

# PDQ: A PLATFORM FOR REFORMULATING QUERIES

Michael Benedikt, George Konstantinidis, Efthymia Tsamoura - Oxford University

## Problem statement

Given a conjunctive query

$$Q := \exists x_1 \dots x_k \wedge_i R_i$$

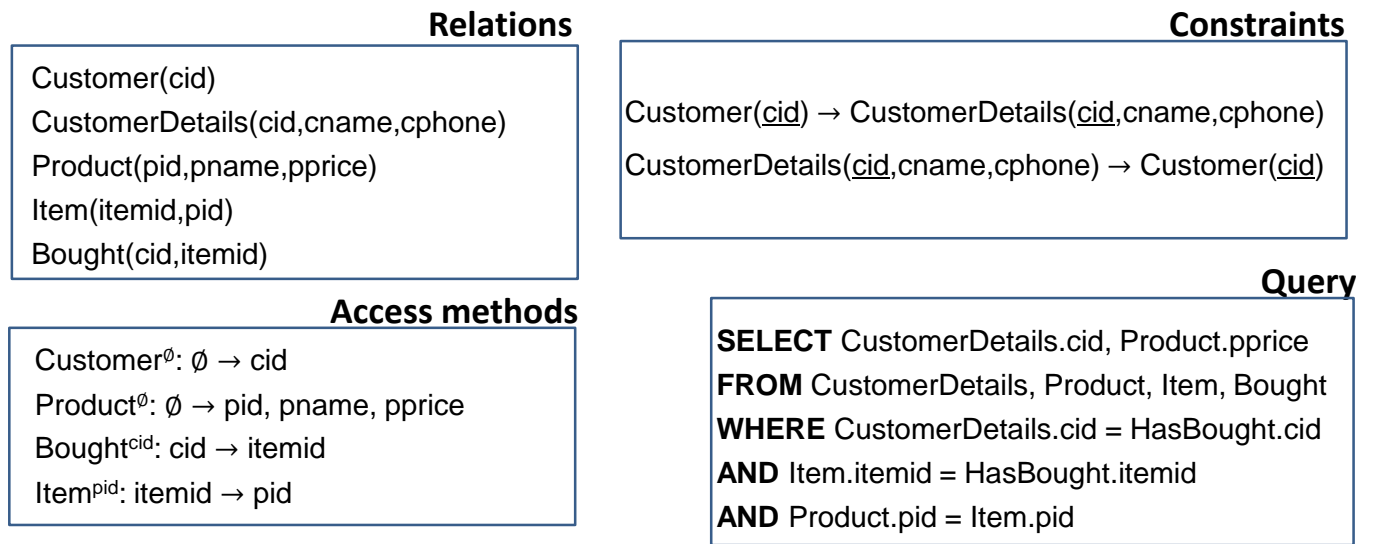
over a schema with **access restrictions** over the schema relations and **integrity constraints**  $\Sigma$  and a **cost function**  $C$  find a plan that **minimises the cost to answer  $Q$  under  $C$** .

Access restrictions over relations require values of certain attributes to be given as inputs to access relations.

The integrity constraints  $\Sigma$  are given by tuple-generating dependencies (TGDs):

$$\forall x_1, \dots, x_k \wedge_i A_i \rightarrow \exists y_1, \dots, y_l \wedge_j B_j$$

## Example



$Q := \exists pid, pname, itemid \text{ CustomerDetails}(cid, cname, cphone), \text{Bought}(cid, itemid), \text{Product}(pid, pname, pprice), \text{Item}(itemid, pid)$

$Q^* := \text{CustomerDetails}(cid_0, cname_0, cphone_0), \text{Bought}(cid_0, itemid_0), \text{Product}(pid_0, pname_0, pprice_0), \text{Item}(itemid_0, pid_0)$

**Hidden database:** CustomerDetails(cid<sub>0</sub>,cname<sub>0</sub>,cphone<sub>0</sub>), Bought(cid<sub>0</sub>,itemid<sub>0</sub>), Product(pid<sub>0</sub>,pname<sub>0</sub>,pprice<sub>0</sub>), Item(itemid<sub>0</sub>,pid<sub>0</sub>), Customer(cid<sub>0</sub>)

### Accessibility axioms

Customer(cid) → InferredAccessibleCustomer(cid) ∧ Accessible(cid)  
 Product(pid,pname,pprice) → InferredAccessibleProduct(pid,pname,pprice) ∧ Accessible(pid) ∧ Accessible(pname) ∧ Accessible(pprice)  
 Accessible(pid) ∧ Item(itemid,pid) → InferredAccessibleItem(itemid,pid) ∧ Accessible(itemid)  
 Accessible(cid) ∧ Bought(cid,itemid) → InferredAccessibleBought(cid,itemid) ∧ Accessible(itemid)

**Status: Unsuccessful, closed**

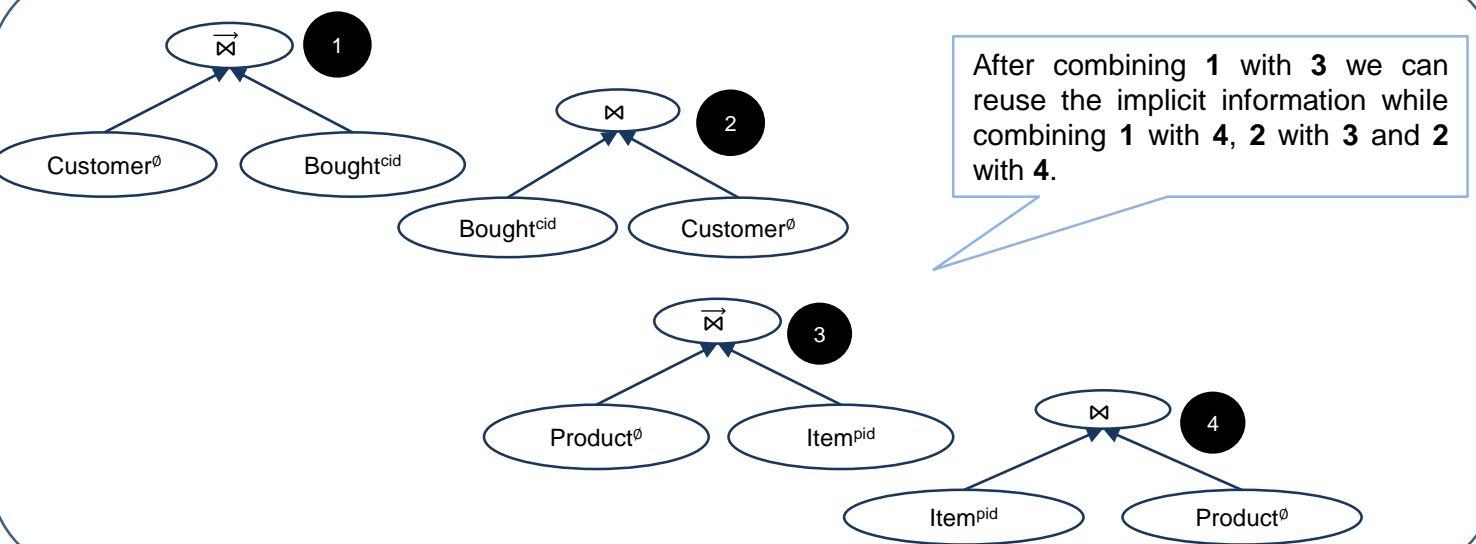
**Inputs:** ∅  
**Implicit information:** Customer(cid<sub>0</sub>), CustomerDetails(cid<sub>0</sub>,cname<sub>0</sub>,cphone<sub>0</sub>), InferredAccessibleCustomer(cid<sub>0</sub>), Accessible={cid<sub>0</sub>}

**Status: Unsuccessful, open**

**Inputs:** cid→cid<sub>0</sub>  
**Implicit information:** Bought(cid<sub>0</sub>,itemid<sub>0</sub>), Customer(cid<sub>0</sub>), CustomerDetails(cid<sub>0</sub>,cname<sub>0</sub>,cphone<sub>0</sub>), InferredAccessibleBought(cid<sub>0</sub>,itemid<sub>0</sub>), InferredAccessibleCustomer(cid<sub>0</sub>), Accessible={cid<sub>0</sub>,itemid<sub>0</sub>}

**Status: Unsuccessful, open**

**Inputs:** cid→cid<sub>0</sub>  
**Implicit information:** Bought(cid<sub>0</sub>,itemid<sub>0</sub>), InferredAccessibleBought(cid<sub>0</sub>,itemid<sub>0</sub>), Accessible={cid<sub>0</sub>,itemid<sub>0</sub>}



## Building up plans

### Access plans

- Plan language
- Access operators
- Dependent join operators
- Joins, selections and projections

Each plan is associated with

- the inputs required to perform the plan
- implicit information found through reasoning

A plan can be

- Open (requires inputs to run) or closed (can be run stand alone)
- Successful (equivalent to query) or unsuccessful

### Deriving implicit information

#### Preprocessing step

- Create "hidden database" of facts by forming canonical database of  $Q^*$  and taking consequences under  $\Sigma$ .
- Augment  $\Sigma$  with accessibility axioms and inferred accessible copies of the constraints in  $\Sigma$ .

#### Implicit information of accesses

- consequences of hidden facts exposed by operators

#### Implicit information of BinaryOperator(Plan<sub>1</sub>, Plan<sub>2</sub>)

- consequences of [implicit information of Plan<sub>1</sub> ∪ implicit information of Plan<sub>2</sub>] under augmented  $\Sigma$

### Pruning out sub-plans

- Keep the best plan within some class
- group plans by input and by implicit information
- discard plans with the same or less implicit information and higher cost

Prune out unpromising sub-plans

- prune out plans with cost higher than the best closed and successful plan found

### Speeding-up planning

- Search the plan space in parallel
- Speeding up reasoning
  - group plans based on their implicit information
  - when combining plans from two groups reason only the first time a composite plan is created from these two groups
- Prune out plans with cost higher than the best plan prior to reasoning

## Prototype with LogicBlox

LogicBlox is a relational database geared toward analytics and predictions.

LogicBlox programs are implemented in the proprietary language LogicQL, derived from Datalog.

A PDQ server instance is started upon the creation or opening of an LB workspace.

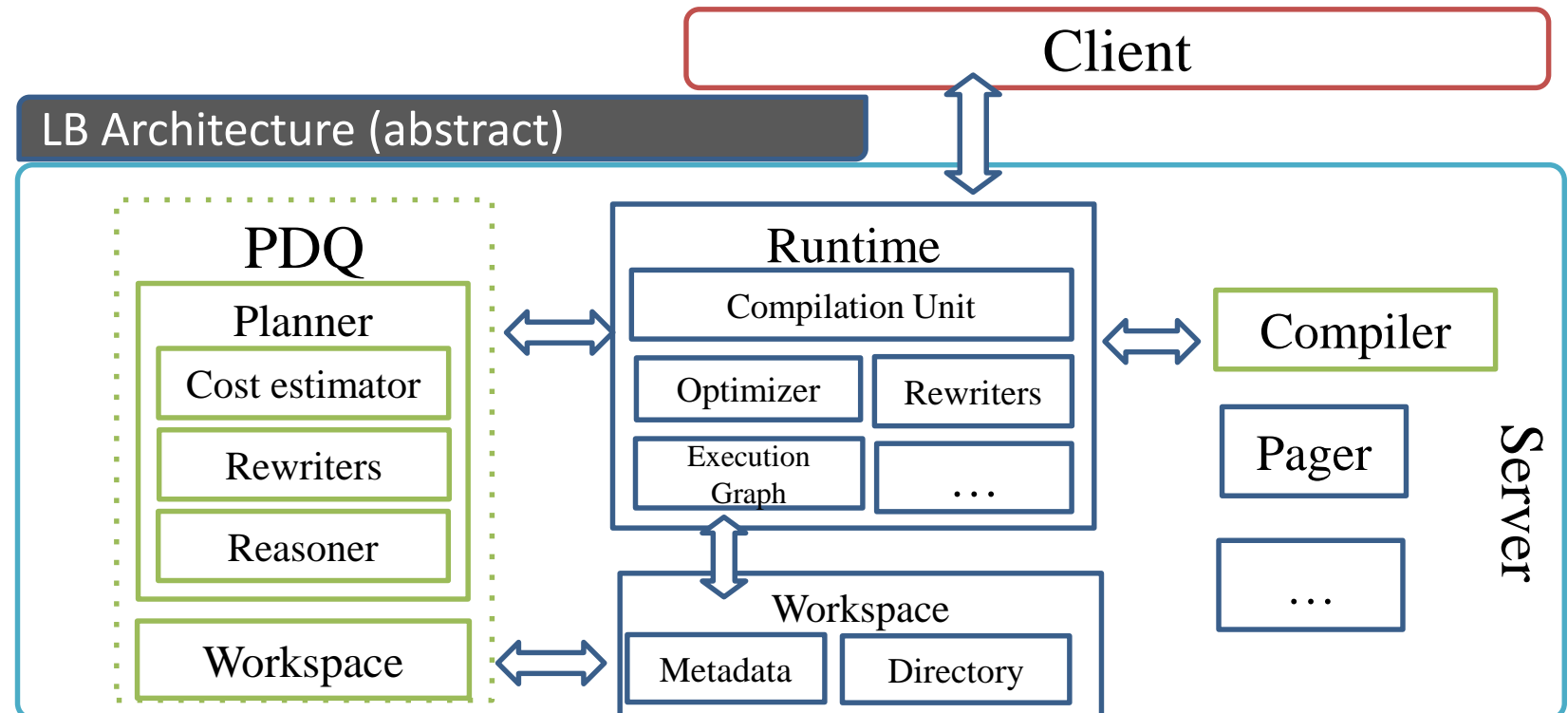
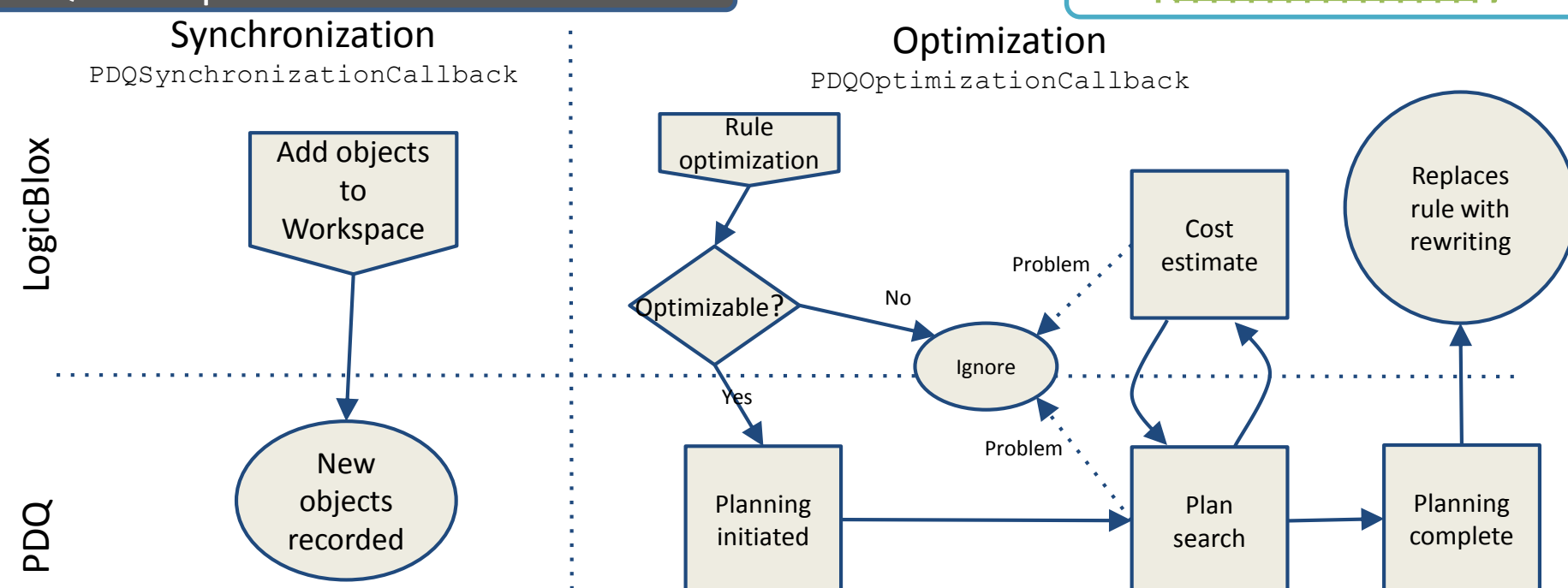
PDQ is initialized with all relevant information from the workspace such as **views** and **constraints**.

PDQ **optimizes** LB rules, by offering equivalent rewritings with different cost.

During planning **PDQ asks LB** for its estimation about the **cost** of a specific subplan.

When a transaction ends successfully, PDQ is updated to account for the objects that have been created or destroyed.

### PDQ – LB Optimization Workflow



### REFERENCES

- Benedikt, M., Leblay, J., Cate, B.T. and Tsamoura, E., 2016. Generating Plans from Proofs: The Interpolation-based Approach to Query Reformulation. *Synthesis Lectures on Data Management*, 8(1), pp.1-205.
- Benedikt, M., Ten Cate, B. and Tsamoura, E., 2016. Generating Plans from Proofs. *ACM Transactions on Database Systems (TODS)*, 40(4), p.22.
- Benedikt, M., Leblay, J. and Tsamoura, E., 2015. Querying with access patterns and integrity constraints. *Proceedings of the VLDB Endowment*, 8(6), pp.690-701.
- Benedikt, M., Leblay, J. and Tsamoura, E., 2014. PDQ: Proof-driven query answering over web-based data. *Proceedings of the VLDB Endowment*, 7(13), pp.1553-1556.
- Benedikt, M., Ten Cate, B. and Tsamoura, E., 2014, June. Generating low-cost plans from proofs. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems* (pp. 200-211). ACM.