

Functorial Data Integration

Categorical Informatics, Inc.

info@catinf.com

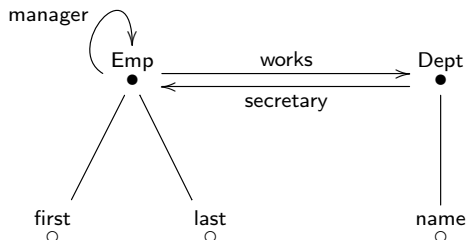
May 1, 2016

\mathcal{C}_i

Outline

- ▶ Review of Functorial Data Migration
 - ▶ Associated with a schema mapping $F : S \rightarrow T$ are three adjoint data migration operations $\Delta_F : T\text{-inst} \rightarrow S\text{-inst}$ (project), $\Pi_F : S\text{-inst} \rightarrow T\text{-inst}$ (join), and $\Sigma_F : S\text{-inst} \rightarrow T\text{-inst}$ (union).
- ▶ Functorial Data Integration
 - ▶ The unary operations Δ, Σ, Π are sufficient to query data and to migrate data, but we require binary operations to *integrate* data from multiple sources.
 - ▶ Integration is expressed category-theoretically using a *pushout* construction.
- ▶ This talk:
 - ▶ Pushouts
 - ▶ How to use pushouts in data warehousing
 - ▶ Extended example
- ▶ The example is discussed in <http://arxiv.org/abs/1503.03571> and is included in FQL as “O Integration”.

Review of the Functorial Data Model



$\text{Emp.manager.works} = \text{Emp.works}$

$\text{Dept.secretary.works} = \text{Dept}$

| Emp | | | | |
|-----|-----|-------|-------|------|
| ID | mgr | works | first | last |
| 101 | 103 | q10 | Al | Akin |
| 102 | 102 | x02 | Bob | Bo |
| 103 | 103 | q10 | Carl | Cork |

| Dept | | |
|------|-----|------|
| ID | sec | name |
| q10 | 102 | CS |
| x02 | 101 | Math |

Review of Functorial Data Migration

- ▶ A functor $F: S \rightarrow T$ is a constraint-respecting mapping:

$$\text{nodes}(S) \rightarrow \text{nodes}(T) \quad \text{edges}(S) \rightarrow \text{paths}(T)$$

and it induces three adjoint data migration functors:

- ▶ $\Delta_F: T\text{-inst} \rightarrow S\text{-inst}$ (like project)

$$\begin{array}{ccc} S & \xrightarrow{F} & T & \xrightarrow{I} & \mathbf{Set} \\ & \searrow & & \nearrow & \\ & & \Delta_F(I) := I \circ F & & \end{array}$$

- ▶ $\Pi_F: S\text{-inst} \rightarrow T\text{-inst}$ (like join)

$$\Delta_F \dashv \Pi_F$$

- ▶ $\Sigma_F: S\text{-inst} \rightarrow T\text{-inst}$ (like outer disjoint union then quotient)

$$\Sigma_F \dashv \Delta_F$$

Data Integration

Exxon
●

Mobil
●

Data Integration



Data Integration

TerroristEvents
●

WeatherData
●

Data Integration



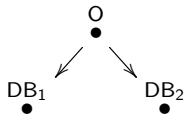
Pushouts

Given two databases:

DB₁ DB₂
● ●

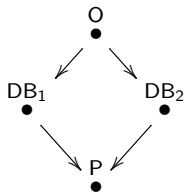
Pushouts

And an overlap between them:



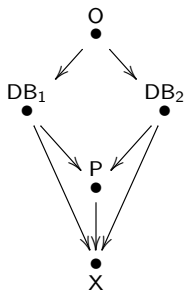
Pushouts

Their pushout is an integrated database:



Pushouts

That is universal among all possible integrated databases:



Pushouts in FQL

- ▶ Pushouts can be defined, but may not exist, in any category.
- ▶ FQL uses automated theorem proving techniques to compute pushouts.
- ▶ Pushouts play a role similar to “the chase” in relational database theory.

Using Pushouts for Data Integration

- Step 1: integrate schemas. Given input schemas S_1, S_2 , an overlap schema S , and mappings F_1, F_2 :

$$S_1 \xleftarrow{F_1} S \xrightarrow{F_2} S_2$$

we propose to use their pushout T as the integrated schema:

$$S_1 \xrightarrow{G_1} T \xleftarrow{G_2} S_2$$

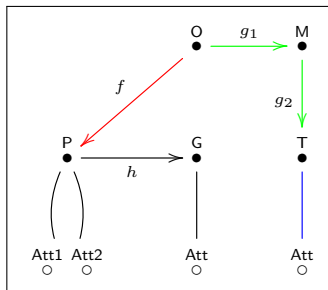
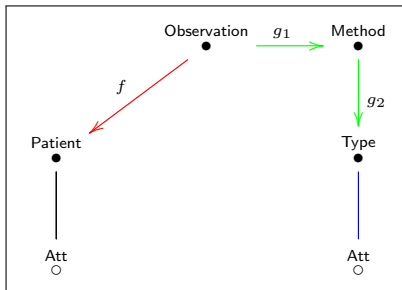
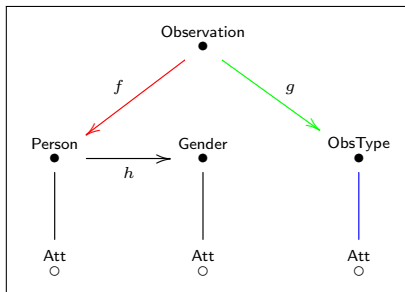
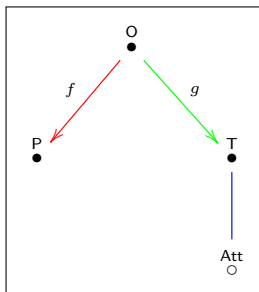
- Step 2: integrate data. Given input S_1 -instance I_1 , S_2 -instance I_2 , overlap S -instance I and transforms $h_1: \Sigma_{F_1}(I) \rightarrow I_1$ and $h_2: \Sigma_{F_2}(I) \rightarrow I_2$, we propose to use the pushout of:

$$\Sigma_{G_1}(I_1) \xleftarrow{\Sigma_{G_1}(h_1)} (\Sigma_{G_1 \circ F_1}(I) = \Sigma_{G_2 \circ F_2}(I)) \xrightarrow{\Sigma_{G_2}(h_2)} \Sigma_{G_2}(I_2)$$

as the integrated T -instance.

Using Pushouts for Data Integration

- ▶ For every two databases to integrate, the strategy just described is a design pattern:
- ▶ The input is one schema (S), two schema mappings (F_1, F_2), one database instance (I), and two database mappings (h_1, h_2).
- ▶ If desired, (S, F_1, F_2) can be constructed automatically using schema matching techniques.
- ▶ If desired, (I, h_1, h_2) can be constructed automatically using data matching techniques.
- ▶ The output is an integrated schema T and a “universal” (best possible) integrated T -instance.



| O | | | P | T | |
|----|---|---|----------|-----------|--------|
| ID | f | g | ID | ID | Att |
| | | | <i>p</i> | <i>bp</i> | BP |
| | | | | <i>wt</i> | Weight |

→

| Gender | | | ObsType | | |
|-------------|-----------|-----------|-----------|-----------|----------|
| ID | f | Att | ID | Att | |
| | <i>f</i> | F | | <i>bp</i> | BP |
| | <i>m</i> | M | | <i>wt</i> | Weight |
| | | | | <i>hr</i> | HR |
| Observation | | | Person | | |
| ID | f | g | ID | Att | h |
| <i>o5</i> | <i>pe</i> | <i>bp</i> | <i>pa</i> | Paul | <i>m</i> |
| <i>o6</i> | <i>pa</i> | <i>hr</i> | <i>p</i> | Peter | <i>m</i> |
| <i>o7</i> | <i>pe</i> | <i>wt</i> | | | |

↓

| Method | | | Type | |
|-------------|-----------|-----------|-----------|--------|
| ID | g2 | | ID | Att |
| <i>m1</i> | <i>bp</i> | | <i>bp</i> | BP |
| <i>m2</i> | <i>bp</i> | | <i>wt</i> | Weight |
| <i>m3</i> | <i>wt</i> | | | |
| <i>m4</i> | <i>wt</i> | | | |
| Observation | | | Patient | |
| ID | f | g | ID | Att |
| <i>o1</i> | <i>p</i> | <i>m1</i> | <i>j</i> | Jane |
| <i>o2</i> | <i>p</i> | <i>m2</i> | <i>p</i> | Pete |
| <i>o3</i> | <i>j</i> | <i>m3</i> | | |
| <i>o4</i> | <i>j</i> | <i>m1</i> | | |

→

| M | | | O | | |
|---------------|--------------------|-------------------|-------------|-----------|---------------|
| ID | g2 | | ID | f | g1 |
| <i>g1(o5)</i> | <i>bp</i> | | <i>o1</i> | <i>p</i> | <i>m1</i> |
| <i>g1(o6)</i> | <i>wt</i> | | <i>o2</i> | <i>p</i> | <i>m2</i> |
| <i>g1(o7)</i> | <i>hr</i> | | <i>o3</i> | <i>j</i> | <i>m3</i> |
| <i>m1</i> | <i>bp</i> | | <i>o4</i> | <i>j</i> | <i>m1</i> |
| <i>m2</i> | <i>bp</i> | | <i>o5</i> | <i>p</i> | <i>g1(o5)</i> |
| <i>m3</i> | <i>wt</i> | | <i>o6</i> | <i>pa</i> | <i>g1(o6)</i> |
| <i>m4</i> | <i>wt</i> | | <i>o7</i> | <i>p</i> | <i>g1(o7)</i> |
| G | | | T | | |
| ID | Att | | ID | Att | |
| <i>f</i> | F | | <i>bp</i> | BP | |
| <i>m</i> | M | | <i>wt</i> | Weight | |
| <i>h(j)</i> | Att(<i>h(j)</i>) | | <i>hr</i> | HR | |
| P | | | h | | |
| ID | Att1 | Att2 | ID | Att | h |
| <i>j</i> | Att1(<i>j</i>) | Jane | <i>h(j)</i> | | |
| <i>pa</i> | Paul | Att2(<i>pa</i>) | <i>m</i> | | |
| <i>p</i> | Peter | Pete | <i>m</i> | | |

Conclusion

- ▶ Pushouts express “union with overlap, then quotient”.
- ▶ We propose a pattern for data integration:
 - ▶ Use pushouts to integrate source schemas.
 - ▶ Use pushouts to integrate source data.
- ▶ The example is discussed in <http://arxiv.org/abs/1503.03571> and is included in FQL as “O Integration”.