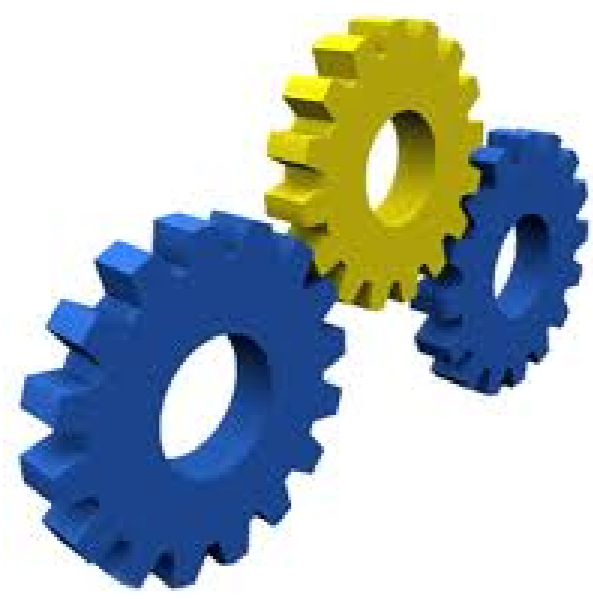


Towards a Standard Development of ETL Processes



Zineb El Akkaoui (promotor Esteban Zimanyi)
 Department of Computer and Decision Engineering (CoDE)
 Universite Libre de Bruxelles
 zelakkao@ulb.ac.be

ETL Process

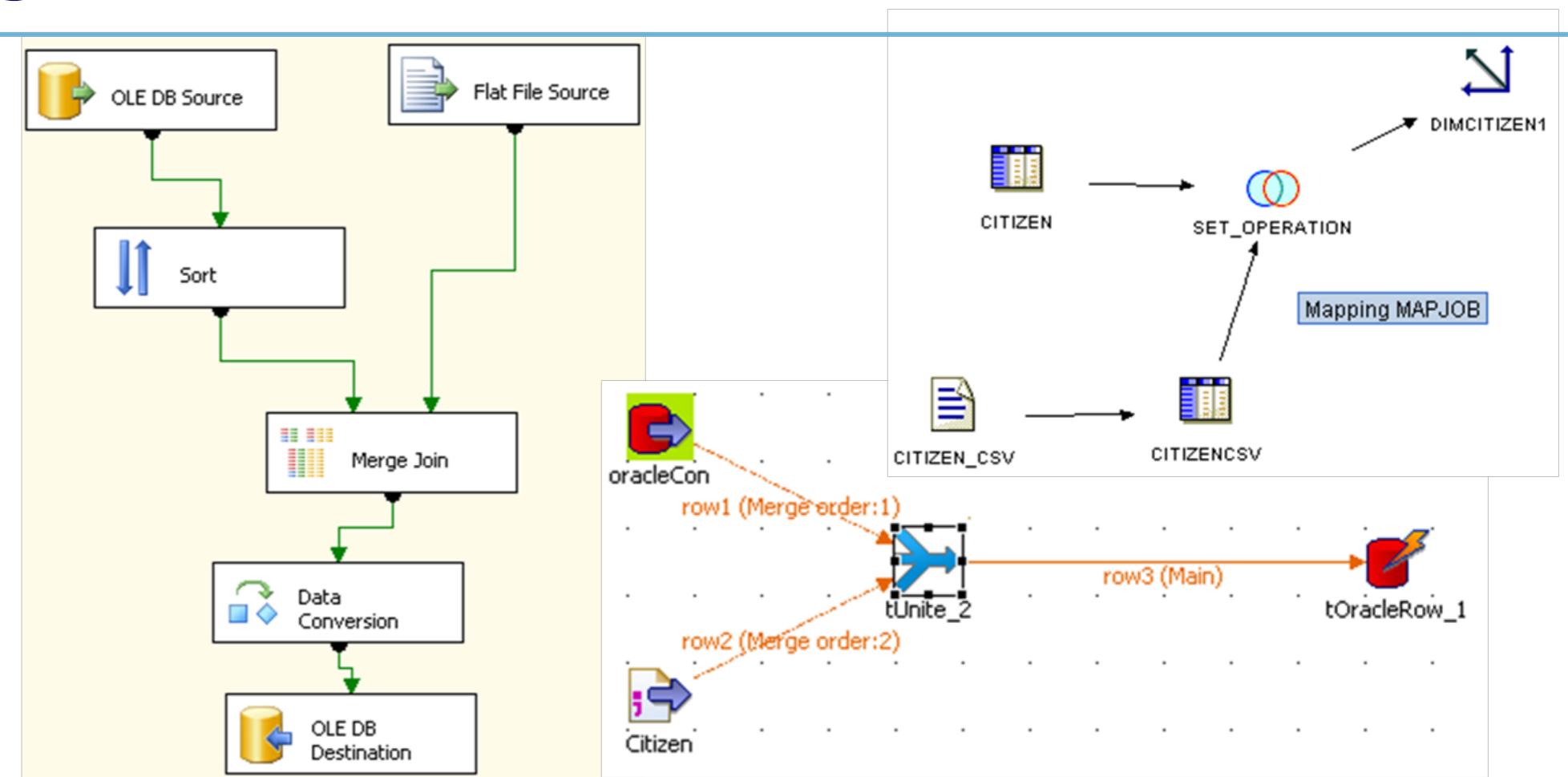
Definition: the backbone component supplying a data warehouse with Integrated and reconciled data, from heterogeneous and distributed data sources

Problem: designed by considering a specific-technology from the beginning of the development process

- ▶ High cost, in both time and resources, part of a data warehouse project
- ▶ No share and reuse methodologies and best practices among projects

Solution: a MDD¹-based framework for ETL development

⁽¹⁾ MDD: Model-Driven Development, a multi-layered software development

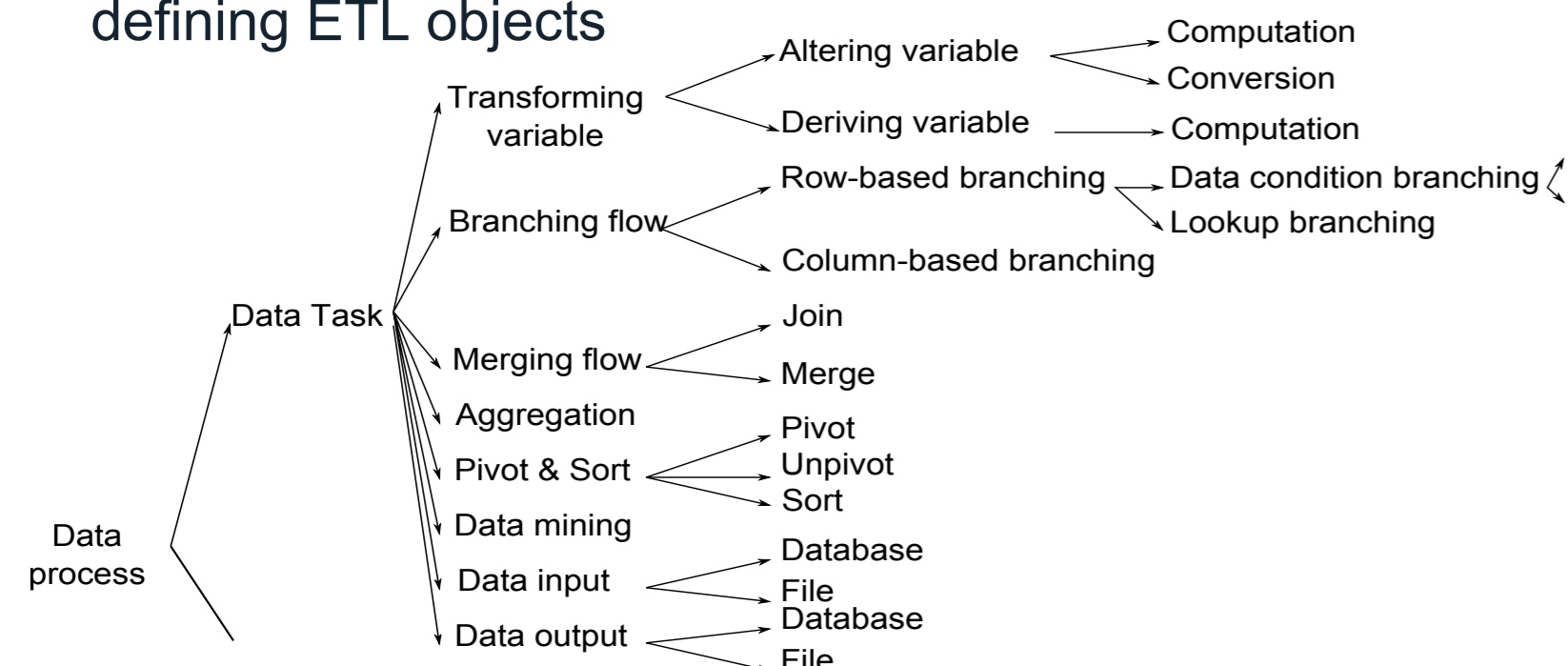


Same ETL process implemented in different ETL tools

Design 1

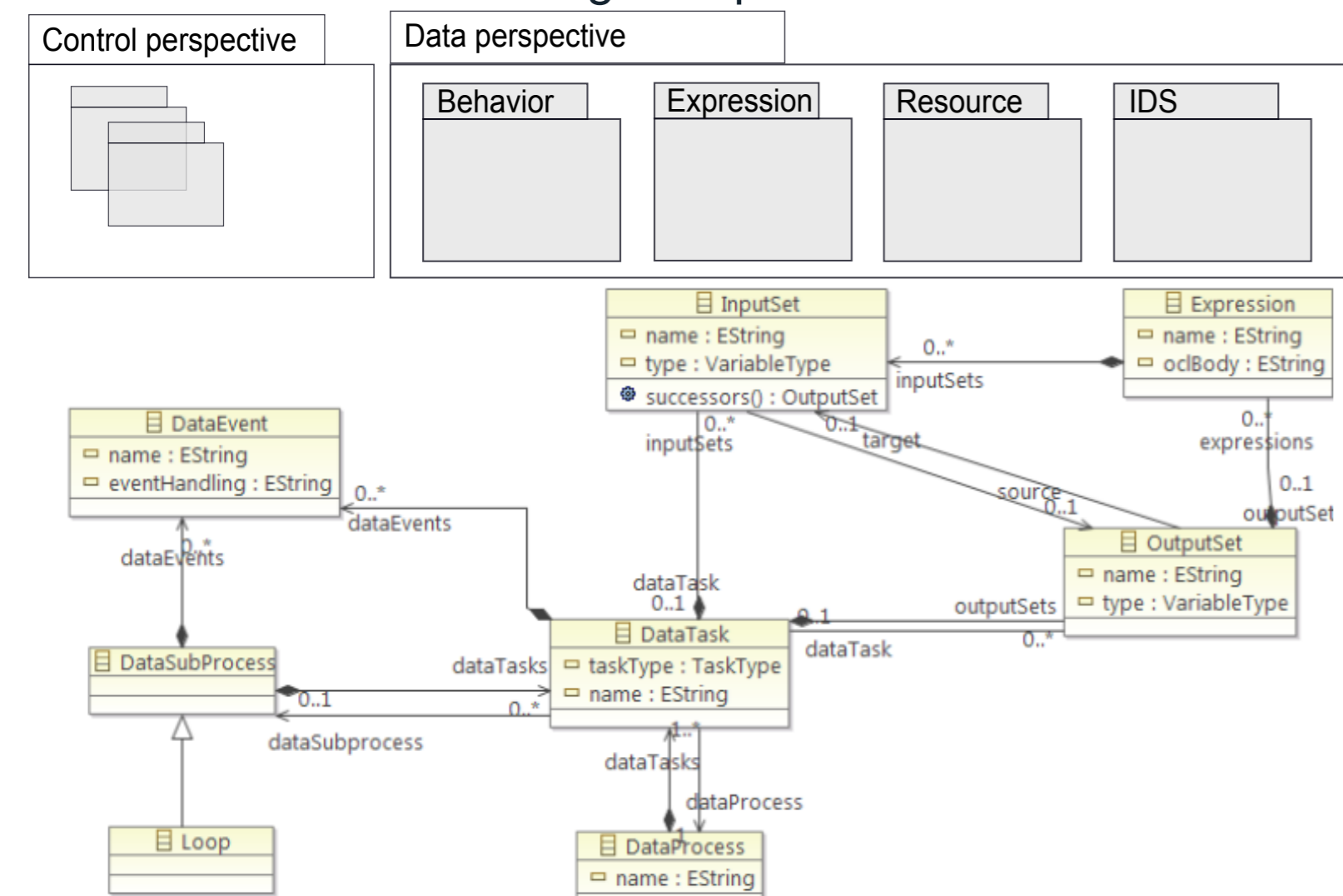
- An ETL process a composition of data & control processes
 - ▶ **A control process:** a sequence of executable programs aiming at loading the whole data warehouse;
 - ▶ **A data process:** an executable program spanning data sources to a part of the data warehouse.

- Towards a full specification of ETL processes:
 - ▶ **Taxonomy** based on a benchmark of ETL tools for defining ETL objects



BPMN-based taxonomy

- Metamodel for defining ETL patterns



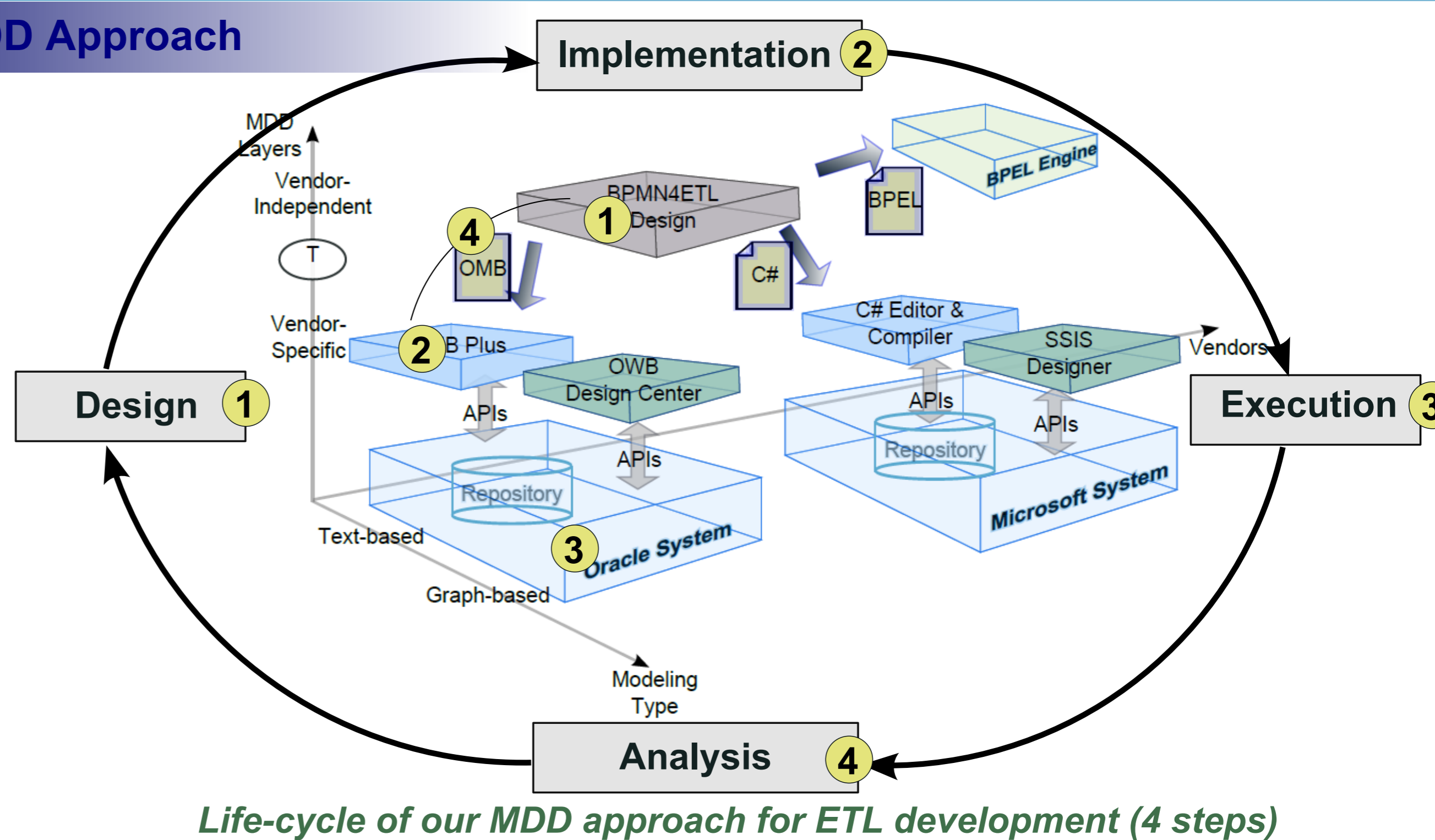
BPMN-based metamodel - excerpt of behavior package

Behavior: data pipeline among ETL tasks
Expression: conditions, computations and queries in tasks
Resource: data storage characteristics
IDS: intermediate Data Storage for non-persistent data

Validation (on progress) 4

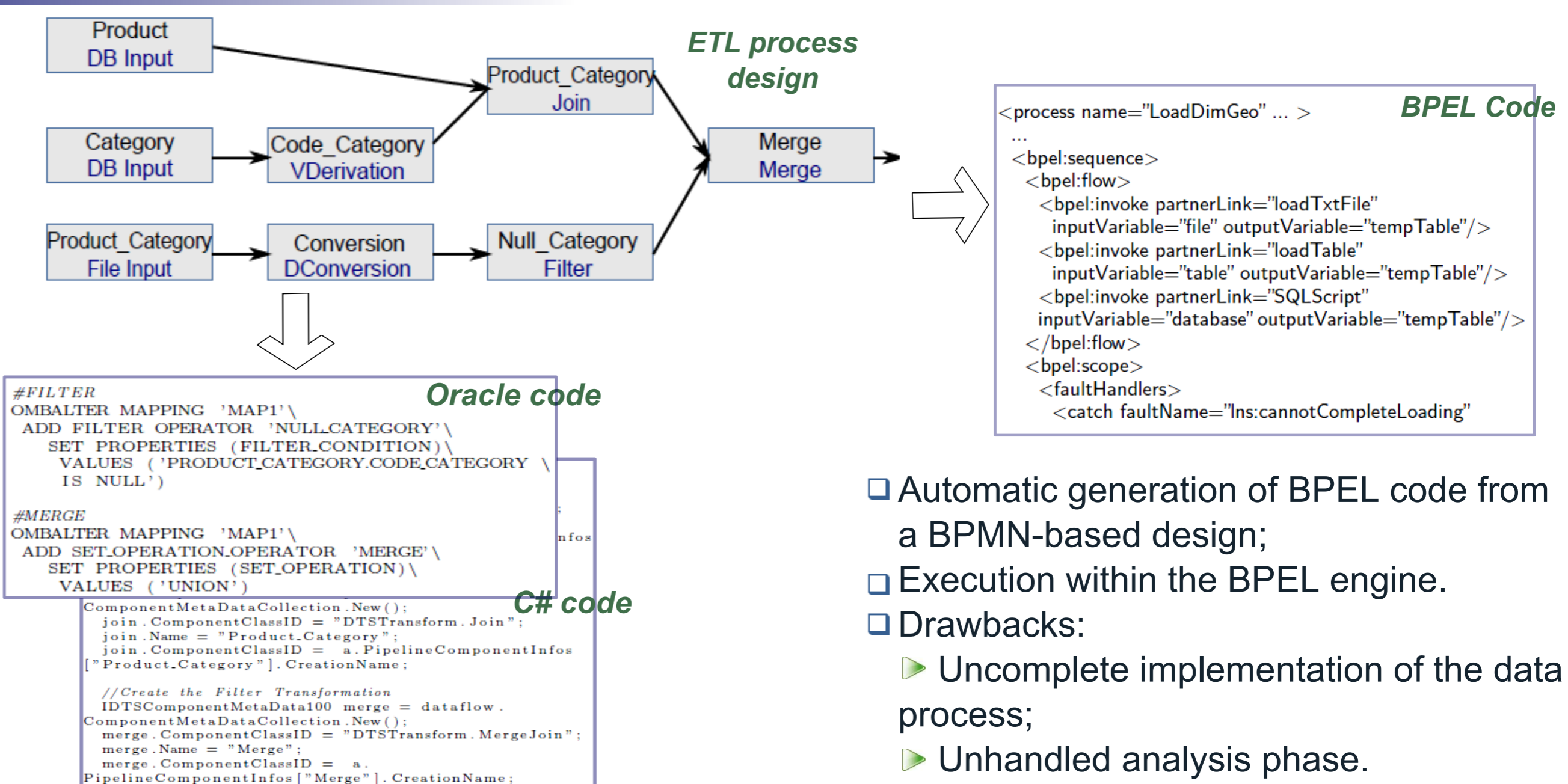
- **Design validation** through:
 - ▶ Design object profiling
 - ▶ Execution metrics
- **Implementation validation** through code metrics
- **Implementation enhancement** through knowledge processing
 - ▶ Code-quality evolution
 - ▶ Detection of most appropriate implementation

MDD Approach



Life-cycle of our MDD approach for ETL development (4 steps)

Implementation 2 3



- Automatic generation of BPEL code from a BPMN-based design;
- Execution within the BPEL engine.
- Drawbacks:
 - ▶ Uncomplete implementation of the data process;
 - ▶ Unhandled analysis phase.
- Useful for distributed data warehouse

- Automatic generation of code for any ETL tool, from the BPMN-based design
- Execution within the ETL engine
- Useful for distributed and non-distributed data warehouses

Conclusion & Future Works

- Our approach comprises:
 - ▶ A **high expressive design** of ETL processes thanks to the use of BPMN
 - ▶ An automatic implementation into code for various environment, which is enhanced by the MDD technology and its corresponding tools (e.g. Acceleo transformations, OCLinEcore...)
- Future works:
 - ▶ A GMF implementation for the BPMN-based design
 - ▶ Continue the validation work.