

---

# COMPUTATIONAL CATEGORY-THEORETIC REWRITING

---

✉ **Kristopher Brown**

Topos Institute  
kris@topos.institute

✉ **Evan Patterson**

Topos Institute  
evan@topos.institute

**Tyler Hanks**

University of Florida  
thanks@ufl.edu

✉ **James Fairbanks**

Department of Computer Science  
University of Florida  
fairbanksj@ufl.edu

## ABSTRACT

We demonstrate how category theory provides specifications that can efficiently be implemented via imperative algorithms and apply this to the field of graph rewriting. By examples, we show how this paradigm of software development makes it easy to quickly write correct and performant code. We provide a modern implementation of graph rewriting techniques at the level of abstraction of finitely-presented  $\mathcal{C}$ -sets and clarify the connections between  $\mathcal{C}$ -sets and the typed graphs supported in existing rewriting software. We emphasize that our open-source library is extensible: by taking new categorical constructions (such as slice categories, structured cospans, and distributed graphs) and relating their limits and colimits to those of their underlying categories, users inherit efficient algorithms for pushout complements and (final) pullback complements. This allows one to perform double-, single-, and sesqui-pushout rewriting over a broad class of data structures.

**Keywords** Double pushout rewriting · category theory · graph rewriting

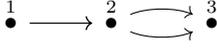
## 1 Introduction and motivation

Term rewriting is a foundational technique in computer algebra systems, programming language theory, and symbolic approaches to artificial intelligence. While classical term rewriting is concerned with tree-shaped terms in a logical theory, the field of graph rewriting extends these techniques to more general shapes of terms, typically simple graphs, digraphs, multigraphs, or typed graphs. Major areas of graph rewriting are graph *languages* (rewriting defines a graph grammar), graph *relations* (rewriting is a relation between input and output graphs), and graph *transition systems* (rewriting evolves a system in time) [15].

When considering the development of software for graph rewriting, it is important to distinguish between studying rewriting systems as mathematical objects and building applications on top of rewriting as infrastructure. The former topic can answer inquiries into confluence, termination, reachability, and whether certain invariants are preserved by rewriting systems. In contrast, we will focus on answering questions that involve the application of concretely specified rewrite systems to particular data.

Category theory is a powerful tool for developing rewriting software, as the numerous and heterogeneous applications and techniques of rewriting are elegantly unified by categorical concepts. Furthermore, the semantics of categorical treatments of graph rewriting are captured by universal properties of limits and colimits, which are easier to reason about than operational characterizations of rewriting. This is an instance of a broader paradigm of *computational applied category theory*, which begins by modeling the domain of interest with category theory, such as using monoidal categories and string diagrams to model processes. One is then free (but not required) to implement the needed categorical structures in a conventional programming language, where the lack of a restrictive type system facilitates a fast software development cycle and enables algorithmic efficiency. For example, arrays can be used to represent finite sets, and union-find data structures can compute equivalence classes.





**Figure 2:** A graph  $G$ , defined by  $G_V = [3]$ ,  $G_E = [3]$ ,  $G_{\text{src}} = [1, 2, 2]$ , and  $G_{\text{tgt}} = [2, 3, 3]$ .

Given two graphs  $G$  and  $H$ , a *graph homomorphism*  $G \xrightarrow{h} H$  consists of a mapping of edges,  $G_E \xrightarrow{h_E} H_E$  and a mapping of vertices,  $G_V \xrightarrow{h_V} H_V$ , that preserve the graph structure, i.e., the following diagrams commute:

$$\begin{array}{ccc}
 G_E & \xrightarrow{G_{\text{src}}} & G_V \\
 h_E \downarrow & & \downarrow h_V \\
 H_E & \xrightarrow{H_{\text{src}}} & H_V
 \end{array}
 \quad
 \begin{array}{ccc}
 G_E & \xrightarrow{G_{\text{tgt}}} & G_V \\
 h_E \downarrow & & \downarrow h_V \\
 H_E & \xrightarrow{H_{\text{tgt}}} & H_V
 \end{array}
 \quad (1)$$

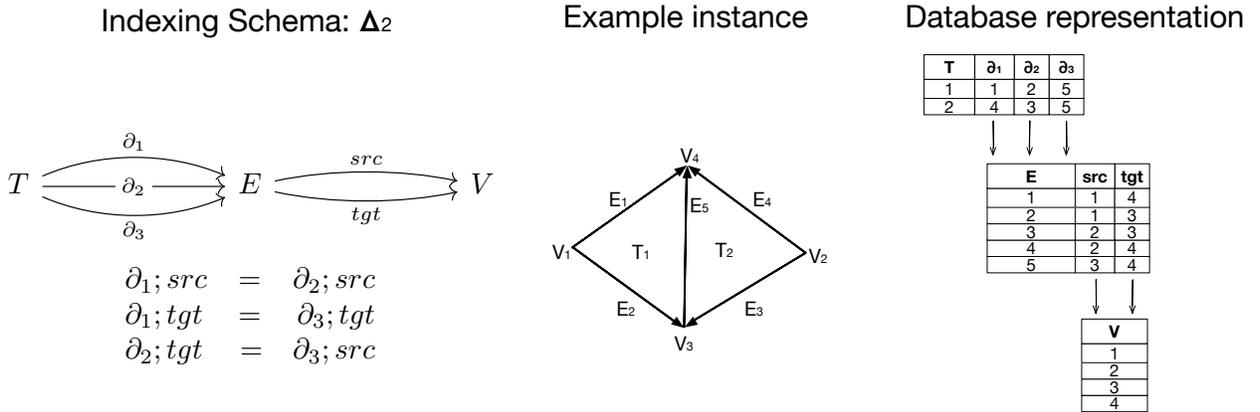
Regarding the source graph as a pattern, the homomorphism describes a pattern match in the target. A graph homomorphism can also be thought of as a typed graph, in which the vertices and edges of  $G$  are assigned types from  $H$ . For a fixed typing graph  $X$ , typed graphs and type-preserving graph homomorphisms form a category, namely the slice category  $\mathbf{Grph}/X$  [10].

## 2.2 $\mathcal{C}$ -sets and their homomorphisms

Graphs are a special case of a class of structures called  $\mathcal{C}$ -sets.<sup>1</sup> Consider the category  $\mathcal{C}$  freely generated by the graph  $E \begin{smallmatrix} \xrightarrow{s} \\ \xrightarrow{t} \end{smallmatrix} V$ . A  $\mathcal{C}$ -set is a functor from the category  $\mathcal{C}$  to  $\mathbf{Set}$ , which by definition assigns to each object a set and to each arrow a function from the domain set to the codomain set. For this choice of  $\mathcal{C}$ , the category of  $\mathcal{C}$ -sets is isomorphic to the category of directed multigraphs. Importantly, we recover the definition of graph homomorphisms between graphs  $G$  and  $H$  as a natural transformation of functors  $G$  and  $H$ .

The category  $\mathcal{C}$  is called the *indexing category* or *schema*, and the functor category  $[\mathcal{C}, \mathbf{Set}]$  is referred to as  $\mathcal{C}\text{-Set}$  or the category of *instances*, *models*, or *databases*. Given a  $\mathcal{C}$ -set  $X$ , the set that  $X$  sends a component  $c \in \text{Ob } \mathcal{C}$  to is denoted by  $X_c$ . Likewise, the finite function  $X$  sends a morphism  $f \in \text{Hom}_{\mathcal{C}}(a, b)$  to is denoted by  $X_f$ . We often restrict to  $[\mathcal{C}, \mathbf{FinSet}]$  for computations.

In addition to graphs,  $\mathbf{Set}$  itself can be thought of as  $\mathcal{C}\text{-Set}_{\text{ew}}$ , where the schema  $\mathcal{C}$  is the terminal category  $\mathbf{1}$ . We can change  $\mathcal{C}$  in other ways to obtain new data structures, as illustrated in Figure 3.  $\mathcal{C}$ -sets can also be extended with a notion of *attributes* to incorporate non-combinatorial data [32, 24], such as symbolic labels or real-valued weights. For simplicity of presentation, we focus on  $\mathcal{C}$ -sets without attributes in our examples.



**Figure 3:** The schema of two-dimensional semi-simplicial sets,  $\Delta_2$ , and an example semi-simplicial set, i.e. an object of  $\Delta_2\text{-Set}$ . The equations enforce the connectivity of edges to be a triangle. Note that MacLane defines  $\Delta$  as our  $\Delta^{op}$ .

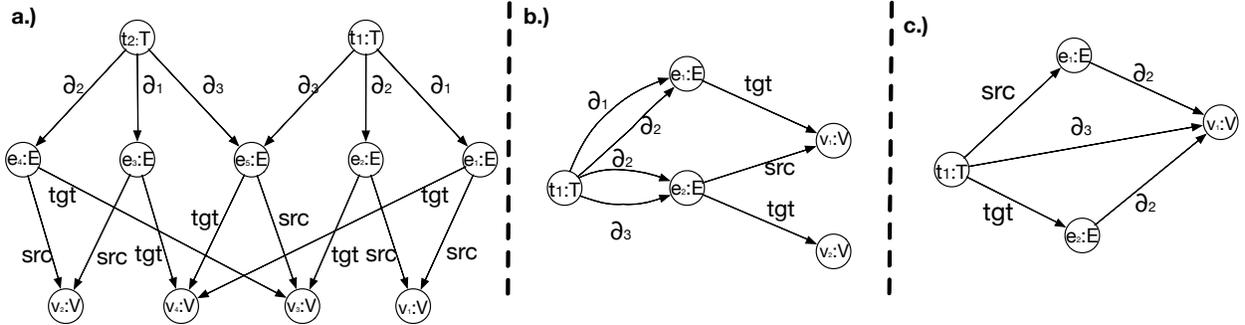
<sup>1</sup> $\mathcal{C}$ -sets are also called *copresheaves* on  $\mathcal{C}$  or *presheaves* on  $\mathcal{C}^{op}$ , and are what Löwe studied as *graph structures* or *unary algebras*.

### 2.3 Relationships between C-sets and typed graphs

One reason to prefer modeling certain domains using typed graphs or  $\mathcal{C}$ -sets rather than graphs is that the domain of interest has regularities that we wish to enforce *by construction*, rather than checking that these properties hold of inputs at runtime and verifying that every rewrite rule preserves them. There are close connections but also important differences between modeling with typed graphs or with  $\mathcal{C}$ -sets.

Every  $\mathcal{C}$ -set instance  $X$  can be functorially transformed into a typed graph. One first applies the category of elements construction,  $\int X : \mathcal{C}\text{-Set} \rightarrow \mathbf{Cat}/\mathcal{C}$ , to produce a functor into  $\mathcal{C}$ . Then the underlying graph functor  $\mathbf{Cat} \rightarrow \mathbf{Grph}$  can be applied to this morphism in  $\mathbf{Cat}$  to produce a graph typed by  $\mathcal{C}$ , i.e., a graph homomorphism into the underlying graph of  $\mathcal{C}$ . Figure 4a shows a concrete example. However, a graph typed by  $\mathcal{C}$  is only a  $\mathcal{C}$ -set under special conditions. The class of  $\mathcal{C}$ -typed graphs representable as  $\mathcal{C}$ -set instances are those that satisfy the path equations of  $\mathcal{C}$  and are, moreover, *discrete opfibrations* over  $\mathcal{C}$ . Discrete opfibrations are defined in full generality in Eq 2.<sup>2</sup>

Given a functor  $F : \mathcal{E} \rightarrow \mathcal{C} : \text{for all } x \xrightarrow{\phi} y \in \text{Hom } \mathcal{C}, \text{ and for all } e_x \in F^{-1}(x),$   
 there exists a unique  $e_x \xrightarrow{e_\phi} e_y \in \text{Hom } \mathcal{E}$  such that  $F(e_\phi) = \phi$  (2)



**Figure 4:** a.) The semi-simplicial set of Figure 3, represented as a typed graph, i.e. a labelled graph with a homomorphism into  $\Delta_2$ . b.) Another valid typed graph which is not a  $\mathcal{C}$ -set for three independent reasons: 1.)  $T_1$  has multiple edges assigned for  $\partial_2$ , 2.)  $e_1$  has no vertices assigned for  $\text{src}$ , and 3.) the last equation of  $\Delta_2$  is not satisfied. c.) A labelled graph which is not well-typed with respect to  $\Delta_2$ , i.e. no labelled graph homomorphism exists into  $\Delta_2$ .

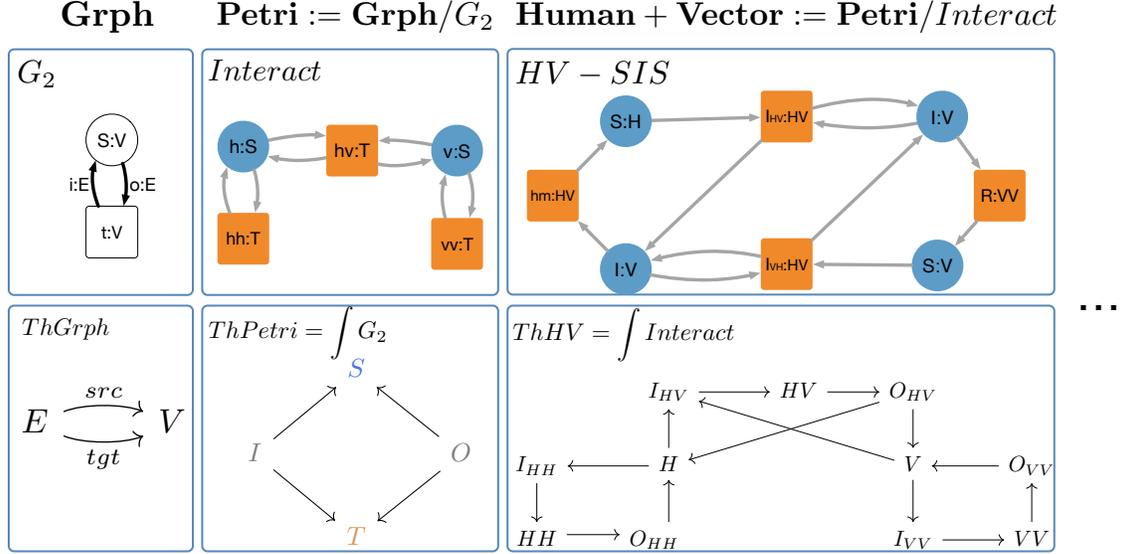
However, there is a sense in which every typed graph is a  $\mathcal{C}$ -set: there exists a schema  $\mathcal{X}$  such that  $\mathcal{X}\text{-Set}$  is equivalent to  $\mathbf{Grph}/\mathcal{X}$ . By the fundamental theorem of presheaf toposes [16],  $\mathcal{X}$  is the category of elements of the graph  $X$ , viewed as a  $\mathcal{C}$ -set on the schema for graphs. Note this procedure of creating a schema to represent objects of a slice category works beyond graphs, which we use to develop a framework of subtype hierarchies for  $\mathcal{C}$ -sets, as demonstrated in Figure 5.

Because every typed graph category is equivalent to a  $\mathcal{C}$ -set category but not the converse,  $\mathcal{C}$ -sets are a more general class of structures. The  $\mathcal{C}$ -set categories equivalent to typed graph categories are those whose instances represent sets and *relations*, in contrast with the general expressive power of  $\mathcal{C}$ -sets to represent sets and *functions*. Concretely for some edge  $a \xrightarrow{f} b$  in a type graph  $X$ , graphs typed over  $X$  can have zero, one, or many  $f$  edges for each vertex of type  $a$ , while  $\mathcal{C}$ -sets come with a restriction of there being exactly one such edge. While functions can represent relations via spans, the converse is not true.

There are practical consequences for this in graph rewriting software, if one is using typed graph rewriting to model a domain that truly has functional relationships. Because rewrite rules could take one out of the class of discrete opfibrations, as in Figure 4b, this becomes a property that one has to verify of inputs and check all rewrite rules preserve. Typed graph rewriting software can allow declaring these constraints and enforce them, but this becomes an additional engineering task outside of the underlying theory. In contrast,  $\mathcal{C}$ -sets are discrete opfibrations by construction.

Path equations are another common means of modeling a domain that are not represented in the theory of typed graph rewriting. This means, for example, that the equation  $\partial_1; \text{tgt} = \partial_2; \text{src}$  in a semi-simplicial set must be checked of all runtime inputs as well as confirmed to be preserved by each rewrite rule. This property is not straightforward to guarantee in the case of sesqui-pushout rewriting. As an upcoming example will demonstrate, it is not sufficient to

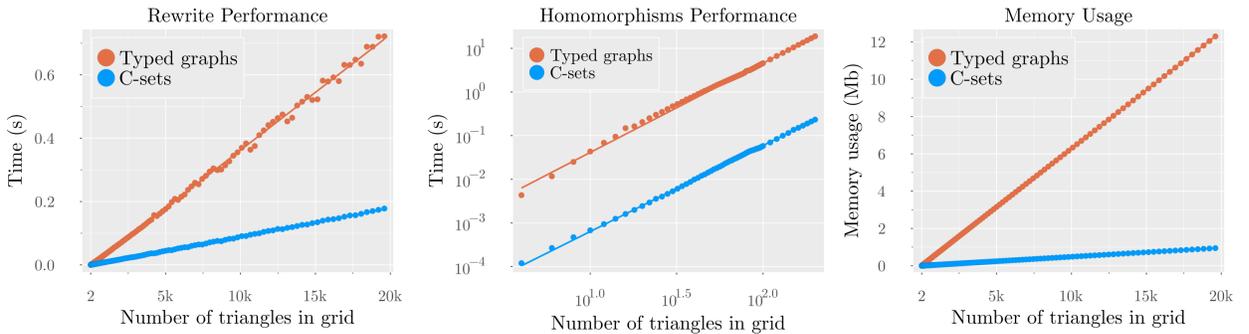
<sup>2</sup>When specialized to typed graphs,  $\mathcal{E} \xrightarrow{F} \mathcal{C}$  is a graph homomorphism and the graphs are regarded as their path categories.



**Figure 5:** Beginning with a theory of graphs, we derive a theory of whole-grain Petri nets (or bipartite graphs) by considering two distinct kinds of vertices (states and transitions) and two kinds of edges (inputs and outputs).  $ThGrph$  is constructed the category of elements of  $G_2$ . Then, taking a slice in **Petri** over an instance, *Interact*, which asserts three kinds of transitions and two kinds of states, we define a type system encoding certain domain knowledge about host-vector interactions, such as the impossibility of a transition which converts a host into a vector. As an example of subtyping, we can interpret hosts as a type of state, implying they are also a type of vertex. This process can be repeated, such as considering SIS disease dynamics for both hosts and vectors. Note that for ease of visualization,  $\mathcal{C}$ -set components at the apex of a span of morphisms (e.g.  $E, I, O$ ) are represented as directed edges.

just check that one's rewrite rule satisfies the path equalities: the rewriting itself must take path equalities into account in order to compute the correct result.

Furthermore, there are performance improvements made possible by working with  $\mathcal{C}$ -sets, rather than typed graphs. Borrowing terminology from relational databases, we first note that data in a  $\mathcal{C}$ -set is organized into distinct tables, so queries over triangles of a semi-simplicial set do not have to consider vertices or edges, for example. Secondly, the uniqueness of foreign keys allows them to be indexed, which is crucial to performance when performing queries that require table joins. This mirrors the well-known performance differences between queries of data organized in relational databases versus knowledge graphs [6]. We compare both representations within the same rewriting tool in a single benchmark experiment, described in Figure 6. This preliminary benchmark evaluates the performance of a single rewrite on semi-simplicial sets in a planar network of tessellated triangles. The rewrite locates a pair of triangles sharing an edge (i.e. a quadrilateral with an internal diagonal edge) and replaces them with a quadrilateral containing the opposite internal diagonal edge. We also chart the performance of finding all quadrilateral instances (homomorphisms) in variously sized grids. The results in Figure 6 demonstrate a lower memory footprint as well as improved rewrite and match searching for  $\mathcal{C}$ -sets.



**Figure 6:** Semisimplicial set edge flip benchmark results. Time was measured on an AMD EPYC 75F3 Milan 3.0 GHz Core with 4GB of allocated RAM.

### 3 Category-theoretic rewriting

#### 3.0.1 Pushout complements

Given a pair of arrows  $A \xrightarrow{f} B \xrightarrow{g} C$ , one constructs a pushout *complement* by finding a pair of morphisms  $A \rightarrow D \rightarrow C$  such that the resulting square is a pushout. While any category of  $\mathcal{C}$ -sets has pushouts, pushout complements are more subtle because they are not guaranteed to exist or be unique [4]. These are both desirable properties to have when using the pushout complement in rewriting, so we will demand that identification and dangling conditions (Eqs 3-4 [20]) hold, which guarantee its existence, and that the first morphism,  $f : A \rightarrow B$ , be monic, which forces it to be unique. [19]

$$\forall X \in \text{Ob } \mathcal{C}, \forall x_1, x_2 \in B_X : \quad (3)$$

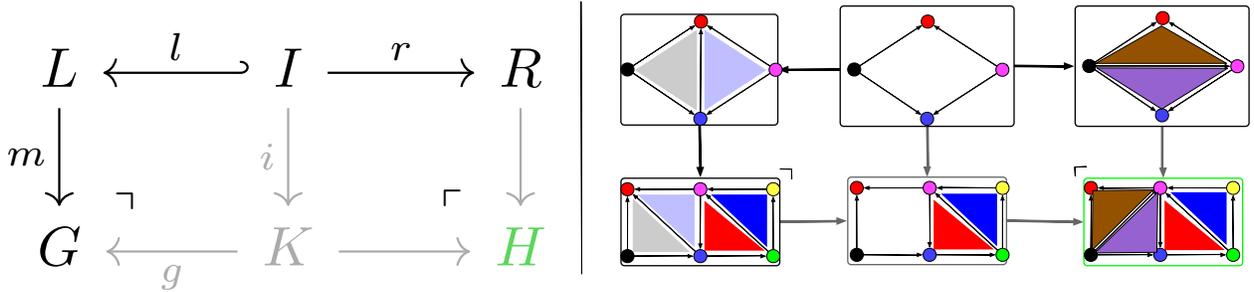
$$g_X(x_1) = g_X(x_2) \implies x_1 = x_2 \vee \{x_1, x_2\} \subseteq f_X(A_X)$$

$$\forall \phi : X \rightarrow Y \in \text{Hom } \mathcal{C}, \forall x \in C_X : \quad (4)$$

$$\phi(x) \in g_Y(B_Y - f_Y(A_Y)) \implies x \in g_X(B_X - f_X(A_X))$$

#### 3.0.2 DPO, SPO, SqPO, PBPO+

The double-pushout (DPO) algorithm [11] formalizes a notion of rewriting a portion of a  $\mathcal{C}$ -set, visualized in Figure 7. The morphism  $m$  is called the *match* morphism. The meaning of  $L$  is to provide a pattern that  $m$  will match to a sub- $\mathcal{C}$ -set in  $G$ , the target of rewriting.  $R$  represents the  $\mathcal{C}$ -set which will be substituted back in for the matched pattern to yield the rewritten  $\mathcal{C}$ -set, and  $I$  indicates what fragment of  $L$  is preserved in the rewrite and its relation to  $R$ . To perform a rewrite, first, a pushout complement computes  $K$ , the original  $\mathcal{C}$ -set with deletions applied. Second, the final rewritten  $\mathcal{C}$ -set is computed via pushout along  $r$  and  $i$ .



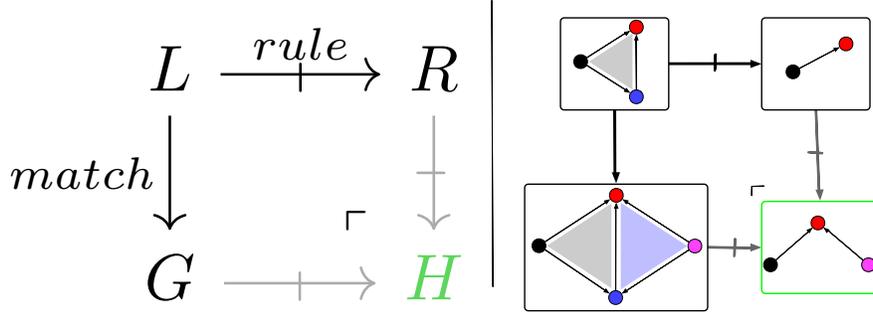
**Figure 7:** **Left:** DPO rewriting. Here and in the following figures, the initial data is in black, intermediate computations in grey, and the final result in green. **Right:** Application of a rewrite rule to flip the internal edge of a quadrilateral in a semi-simplicial set with two adjacent quadrilaterals. Here and in the following figures, colors are used to represent homomorphism data.

Single-pushout (SPO) rewriting [20] generalizes DPO rewriting, as every DPO transformation can be expressed as a SPO transformation. The additional expressivity allows us to delete in an unknown context, as demonstrated in Figure 8. The name comes from the construction being a single pushout in the category of *partial*  $\mathcal{C}$ -set morphisms,  $\mathcal{C}\text{-Par}$ .

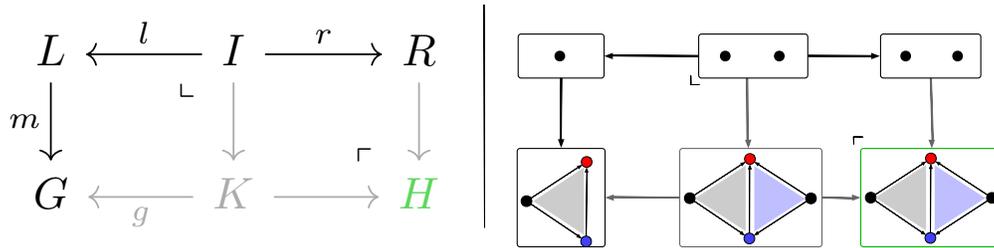
A partial  $\mathcal{C}$ -set morphism is a span  $L \xleftarrow{l} I \xrightarrow{r} R$  where  $l$  is monic. Sesqui-pushout (SqPO) rewriting [9] is a more recent technique which generalizes the previous two. It is defined in terms of the notions of partial map classifiers and final pushout complements, and it further generalizes SPO by allowing both deletion and addition in an unknown context, as demonstrated in Figure 9. Lastly, Pullback-pushout+ (PBPO+) rewriting [22] is the most recent of the four paradigms we have implemented. As shown in Figure 10, each PBPO+ rule has its own type graph,  $L'$ , which allows it to control rewriting of both the explicit matched pattern (described by  $L$ ) as well as *all* elements in the input graph  $G$  which interact with the boundary of the matched pattern. This means the notion of a match must be generalized from a match morphism  $L \rightarrow G$  to include an adherence morphism  $G \rightarrow L'$  which is an interpretation of  $G$  as typed over  $L'$ .

## 4 Design and implementation of generic categorical rewriting

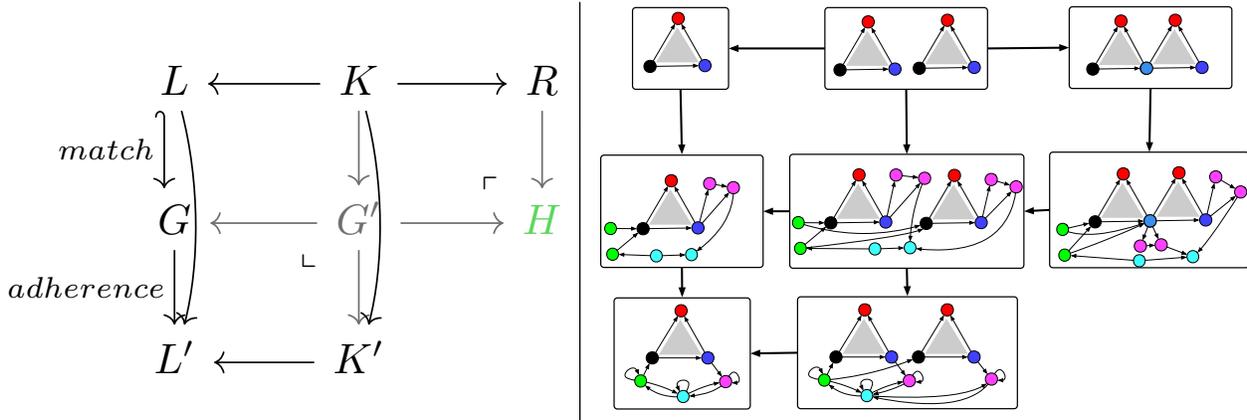
Within the paradigm of computational category theory, Catlab.jl is an open source framework for applied category theory at the center of an ecosystem of software packages called AlgebraicJulia [24, 12]. We have recently added



**Figure 8: Left: SPO rewriting Right:** An instance of deletion in an unknown context.



**Figure 9: Left: SqPO rewriting Right:** an instance of creation in an unknown context. Note that there are multiple possible pushout complements because  $l$  is not monic, but performing DPO using any of these would leave the original graph unchanged. Also note that enforcing the  $\Delta_2$  equations (in Figure 3) when computing the partial object classifier affects the results: without equations, there are four resulting ‘triangle’ objects, although two of these clearly do not form triangles.



**Figure 10: Left: PBPO+ rewriting Right:** an instance of rewriting where we explicitly control how the boundary of our matched triangular pattern is treated. The rule’s type graph  $L'$  says that, besides the matched pattern, we consider three other types of vertices: those that point at the black vertex (in green), those that are pointed at by the blue vertex (in pink) and the rest of the graph (light blue). The self loops on those extra vertices allow entire subgraphs to be mapped onto them, rather than just vertices. In  $K'$ , the rule indicates that we wish to duplicate the part of the graph that gets classified as pink (by the adherence map which assigns types to  $G$ ), while only the *edges* from the green part of the graph will get copied when we duplicate the triangle.  $L'$  has no notion of edges which are incident to the red vertex, so any input graph that has such an edge cannot be matched by this rule.

AlgebraicRewriting.jl to this ecosystem to support the categorical rewriting paradigms described above for  $\mathcal{C}$ -sets on finitely presented schemas  $\mathcal{C}$ . This class of structures balances expressivity and efficiency of manipulation, given that  $\mathcal{C}$ -sets are representable in the concrete language of relational databases [32], modulo equations in  $\mathcal{C}$ . In Catlab, each  $\mathcal{C}$ -set is automatically specialized to an efficient Julia data type; for example, when specialized to graphs, Catlab’s implementation of  $\mathcal{C}$ -sets, performs competitively against libraries optimized for graphs [24]. Catlab now occupies a unique point in the space of rewriting software tools (Table 1). For performance in pattern matching (often the typical bottleneck of rewriting), Catlab outperforms ReGraph, the nearest alternative in terms of expressive capabilities (SqPO) and usability (Table 2).

Software	Typed Graphs	$\mathcal{C}$ -sets	Rewrite type	CT Env	Last update	GUI	Scripting Env	Library vs. App
AGG[34]	Y	N	S	N	2017	Y	N	Both
Groove[27]	Y	N	S	N	2021	Y	N	App
Kappa[14]	N	N		N	2021	Y	Y	App
VeriGraph[1]	Y	N	D	Y	2017	N	Y	Lib
ReGraph[13]	Y	N	Q	N	2018	N	Y	Lib
AlgebraicRewriting	Y	Y	D,S,Q,P	Y	2022	N	Y	Lib

**Table 1:** High-level comparison with contemporary graph rewriting software packages. *Rewrite type* refers to whether DPO (D), SPO (S), SqPO (Q), and PBPO+ (P) are explicitly supported. *CT Env* refers to whether the software was implemented within a general environment of categorical abstractions beyond those immediately useful for graph rewriting. *Last update* refers to the year of the last minor version release (i.e. X.Y.0).

Mesh size	Catlab (s)	ReGraph (s)
2 by 2	$1.2 \times 10^{-4}$	$5.3 \times 10^{-3}$
2 by 3	$2.7 \times 10^{-4}$	8.0
2 by 4	$4.7 \times 10^{-4}$	1313.3
2 by 5	$6.7 \times 10^{-4}$	44979.8

**Table 2:** Catlab  $\mathcal{C}$ -set homomorphism search compared to ReGraph typed graph homomorphism search. The task was to find all quadrilateral patterns in meshes of increasing size. Tests were conducted on a single AMD EPYC 75F3 Milan 3.0 GHz Core with 4GB of RAM.

The development of Catlab has emphasized the separation of syntax and semantics when modeling a domain. This facilitates writing generic code, as diverse applications can share syntactic features, e.g. representability through string diagrams and hierarchical operad composition, with different semantic interpretations of that syntax for diverse applications. One result of this is that library code becomes very reusable, such that new features can be built from the composition of old parts with minimal additions, which reduces both developer time and the surface area for new bugs.

This point is underscored by the developer experience of implementing the above rewriting algorithms: because limits and colimits already existed for  $\mathcal{C}$ -sets, PBPO+ required no serious code writing, and the implementation of DPO only required pushout complements. Like limits and colimits, pushout complements are computed component-wise for  $\mathcal{C}$ -sets, meaning that only basic code related to pushout complements of finite sets was required. More work was needed to implement SPO because no infrastructure for the category  $\mathcal{C}\text{-Par}$  existed at the time. However, with a specification of partial morphism pushouts in terms of pushouts and pullback complements of total morphisms [17, Theorem 3.2], the only engineering required for this feature was an efficient pullback complement for  $\mathcal{C}$ -sets. Lastly, for SqPO, an algorithm for final pullback complements for  $\mathcal{C}$ -sets was the only nontrivial component that needed to be implemented, based on [8, Theorem 1] and [2, Theorem 2]. This required generalizing examples of partial map classifiers from graphs to  $\mathcal{C}$ -sets. Because the partial map classifier can be infinite for even a finitely presented  $\mathcal{C}$ -set, this type of rewriting is restricted to acyclic schemas, which nevertheless includes graphs, Petri nets, semi-simplicial sets, and other useful examples.

Because AlgebraicJulia is a collection of libraries rather than a standalone application, users have a great deal of freedom in defining their own abstractions and automation techniques, using the full power of the Julia programming language. A great deal of convenience follows from having the scripting language and the implementation language be the same: we can specify the pattern of a rewrite rule via a pushout, or we can programmatically generate repetitive rewrite rules based on structural features of a particular graph. Providing libraries rather than standalone black-box software makes integration into other projects (in the same programming language) trivial, and in virtue of being open-source library, individuals can easily extend the functionality. By making these extensions publicly available, all members of the AlgebraicJulia ecosystem can mutually benefit from each other’s efforts. As examples of this, the following additional features that have been contributed to AlgebraicRewriting.jl all serve to extend its utility as a general rewriting tool:

#### 4.1 Computation of homomorphisms and isomorphisms of $\mathcal{C}$ -sets

For rewriting algorithms to be of practical use, morphisms matching the left-hand-side of rules must somehow be supplied. The specification of a  $\mathcal{C}$ -set morphism requires a nontrivial amount of data that must satisfy the naturality condition. Furthermore, in confluent rewriting systems, manually finding matches is an unreasonable request to make of the end user, as the goal is to apply all rewrites possible until the term reaches a normal form. For this reason, DPO rewriting of  $\mathcal{C}$ -sets benefits from a generic algorithm to find homomorphisms, analogous to structural pattern matching in the tree term rewriting case.

The problem of finding a  $\mathcal{C}$ -set homomorphism  $X \rightarrow Y$ , given a finitely presented category  $\mathcal{C}$  and two finite  $\mathcal{C}$ -sets  $X$  and  $Y$ , is generically at least as hard as the graph homomorphism problem, which is NP-complete. On the other hand, the  $\mathcal{C}$ -set homomorphism problem can be framed as a constraint satisfaction problem (CSP), a classic problem in computer science for which many algorithms are known [31, Chapter 6]. Since  $\mathcal{C}$ -sets are a mathematical model of relational databases [33], the connection between  $\mathcal{C}$ -set homomorphisms and constraint satisfaction is a facet of the better-known connection between databases and CSPs [35].

To make this connection precise, we introduce the slightly nonstandard notion of a typed CSP. Given a finite set  $T$  of types, the slice category  $\mathbf{FinSet}/T$  is the category of  $T$ -typed finite sets. A typed CSP then consists of  $T$ -typed finite sets  $V$  and  $D$ , called the variables and the domain, and a finite set of constraints of form  $(\mathbf{x}, R)$ , where  $\mathbf{x} = (x_1, \dots, x_k)$  is a list of variables and  $R \subseteq D^{-1}(V(x_1)) \times \dots \times D^{-1}(V(x_k))$  is a compatibly typed  $k$ -ary relation. An assignment is a map  $\phi : V \rightarrow D$  in  $\mathbf{FinSet}/T$ . The objective is to find a solution to the CSP, namely an assignment  $\phi$  such that  $(\phi(x_1), \dots, \phi(x_k)) \in R$  for every constraint  $(\mathbf{x}, R)$ .

The problem of finding a  $\mathcal{C}$ -set morphism  $X \rightarrow Y$  translates to a typed CSP by taking the elements of  $X$  and  $Y$  to be the variables and the domain of the CSP, respectively. To be precise, let the types  $T$  be the objects of  $\mathcal{C}$ . The variables  $V : \{(c, x) : c \in \mathcal{C}, x \in X(c)\} \rightarrow \mathbf{Ob} \mathcal{C}$  are given by applying the objects functor  $\mathbf{Ob} : \mathbf{Cat} \rightarrow \mathbf{Set}$  to  $\int X \rightarrow \mathcal{C}$ , the category of elements of  $X$  with its canonical projection. Similarly, the domain is  $D := \mathbf{Ob}(\int Y \rightarrow \mathcal{C})$ . Finally, for every generating morphism  $f : c \rightarrow c'$  of  $\mathcal{C}$  and every element  $x \in X(c)$ , introduce a constraint  $((x, x'), R)$  where  $x' := X(f)(x)$  and  $R := \{(y, y') \in Y(c) \times Y(c') : Y(f)(y) = y'\}$  is the graph of  $Y(f)$ . By construction, an assignment  $\phi : V \rightarrow D$  is the data of a  $\mathcal{C}$ -set transformation (not necessarily natural) and  $\phi$  is a solution if and only if the transformation is natural. Thus, the solutions of the typed CSP are exactly the  $\mathcal{C}$ -set homomorphisms  $X \rightarrow Y$ .

With this reduction, CSP algorithms are straightforwardly ported to algorithms for finding  $\mathcal{C}$ -set morphisms, where the types and special structure permits optimizations, one example being the use of the discrete opfibration condition to accelerate the search. We only consider assignments that satisfy the typing relations. We have adapted backtracking search [31, Section 6.3], a simple but fundamental CSP algorithm, to find  $\mathcal{C}$ -set homomorphisms. By also maintaining a partial inverse assignment, this algorithm is easily extended to finding  $\mathcal{C}$ -set monomorphisms, an important constraint when matching for rewriting. Since a monomorphism between finite  $\mathcal{C}$ -sets  $X$  and  $Y$  is an isomorphism if and only if  $X(c)$  and  $Y(c)$  have the same cardinality for all  $c \in \mathcal{C}$ , this extension also yields an algorithm for isomorphism testing, which is useful for checking the correctness of rewrites.

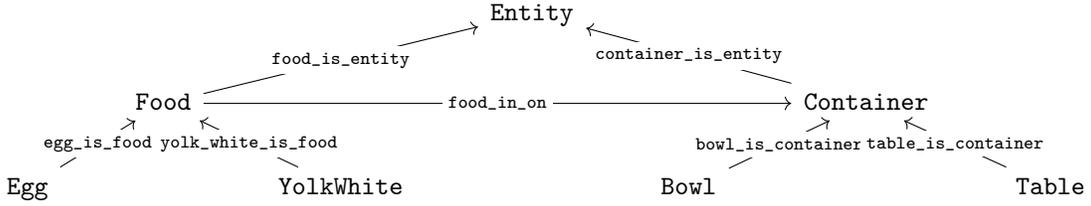
#### 4.2 Diagrammatic syntax

Specifying DPO rewrite rules can be cumbersome as a significant amount of combinatorial data is contained in a span of  $\mathcal{C}$ -sets. To make our system more user-friendly, we have developed a symbolic domain-specific language (DSL) to specify rewrite rules, based on the idea of assembling  $\mathcal{C}$ -sets from the atomic ones known as *representables*. This involves no loss of generality since every  $\mathcal{C}$ -set can be expressed as a colimit of representable  $\mathcal{C}$ -sets [28, Theorem 6.5.7]. For instance, in the category of graphs, the two representables are the graphs with one isolated vertex and with one edge between two distinct vertices, and clearly every graph is a colimit of copies of these two graphs. An example of specifying a rewrite rule in this manner, using a much more elaborate schema, is shown in Figure 11.

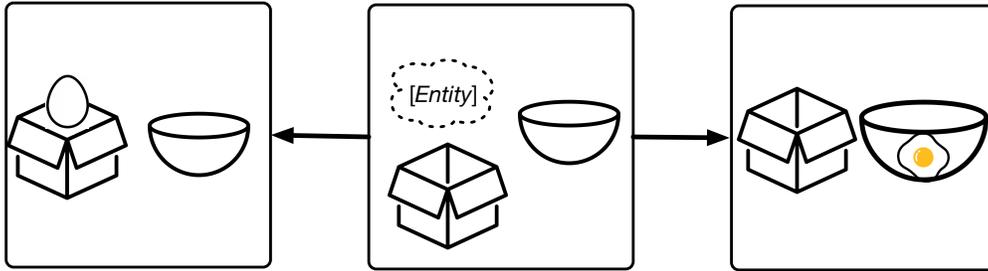
The mathematics behind our DSL uses the underappreciated fact that the diagrams in a given category are themselves the objects of a category; as described in [26, 25, 23] and references therein. Given a category  $\mathcal{S}$ , the *diagram category*  $\mathbf{Diag}(\mathcal{S})$  has, as objects, diagrams  $D : \mathcal{J} \rightarrow \mathcal{S}$  in  $\mathcal{S}$ , and as morphisms  $(\mathcal{J}, D) \rightarrow (\mathcal{J}', D')$ , a functor  $R : \mathcal{J} \rightarrow \mathcal{J}'$  along with a natural transformation  $\rho : D \Rightarrow D' \circ R$ . Another diagram category  $\mathbf{Diag}^{\text{co}}(\mathcal{S})$  is defined similarly, except that the natural transformation in a morphism  $(R, \rho)$  goes in the opposite direction:  $\rho : D' \circ R \Rightarrow D$ .

We now show that a span in  $\mathbf{Diag}^{\text{co}}(\mathcal{C})$  presents a span in  $\mathcal{C}\text{-Set}$ , i.e., a DPO rewrite rule for  $\mathcal{C}$ -sets, as colimits of representables and morphisms between them. The category  $\mathbf{Diag}^{\text{co}}(\mathcal{C})$  has the advantage of referring only to the schema  $\mathcal{C}$  and so can be described syntactically given a finite presentation of  $\mathcal{C}$ .

**Proposition 1.** *By applying the Yoneda embedding and taking colimits, a span in the category  $\mathbf{Diag}^{\text{co}}(\mathcal{C})$  induces a span of  $\mathcal{C}$ -sets.*



(a) Fragment of a schema that models recipes for cooking breakfast



(b) Cartoon visualization of egg cracking rule. Notably we require an abstract entity in the interface, mapping to both the egg and yolk+white, to reflect that they are the same entity.

```

crack_egg_in_bowl = @migration SchCospan SchBreakfastKitchen begin
  L => @join begin # left-hand side of rule
    bowl::Bowl
    egg::Egg
  end
  I => @join begin # intermediate state of rule
    bowl::Bowl
    egg_entity::Entity # entity underlying egg and yolk-white
    old_container::Container # original container of food
  end
  R => @join begin # right-hand side of rule
    bowl::Bowl
    yolk_white::YolkWhite
    food_in_on(yolk_white_is_food(yolk_white)) == bowl_is_container(bowl)
    old_container::Container
  end
  l => begin # left map in rule
    bowl => bowl
    egg_entity => food_is_entity(egg_is_food(egg))
    old_container => food_in_on(egg_is_food(egg))
  end
  r => begin # right map in rule
    bowl => bowl
    egg_entity => food_is_entity(yolk_white_is_food(yolk_white))
    old_container => old_container
  end
end
end
    
```

(c) DPO rewrite rule specified using diagrammatic syntax. This syntax allows us to avoid explicitly treating the underlying entity of the container, for example.

**Figure 11:** Example of a DPO rewrite rule specified using the diagrammatic syntax, adapted from a planning system for the cooking domain.

*Proof.* It is enough to define a functor  $\text{Diag}^{\text{co}}(\mathcal{C}) \rightarrow \mathcal{C}\text{-Set}$ , which we do as the following composite

$$\text{Diag}^{\text{co}}(\mathcal{C}) \xrightarrow{\text{op}} \text{Diag}(\mathcal{C}^{\text{op}}) \xrightarrow{\text{Diag}(y)} \text{Diag}(\mathcal{C}\text{-Set}) \xrightarrow{\text{colim}} \mathcal{C}\text{-Set},$$

where  $\text{op} : \mathbf{Cat}^{\text{co}} \rightarrow \mathbf{Cat}$  is the oppositization 2-functor and  $y : \mathcal{C}^{\text{op}} \rightarrow \mathcal{C}\text{-Set}$  is the Yoneda embedding for  $\mathcal{C}$ . We are using the facts that the diagram construction extends to a (2-)functor  $\text{Diag} : \mathbf{Cat} \rightarrow \mathbf{Cat}$  in which morphisms act by postcomposition [25, §2.1] and that taking colimits is functorial with respect to the category  $\text{Diag}(\mathcal{S})$  whenever  $\mathcal{S}$  is cocomplete [25, §5.1].  $\square$

### 4.3 Typed graph rewriting with slice categories

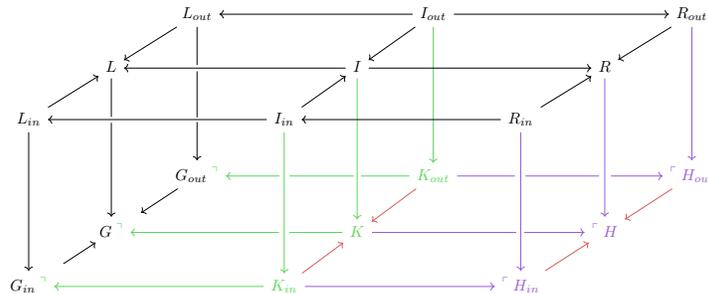
Slice categories offer a form of constraining  $\mathcal{C}$ -sets without altering the schema. Consider the example of rewriting string diagrams encoded as hypergraph cospans [3]. These can be used to represent terms in a symmetric monoidal theory, where it is important to restrict diagrams to only those which draw from a fixed set of boxes with particular arities, given by a monoidal signature  $\Sigma$ , which induces the unique hypergraph  $H\Sigma$  which has all box types from  $\Sigma$  and a single vertex. Working within the slice category  $\mathbf{Hyp}/H\Sigma$  prevents us from performing rewrites which violate the arities of the operations specified by  $\Sigma$ .

There are two ways to implement rewriting in  $\mathcal{C}\text{-Set}/X$  for a particular  $\mathcal{C}$ : the computation can be performed with the objects  $L, I, R, G$  being  $\mathcal{C}$ -set morphisms, or it can be performed in  $[f X, \mathbf{Set}]$ . Programming with generic categorical abstraction greatly lowered the barrier to implementing both of these: for the former, what was needed was to relate the pushout and pushout complement of  $\mathcal{C}\text{-Set}/X$  to the corresponding computations in  $\mathcal{C}\text{-Set}$ . The barrier to the latter was to compute the category of elements and migrate data between the two representations, code which had already been implemented. As the former strategy requires less data transformation, it is preferred.

### 4.4 Open system rewriting with structured cospans

The forms of rewriting discussed up to this point have concerned rewriting closed systems. Structured cospans are a general model for open systems, which formalize the notion of gluing together systems which have designated inputs and outputs. Open systems are modeled as cospans of form  $La \rightarrow x \leftarrow Lb$ , where the apex  $x$  represents the system itself and the feet  $La$  and  $Lb$  represent the inputs and outputs, typically discrete systems such as graphs without edges. Here,  $L : A \rightarrow X$  is a functor that maps from the system category  $A$  to the system interface category  $X$ , and  $L$  must be a left adjoint between categories with finite colimits.<sup>3</sup> Larger systems are built up from smaller systems via pushouts in  $X$ , which glue systems together along a shared interface:  $(La \rightarrow x \leftarrow Lb \rightarrow y \leftarrow Lc) \mapsto (La \rightarrow x +_{Lb} y \leftarrow Lc)$ .

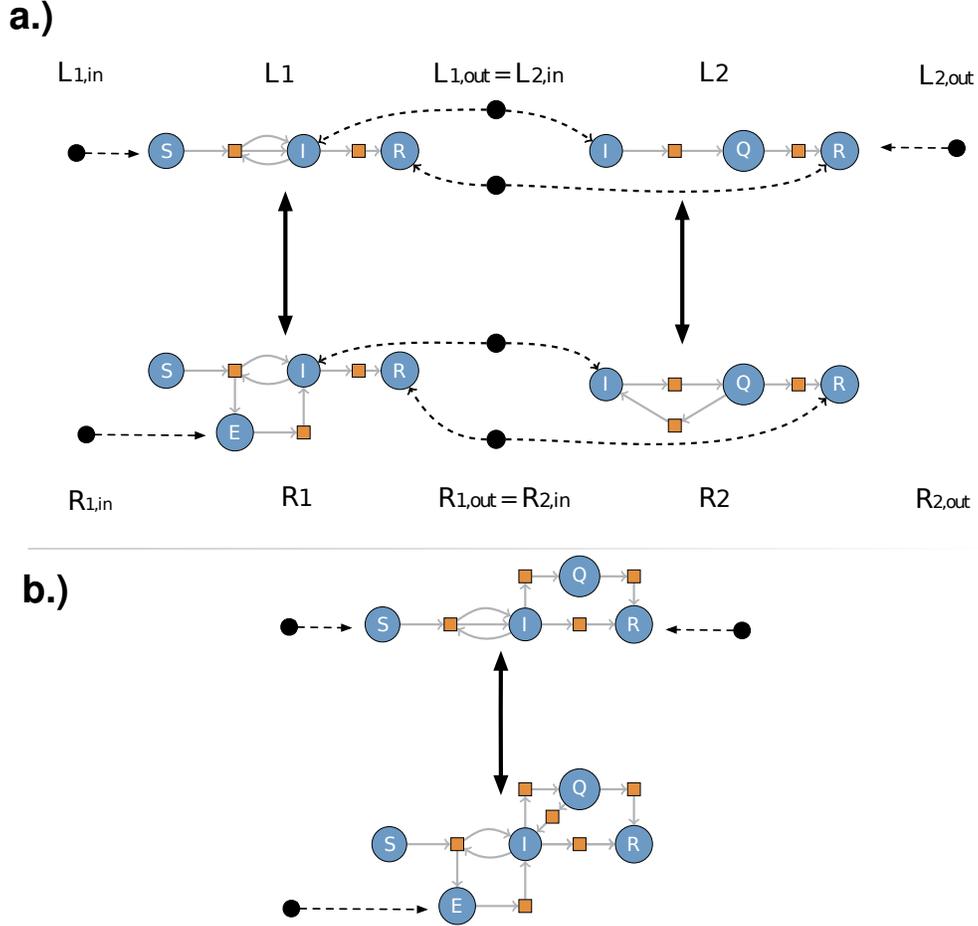
When  $L, I$ , and  $R$  are each structured cospans, there is extra data to consider when rewriting, as shown in Figure 12. In ordinary DPO rewriting, if the  $R$  of one rewrite rule equals the  $L$  of another, a composite rewrite rule can be constructed, which could be called *vertical* composition. In the case of structured cospans, *horizontal* composition emerges from composing the  $L, I$ , and  $R$  of two structured cospan rules pairwise, visualized in Figure 13. These two forms of composition together yield a double category of structured cospan rewrites, where horizontal arrows are in correspondence with structured cospans and squares are in correspondence with all possible rewrites [7].



**Figure 12:** Applying a structured cospan rewrite rule.  $\mathcal{C}$ -sets and morphisms in black are the initial data: the upper face represents the open rewrite rule, the upper left edge represents the open pattern to be matched, and the left face represents the matching. Green morphisms are computed by pushout complement in  $\mathcal{C}\text{-Set}$ . The purple morphisms are computed by the rewriting pushouts and red morphisms are computed by the structured cospan pushouts. Figure adapted from [7, Section 4.2].

While this compositional approach to building open systems can be an illuminating way to organize information about a complex system, there can also be computational benefits. When searching for a match in a large  $\mathcal{C}$ -set, the search

<sup>3</sup>The  $L$  of structured cospans should not be confused with the  $L$  of the rewrite rule  $L \leftarrow I \rightarrow R$ .



**Figure 13: a.)** Example of horizontal composition of structured cospan rewrite rules. The  $L$  and  $R$  structured cospans are positioned on the top and bottom, respectively. For clarity,  $I$  cospans are omitted. **b.)** The result of composition.

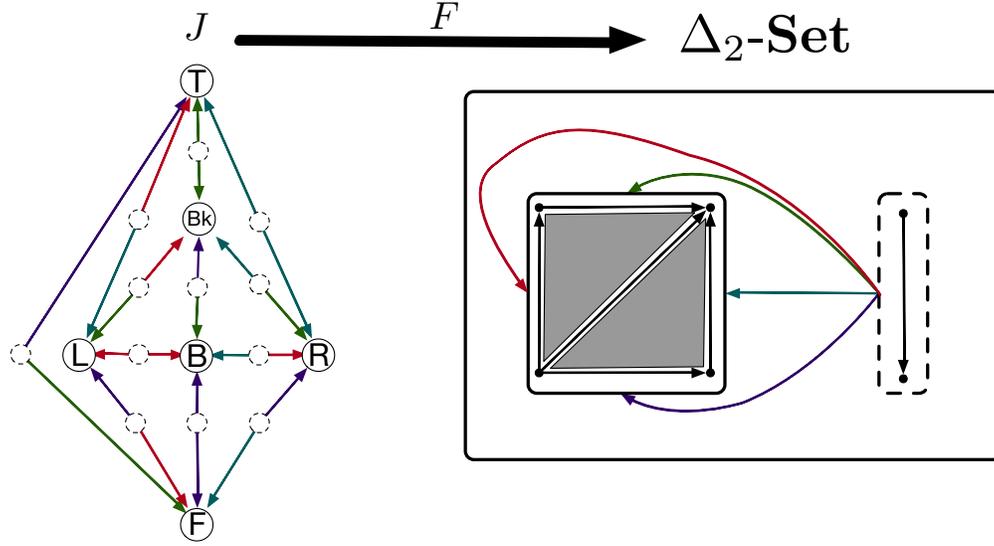
space grows as  $O(n^k)$  where  $k$  is the size of the pattern  $L$  and  $n$  is the size of  $G$ . However, after decomposing  $G$  into a composite of substructures and restricting matches to homomorphisms into a specific substructure, the search space is limited by  $O(m^k)$  where  $m < n$  is the size of the substructure. Not only does this accelerate the computation, but it can be semantically meaningful to restrict matches to those which do not cross borders.

#### 4.5 Distributed graph rewriting

Distributed graphs offer an alternative formalism that allows one to decompose a large graph into smaller ones while maintaining consistency at the boundaries, and thus it is another strategy for parallelizing computations over graphs. The content of a distributed graph can be succinctly expressed in the language of category theory as a diagram in **Grph**. Because Catlab has sophisticated infrastructure in place for manipulating categories of diagrams, it merely takes specializing the codomain of the Diagram datatype to **Grph** to represent distributed graphs and their morphisms. Note that we can easily generalize to distributed semi-simplicial sets or other  $\mathcal{C}$ -sets (Figure 14). Colimits in the category of diagrams (in a cocomplete category) are defined in terms of left Kan extensions [26], and with our implementation [5] it is possible to develop a rewriting tool for distributed graphs.

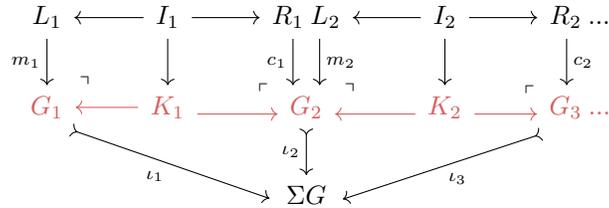
#### 4.6 Graph processes

Given a concrete sequence of rewrites, perhaps representing a sequence of actions required to take one from an initial state to some desired state, it is of practical importance to represent the steps taken in a maximally-parallel manner that has only the necessary dependencies, such as one rewrite step creating an element that another rewrite step deletes. Graph processes [10] are a construction which exposes the causal dependencies between rewrites as a partially-ordered



**Figure 14:** Constructing the surface of a cube compositionally with a distributed graph.  $F$  sends the solid circles to the square face graph and the dashed circles to the edge graph. Colors indicate which morphism from the edge to the face which controls how the faces are being glued together. We construct the assembled cube as a  $\mathcal{C}$ -set simply by taking the colimit of the diagram.

set. The construction of this partial order is expressed as a colimit of a certain bipartite diagram, as shown in Figure 15. Colimits of diagrams being readily computable in Catlab led to this extension requiring only a small amount of programmer effort.



**Figure 15:** The graph processes construction from a sequence of rewrites with match morphisms  $m_i$  and co-match morphisms  $c_i$  labeled.  $\Sigma G$  is constructed as the colimit of the red subdiagram, and its role is to identify the same elements across time, if we interpret  $G_i$  as a temporal sequence. Therefore, given a notion of element production, deletion, and preservation, if  $i$  produces some element that  $j$  preserves or deletes, there must be a causal dependency  $i < j$ .

#### 4.7 Further extensions

Examples of further features, such as negative application conditions, parallel rewriting, rewriting with functions applied to attributes, matching variables on attributes, (e.g. one rule which can identify any triangle that has exactly two edges with an equal length attribute and rewrite to make all three edges have that length) are found in AlgebraicRewriting documentation or tests.

## 5 Conclusions and Future Work

There are many desiderata for software development in academic and industrial settings alike, such as velocity of development, robustness to future changes in design, and correctness. We demonstrated how designing software with category-theoretic abstractions facilitates the achievement all three of these, using the mature field of graph rewriting software as a case study.

While current graph transformation software in use is often very specialized to particular domains, such as chemistry, we show that DPO, SPO, SqPO, and PBPO+ rewriting can be efficiently performed on  $\mathcal{C}$ -sets, which are viewed as

a subset of typed graphs (discrete opfibrations) with desirable theoretical and performance characteristics, and we have presented the first practical implementation for this. This result allows generic rewrite operations to be used in a variety of contexts, when it would otherwise be time-consuming and error-prone to develop custom rewrite algorithms for such a multitude of data structures or to work with typed graphs and enforce the discrete opfibration condition by other means. We also extended these implementations to the first practical implementations of homomorphism search, structured cospan rewriting, and distributed graphs for arbitrary  $\mathcal{C}$ -sets. Our internal benchmark showed that  $\mathcal{C}$ -set rewriting can leverage the discrete opfibration condition to outperform typed graphs in memory and speed, and an external benchmark showed a significant speedup relative to comparable graph rewriting software.

Catlab and AlgebraicRewriting could be extended to a tool for graph transformation researchers to computationally validate and explore new ideas. Researchers interested developing tools to be directly consumed by others could produce a performant and easily interoperable instantiation of their work. Even those interested in rewriting systems as mathematical objects can benefit from this process by gaining intuition and empirically testing conjectures about their constructions. However, many useful concepts from graph rewriting have yet to be added, such as rule control mechanisms and rule algebras, but the extensibility of Catlab allows researchers to do this on their own or with the support of Catlab's active user community.

To create tools for practicing scientists and engineers, our future work involves building practical scientific software that applies rewriting in each its main areas, i.e. *graph relations*, *languages*, and *transition systems*: respectively, a theorem prover for symmetric monoidal categories by performing e-graph equality saturation [36] with rewriting, a tool for defining and exploring a language of open epidemiological models, and a general agent-based model simulator.

## References

- [1] Azzi, G.G., Bezerra, J.S., Ribeiro, L., Costa, A., Rodrigues, L.M., Machado, R.: The verigraph system for graph transformation. In: Graph Transformation, Specifications, and Nets, pp. 160–178. Springer (2018)
- [2] Behr, N., Harmer, R., Krivine, J.: Concurrency theorems for non-linear rewriting theories. In: International Conference on Graph Transformation. pp. 3–21. Springer (2021)
- [3] Bonchi, F., Gadducci, F., Kissinger, A., Sobocinski, P., Zanasi, F.: String diagram rewrite theory i: Rewriting with frobenius structure. arXiv preprint arXiv:2012.01847 (2020)
- [4] Braatz, B., Golas, U., Soboll, T.: How to delete categorically—two pushout complement constructions. Journal of Symbolic Computation **46**(3), 246–271 (2011)
- [5] Brown, K., Hanks, T., Fairbanks, J.: Compositional exploration of combinatorial scientific models (2022). doi:10.48550/ARXIV.2206.08755, <https://arxiv.org/abs/2206.08755>
- [6] Cheng, Y., Ding, P., Wang, T., Lu, W., Du, X.: Which category is better: benchmarking relational and graph database management systems. Data Science and Engineering **4**(4), 309–322 (2019)
- [7] Cicala, D.: Rewriting structured cospans: A syntax for open systems. arXiv preprint arXiv:1906.05443 (2019)
- [8] Corradini, A., Duval, D., Echahed, R., Prost, F., Ribeiro, L.: Agree–algebraic graph rewriting with controlled embedding. In: International Conference on Graph Transformation. pp. 35–51. Springer (2015)
- [9] Corradini, A., Heindel, T., Hermann, F., König, B.: Sesqui-pushout rewriting. In: International Conference on Graph Transformation. pp. 30–45. Springer (2006)
- [10] Corradini, A., Montanari, U., Rossi, F.: Graph processes. Fundamenta Informaticae **26**(3, 4), 241–265 (1996)
- [11] Ehrig, H., Pfender, M., Schneider, H.J.: Graph-grammars: An algebraic approach. In: 14th Annual Symposium on Switching and Automata Theory (swat 1973). pp. 167–180. IEEE (1973)
- [12] Halter, M., Patterson, E., Baas, A., Fairbanks, J.: Compositional scientific computing with catlab and semantic-models. arXiv preprint arXiv:2005.04831 (2020)
- [13] Harmer, R., Oshurko, E.: Reversibility and composition of rewriting in hierarchies. arXiv preprint arXiv:2012.01661 (2020)
- [14] Hayman, J., Heindel, T.: Pattern graphs and rule-based models: The semantics of kappa. In: International Conference on Foundations of Software Science and Computational Structures. pp. 1–16. Springer (2013)
- [15] Heckel, R., Lambers, L., Saadat, M.G.: Analysis of graph transformation systems: Native vs translation-based techniques. arXiv preprint arXiv:1912.09607 (2019)
- [16] Kashiwara, M., Schapira, P.: Categories and Sheaves. Springer Berlin Heidelberg (2006). doi:10.1007/3-540-27950-4, <https://doi.org/10.1007/3-540-27950-4>

- [17] Kennaway, R.: Graph rewriting in some categories of partial morphisms. In: International Workshop on Graph Grammars and their Application to Computer Science. pp. 490–504. Springer (1990)
- [18] Lack, S., Sobociński, P.: Adhesive categories. In: International Conference on Foundations of Software Science and Computation Structures. pp. 273–288. Springer (2004)
- [19] Lack, S., Sobociński, P.: Adhesive and quasiadhesive categories. *RAIRO-Theoretical Informatics and Applications* **39**(3), 511–545 (2005)
- [20] Löwe, M.: Algebraic approach to single-pushout graph transformation. *Theoretical Computer Science* **109**(1-2), 181–224 (1993)
- [21] Minas, M., Schneider, H.J.: Graph transformation by computational category theory. In: Graph Transformations and Model-Driven Engineering, pp. 33–58. Springer (2010)
- [22] Overbeek, R., Endrullis, J., Rosset, A.: Graph rewriting and relabeling with pbpo+: A unifying theory for quasitoposes (2022). doi:10.48550/ARXIV.2203.01032
- [23] Patterson, E., Baas, A., Hosgood, T., Fairbanks, J.: A diagrammatic view of differential equations in physics. *Mathematics in Engineering* **5**(2), 1–59 (2023). doi:10.3934/mine.2023036
- [24] Patterson, E., Lynch, O., Fairbanks, J.: Categorical data structures for technical computing. arXiv preprint arXiv:2106.04703 (2021)
- [25] Perrone, P., Tholen, W.: Kan extensions are partial colimits. *Applied Categorical Structures* (2022). doi:10.1007/s10485-021-09671-9
- [26] Peschke, G., Tholen, W.: Diagrams, fibrations, and the decomposition of colimits. arXiv preprint arXiv:2006.10890 (2020)
- [27] Rensink, A., Boneva, I., Kastenbergh, H., Staijen, T.: User manual for the groove tool set. Department of Computer Science, University of Twente, The Netherlands (2010)
- [28] Riehl, E.: Category theory in context. Courier Dover Publications (2016), <http://www.math.jhu.edu/~eriehl/context.pdf>
- [29] Ringer, T., Palmkog, K., Sergey, I., Gligoric, M., Tatlock, Z.: Qed at large: A survey of engineering of formally verified software. arXiv preprint arXiv:2003.06458 (2020)
- [30] Rushby, J.: Automated test generation and verified software. In: Working Conference on Verified Software: Theories, Tools, and Experiments. pp. 161–172. Springer (2005)
- [31] Russell, S., Norvig, P.: Artificial intelligence: a modern approach (2010)
- [32] Schultz, P., Spivak, D.I., Vasilakopoulou, C., Wisnesky, R.: Algebraic databases. arXiv preprint arXiv:1602.03501 (2016)
- [33] Spivak, D.I.: Functorial data migration. *Information and Computation* **217**, 31–51 (2012). doi:10.1016/j.ic.2012.05.001
- [34] Taentzer, G.: Agg: A graph transformation environment for modeling and validation of software. In: International Workshop on Applications of Graph Transformations with Industrial Relevance. pp. 446–453. Springer (2003)
- [35] Vardi, M.Y.: Constraint satisfaction and database theory: a tutorial. In: Proceedings of the nineteenth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems. pp. 76–85 (2000). doi:10.1145/335168.335209
- [36] Willsey, M., Nandi, C., Wang, Y.R., Flatt, O., Tatlock, Z., Panckheha, P.: egg: fast and extensible equality saturation. *Proceedings of the ACM on Programming Languages* **5**(POPL), 1–29 (2021)