CoAl : Incremental requirement-driven design and deployment of data intensive flows

Petar Jovanovic, Oscar Romero, Alberto Abelló

Universitat Politècnica de Catalunya, BarcelonaTech [petar | oromero | aabello]@essi.upc.edu

Alkis Simitsis

HP Labs, Palo Alto, CA, USA alkis@hp.com

Problem under study

- Incrementally integrating individual data flows into a unified flow satisfying the entire set of inf. requirements
- \succ Maximizing the reuse of the existing data flows
- > Lowering overall execution time by sharing data and computation
- > Considering execution costs



year(o_date) = 2013

l orderkey =

lineitem

line-

item

order

custo

mer

natio

Logical representation of data intensive flows

revenue =



Challenges & Solution

1) Incremental advancement

- \succ To guarantee semantical overlapping when comparing two operations
- > We must ensure that all the predeceasing operations of both operations coincide
- Bottom-up (sources-to-target) algorithm for consolidating data flows

2) Operation reordering

- Swap, distribute/factorize, merge/split, association
- Generic equivalence rules to guarantee the equivalence of reordering transformations
- > In terms of properties of data flow operations (i.e., Schema, Value, Order)
 - $o = (\mathbb{I}, \mathbb{O}, \mathbb{S}, \operatorname{Pre}, \operatorname{Post})$
 - $\mathbf{Pre} = (S_{pre}, V_{pre}, O_{pre})$
 - $\mathbf{Post} = (S_{post_gen}, S_{post_rem}, V_{post}, O_{post})$
- \succ Check for the conflicts

✓ Schema
$$S_{post_rem_B} \cap S_{pre_A} = \emptyset \land S_{post_gen_A} \cap S_{pre_B} = \emptyset$$

✓ Value $V_{post_B} \cap V_{pre_A} = \emptyset \land V_{post_A} \cap V_{pre_B} = \emptyset$

O_A

$$\checkmark Order O_{post_B} \Rightarrow \neg O_{pre_A} \land O_{post_A} \Rightarrow \neg O_{pre_B}$$
Operation comparison

- Generic operation comparison
- Full and partial operation



Case 2: operation reordering & **verial overlap**







Further reading:

- Oscar Romero, Alkis Simitsis, Alberto Abelló: GEM: Requirement-Driven Generation of ETL and 1. Multidimensional Conceptual Designs. DaWaK 2011: 80-95
- Petar Jovanovic, Oscar Romero, Alkis Simitsis, Alberto Abelló: Integrating ETL Processes from 2. Information Requirements. DaWaK 2012: 65-80
- Petar Jovanovic, Oscar Romero, Alkis Simitsis, Alberto Abelló: Requirement-Driven Creation and 3. Deployment of Multidimensional and ETL Designs. ER Workshops 2012: 391-395
- Petar Jovanovic, Alkis Simitsis, Kevin Wilkinson: Engine independence for logical analytic flows. 4. ICDE 2014: 1060-1071
- Petar Jovanovic, Alkis Simitsis, Kevin Wilkinson: BabbleFlow: a translator for analytic data flow 5. programs. SIGMOD Conference 2014: 713-716
- Petar Jovanovic, Oscar Romero, Alkis Simitsis, Alberto Abelló, Daria Mayorova: A requirement-6. driven approach to the design and evolution of data warehouses. Inf. Syst. 44: 94-119 (2014)

Fourth European Business Intelligence Summer School (eBISS 2014) – Berlin, Germany