Database Completeness Reasoning and Completeness Verification over Processes



Simon Razniewski Joint work with Marco Montali and Werner Nutt Free University of Bozen-Bolzano

Free University of Bozen-Bolzano, Italy



- founded in 1997
- 3500 students



Background: School Management in the Province of Bolzano

- Province has central database about pupils, teachers, etc.
- Would like to answer statistical queries
- Problem: Data often entered with delays or not at all
- Thus, administrators would like to know whether a query is currently reliable (complete)

Completeness Problem



Relational Database Setting

- Completeness reasoning can be reduced to query containment [VLDB 2011]
- Null values increase complexity due to ambiguity (Unknown values versus non-existing values) [CIKM 2012]

Linked Data/Semantic Web

- Generally, the semantic web is assumed to be open-world
- Completeness statements can allow closed-world reasoning [ISWC 2013]



Geographic Data

- Spatial predicates allow to conclude query completeness in certain areas *"All hotels near a train station and a museum"*
- Completeness statements exist in OpenStreetMap [BNCOD 2013]

🍓 Oxford - OpenStreetMa	p Wiki - Mozilla Firefox		Addres, Second apply	- Marriell PressPect] 🛛 🗙	<u> </u>
<u>Datei B</u> earbeiten <u>A</u> nsich	nt <u>C</u> hronik <u>L</u> esezeichen E <u>x</u> tras <u>H</u> ilfe	an Annual State Annual	Annalise 1978						
Dxford - OpenStreetMa	p Wiki 🛛 🗴 🕅 Map Compare Geofabrik To 🗡	Venue 29th British Nationa	al × 😻 St Ann	e's College, Oxford H	× 🔝 St Anne's College, South Para × 占	Sec.	-		
Image: State					 bncod 	م	Ŧ	俞	
wiki.openstreet Printable version Permanent link Cite this page	Description The University City of Oxford lies to the north west or short distance from the M40. The River Thames runs through the centre of the city OSM Coverage Oxford is beginning to have excellent coverage in OS Kennington are almost certainly complete. If we've m See Template En:Map_status for an explanation of version of version are almost certainly complete. If we've m See Template En:Map_status for an explanation of version of version of version of version of the status for an explanation of version of	242 Cranepont OpenStreetMap images (and underlying CopenStreetMap images (and underlying f London and Reading, and to the so r, the Oxford Canal comes in from th SM. All the major roads and probably issed the road you live on, please v what the symbols and colours mean Status Rer Status Res Status Rer Status Rer Statu	t Line of the minor ror ror ror ror ror ror ror ror ror	Valey 10 / 10 / 10 / 10 / 10 / 10 / 10 /	Groad is the A34, and the northern edge is the A4 ast. Blackbird Leys, Barton, Sandhills, Dean Court, Ci dd a red dot.		ک streets with no-names & ck for OSM & with more maps nks: j list & j list archives & nspector & listeeBrowser & lis		
	26. Headington Quarry, Wood Farm		achadwick d	necked 2009-05-31	unibzit				1

7

Verifying Completeness over Processes

- Data often created following processes
- Many processes are executed only partially formal (pen&paper, email, phone, ...)
- Valid information may be stored in databases with delays

➔ Database content is of questionable completeness

N.B.: Completeness and timeliness are equivalent problems, if database is monotonic

Enrolment Process in a School



Observation



 At some points, new facts in the real world have not yet been stored

→ queries may give wrong answers

- At other points, all facts that hold in the real world have been stored
 - → queries give correct answers

Formalization: Two Databases

Conceptually, there are

- the state of the information system
- the state of the real world

We model

- each state as a database
- the process interacting with both



Two databases: Example



- **Deciding** about enrolments:
 - read from and write into real-world world database
- Recording accepted enrolments into the information system:
 - read from real-world database
 - write into information system database

Completeness Problem



Is the information system up-to-date wrt. the state of the real world?

Research Questions

- How can we express which data a process generates?
- What does completeness mean?
- How can we find out whether a query is complete?

Model: Quality-aware Transition Systems (QATS)

- Goal: Technique applicable to different modeling languages
- Low-level formalism for process instances: Transition systems
 - Petri nets can be encoded using their reachability graphs (possibly exponential encoding due to parallelism)
- Actions in a QATS can be labeled with two kinds of effects:
 - Real-world effects: allow to create new data in the real world
 - Copy effects: store information that holds in the real world into the information system

QATS are data-monotonic

Example Revisited



Real-world effect: $pupil^{rw}(n, s) \nleftrightarrow request^{rw}(n, s)$

Copy effect: $pupil^{rw}(n, s) \rightarrow pupil^{is}(n, s)$

Real-world and Copy Effects

Real-world effect: $pupil^{rw}(n, s) \leftrightarrow request^{rw}(n, s)$ Copy effect: $pupil^{rw}(n, s) \rightarrow pupil^{is}(n, s)$

In general, a real-world effect has the form

 $R^{rw}(X,Y) \nleftrightarrow G^{rw}(X,Z)$

where G is a condition, X are bound variables and Y are unbound variables.

It allows to introduce new facts R^{rw}(X,Y), if G^{rw}(X,Z) holds for some Z

A copy effect has the form

 $R^{rw}(X), G^{rw}(X,Y) \rightarrow R^{is}(X)$

It copies all facts in R^{rw} that satisfy G^{rw} into R^{is}

Real-world effects are nondeterministic, copy effects are deterministic

Transition Systems for Process Instances



Transition Systems for Process Instances

Two concurrent process instances:

- Middle School A
- High School B



Completeness Verification

Given

- Process description
- State S
- Query Q



Question

Is it safe to pose the query Q in state S against the information system database?

Completeness Verification (2)



Compliance

When does a development ((D^{rw}₀, D^{is}₀), .., (D^{rw}_n, D^{is}_n)) comply to a sequence of real-world and copy effects?

Compliance to Real-world Effects

Real-world database

request(John,HS) request(Mary,HS)

Real-world effect: pupil^{rw}(n, HS) <--- request^{rw}(n, HS)

Because "↔⊷" is nondeterministic

request(John,HS) request(Mary,HS) pupil(John,HS) Pupil(Mary,HS)

Possible successive real-world databases:

request(John,HS) request(Mary,HS) Pupil(Mary,HS)

request(John,HS) request(Mary,HS) pupil(John,HS)

request(John,HS) request(Mary,HS)

Compliance to Copy Effects

Real-world database

pupil(John,HS) Pupil(Mary,HS)

Copy effect: $pupil^{rw}(n, HS) \rightarrow pupil^{is}(n, HS)$

Resulting information system database



Results – Completeness over Paths

• A real-world effect is risky wrt. a query, if it has the potential to change the query result

Adding pupils in class 1A is risky wrt. a query for all pupils, but not wrt. a query for all pupils in level 2

- Copy effects can repair a risky effect, if they copy all data that has the potential to change the query result Copying all pupils in level 1 into the information system repairs the risky effect.
- Result: A query is complete over all developments of a path, if all risky effects in the path are repaired

Theorem: Repair checking can be reduced to query containment

 Query containment for conjunctive queries (SELECT ... FROM ... WHERE ...) has been well studied in database research

Results – Completeness in States

- Completeness holds in a state, if it holds for all paths that lead to that state
- A priori, infinitely many paths (due to cycles)

Theorem: Repeated actions can be ignored

- Thus, only finitely many paths to consider
- Still, number of paths can be exponential wrt. the QATS

Example Revisited



Complexity

Query and effect language	Complexity of completeness checking for a path	Complexity of completeness checking for a state
Arbitrary conjunctive queries (CQ)	Π ^P ₂ -complete	Π ^P ₂ -complete
CQs without <, \leq	NP-complete	In Π ^P ₂
CQs without selfjoins	coNP-complete	coNP-complete
CQs without selfjoins and without <, ≤	ΡΤΙΜΕ	in coNP

Applications

- Annotation of statistics and KPI with completeness information (see next slide)
- Process mining (trace analysis) to validate whether queries over traces return the real state of the process
- Auditing to verify whether the information about the real-world is properly stored

Possible Use: Statistical Reports



Conclusion

- Introduced the problem of query completeness due to delays between real-world events and their recording in a database
- Modelling of the problem using quality-aware transition systems that interact both with the real world and with an information system
- Showed how to verify query completeness over such models
- Future work: Demo for a high-level process language (BPMN or YAWL)

Thank you!



Acknowledgment

This research has been supported by the project "MAGIC", funded by the Province of Bozen-Bolzano