



MAGIK: Managing Completeness of Data

Try out the demo at:

<http://magik-demo.inf.unibz.it>

Ognjen Savkovic

Free University of Bozen-Bolzano, Italy

savkovic@inf.unibz.it

joint work with Sergey Paramonov, Mirza Paramita, and Werner Nutt



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

Data Quality and Data Completeness

What is Data Quality?

- Data is of a high quality if it is fit for intended uses
- Data quality (DQ) has different aspects: **completeness**, correctness, accuracy, etc.
- Little work has been done on **data completeness**

What is Data Completeness?

- A database is complete for a domain if it contains all facts that are true in the domain
- In practice **no DB is complete**, but a database can be sufficiently complete for a given query
e.g., "IMDb does not contain all movies but it contains all movies by Charlie Chaplin"

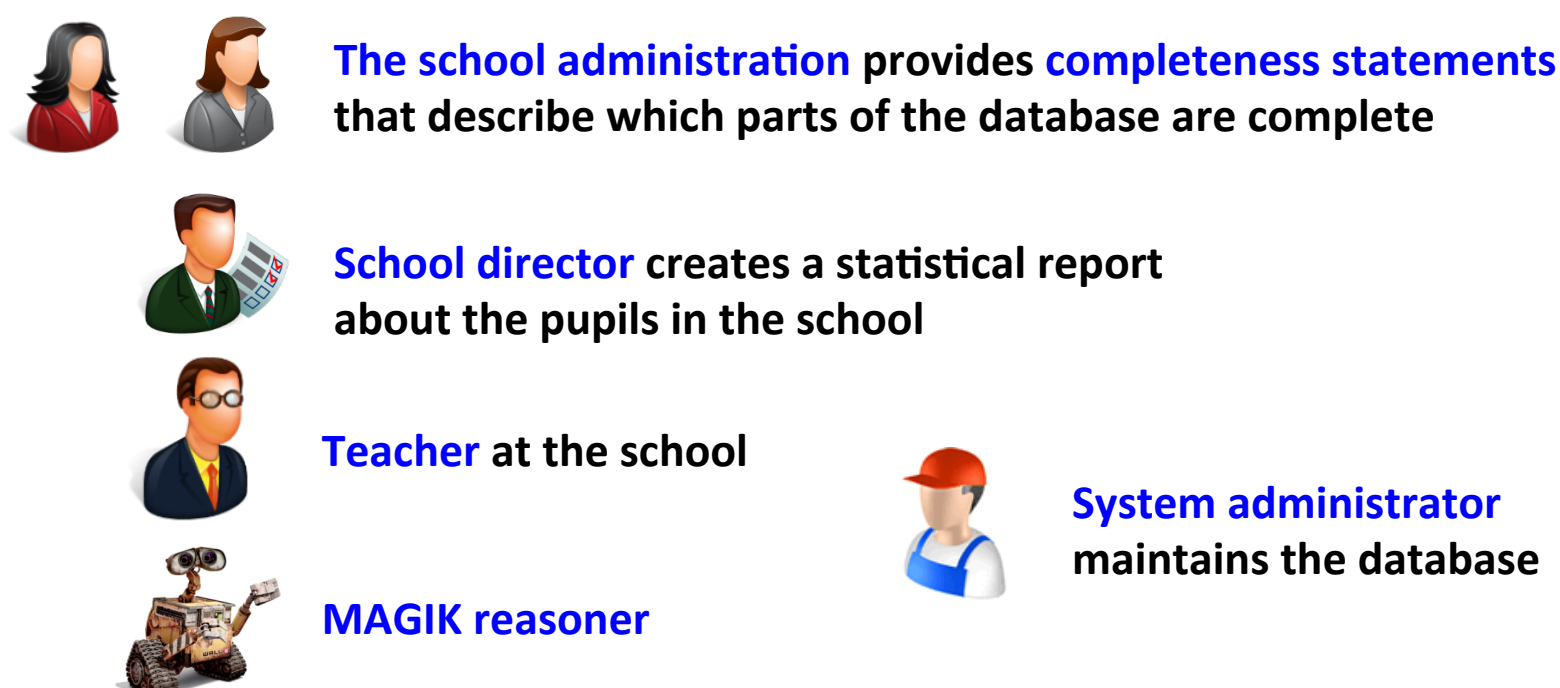
Meta-information about completeness

- Completeness cannot be checked by inspecting the database
 - One cannot see what is missing
- We need information about the database completeness state – **meta-information**
- Often **meta-information** about the data completeness is available
- Information about partial completeness can come from:
 - Business Processes that manipulate the data
 - Humans assertions (e.g., school administration)
 - Origin of the data (**data provenance**)
 - ETL processes that integrate the data, etc.



MAGIK at Work: School Database

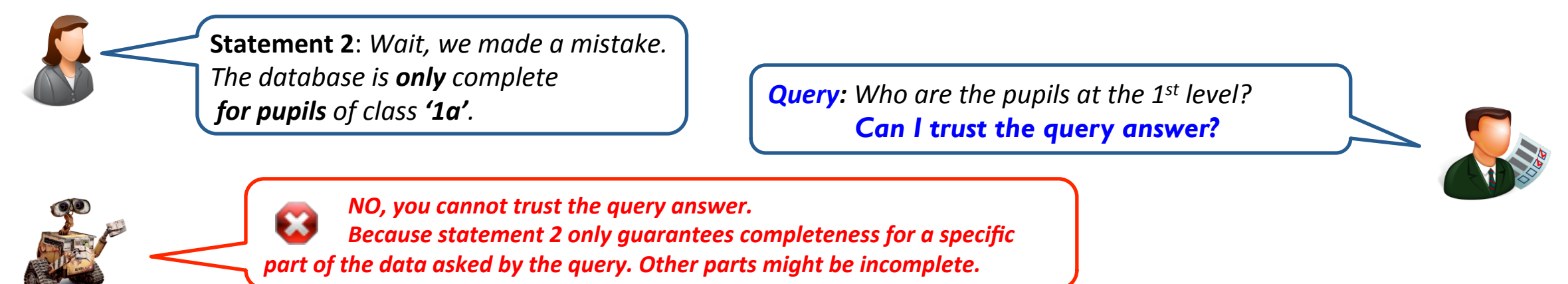
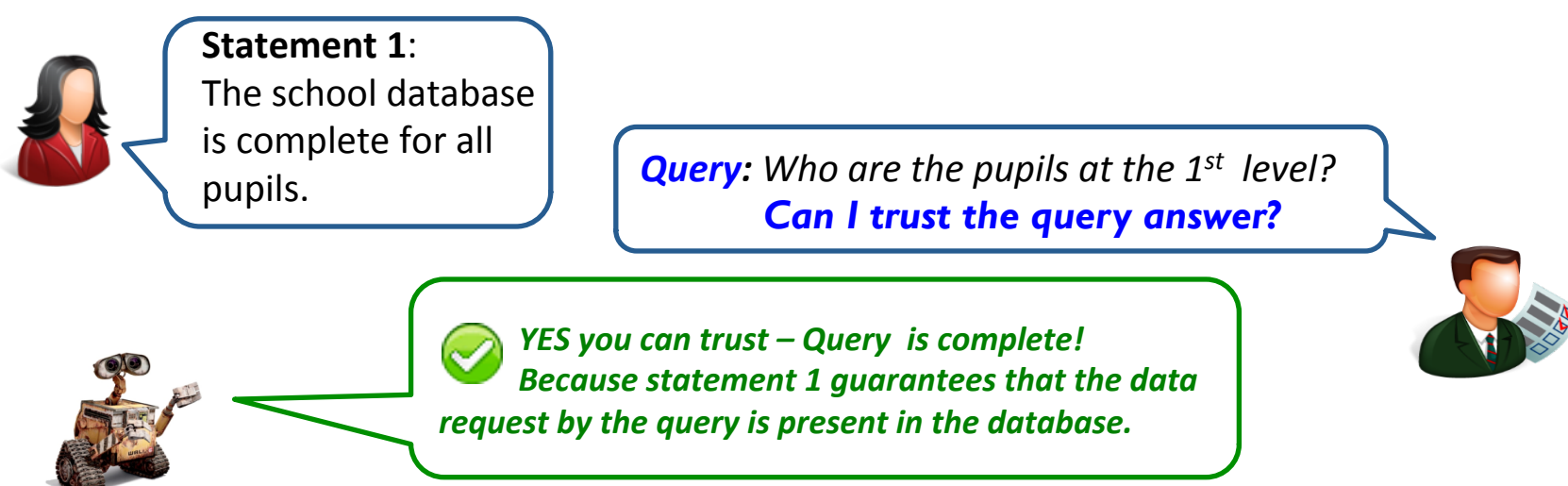
Scenario



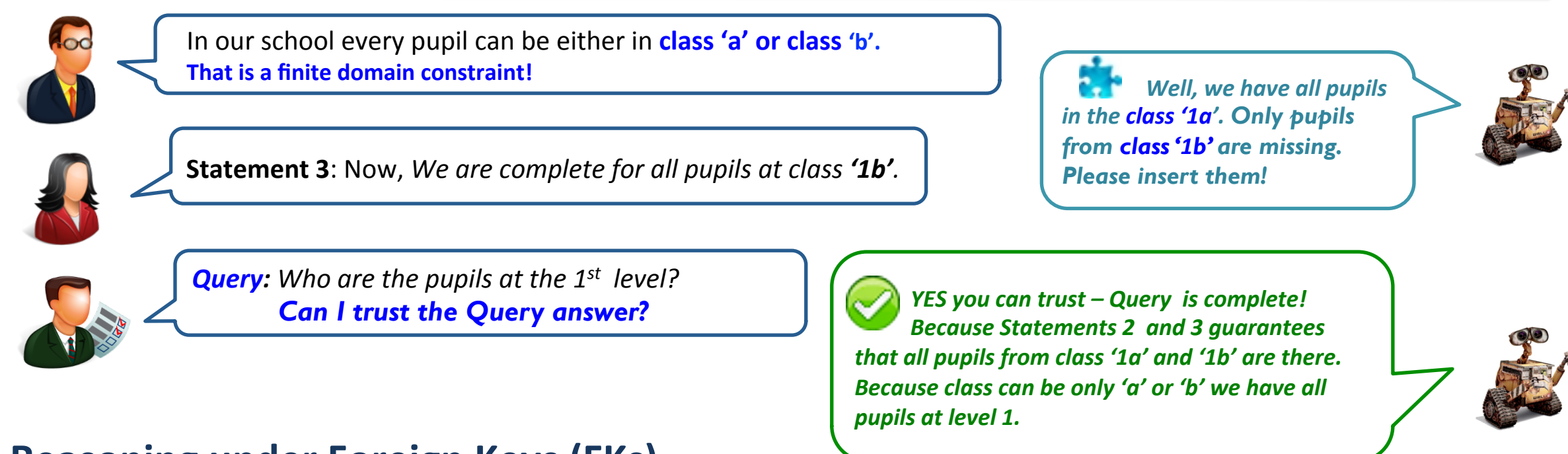
School Database

pupil(name, level, code) ... a pupil belongs to a class of certain level and code
class(level, code, branch) ... every class belongs to some branch
learns(name, language) ... a pupil learns a language

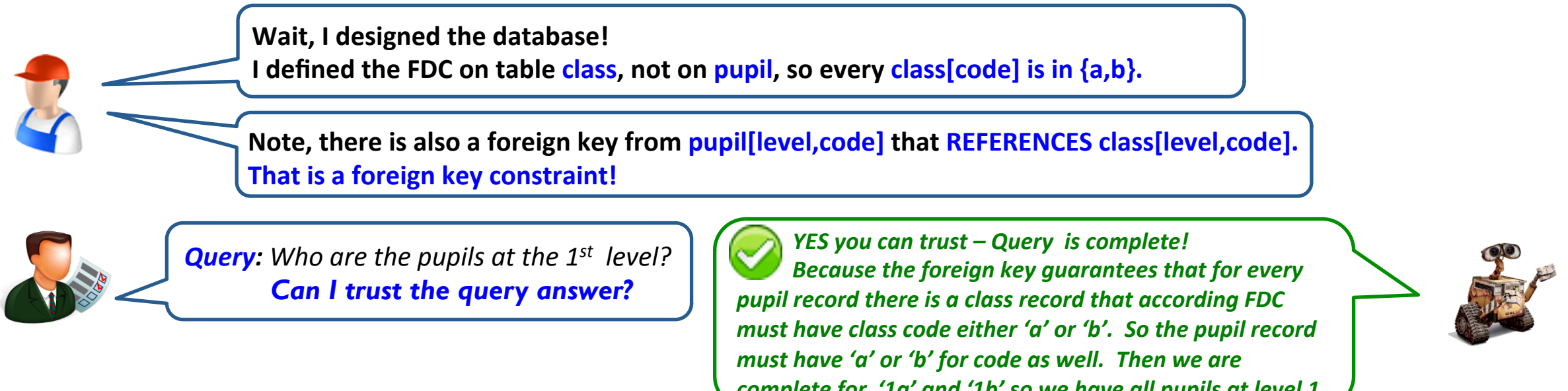
Plain Reasoning



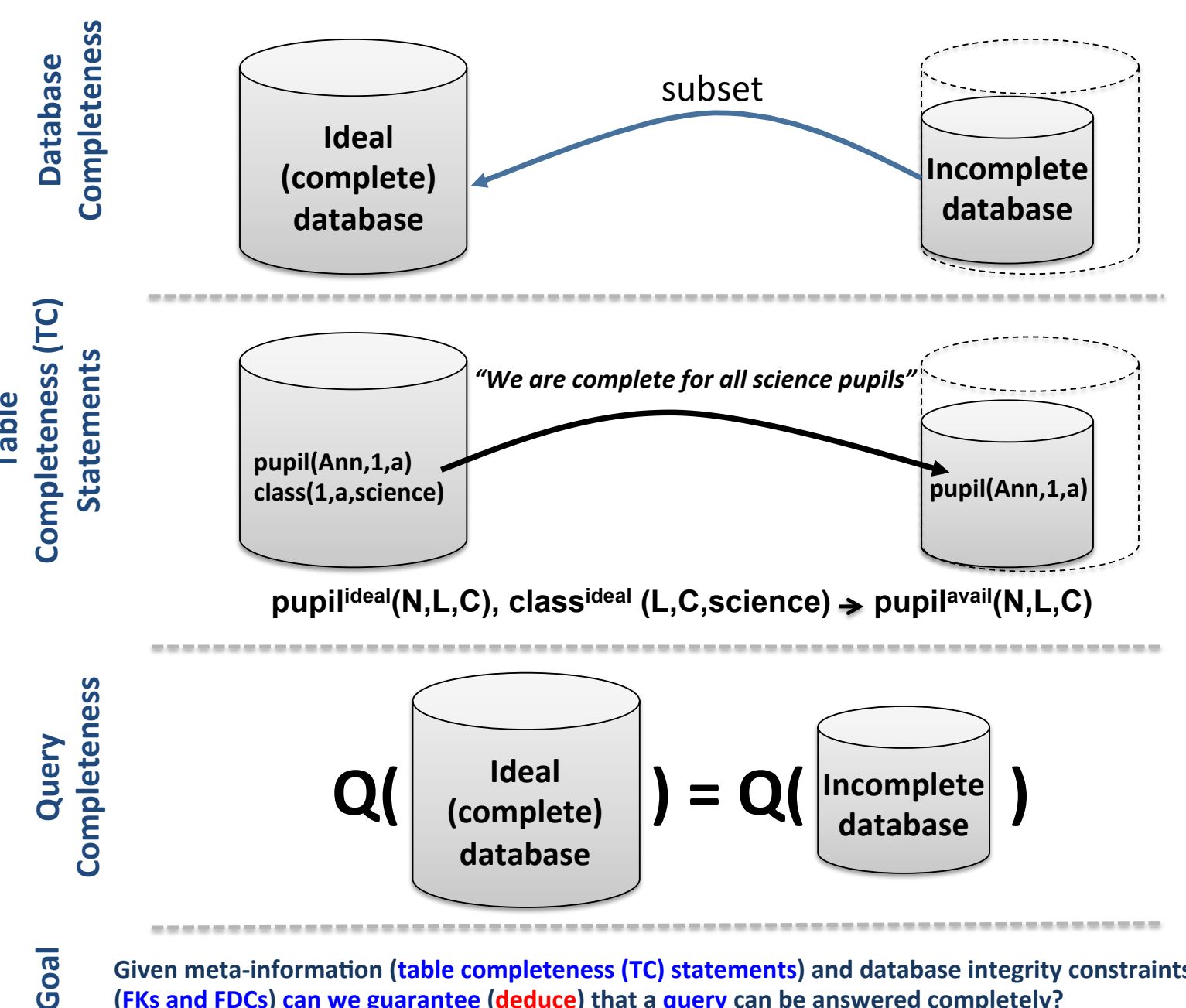
Reasoning under Finite Domain Constraints (FDCs)



Reasoning under Foreign Keys (FKs)

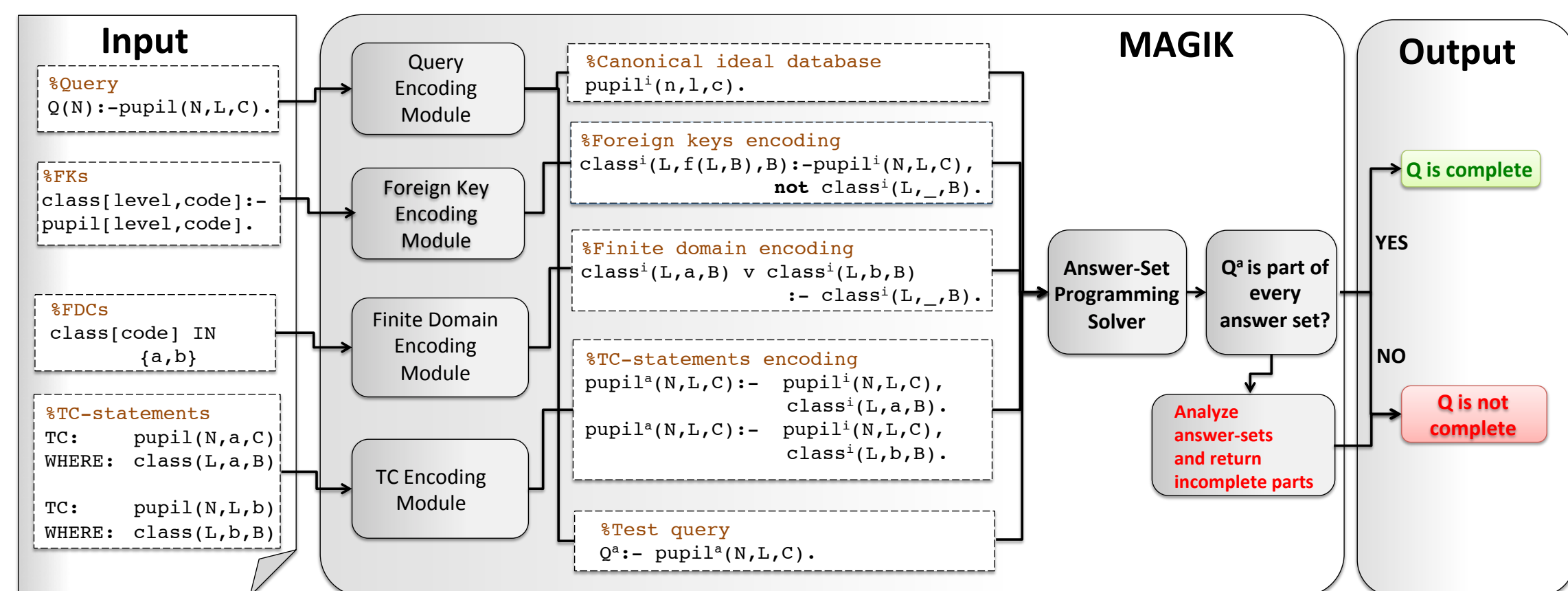


Formalization of the Problem



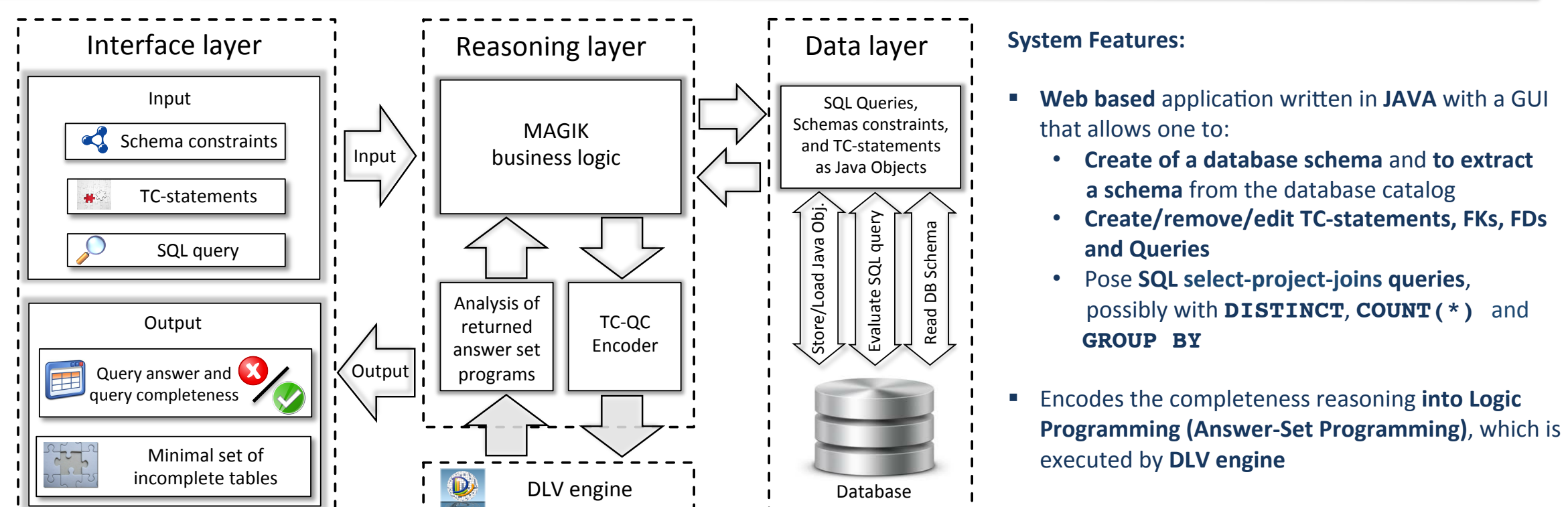
Implementation

Query is Complete iff Query is Complete wrt Canonical Ideal Database



Completeness of Q follows from a set of TCs, FKs and FDCs iff the fact Q^a is in every answer-set of the encoding answer-set program

System Architecture



Summary

- MAGIK checks completeness of queries over incomplete databases given information about partially complete tables
- MAGIK reasons taking into account schema constraints: foreign keys and finite domains constraints
- MAGIK explains its answer and suggests which data to add to make a database sufficiently complete for a query