A Pattern Language for Machine Learning Tasks

Anonymous authors

Paper under double-blind review

Abstract

We formalise the essential data of objective functions as equality constraints on composites of learners. We call these constraints "tasks", and we investigate the idealised view that such tasks determine model behaviours. We develop a flowchart-like graphical mathematics for tasks that allows us to; (1) design and optimise desired behaviours model-agnostically; (2) offer a unified perspective of approaches in machine learning across domains; (3) import insights from theoretical computer science into practical machine learning. As a proof-of-concept of the potential practical impact of our theoretical framework, we exhibit and implement a novel "manipulator" task that minimally edits input data to have a desired attribute. Our modelagnostic approach achieves this without the need for custom architectures, adversarial training, random sampling, or interventions on the data, hence enabling capable, small-scale, and training-stable models.

1 INTRODUCTION

The primary instrument for controlling the training of machine learning (ML) models is the objective function, which can be broken into three modular parts. Let Θ_e, Θ_d be the parameter-spaces of a model enc and dec respectively. Then the reconstruction loss that characterises an autoencoding task amounts to minimising (for $\theta_e \in \Theta_e, \theta_d \in \Theta_d$) the following objective function:

030

000

001 002 003

005 006 007

009

010

011

012

013

014

015

016

017

018

019

020

021

023

024

031

$\operatorname{argmin}_{\theta_e,\theta_d} \left(\mathbb{E}_{x \sim \mathcal{X}} [\mathbf{D} \left(\operatorname{dec}_{\theta_d} (\operatorname{enc}_{\theta_e}(x)), x \right)] \right)$

We take the expected value over a data distribution \mathcal{X} of a measure of statistical divergence D (such as cross-entropy or log-likelihood) of two expressions that, under ideal conditions, should be equal: the decoding of the encoding of some data x, and the original x. In this work, we focus on the two expressions that want to be equal, and we call this equational constraint a *task*.

In practice, designing a good objective function incorporates many technical choices, such as choice of architecture, measure of statistical divergence, and training data (Ciampiconi et al., 2023; Richardson, 2022; Terven et al., 2023). However, these choices are often made by heuristics, or rationalised post hoc. While such choices are sometimes required to make training tractable, they are not always relevant to understanding the final behaviour of the trained model. Instead, we propose and investigate the idealised view that:

- 043
- 044

051

To reason about tasks compactly, we use flowchart-like string diagram notation.

Example 1.1. For a hyperparameter choice of divergence **D**, where *X* is an input datatype, and *LAT* is the datatype of the latent space, the two components of the empirical risk minimisation of the autoencoder *task* consist of (1) applying the encoder enc (typed $X \rightarrow LAT$) followed by the decoder dec (typed $LAT \rightarrow X$) to inputs *x* drawn from a source of data \mathcal{X} over the datatype *X*, *which should be equal to* (2) the original *x*. We depict this as:

Tasks determine model behaviour.

$$\mathbb{E}_{x \sim \mathcal{X}} \left[\mathbf{D} \left(\underbrace{\operatorname{dec}_{\theta_d}(\operatorname{enc}_{\theta_e}(x))}_{(1)}, \underbrace{x}_{(2)} \right) \right] \quad \Leftrightarrow \quad x \underbrace{ \underbrace{ \underset{dec}{\overset{X}{\underset{dec}}}_{dec} \underbrace{\operatorname{dec}}_{d} \underbrace{\operatorname{dec}}_{x}}_{\operatorname{dec}_{\theta_e}(x))} \stackrel{\operatorname{dec}}{\xleftarrow} x \underbrace{ \underset{x}{\overset{X}{\underset{dec}}}_{x}}_{x}$$

The diagrammatic notation is formally equivalent to the traditional symbolic notation with the addition of type-information about inputs and outputs. While the formal details (Appendix A) involve category theory, the power of string diagrams lies in their intuitive visual nature: by reading the diagrams as flowcharts from left to right, practitioners can leverage these diagrams to reason about ML tasks without needing to fully grasp the underlying mathematical formalism. In summary, nodes depicted as various shapes are functions, and wires are datatypes which can be understood as carrying information.

To explore the expressive power of tasks, we abstract away implementation details such as
 architecture and training by idealising models to be *universal function approximators* that can, in principle, perfectly optimise objective functions. Thus each task can be viewed
 purely as an equational constraint on the behaviour of the learners, comparable to equational
 constraints on the possible values of variables in algebra. This perspective allows us two
 ways to specialise tasks by imposing structural inductive biases, by specifying architectural
 choices, or by adding additional objectives.

Example 1.2 (Residuation as an architectural choice). Diagrammatically, choosing an architecture means substituting a "black-box" universal function approximator with another diagram with matching input-output type constraints; intuitively, since a universal function approximator can be any function, it can *in particular* be a specific function of the same input-output type if necessary. The formal semantics for such substitutions are provided in Appendix A.2. A simple example is residuation, where a learner *N* of type $X \to X$ is transformed into $N^{res} := x \mapsto N(x) + x$, depicted as:



Example 1.3 (Perceptual losses as multi-objective learning). A common form of multiobjective learning is to add a normalisation or regularisation term to an objective function, for instance, a perceptual loss $\mathbf{L} : Y \to \mathbb{R}^{\geq 0}$ (which we depict with a "white box", because there are no learnable parameters). Multiple tasks may be combined into single objective functions by means of weighted summation with positive hyperparameter-coefficients α , β .



Specialisation of tasks obtains desirable properties without compromising basic behaviours, and allows us to explain the behaviour of models in terms of their constituent tasks. As a canonical example, we may obtain a variational autoencoder (VAE) (Kingma & Welling, 2022) from a regular one by the two kinds of specialisation described.

095 Example 1.4 (VAE). Setting the output type of the encoder to be a space (mean, variance) of 096 parameters for Gaussians; declaring the decoder to be internally structured as the sequential 097 composite of sampling from a Gaussian of the input (mean, variance) followed by a learner; 098 adding an additional normalisation task which requires the outputs of the encoder to be 099 close to (0, 1) — the parameters of the unit Gaussian. enc maps the input space X to (x, σ) , 100 the parameter-space of Gaussians over the latent space *LAT*. The decoder is the composite 101 of a sampling function samp. : $(\mu, \sigma) \to LAT$ and a learner dec : $LAT \to X$. Note the additional NORMALISE task that enc must satisfy. This presentation of VAEs is equivalent to 102 the traditional probabilistic form when the statistical divergence is KL (Rocca, 2021). 103



1.1 Contributions

075

076 077 078

079

080

081

082

083 084 085

088

090 091

092

093

094

104

105

We make three primary contributions in this work.

Our first, theoretical contribution is **the formalisation of a common but informal standard procedure in deep learning**, which may be summarised as the following recipe:

112 113

114

115

- 1. Characterise desired behavior via equational constraints (tasks) between learners
- 2. Implement tasks by treating neural networks as universal approximators
- 3. Convert equations to loss function by a choice of hyperparameters, namely weighted sum of constituent losses and choice of divergences
- 116 117

118 Our second, theoretical contribution is a demonstration that **using our framework**, we may 119 analyse, predict, and design behaviours of models. Using specialisation and algebraic reasoning, we can analyse the behaviour of complex models in terms of simple, well-120 understood tasks we call *patterns*. Moreover, we can define a novel synthetic relationship 121 between tasks called *refinement* which describes when one optimally trained set of tasks 122 entails satisfaction of a different set of tasks. Altogether, we may use these techniques 123 to understand and compare the behaviour of models before committing to potentially 124 costly training. Example 2.6 and Propositions 2.9, 3.2, and 3.3 illustrate the kinds of formal 125 reasoning our language enables.

Our third, practical contribution and proof-of-concept is the **implementation of a novel** 127 problem class we call manipulation (Section 3), which formalises (Bancilhon & Spyratos, 128 1981) the problem of viewing and editing a targeted attribute of data while "leaving other 129 aspects the same". Notably, this task represents the first problem class to our knowledge that 130 appears to be naturally solved by the relatively underexplored technique of "multi-learner 131 multi-objective learning", which our framework naturally accommodates. As examples in 132 image domains, we; change only the colour of a shape (Figure 1); change the value of a 133 handwritten digit without affecting other stylistic properties (Figure 2); and change only 134 whether a person is smiling in an image (Figure 6). Even in these toy domains, we observe a 135 range of benefits we expect to scale: by following our recipe we obtain architecture-agnostic 136 (Table 1) style-transference models without the need for randomness, adversarial training, or modality- and architecture-specific interventions, with good interpolation properties 137 (Figure 6). 138

We conclude by discussing relations to similar approaches in the literature, along with avenues and prospects for further development.

141 142 143

144

145

2 TASKS AND PATTERNS

2.1 Tasks

146 We assume the following contextual data, omitted if there is no confusion. Let X, Y denote 147 datatypes; Σ a set of processes f, each of which has (possibly empty) learnable parameter 148 datatypes \mathfrak{p}_f . An atomic task is a process-theoretic equational constraint on learners 149 specifying that *f* should behave like *g* on all inputs. The objective function of an atomic task 150 φ of type $X \to Y$ equipped with distribution \mathcal{X} corresponds to a map $\mathfrak{p}_{\varphi} \to \mathbb{R}$ that sends 151 $\pi \mapsto \mathbb{E}_{x \sim \mathcal{X}}[\mathbf{D}(\operatorname{sys}_{\varphi;\pi}(x), \operatorname{spec}_{\varphi;\pi}(x))]$ for some choice of statistical divergence **D**. A learner 152 can optimise multiple atomic tasks simultaneously by optimising a combination α of the 153 atomic objective functions (commonly obtained by taking a weighted sum). 154

Definition 2.1 (Tasks). An *atomic task* φ is a tuple $(f, g, \mathcal{X}, \mathfrak{p})$, where $f, g : X \to Y$ are composite processes of Σ , \mathcal{X} is a distribution over X, and $\mathfrak{p} \subseteq \mathfrak{p}_f \oplus \mathfrak{p}_g$ is a space of trainable parameters. We indicate the *system* f and *specification* g as sys_{φ} and $spec_{\varphi}$, and similarly the *domain* X and *codomain* Y as dom(φ) and cod(φ). A *compound task* Φ (or just *task*) is a non-empty set of tasks. As we have seen in previous examples, we notate a task Φ as a collection of atomic tasks $sys_{\varphi} \rightleftharpoons spec_{\varphi}$, where superscripts on the harpoons indicate which learnable parameters are governed by each atomic task. Tasks become concrete objective functions by a hyperparametric choice of divergences for atomic tasks, followed by a combination via a weighted sum with hyperparameter coefficients, or more generally a *compound function*. As such, a particular objective function that instantiates a task is one where the choices for the measure of statistical divergence and compound function have been made. Beyond these hyperparameters, tasks and objective functions may be viewed as informationally equivalent.

2.2 PATTERNS ARE "NICE" TASKS

We do not expect there to be a general and systematic method to translate between natural-language behavioural specifications and tasks; if such a method existed, then all of deep learning would be reduced to hyperparameter search. There are often many different ways to approach a problem in ML, much like there is no single correct way to write software, or design a building. This suggests to us the view of *patterns*: some tasks are well understood, usable modularly and in many contexts, and easily modifiable, and such tasks can be viewed as *design patterns* – borrowing a term from software engineering (Beck & Cunningham, 1987)¹. In this section, we suggest some examples of patterns that correspond to well-studied methods and paradigms in ML, and we show how to use **patterns as an accessible basis to analyse models**.

Pattern 2.3 (classification).



Pattern 2.4 (autoencoding).



Given a data-label pair (d, l) drawn from \mathcal{X} with labels $l \in L$, a *classifier* cls : $D \to L$ is a function that solves the classification task, in which it seeks to reconstruct the label from the data. This can be done by minimising the corresponding objective function $\mathbb{E}_{(d,l)\sim\mathcal{X}}[\mathbf{D}(cls(d), l)]$ for some measure of statistical divergence **D** on the label space, which may be continuous to encompass regression.

As we have seen, given a data distribution \mathcal{X} over X and some latent space *LAT*, an *autoencoder* consists of an *encoder* enc : $X \rightarrow LAT$ and a *decoder* dec : $LAT \rightarrow X$ which cooperate to reconstruct the identity over the observed distribution, by minimising $\mathbb{E}_{x \sim \mathcal{X}}[\mathbf{D}(\det(\operatorname{enc}(x)), x)].$

Pattern 2.5 (GAN). Given a data distribution \mathcal{X} over X and noise distribution \mathcal{N} over N, a *generative adversarial network* (GAN) consists of a *generator* gen : $N \to X$ and a *discriminator* dsc : $X \to [0, 1]$. The prosaic explanation that "the discriminator seeks to distinguish real data from fake data while the generator aims to fool the discriminator" translates directly into a task description: where 1 indicates "real" and 0 indicates "fake", the discriminator dsc seeks to minimise some positive combination of the terms $\mathbb{E}_{x \sim \mathcal{X}}[\mathbf{D}(dsc(gen(x)), 0)]$, while the generator gen seeks to minimise $\mathbb{E}_{x \sim \mathcal{X}}[\mathbf{D}(dsc(gen(x)), 1)]$.



¹And before that, architecture and urban design (Alexander et al., 1977).

216 2.3 Analysing complex tasks

We can analyse the intended functions of models by viewing them as composites of simple patterns. We have already seen in Example 1.4 how a VAE is analysable by inspection as a specialised regular autoencoder. As a second example, a CycleGAN is two GANs on different distributions, whose generators are mutually autoencoders by a *cycle consistency loss* (Zhu et al., 2017). This suggests that the generators encode the distributions into each other in a reversible manner, and indeed this kind of style transfer between distributions is what a CycleGAN does in practice.

Example 2.6 (cycleGAN). Below, *i*, *j* are nonequal indices taking values in {0, 1}, where X_0 and X_1 are different distributions on the same space *X*: typically these are two classes of images.



Tasks allow us to reason "legalistically" to rule out undesirable model behaviours. For instance, in Example 1.4, we can immediately determine that the normalisation term imposes a nontrivial constraint: without the normalisation task, enc and dec could collude against us by only using variance 0 (i.e. deterministic) representations that are far apart, which would lose the ease of sampling and robustness of representations. We further elaborate on the use of this form of reasoning for informing task design in Appendix B.

We can also upgrade informal intuitions into formal derivations. For example, on the account of Ranzato et al. (2007), a broad class of unsupervised learning techniques — including PCA and *k*-means — are specialisations of the energy minimisation task, which may be considered an autoencoding task from a latent "code" space subject to the representations minimising a measure of "energy". We can in fact formalise the relationship between these forms of unsupervised learning and autoencoding, showing that under mild assumptions, they are the same in the computational limit.

Pattern 2.7 (energy minimisation).



energy minimisation consists of; three types of systems: D(ata), C(ode), and $\mathbb{R}^{\geq 0}$; two learnable processes: an encoder enc : $D \rightarrow C$ and a decoder dec : $C \rightarrow D$; two user-supplied *energy functions*: \mathbf{E}_{enc} : $C \times C \rightarrow \mathbb{R}^{\geq 0}$ and \mathbf{E}_{dec} : $D \times D \rightarrow$ $\mathbb{R}^{\geq 0}$; and a user-supplied constant $\gamma \in \mathbb{R}^{\geq 0}$. Provided a distribution of inputs \mathcal{Y} on D, and a distribution \mathcal{Z} on C the system seeks to minimise $\mathbb{E}_{y \sim \mathcal{Y}, \mathbb{Z} \sim \mathcal{Z}}[\gamma \mathbf{E}_{\text{enc}}(\text{enc}(y), z) + \mathbf{E}_{\text{dec}}(y, d(z))]$. Such a task is called *code-extracting* when $\mathcal{Z} = \text{enc}(\mathcal{Y})$.

To formally relate code-extracting energy minimisation and autoencoding, we introduce a relationship between tasks called *refinement*, which states that perfectly solving Φ allows one to construct perfect solutions for Ψ . When Φ and Ψ refine each other, the tasks are "the same in the computational limit"; a perfect autoencoder is a perfect code-extracting energy minimiser, and vice versa. These relationships are a proxy for the behaviour and relative power of concrete implementations of tasks.

Definition 2.8 (Refinement and equivalence of tasks). Task Φ *refines* task Ψ if, by treating the atomic tasks as equations, the processes of Φ may be composed to satisfy the equations of Ψ . Φ and Ψ are *equivalent* if they refine one another, which we denote $\Phi \equiv \Psi$.

Proposition 2.9. If enc and dec are deterministic and E_{enc} and E_{dec} are positive (e.g. metrics or statistical divergences), then energy minimisation \equiv autoencoding. (proof in Section C.1)

270 3 EXPERIMENT: THE MANIPULATION TASK 271

272 In this section, we design and experimentally validate manipulation, a novel task which 273 aims to view and edit a property of data without explicit guidance. There are many 274 models that behave like manipulators, e.g. unsupervised sentiment translation (Li et al., 275 2018; Sudhakar et al., 2019) in the text domain and prompt-based photo editing (Hertz 276 et al., 2022; Kawar et al., 2023) in the image domain. To constrain the behaviour of such a 277 manipulation, we demonstrate another ability of our framework: importing insights from 278 computer science more broadly via the process-theoretic perspective. For manipulation, we reference the field of Bidirectional Transformations, which studies consistency between 279 different overlapping representations of data (Abou-Saleh et al., 2018). A special case is 280 the view-update problem (Bancilhon & Spyratos, 1981) originally proposed for databases: how do we algebraically characterise reading-out and updating an attribute $a \in A$ from 282 some data $d \in D$? There are a family of solutions called *lenses* which are parameterised by 283 algebraic laws of varying strength (Nakano, 2021), which we take inspiration from below: 284

286

287

289

290

291

293

295

296 297 298

299

300 301

302 303

304

306

307

308

310

311

281

Task 3.1 (manipulation). Let \mathcal{X} : (d, a) be a distribution over some data $d \in D$, each labelled with an attribute $a \in A$ and let \mathcal{A} be a distribution over the attributes. A manipulation consists of a pair of operations (get : $D \rightarrow A$, put : $D \times A \rightarrow D$) which can be understood as reading and writing, respectively. In particular, the put edits a reference data point to exhibit the specified attribute. The two operations have to obey the following tasks, with respect to the modeller's choice of distribution \mathcal{A} on A.



To provide an intuition for each of the tasks, assume that the data consists of images each containing a single shape each labeled with the colour of the shape. CLASSIFY allows us to use get to read out the colour of a shape. PutGet says that first editing the colour of a shape (say, from red to blue) and then immediately reading out that colour will return the edited colour (blue). GETPUT says that reading out the colour of a shape (say, red) followed by editing the shape to have the same colour (i.e., an edit that leaves red unchanged) is the same as doing nothing. UNDOABILITY says that edits can be undone; using the first put to change the colour of a shape (say from red to blue), and then editing again with a second put to restore the original, read-out colour of the shape (red) must restore the original image.

312 For completeness, and to illustrate the ergonomic necessity of our diagrammatic notation, we 313 display the formulaically obtained hyperparameterised objective function of manipulation 314 in standard notation below: 315

$$\operatorname{argmin}_{\boldsymbol{\phi},\psi} \left(\underbrace{\alpha \mathbb{E}_{(d,a)\sim\mathcal{X}} \left(\mathbf{D}_{1}^{A} \left(\operatorname{get}_{\boldsymbol{\phi}} \left(d \right), a \right) \right)}_{\operatorname{CLASSIFY}} + \underbrace{\gamma \mathbb{E}_{d\sim\mathcal{X}|d} \left(\mathbf{D}_{1}^{D} \left(\operatorname{put}_{\boldsymbol{\psi}} \left(d, \operatorname{get}_{\boldsymbol{\phi}} \left(d \right), a \right) \right)}_{\operatorname{GetPur}} + \underbrace{\delta \mathbb{E}_{(d,a')\sim\mathcal{X}|d\times\mathcal{A}} \left(\mathbf{D}_{2}^{A} \left(\operatorname{get}_{\boldsymbol{\phi}} \left(\operatorname{put}_{\boldsymbol{\psi}} \left(d, a' \right), a' \right) \right)}_{\operatorname{UNDOABILITY}} \right) \right)$$

Where; ϕ , ψ are the parameters of put and get to be learnt; α , β , γ , δ are hyperparametric pos-323 itive weighting coefficients for each task; \mathbf{D}_1^A , \mathbf{D}_2^A are hyperparametric statistical divergences

for *A*; \mathbf{D}_1^D , \mathbf{D}_2^D are statistical divergences for *D*; and \mathcal{A} is a hyperparametric distribution over *A* such that supp(\mathcal{A}) contains the attributes the modeller intends to have as targets.

The following pair of theoretical results allow us to anticipate and justify the efficacy of manipulation, by analysing its behaviour with respect to an additional regularisation term we detail in Appendix B; in the same way that a VAE may be viewed as a nicely regularised autoencoder, we may learn some things indirectly about manipulation by considering a nicely-regularised variant. First, we may relate the behaviour of manipulation to style-transfer.

Proposition 3.2 (manipulation vs. style-transfer (Take 1)). There exists a regularisation term for manipulation such that an "optimally-trained" manipulation yields an "optimally-trained" CycleGAN. (Elaboration of conditions and proof in Section C.2)

Second, we have a result that informs us that in the computational limit, manipulation
 allows us to transform any classifier of attributes into a "minimally invasive" editor of those
 attributes, without providing any additional information in the form of augmented data or
 inductive biases.

Proposition 3.3 (put as Bayesian inverse of a classifier). There exists a regularisation term for manipulation such that by identifying a given "well-behaved" classifier cls as get, satisfaction of the manipulation task allows put to induce the Bayesian inversion cls[†]. (Elaboration of conditions and proof in Section C.3)

Indeed, as we proceed to report, we obtain such minimal-editors in practice. Moreover, we
do so in a manner that is deterministic (cf. VAEs, Kingma & Welling (2022)), non-adversarial
(cf. GANs, Goodfellow et al. (2020)), and without hardcoding or handcrafting features in
latent spaces (cf. CVAEs, Shaikh et al. (2022)).

3.1 Experimental Results I: Simple attributes of synthetic and real-world data

N.B. For space, only results are reported in the main body, while all methodological details are reported in Appendix D; as we follow our recipe, these details consist chiefly of hyperparameter choices.

We demonstrate initial proofs-of-concept of the manipulation task on a simple analytic dataset (Figure 1) inspired by Spriteworld (Watters et al., 2019), and on MNIST (Figure 2). In the former, each image depicts a single shape with varying properties, and is labelled by two attributes: shape – circle, square or triangle – and colour – red, green or blue. For each attribute, we train a get/put pair according to the manipulation task specification.



Figure 1: An input Spriteworld image alongside a spectrum of outputs exhibiting the ability of the put to manipulate a single attribute of the input while preserving its other properties. Additionally, the model is able to generalise by interpolating to attribute values unseen during training, in this case producing orange and cyan shapes, whereas during training, it only sees red, green or blue shapes. (further details in Section D.1)

372 373 374

349

350 351

352

354

355

356

357

358

360 361 362

368

370

3.2 Experimental Results II: Derived attributes of synthetic data

Often in practice we are interested in complex, non-explicit attributes that are derived
from those labelled in the data: for example, "eligibility for a loan" may be derived from
other explicit attributes of people in a database by an operationally opaque classifier, with

 Original 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 4
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 4
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 1
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 6
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 7
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 5
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 5
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 5
 0
 1
 2
 3
 4
 5
 6
 7
 8
 9

 4
 0
 1
 2
 3
 4
 5
 6
 7
 8
 <td

Figure 2: Outputs of a put trained against an MNIST classifier. The put preserves several graphological aspects, such as stroke weight, slant, and angularity. This represents qualitative evidence to support our prediction that put as a class-conditioned generative model behaves as a style-preserving edit.

unknown range, distribution, and dependencies on other attributes. A known challenge in manipulating derived attributes is unequal entropy in attribute classes (Chu et al., 2017), which may cause models such as CycleGANs to hide data imperceptibly, making them particularly vulnerable to adversarial attacks. Various solutions have been proposed, including masks (Wu et al., 2024), blurring (Fu et al., 2019) and compression (Dziugaite et al., 2016). We demonstrate via a modification of manipulation (Task 3.1, Figure 4) that our framework permits the design and implementation of end-to-end approaches to editing complex attributes without interventions on the data.



Figure 3: To illustrate the concepts of derived attributes and unequal entropy, consider an attribute on the Spriteworld data called *blue-circleness*, which broadly measures how similar a shape is to a blue circle. We define *blue-circleness* (*bc*) as a function of explicit attributes *shape* and *colour*; we assign a continuous colour score $cs \in [0, 1]$ based on the hue, where red = 0 and blue = 1. To illustrate unequal entropy in this example, the class 0 has higher entropy than 0.4 because there are more shapes that have *bc*-value 0. So manipulating a shape with *bc*-value 0 to 0.4 must lose information.



 Inspired by the complement of symmetric lenses (Hofmann et al., 2011), we introduce a complement *C* to put, changing its type to $S \times L \times C \rightarrow S \times C$. Let $(d, a) \sim \mathcal{X}$ be a distribution over some data $d \in D$, each labelled with an attribute $a \in A$. complement manipulation consists of a pair (get $: D \rightarrow A$, $\overline{put} : D \times A \times C \rightarrow \mathcal{D} \times C$) fulfilling the rules on the left. The idea of the complement is that it provides the manipulation with a scratchpad *C* to keep track of additional data. As none of the tasks check the output of the complement, the put and get can use it freely to store relevant data.



Figure 4: Complement manipulators (Task 3.4) can manipulate derived attributes such as *blue-circleness*, by using the complement as a scratchpad to record a correspondence between data points (further details in Section D.1) while preserving attributes such as position and size.

3.3 Experimental results III: Interpretability applications on real-world data

As a further test of the robustness of tasks to implementation choices, we specialised the put to be a simple vector addition in the latent space of an autoencoder (Figure 5), for the relatively complex *Large-scale CelebFaces Attributes* dataset. We found that restricting put to be linear in this way increased training stability. Moreover, in the same way that we would expect a latent space "filtered through" the probabilistic structure of Gaussians to yield good sampling properties (Example 1.4), we would predict that enforcing linear structure on the latent space would yield "linear" properties. Indeed, we exhibit continuous interpolation in generated outputs between normally discrete class labels (Figure 6), and class-sensitive separation of latent space embeddings in the autoencoder. We consider this to be compelling evidence that **since our framework is agnostic, implementation details may be engineered to obtain additional desirable properties without compromising behavioural specifications.**



Figure 5: Recalling that architectural choices are a form of specialisation by diagrammatic
substitution, the linear put is a specialisation of a generic put as an autoencoder task-bound
pair of learners enc and dec, along with a put' that computes single shift vector to be
added into the latent space, depending only on the label value. enc, dec, and put' are
trained simultaneously along with the manipulation tasks, and intuitively this pressures
the autoencoder pair to adapt their latent representations to fit the needs of the broader
manipulation task.



Figure 6: *Left*: Outputs of a linear manipulator trained on face image data paired with a binary "smile"/"no-smile" label. The remarkable aspect of this experiment is that the original data only carried binary smile/no-smile labels, and that the linear structure in the specialisation of the put admits continuous interpolation. *Right*: A comparison of the spread of the latent embeddings of images from the validation set when pre-training an autoencoder and then training a (linear) classifier on the latent space (top), vs. when trained with a linear manipulator (bottom). We find that linear put automatically separates latent space embeddings of classes : the graphs depict the relative density of embeddings along the direction of the classifiers' weight vectors, normalised so that each combined data spread is centred and has unit variance (details in Section D.2).

4 Concluding discussion

4.1 SUMMARY

512 513

499

500

501

503

505

506

507

509 510

511

514 515

We introduced a diagrammatic language for representing and reasoning about the
behaviour of machine learning models in terms of tasks, viewed as the essential data of
objective functions. By leveraging category theory and string diagrams, our work establishes
a cross-disciplinary formal bridge between theoretical computer science and practical
machine learning, providing new conceptual tools for analyzing ML systems and permitting
the transfer of insights between traditionally separate fields.

The proposed framework allows capturing existing tasks in machine learning, providing intuitive insights rooted in mathematical rigour. We identify a set of widespread and well-understood tasks, which we call patterns. We can analyse some tasks as composites of patterns (Example 2.6) while other tasks can be understood as specialisations of patterns (Example 1.4). The rewrite system inherent to string diagrams, allows us to identify relationships between different tasks and formalise intuitions (Proposition 2.9, Proposition 3.2).

Beyond theoretical insights, the proposed language also allows the creation of new training 528 **paradigms.** As preliminary empirical validation of our theory's utility and potential, we 529 introduced a novel task type called manipulator that produces a class-conditioned and style-530 preserving generative model counterpart for a given classifier. In the image domain, we were able to verify predicted behaviours (Section 3.1), and we demonstrated the ability to design 532 novel end-to-end capabilities, such as end-to-end editing of complex attributes (Section 3.2) and the imposition of linear structure on latent space representations (Section 3.3), which 534 allowed continuous interpolation between discrete class labels on real data, and separated latent space embeddings of different classes. Notably, this was achieved without adversarial training conditions, random sampling, preprocessing of data, or hardcoded interventions in the architectures, i.e.: 537

538 539

Our framework enables capable, small-scale, and training-stable models.

540 4.2 Relation to extant work

542 Regarding our nascent theoretical framework as a whole, the style of engineering beginning 543 from tasks is already common practice in many fields of ML, and we sought here to place these practices on a more rigorous footing, and to probe their strengths and limitations. Our 544 mathematical approach draws broadly from the field of Applied Category Theory (Fong & Spivak, 2019a), particularly in the use of string diagrams for the higher-algebraic data 546 of concurrently and serially composed functions, which enables compact representation 547 and reasoning with otherwise cumbersome symbolic equivalents when dealing with 548 multiple learners in tandem. To our knowledge, our concern with the composition of tasks 549 among many learners distinguishes our aims and formal choices from approaches that 550 employ similar mathematical formulations, both within the category-theoretic literature (cf. Gavranović et al. (2024)), and without (cf. the variational generalisation of Bayesian 552 inference presented in Knoblauch et al. (2022)).

553 554

While our approach is essentially neurosymbolic in spirit, it does not fit neatly into the mainstream triad of neurosymbolic approaches (d'Avila Garcez & Lamb, 2020); we do not encode symbolic data for neural operations, nor do we interface neural approaches with symbolic engines, nor are we hardcoding expert knowledge representations. Moreover, our aims differ: while neurosymbolic approaches often seek to manipulate symbolic data systematically by neural means, our framework operates at a higher level of abstraction, seeking to use the systematicity of higher-algebraic equational characterisations as a means to shape the neural ends. Hence our perspective may complement existing approaches to structure in machine learning.

563

Regarding manipulation in particular, this was to our knowledge the first practically demonstrated synthesis of insights from Bidirectional Transformations as a subfield of database
theory (Abou-Saleh et al., 2018) with ML. While explicitly neurosymbolic approaches have
been tried for similar editing tasks before (see e.g. Smet et al. (2023)), owing to the influence
of database theory in our approach, to our knowledge our statement and execution of this
task enjoys the maximal permissible generality and implementation agnosticism among
similar attribute-editing tasks, without sacrificing rigour and systematicity.

571 572

573

4.3 Limitations and Prospects

Concerning manipulation in particular, an immediately evident limitation of this practical 574 demonstration is a lack of exploration of how the difficulty of training such ensembles of learners behaves at scale, with respect to more complex and multimodal datasets, and 576 with a wider range of architectures. Concerning scale, while none of the products of our 577 experimentation are state-of-the-art with respect to specific applications, we believe the 578 variety and promise of these results serve as a compelling validation of our theoretical 579 framework's utility and potential. Concerning applications of manipulator beyond the image domain, we report on some sketch experiments in sentiment-manipulation on text in Appendix F, where we also comment on the nature of technical difficulties to be overcome 582 in the application of our framework to complex domain data, and offer an explanation for 583 mode collapses observed during training by empirically relating manipulator to other 584 generative classification approaches. Concerning a wider range of architectures, we report on some specialisations of the learners. However, we leave exploring manipulation in combination with state-of-the-art architectures, such as diffusion, for future work.

587

Addressing the theoretical framework of tasks more broadly, our reliance on equational characterisations is double-edged. On one hand, it is uncommon to find such characterisations of mathematical systems of interest as they are usually defined by more direct means, and this presents a theoretical limitation. On the other hand, it appears that the strength of equational characterisations when applied to ML lies in imposing structure on "the way to learn to solve a problem" rather than on the solutions or problems themselves (Sutton, 2019). This suggests promising future possibilities of our mathematical framework in bridging structural-symbolic approaches from computer science more broadly with methods that can
 effectively leverage computation; we suggest "neural data structures" as a sketch experiment
 in Appendix E.

While we have demonstrated that, in certain cases, tasks can determine behaviours, there is a theoretical gap in the converse analysis of behaviours of trained models in terms of their basic tasks. The problem is typified by generative models where it is impossible a priori for all of the constituent tasks to be simultaneously perfectly optimised. The prototypical illustrating example is the GAN, where it is impossible for both the generator and discriminator to be perfect. Conceptually, the gap is that we have only dealt here with "static" ensembles which in principle admit idealised loss-minimisations, whereas some generative models use adversarial tasks to enforce "dynamic" training forces for a variety of purposes, such as representation regularisation for easier sampling in VAEs. While we do offer some initial methods of analysis in Section C.2, a more thorough and encompassing analysis is beyond the scope of this paper.

Future theoretical developments will seek to incorporate other aspects of ML: for example, relating to work that focuses on the choice of model architecture (Khatri et al., 2024) and interactions with the underlying data distribution (Bronstein et al., 2021). While our current experiments focus on demonstrating our framework's validity, future practical developments will explore applications to more complex, real-world ML challenges, where we envision our approach informing areas such as AutoML, interpretable AI, and formal verification of ML systems: the compositional nature of our task-based framework naturally aligns with neural architecture search by potentially informing principled search strategies for optimal model architectures, and the explicit representation of model behaviours as equational constraints could enhance interpretability and facilitate formal verification.

648 ETHICS STATEMENT 649

We recognise the potential impact of ML, both positive and negative, on global *society* and
 nature as ethical stakeholders.

Regarding society, we recognise that expertise in, and steering of, the development and implementation of ML systems is concentrated in the hands of relatively few. Hence, we believe that broad democratisation of methods and understanding would be an improvement upon current circumstances. Accordingly, in the development of our theory and exposition, we were biased away from formal rigidity in favour of accessibility and expressive possibility. Our "pattern language" approach seeks to be a minimally constraining mode of communication, reasoning, and design.

Regarding nature, we recognise that practical ML carries a growing ecological footprint,
 primarily due to the energy-intensive nature of training large models and running extensive
 experiments. We believe that a compositional approach could enable efficient and smaller
 model design, reducing environmental impact.

By formalising ML tasks in a way that allows for more principled and targeted experimentation, we aim to encourage a more thoughtful, resource-conscious, and participatory approach to ML research and development. We acknowledge that there is still much work to be done in this regard, and we hope that our framework will inspire further research into sustainable and harm-free ML practices.

669

Reproducibility Statement

670 671 672

Experiment methodology and details such as model parameters are reported in the Appendix. All code is available in the supplementary materials, and will be made public.

673 674 675

676

682

684 685

686

687

688

689

690

691 692

693

694

References

- Faris Abou-Saleh, James Cheney, Jeremy Gibbons, James McKinna, and Perdita Stevens. *Introduction to Bidirectional Transformations*, pp. 1–28. Springer International Publishing, Cham, 2018. ISBN 978-3-319-79108-1. doi: 10.1007/978-3-319-79108-1_1. URL https: //doi.org/10.1007/978-3-319-79108-1_1.
 - Christopher Alexander, Sara Ishikawa, and Murray Silverstein. *A Pattern Language: Towns, buildings, construction,* volume 2 of *Center for Environmental Structure Series.* Oxford University Press, New York, 1977.
 - Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. Stochastic Variational Video Prediction, March 2018. URL https://doi.org/10. 48550/arXiv.1710.11252.
 - John C. Baez and Jade Master. Open Petri Nets. *Mathematical Structures in Computer Science*, 30(3):314–341, March 2020. ISSN 0960-1295, 1469-8072. doi: 10.1017/S0960129520000043. URL http://arxiv.org/abs/1808.05415.
 - Francois Bancilhon and Nicolas Spyratos. Update semantics of relational views. *ACM Trans. Database Syst.*, 6(4):557–575, dec 1981. ISSN 0362-5915. doi: 10.1145/319628.319634. URL https://doi.org/10.1145/319628.319634.
- Kent Beck and Ward Cunningham. Using Pattern Languages for Object-Oriented Programs.
 Technical Report CR-87-43, September 1987.
- Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips. cc/paper_files/paper/2015/file/e995f98d56967d946471af29d7bf99f1-Paper.pdf.

702 703 704	Guillaume Boisseau and Paweł Sobociński. String Diagrammatic Electrical Circuit Theory. <i>Electronic Proceedings in Theoretical Computer Science</i> , 372:178–191, November 2022. ISSN 2075-2180. doi: 10.4204/EPTCS.372.13. URL http://arxiv.org/abs/2106.07763.
705 706 707 708	Filippo Bonchi, Pawel Sobociński, and Fabio Zanasi. A Categorical Semantics of Signal Flow Graphs. In CONCUR 2014 - Concurrency Theory - 25th International Conference, Rome, Italy, September 2014. URL https://hal.science/hal-02134182.
709 710 711	Filippo Bonchi, Pawel Sobocinski, and Fabio Zanasi. Interacting Hopf Algebras. Journal of Pure and Applied Algebra, 221(1):144–184, January 2017. ISSN 00224049. doi: 10.1016/j.jpaa. 2016.06.002. URL http://arxiv.org/abs/1403.7048.
712 713 714 715 716	Filippo Bonchi, Robin Piedeleu, Pawel Sobociński, and Fabio Zanasi. Graphical Affine Algebra. In 2019 34th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS), pp. 1–12, June 2019. doi: 10.1109/LICS.2019.8785877. URL https://doi.org/10.1109/ LICS.2019.8785877.
717 718 719 720 721	Sam Bond-Taylor, Adam Leach, Yang Long, and Chris G. Willcocks. Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 44(11):7327–7347, November 2022. ISSN 0162-8828, 2160-9292, 1939-3539. doi: 10.1109/ TPAMI.2021.3116668. URL https://doi.org/10.48550/arXiv.2103.04922.
722 723 724 725 726	Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. Generating Sentences from a Continuous Space. In <i>Proceedings of The 20th</i> <i>SIGNLL Conference on Computational Natural Language Learning</i> , pp. 10–21, Berlin, Germany, 2016. Association for Computational Linguistics. doi: 10.18653/v1/K16-1002. URL https://doi.org/10.48550/arXiv.1511.06349.
727 728 729	Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges, May 2021. URL https: //doi.org/10.48550/arXiv.2104.13478.
730 731 732 733 734	Kenta Cho and Bart Jacobs. Disintegration and Bayesian Inversion via String Diagrams. <i>Mathematical Structures in Computer Science</i> , 29(7):938–971, August 2019. ISSN 0960-1295, 1469-8072. doi: 10.1017/S0960129518000488. URL https://doi.org/10.48550/arXiv. 1709.00322.
735 736 737	Casey Chu, Andrey Zhmoginov, and Mark Sandler. Cyclegan, a master of steganography. <i>arXiv preprint arXiv:1712.02950</i> , 2017. URL https://doi.org/10.48550/arXiv.1712. 02950.
738 739 740 741	Lorenzo Ciampiconi, Adam Elwood, Marco Leonardi, Ashraf Mohamed, and Alessandro Rozza. A survey and taxonomy of loss functions in machine learning, January 2023. URL https://doi.org/10.48550/arXiv.2301.05579.
742 743 744	Bob Coecke and Ross Duncan. Interacting quantum observables: categorical algebra and diagrammatics. <i>New Journal of Physics</i> , 13(4):043016, April 2011. doi: 10.1088/1367-2630/13/4/043016. URL https://dx.doi.org/10.1088/1367-2630/13/4/043016.
745 746 747 748	Bob Coecke and Aleks Kissinger. <i>Picturing Quantum Processes: A First Course in Quantum Theory and Diagrammatic Reasoning</i> . Cambridge University Press, Cambridge, 2017. ISBN 978-1-107-10422-8. doi: 10.1017/9781316219317. URL https://www.cambridge.org/core/books/picturing-quantum-processes/1119568B3101F3A685BE832FEEC53E52.
749 750 751 752	Bob Coecke, Mehrnoosh Sadrzadeh, and Stephen Clark. Mathematical Foundations for a Compositional Distributional Model of Meaning, March 2010. URL http://arxiv.org/abs/1003.4394.
753 754 755	Geoffrey SH Cruttwell, Bruno Gavranović, Neil Ghani, Paul Wilson, and Fabio Zanasi. Categorical foundations of gradient-based learning. In <i>Programming Languages and</i> <i>Systems: 31st European Symposium on Programming, ESOP 2022, Held as Part of the European</i> <i>Joint Conferences on Theory and Practice of Software, ETAPS 2022, Munich, Germany, April</i>

2-7, 2022, Proceedings, pp. 1–28. Springer International Publishing Cham, 2022. URL
 https://doi.org/10.48550/arXiv.2103.01931.

- Artur d'Avila Garcez and Luis C. Lamb. Neurosymbolic AI: The 3rd Wave. https://arxiv.org/abs/2012.05876v2, December 2020.
- Gintare Karolina Dziugaite, Zoubin Ghahramani, and Daniel M. Roy. A study of the effect of JPG compression on adversarial images, August 2016. URL https://doi.org/10.48550/arXiv.1608.00853.
- Brendan Fong and David I. Spivak. An Invitation to Applied Category Theory: Seven Sketches in Compositionality. Cambridge University Press, 1 edition, July 2019a. ISBN 978-1-108-66880-4 978-1-108-48229-5 978-1-108-71182-1. doi: 10.1017/9781108668804.
- Brendan Fong and David I. Spivak. Supplying bells and whistles in symmetric monoidal categories., August 2019b. URL https://doi.org/10.48550/arXiv.1908.02633.
- 771
 Thomas Fox. Coalgebras and cartesian categories. Communications in Algebra, 4(7):665–667,

 772
 January 1976a. ISSN 0092-7872. doi: 10.1080/00927877608822127.
- Thomas Fox. Coalgebras and cartesian categories. *Communications in Algebra*, 4(7):665–667,
 January 1976b. ISSN 0092-7872, 1532-4125. doi: 10.1080/00927877608822127. URL
 https://doi.org/10.1080/00927877608822127.
- Tobias Fritz. A synthetic approach to Markov kernels, conditional independence and theorems on sufficient statistics. *Advances in Mathematics*, 370:107239, August 2020. ISSN 00018708. doi: 10.1016/j.aim.2020.107239. URL https://doi.org/10.1016/j.aim.2020. 107239.
- Tobias Fritz, Tomáš Gonda, and Paolo Perrone. De Finetti's Theorem in Categorical
 Probability. *Journal of Stochastic Analysis*, 2(4), November 2021. ISSN 2689-6931. doi:
 10.31390/josa.2.4.06. URL http://arxiv.org/abs/2105.02639.
- Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, Kun Zhang, and
 Dacheng Tao. Geometry-consistent generative adversarial networks for one-sided unsupervised domain mapping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Bruno Gavranović, Paul Lessard, Andrew Dudzik, Tamara von Glehn, João G. M. Araújo,
 and Petar Veličković. Categorical deep learning: An algebraic theory of architectures,
 February 2024. URL https://doi.org/10.48550/arXiv.2402.15332.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. URL https://doi.org/10.1145/3422622.
- Nathan Haydon and Paweł Sobociński. Compositional Diagrammatic First-Order Logic. In Ahti-Veikko Pietarinen, Peter Chapman, Leonie Bosveld-de Smet, Valeria Giardino, James Corter, and Sven Linker (eds.), *Diagrammatic Representation and Inference*, Lecture Notes in Computer Science, pp. 402–418, Cham, 2020. Springer International Publishing. ISBN 978-3-030-54249-8. doi: 10.1007/978-3-030-54249-8_32. URL https://doi.org/10.1007/ 978-3-030-54249-8_32.
- Jules Hedges. String diagrams for game theory, March 2015. URL http://arxiv.org/abs/ 1503.06072.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*, 2022. URL https://doi.org/10.48550/arXiv. 2208.01626.
- 809 Martin Hofmann, Benjamin Pierce, and Daniel Wagner. Symmetric lenses. *ACM SIGPLAN Notices*, 46(1):371–384, 2011.

823

835

836

- 810 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks 811 are universal approximators. *Neural networks*, 2(5):359–366, 1989. 812
- Bart Jacobs, Aleks Kissinger, and Fabio Zanasi. Causal inference by string diagram surgery. 813 In Foundations of Software Science and Computation Structures: 22nd International Conference, 814 FOSSACS 2019, Held as Part of the European Joint Conferences on Theory and Practice of Software, 815 ETAPS 2019, Prague, Czech Republic, April 6–11, 2019, Proceedings 22, pp. 313–329. Springer, 816 2019. URL https://doi.org/10.1007/978-3-030-17127-8_18. 817
- 818 Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, 819 and Michal Irani. Imagic: Text-based real image editing with diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 820 pp. 6007–6017, June 2023. 821
- 822 Nikhil Khatri, Tuomas Laakkonen, Jonathon Liu, and Vincent Wang-Maścianica. On the anatomy of attention, 2024. 824
- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes, December 2022. 825 URL https://doi.org/10.48550/arXiv.1312.6114.
- 827 Jeremias Knoblauch, Jack Jewson, and Theodoros Damoulas. An Optimization-centric View 828 on Bayes' Rule: Reviewing and Generalizing Variational Inference. Journal of Machine 829 *Learning Research*, 23:1–109, 2022.
- 830 Juncen Li, Robin Jia, He He, and Percy Liang. Delete, retrieve, generate: a simple approach 831 to sentiment and style transfer. In Proceedings of the 2018 Conference of the North American 832 Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 833 1 (Long Papers), pp. 1865–1874, 2018. URL https://doi.org/10.18653/v1/N18-1169. 834
 - 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), December 2015.
- 837 Robin Lorenz and Sean Tull. Causal models in string diagrams, April 2023. URL https: 838 //doi.org/10.48550/arXiv.2304.07638. 839
- Martin Markl, Steve Shnider, and Jim Stasheff. Operads in Algebra, Topology and Physics. 840 American Mathematical Society, Providence, RI, July 2007. ISBN 978-0-8218-4362-8. 841
- 842 Samuel G. Müller and Frank Hutter. Trivialaugment: Tuning-free yet state-of-the-art data 843 augmentation, 2021. 844
- Keisuke Nakano. A Tangled Web of 12 Lens Laws. In Reversible Computation: 13th 845 International Conference, RC 2021, Virtual Event, July 7–8, 2021, Proceedings, pp. 185– 846 203, Berlin, Heidelberg, July 2021. Springer-Verlag. ISBN 978-3-030-79836-9. doi: 847 10.1007/978-3-030-79837-6_11. 848
- 849 Andrew Ng and Michael Jordan. On Discriminative vs. Generative Classifiers: A comparison 850 of logistic regression and naive Bayes. In Advances in Neural Information Processing Systems, 851 volume 14. MIT Press, 2001.
- nLab authors. Bayesian inversion. https://ncatlab.org/nlab/show/Bayesian+ 853 inversion, June 2024a. Revision 9. 854
- 855 nLab authors. Markov category. https://ncatlab.org/nlab/show/Markov+category, July 856 2024b. Revision 56.
- 857 Prakash Panangaden. The Category of Markov Kernels. Electr. Notes Theor. Comput. Sci., 22: 858 171–187, December 1999. doi: 10.1016/S1571-0661(05)80602-4. URL https://doi.org/ 859 10.1016/S1571-0661(05)80602-4.
- Dusko Pavlovic. Monoidal computer i: Basic computability by string diagrams. Infor-861 mation and Computation, 226:94-116, 2013. ISSN 0890-5401. doi: https://doi.org/10. 862 1016/j.ic.2013.03.007. URL https://www.sciencedirect.com/science/article/pii/ 863 S0890540113000254. Special Issue: Information Security as a Resource.

864 Dusko Pavlovic. Programs as Diagrams: From Categorical Computability to Computable Categories. 865 Springer, 1st ed. 2023 edition edition, September 2023. ISBN 978-3-031-34826-6. URL 866 https://doi.org/10.48550/arXiv.2208.03817. 867 Boldizsár Poór, Quanlong Wang, Razin A. Shaikh, Lia Yeh, Richie Yeung, and Bob Coecke. 868 Completeness for arbitrary finite dimensions of ZXW-calculus, a unifying calculus, April 869 2023. URL http://arxiv.org/abs/2302.12135. 870 871 Alec Radford and Karthik Narasimhan. Improving language understanding by generative 872 pre-training. 2018. URL https://api.semanticscholar.org/CorpusID:49313245. 873 Marc'Aurelio Ranzato, Y-Lan Boureau, Sumit Chopra, and Yann LeCun. A unified energy-874 based framework for unsupervised learning. In Marina Meila and Xiaotong Shen (eds.), 875 *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics,* 876 volume 2 of Proceedings of Machine Learning Research, pp. 371–379, San Juan, Puerto Rico, 877 21-24 Mar 2007. PMLR. URL https://proceedings.mlr.press/v2/ranzato07a.html. 878 Oliver E. Richardson. Loss as the inconsistency of a probabilistic dependency graph: Choose 879 your model, not your loss function. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), Proceedings of The 25th International Conference on Artificial Intelligence 881 and Statistics, volume 151 of Proceedings of Machine Learning Research, pp. 2706–2735. PMLR, 882 28-30 Mar 2022. URL https://proceedings.mlr.press/v151/richardson22b.html. 883 884 Joseph Understanding Variational Autoencoders Rocca. (VAEs). 885 https://towardsdatascience.com/understanding-variational-autoencoders-vaes-886 f70510919f73, March 2021. Peter Selinger. A Survey of Graphical Languages for Monoidal Categories. In Bob Coecke 888 (ed.), New Structures for Physics, pp. 289–355. Springer Berlin Heidelberg, Berlin, Heidelberg, 889 2011. ISBN 978-3-642-12821-9. doi: 10.1007/978-3-642-12821-9_4. URL https://doi. 890 org/10.1007/978-3-642-12821-9_4. 891 Razin A. Shaikh, Sara Sabrina Zemljic, Sean Tull, and Stephen Clark. The Conceptual VAE, 892 March 2022. URL https://doi.org/10.48550/arXiv.2203.11216. 893 894 Lennert De Smet, Pedro Zuidberg Dos Martires, Robin Manhaeve, Giuseppe Marra, Angelika 895 Kimmig, and Luc De Readt. Neural probabilistic logic programming in discrete-continuous 896 domains. In Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence, 897 pp. 529-538. PMLR, July 2023. 898 Paweł Sobociński. Graphical Linear Algebra, June 2015. URL https:// 899 graphicallinearalgebra.net/. 900 901 Akhilesh Sudhakar, Bhargav Upadhyay, and Arjun Maheswaran. "Transforming" delete, 902 retrieve, generate approach for controlled text style transfer. In Proceedings of the 2019 903 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint 904 Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 3269–3279, 2019. URL 905 https://doi.org/10.48550/arXiv.1908.09368. 906 Richard Sutton. The Bitter Lesson, 2019. 907 908 Juan Terven, Diana M. Cordova-Esparza, Alfonso Ramirez-Pedraza, and Edgar A. Chavez-909 Urbiola. Loss Functions and Metrics in Deep Learning, September 2023. URL https: 910 //doi.org/10.48550/arXiv.2307.0269. 911 Vladimir N. Vapnik. Statistical Learning Theory. Wiley-Interscience, New York, 1st edition 912 edition, September 1998. ISBN 978-0-471-03003-4. 913 914 V. Wang-Maścianica. String Diagrams for Text. http://purl.org/dc/dcmitype/Text, University 915 of Oxford, 2023. 916

917 Vincent Wang-Mascianica, Jonathon Liu, and Bob Coecke. Distilling Text into Circuits, January 2023. URL http://arxiv.org/abs/2301.10595.

918 919 920	Nicholas Watters, Loic Matthey, Sebastian Borgeaud, Rishabh Kabra, and Alexander Lerchner. Spriteworld: A flexible, configurable reinforcement learning environment, 2019. URL https://github.com/deepmind/spriteworld/.
921 922 923 924	Ronald J. Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. <i>Neural Computation</i> , 1(2):270–280, 1989. doi: 10.1162/neco. 1989.1.2.270. URL https://doi.org/10.1162/neco.1989.1.2.270.
925 926 927	Sidi Wu, Yizi Chen, Samuel Mermet, Lorenz Hurni, Konrad Schindler, Nicolas Gonthier, and Loic Landrieu. Stegogan: Leveraging steganography for non-bijective image-to-image translation, March 2024. URL https://doi.org/10.48550/arXiv.2403.20142.
928 929 930	Donald Yau. Higher dimensional algebras via colored PROPs, September 2008. URL https://doi.org/10.48550/arXiv.0809.2161.
931 932 933	Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks, 2017.
934	
935	
930	
937	
930	
940	
941	
942	
943	
944	
945	
946	
947	
948	
949	
950	
951	
952	
953	
954	
956	
957	
958	
959	
960	
961	
962	
963	
964	
965	
966	
967	
968	
909	
970	
3/1	

972 STRING DIAGRAMS FOR TASKS А 973

974

975

976

977

String diagrams are a formal diagrammatic syntax that take semantics in symmetric monoidal categories, and they find usage in a broad variety of fields². Our string diagrams are built using sequential and parallel composition from the following generators, along with a stock of function labels:

978					\mathbb{R}^{B}
979	$\mathbb{R}^M - f - \mathbb{R}^N$	{*}	\mathbb{R}^M — \bullet \cdots { \star }	$\cdots x \cdots = \mathbb{R}^M$	$\mathbb{R}^A - f - g - \mathbb{R}^C$
980	Function	Singleton space	Delete	Data from $\mathcal X$	Sequential composition
981	$f: \mathbb{R}^M \to \mathbb{R}^N$	$\mathbb{R}^0 \simeq \{\star\}$	$\epsilon: \mathbb{R}^M \to \{\star\}$	$\mathcal{X}: \{\star\} \to \mathbb{R}^M$	$(g \circ f): \mathbb{R}^A \to \mathbb{R}^C$
982	J · · · · · · ·	12 <u>(</u> /)	61 Ha (A)	561 (A) - 54	(8-)), 14
983		D N D M	D M	•	
984	$\{\star\}\cdots \oslash \vdash \mathbb{R}^N$		\mathbb{R}^{M} —	$\cdots x \cdots = \mathbb{R}^m$	$\mathbb{R}^{A} - f = \mathbb{R}^{B}$
985	7	$\mathbb{R}^M \longrightarrow \mathbb{R}^N$	$\subseteq \mathbb{R}^M$	$\cdots \mathcal{Y} \cdots \blacksquare \mathbb{R}^N$	$\mathbb{R}^{C} \longrightarrow \mathbb{R}^{D}$
000	Vector	Swap	Сору	Independent \mathcal{X}, \mathcal{Y}	Parallel composition
900		(N+M) = (M+N)		$(\mathcal{X} \times \mathcal{U}) \cdot (\mathcal{U}) \to \mathbb{D}(M+N)$	$(f \oplus a) : \mathbb{D}(A+C) : \mathbb{D}(B+D)$
987	$v \in \mathbb{R}^{n}$	$A: \mathbb{K}_{(1,1,1)} \to \mathbb{K}_{(1,1,1)}$	$\Delta \colon \mathbb{R}^m \to \mathbb{R}^{(m+m)}$	$(\mathcal{I} \times \mathcal{J}) \cdot \{\star\} \to \mathbb{R}^{(m+1)}$	$(\uparrow \oplus \&) \cdot \bowtie (\to) \to \mathbb{K}_{(-,+)}$

988 For conventional reasons that were not by our choice, vectors are depicted as triangular 989 nodes with only output wires, reminiscent of bra-ket notation. (co)associative (co)monoids, 990 such as copy-delete and add-zero, are specially depicted as circular nodes as is common in 991 applied category theory. In this work, encoders and decoders are sometimes depicted as 992 "bottlenecking" trapezia, as is common in ML, and distributional states are given their own 993 notation as thick bars.

994 An attractive characteristic of string diagrams is that visually intuitive equivalences between 995 information flows are guaranteed to correspond to symbolic derivations of behavioural equiv-996 alence: tedious algebraic proofs of equality between sequentially- and parallel-composite 997 processes are suppressed and absorbed by (processive) isotopies of diagrams. In the 998 diagrammatic syntax it is conventional to notate such isomorphisms as plain equalities. 999 Interested readers are referred to (Selinger, 2011) for the relevant mathematical foundations.



1006 A.1 CATEGORICAL SEMANTICS OF TASK DIAGRAMS

The functional effect of the construction below is to extend the category of continuous maps 1008 between Euclidean spaces with global elements that behave as probability distributions 1009 instead of points. We presume familiarity with symmetric monoidal categories and their 1010 graphical calculi (Selinger, 2011). 1011

Let CartSp denote the coloured PROP (Yau, 2008) of continuous maps between Euclidean 1012 spaces, where the tensor product is the cartesian product — i.e. **CartSp** is cartesian monoidal. 1013

1014 Let BorelStoch denote the Markov category (Cho & Jacobs, 2019; Fritz, 2020) of stochastic 1015 kernels (Panangaden, 1999) between Borel-measurable spaces. Stochastic kernels in particular 1016 subsume the continuous maps between Euclidean spaces.

1017 As a Markov category, in the terminology of (Fong & Spivak, 2019b), BorelStoch supplies 1018 cocommutative comonoids. By Fox's theorem (Fox, 1976b) cartesian monoidal categories 1019

²Including linear and affine algebra (Sobociński, 2015; Bonchi et al., 2017; 2019), first order logic 1021 (Haydon & Sobociński, 2020), causal models (Lorenz & Tull, 2023; Jacobs et al., 2019), signal flow graphs (Bonchi et al., 2014), electrical circuits (Boisseau & Sobociński, 2022), game theory (Hedges, 2015), 1023 petri nets (Baez & Master, 2020), probability theory (Fritz et al., 2021), formal linguistics (Coecke et al., 2010; Wang-Mascianica et al., 2023; Wang-Maścianica, 2023), quantum theory (Coecke & Duncan, 2011; Coecke & Kissinger, 2017; Poór et al., 2023), and aspects of machine learning such as backpropagation 1025 (Cruttwell et al., 2022).

are precisely those isomorphic to their own categories of cocommutative comonoids. Hence there is a (semicartesian) functorial embedding of **CartSp** into **BorelStoch** sending \mathbb{R}^N to \mathbb{R}^N (equipped with the usual Borel measure), and continuous maps to deterministic continuous maps. **We declare our semantics** to be taken in the category to be generated by the image of this embedding along with the probability distributions $\mathcal{X} : \{\star\} \to \mathbb{R}^N$, where $\{\star\}$ is the singleton monoidal unit of **BorelStoch**.

1033 A.2 CATEG

A.2 CATEGORICAL SEMANTICS FOR IDEALISED UNIVERSAL APPROXIMATORS

As we are concerned with behaviour, not implementation details, we idealise all neural networks as perfect universal approximators, which we may formulate string-diagrammatically in a monoidal closed category, borrowing evaluator-notation from (Pavlovic, 2013; 2023). In essence, we are assuming that architectures are sufficiently expressive to optimise whatever tasks we give them; in practice, the conditions under which architectures become universal approximators can be mild (Hornik et al., 1989), and the idealisation is increasingly true-in-practice in the contemporary context of increasing data and compute.

Definition A.1 (Learner). Let *X*, *Y* denote input and output types. A process $\Omega : \mathfrak{p} \oplus X \to Y$ with parameters in para(Ω) = $\mathfrak{p} = \mathbb{R}^n$ (for sufficiently large *n*) is a *universal approximator* or *learner* when³:

1034

1045

 $\forall f_{X \to Y} \exists \hat{f}_{\in \mathfrak{p}} : \underline{X f Y} = \underbrace{ \begin{array}{c} & & \\ & \\ & & & \\ & &$

1047 The parameter space could represent e.g. the phase space of weights and biases of a neural network.

Example A.2. By visual convention, we use colours to indicate different data types of wires. We depict processes with no free parameters as white boxes, and learners as black-boxes with variable labels to indicate distinct or shared parameters. The following composite process has one function f with no learnable parameters, and three neural nets: the two labelled α share a parameter in the space \mathfrak{p} , and the one labelled β takes a parameter in \mathbb{Q} . In this paper, we favour the shorthand on the right.

1055

1056

1058

1067



The universal approximation theorem, suitably idealised, manifests as the capacity for a black-box learner to be diagrammatically substituted for any other composite diagram with equal input and output, including those composites that contain other learners. For example, recall the linear put below, which may be viewed as substituting a particular composite of put', enc, dec, and addition in place of put:



This ability is referred to in this work intermittently as *specialisation*, and as *expressive reduction* in (Khatri et al., 2024) where the concept first appeared. For the sake of completeness, we reproduce the relevant construction that gives category-theoretic semantics to universal approximators and specialisation below, along with standard definitions, with the authors' permission.

Definition A.3 (PROP). A PROP is a strict symmetric monoidal category generated by a single object *x*: every object is of the form

$$\begin{array}{c}
n \\
x = x \underbrace{\otimes \cdots \otimes x}_{n} \\
\end{array}$$

³We adapt the shape of the universal approximators to clearly indicate the parameters.

1080 PROPs may be generated by, and presented as *signatures* (Σ , E) consisting of generating 1081 morphisms Σ with arity and coarity in \mathbb{N} , and equations *E* relating symmetric monoidal 1082 composites of generators.

1083 **Definition A.4** (Coloured PROP). A *multi-sorted* or *coloured* PROP with set of colours **C** has 1084 a monoid of objects generated by **C**.

1085 **Definition A.5** (Cartesian PROP). By Fox's theorem (Fox, 1976a), a cartesian PROP is one in 1086 which every object (wire) is equipped with a cocommutative comonoid (copy) with counit 1087 (delete) such that all morphisms in the category are comonoid cohomomorphisms. 1088

Definition A.6 ((symmetric, unital) coloured operad). Where $(\mathcal{V}, \boxtimes, J)$ is a symmetric 1089 monoidal category and \mathfrak{C} denotes a set of *colours* c_i , a coloured operad \mathcal{O} consists of: 1090

- For each $n \in \mathbb{N}$ and each (n + 1)-tuple $(c_1, \dots, c_n; c)$, an object $\mathcal{O}(c_1, \dots, c_n; n) \in \mathcal{V}$
- For each $c \in \mathfrak{G}$, a morphism $1_c : J \to \mathcal{O}(c; c)$ called the *identity of c*
- For each (n + 1)-tuple $(c_1 \cdots c_n; c)$ and *n* other tuples $(d_1^1 \cdots d_{k_1}^1) \cdots (d_1^n \cdots d_{k_n}^n)$ a composition morphism

 $\mathcal{O}(c_1, \cdots, c_n; c) \boxtimes \mathcal{O}(d_1^1 \cdots d_{k_1}^1) \boxtimes \cdots \boxtimes \mathcal{O}(d_1^n \cdots d_{k_n}^n) \to \mathcal{O}(d_1^1 \cdots d_{k_1}^1 \cdots d_1^n \cdots d_{k_n}^n; c)$

• for all $n \in \mathbb{N}$, all tuples of colours, and each permutation $\sigma \in S_n$ the symmetric group on *n*, a morphism:

 $\sigma^*: \mathcal{O}(c_1 \cdots c_n; c) \to \mathcal{O}(c_{\sigma^*(1)} \cdots c_{\sigma^*(n)}; c)$

1104 The σ^* must represent S_n , and composition must satisfy associativity and unitality in a 1105 S_n -invariant manner.

1106 **Construction A.7** (Hom-Operad of coloured PROP). Where $(\mathcal{P}, \otimes, I)$ is a coloured PROP 1107 with colours $\mathfrak{G}_{\mathcal{P}}$, we construct $\mathcal{O}_{\mathcal{P}}$, the *hom-operad* of \mathcal{P} . We do so in two stages, by first 1108 defining an ambient operad, and then restricting to the operad obtained by a collection of 1109 generators. Let the ambient symmetric monoidal category be (**Set**, \times , { \star }). Let the colours $\mathfrak{C}_{\mathcal{O}}$ 1110 be the set of all tuples (**A**, **B**), each denoting a pair of tuples $(A_1 \otimes A_n, B_1 \otimes B_n)$ of $A_i, B_i \in \mathfrak{C}_{\mathcal{P}}$.

- The tuple $((\mathbf{A}^1, \mathbf{B}^1) \cdots (\mathbf{A}^n, \mathbf{B}^n); (\mathbf{A}, \mathbf{B}))$ is assigned the set $[\mathcal{P}(\mathbf{A}^1, \mathbf{B}^1) \times \cdots \times$ $\mathcal{P}(\mathbf{A}^n, \mathbf{B}^n) \to \mathcal{P}(\mathbf{A}, \mathbf{B}) \in \mathbf{Set}$; the set of all *generated* functions from the product of homsets $\mathcal{P}(\mathbf{A}^i, \mathbf{B}^i)$ to the homset $\mathcal{P}(\mathbf{A}, \mathbf{B})$.
- $1_{(\mathbf{A},\mathbf{B})}$: $\{\star\} \to [\mathcal{P}(\mathbf{A},\mathbf{B}) \to \mathcal{P}(\mathbf{A},\mathbf{B})]$ is the identity functional that maps $f : \mathbf{A} \to \mathbf{B}$ in $\mathcal{P}(\mathbf{A}, \mathbf{B})$ to itself.
- The composition operations correspond to function composition in **Set**, where $[X \to Y] \times [Y \to Z] \to [X \to Z]$ sends $(f_{:X \to Y}, g_{:Y \to Z}) \mapsto (g \circ f)_{:X \to Z}$; appropriately generalised to the multi-argument case. The permutations are similarly defined, inheriting their coherence conditions from the commutativity isomorphisms of the categorical product ×.

1123 The generators are: 1124

1093

1095

1097

1099

1100

1101 1102

1103

1111

1112 1113

1114

1115

1116

1117 1118

1119

1120

1121

1122

1127

1130

1131

- 1125 • For every $f \in \mathcal{P}(\mathbf{A}, \mathbf{B})$ that is a generator of \mathcal{P} , define a corresponding generator of 1126 type $\{\star\} \to [\mathcal{P}(I,I) \to \mathcal{P}(\mathbf{A},\mathbf{B})]$, which is the functional $(-\mapsto (f \otimes -))$ that sends endomorphisms of the monoidal unit of \mathcal{P} to their tensor with f, viewed as an 1128 element of the set $[\mathcal{P}(I, I) \rightarrow \mathcal{P}(\mathbf{A}, \mathbf{B})]$. 1129
 - For every pair of tuples $((\mathbf{X}^1, \mathbf{Y}^1) \cdots (\mathbf{X}^m, \mathbf{Y}^m); (\mathbf{A}, \mathbf{B}))$ and $((\mathbf{J}^1, \mathbf{K}^1) \cdots (\mathbf{J}^n, \mathbf{K}^n); (\mathbf{B}, \mathbf{C}))$ in $\mathfrak{C}_{\mathcal{O}}$, a corresponding *sequential composition* operation of type:

$$[\prod_{i\leqslant m}\mathcal{P}(\mathbf{X}^{i},\mathbf{Y}^{i})\to\mathcal{P}(\mathbf{A},\mathbf{B})]\times[\prod_{j\leqslant n}\mathcal{P}(\mathbf{J}^{j},\mathbf{K}^{j})\to\mathcal{P}(\mathbf{B},\mathbf{C})]$$



1136 1137 1138

1139

1140 1141

1142

1143

1162

- Which maps pairs of functionals $(F_{:\prod_{i \leq m} \mathcal{P}(\mathbf{X}^{i}, \mathbf{Y}^{i}) \to \mathcal{P}(\mathbf{A}, \mathbf{B})}, G_{:\prod_{j \leq n} \mathcal{P}(\mathbf{J}^{j}, \mathbf{K}^{j}) \to \mathcal{P}(\mathbf{B}, \mathbf{C})})$ to the functional which sends $p^{i} : \mathbf{X}^{i} \to \mathbf{Y}^{i}$ and $q^{j} : \mathbf{X}^{j} \to \mathbf{Y}^{j}$ to $G(p_{1} \cdots p_{m}) \circ F(q_{1} \cdots q_{n})$.
- An analogous *parallel composition* for every pair of tuples, which sends pairs of functionals (F, G) to $G(p_1 \cdots p_m) \otimes F(q_1 \cdots q_n)$.

Remark A.8. For technical reasons involving scalars (the endomorphisms of the monoidal unit), this construction only works in semicartesian settings, i.e. where the monoidal unit is also terminal, but that is sufficiently general to admit our use cases, which are primarily in cartesian monoidal settings (Fox, 1976a) and semicartesian Markov categories for probabilistic settings (nLab authors, 2024b).

Example A.9. Construction A.7 can be thought of as bridging diagrams with their specific algebraic descriptions using just the basic constructors \circ , \otimes ; the hom-operad (when notated suggestively in the usual tree-notation, found e.g. in Markl et al. (2007)) plays the role of the syntactic tree of \circ , \otimes operators. For instance, given the composite morphism ($g \otimes 1_E$) \circ ($1_A \otimes f$) in PROP \mathcal{P} , the corresponding diagram and operad-state in $\mathcal{O}_{\mathcal{P}}$ is:



Since the PROPs **CartSp** and its free tensoring are cartesian, $\mathcal{P}(I, I)$ is a singleton containing 1163 only the identity of the monoidal unit, so in the settings we are concerned with, we 1164 may simplify colours of the form $[\mathcal{P}(I,I)^N \to \mathcal{P}(\mathbf{A},\mathbf{B})]$ to just $\mathcal{P}(\mathbf{A},\mathbf{B})$, and operad-states 1165 $\{\star\} \to \mathcal{P}(\mathbf{A}, \mathbf{B})$ are in bijective correpondence with morphisms $f : \mathbf{A} \to \mathbf{B}$ of \mathcal{P} ; the fact that 1166 all $f : \mathbf{A} \to \mathbf{B}$ are representable as operad states follows by construction, since any f in 1167 \mathcal{P} is by definition expressible in terms of the generators of \mathcal{P} , and sequential and parallel 1168 composition \circ , \otimes . As we assume homsets are already quotiented by the equational theory 1169 of \mathcal{P} and the symmetric monoidal coherences, our operadic representations inherit them: 1170 for example, we obtain interchange equalities such as the one below for free:



1180 **Definition A.10** (Universal approximators and specialisation). A morphism of a coloured 1181 PROP \mathcal{P} of type (**A**, **B**) containing universal approximators as black-boxes of types $\mathbf{A}^{i \leq n} \rightarrow$ 1182 $\mathbf{B}^{i \leq n}$ is a morphism (($\mathbf{A}^1, \mathbf{B}^1$) ... ($\mathbf{A}^n, \mathbf{B}^n$); (**A**, **B**)) of $\mathcal{O}_{\mathcal{P}}$, and by construction, vice versa. 1183 Specialisation corresponds to precomposition in $\mathcal{O}_{\mathcal{P}}$.





1188 composites containing other universal approximators.

Remark A.12. The extension of the current theory to accommodate parameter sharing between universal approximators is conceptually straightforward but technically involved. Parameter sharing corresponds to the ability to reuse – i.e. copy – data between open wires in the operad $\mathcal{O}_{\mathcal{P}}$, which amounts to having a cartesian operad.

- **B** Strong manipulation
- 1213 1214 1215

1206

1207

1208

1209

1210 1211 1212

The basic manipulation admits pathological counterexamples, which we may block by imposing additional tasks. We can treat these new tasks as additional regularisation terms. This section further illustrates a form of "legalistic" thinking using tasks: by thinking of ways that "noncooperative" or "naughty" learners might seek to satisfy tasks without exhibiting the behaviour that the modeller desires. By identifying these counterexamples and constructing additional tasks that block them, the modeller may iteratively improve the behaviour of the model. For illustration, consider the following examples of pathological behaviour that satisfy basic manipulation, again in the setting of editing the colour of a shape.

Example B.1 (Flipping). Consider a put that changes the colour of a shape as desired, but then vertically flips the shape. If the classifier get is insensitive to the position and orientation of the shape, then CLASSIFY, PUTGET, and GETPUT are satisfied. Moreover, since a vertical flip is its own inverse, composing two puts as in the UNDOABILITY task will not detect this aberration. Speaking in more general terms, if there are properties that get is insensitive to, there must be additional guardrails to ensure that put preserves these other properties as the identity, rather than one of potentially very many self-inverse symmetries.

1231

Example B.2 (Adversarial decorations). While the classifier get may be perfect in-distribution, there are no guarantees about its behaviour out of distribution, for example, when given images with multiple shapes, where it might only classify the leftmost shape. So, it is possible that put learns to make edits that take the resulting image out-of-distribution: for example, by adding a red circle next to a blue square to fool the classifier into outputting "red". This would satisfy CLASSIFY, PUTGET, and GETPUT. If the put can recognise and undo its own decorations, then UNDOABILITY will also be satisfied. Speaking more generally, we require additional guardrails to ensure that put returns something in-distribution.

19/0

1241 The strong manipulation adds additional tasks to the manipulator. A strong manipulator has to satisfy the original four tasks plus the following four:



The PutPut task (which is strictly stronger than UNDOABILITY in that it is algebraically implied) says that the effect of putting twice is the same as discarding the effect of the first edit and only keeping the last edit. In conjunction with PutGet and GetPut, this creates what is known in the literature as a *very well-behaved lens*, which blocks Example B.1 and similar modifications of the data get is insensitive to. The TRUE, FAKE, and FOOL tasks introduce a discriminator component dsc, which forms a GAN pattern with respect to put as the generator. When well-trained, this forces the outputs of put to lie in-distribution. As in general there are no algebraic or equational laws that characterise arbitrary distributions of data, using GANs in this way is a generic recipe for shaping outputs of generators to behave well in-distribution.

Remark B.4 (Why basic manipulation is preferable in practice). We have observed informally that conditions such as those in basic manipulation where learners are cooperative and 1261 there are only learners on the LHS appear to be more stable during training. We suggest 1262 a sketch reason why: in the tasks of strong manipulation, PutPut has put occur on both 1263 the LHS and RHS, which establishes a nontrivial dependence on the current position on 1264 parameter-space of put in the process of finding a solution. Similarly, the GAN rules of strong 1265 manipulation establishes adversarial mutual dependencies in the parameters of dsc and put. 1266 Conceptually, these dependencies create dynamical systems on the paths that the learners 1267 take over the course of training in parameter-space, which may for instance include stable 1268 orbits and chaotic behaviour, and may be highly sensitive to initial conditions. A further 1269 elaboration of "static" versus "dynamic" tasks in tandem with the ability to express equivalent tasks is potentially useful for creating train-stable models with equivalent behaviour, but 1270 this is beyond the scope of this paper, and left for future work. 1271

1272 1273

C Deferred proofs

1274 1275

1279 1280

1281 1282

1283

1284 1285 1286

In this section, we refer to tasks and their realisations interchangeably: so instead of "manipulators" realising the manipulation task, we speak just of a manipulator, disambiguating when necessary.

C.1 Proof of Proposition 2.9

Lemma C.1. For all well-typed *f* , *g* , and for any positive linear combination α : $\mathbb{R}^{\geq 0} \times \mathbb{R}^{\geq 0} \rightarrow \mathbb{R}^{\geq 0}$:

$$x \vdash \left(\begin{array}{c} f \\ g \end{array} \right) = \alpha \quad \Rightarrow \quad x \vdash \bullet \bullet - \quad \sim \quad x \vdash f \quad \Rightarrow \quad x \vdash \bullet \bullet + \quad x \vdash g \quad = \quad$$

1287 1288

Proof. For the forward refinement, as α is a positive linear combination, we have for all $x \in \mathcal{X}$ that $\alpha(f(x), g(x)) = \alpha_1 \cdot f(x) + \alpha_2 \cdot g(x) = 0$. If f and g are constant-functions 0, we are done. Otherwise, positivity implies that $\alpha_1 \cdot f(x) = \alpha_2 \cdot g(x) = 0$, and since $\alpha_1, \alpha_2 \in \mathbb{R}^{\geq 0}$, then f(x) = 0 = g(x), which is the desired task. For the backwards refinement, if f(x) = 0 = g(x)for all $x \in \mathcal{X}$, then $\alpha(f(x), g(x)) = \alpha_1 \cdot f(x) + \alpha_2 \cdot g(x) = 0$. **Lemma C.2.** For a real-valued pairwise measure $\mathbf{D} : \mathcal{D}(Y) \times \mathcal{D}(Y) \to \mathbb{R}^{\geq 0}$ on the space of distributions over *Y*, the positivity axiom $\mathbf{D}(\mathcal{Y}_1, \mathcal{Y}_2) = 0 \iff \mathcal{Y}_1 = \mathcal{Y}_2$ implies, for (almost⁴) all $f, g : X \to Y$ and \mathcal{X} :

$$x \vdash \underbrace{f}_{g} \vdash \mathbf{D}_{-} \Leftrightarrow x \vdash \bullet \diamondsuit_{-} \sim x \vdash f_{-} :\Leftrightarrow x \vdash g_{-}$$

Proof. For the forward refinement, we assume that $\mathbf{D}(f(\mathcal{X}, g(\mathcal{X}))) = 0$. By positivity of **D**, $f(\mathcal{X}) = g(\mathcal{X})$, which is the right hand task. For the backward refinement, if $f(\mathcal{X}) = g(\mathcal{X})$, then, by positivity, $\mathbf{D}(f(\mathcal{X}, g(\mathcal{X}))) = 0$.

Proposition C.3. If enc and dec are deterministic and \mathbf{E}_{enc} and \mathbf{E}_{dec} are positive (e.g. metrics or statistical divergences), then energy minimisation \equiv autoencoding.

Proof. As enc and dec are functions, we may copy them through their outputs to re-express energy minimisation as:

1313 1314 1315

1316

1318 1319

1320

1324

1325

1309

1300 1301 1302

 $y \vdash \underbrace{enc-dec}_{enc} = \underbrace{E_{dec}}_{E_{enc}} + \underbrace{xy}_{enc,dec} + \underbrace{enc,dec}_{enc}$

¹³¹⁷ By Lemma C.1, this is equivalent to:



Recall that \mathbf{E}_{enc} and \mathbf{E}_{dec} are positive by definition, hence Lemma C.2 allows us to express the two minimisations as:



The left task is autoencoding, so we have that energy minimisation refines autoencoding.
 For the other direction, we observe that autoencoding refines the right task by postcomposing both sides with enc, and as the left and right tasks together are equivalent to energy minimisation, we have the claim.

1330 1331

C.2 Proof of Proposition 3.2

CycleGANs (Example 2.6) solve a similar task as the manipulator, translating between two distributions. In fact, with the additional regularisation terms, the strong manipulator (see Appendix B) is a refinement of CycleGANs, giving us more guarantees by avoiding certain failure cases. We can use the pattern language to show this formally. However, in both cases, the tasks are not perfectible, as, by design, it is impossible for the generator and the discriminator to have a loss of zero at the same time. Therefore, we have to generalise the definition of *refinement* for *partially perfectible tasks*.

Definition C.4 (Partially perfectible task). A *partially perfectible task* $\{(f_i, g_i, \mathcal{X}_i, \mathfrak{p}_i)\}_i$ with learnable functions Σ_l is a compound task with excluded learners $E \subseteq \Sigma_l$ such that:

- 1. no two $e \in E$ share atomic tasks, i.e. $\forall e, e' \in E.\{(f_i, g_i, \mathcal{X}_i, \mathfrak{p}_i) | e \in f_i \lor e \in g_i\} \cap \{(f_i, g_i, \mathcal{X}_i, \mathfrak{p}_i) | e' \in f_i \lor e' \in g_i\} = \emptyset$ and
 - 2. all tasks not involving learners from *E* are perfectible, i.e. $\exists \pi \in \mathfrak{p} . \forall x \in X . f_{i;\pi}(x) = g_{i;\pi}(x)$

1342

1344

¹³⁴⁷

¹³⁴⁸ ⁴When the function space containing f and g is large, the edge case where f and g are nonconstant **1349** and f = -g is negligible. Since we are concerned with behaviour in the computational limit, this is an acceptable assumption for our purposes.

1350 Importantly, there does not need to exist a perfect solution for all tasks. As condition (1) 1351 explicitly forbids excluded learners to share tasks, we do not need to make restrictions 1352 on their composite behaviour. An example of a partially perfectible task is CycleGAN 1353 with excluded learners dsc₁, dsc₂. Assuming an appropriate data distribution (Wu et al., 1354 2024), the autoencoding tasks, i.e. the reconstruction losses, are perfectible. However, the generator-discriminator tasks are not. GAN with excluded learner dsc is another example of 1355 a partially perfectible tasks. In this case no tasks are required to be perfectible, as all of them 1356 involve the discriminator. 1357

When we have partially perfectible tasks, we have to define what we mean by an optimalsolution.

1360 Definition C.5 (Optimal partially pefect solution). Let $\Phi = \{(f_i, g_i, \mathcal{X}_i, \mathfrak{p}_i)\}_i$ be a partially **1361** perfectible tasks with learners Σ and excluded learners $E \subseteq \Sigma$. Then an *optimal partially* **1362** *perfect solution* for the learners $L = \Sigma \setminus E$ is, if it exists, a set of parameters $\pi_L \in \mathfrak{p}_L$ such that **1363** for all α_L , \mathbf{D}_{ϕ} they perfect the perfectible tasks:

$$\forall \{(f_i, g_i, \mathcal{X}_i, \mathfrak{p}_i) \in \Phi | \forall e \in E.e \notin f_i \land e \notin g_i\} \\ \forall x \in \mathcal{X}_i \\ f_{i;\pi_L}(x) = g_{i;\pi_L}(x)$$

and for all $l \in L$ and all tasks Φ_l with para $(l) \cap \mathfrak{p}_i \neq \emptyset$:

$$\alpha_{l}(\{\mathop{\mathbb{E}}_{x\sim\mathcal{X}_{i}}(\mathbf{D}_{\phi}(\operatorname{sys}_{\phi;(\pi_{L}\cup\mathfrak{p}_{E}(\pi_{L}))}(x),\operatorname{spec}_{\phi;\pi_{L}\cup\mathfrak{p}_{E}(\pi_{L})}(x))) \mid \phi \in \Phi_{l}\})$$

is minimal over all possible parameter combinations, where

$$\mathfrak{p}_{E}(\pi_{L}) = \bigcup_{e \in E} \inf_{\pi_{e} \in \mathfrak{p}_{e}} (\alpha_{e}(\{ \mathop{\mathbb{E}}_{x \sim \mathcal{X}_{i}} (\mathbf{D}_{\psi}(\operatorname{sys}_{\psi;(\pi_{L} \cup \pi_{e})}(x), \operatorname{spec}_{\psi;(\pi_{L} \cup \pi_{e})}(x))) \mid \psi \in \Phi_{e}\}))$$

In words, we consider a set of parameters *optimal partially perfect*, if they are perfect for the perfectible tasks and have a minimal objective function for the non-perfectible tasks, even if the excluded learners only optimise for themselves. This may not always exist, either because there is no solution that is optimal for all α_L , \mathbf{D}_{ϕ} or as there may be more optimal solutions for the non-perfectible tasks that do not satisfy the perfectible tasks.

Given this, we can generalise the definition of *refinement* for partially perfectible tasks.

Definition C.6 (Optimal refinement of partially-perfectible tasks). Given two partiallyperfectible tasks Ψ , Φ with excluded learners E_{Ψ} , E_{Φ} , we say that Ψ optimally refines Φ if optimal partially perfect solutions for $\Sigma_{l_{\Psi}} \setminus E_{\Psi}$ can be composed to form optimal partially perfect solutions for $\Sigma_{l_{\Phi}} \setminus A_{\Phi}$.

¹³⁸⁵ For A_{Ψ} , $A_{\Phi} = \emptyset$, optimal refinement is equivalent to refinement.

To show that manipulation optimally refines cycleGAN, we need one assumption: we assume that the measure of statistical divergence and compound function have been chosen such that a generator in a generative-adversarial setting is optimal if and only if its output distribution is equal to the original distribution. As the goal of such a generative setting is to approximate the original distribution, this is quite a natural assumption. One possible choice of measure of statistical divergence and compound function is given by Goodfellow et al. (Goodfellow et al., 2020, Theorem 1) and has been proven to fulfil this assumption.

Proposition C.7 (Manipulation vs. Style-transfer). Given appropriately chosen measure of statistical divergence and compound functions, the strong manipulator with excluded learner dsc optimally refines the CycleGAN with excluded learners dsc₁, dsc₂.

1396

1364 1365

1367

1369

1371 1372 1373

Proof. Let $\mathcal{X}_1, \mathcal{X}_2$ be two distributions over the same type. We create $\mathcal{X} = (x, i)$ for $x \in \mathcal{X}_i$ where the label *i* indicates the distribution the data point came from. Assuming that we have a strong manipulator (get, put, disc) with an optimal partially perfect put, get, we can construct optimal partially perfect generators G_1 and G_2 for CycleGAN.

1401 For $i \in 0, 1, j = 1 - i$, we define:

1403

 $-G_i$:= $-f_i$ put -

First, we show that the perfectible tasks are indeed perfected, i.e. that both the autoencoder
 tasks are fulfilled. We have:



1414 Next, we will show that these generators are indeed globally optimal. By assumption, a generator is optimal if and only if its output distribution is indistinguishable from the training distribution. Thus, assuming that the manipulator is optimal with respect to the generative-adversarial tasks, put output distribution approaches $\mathcal{X} = (x, i)$ for $x \in \mathcal{X}_i$. But by PutGet, we have:



Therefore, by CLASSIFY, G_1 and G_2 can only return values that are akin to values in \mathcal{X}_1 and \mathcal{X}_2 respectively. As these two labels make up the entirety of \mathcal{X} , G_1 's and G_2 's output distributions are indistinguishable from \mathcal{X}_1 and \mathcal{X}_2 respectively. Therefore, they are indeed optimal generators for their respective distribution.

But then we have shown that G_1 and G_2 are indeed optimal solutions to the CycleGAN pattern and therefore that the concept manipulator is a specialisation of the CycleGAN⁵.

1430 In turn, the CycleGAN, however, is not a refinement of the strong manipulator. This 1431 means there exist solutions to CycleGAN, which violate strong manipulator. In these image 1432 translation tasks, we expect the translators to change as little as possible to go from one 1433 distribution to the other, i.e. preserving as much information from the original as possible. 1434 However, the generators of the CycleGAN could, for example, flip the images horizontally. 1435 As the autoencoder tasks have an even number of generators on each side, this would be a 'perfect solution' to the outlined task, yet not the behaviour we would want or expect. In 1436 contrast, the PutPut rule of the strong manipulator does not allow this behaviour. As 1437 such, the strong manipulator gives us more guarantees than the CycleGAN. 1438

1439 Despite the additional guarantees, the strong manipulator does not guarantee that it 1440 indeed only changes as little as necessary to go from one distribution to the other. This can 1441 be seen when considering a limited toy scenario: the data has the four objects *circle*, square, red and green. There is no unique mapping between shapes and colours. Without further 1442 information, it is therefore completely impossible to say what it would mean to change 1443 as little as possible to go from shapes to colours. When considering this toy scenario, it 1444 becomes obvious that, without specifying further restrictions on all remaining attributes, it 1445 is impossible to know which mapping is correct. Yet, despite these theoretical concerns, in 1446 practice, even the weaker manipulator often converges to the desired behaviour, similar to 1447 CycleGANs, which have even fewer guarantees. 1448

1448 1449 1450

1455

1407

1408

1413

1419

1420 1421

1422

C.3 Proof of Proposition 3.3

1451 Definition C.8 (Balanced entropy of an attribute). Given a distribution of data \mathcal{D} on space *D*, we say that a distribution \mathcal{A} over space *A* represents an *entropy-balanced attribute* of \mathcal{D} if there exists a complement type *C* with distribution \mathcal{C} such that we have the equality in distributions $\mathcal{D}=\mathcal{A} \times \mathcal{C}$ up an isomorphism of the underlying spaces $D \simeq (A \times C)$.

⁵In (Zhu et al., 2017), when doing style transfer, they additionally add the *identity loss* which enforces that generating an image in X_i given an image from X_1 should return the identity. This perfectible task is also guaranteed by the strong manipulator by cls and GetPut 1461

1477

1478 1479 1480



Lemma C.10 (Perfect-task representation of equal entropy). Balanced entropy in the attributes is equivalent to the perfectability of entropy balancer.

1465*Proof.* If there is balanced entropy in the attributes, Task C.9 may be satisfied up to identity by
the enc and dec witnessing the isomorphism $D \simeq (A \times C)$ that realises $\mathcal{D} = \mathcal{A} \times \mathcal{C}$. Conversely,
if the tasks above are equalities, we recover the definition of balanced entropy.1468

Contextually, we will assume that distributions have full support over spaces; this can always
 morally be the case by restrictions to subspaces. This assumption strengthens Lemma C.10
 to yield the following corollary.

1472 **Corollary C.11.** If entropy balancer is perfect, enc and dec witness an isomorphism $D \simeq (A \times C)$.

Lemma C.12 (Hardcoding latent-spaces). Given a perfected entropy balancer, we may construct the following composites, which are the put and get of a strong manipulator⁶.



Proof. First, we argue that the lens laws are satisfied. PutGet, GetPut, UNDOABILITY, and
 PUTPut follow from pure diagrammatic reasoning and applying Corollary C.11. For PutGet:



⁶modulo cLASSIFY, which in this setting is trivially obtained as we assume the attribute is derived from get.

1512By construction, putting a by independently sampling the marginal on A is indistinguishable1513in distribution from the original distribution, hence the GAN laws of strong manipulator1514are optimally satisfied, and we are done.

Definition C.13 (Bayesian Inversion in Markov Categories). In a Markov category, the Bayesian inversion (Cho & Jacobs, 2019; nLab authors, 2024a) of a stochastic map $f : X \to Y$ with respect to a distribution \mathcal{P} on X is a stochastic map $f^{\dagger} : Y \to X$ such that, in distribution:

$$\mathcal{P} \vdash X f = \mathcal{P} \vdash X f = \mathcal{P} \vdash X f = \mathcal{P}$$

Proposition C.14 (strong manipulator as Bayesian inversion). If a discriminative classifier cls : $D \rightarrow A$ induces a balanced attribute cls(D) over A with respect to D, there exists a strong manipulator for which the put induces the Bayesian inversion cls[†] : $A \rightarrow D$.

Proof. We show that, under our premises, the put composed with an independent copy of the data source \mathcal{X} is the Bayesian inversion of cls.



1550 D.1 Spriteworld

1519 1520 1521

1526

1527

1551

1552

1553

For the Spriteworld experiment, we procedurally generate 32x32 images containing a single shape with the following attributes:

1554		
1555	Attributes	Possible Values
1556	Shape	{ Ellipse, Rectangle, Triangle }
1557	Hue	{ Red : 0±8, Green : 85±8, Blue : 170±8 }
1558	Saturation	64-255
1559	Value	64-255
1560	Background Color	Black
1561	Width & Height	5-27
1562	X & Y position	5-27
1500		

Only the first two attributes, shape and hue, are changed in the manipulation task, but all unchanged properties are intended to be preserved by the transformation. We use an autoencoder with a CNN/DCNN architecture to embed each image into a latent space:

1566		
1567	Parameter	Value
1568	Latent Size	32
1569	Layers	4
1570	Hidden Channels	64
1571	Kernel Size	5x5
1572	Stride	2 Lealer Dat U(0.1) faller a har Datab Norm
1570	Activation Function	Leaky ReLU(0.1) Ionowed by BatchNorm

1574 We train separate get/put models for each of the three concepts: shape, colour, blue-1575 circleness. Each model uses the encoder of the autoencoder to embed input images into latent space, and only sees the labels for the particular attribute it is manipulating.

1577 For shape and colour, the get model uses a linear classifier from the latent space (of size 1578 32) to 3 output logit values, one for each possible value. The put model maps a one-hot 1579 vector of the input value to a vector in latent space that is added to the embedding. This 1580 new embedding is then decoded by the autoencoder. 1581

For blue-circleness, the get model uses a linear classifier from the latent space to a single output value from zero to one (we do not use a sigmoid output layer to restrict the output). 1583 The put model uses a complement of size 8. It concatenates the one-hot value vector with the image embedding and the complement vector (using a default trainable complement 1585 vector if one is not provided) and passes that through a linear layer to get a new embedding 1586 vector (which is then decoded) and complement vector. 1587

In total, these models contain 644,130 parameters. All models are trained simultaneously according to the autoencoding and manipulation rules, along with PutPut. At each step, a batch of images is generated, along with four batches of random values, containing 1590 random labels for the shape and colour manipulators, and random real numbers for the 1591 blue-circleness manipulator, uniformly sampled from [-0.1, 1.1] and then clamped to [0, 1]. 1592 The loss function is a weighted sum of the losses from each atomic task in order to balance 1593 the signal from the image loss with the signal from the classifier loss: 1594

6	Hyper-parameter	Value	Task	Weight
7	Steps	100,000	AUTOENCODING	100
8	Batch Size	512	GetPut	1
9	Optimiser	AdamW	ΡυτΡυτ	1
0	Learning Rate	10^{-3}	Undo	10
1	Weight Decay	10^{-2}	PutGet	
	Gradient Clipping	1 (element-wise)	(blue-circleness)	10
	Image Loss	$L_{2} + 0.25 \cdot L_{1}$	(shape and colour)	1
	Discrete Value Loss	Binary cross-entropy	CLASSIFICATION	
	Continuous Value Loss	Mean squared error	(blue-circleness)	10
	Seed	0	(shape and colour)	1

D.2 Faces

1607

For the faces experiment, we use the *CelebFaces Attributes* dataset (Liu et al., 2015), with an 1610 off-the-shelf data augmentation method called "TrivialAugment" (Müller & Hutter, 2021). 1611 Again, we use an autoencoder with a CNN/DCNN architecture to embed each image: 1612

1613		
1614	Parameter	Value
1615	Latent Size	128
1616	Layers	5
1617	Hidden Channels	8, 16, 32, 64, 128
1618	Kernel Size	5
1619	Stride Activation Function	2 LeakyReLU(0.1) followed by BatchNorm

We train linear get/put models for the binary concept of "Smiling", resulting in a total of 1,071,749 parameters. The loss function is a weighted sum of the losses from each atomic task:

Hyper-parameter	Value		Task	Weight
Steps	100,000	-	AUTOENCODING	10
Batch Size	64		GetPut	1
Optimiser	AdamW		ΡυτΡυτ	1
Learning Rate	10^{-3}		Undo	1
Weight Decay	10^{-2}		PutGet	1
Gradient Clipping	1 (element-wise)		CLASSIFICATION	1
Image Loss	$L_2 + 0.2 \cdot L_1 + SSIM$	-		
Value Loss	Binary cross-entropy			
Seed	0			

D.3 MNIST

We trained the manipulator pattern on the MNIST dataset, using the digit label as the property. The get operated directly on images, while put was trained to act on the latent space of an autoencoder, as in option (1) of Section 3.3. The images are input as 28×28 matrices, flattened to 784-dimensional vectors, and the labels are provided as 10-dimensional vectors with one-hot encoding. All of the components were structured as multilayer perceptrons. The hyperparameters are given below:

1649					
1650		enc	dec	put	get
1651	Input Dimension	$784 = 28 \times 28$	32	42 = 32 + 10	$784 = 28 \times 28$
1652	Output Dimension	32	$784 = 28 \times 28$	32	10
1653	Hidden Dimensions	$\{128, 128, 64\}$	$\{64, 128, 128\}$	$\{128, 128, 128\}$	{64, 64}
1654	Hidden Activations	ReLU	ReLU	ReLU	ReLU
1655	Final Activation	Sigmoid	Sigmoid	Sigmoid	Softmax

With these architectural components, we trained four tasks: (a) training get supervised, (b) training get given pre-trained put, enc, and dec, (c) training put given pre-trained get', enc, and dec, and (d) training get to match a previous get'. These were trained using the manipulation rules, as well as PutPut, and additional regularization term we denote as ENTROPY. The loss function of ENTROPY is given by

$$\mathcal{L}_{\text{Entropy}} = \mathbb{E}[H(\texttt{get}(\texttt{enc}(x)))] - H(\mathbb{E}[\texttt{get}(\texttt{enc}(x))])$$

where x is a batch of input images, $H(\cdot)$ is the entropy of a categorical distribution, and the expectation is approximated by the mean over each batch. The idea behind ENTROPY is to encourage the output of get to be well-distributed across labels (by maximizing the entropy of the mean distribution) but to be sure of each label (by minimizing the entropy for each specific input). For task (d), labels were generated for the CLASSIFY rule using get'. Each task is trained by minimizing a weighted linear combination of the rules. We give the hyperparameters, rule weights, and loss functions for each of these below.

	(a))	((b)	(c))	(d))
Optimizer Learning Rate Epochs	Adam 0.001 20		Adam 0.001 20		Adam 0.0001 20		Adam 0.001 20	
	Weight	Loss	Weight	Loss	Weight	Loss	Weight	Loss
CLASSIFY	1	CE	_	_	_	_	1	CE
PutGet	_	_	10	CE	10	CE	_	_
GetPut	_	_	10	L2	10	L2	_	_
ΡυτΡυτ	_	_	10	L2	10	L2	_	_
Undoability	_	_	10	L2	10	L2	_	_
Entropy	_	_	1	$\mathcal{L}_{\text{Entropy}}$	_	_	-	_

We also have an additional task (e) of training enc and dec unsupervised. This was done using the Adam optimizer, with a learning rate of 0.001 for 80 epochs. The reconstruction loss was given by

 $\mathcal{L}(x, \hat{x}) = L2(x, \hat{x}) + (1 - SSIM(x, \hat{x}))$

where SSIM is the structure similarity image metric.

To produce Figure 2, an enc, dec, put, and get were trained using (a) \rightarrow (e) \rightarrow (c), followed by training (c) for an additional 40 epochs. A slightly larger (but still MLP-based) model, where both put and get act on the latent space of the autoencoder, was used to achieve better visual quality. The hyperparameters are detailed below:

	enc	dec	put	get
Input Dimension	$784 = 28 \times 28$	32	42 = 32 + 10	32
Output Dimension	32	$784 = 28 \times 28$	32	10
Hidden Dimensions	{128, 128, 128, 32, 32}	{32, 32, 128, 128, 128}	{256, 256}	{256}
Hidden Activations	ReLU	ReLU	ReLU	ReLU
Final Activation	Sigmoid	Sigmoid	Sigmoid	Softmax

enc, dec, and put were used to manipulate six examples picked from the dataset, putting each of the ten classes onto each example. The examples were cherrypicked to provide maximum stylistic contrast across the sample but were not selected for maximum style transfer accuracy - a similar level was observed across the entire dataset. Code for all of these models, as well as the training schedules of tasks (a)-(e), are provided in the supplementary material.

E STACKS

Here we define tasks that obtain learned operators that implement the well-known data structure stack. There are two interacting operations: *push* and *pop*. By itself, the stack is not particularly interesting, as it can algorithmically be implemented rather easily, but a fully differentiable and automatically learned stack may be useful for other downstream tasks.



Given a data distribution \mathcal{X} over X and a pretrained autoencoder (enc: $X \rightarrow LAT$, dec: $LAT \rightarrow X$), a stack consists of a *empty stack* \bot : $\star \rightarrow S$, a push operation *psh* : $S \times X \rightarrow S$ and a pop operation *pop* : $S \rightarrow S \times X$. We recursively define the distribution:

 $\alpha := \bot \otimes \operatorname{enc}(\mathcal{X}) \mid \operatorname{psh}(\alpha) \otimes \operatorname{enc}(\mathcal{X})$

meaning α is a stack of arbitrary size and an encoded element from \mathcal{X} .

Then psh and pop have to obey the rules on the left.

Experimentally in the Spriteworld setting, we observe that the stack works exactly as 1741 expected. The stack space *S* has to be *n* times larger than the latent space **Lat** for the stack to 1742 be able to hold up to *n* items. At the moment, we do not observe any additional information 1743 compression in this basic setup. The experiment (Figure 7) uses an autoencoder architecture 1744 for processing 32×32 RGB images, with a latent space size of 16 dimensions. The encoder 1745 consists of four convolutional layers, each with 64 channels and a channel multiplier of 1. 1746 Additionally, a stack of 64 latent features is processed through an MLP with 256 hidden 1747 units. Training is performed on a GPU (if available), with a batch size of 64, learning rate of 1×10^{-4} , weight decay of 1×10^{-2} , and gradient clipping at 1. The model trains for 1748 100,000 steps, logging every 10,000 steps. Input images are converted to tensors and scaled 1749 to floating-point precision, and a fixed random seed of 0 ensures reproducibility. 1750



Figure 7: In this example, we train a stack (alongside an autoencoder) to store the latent vectors of Spriteworld shapes. With an image latent size 16 and stack vector size 64, it is able to retain information to faithfully restore up to 4 shapes.

1773 1774

1767

1768 1769

1751 1752 1753

1754

1755 1756

- 1775 F Examining Manipulation in complex domains
- 1776 1777 1778

F.1 Manipulation for text sentiment

We attempted to fine-tune an extant strong solution for text-sentiment modification by additionally imposing the constraints of manipulation on top of the original objective function. Our findings suggest that additionally imposing the constraints of manipulation on architectures that are already performant does not make an appreciable difference

1782 (Table 1). We pretrained the Blind Generative Style Transformer (B-GST) model of 1783 (Sudhakar et al., 2019) which takes in the non-stylistic components of a sentence and the 1784 target sentiment, and outputs the sentence generated in the target style. This was done until 1785 we achieved baselines higher than the original ones reported by the authors of the model. Afterwards, we continued training under three different conditions: (1) resuming training with only the original objective, (2) only using objective functions from manipulation, 1787 and (3) using both. For (1) and (3), we reused the same objective in the original paper. 1788 In all experiments, we used the YELP dataset used by (Li et al., 2018), reusing the same 1789 train-dev-test split they used. It consists of 270K positive and 180K negative sentences for 1790 the training set, 2000 sentences each for the dev set, and 500 sentences each for the test set. 1791 Furthermore, we used the human gold standard references they provided for their test set. 1792 **B-GST** uses a sequence length of 512, 12 attention blocks each with 12 attention heads. We 1793 used 768-dimensional internal states (keys, queries, values, word embeddings, positional 1794 embeddings). We tokenized the input text using Byte-Pair Encoding (BPE).

1795

1796 We used the same input autoencoding and output decoding used in (Sudhakar et al., 2019) 1797 across all experiments. For the get of the manipulation task, we used the PyTorch version of 1798 the pretrained Transformer by HuggingFace, which uses the OpenAI GPT model pretrained 1799 by (Radford & Narasimhan, 2018) on the BookCorpus dataset which contains over 7000 books with approximately 800M words. We trained it on a sentiment classification task 1801 using the YELP dataset reaching 98% accuracy on the test set. The get was fixed for the 1802 entire duration of training conditions (2) and (3) above. For the put of manipulation, we used the **B-GST** model to generate text with a specified sentiment. This was a computational 1803 bottleneck for training conditions (2) and (3) as autoregressive decoding is required to generate model inputs for PUTGET and UNDOABILITY in manipulation. We used 'teacher forcing' or 'guided approach' (Bengio et al., 2015; Williams & Zipser, 1989) whenever we 1806 computed the reconstruction loss of put. Additionally, for training conditions (2) and (3), 1807 we only used PUTGET, GETPUT, and UNDOABILITY from manipulation. We used a weighted sum of the losses computed for each of these and the original reconstruction loss if present the weights can be considered as training hyperparameters. For (2), we used (PutGet=5, 1810 GetPut=20, Undoability=20), while for (3), we used (PutGet=5, GetPut=10, Undoability=25, B-GST=30). Code for all of the models, training schedules, and hyperparameter values for 1812 training conditions (1)-(3) are also provided in the supplementary material. 1813

Model	GLEU	BLEU _{SRC}	BLEU _{REF}	ACC (fasttext)
B-GST-pretrained	11.869	74.563	52.770	84.6
B-GST-only	11.426	74.876	52.549	85.7
manipulation-only	11.712	74.428	52.646	84.1
B-GST+manipulator	11.338	74.608	52.836	85.1
Human Reference	100.00	58.158	100.00	67.6

Table 1: We pretrained the Blind Generative Style Transformer (B-GST) model (Sudhakar et al., 2019) based on the Delete-Retrieve-Generate (Li et al., 2018) framework for sentiment modification until we recovered higher baselines than reported by the authors of the model (GLEU=11.6, BLEU_{SRC}=71.0), and we continued training in three different conditions: (1) keeping the original objective, (2) only using objective functions obtained from manipulation, and (3) using both. We report no statistically significant differences in scores, even under continued training. Below we report the results of training conditions (1) and (2) for an additional epoch, and (3) for two epochs (details in Appendix F).

1829

1830 F.2 Characterising Manipulation as generative classification

Training manipulation autoregressively (e.g. when instantiating the learners as transformers for sequential data) is slow due to autoregressing twice for PutPut and Undoability.
Moreover, our attempts to autoregressively manipulate the sentiment of IMDB reviews often resulted in a form of posterior collapse where put ignored the attribute and behaved as the identity function on text. Conceptually, this is because the identity function satisfies

1836 GETPUT, PUTPUT and UNDOABILITY, and while the identity fails on PUTGET, failing on one component of the combined loss function does not provide a strong enough incentive to move away from the identity in parameter-space.

1839

1840 Notably, these shortcomings mirror that of VAEs, which also suffer from posterior collapse 1841 (Bond-Taylor et al., 2022) in highly structured domains such as video (Babaeizadeh et al., 2018) and text (Bowman et al., 2016). This suggested to us that **puts may be generative** 1843 classifiers, which could potentially explain why mode collapse was occurring in complex 1844 domains. Borrowing terminology from (Ng & Jordan, 2001), classifiers are *discriminative* if they seek to learn the conditional distribution p(a|d) of attributes given data (as in the classification pattern), and otherwise they are *generative* if they seek to learn the joint distribution p(d, a) (as, for example, a VAE). The tradeoffs between the two types are well 1847 studied, e.g. performance-wise, generative models may converge faster with limited data, but discriminative models often achieve lower asymptotic error, and it is well known that 1849 learning generative models is harder (Vapnik, 1998). 1850

1851

As generative classifiers, puts in the manipulation task appear to approximate the 1853 **Bayesian inverse of a discriminative classifier cls, approximately as well as VAEs do.** To demonstrate this, we tried three ways to train an "informationally identical" cls' given an initial cls, provided access to unlabelled data to obtain a distribution of pairs (d, cls(d)) by: 1855 (a) directly training cls by classification, (b) training a generative classifier, in our case a 1856 VAE, and (c) training manipulation around cls-as-get to obtain a put, and then train cls' 1857 to satisfy the tasks of manipulation except for CLASSIFY. Repeating this process several times, we would expect to see some loss of accuracy due to imperfect Bayesian inversion. Indeed, we see in Figure 8 that (b) and (c) have similar decays in accuracy, indicating that manipulation 1860 and VAEs have similar performance in this case. This experiment was performed using the 1861 trained components obtained from the MNIST experiment (Section D.3), and in addition to the tasks (a-e) we (f) trained a VAE to learn the joint distribution of images and labels produced by get', and we (g) trained a get supervised using labels and images generated from the VAE. The VAE encoder and decoder are also based on multilayer perceptrons each is comprised of an MLP trunk and two linear heads for generating the means and 1866 log-variances of the latent space, or the image and labels, respectively. The latent space is comprised of independent normally distributed variables as in (Kingma & Welling, 2022), and is sampled using the standard reparameterization trick. The hyperparameters of the 1868 architecture are given below:

VAE Encoder	VAE Decoder
$794 = 28 \times 28 + 10$	32
$\{128, 128, 64\}$	$\{64, 128, 128\}$
ReLU	ReLU
32	$784 = 28 \times 28$
-	Sigmoid
32	10
-	Softmax
	VAE Encoder 794 = 28 × 28 + 10 {128, 128, 64} ReLU 32 - 32 -

1883

1871 1872

1884

1885

Three loss functions were used for tasks (f) and (g) — the reconstruction loss of the autoencoder, which can be separated into a label loss and an image loss, the K-L divergence regularization term \mathcal{L}_{KL} of the VAE (Kingma & Welling, 2022), and the CLASSIFY loss of get. The training hyperparameters are given as follows:

1000					
1890		(f)		(g)	
1891					
1892	Optimizer	Adam		Adam	
1893	Learning Rate	0.001		0.001	
1894	Epochs	40		20	
1895		Weight	Loss	Weight	Loss
1896				1	CE
1897	Image Reconstruction	100	12	1	
1898	Label Reconstruction	100	CE	_	_
1899	K-L Divergence	0.5	\mathcal{L}_{KL}	_	_
1900	0		nn nn		

In order to evaluate the three methods, we trained the tasks in the following order. At each step, the component being trained (e.g. get, put, etc) was initialized randomly (the previous weights were discarded). Measurements of the test accuracy were made after each (a), (b), (d), or (g) training run, and used to produce Figure 8.



1934

1935 Figure 8: We tried three ways to train an "informationally identical" cls' given an initial cls 1936 - depicted are the results of training successive MNIST classifiers using methods (a), (b) and (c) given above. 'steps' refers to the number of times this process was repeated. We observe 1937 that the degradation of accuracy is approximately the same for both (b) and (c), which we 1938 consider evidence that manipulator and VAEs have similar performance characteristics. 1939 Both models had roughly the same number of parameters (300K) and were based on the 1940 same MLP architecture. We ran 20 repetitions of each method, the shaded regions represent 1941 one standard deviation (method (a) had a standard deviation of less than 1%). 1942

1943