

# Cyber Kittens, or Some First Steps Towards Categorical Cybernetics

Toby St. Clere Smithe

Department of Experimental Psychology,  
University of Oxford  
arxiv@tsmithe.net

We define a categorical notion of cybernetic system as a dynamical realisation of a generalized open game, along with a coherence condition. We show that this notion captures a wide class of cybernetic systems in computational neuroscience and statistical machine learning, exposes their compositional structure, and gives an abstract justification for the bidirectional structure empirically observed in cortical circuits. Our construction is built on the observation that Bayesian updates compose optically, a fact which we prove along the way, via a fibred category of state-dependent stochastic channels.

## 1 Introduction

Those systems that we might classify as living, adaptive, or somehow intelligent all display a fundamental property: they resist or avoid perturbations that would result in their existence becoming unsustainable. This means that they must somehow be able to sense their current state of affairs (*perception*) and respond appropriately (*action*). In particular, an adaptive system should sense the relevant aspects of its current environmental state, and form expectations about the consequences of that state. In general, the interaction with the environment will be stochastic, and the statistically optimal method of ‘sensing’ and prediction is Bayesian inference.

Typically, however, the system has no direct access to the external state, only to sense data that indirectly have external causes. Moreover, sense data are often very high-dimensional, and predicting their consequences is underdetermined. As a result, it is common to assume that successful organisms are imbued with some kind of *generative model* of the process by which external causes generate their sense data. They can then use this model to infer those actions will bring (their beliefs about) their current state closer to those expectations: a process called *active inference*.

Systems such as these are inherently open, and often their internal models and beliefs are supposed to be structured hierarchically—that is, compositionally. The processes of prediction and action sketched here are naturally bidirectional, and indeed our first contribution in the present work is to show that Bayesian inference is abstractly structured as a *category of optics* [21, 6], the emerging canonical formalism for (open) bidirectionally structured compositional systems.

The compositional framework of *open games* [2, 13] builds on categories of optics to describe systems of motivated interacting agents, but it is substantially more general than needed for classical game theory: generalized open games naturally describe any bidirectionally structured open systems that can be associated with a measure of fitness. Consequently, such generalized open games provide a natural home for a compositional theory of interacting cybernetic systems, and using our notion of *Bayesian lens*, we characterize a number of canonical statistical models as *statistical games*.

However, mere open games themselves supply no notion of *dynamics* mediating the interactions. We therefore introduce the concept of *dynamical realisation* of an open game (Definition 4.7), as well as a

coherence condition that ensures such a realisation behaves as we would expect from a cybernetic system (Definition 4.8). We use these concepts to show that two prominent frameworks for active inference instantiate such categories of cybernetic systems.

**Acknowledgements** We thank the organizers of *Applied Category Theory 2020* for the opportunity to present this work, and the anonymous reviewers for helpful comments and questions. We also thank Bruno Gavranović, Jules Hedges, and Neil Ghani for stimulating and insightful conversations, and credit Jules Hedges for observing the correct form of the Bayesian update map in discussions at SYCO 6.

## 2 Bayesian Updates Compose Optically

We begin by proving that Bayesian updates compose according to the ‘lens’ pattern [9] that sits at the heart of categories of open games and other ‘bidirectional’ structures. We first show that Bayesian inversions are ‘vertical’ maps in a fibred category of state-dependent channels. The Grothendieck construction of this structure gives a category of lenses. Open games are commonly defined using the more general ‘optics’ pattern [2], and so we also show that, under the Yoneda embedding, our category of lenses is equivalently a category of optics.

Throughout the paper, we work in a general category of stochastic channels; abstractly, this corresponds to a *Markov category* [12] or *copy-delete category* [5]. Familiar examples of such categories include  $\mathcal{Kl}(\mathcal{D})$ , the Kleisli category of the finitely-supported distribution monad  $\mathcal{D}$ , and, for ‘continuous’ probability,  $\mathcal{Kl}(\mathcal{G})$ , the Kleisli category of the Giry monad. We will write  $c_\pi^\dagger := c_{(\cdot)}^\dagger(\pi)$  to indicate the **Bayesian inversion** of the channel  $c$  with respect to a state  $\pi$ . Then, given some  $y \in Y$ ,  $c_\pi^\dagger(y)$  is a new ‘posterior’ distribution on  $X$ . We will call  $c_\pi^\dagger(y)$  the **Bayesian update** of  $\pi$  along  $c$  given  $y$ .

For a substantially expanded version of this section, including proofs and background exposition with precise definitions of Bayesian inversion, see the author’s [25]. We will occasionally here refer to definitions or results in that paper.

**Definition 2.1** (State-indexed categories). Let  $(\mathcal{C}, \otimes, I)$  be a monoidal category enriched in a Cartesian closed category  $\mathbf{V}$ . Define the  $\mathcal{C}$ -state-indexed category  $\text{Stat} : \mathcal{C}^{\text{op}} \rightarrow \mathbf{V}\text{-Cat}$  as follows.

$\text{Stat} : \mathcal{C}^{\text{op}} \rightarrow \mathbf{V}\text{-Cat}$

$$X \mapsto \text{Stat}(X) := \left( \begin{array}{lll} \text{Stat}(X)_0 & := & \mathcal{C}_0 \\ \text{Stat}(X)(A, B) & := & \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(A, B)) \\ \text{id}_A : \text{Stat}(x)(A, A) & := & \begin{cases} \text{id}_A : \mathcal{C}(I, X) \rightarrow \mathcal{C}(A, A) \\ \rho \mapsto \text{id}_A \end{cases} \end{array} \right) \quad (1)$$

$$f : \mathcal{C}(Y, X) \mapsto \left( \begin{array}{lll} \text{Stat}(f) : \text{Stat}(X) & \rightarrow & \text{Stat}(Y) \\ & \text{Stat}(X)_0 = & \text{Stat}(Y)_0 \\ \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(A, B)) & \rightarrow & \mathbf{V}(\mathcal{C}(I, Y), \mathcal{C}(A, B)) \\ \alpha & \mapsto & f^* \alpha : (\sigma : \mathcal{C}(I, Y)) \mapsto (\alpha(f \bullet \sigma) : \mathcal{C}(A, B)) \end{array} \right)$$

Composition in each fibre  $\text{Stat}(X)$  is given by composition in  $\mathcal{C}$ ; that is, by the left and right actions of the profunctor  $\text{Stat}(X)(-, =) : \mathcal{C}^{\text{op}} \times \mathcal{C} \rightarrow \mathbf{V}$ . Explicitly, given  $\alpha : \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(A, B))$  and  $\beta : \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(B, C))$ , their composite is  $\beta \circ \alpha : \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(A, C)) := \rho \mapsto \beta(\rho) \bullet \alpha(\rho)$ . Since  $\mathbf{V}$

is Cartesian, there is a canonical copier  $\gamma : x \mapsto (x, x)$  on each object, so we can alternatively write  $(\beta \circ \alpha)(\rho) = (\beta(-) \bullet \alpha(-)) \circ \gamma \circ \rho$ . Note that we indicate composition in  $\mathcal{C}$  by  $\bullet$  and composition in the fibres  $\text{Stat}(X)$  by  $\circ$ .

**Example 2.2.** Let  $\mathbf{V} = \mathbf{Meas}$  be a ‘convenient’ (i.e., Cartesian closed) category of measurable spaces, such as the category of quasi-Borel spaces [14], let  $\mathcal{P} : \mathbf{Meas} \rightarrow \mathbf{Meas}$  be a probability monad defined on this category, and let  $\mathcal{C} = \mathcal{Kl}(\mathcal{P})$  be the Kleisli category of this monad. Its objects are the objects of  $\mathbf{Meas}$ , and its hom-spaces  $\mathcal{Kl}(\mathcal{P})(A, B)$  are the spaces  $\mathbf{Meas}(A, \mathcal{P}B)$  [12]. This  $\mathcal{C}$  is a monoidal category of stochastic channels, whose monoidal unit  $I$  is the space with a single point. Consequently, states of  $X$  are just measures (distributions) in  $\mathcal{P}X$ . That is,  $\mathcal{Kl}(\mathcal{P})(I, X) \cong \mathbf{Meas}(1, \mathcal{P}X)$ . Instantiating  $\text{Stat}$  in this setting, we obtain:

$$\begin{aligned} \text{Stat} &: \mathcal{Kl}(\mathcal{P})^{\text{op}} \rightarrow \mathbf{V}\text{-Cat} \\ X \mapsto \text{Stat}(X) &:= \left( \begin{array}{lcl} \text{Stat}(X)_0 & := & \mathbf{Meas}_0 \\ \text{Stat}(X)(A, B) & := & \mathbf{Meas}(\mathcal{P}X, \mathbf{Meas}(A, \mathcal{P}B)) \\ \text{id}_A : \text{Stat}(X)(A, A) & := & \begin{cases} \text{id}_A : \mathcal{P}X \rightarrow \mathbf{Meas}(A, \mathcal{P}A) \\ \rho \mapsto \eta_A \end{cases} \end{array} \right) \end{aligned} \quad (2)$$

$$c : \mathcal{Kl}(\mathcal{P})(Y, X) \mapsto \text{Stat}(c) :=$$

$$\left( \begin{array}{lcl} \text{Stat}(c) : & \text{Stat}(X) & \rightarrow \text{Stat}(Y) \\ & \text{Stat}(X)_0 & = \text{Stat}(Y)_0 \\ \left( \begin{array}{lcl} d^\dagger : \mathcal{P}X & \rightarrow & \mathcal{Kl}(\mathcal{P})(A, B) \\ \pi & \mapsto & d_\pi^\dagger \end{array} \right) & \mapsto & \left( \begin{array}{lcl} c^* d^\dagger : \mathcal{P}Y & \rightarrow & \mathcal{Kl}(\mathcal{P})(A, B) \\ \rho & \mapsto & d_{c \bullet \rho}^\dagger \end{array} \right) \end{array} \right)$$

Each  $\text{Stat}(X)$  is a category of stochastic channels with respect to measures on the space  $X$ . We can write morphisms  $d^\dagger : \mathcal{P}X \rightarrow \mathcal{Kl}(\mathcal{P})(A, B)$  in  $\text{Stat}(X)$  as  $d_{(\cdot)}^\dagger : A \xrightarrow{(\cdot)} B$ , and think of them as generalized Bayesian inversions: given a measure  $\pi$  on  $X$ , we obtain a channel  $d_\pi^\dagger : A \xrightarrow{\pi} B$  with respect to  $\pi$ . Given a channel  $c : Y \rightarrow X$  in the base category of priors, we can pull  $d^\dagger$  back along  $c$ , to obtain a  $Y$ -dependent channel in  $\text{Stat}(Y)$ ,  $c^* d^\dagger : \mathcal{P}Y \rightarrow \mathcal{Kl}(\mathcal{P})(A, B)$ , which takes  $\rho : \mathcal{P}Y$  to the channel  $d_{c \bullet \rho}^\dagger : A \xrightarrow{c \bullet \rho} B$  defined by pushing  $\rho$  through  $c$  and then applying  $d^\dagger$ .

**Remark 2.3.** Note that by taking  $\mathbf{Meas}$  to be Cartesian closed, we have  $\mathbf{Meas}(\mathcal{P}X, \mathbf{Meas}(A, \mathcal{P}B)) \cong \mathbf{Meas}(\mathcal{P}X \times A, \mathcal{P}B)$  for each  $X, A$  and  $B$ , and so a morphism  $c^\dagger : \mathcal{P}Y \rightarrow \mathcal{Kl}(\mathcal{P})(X, Y)$  equivalently has the type  $\mathcal{P}Y \times X \rightarrow \mathcal{P}Y$ . Paired with a channel  $c : Y \rightarrow \mathcal{P}X$ , we have something like a Cartesian lens; and to compose such pairs, we can use the Grothendieck construction [20, 26].

**Definition 2.4 (GrLens<sub>Stat</sub>).** Instantiating the category of Grothendieck  $F$ -lenses  $\mathbf{GrLens}_F$  (see [26]) with  $F = \text{Stat} : \mathcal{C}^{\text{op}} \rightarrow \mathbf{V}\text{-Cat}$ , we obtain the category  $\mathbf{GrLens}_{\text{Stat}}$  whose objects are pairs  $(X, A)$  of objects of  $\mathcal{C}$  and whose morphisms  $(X, A) \rightarrow (Y, B)$  are elements of the set

$$\mathbf{GrLens}_{\text{Stat}}((X, A), (Y, B)) \cong \mathcal{C}(X, Y) \times \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(B, A)). \quad (3)$$

The identity  $\text{Stat}$ -lens on  $(Y, A)$  is  $(\text{id}_Y, \text{id}_A)$ , where by abuse of notation  $\text{id}_A : \mathcal{C}(I, Y) \rightarrow \mathcal{C}(A, A)$  is the constant map  $\text{id}_A$  defined in (1) that takes any state on  $Y$  to the identity on  $A$ . The sequential composite of  $(c, c^\dagger) : (X, A) \rightarrow (Y, B)$  and  $(d, d^\dagger) : (Y, B) \rightarrow (Z, C)$  is the  $\text{Stat}$ -lens  $((d \bullet c), (c^\dagger \circ c^* d^\dagger)) : (X, A) \rightarrow (Z, C)$  with  $(d \bullet c) : \mathcal{C}(X, Z)$  and where  $(c^\dagger \circ c^* d^\dagger) : \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(C, A))$  takes a state  $\pi : \mathcal{C}(I, X)$  on  $X$  to the

channel  $c_\pi^\dagger \bullet d_{c \bullet \pi}^\dagger$ . If we think of the notation  $(\cdot)^\dagger$  as denoting the operation of forming the Bayesian inverse of a channel (in the case where  $A = X$ ,  $B = Y$  and  $C = Z$ ), then the main result of this section is to show that  $(d \bullet c)_\pi^\dagger \stackrel{d \bullet c \bullet \pi}{\sim} c_\pi^\dagger \bullet d_{c \bullet \pi}^\dagger$ , where  $\stackrel{d \bullet c \bullet \pi}{\sim}$  denotes  $(d \bullet c \bullet \pi)$ -almost-equality [25, Definition 2.5].

In order to give an optical form for  $\mathbf{GrLens}_{\text{Stat}}$ , we need to find two  $\mathcal{M}$ -actegories with a common category of actions  $\mathcal{M}$ . Let  $\hat{\mathcal{C}}$  and  $\check{\mathcal{C}}$  denote the categories  $\hat{\mathcal{C}} := \mathbf{V}\text{-Cat}(\mathcal{C}^{\text{op}}, \mathbf{V})$  and  $\check{\mathcal{C}} := \mathbf{V}\text{-Cat}(\mathcal{C}, \mathbf{V})$  of presheaves and copresheaves on  $\mathcal{C}$ , and consider the following natural isomorphisms.

$$\begin{aligned} \mathbf{GrLens}_{\text{Stat}}((X, A), (Y, B)) &\cong \mathcal{C}(X, Y) \times \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(B, A)) \\ &\cong \int^{M: \mathcal{C}} \mathcal{C}(X, Y) \times \mathcal{C}(X, M) \times \mathbf{V}(\mathcal{C}(I, M), \mathcal{C}(B, A)) \\ &\cong \int^{\hat{M}: \hat{\mathcal{C}}} \mathcal{C}(X, Y) \times \hat{M}(X) \times \mathbf{V}(\hat{M}(I), \mathcal{C}(B, A)) \end{aligned} \quad (4)$$

The second isomorphism follows by Yoneda reduction [17, 23], and the third follows by the Yoneda lemma. We take  $\mathcal{M}$  to be  $\mathcal{M} := \hat{\mathcal{C}}$ , and define an action  $\odot$  of  $\hat{\mathcal{C}}$  on  $\check{\mathcal{C}}$  as follows.

**Definition 2.5** ( $\odot$ ). We give only the action on objects; the action on morphisms is analogous.

$$\begin{aligned} \odot : \hat{\mathcal{C}} &\rightarrow \mathbf{V}\text{-Cat}(\check{\mathcal{C}}, \check{\mathcal{C}}) \\ \hat{M} &\mapsto \left( \begin{array}{ccc} \hat{M} \odot - & : & \check{\mathcal{C}} \rightarrow \check{\mathcal{C}} \\ P & \mapsto & \mathbf{V}(\hat{M}(I), P) \end{array} \right) \end{aligned} \quad (5)$$

Functoriality of  $\odot$  follows from the functoriality of copresheaves.  $\square$

**Proposition 2.6.**  $\odot$  equips  $\check{\mathcal{C}}$  with a  $\hat{\mathcal{C}}$ -actegory structure: unitor isomorphisms  $\lambda_F^\odot : 1 \odot F \xrightarrow{\sim} F$  and associator isomorphisms  $a_{\hat{M}, \hat{N}, F}^\odot : (\hat{M} \times \hat{N}) \odot F \xrightarrow{\sim} \hat{M} \odot (\hat{N} \odot F)$  for each  $\hat{M}, \hat{N}$  in  $\hat{\mathcal{C}}$ , both natural in  $F : \mathbf{V}\text{-Cat}(\mathcal{C}, \mathbf{V})$ .

We are now in a position to define the category of abstract Bayesian lenses, and show that this category coincides with the category of Stat-lenses.

**Definition 2.7** (Bayesian lenses). Denote by  $\mathbf{BayesLens}$  the category of optics  $\mathbf{Optic}_{\times, \odot}$  for the action of the Cartesian product on presheaf categories  $\times : \hat{\mathcal{C}} \rightarrow \mathbf{V}\text{-Cat}(\hat{\mathcal{C}}, \hat{\mathcal{C}})$  and the action  $\odot : \hat{\mathcal{C}} \rightarrow \mathbf{V}\text{-Cat}(\check{\mathcal{C}}, \check{\mathcal{C}})$  defined in (5). Its objects  $(\hat{X}, \check{Y})$  are pairs of a presheaf and a copresheaf on  $\mathcal{C}$ , and its morphisms  $(\hat{X}, \check{A}) \rightarrow (\hat{Y}, \check{B})$  are abstract *Bayesian lenses*—elements of the type

$$\mathbf{Optic}_{\times, \odot}((\hat{X}, \check{A}), (\hat{Y}, \check{B})) = \int^{\hat{M}: \hat{\mathcal{C}}} \mathcal{C}(\hat{X}, \hat{M} \times \hat{Y}) \times \check{\mathcal{C}}(\hat{M} \odot \check{B}, \check{A}) \quad (6)$$

Given  $v : \mathcal{C}(X, Y)$  and  $u : \mathbf{V}(\mathcal{C}(I, X), \mathcal{C}(B, A))$ , we denote the corresponding element of this type by  $\langle v | u \rangle$ . A Bayesian lens  $(\hat{X}, \check{X}) \rightarrow (\hat{Y}, \check{Y})$  is called a **simple** Bayesian lens.

**Proposition 2.8.**  $\mathbf{BayesLens}$  is a category of lenses; a definition is given in [25, §2.2.1].

**Proposition 2.9** (Stat-lenses are Bayesian lenses). Let  $(\hat{\cdot}) : \mathcal{C} \hookrightarrow \mathbf{V}\text{-Cat}(\mathcal{C}^{\text{op}}, \mathbf{V})$  denote the Yoneda embedding and  $(\check{\cdot}) : \mathcal{C} \hookrightarrow \mathbf{V}\text{-Cat}(\mathcal{C}, \mathbf{V})$  the coYoneda embedding. Then

$$\mathbf{Optic}_{\times, \odot}((\hat{X}, \check{A}), (\hat{Y}, \check{B})) \cong \mathbf{GrLens}_{\text{Stat}}((X, A), (Y, B)) \quad (7)$$

so that  $\mathbf{GrLens}_{\text{Stat}}$  is equivalent to the full subcategory of  $\mathbf{Optic}_{\times, \odot}$  on representable (co)presheaves.

**Remark 2.10.** We will often abuse notation by indicating representable objects in **BayesLens** by their representations in  $\mathcal{C}$ . That is, we will write  $(X, A)$  instead of  $(\hat{X}, \check{A})$  where this would be unambiguous.

**Proposition 2.11.** **BayesLens** is a symmetric monoidal category. The monoidal product  $\otimes$  is inherited from  $\mathcal{C}$ ; the unit object is the pair  $(I, I)$  where  $I$  is the unit object in  $\mathcal{C}$ . For more details on the structure, see [21] or [19].

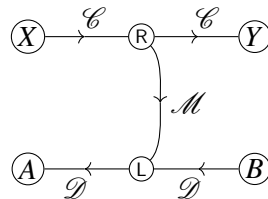
**Definition 2.12** (Exact and approximate Bayesian lens). Let  $\langle c | c^\dagger \rangle : (X, X) \rightarrow (Y, Y)$  be a simple Bayesian lens. We say that  $\langle c | c^\dagger \rangle$  is **exact** if  $c$  admits Bayesian inversion and, for each  $\pi : I \rightarrow X$  such that  $c \bullet \pi$  has non-empty support,  $c^\dagger_\pi$  is the Bayesian inversion of  $c$  with respect to  $\pi$ . Simple Bayesian lenses that are not exact are said to be **approximate**.

**Lemma 2.13.** Let  $\langle c | c^\dagger \rangle$  and  $\langle d | d^\dagger \rangle$  be sequentially composable exact Bayesian lenses. Then the contravariant component of the composite lens  $\langle d | d^\dagger \rangle \circ \langle c | c^\dagger \rangle \cong \langle d \bullet c | c^\dagger \circ c^\dagger d^\dagger \rangle$  is, up to  $d \bullet c \bullet \pi$ -almost-equality, the Bayesian inversion of  $d \bullet c$  with respect to any state  $\pi$  on the domain of  $c$  such that  $c \bullet \pi$  has non-empty support. That is to say, *Bayesian updates compose optically*:  $(d \bullet c)^\dagger_\pi \stackrel{d \bullet c \bullet \pi}{\sim} c^\dagger_\pi \bullet d^\dagger_{c \bullet \pi}$ .

### 3 Open Games for General Optics

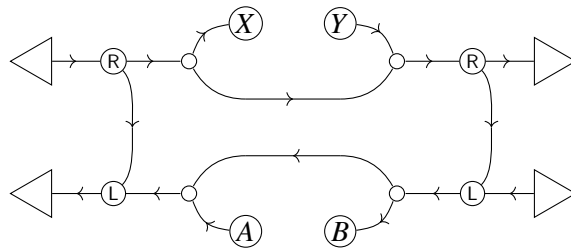
In this section, we supply mild generalizations of the structures underlying open games, building on those in [2]; at first, then, we consider games over arbitrary categories of optics  $\mathbf{Optic}_{\otimes, \circ}$ . Subsequently, we use games over Bayesian lenses (in the category of optics **BayesLens** introduced above) to exemplify a number of canonical statistical concepts, such as maximum likelihood estimation and the variational autoencoder, and clarify their compositional structure using the notion of *optimization game* (Definition 3.22). Owing to space constraints, we omit most proofs in this section; they will appear in a full paper expanding the present abstract, and can be supplied at the request of the reader.

**Observation 3.1.** In the graphical calculus for the compact closed bicategory of profunctors **Prof** [22], the hom object  $\mathbf{Optic}_{\otimes, \circ}((X, A), (Y, B))$  has the depiction



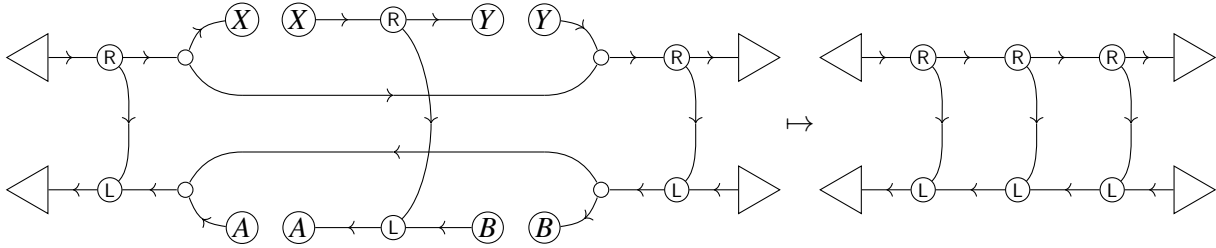
where the types on the wires are the 0-cells of **Prof**, the monoidal actions  $\otimes$  and  $\circ$  are depicted as (co)monoids, and the states and effects are (co)representable functors on the objects  $X, A, Y, B$ , treated as profunctors.

**Definition 3.2** (Generalized context). The context functor  $C : \mathbf{Optic}_{\otimes, \circ}^{\text{op}} \times \mathbf{Optic}_{\otimes, \circ} \rightarrow \mathbf{Set}$  takes the pair of optical objects  $((X, A), (Y, B))$  to the type with depiction



The triangles depict the (co)presheaves on the monoidal unit  $I$  in the underlying actegories. The action on morphisms (*i.e.*, optics) is by precomposition on the left and postcomposition on the right. Functoriality follows accordingly.  $\square$

We can compose a context with an optic to obtain a ‘closed’ system, as follows:



**Conjecture 3.3.** It is easy to show that a context on  $((X,A),(Y,B))$  is equivalently a state  $(I,I) \rightarrow ((X,A),(Y,B))$  in the monoidal category of ‘double lenses’,  $\mathbf{Lens}_{\mathbf{Optic}_{\otimes,\otimes}}$  [2]. Rendering this graphically leads us to the following conjecture: categories of double optics are instances of the *doubling* or *CP* construction from categorical quantum mechanics (*cf.* [8, 7]).

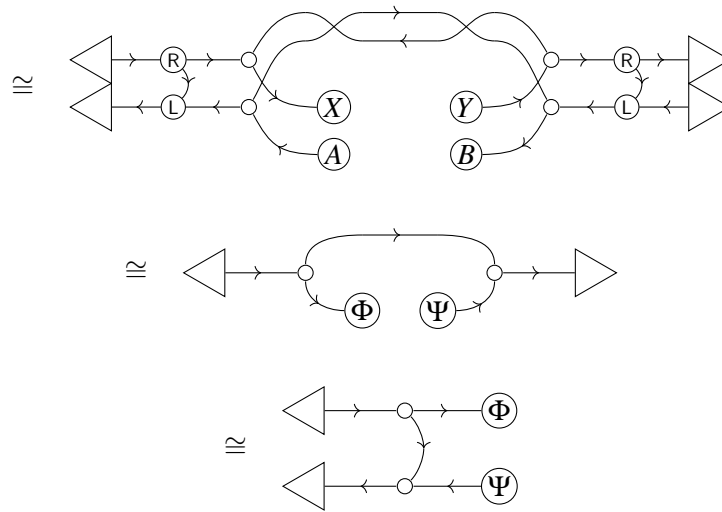
**Observation 3.4** (Justification for Conjecture 3.3). Assuming  $\mathbf{Optic}_{\otimes,\otimes}$  is monoidal, a context on  $((X,A),(Y,B))$  is equivalently a monoidal lens  $((I,I),(I,I)) \rightarrow ((X,A),(Y,B))$  over  $\mathbf{Optic}_{\otimes,\otimes}$ :

$$\int^{(M,N):\mathbf{Optic}_{\otimes,\otimes}} \mathbf{Optic}_{\otimes,\otimes}((I,I),(M,N) \otimes (X,A)) \times \mathbf{Optic}_{\otimes,\otimes}((M,N) \otimes (Y,B),(I,I)) \\ \cong \mathbf{Optic}_{\otimes,\otimes} \left( ((I,I),(I,I)), ((X,A),(Y,B)) \right)$$

where  $\otimes$  is here the induced monoidal product on  $\mathbf{Optic}_{\otimes,\otimes}$ . Graphically, we can represent this isomorphism as follows:

$$\int^{(M,N):\mathbf{Optic}_{\otimes,\otimes}} \mathbf{Optic}_{\otimes,\otimes}((I,I),(M,N) \otimes (X,A)) \times \mathbf{Optic}_{\otimes,\otimes}((M,N) \otimes (Y,B),(I,I)) \\ \cong \mathbf{Optic}_{\otimes,\otimes} \left( ((I,I),(I,I)), ((X,A),(Y,B)) \right)$$

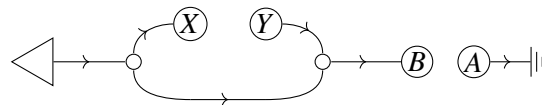
The diagram illustrates the isomorphism between the integral expression and the simplified form. The top part shows the integral expression with a complex network of nodes and arrows. The middle part shows a simplified version of the same system. The bottom part shows another simplified version, with nodes arranged in two rows and arrows connecting them in a more regular pattern.



$$\cong \mathbf{Optic}_{\otimes, \otimes} \left( ((I, I), (I, I)), ((X, A), (Y, B)) \right)$$

where for the second isomorphism we used the symmetry of the monoidal product assumed on the  $\mathbb{R}$  category. The fourth isomorphism suggests that the symmetric monoidal product  $\otimes$  on  $\mathbf{Optic}_{\mathbb{R}, \mathbb{Q}}$  is a ‘quantum spider’ in the terminology of (author?) [7].

**Proposition 3.5.** Let  $\mathcal{C}$  and  $\mathcal{D}$  be the (monoidal) categories underlying  $\mathbf{Optic}_{\mathbb{R}, \mathbb{Q}}$ , and denote their respective monoidal units by  $I_{\mathcal{C}}$  and  $I_{\mathcal{D}}$ . If these unit objects are terminal in their respective categories, then the contexts  $C((X, A), (Y, B))$  simplify to



where we have depicted the representable presheaf on  $I_{\mathcal{D}}$  as  $\bar{\top}$  to indicate that  $A$  is just discarded. Consequently, in this case, a context is just an optic  $(I, B) \rightarrow (X, Y)$ .

*Proof of Proposition 3.5.* The proof that contexts simplify when the underlying monoidal units are terminal is a simple exercise in the graphical calculus: just throw away any hom-objects with terminal codomains, and simplify the resulting diagrams.  $\square$

**Definition 3.6** (Generalized open game). Let  $(X, A)$  and  $(Y, B)$  be objects in any symmetric monoidal category of optics  $\mathbf{Optic}_{\mathbb{R}, \mathbb{Q}}$ . Let  $\Sigma$  be a  $\mathcal{U}$ -category, for any base of enrichment  $\mathcal{U}$  such that  $\mathcal{U}$ -**Prof** is compact closed. An **open game** from  $(X, A)$  to  $(Y, B)$  with strategies in  $\Sigma$ , denoted  $G : (X, A) \xrightarrow{\Sigma} (Y, B)$ , is given by:

- (a) a play function  $P : \Sigma_0 \rightarrow \mathbf{Optic}_{\mathbb{R}, \mathbb{Q}}((X, A), (Y, B))$ ; and
- (b) a best response function  $B : C((X, A), (Y, B)) \rightarrow \mathcal{U}\text{-Prof}(\Sigma, \Sigma)$ .

Given a strategy  $\sigma : \Sigma$ , we will often write  $\langle v | u \rangle_{\sigma}$  or similar to denote its image under  $P$ . A strategy is an **equilibrium** in a context  $\langle \pi | k \rangle$  if it is a fixed point of  $B(\langle \pi | k \rangle)$ .



Roughly speaking, the ‘best responses’ to a strategy  $\sigma$  in a context is are those strategies  $\tau$  such that choosing  $\tau$  would result in performance at the game at least as good as choosing  $\sigma$ ; equilibrium strategies are those for which such deviation would not improve performance.

**Remark 3.7.** Note: whereas classic open games use a best-response relation, we categorify that here to a best-response *relator* (in the terminology of [17]; *i.e.*, a ‘proof-relevant’ relation), so that we can describe the trajectories witnessing the computation of equilibria, rather than their mere existence.

**Proposition 3.8.** Generalized open games over the symmetric monoidal category of optics  $\mathbf{Optic}_{\mathbb{R}, \mathbb{C}}$  with strategies enriched in  $\mathcal{U}$  form a symmetric monoidal category denoted  $\mathbf{Game}(\mathcal{U}, \mathbb{R}, \mathbb{C})$ .

*Proof of Proposition 3.8.* Since we assume  $\mathcal{U}$ -**Prof** to be compact closed, the proof is a straightforward generalization of that given by (author?) [2], replacing best-response relations by best-response *relators* [17] (*i.e.*, ‘proof-relevant’ relations). The play function of a sequentially composite game is given by the sequential composition of the corresponding optics, and the associated best response function is given by relators on the product of the strategies using a natural notion of ‘local context’. The symmetric monoidal structure is given similarly. We refer the reader to (author?) [2] for the details.  $\square$

Since our games are only a mild generalization of those of [2], we refer the reader to §3.10 of that paper for an idea of the proof of the foregoing proposition, which goes through analogously. The sequential composition of games is given by the sequential composition of optics, with the best response to the composite being the product of the best responses to the factors. Similarly, parallel composition is given by the monoidal product of optics, and the best response to the composite is again the product of the best responses to the factors.

We now consider some games over **BayesLens** that supply the building blocks of the archetypal cybernetic systems to be considered in §4. For now, we will take the strategies simply to be discrete categories (*i.e.*, sets), as in the standard formulation of open games. Consequently, we will take the codomain of the best response function to be  $\mathbf{Set}(\Sigma, \mathbf{Set}(\Sigma, 2))$ , for each strategy type  $\Sigma$ . We assume the ambient category of stochastic channels is semicartesian, so that the monoidal unit is the terminal object.

**Remark 3.9.** All the games we will consider henceforth will have play functions whose codomains restrict to the representable subcategory  $\mathbf{GrLens}_{\text{Stat}}$  of **BayesLens**; in this work, we do not use the extra generality afforded by **BayesLens**, except insofar as it grants us the use of string diagrams in **Prof**, which we find helpful for reasoning intuitively about these systems. The generality of optics *is* however used in the ‘game-theoretic’ games of [2], and in future work we hope to relate the cybernetic systems of this paper to the game-theoretic setting of that earlier work.

**Remark 3.10.** All the statistical games considered in this paper will be ‘atomic’ in the sense of [2]: in particular, the best response functions we consider will be constant, meaning that, in any context, the set of best strategies does not depend on the ‘current’ choice of strategy. Permitting such dependence will be important in future work, however, when we consider how cybernetic systems interact, and hence respond to each other.

**Example 3.11.** A Bayesian lens of the form  $(I, I) \leftrightarrow (X, X)$  is fully specified by a state  $\pi : I \rightarrow X$ . A context for such a lens is given by a lens  $\langle ! | k \rangle : (I, X) \leftrightarrow (X, X)$  where  $! : I \rightarrow I$  is the unique map and  $k : X \rightarrow X$  is any endochannel on  $X$ . A **maximum likelihood game** is any game whose play function has codomain in Bayesian lenses of this form  $(I, I) \leftrightarrow (X, X)$  for any  $X : \mathcal{C}$ , and whose best response function is isomorphic to

$$B(\langle ! | k \rangle) = \langle \rho | ! \rangle_{\sigma} \mapsto \left\{ \langle \pi | ! \rangle_{\tau} \mid \pi \in \arg \max_{\pi : I \rightarrow X} \mathbb{E}_{k \bullet \pi} [\pi] \right\}$$



where  $\mathbb{E}$  is the canonical expectation operator (*i.e.* algebra evaluation) associated to states in  $\mathcal{C}$ , and where we have written  $\langle \rho \mid ! \rangle_\sigma$  and  $\langle \pi \mid ! \rangle_\tau$  to denote the images of the strategies  $\sigma$  and  $\tau$  under the play function. Intuitively, then, the best response is given by the strategy that maximises the likelihood of the state obtained from the context  $k$ .

**Remark 3.12.** In what follows, we assume that the underlying category  $\mathcal{C}$  of stochastic channels *admits density functions*. Informally, a density function for a stochastic channel  $c : X \rightarrow Y$  is a measurable function  $p_c : Y \times X \rightarrow [0, 1]$  whose values are the probabilities (or probability densities)  $p_c(y|x)$  at each pair  $(y, x) : Y \times X$ . We say that the value  $p_c(y|x)$  is the probability (or probability density) of  $y$  *given*  $x$ . In a category such as  $\mathcal{Kl}(\mathcal{D}_{\leq 1})$ , whose objects are sets and whose morphisms  $X \rightarrow Y$  are functions  $X \rightarrow \mathcal{D}(Y + 1)$ , a density function for  $c : X \rightarrow Y$  is a morphism  $Y \otimes X \rightarrow I$ ; note that in  $\mathcal{Kl}(\mathcal{D}_{\leq 1})$ ,  $I$  is not terminal. In the finitely-supported case, density functions are effectively equivalent to channels, but this is not the case in the continuous setting, where they are of most use. For more on this, see [25, §2.1.4].

A natural first generalization of maximum likelihood games takes us from states  $I \rightarrow X$  to channels  $Z \rightarrow X$ ; that is, from ‘elements’ to ‘generalized elements’ in the covariant (forwards) part of the lens. Unlike Bayesian lenses  $(I, I) \rightarrow (X, X)$ , lenses  $(Z, Z) \rightarrow (X, X)$  admit nontrivial contravariant components, which we think of as generalized Bayesian inversions. Consequently, our first generalization is a notion of ‘Bayesian inference game’. A context  $\langle \pi \mid k \rangle : (I, X) \rightarrow (Z, X)$  for a Bayesian lens  $(Z, Z) \rightarrow (X, X)$  then constitutes a ‘prior’ state  $\pi : I \rightarrow Z$  and a ‘continuation’ channel  $k : X \rightarrow X$  which together witness the closure of the otherwise open system.

**Example 3.13.** Fix a channel  $c : Z \rightarrow X$  with associated density function  $p_c : X \times Z \rightarrow \mathbb{R}_+$  and a measure of divergence between states on  $Z$ ,  $D : \mathcal{C}(I, Z) \times \mathcal{C}(I, Z) \rightarrow \mathbb{R}$ . A corresponding (generalized) **simple Bayesian inference game** is any game whose play function has codomain  $\mathbf{BayesLens}((Z, Z), (X, X))$  and whose best response function is isomorphic to

$$\begin{aligned} B(\langle \pi \mid k \rangle) &= \langle d \mid d' \rangle_\sigma \mapsto \left\{ \langle c \mid c' \rangle_\tau \mid c' \in \arg \min_{c' : \mathbf{V}(\mathcal{C}(I, Z), \mathcal{C}(X, Z))} \mathbb{E}_{x \sim k \bullet c \bullet \pi} \left[ \mathbb{E}_{z \sim c'_\pi(x)} [-\log p_c(x|z)] + D(c'_\pi(x), \pi) \right] \right\} \\ &= \langle d \mid d' \rangle_\sigma \mapsto \left\{ \langle c \mid c' \rangle_\tau \mid c' \in \arg \min_{c' : \mathbf{V}(\mathcal{C}(I, Z), \mathcal{C}(X, Z))} \left( \mathbb{E}_{z \sim c'_\pi \bullet k \bullet c \bullet \pi} \left[ - \int_X \log p_c(dk \bullet c \bullet \pi|z) \right] \right. \right. \\ &\quad \left. \left. + D(c'_\pi \bullet k \bullet c \bullet \pi, \pi) \right) \right\} \end{aligned}$$

where  $\pi : I \rightarrow Z$  and  $k : X \rightarrow X$ , and where the notation  $z \sim \pi$  means “ $z$  distributed according to the state  $\pi$ ”. Note that the second line follows from the first by linearity of expectation.

**Proposition 3.14** ([16, Thm. 1]). When  $D$  is chosen to be the Kullback-Leibler divergence  $D_{KL}$ , minimizing the objective function defining a simple Bayesian inference game is equivalent to computing an (exact) Bayesian inversion.

**Corollary 3.15.** Given two Bayesian inference games  $G : (Z, Z) \rightarrow (Y, Y)$  and  $H : (Y, Y) \rightarrow (X, X)$ , we can compose them sequentially to obtain a game  $H \circ G : (Z, Z) \rightarrow (X, X)$ , which we will call a **hierarchical Bayesian inference game**. It is then an immediate consequence of Lemma 2.13 that, in any given context for which the forwards channels admit Bayesian inversion, the best response to the composite game  $H \circ G$  (that is, the optimal inversion of the composite channel) is given simply by (the composition of) the best responses to the factors  $H$  and  $G$ . Consequently, Bayesian inference games are closed under composition.

Similarly, given a channel  $c : Z \otimes Y \rightarrow X$ , we can consider the **marginal Bayesian inference game** in which the objective is to compute the inversion of the channel onto just one of the factors  $Z$  or  $Y$  in the domain.

**Example 3.16** (Variational autoencoder game). Fix a family  $\mathcal{F} \hookrightarrow \mathcal{C}(Z, X)$  of forward channels and a family  $\mathcal{P} \hookrightarrow \mathcal{C}(X, Z)$  of backward channels such that each  $c : \mathcal{F}$  admits a density function  $p_c : X \otimes Z \rightarrow \mathbb{R}_+$  and each  $d : \mathcal{P}$  admits a density function  $q : Z \otimes X \rightarrow \mathbb{R}_+$ ; think of these families as determining parameterizations of the channels. We take our strategy type to be  $\Sigma = \mathcal{F} \times \mathcal{P}$ . A **simple variational autoencoder game**  $(Z, Z) \xrightarrow{\Sigma} (X, X)$  is any game with play function  $P : \Sigma \rightarrow \mathbf{BayesLens}((Z, Z), (X, X))$  and whose best response function is isomorphic to

$$B(\langle \pi | k \rangle) = \langle d | d' \rangle_{\sigma} \mapsto \left\{ \langle c | c' \rangle_{\tau} \mid (c, c') \in \underset{\substack{c \in \mathcal{F}, \\ c' \in \mathbf{V}(\mathcal{C}(I, Z), \mathcal{P})}}{\operatorname{argmin}}} \mathbb{E}_{x \sim k \bullet c \bullet \pi} \mathbb{E}_{z \sim c'_\pi(x)} \left[ \log \frac{q(z|x)}{p_c(x|z)p_\pi(z)} \right] \right\}$$

where  $\pi : I \rightarrow Z$  admits a density function  $p_\pi : Z \rightarrow \mathbb{R}_+$ ,  $q : Z \otimes X \rightarrow \mathbb{R}_+$  is a density function associated to  $c'_\pi$ , and  $k$  has type  $X \rightarrow X$ .

**Proposition 3.17.** A best response to a variational autoencoder game is a stochastic channel  $c : \mathcal{F}$  that maximises the likelihood of the state observed through the continuation  $k$  under the assumption that the generative process is in  $\mathcal{F}$ , along with an inverse channel  $c'_\pi : \mathcal{P}$  that best approximates the exact Bayesian inverse  $c_\pi^\dagger$  under the constraint of being in  $\mathcal{P}$ .

*Proof of Proposition 3.17.* We seek to show that giving a best response to a simple variational autoencoder game is equivalent to maximizing the likelihood of the observed data under the assumption that the data is generated by the process  $k \bullet c \bullet \pi$ , with  $c$  constrained to lie in  $\mathcal{F}$  and with a (possibly approximate) inverse  $c'_\pi$  constrained to lie in  $\mathcal{P}$ . The intuition is that an agent responding to such a game builds a model of the process generating the data such that the agent can invert the model to infer the latent causes of the observed data. Typically, therefore, one chooses families  $\mathcal{F}$  and  $\mathcal{P}$  that are computationally tractable at the cost of some approximation.

As in the maximum likelihood game (Example 3.11), maximising the likelihood of data assumed to be generated by  $k \bullet c \bullet \pi$  means maximising  $\mathbb{E}_{k \bullet c \bullet \pi} [c \bullet \pi]$ . We assume that  $c$  is represented by the density function  $p_c : X \times Z \rightarrow \mathbb{R}_+$ ,  $\pi$  is represented by  $p_\pi : Z \rightarrow \mathbb{R}_+$ , and  $c'_\pi$  is represented by  $q : Z \times X \rightarrow \mathbb{R}_+$ . The density function for the composite state  $c \bullet \pi : I \rightarrow X$  is then given by  $p_{c \bullet \pi}(x) = \int_Z p_c(x|z)p_\pi(z)dz$ . Consequently, we have  $\mathbb{E}_{k \bullet c \bullet \pi} [c \bullet \pi] = \mathbb{E}_{k \bullet c \bullet \pi} [p_{c \bullet \pi}]$ . Moreover, since  $\log$  is monotonic, we can just as well minimize  $\mathbb{E}_{k \bullet c \bullet \pi} [-\log p_{c \bullet \pi}]$ .

Let  $p_\omega : Z \times X \rightarrow \mathbb{R}_+$  be given by  $p_\omega(z, x) = p_c(x|z)p_\pi(z)$ , and suppose the exact Bayesian inversion of  $c$  with respect to  $\pi$ , denoted  $c_\pi^\dagger$ , has density function  $p_{c_\pi^\dagger} : Z \times X \rightarrow \mathbb{R}_+$ .

We then have the following equalities:

$$\begin{aligned} -\log p_{c \bullet \pi}(x) &= \mathbb{E}_{z \sim c'_\pi(x)} [-\log p_{c \bullet \pi}(x)] \\ &= \mathbb{E}_{z \sim c'_\pi(x)} \left[ -\log \frac{p_\omega(z, x)}{p_{c_\pi^\dagger}(z|x)} \right] \\ &= \mathbb{E}_{z \sim c'_\pi(x)} \left[ -\log \frac{p_\omega(z, x)}{q(z|x)} \frac{q(z|x)}{p_{c_\pi^\dagger}(z|x)} \right] \\ &= - \mathbb{E}_{z \sim c'_\pi(x)} \left[ \log \frac{p_\omega(z, x)}{q(z|x)} \right] - D_{KL} [c'_\pi(x), c_\pi^\dagger(x)] \end{aligned}$$

so that

$$\mathbb{E}_{z \sim c'_\pi(x)} \left[ \log \frac{q(z|x)}{p_c(x|z)p_\pi(z)} \right] = D_{KL} [c'_\pi(x), c_\pi^\dagger(x)] - \log p_{c \bullet \pi}(x). \quad (8)$$

This quantity is called the *free energy*, from an analogy with statistical mechanics, or *evidence upper bound* [15], from machine learning. Minimizing the free energy entails both minimizing the divergence between the model posterior  $c'_\pi(x)$  and the exact Bayesian posterior  $c_\pi^\dagger(x)$  (given an observation  $x$ ), as well as maximizing the likelihood of the observation  $x$  under the assumed generative model  $c$ . If the exact posterior  $c_\pi^\dagger$  is in  $\mathcal{F}$ , then the divergence term has a minimum at 0 where the two posteriors are equal almost everywhere; in this case, the objective absolutely maximises the likelihood of the data under the model  $c$ .

In general, we are not just interested in a single observation  $x : X$ , but in a distribution (*i.e.* state) induced on  $X$  by some external process, represented here by the continuation  $k$ . Thus, taking the free energy under the expectation induced by  $k$  gives precisely the objective of the simple variational autoencoder game.  $\square$

**Proposition 3.18.** Variational autoencoder games generalize inference games for the Kullback-Leibler divergence. More precisely, the objective function defining autoencoder games is of the same form as that defining inference games (3.13) when  $D = D_{KL}$ .

*Proof of Proposition 3.18.* We start from the free energy (8):

$$\begin{aligned} \mathbb{E}_{z \sim c'_\pi(x)} \left[ \log \frac{q(z|x)}{p_c(x|z)p_\pi(z)} \right] &= \mathbb{E}_{z \sim c'_\pi(x)} [\log q(z|x) - \log p_c(x|z) - \log p_\pi(z)] \\ &= \mathbb{E}_{z \sim c'_\pi(x)} [-\log p_c(x|z)] + D_{KL} [c'_\pi(x), \pi]. \end{aligned}$$

Taking the expectation induced by the continuation  $k$  gives a generalization of the objective of a simple Bayesian inference game, with  $D$  taken to be the Kullback-Leibler divergence  $D_{KL}$ , and with  $c'$  constrained to lie in  $\mathcal{P}$  and  $c$  allowed to vary in  $\mathcal{F}$ .  $\square$

This prompts the following generalization:

**Example 3.19** (Generalized autoencoder game). Fix two families of channels  $\mathcal{F}, \mathcal{P}$  and a strategy type  $\Sigma$  as in Example 3.16. Then a (generalized) **simple autoencoder game**  $(Z, Z) \xrightarrow{\Sigma} (X, X)$  is any game with play function  $P : \Sigma \rightarrow \mathbf{BayesLens}((Z, Z), (X, X))$  and whose best response function is isomorphic to

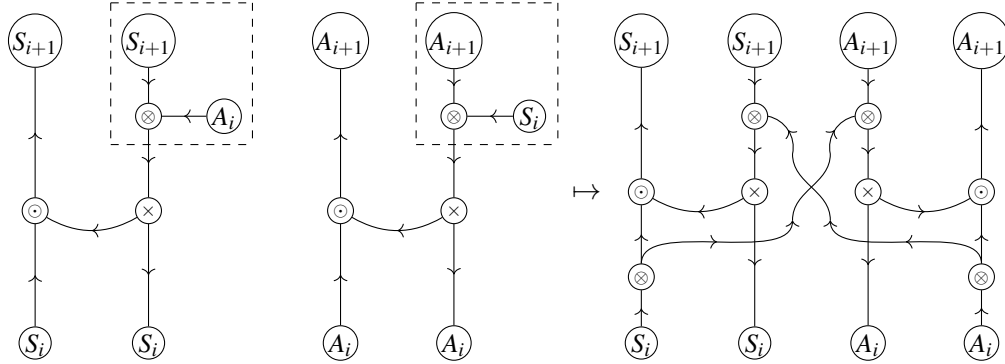
$$B(\langle \pi | k \rangle) = \langle d | d' \rangle_\sigma \mapsto \left\{ \langle c | c' \rangle_\tau \left| (c, c') \in \arg \min_{\substack{c \in \mathcal{F}, \\ c' \in \mathbf{V}(\mathcal{C}(I, Z), \mathcal{P})}} \left( \mathbb{E}_{z \sim c'_\pi \bullet k \bullet c \bullet \pi} \left[ - \int_X \log p_c(dk \bullet c \bullet \pi | z) \right] + D(c'_\pi \bullet k \bullet c \bullet \pi, \pi) \right) \right. \right\}$$

where  $\pi$  and  $k$  have respective types  $I \rightarrow Z$  and  $X \rightarrow X$ , and  $D$  is any measure of divergence between states.

As with Bayesian inference games, we can generalize simple autoencoder games to **hierarchical** and **marginal** autoencoder games via the corresponding sequential and parallel compositions.

The foregoing games have been purely statistically formulated, without capturing the motivating feature of an open system as something in interaction with an external environment. Nonetheless, we can model a simple open system of hierarchical **active inference** that receives stochastic inputs from an environment and emits actions stochastically into the environment, as follows.

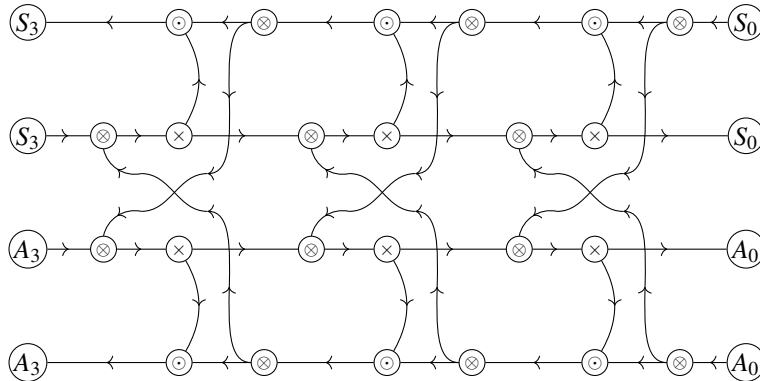
**Example 3.20** (Active inference game). Let  $\{S_i\}_i$  be set of spaces of sensory data indexed by hierarchical levels of abstraction  $i$  (for instance, the levels of abstraction might range from representations of whole objects to fine details about their texture); similarly, let  $\{A_i\}_i$  be a set of spaces of possible actions similarly hierarchically organized. Consider the marginal autoencoder games  $(S_{i+1} \otimes A_i, S_{i+1}) \rightarrow (S_{i+1}, S_{i+1})$  and  $(A_{i+1} \otimes S_i, A_{i+1}) \rightarrow (A_{i+1}, A_{i+1})$  coupled via the symmetric monoidal structure  $\otimes$  of  $\mathcal{C}$ :



giving a composite game  $(S_{i+1} \otimes A_{i+1}, S_{i+1} \times A_{i+1}) \rightarrow (S_i \otimes A_i, S_i \times A_i)$ . Recall from [2, §§3.7-3.8] that a composite game is given by the (sequential and parallel) composition of optics, with best-response given by the product of the best-responses of the factors.

Note that the Bayesian posterior inferred by such a game has independent factors on  $S_{i+1}$  and  $A_{i+1}$ . This is not merely a diagrammatic convenience, but coincides with a common ‘mean field’ simplification in the modelling literature [3, 15]. The dashed box is a functorial box [18] depicting the Yoneda embedding; recall that optics in **BayesLens** were defined over (co)presheaves, and so here we needed to lift the monoidal product on  $\mathcal{C}$  into a diagram over its presheaf category  $\mathbf{Cat}(\mathcal{C}^{\text{op}}, \mathbf{Set})$ .

Next, compose these games along the hierarchy indexed by  $i$ , to obtain a game  $(S_N \otimes A_N, S_N \times A_N) \rightarrow (S_0 \otimes A_0, S_0 \times A_0)$ , such as an element of the following object:



Given a context with a strong prior about expected sensory states and a continuation that responds to an action of type  $A_0$  by feeding back a state on  $S_0$ , the best response can be shown to be that which selects actions that, under the current state, maximize the likelihood of obtaining the expected ‘goal’ state [3, 11].

**Remark 3.21.** We have framed each of these statistical procedures as optimization problems not only to suggest a link to the utility-maximising agents of game theory, but also because it suggests the use of iterative methods to compute best responses; note that computational tractability is an important motivation in the proof of Proposition 3.17.

The question of providing such dynamical or, thinking of game composition as an algebra for building complex systems, ‘coalgebraic’ semantics for (generalized) optimization games is the topic of the next section. We first formalize this notion.

**Definition 3.22.** An **optimization game** is any open game whose best response function can be defined by a function of the form  $\Sigma \times C \xrightarrow{\pi} M \xrightarrow{\varphi} P$ , where  $\Sigma$  is a strategy type,  $C$  a context type,  $M$  any space, and  $P$  a poset. We call  $\varphi$  the **fitness function**, and think of  $\pi$  as projecting systems into a space whose points can be assigned a fitness. The best response function of an optimization game can then be defined by giving the subset of strategies contextually maximizing fitness, for each context  $c : C$ .

## 4 Cybernetic Systems and Dynamical Realisation

In this section, we begin to answer the question of precisely how the optimization games of the previous section may be realized in physical systems, such as brains or computers. More formally, this means we seek open dynamical systems whose input and output types correspond to the domain and codomain types of the foregoing games, such that there is a correspondence between the behaviours of the abstract games and their dynamical realisations, and such that the evolutions of the internal states of the dynamical systems correspond to strategic improvements in game-playing: by concentrating on optimization games, a natural measure of such improvement is encoded in the fitness function underlying the best-response relator.

We do not require that there is a correspondence between internal states of the realisations and strategies for the corresponding games, but we do require that the fitness functions extend to the total state spaces of the closure of a realisation induced by the context. When there *is* a correspondence between internal states and strategies, we can take advantage of Definition 3.6 and interpret trajectories over the state space as trajectories over strategies witnessing the strategic improvement.

We begin by sketching categories of dynamical games, and then use these ideas to define preliminary notions of open cybernetic systems and categories thereof. We consider principally single systems whose underlying games are atomic (in the sense of Remark 3.10), and leave the study of the behaviour of interacting cybernetic systems to future work. Once more, we omit proofs in this section; they will appear in a paper to follow.

**Definition 4.1** (Discrete-time dynamical system over  $\mathcal{C}$ ; after [24, 6]). A **discrete-time dynamical system** over  $\mathcal{C}$  with state space  $S : \mathcal{C}$ , input type  $A : \mathcal{C}$  and output type  $B : \mathcal{C}$  is a lens  $(S, S) \rightarrow (B, A)$  over  $\mathcal{C}$ , *i.e.* in the following optical hom object:

$$\int^{M:\mathcal{C}} \mathbf{Comon}(\mathcal{C})(S, M \otimes B) \times \mathcal{C}(M \otimes A, S) \cong \mathbf{Comon}(\mathcal{C})(S, B) \times \mathcal{C}(S \otimes A, S)$$

where the isomorphism follows by Yoneda reduction. Note that this requires that the ‘output’ map of the dynamical system is a comonoid homomorphism in  $\mathcal{C}$  and hence deterministic in a category of stochastic channels.

**Definition 4.2** (Category of discrete-time dynamical systems). We define a category  $\mathbf{Dyn}_{\mathcal{C}}$  whose objects are the objects of  $\mathcal{C}$  and whose morphisms, denoted  $A \xrightarrow{S} B$ , are discrete-time dynamical systems; the symbol above the arrow denotes the internal state space. Hom objects are given by

$$\mathbf{Dyn}_{\mathcal{C}}(A, B) = \sum_{S:\mathcal{C}} \mathbf{Comon}(\mathcal{C})(S, B) \times \mathcal{C}(S \otimes A, S).$$

Identity dynamical systems on each  $A : \mathcal{C}$  are the ‘no-op’ dynamical systems  $A \xrightarrow{A} A$  given by identity optics  $\text{id}_A : (A,A) \rightarrow (A,A)$ . Associativity and unitality of composition is inherited from the category of optics underlying Definition 4.1; a symmetric monoidal structure is similarly inherited.  $\square$

**Definition 4.3** (Lenses over dynamical systems; after [21]). The category of (monoidal) lenses over  $\mathcal{C}$ -dynamical systems has as objects pairs  $(X,A)$  of objects in  $\mathcal{C}$  and as morphisms, **dynamical lenses**  $(X,A) \rightarrow (Y,B)$ , elements of the type

$$\begin{array}{c}
 \begin{array}{c}
 \textcircled{X} \rightarrow \textcircled{\phantom{X}} \rightarrow \textcircled{Y} \\
 \textcircled{A} \leftarrow \textcircled{\phantom{A}} \leftarrow \textcircled{B}
 \end{array} \\
 \int^{M:\mathcal{C}} \mathbf{Dyn}_{\mathcal{C}}(X, M \otimes Y) \times \mathbf{Dyn}_{\mathcal{C}}(M \otimes B, A) \\
 \cong \\
 \sum_{P,Q:\mathcal{C}} \int^{M:\mathcal{C}} \mathcal{C}(P \otimes X, P) \times \mathbf{Comon}(\mathcal{C})(P, M \otimes Y) \times \mathcal{C}(Q \otimes M \otimes B, Q) \times \mathbf{Comon}(\mathcal{C})(Q, A) \\
 \sum_{P,Q:\mathcal{C}} \begin{array}{c}
 \textcircled{P} \\
 \downarrow \\
 \textcircled{X} \rightarrow \textcircled{\phantom{X}} \rightarrow \textcircled{P} \quad \textcircled{P} \rightarrow \textcircled{\phantom{P}} \rightarrow \textcircled{Y} \\
 \textcircled{A} \leftarrow \textcircled{\phantom{A}} \leftarrow \textcircled{Q} \quad \textcircled{Q} \leftarrow \textcircled{\phantom{Q}} \leftarrow \textcircled{B} \\
 \textcircled{Q} \uparrow
 \end{array}
 \end{array}$$

That is, a dynamical lens is a pair of dynamical systems coupled along some ‘residual’ type.

**Remark 4.4.** At this point we begin to run into sizes issues. However, for the purposes of this paper, we will simply assume that a satisfactory resolution of these matters is at hand; for instance, that there is a hierarchy of Grothendieck universes such that the coends over (‘large’) sums in the preceding definition constitute accessible objects.

We now expand the definition of context in the dynamical setting. We will see that a dynamical context is simply a closure of an open dynamical system: that is, a ‘larger’ system into which a ‘smaller’ open dynamical system can plug such that the composite is a closed (but still uninitialized) system.

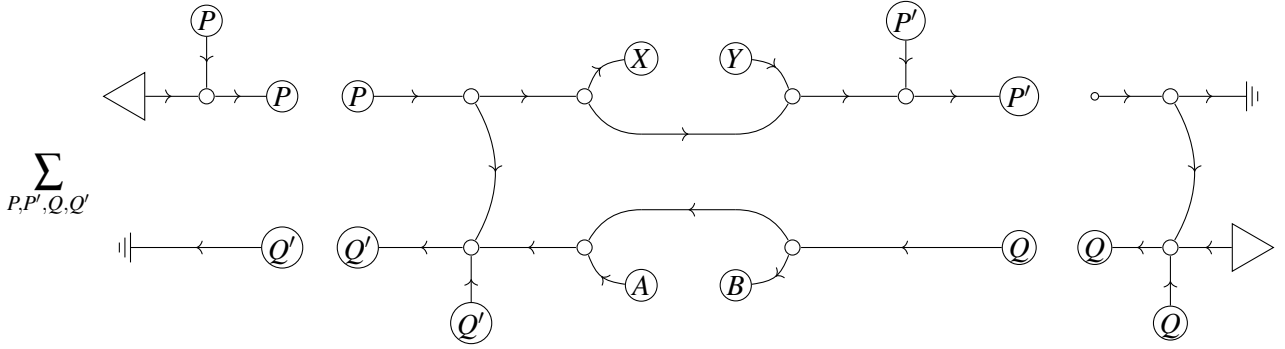
**Proposition 4.5.** If  $I$  is terminal in  $\mathcal{C}$ , a context for a dynamical lens  $(X,A) \rightarrow (Y,B)$  is an element of the following type, denoted  $\tilde{\mathcal{C}}((X,A), (Y,B))$ :

$$\sum_{P,Q:\mathcal{C}} \textcircled{P} \rightarrow \textcircled{P} \quad \textcircled{P} \rightarrow \textcircled{\phantom{P}} \rightarrow \textcircled{X} \quad \textcircled{Y} \rightarrow \textcircled{\phantom{Y}} \rightarrow \textcircled{Q} \quad \textcircled{Q} \rightarrow \textcircled{B} \quad \textcircled{A} \rightarrow \textcircled{\phantom{A}}$$

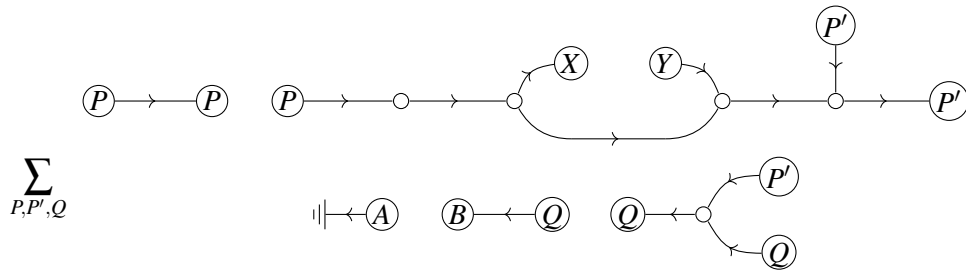
Interpreting this diagram, a context for a dynamical lens  $(X,A) \rightarrow (Y,B)$  amounts to an autonomous dynamical system with output type of the form  $X \otimes M$  (for some residual type  $M$ ), coupled along the

residual  $M$  to an open dynamical system with input type  $Y \otimes M$  and output type  $B$ ; and the  $A$  type is discarded. This is precisely what we should expect from a dynamical analogue of Proposition 3.5.

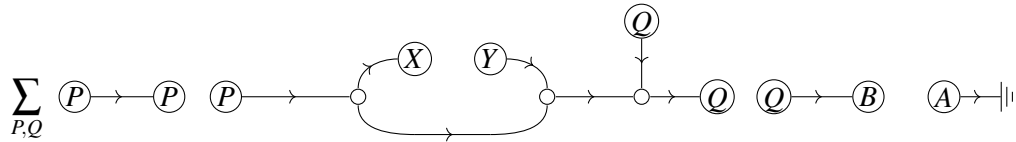
*Proof of Proposition 4.5.* Writing it out in full, we find that the context  $\tilde{C}((X,A),(Y,B))$  is an element of the following type:



By discarding the factors that output into the terminal object, the context reduces to



Since  $P'$  is an arbitrary choice of internal state space type, we can without loss of generality set  $P' = Q$ , which gives us



as required. □

**Definition 4.6.** A **dynamical game** is just a generalized open game (3.6) over the category of dynamical lenses. We write  $(X,A) \xrightarrow{\tilde{\Sigma}, S} (Y,B)$  to indicate both the strategy type  $\tilde{\Sigma}$  and state space  $S$ . Dynamical games form a symmetric monoidal category in the corresponding way. For notational clarity, we will write  $\tilde{G}$  for a dynamical game,  $\tilde{P}$  for its play function, and  $\tilde{B}$  for its best response function.

**Definition 4.7** (Dynamical realisation of an open game). Let  $G : (X,A) \xrightarrow{\Sigma} (Y,B)$  be an open game with  $X,A,Y,B$  all objects of some symmetric monoidal category  $\mathcal{C}$ . A **dynamical realisation** of  $G$  is a choice of dynamical game  $\tilde{G} : (X,A) \xrightarrow{\tilde{\Sigma}, S} (Y,B)$  on the same objects, along with a function  $[[\cdot]] : C((X,A),(Y,B)) \rightarrow \tilde{C}((X,A),(Y,B))$  lifting static contexts to dynamical contexts. Given a context  $\langle \pi | k \rangle : C((X,A),(Y,B))$ , we choose a representative  $\langle [[\pi]] | [[k]] \rangle \cong [[\langle \pi | k \rangle]] : \tilde{C}((X,A),(Y,B))$  for its realisation.



A ‘dynamical context’ is an element of the type given in Proposition 4.5: a context for a dynamical lens. A ‘static context’ is simply a context for the ‘static’ game that is being dynamically realized. At this stage, we impose no particular requirements on the context realisation function  $\llbracket \cdot \rrbracket$ , except to say that in the intended semantics,  $\llbracket \langle \pi | k \rangle \rrbracket$  is a (coupled, open) dynamical system that constantly emits the state  $\pi$  and (by some mechanism) realizes the channel  $k$ . We call such a context *stationary* as neither  $\pi$  nor  $k$  vary in time; future work will generalize the results of this section to *non-stationary* contexts.

**Definition 4.8** (Open cybernetic systems). An open **cybernetic system** is defined by the data:

- an open optimization game (Def. 3.22)  $G : (X, A) \xrightarrow{\Sigma} (Y, B)$  with  $X, A, Y, B$  all objects of some symmetric monoidal category  $\mathcal{C}$ ,
- a fitness function  $\varphi_G : \Sigma \times C \rightarrow M \xrightarrow{\varphi} F$  where  $C = C((X, A), (Y, B))$ ,
- a dynamical realisation  $(\tilde{G} : (X, A) \xrightarrow{\tilde{\Sigma}, \tilde{S}} (Y, B), \llbracket \cdot \rrbracket : C((X, A), (Y, B)) \rightarrow \tilde{C}((X, A), (Y, B)))$  of  $G$ ,

satisfying the following condition for each context  $\langle \pi | k \rangle : C((X, A), (Y, B))$ :

- there exists a dynamical strategy  $\tilde{\sigma} : \tilde{\Sigma}$ , such that
- writing  $Z$  for the total state space of the autonomous dynamical system  $\llbracket \langle \pi | k \rangle \rrbracket \circ \tilde{P}(\tilde{\sigma})$  induced by the context, there exists a function  $\nu : Z \rightarrow M$  projecting  $Z$  into the ‘fitness landscape’  $M$ , such that
- there exists a fitness-maximising fixed point  $\zeta^* : Z$ , in the sense that
- for some equilibrium strategy of the static system  $\sigma^* : \text{fix } B(\langle \pi | k \rangle)$ ,  $\varphi(\nu(\zeta^*)) \leq \varphi_G(\sigma^*, \langle \pi | k \rangle)$ .

A **category of open cybernetic systems** is a category of (generalized) open games such that each game is an open cybernetic system with dynamics realised in the same category  $\mathcal{C}$ , and such that the composite of games is a cybernetic system whose fitness-maximising fixed point projects onto fitness-maximising fixed points of each of the factors in their corresponding local contexts. (See [2, §3.7] for the definition of local context.)

The idea here is that, by using the fitness function of the underlying optimization game, the cybernetic condition forces the behaviour of the dynamical realisation to coincide with the process of iteratively improving the strategies deployed by the system in playing the game. We summarize the condition in the diagram

$$\begin{array}{ccccc} \Sigma \times C & \longrightarrow & M & \xrightarrow{\varphi} & F \\ \llbracket \cdot \rrbracket \downarrow & & \uparrow & & \\ \tilde{\Sigma} \times \tilde{C} & \xrightarrow{\text{fix}} & Z & & \end{array}$$

though this is in general ill-defined: we do not require a function  $\llbracket \cdot \rrbracket : \Sigma \rightarrow \tilde{\Sigma}$ , and nor do we require that the best response to  $\tilde{G}$  coincides in any way with the best response to  $G$ . Investigating such conditions is the subject of future work; for instance, we may be interested in nested cybernetic systems, such as characterize evolution by natural selection, and how their fitness functions constrain one another. For similar reasons, we are also interested in the case where the fitness function is itself non-stationary.

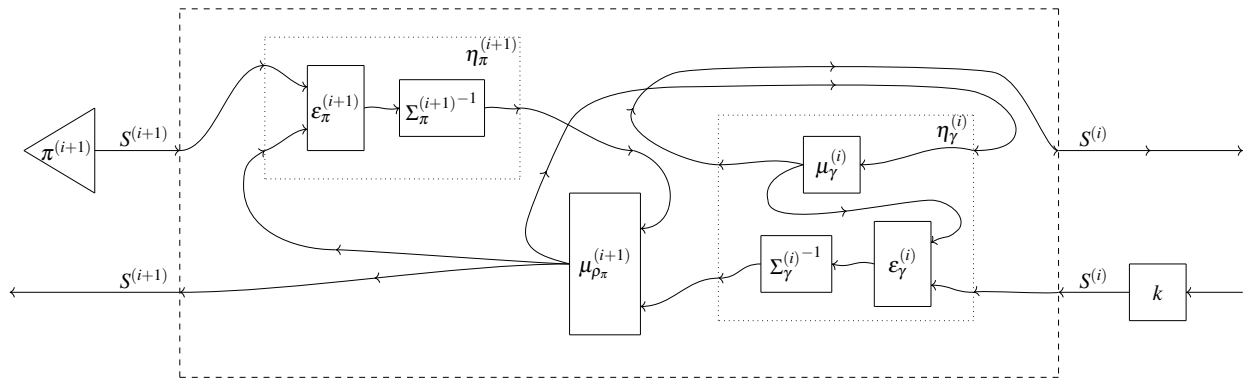
**Remark 4.9.** The codomain category of the cybernetic realisation functor is in general much larger than the domain category of static games, and often it makes sense to consider dynamical games in this codomain category as if they were dynamical realisations of static games, even if in fact there is no static game to which they could correspond. For instance, adaptive systems in physical environments are in general not realisations of static games because their contexts are irreducibly dynamical and thus not

the dynamical realisation of a static context; but over short time intervals, it can be productive to treat such systems as realisations of static games. In continuous time (not treated here), it is even possible to consider dynamical games that are indeed realisations of games that are static when represented in a smoothly varying coordinate system. The free-energy framework of Theorem 4.10 is an example of a category of cybernetic systems with a rich underlying category of dynamic games.

A classic category of open cybernetic systems is found in the computational neuroscience literature, as summarized in the following theorem.

**Theorem 4.10.** Consider the subcategory of **BayesLens** spanned by finite-dimensional Euclidean spaces, with morphisms generated (under sequential and parallel composition) by the (variational) autoencoder and inference games whose forwards and backwards channels emit Gaussian measures with high-precision. The (discrete-time) free-energy framework for action and perception [3] instantiates a category of open cybernetic systems realising games over this subcategory.

*Proof of Theorem 4.10.* By generalizing the exposition given by (author?) [3] and rendering it compositional, we show that every such game is realised by a dynamical system of the following form, depicted as a wiring diagram [24], and that the dynamics has a fixed point corresponding to a best-response to the underlying game:



While it is possible to give a dynamical realisation by a system of the form in Definition 4.3, it is more parsimonious to allow the ‘forwards’ view system to receive inputs from the ‘backwards’ update system, as well as vice versa. In a sense, we want the coupling along the residual to be bidirectional. We can depict this as the following type:

$$\begin{array}{c}
 \textcircled{X} \rightarrow \textcircled{\phantom{X}} \rightarrow \textcircled{Y} \\
 \textcircled{A} \leftarrow \textcircled{\phantom{A}} \leftarrow \textcircled{B}
 \end{array}
 \cong
 \int^{M,N:\mathcal{C}} \mathbf{Dyn}_{\mathcal{C}}(M \otimes X, N \otimes Y) \times \mathbf{Dyn}_{\mathcal{C}}(N \otimes B, M \otimes A)$$

We can choose residuals of type  $N = A$  and  $M = B$ , so that representatives are of the type  $\mathbf{Dyn}_{\mathcal{C}}(X \otimes B, Y \otimes A) \times \mathbf{Dyn}_{\mathcal{C}}(A \otimes B, A \otimes B)$ . We can then choose the second factor of the representative to be the identity system, leaving us only needing to construct a system of type  $\mathbf{Dyn}_{\mathcal{C}}(X \otimes B, Y \otimes A)$ , for some  $X, A, Y, B$ . We call such a realisation *bidirectionally coupled*. Reconciling this type with the type given in Definition 4.3 is the subject of ongoing work. Note that we retain the form for the dynamical context given in Proposition 4.5.

We will work in the ‘convenient’ setting of Example 2.2, in which  $\mathcal{C} = \mathcal{Kl}(\mathcal{P})$  for some probability monad  $\mathcal{P} : \mathbf{Meas} \rightarrow \mathbf{Meas}$  on a Cartesian closed category of measurable spaces, such as the category of quasi-Borel spaces [14]. In this setting, a Bayesian lens  $(X, A) \leftrightarrow (Y, B)$  is given by a pair of maps of type  $\mathbf{Meas}(X, \mathcal{P}Y) \times \mathbf{Meas}(\mathcal{P}X \times B, \mathcal{P}A)$ . A bidirectionally coupled dynamical realisation is then of the type  $\mathbf{Dyn}_{\mathcal{Kl}(\mathcal{P})}(X \otimes B, Y \otimes A)$ , or equivalently

$$\sum_{Z: \mathcal{Kl}(\mathcal{P})} \mathcal{Kl}(\mathcal{P})(Z \otimes X \otimes B, Z) \times \mathbf{Comon}(\mathcal{Kl}(\mathcal{P}))(Z, Y \times A).$$

A corresponding dynamical context has the type:

$$\sum_{Z_\pi, Z_k: \mathcal{Kl}(\mathcal{P})} \int^{M: \mathcal{Kl}(\mathcal{P})} \mathcal{Kl}(\mathcal{P})(Z_\pi, Z_\pi) \times \mathbf{Comon}(\mathcal{Kl}(\mathcal{P}))(Z_\pi, M \times X) \\ \times \mathcal{Kl}(\mathcal{P})(M \otimes Z_k \otimes Y, Z_k) \times \mathbf{Comon}(\mathcal{Kl}(\mathcal{P}))(Z_k, B)$$

We will consider first simple variational autoencoder games (as in Example 3.16), of type  $\langle \gamma^{(i)} \mid \rho^{(i+1)} \rangle : (S^{(i+1)}, S^{(i+1)}) \leftrightarrow (S^{(i)}, S^{(i)})$ , in which the spaces  $S^{(i)}$  are Euclidean, and we restrict to Gaussian measures over them. Let  $\pi^{(i+1)} : \mathcal{P}S^{(i+1)}$  be a prior on  $S^{(i+1)}$ . We will write  $p_\gamma^{(i)} : S^{(i)} \times S^{(i+1)} \rightarrow [0, 1]$ ,  $p_{\rho_\pi}^{(i+1)} : S^{(i+1)} \times S^{(i)} \rightarrow [0, 1]$ , and  $p_\pi^{(i+1)} : S^{(i+1)} \rightarrow [0, 1]$  for the corresponding density functions. The objective of such a variational autoencoder game is then

$$\arg \min_{\substack{\gamma^{(i)} \in \Gamma^{(i)}, \\ \rho^{(i+1)} \in \mathbf{Meas}(\mathcal{P}S^{(i+1)}, \mathbf{P}^{(i+1)})}} \mathbb{E}_{s^{(i)} \sim k \bullet \gamma^{(i)} \bullet \pi^{(i+1)}} \mathbb{E}_{s^{(i+1)} \sim \rho_\pi(s^{(i)})} \left[ \log \frac{p_{\rho_\pi}(s^{(i+1)} | s^{(i)})}{p_\gamma(s^{(i)} | s^{(i+1)}) p_\pi(s^{(i+1)})} \right]$$

where  $\Gamma^{(i)} \subseteq \mathcal{Kl}(\mathcal{P})(S^{(i+1)}, S^{(i)})$  and  $\mathbf{P}^{(i+1)} \subseteq \mathcal{Kl}(\mathcal{P})(S^{(i)}, S^{(i+1)})$  are the subtypes of stochastic channels emitting the corresponding Gaussian measures.

Note that these types are families of *parameterized* stochastic channels, in the sense that each  $\gamma^{(i)}(s^{(i+1)})$ ,  $\rho^{(i+1)}(s^{(i)})$  and  $\pi^{(i+1)}$  can be completely specified by parameter vectors  $\vartheta_\gamma^{(i)}(s^{(i+1)}) : \Theta_\gamma^{(i)}$ ,  $\vartheta_{\rho_\pi}^{(i+1)}(s^{(i)}) : \Theta_{\rho_\pi}^{(i+1)}$ , and  $\vartheta_\pi^{(i+1)} : \Theta_\pi^{(i+1)}$ , where we have denoted the parameter spaces by  $\Theta$ . Consequently, the channels  $\gamma^{(i)}$  and  $\rho_\pi^{(i+1)}$  themselves are completely specified by functions  $\vartheta_\gamma^{(i)} : \mathbf{Meas}(S^{(i+1)}, \Theta_\gamma^{(i)})$  and  $\vartheta_{\rho_\pi}^{(i+1)} : \mathbf{Meas}(S^{(i)}, \Theta_{\rho_\pi}^{(i+1)})$ . Indeed, we can factor the channels as  $\gamma^{(i)} : S^{(i+1)} \xrightarrow{\vartheta} \Theta_\gamma^{(i)} \hookrightarrow \mathcal{P}S^{(i)}$  and  $\rho_\pi^{(i+1)} : S^{(i)} \xrightarrow{\vartheta} \Theta_{\rho_\pi}^{(i+1)} \hookrightarrow \mathcal{P}S^{(i+1)}$ , where the inclusions map a parameter vector onto the corresponding measure.

Since each measure is Gaussian, the relevant parameters are means and covariances, which we denote by  $\mu$  and  $\Sigma$ , so that  $\vartheta_\gamma^{(i)}(s^{(i+1)}) = (\mu_\gamma^{(i)}, \Sigma_\gamma^{(i)})(s^{(i+1)}) : (S^{(i)} \otimes \mathbf{FVect}_\mathbb{R}(S^{(i)}, S^{(i)})) = \Theta_\gamma^{(i)}$ , and similarly for  $\rho_\pi$  and  $\pi$ . Write  $|S^{(i)}|$  to denote the dimension of the space  $|S^{(i)}|$ , and  $\langle \cdot, \cdot \rangle$  for the inner product of vectors. We can write down the (log) density functions explicitly as

$$\log p_\gamma^{(i)}(s^{(i)} | s^{(i+1)}) = \frac{1}{2} \left\langle \varepsilon_\gamma^{(i)}, \Sigma_\gamma^{(i)-1} \varepsilon_\gamma^{(i)} \right\rangle - \log \sqrt{(2\pi)^{|S^{(i)}|} \det \Sigma_\gamma^{(i)}} \\ \log p_{\rho_\pi}^{(i+1)}(s^{(i+1)} | s^{(i)}) = \frac{1}{2} \left\langle \varepsilon_{\rho_\pi}^{(i+1)}, \Sigma_{\rho_\pi}^{(i+1)-1} \varepsilon_{\rho_\pi}^{(i+1)} \right\rangle - \log \sqrt{(2\pi)^{|S^{(i+1)}|} \det \Sigma_{\rho_\pi}^{(i+1)}} \\ \log p_\pi^{(i+1)}(s^{(i+1)}) = \frac{1}{2} \left\langle \varepsilon_\pi^{(i+1)}, \Sigma_\pi^{(i+1)-1} \varepsilon_\pi^{(i+1)} \right\rangle - \log \sqrt{(2\pi)^{|S^{(i+1)}|} \det \Sigma_\pi^{(i+1)}}$$

where for clarity we have omitted the dependence of  $\Sigma_\gamma^{(i)}$  on  $s^{(i+1)}$  and  $\Sigma_{\rho_\pi}^{(i+1)}$  on  $s^{(i)}$ , and where

$$\begin{aligned}\varepsilon_\gamma^{(i)} &:= s^{(i)} - \mu_\gamma^{(i)}(s^{(i+1)}) \\ \varepsilon_{\rho_\pi}^{(i+1)} &:= s^{(i+1)} - \mu_{\rho_\pi}^{(i+1)}(s^{(i)}) \\ \varepsilon_\pi^{(i+1)} &:= s^{(i+1)} - \mu_\pi^{(i+1)}.\end{aligned}$$

We now construct a dynamical system performing gradient descent on the free energy

$$\begin{aligned}\varphi^{(i)} : \mathcal{S}^{(i)} \rightarrow \mathbb{R}_+ &:= s^{(i)} \mapsto \mathbb{E}_{s^{(i+1)} \sim \rho_\pi(s^{(i)})} \left[ \log \frac{p_{\rho_\pi}(s^{(i+1)} | s^{(i)})}{p_\gamma(s^{(i)} | s^{(i+1)}) p_\pi(s^{(i+1)})} \right] \\ &= \mathbf{H}[\rho_\pi^{(i+1)}(s^{(i)})] - \mathbb{E}_{s^{(i+1)} \sim \rho_\pi(s^{(i)})} \left[ \log p_\gamma(s^{(i)} | s^{(i+1)}) + \log p_\pi(s^{(i+1)}) \right] \\ &= \mathbf{H}[\rho_\pi^{(i+1)}(s^{(i)})] + \mathbb{E}_{s^{(i+1)} \sim \rho_\pi(s^{(i)})} \left[ E^{(i)}(s^{(i+1)}, s^{(i)}) \right].\end{aligned}$$

where  $\mathbf{H}[\cdot]$  is the Shannon entropy functor and  $E^{(i)}(s^{(i+1)}, s^{(i)})$  is called the *energy* of the system. Note that this defines a fitness function  $\varphi^{(i)}$  for the variational autoencoder game, in the sense of Definitions 3.22 and 4.8, by

$$\varphi_G : \Gamma \times \mathbf{P} \times \mathcal{C} \rightarrow \mathbb{R}_+ := (\gamma, \rho, \langle \pi | k \rangle) \mapsto \mathbb{E}_{k \bullet \gamma \bullet \pi} \left[ \varphi^{(i)} \right].$$

Note also that under the assumption that the posterior  $\rho_\pi^{(i+1)}(s^{(i)})$  has small variance, we can approximate the expected energy by its second-order Taylor expansion around the mean  $\mu_{\rho_\pi}^{(i+1)}(s^{(i)})$ , so that

$$\begin{aligned}\varphi^{(i)}(s^{(i)}) &\approx E^{(i)}(\mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)}) + \frac{1}{2} \left\langle \varepsilon_{\rho_\pi}^{(i+1)}, \left( \partial_{s^{(i+1)}}^2 E^{(i)} \right) \left( \mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)} \right) \varepsilon_{\rho_\pi}^{(i+1)} \right\rangle \\ &\quad + \mathbf{H}[\rho_\pi^{(i+1)}(s^{(i)})].\end{aligned}$$

where  $\left( \partial_{s^{(i+1)}}^2 E^{(i)} \right) \left( \mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)} \right)$  is the Hessian of  $E^{(i)}$  with respect to  $s^{(i+1)}$  evaluated at  $(\mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)})$ . Note that

$$\left\langle \varepsilon_{\rho_\pi}^{(i+1)}, \left( \partial_{s^{(i+1)}}^2 E^{(i)} \right) \left( \mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)} \right) \varepsilon_{\rho_\pi}^{(i+1)} \right\rangle = \text{tr} \left[ \left( \partial_{s^{(i+1)}}^2 E^{(i)} \right) \left( \mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)} \right) \Sigma_{\rho_\pi}^{(i+1)}(s^{(i)}) \right],$$

that the entropy of a Gaussian measure depends only on its covariance,

$$\mathbf{H}[\rho_\pi^{(i+1)}(s^{(i)})] = \frac{1}{2} \log \det \left( 2\pi e \Sigma_{\rho_\pi}^{(i+1)}(s^{(i)}) \right),$$

and that the energy  $E^{(i)}(\mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)})$  does not depend on  $\Sigma_{\rho_\pi}^{(i+1)}(s^{(i)})$ . We can therefore write down directly the covariance  $\Sigma_{\rho_\pi}^{(i+1)*}(s^{(i)})$  minimizing  $\varphi^{(i)}(s^{(i)})$  as a function of  $s^{(i)}$ . We have

$$\partial_{\Sigma_{\rho_\pi}^{(i+1)}} \varphi^{(i)}(s^{(i)}) \approx \frac{1}{2} \left( \partial_{s^{(i+1)}}^2 E \right) \left( \mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)} \right) + \frac{1}{2} \Sigma_{\rho_\pi}^{(i+1)-1}.$$

Setting  $\partial_{\Sigma_{\rho_\pi}^{(i+1)}} \varphi^{(i)}(s^{(i)}) = 0$ , we find the optimum

$$\Sigma_{\rho_\pi}^{(i+1)*}(s^{(i)}) = \left( \partial_{s^{(i+1)}}^2 E \right) \left( \mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)} \right)^{-1}. \quad (9)$$

Next, we optimize the mean parameter  $\mu_{\rho_\pi}^{(i+1)}$ . It suffices to consider only the energy term  $E^{(i)}(\mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)})$ . We have

$$\begin{aligned} E^{(i)}(\mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)}) &= -\log p_\gamma(s^{(i)}|s^{(i+1)}) - \log p_\pi(s^{(i+1)}) \\ &= -\frac{1}{2} \left\langle \boldsymbol{\varepsilon}_\gamma^{(i)}, \boldsymbol{\Sigma}_\gamma^{(i-1)} \boldsymbol{\varepsilon}_\gamma^{(i)} \right\rangle - \frac{1}{2} \left\langle \boldsymbol{\varepsilon}_\pi^{(i+1)}, \boldsymbol{\Sigma}_\pi^{(i+1)-1} \boldsymbol{\varepsilon}_\pi^{(i+1)} \right\rangle \\ &\quad + \log \sqrt{(2\pi)^{|S^{(i)}|} \det \boldsymbol{\Sigma}_\gamma^{(i)}} + \log \sqrt{(2\pi)^{|S^{(i+1)}|} \det \boldsymbol{\Sigma}_\pi^{(i+1)}} \end{aligned}$$

and a simple computation shows that

$$\partial_{\mu_{\rho_\pi}^{(i+1)}} E^{(i)}(\mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)}) = - \left( \partial_{s^{(i+1)}} \boldsymbol{\mu}_\gamma^{(i)} \right) (\boldsymbol{\mu}_{\rho_\pi}^{(i+1)})^T \boldsymbol{\Sigma}_\gamma^{(i)-1} \boldsymbol{\varepsilon}_\gamma^{(i)} + \boldsymbol{\Sigma}_\pi^{(i+1)-1} \boldsymbol{\varepsilon}_\pi^{(i+1)}.$$

Define the precision-weighted error functions

$$\begin{aligned} \boldsymbol{\eta}_\gamma^{(i)} &: S^{(i)} \times S^{(i)} \xrightarrow{\boldsymbol{\varepsilon}_\gamma^{(i)}} S^{(i)} \xrightarrow{\boldsymbol{\Sigma}_\gamma^{(i)-1}} S^{(i)} \quad \text{and} \\ \boldsymbol{\eta}_\pi^{(i+1)} &: S^{(i+1)} \times S^{(i+1)} \xrightarrow{\boldsymbol{\varepsilon}_\pi^{(i+1)}} S^{(i+1)} \xrightarrow{\boldsymbol{\Sigma}_\pi^{(i+1)-1}} S^{(i+1)} \end{aligned}$$

so that

$$\partial_{\mu_{\rho_\pi}^{(i+1)}} E^{(i)}(\mu_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)}) = - \left( \partial_{s^{(i+1)}} \boldsymbol{\mu}_\gamma^{(i)} \right) (\boldsymbol{\mu}_{\rho_\pi}^{(i+1)})^T \boldsymbol{\eta}_\gamma^{(i)} + \boldsymbol{\eta}_\pi^{(i+1)}.$$

We can therefore define a discrete-time dynamical system with update function

$$\begin{aligned} \boldsymbol{\mu}_{\rho_\pi}^{(i+1)} &: S^{(i+1)} \times S^{(i)} \rightarrow S^{(i+1)} \\ \boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(t) \times s^{(i)} &\mapsto \boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(t+1) \end{aligned} \tag{10}$$

where

$$\begin{aligned} \boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(t+1) &= \boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(t) - \lambda \left( \partial_{s^{(i+1)}} \boldsymbol{\mu}_\gamma^{(i)} \right) (\boldsymbol{\mu}_{\rho_\pi}^{(i+1)})^T \boldsymbol{\eta}_\gamma^{(i)} \left( s^{(i)}, \boldsymbol{\mu}_\gamma^{(i)}(\boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(t)) \right) + \lambda \boldsymbol{\eta}_\pi^{(i+1)} \left( \boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(t), \boldsymbol{\mu}_\pi^{(i+1)} \right) \\ &= \boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(t) - \lambda \partial_{\boldsymbol{\mu}_{\rho_\pi}^{(i+1)}} E^{(i)} \left( \boldsymbol{\mu}_{\rho_\pi}^{(i+1)}(s^{(i)}), s^{(i)} \right) \end{aligned}$$

for some ‘learning rate’  $\lambda : \mathbb{R}_+$ , and identity output function. Note that these dynamics depend on the parameters of the prior  $\pi$ , thereby witnessing the exponential structure of the Bayesian inversion operation as a ‘state-dependent stochastic channel’ (Definition 2.4); and, less subtly, the dynamics for the inverse channel  $\rho_\pi^{(i+1)}$  depend on the forwards channel  $\gamma^{(i)}$ .

Equations (10) and (9) together give a dynamical system with input space  $S^{(i)}$  and state space  $\Theta_{\rho_\pi}^{(i+1)} = S^{(i+1)} \times \mathbf{FVect}_{\mathbb{R}}(S^{(i+1)}, S^{(i+1)})$ . Given an identity output function and since  $\Theta_{\rho_\pi}^{(i+1)} \hookrightarrow \mathcal{P}S^{(i+1)}$ , we have a stochastic dynamical system of type  $\mathbf{Dyn}_{\mathcal{X}\mathcal{V}(\mathcal{P})}(S^{(i)}, S^{(i+1)})$  realising the update factor (or ‘recognition model’)  $\rho_\pi^{(i+1)}$  of the variational autoencoder game. **We now give a corresponding dynamical system for the view factor (or ‘generative model’)  $\gamma^{(i)}$ .**  $\square$

**Remark 4.11.** Typical presentations of ‘active inference’ under the free-energy principle are excessively complicated by the lack of attention paid to compositionality. Because the free-energy framework instantiates a *category* of open cybernetic systems, a radically simplified compositional presentation is possible. Such a presentation forms a companion to the present work.

**Corollary 4.12.** The free-energy framework has been used to supply a computational explanation for the pervasive bidirectionality of cortical circuits in the mammalian brain [1, 10]. A corollary of Theorem 4.10 is that this bidirectionality is furthermore justified by the abstract structure of Bayesian inference and its dynamical realisation: because Bayesian updates compose optically, a cybernetic system realising Bayesian inference compositionally must instantiate this structure. We note also that the parallel interacting bidirectional structure of the active inference game (Example 3.20) is reproduced in the cortex.

The free-energy framework realisation of autoencoder games is not unique; an alternative is found in machine learning.

**Theorem 4.13.** Consider the subcategory of **BayesLens** spanned by finite-dimensional Euclidean spaces, with morphisms generated (under sequential and parallel composition) by the (variational) autoencoder and inference games whose forwards and backwards channels emit exponential-family measures. The deep (variational) autoencoder framework [15] instantiates a category of open cybernetic systems realising games over this subcategory.

Increasingly, the variational autoencoder framework is used to model complete agents in machine learning, rather than merely dynamically realise static inference or learning problems. Indeed, thinking of the ‘free-energy framework’ as a collection of cybernetic realisations of autoencoder and active-inference games, the demonstration of the following corollary of Theorem 4.13 is unsurprising:

**Corollary 4.14.** The “deep active inference agent” [27] is a cybernetic system realising an active inference game in the variational autoencoder framework.

We have heretofore concentrated on ‘variational Bayesian’ realisations of the games introduced in §3, as they most strikingly fit the language of optimization used there. But we expect any other family of approximate inference methods to supply a corresponding category of cybernetic systems. We thus make the following conjecture.

**Conjecture 4.15.** Consider the subcategory of **BayesLens** spanned by finite-dimensional smooth manifolds, with morphisms generated (under sequential and parallel composition) by the generalized autoencoder and inference games. We expect sampling algorithms, such as Markov chain Monte Carlo, to supply a corresponding category of open cybernetic systems of interest.

Finally, we provide further justification for Remark 3.7.

**Observation 4.16.** Consider a variational autoencoder, realised as in Theorem 4.13. By choosing the parameterizations  $\mathcal{F}, \mathcal{P}$  of the forwards and backwards channels to coincide with the state spaces of their dynamical realisations, and the (static) play function  $P$  to take a parameter vector to the corresponding channel, the dynamical realisation induces a trajectory over the strategy space. Such trajectories organize into sheaf whose sections are trajectories of arbitrary length [24], spans of which are again just (generalized) dynamical systems; these spans are equivalently profunctors [4]. We can thus define a best-response function valued in profunctors whose elements are trajectories witnessing deviations of strategies to ‘better’ strategies, and whose dynamical equilibria correspond precisely to the equilibria of the ‘static’ best response function.

**On-going and Future Work** The structures sketched in this paper are merely first steps towards a categorical theory of cybernetics. In particular, since the first draft of this work was written, we have come to believe that the preliminary notions presented here of dynamical realisation, and by extension of open cybernetic system, are substantially less elegant than they could be. On-going work is focusing on this issue. We hope that a consequence of this refinement will be that the treatment of *interacting* cybernetic

systems is simplified. In this new setting, we will also treat non-stationary systems in dynamical contexts and in continuous time, thereby supplying a general compositional treatment of (amongst other things) the ‘free-energy’ framework.

Finally, with respect to applications, we are interested in using these tools to realise game-theoretic games and to investigate the connections between repeated games and dynamical realisation. There are deep links with reinforcement learning to be explored, and we seek a setting for the study of nested and multi-agent (‘ecological’) systems.

## References

- [1] A. M. Bastos, W. M. Usrey, R. A. Adams, G. R. Mangun, P. Fries & K. J. Friston (2012): *Canonical microcircuits for predictive coding*. *Neuron* 76(4), pp. 695–711, doi:[10.1016/j.neuron.2012.10.038](https://doi.org/10.1016/j.neuron.2012.10.038).
- [2] Joe Bolt, Jules Hedges & Philipp Zahn (2019): *Bayesian open games*. Available at <http://arxiv.org/abs/1910.03656v1>.
- [3] Christopher L Buckley, Chang Sub Kim, Simon McGregor & Anil K Seth (2017): *The free energy principle for action and perception: A mathematical review*. *Journal of Mathematical Psychology* 81, pp. 55–79, doi:[10.1016/j.jmp.2017.09.004](https://doi.org/10.1016/j.jmp.2017.09.004). Available at <http://arxiv.org/abs/1705.09156v1>.
- [4] Jean Bénabou (2000): *Distributors at work*. Available at <http://www.mathematik.tu-darmstadt.de/~streicher/FIBR/DiWo.pdf>. Lecture notes written by Thomas Streicher.
- [5] Kenta Cho & Bart Jacobs (2017): *Disintegration and Bayesian Inversion via String Diagrams*. *Math. Struct. Comp. Sci.* 29 (2019) 938–971, doi:[10.1017/S0960129518000488](https://doi.org/10.1017/S0960129518000488). Available at <http://arxiv.org/abs/1709.00322v3>.
- [6] Bryce Clarke, Derek Elkins, Jeremy Gibbons, Fosco Loregian, Bartosz Milewski, Emily Pillmore & Mario Román (2020): *Profunctor optics, a categorical update*. Available at <http://arxiv.org/abs/2001.07488v1>.
- [7] Bob Coecke & Aleks Kissinger (2016): *Categorical Quantum Mechanics II: Classical-Quantum Interaction*. doi:[10.1142/S0219749910006502](https://doi.org/10.1142/S0219749910006502). Available at <http://arxiv.org/abs/1605.08617v1>.
- [8] Bob Coecke & Aleks Kissinger (2017): *Categorical Quantum Mechanics I: Causal Quantum Processes*. In Elaine Landry, editor: *Categories for the Working Philosopher*, chapter 12, Oxford University Press, pp. 286–328. Available at <https://arxiv.org/abs/1510.05468v3>.
- [9] J. Nathan Foster, Michael B. Greenwald, Jonathan T. Moore, Benjamin C. Pierce & Alan Schmitt (2007): *Combinators for bidirectional tree transformations*. *ACM Transactions on Programming Languages and Systems* 29(3), p. 17, doi:[10.1145/1232420.1232424](https://doi.org/10.1145/1232420.1232424).
- [10] Karl Friston (2010): *The free-energy principle: a unified brain theory?* *Nature Reviews Neuroscience* 11(2), pp. 127–138, doi:[10.1038/nrn2787](https://doi.org/10.1038/nrn2787).
- [11] Karl Friston, Francesco Rigoli, Dimitri Ognibene, Christoph Mathys, Thomas Fitzgerald & Giovanni Pezzulo (2015): *Active inference and epistemic value*. *Cognitive Neuroscience* 6(4), pp. 187–214, doi:[10.1080/17588928.2015.1020053](https://doi.org/10.1080/17588928.2015.1020053).
- [12] Tobias Fritz (2019): *A synthetic approach to Markov kernels, conditional independence and theorems on sufficient statistics*. Available at <http://arxiv.org/abs/1908.07021v8>.



- [13] Neil Ghani, Jules Hedges, Viktor Winschel & Philipp Zahn (2016): *Compositional game theory*. *Proceedings of Logic in Computer Science (LiCS) 2018*, doi:10.1145/3209108.3209165. Available at <http://arxiv.org/abs/1603.04641v3>.
- [14] Chris Heunen, Ohad Kammar, Sam Staton & Hongseok Yang (2017): *A Convenient Category for Higher-Order Probability Theory*. doi:10.1109/lics.2017.8005137.
- [15] Diederik P. Kingma (2017): *Variational Inference & Deep Learning*. Available at <https://hdl.handle.net/11245.1/8e55e07f-e4be-458f-a929-2f9bc2d169e8>.
- [16] Jeremias Knoblauch, Jack Jewson & Theodoros Damoulas (2019): *Generalized Variational Inference*. Available at <http://arxiv.org/abs/1904.02063v4>.
- [17] Fosco Loregian (2015): *This is the (co)end, my only (co)friend*. Available at <http://arxiv.org/abs/1501.02503v4>.
- [18] Paul-André Melliès (2006): *Functorial Boxes in String Diagrams*. In: *Computer Science Logic*, Springer Berlin Heidelberg, pp. 1–30, doi:10.1007/11874683\_1.
- [19] Joe Moeller & Christina Vasilakopoulou (2018): *Monoidal Grothendieck Construction*. Available at <http://arxiv.org/abs/1809.00727v2>.
- [20] nLab authors (2020): *Grothendieck construction*. Available at <http://ncatlab.org/nlab/show/Grothendieck+construction>. Revision 62.
- [21] Mitchell Riley (2018): *Categories of Optics*. Available at <http://arxiv.org/abs/1809.00738v2>.
- [22] Mario Román (2020): *Open Diagrams via Coend Calculus*. Available at <http://arxiv.org/abs/2004.04526v2>.
- [23] Mario Román (2020): *Profunctor optics and traversals*. Available at <http://arxiv.org/abs/2001.08045v1>.
- [24] Patrick Schultz, David I Spivak & Christina Vasilakopoulou (2019): *Dynamical Systems and Sheaves*. *Applied Categorical Structures*, pp. 1–57, doi:10.1007/s10485-019-09565-x. Available at <http://arxiv.org/abs/1609.08086v4>.
- [25] Toby St. Clere Smithe (2020): *Bayesian Updates Compose Optically*. Available at <http://arxiv.org/pdf/2006.01631v1>.
- [26] David I. Spivak (2019): *Generalized Lens Categories via functors  $\mathcal{C}^{op} \rightarrow \mathbf{Cat}$* . Available at <http://arxiv.org/abs/1908.02202v2>.
- [27] Kai Ueltzhöffer (2018): *Deep Active Inference*. *Biological Cybernetics* 112(6), pp. 547–573, doi:10.1007/s00422-018-0785-7. Available at <http://arxiv.org/abs/1709.02341v5>.