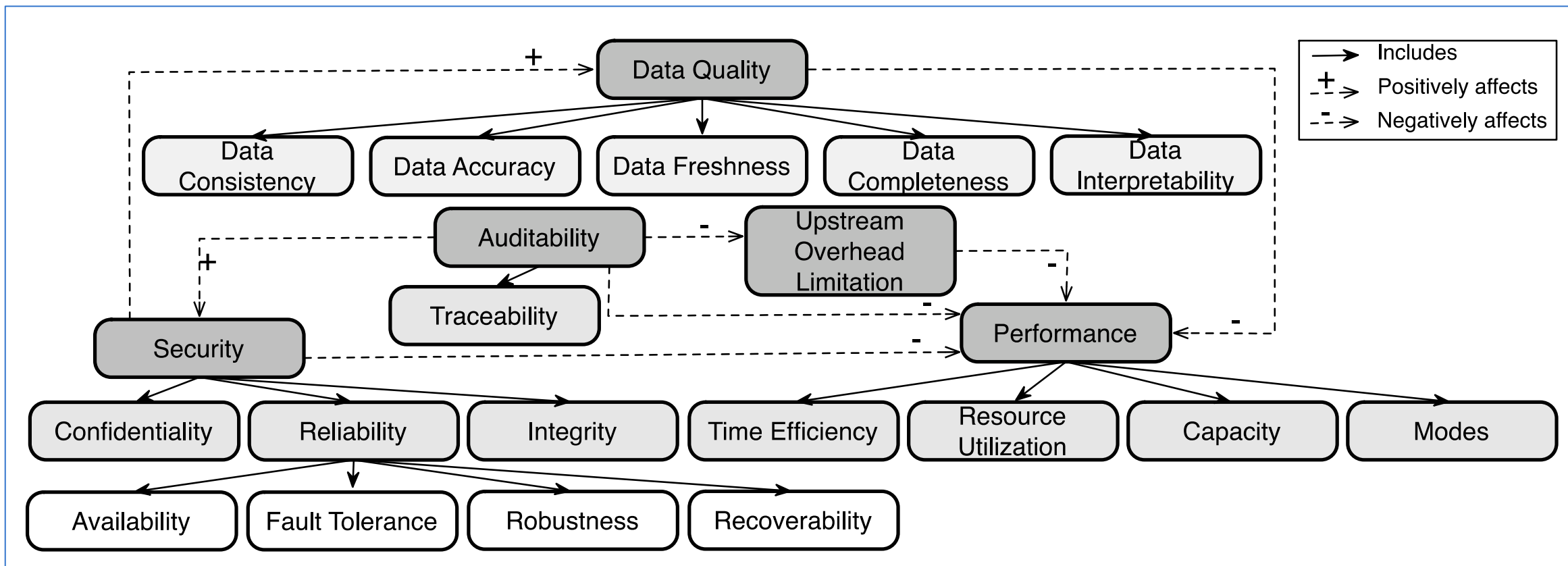# Automating User-Centered Design of Data-IntensiveProcesses
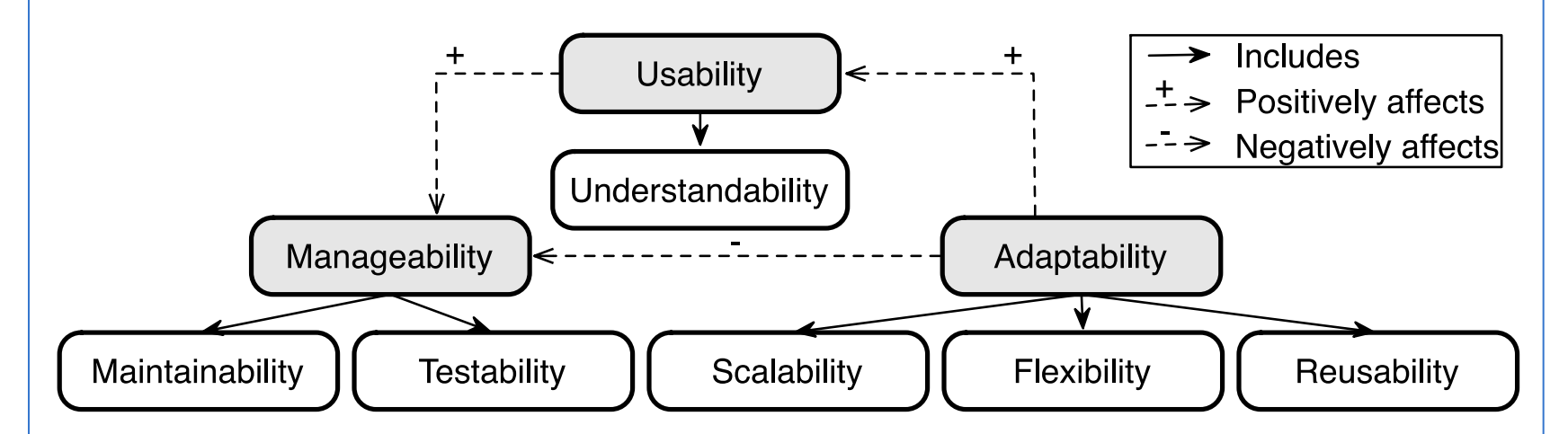
**Vasileios Theodorou, Alberto Abelló**
**Universitat Politècnica de Catalunya, BarcelonaTech**
**[vasileios | aabello]@essi.upc.edu**

**Wolfgang Lehner, Maik Thiele**
**Technische Universität Dresden**
**[wolfgang.lehner | maik.thiele]@tu-dresden.de**

## Dependencies among process characteristics with construct implications



## Dependencies among characteristics for design evaluation



**POIESIS:** A tool to automatically generate quality patterns over existing ETL processes in an iterative, dynamic fashion, with high-level user interaction and based on pursued goals.

## ETL Generation and Improvement

### • Functionality-Based Design

ETL Process Designer: Semi-automatically designs ETL process model that implements basic ETL functionality. Input: domain metadata & business requirements.

### • Quality Enhancement

- iterative, incremental, user-centered
- ETL flow is represented as logical model that can be visualized for user as a BPMN process
- iterations are terminated at any point once user approves the model as adequate

Process Simulator : Simulates ETL processes and produces meaningful simple and aggregate analytics according to user's interest.
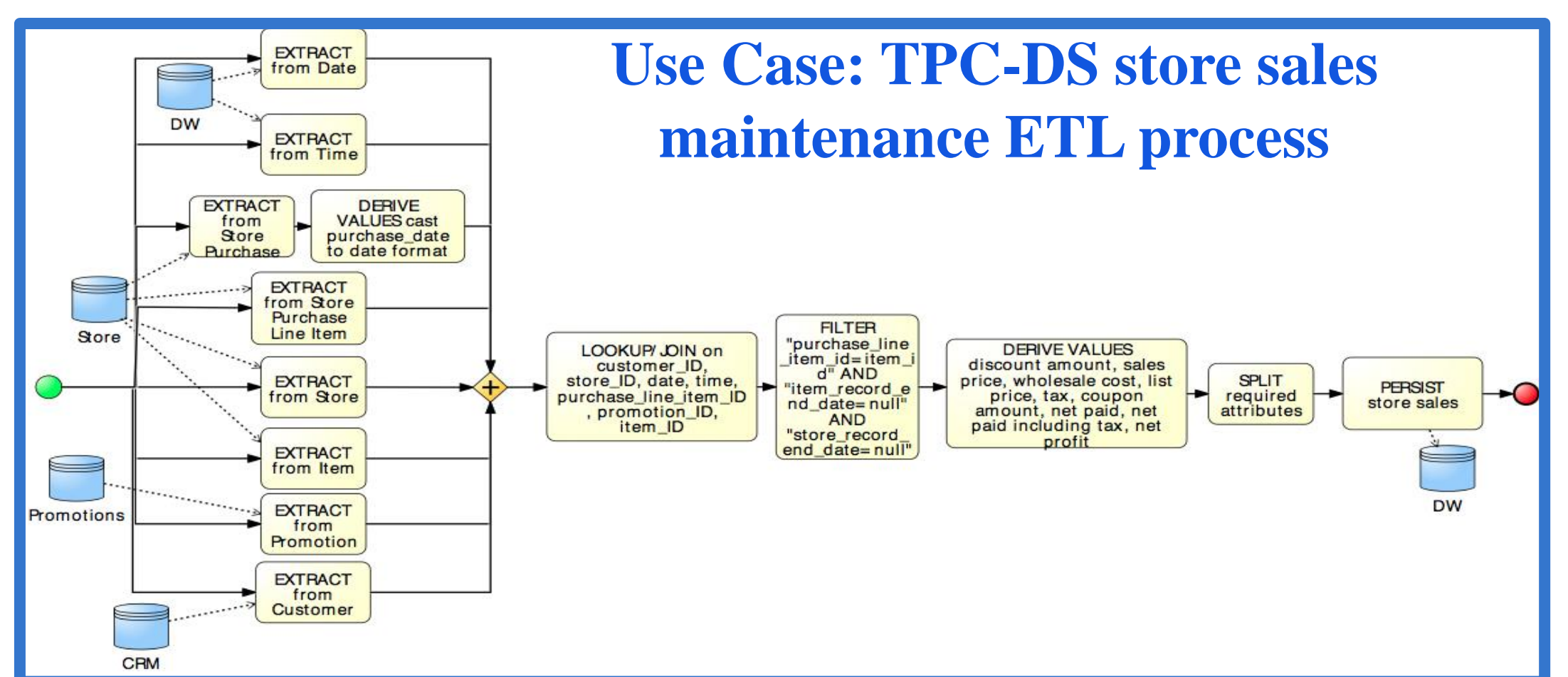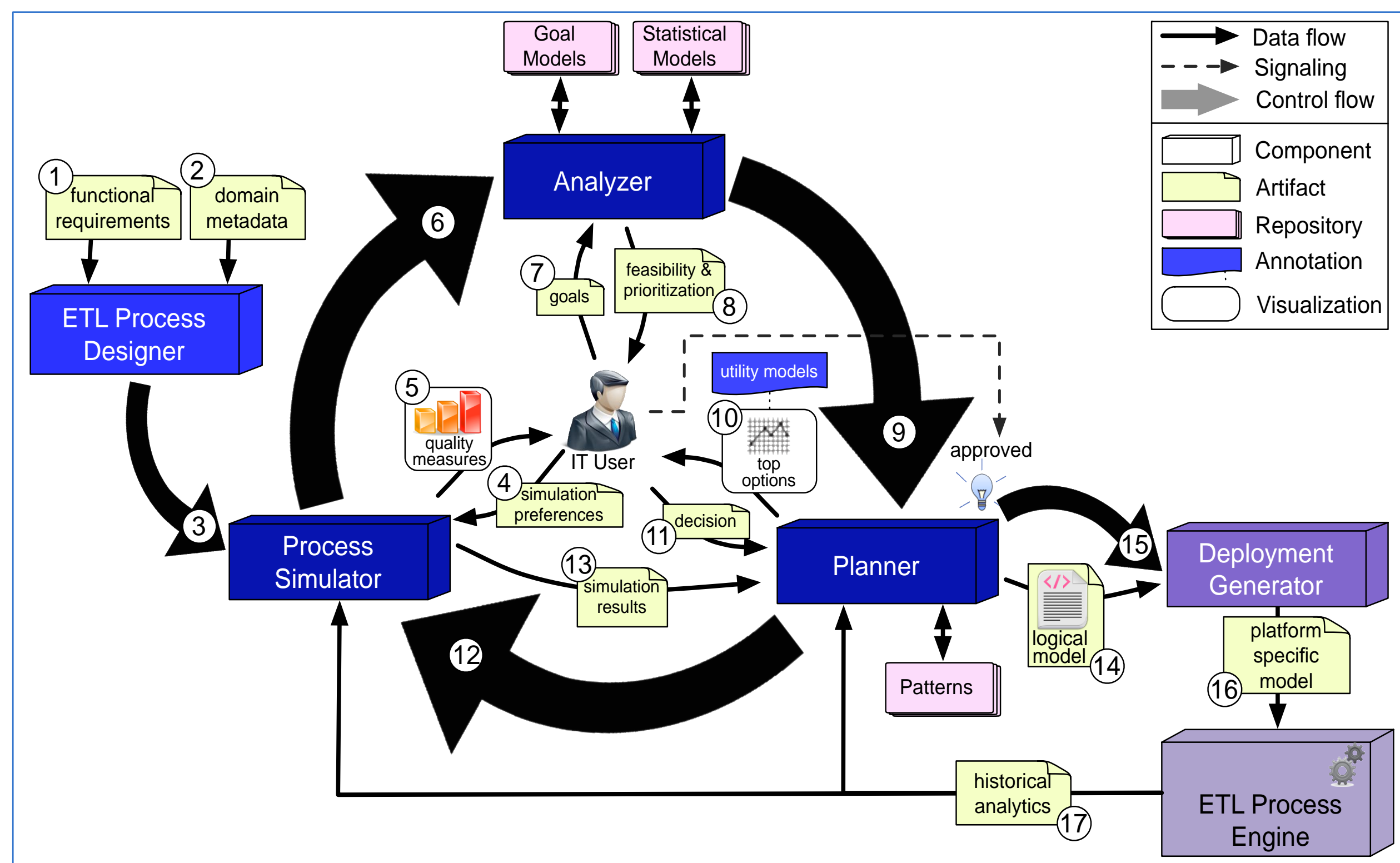
Analyzer: Performs feasibility analysis and prioritization of tasks about the quality patterns that can be integrated on the ETL process, using as input user-defined goals.

Planner: Using a set of available patters, it conducts a pre-selection of highest ranked pattern combinations, based on heuristics and cost models, as adjusted from real execution and simulation. The user selects one combination.
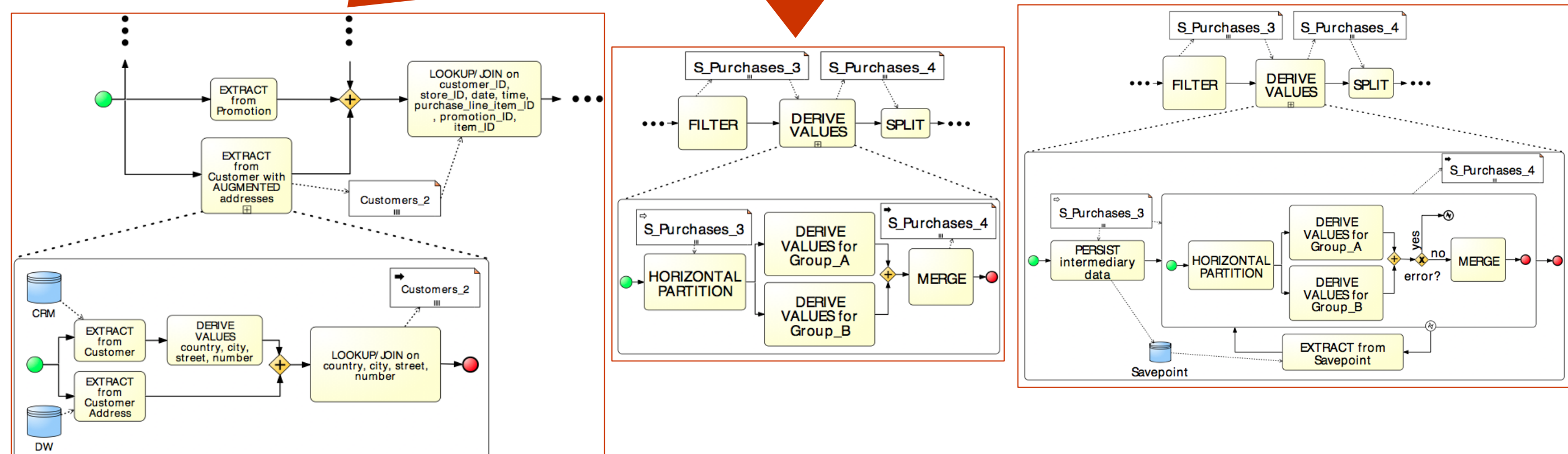
### • Deployment & Execution

Deployment Generator : Translates logical model to platform specific model.

ETL Process Engine : Executes ETL and keeps traces for providing historical analytics.





**Use Case: TPC-DS store sales maintenance ETL process**

improve Data Quality

improve Performance

improve Reliability

| Characteristic | Sub-characteristic | Measure |
|---|---|---|
| performance | time efficiency | • Process cycle time<br>• Average latency per tuple |
| | capacity | • Throughput of regular execution |
| data quality | data consistency | • % of tuples that violate business rules<br>• % of duplicates |
| | data freshness | • Request time - Time of last update<br>• 1 / (1 - age * Frequency of updates) |
| reliability | availability | • Mean Time Between Failures (MTBF)<br>• Uptime of ETL process |
| | recoverability | • Number of recovery points used<br>• % of successfully resumed workflow executions<br>• Mean time to repair (MTTR) |
| manageability | maintainability | • Length of process workflow's longest path<br>• Coupling of process workflow<br># of merge elements in the process model |
| | testability | • Cyclomatic Complexity of the ETL process workflow |