

Algebraic Databases

Patrick Schultz, David Spivak
MIT

Ryan Wisnesky
Categorical Informatics

LambdaConf 2017

Introduction

- ▶ This talk describes a new “algebraic” way to formalize databases based on category theory.
- ▶ Category theory was designed to migrate **theorems** from one **area of mathematics** to another, but researchers at MIT developed a way to use it to migrate **data** from one **computer system** to another.
- ▶ Research has culminated in an open-source prototype ETL and data integration tool, AQL (Algebraic Query Language), available at categoricaldata.net/aql.html. (These slides are also there.)
- ▶ Outline:
 - ▶ Review of basic category theory.
 - ▶ Introduction to AQL.
 - ▶ AQL demo.
 - ▶ Optional interlude: additional AQL constructions.
 - ▶ How AQL instances model the simply-typed λ -calculus.

Category Theory

- ▶ A category \mathcal{C} consists of
 - ▶ a set of *objects*, $\text{Ob}(\mathcal{C})$
 - ▶ for all $X, Y \in \text{Ob}(\mathcal{C})$, a set $\mathcal{C}(X, Y)$ of *morphisms* a.k.a *arrows*
 - ▶ for all $X \in \text{Ob}(\mathcal{C})$, a morphism $id \in \mathcal{C}(X, X)$
 - ▶ for all $X, Y, Z \in \text{Ob}(\mathcal{C})$, a function $\circ: \mathcal{C}(Y, Z) \times \mathcal{C}(X, Y) \rightarrow \mathcal{C}(X, Z)$ s.t.

$$f \circ id = f \quad id \circ f = f \quad (f \circ g) \circ h = f \circ (g \circ h)$$

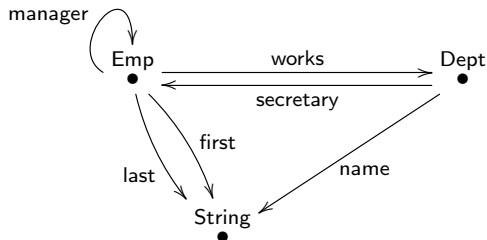
- ▶ The category **Set** has sets as objects and functions as arrows, and the “category” **Haskell** has types as objects and programs as arrows.
-

- ▶ A functor $F: \mathcal{C} \rightarrow \mathcal{D}$ between categories \mathcal{C}, \mathcal{D} consists of
 - ▶ a function $\text{Ob}(\mathcal{C}) \rightarrow \text{Ob}(\mathcal{D})$
 - ▶ for all $X, Y \in \text{Ob}(\mathcal{C})$, a function $\mathcal{C}(X, Y) \rightarrow \mathcal{D}(F(X), F(Y))$ s.t.

$$F(id) = id \quad F(f \circ g) = F(f) \circ F(g)$$

- ▶ The functor $\mathcal{P}: \mathbf{Set} \rightarrow \mathbf{Set}$ takes each set to its power set, and the functor $\text{List}: \mathbf{Haskell} \rightarrow \mathbf{Haskell}$ takes each type t to the type $\text{List } t$.

Schemas and Instances



$\text{works} \circ \text{manager} = \text{works}$

$\text{works} \circ \text{secretary} = id$

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	101	CS
x02	102	Math

String	
ID	
Al	
Bob	
...	

Categorical Semantics of Schemas and Instances

- ▶ The meaning of a schema S is a category $\llbracket S \rrbracket$.
 - ▶ $\text{Ob}(\llbracket S \rrbracket)$ is the nodes of S .
 - ▶ For all nodes X, Y , $\llbracket S \rrbracket(X, Y)$ is the set of finite paths $X \rightarrow Y$, modulo the path equivalences in S .
 - ▶ Path equivalence in S may not be decidable! (“the word problem”)
- ▶ A morphism of schemas (a “**schema mapping**”) $S \rightarrow T$ is a functor $\llbracket S \rrbracket \rightarrow \llbracket T \rrbracket$.
 - ▶ It can be defined as an equation-preserving function:

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

- ▶ An S -instance is a functor $\llbracket S \rrbracket \rightarrow \mathbf{Set}$.
 - ▶ It can be defined as a set of tables, one per node in S and one column per edge in S , satisfying the path equivalences in S .
- ▶ A morphism of S -instances $I \rightarrow J$ (a “**data mapping**”) is a natural transformation $I \rightarrow J$.
 - ▶ Instances on S and their mappings form a category, written $S\text{-inst}$.

Functorial Data Migration

A schema mapping $F: S \rightarrow T$ induces three data migration functors:

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

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

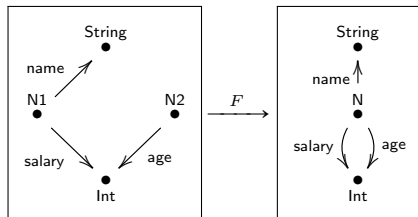
-
- ▶ $\Pi_F: S\text{-inst} \rightarrow T\text{-inst}$ (right adjoint to Δ_F ; like join)

$$\forall I, J. \quad S\text{-inst}(\Delta_F(I), J) \cong T\text{-inst}(I, \Pi_F(J))$$

-
- ▶ $\Sigma_F: S\text{-inst} \rightarrow T\text{-inst}$ (left adjoint to Δ_F ; like outer union then merge)

$$\forall I, J. \quad S\text{-inst}(J, \Delta_F(I)) \cong T\text{-inst}(\Sigma_F(J), I)$$

Δ (Project)



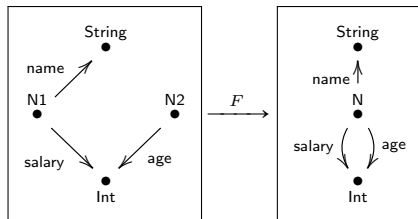
N1		
ID	name	salary
1	Alice	\$100
2	Bob	\$250
3	Sue	\$300

N2	
ID	age
4	20
5	20
6	30

Δ_F

N			
ID	name	salary	age
a	Alice	\$100	20
b	Bob	\$250	20
c	Sue	\$300	30

II (Product)



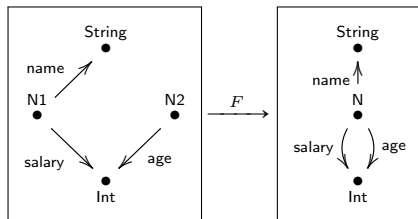
N1		
ID	name	salary
1	Alice	\$100
2	Bob	\$250
3	Sue	\$300

N2	
ID	age
4	20
5	20
6	30

Π_F

N			
ID	name	salary	age
a	Alice	\$100	20
b	Alice	\$100	20
c	Alice	\$100	30
d	Bob	\$250	20
e	Bob	\$250	20
f	Bob	\$250	30
g	Sue	\$300	20
h	Sue	\$300	20
i	Sue	\$300	30

Σ (Outer Union)



N1		
ID	Name	Salary
1	Alice	\$100
2	Bob	\$250
3	Sue	\$300

N2	
ID	Age
4	20
5	20
6	30

Σ_F

N			
ID	Name	Salary	Age
a	Alice	\$100	$null_1$
b	Bob	\$250	$null_2$
c	Sue	\$300	$null_3$
d	$null_4$	$null_5$	20
e	$null_6$	$null_7$	20
f	$null_8$	$null_9$	30

Unit of $\Sigma_F \dashv \Delta_F$

N1			N2	
ID	Name	Salary	ID	Age
1	Alice	\$100	4	20
2	Bob	\$250	5	20
3	Sue	\$300	6	30

 $\Sigma_F \rightarrow$

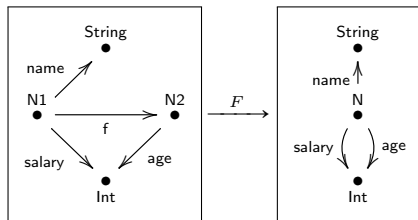
N			
ID	Name	Salary	Age
a	Alice	\$100	<i>null</i> ₁
b	Bob	\$250	<i>null</i> ₂
c	Sue	\$300	<i>null</i> ₃
d	<i>null</i> ₄	<i>null</i> ₅	20
e	<i>null</i> ₆	<i>null</i> ₇	20
f	<i>null</i> ₈	<i>null</i> ₉	30

 $\Delta_F \swarrow$

η

N1			N2	
ID	Name	Salary	ID	Age
a	Alice	\$100	a	<i>null</i> ₁
b	Bob	\$250	b	<i>null</i> ₂
c	Sue	\$300	c	<i>null</i> ₃
d	<i>null</i> ₄	<i>null</i> ₅	d	20
e	<i>null</i> ₆	<i>null</i> ₇	e	20
f	<i>null</i> ₈	<i>null</i> ₉	f	30

A Foreign Key



N1			
ID	name	salary	f
1	Alice	\$100	4
2	Bob	\$250	5
3	Sue	\$300	6

N2	
ID	age
4	20
5	20
6	30

Δ_F
 Π_F, Σ_F

N			
ID	name	salary	age
a	Alice	\$100	20
b	Bob	\$250	20
c	Sue	\$300	30

AQL Demo

- ▶ AQL implements Δ, Σ, Π , and more in software.
 - ▶ catinf.com
- ▶ The AQL “execution engine” is an automated theorem prover.
 - ▶ Value proposition: AQL catches mistakes at compile time that existing ETL / data integration tools catch at runtime – if at all.
 - ▶ Data import and export by JDBC-SQL and CSV.
- ▶ We are looking for collaborators for a “real-world pilot project”.

Interlude - Additional Constructions

- ▶ What is “algebraic” here?
- ▶ AQL vs SQL.
- ▶ Pivot.
- ▶ Non-equational data integrity constraints.
- ▶ Data integration via pushouts.
- ▶ AQL vs comprehension calculi.

Why “Algebraic”?

- ▶ A schema can be identified with an algebraic (equational) theory.

$\text{Emp Dept String} : \text{Type} \quad \text{first last} : \text{Emp} \rightarrow \text{String} \quad \text{name} : \text{Dept} \rightarrow \text{String}$

$\text{works} : \text{Emp} \rightarrow \text{Dept} \quad \text{mgr} : \text{Emp} \rightarrow \text{Emp} \quad \text{secr} : \text{Dept} \rightarrow \text{Emp}$

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

- ▶ This perspective makes it easy to add functions such as $+$: $\text{Int}, \text{Int} \rightarrow \text{Int}$ to a schema. See *Algebraic Databases*.

-
- ▶ An S -instance can be identified with the initial algebra of an algebraic theory extending S .

$101 \ 102 \ 103 : \text{Emp} \quad \text{q10} \ \text{x02} : \text{Dept}$

$\text{mgr}(101) = 103 \quad \text{works}(101) = \text{q10} \quad \dots$

- ▶ Treating instances as theories allows instances that are infinite or inconsistent (e.g., $\text{Alice} = \text{Bob}$).

AQL vs SQL

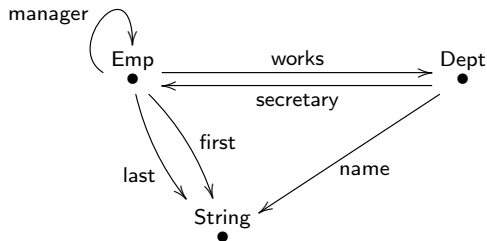
- ▶ Data migration triplets of the form

$$\Sigma_F \circ \Pi_G \circ \Delta_H$$

can be expressed using relational algebra and keygen, provided:

- ▶ F is a discrete op-fibration (ensures union compatibility).
- ▶ G is surjective on attributes (ensures domain independence).
- ▶ All categories are finite (ensures computability).
- ▶ The difference-free fragment of relational algebra can be expressed using such triplets. See *Relational Foundations*.
- ▶ Such triplets can be written in “foreign-key aware” SQL-ish syntax.

Select-From-Where Syntax



Find the name of every manager's department:

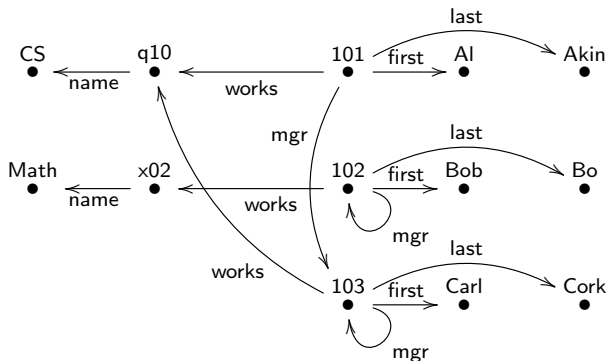
AQL

```
select e.manager.works.name  
from Emp as e
```

SQL

```
select d.name  
from Emp as e1, Emp as e2, Dept as d  
where e1.manager = e2.ID and  
e2.works = d.ID
```


Pivot (Instance \Leftrightarrow Schema)

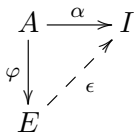


Emp				
ID	mgr	works	first	last
101	103	q10	Al	Akin
102	102	x02	Bob	Bo
103	103	q10	Carl	Cork

Dept	
ID	name
q10	CS
x02	Math

Richer Constraints

- ▶ Not all data integrity constraints are equational (e.g., keys).
- ▶ A data mapping $\varphi : A \rightarrow E$ defines a constraint: instance I satisfies φ if for every $\alpha : A \rightarrow I$ there exists an $\epsilon : E \rightarrow I$ s.t $\alpha = \epsilon \circ \varphi$.



- ▶ Most constraints used in practice can be captured the above way. E.g.,

$$\forall d_1, d_2 : \text{Dept. name}(d_1) = \text{name}(d_2) \rightarrow d_1 = d_2$$

is captured as

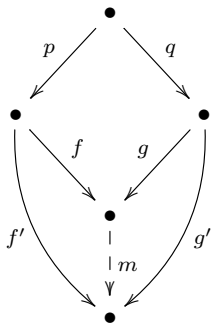
$$A(\text{Dept}) = \{d_1, d_2\} \quad A(\text{name})(d_1) = A(\text{name})(d_2)$$

$$E(\text{Dept}) = \{d\} \quad \varphi(d_1) = \varphi(d_2) = d$$

- ▶ See *Database Queries and Constraints via Lifting Problems* and *Algebraic Model Management*.

Pushouts

- ▶ A pushout of p, q is f, g s.t. for every f', g' there is a unique m s.t.:



- ▶ The category of schemas has all pushouts.
- ▶ For every schema S , the category S -inst has all pushouts.
- ▶ Pushouts of schemas, instances, and Σ are used together to integrate data - see *Algebraic Data Integration*.

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:

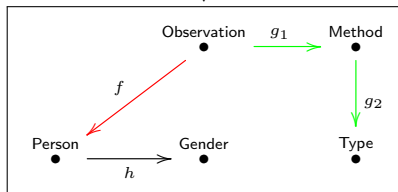
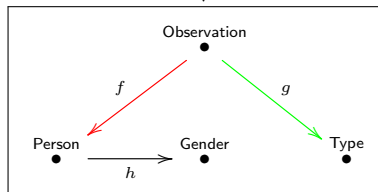
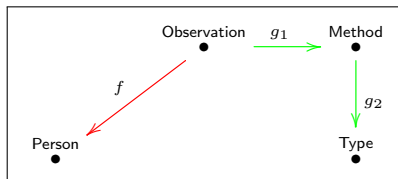
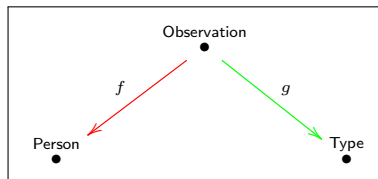
$$S_1 \xrightarrow{G_1} T \xleftarrow{G_2} S_2$$

- Step 2: integrate data. Given input S_1 -instance I_1 , S_2 -instance I_2 , overlap S -instance I and data mappings $h_1: \Sigma_{F_1}(I) \rightarrow I_1$ 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.

Schema Integration



Data Integration

Observation			Person		Type
ID	f	g	ID		ID
			<i>p</i>		BP
					Wt

→

Gender			Type	
ID			ID	
				BP
				Wt
				HR

Observation			Person	
ID	f	g	ID	h
<i>o5</i>	Peter	BP	Paul	M
<i>o6</i>	Paul	HR	<i>Peter</i>	M
<i>o7</i>	Peter	Wt		

↓

Method			Type
ID	g2		ID
<i>m1</i>	BP		BP
<i>m2</i>	BP		Wt
<i>m3</i>	Wt		
<i>m4</i>	Wt		

Observation			Person
ID	f	g1	ID
<i>o1</i>	Pete	<i>m1</i>	Jane
<i>o2</i>	Pete	<i>m2</i>	<i>Pete</i>
<i>o3</i>	Jane	<i>m3</i>	
<i>o4</i>	Jane	<i>m1</i>	

→

Method			Observation		
ID	g2		ID	f	g1
<i>null1</i>	BP		<i>o1</i>	Peter	<i>m1</i>
<i>null2</i>	Wt		<i>o2</i>	Peter	<i>m2</i>
<i>null3</i>	HR		<i>o3</i>	Jane	<i>m3</i>
<i>m1</i>	BP		<i>o4</i>	Jane	<i>m1</i>
<i>m2</i>	BP		<i>o5</i>	Peter	<i>null1</i>
<i>m3</i>	Wt		<i>o6</i>	Paul	<i>null2</i>
<i>m4</i>	Wt		<i>o7</i>	Peter	<i>null3</i>

Gender	Type	Person	
ID	ID	ID	h
	BP	Jane	<i>null4</i>
	Wt	Paul	M
<i>null4</i>	HR	<i>Peter</i>	M

AQL vs LINQ

- ▶ Treating entity sets as types rather than terms makes AQL a conceptual dual to comprehension calculi (e.g., LINQ). See *QINL: Query-Integrated Languages*.
- ▶ LINQ enriches programs with (schemas, queries and instances).
 - ▶ Collections are *terms*

Employee: Set Int manager: Set (Int × Int)

- ▶ e: Employee *is not* a judgment.
 - ▶ There *is* a term $\epsilon: \text{Int} \times \text{Set Int} \rightarrow \text{Bool}$.
- ▶ AQL enriches (schemas, queries and instances) with programs.
 - ▶ Collections are *types*

Employee: Type manager: Employee \rightarrow Employee

- ▶ e: Employee *is* a judgment.
 - ▶ There *is not* a term $\epsilon: \text{Employee} \times \text{Type} \rightarrow \text{Bool}$.
- ▶ LINQ is more popular, but AQL's style is common in Coq, Agda, etc.

AQL is “one level up” from LINQ

▸ LINQ

- Schemas are collection types over a base type theory

$$\text{Set } (\text{Int} \times \text{String})$$

- Instances are terms

$$\{(1, \text{CS})\} \cup \{(2, \text{Math})\}$$

- Data migrations are functions

$$\pi_1 : \text{Set } (\text{Int} \times \text{String}) \rightarrow \text{Set Int}$$

▸ AQL

- Schemas are type theories over a base type theory

$$\text{Dept, name: Dept} \rightarrow \text{String}$$

- Instances are term models (initial algebras) of theories

$$d_1, d_2 : \text{Dept, name}(d_1) = \text{CS, name}(d_2) = \text{Math}$$

- Data migrations are functors

$$\Delta_{\text{Dept}} : (\text{Dept, name: Dept} \rightarrow \text{String})\text{-inst} \rightarrow (\text{Dept})\text{-inst}$$

Part 2

- ▶ For every schema S , S -inst models simply-typed λ -calculus (STLC).
- ▶ The STLC is the core of typed functional languages ML, Haskell, etc.
- ▶ We will use the internal language of a cartesian closed category, which is equivalent to the STLC.
- ▶ Lots of “point-free” functional programming ahead.
- ▶ The category of schemas and mappings is also cartesian closed - see talk at Boston Haskell.

Categorical Abstract Machine Language (CAML)

- Types t :

$$t ::= 1 \mid t \times t \mid t^t$$

- Terms f, g :

$$id_t : t \rightarrow t \quad ()_t : t \rightarrow 1 \quad \pi_{s,t}^1 : s \times t \rightarrow s \quad \pi_{s,t}^2 : s \times t \rightarrow t$$

$$eval_{s,t} : t^s \times s \rightarrow t \quad \frac{f : s \rightarrow u \quad g : u \rightarrow t}{g \circ f : s \rightarrow t} \quad \frac{f : s \rightarrow t \quad g : s \rightarrow u}{(f, g) : s \rightarrow t \times u}$$

$$\frac{f : s \times u \rightarrow t}{\lambda f : s \rightarrow t^u}$$

- Equations:

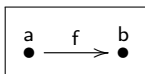
$$id \circ f = f \quad f \circ id = f \quad f \circ (g \circ h) = (f \circ g) \circ h \quad () \circ f = ()$$

$$\pi^1 \circ (f, g) = f \quad \pi^2 \circ (f, g) = g \quad (\pi^1 \circ f, \pi^2 \circ f) = f$$

$$eval \circ (\lambda f \circ \pi^1, \pi^2) = f \quad \lambda(eval \circ (f \circ \pi^1, \pi^2)) = f$$

Programming AQL in CAML

- ▶ For every schema S , the category $S\text{-inst}$ is cartesian closed.
 - ▶ Given a type t , you get an S -instance $[t]$.
 - ▶ Given a term $f : t \rightarrow t'$, you get a data mapping $[f] : [t] \rightarrow [t']$.
 - ▶ All equations obeyed.
- ▶ $S\text{-inst}$ is further a topos (model of higher-order logic / set theory).
- ▶ We consider the following schema in the examples that follow:



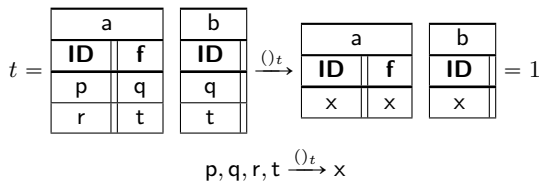
Programming AQL in CAML: Unit

- ▶ The unit instance 1 has one row per table:

a	
ID	f
x	x

b	
ID	
x	

- ▶ The data mapping $(\cdot)_t : t \rightarrow 1$ sends every row in t to the only row in 1. For example,



Programming AQL in CAML: Products

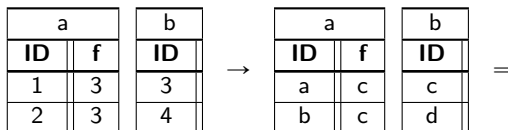
- Products $s \times t$ are computed row-by-row, with evident projections $\pi^1 : s \times t \rightarrow s$ and $\pi^2 : s \times t \rightarrow t$. For example:

a		b		×	a		b		=	a		b	
ID	f	ID			ID	f	ID			ID	f	ID	
1	3	3		a	c	c		(1,a)	(3,c)	(3,c)			
2	3	4		b	c	d		(1,b)	(3,c)	(3,d)			
								(2,a)	(3,c)	(4,c)			
								(2,b)	(3,c)	(4,d)			

- Given data mappings $f : s \rightarrow t$ and $g : s \rightarrow u$, how to define $(f, g) : s \rightarrow t \times u$ is left to the reader.
 - hint: try it on π^1 and π^2 and verify that $(\pi^1, \pi^2) = id$.

Programming AQL in CAML: Exponentials

- Exponentials t^s are given by finding all data mappings $s \rightarrow t$:



a	
ID	f
$1 \mapsto a, 2 \mapsto b, 3 \mapsto c, 4 \mapsto d$	
$1 \mapsto b, 2 \mapsto a, 3 \mapsto c, 4 \mapsto d$	
$1 \mapsto a, 2 \mapsto a, 3 \mapsto c, 4 \mapsto d$	
$1 \mapsto b, 2 \mapsto b, 3 \mapsto c, 4 \mapsto d$	
$1 \mapsto a, 2 \mapsto b, 3 \mapsto d, 4 \mapsto c$	
$1 \mapsto b, 2 \mapsto a, 3 \mapsto d, 4 \mapsto c$	
$1 \mapsto a, 2 \mapsto a, 3 \mapsto d, 4 \mapsto c$	
$1 \mapsto b, 2 \mapsto b, 3 \mapsto d, 4 \mapsto c$	

b	
ID	f
$3 \mapsto c, 4 \mapsto c$	
$3 \mapsto c, 4 \mapsto d$	
$3 \mapsto d, 4 \mapsto c$	
$3 \mapsto d, 4 \mapsto d$	

- Defining *eval* and λ are left to the reader.

Conclusion

- ▶ We described a new “algebraic” approach to databases based on category theory.
 - ▶ Schemas are categories, instances are set-valued functors.
 - ▶ Three adjoint data migration functors, Σ, Δ, Π manipulate data.
 - ▶ Instances on a schema model the simply-typed λ -calculus.
- ▶ Our approach is implemented in AQL, an open-source project, available at catinf.com.
- ▶ Collaborators welcome!
 - ▶ We are looking for a “real-world pilot project”.