

From Graphs to Categories

And Algebraic Property Graphs Too

Ryan Wisnesky
Conexus AI

November 10, 2020

$$\Sigma \dashv \Delta \dashv \Pi$$

Introduction

- ▶ These slides describe new ways to formalize databases and migrate data based on category theory.
- ▶ Category theory was designed to migrate theorems from one area of mathematics to another, so it is a very natural language with which to describe migrating data / knowledge from one schema / ontology to another.
- ▶ Research has culminated in an open-source language, CQL, available at categoricaldata.net, being commercialized by Conexus AI, conexus.com.
- ▶ Categories are graphs with extra structure, and so category theory has deep connections to property graphs (joint work with Joshua Shinavier), also described in these slides.

Category Theory

A category \mathcal{C} consists of

- ▶ objects $A, B, C \dots$ and arrows (also called *morphisms*) $f, g, h \dots$ such that:
- ▶ For every arrow f there is an object $\text{src}(f)$ called the *source* of f and an object $\text{tgt}(f)$ called the *target* of f . When $S = \text{src}(f)$ and $T = \text{tgt}(f)$, we may write $f : S \rightarrow T$. Visually:

$$S \xrightarrow{f} T$$

- ▶ For every arrow $f : A \rightarrow B$ and arrow $g : B \rightarrow C$ there is an arrow $g \circ f : A \rightarrow C$ called the *composite* of f and g :

$$\begin{array}{ccccc} A & \xrightarrow{f} & B & \xrightarrow{g} & C \\ & \searrow & & \nearrow & \\ & & g \circ f & & \end{array}$$

- ▶ Composition is *associative*, i.e. $h \circ (g \circ f) = (h \circ g) \circ f$ for arbitrary f, g , and h .
- ▶ For every object A there is an *identity arrow* $\text{id}_A : A \rightarrow A$:

$$\begin{array}{c} \text{id}_A \\ \curvearrowright \\ A \end{array}$$

- ▶ Furthermore, for any arrow $f : A \rightarrow B$, $f \circ \text{id}_A = f = \text{id}_B \circ f$.

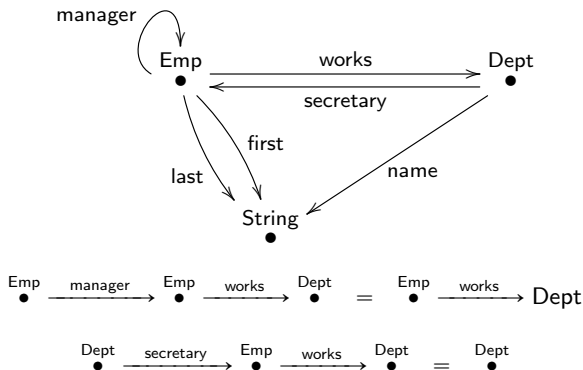
A *functor* $F : \mathcal{C} \rightarrow \mathcal{D}$ is a function from \mathcal{C} 's objects to \mathcal{D} 's objects and \mathcal{C} 's arrows to \mathcal{D} 's arrows that preserves composition and identity:

$$F(\text{id}_C) = \text{id}_{F(C)} \quad F(f \circ g) = F(f) \circ F(g).$$

Categorical Schemas and Databases

- ▶ A *schema* S is a directed multi-graph and a set of paths through the graph called “equivalent”.
- ▶ A schema S denotes a category $\llbracket S \rrbracket$:
 - ▶ The objects of $\llbracket S \rrbracket$ are the nodes of S .
 - ▶ The arrows of $\llbracket S \rrbracket$ are the paths through S , modulo the path equivalences in S .
- ▶ An S -instance (database on schema S) is a collection of sets, one per node in S , and a collection of (unary) functions, one per edge in S , satisfying the path equivalences in S .
- ▶ For example, these sets and functions may be represented as a collection of SQL tables, one per node in S , each with columns for edges out of that node.
- ▶ An S -instance denotes a functor $\llbracket S \rrbracket \rightarrow \mathbf{Set}$, where \mathbf{Set} , the category of sets, has for objects all sets and for arrows all (unary) functions.

Example Categorical Schema and Database

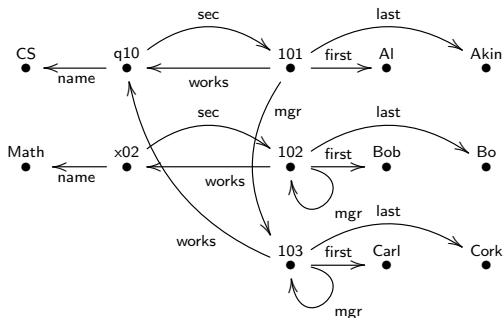


Emp				
ID	mgr	works	first	last
101	103	q10	Al	Akin
102	102	x02	Bob	Bo
103	103	q10	Carl	Cork

Dept		
ID	sec	name
q10	101	CS
x02	102	Math

String	
ID	
Al	
Bob	
...	

Categorical Databases to Triples (Graphs) and Back

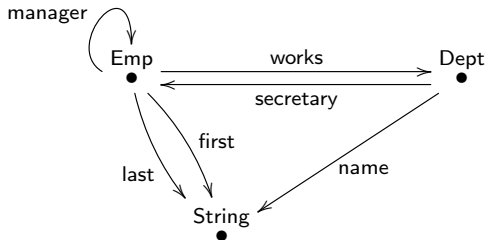


Emp				
ID	mgr	works	first	last
101	103	q10	Al	Akin
102	102	x02	Bob	Bo
103	103	q10	Carl	Cork

Dept		
ID	sec	name
q10	101	CS
x02	102	Math

String	
ID	
Al	
Bob	
...	

Categorical Select-From-Where/For-Where-Return Syntax



Find the name of every manager's department:

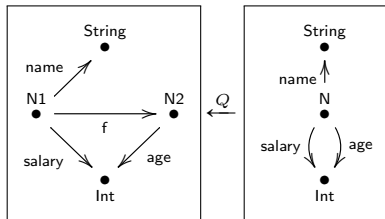
CQL

```
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
```

Query Evaluation and Co-evaluation



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

$\xleftarrow{eval_Q}$
 $\xrightarrow{coeval_Q}$

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

Example Round Trip (column f removed)

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

N2	
ID	Age
4	20
5	20
6	30

$\xrightarrow{\text{coeval}_Q}$

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

$\nwarrow \text{eval}_Q$

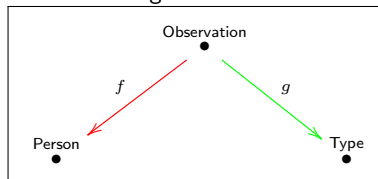
$\downarrow \eta$

N1		
ID	Name	Salary
a	Alice	\$100
b	Bob	\$250
c	Sue	\$300
d	null_4	null_5
e	null_6	null_7
f	null_8	null_9

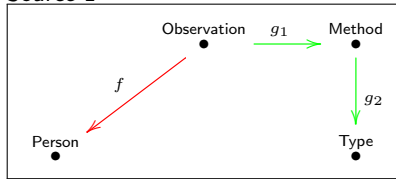
N2	
ID	Age
a	null_1
b	null_2
c	null_3
d	20
e	20
f	30

Schema Integration

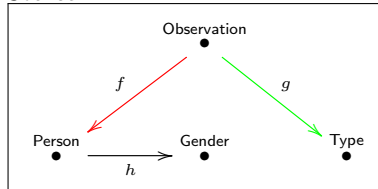
Column-Linkages



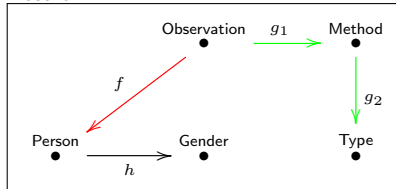
Source-1



Source-2



Result



Data Integration

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

→

Method			Type
ID	g2		ID
<i>m</i> ₁	BP		BP
<i>m</i> ₂	BP		Wt
<i>m</i> ₃	Wt		
<i>m</i> ₄	Wt		

Observation			Person
ID	f	g1	ID
<i>o</i> ₁	Pete	<i>m</i> ₁	Jane
<i>o</i> ₂	Pete	<i>m</i> ₂	<i>Pete</i>
<i>o</i> ₃	Jane	<i>m</i> ₃	
<i>o</i> ₄	Jane	<i>m</i> ₁	

↓

Gender		Type
ID		ID
F		BP
M		Wt
		HR

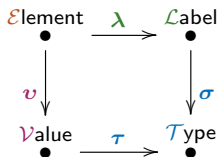
Observation			Person	
ID	f	g	ID	h
<i>o</i> ₅	Peter	BP	Paul	M
<i>o</i> ₆	Paul	HR	Peter	M
<i>o</i> ₇	Peter	Wt		

→

Method		Observation		
ID	g2	ID	f	g1
<i>null</i> ₁	BP	<i>o</i> ₁	Peter	<i>m</i> ₁
<i>null</i> ₂	Wt	<i>o</i> ₂	Peter	<i>m</i> ₂
<i>null</i> ₃	HR	<i>o</i> ₃	Jane	<i>m</i> ₃
<i>m</i> ₁	BP	<i>o</i> ₄	Jane	<i>m</i> ₁
<i>m</i> ₂	BP	<i>o</i> ₅	Peter	<i>null</i> ₁
<i>m</i> ₃	Wt	<i>o</i> ₆	Paul	<i>null</i> ₂
<i>m</i> ₄	Wt	<i>o</i> ₇	Peter	<i>null</i> ₃

Gender		Type	Person	
ID		ID	ID	h
F		BP	Jane	<i>null</i> ₄
M		Wt	Paul	M
<i>null</i> ₄		HR	<i>Peter</i>	M

An Overly-General Schema for Algebraic Property Graphs



$$\mathcal{E} \xrightarrow{v} \mathcal{V} \xrightarrow{\tau} \mathcal{T} = \mathcal{E} \xrightarrow{\lambda} \mathcal{L} \xrightarrow{\sigma} \mathcal{T}$$

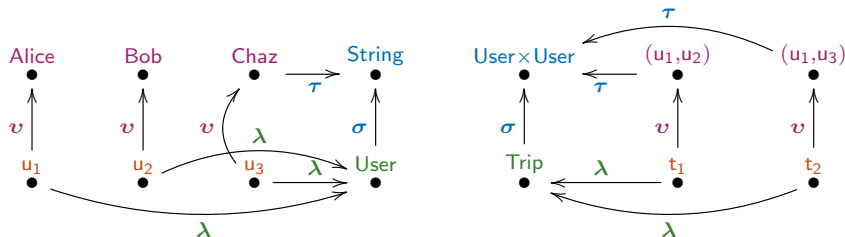
Element		
ID	λ	v
t_1	Trip	(u_1, u_2)
t_2	Trip	(u_1, u_3)
u_1	User	Alice
u_2	User	Bob
u_3	User	Chaz

Value	
ID	τ
Alice	String
Bob	String
Chaz	String
(u_1, u_2)	User \times User
(u_1, u_3)	User \times User

Label	
ID	σ
User	String
Trip	User \times User

Type	
ID	
String	
User \times User	

Algebraic Property Graphs as 4-sorted Triples



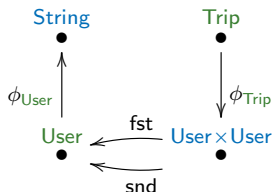
Element		
ID	λ	v
t_1	Trip	(u_1, u_2)
t_2	Trip	(u_1, u_3)
u_1	User	Alice
u_2	User	Bob
u_3	User	Chaz

Value	
ID	τ
Alice	String
Bob	String
Chaz	String
(u_1, u_2)	User \times User
(u_1, u_3)	User \times User

Label	
ID	σ
User	String
Trip	User \times User

\mathcal{T} type	
ID	
String	
User \times User	

Algebraic Property Graphs with Product Schemas



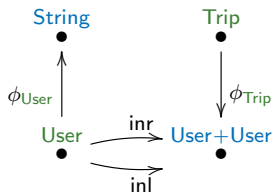
User × User		
ID	fst	snd
(u_1, u_1)	u_1	u_1
(u_1, u_2)	u_1	u_2
(u_1, u_3)	u_1	u_3
(u_2, u_1)	u_2	u_1
(u_2, u_2)	u_2	u_2
(u_2, u_3)	u_2	u_3
(u_3, u_1)	u_3	u_1
(u_3, u_2)	u_3	u_2
(u_3, u_3)	u_3	u_3

User	
ID	ϕ_{User}
u_1	Alice
u_2	Bob
u_3	Chaz

Trip	
ID	ϕ_{Trip}
t_1	(u_1, u_2)
t_2	(u_1, u_3)

String
ID
Alice
Bob
Chaz

Algebraic Property Graphs with Sum Schemas



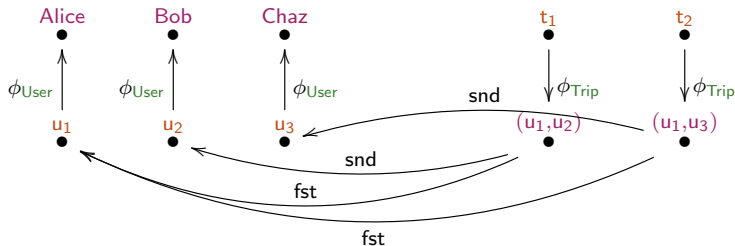
User + User	
ID	
	$\text{inl}(u_1)$
	$\text{inl}(u_2)$
	$\text{inl}(u_3)$
	$\text{inr}(u_1)$
	$\text{inr}(u_2)$
	$\text{inr}(u_3)$

User			
ID	ϕ_{User}	inl	inr
u_1	Alice	$\text{inl}(u_1)$	$\text{inr}(u_1)$
u_2	Bob	$\text{inl}(u_2)$	$\text{inr}(u_2)$
u_3	Chaz	$\text{inl}(u_3)$	$\text{inr}(u_3)$

Trip	
ID	ϕ_{Trip}
t_1	$\text{inl}(u_1)$
t_2	$\text{inr}(u_2)$

String	
ID	
	Alice
	Bob
	Chaz

Algebraic Property Graphs as Typed-sorted Triples



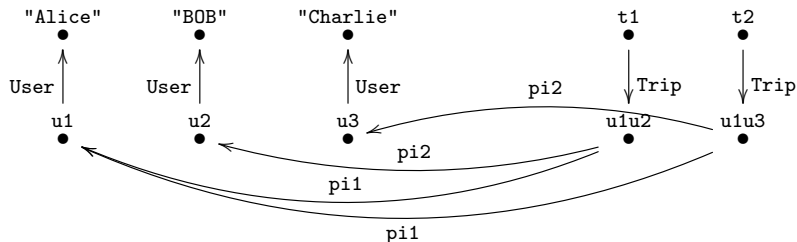
User \times User		
ID	fst	snd
(u_1, u_2)	u_1	u_2
(u_1, u_3)	u_1	u_3

User	
ID	ϕ_{User}
u_1	Alice
u_2	Bob
u_3	Chaz

Trip	
ID	ϕ_{Trip}
t_1	(u_1, u_2)
t_2	(u_1, u_3)

String
ID
Alice
Bob
Chaz

Algebraic Property Graphs as RDF



User \times User			
ID	fst	snd	Resource
(u_1, u_2)	u_1	u_2	u1u2
(u_1, u_3)	u_1	u_3	u1u3

User		
ID	ϕ_{User}	Resource
u_1	Alice	u1
u_2	Bob	u2
u_3	Charlie	u3

Trip		
ID	ϕ_{Trip}	Resource
t_1	(u_1, u_2)	t1
t_2	(u_1, u_3)	t2

String	
ID	Resource
Alice	"Alice"
Bob	"BOB"
Chaz	"Charlie"