

UNIVERSITY OF ZAGREB  
FACULTY OF ELECTRICAL ENGINEERING AND COMPUTING

arXiv:1907.08292v1 [cs.LG] 16 Jul 2019

MASTER THESIS num. 2034

# **Compositional Deep Learning**

Bruno Gavranović

Zagreb, July 2019.

*I've had an amazing time these last few years. I've had my eyes opened to a new and profound way of understanding the world and learned a bunch of category theory. This thesis was shaped with help of a number of people who I owe my gratitude to. I thank my advisor Jan Šnajder for introducing me to machine learning, Haskell and being a great advisor throughout these years. I thank Martin Tutek and Siniša Šegvić who have been great help for discussing matters related to deep learning and for proofchecking early versions of these ideas. David Spivak has generously answered many of my questions about categorical concepts related to this thesis. I thank Alexander Poddubny for stimulating conversations and valuable theoretic insights, and guidance in thinking about these things without whom many of the constructions in this thesis would not be in their current form. I also thank Mario Roman, Carles Sáez, Ammar Husain, Tom Gebhart for valuable input on a rough draft of this thesis.*

*Finally, I owe everything I have done to my brother and my parents for their unconditional love and understanding throughout all these years. Thank you.*

# CONTENTS

<b>1. Introduction</b>	<b>1</b>
<b>2. Background</b>	<b>3</b>
2.1. Neural Networks . . . . .	3
2.2. Category Theory . . . . .	4
2.3. Database Schemas as Neural Network Schemas . . . . .	6
2.4. Outline of the Main Results . . . . .	9
<b>3. Categorical Deep Learning</b>	<b>10</b>
3.1. Model Schema . . . . .	10
3.2. What Is a Neural Network? . . . . .	11
3.3. Parametrization and Architecture . . . . .	13
3.4. Architecture Parameter Space . . . . .	16
3.5. Data . . . . .	19
3.6. Optimization . . . . .	21
3.7. From Categorical Databases to Deep Learning . . . . .	27
<b>4. Examples</b>	<b>28</b>
4.1. Existing Architectures . . . . .	28
4.2. Product Task . . . . .	29
<b>5. Experiments</b>	<b>33</b>
5.1. Circles . . . . .	33
5.2. CelebA . . . . .	36
5.3. StyleGAN . . . . .	39
5.4. Experiment Summary . . . . .	46
<b>6. Conclusion</b>	<b>47</b>
<b>Bibliography</b>	<b>48</b>

# 1. Introduction

Our understanding of intelligent systems which augment and automate various aspects of our cognition has seen rapid progress in recent decades. Partially prompted by advances in hardware, the field of study of multi-layer artificial neural networks – also known as *deep learning* – has seen astonishing progress, both in terms of theoretical advances and practical integration with the real world. Just as mechanical muscles spawned by the Industrial revolution automated significant portions of human manual labor, so are mechanical minds brought forth by modern deep learning showing the potential to automate aspects of cognition and pattern recognition previously thought to have been unique only to humans and animals.

In order to design and scale such sophisticated systems, we need to take extra care when managing their complexity. As with all modern software, their design needs to be done with the principle of *compositionality* in mind. Although at a first glance it might seem like an interesting yet obscure concept, the notion of compositionality is at the heart of modern computer science, especially type theory and functional programming.

Compositionality describes and quantifies how complex things can be assembled out of simpler parts. It is a principle which tells us that the design of abstractions in a system needs to be done in such a way that we can intentionally forget their internal structure (Hedges, 2017). This is tightly related to the notion of a leaky abstraction (Spolsky, 2002) – a system whose internal design affects its users in ways not specified by its interface.

Indeed, going back to deep learning, we observe two interesting properties of neural networks related to compositionality: (i) they are compositional – increasing the number of layers tends to yield better performance, and (ii) they are discovering (compositional) structures in data. Furthermore, an increasing number of components of a modern deep learning system is learned. For instance, Generative Adversarial Networks (Goodfellow et al., 2014) learn the *cost function*. The paper *Learning to Learn by gradient descent by gradient descent* (Andrychowicz et al., 2016) specifies networks that learn the *optimization function*. The paper *Decoupled Neural Interfaces using Synthetic Gradients* (Jaderberg et al., 2016) specifies how gradients themselves can be learned. The neural network system in (Jaderberg et al., 2016) can be thought of as a cooperative multi-player game, where some players depend on other ones to learn but can be trained in an asynchronous manner.



These are just rough examples, but they give a sense of things to come. As more and more components of these systems stop being fixed throughout training, there is an increasingly larger need for more precise formal specification of the things that *do* stay fixed. This is not an easy task; the invariants across all these networks seem to be rather abstract and hard to describe.

In this thesis we explore the hypothesis that category theory – a formal language for describing general abstract structures in mathematics – could be well suited to describe these systems in a precise manner. In what follows we lay out the beginnings of a formal compositional framework for reasoning about a number of components of modern deep learning architectures. As such the general aim of this thesis is twofold. Firstly, we translate a collection of abstractions known to machine learning practitioners into the language of category theory. By doing so we hope to uncover and disentangle some of the rich conceptual structure underpinning gradient-based optimization and provide mathematicians with some interesting new problems to solve. Secondly, we use this abstract category-theoretic framework to conceive a new and practical way to train neural networks and to perform a novel task of object deletion and insertion in images with unpaired data.

The rest of the thesis is organized as follows. In Chapter 2 we outline some recent work in neural networks and provide a sense of depth to which category theory is used in this thesis. We also motivate our approach by noting a surprising correspondence between a class of neural network architectures and database systems. Chapter 3 contains the meat of the thesis and most of the formal categorical structure. We provide a notion of generalization of parametrization using the construction we call **Para**. Similarly, provision of categorical analogues of neural network architectures as functors allows us to generalize parameter spaces of these network architectures to parameter space of *functors*. This chapter concludes with a description of the optimization process in such a setting. In Chapter 4 we show how existing neural network architectures fit into this framework and we conceive a novel network architecture. In the final chapter we report our experiments of this novel neural network architecture on some concrete datasets.

## 2. Background

In this chapter we give a brief overview of the necessary background related to neural networks and category theory, along with an outline of the used notation and terminology. In Section 2.3 we motivate our approach by informally presenting categorical database systems as they pertain to the topic of this thesis. Lastly, we outline the main results of the thesis.

### 2.1. Neural Networks

Neural networks have become an increasingly popular tool for solving many real-world problems. They are a general framework for differentiable optimization which includes many other machine learning approaches as special cases.

Recent advances in neural networks describe and quantify the process of discovering high-level, abstract structure in data using gradient information. As such, learning inter-domain mappings has received increasing attention in recent years, especially in the context of *unpaired data* and image-to-image translation (Zhu et al., 2017; Almahairi et al., 2018). *Pairedness* of datasets  $X$  and  $Y$  generally refers to the existence of some invertible function  $X \rightarrow Y$ . Note that in classification we might also refer to the input dataset as being paired with the dataset of labels, although the meaning is slightly different as we cannot obviously invert a label  $f(x)$  for some  $x \in X$ .

Obtaining this pairing information for datasets usually requires considerable effort. Consider the task of object removal from images; obtaining pairs of images where one of them lacks a certain object, with everything else the same, is much more difficult than the mere task of obtaining two sets of images: one that contains that object and one that does not, with everything in these images possibly varying. Moreover, we further motivate this example by the reminiscence of the way humans reason about the missing object: simply by observing two unpaired sets of images, where we are told one set of images lack an object, we are able to learn how the missing object looks like.

There is one notable neural network architecture related to generative modelling – Generative Adversarial Networks (GANs) (Goodfellow et al., 2014). Generative Adversarial Networks present a radically different approach to training neural networks. A GAN is a

generative model which is a composition of two networks, one called *the generator* and one called *the discriminator*. Instead of having the cost function fixed, GAN *learns* the cost function using the discriminator. The generator and the discriminator have opposing goals and are trained in an alternating fashion where both continue improving until the generator learns the underlying data distribution. GANs show great potential for the development of accurate generative models for complex distributions, such as the distribution of images of written digits or faces. Consequently, in just a few years, GANs have grown into a major topic of research in machine learning.

Motivated by the success of Generative Adversarial Networks in image generation, existing unsupervised learning methods such as CycleGAN (Zhu et al., 2017) and Augmented CycleGAN (Almahairi et al., 2018) use adversarial losses to learn the true data distribution of given domains of natural images and *cycle-consistency* losses to learn *coherent* mappings between those domains. CycleGAN is an architecture which learns a one-to-one mapping between two domains. Each domain has an associated *discriminator*, while the mappings between these domains correspond to *generators*. The generators in CycleGAN are a collection of neural networks which is closed under composition, and whose inductive bias is increased by enforcing composition invariants, i.e. cycle-consistencies.

Augmented CycleGAN (Almahairi et al., 2018) notices that most relationships across domains are more complex than simple isomorphisms. It is a generalization of CycleGAN which learns *many-to-many* mappings between two domains. Augmented CycleGAN augments each domain with an auxiliary latent variable and extends the training procedure to these augmented spaces.

## 2.2. Category Theory

This work builds on the foundations of a branch of mathematics called Category theory. Describing category theory in just a few paragraphs is not an easy task as there exist a large number of equally valid vantage points to observe it from (Sobocinski, 2017). Rather, we give some intuition and show how it is becoming an unifying force throughout sciences (Baez and Stay, 2009), in all the places in which we need to reason about compositionality.

First and foremost, category theory is a language - a rigorous and a formal one. We mean this in the full definition of the word *language* – it enables us to specify and communicate complex ideas in a succinct manner. Just as any language – it *guides and structures thought*.

It is a toolset for describing general abstract structures in mathematics. Called also “the architecture of mathematics” (Cheng, 2000), it can be regarded as the *theory of theories*, a tool for organizing and layering abstractions and finding formal connections between seemingly disparate fields (Fong and Spivak, 2018). Originating in algebraic topology, it has not

been designed with the compositionality in mind. However, category theory seems to be deeply rooted in all the places we need to reason about composition.

As such, category theory is slowly finding applications outside of the realm of pure mathematics. Some categorical structures have been emerging across the sciences: in Bayesian networks (Fong, 2013; Culbertson and Sturtz, 2013), database systems (Fleming et al., 2003; Rosebrugh and Wood, 1992; Spivak, 2010; Schultz et al., 2016), version control systems (Mimram and Di Giusto, 2013), type theory (Jacobs, 1999), electric circuits (Baez and Fong, 2015), natural language processing (Coecke et al., 2010; Bradley et al., 2018), game theory (Ghani et al., 2016; Hedges and Lewis, 2018; Hedges, 2017b,a), and automatic differentiation (Elliott, 2018), not to mention its increased use in quantum physics (Coecke and Kissinger, 2017; Abramsky and Coecke, 2004, 2008).

In the context of this thesis we focus on categorical formulations of neural networks (Fong et al., 2017; Harris, 2019; Elliott, 2018) and databases (Spivak, 2010). Perhaps the most relevant to this thesis is compositional formulation of supervised learning found in (Fong et al., 2017), whose construction **Para** we generalize in Section 3.2.

### 2.2.1. Assumptions, Notation, and Terminology

We assume a working knowledge of fundamental category theory. Although most of the notation we use is standard, we outline some of the conventions here.

For any category  $\mathcal{C}$  we denote the set of its objects with  $Ob(\mathcal{C})$  and individual objects using uppercase letters such as  $A$ ,  $B$ , and  $C$ . The hom-set of morphisms between two objects  $A$  and  $B$  in a category  $\mathcal{C}$  is written as  $\mathcal{C}(A, B)$ . When we want to consider a discretization of a category  $\mathcal{C}$  such that the only morphisms are the identity morphisms we will write  $|\mathcal{C}|$ . We use  $A \subseteq B$  to denote  $A$  is a subset of  $B$ , but also more generally to denote  $A$  is a subobject of  $B$ . Given some function  $f : P \times A \rightarrow B$  and a  $p \in P$ , we write a partial application of the  $p$  to the first argument of  $f$  as  $f(p, -) : A \rightarrow B$ .

Given categories  $\mathcal{C}$  and  $\mathcal{D}$ , we write  $\mathcal{D}^{\mathcal{C}}$  for the functor category whose objects are functors  $\mathcal{C} \rightarrow \mathcal{D}$  and morphisms are natural transformations between such functors.

When talking about a monoidal category  $\mathcal{C}$  we will use  $\otimes : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$  for the monoidal product,  $I \in Ob(\mathcal{C})$  for the unit object,  $\alpha_{A,B,C} : (A \otimes B) \otimes C \rightarrow A \otimes (B \otimes C)$  for the associator, and  $\lambda_A : I \otimes A \rightarrow A$  for the left unitor.

A notable category we use is **Euc**, the strict symmetric monoidal category whose objects are finite-dimensional Euclidean spaces and morphisms are differentiable maps. A monoidal product on **Euc** is given by the Cartesian product.

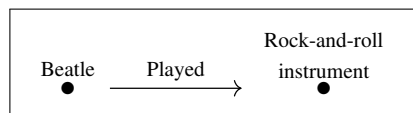
Given a directed multigraph  $G$ , we will write **Free**( $G$ ) for the free category on that graph  $G$  and **Free**( $G$ )/ $\sim$  for its quotient category by some congruence relation  $\sim$ . Lastly, given a category  $\mathcal{C}$  with generators  $G$ , we write  $Gen_{\mathcal{C}}$  for the set of generating morphisms  $G$  in  $\mathcal{C}$ .

## 2.3. Database Schemas as Neural Network Schemas

In this section we motivate our approach by highlighting a remarkable correspondence between database systems, as defined by Spivak (2010), and a class of neural network architectures, here exemplified by CycleGAN (Zhu et al., 2017), but developed in generality in Chapter 3.

Our aim in this section is to present an informal, high-level overview of this correspondence. We hope the emphasized intuition in this section serves as a guide for Chapter 3 where these structures are described in a formal manner.

The categorical formulation of databases found in (Spivak, 2010) can roughly be described as follows. A database is modelled as a category which holds just the abstract relationship between concepts, and a structure-preserving map into another category which holds the actual data corresponding to those concepts. That is, a database *schema*  $\mathcal{C}$  specifies a reference structure for a given database instance. An example, adapted from Fong and Spivak (2018), is shown in Figure 2.1.



**Figure 2.1:** Toy example of a database schema

Actual data corresponding to such a schema is a functor  $\mathcal{C} \rightarrow \mathbf{Set}$ , shown in Figure 2.1 as a system of interlocking tables.

<b>Beatle</b>	<b>Played</b>	<b>Rock-and-roll instrument</b>
George	Lead guitar	Bass guitar
John	Rhythm guitar	Drums
Paul	Bass guitar	Keyboard
Ringo	Drums	Lead guitar
		Rhythm guitar

**Table 2.1:** Toy example of a database instance corresponding to the schema in Figure 2.1

Observe the following: The actual data – sets and functions between them – are available or known beforehand. There might be some missing data, but all functions usually have well-defined implementations.

We contrast this with *machine learning*, where we might have plenty of data, but no known implementation of functions that map between data samples. Table 2.2 shows an example from the setting of supervised learning. In this example, samples are paired: for every input we have an expected output. Thus, given a trained model and a new sample from a test set – say “DataSample4717” – we hope our model has learned to generalize well and assign a corresponding output to this input.

<b>Input</b>	<b>Corresponding output</b>	<b>Output</b>
DataSample1	ExpectedOutput1	ExpectedOutput1
DataSample17	ExpectedOutput17	ExpectedOutput17
DataSample30	ExpectedOutput30	ExpectedOutput30
DataSample400	ExpectedOutput400	ExpectedOutput400
...		...

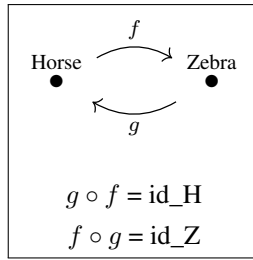
**Table 2.2:** Example of paired datasets in a setting of supervised learning

Moreover, we point out that the luxury of having paired data is not always at our disposal: real life data is mostly *unpaired*. We might have two datasets that are related *somehow*, but without knowing if any inputs match any of the outputs. Consider the example shown in Table 2.3.

<b>Horse image</b>	<b>Horse → Zebra</b>	<b>Zebra image</b>	<b>Zebra → Horse</b>
HorseImg1	?	ZebraImg10	?
HorseImg24	?	ZebraImg430	?
⋮		⋮	
HorseImg303	?	ZebraImg566	?
HorseImg2392	?	ZebraImg637	?
		ZebraImg700	?
		⋮	

**Table 2.3:** An example of two unpaired datasets

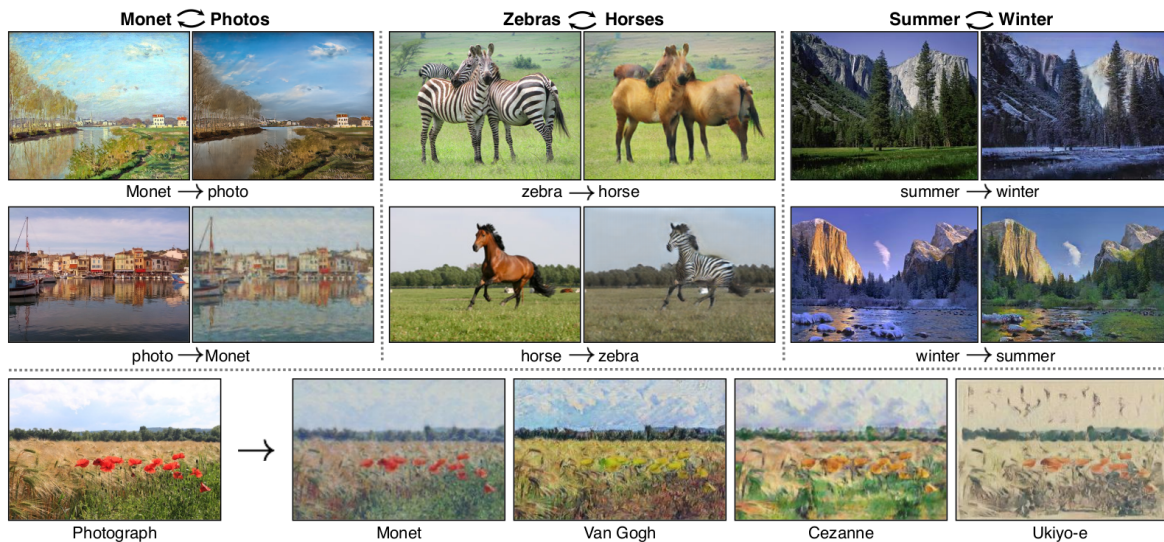
The example in Table 2.3 is given by the schema in Figure 3.1b. They depict the scenario where we have two image datasets: of images of horses and of images of zebras. Consider those images as photographs of these animals in nature, in various positions and from various angles (Figure 2.3). We just have *some* images and do not necessarily have any pairs of images in our dataset.



**Figure 2.2:** We might hypothesize mappings  $f$  and  $g$  exist – without knowing anything about them other than their composition invariants.

Although our dataset does not contain any explicit Horse – Zebra pairs – we still might hypothesize that these pairs exist. In other words, we could think that it should be possible to map back-and-forth between images of Horses and Zebras, just by changing the texture of animal in such an image. In other words, we posit there exists a specific relationship between the datasets. Compared to *databases*, where we have the data and well-defined function implementations, here all we have is data and a *rough idea* of which mappings exist, without known implementations. The issue is that our dataset does not contain explicit pairs usable in the context of supervised learning.

What is described here is first introduced in a paper by Zhu et al. (2017). They introduce a model CycleGAN which is a generalization of Generative Adversarial Networks (Goodfellow et al., 2014). Figure 2.3 is adapted from their paper and showcases the main results.



**Figure 2.3:** Given any two unordered image collections  $X$  and  $Y$ , CycleGAN learns to automatically “translate” an image from one into the other and vice versa. Figure taken from Zhu et al. (2017).

This shows us that, at least at a first glance, CycleGAN and categorical databases are related in some abstract way. After developing the necessary categorical tools to reason about CycleGAN, we elaborate on this correspondence in Section 3.7.

## 2.4. Outline of the Main Results

This thesis aims to bridge two seemingly distinct ideas: category theory and deep learning. In doing so we take a collection of abstractions in deep learning and formalize the notation in categorical terms. This allows us to begin to consider a formal theory of gradient-based learning in the *functor space*.

We package a notion of the interconnection of networks as a free category  $\mathbf{Free}(G)$  on some graph  $G$  and specify any equivalences between networks as relations between morphisms as a quotient category  $\mathbf{Free}(G)/\sim$ . Given such a category – which we call a *schema*, inspired by Spivak (2010) – we specify the architectures of its networks as a functor  $\text{Arch}$ . We reason about various other notions found in deep learning, such as datasets, embeddings, and parameter spaces. The training process is associated with an indexed family of functors  $\{H_{p_i} : \mathbf{Free}(G) \rightarrow \mathbf{Set}\}_{i=1}^T$ , where  $T$  is the number of training steps and  $p_i$  is some choice of a parameter for that architecture at the training step  $i$ .

Analogous to standard neural networks – we start with a randomly initialized  $H_p$  and iteratively update it using gradient descent. Our optimization is guided by *two* objectives. These objectives arise as a natural generalization of those found in (Zhu et al., 2017). One of them is the adversarial objective – the minmax objective found in any Generative Adversarial Network. The other one is a generalization of the cycle-consistency loss which we call *path-equivalence loss*.

Although mathematically abstract, this approach yields useful insights. Our formulation provides maximum generality: (i) it enables learning with unpaired data as it does not impose any constraints on ordering or pairing of the sets in a category, and (ii) although specialized to generative models in the domain of computer vision, the approach is domain-independent and general enough to hold in any domain of interest, such as sound, text, or video.

We show that for specific choices of  $\mathbf{Free}(G)/\sim$  and the dataset we recover GAN (Goodfellow et al., 2014) and CycleGAN (Zhu et al., 2017). Furthermore, a novel neural network architecture capable of learning to remove and insert objects into an image with unpaired data is proposed. We outline its categorical perspective and show it in action by testing it on three different datasets.



## 3. Categorical Deep Learning

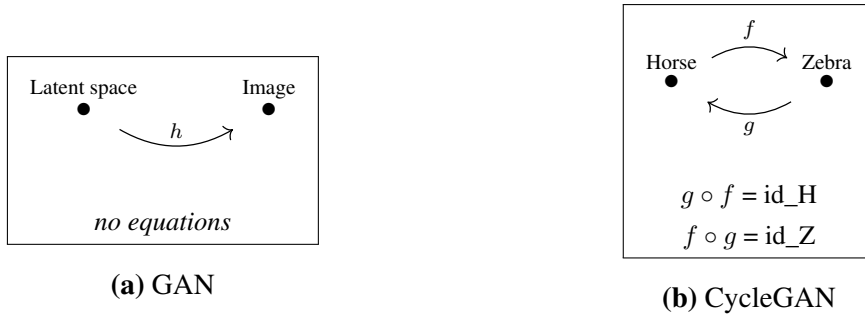
Modern deep learning optimization algorithms can be framed as a gradient-based search in some function space  $Y^X$ , where  $X$  and  $Y$  are sets that have been endowed with extra structure. Given some sets of data points  $D_X \subseteq X$ ,  $D_Y \subseteq Y$ , a typical approach for adding inductive bias relies on exploiting this extra structure associated to the data points embedded in those sets, or those sets themselves. This structure includes domain-specific features which can be exploited by various methods – convolutions for images, Fourier transform for audio, and specialized word embeddings for textual data.

In this chapter we develop the categorical tools to increase inductive bias of a model without relying on any extra such structure. We build on top of the work of (Fong et al., 2017; Spivak, 2010) and, very roughly, define our model as a *collection of networks* and increase its inductive bias by *enforcing their composition invariants*.

### 3.1. Model Schema

Many deep learning models are complex systems, some comprised of several neural networks. Each neural network can be identified with domain  $X$ , codomain  $Y$ , and a *differentiable parametrized function*  $X \rightarrow Y$ . Given a *collection* of such networks, we use a directed multigraph to capture their interconnections. We use vertices to represent the domains and codomains, and edges to represent differentiable parametrized functions. Observe that an ordinary graph will not suffice, as there can be two *different* differentiable parametrized functions with the same domain and codomain.

Each directed multigraph  $G$  gives rise to a corresponding free category on that graph  $\mathbf{Free}(G)$ . Based on this construction, Figure 3.1 shows the interconnection pattern for generators of two popular neural network architectures: GAN (Goodfellow et al., 2014) and CycleGAN (Zhu et al., 2017).



**Figure 3.1:** Bird's-eye view of two popular neural network models

Observe that CycleGAN has some additional properties imposed on it, specified by equations in Figure 3.1 (b). These are called CycleGAN cycle-consistency conditions and can roughly be stated as follows: given domains  $A$  and  $B$  considered as sets,  $a \approx g(f(a)), \forall a \in A$  and  $b \approx f(g(b)), \forall b \in B$ .

Our approach involves a realization that cycle-consistency conditions *can be generalized* to path equivalence relations, or, in formal terms – a congruence relation. The condition  $a \approx g(f(a)), \forall a \in A$  can be reformulated such that it does not refer to the elements of the set  $a \in A$ . By *eta-reducing* the equation we obtain  $id_a = g \circ f$ . Similar reformulation can be done for the other condition:  $id_b = f \circ g$ .

This allows us to package the newly formed equations as equivalence relations on the sets  $\mathbf{Free}(G)(A, A)$  and  $\mathbf{Free}(G)(B, B)$ , respectively. This notion can be further packaged into a quotient category  $\mathbf{Free}(G)/\sim$ , together with the quotient functor  $\mathbf{Free}(G) \xrightarrow{Q} \mathbf{Free}(G)/\sim$ .

This formulation – as a free category on a graph  $G$  – represents the cornerstone of our approach. These schemas allow us to reason solely about the interconnections between various concepts, rather than jointly with functions, networks or other some other sets. All the other constructs in this thesis are structure-preserving maps between categories whose domain, roughly, can be traced back to  $\mathbf{Free}(G)$ .

## 3.2. What Is a Neural Network?

In computer science, the idea of a *neural network* colloquially means a number of different things. At a most fundamental level, it can be interpreted as a system of interconnected units called neurons, each of which has a firing threshold acting as an information filtering system. Drawing inspiration from biology, this perspective is thoroughly explored in literature. In many contexts we want to focus on the mathematical properties of a neural network and as such identify it with a function between sets  $A \xrightarrow{f} B$ . Those sets are often considered to have extra structure, such as those of Euclidean spaces or manifolds. Functions are then considered to be maps of a given differentiability class which preserve such structure. We

also frequently reason about a neural network jointly with its parameter space  $P$  as a function of type  $f : P \times A \rightarrow B$ . For instance, consider a classifier in the context of supervised learning. A convolutional neural network whose input is a  $32 \times 32$  RGB image and output is real number can be represented as a function with the following type:  $\mathbb{R}^n \times \mathbb{R}^{32 \times 32 \times 3} \rightarrow \mathbb{R}$ , for some  $n \in \mathbb{N}$ . In this case  $\mathbb{R}^n$  represents the parameter space of this network.

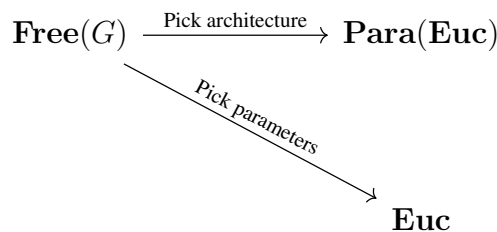
The former ( $A \rightarrow B$ ) and the latter ( $P \times A \rightarrow B$ ) perspective on neural networks are related. Namely, consider some function space  $B^A$ . Any notion of smoothness in such a space is not well defined without any further assumptions on sets  $A$  or  $B$ . This is the reason deep learning employs a gradient-based search in such a space via a proxy function  $P \times A \rightarrow B$ . This function specifies an entire *parametrized family* of functions of type  $A \rightarrow B$ , because partial application of each  $p \in P$  yields a function  $f(p, -) : A \rightarrow B$ . This choice of a parametrized family of functions is part of the *inductive bias* we are building into the training process. For example, in computer vision it is common to restrict the class of functions to those that can be modeled by convolutional neural networks.

In more general terms, by currying  $f : P \times A \rightarrow B$  we obtain an important construction in the literature as a parameter-function map  $\mathcal{M} : P \rightarrow B^A$  (Valle-Perez et al., 2019).

The parameter-function map is important as it allows us to map behaviors in the parameter space to the behavior of functions in the function space. For example, partial differentiation of  $f$  with respect to  $p$  allows us to use gradient information to search the parameter space – but also the space of a particular family of functions of type  $A \rightarrow B$  specified by  $f$ .

With this in mind, we recall the model schema. For each morphism  $A \rightarrow B$  in  $\mathbf{Free}(G)$  we are interested in specifying a parametrized function  $f : P \times A \rightarrow B$ , i.e. a parametrized *family of functions* in  $\mathbf{Set}$ . The construction  $f$  describes a neural network architecture, and a choice of a partially applied  $p \in P$  to  $f$  describes a choice of some parameter value for that specific architecture.

We capture the notion of parametrization with a construction  $\mathbf{Para}$ . We package both of these notions – choosing an architecture and choosing parameters – into functors. Figure 3.2 shows a high-level overview of these constructions – including a notion that will be central to this thesis –  $\mathbf{Para}(\mathbf{Euc})$ .



**Figure 3.2:** High-level structure of architecture and parameter selection

### 3.3. Parametrization and Architecture

In this section we begin to provide a rigorous categorical framework for reasoning about neural networks.

#### 3.3.1. Parametrization

We now turn our attention to a construction  $\mathbf{Para}$ , which allows us to compose parametrized functions of type  $f : P \times A \rightarrow B$  in such a way that it abstracts away the notion of a parameter. This is a generalization of the construction  $\mathbf{Para}$  found in (Fong et al., 2017).

**Definition 1** ( $\mathbf{Para}$ ). *Given any small symmetric monoidal category  $(\mathcal{V}, I, \otimes)$ , we can construct another small symmetric monoidal category  $(\mathbf{Para}(\mathcal{V}), I, \otimes)$  given by the following:*

- **Objects.**  $\mathbf{Para}(\mathcal{V})$  has the same objects as  $\mathcal{V}$ ;
- **Morphisms.**  $\mathbf{Para}(\mathcal{V})(A, B) := \{f : P \otimes A \rightarrow B \mid P \in \text{Ob}(\mathcal{V})\} / \sim$ , where  $f \sim f'$  if there exists an isomorphism  $g \in \mathcal{V}(P, P')$  such that  $f' \circ (g \otimes \text{id}_A) = f$ ;
- **Identity.** Identity of any object  $A \in \text{Ob}(\mathbf{Para}(\mathcal{V}))$  is the left unitor  $\lambda_A : I \otimes A \rightarrow A$ ;
- **Composition.** For every three objects  $A, B, C$  and morphisms  $f : P \otimes A \rightarrow B \in \mathbf{Para}(\mathcal{V})(A, B)$  and  $g : Q \otimes B \rightarrow C \in \mathbf{Para}(\mathcal{V})(B, C)$  for some  $P, Q \in \text{Ob}(\mathcal{V})$ , we specify a morphism  $g \circ f \in \mathbf{Para}(\mathcal{V})(A, C)$  as follows:

$$\begin{aligned} g \circ f &: (P \otimes Q) \otimes A \rightarrow C \\ g \circ f &= \lambda((p, q), a) \rightarrow g(q, f(p, a)) \end{aligned}$$

*Monoidal structure is inherited from  $\mathcal{V}$ .*

*Proof.* To prove  $\mathbf{Para}(\mathcal{V})$  is indeed a category, we need to show associativity and unitality of composition strictly holds. Observe that composition in  $\mathbf{Para}(\mathcal{V})$  is defined in terms of the monoidal product in  $\mathcal{V}$ . Consider the two different ways of composing following morphisms in  $\mathbf{Para}(\mathcal{V})$ :

$$\begin{aligned} f &: P \otimes A \rightarrow B \in \mathbf{Para}(\mathcal{V})(A, B), \\ g &: Q \otimes B \rightarrow C \in \mathbf{Para}(\mathcal{V})(B, C), \\ h &: R \otimes C \rightarrow D \in \mathbf{Para}(\mathcal{V})(C, D). \end{aligned}$$

Depending on the order we compose them, we end up with one of the two following morphisms:

$$\begin{aligned} h \circ (g \circ f) &: ((P \otimes Q) \otimes R) \otimes A \rightarrow D, \\ (h \circ g) \circ f &: (P \otimes (Q \otimes R)) \otimes A \rightarrow D. \end{aligned}$$

Even though it might seem strictness of the associativity of composition in  $\mathbf{Para}(\mathcal{V})$  depends on the strictness of the monoidal product in  $\mathcal{V}$ , we note that morphisms in  $\mathbf{Para}(\mathcal{V})$  are actually equivalence classes. Namely, because every monoidal category comes equipped with an associator  $\alpha_{P,Q,R} : (P \otimes Q) \otimes R \cong P \otimes (Q \otimes R)$ , both  $h \circ (g \circ f)$  and  $(h \circ g) \circ f$  fall into the same equivalence class, making composition in  $\mathbf{Para}(\mathcal{V})$  strictly associative.

Similar argument can be made for the unitality condition, thus showing  $\mathbf{Para}(\mathcal{V})$  is a category.  $\square$

We will call morphisms in  $\mathbf{Para}(\mathcal{V})$  *parametrized morphisms*, *neural networks*, or *neural network architectures* depending on the context.<sup>1</sup> Abusing our notation slightly, we will refer to a morphism in  $\mathbf{Para}(\mathcal{V})$  by some elements from the corresponding equivalence class.

The composition of morphisms in  $\mathbf{Para}(\mathcal{V})$  is defined in such a way that it explicitly keeps track of parameters. Namely, when we sequentially compose two morphisms  $A \xrightarrow{f} B$  and  $B \xrightarrow{g} C$  in  $\mathbf{Para}(\mathcal{V})$ , we are actually composing morphisms  $P \otimes A \rightarrow B$  and  $Q \otimes B \rightarrow C$  in  $\mathcal{V}$  such that the composition  $(P \otimes Q) \otimes A \rightarrow C$  keeps track of parameters separately.<sup>2</sup>

This construction  $\mathbf{Para}(\mathcal{V})$  generalizes the category  $\mathbf{Para}$  as originally defined in Fong et al. (2017). Namely, by setting  $\mathcal{V} := \mathbf{Euc}$ , we recover the notion  $\mathbf{Para}$  as it is described in the aforementioned paper. As  $\mathbf{Para}(\mathbf{Euc})$  will make continued appearance in this thesis, we describe some of its properties here.

$\mathbf{Para}(\mathbf{Euc})$  is a strict symmetric monoidal category whose objects are Euclidean spaces and morphisms  $\mathbb{R}^n \rightarrow \mathbb{R}^m$  are equivalence classes of differentiable function of type  $\mathbb{R}^p \times \mathbb{R}^n \rightarrow \mathbb{R}^m$ , for some  $p \in \mathbb{N}$ . We refer to  $\mathbb{R}^p$  as the parameter space.

Monoidal product in  $\mathbf{Para}(\mathbf{Euc})$  is the Cartesian product inherited from  $\mathbf{Euc}$ . As maps in  $\mathbf{Euc}$  are differentiable, so are maps in  $\mathbf{Para}(\mathbf{Euc})$  thus enabling us to consider gradient-based optimization in a more abstract setting.

## Parameter selection in a monoidal category

Previously, in the context of functions  $f : P \times A \rightarrow B$  between sets, we have considered the partial application  $f(p, -)$ . Now, we are interested in doing the same in  $\mathbf{Para}(\mathcal{V})$ , given some small symmetric monoidal category  $\mathcal{V}$ .

There are two issues with this statement.

---

<sup>1</sup> Note that  $\mathbf{Para}$  is not a functor between categories  $\mathcal{V}$  and  $\mathbf{Para}(\mathcal{V})$ , but rather an endofunctor on  $\mathbf{SMC}_{\text{str}}$ , the category of all small symmetric monoidal categories. We outline this for completeness but do not prove or explore this direction further.

<sup>2</sup> Even though objects in  $\mathbf{Para}(\mathcal{V})$  generally do not have elements, in Definition 1 composition is stated in terms of elements to supply intuition. We have also bracketed the monoidal product even though  $\otimes$  is strict for a similar reason – to show we think of  $P \otimes Q$  as the parameter of the composite  $g \circ f$ . We invite the reader to (Fong et al., 2017), which contains a particularly clean interpretation of  $\mathbf{Para}(\mathcal{V})$  in terms of string diagrams.

1. Internal structure of objects in  $\mathcal{V}$  is unknown – they might not be sets with elements
2. It assumes a specific notion of completeness: for every  $f : P \otimes A \rightarrow B$  in  $\mathbf{Para}(\mathcal{V})$  we assume  $f(p, -) \in \mathcal{V}(A, B), \forall p \in P$ .

We solve both of the issues by noting that picking a parameter  $p \in P$  in a monoidal category without any assumption on the internals of its objects amounts to picking a morphism  $I \xrightarrow{p} P$ . Then the specific notion of completeness *is already given to us* by the monoidal product on  $\mathcal{V}$ . Indeed,  $f(p, -)$  amounts to the composition  $f \circ (p \otimes id_A) \circ \lambda_A^{-1} : \mathcal{V}(A, B)$ , where  $\lambda^{-1}$  is the inverse of left unitor of the monoidal category.

However, most of our consideration in this thesis will be where  $\mathcal{V} := \mathbf{Euc}$ . In these cases the notation  $p \in P$  is well-justified, as  $\mathbf{Euc}$  comes equipped with a forgetful functor  $\mathbf{Euc} \xrightarrow{U} \mathbf{Set}$ .

### 3.3.2. Model Architecture

We now formally specify *model architecture* as a functor. We chose  $\mathbf{Free}(G)$  as the domain of the functor, rather than  $\mathbf{Free}(G)/\sim$ , for reasons that will be explained in Remark 3. As such, observe that the action on morphisms is defined on the generators in  $\mathbf{Free}(G)$ .

**Definition 2.** *Architecture of a model is a functor  $\mathbf{Arch} : \mathbf{Free}(G) \rightarrow \mathbf{Para}(\mathbf{Euc})$ .*

- For each  $A \in \mathbf{Ob}(\mathbf{Free}(G))$ , it specifies an Euclidean space  $\mathbb{R}^a$ ;
- For each generating morphism  $A \xrightarrow{f} B$  in  $\mathbf{Free}(G)$ , it specifies a morphism  $\mathbb{R}^a \xrightarrow{\mathbf{Arch}(f)} \mathbb{R}^b$  which is a differentiable parametrized function of type  $\mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b$ .

Given a non-trivial composite morphism  $f = f_n \circ f_{n-1} \circ \dots \circ f_1$  in  $\mathcal{C}$ , the image of  $f$  under  $\mathbf{Arch}$  is the composite of the image of each constituent:  $\mathbf{Arch}(f) = \mathbf{Arch}(f_n) \circ \mathbf{Arch}(f_{n-1}) \circ \dots \circ \mathbf{Arch}(f_1)$ .  $\mathbf{Arch}$  maps identities to the projection  $\pi_2 : I \times A \rightarrow A$ .

**Remark 3.** *The reason the domain of  $\mathbf{Arch}$  is  $\mathbf{Free}(G)$ , rather than  $\mathbf{Free}(G)/\sim$  can be illustrated with the following example. Consider two morphisms  $id_A : A \rightarrow A$  and  $g \circ f : A \rightarrow A$  in some  $\mathbf{Free}(G)/\sim$ . Suppose  $id_A = g \circ f$ . The value image of architecture at  $id_A$  is already given:  $\mathbf{Arch}(id_A) := I \otimes A \rightarrow A$ , but for  $g \circ f$  it is defined as  $\mathbf{Arch}(g) \circ \mathbf{Arch}(f)$ . Hence, there exists a choice  $\mathbf{Arch}(g)$  and  $\mathbf{Arch}(f)$  such that  $\mathbf{Arch}(id_A) \neq \mathbf{Arch}(g) \circ \mathbf{Arch}(f)$ , rendering this structure defined on  $\mathbf{Free}(G)/\sim$  not a functor. However, this is not an issue; we will show it will be possible to learn those relations.*

The choice of architecture  $\mathbf{Free}(G) \xrightarrow{\mathbf{Arch}} \mathbf{Para}(\mathbf{Euc})$  goes hand in hand with the choice of an *embedding*.

**Proposition 4.** *An embedding is a functor  $|\mathbf{Free}(G)| \xrightarrow{E} \mathbf{Set}$  which agrees with  $\mathbf{Arch}$  on objects.*

Observe that the codomain of  $E$  is  $\mathbf{Set}$ , rather than  $\mathbf{Para}(\mathbf{Euc})$ . This can be shown in two steps: (i)  $\mathbf{Para}(\mathbf{Euc})$  and  $\mathbf{Euc}$  have the same objects, and (ii) objects in  $\mathbf{Euc}$  are just sets with extra structure.

Embedding  $E$  and  $\mathbf{Arch}$  come up in two different scenarios. Sometimes we start out with a choice of architecture which then induces the embedding. In other cases, the embedding is given to us beforehand and it restricts the possible choice of architectures. The embedding construction will prove to be important later in this thesis.

Having defined  $\mathbf{Free}(G) \xrightarrow{\mathbf{Arch}} \mathbf{Para}(\mathbf{Euc})$ , we shift our attention to the notion of parameter specification. For a given a differentiable parametrized function  $\mathbf{Arch}(f) : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b$  the training process involves repeatedly updating the chosen  $p \in \mathbb{R}^n$ . We might suspect this process of choosing parameters might be made into a functor  $\mathbf{Para}(\mathbf{Euc}) \rightarrow \mathbf{Euc}$ , mapping each  $\mathbf{Arch}(f) : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b$  into  $\mathbf{Arch}(f)(p, -) : \mathbb{R}^a \rightarrow \mathbb{R}^b$ . However, it can be seen that this is not the case by considering fully-connected neural networks; we want to specify parameters for an  $N$ -layer neural network by specifying parameters for each of its layers. In categorical terms, this means that we want to specify the action of this functor on generators in  $\mathbf{Para}(\mathbf{Euc})$ . Since  $\mathbf{Para}(\mathbf{Euc})$  might also have arbitrary relations between morphisms, we cannot be sure this recursive approach will satisfy any such relations. This, in turn, stops us from considering this construction as a functor.

Moreover, even if this construction could be a functor, we show that it might not be the construction we are interested in. Observe that  $\mathbf{Arch}$  is not necessarily faithful. Suppose two different arrows  $A \xrightarrow[f]{g} B$  in  $\mathbf{Free}(G)$  are mapped to the same neural network architecture  $\mathbf{Arch}(f) = \mathbf{Arch}(g) : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b : \mathbf{Para}(\mathbf{Euc})(\mathbf{Arch}(A), \mathbf{Arch}(B))$ . Even though images of these arrows are the same, it is beneficial and necessary to keep in mind that those two are separate neural networks, each of which could have a different parameter assigned to it during training. Any such parameter specification functor whose domain is  $\mathbf{Para}(\mathbf{Euc})$  would have to specify *one* parameter value for such a morphism. This, in turn, means that this construction would not allow us to have two distinct parameters for what we consider to be two distinct networks.

Before coming back to this issue, we take a slight detour and consider various notions of *parameter spaces*.

### 3.4. Architecture Parameter Space

Each network architecture  $f : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b$  comes equipped with its parameter space  $\mathbb{R}^n$ . Just as  $\mathbf{Free}(G) \xrightarrow{\mathbf{Arch}} \mathbf{Para}(\mathbf{Euc})$  is a categorical generalization of architecture, we now show there exists a categorical generalization of a parameter space. In this case – it is the parameter space of a functor. Before we move on to the main definition, we package the

notion of parameter space of a function  $f : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b$  into a simple function  $\mathfrak{p}(f) = \mathbb{R}^n$ .

**Definition 5** (Functor parameter space). *Let  $\text{Gen}_{\mathbf{Free}(G)}$  the set of generators in  $\mathbf{Free}(G)$ . The total parameter map  $\mathcal{P} : \text{Ob}(\mathbf{Para}(\mathbf{Euc})^{\mathbf{Free}(G)}) \rightarrow \text{Ob}(\mathbf{Euc})$  assigns to each functor Arch the product of the parameter spaces of all its generating morphisms:*

$$\mathcal{P}(\text{Arch}) = \prod_{f \in \text{Gen}_{\mathbf{Free}(G)}} \mathfrak{p}(\text{Arch}(f))$$

Essentially, just as  $\mathfrak{p}$  returns the parameter space of a function,  $\mathcal{P}$  does the same for a functor.

We are now in a position to talk about parameter specification. Recall the non-categorical setting: given some network architecture  $f : P \times A \rightarrow B$  and a choice of  $p \in \mathfrak{p}(f)$  we can partially apply the parameter  $p$  to the network to get  $f(p, -) : A \rightarrow B$ . This admits a straightforward generalization to the categorical setting.

**Definition 6** (Parameter specification). *Parameter specification  $\text{PSpec}$  is a dependently typed function with the following signature:*

$$\text{PSpec} : (\text{Arch} : \text{Ob}(\mathbf{Para}(\mathbf{Euc})^{\mathbf{Free}(G)})) \times \mathcal{P}(\text{Arch}) \rightarrow \text{Ob}(\mathbf{Euc}^{\mathbf{Free}(G)})$$

*Given an architecture Arch and a parameter choice  $(p_f)_{f \in \text{Gen}_{\mathbf{Free}(G)}} \in \mathcal{P}(\text{Arch})$  for that architecture, it defines a choice of a functor in  $\mathbf{Euc}^{\mathbf{Free}(G)}$ . This functor acts on objects the same as Arch. On morphisms, it partially applies every  $p_f$  to the corresponding morphism  $\text{Arch}(f) : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b$ , thus yielding  $f(p_f, -) : \mathbb{R}^a \rightarrow \mathbb{R}^b$  in  $\mathbf{Euc}$ .*

Elements of  $\mathbf{Euc}^{\mathbf{Free}(G)}$  will play a central role later on in the thesis. These elements are functors which we will call *Models*. Given some architecture Arch and a parameter  $p \in \mathcal{P}(\text{Arch})$ , a model  $\mathbf{Free}(G) \xrightarrow{\text{Model}_p} \mathbf{Euc}$  generalizes the standard notion of a model in machine learning – it can be used for inference and evaluated.

Analogous to database instances in Spivak (2010), we call a *network instance*  $H_p$  the composition of some  $\text{Model}_p$  with the forgetful functor  $\mathbf{Euc} \xrightarrow{U} \mathbf{Set}$ . That is to say, a network instance is a functor  $\mathbf{Free}(G) \xrightarrow{H_p} \mathbf{Set} := U \circ \text{Model}_p$ .

**Corollary 7.** *Given an architecture  $\mathbf{Free}(G) \xrightarrow{\text{Arch}} \mathbf{Para}(\mathbf{Euc})$ , for each  $p \in \mathcal{P}(\text{Arch})$  all network instances  $\mathbf{Free}(G) \xrightarrow{H_p} \mathbf{Set}$  act the same on objects.*

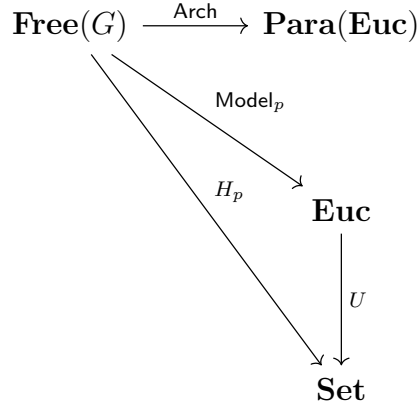
*Proof.* As  $H_p : U \circ \text{Model}_p$ , it is sufficient to prove that for all  $p \in \mathcal{P}(\text{Arch})$  all models act the same on objects. This follows from Definition 6 which tells us that action of  $\text{Model}_p$  on objects is independent of  $p$ .  $\square$

This means that the only thing different between any two network instances during training is the choice of a parameter they partially apply to morphisms in the image of



$\mathbf{Para}(\mathbf{Euc})$ . Objects – the domain of functions resulting from partial applications – stay the same. This is evident in standard machine learning models as well, where we obviously do not change their type of input during architecture parameter updates.

We shed some more light on these constructions using Figure 3.3.



**Figure 3.3:**  $\mathbf{Free}(G)$  is the domain of three types of functors of interest:  $\text{Arch}$ ,  $\text{Model}_p$  and  $H_p$ .

### 3.4.1. Parameter-functor Map.

In this section we show the parameter-function map detailed in Section 3.2 admits a natural generalization to a categorical setting. Recall that given some architecture – a differentiable function between sets  $f : P \times A \rightarrow B$  – we can obtain the *parameter-function map*  $P \rightarrow B^A$ .

In a categorical setting, given an architecture  $\mathbf{Free}(G) \xrightarrow{\text{Arch}} \mathbf{Para}(\mathbf{Euc})$  we can obtain the *parameter-functor map* by partially applying  $\text{Arch}$  to  $\text{PSpec}$ .

**Definition 8** (Parameter-functor map). *Let  $\mathbf{Free}(G) \xrightarrow{\text{Arch}} \mathbf{Para}(\mathbf{Euc})$  be the architecture. Then the parameter-functor map  $\mathcal{M}_{\mathbf{Euc}}$  is the partial application  $\mathcal{M}_{\mathbf{Euc}} := \text{PSpec}(\text{Arch}, -)$  which has the type:*

$$\mathcal{M}_{\mathbf{Euc}} : \mathcal{P}(\text{Arch}) \rightarrow \text{Ob}(\mathbf{Euc}^{\mathbf{Free}(G)}).$$

*More precisely, its values in  $\text{Ob}(\mathbf{Euc}^{\mathbf{Free}(G)})$  are models  $\mathbf{Free}(G) \xrightarrow{\text{Model}_i} \mathbf{Euc}$ . The map  $\mathcal{M}_{\mathbf{Euc}}$  map naturally extends to  $\mathcal{M}_{\mathbf{Set}} : \mathcal{P}(\text{Arch}) \rightarrow \mathbf{Set}^{\mathbf{Free}(G)}$  via the forgetful functor  $\mathbf{Euc} \rightarrow \mathbf{Set}$ . When considering  $\mathcal{M}_{\mathbf{Set}}$  we will simply write  $\mathcal{M}$ .*

**Example 9.** *We show how the parameter-function map arises as a special case of the parameter-functor map. Let  $\mathbf{Free}(G) = \mathbf{A} \xrightarrow{f} \mathbf{B}$ . Consider the architecture  $\mathbf{A} \xrightarrow{f} \mathbf{B} \xrightarrow{\text{Arch}}$*

*$\mathbf{Para}(\mathbf{Euc})$  such that  $\text{Arch}(f) : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b$ . Then  $\mathcal{P}(\text{Arch}) = \mathbb{R}^n$ . The parameter-functor map partially applied to  $\text{Arch}$  has the following type  $\text{PSpec}(\text{Arch}, -) : \mathbb{R}^n \rightarrow$*

$Ob(\mathbf{Euc} \begin{matrix} \bullet & \xrightarrow{f} & \bullet \\ \text{A} & & \text{B} \end{matrix})$ . For each  $p \in \mathbb{R}^p$  it specifies a choice of a functor which can be identified with a function  $\text{Arch}(f)(p, -) : \mathbb{R}^a \rightarrow \mathbb{R}^b$ . This means  $\text{PSpec}(\text{Arch}, -)$  can be identified with a function  $\mathbb{R}^n \rightarrow \mathbb{R}^{b\mathbb{R}^a}$ , thus reducing to the parameter-function map.

Parameter-functor map  $\mathcal{M} : \mathcal{P}(\text{Arch}) \rightarrow Ob(\mathbf{Set}^{\text{Free}(G)})$  is an important construction because of the following reason. It allows us to consider its image in  $\mathbf{Set}^{\text{Free}(G)}$  as a *space* which inherits all the properties of  $\mathbf{Euc}$ . We explore this statement in detail in Section 3.6.4.

### 3.5. Data

We have described constructions which allow us to pick an architecture for a schema and consider its different models  $\text{Model}_p$ , each of them identified with a choice of a parameter  $p \in \mathcal{P}(\text{Arch})$ . In order to understand how the optimization process is steered in updating the parameter choice for an architecture, we need to understand a vital component of any deep learning system – datasets themselves.

Understanding how datasets fit into the broader picture necessitates that we also understand their relationship to the space they are embedded in. Recall the *embedding* functor and the notation  $|\mathcal{C}|$  for the discretization of some category  $\mathcal{C}$ .

**Definition 10.** Let  $|\text{Free}(G)| \xrightarrow{E} \mathbf{Set}$  be the embedding. A **dataset** is a subfunctor  $D_E : |\text{Free}(G)| \rightarrow \mathbf{Set}$  of  $E$ .  $D_E$  maps each object  $A \in Ob(\text{Free}(G))$  to a dataset  $D_E(A) := \{a_i\}_{i=1}^N \subseteq E(A)$ .

Note that we refer to this functor in the singular, although it assigns a dataset to *each* object in  $\text{Free}(G)$ . We also highlight that the domain of  $D_E$  is  $|\text{Free}(G)|$ , rather than  $\text{Free}(G)$ . We generally cannot provide an action on morphisms because datasets might be incomplete. Going back to the example with Horses and Zebras – a dataset functor on  $\text{Free}(G)$  in Figure 3.1 (b) maps Horse to the set of obtained horse images and Zebra to the set of obtained zebra images.

The subobject relation  $D_E \subseteq E$  in Proposition 10 reflects an important property of data; we cannot obtain some data without it being in some shape or form, embedded in some larger space. Any obtained data thus implicitly fixes an embedding.

Observe that when we have a dataset in standard machine learning, we have a dataset *of something*. We can have a dataset of historical weather data, a dataset of housing prices in New York or a dataset of cat images. What ties all these concepts together is that each element  $a_i$  of some dataset  $\{a_i\}_{i=1}^N$  is an instance of a more general concept. As a trivial example, every image in the dataset of horse images is a *horse*. The word *horse* refers to a more general concept and as such could be generalized from some of its instances which we

*do not possess*. But all the horse images we possess are indeed an example of a horse. By considering everything to be embedded in some space  $E(A)$  we capture this statement with the relation  $\{a_i\}_{i=1}^N \subseteq \mathfrak{C}(A) \subseteq E(A)$ . Here  $\mathfrak{C}(A)$  is the set of all instances of some notion  $A$  which are embedded in  $E(A)$ . In the running example this corresponds to all images of horses in a given space, such as the space of all  $64 \times 64$  RGB images. Obviously, the precise specification of  $\mathfrak{C}(A)$  is unknown – as we cannot enumerate or specify the set of *all* horse images.

We use such calligraphy to denote this is an abstract concept. Despite the fact that its precise specification is unknown, we can still reason about its relationship to other structures. Furthermore, as it is the case with any abstract notion, there might be some edge cases or it might turn out that this concept is ambiguously defined or even inconsistent. Moreover, it might be possible to identify a dataset with multiple concepts; is a dataset of male human faces associated with the concept of male faces or with with all faces in general? We ignore these concerns and assume each dataset is a dataset of some well-defined, consistent and unambiguous concept. This does not change the validity of the rest of the formalism in any way as there exist plenty of datasets satisfying such a constraint.

Armed with intuition, we show this admits a generalization to the categorical setting. Just as  $\{a_i\}_{i=1}^N \subseteq \mathfrak{C}(A) \subseteq E(A)$  are all subsets of  $E(A)$  we might hypothesize  $D_E \subseteq \mathfrak{C} \subseteq E$  are all subfunctors of  $E$ . This is quite close to being true. It would mean that the domain of  $\mathfrak{C}$  is the discrete category  $|\mathbf{Free}(G)|$ . However, just as we assign a set of all concept instances to *objects* in  $\mathbf{Free}(G)$ , we also assign a function between these sets to *morphisms* in  $\mathbf{Free}(G)$ . Unlike with datasets, this can be done because, by definition, these sets are not incomplete.

We make this precise as follows. Recall the inclusion functor  $|\mathbf{Free}(G)| \xrightarrow{I} \mathbf{Free}(G) / \sim$ .

**Definition 11.** *Given a schema  $\mathbf{Free}(G) / \sim$  and a dataset  $|\mathbf{Free}(G)| \xrightarrow{D_E} \mathbf{Set}$ , a **concept** associated with the dataset  $D_E$  embedded in  $E$  is a functor  $\mathfrak{C} : \mathbf{Free}(G) / \sim \rightarrow \mathbf{Set}$  such that  $D_E \subseteq \mathfrak{C} \circ I \subseteq E$ . We say  $\mathfrak{C}$  picks out sets of concept instances and functions between those sets.*

Another way to understand a concept  $\mathbf{Free}(G) / \sim \xrightarrow{\mathfrak{C}} \mathbf{Set}$  is that it is required that a human observer can tell, for each  $A \in \mathit{Ob}(\mathbf{Free}(G))$  and some  $a \in E(A)$  whether  $a \in \mathfrak{C}(A)$ . Similarly for morphisms, a human observer should be able to tell if some function  $\mathfrak{C}(A) \xrightarrow{f} \mathfrak{C}(B)$  is an image of some morphism in  $\mathbf{Free}(G) / \sim$  under  $\mathfrak{C}$ .

For instance, consider the GAN schema in Figure 3.1 (a) where  $\mathfrak{C}(\mathit{Image})$  is a set of all images of human faces embedded in some space such as  $\mathbb{R}^{64 \times 64 \times 3}$ . For each image in this space, a human observer should be able to tell if that image contains a face or not. We cannot enumerate such a set  $\mathfrak{C}(\mathit{Image})$  or write it down explicitly, but we can easily tell if an image contains a given concept. Likewise, for a morphism in the CycleGAN schema (Figure 3.1

(b)), we cannot explicitly write down a function which transforms a horse into a zebra, but we can tell if some function did a good job or not by testing it on different inputs.

The most important thing related to this concept is that this represents the goal of our optimization process. Given a dataset  $|\mathbf{Free}(G)| \xrightarrow{D_E} \mathbf{Set}$ , want to extend it into a functor  $\mathbf{Free}(G)/\sim \xrightarrow{\mathcal{C}} \mathbf{Set}$ , and actually *learn* its implementation.

## 3.6. Optimization

We now describe how data guides the search process. We identify the goal of this search with the concept functor  $\mathbf{Free}(G)/\sim \xrightarrow{\mathcal{C}} \mathbf{Set}$ . This means that given a schema  $\mathbf{Free}(G)/\sim$  and data  $|\mathbf{Free}(G)| \xrightarrow{D_E} \mathbf{Set}$  we want to train some architecture  $\text{Arch}$  and find a functor  $\mathbf{Free}(G)/\sim \xrightarrow{H'} \mathbf{Set}$  that can be identified with  $\mathcal{C}$ . Of course, unlike in the case of the concept  $\mathcal{C}$ , the implementation of  $H'$  is something that will be known to us.

We now define the notion of a *task*.

**Definition 12.** *Let  $G$  be a directed multigraph and  $\sim$  a congruence relation on  $\mathbf{Free}(G)$ . A task is a triple  $(G, \sim, |\mathbf{Free}(G)| \xrightarrow{D_E} \mathbf{Set})$ .*

In other words, a graph  $G$  and  $\sim$  specify a schema  $\mathbf{Free}(G)/\sim$  and a functor  $D_E$  specifies a dataset for that schema. Each dataset is a dataset *of something* and thus can be associated with a functor  $\mathbf{Free}(G)/\sim \xrightarrow{\mathcal{C}} \mathbf{Set}$ . Moreover, recall that a dataset fixes an embedding  $E$  too, as  $D_E \subseteq E$ . This in turn also narrows our choice of architecture  $\mathbf{Free}(G) \xrightarrow{\text{Arch}} \mathbf{Para}(\mathbf{Euc})$ , as it has to agree with the embedding on objects. This situation fully reflects what happens in standard machine learning practice – a neural network  $P \times A \rightarrow B$  has to be defined in such a way that its domain  $A$  and codomain  $B$  embed the datasets of all of its inputs and outputs, respectively.

Even though for the same schema  $\mathbf{Free}(G)/\sim$  we might want to consider different datasets, we will always assume a chosen dataset corresponds to a single training goal  $\mathcal{C}$ .

### 3.6.1. Optimization Objectives

As it is the general theme of this thesis – we generalize an already established construction to a categorical setting in a natural way, free of ad-hoc choices. This time we focus on the training procedure as described in Zhu et al. (2017).

Suppose we have a task  $(G, \sim, |\mathbf{Free}(G)| \xrightarrow{D_E} \mathbf{Set})$ . After choosing an architecture  $\mathbf{Free}(G) \xrightarrow{\text{Arch}} \mathbf{Para}(\mathbf{Euc})$  consistent with the embedding  $E$  and, hopefully, with the right inductive bias, we start with a randomly chosen parameter  $\theta_0 \in \mathcal{P}(\text{Arch})$ . Via the parameter-functor map (Definition 8), this amounts to the choice of a specific  $\mathbf{Free}(G) \xrightarrow{\text{Model}_{\theta_0}} \mathbf{Euc}$ . Using the loss function defined further down in this section, we partially differentiate each

$f : \mathbb{R}^n \times \mathbb{R}^a \rightarrow \mathbb{R}^b \in \text{Gen}_{\mathbf{Free}(G)}$  with respect to the corresponding  $p_f$ . We then obtain a new parameter value for that function using some update rule, such as Adam (Kingma and Ba, 2015). The product of these parameters for each of the generators  $(p_f)_{f \in \text{Gen}_{\mathbf{Free}(G)}}$  (Definition 5) defines a new parameter  $\theta_1 \in \mathcal{P}(\text{Arch})$  for the model  $\text{Model}_{\theta_1}$ . This procedure allows us to iteratively update a given  $\text{Model}_{\theta_i}$  and as such fixes a linear order  $\{\theta_0, \theta_1, \dots, \theta_T\}$  on some subset of  $\mathcal{P}(\text{Arch})$ .

The optimization objective for a model  $\mathbf{Free}(G) \xrightarrow{\text{Model}_\theta} \mathbf{Euc}$  and a task  $(G, \sim, |\mathbf{Free}(G)| \xrightarrow{D_E} \mathbf{Set})$  is twofold. The total loss will be stated as a sum of the *adversarial loss* and a *path-equivalence loss*. We now describe both of these losses.

Adversarial loss is tightly linked to an important, but orthogonal concept to this thesis – Generative Adversarial Networks (GANs). As we slowly transition to standard machine learning terminology, we note that some of the notation here will be untyped due to the lack of the proper categorical understanding of these concepts.<sup>3</sup>

We now roughly describe how discriminators fit into the story so far, assuming a fixed architecture. We assign a discriminator to each object  $A \in \text{Ob}(\mathbf{Free}(G))$  using the following function:

$$\mathbf{D} : (A : \text{Ob}(\mathbf{Free}(G))) \rightarrow \mathbf{Para}(\mathbf{Euc})(\text{Arch}(A), \mathbb{R}) \quad (3.1)$$

This function assigns to each object  $A \in \text{Ob}(\mathbf{Free}(G))$  a morphism in  $\mathbf{Para}(\mathbf{Euc})$  such that its domain is that given by  $\text{Arch}(A)$ . This will allow us to compose compatible generators and discriminators. For instance, consider  $\text{Arch}(A) = \mathbb{R}^a$ . Discriminator  $\mathbf{D}(A)$  is then a function of type  $\mathbb{R}^q \times \mathbb{R}^a \rightarrow \mathbb{R} : \mathbf{Para}(\mathbf{Euc})(\mathbb{R}^a, \mathbb{R})$ , where  $\mathbb{R}^q$  is discriminator’s parameter space. As a slight abuse of notation – and to be more in line with machine learning notation – we will call  $\mathbf{D}_A$  discriminator of the object  $A$  with some partially applied parameter value  $\mathbf{D}(A)(p, -)$ .

In the context of GANs, when we refer to a generator we refer to the image of a generating morphism in  $\mathbf{Free}(G)$  under  $\text{Arch}$ . Similarly as with discriminators, a generator corresponding to a morphism  $\mathbb{R}^a \xrightarrow{f} \mathbb{R}^b$  in  $\mathbf{Para}(\mathbf{Euc})$  with some partially applied parameter value will be denoted using  $\mathbf{G}_f$ .

The GAN minimax objective  $\mathcal{L}_{GAN}^B$  for a generator  $\mathbf{G}_f$  and a discriminator  $\mathbf{D}_B$  is stated in Eq. (3.2). In this formulation we use Wasserstein distance (Arjovsky et al., 2017).

---

<sup>3</sup> In this thesis this adversarial component is used in the optimization procedure, but to the best of our knowledge, it has not been properly framed in a categorical setting yet and is still an open problem. It seems to require nontrivial reformulations of existing constructions (Fong et al., 2017) and at least a partial integration of Open Games (Ghani et al., 2016) into the framework of gradient-based optimization. In this thesis we do not concern ourselves with these matters and they present no practical issues for training these networks.

$$\begin{aligned} \mathcal{L}_{GAN}^B(\mathbf{G}_f, \mathbf{D}_B) &:= \mathbb{E}_{b \sim D_E(B)} [\mathbf{D}_B(b)] \\ &\quad - \mathbb{E}_{a \sim D_E(A)} [\mathbf{D}_B(\mathbf{G}_f(a))] \end{aligned} \quad (3.2)$$

The generator is trained to minimize the loss in the Eq. (3.2), while the discriminator is trained to maximize it.

The second component of the total loss is a generalization of *cycle-consistency loss* in CycleGAN (Zhu et al., 2017), analogous to the generalization of the cycle-consistency condition in Section 3.1.

**Definition 13.** Let  $A \xrightleftharpoons[f]{g} B$  and suppose there exists a path equivalence  $f = g$ . For the equivalence  $f = g$  and the model  $\mathbf{Free}(G) \xrightarrow{\text{Model}_i} \mathbf{Euc}$  we define a **path equivalence loss**  $\mathcal{L}_{\sim}^{f,g}$  as follows:

$$\mathcal{L}_{\sim}^{f,g} := \mathbb{E}_{a \sim D_E(A)} [\| \text{Model}_i(f)(a) - \text{Model}_i(g)(a) \|_1]$$

This enables us to state the total loss simply as a weighted sum of adversarial losses for all generators and path equivalence losses for all equations.

**Definition 14.** The **total loss** is given as the sum of all adversarial and path equivalence losses:

$$\mathcal{L}_i := \sum_{A \xrightarrow{f} B \in \text{Gen}_{\mathbf{Free}(G)}} \mathcal{L}_{GAN}^B(\mathbf{G}_f, \mathbf{D}_B) + \gamma \sum_{f=g \in \sim} \mathcal{L}_{\sim}^{f,g} \quad (3.3)$$

where  $\gamma$  is a hyperparameter that balances between the adversarial loss and the path equivalence loss.

Observe that identity morphisms are not elements of  $\text{Gen}_{\mathbf{Free}(G)}$ . Even if they were, they would cause the discriminator to be steered in two directions. They would incentivize the discriminator to classify samples of  $D_E(A)$  as real, but also to classify samples of  $D_E(A)$  under  $\mathbf{G}_{id_A} : A \rightarrow A$  as fake. But notice that  $\mathbf{G}_{id_A}$  does not change the samples, so it would be pushed to classify the same samples as both real and fake – making the net result is zero as these effects cancel each other out.

### 3.6.2. Restriction of Network Instance to the Dataset

We have seen how data relates to the architecture and how model  $\text{Model}_{p_i}$  corresponding to a parameter  $p_i \in \mathcal{P}(\text{Arch})$  is updated. Observe that network instance  $H_p$  maps each object  $A \in \text{Ob}(\mathbf{Free}(G))$  to the entire embedding  $H_{p_i}(A) = E(A)$ , rather than just the

concept  $\mathfrak{C}(A)$ . Even though we started out with an embedding  $E(A)$ , we want to narrow that embedding down just to the set of instances corresponding to some concept  $A$ .

For example, consider a diagram such as the one in Figure 3.1 (a). Suppose the result of a successful training was a functor  $\mathbf{Free}(G) \xrightarrow{H} \mathbf{Set}$ . Suppose that the image of  $h$  is  $H(h) : [0, 1]^{100} \rightarrow [0, 1]^{64 \times 64 \times 3}$ . As such, our interest is mainly the restriction of  $[0, 1]^{64 \times 64 \times 3}$  to  $\mathfrak{C}(\text{Image})$ , the image of  $[0, 1]^{100}$  under  $H(h)$ , rather than the entire  $[0, 1]^{64 \times 64 \times 3}$ . In the case of horses and zebras in Figure 3.1 (b), we are interested in a map  $\mathfrak{C}(\text{Horse}) \rightarrow \mathfrak{C}(\text{Zebra})$  rather than a map  $[0, 1]^{64 \times 64 \times 3} \rightarrow [0, 1]^{64 \times 64 \times 3}$ . In what follows we show a construction which restricts some  $H_p$  to its smallest subfunctor which contains the dataset  $D_E$ . Recall the previously defined forgetful functor  $\mathbf{Euc} \xrightarrow{U} \mathbf{Set}$  and the inclusion  $|\mathbf{Free}(G)| \xrightarrow{I} \mathbf{Free}(G)$ .

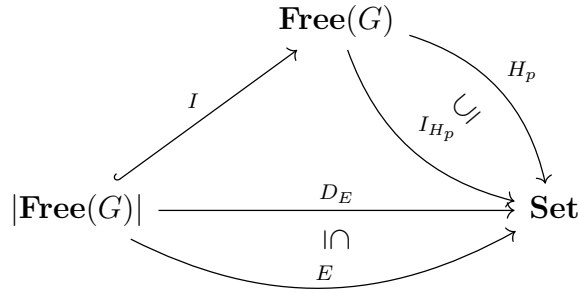
**Definition 15.** Let  $D_E : |\mathbf{Free}(G)| \rightarrow \mathbf{Set}$  be the *dataset*. Let  $\mathbf{Free}(G) \xrightarrow{H_p} \mathbf{Set}$  be the network instance on  $\mathbf{Free}(G)$ . The *restriction* of  $H_p$  to  $D_E$  is a subfunctor of  $H_p$  defined as follows:

$$I_{H_p} := \bigcap_{\{G \in \text{Sub}(H_p) \mid D_E \subseteq G \circ I\}} G$$

where  $\text{Sub}(H_p)$  is the set of subfunctors of  $H_p$ .

This definition is quite condensed so we supply some intuition. We first note that the meet is well defined because each  $G$  is a subfunctor of  $H$ .

In Figure 3.4 we depict the newly defined constructions using a commutative diagram.



**Figure 3.4:** The functor  $I_{H_p}$  is a subfunctor of  $H_p$ , and  $D_E$  is a subfunctor of  $I_{H_p} \circ I$ .

It is useful to think of  $I_H$  as a restriction of  $H$  to the *smallest* functor which fits all data and mappings between the data. This means that  $I_{H_p}$  contains all data samples specified by  $D_E$ .

**Corollary 16.**  $D_E$  is a subfunctor of  $I_{H_p} \circ I$ :

*Proof.* This is straightforward to show, as  $I_{H_p}$  is the intersection of all subobjects of  $H$  which, when composed with the inclusion  $I$  contain  $D_E$ . Therefore  $I_{H_p} \circ I$  contains  $D_E$  as well.  $\square$

Even though Corollary 7 tells us that all  $H_p$  act the same on objects, we can see that this is not the case with  $I_{H_p}$ .

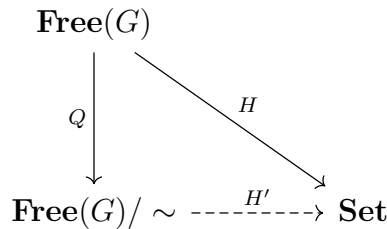
Consider some  $A \xrightarrow{f} B$  in  $\mathbf{Free}(G)$ . Suppose we have two different network instances,  $\mathbb{R}^a \xrightarrow{H_p(f)} \mathbb{R}^b$  and  $\mathbb{R}^a \xrightarrow{H_q(f)} \mathbb{R}^b$ . Even though these instances act the same on objects, their restrictions to data  $D_E$  might not. Depending on the implementations of  $H_p(f)$  and  $H_q(f)$ , images of  $\mathbb{R}^a$  under  $H_p(f)$  and  $H_q(f)$  might be different, and as such impose different constraints on the minimum size of their restriction to  $D_E$  on the object  $B$ . This in turn means that the sets  $I_{H_p}(B)$  and  $I_{H_q}(B)$  might be vastly different and practically agree only on the  $D_E(B)$ . More generally, for any  $X \in \mathit{Ob}(\mathbf{Free}(G))$  and  $X \xrightarrow{g} A$ ,  $I_{H_p}(A)$  contains the image of  $D_E(X)$  under  $H(X) \xrightarrow{H(g)} H(A)$ .

### 3.6.3. Path Equivalence Relations

There is one interesting case of the total loss – when the total path-equivalence loss is zero:  $\sum_{f=g \in \sim} \mathcal{L}_{\sim}^{f,g} = 0$ . This tells us that  $H(f) = H(g)$  for all  $f = g$  in  $\sim$ . In what follows we elaborate on what this means by recalling how  $\mathbf{Free}(G)$  and  $\mathbf{Free}(G)/\sim$  are related.

So far, we have been only considering schemas given by  $\mathbf{Free}(G)$ . This indeed is a limiting factor, as it assumes the categories of interest are only those without any imposed relations  $R$  between the generators  $G$ . One example of a schema *with* relations is the CycleGAN schema 3.1 (b). As briefly mentioned in Remark 3, fixing a functor  $\mathbf{Free}(G)/\sim \rightarrow \mathbf{Set}$  requires its image to satisfy any relations imposed by  $\mathbf{Free}(G)/\sim$ . As neural network parameters usually are initialized randomly, any such image in  $\mathbf{Set}$  will most surely not preserve such relations and thus will not be a proper functor whose domain is  $\mathbf{Free}(G)/\sim$ .

However, this construction is a functor if we consider its domain to be  $\mathbf{Free}(G)$ . Furthermore, assuming a successful training process whose end result is a path-equivalence relation preserving functor  $\mathbf{Free}(G) \rightarrow \mathbf{Set}$ , we show this induces a unique  $\mathbf{Free}(G)/\sim \rightarrow \mathbf{Set}$  (Figure 3.5).



**Figure 3.5:** Functor  $H$  which preserves path-equivalence relations factors uniquely through  $Q$ .

**Theorem 17.** *Let  $Q : \mathbf{Free}(G) \rightarrow \mathbf{Free}(G)/\sim$  be the quotient functor and let  $H : \mathbf{Free}(G) \rightarrow \mathbf{Set}$  be a path-equivalence relation preserving functor. Then there exists a unique functor  $H' : \mathbf{Free}(G)/\sim \rightarrow \mathbf{Set}$  such that  $H' \circ Q = H$ .*



*Proof.* We refer the reader to MacLane (1971), Section 2.8., Proposition 1.  $\square$

Finding such a functor  $H$  is no easier task than finding a functor  $H'$ . However, considering the domain to be  $\mathbf{Free}(G)$ , rather than  $\mathbf{Free}(G)/\sim$  allows us to initially just guess a functor  $H_0$ , since our initial choice will not have to preserve the equivalence  $\sim$ . As training progresses and the path-equivalence loss of a network instance  $\mathbf{Free}(G) \xrightarrow{H_p} \mathbf{Set}$  converges to zero, by Theorem 17 we show  $H_p$  factors uniquely through  $\mathbf{Free}(G) \xrightarrow{Q} \mathbf{Free}(G)/\sim$  via  $\mathbf{Free}(G)/\sim \xrightarrow{H'} \mathbf{Set}$ .

### 3.6.4. Functor Space

Recall the parameter-functor map (Definition 8)  $\mathcal{M}_{\mathbf{Euc}} : \mathcal{P}(\text{Arch}) \rightarrow \text{Ob}(\mathbf{Euc}^{\mathbf{Free}(G)})$ . Each choice of  $p \in \mathcal{P}(\text{Arch})$  specifies a functor of type  $\mathbf{Free}(G) \rightarrow \mathbf{Set}$ . In this way exploration of the parameter space amounts to exploration of part of the functor category  $\mathbf{Set}^{\mathbf{Free}(G)}$ . Roughly stated, this means that a choice of an architecture  $\mathbf{Free}(G) \xrightarrow{\text{Arch}} \mathbf{Para}(\mathbf{Euc})$  adjoins a notion of *space* to the image of  $\mathcal{P}(\text{Arch})$  under  $\mathcal{M}_{\mathbf{Euc}}$  in the functor category  $\mathbf{Euc}^{\mathbf{Free}(G)}$ . This space inherits all the properties of  $\mathbf{Euc}$ .

By using gradient information to search the parameter space  $\mathcal{P}(\text{Arch})$ , we are effectively using gradient information to search part of the functor space  $\mathbf{Set}^{\mathbf{Free}(G)}$ . Although we cannot explicitly explore just  $\mathbf{Set}^{\mathbf{Free}(G)/\sim}$ , we penalize the search method for veering into the parts of this space where the specified path equivalences do not hold. As such, the inductive bias of the model is increased without special constraints on the datasets or the embedding space – we merely require that the space is differentiable and that it has a sensible notion of distance. Note that we do not claim inductive bias is *sufficient* to guarantee training convergence, merely that it is a useful regularization method applicable to a wide variety of situations.

As categories can encode complex relationships between concepts and as functors map between categories in a structure-preserving way – this enables *structured learning* of concepts and their interconnections in a very general fashion. Of course, this does not tell us which such structures are practical to learn or what other inductive biases we need to add, but simply disentangles those structures such that they can be studied in their own right.

To the best of our knowledge, this represents a considerably different approach to learning. Rather than employing a gradient-based search in *one* function space  $Y^X$  this formulation describes a search method in a *collection* of such function spaces and then regularizes it to satisfy any composition rules that have been specified. Even though it is just a rough sketch, this concludes our reasoning which shows functors of the type  $\mathbf{Free}(G) \rightarrow \mathbf{Set}$  can be *learned* using gradient-descent.

	$\mathbf{Free}(G)/\sim \rightarrow \mathbf{Set}$	$\mathbf{Free}(G) \rightarrow \mathbf{Free}(G)/\sim$
Functorial Data Migration	Fixed	Data integrity
Compositional Deep Learning	Learned	Cycle-consistency

**Table 3.1:** Pocket version of the Rosetta stone between Functorial Data Migration and Compositional Deep Learning

### 3.7. From Categorical Databases to Deep Learning

The formulation presented in this thesis bears a striking and unexpected similarity to Functorial Data Migration (FDM) (Spivak, 2010). Given a *categorical schema*  $\mathcal{C}$  on some directed multigraph  $G$ , FDM specifies a functor category  $\mathbf{Set}^{\mathcal{C}}$  of database instances on that schema. The notion of *data integrity* is captured by *path equivalence relations* which constrain the schema by demanding specific paths in the graph be the same. The analogue of data integrity is captured by cycle-consistency conditions, first introduced in CycleGAN (Zhu et al., 2017). In Table 3.1 we outline the main differences between the approaches. Namely, our approach requires us not to specify functors into  $\mathbf{Set}$ , but rather to *learn* them.

This shows that the underlying structures used for specifying data semantics and concrete instances for a given database system are equivalent to the structures used for describing network semantics and their concrete instances.

## 4. Examples

We have seen how specific concepts in deep learning fit into a rigorous categorical framework. We shed further light on this framework by considering some of its interesting special cases. Namely, we show how different instances of a task  $(G, \sim, |\mathbf{Free}(G)| \xrightarrow{D_E} \mathbf{Set})$  give rise not just to two existing architectures, but also to a completely new system.

### 4.1. Existing Architectures

We first shift our attention to two existing architectures: GAN and CycleGAN. We specify the choice of  $G$ ,  $\sim$  and the dataset  $\mathbf{Free}(G) \xrightarrow{D_E} \mathbf{Set}$  and show this determines our interpretation of the learned semantics.

#### 4.1.1. GAN Task

The choice of a graph and path equivalence relations corresponding to the GAN task can be identified with the schema in Figure 3.1 (a) and a choice of a dataset. We will abbreviate “Latent space” with  $LS$  and “Image space” with  $IS$ . A choice of a dataset can be understood as follows. We think of  $D_E(LS) = [0, 1]^{100}$  as the choice of uniform distribution in some high-dimensional space and  $D_E(IS) \subseteq [0, 1]^{64 \times 64 \times 3}$  as the choice of some image dataset, such as faces, cars, or handwritten digits. The choice of a dataset fixes the embedding  $|\mathbf{Free}(G)| \xrightarrow{E} \mathbf{Set}$ , but also the choice of a concept  $\mathfrak{C}$ . In our case  $\mathfrak{C}(IS)$  is the set of  $64 \times 64$  RGB images human faces and  $\mathfrak{C}(LS)$  is some notion of an indexing set. As the corresponding dataset  $D_E(IS)$  does not have any meaningful interpretation, neither does  $\mathfrak{C}(IS)$ .

There are no path equivalences here, so the total loss reduces just to the adversarial loss. Moreover, since the only generator in this schema is  $LS \xrightarrow{h} IS$ , the total adversarial loss reduces just to the standard minimax loss. As per Example 9, the parameter space of any chosen model  $\mathbf{Free}(G) \xrightarrow{\text{Model}} \mathbf{Para}(\mathbf{Euc})$  reduces to the parameter space of the function  $\text{Arch}(h)$ .

In a similar fashion, the network instance  $H_p$  can be identified with a choice of a function  $H_p(h)$ . Observe that the image some function  $[0, 1]^{100} \xrightarrow{H_p(h)} [0, 1]^{1024 \times 1024 \times 3}$  might only

partially overlap the dataset  $D_E(IS)$ . Nevertheless, the induced  $I_{H_p}$  allows us to consider the union of the dataset and the image of this function.

### 4.1.2. CycleGAN Task

CycleGAN task corresponds to the choice of schema in Figure 3.1 (b) and a choice of two datasets for each of the objects whose corresponding concepts are isomorphic. As previously mentioned, these could be sets of images of *apples and oranges*, *paintings and photos*, *horses and zebras* etc.

The parameter space of some chosen architecture  $\mathbf{Free}(G) \xrightarrow{\text{Arch}} \mathbf{Para}(\mathbf{Euc})$  for this schema is the product  $\mathcal{P}(\text{Arch}) = \text{Arch}(f) \times \text{Arch}(g)$ . There are two discriminators – one which distinguishes between real and fake horse images and the one which distinguishes between real and fake zebra images. Unlike in the GAN task, the total loss does not reduce just to the adversarial one. Furthermore, the adversarial component of the loss is the sum of *two* minimax losses, for  $(G_f, D_B)$  and  $(G_g, D_A)$ .

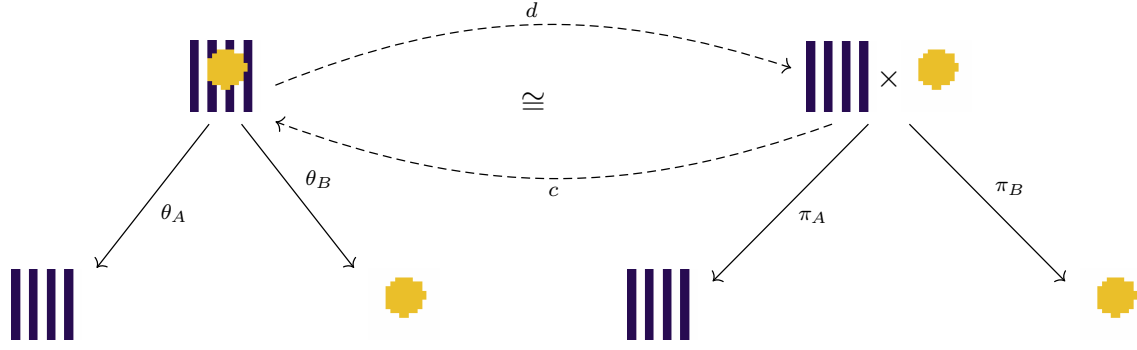
## 4.2. Product Task

We have given two examples of a categorical *schema*. Each of them has specific semantics which can be interpreted. CycleGAN posits two objects  $A$  and  $B$  are *the same* in a certain way by learning an isomorphism  $A \cong B$ . GAN has a simple semantic which tells us that one object  $A$  indexes at least part of another object  $B$  in some way – i.e. that there exists a mapping  $A \rightarrow B$ . Together with the choice of a dataset, they constitute a choice of a task.

We now present a different choice of a dataset for the CycleGAN schema which makes up a task we will call *the product task*. The interpretation of this task comes in two flavors. From one perspective, it is a simple change of dataset for the CycleGAN schema. The other perspective justifies the name *product task* by grounding the explanation in categorical terms.

Just as we can take the product of two real numbers  $a, b \in \mathbb{R}$  with a multiplication function  $(a, b) \mapsto ab$ , we show we can take a product of some two sets of images  $A, B \in \text{Ob}(\mathbf{Set})$  with a neural network of type  $A \times B \rightarrow C$ . We will show  $C \in \text{Ob}(\mathbf{Set})$  is a set of images which possesses all the properties of a categorical product.

In a Cartesian category such as  $\mathbf{Set}$  there already exists a notion of a categorical product – the Cartesian product. Recall that it is unique up to unique isomorphism. This means that any other object  $X$  which satisfies the universal property of the categorical product of  $A$  and  $B$  is uniquely isomorphic to  $A \times B$ . This isomorphism will be central to the notion of the product task. Before continuing, we supply some intuition with Figure 4.1.



**Figure 4.1:** Two different notions of a categorical product of  $A$  and  $B$  in  $\mathbf{Set}$ , illustrated with examples from some choice of those sets. Right side depicts the Cartesian product  $A \times B$ , while the left side depicts some different product which we call  $AB$ . Both of them have a notion of projections to their constituents. By the universal property of the categorical product, they are uniquely isomorphic. This central isomorphism represents the CycleGAN perspective of the product task. The decomposition map  $d : AB \rightarrow A \times B$  and the composition map  $c : A \times B \rightarrow AB$  are things we want to learn.

By *learning* the model  $\mathbf{Free}(G) \xrightarrow{\text{Model}} \mathbf{Euc}$  corresponding to the isomorphism  $AB \cong A \times B$  we are also learning the projection maps  $\theta_A$  and  $\theta_B$ . This follows from the universal property of the categorical product:  $\pi_A \circ d = \theta_A$  and  $\pi_B \circ d = \theta_B$ . Note how  $AB$  differs from a Cartesian product. The domain of the corresponding projections  $\theta_A$  and  $\theta_B$  is not a simple pair of objects  $(a, b)$  and thus these projections cannot merely discard an element.  $\theta_A$  needs to learn to remove  $A$  from a potentially complex domain. As such, this can be any complex, highly non-linear function which satisfies coherence conditions of a categorical product.

We will be concerned with supplying this new notion of the product  $AB$  with a dataset and learning the image of the isomorphism  $AB \cong A \times B$  under  $\mathbf{Free}(G) \xrightarrow{\text{Model}} \mathbf{Set}$ . We illustrate this on a concrete example. Consider a dataset  $A$  of images of human faces, a dataset  $B$  of images of sunglasses, and a dataset  $AB$  of people *wearing* glasses. Learning this isomorphism amounts to learning two things: (i) learning how to decompose an image of a person wearing glasses  $(ab)_i$  into an image of a person  $a_j$  and image  $b_k$  of these glasses, and (ii) learning how to map this person  $a_j$  and perhaps some other glasses  $b_l$  into an image of a person  $a_j$  wearing glasses  $b_l$ . Generally,  $AB$  represents some sort of composition of objects  $A$  and  $B$  in the image space such that all information about  $A$  and  $B$  is preserved in  $AB$ . Of course, this might only be approximately true. Glasses usually cover a part of a face and sometimes its dark shades cover up the eyes – thus losing information about the eye color in the image and rendering the isomorphism invalid. However, in this thesis we ignore such issues and assume that the networks  $\text{Arch}(d)$  can learn to unambiguously fill part of the face where the glasses were and that  $\text{Arch}(c)$  can learn to generate and superimpose the glasses on the relevant part of the face.

This example shows us that product task is tightly linked with composition and decomposition of images. Although we use the same CycleGAN schema from Figure 3.1 (b), we label one of its objects as  $AB$  and the other one as  $A \times B$ . Note that this does not change the schema itself, the labeling is merely for our convenience. We highlight that the notion of a product or its projections is not captured in the schema itself. As schemas are merely categories presented with generators  $G$  and relations  $R$ , they lack the tools needed to encode a complex abstraction such as a universal construction. So how do we capture the notion of a product?

In this thesis we frame this simply as a specific dataset functor, which we now describe. A dataset functor corresponding to the product task maps the object  $A \times B$  in CycleGAN schema to a *Cartesian product of two datasets*,  $D_E(A \times B) = \{a_i\}_{i=0}^N \times \{b_j\}_{j=0}^M$ . It maps the object  $AB$  to a dataset  $\{(ab)_i\}_{i=0}^N$ . In this case  $ab$ ,  $a$ , and  $b$  are free to be any elements of datasets of a well-defined concept  $\mathcal{C}$ . Although the difference between the product task and the CycleGAN task boils down to a different choice of a dataset functor, we note this is a key aspect which allows for a significantly different interpretation of the task semantics.<sup>1</sup>

We now describe two possible classes of choice of datasets for the product task. One is already given in the example with glasses and faces. This example includes three distinct datasets in its consideration. However, observe that sometimes three such related datasets might not always be available. For instance, consider the scenario where we only possess the dataset of faces with glasses  $AB$  and a dataset of just faces  $A$ , but not a dataset of images of glasses. This prevents us from generating images of glasses which some person is wearing. But it does not prevent us from using the product task. Namely, we know that we can decompose the point in some image space corresponding to a person with glasses  $(ab) \in D_E(AB)$  into point in some image space of a person  $a$  and of a point in some different space which captures whatever the missing information between  $(ab)$  and  $a$  is. In this case, this point would parametrize the color, shape, and type of glasses.

Put another way, one possible choice for the action of a dataset on  $A \times B$  is  $D_E(A \times B) = \{\text{Faces}_i\}_{i=0}^N \times [0, 1]^{100}$ . Similarly, as a regular GAN does not allow us to know in what way the object  $LS$  indexes  $IS$ , here do we not know how  $[0, 1]^{100}$ , together with  $D_E(B)$ , indexes  $D_E(AB)$ . This represents one perspective on the many-to-many mappings introduced by Almahairi et al. (2018). This is a map of a face to a face with glasses as a one-to-many map.

The product schema admits a straightforward generalization to  $n$  objects. For example, the isomorphism  $ABC \cong A \times B \times C$  can be understood as a one-to-many map between  $A \times B$  and  $ABC$ , but also as a one-to-many map between  $A$  and  $ABC$ .

---

<sup>1</sup> Just as we can consider a joint probability distribution of two random variables  $X$  and  $Y$ , here we consider a joint *dataset* of two datasets  $\{a_i\}_{i=0}^N$  and  $\{b_j\}_{j=0}^M$ . In a similar fashion we consider the Cartesian product  $\mathbf{D}(A) \times \mathbf{D}(B)$  of two discriminators  $\mathbf{D}(A)$  and  $\mathbf{D}(B)$  to be a discriminator itself.

With this in mind, we now envision potential applications of the product schema.

Consider the task of composing a car  $A$  with a road  $B$ . In most cases there is certainly more than one way to do so – by putting the car on different parts of this road. By generalizing the product schema to more than two datasets we begin to explore its full potential.

We could consider an isomorphism  $ABC \cong A \times B \times C$ , where  $D_E(A)$  is a set of car images,  $D_E(B)$  is a set of road images and  $D_E(C) = [0, 1]^{100}$  is some choice of uniform distribution, similarly as with the GAN task. Now, given a car  $a \in \mathcal{C}(A)$ , a road  $b \in \mathcal{C}(B)$  and some parametrization the position of car on the road  $c \in \mathcal{C}(C)$ , we can learn an unique composition and decomposition function. Observe that the purpose of  $C$  is to add information on *whatever is missing* from  $A \times B$  so that the product of  $A \times B$  with this missing information is isomorphic to  $ABC$ .

By considering  $A$  as some *image background* and  $B$  as the *object* which will be inserted, this allows us to interpret  $d$  and  $c$  as maps which *remove an object from the image* and *insert an object* in an image, respectively. They do it in a coherent, parametric way, such that  $B$  parametrizes the type of object which will be inserted into the domain  $A$  and  $C$  all the extra information needed to make insertion and removal inverses. Furthermore, we restate that is all done on unpaired data. As such, no extra structure is imposed on the *sets*  $A$ ,  $B$  and  $AB$ , easing the data collection process.

This seems like a novel method of object generation and deletion with unpaired data, though a more thorough literature survey should be done to confirm this.

## 5. Experiments

In this chapter, we put on our engineering hat and test whether the product task described in Section 4.2 can be trained in practice. We perform experiments on three datasets. The first experiment is done on a dataset created by us and the second is done on a well-known image dataset of celebrities. The third experiment is significantly different and involves the usage of an already pretrained generative model. As such the choices of architecture and hyperparameters for first two experiments differs from the choices for the third one.

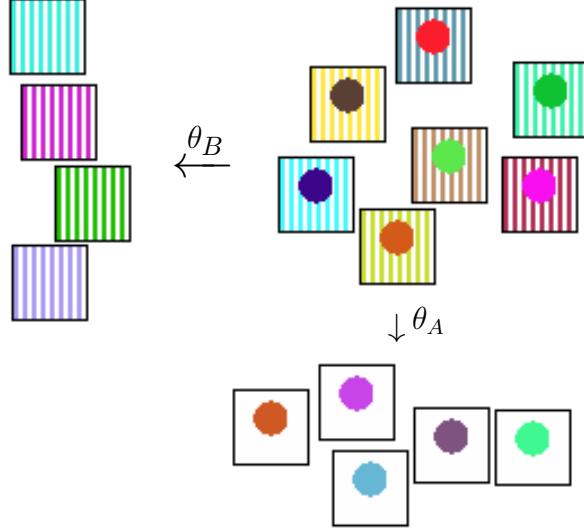
In the first two experiments we have used optimizer Adam (Kingma and Ba, 2015) and the Wasserstein GAN with gradient penalty (Gulrajani et al., 2017). We used the suggested choice of hyperparameters in Gulrajani et al. (2017). The parameter  $\gamma$  is set to 20 and as such weighted the optimization procedure towards the path-equivalence, rather than the cycle-consistency loss. All weights were initialized from a Gaussian distribution  $\mathcal{N}(0, 0.01)$ . As suggested in (Gulrajani et al., 2017), we always gave the discriminator a head start and trained it more, especially in the beginning. We set  $n_{critic} = 50$  for the first 50 time steps and  $n_{critic} = 5$  for all other time steps.

Discriminator  $\mathbf{D}(AB)$  and the discriminators for each  $A$  and  $B$  in  $\mathbf{D}(A \times B)$  in first two experiments were 5-layer ReLU convolutional neural networks of type  $\mathbb{R}^q \times \mathbb{R}^{32 \times 32 \times 3} \rightarrow \mathbb{R}$ . Kernel size was set to 5 and padding to 2. We used stride 2 to halve image size in all layers except the second, where we used stride 1. The only things varying in Sections 5.1 and 5.2 were the filter sizes. We used a fully-connected layer without any activations at the end of the convolutional network to reduce the output size to 1.

### 5.1. Circles

We created datasets whose samples are shown in a suggestive manner in Figure 5.1.





**Figure 5.1:** A dataset of stripes and circles can be decomposed sensibly into a dataset of circles and a dataset of stripes.

Each dataset consists of  $32 \times 32$  RGB images of either a circle (corresponding to  $A$  in the schema), stripes (to  $B$ ) or circle on stripes (to  $AB$ ), all normalized to a  $[0, 1]$  range. The background dataset  $D_E(B)$  consists of unicolored stripes on a white background. The object dataset  $D_E(A)$  consists of unicolored circles in a fixed position on a white image. The dataset  $D_E(AB)$  consist of a unicolored circles superimposed over unicolored stripes on a white background. Each element of both  $D_E(A)$  and  $D_E(B)$  can be identified with a 3-dimensional vector describing the choice of color of the circle and the stripes, respectively. Likewise, each element of  $D_E(AB)$  can be identified with a 6-dimensional vector, describing the color of both the circle and the stripes.

More precisely, we define the dataset  $D_E$  as follows:  $D_E(A \times B) = \{\text{Circles}_i\}_{i=0}^N \times \{\text{Stripes}_j\}_{j=0}^M$  and  $D_E(B) = \{\text{Circles\_on\_stripes}_k\}_{k=0}^K$ .

The task of some network architecture is then to learn to decompose an image into its two constituents. We highlight two important things. We do not provide explicit dataset pairings to the network. We also do not tell the networks *what* the images contain. No information was specified about the circle location or even the fact that there some kind of an object in the image.

We now specify an architecture for this task. It is given by a choice of morphisms in  $\text{Para}(\mathbf{Euc})$  for two generators  $d : AB \rightarrow A \times B$  and  $c : A \times B \rightarrow AB$  in  $\mathbf{Free}(G)$ . For both of them, we choose fully-connected networks and thus flatten the images into a 3072-dimensional vector before using them as inputs to  $\text{Model}(c)$  and  $\text{Model}(d)$ . This fixes the type of possible network architectures  $-\mathbb{R}^p \times \mathbb{R}^{3072+3072} \xrightarrow{\text{Arch}(c)} \mathbb{R}^{3072}$  and  $\mathbb{R}^q \times \mathbb{R}^{3072} \xrightarrow{\text{Arch}(d)} \mathbb{R}^{3072+3072}$ .

We implement both networks as simple autoencoders with a bottleneck of dimension 6.

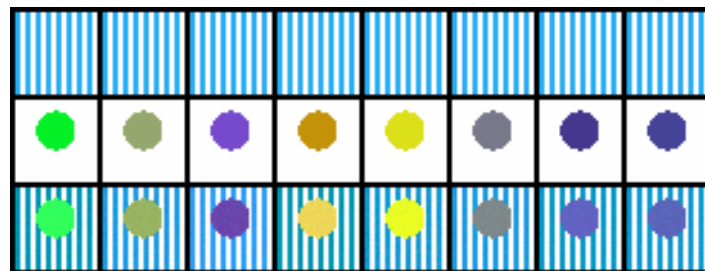
This means that both of them are incentivized to extract a 6-dimensional summary of the image – the color of the circle and the color of the stripes. Even though the path-equivalence relations regularize the networks to be inverses of each other, there could actually exist *many* isomorphisms between the objects  $\mathcal{C}(A)$  and  $\mathcal{C}(B)$ . For instance, the networks could learn to preserve color information by encoding it as a *different* color. In this experiment we remedy this issue by adding an auxiliary *identity mapping* loss function, as outlined in (Zhu et al., 2017). It incentivizes the training procedure to pick out a *specific* isomorphism which preserves the color information in a canonical way.

The architectures of the discriminators are described in the beginning of this chapter, with a sequence of filter sizes for each layer as follows: [3, 8, 16, 32, 64, 64].

### 5.1.1. Evaluation of Trained Networks

With this in mind, we show the results of the training. Given that this is a toy task and that the architecture parameter count is quite low compared to standard deep learning architectures, both networks quickly solve the task. The adversarial loss steers the generators in generating realistic looking images, while the path-equivalence loss steers the generators in preserving color information.

We test the results in three distinct ways – by changing either  $A$ ,  $B$ , or  $AB$  while keeping everything else fixed and evaluating the results. In Figure 5.2 we test if the network learned to preserve background color while the circle color is changing. Figure 5.3 tests if the network learned to preserve background color while the circle color is changing. And lastly, Figure 5.4 test the decomposition network.



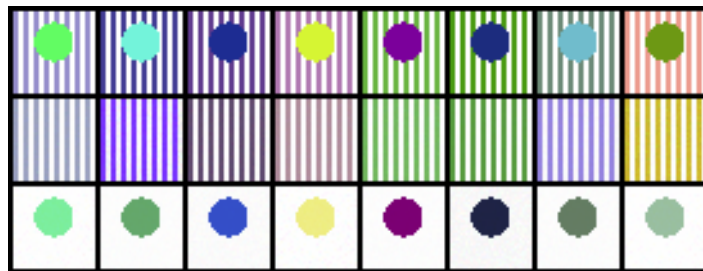
**Figure 5.2:** We evaluate the composition network by keeping the *stripe color fixed* and circle color changing. First two rows are the inputs to the network and third row shows the results. We can see that both the stripe and circle color in the generated image match the expected one.

We observe a general pattern throughout these tests. Although the colors match the expected ones in most cases, there exists a slight discrepancy between the obtained and the expected results. Most of the experiments we ran exhibited similar symptoms, but we did

not look further into this issue. Nevertheless, they do manage to quickly learn and solve this simple task and as such prove the viability of the product task.



**Figure 5.3:** We evaluate the composition network by keeping the *circle color fixed* and stripe color changing. First two rows are the inputs to the network and third row shows the results. We observe that although the stripe color matches the expected one, circle color varies slightly in brightness.



**Figure 5.4:** We evaluate the decomposition network whose input is shown in the first row. The second and third row show the generated stripe and circle image, respectively. We note that some discrepancy between the expected and generated colors exists.

## 5.2. CelebA

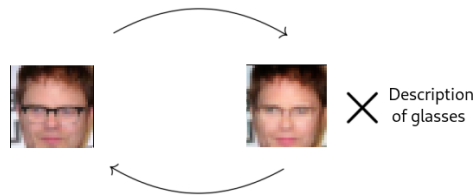
The circle dataset was made with the product task in mind and as such was its canonical example. We had a notion of a dataset with circles, a dataset with stripes and a dataset which contains both. However, in the real world, the luxury of obtaining three such datasets related in a special way is not always there. We now turn our attention to an existing image dataset which we use as a more realistic test of the feasibility of the product task.

CelebFaces Attributes Dataset (CelebA) (Liu et al., 2015) is a large-scale face attributes dataset with more than 200000 celebrity images, each with 40 attribute annotations. Frequently used for image generation purposes, it fits perfectly into the proposed paradigm of the product task. The images in this dataset cover large pose variations and background clutter. The attribute annotations include “eyeglasses”, “bangs”, “pointy nose”, “wavy hair” etc., as boolean flags for every image.

We used the attribute annotations to separate CelebA into two datasets. The dataset  $D_E(AB)$  consists of images with the attribute “Eyeglasses”, while the dataset  $D_E(A)$  consists of all the other images.

Given that we have no dataset of just glasses, we set  $D_E(B) = [0, 1]^{100}$ . As outlined in section 4.2, this is a parametrization of all the missing information from  $A$  such that  $A \times B \cong AB$ . We refer to an element  $z \in [0, 1]^{100}$  as a *latent vector*, in line with machine learning terminology. We will also use  $Z$  instead of  $B$ , as to make it more clear we are not generating images of this object.

As before, we highlight that there are no explicit pairs in this dataset – there does not exist an image of the same person with and without glasses in completely the same position. Most of the people in those images appear only once and as such this dataset provides an interesting challenge for generalization in the context of the product task (Figure 5.5).



**Figure 5.5:** We trained two networks, a *composition* and *decomposition* network which learned to insert and remove glasses from images, respectively. We hypothesize this task of object removal and insertion works in more general scenarios, as nothing in the training process is specific to glasses or people.

We resize all images to  $32 \times 32$  to make it possible to train such network with available compute resources. The network architecture of  $\text{Arch}(f)$  and  $\text{Arch}(g)$  are combinations of convolutional layers and residual blocks. Unlike with the circle dataset, we do not think of these networks as autoencoders because we have no guarantees on the exact bottleneck size. The architectures of the discriminators are described in the beginning of this chapter, with a sequence of filter sizes for each layer as follows:  $[3, 64, 128, 256, 512, 512]$ .

The task for the decomposition neural network is then to transform an image of a person wearing glasses into an image satisfying the following: (i) it needs to be a realistic looking image of a person, and (ii) the identity of the person needs to be preserved when removing glasses. The first condition depends on the adversarial loss, while the second one depends on the path-equivalence loss. The composition neural network has a task to transform an image of a person  $a_i$  and some parametrization of glasses  $z_k$  into an image such that: (i) the image contains a realistic looking person, (ii) the person in the image is wearing glasses, (iii) the input of any other person  $a_j$  with that vector  $z$  would produce an image of the other person wearing *the same glasses*, and (iv) input of any other glasses  $z_l$  would produce the same

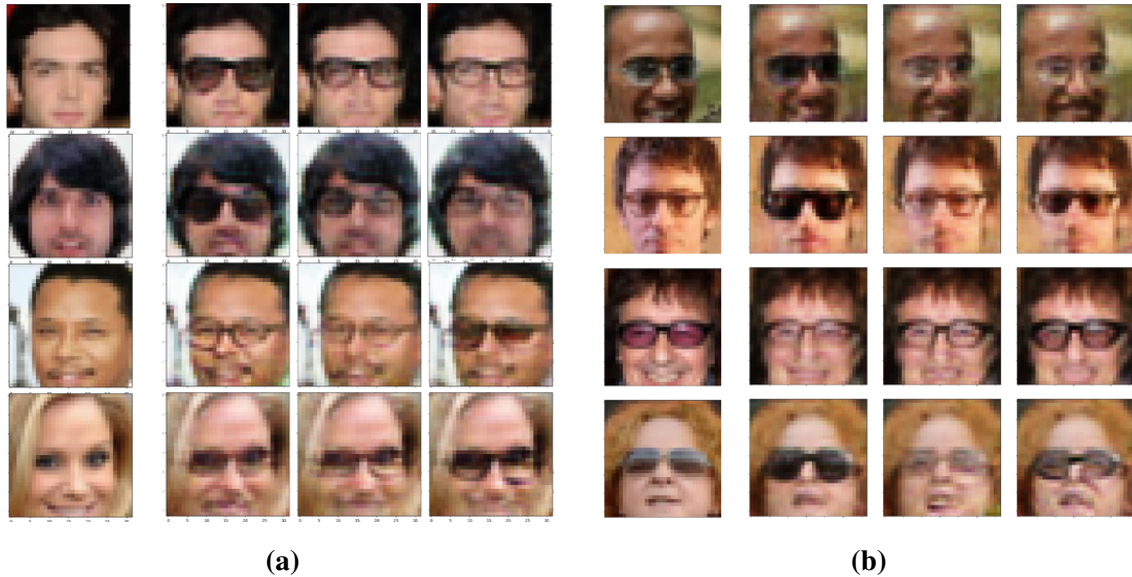
person wearing *different* glasses. There are many different conditions here, all enforced by just two path-equivalence relations.

We highlight that none of the networks were told that images contain people, glasses, or objects of any kind.

### 5.2.1. Evaluation of Trained Networks

Similarly as with the circle dataset, we investigate if these networks generalize well. We examine three things: (i) whether it is possible to *generate an image of a specific person wearing specific glasses*, (ii) whether we can *change* glasses a person wears by changing the corresponding latent vector, and (iii) whether the same latent vector corresponds to the same glasses, independently of the person we pair it with.

We recall that  $AB$  refers to the composition,  $A$  to faces and  $Z$  to parametrization of glasses.



**Figure 5.6:** Parametrically *adding* glasses (a) and *changing* glasses (b) on a person’s face. (a): the leftmost column shows a sample from the dataset  $a_i \in D_E(A)$ . Three rightmost columns show the result of  $c(a_i, z_j)$ , where  $z_j \in [0, 1]^{100}$  is a randomly sampled latent vector. (b): leftmost column shows a sample from the dataset  $(ab)_i \in D_E(AB)$ . Three rightmost columns show the image  $c(\pi_A(d((ab)_i)), z_j)$  which is the result of changing the glasses of a person. The latent vector  $z_j \in [0, 1]^{100}$  is randomly sampled.

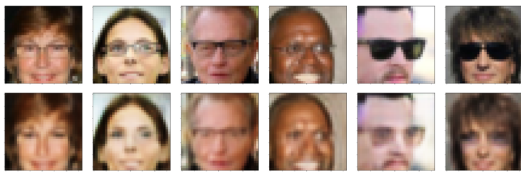
In Figure 5.6 (left) we show the model learns the task (i): generating image of a specific person wearing glasses. Glasses are parametrized by the latent vector  $z \in D_E(Z)$ . The model learns to warp the glasses and put them in the right angle and size, based on the shape of the face. This can especially be seen in Figure 5.8, where some of the faces are seen from

an angle, but glasses still blend in naturally. Figure 5.6 (right) shows the model learning task (ii): *changing* the glasses a person wears.

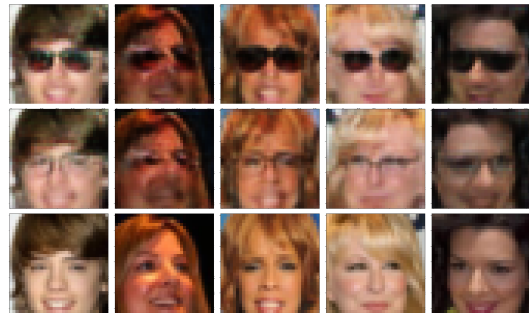
Figure 5.7 shows the model can learn to *remove* glasses. Observe how in some cases the model did not learn to remove the glasses properly, as a slight outline of glasses can be seen.

An interesting test of the learned semantics can be done by checking if a specific randomly sampled latent vector  $z_j$  is consistent across different images. Does the resulting image of the application of  $g(a_i, z_j)$ , contain the same glasses as we vary the input image  $a_i$ ? The results for that third task (iii) are shown in Figure 5.8. It shows how the network has learned to associate a specific vector  $z_j$  to a specific type of glasses and insert it in a natural way.

We note low diversity in generated glasses and a slight loss in image quality, which is due to suboptimal architecture choice for neural networks. It can be seen that the network has learned to preserve all the main facial features.



**Figure 5.7:** Top row shows samples  $(ab)_i \in D_E(AB)$ . Bottom row shows the result of a function  $\pi_A \circ d : AB \rightarrow A$  which removes the glasses from the person.



**Figure 5.8:** Bottom row shows true samples  $a_i \in D_E(A)$ . Top two rows show the image  $c(a_i, z_j)$  of adding glasses with a specific latent vector  $z_1$  for the topmost row and  $z_2$  for the middle row. Observe how the general style of the glasses stays the same in a given row, but gets adapted for every person that wears them.

### 5.3. StyleGAN

In previous sections we have used a dataset  $D_E$  of a concept  $\mathcal{C}$  in some *image space* to train neural networks. We now consider the same schema in Figure 3.1 (b), but conceive a different notion of a dataset to test the product task on – dataset embedded in some *latent space* of a trained generator.

We turn our attention to a GAN architecture called StyleGAN (Karras et al., 2018). We do not attempt to train it but merely use its pretrained generator. This generator was trained

on a dataset of human faces and is capable of producing high-resolution images of stunning quality, offering a lot of variety in terms of pose, age, ethnicity, lightning and image background. The architecture of the generator – a complex neural network with a total of 26.2 million parameters – will be referred to as a function  $g : [0, 1]^{512} \rightarrow [0, 1]^{1024 \times 1024 \times 3}$ .

Given the large variety of generated images, we turn our attention to two specific classes of images. Observe that there exist two subsets of generated people: those that are wearing some sort of headwear (such as glasses or a hat) and those that are not. Just as in Section 5.2, we focus to glasses and make these statements more precise. A subset  $\mathcal{C}(AB) \subseteq [0, 1]^{512}$  of latent vectors are mapped to the set of images of a people *wearing glasses* and another subset  $\mathcal{C}(A) \subseteq [0, 1]^{512}$  is mapped to a set of images of people *not wearing glasses*.

In Section 5.2 we have shown that learning from samples in the *image space* yields promising results. In this section we hypothesize that *learning from samples in the latent space* of a trained generator can produce similar results. As long as the dataset corresponding to the product task is a dataset of any well-defined concept  $\mathcal{C}$ , we observe that the product task is *independent* of the semantic interpretation of such a dataset. Namely, we conjecture that any structure between points in the *image space* is the image of this structure in the *latent space*. For example, it is well-known that interpolation between two points in the latent space yields realistic looking results when mapped to the image space.

As such we embark on a task of collection of data points embedded in the *latent space*. These data points correspond to natural images of a chosen class after they are put through the trained generator. In other words, we collect two datasets of points which are embedded *in the latent space*. The datasets are collected as follows. We randomly sample a vector  $z$  from the latent space, put it through the trained StyleGAN generator and display the result  $g(z)$  on the screen as a  $1024 \times 1024$  RGB image. By creating a script which does that for many points at once and prompts us to classify the result as an image of a person *wearing glasses* or a person *not wearing glasses*, we were able to create two such datasets,  $D_E(AB)$  and  $D_E(A)$ . Sizes of these datasets were 356 and 1576, respectively. This class imbalance stems from the fact that for a randomly sampled vector  $z$ , StyleGAN is about 5 times more likely to produce an image of a person without, than with glasses.

We highlight that, unlike regular datasets, the elements of these datasets only have an interpretation *given an already trained generator*  $g$ . In contrast to collected datasets embedded in  $\mathbb{R}^{1024 \times 1024 \times 3}$  these collected points in  $\mathbb{R}^{512}$  do not have a canonical interpretation as an image of a person, but rather depend on the trained generator  $g : [0, 1]^{512} \rightarrow [0, 1]^{1024 \times 1024 \times 3}$ .

Usage of a pretrained generator allows us to separate two notions: (i) learning the product task, and (ii) learning to generate realistic images of a given class. By using a large pretrained network such as StyleGAN to do the actual image generation, the available computing power can be used solely for learning the isomorphism in Figure 3.1 (b). As a result, all generated



images *look realistic from the start*. The composition and decomposition networks quickly learn to constrain their image to the part of the latent space StyleGAN generator is trained on, thus making the output as realistic as an average sample of StyleGAN. The rest of the time is spent learning *which subset* of the latent space they should focus on. This translates to choosing what type of person and glasses will be generated. Before describing the details of this experiment, we show some a selection of results in Figure 5.9.



**Figure 5.9:** Top row: input image of a person with glasses; Bottom row: output after the glasses have been removed. The entire learning is done in the *latent space* of StyleGAN and these images each depict a latent vector  $z$  after it has been applied to the generator  $g$ . Note how the network has learned to preserve facial features, without being specialized in any way for faces.

We now restate the problem setting. Given the isomorphism in Figure 3.1 (b) we supply two datasets  $D_E(AB)$  and  $D_E(A)$  of latent vectors of people with and without glasses, respectively. This dataset fixes the embedding  $\text{Free}(G) \xrightarrow{E} \text{Set}$  and thus the possible choice of architectures. The *composition* network architecture is an element of  $\text{Para}(\text{Euc})(\mathbb{R}^{512} \times \mathbb{R}^{512}, \mathbb{R}^{512})$ , while the *decomposition* network architecture is an element of  $\text{Para}(\text{Euc})(\mathbb{R}^{512}, \mathbb{R}^{512} \times \mathbb{R}^{512})$ .

As these networks are generally orders of magnitude smaller than networks in Section 5.2, we had enough computing power at our disposal to try out a variety of different network architectures and hyperparameter choices. For both the generators and discriminators we have tested fully-connected, convolutional and mixed architectures, ranging from just a couple of thousand parameters to about half a million parameters. We have also tried a range of choices of hyperparameters. That includes the learning rate, batch size,  $\gamma$  (Definition 14) and  $n_{critic}$ .

Using this dataset we have obtained interesting, but generally limited and inconsistent results. Namely, we have not succeeded in finding a model that solves the task in its entirety. All the trained models seem to have trouble learning to optimize for both the adversarial and



path-equivalence loss. Either the models learn to preserve facial features or they learn to remove and insert glasses, but never both. We especially note difficulties with the insertion of glasses into the image.

In the Figures 5.11, 5.12, and 5.10 we use  $c : \mathbb{R}^{512} \times \mathbb{R}^{512} \rightarrow \mathbb{R}^{512}$  and  $d : \mathbb{R}^{512} \rightarrow \mathbb{R}^{512} \times \mathbb{R}^{512}$  to refer to a trained decomposition and composition network, respectively. Even though in some samples we can see the model has learned to do an almost perfect job, we have not been able to obtain uniformly positive results.

We hypothesize that the lack of these results could be for three reasons. The first is that the product task regularization condition is not always sufficient to successfully train the networks. Given some coherent choice of datasets, perhaps it is necessary to impose extra inductive biases to steer the training into the right direction. The second reason could be that any structure between points in the *image space* is not a preservation of this structure in the latent space. Even though the generator produces realistic-looking images of people with and without glasses, perhaps the coherence of this mapping is too strong of an assumption. We hypothesize that initial training of the generator *with* this regularization condition could improve the results. The third possible reason is that we simply have not stumbled upon the right choice of hyperparameters. This would be a sign that the proposed network is not as robust as expected.



**Figure 5.10:** Top row: input image  $a_i$  of a person *without* glasses. Every other row: the image  $c(a_i, z_j)$ , where  $z_j$  is a per-row randomly sampled latent vector of glasses. Observe the failure of the model to learn to insert glasses on many of the faces. Even when it does learn to insert the glasses, it generates only minimal variation in glass shape, size and color, contrary to the expected. Furthermore, even though  $z_j$  is fixed for each of the last three rows, the generated images contain substantially *different* glasses for each person in a row.





(a)



(b)

**Figure 5.11:** All subfigures depict the same mode on different input faces. Top row: input image  $(ab)_i$  of a person with glasses; 2nd row: image  $\pi_A(d((ab)_i))$  of that person after the glasses have been removed. Every other row:  $c(\pi_A(d((ab)_i)), z_j)$  where  $z_j$  is a per-row randomly sampled latent vector of glasses. Observe the failure of the model to coherently vary the generated glasses on a person as we change the latent vector  $z_j$ . Furthermore, even though  $z_j$  is fixed for each of the last three rows, the generated images contain substantially *different* glasses for each person in a row.





(a)



(b)

**Figure 5.12:** All subfigures depict the same mode on different input faces. Top row: input image  $(ab)_i$  of a person with glasses; 2nd row: image  $\pi_A(d((ab)_i))$  of that person after the glasses have been removed. Every other row:  $c(\pi_A(d((ab)_i)), z_j)$  where  $z_j$  is a per-row randomly sampled latent vector of glasses. Observe the failure of the model to coherently vary the generated glasses on a person as we change the latent vector  $z_j$ . Furthermore, even though  $z_j$  is fixed for each of the last three rows, the generated images contain substantially *different* glasses for each person in a row.

## 5.4. Experiment Summary

With these three experiments, we have demonstrated the practical results obtained from careful ascend into abstract categorical structures. In two out of the three experiments we obtained promising results. Even though the last experiment contains some samples of high quality in which careful removal and insertion of glasses can be seen, we highlight that we have not been able to obtain uniformly positive results.

With these experiments in mind, we believe the practical potential of the product task has been shown. Simply by enforcing high-level rules about what “removal” and “insertion” mean in the product task, we have been able to train networks to remove and insert objects in images without ever being told anything about the nature of these images or objects in them. Moreover, we highlight this is done with *unpaired data*. As such, the high-level regularization given by these path-equivalence losses proves to be a useful training signal applicable to a wide variety of situations.

### **Note on categorical schemas**

In this thesis our focus was on two simple schemas already given to us by existing neural networks in Figure 3.1. In both cases there exist clearly interpretable semantics and trainable network instances. Even though the set of all possible schemas seems vast and potentially interesting, we did not explore any other specific schemas. It remains to be seen whether any other schemas capture new and interesting semantics whose network instances can be efficiently trained.

## 6. Conclusion

In this thesis we have begun to draw an outline of the rich categorical structure underpinning deep learning. Even though the category theory used is elementary, it is a versatile tool which we use to do many things, including the provision of a tangible result. As such category theory did not just help with formality but has guided our thought toward interesting questions.

We highlight that this endeavour is far from finished, for two reasons. The first reason is that, even this narrow domain, there are plenty of structures left to be put in their rightful place. Many of the constructions – especially those towards the end of the thesis – lose their pure categorical flavor. The second reason is that we focus on a narrow subfield of deep learning, namely generative adversarial models in the domain of computer vision. Other areas such as recurrent neural networks, variational autoencoders, optimization, and meta-learning are beyond the scope of this thesis and provide plausible directions in which this study could be continued.

We have not fully reaped the payoff of the categorical approach. It is our hope that in the coming years this thesis can serve as a part of the foundation of further work which uses category theory to reason about the way we can teach machines to reason.

# BIBLIOGRAPHY

- Samson Abramsky and Bob Coecke. A categorical semantics of quantum protocols. *arXiv e-prints*, art. quant-ph/0402130, Feb 2004.
- Samson Abramsky and Bob Coecke. Categorical quantum mechanics. *arXiv e-prints*, art. arXiv:0808.1023, Aug 2008.
- Amjad Almahairi, Sai Rajeswar, Alessandro Sordani, Philip Bachman, and Aaron C. Courville. Augmented cyclegan: Learning many-to-many mappings from unpaired data. *CoRR*, abs/1802.10151, 2018. URL <http://arxiv.org/abs/1802.10151>.
- Marcin Andrychowicz, Misha Denil, Sergio Gomez Colmenarejo, Matthew W. Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. *CoRR*, abs/1606.04474, 2016. URL <http://arxiv.org/abs/1606.04474>.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *arXiv e-prints*, art. arXiv:1701.07875, January 2017.
- John C. Baez and Brendan Fong. A Compositional Framework for Passive Linear Networks. *arXiv e-prints*, art. arXiv:1504.05625, April 2015.
- John C. Baez and Mike Stay. Physics, Topology, Logic and Computation: A Rosetta Stone. *arXiv e-prints*, art. arXiv:0903.0340, Mar 2009.
- Tai-Danae Bradley, Martha Lewis, Jade Master, and Brad Theilman. Translating and Evolving: Towards a Model of Language Change in DisCoCat. *arXiv e-prints*, art. arXiv:1811.11041, November 2018.
- Eugenia Cheng. Higher-dimensional category theory. 2000. URL <http://cheng.staff.shef.ac.uk/misc/4000.pdf>.
- Bob Coecke and Aleks Kissinger. *Picturing Quantum Processes: A First Course in Quantum Theory and Diagrammatic Reasoning*. Cambridge University Press, 2017. doi: 10.1017/9781316219317.

- Bob Coecke, Mehrnoosh Sadrzadeh, and Stephen Clark. Mathematical Foundations for a Compositional Distributional Model of Meaning. *arXiv e-prints*, art. arXiv:1003.4394, March 2010.
- Jared Culbertson and Kirk Sturtz. Bayesian machine learning via category theory. *arXiv e-prints*, art. arXiv:1312.1445, December 2013.
- Conal Elliott. The simple essence of automatic differentiation (differentiable functional programming made easy). *CoRR*, abs/1804.00746, 2018. URL <http://arxiv.org/abs/1804.00746>.
- Michael Fleming, Ryan Gunther, and Robert Rosebrugh. A database of categories. *Journal of Symbolic Computation*, 35(2):127–135, 2003.
- Brendan Fong. Causal Theories: A Categorical Perspective on Bayesian Networks. *arXiv e-prints*, art. arXiv:1301.6201, January 2013.
- Brendan Fong and David I Spivak. Seven Sketches in Compositionality: An Invitation to Applied Category Theory. *arXiv e-prints*, art. arXiv:1803.05316, March 2018.
- Brendan Fong, David I. Spivak, and Rémy Tuyéras. Backprop as functor: A compositional perspective on supervised learning. *CoRR*, abs/1711.10455, 2017. URL <http://arxiv.org/abs/1711.10455>.
- Neil Ghani, Jules Hedges, Viktor Winschel, and Philipp Zahn. Compositional game theory. *arXiv e-prints*, art. arXiv:1603.04641, March 2016.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. in Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, urednici, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014. URL <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved Training of Wasserstein GANs. *arXiv e-prints*, art. arXiv:1704.00028, March 2017.
- Kenneth D. Harris. Characterizing the invariances of learning algorithms using category theory. *arXiv e-prints*, art. arXiv:1905.02072, May 2019.
- Jules Hedges. Coherence for lenses and open games. *arXiv e-prints*, art. arXiv:1704.02230, April 2017a.



- Jules Hedges. Morphisms of open games. *arXiv e-prints*, art. arXiv:1711.07059, November 2017b.
- Jules Hedges. On compositionality. 2017. URL <https://julesh.com/2017/04/22/on-compositionality/>.
- Jules Hedges and Martha Lewis. Towards Functorial Language-Games. *arXiv e-prints*, art. arXiv:1807.07828, July 2018.
- B. Jacobs. *Categorical Logic and Type Theory*. Number 141 u Studies in Logic and the Foundations of Mathematics. North Holland, Amsterdam, 1999.
- Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, and Koray Kavukcuoglu. Decoupled neural interfaces using synthetic gradients. *CoRR*, abs/1608.05343, 2016. URL <http://arxiv.org/abs/1608.05343>.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *CoRR*, abs/1812.04948, 2018. URL <http://arxiv.org/abs/1812.04948>.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. in *ICLR*, 2015.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. in *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- Saunders MacLane. *Categories for the Working Mathematician*. Springer-Verlag, New York, 1971. Graduate Texts in Mathematics, Vol. 5.
- Samuel Mimram and Cinzia Di Giusto. A Categorical Theory of Patches. *arXiv e-prints*, art. arXiv:1311.3903, November 2013.
- Robert Rosebrugh and RJ Wood. Relational databases and indexed categories. in *Proceedings of the International Category Theory Meeting 1991, CMS Conference Proceedings*, volume 13, pages 391–407, 1992.
- Patrick Schultz, David I. Spivak, Christina Vasilakopoulou, and Ryan Wisnesky. Algebraic Databases. *arXiv e-prints*, art. arXiv:1602.03501, February 2016.
- Pawel Sobocinski. ..., a monoid is a category, a category is a monad, a monad is a monoid, .... 2017. URL <https://graphicallinearalgebra.net/2017/04/16/a-monoid-is-a-category-a-category-is-a-monad-a-monad-is-a-monoid/>.

David I. Spivak. Functorial data migration. *CoRR*, abs/1009.1166, 2010. URL <http://arxiv.org/abs/1009.1166>.

Joel Spolsky. The law of leaky abstractions. 2002. URL <https://www.joelonsoftware.com/2002/11/11/the-law-of-leaky-abstractions/>.

Guillermo Valle-Perez, Chico Q. Camargo, and Ard A. Louis. Deep learning generalizes because the parameter-function map is biased towards simple functions. in *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rye4g3AqFm>.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *CoRR*, abs/1703.10593, 2017. URL <http://arxiv.org/abs/1703.10593>.

## Kompozicijsko duboko učenje

### Sažetak

Neuronske mreže postaju sve popularniji alat za rješavanje mnogih problema iz stvarnoga svijeta. Neuronske mreže općenita su metoda za diferencijabilnu optimizaciju koja uključuje mnoge druge algoritme strojnog učenja kao specijalne slučajeve. U ovom diplomskom radu izloženi su početci formalnog kompozicijskog okvira za razumijevanje različitih komponenta modernih arhitektura neuronskih mreža. Jezik teorije kategorija korišten je za proširenje postojećeg rada o kompozicijskom nadziranom učenju na područja nenadziranog učenja i generativnih modela. Prevođenjem arhitektura neuronskih mreža, skupova podataka, mape parameter-funkcija i nekolicine drugih koncepata iz dubokog učenja u kategorijski jezik, pokazano je da se optimizacija može raditi u prostoru funktora između dvije fiksne kategorije umjesto u prostoru funkcija između dva skupa. Dajemo pregled znakovite poveznice između formulacije dubokog učenja u ovom diplomskom radu i formulacije kategorijskih baza podataka. Nadalje, koristimo navedenu kategorijsku formulaciju kako bismo osmislili novu arhitekturu neuronskih mreža kojoj je cilj naučiti umetanje i brisanje objekata iz slike sa neuparenim podacima. Ispitivanjem te arhitekture na tri skupa podataka dobivamo obećavajuće rezultate.

**Ključne riječi:** Neuronske mreže, duboko učenje, teorija kategorija, generativne suparničke mreže

# Compositional Deep Learning

## Abstract

Neural networks have become an increasingly popular tool for solving many real-world problems. They are a general framework for differentiable optimization which includes many other machine learning approaches as special cases. In this thesis we lay out the beginnings of a formal compositional framework for reasoning about a number of components of modern neural network architectures. The language of category theory is used to expand existing work on compositional supervised learning into territories of unsupervised learning and generative models. By translating neural network architectures, datasets, parameter-function map, and a number of other concepts to the categorical setting, we show optimization can be done in the functor space between two fixed categories, rather than functions between two sets. We outline a striking correspondence between the deep learning formulation in this thesis and that of categorical database systems. Furthermore, we use the category-theoretic framework to conceive a novel neural network architecture whose goal is to learn the task of object insertion and object deletion in images with unpaired data. We test the architecture on three different datasets and obtain promising results.

**Keywords:** Neural networks, deep learning, category theory, generative adversarial networks