

Completeness Statements about RDF Data Sources and Their Use for Query Answering + Research Story

Fariz Darari

Werner Nutt
Giuseppe Pirrò
Simon Razniewski



FREIE UNIVERSITÄT BOZEN
LIBERA UNIVERSITÀ DI BOLZANO
FREE UNIVERSITY OF BOZEN · BOLZANO

EMCL Workshop 2014, Vienna

Slides about Research

In this color!

Slides about Research Story (Meta-Research)

In this color!

Story: Motivation

- I want to do a project and a thesis.
- I like Semantic Web.
- What can be a good research topic? Finding a good research topic (and good supervisor) is also part of research!
- Werner: “Hi, we are doing research about completeness reasoning on databases”
- Me: “Why not also on the Semantic Web?”

Motivation



IMDb > Reservoir Dogs



Full cast and crew for

Reservoir Dogs (1992) [More at IMDbPro »](http://www.imdb.com/title/tt0105236/fullcredits?ref_tt_ov_st_sm#cast)

http://www.imdb.com/title/tt0105236/fullcredits?ref_tt_ov_st_sm#cast



IMDbPro.com offers representation listings for over 120,000 individuals, including actors, directors, and producers, as well as company and employee contact details for over 50,000 companies in the entertainment industry.

[Click here for a free trial!](#)

Directed by

Quentin Tarantino

Writing credits

Quentin Tarantino (written by)

Roger Avary (background radio dialog) &

Quentin Tarantino (background radio dialog)

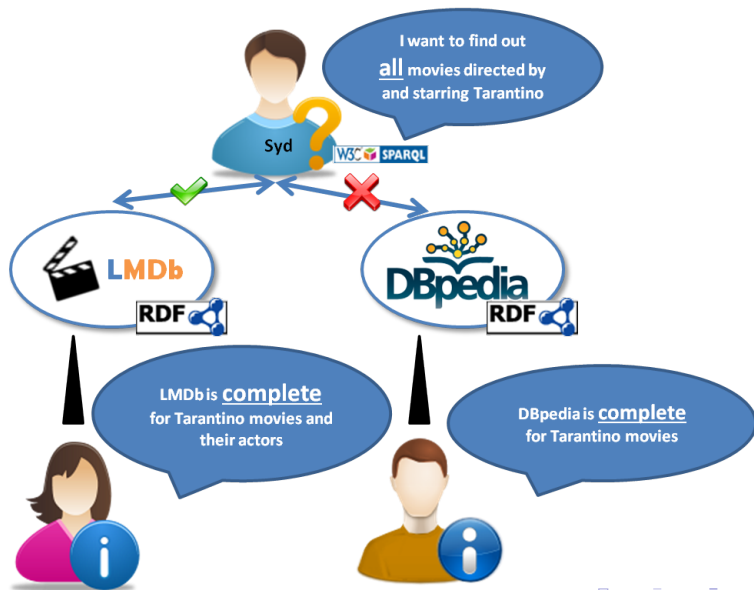
Cast (in credits order) verified as complete

	Harvey Keitel	...	Mr. White - Larry Dimmick
.....			
	Edward Bunker	...	Mr. Blue (as Eddie Bunker)
	Quentin Tarantino	...	Mr. Brown
.....			
.....			

Completeness statement about the IMDB data source

Quentin Tarantino was the character Mr. Brown

Motivation




Story: Literature Study

- Hmm..for realizing my goal, what should I learn?
- Aha, I need to learn Semantic Web, I think the Semantic Web Technologies course offered by FUB would be useful for me
- Aha, I also need to learn existing work on completeness reasoning for databases, I guess this paper¹ is worth to read!

¹Completeness of Queries over Incomplete Databases by Simon and Werner

Story: Google (and your supervisors) are your Googles, ask them!

completeness database 

Web Images Videos News Shopping More Search tools




About 15,000,000 results (0.27 seconds)

Database completeness - IBM
pic.dhe.ibm.com/infocenter/.../c_sysadm_db_completeness.html ▾ IBM ▾
Database completeness. The standard backup and restore by using the nzbackup and nzrestore commands provide transactionally consistent, automated ...

[PDF] Completeness of Queries over Incomplete Databases - VLDB ...
www.vldb.org/pvldb/vol4/p749-razniewski.pdf ▾
by S Razniewski - 2011 - Cited by 15 - Related articles
We develop techniques to conclude the **completeness** of query an- swers from information about the **completeness** of parts of a gen- erally incomplete database.

What is completeness constraint - Wiki Answers
wiki.answers.com > ... > Computer Programming > Database Programm
The **Completeness** Constraint addresses the issue of whether or not an ... concurrency control and why you think its important in **database** environme

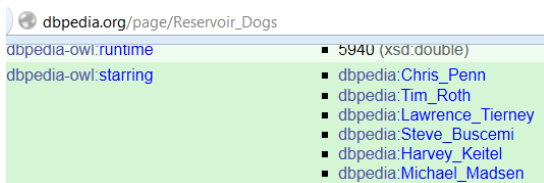
[PDF] Integrity = Validity + Completeness - Department of Com
cs.gmu.edu/~ami/research/.../tods89.pdf ▾ George Mason University ▾
by A MOTRO - 1989 - Cited by 164 - Related articles



Story: Start Doing Real Research – Problem Understanding

- Completeness reasoning on databases has been investigated
- Hmm..but **databases** are different than **Linked Data**, what should be adjusted?
- Well, in Linked Data, instead of **databases**, we have **RDF data sources**
- Well, in Linked Data, instead of **SQL**, we have **SPARQL** as a query language
- Well, Linked Data also is **more open, more heterogeneous and federated**
- Well, Linked Data also has **ontologies**
- Then, I have to discuss these issues with my supervisors!

Incomplete Data Source



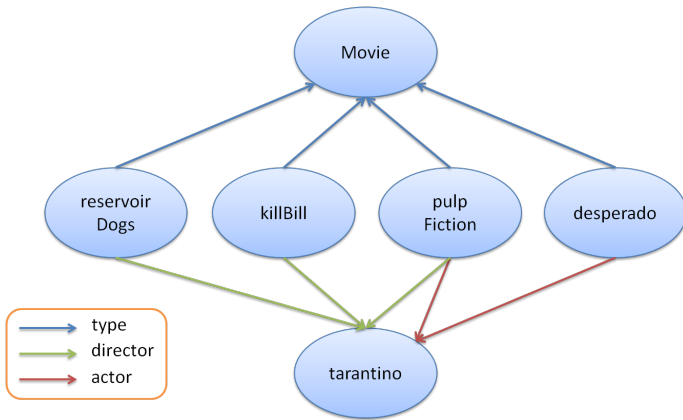
dbpedia.org/page/Reservoir_Dogs	
dbpedia-owl:runtime	▪ 5940 (xsd:double)
dbpedia-owl:starring	▪ dbpedia:Chris_Penn ▪ dbpedia:Tim_Roth ▪ dbpedia:Lawrence_Tierney ▪ dbpedia:Steve_Buscemi ▪ dbpedia:Harvey_Keitel ▪ dbpedia:Michael_Madsen

Quentin Tarantino is missing..

Incomplete Data Source

An incomplete data source of Tarantino movies,

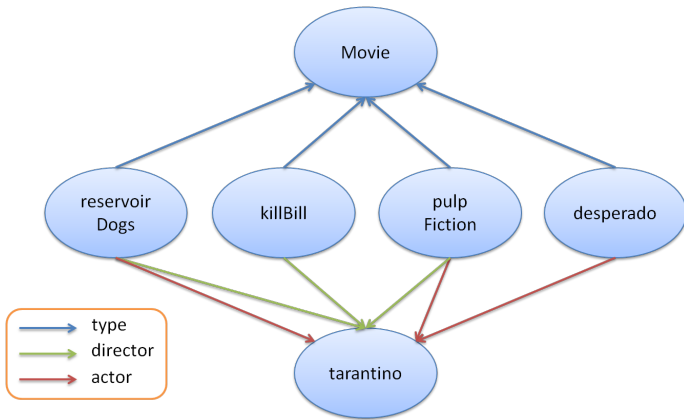
$$\mathcal{G}_{dbp} = (\mathcal{G}_{dbp}^a, \mathcal{G}_{dbp}^i):$$



Incomplete Data Source

An incomplete data source of Tarantino movies,

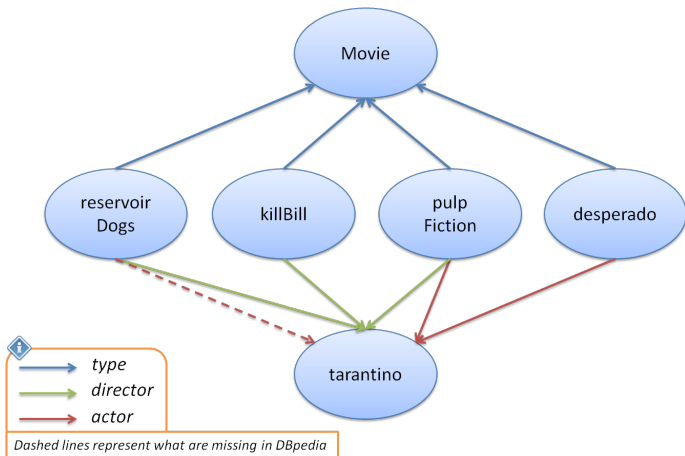
$$\mathcal{G}_{dbp} = (\mathcal{G}_{dbp}^a, \mathcal{G}_{dbp}^i):$$



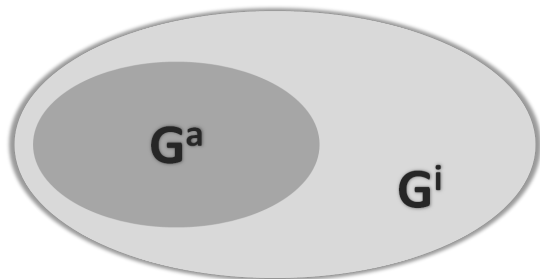
Incomplete Data Source

An incomplete data source of Tarantino movies,

$$\mathcal{G}_{dbp} = (\mathcal{G}_{dbp}^a, \mathcal{G}_{dbp}^i):$$



Incomplete Data Source



Incomplete Data Source

An **incomplete data source** is a pair of two graphs,

$$\mathcal{G} = (G^a, G^i), \text{ where } G^a \subseteq G^i.$$

We call G^a the **available** graph and G^i the **ideal** graph.

Completeness Statements: Examples

To express that a source is complete
for all the triples about movies directed by Tarantino,
we use the statement

$$C_{dir} = Compl((?m, type, Movie), (?m, director, tarantino) \mid \emptyset),$$

Completeness Statements: Examples

To express that a source is complete
for all the triples about movies directed by Tarantino,
we use the statement

$$C_{dir} = Compl((?m, type, Movie), (?m, director, tarantino) \mid \emptyset),$$

To express that a source is complete
for all triples about actors in movies directed by Tarantino,
we use

$$C_{act} =$$

$$Compl((?m, actor, ?a) \mid (?m, type, Movie), (?m, director, tara))$$

Completeness Statement: Definition

Let P_1 be a non-empty BGP (Basic Graph Pattern) and P_2 a BGP.

A **completeness statement** is defined as

$$\text{Compl}(P_1 \mid P_2)$$

where we call P_1 the **pattern** and P_2 the **condition** of the statement.

Satisfaction of Completeness Statements

To the statement

$$C = \text{Compl}(P_1 \mid P_2),$$

we associate the CONSTRUCT query

$$Q_C = (P_1, P_1 \cup P_2).$$

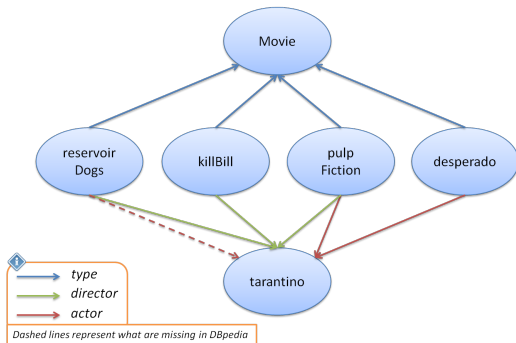
Then, we say:

C is **satisfied** by an IDS $\mathcal{G} = (G^a, G^i)$, written $\mathcal{G} \models C$, if

$$\llbracket Q_C \rrbracket_{G^i} \subseteq G^a.$$

Completeness Statements

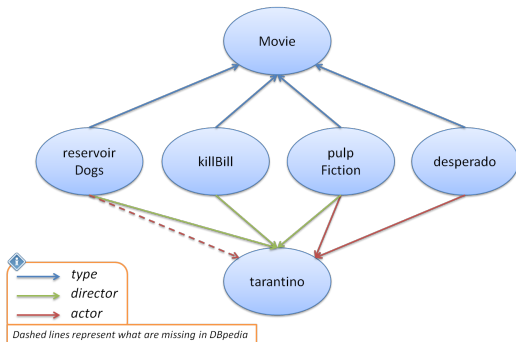
$$C_{dir} = Compl((?m, type, Movie), (?m, director, tarantino) \mid \emptyset)$$



Question: $\mathcal{G}_{dbp} \models C_{dir}$?

Completeness Statements

$$C_{dir} = Compl((?m, type, Movie), (?m, director, tarantino) \mid \emptyset)$$

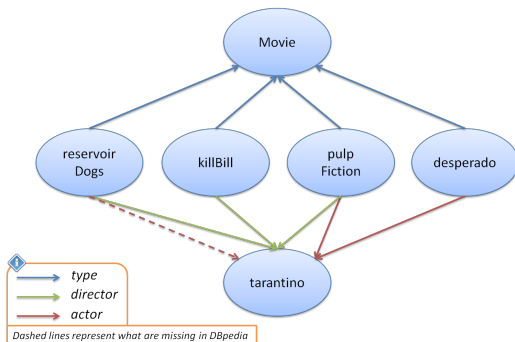


Question: $\mathcal{G}_{dbp} \models C_{dir}$?

Yes, because $\llbracket Q_{C_{dir}} \rrbracket \mathcal{G}_{dbp}^i \subseteq \mathcal{G}_{dbp}^a$.

Completeness Statements

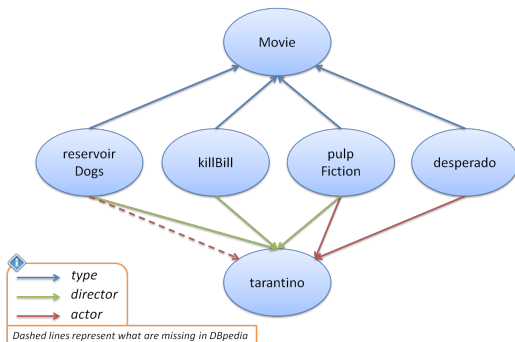
$$C_{act} = Compl((?m, actor, ?a) \mid (?m, type, Movie), (?m, director, tara$$



Question: $\mathcal{G}_{dbp} \models C_{act}$?

Completeness Statements

$C_{act} = Compl((?m, actor, ?a) \mid (?m, type, Movie), (?m, director, tara$

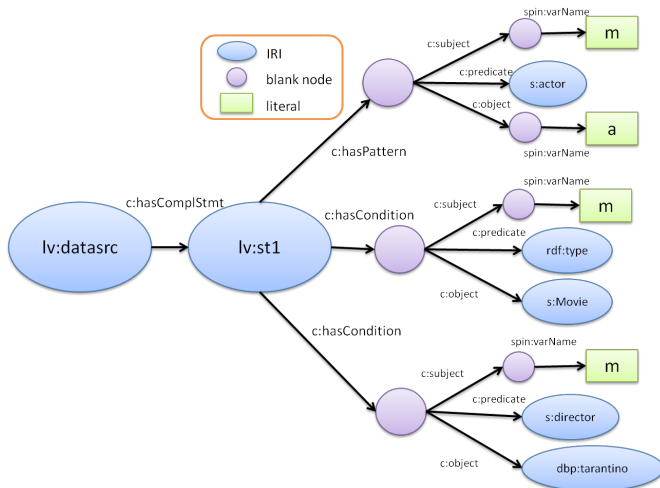


Question: $G_{dbp} \models C_{act}$?

No, because among the results of $\llbracket Q_{C_{act}} \rrbracket_{G_{dbp}^i}$, there is $(reservoirDogs, actor, tarantino)$ not in G_{dbp}^a .

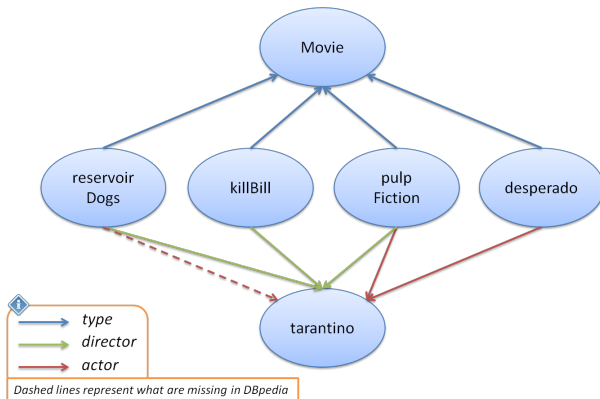
Completeness Statements in RDF

$$C_{act} = Compl((?m, actor, ?a) \mid (?m, type, Movie), (?m, director, tara))$$



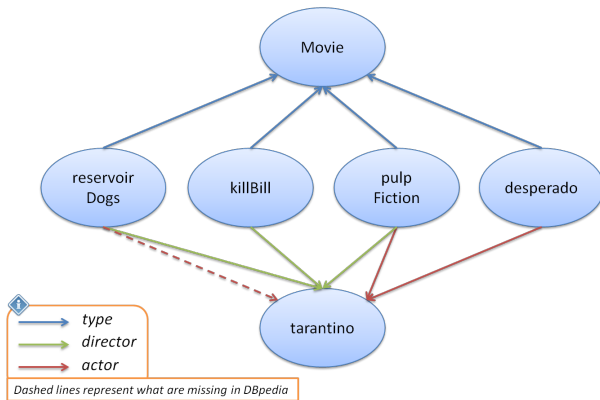
Query Completeness: Example

$$Q_{dir} = (\{?m\}, \{(?m, type, Movie), (?m, director, tarantino)\})$$



Query Completeness: Example

$$Q_{dir} = (\{?m\}, \{(?m, type, Movie), (?m, director, tarantino)\})$$



Then, $\llbracket Q_{dir} \rrbracket_{G_{dbp}^i} = \llbracket Q_{dir} \rrbracket_{G_{dbp}^a}$, and therefore $G_{dbp} \models Compl(Q_{dir})$.

Query Completeness: Definition

Definition

Let Q be a `SELECT` query. We write

$$\text{Compl}(Q)$$

to say that Q is **complete**.

An incomplete data source $\mathcal{G} = (G^a, G^i)$ **satisfies** $\text{Compl}(Q)$, written

$$\mathcal{G} \models \text{Compl}(Q) \quad \text{if and only if} \quad \llbracket Q \rrbracket_{G^i} = \llbracket Q \rrbracket_{G^a}.$$

Story: Defining Theorems to Prove

- Theorems are formal representations of goals you want to achieve (remember the first slides about motivation)
- In fact..all the definitions and framework introduced before are actually to well-define these theorems
- Hmm..I am sure these theorems should hold!
- Then, I have to discuss these issues with my supervisors!

Completeness Entailment

Let \mathcal{C} be a set of completeness statements and Q a `SELECT` query.

We say that \mathcal{C} **entails the completeness of Q** , written

$$\mathcal{C} \models \text{Compl}(Q),$$

if any incomplete data source satisfying \mathcal{C} also satisfies $\text{Compl}(Q)$.

Transfer Operator

For any set \mathcal{C} of completeness statements and a graph G , we define the **transfer operator** $T_{\mathcal{C}}$:

$$T_{\mathcal{C}}(G) = \bigcup_{\mathcal{C} \in \mathcal{C}} \llbracket Q_{\mathcal{C}} \rrbracket_G$$

Completeness Reasoning

Transfer Operator

For any set \mathcal{C} of completeness statements and a graph G , we define the **transfer operator** $T_{\mathcal{C}}$:

$$T_{\mathcal{C}}(G) = \bigcup_{C \in \mathcal{C}} \llbracket Q_C \rrbracket_G$$

Prototypical Graph

Let $Q = (W, P)$ be a `SELECT` query.

The **prototypical graph** \tilde{P} is the graph resulting from the mapping of variables in P to fresh, unique IRIs.

Completeness of Basic Queries

Theorem

Let \mathcal{C} be a set of completeness statements and $Q = (W, P)$ a basic query. Then,

$$\mathcal{C} \models \text{Compl}(Q) \quad \text{if and only if} \quad \tilde{P} = T_{\mathcal{C}}(\tilde{P}).$$

Story: Start Doing The Hardcore Part - Proving Theorems

- Argh, I got stuck
- This proving theorems stuff always haunts me!
- But:



Because you are Excellent Marvelous Champion
Legendary (aka EMCL)

Story: Start Doing The Hardcore Part - Proving Theorems (2)

- Get the proof idea first
- Generate some concrete examples of the problems with their solutions from the theorem you want to prove
- Discuss with your supervisors
- Now, start to prove in detail
- There is a problem on this, and that, then refine your proof idea, go back to the first point
- Start doing the above points until (hopefully :) you are able to prove it

Completeness Reasoning: Example

Consider the set of completeness statements

$$C_{dir,act} = \{ C_{dir}, C_{act} \}$$

and the query

$$Q_{dir+act} = (\{ ?m \}, P_{dir+act})$$

where

$$P_{dir+act} = \{ (?m, type, Movie), (?m, director, tarantino), (?m, actor, tarantino) \}$$

Completeness Reasoning: Example

Consider the set of completeness statements

$$C_{dir,act} = \{ C_{dir}, C_{act} \}$$

and the query

$$Q_{dir+act} = (\{ ?m \}, P_{dir+act})$$

Then,

$$\tilde{P}_{dir+act} = \{ (C_{?m}, type, Movie), (C_{?m}, director, tarantino), (C_{?m}, actor, tarantino) \}$$

Completeness Reasoning: Example

$$\mathcal{C}_{dir,act} = \{ \mathbf{C}_{dir}, \mathbf{C}_{act} \}$$

$$\mathbf{Q}_{dir+act} = (\{ ?m \}, \mathbf{P}_{dir+act})$$

Therefore,

$$\mathcal{T}_{\mathcal{C}_{dir,act}}(\tilde{\mathbf{P}}_{dir+act})$$

Completeness Reasoning: Example

$$\mathcal{C}_{dir,act} = \{ \mathbf{C}_{dir}, \mathbf{C}_{act} \}$$

$$Q_{dir+act} = (\{ ?m \}, P_{dir+act})$$

Therefore,

$$\begin{aligned} & \mathcal{T}_{\mathcal{C}_{dir,act}}(\tilde{P}_{dir+act}) \\ = & \llbracket Q_{\mathbf{C}_{dir}} \rrbracket_{\tilde{P}_{dir+act}} \cup \llbracket Q_{\mathbf{C}_{act}} \rrbracket_{\tilde{P}_{dir+act}} \end{aligned}$$

Completeness Reasoning: Example

$$\mathcal{C}_{dir,act} = \{ \mathbf{C}_{dir}, \mathbf{C}_{act} \}$$

$$Q_{dir+act} = (\{ ?m \}, P_{dir+act})$$

Therefore,

$$\begin{aligned} & \mathcal{T}_{\mathcal{C}_{dir,act}}(\tilde{P}_{dir+act}) \\ &= \llbracket Q_{\mathbf{C}_{dir}} \rrbracket \tilde{P}_{dir+act} \cup \llbracket Q_{\mathbf{C}_{act}} \rrbracket \tilde{P}_{dir+act} \\ &= \{ (\mathbf{c}_{?m}, type, Movie), (\mathbf{c}_{?m}, director, tara), \end{aligned}$$

Completeness Reasoning: Example

$$\mathcal{C}_{dir,act} = \{ \mathcal{C}_{dir}, \mathcal{C}_{act} \}$$

$$Q_{dir+act} = (\{ ?m \}, P_{dir+act})$$

Therefore,

$$\begin{aligned} & \mathcal{T}_{\mathcal{C}_{dir,act}}(\tilde{P}_{dir+act}) \\ &= \llbracket Q_{\mathcal{C}_{dir}} \rrbracket \tilde{P}_{dir+act} \cup \llbracket Q_{\mathcal{C}_{act}} \rrbracket \tilde{P}_{dir+act} \\ &= \{ (c_{?m}, type, Movie), (c_{?m}, director, tara), (c_{?m}, actor, tara) \} \end{aligned}$$

Completeness Reasoning: Example

$$\mathcal{C}_{dir,act} = \{ \mathcal{C}_{dir}, \mathcal{C}_{act} \}$$

$$Q_{dir+act} = (\{ ?m \}, P_{dir+act})$$

Therefore,

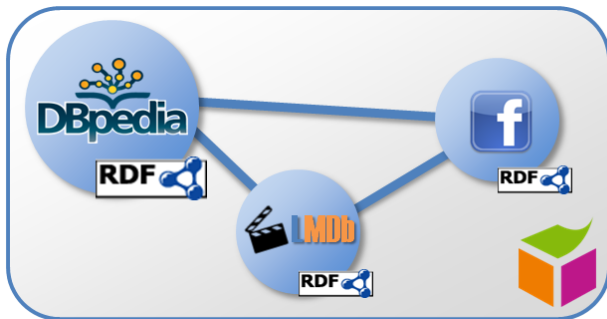
$$\begin{aligned} & \mathcal{T}_{\mathcal{C}_{dir,act}}(\tilde{P}_{dir+act}) \\ &= \llbracket Q_{\mathcal{C}_{dir}} \rrbracket \tilde{P}_{dir+act} \cup \llbracket Q_{\mathcal{C}_{act}} \rrbracket \tilde{P}_{dir+act} \\ &= \{ (c_{?m}, type, Movie), (c_{?m}, director, tara), (c_{?m}, actor, tara) \} \\ &= \tilde{P}_{dir+act}. \end{aligned}$$

Thus, $\mathcal{C}_{dir,act} \models Compl(Q_{dir+act})$

The framework can also be applied to:

- DISTINCT Queries: with set semantics
- OPT Queries: eg. get all people and in case they are not single, get also their spouse
- Queries with RDFS Data Sources: incorporating ontologies

Federated Framework



Can I get a complete answer?

```
SELECT ?m ?l
WHERE { ?m rdf:type s:Movie .
        ?m s:director dbp:tarantino .
        ?m fbo:likes ?l }
```

CoRner: Completeness Reasoner

- Can check the completeness of a subset of SPARQL queries depending on the completeness statements of a single data source.
- Developed in Java using the Apache Jena library.
- Takes three inputs:
 - Completeness statements in RDF format
 - A SPARQL query
 - (optional) an RDFS ontology
- Available at <http://rdcorner.wordpress.com/>.

Story: Cooling Down, or is it?

- I have got all the results
- But, I haven't started writing the thesis yet
- Nooo, I will miss the deadline :(

Story: Cooling Down - Alternative

- I have got all the results
- And, I have started writing the thesis since the very beginning (right after I have studied the literature)
- I think I can manage to meet the deadline :)
- Aha, I think our research results can be also useful if shown to the Semantic Web community, why not to try submit it at ISWC 2013? – Another research story :)

Conclusions and Future Work

- We developed a theoretical framework for the expression of completeness statements about data sources, from which we can ensure query completeness.
- We provide a compact RDF representation of completeness statements, which can be embedded in `VOID` descriptions, and implemented the framework in `CoRner`.
- We are interested in studying completeness reasoning with more expressive queries and `OWL 2`.

Terima kasih, grazie, danke!!



[Link to References](#)