# Composing knowledge graphs, inside and out

Spencer Breiner[1]

`spencer.breiner@nist.gov`

Joint with Blake Pollard[1], Peter Denno[1]

and Eswaran Subrahmanian[1,2]

[1]NIST        [2]CMU

March 18, 2020

**NIST**
National Institute of
Standards and Technology
U.S. Department of Commerce

**Carnegie
Mellon
University**

# About me (Spencer Breiner)

National Institute of Standards and Technology

- Information Technology Lab - Software & Systems Division
- Ph.D., CMU, 2013 - Category theory (CT) and logic

Current work: *Applied* CT for systems modeling

- Knowledge representation
- Knowledge integration
- Multiple semantics

Outline for today:

- Graphs & categories
- Why not (just) graphs?
- Knowledge graphs as categories and functors

# What's beneath a knowledge graph?

Knowledge Graphs:

> *"structured representations of semantic*
>
> *knowledge that are stored in a graph"*

What structure? Stored how?

Today, some possible answers from category theory.

Some themes:

Bite-size ontologies
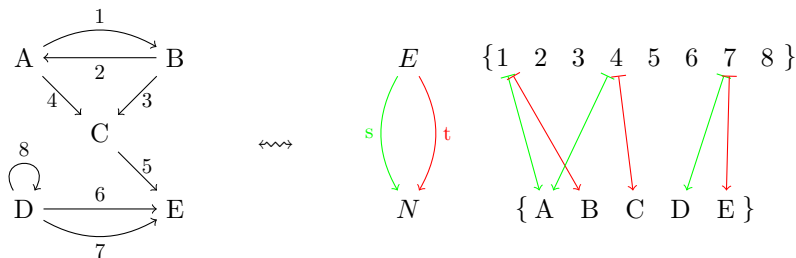
Data/concept duality

Internalized computation/proof

# Graphs

For today, graphs are *directed* and (optionally) *multi-*.

Any graph can be represented as

- A pair of sets $N = \mathtt{Node}$ and $E = \mathtt{Edge}$, and

- A pair of functions $s = \mathtt{src}, t = \mathtt{tgt} : E \rightrightarrows N$.

For example,

# Categories

A category is a graph $G$ together with

- Version 1: A partial associative operation (with identities)

$$E \times E$$
$$\cup|$$
$$\{f.\mathtt{tgt} = g.\mathtt{src}\} \xrightarrow{\quad f.g \quad} E$$

  Semantic categories: **Sets**, **Graph**, **Vect**, **Type**

- Version 2: A (concat-stable) equivalence relation over paths

$$\{\langle f_i \rangle \sim \langle g_j \rangle\} \subseteq \mathbf{Path}(G) \times \mathbf{Path}(G)$$

  Schemas: $\mathcal{G} = \langle E \rightrightarrows N \rangle$, $\mathcal{P} = \mathbf{OSProb}$, $\mathcal{S} = \mathbf{OSSoln}$

# Free categories (!)

Upshot: Any graph $G$ already "is" a category.

The relationship is mediated by a free/forgetful *adjunction*



$$\frac{\mathbf{Free}(G) \longrightarrow \mathbf{C} \quad \in \mathbf{Cat}}{G \longrightarrow \mathbf{Forget}(\mathbf{C}) \quad \in \mathbf{Graph}}$$

Two round trips:

A *monad* $\eta_G : G \to \mathbf{Path}(G)$ (concat)

A *comonad* $\epsilon_{\mathbf{C}} : \mathbf{Factor}(\mathbf{C}) \to \mathbf{C}$ (compute)

# A bite-sized example

Open-shop scheduling

<table>
<tr><th colspan="6" align="center">Problem</th></tr>
<tr><th></th><th>Jobs</th><th>$j_1$</th><th>$j_2$</th><th>$j_3$</th><th>$j_4$</th></tr>
<tr><td rowspan="4">Machines</td><td>saw</td><td>2 hr</td><td>2 hr</td><td>2 hr</td><td>1 hr</td></tr>
<tr><td>drill</td><td>2 hr</td><td>3 hr</td><td>0</td><td>3 hr</td></tr>
<tr><td>lathe</td><td>2 hr</td><td>3 hr</td><td>3 hr</td><td>0</td></tr>
<tr><td>mill</td><td>2 hr</td><td>2 hr</td><td>1 hr</td><td>3 hr</td></tr>
</table>

<table>
<tr><th colspan="11" align="center">Schedule</th></tr>
<tr><th></th><th>1</th><th>2</th><th>3</th><th>4</th><th>5</th><th>6</th><th>7</th><th>8</th><th>9</th><th>10</th></tr>
<tr><td>saw</td><td colspan="2">$j_1$</td><td></td><td colspan="2">$j_4$</td><td colspan="2">$j_2$</td><td></td><td colspan="2">$j_3$</td></tr>
<tr><td>drill</td><td colspan="2">$j_2$</td><td colspan="2">$j_1$</td><td></td><td></td><td colspan="3">$j_4$</td><td></td></tr>
<tr><td>lathe</td><td colspan="2">$j_3$</td><td></td><td></td><td></td><td colspan="2">$j_1$</td><td colspan="3">$j_2$</td></tr>
<tr><td>mill</td><td colspan="2">$j_4$</td><td colspan="2">$j_2$</td><td colspan="2">$j_3$</td><td></td><td colspan="2">$j_1$</td><td></td></tr>
</table>

Schematically,

$$\mathcal{P} = \langle \tau : J \times M \longrightarrow \mathbb{R}^+ \rangle \qquad \mathcal{S} = \langle s, t : J \times M \rightrightarrows \mathbb{R}^+ \mid \text{ax.} \rangle$$

The two are related via a functor

$$F : \mathcal{P} \longrightarrow \mathcal{S}$$
$$\tau \longmapsto t - s$$

## Functorial semantics & duality

$F : \mathcal{P} \to \mathcal{S}$ encodes: "every schedule solves *some* problem."

Any concrete schedule (instance) defines a functor $P : \mathcal{P} \to \mathbf{Sets}$

Nodes map to sets: $\qquad P(J) = \{j_1, j_2, j_3, j_4\}$

Edges map to functions: $\quad P(\tau) : (j_2, \mathtt{lathe}) \mapsto 3 \text{ hr}$

Every schema functor defines a *dual* transformation on instances

$$\mathbf{Inst}(\mathcal{S}) \xrightarrow{\;\;F^*\;\;} \mathbf{Inst}(\mathcal{P})$$
$$S(s),\; S(t) \longmapsto S(t) - S(s)$$

Duality is just precomposition:

$$
\begin{array}{ccc}
\mathcal{S} & \xleftarrow{\;\;F\;\;} & \mathcal{P} \\
{\scriptstyle S}\downarrow & & \downarrow{\scriptstyle F^*(S) = F.S} \\
\mathbf{Sets} & =\!=\!=\!=\!= & \mathbf{Sets}
\end{array}
$$

# Why not (just) graphs?

- Structured nodes/edges: $J \times M$

- Built-in elements (libraries): $\mathtt{diff} : \mathbb{R}^+ \times \mathbb{R}^+ \longrightarrow \mathbb{R}^+$

- First-class axioms/proofs: $s_{jm} \leq t_{jm} \ \vdash \ F(\tau)_{jm} \in \mathbb{R}^+$

- Not a graph homomorphism: $F(\tau) = \ell.p$

# Structure in a category

The *Cartesian product* of two objects $X$ and $Y$ is a diagram $X \xleftarrow{\pi_1} P \xrightarrow{\pi_2} Y$ such that, for any object $Z$ and any pair of arrows $x : Z \to X$ and $y : Z \to Y$, there exists a unique map $p = \langle x, y \rangle$ such that $p.\pi_1 = x$ and $p.\pi_2 = y$.

## Generalized elements

A suggestive notation:

$$x : Z \to X \iff x \underset{Z}{\in} X$$

Compare:

| **In set theory** | **In category theory** |
|---|---|
| $p \in X \times Y$ | $p \underset{Z}{\in} X \times Y$ |
| $x \in X,\ y \in Y,\ p = \langle x, y \rangle$ | $x \underset{Z}{\in} X,\ y \underset{Z}{\in} Y,\ \underbrace{p = \langle x, y \rangle}_{p.\pi_1 = x,\ p.\pi_2 = y}$ |

Why generalize? In **Sets**, arrows $\{*\} \to X$ "see" everything in $X$, but...

  In **Graph**, $\{*\}$ can't distinguish $\{*\quad *\}$ from $\{* \to *\}$.

  In **Vect**, $\{*\} = \mathbb{R}^0$ can't see *anything* (zero object).

# More structure

In programming, a function `f(x : X) : Y` is *pure* if

- It has no side effects (e.g., no IO, non-local variable mutation)

- It has consistent return values (e.g., no non-local variable dependence)

The pure fragment of a programming language defines a category **Type**.

The *exponential* adjunction mediates global/generalized elements

$$
\begin{array}{cc}
\textbf{Type} & \qquad f: \quad X \times Z \longrightarrow Y \\
{\scriptstyle -\times Z}\Big\uparrow\Big\downarrow{\scriptstyle (-)^Z} & \overline{\qquad\qquad\qquad\qquad\qquad} \\
& \ulcorner f(x,-)\urcorner: \quad X \longrightarrow Y^Z \\
\textbf{Type} & \overline{\qquad\qquad\qquad\qquad\qquad} \\
& \ulcorner f \urcorner: \quad \{*\} \longrightarrow (Y^Z)^X
\end{array}
$$

Round trips:   `eval`:   $Y^Z \times Z \longrightarrow Y$

　　　　　　`coeval`:   $X \longrightarrow (X \times Z)^Z$

## Types in a schema

We can think of schema libraries as

- A subschema $\mathcal{P}_0 \subseteq \mathcal{P}$, together with
- A fixed implementation functor $\mathbf{impl} : \mathcal{P}_0 \to \mathbf{Type}$

An instance should respect the behavior of the implementation:

$$
\begin{array}{ccc}
\mathcal{P} & \xrightarrow{\quad P \quad} & \mathbf{Sets} \\
\cup| & & \big\uparrow \mathbf{glob}=\mathrm{Hom}(\{*\},-) \\
\mathcal{P}_0 & \xrightarrow[\quad \mathbf{impl} \quad]{} & \mathbf{Type}
\end{array}
$$

Problem:    We want $\tau \in \mathbb{R}^+$, but    $s, t \in \mathbb{R}^+ \not\Rightarrow t - s \in \mathbb{R}^+$.

## Logic in a schema

In general, formulas define subobjects, and inferences define sub-sub-objects:

$$\frac{\varphi(x) \vdash \psi(x)}{}$$

$$[\![\varphi]\!] \dashrightarrow [\![\psi]\!]$$
$$X$$

Interpretations are defined recursively:

| $x = y$ | $\varphi \wedge \psi$ | $\varphi \vee \psi$ | $\exists y.\varphi$ |
|---------|----------------------|---------------------|---------------------|
| Diagonal | Pullback | Pushout | Image |

$X$
$\langle \text{id},\text{id} \rangle \downarrow$
$X \times X$

$[\![\varphi \wedge \psi]\!] \longrightarrow [\![\psi]\!]$
$\downarrow \qquad \downarrow$
$[\![\varphi]\!] \longrightarrow X$

$[\![\varphi \wedge \psi]\!] \longrightarrow [\![\psi]\!]$
$\downarrow \qquad \downarrow$
$[\![\varphi]\!] \longrightarrow [\![\varphi \vee \psi]\!]$

$[\![\varphi]\!] \longrightarrow [\![\exists y.\varphi]\!]$
$\downarrow \qquad \downarrow$
$X \times Y \longrightarrow X$

## Proofs as diagrams

Formulas are (sub)objects, inferences & proofs are arrows:

An axiom:

$$\vdash s_{jm} \leq t_{jm}$$



An inference:

$$x \leq y \vdash (y - x) \in \mathbb{R}^+$$



The cut rule corresponds to concatenation of diagrams

e.g., $\vdash F(\tau)_{jm} \in \mathbb{R}^+$

## Functors between graphs

Functors are more flexible than graph homomorphisms:

Nodes map to nodes, but edges map to *paths*.

$$
\begin{array}{ccc}
\mathcal{P} & & J \times M \xrightarrow{\quad\tau\quad} \mathbb{R}^+ \\
{\scriptstyle F}\downarrow & & \quad\downarrow\qquad\qquad\qquad\downarrow \\
\mathcal{S} & & J \times M \xrightarrow[\ell]{} \bullet \xrightarrow[p]{} \mathbb{R}^+
\end{array}
$$

Usually interested in *structure-preserving* functors (instances, too!)

Constructions: $\quad F(J \times M) \cong F(J) \times F(M)$

Types: $\quad \mathcal{P}_0 \dashrightarrow^{F(\mathcal{P}_0) \subseteq \mathcal{S}_0} \mathcal{S}_0$

$$
\mathcal{P}_0 \searrow_{\textbf{impl}} \quad \swarrow_{\textbf{impl}} \mathcal{S}_0
$$

$$\textbf{Type}$$

## Solutions as functors

Any solution algorithm $a$ defines a matrix endomorphism

$$(\mathbb{R}^+)^{J \times M} \xrightarrow{\quad a \quad} (\mathbb{R}^+)^{J \times M}$$

$$(\tau_{jm}) \longmapsto (s_{jm})$$

From this, we can define a functor $A : \mathcal{S} \to \mathcal{P}$

$$A(s)_{jm} = \texttt{eval}\big(a(\ulcorner \tau \urcorner), (j, m)\big)$$
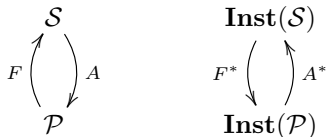
$$A(t)_{jm} = A(s)_{jm} + \tau_{jm}$$

$$J \times M \xrightarrow{\langle \mathrm{id}, \ulcorner \tau \urcorner \rangle} (J \times M) \times (\mathbb{R}^+)^{J \times M} \xrightarrow{\mathrm{id} \times a} (J \times M) \times (\mathbb{R}^+)^{J \times M} \xrightarrow{\texttt{eval}} \mathbb{R}^+$$

Defining $A$ requires proof: $a$ satisfies the axioms of $\mathcal{S}$.

Note: Equivalence $(\mathbb{R}^+)^{J \times M} \cong \mathrm{Mat}_{\mathbb{R}^+}(|J|, |M|)$ requires a *labeling*.

# A knowledge "graph"

By duality, every problem $P \in \mathbf{Inst}(\mathcal{P})$ defines a solution $A^*(P) \in \mathbf{Inst}(\mathcal{S})$.

$$
\begin{array}{cc}
\mathcal{S} & \mathbf{Inst}(\mathcal{S}) \\
F \big( \quad \big) A & F^* \big( \quad \big) A^* \\
\mathcal{P} & \mathbf{Inst}(\mathcal{P})
\end{array}
$$

The functors should satisfy $F.A = \mathrm{id}_\mathcal{P}$ $(\Rightarrow A^*.F^* = \mathrm{id}_{\mathbf{Inst}(\mathcal{P})})$:

$$
\begin{aligned}
A(F(\tau)) &= A(t - s) \\
&= A(t) - A(s) \\
&= A(s) + \tau - A(s) \\
&= \tau
\end{aligned}
$$

# Variation I

What's the difference?

$$\mathcal{P}' := \left\langle \begin{array}{c} T \xrightarrow{\tau'} \mathbb{R}^+ \\ {}^{j}\swarrow \quad {}^{m}\searrow \\ J \qquad\qquad M \end{array} \right\rangle \qquad \mathcal{S}' := \left\langle \begin{array}{c} T \underset{t'}{\overset{s'}{\rightrightarrows}} \mathbb{R}^+ \\ {}^{j}\swarrow \quad {}^{m}\searrow \\ J \qquad\qquad M \end{array} \right\rangle$$

What's the same?

Problem generalization, functorially:

$$
\begin{array}{ccc}
\mathcal{S}' & \dashrightarrow^{\exists! \ \overline{K}} & \mathcal{S} \\
{\scriptstyle F'}\big\uparrow & & \big\uparrow{\scriptstyle F} \\
\mathcal{P}' & \xrightarrow[K:T \mapsto J \times M]{} & \mathcal{P}
\end{array}
$$

The other direction(s)?

# Variation II

Duplicate machines ($C=$"capability", $a=$"assignment")

$$\mathcal{P}^d := \left\langle \begin{array}{c} J \times C \xrightarrow{\tau^c} \mathbb{R}^+ \\ \\ M \xrightarrow{\quad c \quad} C \end{array} \right\rangle \qquad \mathcal{S}^d := \left\langle \begin{array}{c} M \xleftarrow{a} \\ c\downarrow \quad J \times C \begin{array}{c} \xrightarrow{s^c} \\ \xrightarrow{t^c} \end{array} \mathbb{R}^+ \\ C \xleftarrow{p_2} \end{array} \right\rangle$$

The arrow (bundle, dep. type) $M \to C$ represents a family of sets $\{M_c\}_{c \in C}$.

This time, we can go both ways (sort of)

$$\begin{array}{ccccc} \mathcal{S}^d & \xrightarrow{\exists! \ \overline{I}:a \mapsto p_2} & \mathcal{S} & \xrightarrow{H:s,t \mapsto s|^a,t|^a} & \mathcal{S}^d \\ \Big\uparrow{F^d} & & \Big\uparrow{F} & \times & \Big\uparrow{F^d} \\ \mathcal{P}^d & \xrightarrow{I:c \mapsto \mathrm{id}_M} & \mathcal{P} & \xrightarrow{G:M \mapsto C} & \mathcal{P}^d \end{array}$$

Here $H : M \mapsto M$ and $s|^a, t|^a$ denote extension by zero along $a$.

## Variation III

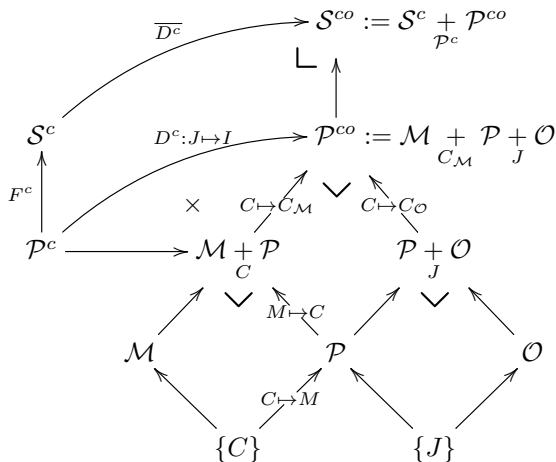Duplicate jobs ($\mathcal{P}$=Process *catalog*, $\mathcal{O}$=Orders database)

$$\mathcal{O} := \left\langle \begin{array}{c} C \xleftarrow{\ c\ } O \xleftarrow{\ o\ } I \\ {\scriptstyle b}\Big\downarrow\Big\downarrow{\scriptstyle r} \quad \swarrow{\scriptstyle d} \qquad \Big\downarrow{\scriptstyle j} \\ \mathbb{R}^+ \xleftarrow[\ k\ ]{} J \end{array} \right\rangle, \qquad b = r + \sum_{o:O_c} \sum_{i:I_o} k(j(i))$$

Extract the daily schedule by mapping to a pushout:

$$\mathcal{P} \longrightarrow \mathcal{O} + \mathcal{P}$$

$$
\begin{array}{ccc}
J \times M & \longmapsto & I \times M \\
& & \Big\downarrow{\scriptstyle j \times \mathrm{id}} \\
{\scriptstyle \tau}\Big\downarrow & & J \times M \\
& & \Big\downarrow{\scriptstyle \tau} \\
\mathbb{R}^+ & \longmapsto & \mathbb{R}^+
\end{array}
$$

# Variation IV

Duplicate jobs and machines ($\mathcal{M}$=Shop floor model)

# Wrapping Up

Recap:

- Bite-sized semantic models & functorial instances

- Built-in logic & computation via (preservation of) structure

- Knowledge graphs as schemas & functors.

More goodies:

- Build-your-own semantics (presheaves)

- Internal concepts generate external schemas (Yoneda/slice cat.)

- Relationships between relationships (Natural transformations)

The bad news...

- Limited tooling

- Steep learning curve

# Thank you!

PS. This talk is based on a paper under review, but a draft is available on request from `spencer.breiner@nist.gov`.