathfrak{D} = (D, \mathcal{L})$

The operator  $App_{\mathfrak{D}} : \mathfrak{L}^c \rightarrow \mathfrak{L}^c$  maps a three-valued structure  $\mathcal{K}$  to a three-valued structure  $\mathcal{K}' = App_{\mathfrak{D}}(\mathcal{K})$  such that, for every predicate  $P$  of  $\mathcal{R}(\Sigma)$  and every tuple  $\bar{a}$ ,

$$P(\bar{a})^{\mathcal{K}'} = \begin{cases} \mathbf{t} & \text{if } P(\bar{a}) \in D, \\ \mathbf{f} & \text{if there exists } \mathcal{LCWA}(P(\bar{x}), \Psi_P[\bar{x}]) \in \mathcal{L} \text{ such that} \\ & \Psi_P[\bar{a}]^{\mathcal{K}} = \mathbf{t} \text{ and } P(\bar{a}) \notin D, \\ \mathbf{u} & \text{otherwise.} \end{cases}$$

Let  $\mathcal{C}_{\mathfrak{D}}$  be the  $\leq_p$ -least fixpoint of  $App_{\mathfrak{D}}$ .

- $\mathcal{C}_{\mathfrak{D}}$  'approximates'  $\mathcal{M}(\mathfrak{D})$ .
- $Cert_{\mathcal{C}_{\mathfrak{D}}}(Q[\bar{x}]) \subseteq Cert_{\mathfrak{D}}(Q[\bar{x}]) \subseteq Poss_{\mathfrak{D}}(Q[\bar{x}]) \subseteq Poss_{\mathcal{C}_{\mathfrak{D}}}(Q[\bar{x}])$ .

## Motivating Fixpoint Queries for the LCWA

$\mathcal{C}_D$  can be constructed in **polynomial time** in  $|D|$ .

But it needs to be recomputed each time the database changes...

- We partially avoid this by using **fixpoint formulas** that **symbolically** describe the construction of  $\mathcal{C}_D$ .
- Certain or possible answers to queries can be computed by transforming the original query into a **fixpoint/datalog query**.
- It suffices to **compute the relations** that are relevant for the **query** rather than computing all the relations in  $\mathcal{C}_D$ .

## Fixpoint Queries for the LCWA: Preliminaries

Given a database vocabulary  $\Sigma$ , we introduce, for each element in  $\mathcal{R}(\Sigma) = \{P_1, \dots, P_n\}$ , four new predicate symbols  $P_i^c, P_i^p, P_i^{c\bar{}}$  and  $P_i^{p\bar{}}$  of the same arity as  $P_i$ .

- $\Phi^c$  is the formula obtained by:
  - substituting  $P_i^c(\bar{t})$  for each **positive occurrence** of  $P_i(\bar{t})$  in  $\Phi$ ,
  - substituting  $\neg P_i^{c\bar{}}(\bar{t})$  for each **negative occurrence** of  $P_i(\bar{t})$  in  $\Phi$ .
  
- $\Phi^p$  is, inversely, the formula obtained by:
  - substituting  $P_i^p(\bar{t})$  for each **positive occurrence** of  $P_i(\bar{t})$  in  $\Phi$
  - substituting  $\neg P_i^{p\bar{}}(\bar{t})$  for each **negative occurrence** of  $P_i(\bar{t})$  in  $\Phi$ .

## Fixpoint Queries

Let  $\mathfrak{D} = (D, \mathcal{L})$  be a locally close database. For a query  $Q[\bar{x}]$  we introduce two new variables  $Q^c$  and  $Q^{c\bar{\neg}}$ , the arity of which is the number of free variables of  $Q[\bar{x}]$ , and define:

$$\Delta_{\mathcal{Q}, \mathcal{L}} = \left\{ \begin{array}{l} Q^c(\bar{x}) \leftarrow Q[\bar{x}]^c \\ Q^{c\bar{\neg}}(\bar{x}) \leftarrow (\neg Q[\bar{x}])^c \end{array} \right\} \cup \cup \left\{ \begin{array}{l} P_i^c(\bar{x}_i) \leftarrow P_i(\bar{x}_i) \\ P_i^{c\bar{\neg}}(\bar{x}_i) \leftarrow \neg P_i(\bar{x}_i) \wedge (\Psi_{P_i}[\bar{x}_i])^c \end{array} \right\},$$

where the right union is over the database predicates  $P_i$ , and  $\Psi_{P_i}$  is the window of expertise of  $P_i$ .

Certain answers for  $Q[\bar{x}]$ :  $[\mathbf{lfp}_{Q^c, \Delta_{\mathcal{Q}, \mathcal{L}}}] (\bar{x})$

Possible answers for  $Q[\bar{x}]$ :  $\neg[\mathbf{lfp}_{Q^{c\bar{\neg}}, \Delta_{\mathcal{Q}, \mathcal{L}}}] (\bar{x})$

$\Rightarrow$  Both of these expressions are evaluated in  $D$ .

## Fixpoint Queries: Example

Consider  $\mathcal{LCWA}(\text{Tel}(x, y), \text{Dept}(x, \text{CS}))$ . Assume that no closure exists for the Dept relation, i.e.  $\mathcal{LCWA}(\text{Dept}(x, y), \mathbf{f})$ .

Let  $Q = \text{Tel}(\text{BD}, 3962836)$

$$\Delta_{Q, \mathcal{L}} = \left\{ \begin{array}{l} Q^c \leftarrow \text{Tel}^c(\text{BD}, 3962836). \\ Q^{c\top} \leftarrow \text{Tel}^{c\top}(\text{BD}, 3962836). \\ \text{Tel}^c(x, y) \leftarrow \text{Tel}(x, y). \\ \text{Tel}^{c\top}(x, y) \leftarrow \neg \text{Tel}(x, y) \wedge \text{Dept}^c(x, \text{CS}). \\ \text{Dept}^c(x, y) \leftarrow \text{Dept}(x, y). \\ \text{Dept}^{c\top}(x, y) \leftarrow \neg \text{Dept}(x, y) \wedge \mathbf{f}. \end{array} \right.$$

As  $\text{lfp}_{Q^{c\top}, \Delta_{Q, \mathcal{L}}}$  is true in  $D$ ,  $\text{Tel}(\text{Bart Delvaux}, 3962836)$  is certainly false.

## Fixpoint queries: Soundness & completeness

Answers obtained by fixpoint queries are sound w.r.t. to  $\mathcal{C}_{\mathfrak{D}}$ .

Given a locally closed database  $(D, \mathcal{L})$  and a query  $Q[\bar{x}]$ .

Let  $(\mathcal{R}_Q^c, \mathcal{R}_Q^{c\bar{c}}, \mathcal{R}_1^c, \mathcal{R}_1^{c\bar{c}}, \dots, \mathcal{R}_n^c, \mathcal{R}_n^{c\bar{c}})$  be the relations defined by  $\Delta_{Q, \mathcal{L}}$  in  $D$ .

Then, for all  $1 \leq i \leq n$ ,

$$\mathcal{R}_i^c = \{\bar{d} \mid P_i(\bar{d})^{\mathcal{C}_{\mathfrak{D}}} = \mathbf{t}\}, \quad \mathcal{R}_i^{c\bar{c}} = \{\bar{d} \mid P_i(\bar{d})^{\mathcal{C}_{\mathfrak{D}}} = \mathbf{f}\},$$

$$\mathcal{R}_Q^c = \{\bar{d} \mid Q(\bar{d})^{\mathcal{C}_{\mathfrak{D}}} = \mathbf{t}\}, \quad \mathcal{R}_Q^{c\bar{c}} = \{\bar{d} \mid Q(\bar{d})^{\mathcal{C}_{\mathfrak{D}}} = \mathbf{f}\}.$$

Answers obtained by fixpoint queries are 'sound' and 'complete' w.r.t. to  $\mathcal{C}_{\mathfrak{D}}$ .

## Fixpoint queries: Soundness & completeness

But query answering in  $\mathcal{C}_{\mathcal{D}}$  is **not always complete** w.r.t  $\mathcal{M}(\mathcal{D})$ .

Let  $D = \emptyset$  and  $\mathcal{L} = \{\mathcal{LCWA}(P, R), \mathcal{LCWA}(Q, R \supset \neg P)\}$ .

In this case  $\mathcal{M}(\mathcal{D}) = (R \supset \neg P) \wedge ((R \supset \neg P) \supset \neg Q)$

**This theory entails  $\neg Q$ .**

The fact that in this case the window of expertise of the second LCWA is exactly the meaning of the first LCWA **is not captured** by  $\mathcal{C}_{\mathcal{D}}$ , and so  $Q^{\mathcal{C}_{\mathcal{D}}} = \mathbf{u}$ .

## Restricting the LCWA databases

The *LCWA dependency graph* of  $\mathcal{L}$  is the directed graph on  $\mathcal{R}(\Sigma)$ :

- A directed edge from predicate  $Q$  to  $P$  iff there exists  $LCWA(P(\bar{x}), \Psi[\bar{x}]) \in \mathcal{L}$  such that  $Q$  occurs negatively in  $\Psi$ .

A *hierarchically closed database*  $\mathfrak{D}$  is a locally closed database in which the LCWA dependency graph is *cycle-free*.

## Completeness

Let  $\mathfrak{D} = (D, \mathcal{L})$  be a hierarchically closed database such that every window of expertise in  $\mathcal{L}$  is a **conjunction of literals**.

If  $Q[\bar{x}]$  is a **conjunction of literals**, then

$$Cert_{\mathcal{C}_{\mathfrak{D}}}(Q[\bar{x}]) = Cert_{\mathfrak{D}}(Q[\bar{x}]).$$

If  $Q[\bar{x}]$  is a **disjunction of literals**, then

$$Poss_{\mathcal{C}_{\mathfrak{D}}}(Q[\bar{x}]) = Poss_{\mathfrak{D}}(Q[\bar{x}]).$$

## Conclusions and Ongoing Work

Tractable methods for approximate querying in locally closed databases, based on standard fixpoint techniques.

- Current research **refining** the class of LCWA for which the method is complete.
- Preliminary results that incorporate **integrity constraints and views**.  $\Rightarrow$  XSB implementation.
- Also, **safety issues**.