#### Query-driven Data Completeness Assessment

Simon Razniewski Free University of Bozen-Bolzano, Italy

Part I joint work with Flip Korn, Werner Nutt, Divesh Srivastava

# Background

Freie Universität Bozen

Libera Università di Bolzano

Università Liedia de Bulsan

unibz

- 2011 2014: PhD (on reasoning about data completeness)
- 2014 now: Assistant professor
- Research visits at UCSD (2012), AT&T Labs-Research (2013), UQ (2015), MPII (now)



# Background (2)

• Research centered on data completeness

- Completeness often a problem in
  - Data integration
  - Complex data-generating processes
  - Large data/knowledge bases

# Outline

#### Part I: Completeness reasoning in databases

#### Part II: Assessing completeness of general-purpose knowledge bases (=Wikidata, Google Knowledge Graph, YAGO, ..)

# Part I: Reasoning in databases

#### • Make data more complete

- Missing value imputation
- Information extraction
- Data fusion
- ...

#### • Reason about (in-)completeness information

- Missing records in relational databases (VLDB 2011)
- Null values (CIKM 2012)
- RDF databases (ISWC 2013)
- Derivations from process states (BPM 2013)
- Spatial data (SIGSPATIAL 2014)
- An algebra for completeness information (SIGMOD 2015)

# Motivation: Data warehouse of a telecommunication company



#### Admin John knows

- Team table is complete (HR says so)
- Maintenance is complete for teams A, B and C
  - their reporting systems export data automatically
- Warnings is complete for all of Week 1, and Monday and Wednesday of Week 2
  - Potential data loss due to a system failure on Tuesday
  - Data further than Wednesday maybe not fully loaded

| Teams |                |  |  |  |  |
|-------|----------------|--|--|--|--|
| name  | specialization |  |  |  |  |
| А     | hardware       |  |  |  |  |
| В     | hardware       |  |  |  |  |
| С     | network        |  |  |  |  |
| С     | software       |  |  |  |  |
| D     | network        |  |  |  |  |

| Maintenance |      |                |  |  |  |  |
|-------------|------|----------------|--|--|--|--|
| ID          | resp | reason         |  |  |  |  |
| tw37        | А    | disk failure   |  |  |  |  |
| tw59        | D    | software crash |  |  |  |  |
| tw83        | В    | unknown        |  |  |  |  |
| tw91        | С    | update failure |  |  |  |  |
| tw91        | С    | network error  |  |  |  |  |

| Warnings |      |      |              |  |  |  |
|----------|------|------|--------------|--|--|--|
| day      | week | ID   | message      |  |  |  |
| Mon      | 1    | tw37 | high voltage |  |  |  |
| Fri      | 1    | tw37 | high voltage |  |  |  |
| Wed      | 2    | tw37 | overheat     |  |  |  |
| Tue      | 1    | tw59 | auto restart |  |  |  |
| Fri      | 1    | tw59 | overheat     |  |  |  |
| Mon      | 2    | tw83 | high voltage |  |  |  |
| Tue      | 2    | tw83 | auto restant |  |  |  |

# John wants to know



| W.Day | W.week | W.ID | W.message    | M.ID | M.resp | M.reason     | T.name | T.specialization |
|-------|--------|------|--------------|------|--------|--------------|--------|------------------|
| Wed   | 2      | tw37 | overheat     | tw37 | А      | disk failure | А      | hardware         |
| Mon   | 2      | tw83 | high voltage | tw83 | В      | unknown      | В      | hardware         |
| Tue   | 2      | tw83 | auto restart | tw83 | В      | unknown      | В      | hardware         |

# John reasons



"Give me all warnings in Week 2 that are generated by objects in maintenance with a hardware team."

- Warnings is complete for Week 1 and Monday and Wednesday of Week 2
- Maintenance is complete for teams A, B and C
- Team is complete

ightarrow The query result definitely contains all warnings from

- Monday for team A
- Monday for team B
- Monday for team C
- Wednesday for team A
- Wednesday for team B
- Wednesday for team C

# John looks at the data

"Give me all warnings in Week 2 that are generated by objects in maintenance with a **hardware team**."

ightarrow The query result definitely contains all warnings from

- Monday for team A
  Monday for team B
  - = All data from Monday
- Monday for team C
- Wednesday for team A
- Wednesday for team B
- Wednesday for team C
- ] = All data from Wednesday

| Teams |                |  |  |  |
|-------|----------------|--|--|--|
| name  | specialization |  |  |  |
| Α     | hardware       |  |  |  |
| В     | hardware       |  |  |  |
| С     | network        |  |  |  |
| С     | software       |  |  |  |
| D     | network        |  |  |  |

• HR says: The team table is complete!

# Problems

"Warnings are complete for Week 1"

(Input)

1. How can we formally describe complete parts of a database?

"The query result contains all warnings from Monday of Week 2 for Team A" (Output)

2. How can we use database completeness information to identify complete parts of query answers?

#### Formalism: Patterns

| We have all warnings from Week 1        | Warnings |      |      |              |  |
|---|----------|------|------|--------------|--|
|   | day      | week | ID   | message      |  |
| We have all warpings from               | Mon      | 1    | tw37 | high voltage |  |
| Monday of Week 2                        | Fri      | 1    | tw37 | high voltage |  |
|   | Wed      | 2    | tw37 | overheat     |  |
|   | Tue      | 1    | tw59 | auto restart |  |
|   | Fri      | 1    | tw59 | overheat     |  |
|   | Mon      | 2    | tw83 | high voltage |  |
|   | Tue      | 2    | tw83 | auto restart |  |
| ``````````````````````````````````````` | *        | 1    | *    | *            |  |
|   | Mon      | 2    | *    | *            |  |

- Less expressive than previously known formalisms (views, Datalog/first-order queries, ..)
- Can be expressed in the same schema as the data
- Efficient reasoning

#### John's knowledge expressed by patterns

#### Team table is complete

Maintenance is complete for teams A, B and C Warnings is complete for all of Week 1, and Monday and Wednesday of Week 2

| Teams |                |  |  |  |  |
|-------|----------------|--|--|--|--|
| name  | specialization |  |  |  |  |
| А     | hardware       |  |  |  |  |
| В     | hardware       |  |  |  |  |
| С     | network        |  |  |  |  |
| С     | software       |  |  |  |  |
| D     | network        |  |  |  |  |
| *     | *              |  |  |  |  |

| Maintenance |      |                |  |  |  |  |
|-------------|------|----------------|--|--|--|--|
| ID          | resp | reason         |  |  |  |  |
| tw37        | А    | disk failure   |  |  |  |  |
| tw59        | D    | software crash |  |  |  |  |
| tw83        | В    | unknown        |  |  |  |  |
| tw91        | С    | update failure |  |  |  |  |
| tw91        | С    | network error  |  |  |  |  |
| *           | А    | *              |  |  |  |  |
| *           | В    | *              |  |  |  |  |
| *           | С    | *              |  |  |  |  |

| Warnings |      |      |              |  |  |  |
|----------|------|------|--------------|--|--|--|
| day      | week | ID   | message      |  |  |  |
| Mon      | 1    | tw37 | high voltage |  |  |  |
| Fri      | 1    | tw37 | high voltage |  |  |  |
| Wed      | 2    | tw37 | overheat     |  |  |  |
| Tue      | 1    | tw59 | auto restart |  |  |  |
| Fri      | 1    | tw59 | overheat     |  |  |  |
| Mon      | 2    | tw83 | high voltage |  |  |  |
| Tue      | 2    | tw83 | auto restart |  |  |  |
| *        | 1    | *    | *            |  |  |  |
| Mon      | 2    | *    | *            |  |  |  |
| Wed      | 2    | *    | *            |  |  |  |

#### John's conclusions expressed by patterns

#### "Give me all warnings in week 2 that are generated by objects in maintenance with a hardware team."

| W.Day | W.week | W.ID | W.message    | M.ID | M.resp | M.reason     | T.name | T.specialization |
|-------|--------|------|--------------|------|--------|--------------|--------|------------------|
| Wed   | 2      | tw37 | overheat     | tw37 | А      | disk failure | А      | hardware         |
| Mon   | 2      | tw83 | high voltage | tw83 | В      | unknown      | В      | hardware         |
| Tue   | 2      | tw83 | auto restart | tw83 | В      | unknown      | В      | hardware         |
| Mon   | *      | *    | *            | *    | А      | *            | А      | *                |
| Mon   | *      | *    | *            | *    | В      | *            | В      | *                |
| Mon   | *      | *    | *            | *    | С      | *            | С      | *                |
| Wed   | *      | *    | *            | *    | A      | *            | A      | *                |
| Wed   | *      | *    | *            | *    | В      | *            | В      | *                |
| Wed   | *      | *    | *            | *    | С      | *            | С      | *                |

ightarrow The query result contains all warnings from

• Monday for team A

•

How to compute the completeness patterns for queries?

Queries are computed by relational algebra Here: Select, project, equijoin



 Apply algebra operators to completeness patterns (analogous to query result computation)

#### Reasoning about selections



Rule 1: Statements with \* survive



#### Reasoning about selections (2)

|     | warnings |       |              |
|-----|----------|-------|--------------|
| day | week     | ID    | message      |
| Mon | 1        | tw37  | high voltage |
| Fri | 1        | tw37  | high voltage |
| Wed | 2        | tw37  | overheat     |
| Tue | 1        | tw59  | auto restart |
| Fri | 1        | tw/59 | overheat     |
| Man | -<br>-   | tw02  | high voltage |
| won | Z        | tw83  | nign voltage |
| Tue | 2        | tw83  | auto restart |
| *   | 1        | *     | *            |
| Mon | 2        | *     | *            |
| Wed | 2        | *     | *            |

Rule 2: Irrelevant constants are ignored

Rule 3: Selected constants survive and are promoted



#### Reasoning about joins

|      | Maintenance |                |  |  |  |  |
|------|-------------|----------------|--|--|--|--|
| ID   | resp        | reason         |  |  |  |  |
| tw37 | А           | disk failure   |  |  |  |  |
| tw59 | D           | software crash |  |  |  |  |
| tw83 | В           | unknown        |  |  |  |  |
| tw91 | С           | update failure |  |  |  |  |
| tw91 | С           | network error  |  |  |  |  |
| *    | А           | *              |  |  |  |  |
| *    | В           | *              |  |  |  |  |
| *    | С           | *              |  |  |  |  |

| M.ID | M.resp | M.reason     | T.name | T.specialization |
|------|--------|--------------|--------|------------------|
| tw37 | А      | disk failure | А      | hardware         |
| tw83 | В      | unknown      | В      | hardware         |
| *    | А      | *            | *      | *                |
| *    | В      | *            | *      | *                |
| *    | С      | *            | *      | *                |
| *    | *      | *            | А      | *                |
| *    | *      | *            | В      | *                |
| *    | *      | *            | С      | *                |

| M.resp=          | T.name      |
|------------------|-------------|
| T.specialization |             |
| hardware         | Rule 1. Con |
|                  |             |

| $\sigma_{spec="hw"}(T)$ |          |  |  |  |
|-------------------------|----------|--|--|--|
| name specializatio      |          |  |  |  |
| А                       | hardware |  |  |  |
| В                       | hardware |  |  |  |
| *                       | *        |  |  |  |

Rule 1: Constants join with equal constants Rule 2: Wildcards join with anything Rule 3: Constants can be promoted

# Algorithmic completeness

#### **Proven**: Extended algebra gives all conclusions that hold on the schema level (reasoning only with the yellow metadata)

• Independent of the algebra tree chosen

## Looking at the data



# So much about the theory, but...

1. How can we implement this?

- 2. How fast is this?
  - In comparison with query evaluation

3. How can we manage large sets of statements?

# How can we implement this?

• Ideally, a plugin inside a DBMS

Promotion procedure benefits from fast access to data

• So far: Separate Java tool

Schema-level algebra can also be encoded in SQL
 → Could compile normal queries into metadata queries

#### How fast is this? (1)

- Synthetic data
- Wikipedia has around 1000 lists declared as complete (using a template or in natural language)

| This is a complete list of the 72 communities in Carmarthenshire. <sup>[Community 1]</sup> |                                  |                                  |  |  |  |
|--|----------------------------------|----------------------------------|--|--|--|
| Abergwili  | <ul> <li>Llanarthney</li> </ul>  | <ul> <li>Llangadog</li> </ul>    |  |  |  |
| Abernant   | <ul> <li>Llanboidy</li> </ul>    | <ul> <li>Llangain</li> </ul>     |  |  |  |
| <ul> <li>Ammanford</li> </ul>  | <ul> <li>Llanddarog</li> </ul>   | <ul> <li>Llangathen</li> </ul>   |  |  |  |
| Betws  | <ul> <li>Llanddeusant</li> </ul> | <ul> <li>Llangeler</li> </ul>    |  |  |  |
| Bronwydd   | <ul> <li>Llanddowror</li> </ul>  | <ul> <li>Llangennech</li> </ul>  |  |  |  |
| Carmarthen   | Llandeilo                        | <ul> <li>Llangunnor</li> </ul>   |  |  |  |
| Cenarth  | <ul> <li>Llandovery</li> </ul>   | <ul> <li>Llangyndeyrn</li> </ul> |  |  |  |
| Cilycwm  | <ul> <li>Llandybie</li> </ul>    | <ul> <li>Llangynin</li> </ul>    |  |  |  |
| - Cilumooplluud  | - Llondyfoolog                   | Llongunog                        |  |  |  |

http://en.wikipedia.org/wiki/List\_of\_places\_in\_Carmarthenshire\_%28categorised%29

## How fast is this? (2)

- Manually extracted some and grouped them by topic
  - Recurrent topics: Sports teams, political assemblies, geographical features, songs, operas and other pieces of art
- Generated one table each about cities, schools and countries (21 statements)

| city |          |                  |                  |  |  |
|------|----------|------------------|------------------|--|--|
| name | country  | state            | county           |  |  |
| *    | USA      | Virginia         | *                |  |  |
| *    | Germany  | *                | *                |  |  |
| *    | Ukraine  | *                | *                |  |  |
| *    | Bulgaria | *                | *                |  |  |
| *    | USA      | New York         | *                |  |  |
| *    | UK       | Carmarthenshire  | *                |  |  |
| *    | USA      | West Virginia    | Hampshire County |  |  |
| *    | Czech    | Moravian-Silesia | Nový Jičín       |  |  |
| *    | Slovenia | *                | *                |  |  |

### How fast is this? (3)

SELECT \*
FROM country, city, school
WHERE country.capital=city.name
AND city.state=school.state

"All schools in capital states"

SQL runtime: 2 seconds (25891 records) Completeness pattern runtime: 0.9 seconds (46 patterns)

Median over 7 join queries:

- SQL runtime: 2 seconds
- Completeness pattern runtime: 0.5 seconds

#### How to manage large sets of patterns?

Redundancies in workflows may lead to redundant patterns

John reports first that all data for Monday of Week 2 is complete, later, that the data for the whole Week 2 is complete

(Monday, 2) (\*, 2)

Redundancies introduce overhead and restrict comprehensibility

 $\rightarrow$  Should be identified and possibly removed

#### Trivial?

- (Monday, \*, hardware)
   (Wedr

   (Tuesday, 2, software)
   (\*, \*, h

   (Monday, 2, \*)
   (\*, 2, s)
  - (Wednesday, \*, software) (\*, \*, hardware) (\*, 2, software)

#### How to manage large sets of patterns? (2)

- Similar problems occur in term indexing for unification in classic AI
- Compared proposed data structures with classic hashing
- Time/space tradeoff:
  - Discrimination trees are most time-efficient (60 seconds for minimizing 800k patterns)
  - Hashing is marginally better in terms of space

# Summary Part I

- Introduced completeness patterns to describe complete parts of databases and query answers
  - Can be expressed in the same schema as the data
- Modified the relational algebra to manipulate completeness patterns
  - Implemented join operator and evaluated scalability
- Limitation: No complete algorithm for data-dependent promotion

#### PART II: ASSESSING COMPLETENESS OF GENERAL-PURPOSE KNOWLEDGE BASES

(=Wikidata, Google Knowledge Graph, YAGO, ..)

#### How complete are knowledge bases?

# KBs are pretty incomplete

DBpedia: contains 6 out of 35 Dijkstra Prize winners 😕





YAGO: the average number of children per person is **0.02** ⊗



Google Knowledge Graph: ``Points of Interest'' – Completeness? 😕

## KBs are pretty complete



DBpedia: 167 out of 199 Nobel laureates in Physics 🙂

| All     Images     Values     News     Elberger     News     Exception       Constant Formation of Movies       Constant formation of Movies <td co<="" colspan="4" th=""><th>G<mark>oogle</mark></th><th colspan="7">tarantino movies</th></td> | <th>G<mark>oogle</mark></th> <th colspan="7">tarantino movies</th> |   |  |  | G <mark>oogle</mark> | tarantino movies |  |  |  |  |  |  |
|---|--|---|--|--|----------------------|------------------|--|--|--|--|--|--|
| Cuentin Tarantho / Movies   |  | All Images Videos News Shopping More + Search tools |  |  |                      |                  |  |  |  |  |  |  |
| - F 🔂 🖄 🐝 🗍 🔭   |  | Quentin Tarantino / Movies                          |  |  |                      |                  |  |  |  |  |  |  |
| 7 📷 🎑 🔛 🎆 🚹   |  |   |  |  |                      |                  |  |  |  |  |  |  |
|   |  |   |  |  |                      |                  |  |  |  |  |  |  |
|   |  |   |  |  |                      |                  |  |  |  |  |  |  |
|   |  |   |  |  |                      |                  |  |  |  |  |  |  |

| child | Malia Obama |                                  |
|-------|-------------|----------------------------------|
|       |             | ▼ 0 references                   |
|       |             | Sasha Obama                      |
|       |             | <ul> <li>0 references</li> </ul> |

Wikidata: 2 out of 2 children of Obama 😳

Google Knowledge Graph: 36 out of 48 Tarantino movies 🙂

# So, how complete are KBs?



| Name                 | # Entity types | # Entity instances | # Relation types | # Confident facts (relation instances) |
|----------------------|----------------|--------------------|------------------|--|
| Knowledge Vault (KV) | 1100           | $45\mathrm{M}$     | 4469             | $271\mathrm{M}$                        |
| DeepDive [32]        | 4              | $2.7\mathrm{M}$    | 34               | $7 \mathrm{M}^{a}$                     |
| NELL [8]             | 271            | $5.19\mathrm{M}$   | 306              | $0.435 \mathrm{M}^{b}$                 |
| PROSPERA [30]        | 11             | N/A                | 14               | $0.1\mathrm{M}$                        |
| YAGO2 [19]           | 350,000        | $9.8\mathrm{M}$    | 100              | $4\mathrm{M}^{c}$                      |
| Freebase [4]         | 1,500          | $40\mathrm{M}$     | 35,000           | $637 \mathrm{M}^{d}$                   |
| Knowledge Graph (KG) | 1,500          | $570\mathrm{M}$    | 35,000           | $18,000\mathrm{M}^e$                   |

KB engineers have only tried to make KBs bigger. The point, however is to **understand** what they are actually trying to approximate.

There are **known knowns**; there are things we know we know. We also know there are **known unknowns**; that is to say we know there are some things we do not know. But there are also **unknown unknowns** – the ones we don't know we don't know.

#### [Dong et al., KDD 2014]



#### http://www.how-complete-are-kbs.today



# How can we get there?

- Find patterns in data that tell about completeness
  - Extended AMIE system to mine rules about completeness
     [WSDM 2017]

```
hockeyPlayer(x) \rightarrow Incomplete(x, hasChild)
scientist(x), hasWonNobelPrize(x) \rightarrow Complete(x, graduatedFrom)
```

- Extract cardinality information from texts
  - Manually created patterns that allow to find information about 178% more children than Wikidata currently contains [ISWC 2016 Poster]

"John lives with his spouse and **five** children on a farm in Alabama."

#### Are we done soon?

#### Challenge 1 – Modelling incompleteness

**Google** Knowledge Graph *knows 6 out of 7 children of President Garfield* 

• RDF blank node for the 7<sup>th</sup> child?

→No identity, thus no way to say that the 7<sup>th</sup> child is different from the previous 6

• Creating a new nameless entity?

→Not known that the 7<sup>th</sup> child is different from all other entities in the KB

Wikidata: New property nr\_of\_children(Garfield, 7)?
 → Semantic relation between ``hasChild'' and ``nr\_of\_children'' lost

# Challenge 2 – Beyond triples

These are facts for humans:

Galileo Galilei:

Contrary to the dogma of the time, postulates that the earthorbits the sun[http://discovermagazine.com/2007/jul/20-things-you-didn2019t-know-about-galileo]

Reinhold Messner:

 First person to climb all mountains >8000mt without

 supplemental oxygen
 [http://www.telegraph.co.uk/travel/lists/reinhold-messner-tribute-quotes-facts/]

Toothbrush? MP3-Player?

These are not KB triples

FirstPersonToClimbAllMountainsAbove8000Without(Supplemental oxygen, Reinhold Messner)

 $\rightarrow$  How incomplete are KB models?

#### Challenge 3 – Aggregating completeness

Can we say whether a KB knows more about Obama than about Trump?

Or about Ronaldo than about Justin Bieber?







 $\rightarrow$  Need to understand relevance, interestingness

https://www.wikidata.org/wiki/User:Ls1g/Recoin

# Outlook

- 1. KBs contain a lot of knowledge, but little is known about how much they actually know
- 2. Completeness can be assessed from within KBs using pattern mining, or using external sources
- 3. Big challenges ahead
  - How to model incompleteness in KBs?
  - How incomplete are KB models?
  - How to aggregate completeness?