

Computational Schemas

Categorical Informatics, Inc.

info@catinf.com

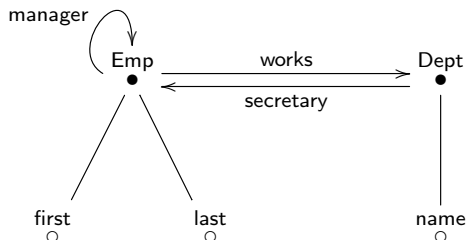
May 1, 2016

\mathcal{C}_i

Outline

- ▶ Review of the Functorial Data Model
 - ▶ Schemas are categories, instances are functors.
- ▶ This talk:
 - ▶ More recent languages in the FQL family can *compute* in schemas, queries, migrations, etc.
 - ▶ Computational primitives are user-defined using equational logic, and, optionally, a java binding.
 - ▶ Examples: boolean algebra, arithmetic, strings, images, machine learning, entity resolution.
 - ▶ Our approach is a kind of dual to LINQ; see <http://arxiv.org/abs/1511.06459>.
 - ▶ See <http://arxiv.org/abs/1503.03571> for more examples.

Review of the Functorial Data Model



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

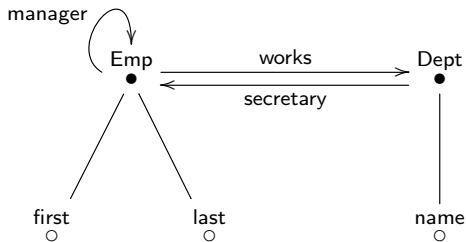
$\text{Dept.secretary.works} = \text{Dept}$

| 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 |

Schemas as Equational Theories

- ▶ We call the graph of a schema its *signature*.



- ▶ In this talk we will write the path equivalences in a schema:

$$\text{Emp.manager.works} = \text{Emp.works} \quad \text{Dept.secretary.works} = \text{Dept}$$

using logical notation:

$$\forall e : \text{Emp. works}(\text{manager}(e)) = \text{works}(e) \quad \forall d : \text{Dept. works}(\text{secretary}(d)) = d$$

- ▶ Attributes and foreign keys are unary function symbols:

$$\text{name} : \text{Dept} \rightarrow \text{String} \quad \text{manager} : \text{Emp} \rightarrow \text{Emp} \quad \textit{etc.}$$

The Type Side

- ▶ A schema exists in an ambient equational theory, called the *type side*.
- ▶ For example, a type side for arithmetic has function symbols:

$$0 : \text{Nat} \quad S : \text{Nat} \rightarrow \text{Nat} \quad + : \text{Nat}, \text{Nat} \rightarrow \text{Nat}$$

and equations:

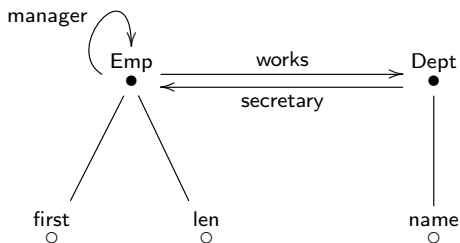
$$\forall n : \text{Nat}. n + 0 = n \quad \forall n, m : \text{Nat}. n + (S(m)) = S(n + m)$$

- ▶ We will assume function symbols for string manipulation, including:

$$\text{length} : \text{String} \rightarrow \text{Nat}$$

Computations in Schemas

- ▶ We can use a type side when we define a schema:



$$\forall e : \text{Emp. works}(\text{manager}(e)) = \text{works}(e) \quad \forall d : \text{Dept. works}(\text{secretary}(d)) = d$$
$$\forall e : \text{Emp. len}(e) = \text{length}(\text{first}(e))$$

| Emp | | | | |
|-----|-----|-------|-------|-----|
| ID | mgr | works | first | len |
| 101 | 103 | q10 | Al | 2 |
| 102 | 102 | x02 | Bob | 3 |
| 103 | 103 | q10 | Carl | 4 |

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

Derived Columns

- ▶ Every “partial” instance on a schema has a unique “completion”.
- ▶ So an equation such as

$$\forall e : \text{Emp. } \text{len}(e) = \text{length}(\text{first}(e))$$

will cause the FQL tool to compute the len attribute from the first attribute whenever the len attribute value is missing in an instance.

- ▶ The FQL tool never requires complete instances, since it can always complete a partial instance.

Example Partial Instance

- Partial Instance (Dept table omitted):

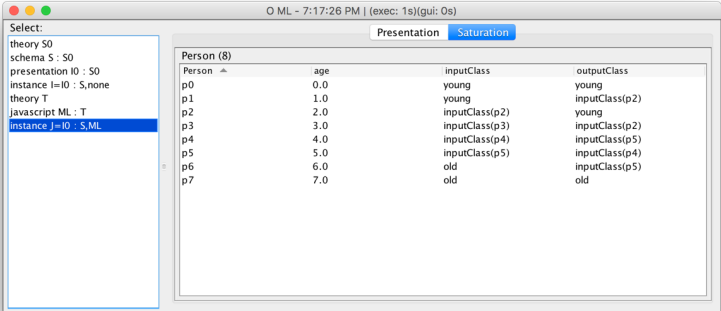
| Emp | | | | |
|-----|-----|-------|-------|------|
| ID | mgr | works | first | last |
| 101 | 103 | q10 | | Akin |
| 102 | 102 | x02 | Bob | Bo |
| 103 | | q10 | Carl | Cork |

- Completion (Dept table omitted):

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

Binding Primitives to Existing Code

- ▶ Function symbols can be bound to user-defined java functions.
- ▶ Here is a screen shot of an FQL instance that invokes a 3rd party machine learning library to classify input ages into old and young:



The screenshot shows a window titled "O ML - 7:17:26 PM | (exec: 1s)(gui: 0s)". On the left, a "Select:" panel lists several items, with "instance J=I0 : S,ML" selected. The main area displays a table titled "Person (8)" with two tabs: "Presentation" and "Saturation". The table has four columns: "Person", "age", "inputClass", and "outputClass".

| Person | age | inputClass | outputClass |
|--------|-----|----------------|----------------|
| p0 | 0.0 | young | young |
| p1 | 1.0 | young | inputClass(p2) |
| p2 | 2.0 | inputClass(p2) | young |
| p3 | 3.0 | inputClass(p3) | inputClass(p2) |
| p4 | 4.0 | inputClass(p4) | inputClass(p5) |
| p5 | 5.0 | inputClass(p5) | inputClass(p4) |
| p6 | 6.0 | old | inputClass(p5) |
| p7 | 7.0 | old | old |

- ▶ This example illustrates the interplay between computation and unknown values.

Conclusion

- ▶ All schemas live in a user-defined computational context.
- ▶ Computational primitives, attributes, and foreign keys can mix freely and can optionally be bound to user-provided java code.
- ▶ Columns to be derived from other columns on the fly.
- ▶ Computation and unknown values work together.