

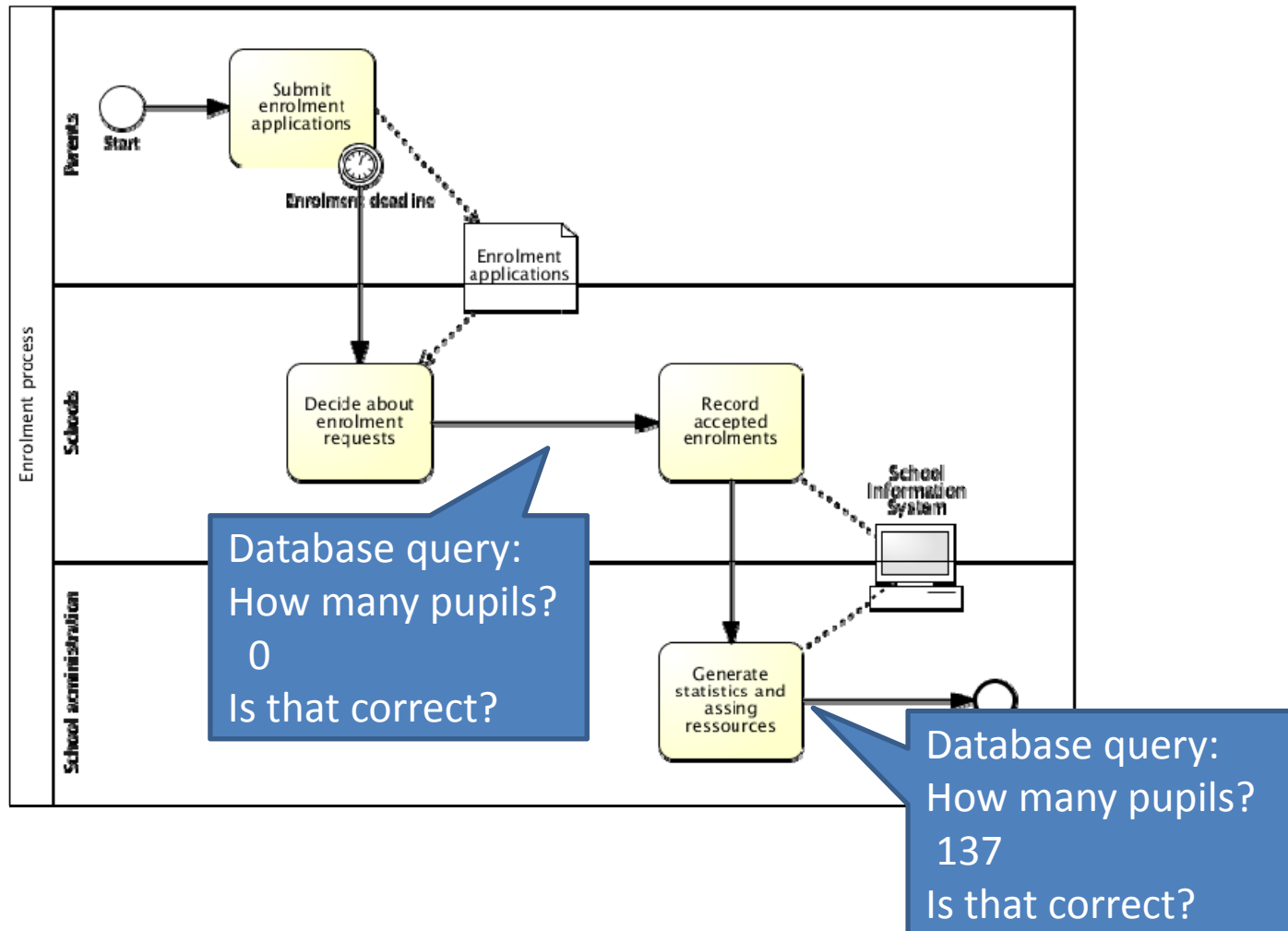
Observation

- Often, process execution is only **partially formal** (pen&paper, email, phone, ...)
- ➔ Valid information may be **stored** in databases **with delays**
- ➔ Database content is of **questionable completeness**

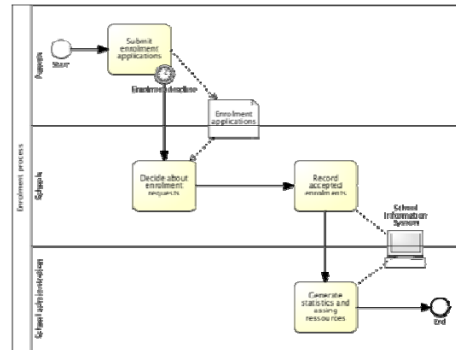
Background: School Management in the Province of Bolzano

- Province has **central database** about pupils, teachers, etc.
- Would like to answer **statistical queries**
- Problem: Data often entered with delays
- Thus, administrators would like to know
whether a query is currently reliable (complete)

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

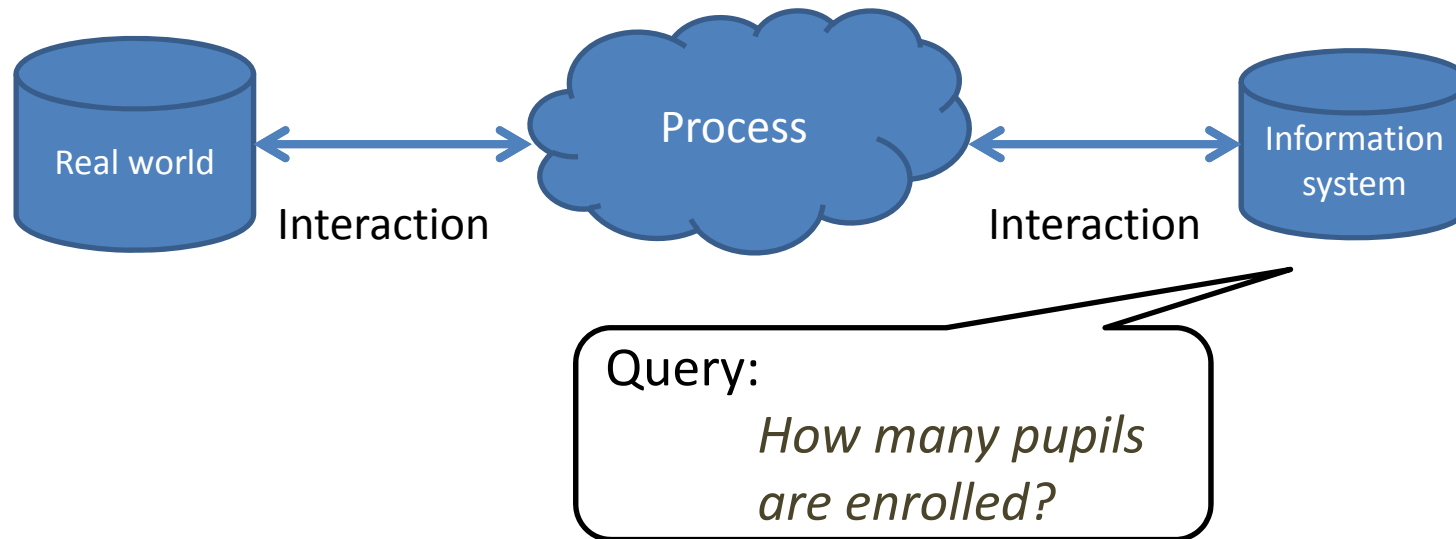


Two databases: Example



- **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



Does the information system contain **all the enrolments**?

Research Questions

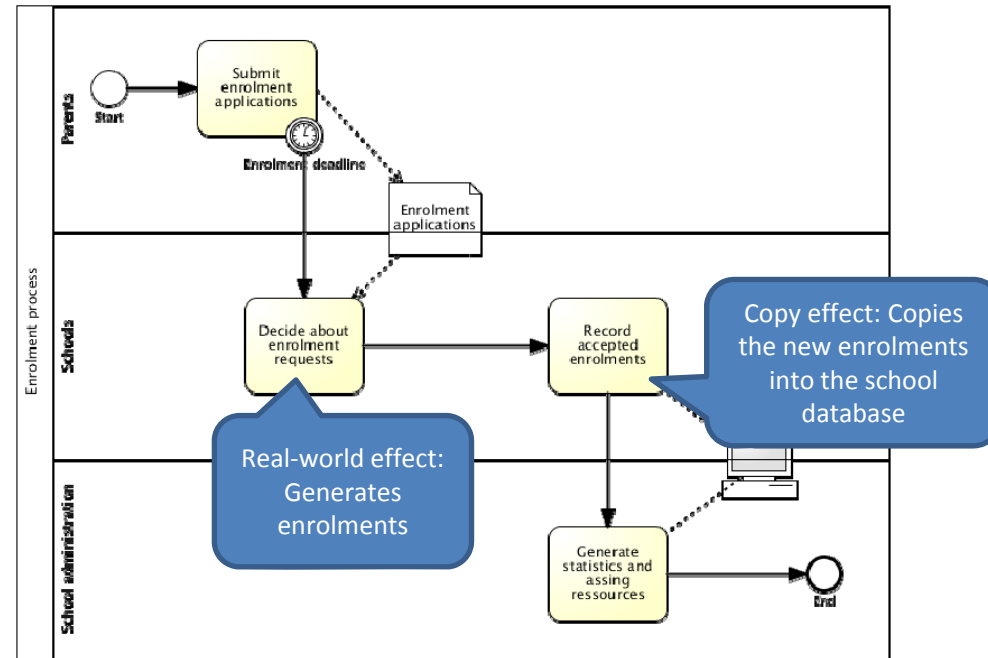
- How can we **express which data a process generates?**
- How can we **express which data are recorded** in the information system?
- How are reads and writes of data related?
- What does **completeness** mean?
- How can we **find out** whether a **query is complete?**

Model: Quality-aware Transition Systems

- Goal: **General technique** applicable to different modeling languages
- Therefore, we use **transition systems** as mathematical formalism
 - *Petri nets can be encoded using their reachability graph (possibly exponential encoding due to parallelism)*
- Actions in our transition systems 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

Thus, our models are data-monotonic

Example Revisited



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

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

Real-world and Copy Effects

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

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

In general, a **real-world effect** has the form

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


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

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

A **copy effect** has the form

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

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



Real-world effects are
nondeterministic,
copy effects are
deterministic

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)

A state S of a QATS
satisfies **completeness**
for a query Q ,

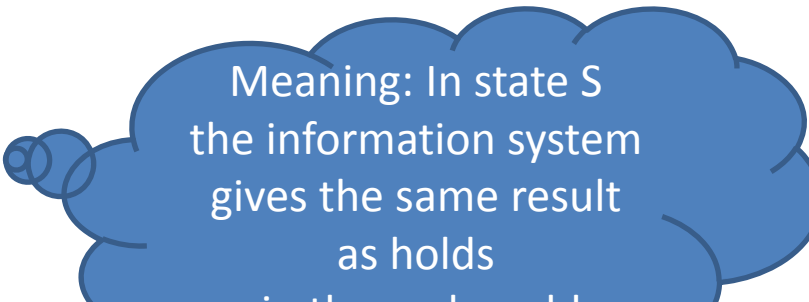
if

for all paths leading to S ,

for all process-compliant database developments

$((D_0^{rw}, D_0^{is}), \dots, (D_n^{rw}, D_n^{is}))$,

$$Q(D_n^{is}) = Q(D_n^{rw})$$



Meaning: In state S
the information system
gives the same result
as holds
in the real world

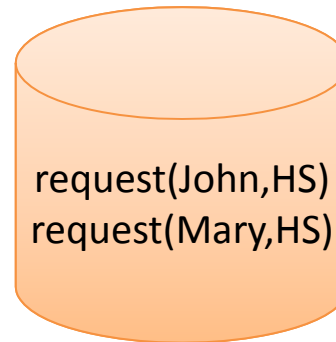
This is what we want to decide!

Compliance

When does a development $((D^{rw}_0, D^{is}_0), \dots, (D^{rw}_n, D^{is}_n))$ comply to a sequence of real-world and copy effects?

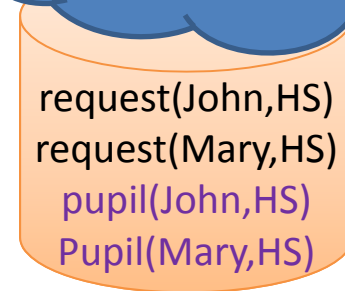
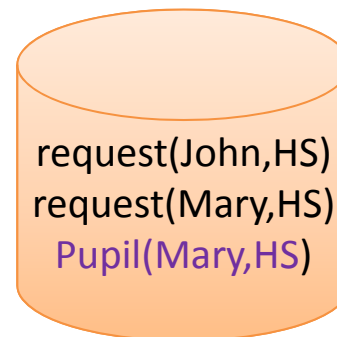
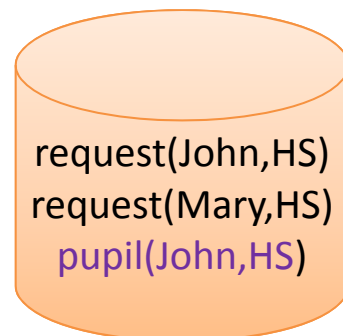
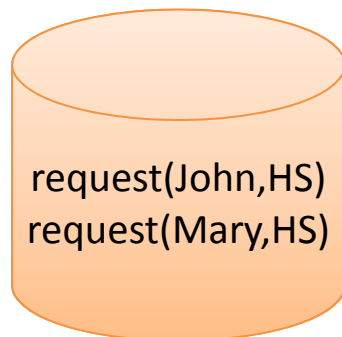
Compliance to Real-world Effects

Real-world database



Real-world effect: $\text{pupil}^{\text{rw}}(n, \text{HS}) \leftarrow^{\text{rw}} \text{request}^{\text{rw}}(n, \text{HS})$

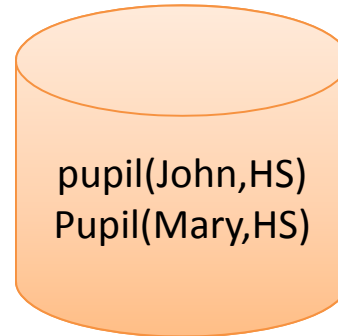
Possible successive real-world databases:



Because " \leftarrow^{rw} " is nondeterministic

Compliance to Copy Effects

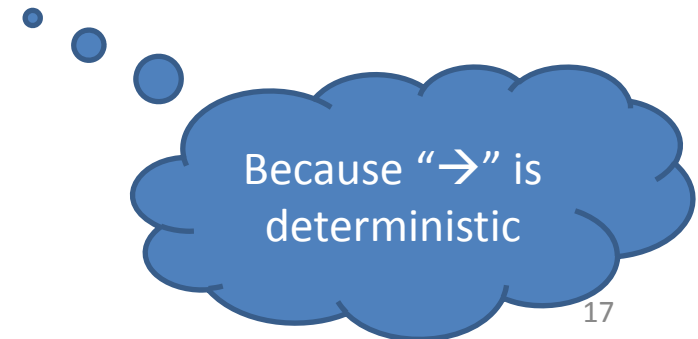
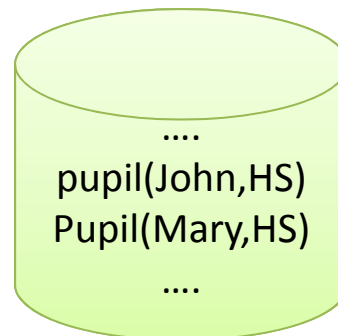
Real-world database



Copy effect:

$$\text{pupil}^{\text{rw}}(n, \text{HS}) \rightarrow \text{pupil}^{\text{is}}(n, \text{HS})$$

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 actions 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

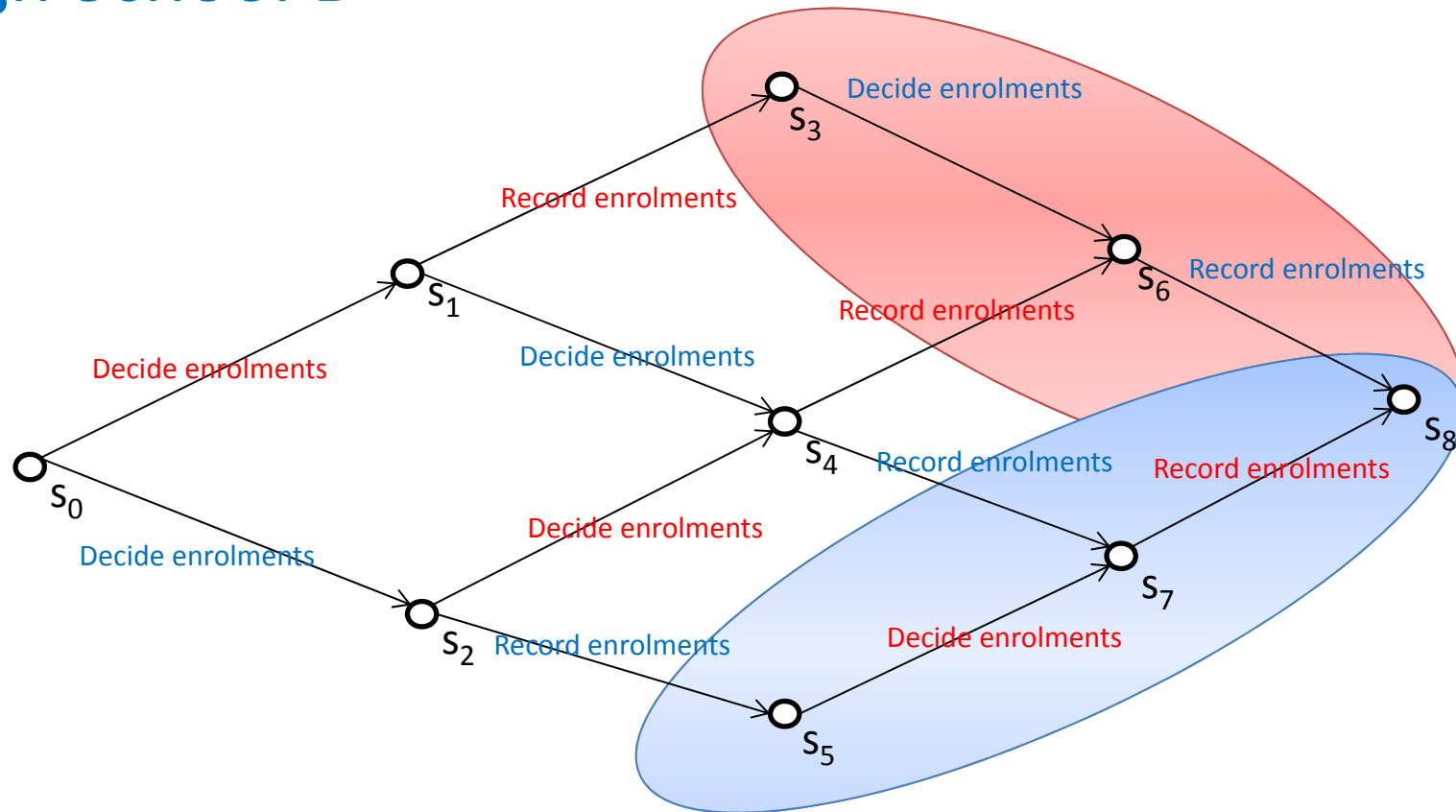
Completeness Checking - Intuition

Middle School A

High School B

How many middle school pupils?

How many high school pupils?



Complexity

Query and effect language	Complexity of completeness checking for a path	Complexity of completeness checking for a state
Arbitrary conjunctive queries (CQ)	Π_2^P -complete	Π_2^P -complete
CQs without $<$, \leq	NP-complete	In Π_2^P
CQs without selfjoins	coNP-complete	coNP-complete
CQs without selfjoins and without $<$, \leq	P TIME	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

School Report 2013

Pupils in primary schools: 548
Pupils in middle schools: 390
Pupils in high schools: 242



Data from the Da Vinci School and the Gherdena School is missing

Pupils taking English: 1157
Pupils taking French: 685
Pupils taking Chinese: 52



The Hofer School did not enter its language course attendance yet

.....

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!

Questions?

