

*To my beloved
Benedetta*

Introduction

English version

Italian version

[Traduzione italiana dell'introduzione.](#)

Contents

Introduction	i
1 Parametric Optics for Gradient-Based Learning	1
1.1 Categorical toolkit	1
1.1.1 Actegories	1
1.1.2 The Para construction	2
1.1.3 Optics	4
1.1.4 Weighted optics	5
1.1.5 Differential categories	6
1.1.6 Parametric lenses	9
1.2 Supervised learning with parametric lenses	9
1.2.1 Model, loss, optimizer, learning rate	9
1.2.2 Weight tying, batching, and the learning iteration	10
1.2.3 Empirical evidence	11
1.3 Future directions and related work	12
1.3.1 Learners	12
1.3.2 Exotic differential categories	13
1.3.3 Functional reverse-mode automatic differentiation	14
2 From Classical Computer Science to Neural Networks	15
2.1 Categorical toolkit	15
2.1.1 (Co)algebras	15
2.1.2 Integral transform	17
2.2 Categorical deep learning	19
2.2.1 From GDL to CDL	19
2.2.2 (Co)inductive definitions for RNNs	20
2.3 Algorithmic alignment: GNNs and dynamic programming	22
2.3.1 Integral transforms for GNNs and dynamic programming	22
2.3.2 Asynchronous algorithmic alignment	24
2.4 Future directions and related work	25
2.4.1 Differentiable causal computations	26
2.4.2 Sheaf neural networks	27
3 Functor Learning	31
3.1 Using functors to separate layers of abstraction	31
3.1.1 Schemas, architectures, models, and concepts	31
3.1.2 Datasets influence the semantics of tasks	34
3.2 Categorical representation learning	35

3.2.1	Unsupervised functorial translation	35
3.2.2	Compositional distributional model for NLP	38
3.3	Equivariant neural networks as functors	40
3.3.1	An equivariant classifier to diagnose Alzheimer’s disease	40
3.3.2	Transporting algebraic structure into embedding spaces	42
4	Detailed Representations of Neural Networks	45
4.1	Neural circuit diagrams	47
4.1.1	Functor string diagrams	47
4.1.2	Neural circuit diagrams	47
4.2	Diagrams with universal approximators	49
4.2.1	The diagrammatic approach of [KLLWM24]	50
4.2.2	Comparing attention mechanisms	52
4.3	Future directions and related work	52
4.3.1	Tensor networks	53
4.3.2	Non-diagrammatic notations	54
4.3.3	A multifaceted approach	55
	Conclusions	57
	Bibliography	63

List of Figures

1.1	Parametric morphisms	3
1.2	Lenses	4
1.3	Parametric lens	9
1.4	Supervised learning with parametric lenses	11
1.5	Weight tying and batching	12
2.1	(Co)algebra homomorphisms	16
2.2	Integral transform	18
2.3	Folding RNN	21
2.4	Lax-commutative diagram for the unrolling of a folding RNN	21
2.5	Neural network cells	22
2.6	Bellman-Ford algorithm and message passing networks	23
2.7	Graph neural network architecture	25
2.8	Tiles and stateful morphisms	27
2.9	Delayed trace and derivatives in $\mathbf{St}_0(\mathcal{C})$	28
3.1	CycleGAN schema	32
3.2	CycleGAN schema	33
3.3	Pregroup representation of grammar	39
3.4	Classifier	41
3.5	Classifier	42
3.6	Transport of algebraic structure	43
4.1	Diagram of a transformer	45
4.2	Functor string diagrams	48
4.3	Neural string diagrams	49
4.4	From functor string diagrams to neural circuit diagrams	50
4.5	Neural string diagram of a transformer	51
4.6	Diagrammatic syntax of [KLLWM24]	52
4.7	Dimensionality, SIMD boxes, and syntactic sugar	53
4.8	Transformer blueprint according to [KLLWM24]	53
4.9	Abstract blueprint for self-attention	54
4.10	Evolutionary tree of attention	55
4.11	GPT-2 architecture	56

Chapter 1

Parametric Optics for Gradient-Based Learning

Despite the unquestionable success that gradient-based deep learning has enjoyed in recent years, the field is still both young and poorly understood. As mentioned in the introduction, the lack of theoretical underpinnings means that good performance is highly dependent on *ad hoc* choices and empirical euristics leading to brittleness and poorly understood phenomena ([CGG⁺22], [SGW21], [Gav24]). The ever-growing complexity of deep learning models poses significant challenges both in terms of optimization ([Ell18]) and architectural design ([GLD⁺24]), and, while there are a number of general purpose deep learning libraries that automatically implement backpropagation and provide tools for designing a wide variety of neural networks, these tools often rely on inelegant machinery difficult to parallelize ([Ell18]). Given the ever-increasing role gradient based learning plays in the sciences, in industry, and in every day life, solving these issues is of the outmost importance.

Hence, it would be auspicious to develop a mathematically structured framework for gradient-based learning able to act as a bridge between low-level automatic differentiation and high level architectural specifications ([Gav24]). The great number of architectures developed in recent years and the inherently modular structure of deep neural networks call for a model which is general (that is, not dependent on a specific differentiation algorithm or a specific optimizer) and compositional (that is, we should be able to predict the behavior of the entire model if the behavior of each part is known). Among the various attempts (see e.g. [FST19] or [SXM23]), a promising combination of differential categories, parametrization and optics has been recently proposed by [CCG⁺19], [CGG⁺22], [Gav24] and more as full-featured gradient-based framework able to challenge established tools. In this chapter, we illustrate such framework and part of its mathematical foundations.

1.1 Categorical toolkit

Learning neural networks have two main properties: they depend on parameters and information flows through them bidirectionally (forward propagation and back propagation): any aspiring categorical model of gradient-based learning must take these two aspects into consideration. A number of authors (see e.g. [Gav24], and [CGG⁺22]) have proposed the **Para** construction as a categorical model of parameter dependence and various categories of optics as the right categorical abstraction for bidirectionality.

1.1.1 Actegories

Before we can deal with parametric maps, we need to find a way to glue input/output spaces to parameter spaces, so that such maps have well-defined domains. One common strategy is to provide the category at hand with a monoidal structure. However, monoidal products can only combine elements within the same underlying category. Since (co)parameters are often taken from spaces that are different in nature from the input and output spaces, a more general mathematical tool is needed: namely, actegories (see the survey [CG22] for a thorough treatment of the subject). Actegories are actions of symmetric monoidal categories on other categories. For brevity's sake, we will only give an incomplete definition (see [CG22] or [Gav24] for further information).

Definition 1 (Actegory). *Let $(\mathcal{M}, I, \otimes)$ be a strict symmetric monoidal category. A \mathcal{M} -actegory is a tuple $(\mathcal{C}, \bullet, \eta, \mu)$, where \mathcal{C} is a category, $\bullet : \mathcal{M} \times \mathcal{C} \rightarrow \mathcal{C}$ is a functor, and η and μ are natural isomorphisms enforcing $I \bullet C \xrightarrow{\eta_C} C$ and $(M \bullet (N \bullet C)) \xrightarrow{\mu_{M,N,C}} (M \otimes N) \bullet C$. The isomorphisms η and μ must also satisfy coherence conditions. If η and μ are identical transformations, we say that the actegory is strict.*

Remark 2. Although the requirement for strictness is somewhat restrictive, we will proceed under the assumption that the actegories we encounter are strict to streamline notation.

We will also be interested in actegories that interact with the monoidal structure of the underlying category.

Definition 3 (Monoidal actegory). *Let $(\mathcal{M}, I, \otimes)$ be a strict symmetric monoidal category and let $(\mathcal{C}, \bullet, \eta, \mu)$ be a strict actegory. Suppose \mathcal{C} has a monoidal structure (J, \boxtimes) . Then we say that (\mathcal{C}, \bullet) is monoidal if the underlying functor \bullet is monoidal.*

We may also be interested in studying the interaction between actegorical structures and endofunctors. This interaction can happen owing to a natural transformation known as strength. We will not provide coherence diagrams in the definition below for the sake of brevity, but the interested reader can find more detail in [GLD⁺24]. The paper also provides a definition of actegorical strong monad, which is a very similar concept.

Definition 4 (Actegorical strong functor). *Let (\mathcal{C}, \bullet) be an \mathcal{M} -actegory. A strong actegorical endofunctor on (\mathcal{C}, \bullet) is a pair (F, σ) where $F : \mathcal{C} \rightarrow \mathcal{C}$ is an endofunctor and σ is a natural transformation with components $\sigma_{P,A} : P \bullet F(A) \rightarrow F(P \bullet A)$ which satisfies a few coherence conditions that we do not list here.*

1.1.2 The Para construction

Suppose we have an \mathcal{M} -actegory (\mathcal{C}, \bullet) . We wish to study maps in \mathcal{C} which are parametrized using objects of \mathcal{M} , that is, maps in the form $P \bullet A \rightarrow B$. We are not just interested in the maps by themselves, but also in their compositional structure. Thus, we abstract away the details by defining a new category $\mathbf{Para}_\bullet(\mathcal{C})$ (first introduced in simplified form in [FST19]). Since we also want to formalize the role of reparametrization, we actually construct $\mathbf{Para}_\bullet(\mathcal{C})$ as a bicategory, so that its 0-cells A can serve as input/output spaces, its 1-cells (P, f) can serve as parametric maps, and, finally, its 2-cells r can serve as reparametrizations.

Definition 5 ($\mathbf{Para}_\bullet(\mathcal{C})$). *Let (\mathcal{C}, \bullet) be an \mathcal{M} -actegory. Then, we define $\mathbf{Para}_\bullet(\mathcal{C})$ as the bicategory whose components are as follows.*

- The 0-cells are the objects of \mathcal{C} .

- The 1-cells are pairs $(P, f) : A \rightarrow B$, where $P : \mathcal{C}$ and $f : P \bullet A \rightarrow B$.
- The 2-cells come in the form $r : (P, f) \Rightarrow (Q, g)$, where $r : P \rightarrow Q$ is a morphism in \mathcal{C} . r must also satisfy a naturality condition.
- The 1-cell composition law is

$$(P, f) \circ (Q, g) = (Q \otimes P, (Q \bullet f) \circ g).$$

- The horizontal and vertical 2-cell composition laws are respectively given by parallel and sequential composition in \mathcal{M} .

It is quite handy to represent the cells of $\mathbf{Para}_\bullet(\mathcal{C})$ using the string diagram notation illustrated in Fig. 1.1. The **Para** construction has a dual **coPara** construction whose 1-cells $f : \mathbf{coPara}_\bullet(\mathcal{C})(A, B)$ take the form (P, f) , where $f : A \rightarrow P \bullet B$. Cells in $\mathbf{coPara}_\bullet(\mathcal{C})$ can also be represented with appropriate string diagrams. The reader can find a complete definition in [Gav24].

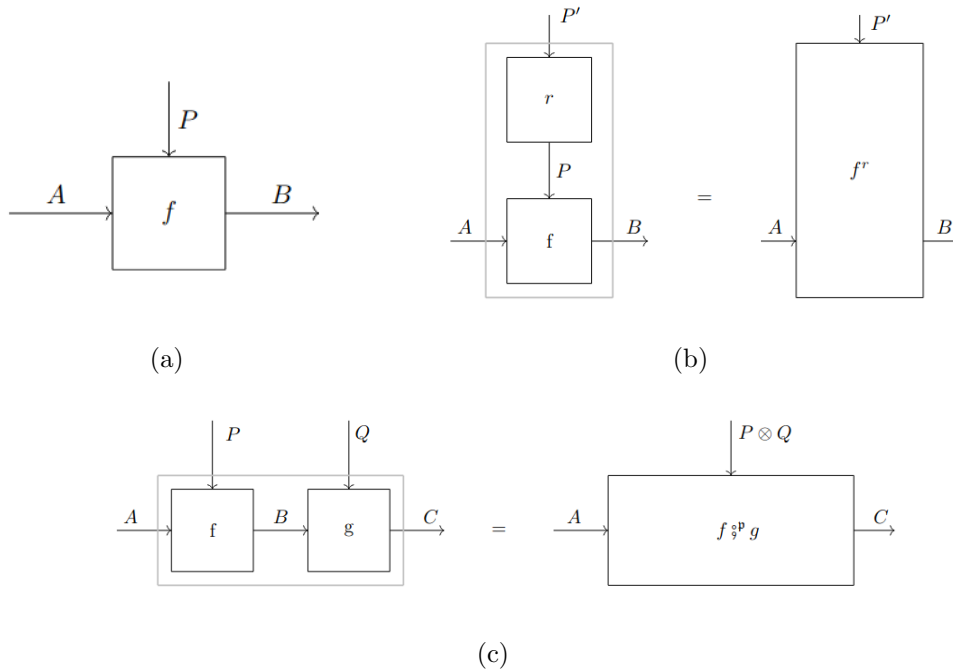


Figure 1.1: String diagrams representing (a) a parametric morphism, (b) a reparametrization of a parametric morphism, (c) a composition of parametric morphisms. (Images taken from [Gav24].)

It is shown in [Gav24] that $\mathbf{Para}_\bullet(\mathcal{C})$ is actually a 2-category if the underlying actegory is strict. Assuming this is the case (as we do in this thesis), we can use a functor $F : \mathbf{Cat} \rightarrow \mathbf{Set}$ to quotient out the 2-categorical structure and turn $\mathbf{Para}_\bullet(\mathcal{C})$ into a 1-category $F_*(\mathbf{Para}_\bullet(\mathcal{C}))$. Here, $F_* : \mathbf{2Cat} \rightarrow \mathbf{Set}$ is the change of enrichment basis functor induced by F . This meaningfully recovers the 1-categorical perspective of [FST19].

Both $\mathbf{Para}_\bullet(\mathcal{C})$ and $\mathbf{coPara}_\bullet(\mathcal{C})$ can be given a monoidal structure if (\mathcal{C}, \bullet) is a monoidal actegory. This is extremely important because it allows us to compose (co)parametric morphisms both in sequence and in parallel. Once again, more detail can be found in [Gav24].

Remark 6. Another way to parametrize morphisms is the coKleinsli construction. As noted by [Gav24], the main difference between **coKl** and **Para** is that the parametrization offered by **coKl**

is global, while the parametrization offered by **Para** is local: all morphisms in $\mathbf{coKl}(X \times -)$ must take a parameter in X , while the parameter space of different morphisms of $\mathbf{Para}(\mathcal{C})$ admit different parameter spaces. Nevertheless, the two constructions are related, and the former can be embedded into the latter.

If we take a parametrized category $\mathbf{Para}_\bullet(\mathcal{C})$ and we restrict our attention to morphisms parametrized with the monoidal identity I , we get back the original category \mathcal{C} . This is expressed by the following proposition ([Gav24]).

Proposition 7. *Let (\mathcal{C}, \bullet) be an \mathcal{M} -actegory. Then, there exists an identity-on-objects pseudofunctor $\gamma : \mathcal{C} \rightarrow \mathbf{Para}_\bullet(\mathcal{C})$ that maps $f \mapsto (I, f)$. If \mathcal{M} is strict, this is a 2-functor.*

1.1.3 Optics

Modelling bidirectional flows of information is not only useful in machine learning, but also in game theory, database theory, and more. As such, categorical tools for bidirectionality have been sought after for a long time: in particular, the greatest deal of efforts has been devoted to developing lens theory. Lenses have then been generalized into optics (see e.g. [Ril18]) to subsume other tools such as prisms and traversals into the same framework. Finally, there have also been various attempts to generalize optics (see e.g. [CEG⁺24] for a definition of mixed optics). We will introduce lenses and optics, and focus on the generalization of optics that appears (to us) to be the most versatile: weighted optics (first introduced in [Gav24]).

As stated in [Gav24], there is no standard definition of lens, and different authors opt for different *ad hoc* definitions that best suit their purposes. We will borrow the perspective of [CGG⁺22] and give the following definition.

Definition 8 (Lenses). *Let \mathcal{C} be a Cartesian category. Then, $\mathbf{Lens}(\mathcal{C})$ is the category defined by the following data:*

- an object of $\mathbf{Lens}(\mathcal{C})$ is a pair $\begin{pmatrix} A \\ A' \end{pmatrix}$ of objects in \mathcal{C} ;
- a $\begin{pmatrix} A \\ A' \end{pmatrix} \rightarrow \begin{pmatrix} B \\ B' \end{pmatrix}$ morphism (or lens) is a pair $\begin{pmatrix} f \\ f' \end{pmatrix}$ of morphisms of \mathcal{C} such that $f : A \rightarrow B$ and $f' : A \times B' \rightarrow A'$. f is known as the forward pass of the lens $\begin{pmatrix} f \\ f' \end{pmatrix}$, whereas f' is known as the backward pass;
- given $\begin{pmatrix} A \\ A' \end{pmatrix} : \mathbf{Lens}(\mathcal{C})$, the associated identity lens is $\begin{pmatrix} 1_A \\ \pi_1 \end{pmatrix}$;
- the composition of $\begin{pmatrix} f \\ f' \end{pmatrix} : \begin{pmatrix} A \\ A' \end{pmatrix} \rightarrow \begin{pmatrix} B \\ B' \end{pmatrix}$ and $\begin{pmatrix} g \\ g' \end{pmatrix} : \begin{pmatrix} B \\ B' \end{pmatrix} \rightarrow \begin{pmatrix} C \\ C' \end{pmatrix}$ is

$$\begin{pmatrix} f \circ g \\ \langle \pi_0, \langle \pi_0 \circ f, \pi_1 \rangle \circ g' \rangle \circ f' \end{pmatrix}.$$

Lenses are best thought of in the helpful language of the string diagrams illustrated in Fig. 1.2.

Lenses are a powerful tool, but they cannot be used to model all situations: for instance, lenses cannot be used if we wish to be able to choose not to interact with the environment depending on the input, or if we wish like to reuse values computed in the forward pass for further computation in the backward pass.

Optics generalize lenses by weakening the link between forward and backward passes, and by replacing the Cartesian structure of the underlying category with a simpler symmetric monoidal structure. In an optic over \mathcal{C} , an object $M : \mathcal{C}$ acts as an inaccessible residual space transferring information

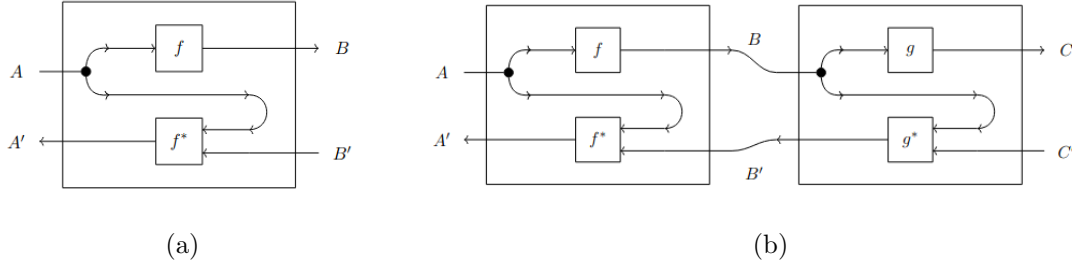


Figure 1.2: String diagrams representing (a) a lens $\left(\begin{smallmatrix} f \\ f^* \end{smallmatrix}\right)$, (b) the composition of two lenses $\left(\begin{smallmatrix} f \\ f^* \end{smallmatrix}\right)$ and $\left(\begin{smallmatrix} g \\ g^* \end{smallmatrix}\right)$. (Images taken from [CGG⁺22].)

between the upper components and the lower component. We will provide the definition given by [Ril18]¹, but we will use the string diagram notation presented by [Gav24] for the sake of consistency.

Definition 9 (Optics). *Let $(\mathcal{C}, I, \otimes, \lambda, \rho, \alpha)$ be a symmetric monoidal category (we make the unitors and associators explicit for later use). Then, $\mathbf{Optic}(\mathcal{C})$ is the category defined by the following data:*

- *an object of $\mathbf{Optic}(\mathcal{C})$ is a pair $\left(\begin{smallmatrix} A \\ A' \end{smallmatrix}\right)$ of objects in \mathcal{C} ;*
- *a $\left(\begin{smallmatrix} A \\ A' \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} B \\ B' \end{smallmatrix}\right)$ morphism (or optic) is a pair $\left(\begin{smallmatrix} f \\ f' \end{smallmatrix}\right)$ of morphisms of \mathcal{C} such that $f : A \rightarrow M \otimes B$ and $f' : M \otimes B' \rightarrow A'$, where $M : \mathcal{C}$ is known as residual space; such pairs $\left(\begin{smallmatrix} f \\ f' \end{smallmatrix}\right)$ are also quotiented by an equivalence relation that allows for reparametrization of the residual space and effectively makes it inaccessible;*
- *the identity on $\left(\begin{smallmatrix} A \\ A' \end{smallmatrix}\right)$ is the optic represented by $(\lambda_A^{-1}, \lambda_A)$.*

See [Ril18] for the more information about how optics compose and about the associated string diagrams.

Lenses come up as a special case of optics ([Ril18]), and optics do solve some of the issues we have with lenses. However, optics are not perfect either: for instance, [Gav24] points out that optics cannot be used in cases where we ask that the forward pass and backward pass are different kind of maps, as they are both forced to live in the same category. Thus, we need a further layer of generalization: namely, weighted optics.

1.1.4 Weighted optics

Before we define weighted optics, we need to introduce a new tool to our toolbox: the category of elements of a functor.

Definition 10 (Elements of a functor). *Let $F : \mathcal{C} \rightarrow \mathbf{Set}$ be a functor. We define $\mathbf{El}(F)$ as the category with the following data: (i) the objects of $\mathbf{El}(F)$ are pairs (C, x) where $C : \mathcal{C}$ and $x : F(C)$; (ii) the $(C, x) \rightarrow (D, y)$ morphisms in $\mathbf{El}(F)$ are the morphisms $f : C \rightarrow D$ in \mathcal{C} such that $F(f)(x) = y$.*

[Gav24] studies \mathcal{B} -actegories (\mathcal{C}, \bullet) , which are then reparametrized so that the acting category becomes $\mathcal{E} = \mathbf{El}(W)$ for some weight functor $W : \mathcal{B}^{\text{op}} \rightarrow \mathbf{Set}$ (which is to be specified). The

¹[Ril18]also provides a more versatile (but more sophisticated) definition of optics that relies on coends. Under the coend formalism,

$$\mathbf{Optic}(\mathcal{C})\left(\left(\begin{smallmatrix} A \\ A' \end{smallmatrix}\right), \left(\begin{smallmatrix} B \\ B' \end{smallmatrix}\right)\right) = \int^{M:\mathcal{C}} \mathcal{C}(A, M \times B) \times \mathcal{C}(M \otimes B', A').$$

reparametrization takes place thanks to the opposite of the forgetful functor $\pi_W : \mathcal{E} \rightarrow \mathcal{B}^{\text{op}}$, which maps $(B, x) \mapsto B$. Hence, we consider the action

$$\bullet^{\pi_W^{\text{op}}} = \mathbf{El}(W)^{\text{op}} \times \mathcal{C} \xrightarrow{\pi_W^{\text{op}} \times \mathcal{C}} \mathcal{B} \times \mathcal{C} \xrightarrow{\bullet} \mathcal{C}.$$

We are finally ready to define weighted optics².

Definition 11 (Weighted **coPara**). *If W is a weight functor as above and (\mathcal{C}, \bullet) is a \mathcal{B} -actegory, we define*

$$\mathbf{coPara}_{\bullet}^W(\mathcal{C}) = \pi_{0*}(\mathbf{coPara}_{\bullet^{\pi_W^{\text{op}}}}(\mathcal{C})),$$

where π_{0*} is the enrichment base change functor generated by the connected component functor $\pi_0 : \mathbf{Cat} \rightarrow \mathbf{Set}$. More explicitly, π_{0*} quotients the connections provided by reparametrizations.

Definition 12 (Weighted optics). *Suppose (\mathcal{C}, \bullet) is an \mathcal{M} -actegory, suppose $(\mathcal{D}, \blacksquare)$ is an \mathcal{M}' -actegory, and suppose $W : \mathcal{M}^{\text{op}} \times \mathcal{M}' \rightarrow \mathbf{Set}$ is a lax monoidal functor. We define the category of W -weighted optics over the product actegory $(\mathcal{C} \times \mathcal{D}^{\text{op}}, (\bullet_{\text{op}}))$ as*

$$\mathbf{Optic}_{(\bullet_{\text{op}})}^W = \mathbf{coPara}_{(\bullet_{\text{op}})}^W(\mathcal{C} \times \mathcal{D}^{\text{op}}).$$

The definition is very dense and deserves some explanation. First of all, we assume that W maps (M, M') to a set of maps $s : M \rightarrow M'$. If that's the case, a $\left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} Y \\ Y' \end{smallmatrix}\right)$ map is a triplet $\left(\left(\begin{smallmatrix} M' \\ M' \end{smallmatrix}\right), s, \left(\begin{smallmatrix} f \\ f' \end{smallmatrix}\right)\right)$, where M is the forward residual, M' is the backward residual, $s : M \rightarrow M'$ links the two residuals, $f : X \rightarrow M \bullet Y$ is the forward pass, and $f' : M' \bullet Y' \rightarrow X'$ is the backward pass. The triplets are also quotiented with respect to reparametrization, which makes the residual spaces effectively inaccessible (as it happens in the case of ordinary optics). We can get a clear "operational" understanding of how a weighted optic works looking at an associated string diagram: data from X flows through the forward map, which computes an output in Y and a forward residual in M . Such forward residual is then converted into a backward residual in M' by the map s , which is provided by the weight functor. Finally, the backward residual is used to compute, together with input from Y' a value in X' . The last step happens thanks to the backward map f' . A full account of the composition law for weighted optics can be found on [Gav24]. As stated in [Gav24], since **coPara** can be given a monoidal structure, we can also give **Optic** _{(\bullet_{op})} ^{W} one such structure as long as the underlying actegories are monoidal and the weight functor W is braided monoidal.

The advantages of weighted optics over ordinary optics are clear: when dealing with weighted optics, we are no longer forced to take reverse maps from the same category as the forward maps. The action on the category of forward spaces is now separated from the action on the category of backward spaces, and the link between the two actions is provided by an external functor. Such modular approach provides a great deal of conceptual clarity and flexibility, more than regular optics or lenses can provide on their own. It is also shown in [Gav24] that weighted optics are indeed a generalization of optics. In particular, it is shown that the lenses in Def. 8 are the specialized weighted optics obtained when $\mathcal{C} = \mathcal{D}$ is Cartesian and the actegories are given by the Cartesian product. More generally, [Gav24] claims that - to the best of the author's knowledge - all definitions of lenses currently used in the literature are subsumed by the definition of weighted optics.

[Gav24] goes on to apply the **Para** construction onto weighted optics, obtaining parametric weighted optics, which are proposed as a full-featured model for deep learning. The author conjectures that "weighted optics provide a good denotational and operational semantics for differentiation". In its full, generality, this is still an unproven conjecture. However, restricting our attention to a special

²Weighted optics also admit a coend definition. Refer to [Gav24] for more information.

class \mathbf{Lens}_A of lenses with an additive backward passes yields a fully formal theory of structural back-propagation, which will be illustrated in the rest of the capter, after a short digression on differential categories.

1.1.5 Differential categories

Modelling gradient-based learning obviously requires a setting where differentiation can take place. Although it is tempting to directly employ smooth functions over Euclidean spaces, recent research has shown that there are tangible advantages in working with generalized differential combinators that extend the notion of derivative to polynomial circuits ([WZ22], [WZ21]), manifolds ([PVM⁺21]), complex spaces ([BQL21]), and so on. Thus, it makes sense to work with an abstract notion of derivative which can then be appropriately implemented depending on the requirements at hand.

One approach to this problem involves the explicit definition of two kinds of differential categories: Cartesian differential categories (first introduced in [BCS06]) and Cartesian reverse differential categories (first introduced by [CCG⁺19]). The former allow for forward differentiation, while the latter allow for reverse differentiation. We will omit the defining axioms for the sake of brevity, but the reader can find complete definitions in [CCG⁺19].

Definition 13 (Cartesian differential category). *A Cartesian differential category (CDC) \mathcal{C} is a Cartesian left-additive category where a differential combinator D is defined. Such differential combinator must take a morphism $f : A \rightarrow B$ and return a morphism $D[f] : A \times A \rightarrow B$, which is known as the derivative of f . The combinator D must satisfy a number of axioms.*

Definition 14 (Cartesian reverse differential category). *A Cartesian reverse differential category (CRDC) \mathcal{C} is a Cartesian left-additive category where a reverse differential combinator R is defined. Such reverse differential combinator must take a morphism $f : A \rightarrow B$ and return a morphism $R[f] : A \times B \rightarrow A$, which is known as the reverse derivative of f . The combinator R must satisfy a number of axioms.*

Example 15. Smooth is both a CDC and a CRDC. In fact, if \mathcal{J}_f is the Jacobian matrix of a smooth morphism f ,

$$D[f] : (x, v) \mapsto \mathcal{J}_f(x)v$$

and

$$R[f] : (x, y) \mapsto \mathcal{J}_f(x)^T y$$

induce well-defined combinators D and R . This is only a partial coincidence, as shown in [CCG⁺19] that CRDCs are always CDCs under a canonical choice of differential combinator. The converse, however, is generally false.

As it turns out, forward differentiation tends to be less efficient when dealing with neural networks that come up in practice ([Ell18]), so CDCs are not extremely useful when studying deep learning. CRDCs, on the other hand, have been applied to great success (see e.g. [CGG⁺22]). As shown in [WZ22], a large supply of CRDCs can be obtained by providing the generators of a finitely presented Cartesian left-additive category with associated reverse derivatives (as long as the choices of reverse derivative are consistent). Moreover, CRDCs have been recently generalized by [Gav24] to coalgebras associated with copointed endofunctors, which could also increase the number of known CRDCs in the future. The rest of this section is devoted to this generalization.

It is shown in [Gav24] that there is a particular class of weighted optics which is useful for reverse differentiation, being able to represent both maps (through forward passes) and the associated reverse derivatives (through backward passes). Moreover, such weighted optics can be represented as lenses

in the sense of Def 8, which means that their inner workings can be pictured in a simple, intuitive way.

Definition 16 (Additively closed Cartesian left-additive category). *A Cartesian left-additive category \mathcal{C} is an additively closed Cartesian left-additive category (ACCLAC) if and only if the following are true:*

- the subcategory $\mathbf{CMon}(\mathcal{C})$ of additive maps has a closed monoidal structure (I, \otimes) ;
- the embedding $\iota : \mathbf{CMon}(\mathcal{C}) \rightarrow \mathcal{C}$ is a lax monoidal functor with respect to the aforementioned structure of $\mathbf{CMon}(\mathcal{C})$ and the Cartesian structure of \mathcal{C} .

Then, we can define the category of lenses with backward passes additive in the second component.

Definition 17. *Let \mathcal{C} be an ACCLAC with Cartesian structure is $(1, \times)$ and whose subcategory $\mathbf{CMon}(\mathcal{C})$ has monoidal structure (I, \otimes) . Then, we define*

$$\mathbf{Lens}_A(\mathcal{C}) = \mathbf{Optic}_{\left(\begin{smallmatrix} \times \\ \otimes \end{smallmatrix}\right)}^{\mathcal{C}(-, \iota(-))}.$$

As argued in [Gav24], the symbol \mathbf{Lens}_A is justified because one such optic of type $\left(\begin{smallmatrix} X \\ Y' \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} Y \\ Y' \end{smallmatrix}\right)$ can be concretely represented as a lens with forward pass $f : \mathcal{C}(X, Y)$ and backward pass $f' : \mathcal{C}(X \times Y', X')$, which is the approach we illustrate in this thesis. Nevertheless, some potential expressiveness is lost when passing from weighted optic composition to concrete lens composition. In particular, if we operated with optics, we would be able to implement backpropagation without resorting to gradient checkpointing, which is not possible if we use lenses ([Gav24]).

The generalization mentioned above is possible because \mathbf{Lens}_A is an endofunctor.

Definition 18. *We defined \mathbf{CLACat} as the category whose objects are Cartesian left-additive categories and whose morphisms are Cartesian left-additive functors (see e.g. [BCS06]).*

Proposition 19. *If $\mathcal{C} : \mathbf{CLACat}$, then $\mathbf{Lens}_A(\mathcal{C}) : \mathbf{CLACat}$.*

Proof. The Cartesian structure on $\mathbf{Lens}_A(\mathcal{C})$ is given by $\left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right) \times \left(\begin{smallmatrix} Y \\ Y' \end{smallmatrix}\right) = \left(\begin{smallmatrix} X \times Y \\ X' \times Y' \end{smallmatrix}\right)$ and by the initial object $\left(\begin{smallmatrix} 1 \\ 1 \end{smallmatrix}\right)$. The monoidal structure on each $\left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right)$ is given by the unit $0_{\left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right)} = \left(\begin{smallmatrix} 0_A \\ 1_{X \times A'} \end{smallmatrix}\right) : \left(\begin{smallmatrix} 1 \\ 1 \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right)$ and by the multiplication $+\left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right) = \left(\begin{smallmatrix} +_A \\ \pi_2 \circ \Delta_{A'} \end{smallmatrix}\right) : \left(\begin{smallmatrix} X \times X \\ X' \times X' \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right)$. \square

Proposition 20. $\mathbf{Lens}_A : \mathbf{CLACat} \rightarrow \mathbf{CLACat}$ *is a functor.*

Proof. Given a Cartesian left-additive functor $F : \mathcal{C} \rightarrow \mathcal{D}$, we can define $\mathbf{Lens}_A(F)$ as the functor that maps $\left(\begin{smallmatrix} X \\ X' \end{smallmatrix}\right) \mapsto \left(\begin{smallmatrix} F(X) \\ F(X') \end{smallmatrix}\right)$ and maps $\left(\begin{smallmatrix} f \\ f' \end{smallmatrix}\right) \mapsto \left(\begin{smallmatrix} F(f) \\ \underline{f} \end{smallmatrix}\right)$, where $\underline{f} = F(X) \times F(Y') \xrightarrow{\cong} F(X \times Y') \xrightarrow{F(f')} F(X')$. It can be shown that $\mathbf{Lens}_A(F)$ is also Cartesian left-additive. \square

Proposition 21. \mathbf{Lens}_A *has a copointed structure*³.

Proof. It suffices to endow \mathbf{Lens}_A with the natural transformation ϵ whose components are the forgetful functors $\epsilon_{\mathcal{C}} : \mathbf{Lens}_A(\mathcal{C}) \rightarrow \mathcal{C}$ which strip away the backward passes. \square

Hence, [Gav24] defines generalized CRDCs as follows.

Definition 22 (Generalized Cartesian reverse differential category). *A generalized Cartesian reverse differential category is a coalgebra for the pointed endofunctor \mathbf{Lens}_A .*

³An endofunctor $F : \mathcal{C} \rightarrow \mathcal{C}$ is copointed if it is endowed with a natural transformation $\epsilon : F(\mathcal{C}) \Rightarrow \text{id}_{\mathcal{C}}$.

Explicitly, a colagebra for \mathbf{Lens}_A is a pair $(\mathcal{C}, \mathbf{R}_\mathcal{C})$ such that $\mathcal{C} : \mathbf{CLACat}$ and $\mathbf{R}_\mathcal{C} : \mathcal{C} \rightarrow \mathbf{Lens}_A(\mathcal{C})$ satisfies $\mathbf{R}_\mathcal{C} \circ \epsilon_\mathcal{C} = \text{id}_\mathcal{C}$. The intuition behind such definition is that $\mathbf{R}_\mathcal{C}$ should map $f \mapsto \left(\begin{smallmatrix} f \\ R[f] \end{smallmatrix} \right)$, where $R[f]$ is a generalized reverse derivative combinator. [Gav24] shows that such a definition of $\mathbf{R}_\mathcal{C}$ does indeed prove that ordinary CRDC fall into the definition of generalized CRDC.

1.1.6 Parametric lenses

We conclude this section discussing the relation between the **Para** construction and the \mathbf{Lens}_A endofunctions. [Gav24] shows that, under an appropriate definition, ‘morphisms of actegories induce morphisms of parametric bicategories’. As a consequence, it can be shown that, if $(\mathcal{C}, \mathbf{R}_\mathcal{C})$ is a generalized CRDC, $\mathbf{R}_\mathcal{C}$ induces a functor $\mathbf{Para}(\mathbf{R}_\mathcal{C}) : \mathbf{Para}_\times(\mathcal{C}) \rightarrow \mathbf{Para}_\bullet(\mathbf{Lens}_A(\mathcal{C}))$, which takes a parametric map $f : P \times A \rightarrow B$ and augments it with its reverse derivative $R[f]$, forming a parametric lens. Parametric lenses behave very similarly to lenses, but we provide a separate stand-alone definition (which we take from [CGG⁺22]) for the reader’s convenience.

Definition 23 (Parametric lenses). *The category of parametric lenses over a Cartesian category $(\mathcal{C}, 1, \times)$ is $\mathbf{Para}_\bullet(\mathbf{Lens}(\mathcal{C}))$, where \bullet is the action on the lenses generated by the Cartesian structure of \mathcal{C} :*

$$\begin{pmatrix} P \\ P' \end{pmatrix} \bullet \begin{pmatrix} A \\ A' \end{pmatrix} = \begin{pmatrix} P \times A \\ P' \times A' \end{pmatrix}.$$

Refer to Fig. 1.3 to see a string diagram that shows the inner workings of a parametric lens.

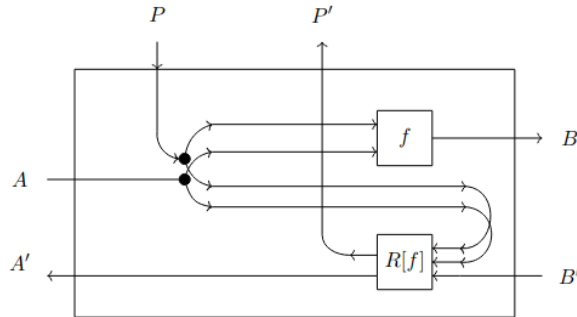


Figure 1.3: String diagrams representing the inner workings of a parametric lens. (Images taken from [CGG⁺22].)

1.2 Supervised learning with parametric lenses

In this section, we show how parametric lenses can be used to model supervised gradient-based learning ([CGG⁺22], [Gav24], [SGW21]). While lenses are not as general as weighted optics, it is shown in [CGG⁺22] that they are powerful enough for most purposes and that there is some empirical evidence of their performance and applicability. The paper also discusses the use of parametric lenses in modeling unsupervised deep learning and deep dreaming, but we do not have the space to discuss them here.

1.2.1 Model, loss, optimizer, learning rate

Supervised gradient-based learning can be modeled using parametric lenses as follows:

1. we can design an architecture (P, Model) as a parametric morphism in $\mathbf{Para}_\bullet(\mathcal{C})$ for some generalized CRDC $(\mathcal{C}, \mathbf{R}_\mathcal{C})$;
2. we can use the functor $\mathbf{R}_\mathcal{C}$ to endow (P, Model) with its reverse derivative $\mathbf{R}[(P, \text{Model})]$, yielding a lens in $\mathbf{Para}_\bullet(\mathbf{Lens}_A(\mathcal{C}))$;
3. we can use 2-categorical machinery of $\mathbf{Para}_\bullet(\mathbf{Lens}_A(\mathcal{C}))$ to provide a loss function, a learning rate, and optimizer, which can be assembled onto $\mathbf{R}_\mathcal{C}(P, \text{Model})$ to yield a parameter update;
4. we can use copy maps from the Cartesian structure of \mathcal{C} to create a learning iteration.

The theory of parametric optics and differential categories does not offer explicit insight with respect to architecture design, so we will assume a good architecture as already been designed (at least for now). Given an architecture (P, Model) , it can be embedded into $\mathbf{Para}_\bullet(\mathbf{Lens}_A(\mathcal{C}))$ as a lens $\left(\begin{smallmatrix} P \times A \\ P \times A \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} B \\ B \end{smallmatrix}\right)$ by breaking it up into its basic components (such as linear layers, convolutional layers, etc.), augmenting such components with their reverse derivatives, and the composing the resulting lenses. The backward pass of the composition is the reverse derivative of its forward pass because $\mathbf{R}_\mathcal{C}$ is a functor⁴. Many examples can be found in [CGG⁺22].

Updating the parameters based on data requires a loss function, an optimizer and a learning rate. Loss functions can be implemented as parametric lenses which take in predictions as input and labels as parameters. The output they produce can be considered the actual loss that needs to be differentiated. Given a model parametric lens $\left(\begin{smallmatrix} \text{Model} \\ \mathbf{R}[\text{Model}] \end{smallmatrix}\right) : \left(\begin{smallmatrix} P \times A \\ P \times A \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} B \\ B \end{smallmatrix}\right)$ and a loss parametric lens $\left(\begin{smallmatrix} \text{Loss} \\ \mathbf{R}[\text{Loss}] \end{smallmatrix}\right) : \left(\begin{smallmatrix} B \times B \\ B \times B \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} L \\ L \end{smallmatrix}\right)$, the composition $\left(\begin{smallmatrix} \text{Model} \\ \mathbf{R}[\text{Model}] \end{smallmatrix}\right) \circ \left(\begin{smallmatrix} \text{Loss} \\ \mathbf{R}[\text{Loss}] \end{smallmatrix}\right)$ takes in features as input and takes model parameters and labels as parameters. Then, this information is used to compute the loss associated with the model predictions. See *Fig. 1.4 (a)* for the associated string diagram.

It can be helpful to think about dangling wires in the diagrams as open slots where other components can be plugged. For instance, the diagram of *Fig. 1.4 (a)* has dangling wires labeled with L on its right. We can use a learning rate lens α to link these wires and allow forward-propagating information to "change direction" and go backwards. α must have domain equal to $\left(\begin{smallmatrix} L \\ L \end{smallmatrix}\right)$ and codomain equal to $\left(\begin{smallmatrix} 1 \\ 1 \end{smallmatrix}\right)$, where 1 is the terminal object of \mathcal{C} . For instance, if $\mathcal{C} = \mathbf{Smooth}$, α might just multiply the loss by some ϵ , which is what machine learning practitioners would ordinarily call learning rate. *Fig. 1.4 (b)* shows how a learning rate can be linked to the loss function and the model using post-composition.

The final element needed for the model Model in *Fig. 1.4 (b)* to learn is an optimizer. It is shown in [CGG⁺22] that optimizers can be represented as reparametrisations in $\mathbf{Para}(\mathbf{Lens}(\mathcal{C}))$. More specifically, we might see an optimizer as a lens $\left(\begin{smallmatrix} P \\ P \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} Q \\ Q \end{smallmatrix}\right)$. In gradient descent, for example, $P = Q$ and the aforementioned lens is $\left(\begin{smallmatrix} 1_P \\ +_P \end{smallmatrix}\right)$. We can plug such reparametrisation on top of the model to obtain the string diagram in *Fig. 1.4 (c)*. The diagram shows how the machinery hidden by the $\mathbf{Para}(\mathbf{Lens}(\mathcal{C}))$ can take care of forward propagation, loss computation, backpropagation and parameter updating in a seamless fashion.

In the supervised learning lens in *Fig. 1.4 (c)*, there are actually has two more dangling wires we do not have any use for: namely, the second input A of the model and the second parameter B of the loss functions. These are not a problem as we can just plug them with delete maps taken from the Cartesian structure of \mathcal{C} . We can also bend the input wires with an appropriate lens to turn them into parameters. We are then left with a $\left(\begin{smallmatrix} 1 \\ 1 \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} 1 \\ 1 \end{smallmatrix}\right)$ parametric lens with parameter space $\left(\begin{smallmatrix} A \\ 1 \end{smallmatrix}\right) \times \left(\begin{smallmatrix} P \\ 1 \end{smallmatrix}\right) \times \left(\begin{smallmatrix} B \\ 1 \end{smallmatrix}\right)$. This lens is pictured in *Fig. 1.4 (d)*.

⁴As highlighted by [SGW21], the diagram for the backward pass of the composition of two lenses looks exactly like the diagram describing the chain rule for reverse derivatives, which is what makes $\mathbf{R}_\mathcal{C}$ a well-defined functor.

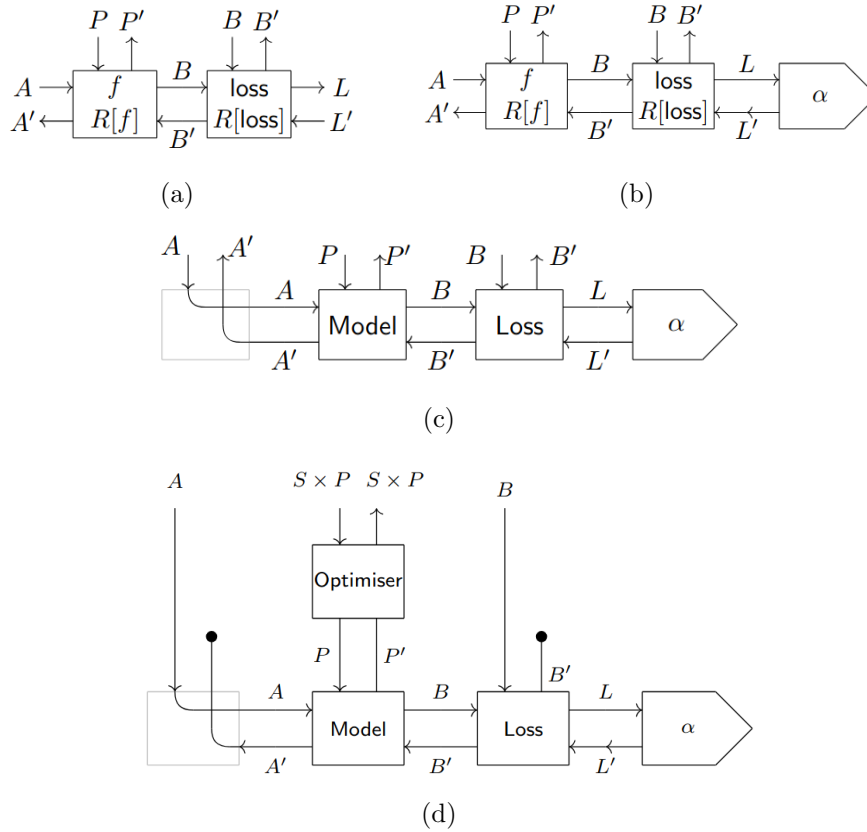


Figure 1.4: String diagrams representing (a) the composition of a model lens and an optimizer lens (b) the composition of a model lens, an optimizer lens, and a learning rate lens, (c) a composition of lenses representing a complete supervised learning algorithm, (d) the same lens of the previous point with useless wires plugged. (Images taken from [CGG⁺22].)

1.2.2 Weight tying, batching, and the learning iteration

Both [CGG⁺22] and [Gav24] emphasize the essential role played by weight tying in deep learning. Weight tying can be implemented within the parametric lens framework as a reparametrization that copies a single parameter to many parameter slots (see Fig. 1.5 (a)): given $(P \times P, f) : \mathbf{Para}(\mathcal{C})(X, Y)$, we can define $(P, f^{\Delta_P}) : \mathbf{Para}(\mathcal{C})(X, Y)$ so that

$$f^{\Delta_P} : P \times X \xrightarrow{\Delta_{P \times X}} P \times P \times X \xrightarrow{f} Y.$$

Weight tying can also be used for batching: batching is implemented by instantiating n different copies of our supervised learning lens (comprised of model, loss function, and learning rate) and tying the parameters to a unique value. Then, it suffices to feed the n data points to the n lenses, and we can optimize across a single parameter (see Fig. 1.5 (b)).

[CGG⁺22] introduces a possible representation for the whole learning iteration of a supervised learning model within the parametric lens framework. The paper suggests extracting the backward pass of the lens in Fig. 1.4 (d) and reframing it as a $P \rightarrow P$ parametric map with parameters $A \times B$. Since this is an endomap, it can be composed n times with itself to obtain a $P \rightarrow P$ map, which is proposed as a model of the learning iteration. While this approach requires breaking lenses apart, it is markedly simple.

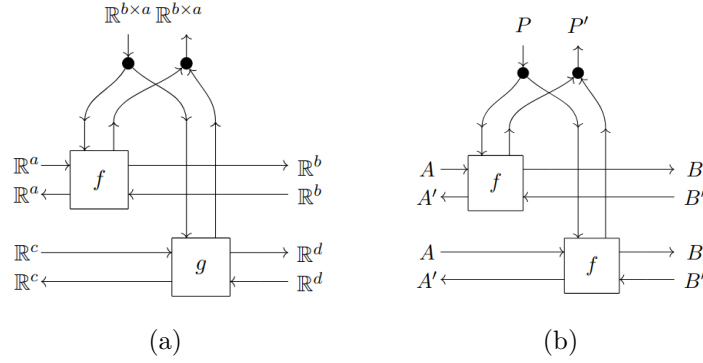


Figure 1.5: String diagrams representing (a) weight trying and (b) batching, both implemented using the copy maps of a Cartesian category. (Images taken from [CGG⁺22].)

1.2.3 Empirical evidence

Empirical evidence for the effectiveness of the parametric lens framework discussed in this section can be found in [CGG⁺22], where the authors go on to implement a Python libraries for gradient-based learning based on these ideas. They use the library to develop a MNIST classifier, obtaining comparable accuracy to models developed using traditional tools.

Since parametric lenses can easily be implemented functionally and without side effects, success stories as the one mentioned above foreshadow a future where popular machine learning libraries also follow elegant functional paradigms informed by category theory. Quoting [CGG⁺22] directly, ‘[the] proposed algebraic structures naturally guide programming practice’.

1.3 Future directions and related work

The parametric optic framework discussed in this chapter is very promising, but there is still a lot of work that needs to be done so that it can reach its full potential. For instance, [Gav24] conjectures that weighted optics can be used in its full generality to model differentiation in cases which are not covered by lenses. For instance, lenses cannot model AD algorithms that do not use gradient checkpointing, while weighted optics are conjectured to be able to do so. [Gav24] also suggests investigating locally graded categories as potential replacements for actegories, and also investigating the applications of parametric optics to meta-learning, that is deep learning where the optimizers themselves are learned. Moreover, [CGG⁺22] conjectures that some of the axioms of CRDC may be used to model higher order optimization algorithms. Finally, as suggested by [CGG⁺22], hopefully future work will allow the parametric optic framework to encompass non-gradient based optimizers such as the ones used in probabilistic learning. See [SGW21] for more on this topic.

We conclude this chapter by discussing three other directions of machine learning research that are closely related to the framework of parametric optics.

1.3.1 Learners

One of the first compositional approaches to training neural networks in the literature can be found in the seminal paper [FST19], which spurred much research in the field, including what is presented in [Gav24] and [CGG⁺22]. The authors introduce a category of learners, objects which are meant to represent components of a neural network and behave similarly to parametric lenses.

Definition 24 (Category of learners). *Let A and B be sets. A learner $A \rightarrow B$ is a tuple (P, I, U, r) where P is a set, and $I : P \times A \rightarrow B$, $U : P \times A \times B \rightarrow P$, and $r : P \times A \times B \rightarrow A$ are functions. P is known as parameter space, I as implement functions, U as update function, and r as request function. Two learners $(P, I, U, r) : A \rightarrow B$ and $(Q, J, V, s) : B \rightarrow C$ compose forming $(P \times Q, I * J, U * V, r * s) : A \rightarrow C$, where*

$$(I * J)(p, q, a) = J(q, I(p, a)),$$

$$(U * V)(p, q, a, c) = (U(p, a, s(q, I(p, a), c)), V(q, I(p, a), c)),$$

$$(r * s)(p, q, a, c) = r(p, a, s(q, I(p, a), c)).$$

*Learners quotiented by an appropriate reparametrization relationship⁵ form a category **Learn**.*

A learner represents an instance of supervised learning: the implement function takes a parameter and implements a function and the update function updates the parameters using a data from a dataset. The request function is necessary to implement backpropagation when optimizing a composition of learners. Suppose we select a learning rate ϵ and an error function $e : \mathbb{R}^2 \rightarrow \mathbb{R}$ such that $y \mapsto \frac{\partial e}{\partial x}(x_0, y)$ is invertible for all y . It is argued in [FST19] that we can define a functor $L_{\epsilon, e} : \mathbf{Para}_\times(\mathbf{Smooth}) \rightarrow \mathbf{Learn}$ which takes a parametric map and yields an associated learner that implements gradient descent.

We do not have the space to talk about learners at length, but we wish to draw a short comparison between parametric weighted optics (and, in particular, parametric lenses) and the approach of [FST19], given the relevant position held by the paper in the machine learning literature. The similarities between learner-based learning and lens-based learning are evident: every learner (P, I, U, r) looks like a parametric lens, where I passes information forward, r passes information backwards and P is the parameter space. Moreover, the role of $L_{\epsilon, e}$ is very similar to the role played by $\mathbf{Para}(\mathbf{R}_\mathcal{C})$ in optic-based learning. Such similarities were even discussed in the original paper [FST19] and have been researched at length: it has been proved in [FJ19] that learners can be functorially and faithfully embedded in a special category of symmetric lenses (as opposed to the lenses of *Def. 8*, which are asymmetric).

Despite the similarities, there is one fundamental difference between the lens-based approach and the learner-based approach: each learner carries its own optimizer, whereas optimization of lenses is usually carried out separately. Moreover, if we compare parametric weighted optics with learners, the latter clearly win in versatility, generality, and (at least from our point of view) conceptual clarity. It is argued in [SGW21] and [CGG⁺22] that the parametric lens framework largely subsumes the learner approach. More information regarding the comparison can also be found in [Gav24].

1.3.2 Exotic differential categories

We have presented the parametric weighted optic approach of [Gav24] and [CGG⁺22] within the context of neural networks for the sake of simplicity, but the framework has been developed with generality in mind and applies to a much wider range of situations. For instance, we can easily replace **Smooth** with any other CRDC \mathcal{C} , yielding a full-feature compositional framework for gradient-based learning over \mathcal{C} .

Switching to a different CRDC is useful because different differential categories can lead to different learning outcomes, both in terms of accuracy of the model and in terms of training computational

⁵As argued in [FST19], learners could be studied from a bicategorical point of view, where reparametrizations would just be 2-cells. We could then use a connected component projection to compress **Learn** into a 1-category **Learn**, as it is done for **coPara** when defining weighted optics.

costs ([WZ22]). For instance, is argued in [WZ22] that polynomial circuits can be used to define and train intrinsically discrete machine learning models. Even ‘radical ’environments such as Boolean circuits - where scalars reside in \mathbb{Z}_2 - seem to be conducive to machine learning under the right choice of architecture and optimizer ([WZ21]). Using such exotic differential categories may be of great advantage because they might be able to better reflect the intrinsic computational limits of computer arithmetics, leading to more efficient learning ([WZ22]).

1.3.3 Functional reverse-mode automatic differentiation

Finally, we wish to highlight the similarities between the formal theory of differential categories illustrated here and the work in [Ell18]. The paper describes the Haskell implementation of a purely functional automatic differentiation library, which is able to handle both forward mode and backward mode AD without resorting to the mutable computational graphs used by most current day libraries.

Among the main insights of [Ell18], it is stated that derivatives should not be treated as simple vectors, but as linear maps, or multilinear maps in the case of uncurried higher-order derivatives. Moreover, the author shows that differentiation can be made compositional by working on pairs (f, Df) , which behaved very similarly to lenses. As noted by [SGW21], however, [CGG⁺22] and other lens-theoretical perspectives do not subsume the work in [Ell18] because of the latter’s programming focus. See [SGW21] for more information regarding this comparison.

Chapter 2

From Classical Computer Science to Neural Networks

Classical computer science focuses on discovering algorithms, that is ordered sequences of steps which operate in precisely set, idealized conditions and have strong guarantees of correctness due to their exact mathematical formulations. Neural networks, on the other hand, are able to work in messy, real-world conditions, but offer very so few guarantees of correctness that their performance is often described as *unreasonably* good. Moreover, whereas algorithms generalize very well (most software engineers will only need a few dozen algorithms in their entire career), neural networks are often completely helpless when pitted against out of distribution inputs. Hence, algorithms and neural networks can be seen as complementary opposites ([VB21], [VBB⁺22]).

Recent attempts going under the label of *neural algorithmic reasoning* (see [VB21] for a very short introduction to the subject) have tried to get the best of both worlds by training neural networks to execute algorithms (see e.g [IKP⁺22]). The CLRS benchmark (introduced by [VBB⁺22]) uses graphs to represent the computations associated with a few classical algorithms from the famous CLRS introductory textbook ([CLRS22]) so that graph neural networks (GNNs) can be trained to learn these algorithms. The benchmark has spurred a large amount of research in this direction, with very promising results.

More generally, linking machine learning to classical computer science may unlock interesting advances. For example, recovering neural networks as parametric versions of known algorithms may classify existing architectures in a conceptually clear manner and may even help develop new neural network architectures by taking inspiration from well-researched classical notions. In this chapter, we illustrate two lines of inquiry which use category theory to build such a bridge: *categorical deep learning* and an interesting categorical approach to *algorithmic alignment*. Before treating such topics, we will go on a short tangent categorical tangent regarding (co)algebras and the integral transform.

2.1 Categorical toolkit

2.1.1 (Co)algebras

Algebras and coalgebras are a categorical formalization of the principles of induction and coinduction. Induction and coinduction are fundamental to computer science because they allow us to give precise definitions for many data structures and to formalize recursive and corecursive algorithms on such structures. We will touch on (co)algebras very briefly but we refer interested readers to [JR97] and [Wis08] for further detail.

Definition 25 ((Co)algebra over an endofunctor). Let $F : \mathcal{C} \rightarrow \mathcal{C}$ be an endofunctor. An algebra over F is a pair (A, a) where $A : \mathcal{C}$ and $a : \mathcal{C}(F(A), A)$. A coalgebra is a pair (A, a) where $A : \mathcal{C}$ and $a : \mathcal{C}(A, F(A))$. In both cases A is known as carrier set and a as structure map.

(Co)algebras can also be defined on monads: the only difference between (co)algebras over an endofunctor and (co)algebras over a monad is that the latter also need to be compatible with the monad structure, that is satisfy commutative diagrams that represent coherence conditions (see [GLD⁺24]). (Co)algebras over the same functor can be given a categorical structure by using the following notion of homomorphism.

Definition 26 (Homomorphisms of (co)algebras over an endofunctor). Let (A, a) and (B, b) be algebras over the same endofunctor $F : \mathcal{C} \rightarrow \mathcal{C}$. An algebra homomorphism $(A, a) \rightarrow (B, b)$ is a map $f : \mathcal{C}(A, B)$ such that the diagram in Fig. 2.1 (a) is commutative.

Now suppose (A, a) and (B, b) are coalgebras. A homomorphism between them is a map $f : \mathcal{C}(A, B)$ such that the diagram in Fig. 2.1 (b) is commutative.

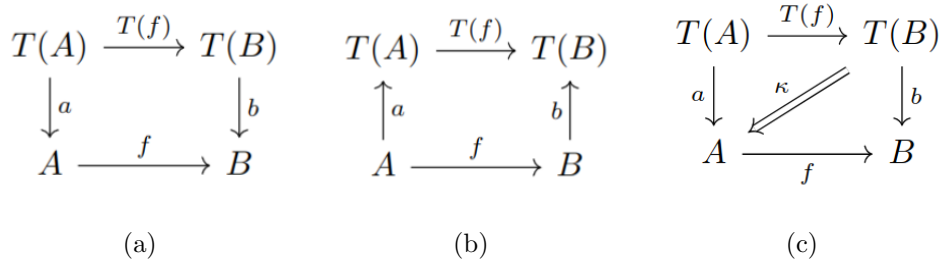


Figure 2.1: Conditions under which $f : (A, a) \rightarrow (B, b)$ is (a) an algebra homomorphism, (b) a coalgebra homomorphism, (c) a lax algebra homomorphism. (Image (c) is taken from [GLD⁺24].)

The main intuition behind the notions of algebra and coalgebra is the following: the underlying functor defines a signature for the (co)algebraic structure; the structure of an algebra is a constructor that takes data from $F(A)$ and uses it to build data from A , whereas the structure of a coalgebra observes data from A and produces an observation in the form of data from $F(A)$; (co)algebra homomorphisms are arrows that preserve the underlying structure. Consider the following clarifying examples from [GLD⁺24].

Remark 27. In the examples below we use polynomial and exponential expressions to define endofunctors over **Set**. In this context, X is the argument of the functor, \times is the Cartesian product, $+$ is the disjoint union, $\langle f, g \rangle$ is the pairing induced by \times , $[f, g]$ is the pairing induced by $+$, and B^A is the set of functions $A \rightarrow B$. The \times operator is assumed to take precedence over the $+$ operator. Similarly, the exponential operator is assumed to take precedence over the \times operator.

Example 28 (Lists). Let A be a set. Consider the endofunctor $1 + A \times X : \mathbf{Set} \rightarrow \mathbf{Set}$. If $\text{List}(A)$ is the set of A -labeled lists, $(\text{List}(A), [\text{Nil}, \text{Cons}])$ is an algebra over $1 + A \times X$. Here, $\text{Nil} : 1 \rightarrow \text{List}(A)$ is the map which takes the unique object of 1 and returns the empty list, while $\text{Cons} : A \times \text{List}(A) \rightarrow \text{List}(A)$ is the map which takes an element $a \in A$ and a list l of elements of A and returns the concatenated list $l \cup \{a\}$. The algebra $(\text{List}(A), [\text{Nil}, \text{Cons}])$ describes lists in $\text{List}(A)$ inductively as object formed by concatenating elements of A to other lists in $\text{List}(A)$. The base case is the empty list.

Example 29 (Mealy machines). Now consider two sets I and O of possible inputs and outputs, respectively. Consider the endofunctor $(O \times X)^I : \mathbf{Set} \rightarrow \mathbf{Set}$. Define $\text{Mealy}_{I,O}$ as the set of Mealy

machines with inputs and outputs in I and O , respectively. Now we can consider the coalgebra $(\text{Mealy}_{I,O}, \text{Next})$, where Next is the map that takes a Mealy machine $m \in \text{Mealy}_{I,O}$ and yields a function which in turn, given $i \in I$, returns the output of m at i and a new machine m' . This is a coinductive description of Mealy machines.

Remark 30. Notice how the description we have given of Mealy machines does not mention the internal states of these objects. This is a recurring aspect of coinductive descriptions: as argued in [JR97], coinduction is best interpreted as a process where an observer tracks the behavior of an object from the outside, with no access to its internal state. This is very useful in machine learning because the internal state of a learning model is often unknown or uninterpretable.

The link between (co)algebras and (co)induction does not stop at the definition level. The example below shows that an algebra homomorphism can model a recursive fold procedure. A similar corecursive unfold procedure can be defined by using a coalgebra homomorphism (see [GLD⁺24] for further detail).

Example 31 (List folds). Consider the algebra $(\text{List}(A), [\text{Nil}, \text{Cons}])$ of lists from *Ex. 28*, and consider a second algebra $(Z, [r_0, r_1])$ over the same functor. A homomorphism $f : \text{List}(A) \rightarrow Z$ from the former into the latter must satisfy

$$\begin{aligned} f(\text{Nil}) &= r_0, \\ f(\text{Cons}(a, l)) &= r_1(a, f(l)). \end{aligned}$$

Hence, f is necessarily a fold over a list with recursive components r_0 and r_1 . Incidentally, this proves that f is unique, making $(\text{List}(A), [\text{Nil}, \text{Cons}])$ an initial object in the category of algebras over the polynomial endofunctor $1 + A \times X$.

The notion of (co)algebra over a functor can be generalized to the sphere 2-categories, defining the notion of (co)algebra over a 2-endofunctor. The basic concepts stay the same but the commutativity of the diagrams defining (co)algebra homomorphisms is relaxed into lax-commutativity. A square diagram of 1-cells is lax-commutative if there exists a 2-cell that carries the composition of the top and the right edge onto the composition of the left and bottom edge, as in *Fig. 2.1 (c)*. Once again we refer to [GLD⁺24] for further information.

We conclude the section with the following proposition ([GLD⁺24]), which shows that actegorical strong endofunctors induce 2-endofunctors over parametric categories. We will use the proposition to study parametric versions of the endofunctors of *Ex. 28* and *Ex. 29*.

Proposition 32. *Suppose (\mathcal{C}, \bullet) is an \mathcal{M} -actegory and $F : \mathcal{C} \rightarrow \mathcal{C}$ is an actegorical endofunctor with strength σ . Then, F induces a 2-endofunctor $\mathbf{Para}(F) : \mathbf{Para}_\bullet(\mathcal{C}) \rightarrow \mathbf{Para}_\bullet(\mathcal{C})$.*

Proof. Define $\mathbf{Para}(F)$ so that:

1. $\mathbf{Para}(F)$ acts like F on objects $A : \mathcal{C}$;
2. $\mathbf{Para}(F)(f) = P \bullet F(A) \xrightarrow{\sigma_{P,A}} F(P \bullet A) \xrightarrow{F(f)} F(B)$ for all $(P, f) : \mathbf{Para}_\bullet(\mathcal{C})(A, B)$;
3. $\mathbf{Para}(F)$ leaves reparametrizations unchanged.

□

2.1.2 Integral transform

Remark 33. In accordance with the notation of [DV22] and [DvGPV24], we use $[A, B]$ to represent the set of $A \rightarrow B$ functions, where A and B are sets.

Suppose (R, \oplus, \otimes) is a commutative semiring. An integral transform is a transformation that carries a function in $[W, R]$ to a function in $[Z, R]$ following a precise chain of steps. Integral transforms¹ have been introduced by [DV22] to provide a single formalism able to describe both dynamic programming and GNNs. Integral transforms can be encoded as polynomial spans.

Definition 34 (Polynomial span). *A polynomial span is a triplet $(i : X \rightarrow W, p : X \rightarrow Y, o : Y \rightarrow Z)$ of morphisms in **FinSet**, the category of finite sets and functions. i is known as input, p as process, o as output. W is known as input set, X as argument set, Y as message set, Z as output set. We also ask that the fibers of p have total orderings². The polynomial span (i, p, o) can be graphically represented as the string diagram in Fig. 2.2 (a).*

Definition 35 (Integral transform). *Let (R, \oplus, \otimes) be a commutative semiring. Let $(i : X \rightarrow W, p : X \rightarrow Y, o : Y \rightarrow Z)$ be a polynomial span. The associated integral transform is the triplet $(i^* : [W, R] \rightarrow [X, R], p_\otimes : [X, R] \rightarrow [Y, R], o_\oplus : [Y, R] \rightarrow [Z, R])$, where:*

1. i^* is the pullback mapping $f \mapsto i \circ f$;
2. p_\otimes is the argument pushforward mapping

$$p_\otimes(a)(u) = \bigotimes_{e \in p^{-1}(u)} a(e);$$

3. o_\oplus is the message pushforward mapping

$$o_\oplus(m)(v) = \bigoplus_{e \in o^{-1}(v)} m(e).$$

The integral transform $(i^*, p_\otimes, o_\oplus)$ can be represented by the diagram in Fig. 2.2 (b).

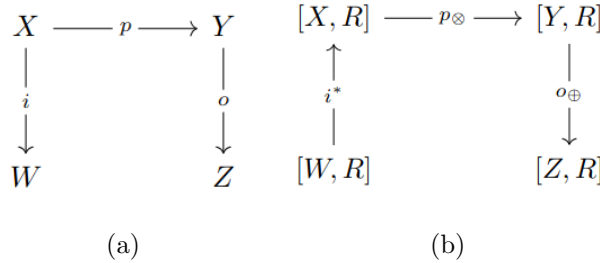


Figure 2.2: (a) A polynomial span and (b) the associated integral transform. (Images taken from [DV22].)

Remark 36. Whereas defining i^* is quite straight-forward, defining p_\otimes and o_\oplus is more difficult because the arrows p and o point in the wrong direction, which implies that the underlying functions must be inverted before considering the associated pullbacks. However, inverting non-invertible functions yields functions into the powersets of the original domains. Moreover, if we want to preserve the multiplicity of arguments and messages, we have to construct inverses that go into the sets of multisets over the original domains. Hence why we need \otimes and \oplus to aggregate results over such multisets. The significance these steps will be clarified later on in this chapter.

¹The label integral transform refers to the fact that similar ideas can be used to write categorical definitions for familiar analytical integral transforms ([Wil10]). A similar construct is also used in physics ([EPWJ80]).

²Neither we nor [DV22] use this requirement but, as stated in the original paper, the requirement is useful to support functions with non-commuting arguments.

2.2 Categorical deep learning

The optic-based framework we presented in the last chapter provides a structured general-purpose compositional framework for gradient-based learning, but its great versatility has a price: optics are unable to guide the architectural design of our models. It has been shown times and times again that a better architecture makes as much of a difference in machine learning as an algorithm with a better asymptotic cost does in classical computer science. Therefore, finding a principled mathematical framework able to guide such architectural choices is of paramount importance.

Designing an architecture is almost synonymous to imposing a set of constraints: for instance, convolutional layers are notoriously equivalent to translationally equivariant linear layers. Hence, it makes sense to inform the choice of architecture by informing the choice of constraint. Among the various attempts to develop a principled theory of machine learning constraints, *geometric deep learning* (see e.g. [BBCV21]), or GDL is particularly relevant. Categorical deep learning, or CDL, evolved as categorical generalization of GDL ([GLD⁺24]).

2.2.1 From GDL to CDL

Geometric deep learning focuses on equivariance constraints defined with respect to group actions, according to the following definition.

Definition 37 (Group action equivariance and invariance). *Let G be a group and let (S, \cdot) and $(T, *)$ be G -actions. A function $f : S \rightarrow T$ is equivariant with respect to the aforementioned actions if $f(g \cdot s) = g * f(s)$ for all $s \in S$ and for all $g \in G$. We say that f is invariant if $*$ is the trivial action on T , and thus $f(g \cdot s) = f(s)$ for all s and g .*

The GDL approach prescribes that a desired constraint should be expressed in terms of equivariance with respect to group action so that associated equations can be derived. Solving such equations usually implies tying weights, which reduces the net number of parameters in the models and provides the various advantages (see e.g. [OCV⁺20]) afforded by equivariance.

GLD has been successfully applied in contexts where data transformations can be expressed in group-theoretical terms (see e.g. [BBCV21]). Regrettably, this is not always the case, as many transformations we can subject the data to are either not invertible or even not compositional ([GLD⁺24]). Moreover, as highlighted in [GLD⁺24], although GDL is effective at describing the constraints that a model should implement, it is not always clear how such constraints can be actually be implemented since solving the equations associated with a constraint might be extremely challenging. Categorical deep learning, introduced by [GLD⁺24], attempts to solve the aforementioned problems using category theory³ to generalize the notion of equivariance and to tie specific constraints to specific architectures.

Remark 38. At the moment, to the best of our knowledge, [GLD⁺24] is the only publicly available paper that discusses the ideas of CDL.

The main insight of CDL is that group actions can be represented as algebras over the group action monads, and that maps that are equivariant with respect to these actions are homomorphisms between these algebras. Hence, GDL can be generalized by taking into consideration (co)algebras over other monads and endofunctors. CDL generalizes GDL into a theory of (co)algebras over endofunctors and monads, yielding a "theory of all architectures" ([GLD⁺24]). This prophecy has not yet been fulfilled,

³As stated in [BBCV21], GDL is inspired by the *Erlangen Programme*, which unified geometry around the notion of invariant at the end of the nineteenth century. Since category theory can be seen as an extension of the *Programme* (as we already remarked in the introduction to this thesis), it is only natural to attempt to generalize GDL by categorical means.

since only a few architectures have been studied within the framework of CDL at the moment of writing this thesis. However, the results shown in [GLD⁺24] are very promising and, hopefully, future work will widen the scope of CDL to include derivations for more architectures.

The following proposition and the subsequent example show how exactly CDL subsumes GDL.

Proposition 39. *Let (G, e, \cdot) be a group. The endofunctor $G \times X : \mathbf{Set} \rightarrow \mathbf{Set}$ can be given a monad structure using the natural transformations η , with components $\eta_S : s \mapsto (e, s)$, and μ , with components $(g, h, s) \mapsto (g \times h, s)$. The monad $(G \times -, \eta, \mu)$ can serve as a signature for G -actions. The actions themselves can be recovered by considering algebras $(S, *)$ for the monad, and, given two actions $(S, *)$ and (T, \star) , an associated equivariant map $f : S \rightarrow T$ is a $(S, *) \rightarrow (T, \star)$ monad algebra homomorphism.*

Proof. It suffices to compare the equations that define group actions and group action invariance with the commutative diagrams in Fig. 2.1. \square

Example 40 (Linear equivariant layer). Consider a carrier set $S = \mathbb{R}^{\mathbb{Z}_2}$, which can be seen as a pair of pixels. Consider the translation action $(i * s)(j) = s(i - j)$ of $G = \mathbb{Z}_2$ on S , which can be seen as swapping the pixels. We want to find a linear map $f : S \rightarrow S$ which is equivariant with respect to the action. Imposing the equivariance constraints as equations on the entries of the matricial representation $W_f \in \mathbb{R}^{2 \times 2}$ of the map, we can prove that f is equivariant if and only if W_f is symmetric ([GLD⁺24]).

2.2.2 (Co)inductive definitions for RNNs

As seen in Ex. 40, the formalism of CDL subsumes the formalism of GDL, but the difference between the two is not a simple matter of notation: CDL is much more general. The most significant piece of novel contribution delineated in [GLD⁺24] is the use of (co)algebras and (co)algebra homomorphisms over parametric categories to (co)inductively define recurrent neural networks (GNNs) and recursive neural networks (TreeRNNs). (Co)algebras are used to define cells, whereas the associated homomorphisms provide the weight-sharing mechanics used to unroll them. Let us build on Ex. 28 and Ex. 31, as is done in [GLD⁺24].

Example 41 (Folding recurrent neural network cell). Consider the endofunctor $1 + A \times X : \mathbf{Set} \rightarrow \mathbf{Set}$ from Ex. 28. Consider the Cartesian action of \mathbf{Set} on itself and associate the following actegorical strength to the functor: $\sigma_{P, X}(p, \text{inl}) = \text{inl}$ and $\sigma_{P, A}(p, \text{inr}(x, x')) = \text{inr}((p, x), (p, x'))$. Now that the functor is actegorical strong, we can use Prop. 32 to construct an endofunctor $\mathbf{Para}(1 + A \times X) : \mathbf{Para}_\bullet(\mathbf{Set}) \rightarrow \mathbf{Para}_\bullet(\mathbf{Set})$. Consider an algebra $(S, (P, \text{Cell}))$ for this functor. Via the isomorphism $P \times (1 + A \times X) \cong P + P \times A \times X$, we deduce that $\text{Cell} = [\text{Cell}_0, \text{Cell}_1]$, where $\text{Cell}_0 : P \rightarrow S$ and $\text{Cell}_1 : P \times A \times S \rightarrow S$. We can interpret Cell_0 and Cell_1 as folding recurrent neural network cells: Cell_0 provides the initial state based on its parameter and Cell_1 takes in the old state, a parameter, and an input, which are then used to return a new state (Fig. 2.3 (a)).

Example 42 (Unrolling of a folding recurrent neural network). Use Prop. 7 to embed the list algebra $(\text{List}(A), [\text{Nil}, \text{Cons}])$ from Ex. 28 as an algebra over the endofunctor $\mathbf{Para}(1 + A \times X)$ define in Def. 41. Now consider an algebra homomorphism $(P, f) : (\text{List}(A), [\text{Nil}, \text{Cons}]) \rightarrow (S, (P, \text{Cell}))$. Since we are working with algebras over a 2-endofunctor, we also need to specify a 2-cell that makes the homomorphism diagram (Fig. 2.1 (c)) lax-commutative. Using the weight-tying reparameterization Δ_P yields the lax commutative diagram in Fig. 2.4, which uniquely identifies f as the fold function which takes a list of inputs in A and unrolls a folding recurrent neural network that reads such inputs. The weight-tying reparameterization makes sure that each cell of the unrolled network uses the same parameters (see Fig. 2.3 (b) for a graphical representation).

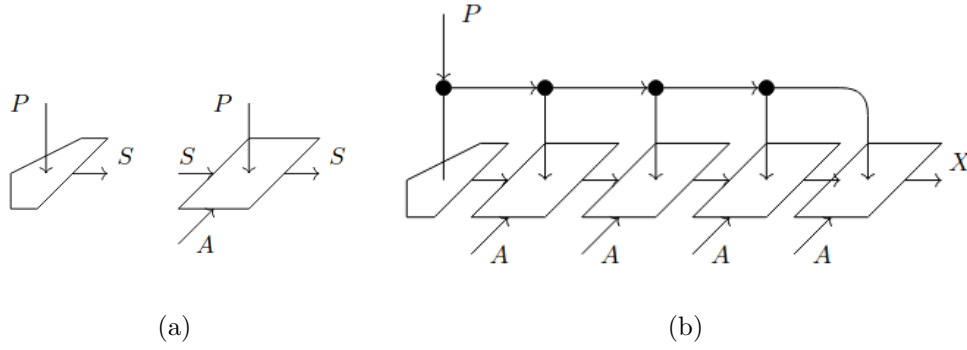


Figure 2.3: (a) Two kinds of folding RNN cell representing the initial state of the network and a generic state update based on input. (b) A folding RNN unrolled using weight tying. (Images taken from [GLD⁺24].)

$$\begin{array}{ccc}
 1 + A \times \text{List}(A) & \xrightarrow{\text{Para}(1 + A \times -)((P, f_{\text{front}}))} & 1 + A \times X \\
 \downarrow \gamma([\text{Nil}, \text{Cons}]) & \swarrow \Delta_P & \downarrow (P, \text{cell}^{\text{front}}) \\
 \text{List}(A) & \xrightarrow{(P, f_{\text{front}})} & X
 \end{array}$$

Figure 2.4: Lax-commutative diagram that makes $(P, f) : (\text{List}(A), [\text{Nil}, \text{Cons}]) \rightarrow (S, (P, \text{Cell}))$ an algebra homomorphisms. (Image taken from [GLD⁺24].)

The construction in *Ex. 41* and *Ex. 42* constitutes a precise mathematical link between the classical data structure of lists and the machine learning construct of folding RNNs. Similarly, [GLD⁺24] recovers recursive neural networks (TreeRNNs) by building upon classical binary trees and, even more interestingly, complete RNNs are recovered from the coalgebra of *Ex. 29*, which reveals an interesting link between RNNs and Mealy machines. This begs the question: if Mealy machines generalize to recurrent neural networks, what do Moore machines generalize to? It is argued in the paper that they generalize to a variant of RNN where different cells (which share the same weights) are used for state update and output production. Hopefully, more work in this direction will lead to new NN architectures inspired from other classical concepts. *Fig. 2.5* shows various kinds of NN cells and the endofunctors used in their (co)algebraic definitions.

Remark 43 (CDL and optic-based learning). In all the examples discussed above, the (co)algebra homomorphisms in question return parametric maps (P, model) , which we can interpret as untrained NN models. We can feed these maps into the \mathbf{R}_C functor associated with a generalized Cartesian reverse differential category⁴ to augment them with their reverse derivative. The framework of parametric lenses described in *Sec. 1.2* can then be used to train these networks. CDL and optic-based learning are thus compatible and even complementary.

⁴The examples illustrated in this section have been developed in **Set**, but we see no reason why they couldn't be specialized to an appropriate CRDC.

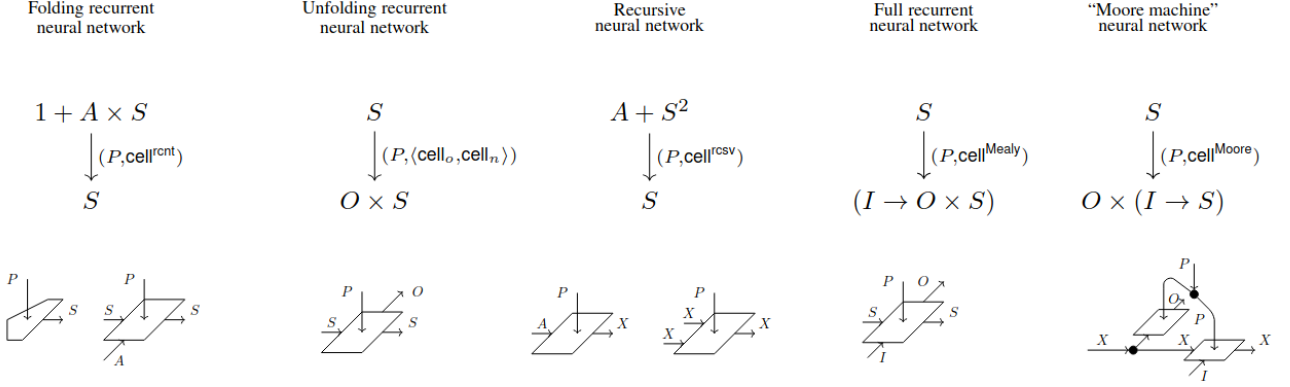


Figure 2.5: Various kinds of NN cells and the endofunctors used in their (co)algebraic definitions. (Image taken from [GLD⁺24].)

2.3 Algorithmic alignment: GNNs and dynamic programming

One of the main tenets of neural algorithmic reasoning is *algorithmic alignment* ([XLZ⁺19]), that is, the presence of structural similarities between the subroutines of a particular algorithm and the architecture of the neural network selected to learn such algorithm. Since [XLZ⁺19] has shown that dynamic programming algorithms align very well with GNNs, and since dynamic programming encompasses a wide variety of techniques algorithms to many different domains, GNNs are at the forefront of neural algorithmic reasoning research ([DV22]). However, the exact link between GNNs and dynamic programming has yet to be fully formalized. In this section we present the work of [DV22], which attempts to derive such a formalization, and the work in [DvGPV24], which studies conditions under which GNNs are invariant with respect to various forms of asynchrony⁵, which is argued to improve algorithmic alignment in some cases.

2.3.1 Integral transforms for GNNs and dynamic programming

The main link between dynamic programming and GNNs is that dynamic programming itself can be interpreted from a graph-theoretical point of view. Dynamic programming breaks up problems into subproblems recursively until trivial base cases are reached. We can thus consider the graph with nodes corresponding to subproblems and edges (y, x) corresponding to the relationships ‘ y is a subproblem of x ’. Then, the solutions of the subproblems are recursively recombined to solve the original problem. This dynamic is very similar to message passing: the simpler cases are solved first, and their solutions are passed as messages along the edges so that they can be used to solve more complex cases. The architecture of a graph neural network usually implements such a message passing dynamic. However, rigorously formulating the link between the architecture of GNN and the structure of an associated dynamic programming algorithm is not easy, the main obstacle being the difference in data type handled by the two mathematical processes: dynamic programming usually deals with tropical objects such as the semiring $(\mathbb{N} \cup \{\infty\}, \min, +)$, while GNNs usually deal with linear algebra over \mathbb{R} ([DV22]).

[DV22] proposes the formalism of integral transforms as the common structure behind both GNNs and dynamic programming. While a full formal proof is not given, the idea is illustrated by showing

⁵While much of the work described in [DvGPV24] does not fall under the umbrella of applied category theory, we still mention it because of its close link with the work of [DV22] and with the idea of algorithmic alignment. Hopefully, future work will explore the intersection between this work and category theory.

that both the Bellman-Ford algorithm and the message passing neural network (MPNN) architecture can be expressed with the help of integral transforms. The difference in data type is overcome by using the weakest common hypothesis: that the data and associated operations form a semiring.

Bellman-Ford algorithm

The Bellman-Ford (BF) algorithm is one of the most popular dynamic programming algorithms and is used to find the shortest paths between a single starting node and every other node in a weighted graph $G = (V, E)$. Since we can see every node of the graph as a subproblem, and since we can see the associated edges as subproblem relationships, the BF algorithm is a very good candidate for a GNN implementation. The algorithm operates within the tropical min-plus semiring $(R = \mathbb{N} \cup \{\infty\}, \oplus = \min, \otimes = +)$, and the data can be provided as a tuple (d, b, w) of three functions into R . Here, $d : V \rightarrow R$ stores the current best distances of the nodes, $b : V \rightarrow R$ stores the weights of the nodes, and $w : E \rightarrow R$ stores the weights of the edges. d is initialized as the function that maps the initial node to 0 and every other node to ∞ . The values of d are updated at each step of the algorithm according to the following formula, where $\mathcal{N}_u = \{v \text{ s. t. } (v, u) \in E\}$ represent the one-hop neighborhood of a node u :

$$d_u \leftarrow \min \left(d_u, \min_{v \in \mathcal{N}_u} (d_v) + w_{v,u} \right).$$

[DV22] propose the integral transform encoded by the polynomial span in *Fig. 2.6 (a)* as the supporting structure of the BF algorithm. The functions i , p , and o are defined as follows:

1. $i : (V + E) + (V + E) \rightarrow V + (V + E)$ acts as the identity on the first V , it maps the edges of the first E to their sources, and it acts as the identity on the second $V + E$ pair;
2. $p : (V + E) + (V + E) \rightarrow V + E$ just collapses the two copies of $V + E$;
3. o acts as the target function on the E and as the identity on V .

It is argued in the paper that the whole integral transform acts as step of the algorithm, carrying the data in (d, b, w) to the updated function d . Let's examine each step: the input pullback extracts the distances the sources of every edge; the argument pushforward computes the lengths of the one-hop extensions of the known shortest paths (the weight of each node is treated as the weight of a self-edge in this case); finally, the message pushforward selects the shortest paths to each node among the ones studied by the argument pushforward. Hence, the simple polynomial span in *Fig. 2.6 (a)* successfully encode the whole BP algorithm without any information loss or ad hoc choice.

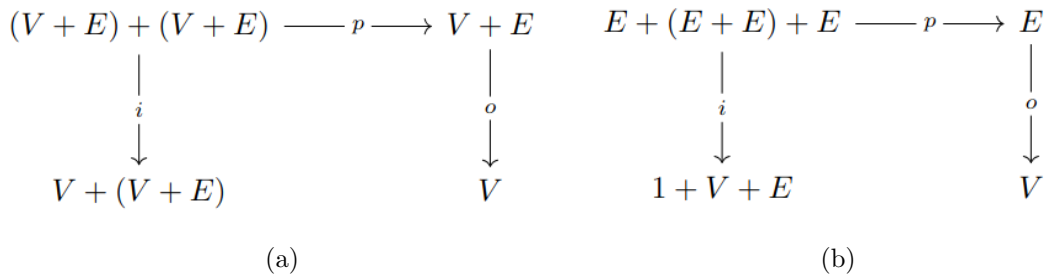


Figure 2.6: Polynomial spans used to by [DV22] to model (a) the Bellman-Ford algorithm and (b) a message passing neural network. (Images taken from [DV22].)

Message passing neural network

A message passing neural network (MPNN) is a graph neural network described by the following equations ([GSR⁺17]):

$$\begin{aligned} m_v^{t+1} &= \sum_{w \in \mathcal{N}_v} M_t(h_v^t, h_w^t, e_{v,w}), \\ h_v^{t+1} &= U_t(h_v^t, m_v^{t+1}), \end{aligned}$$

where t represents the time step, and M_t and U_t are learned differentiable functions.

[DV22] argues that an MPNN layer can be implemented as the integral transform associated with the polynomial span of *Fig. 2.6 (b)*, with an extra MLP. Here,

1. i sends the first E to 1, it acts as source and target on the second E and third E respectively, and it acts as the identity on the fourth E ;
2. p collapses four E 's into one;
3. o acts as the target function.

In the associated integral transform, i^* gathers graph, features, node features and edge features, p_\otimes "projects" the such features on the edges, the MLP combines them, and, finally, o_\oplus sends them to the right target.

Although not a perfect representation of the target architecture (due to the extra MLP), the polynomial span in *Fig. 2.6 (b)* can be used to inform the design of new architectures which are obtained by simple manipulations of the arrows or objects in the diagram. For instance, [DV22] uses the integral transform formalism to investigate possible performance improvements on CLRS benchmark tasks ([VBB⁺22]). They consider messages that reduce over intermediate nodes they show these architectures lead to better average performance on these tasks, likely a result of better algorithmic alignment.

2.3.2 Asynchronous algorithmic alignment

[DvGPV24] derives conditions under which synchronous GNNs are invariant under asynchronous execution. This is relevant because, as stated in the paper, in many dynamic programming tasks modeled by graphs, only small parts of the aforementioned graphs are changed at each step. A synchronously executed GNN that is trained on these tasks must learn the identity function many times over, which leads to brittleness and wasted computational resources. However, implementing many algorithms that are invariant under asynchronous execution as synchronous GNNs that are also invariant under asynchronous execution leads to better algorithmic alignment and thus more robust and efficient learning.

The authors of [DvGPV24] revise the model explored in [DV22] so that it includes a message function ψ that generates messages based on gathered arguments (see *Fig. 2.7 (a)* for the update diagram). Moreover, the authors argue that it is best to consider GNNs where every graph component that has a persistent state is elevated to the status of node, whereas transient computations are carried out along edges. The resulting GNN can be described by the diagram in *Fig. 2.7 (b)*, where ϕ is the transit function that updates the persistent state of each node, and δ is the function that computes the arguments needed to generate the next messages.

It is argued in [DvGPV24] that invariance under asynchrony can be modeled by giving both arguments and messages monoidal structures. For instance, let $(M, 1, \cdot)$ be the message monoid and let $(A, 0, +)$ be the argument monoid. Then, if S is the set of persistent states, state update and argument generation can be modeled as a function $M \times S \rightarrow S \times A$ which maps $(m, s) \mapsto (m \bullet s, \delta_m(s))$.

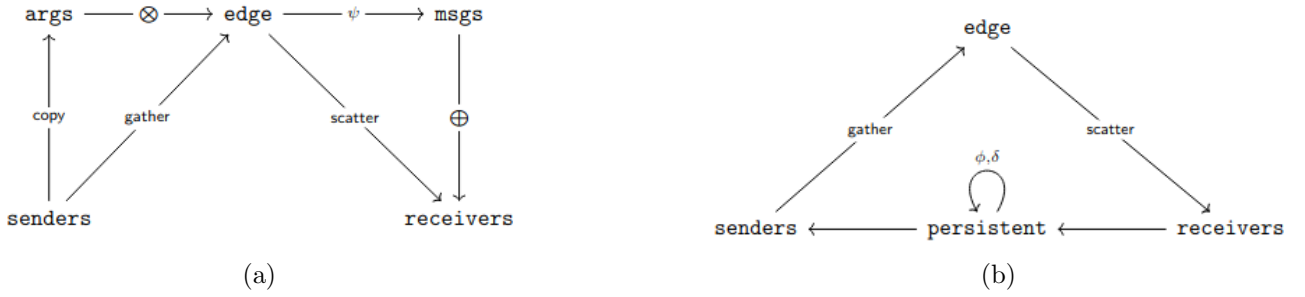


Figure 2.7: Diagrams describing the GNN architecture described in [DvGPV24]. (Images taken from [DvGPV24].)

Invariance under asynchronous message aggregation is obtained by defining \bullet as a monoidal action of M on S . However, [DvGPV24] shows that this is meaningful if and only if the argument generation function δ is compatible with the unitality and associativity equations of the action. This can only happen if $\delta : m \mapsto \delta_m$ is a 1-cocycle.

Definition 44 (1-cocycle). *A map $\delta : M \rightarrow [S, A]$ is a 1-cocycle if and only if the following are satisfied:*

1. $\delta_1(s) = 0$ for all $s \in S$;
2. $\delta_{n \cdot m}(s) = \delta_n(m \cdot s) + \delta_m(s)$ for all $s \in S$.

Proposition 45. *The state update function δ described above is asynchronous with respect to message passing if and only if it is a 1-cocycle.*

[DvGPV24] also proves the following.

Proposition 46. *Under the hypotheses described above, a single-input message function ψ supports asynchronous invocation if and only if ψ is a homomorphism of monoids.*

We will not describe the whole formalism of [DvGPV24], but we will show (without proof) its implications on GNN architecture design.

Example 47. Consider the GNN architecture:

$$x'_u = \phi \left(x_u, \bigoplus_{v \in \mathcal{N}_u} \psi(x_u, x_v) \right),$$

where \oplus is a message aggregator function. The authors of [DvGPV24] derive conditions under which this architecture is invariant under asynchronies in message aggregation, node update, and argument generation: the GNN is trivially invariant under asynchronous message aggregation if messages $(M, 0, \oplus)$ are given a commutative monoidal structure; invariance under asynchronies in node updates is obtained by selecting an update function ϕ which satisfies the associative law $\phi(s, m \oplus n) = \phi(\phi(s, m), n)$ for all $m, n \in M$ and for all $s \in S$; finally, invariance under argument generation is obtained if ϕ satisfies the 1-cocycle equations (Def. 44). [DvGPV24] also gives sufficient conditions that ensure that the 1-cocycle equations as well as the associative law are true: it suffices that $M = A = S$ and that $\phi = \oplus$ is idempotent.

2.4 Future directions and related work

In this section we provide a brief introduction to the theory of differentiable causal computations and the theory of sheaf neural networks. These two lines of work are adjacent to the main theme of this chapter - relating classical computer science to modern machine learning - and they highlight possible directions for future research into categorical deep learning and the application of integral transforms to neural networks.

2.4.1 Differentiable causal computations

[SK19] studies the differential properties of causal computations offering valuable insight into the formal properties of recurrent neural networks (RNNs). The paper focuses on sequences $\{f_k\}_{k=0,1,\dots}$ of functions $f_k : S_k \times X_k \rightarrow S_{k+1} \times Y_k$ which represent computations executed in discrete time k , where, at each tick k of the clock, f_k takes an input x_k and the previous state s_k , and uses this data to compute an output y_k and a new state s_{k+1} . In symbols, $f_k(s_k, x_k) = (s_{k+1}, y_k)$. Such sequences are given a nice compositional structure using the formalism of double categories.

Definition 48 (Category of tiles). *Let \mathcal{C} be a Cartesian category. Define $\mathbf{Db}(\mathcal{C})$ as the double category with the following data:*

1. *there is only one 0-cell, which we represent with the symbol \cdot ;*
2. *the horizontal and vertical 1-cells are the objects of \mathcal{C} ;*
3. *a 2-cell (tile) with horizontal source S , horizontal target S' , vertical source X , and vertical target Y is a morphism $f : S \times X \rightarrow S' \times Y$ which we represent with the symbol $f : X \xrightarrow[S']{S} Y$.*

It is handy to also represent 2-cells f as the tile string diagrams in Fig. 2.8 (a). The horizontal and vertical composition laws for 2-cells are given consistently with the tile diagrams.

Definition 49 (Category of stateful morphism sequences). *Let \mathcal{C} be a Cartesian category. Define $\mathbf{St}(\mathcal{C})$ as the category with the following data:*

1. *the objects of $\mathbf{St}(\mathcal{C})$ are sequences $\mathbf{X} = \{X_k\}_{k=0,1,\dots}$ of objects of \mathcal{C} ;*
2. *the morphisms $\mathbf{X} \rightarrow \mathbf{Y}$ are pairs (\mathbf{f}, i) , where $\mathbf{f} = \{f_k\}_{k=0,1,\dots}$ is a sequence of tiles in $\mathbf{Db}(\mathcal{C})$ such that $f_k : X_k \xrightarrow[S_{k+1}]{S_k} Y_k$, for some sequence \mathbf{S} of states, and $i : 1 \rightarrow S_0$ selects an initial state.*

The morphisms of $\mathbf{St}(\mathcal{C})$ are known as stateful morphism sequences and are represented using string diagrams as in Fig. 2.8 (b).

Stateful morphism sequences can be truncated and unrolled as one would expect, and it is proved in [SK19] that there is a bijection between stateful morphism sequences in $\mathbf{St}(\mathbf{Set})([A], [B])$ and causal functions $A^{\mathbb{N}} \rightarrow B^{\mathbb{N}}$ (here $[A]$ is the constant sequence of objects $\{A, A, A, \dots\}$). More generally, given any $\mathbf{St}(\mathcal{C})$, we can restrict our attention to constant sequences $[A]$ and stateful sequences of morphisms in the form $([f], i)$, where $f : X \xrightarrow[S]{S} Y$ is a tile in $\mathbf{Db}(\mathcal{C})$. This yields a subcategory $\mathbf{St}_0(\mathcal{C})$ whose morphisms can be thought of as Mealy machines that take in an input and produce an output based on an internal state which is updated after every computation. The new state is in a sense fed back to the machine after the computation so that a new computation can take place. This is represented by the diagram in Fig. 2.8 (b).

The authors of [SK19] go on to define a delayed trace operation which provides a rigorous formalization of feedback loops such as the one in Fig. 2.9 (a). As stated in the paper, the delayed trace operator is closely related to the more popular trace operator ([JSV96]) and shares many of the same properties. Finally, the authors of [SK19] show how both $\mathbf{St}(\mathcal{C})$ and $\mathbf{St}_0(\mathcal{C})$ can be given the structure of a CDC (Def. 13), as long as \mathcal{C} is itself a CDC. This differential structure is conceptually clear,

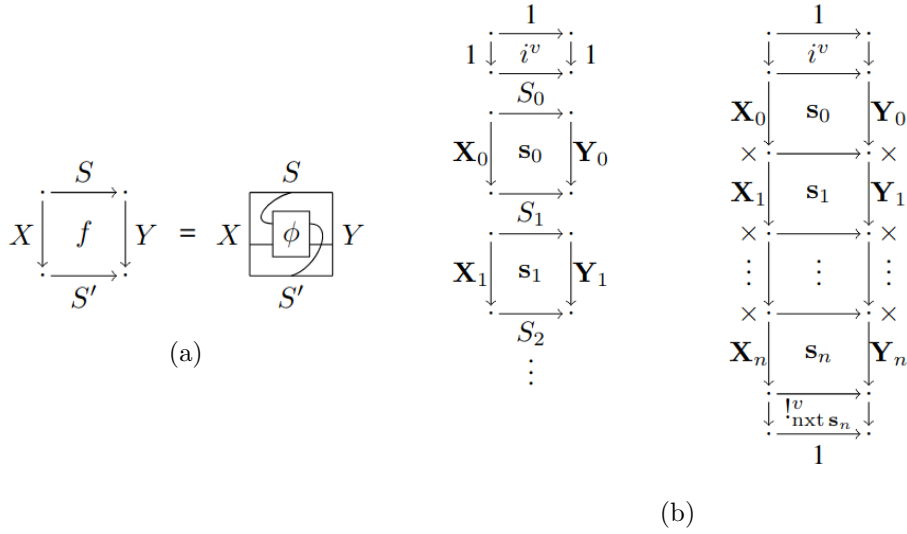


Figure 2.8: (a) String diagrams for a tile in $\mathbf{Dbl}(\mathcal{C})$. (b) String diagrams for a stateful morphism in $\mathbf{St}(\mathcal{C})$ and one of its truncation. (Images taken from [SK19].)

rigorously defined, and compatible with the dealed trace operator. We do not have space to describe the details of these definitions, but we report the relevant string diagrams in *Fig. 2.9 (b),(c)*.

The work in [SK19] provides a theoretical foundation for the technique of backpropagation through time (BPTT), which consists in computing the gradient of the k -th unrolling of a RNN in place of the gradient of the such RNN at discrete time k . Despite the alleged ad hoc nature of BPTT, [SK19] proves that the technique does not just ‘involve differentiation’ but is an actual ‘form of differentiation’ that can be reasoned about in the formalism of CDCs. Nevertheless, as stated in the paper, the differential operator of $\mathbf{St}(\mathcal{C})$ does not compute explicit gradients, and deriving the latter would from the former would be computationally intractable when there are millions of parameters.

It is interesting to compare the approach of [SK19] with the framework of categorical deep learning: both CDL and the work in [SK19] synthetically describe RNN architectures, but, while CDL focuses on weight sharing mechanics and the coinductive nature of the definition, [SK19] focuses the differential properties of these architectures. However, neither categorical framework deals with the problems that come up when computing gradients of unrolled RNNs, such as vanishing or exploding gradients (see e.g. [Han18]).

2.4.2 Sheaf neural networks

The theory of sheaf neural networks ([HG20], [BDGC⁺22], [Zag24]), or SNNs, is informed by both topology and category theory and aims to improve the GNN architecture by endowing graphs with cellular sheaf structures. In particular, SNNs are designed to solve two main issues that are encountered when training GNNs: oversmoothing, which is the tendency of deep GNNs to spread information too far in the graph to be able to effectively classify nodes, and poor performance on heterophilic graphs, that is graphs where the nodes features are diverse in structure and attributes.

Definition 50 (Cellular sheaf). *A cellular sheaf \mathcal{F} associated with a graph $G = (E, V)$ consists of the following data:*

1. a vector space $\mathcal{F}(v)$ for every node $v \in V$;
2. a vector space $\mathcal{F}(e)$ for every edge $e \in E$;

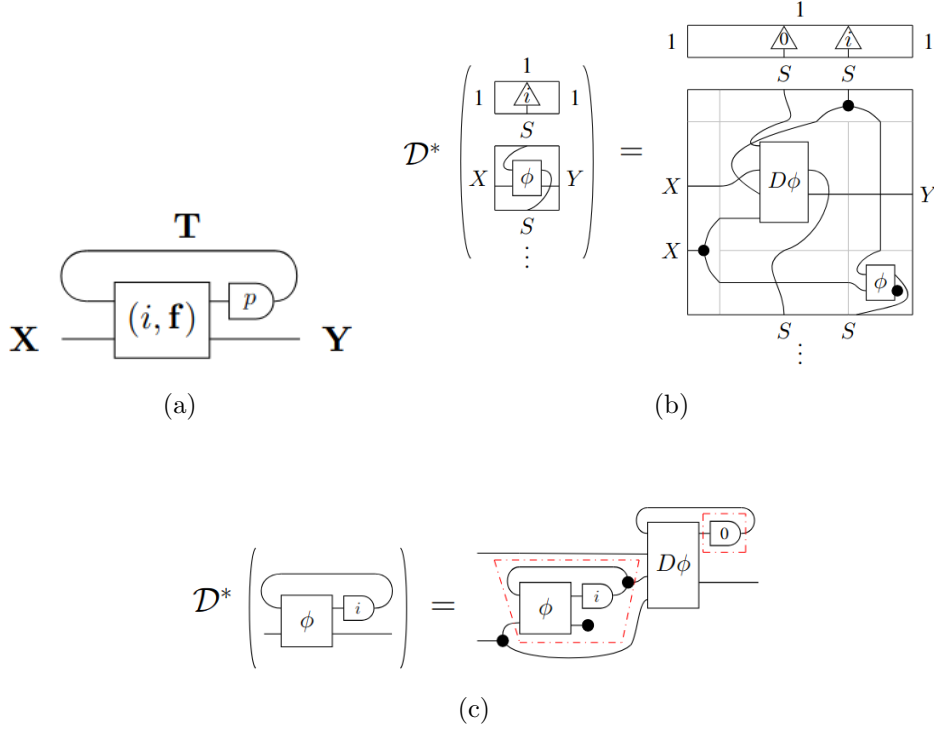


Figure 2.9: String diagrams (a) representing a delayed trace in $\mathbf{St}(\mathcal{C})$ as a delayed feedback loop, (b) representing the derivative of a morphism in $\mathbf{St}_0(\mathcal{C})$, (c) representing the derivative of a delayed trace in $\mathbf{St}_0(\mathcal{C})$. (Images taken from [SK19].)

3. a linear map $\mathcal{F}_{v \leq e} : \mathcal{F}(v) \rightarrow \mathcal{F}(e)$ for each incident node-edge pair $v \leq e$.

The vector spaces associated to nodes and edges are known as *stalks*. The linear maps associated to incident node-edge pairs are known as *restriction maps*. The direct sum $C^0(G, \mathcal{F})$ of all node stalks is known as *space of 0-cochains*, and the direct sum $C^1(G, \mathcal{F})$ of all edge stalks is known as *space of 1-cochains*.

As stated in [Zag24], the node stalks assigned by \mathcal{F} serve as spaces for node features, while the restriction maps allow the data that resides on adjacent nodes to interact on edge stalks. Given a cellular sheaf \mathcal{F} , we can define a coboundary map δ which measures the amount ‘disagreement’⁶ between nodes. The coboundary map can then be used to define a sheaf Laplacian which can be used to propagate information in the graph.

Definition 51 (Coboundary map). *Let \mathcal{F} be a cellular sheaf on a directed graph $G = (E, V)$. The coboundary map associated with \mathcal{F} is the linear map $\delta : C^0(G, \mathcal{F}) \rightarrow C^1(G, \mathcal{F})$ that maps $\delta(\mathbf{x})_e = \mathcal{F}_{v \leq e}(x_v) - \mathcal{F}_{u \leq e}(x_u)$ for each edge $e : u \rightarrow v$.*

Definition 52 (Sheaf Laplacian). *Let \mathcal{F} be a cellular sheaf on a directed graph $G = (E, V)$ and let δ be the associated coboundary map. The sheaf Laplacian associated with \mathcal{F} is the linear map $L_{\mathcal{F}} = \delta^T \circ \delta : C^0(G, \mathcal{F}) \rightarrow C^0(G, \mathcal{F})$. The normalized sheaf Laplacian associated with the sheaf is the linear map $\Delta_{\mathcal{F}} = D^{-\frac{1}{2}} \circ L_{\mathcal{F}} \circ D^{\frac{1}{2}}$, where D is the diagonal of $L_{\mathcal{F}}$.*

⁶There is a close link between the theory of SNNs and the theory of opinion dynamics. See [Zag24] for further information.

Remark 53. The coboundary map and the sheaf Laplacian associated with a cellular sheaf \mathcal{F} are generalizations of the more commonly known incidence matrix and Laplacian associated to a graph (see e.g. [WJL⁺22]).

There are many kinds of SNN architectures ([HG20], [BDGC⁺22], [Zag24]). Due to space constraints, we only give a short description of the first one to appear in the literature: the Hansen-Gebhart SNN proposed by [HG20].

Definition 54 (Sheaf neural network). *Suppose $G = (E, V)$ is a directed graph and \mathcal{F} is a cellular sheaf on it. Suppose the stalks of \mathcal{F} are all equal to $\mathbb{R}^{f \times d}$, where f is the number of channels and d is the dimension of each feature. Then, $C^0(G, \mathcal{F})$ is isomorphic to $\mathbb{R}^{nd \times f}$, where n is the number of nodes, and its elements can be represented as matrices \mathbf{X} . The sheaf neural network proposed by [HG20] uses the following transition function to update this features:*

$$\mathbf{Y} = \sigma((I_{nd} - \Delta_{\mathcal{F}})(I_n \otimes W_1)\mathbf{X}W_2),$$

where σ is a non-linearity, \otimes is the Kronecker product, $W_1 \in \mathbb{R}^{d \times d}$ and $W_2 \in \mathbb{R}^{f \times f}$ are weight matrices.

Remark 55. The values of f_1 , f_2 , n , d , and the restriction maps are all hyperparameters. Choosing $W_2 \in \mathbb{R}^{f_1 \times f_2}$ allows the SNN layer described above to change the number of features from f_1 to f_2 .

As observed by [BDGC⁺22], the SNN architecture proposed by [HG20] can be seen as a discretization of the differential equation

$$\dot{\mathbf{X}}(t) = -\Delta_{\mathcal{F}}\mathbf{X}(t),$$

which is known as sheaf diffusion equation and is analogous to the heat diffusion equation used in graph convolutional networks ([BDGC⁺22]). Studying the time limit of the sheaf diffusion equation yields important results about the diffusion of information through the graph after repeated application of the transformation in Def. 54. In particular, [BDGC⁺22] argues that, in the time limit, node features tend to values that ‘agree’ on the edges. Hence, ‘sheaf diffusion can be seen as a synchronization process over the graph’. [BDGC⁺22] goes on to study the discriminative power of different classes of cellular sheaves and proposes strategies to learn the restriction maps themselves. [Zag24] extends the work of [BDGC⁺22] by analyzing non-linear sheaf Laplacians and the associated sheaf diffusion process.

It is important to notice that SNNs can be considered a strict generalization of GNNs since the latter are nothing but instances of the former where the sheaf structure is trivial ([BDGC⁺22]). Thus, although we are not aware of any work applying SNNs to neural algorithmic reasoning, given the success enjoyed by GNNs in this area of research, we can only imagine that SNNs would be even more effective at executing algorithms. To the best of our knowledge, no one has ever explicitly described SNN architectures using integral transforms either. However, a passing remark in [DvGPV24] hints that the message passing dynamic of GNNs is very similar to sheaf diffusion and thus such a generalization should be all but impossible. Hopefully, future research will shed light on these conjectures.

Chapter 3

Functor Learning

All categorical machine learning frameworks examined in the previous chapters represent machine learning models, both trained and untrained, as morphisms in some category. Morphisms capture the core idea of compositionality and are thus a very good choice in many contexts; nevertheless, there are also cases where a simple morphism is unable to capture the structure that one might want preserved. For instance, different datasets might be linked using morphisms in an appropriate category (see e.g. [Spi12], [Gav19]) or sets containing machine learning data could be given a categorical structure, where elements are objects and morphisms capture relations between such objects (see e.g. [Lam99]). In these and other cases, learning functors instead of morphisms is advantageous as it allows us to preserve the aforementioned structure during the learning process. In this chapter we examine various approaches that testify to the usefulness of this insight: we will see how functors can be used to separate different layers of abstraction in the machine learning process ([Gav19]), to embed data in vector spaces ([SY21], [CSC10],[Lew19]), to carry out unsupervised translation ([SY21]), and to impose equivariance constraints and pool data effectively ([CLLS24]). We will illustrate how functors can be learned by gradient descent, just like morphisms, as long as we have appropriate parametrizations ([Gav19]) or specialized objective functions ([SY21], [CLLS24]).

3.1 Using functors to separate layers of abstraction

The author of [Gav19] takes inspiration from the field of categorical data migration ([Spi12]) to create a categorical framework for deep learning that separates the development of a machine learning model into a number of key steps. The different steps concern different levels of abstraction and are linked by functors.

3.1.1 Schemas, architectures, models, and concepts

The first step in the learning pipeline proposed by [Gav19] is to write down the bare-bones structure of the model in question. This can be done by using a directed multigraph G , where nodes represent data and edges represent neural networks interacting with such data. Constraints can be added at this level in the form of a set \mathcal{X} of equations that identify parallel paths (see e.g. *Fig. 3.1*).

Definition 56 (Model schema). *The schema of a model represented by a multigraph G is the freely generated category $\mathbf{Free}(G)$.*

The schema of a model does not contain any data nor does it do any computation, but it encodes the bare-bones structure of such model. If \sim is the path congruence relation induced by the equations in \mathcal{X} , we can take the module $\mathbf{Free}(G)/\sim$ to impose the constraints represented by \mathcal{X} . Depending on

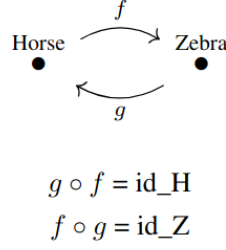


Figure 3.1: Schema for the cycleGAN architecture. (Image taken from [Gav19].)

the context, the word schema will refer either to $\mathbf{Free}(G)$ or to $\mathbf{Free}(G)/\sim$. We will see later that each equation in \mathcal{X} can be associated with a specific loss function and these losses are used to teach models to abide by \sim .

Given a schema $\mathbf{Free}(G)$, we can choose an architecture for the model, that is, we can assign to each node a Euclidean space and to each morphism a parametric map. This procedure yields an untrained neural network.

Definition 57 (Model architecture). *Let $\mathbf{Free}(G)$ be a model schema. An architecture for such schema is a functor $\mathbf{Arch} : \mathbf{Free}(G) \rightarrow \mathbf{Para}_\times(\mathbf{Smooth})$ (see Def. 5 and see Ex. 15).*

The Euclidean spaces \mathbf{Arch} maps objects to might be intuitively interpreted as the spaces the data will live in, but it is wiser to put data outside the \mathbf{Para} machinery and in the simpler \mathbf{Set} category, as this allows for better compartmentalisation. Thus, [Gav19] also defines an embedding functor, which agrees with \mathbf{Arch} on objects but exists independently of it.

Definition 58 (Model embedding). *Let $\mathbf{Free}(G)$ be a model schema and let \mathbf{Arch} be a chosen architecture. An embedding for such schema is a functor $E : |\mathbf{Free}(G)| \rightarrow \mathbf{Set}$ which agrees with \mathbf{Arch} on objects¹.*

Now that we have a model, we must find a way to assign specific values to the parameters, so that the model can be optimized by gradient descent. Consider the function $\mathbf{p} : (P, f) \mapsto P$, which takes the parameter space out of a parametric map in $\mathbf{Para}_\times(\mathbf{Smooth})$. We can use it to define the function

$$\mathfrak{P} : \mathbf{Arch} \mapsto \prod_{f : \mathbf{Gen}_{\mathbf{Free}(G)}} \mathbf{p}(\mathbf{Arch}(f)),$$

where $\mathbf{Gen}_{\mathbf{Free}(G)}$ is the set of generating morphisms of the free category on the multigraph G . The function \mathfrak{P} takes an architecture and returns the parameter space. Given \mathfrak{P} , we can define the notion of parameter specification function.

Definition 59 (Parameter specification function). *Let $\mathbf{Free}(G)$ be a model schema and let \mathbf{Arch} be a chosen architecture. A parameter specification function is a function \mathbf{PSpec} which maps a pair (\mathbf{Arch}, p) - comprised of an architecture \mathbf{Arch} and some $p \in \mathfrak{P}(\mathbf{Arch})$ - to a functor $\mathbf{Model}_p : \mathbf{Free}(G) \rightarrow \mathbf{Smooth}$. The functor \mathbf{Model}_p takes the model schema and returns its implementation according to \mathbf{Arch} , partially applying p_f to each $\mathbf{Arch}(f)$, so that we obtain an actual smooth map.*

The functor \mathbf{Model}_p takes a schema and implements it replacing nodes with Euclidean spaces and arrows with appropriate smooth functions. Notice that, although we need to choose an architecture

¹The reason why the domain E is the discretized schema $|\mathbf{Free}(G)|$ instead of the original schema $\mathbf{Free}(G)$ is clarified in Rem. 61.

and specific values for its parameter in order to define Model_p , the latter acts on the schema directly and does not pass through the architecture. See *Fig. 3.2* for a string diagram that depicts the relationship between Arch and Model_p .

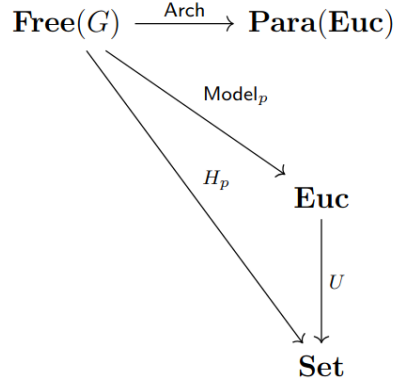


Figure 3.2: Functors are used to link different layers of abstraction in the framework of [Gav19]. (Image taken from [Gav19].)

Now, if we hope to train the model we have defined, we will need a dataset. [Gav19] suggests that a dataset should be represented as a subfunctor of the model embedding functor.

Definition 60 (Dataset). *Let E be a model embedding. Then, a dataset is a subfunctor $D_E : |\mathbf{Free}(G)| \rightarrow \mathbf{Set}$ which maps every object A of the discretized free category $|\mathbf{Free}(G)|$ to a finite subset $D_E(A) \subseteq E(A)$.*

Remark 61. The reason why [Gav19] defines E and D_E on discretized categories is because it often happens in practical machine learning that the available data is not paired. In these cases, it would be meaningless to provide an action on morphisms because they would end up being incomplete maps.

Given a node A of G , we have associated to A a Euclidean space $E(A)$ and a dataset $D_E(A)$. A dataset may be considered a collection of instances of something more specific than just vectors; for instance, if we have a finite dataset of pictures of horses, we are clearly interested in the concept of horse, that is in the set of all possible pictures of horses, which is much larger than our dataset but still much smaller than the vector space used to host such pictures. It makes thus sense to define another set $\mathfrak{C}(A)$ representing such concept: this set will satisfy the inclusion relations $D_E(A) \subseteq \mathfrak{C}(A) \subseteq E(A)$. Moreover, since concepts are assumed to be complete, we can extend \mathfrak{C} to a functor.

Definition 62 (Concept functor). *Given a schema $\mathbf{Free}(G)/\sim$, an embedding E and a dataset D_E , a concept associated with this information is a functor $\mathfrak{C} : \mathbf{Free}(G)/\sim \rightarrow \mathbf{Set}$ such that, if $I : |\mathbf{Free}(G)| \rightarrow \mathbf{Free}(G)/\sim$ is the inclusion functor, $D_E \subseteq I \circ \mathfrak{C} \subseteq E$.*

As [Gav19] states, \mathfrak{C} is an idealization, but it is a useful idealization as it represents the goal of the optimization process: given a dataset $D_E : |\mathbf{Free}(G)| \rightarrow \mathbf{Set}$, we wish to learn the concept functor $\mathfrak{C} : \mathbf{Free}(G)/\sim \rightarrow \mathbf{Set}$. More concretely, we want to train a model Model_p which is as close of an approximation of \mathfrak{C} as possible². Total achievement of such goal is clearly impossible as, even in the

²Notice that Model_p and \mathfrak{C} have different domains and codomains. The difference of codomains is not an issue: we can just forget the differential structure of the codomain of Model_p , which yields a functor H_p with image in \mathbf{Set} (see *Fig. 3.2*). The difference of domains, on the other hand, is by design as Model_p can only approximate the constraints imposed by \sim . This is not an issue in practice as we are only interested in the performance of the image of the trained model Model_p .

simplest of cases (such as linear regression on synthetic linearly generated data), finite arithmetics and the finite nature of the learning iteration prevent us from obtaining a perfect copy of the generating function. Nevertheless, we will hopefully design an optimization process which makes the learning iteration converge towards the ideal goal.

Now that we know what the optimization goal is, we can define tasks. The task formalism brings together what has been defined in this section in an integrated fashion.

Definition 63 (Task). *Let G be a directed multigraph, let \sim be a congruence relation on $\mathbf{Free}(G)$ and let $D_E : |\mathbf{Free}(G)| \rightarrow \mathbf{Set}$ be a dataset. Then, we call the triple (G, \sim, D_E) a task.*

Once we are assigned a machine learning task (G, \sim, D_E) , we have to choose an architecture, an embedding and a concept compatible with the given multigraph, equations and dataset. Then, we specify a random initial parameter with an appropriate parameter specification function. Now we can choose and optimizer, but we must be careful to choose an appropriate loss function. The loss function should incorporate both an architecture specific loss and a path equivalence loss. The former penalizes wrong predictions while the latter penalizes violations of the constraints embodied by \sim .

Definition 64 (Path equivalence loss). *Let (G, \sim, D_E) be a task. Let Model_p be an associated model. Then, if $f \sim g : A \rightarrow B$ in G , we define the path equivalence loss associated with f , g and Model_p as*

$$\mathcal{L}_{\sim}^{f,g} = \mathbb{E}_{a \sim D_E(A)} [\|\text{Model}_p(f)(a) - \text{Model}_p(g)(a)\|].$$

Definition 65 (Total loss). *Let (G, \sim, D_E) be a task. Let Arch be an associated architecture, let Model_p be an associated model, and let \mathcal{L}' be an architecture specific loss. Suppose γ is a non-negative hyperparameter. Then, we define the total loss associated with the task, the architecture, the model, and the hyperparameter as*

$$\mathcal{L} = \mathcal{L}' + \gamma \sum_{f \sim g} \mathcal{L}_{\sim}^{f,g} \quad (3.1)$$

We can now proceed as usual, computing the loss on the dataset for a number of epochs and updating the parameter p each time. Notice that \mathcal{L} implicitly depends on p because each $\mathcal{L}_{\sim}^{f,g}$ does.

It is important to notice that, while the learning iteration employed by [Gav19] is nothing new, the functor approach is actually novel, in that the usual optimization process is used to explore a functor space instead of a simple morphism space. This point of view offers two main advantages: on one hand, it separates different layers of abstraction are separated providing much needed conceptual clarity; on the other hand, it offers an explicit treatment of constraints, which are often only described implicitly and hidden away in the architecture of the model or in the loss function.

3.1.2 Datasets influence the semantics of tasks

[Gav19] centers its investigation around the powerful cycleGAN architecture ([ZPIE17]), which is described using an appropriate task consisting of the cycleGAN schema (see *Fig. 3.1*), the cycle consistency equations, and a cycleGAN dataset, that is, a dataset with two sets A and B containing data which is essentially isomorphic, such as pictures of horses and zebras. This description is particularly insightful because it can be used to prove a very important point: the choice of dataset influences the semantics of the learned task in non-trivial ways. In other words, changing the dataset functor can result in semantically different tasks even if the same schema and equations are retained.

The example shown in the paper is the [Gav19]: combining the cycleGAN schema with the cycleGAN equations and a cycleGAN dataset yields a task whose semantics can be described as *learn maps that turn horses into zebras and vice versa*. Now replace the cycleGAN dataset with a new dataset consisting of two sets A and B , where A contains pictures depicting two elements X and Y

together, and B contains separate images of X and Y . The resulting task has very different semantics: *learn how to separate X from Y* . For instance, [Gav19] shows how to use the CelebA dataset to train a neural network able to remove glasses from pictures of faces, or even insert them.

This example is especially relevant to the present discussion because it shows how important the categorical structure of **Set** can be to machine learning. We can interpret pairs *(face, glasses)* as elements of the Cartesian product of the set of faces and the set of glasses. On the other hand, the set of pictures of faces with glasses can instead be considered another categorical product of the aforementioned sets. Since categorical products are unique, we know that there must be a unique isomorphism between the two products. Then, one interpretation of the task becomes: *find the canonical isomorphism*. [Gav19] labels this task as the product task and shows empirical evidence of the validity of the approach.

3.2 Categorical representation learning

We now discuss frameworks where the categorical structure of data is not only used to integrate data distributed among different sets (as in [Gav19]) but also to model relationships within the data. Preserving and exploiting this structure functorially can lead to marked improvements in effectiveness and efficiency of the models in question because the structure of data is often closely linked with its semantics (see e.g. [SY21]), especially in the field of natural language processing, or NLP (see e.g. [CSC10], [Lew19]).

The field of *categorical representation learning* ([CSC10], [Lew19], [SY21]) aims to learn vectorial representations of data endowed with a categorical structure mirroring the one of the original data, so that models trained on these representations can use the structure efficiently. This is often a symbolic-to-vectorial transformation, which poses conceptual and computational challenges (see e.g. [CSC10]). We will examine two examples of categorical representation learning: the approach of [SY21], which uses the obtained representations to carry out (partially) unsupervised translation, and the approaches of [CSC10] and [Lew19], which develop a compositional distributional model of meaning for NLP.

3.2.1 Unsupervised functorial translation

We now illustrate the approach of [SY21] showing how it is applied to the same example described in the original paper: unsupervised translation of the names of chemical elements from English to Chinese. The authors are tasked with converting between two identical chemical compound datasets, one labeled in English and the other in Chinese. Both datasets can be given a categorical structure by considering elements as objects and chemical bonds as morphisms. Suppose \mathcal{C} and \mathcal{D} are the resulting categories. The authors of [SY21] leverage the share structure by training a model to learn a translation functor $\mathcal{F} : \mathcal{C} \rightarrow \mathcal{D}$. To this aim, they then functorially embed these categories into the vector space category defined below, which will contain representations as objects.

Definition 66 (Vector space category). *Let $n \in \mathbb{N}$. Let \mathcal{R} be the category whose objects are the vectors in \mathbb{R}^n and such that, for all $u, v \in \mathcal{R}$,*

$$\mathcal{R}(u, v) = \{M \in \mathbb{R}^{n \times n} \text{ s.t. } v = Mu\}.$$

Composition is ordinary matrix multiplication and the identity on v is $\text{id}_v = \frac{vv^T}{|v|^2}$.

The embedding of \mathcal{C} into \mathcal{R} is $\mathcal{C} \rightarrow \mathcal{R}$ functor which maps each object a to a vector v_a and each morphism f to a matrix M_f . The actual mapping can be learned with a neural network consisting of

two separate embedding layers: one mapping objects to vectors and one mapping relations to matrices. The authors of [SY21] train the embedding layers using co-occurrence statistics and negative sampling to make sure that the embedded morphisms actually represent the same relations as the original morphisms. Here, the authors of [SY21] take on a distributional point of view, positing that co-occurrence must encode such relations. In the authors’ words, ‘co-occurrence does not happen for no reason’. The training strategy used in the paper is the following: given two embedded words a and b , model the probability of co-occurrence as $P(a \rightarrow b) = \text{sigmoid}(z(a \rightarrow b))$, where the logit $z(a \rightarrow b)$ is defined as

$$z(a \rightarrow b) = F \left(\bigoplus_f v_a^T M_f v_b \right).$$

Here, F is non-linear and \bigoplus_f represents concatenation (or more sophisticated forms of aggregation) over all morphisms in \mathcal{C} . Taking on a probabilistic perspective allows us to represent non-strict relations, which is often necessary in machine learning. From this point of view, the likelihood of $f : a \rightarrow b$ is proportional to the degree of alignment between v_b and $M_f v_a$ are closely aligned. The latter can be computed as $v_a^T M_f v_b$, and the non-linearity F reads out these measurements into a single value. Now, the actual co-occurrence probability $p(a, b)$ of two objects $a, b : \mathcal{C}$ can be approximated directly from the dataset. Given a negative sampling distribution p_N on objects unrelated to a , we can implement the negative sampling objective

$$\mathcal{L} = \mathbb{E}_{(a,b) \sim p(a,b)} (\log P(a \rightarrow b)) + \mathbb{E}_{b' \sim p(b')} (\log(1 - P(a \rightarrow b'))).$$

The embedding network can then be trained by maximizing this objective function. The embedding of \mathcal{D} into \mathcal{R} can be defined and trained in the same way.

Remark 67. Strictly speaking, morphisms in category should have a specific domain and a specific codomain but, in the case of the vector space category of *Def. 66*, a single matrix M can be in many hom-sets of \mathcal{R} . This is mirrored by \mathcal{C} and \mathcal{D} , where morphisms correspond to classes of chemical bond and the same morphism can link different objects. From this point of view, morphisms in \mathcal{R} , in \mathcal{C} , and in \mathcal{D} behave more similarly to relations on sets than to morphisms in categories. Although this line of thinking does not strictly adhere to the definition of category, it helps our intuition and greatly reduces the size of the model by limiting the number of parameters that need to be learned. We will see a possible formal solution to this problem later on in this section, when we discuss the approach of [CSC10] and [Lew19].

The procedure described above can be applied to both \mathcal{C} and \mathcal{D} to obtain meaningful vectorial representations of both datasets. An English to Chinese translation functor $\mathcal{F} : \mathcal{C} \rightarrow \mathcal{D}$ can be similarly embedded as a $\mathcal{R} \rightarrow \mathcal{R}$ functor by precomposition with the $\mathcal{C} \rightarrow \mathcal{R}$ embedding and postcomposition with the inverse of the $\mathcal{D} \rightarrow \mathcal{R}$ embedding. Such functor must equate chemical bonds of the same kind, for example, if f is a covalent bond so is $\mathcal{F}(f)$. It is argued in [SY21] that this functor can be represented by a matrix $V_{\mathcal{F}}$ so that $v_{\mathcal{F}(a)} = V_{\mathcal{F}} v_a$ and $M_{\mathcal{F}(f)} = V_{\mathcal{F}} M_f$. This representation is only meaningful if (i) $V_{\mathcal{F}} M_f = M_{\mathcal{F}(f)} V_{\mathcal{F}}$ for all f , (ii) $V_{\mathcal{F}} M_{\text{id}_a} = M_{\text{id}_{\mathcal{F}(a)}} V_{\mathcal{F}}$ for all a , and (iii) $V_{\mathcal{F}} M_{f \circ g} = V_{\mathcal{F}} M_f V_{\mathcal{F}} M_g$ for all f, g . This is not true for all choices of $V_{\mathcal{F}}$ but, if we choose every v_a to be a unit vector, and if we constrain $V_{\mathcal{F}}$ to be orthogonal, (ii) and (iii) are trivially satisfied. Requirement (i), on the other hand, can be learned by minimizing the following structure loss:

$$\mathcal{L}_{\text{struc}} = \sum_f \|V_{\mathcal{F}} M_f - M_{\mathcal{F}(f)} V_{\mathcal{F}}\|^2.$$

As the authors of [SY21] remark, this loss is universal, in the sense that it does not depend on any specific object, but acts on the morphisms themselves. While this approach is very elegant and does

indeed return a functor, it might not produce the functor we expect because $V_{\mathcal{F}}$ is not unique if the M_f happen to be singular. Thus, it is better to integrate the structure loss with a second alignment loss that introduces some supervision to the unsupervised translation task. For instance, if the value of $\mathcal{F}(a)$ is known for a set A for objects, we can define

$$\mathcal{L}_{\text{align}} = \sum_{a \in A} \|V_{\mathcal{F}}v_a - v_{\mathcal{F}(a)}\|.$$

Then, the total loss can be written as a weighted sum

$$\mathcal{L} = \mathcal{L}_{\text{align}} + \lambda \mathcal{L}_{\text{struct}}, \quad (3.2)$$

where λ is a hyperparameter that regulates the relative importance of the two losses.

Remark 68. It is interesting to compare the total loss in *Eq. 3.2* to the one in *Eq. 3.1*. In both cases, the total loss is obtained as a linear combination of a model specific loss and a second more abstract loss. Both abstract losses enforce equality between morphisms but the abstract loss in *Eq. 3.2* does not enforce functoriality because Model_p is already a functor by definition in the framework of [Gav19]. This highlights the main difference between the approach of [Gav19] and the approach of [SY21]: the latter requires functoriality to be learned from structure of data, while the former takes functoriality as a given.

According to [SY21], the categorical approach described above can also be strengthened by endowing the categories \mathcal{C} , \mathcal{D} , and \mathcal{R} with monomial products that make combining objects into higher-level structures possible. In particular, the monomial structure \otimes' of \mathcal{R} can be given as

$$\begin{aligned} v_a \otimes' v_b &= \Theta(v_a \otimes v_b), \\ M_f \otimes' M_g &= \Theta(M_f \otimes M_g), \end{aligned}$$

where \otimes is the Kronecker tensor product and Θ is a learned operator that sends the products back into the original spaces. This monomial structure can then be used to mine categorical structure on multiple levels: for instance, in the running example of unsupervised translation of chemicals, it can be used to derive vectorial representations for functional groups or even whole compounds. If the data has a non-obvious high-level structure, [SY21] suggests a bootstrap approach where different possible links are randomly tested and the stronger ones are selected.

It is shown in [SY21] that the categorical representation learning framework described in the paper and the associated functor learning architecture for unsupervised translation can be successfully implemented. Benchmark tests against traditional sequence-to-sequence models show that the functorial paradigm leads to marked improvements in efficiency of learning. In particular, the authors compare a GRU cell model of similar performance as the functorial model described in the paper, noting that the former needs 17 times more parameters than the latter to learn to translate element names with a similar accuracy. The authors also compare their approach with the multi-head attention approach first introduced by [VSP⁺17], arguing that their categorical approach is an improvement over multi-head attention as the matrices M_f are essentially equivalent to the products $Q_f^T K_f$, where Q_f is the query matrix associated to f and K_f is the key matrix associated to f . Keeping M_f united emphasizes the important role of functoriality and provides intuition concerning the nature of the vectorial representations of objects: in fact, it makes sense to interpret each M_f as a metric that distorts the space \mathbb{R}^n and makes two vectors v_a and v_b closer if and only if there is a high likelihood that $a \xrightarrow{f} b$.

3.2.2 Compositional distributional model for NLP

There are two general approaches to natural language processing (NLP): the *compositional* approach and the *contextual* approach. The former aims to exploit the symbolic structure of grammar to understand text, while the latter posits that many aspects of text can be understood from distributional properties of the words ([AACFM22]). While compositionality has been very popular with theoreticians due to its elegance and its abstract structure, contextuality has been successfully applied by practitioners, yielding remarkably effective NLP models that learn the meaning of words and utterances by embedding them in vector spaces. These embeddings are mostly based on contextual information. The fact that both compositionality and contextuality are successful in their domains hints that the two approaches should be seen as complementary and that it would be auspicious to integrate them both in a single *compositional distributional framework*³ (see e.g. [AACFM22], [Mar19]).

There have been numerous attempts to develop such a framework in the last few decades. For instance, [CP07] proposes to represent both text and the structural roles of its constituents as vectors, which can then be combined using sums and tensor products. We will focus on the approach of [CSC10], whose authors take inspiration from the ideas in [CP07] and propose a framework that uses the language of compact closed categories to unite symbolic grammar rules with concrete representation vectors. Grammar is represented as a pregroup category, whereas vectorial representations are learned on vector spaces which are then organized within **FVect**, the category of finite dimensional real vector spaces. We start by defining pregroup categories⁴.

Definition 69 (Pregroup). *A pregroup is a monoidal poset category P where each object a has a left adjoint a^l and a right adjoint a^r . The adjoints must be such that the following are true:*

$$\begin{aligned} a^l a &\leq 1 \leq a a^l, \\ a a^r &\leq 1 \leq a^r a, \end{aligned}$$

where 1 is the monoidal unit of P .

The use of pregroups to model grammatical structures dates back to [Lam99]. The key idea is that we can define a grammar category as the free pregroup P generated by an initial set of basic grammatical roles. Then, we can represent various parts of speech using basic objects, their adjoints, and their monomial products. The poset structure of the pregroup serves as a system of reduction rules which can be used to prove if a sentence is grammatical. Consider the following example (taken from [Lew19]).

Example 70. Suppose n represents the role of a noun and s represents the role of a sentence. Then, an English transitive verb can be represented by the product $n^r s n^l$ in the pregroup P generated by s and n . Thus, the grammatical structure of the sentence “Dragons breathe fire.” is

$$n(n^r s n^l)n = n n^r s n^l n = (n n^r) s (n^l n) \leq 1 s 1 \leq s. \quad (3.3)$$

The pregroup axioms show that the structure of “Dragons breathe fire.” can be reduced to the sentence type s . This is proof that the sentence is grammatical. More complex pattern can be modeled by adding more grammatical types or by combining types in other ways. These kinds of computations have a nice representation in the form of string diagrams such as the ones of *Fig. 3.3*.

³Cognitive science faces a similar challenge in integrating the competing connectionists and symbolic models of the human mind. This is the central theme of [SL06], whose approach inspired [CP07] (discussed in this section). The important role of language in cognition suggests that advances in the development of a compositional distributional model of language may yield advances in cognitive science and vice versa.

⁴We work under the same assumption as [CSC10], namely, that each vector space in **FVect** is endowed with a scalar product $\langle -, - \rangle$.

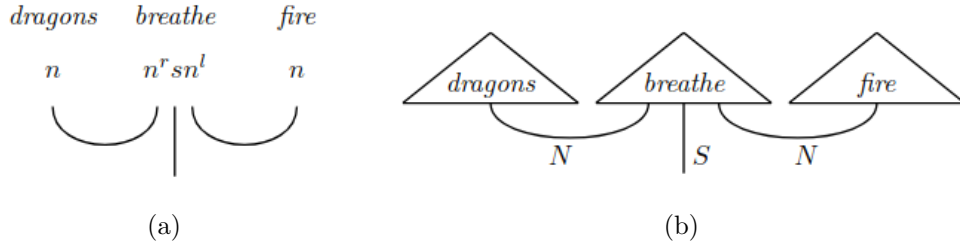


Figure 3.3: Two string diagram representations of the grammatical structure of the sentence “Dragons breathe fire.” (Images taken from [Lew19].)

The pregroup grammar framework is useful to analyze the basic structure of language but gives no access to the actual meaning of words and sentences. The compositional distributional model of [CSC10] solves this problem by building a functor $\Phi : P \rightarrow \mathbf{FVect}$. Given a certain pregroup grammar category P , for every grammatical type $x : \mathcal{P}$, we can associate a vector space $X : \mathbf{FVect}$ containing vectorial representations of utterances of type x . Φ can be defined so that it preserves the pregroup structure of P thanks to the following proposition ([CSC10]).

Proposition 71. *Consider the category \mathbf{FVect} of finite-dimensional real vector spaces. Suppose \otimes is the tensor product of vector spaces. Then, $(\mathbf{FVect}, \mathbb{R}, \otimes)$ is a monoidal category. For all $V : \mathbf{FVect}$, define $V^l = V^r = V$ and consider the (co)unit linear maps:*

$$\begin{aligned} \eta_V^l = \eta_V^r : \mathbb{R} &\rightarrow V \otimes V, \quad 1 \mapsto \sum_i e_i \otimes e_i, \\ \epsilon_V^l = \epsilon_V^r : V \otimes V &\rightarrow \mathbb{R}, \quad \sum_{i,j} c_{ij} e_i \otimes e_j \mapsto \sum_{i,j} c_{ij} \langle e_i, e_j \rangle, \end{aligned}$$

where e_1, \dots, e_n is a basis for V . Then, $(\mathbf{FVect}, \mathbb{R}, \otimes, (-)^l, (-)^r, \eta^l, \eta^r, \epsilon^l, \epsilon^r)$ is a compact closed monoidal category.

Now, pregroups have a trivial compact-closed structure as well. Thus, the functor Φ can be defined so that it preserves the shared compact-closed structure just by mapping the poset relations of P to the appropriate linear maps of \mathbf{FVect} . In practice, Φ can be used to take a formal reduction in P and implement it in \mathbf{FVect} as a series of operations on vectorial representations. The goal of this process is to take vectorial representations of single words in a sentence and obtain the vectorial representation of the meaning of the sentence. The compositional distributional model is thus compositional because it is build around the compositional structure of grammar and distributional because it allows us to encode distributional information in vectorial representations of words and sentences.

Example 72. The reduction of Eq. 3.3 can be seen as a morphism $nn^r sn^l n \rightarrow s$. The image of this morphism through Φ is a linear map $N \otimes N \otimes S \otimes N \otimes N \rightarrow S$.

Remark 73. Differently from [SY21], [CSC10] does not prescribe a specific method for learning the representations of individual words (some possible strategies are listed in the original paper), but specifies a strategy to combine existing representations to generate representations of more complex utterances up to the level of whole sentences. It is interesting to notice that both papers suggest using tensor products in order to combine representations of lower-level structures into representations of higher-level structures.

String diagrams such as the ones in *Fig. 3.3* can be written down as binary trees thanks to the associativity of the monomial product of P and the associativity of the tensor product of **FVect**. [Lew19] proposes to exploit this fact and augment the compositional distributional model of [CSC10] by implementing the (co)unit maps of **FVect** as one or more recursive neural network (TreeRNN) cells. The paper argues that the use of TreeRNN cells should make the model easier to train and more robust. To our knowledge, such proposal has not yet been implemented, although the use of TreeRNNs in NLP has already been applied in practice (see e.g. [ASM19]).

3.3 Equivariant neural networks as functors

Warning 74. We warn the reader that the paper [CLLS24], whose contents are discussed in this section, has been *withdrawn by the authors* citing issues with the references. We decided against omitting the paper from this thesis because we think its research direction is intriguing and potentially highly valuable. Hopefully, future research will soon clarify the epistemological status of the ideas hereby discussed.

[CLLS24] presents a categorical framework that describes neural network layers as functors and uses their functorial nature to impose invariance and equivariance constraints. In particular, equivariance constraints are used to effectively transport the unary operators from the dataset into a latent embedding layer. It is shown that this enables the creation of models that account for shifts and imbalance in covariates when being trained on pooled medical image datasets. In this section, we describe the framework of [CLLS24] and we draw a short comparison with the work of [PAS24], which also deals with the transport of algebraic structure into embedding layers.

3.3.1 An equivariant classifier to diagnose Alzheimer’s disease

The authors of [CLLS24] consider a data category \mathcal{S} whose objects s are data points and whose morphisms $f : s_1 \rightarrow s_2$ represent differences in covariates. An example considered in the paper is the following: suppose the objects s are comprised of brain scans and associated information concerning patient age and other covariates. The goal is to develop a model trained to diagnose Alzheimer’s disease from the scans. An example of morphism in such data category is $f_x : s_1 \rightarrow s_2$, which indicates a difference of x years in age⁵: $s_2.\text{age} = s_1.\text{age} + x$. Since we are dealing with a classification task, the dataset here has labels. It is important not to include the labels in the data category, else any classifier model would just read such labels instead of learning to predict them. Use the notation \mathbf{y}_s to represent the label associated to s . Since the dataset has labels while the data category must not have them, we feel justified in using the phrase data category in place of the phrase dataset category employed above.

Now, if we consider another category \mathcal{T} , we can use a functor $F : \mathcal{S} \rightarrow \mathcal{T}$ to project \mathcal{S} onto \mathcal{T} . Learning such functor instead of a simple map between objects is advantageous because the functoriality axioms imply that F automatically satisfies equivariance constraints: if $g : s_1 \rightarrow s_2$,

$$F(g) : F(s_1) \rightarrow F(s_2). \quad (3.4)$$

Invariance with respect to g is not much harder to define: it suffices to impose $F(s_1) = F(s_2)$ and $F(g) = \text{id}_{F(s_1)}$.

Remark 75. *Eq. 3.4* is a categorical generalization of the more usual group-theoretical notion of invariance, as defined in *Def. 37*. To be precise, *Def. 37* is equivalent to *Eq. 3.4* if \mathcal{S}, \mathcal{T} are Borel spaces and $g, F(g)$ are group actions, as highlighted in [CLLS24].

⁵Notice that this kind of operator clearly possesses an algebraic structure as $f_x \circ f_y = f_{x+y}$.

The authors of [CLLS24] solve the Alzheimer’s disease diagnosis task using a classifier that actively exploits the equivariant nature of functors. The proposed architecture consists of two modules: an invertible functor $F : \mathcal{S} \rightarrow \mathcal{T}$ that embeds the data into a latent space, and a functor $C : \mathcal{T} \rightarrow \mathbf{Free}(\mathbb{N})$ that does the actual classification (see *Fig. 3.4* for a diagrammatic representation). Here, \mathcal{T} acts as a latent space and the equivariance of F forces the representations in \mathcal{T} to be robust with respect to covariate shifts and imbalances, so that the actual classification operated by C can be more effective at diagnosing the disease. The whole model can be compactly represented as $F \circ C : \mathcal{S} \rightarrow \mathbf{Free}(\mathbb{N})$. Although the inverse of F does not appear in the formula, its existence is needed to ensure that the latent space accurately represents the data.

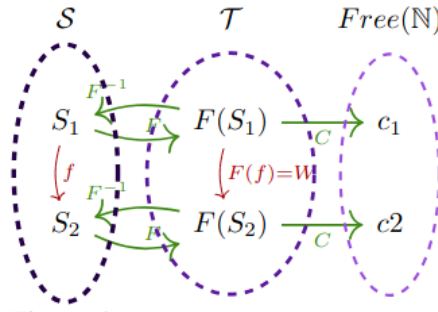


Figure 3.4: Classifier model proposed by [CLLS24] (Image taken from [CLLS24].)

The classifier described above is implemented using similar tools as the ones employed by [SY21]. The latent space \mathcal{T} can be defined as a vector space category as in *Def. 66*. The functor F is realized as the encoder of an appropriate autoencoder, and the decoder portion of this autoencoder acts as the inverse F^{-1} . The maps $f : \mathcal{S}(s_1, s_2)$ can be embedded as matrices as described in the previous section. In particular, [CLLS24] suggests using orthogonal matrices $F(f) = W_f$ so that the resulting transformations can be efficiently inverted. Finally, C can be implemented as a MLP classifier.

Now that the architecture has been specified in terms of neural networks, the resulting model can be trained using gradient descent and a linear combination of three separate losses: $\mathcal{L} = \gamma_1 \mathcal{L}_r + \gamma_2 \mathcal{L}_p + \gamma_3 \mathcal{L}_s$, where γ_1 , γ_2 , and γ_3 are hyperparameters. Here, \mathcal{L}_r is a reconstruction loss, which makes sure that F is invertible and that its inverse accurately reconstructs the original data; \mathcal{L}_p is a prediction loss, which makes sure that $F \circ C$ accurately predicts the labels of the data; \mathcal{L}_s a structure loss, which makes sure that F acts as a functor and not just a map. In formulas,

$$\begin{aligned} \mathcal{L}_r &= \sum_{s: \mathcal{S}} \|s - (F^{-1} \circ F)(s)\|_2^2, \\ \mathcal{L}_p &= \sum_{s: \mathcal{S}} \text{crossentropy}(\mathbf{y}_s, (C \circ F)(s)), \\ \mathcal{L}_s &= \sum_{\substack{s_1, s_2: \mathcal{S} \\ f: s_1 \rightarrow s_2}} \|W_f F(s_1) - F(s_2)\|_2^2. \end{aligned}$$

Remark 76. Notice that the amount of equivariance constraints imposed can vary with no substantial changes to the proposed architecture: it suffices to add or remove terms from \mathcal{L}_s . This is in stark contrast with approaches such as [LCRS22], where handling multiple covariates requires a complicate multi-stage model.

The authors of [CLLS24] test the validity of the proposed approach with two interesting experiments: a proof of concept trained on the MNIST dataset, and a working classifier trained on the ADNI brain imaging dataset. The proposed MNIST model implements equivariance with respect to increments, rotations, and zooming. It is shown in the paper that representing the associated morphisms with orthogonal matrices allows such morphisms to be inverted and combined in the latent space. A subsequent application of F^{-1} shows the results of the aforementioned manipulation in human-understandable form. Such results are indeed very promising: the authors are able to combine rotations and scaling successfully, even though the network was only trained to apply them separately (see Fig. 3.5) for an example. The ADNI classifier model also shows promising results which are on par with state-of-the-art models that do not use categorical tools. The comparison takes place according to accuracy of prediction, maximum mean discrepancy, and adversarial validation.

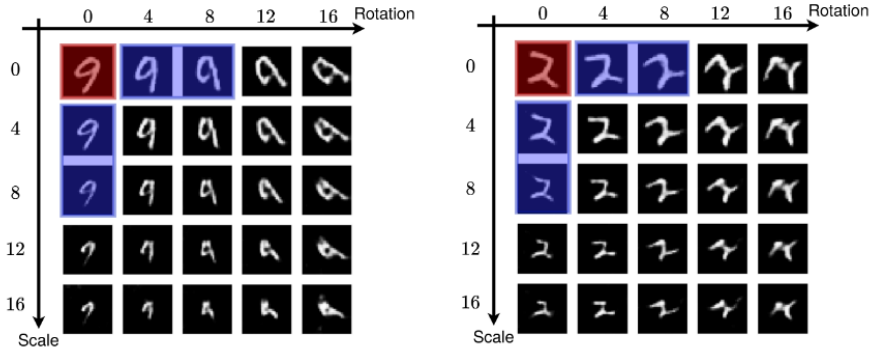


Figure 3.5: The authors of [CLLS24] applied rotations and scaling on MNIST images using the operators learned in the latent space. Despite being learned separately, the two transformations were successfully combined. (Image taken from [CLLS24].)

3.3.2 Transporting algebraic structure into embedding spaces

Although stated in terms of equivariance conditions, the work of [CLLS24] can be seen from an algebraic point of view as a framework that allows to transport unary operators from a space of data to an embedding space. This line of thinking has also been explored by [PAS24]. While [PAS24] does not use category theory, the authors themselves state that their approach could be easily categorified using the language of categorical algebras (see Def. 25). The main innovation of [PAS24] is the use of *mirrored algebras* to fill the gaps between original data and the embeddings.

Consider a pre-trained autoencoder consisting of an encoder $E : \mathcal{S} \rightarrow \mathcal{L}$ and a decoder $D : \mathcal{L} \rightarrow \mathcal{S}$, where \mathcal{S} is a data space endowed with an algebraic structure, such as a monomial product $\bullet^{\mathcal{S}}$, and where $\mathcal{L} = \mathbb{R}^n$ is a latent space. Suppose we want to transport this algebraic structure to \mathbb{R}^n through E , that is, we want to find a monomial product $\bullet^{\mathcal{L}}$ over \mathcal{L} so that E acts as a $(\mathcal{S}, \bullet^{\mathcal{S}}) \rightarrow (\mathcal{L}, \bullet^{\mathcal{L}})$ homomorphism. While finding a product $*$ that makes $(\mathcal{S}, \bullet^{\mathcal{S}})$ and $(\mathcal{L}, \bullet^{\mathcal{L}})$ homomorphic might be easy, there is no guarantee that E is an actual homomorphism with respect to this choice of algebraic structure. [CLLS24] solves this problem for unary operators by learning appropriate orthogonal matrices, but it is not clear how this approach could be scaled to operators with arities greater than 1. [PAS24] suggests an alternative approach consisting in learning a bijection ϕ between the latent space and a copy $\mathcal{M} = \mathbb{R}^n$ of this space. Endow $\mathcal{M} = \mathbb{R}^n$ with a product $\bullet^{\mathcal{M}}$ which enjoys similar properties as $\bullet^{\mathcal{S}}$ (see Fig. ??). We will say that $(\mathcal{M}, \bullet^{\mathcal{M}})$ is the mirrored algebra of $(\mathcal{S}, \bullet^{\mathcal{S}})$.

Now train ϕ so that E is $(\mathcal{S}, \bullet^{\mathcal{S}}) \rightarrow (\mathcal{L}, \bullet^{\mathcal{L}})$ homomorphism, where $\bullet^{\mathcal{L}}$ is defined by

$$l_1 \bullet^{\mathcal{L}} l_2 = \phi^{-1}(\phi(l_1) \bullet^{\mathcal{M}} \phi(l_2)).$$

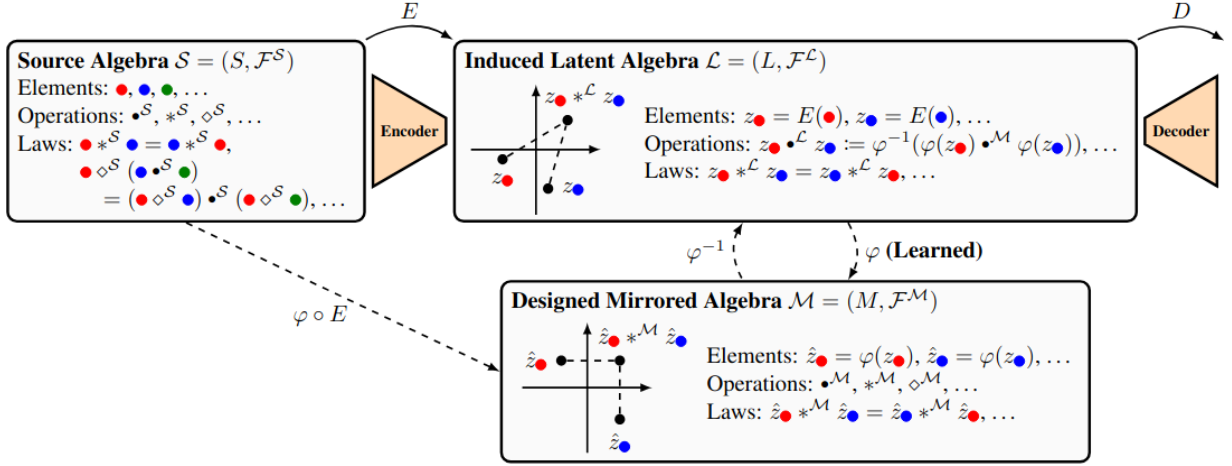


Figure 3.6: Method proposed by [PAS24] to transport algebraic structure. (Image taken from [PAS24].)

As stated in [PAS24], the bijection ϕ can be implemented as an invertible neural network and can be trained using a loss in the form

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N d(D(p_i^{\mathcal{L}}(E(s_1), \dots, E(s_{n_i}))), p_i^{\mathcal{S}}(s_1, \dots, s_{n_i}))$$

where the $p_i^{\mathcal{S}}$ are algebraic expressions sampled from the set of expressions over $(\mathcal{S}, \bullet^{\mathcal{S}})$, the s_j are sampled from \mathcal{S} , $p_i^{\mathcal{L}}(t_1, \dots, t_{n_i})$ is defined as $\phi^{-1}(p_i^{\mathcal{M}}(\phi(t_1), \dots, \phi(t_{n_i})))$, $p_i^{\mathcal{M}}$ is the translation of $p_i^{\mathcal{S}}$ in terms of the algebraic structure of $(\mathcal{M}, \bullet^{\mathcal{M}})$, and d is an appropriate distance function.

We only presented a sketch of the work in [PAS24]: the original paper presents the framework in terms of any general algebraic structure, deals with laws that regulate the interaction of different operations, proves theoretical guarantees, and explores limitations of the proposed approach. It is shown in the paper that it is not always possible to find a mirrored algebra whose operations satisfy the same laws of the source algebra, which means that it is not always possible to transport this algebraic structure in its entirety. Nevertheless, the authors carry out experiments that prove that there are tangible advantages even in transporting a subset of the algebraic laws. The experiments support the following conjecture: “learned latent space operations will achieve higher performance if they are constructed to satisfy the laws of the underlying source algebra” ([PAS24]).

Chapter 4

Detailed Representations of Neural Networks

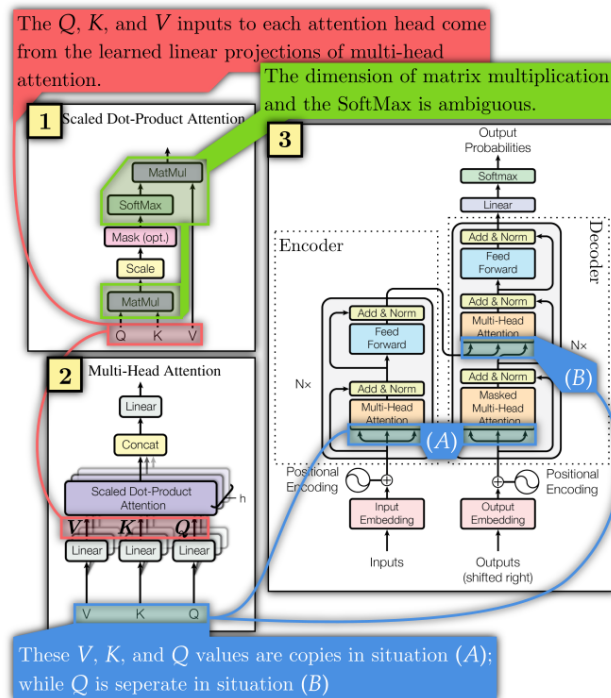


Figure 4.1: Schematics for a transformer as they were presented in the original paper [VSP⁺17]. The annotation are from the author of [Abb24]. (Image taken from [Abb23].)

One of the main concerns arising from machine learning research is low reproducibility. This is not a simple issue and investigation by many authors (see e.g. [?], [?], [?]) have suggested that many factors might be at play (such as bad incentives, unsatisfying experiment design, incorrect evaluation practices, and so on). Source code availability is a recurring theme in these discussions: for instance [?] observes that “Code describes implementation details perfectly and is required for outcome reproducibility for experiments of some complexity”. Similarly, [?] describes efforts of the NeurIPS19 conference to improve reproducibility by requiring authors to provide source code. The paper justifies

this requirement by speculating that the reproducibility issues affecting machine learning research might be partly caused by a reporting problem within the community, as machine learning papers often do not report all the details needed to reimplementing the models they describe. For instance, it is argued that hyperparameters are often missing, together descriptions of entire components of the models at hand and descriptions of the pipeline and external tools employed.

While making code and datasets as available as possible can certainly be considered a good practice for machine learning researchers, it is argued in [Abb24], [KLLWM24], and [?] that over-reliance on code is detrimental to research. In fact, while it is tautologically true that code perfectly describes implementation, code is hardly effective at capturing the essential structure of a machine learning model. Even high-level code (like Python or R code) tends to be excessively verbose and drowns important details in a sea of instructions that have little to do with the theory at play. Hence, the question remains: how should machine learning researches represent their models? This question is extremely relevant in deep learning, where the size and complexity of models is ever increasing. For instance, even the ubiquitous transformer introduced by [VSP⁺17] is a complex piece of machinery with many components, and it is very challenging to provide a satisfying description.

Many authors describe their models through equations, but this becomes more and more difficult as the size of each model grows, especially if complex tensor operations need to be carried out as ordinary tensor notation is not very readable ([XKM23]). Interestingly, the number of equations contained in a paper has been found to be negatively correlated to reproducibility ([?]). Alternatively, it is becoming more and more common to describe models using specially crafted diagrams. For instance, *Fig. 4.1* shows the transformer blueprint as provided in the original paper [VSP⁺17] (the annotations are not from the original authors). These diagrams are intuitive but, as argued in [Abb24] and [KLLWM24], they suffer from a lack of systematicity: they are usually *ad hoc* and do not have significant mathematical properties. Moreover, important information is often omitted from the diagrams and the reader has to carefully parse the text to (hopefully) find or deduce it. For instance, the diagrams in *Fig. 4.1* do not provide any information concerning the size of the matrices Q , K , or V , nor is it clear from them across which dimension the softmax should be computed.

The *ad hoc* diagrams often employed by machine learning researchers should be compared with the monoidal string diagrams that are almost standardized across applied category theory research (see e.g. [?]). We have already seen many such diagrams in the previous chapters of this thesis, especially in the description of optics (see e.g. *Fig. 1.2*). On top of being standardized, monoidal string diagrams have precise categorical semantics and can thus be used to write down rigorous mathematical proofs and computations. Given these good properties and given the promising results described in *Chpt. 1*, it is reasonable to ask whether the systematic use of monoidal string diagrams to represent machine learning models might be helpful. Regrettably monoidal string diagrams are also very inefficient at representing the details of tensor operations. The reader can appreciate this from *Fig. 1.4*, where dimensionality is all but hidden. This is the price that needs to be paid for generality, as monoidal string diagrams can represent mathematical objects that go well-beyond tensors. Nevertheless, there have been attempts by [Abb23] and by [KLLWM24] to adapt the box-wire notation of monoidal string diagrams so that it can perfectly represent a neural network architecture without leaving out implementational details. Category theory can then be used to provide rigorous semantics to these diagrams. Hopefully, such diagrams will be able to convey the information that is necessary for implementation and theoretical analysis, leading to improved insight and better reproducibility.

In this chapter, we discuss the two aforementioned categorical diagrammatic approaches and we compare them with other non-categoric and even non-graphical approaches to the problem of representing a neural network. The papers we discuss measure the validity of their approaches against the problem of representing the transformer architecture, which is both well-known and highly non-trivial. We will reproduce some of the diagrams in questions to give an idea of how the different approaches

compare. Regrettably, space concerns prevent us from fully describing the categorical mathematics behind the approaches in question, but the reader can find complete accounts of each in the respective papers.

4.1 Neural circuit diagrams

[Abb23] proposes to compensate for the shortcomings of monoidal string diagrams by introducing a novel kind of string diagrams: neural circuits diagrams. neural circuits diagrams are a specialization of a broader class of powerful string diagrams known as functor string diagrams, developed in order to represent objects, morphisms, functors, natural transformations, and products in a single diagram. In this section, we provide a brief description of both classes of diagrams, and we show how neural circuits diagrams can be used to describe the sophisticated transformer architecture.

4.1.1 Functor string diagrams

Despite their unquestionable usefulness, monoidal string diagrams have are fundamentally limited by their inability to represent functors and natural transformations. [Mar14] and [Nak23] solve this relying on colors and symbols, but this only shifts the problem as their diagrams are unable to effectively represent products ([AZ24]). Functor string diagrams follow a different approach: they represent both objects and functors as wires and both morphisms and natural transformations as boxes¹. These diagrams as built around two fundamental principles, which are already observed in traditional string diagrams. Following these principles ensures that the diagrammatic syntax stays consistent and coherent with the underlying categorical structure as notation becomes more complicated and more features as added.

Principle 77 (Vertical section decomposition). *A string diagram can be divided into vertical columns. A single vertical columns must either contain only objects or only morphisms. Object columns and morphism columns must alternate each other.*

Principle 78 (Equivalent expression). *Every newly introduced piece of graphical notation must be expressible using the notation already present in a compatible manner.*

In accordance with the above principles, [Abb23] introduces notation that represents a functors as wires that lie above objects wires and natural transformations as boxes that lie above morphism boxes (see *Fig. 4.2 (a), (b)*). In standard string diagrams, a wire running through a morphism column is interpreted as an identity morphism. Similarly, in the notation of [Abb23], a wire running where a natural transformation is supposed to be is interpreted as an identical natural transformation. This diagrammatic notation is already effective enough to write a graphical proof of the Yoneda lemma, and can be easily adapted to represent product bifunctors: it suffices to stack two rows, one per category involved in the product, and seprate them using a double dashed line (as in *Fig. 4.2 (c)*). Finally, monomial products can be represented by using a single dashed separation line or no line at all (as in *Fig. 4.2 (d)*). These piece of diagrammatical notation is also in perfect compliance with the vertical section decompositionl principle and the equivalent expression principle.

4.1.2 Neural circuit diagrams

Neural circuit diagrams are a specialization of functor string diagrams designed to represent precise schematics for deep neural networks architectures. The idea is to capture all information relevant to

¹In simple string diagrams, the actual box is often not drawn for the sake of simplicity. In these cases, the box is replaced by a single letter, but we still use the word “box” in lack of a better word.

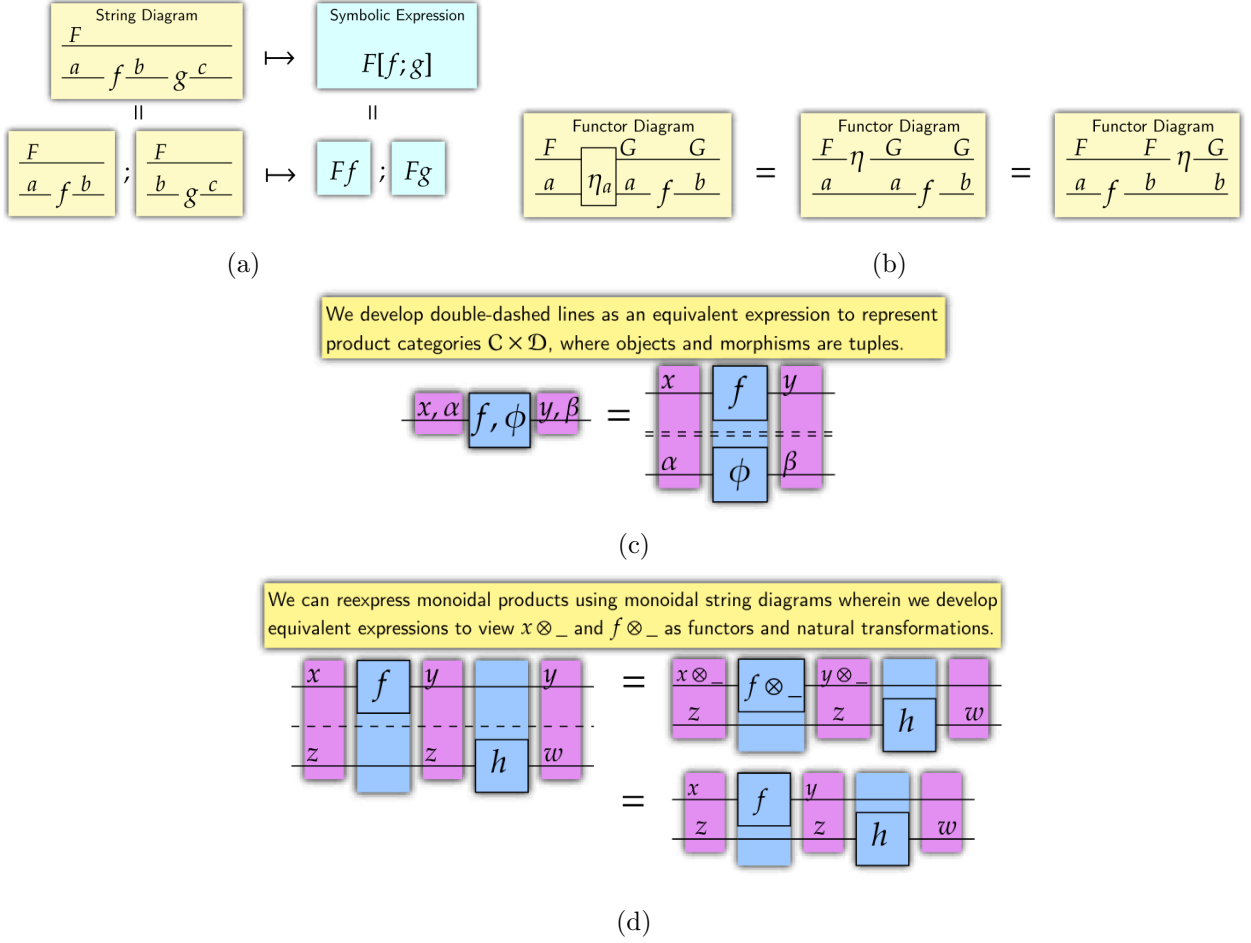


Figure 4.2: Functor string diagrams associated with (a) a functor, (b) a natural transformation, (c) a product bifunctor, (d) a monoidal product. (Images and annotations taken from [AZ24].)

implement and analyze the described architecture without using specific programming frameworks. In this sense, neural circuit diagrams can be compared to very detailed pseudocode. Although neural circuit diagrams can be learned and applied without any reference to category theory (this is the perspective adopted in [Abb24]), the fact that these diagrams are just an instance of the wider class of functor string diagrams is not just a curiosity, but proof that they are categorically sound. In other words, every well-formed neural circuit diagrams corresponds to a neural networks and neural circuit diagrams can be composed intuitively.

Differently from the string diagrams seen in e.g. [Gav24] or [WZ22], neural circuit diagrams explicitly keep track of dimensionality and indexing. This allows them to represent complex architectures that involve convolutions, residuality, and so on. Crucially, neural string diagrams provide convenient notation to represent broadcasting, which is of fundamental importance in deep learning but rarely adequately represented by ad-hoc diagrams often seen in papers.

Neural circuit diagrams represent axes of tensors as parallel wires and decorate each wire with the associated number of dimensions (see Fig. 4.3 (d)). Indexing is implemented by assigning a specific value to the wire representing the axis with respect to which we want to index (see Fig. 4.3 (b)). Finally, broadcasting is represented by running a wire that represents the additional axis over the morphism that is to be broadcasted (see Fig. 4.3 (c)). Mathematically, tensors are constructed as

hom-functors over the real line and broadcasting is just the action of such a functor on morphisms (*Fig. 4.4*). This is all representable using functor string diagram: neural circuit diagrams are just the result of omitting unnecessary detail from the latter. See [Abb24] for more information regarding neural circuit diagrams and see [Abb23] for more information regarding the relationship between neural circuit diagrams and functor string diagrams.

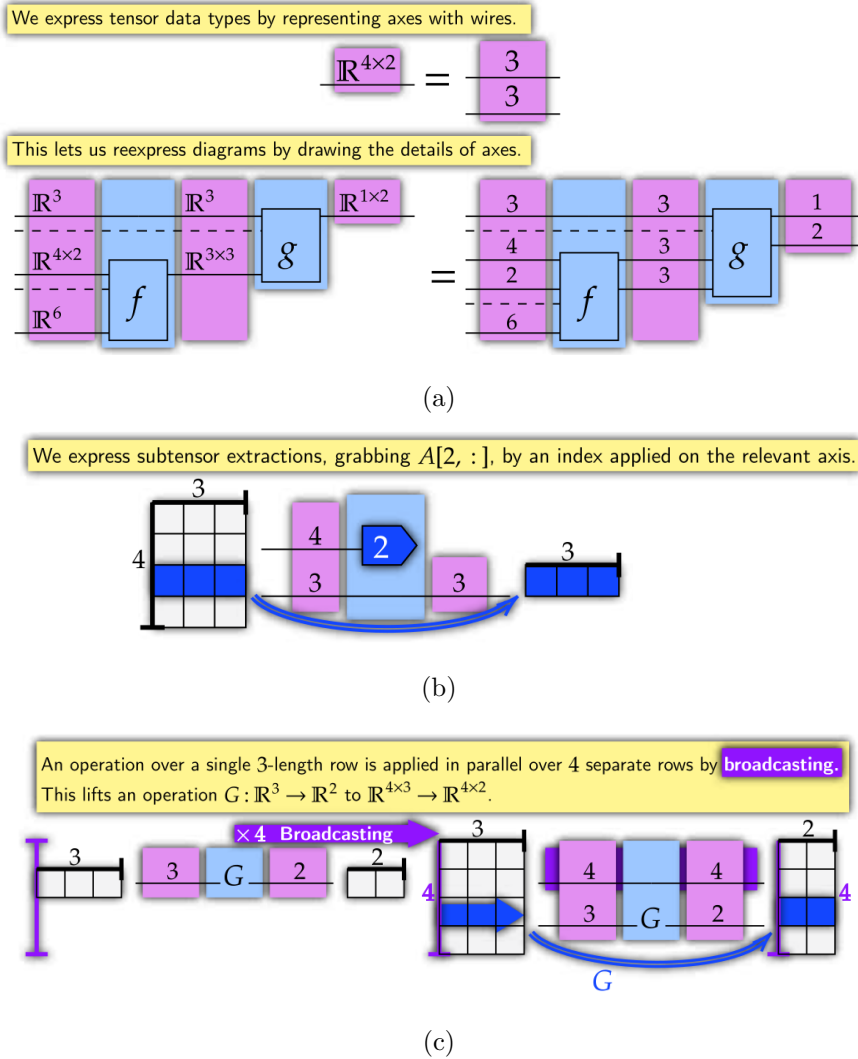


Figure 4.3: Neural string diagram formalisms for (a) representation of tensors, (b) indexing of tensors, (c) broadcasting. (Images and annotations taken from [Abb24].)

[Abb24] proves the usefulness of neural circuit diagrams by drawing schematics (*Fig. 4.5*) for the powerful transformer architecture proposed by [VSP⁺17]. Comparing these schematics with the ones in *Fig. 4.1* presented in the original paper shows the wealth of detail that the former easily encode but are completely absent from the latter. For instance, *Fig. 4.1* provides no information regarding the dimensions of Q , K , and V , despite the fact that these matrices play a central role in the architecture. Neural circuit diagrams for additional architectures can be found in [Abb24].

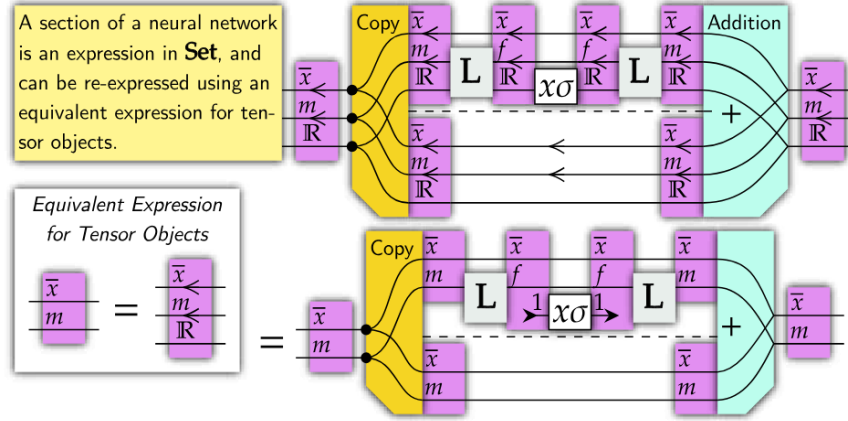


Figure 4.4: From functor string diagrams to neural circuit diagrams. (Image taken from [AZ24].)

4.2 Diagrams with universal approximators

[KLLWM24] also develops a string-diagrammatic language capable of giving detailed unambiguous descriptions of neural network architectures. It is only natural to compare the approach of [KLLWM24] with the neural circuit diagrams of [Abb23] discussed in the previous section. Despite their unquestionable similarity, the two approaches differ on one key point: while neural circuit diagrams leave no space for abstraction and only allow equational rewrites, [KLLWM24] develops a semantics based on the notion of universal approximator that allows to abstract away details and compare different architectures based on share abstractions. In this section, we give a brief overview of the approach of [KLLWM24] and we discuss interesting experimental results presented by the paper. We do not discuss the precise categorical semantics of the diagrams introduced [KLLWM24] because that would require too much space, but the reader can find a full account of them in the appendices of the original paper.

4.2.1 The diagrammatic approach of [KLLWM24]

Neural circuit diagrams effectively describe neural network architectures, but the abundance of detail they convey also means that they can make it challenging to extract the fundamental structure of an architecture from its complete schematics. [KLLWM24] attacks this problem by drawing different diagrams for different levels of abstraction so that complete schematics like the ones encoded in neural circuit diagrams become only the end of a spectrum that begins with a purely abstract black box. The same fundamental grammar underpins all these diagrams, and appropriate rewrites can be used to seamlessly move between layers of abstraction.

Let us begin by discussing the most concrete end of the aforementioned spectrum. [KLLWM24] describes computation in the Cartesian category **CartSp** of Euclidean spaces and continuous functions, and follows the basic conventions of monoidal string diagrams for Cartesian categories. A sketch of the language can be found in Fig. 4.6. Unlike neural circuit diagrams, the diagrams of [KLLWM24] represent broadcasting using *SIMD* (Same Instruction applied to Multiple Data) boxes, as in Fig. 4.6 (b), which are also used to represent tensor contractions, as in Fig. 4.6 (c). Finally, syntactic sugar is offered to represent reshaping operations (see Fig. 4.6 (c)).

The reader can compare the upper portion of Fig. 4.5 against Fig. 4.8 to see how both neural circuit diagrams and the approach of [KLLWM24] can be used to describe the same transformer

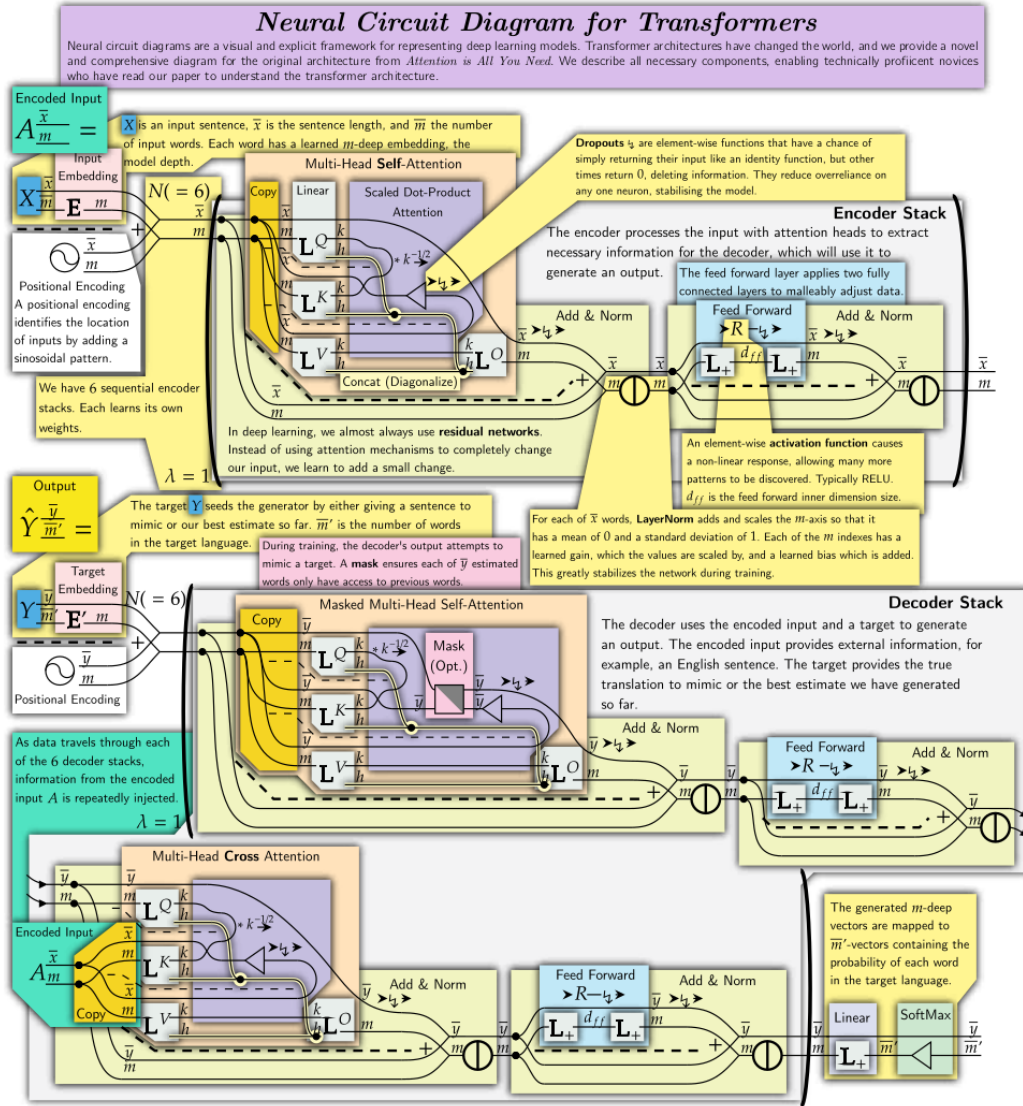


Figure 4.5: Neural string diagram of a transformer. (Image taken from [Abb23].)

architecture in detail. While both formalisms shine in accuracy of representation, they also both suffer from low readability: one might ask whether adopting such systems is even worth the trouble given that simpler diagrams such as the one in Fig. 4.1 seem to be better at communicating basic ideas despite their lack of detail. While [Abb23] compensates for this by annotating its diagrams and by encapsulating portions of the architectures at hand into dedicated boxes, [KLLWM24] proposes to abstract away detail by replacing portions of the graph representing specific implementations with filled boxes that represent universal approximators with the same inputs and outputs.

Definition 79 (Universal approximators). *A universal approximator for a class \mathcal{C} of functions is a neural network architecture able to approximate any element of \mathcal{C} to any desired precision.*

Many well known results concern the existence and characterization of such approximators (see e.g. [Cyb89]). [KLLWM24] uses universal approximators to rigorously represent “typed holes” in a diagram that need to be filled with an appropriate operator. Replacing detailed structures with

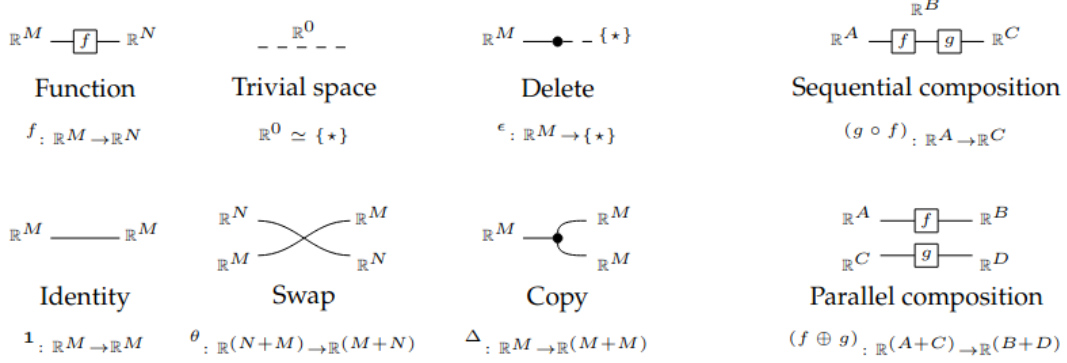


Figure 4.6: Basic elements of the diagrammatic syntax proposed by [KLLWM24]. (Image taken from [KLLWM24].)

abstract universal approximators allows us to climb the ladder of abstraction. This can be combined with the omission of dimensional information that is not necessary to understand the architecture at hand, but only to implement it. [KLLWM24] showcases this process in the diagram of *Fig. 4.9*, which highlights the fundamental structure of the attention mechanism proposed by [VSP⁺17]. Notice how colors are used to distinguish between dimensions that are fixed as hyperparameters and dimensions that vary with the length of the sequence that is being processed.

4.2.2 Comparing attention mechanisms

Not only does the abstraction allowed by universal approximators highlight the fundamental structure of architectures, but it also allows comparisons between architectures. Two architectures that only differ at a point can be considered as descendants of a common ancestor architecture, and this common link allows us to compare them. For instance, [KLLWM24] derives all known attention mechanisms as descendants of a prototype “primordial attention”, forming a sort of evolutionary tree of attention (see *Fig. 4.10*). The process of filling in a universal approximator with a more specialized structure is called by [KLLWM24] an “expressivity reduction”, as it reduces number of architectures represented by the given schematics. Expressivity reductions can be formalized as string diagrammatic rewrites. Hence, in the example discussed above, rewrites are employed to derive descendent attention mechanisms from the primordial attention.

Just like biological taxonomy classifies species based on their evolutionary history, we can use the tree in *Fig. 4.10* as a blueprint for the classification of attention mechanisms. This undertaking is analogous to the one of [GLD⁺24], where recurrent neural network cells are classified based on the classical data structures and automata they are related to. Both taxonomies serve to systematize the scattered panorama of machine learning research, but may also inspire new architectures.

The authors of [KLLWM24] do pursue this line of thinking but they do explore a space of possible attention mechanisms freely constructed from five generators (represented as abstract diagrams). This large space is reduced in size by rewriting sections that would lead to overparameterization. This results in a selection of eleven possible attention mechanisms, that the authors of [KLLWM24] go on to test: the mechanisms are embedded in a decoder-only transformer and trained on a Penn-Treebank task. Interestingly, all the 11 mechanisms (among which there are known architectures such as the one proposed by [VSP⁺17]) perform comparably well (with only one having a slightly better performance than the others). The authors entertain the conjecture that the aforementioned results might be a sign that the choice of attention mechanism matters little, as long as data is exchanged between tokens.

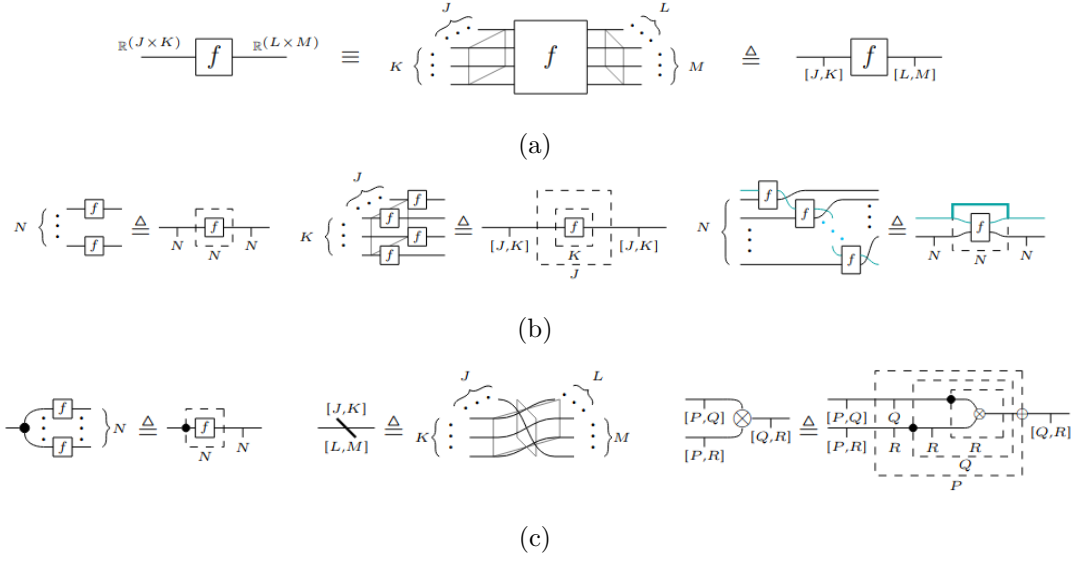


Figure 4.7: Syntax proposed by [KLLWM24] for (a) dimensionality, (b) SIMD boxes, (c) copys, reshape operators, and contractions. (Images taken from [KLLWM24].)

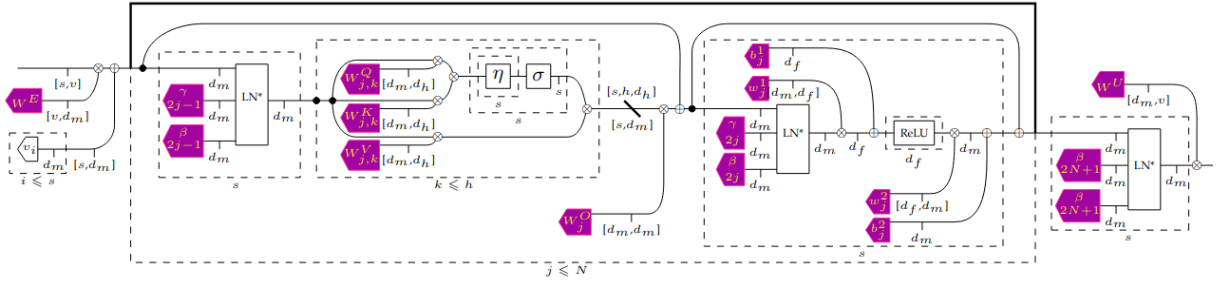


Figure 4.8: Blueprint of the encoder part of a transformer according to [KLLWM24]. (Image taken from [KLLWM24].)

While there is no uncontroversial evidence that this is the case, these experiments prove the utility of the approach of [KLLWM24] and, more generally, of the use of diagrams in machine learning.

4.3 Future directions and related work

The amount of implementational detail conveyed by neural circuit diagrams and the diagrams in [KLLWM24] suggests that it should be possible to find a 1-to-1 correspondence between these diagrams and code. It is even argued in [Abb23] that neural circuit diagrams are expressive enough to be considered a high-level language with explicit memory management. Hopefully, future research will lead to the creations of tools that convert neural circuit diagrams into working models, not unlike the library developed by [CGG⁺22] that allows programmers to code using parametric lenses. Another interesting direction of research is the application of methods analogous to the ones of [KLLWM24] to explore other classes of architectures.

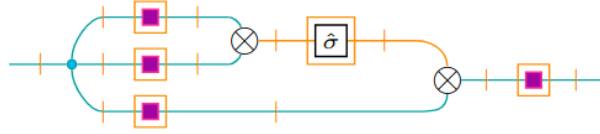


Figure 4.9: Abstract blueprint for the self-attention mechanism for the transformer architecture of [VSP⁺17]. (Image taken from [KLLWM24].)

We conclude the chapter by briefly discussing a few alternative approaches to the problem of unambiguous representation of neural networks. While these approaches are not categorical (yet), we believe there are fruitful comparisons to be drawn between the latter and the categorical approaches discussed in the previous sections. We also believe that these all these methods can be fruitfully integrated on various level, as discussed below.

4.3.1 Tensor networks

The enormous importance of tensors in quantum physics lead to the development of tensor networks (pioneered by [P⁺71]), in which tensor operations are represented as graphs, where tensors are nodes and modes are represented as edges coming out of the nodes. Since tensors are the main data structure used in deep learning, it makes sense to try to adapt tensor networks to represent neural networks. This is done in [XKM23] and [Tay24]: the two papers introduce similar conventions that endow tensor networks with dedicated notation for non-linear activation functions and copy maps, which are not present in ordinary tensor networks.

The main difference between tensor networks and the diagrams described above is that the latter represent data as boxes and processes as wires, whereas the former do the opposite. In tensor network an edge linking two tensors represents a tensor contraction: this choice yields intuitive representations for sophisticated tensor operations such as the Tucker product ([XKM23]) and tensor decompositions techniques ([Tay24]), but it also makes it much more difficult to conceptualize how data moves through a model. For instance, the reader can observe *Fig. 4.11* to understand how much more difficult it is to represent copy maps in the notation of [Tay24]. Moreover, while the diagrams of [Abb23] and [KLLWM24] have sophisticated categorical semantics, the notation of [XKM23] and [Tay24] mainly relies on human intuition².

Despite their weaknesses, tensor networks can be of great utility to deep learning practitioners as they can offer clear insight into tensorization processes. Tensorization is an architecture design pattern that aims to reduce the number of parameters of a neural network without compromising its expressivity: this can be done by replacing large matrices with appropriately decomposed low-rank tensors. Thanks to the redundancy of ordinary fully connected layers (see e.g. [?]), this strategy leads to remarkable levels of compression ([?], [XKM23]). Tensorization yields improved efficiency at the price of human intuition, as it is much more difficult to think in terms of tensor operations across many modes than in terms of matrix-vector multiplication ([XKM23]). Tensor networks such as the ones discussed in [XKM23] and [Tay24] can bridge this gap by helping practitioners to visualize tensor-based architectures. We refer to the original papers for examples and further information.

²We see no reason why the diagrams in [XKM23] and [Tay24] should not be interpreted under the lens of category theory. Refer to [BCJ11] for a categorical approach to tensor networks.

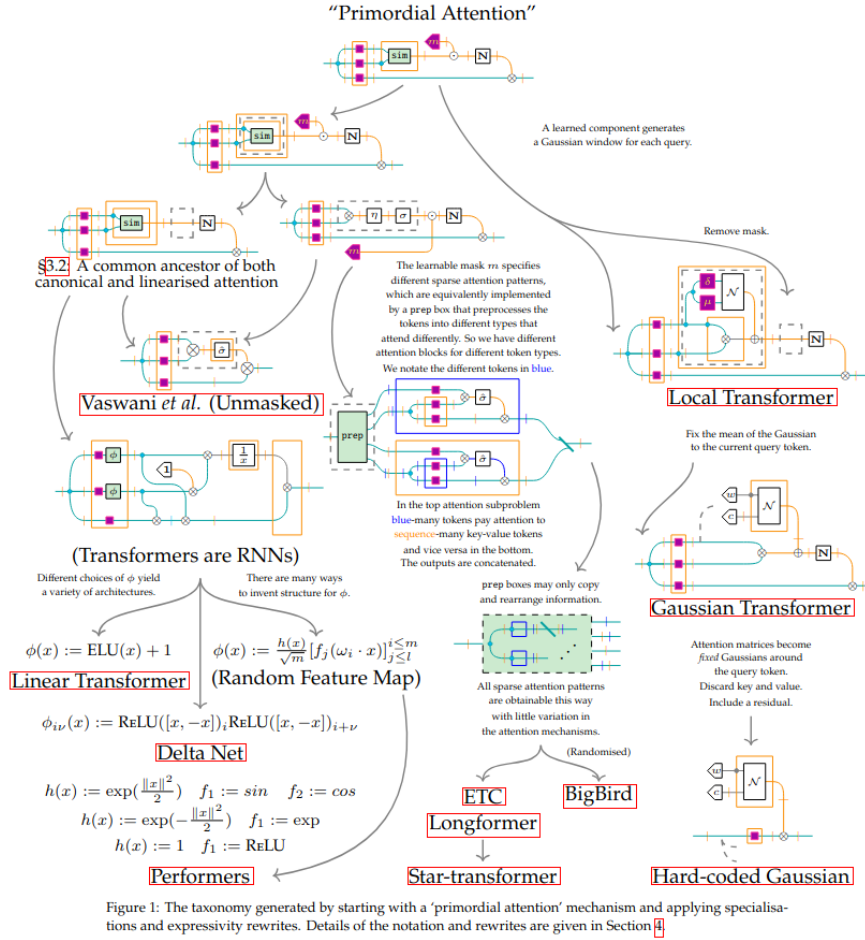


Figure 4.10: Evolutionary tree of attention mechanisms according to [KLLWM24]. Child nodes are obtained by applying rewrites to parent nodes. (Image taken from [KLLWM24].)

4.3.2 Non-diagrammatic notations

An alternative to diagrammatic notation is the enhancement of algebraic notation. For instance, [?] proposes to give English labels to tensor axes to make it easier to keep track of them. Following this convention, known as *named tensor notation*, the dot-product self-attention mechanism of [VSP⁺17] can be represented with the equation (which we lift from [?]):

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \left(\text{softmax}_{\text{seq}} \frac{\mathbf{Q} \odot \mathbf{K}}{\sqrt{|\text{key}|}} \right) \odot \mathbf{V}, \quad (4.1)$$

where **key** represents the axes dedicated to query and key features, and **seq** represents the axis of sequence tokens. The authors argue that their named tensor notation is an improvement over the popular Einstein notation for sums as the axes having English names allows the reader to immediately understand what kind of information they encode; moreover, they claim that named tensor notation limits repetition of indices and thus makes equations cleaner.

Another viable alternative approach is described in [?], where it is suggested that neural network architectures should be represented using pseudocode analogous to the one used in algorithm design

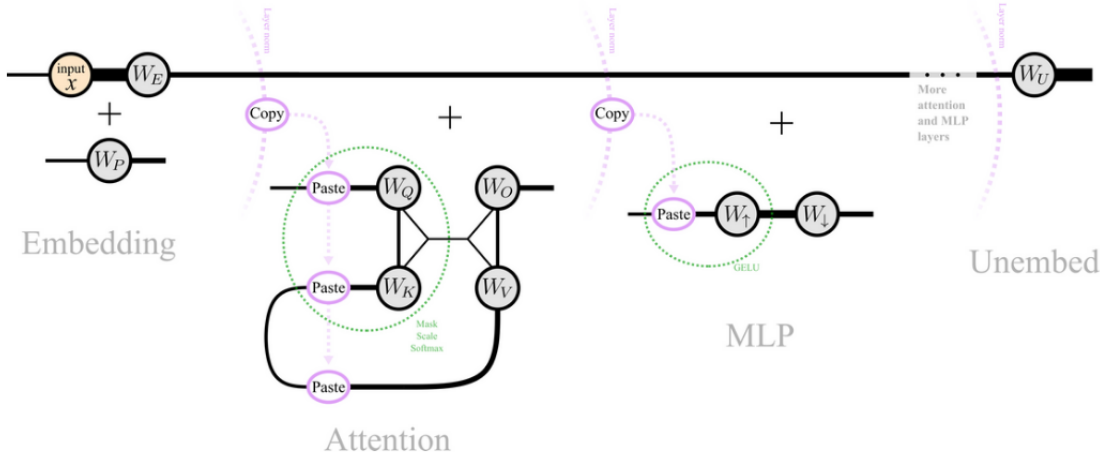


Figure 4.11: Tensor network representing the GPT-2 architecture. (Image taken from [Tay24].)

and analysis (see e.g. [CLRS22]). The authors claim that writing pseudocode distills each architecture down to its fundamental components and remove incidental distractions, while preserving all the detail necessary for implementation. One disadvantage of using pseudocode is that the structure of the involved tensors and the flow of information are partially obscured, which is detrimental to intuition. More generally, one can argue that pseudocode is readily implementable but hides away mathematics and structure.

4.3.3 A multifaceted approach

In our opinion, the problem of representing neural network architectures is best solved by using a combination of methods. For instance, there is no reason (aside maybe from space concerns rendered null by the use of appendices) why the same architecture cannot be described with a string diagram, a tensor network, named tensor equations, pseudocode, and even code in the same paper. One or more string diagrams might be designed using either of the two approaches described in this chapter; tensor networks and named tensor equations might be employed to clarify the dynamics of complex tensor operations³; pseudocode could be provided so that practitioners can readily implement the architecture; finally, providing source code might clarify any ambiguity left by the previous representations. It is argued in [?] that writing pseudocode might highlight mistakes in the actual source code that remained unnoticed by coincidence. We believe that the same is true for every other representation discussed in this chapter.

³Since the edges of tensor networks represent axes, we see no reason why they cannot be themselves labelled, effectively integrating the approaches of [XKM23] and [?].

Conclusions

Bibliography

- [AACFM22] Enrique Amigó, Alejandro Ariza-Casabona, Victor Fresno, and M Antònia Martí. Information theory–based compositional distributional semantics. *Computational Linguistics*, 48(4):907–948, 2022.
- [Abb23] Vincent Abbott. *Robust Diagrams for Deep Learning Architectures: Applications and Theory*. PhD thesis, Honours Thesis, The Australian National University, Canberra, 2023.
- [Abb24] Vincent Abbott. Neural circuit diagrams: Robust diagrams for the communication, implementation, and analysis of deep learning architectures. *arXiv preprint arXiv:2402.05424*, 2024.
- [ASM19] Mahtab Ahmed, Muhammad Rifayat Samee, and Robert E Mercer. Improving tree-lstm with tree attention. In *2019 IEEE 13th international conference on semantic computing (ICSC)*, pages 247–254. IEEE, 2019.
- [AZ24] Vincent Abbott and Gioele Zardini. Functor string diagrams: A novel approach to flexible diagrams for applied category theory. *arXiv preprint arXiv:2404.00249*, 2024.
- [BBCV21] Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- [BCJ11] Jacob D Biamonte, Stephen R Clark, and Dieter Jaksch. Categorical tensor network states. *AIP Advances*, 1(4), 2011.
- [BCS06] Richard F Blute, J Robin B Cockett, and Robert AG Seely. Differential categories. *Mathematical structures in computer science*, 16(6):1049–1083, 2006.
- [BDGC⁺22] Cristian Bodnar, Francesco Di Giovanni, Benjamin Chamberlain, Pietro Lio, and Michael Bronstein. Neural sheaf diffusion: A topological perspective on heterophily and oversmoothing in gnns. *Advances in Neural Information Processing Systems*, 35:18527–18541, 2022.
- [BQL21] Joshua Bassey, Lijun Qian, and Xianfang Li. A survey of complex-valued neural networks. *arXiv preprint arXiv:2101.12249*, 2021.
- [CCG⁺19] Robin Cockett, Geoffrey Cruttwell, Jonathan Gallagher, Jean-Simon Pacaud Lemay, Benjamin MacAdam, Gordon Plotkin, and Dorette Pronk. Reverse derivative categories. *arXiv preprint arXiv:1910.07065*, 2019.

- [CEG⁺24] Bryce Clarke, Derek Elkins, Jeremy Gibbons, Fosco Loregian, Bartosz Milewski, Emily Pillmore, and Mario Román. Profunctor optics, a categorical update. *Compositionality*, 6, 2024.
- [CG22] Matteo Capucci and Bruno Gavranović. Actegories for the working amthematician. *arXiv preprint arXiv:2203.16351*, 2022.
- [CGG⁺22] Geoffrey SH Cruttwell, Bruno Gavranović, Neil Ghani, Paul Wilson, and Fabio Zanasi. Categorical foundations of gradient-based learning. In *European Symposium on Programming*, pages 1–28. Springer International Publishing Cham, 2022.
- [CLLS24] Sotirios Panagiotis Chytas, Vishnu Suresh Lokhande, Peiran Li, and Vikas Singh. Pooling image datasets with multiple covariate shift and imbalance, 2024.
- [CLRS22] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2022.
- [CP07] Stephen Clark and Stephen Pulman. Combining symbolic and distributional models of meaning. 2007.
- [CSC10] Bob Coecke, Mehrnoosh Sadrzadeh, and Stephen Clark. Mathematical foundations for a compositional distributional model of meaning. *arXiv preprint arXiv:1003.4394*, 2010.
- [Cyb89] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- [DV22] Andrew J Dudzik and Petar Veličković. Graph neural networks are dynamic programmers. *Advances in neural information processing systems*, 35:20635–20647, 2022.
- [DvGPV24] Andrew Joseph Dudzik, Tamara von Glehn, Razvan Pascanu, and Petar Veličković. Asynchronous algorithmic alignment with cocycles. In *Learning on Graphs Conference*, pages 3–1. PMLR, 2024.
- [Ell18] Conal Elliott. The simple essence of automatic differentiation. *Proceedings of the ACM on Programming Languages*, 2(ICFP):1–29, 2018.
- [EPWJ80] Michael G Eastwood, Roger Penrose, and Raymond O Wells Jr. Cohomology and massless fields. 1980.
- [FJ19] Brendan Fong and Michael Johnson. Lenses and learners. *arXiv preprint arXiv:1903.03671*, 2019.
- [FST19] Brendan Fong, David Spivak, and Rémy Tuyéras. Backprop as functor: A compositional perspective on supervised learning. In *2019 34th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS)*, pages 1–13. IEEE, 2019.
- [Gav19] Bruno Gavranović. Compositional deep learning. *arXiv preprint arXiv:1907.08292*, 2019.
- [Gav24] Bruno Gavranović. Fundamental components of deep learning: A category-theoretic approach, 2024.
- [GLD⁺24] Bruno Gavranović, Paul Lessard, Andrew Joseph Dudzik, Tamara von Glehn, João Guilherme Madeira Araújo, and Petar Veličković. Position: Categorical deep learning is an algebraic theory of all architectures. In *Forty-first International Conference on Machine Learning*, 2024.

- [GSR⁺17] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1263–1272. PMLR, 06–11 Aug 2017.
- [Han18] Boris Hanin. Which neural net architectures give rise to exploding and vanishing gradients? *Advances in neural information processing systems*, 31, 2018.
- [HG20] Jakob Hansen and Thomas Gebhart. Sheaf neural networks. *arXiv preprint arXiv:2012.06333*, 2020.
- [IKP⁺22] Borja Ibarz, Vitaly Kurin, George Papamakarios, Kyriacos Nikiforou, Mehdi Bennani, Róbert Csordás, Andrew Joseph Dudzik, Matko Bošnjak, Alex Vitvitskyi, Yulia Rubanova, et al. A generalist neural algorithmic learner. In *Learning on graphs conference*, pages 2–1. PMLR, 2022.
- [JR97] Bart Jacobs and Jan Rutten. A tutorial on (co) algebras and (co) induction. *Bulletin-European Association for Theoretical Computer Science*, 62:222–259, 1997.
- [JSV96] André Joyal, Ross Street, and Dominic Verity. Traced monoidal categories. In *Mathematical proceedings of the cambridge philosophical society*, volume 119, pages 447–468. Cambridge University Press, 1996.
- [KLLWM24] Nikhil Khatri, Tuomas Laakkonen, Jonathon Liu, and Vincent Wang-Maścianica. On the anatomy of attention. *arXiv preprint arXiv:2407.02423*, 2024.
- [Lam99] Joachim Lambek. Type grammar revisited. In *Logical Aspects of Computational Linguistics: Second International Conference, LACL’97 Nancy, France, September 22-24, 1997 Selected Papers 2*, pages 1–27. Springer, 1999.
- [LCRS22] Vishnu Suresh Lokhande, Rudrasis Chakraborty, Sathya N Ravi, and Vikas Singh. Equivariance allows handling multiple nuisance variables when analyzing pooled neuroimaging datasets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10432–10441, 2022.
- [Lew19] Martha Lewis. Compositionality for recursive neural networks. *arXiv preprint arXiv:1901.10723*, 2019.
- [Mar14] Daniel Marsden. Category theory using string diagrams. *arXiv preprint arXiv:1401.7220*, 2014.
- [Mar19] Yoshihiro Maruyama. Compositionality and contextuality: the symbolic and statistical theories of meaning. In *Modeling and Using Context: 11th International and Interdisciplinary Conference, CONTEXT 2019, Trento, Italy, November 20–22, 2019, Proceedings 11*, pages 161–174. Springer, 2019.
- [Nak23] Kenji Nakahira. Diagrammatic category theory. *arXiv preprint arXiv:2307.08891*, 2023.
- [OCV⁺20] Chris Olah, Nick Cammarata, Chelsea Voss, Ludwig Schubert, and Gabriel Goh. Naturally occurring equivariance in neural networks. *Distill*, 5(12):e00024–004, 2020.
- [P⁺71] Roger Penrose et al. Applications of negative dimensional tensors. *Combinatorial mathematics and its applications*, 1:221–244, 1971.

- [PAS24] Samuel Pfrommer, Brendon G Anderson, and Somayeh Sojoudi. Transport of algebraic structure to latent embeddings. *arXiv preprint arXiv:2405.16763*, 2024.
- [PVM⁺21] Wei Peng, Tuomas Varanka, Abdelrahman Mostafa, Henglin Shi, and Guoying Zhao. Hyperbolic deep neural networks: A survey. *IEEE Transactions on pattern analysis and machine intelligence*, 44(12):10023–10044, 2021.
- [Ril18] Mitchell Riley. Categories of optics. *arXiv preprint arXiv:1809.00738*, 2018.
- [SGW21] Dan Shiebler, Bruno Gavranović, and Paul Wilson. Category theory in machine learning. *arXiv preprint arXiv:2106.07032*, 2021.
- [SK19] David Sprunger and Shin-ya Katsumata. Differentiable causal computations via delayed trace. In *2019 34th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS)*, pages 1–12. IEEE, 2019.
- [SL06] Paul Smolensky and Géraldine Legendre. *The harmonic mind: From neural computation to optimality-theoretic grammar*. MIT, 2006.
- [Spi12] David I Spivak. Functorial data migration. *Information and Computation*, 217:31–51, 2012.
- [SXM23] Eli Sennesh, Tom Xu, and Yoshihiro Maruyama. Computing with categories in machine learning. In *International Conference on Artificial General Intelligence*, pages 244–254. Springer, 2023.
- [SY21] Artan Sheshmani and Yi-Zhuang You. Categorical representation learning: morphism is all you need. *Machine Learning: Science and Technology*, 3(1):015016, 2021.
- [Tay24] Jordan K Taylor. An introduction to graphical tensor notation for mechanistic interpretability. *arXiv preprint arXiv:2402.01790*, 2024.
- [VB21] Petar Veličković and Charles Blundell. Neural algorithmic reasoning. *Patterns*, 2(7), 2021.
- [VBB⁺22] Petar Veličković, Adrià Puigdomènech Badia, David Budden, Razvan Pascanu, Andrea Banino, Misha Dashevskiy, Raia Hadsell, and Charles Blundell. The clsr algorithmic reasoning benchmark. In *International Conference on Machine Learning*, pages 22084–22102. PMLR, 2022.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [Wil10] Simon Willerton. Integral transforms and the pull-push perspective, 2010.
- [Wis08] Robert Wisbauer. Algebras versus coalgebras. *Applied Categorical Structures*, 16(1):255–295, 2008.
- [WJL⁺22] Isaac Ronald Ward, Jack Joyner, Casey Lickfold, Yulan Guo, and Mohammed Ben-namoun. A practical tutorial on graph neural networks. *ACM Computing Surveys (CSUR)*, 54(10s):1–35, 2022.

- [WZ21] Paul Wilson and Fabio Zanasi. Reverse derivative ascent: A categorical approach to learning boolean circuits. *arXiv preprint arXiv:2101.10488*, 2021.
- [WZ22] Paul Wilson and Fabio Zanasi. Categories of differentiable polynomial circuits for machine learning. In *International Conference on Graph Transformation*, pages 77–93. Springer, 2022.
- [XKM23] Yao Lei Xu, Kriton Konstantinidis, and Danilo P Mandic. Graph tensor networks: An intuitive framework for designing large-scale neural learning systems on multiple domains. *arXiv preprint arXiv:2303.13565*, 2023.
- [XLZ⁺19] Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? *arXiv preprint arXiv:1905.13211*, 2019.
- [Zag24] Olga Zaghen. Nonlinear sheaf diffusion in graph neural networks. *arXiv preprint arXiv:2403.00337*, 2024.
- [ZPIE17] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

Acknowledgements

I wish to acknowledge the essential role my advisor Professor F. Zanasi had in guiding me through the process of writing this thesis. He introduced me to applied category theory, to machine learning, and to the world of academic research, and I will be forever grateful for it. I also wish to thank my family and my friends, who supported me throughout my academic journey, and without whom I would not have been here writing. To all of you, my most sincere gratitude.