

A PRESCRIPTIVE THEORY FOR BRAIN-LIKE INFERENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

The Evidence Lower Bound (ELBO) is a widely used objective for training deep generative models, such as Variational Autoencoders (VAEs). In the neuroscience literature, an identical objective is known as the Free Energy Principle (FEP), hinting at a potential unified framework for brain function and machine learning. Despite its utility in interpreting generative models, including diffusion models, ELBO maximization is often seen as too broad to offer prescriptive guidance for specific architectures in neuroscience or machine learning. In this work, we show that maximizing ELBO under Poisson assumptions for general sequences leads to a spiking neural network that performs Bayesian posterior inference through its membrane potential dynamics. The resulting model, the iterative Poisson VAE ($i\mathcal{P}$ -VAE), has a closer connection to biological neurons than previous brain-inspired predictive coding models based on Gaussian assumptions. Compared to amortized and iterative VAEs, $i\mathcal{P}$ -VAE learns sparser representations and exhibits superior generalization to out-of-distribution samples. These findings suggest that optimizing ELBO, combined with Poisson assumptions, provides a solid foundation for developing prescriptive theories in NeuroAI.

1 INTRODUCTION

Optimizing the Evidence Lower Bound (ELBO) serves as a unifying objective for training deep generative models (Hinton et al., 1995; Dayan et al., 1995; Kingma & Welling, 2014; Rezende et al., 2014; Luo, 2022). Even when models don’t explicitly reference ELBO, they’re often optimizing objectives closely related to it (Luo, 2022; Kingma & Gao, 2023). This is directly paralleled by the Free Energy Principle (FEP) in neuroscience, which absorbs previous theoretical frameworks like Predictive Coding, Bayesian Brain, and Active Learning (Friston, 2005; 2009; 2010). FEP states that a single objective, the minimization of variational free energy, is all that is needed. Because this is equivalent to maximizing ELBO, it suggests a powerful unifying theoretical framework for neuroscience and machine learning (Friston, 2010).

However, in many ways, Free Energy (and by proxy, ELBO) is too general to be useful as a theory (Gershman, 2019; Andrews, 2021). In practice, the specific implementations of FEP predictive coding have been difficult to map directly onto neural circuits (Millidge et al., 2021a; 2022), struggling with negative rates and prediction signals that have not been observed empirically (Walsh et al., 2020; Millidge et al., 2022). Similarly, in machine learning, it is often discovered after the fact that a new objective is actually ELBO maximization (or KL minimization; Hobson (1969)) masquerading as something else (Kingma & Gao, 2023)—and not the other way around. If ELBO is “all you need,” then why is ELBO not prescriptive?

One possibility, at least in neuroscience, is that ELBO’s lack of prescriptive theory results from incorrect approximating distributions. In fact, most of the difficulty mapping predictive coding onto neural circuits has to do with terms that result from the Gaussian assumption (Millidge et al., 2022). In contrast, biological neurons are largely modeled as conditionally Poisson (Goris et al., 2014).

Recent work provides a potential prescriptive route: replacing Gaussians with Poisson distributions. To this end, Vafaii et al. (2024) introduced a reparameterization algorithm for training Poisson Variational Autoencoders (\mathcal{P} -VAE). They observed that replacing Gaussians in ELBO reduces to an amortized version of sparse coding, an influential model inspired by the brain that captures many features of the selectivity in early visual cortex (Olshausen & Field, 1996; 2004). \mathcal{P} -VAE learns sparse representations, avoids posterior collapse, and performs better on downstream classification

tasks. However, the authors identified a large amortization gap in \mathcal{P} -VAE (Vafaii et al., 2024), adding to a growing body of work that highlights limitations of amortized inference Cremer et al. (2018); Kim & Pavlovic (2021). A potential solution is to develop more general iterative inference solutions, or hybrid iterative-amortized ones (Marino et al., 2018; Kim et al., 2018).

Here, we extend the Poisson VAE to include iterative inference (“iterative \mathcal{P} -VAE,” or $i\mathcal{P}$ -VAE). This results in a generalization of predictive coding that maps well onto biological neurons. $i\mathcal{P}$ -VAE implements Bayesian posterior inference via private membrane potential dynamics, resembling a spiking version of the Locally Competitive Algorithm (LCA) for sparse coding (Rozell et al., 2008). This solution avoids the major problems with predictive coding: there is no explicit prediction, neurons communicate through spikes, and feedback is modulatory—all consistent with real neurons (Gilbert & Li, 2013; Kandel et al., 2000). But how effective is $i\mathcal{P}$ -VAE as a machine learning model?

We evaluate $i\mathcal{P}$ -VAE in terms of convergence, reconstruction performance, efficiency, and out-of-distribution (OOD) generalization. We find that $i\mathcal{P}$ -VAE converges to sparse posterior representations, outperforming other iterative VAEs (Kim et al., 2018; Marino et al., 2018).

Contributions. We introduce a new architecture, $i\mathcal{P}$ -VAE, that accomplishes the following:

- Deriving the ELBO for sequences with Poisson-distributed latents results in a neural network that spikes, and performs predictive coding in the dynamics of the membrane potential.
- By reusing the same set of weights across iterations and utilizing sparse, integer spike counts, $i\mathcal{P}$ -VAE is well-suited for hardware implementations and energy-efficient deployment.
- $i\mathcal{P}$ -VAE demonstrates robust out-of-distribution generalization, excelling in both within-dataset perturbations and cross-dataset generalization.

Taken together, $i\mathcal{P}$ -VAE is a powerful brain-inspired architecture that tightly maps onto biological neurons while outperforming much larger models in key objectives such as performance, parameter count, sparsity, and out-of-distribution generalization.

2 BACKGROUND AND RELATED WORK

Generative models and ELBO. Generative models learn to represent the data distribution, $p(\mathbf{x})$, typically by invoking latent variables \mathbf{z} , such that $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ (Bishop & Nasrabadi, 2006). The key challenge is computing, $p(\mathbf{z}|\mathbf{x})$, the posterior distribution of these latent variables given the data, which is typically intractable except for simple cases.

Variational inference offers a practical solution by introducing an approximate posterior $q_\phi(\mathbf{z}|\mathbf{x})$ parameterized by ϕ (Blei et al., 2017). The goal is to make this approximation as close as possible to the true posterior $p(\mathbf{z}|\mathbf{x})$. Ideally, one would minimize the KL divergence between $q_\phi(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z}|\mathbf{x})$, but since we cannot compute $p(\mathbf{z}|\mathbf{x})$ exactly, direct minimization is not feasible.

The Evidence Lower Bound (ELBO) provides a tractable objective that indirectly minimizes the KL divergence between the approximate and true posteriors. Specifically, the relationship is:

$$\log p(\mathbf{x}) = \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]}_{\text{ELBO}} + \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})) \quad (1)$$

Since $\log p(\mathbf{x})$ does not depend on ϕ and the KL divergence is non-negative, maximizing the ELBO effectively minimizes the intractable KL divergence (Hinton et al., 1995; Kingma & Welling, 2014; Rezende et al., 2014). Interestingly, even when generative models seem to optimize a different loss function, like diffusion models (Chan, 2024; Ho et al., 2020), they are often still performing KL minimization through the ELBO (Kingma & Gao, 2023; Luo, 2022).

ELBO in Neuroscience. The Evidence Lower Bound (ELBO) has an identical formulation in neuroscience, where it is referred to as the Free Energy (Friston, 2005; 2009; 2010). The Free Energy Principle (FEP) extends the framework of perception as inference (Alhazen, 1011–1021 AD;

Von Helmholtz, 1867; Mumford, 1992), drawing concepts from predictive coding (PC; Srinivasan et al. (1982); Rao & Ballard (1999)). Extensive research has explored how PC might be implemented by neurons (Boerlin et al., 2013; Millidge et al., 2021a), and PC has been applied in machine learning for predictive models (Lotter et al., 2017; Wen et al., 2018; Millidge et al., 2024).

Despite their neural inspiration, FEP is challenging to map directly onto neuronal circuits (Kogo & Trengove, 2015; Aitchison & Lengyel, 2017; Millidge et al., 2022). This difficulty results from assuming Gaussian for the approximate posterior and prior (Millidge et al., 2022). The Gaussian assumption results in models with explicit predictions or prediction errors, which have not been observed empirically (Mikulasch et al., 2023). Solutions also struggle with how to avoid negative firing rates due to subtraction operations (Bastos et al., 2012; Keller & Mrsic-Flogel, 2018). While leaky integrate-and-fire (LIF) circuits can be engineered to perform predictive coding (Boerlin et al., 2013), these implementations do not naturally arise from ELBO maximization, making the theory more postdictive than prescriptive. The related framework of sparse coding can be thought of as a form of predictive coding with a sparse prior Olshausen & Field (1996; 2004). A biologically plausible implementation of sparse coding, known as the locally competitive algorithm (LCA; Rozell et al. (2008)), results naturally in a dynamic update rule that resembles neural circuits. However, LCA relies on maximum a posteriori inference, which is restrictive if we aim to sample from the full posterior distribution.

Bayesian posterior inference: iterative versus amortized. In contrast to predictive coding, Variational Autoencoders (VAEs) introduced a computationally-efficient solution to maximize ELBO through *amortized* inference (Kingma & Welling, 2014; Rezende et al., 2014). Amortized inference uses a parameterized neural network (the “encoder” or “recognition” network) to produce the parameters of an approximate posterior, $q_\phi(\mathbf{z}|\mathbf{x})$, in one shot. The term “amortized” reflects that the computational cost of inference is paid during training, not at test time, similar to cost distribution in accounting (Gershman & Goodman, 2014). While amortized inference is considered efficient, it can suffer from an *amortization gap*—the discrepancy between the approximate posterior provided by the encoder and the optimal variational parameters—which can be significant (Cremer et al., 2018).

To address the amortization gap, hybrid approaches have been developed that introduce iterative elements into the VAE framework (Marino et al., 2018; Kim et al., 2018; Marino et al., 2021). For example, Marino et al. (2018) proposed a method where the encoder network takes as input both the data sample \mathbf{x} and the gradients of the loss with respect to the variational parameters $\nabla_\lambda \mathcal{L}$, with $\lambda = \{\mu, \sigma^2\}$. Alternatively, semi-amortized inference (Kim et al., 2018) starts with an amortized initial estimate and refines it using stochastic variational inference (SVI) updates (Hoffman et al., 2013). Our method is closely related to these approaches, and we compare to them in the results.

Although VAEs and predictive coding are related through their optimization of ELBO (Marino, 2022), recent work has made that connection more explicit, demonstrating that classical predictive coding networks can be seen as a subclass of iterative inference in VAEs (Boutin et al., 2020). A key difference between our work and Boutin et al. (2020) is that they show the Rao & Ballard (1999) loss function arises from assuming a delta-function posterior in the ELBO. In our work, predictive coding naturally emerges in the dynamics of the log spike rates, which comes from a fairly general assumption of Poisson distributions.

Poisson VAE. A large body of literature in neuroscience has demonstrated that neuron spike counts are well described by a Poisson process over short counting windows (Goris et al., 2014). Building on this, Vafaii et al. (2024) introduced the Poisson Variational Autoencoder (P-VAE), which performs posterior inference using discrete spike counts. They developed a Poisson reparameterization trick and derived the ELBO for Poisson-distributed VAEs (\mathcal{P} -VAE).

In \mathcal{P} -VAE, the KL term penalizes firing rates, similar to sparse coding, and the ELBO, when paired with a linear generative model, reduces to amortized sparse coding. When trained on natural image patches, \mathcal{P} -VAE learns sparse solutions with Gabor-like basis vectors and latent sparsity, similar to sparse coding. While \mathcal{P} -VAE outperformed Gaussian VAEs in sparsity and downstream classification, the authors noted a significant performance gap with traditional sparse coding, likely arising from an amortization gap due to the lack of iterative updates. Our work builds upon \mathcal{P} -VAE, suggesting that Poisson is the right choice for parameterizing the distributions in ELBO (see Appendix B for a discussion).

3 INTRODUCING THE ITERATIVE POISSON VAE (IP-VAE)

In this section, we derive the ELBO for sequences with Poisson distributions. We show the resulting architecture (iP-VAE) implements iterative Bayesian posterior inference with dynamics on the log rates. We relate this directly to membrane potential dynamics in a spiking neural network and show that it solves many of the implementation limitations of classic predictive coding.

General setup. We conceptualize iterative inference by starting with the more general framework of inference over a sequence (Chung et al., 2015). From there, we can treat iterative inference for images as a sequences of the same image repeated at all time points. This approach is appealing because dynamics emerge necessarily, and it builds a foundation for future work on dynamic sequences.

Consider a sequence of $T + 1$ observed data points, $\vec{x} = \{\mathbf{x}_t : t = 0, \dots, T\}$ where $\mathbf{x}_t \in \mathbb{R}^M$, and corresponding latent variables, $\vec{z} = \{\mathbf{z}_t : t = 0, \dots, T\}$, where each \mathbf{z}_t is K -variate. We denote the full probabilistic generative model as the joint distribution, $p(\vec{x}, \vec{z})$. A reasonable starting assumption for modeling the physical world is Markovian dependence between consecutive data points (Van Kampen, 1992), resulting in the marginal distribution:

$$p(\vec{x}) = \int p(\vec{x}, \vec{z}) d\vec{z} = p(\mathbf{x}_0) \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{x}_{t-1}), \quad (2)$$

where $p(\mathbf{x}_0) = \int p(\mathbf{x}_0 | \mathbf{z}_0) p(\mathbf{z}_0) d\mathbf{z}_0$, and $p(\mathbf{x}_t | \mathbf{x}_{t-1}) = \int p(\mathbf{x}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{x}_{t-1}) d\mathbf{z}_t$. For our sequence data, the ELBO can be written as follows:

$$\begin{aligned} \log p_\theta(\vec{x}) &\geq \mathbb{E}_{q_\phi(\vec{z}|\vec{x})} \left[\log \frac{p_\theta(\vec{x}, \vec{z})}{q_\phi(\vec{z}|\vec{x})} \right] \\ &= \mathbb{E}_{q_\phi(\vec{z}|\vec{x})} \left[\log p_\theta(\vec{x}|\vec{z}) \right] - \mathcal{D}_{\text{KL}}(q_\phi(\vec{z}|\vec{x}) \parallel p_\theta(\vec{z})) \\ &= \mathcal{L}_{\text{ELBO}}(\vec{x}; \theta, \phi), \end{aligned} \quad (3)$$

where $p_\theta(\vec{z})$ is a prior (either learned or fixed) over latents, and $p_\theta(\vec{x}|\vec{z})$ is the conditional likelihood distribution, which is computed via a decoder network. Model parameters, (ϕ, θ) —corresponding to the encoder and decoder networks of a VAE, respectively—are jointly optimized. Below we will express the ELBO for sequences when using the Poisson Variational Autoencoder framework.

Iterative Poisson VAE. To extend the \mathcal{P} -VAE to sequences, iP-VAE needs to make explicit how the prior and posterior distributions update with each sample. The simplest starting point is assuming stationarity, implying that the posterior over the previous stimulus should act as a prior for the current one (although future extensions could extend to nonstationary signals such as videos with a more sophisticated update rule). Because of the Markovian assumption, the prior, $p(\vec{z})$, then factorizes into the initial prior, $p(\mathbf{z}_0)$ and a product over all time steps:

$$p(\vec{z}) = p(\mathbf{z}_0) \prod_{t=1}^T p(\mathbf{z}_t | \mathbf{x}_{t-1}) \quad (4)$$

The initial prior, $p(\mathbf{z}_0) = \mathcal{Pois}(\mathbf{z}_0; \mathbf{r}_0)$, is Poisson with learned prior rates, $\mathbf{r}_0 \in \mathbb{R}_{>0}^K$. Subsequent time steps have prior rates that depend on the stimulus from the previous time step, $p(\mathbf{z}_t | \mathbf{x}_{t-1}) = \mathcal{Pois}(\mathbf{z}_t; \mathbf{r}_t(\mathbf{x}_{t-1}))$. The approximate posterior factorizes as well:

$$q(\vec{z}|\vec{x}) = q(\mathbf{z}_0 | \mathbf{x}_0) \prod_{t=1}^T q(\mathbf{z}_t | \mathbf{x}_t, \mathbf{x}_{t-1}), \quad (5)$$

with initial posterior, $q(\mathbf{z}_0 | \mathbf{x}_0) = \mathcal{Pois}(\mathbf{z}_0; \mathbf{r}_0 \odot \delta \mathbf{r}(\mathbf{x}_0))$, and time-dependent posterior, $q(\mathbf{z}_t | \mathbf{x}_t, \mathbf{x}_{t-1}) = \mathcal{Pois}(\mathbf{z}_t; \mathbf{r}_t(\mathbf{x}_{t-1}) \odot \delta \mathbf{r}(\mathbf{x}_t))$, both parameterized as Poisson distributions. We follow the formulation in Vafaii et al. (2024), and define the posterior rates via an element-wise multiplicative interaction between \mathbf{r} and some gain modulator, $\delta \mathbf{r} \in \mathbb{R}_{>0}^K$. This is a natural choice because rates must be positive, and without loss of generality, the relationship between two positive variables can be written in terms of a base rate, and a multiplicative gain on that base rate.

The conditional log-likelihood for $i\mathcal{P}$ -VAE factorizes into a sum over individual sample likelihoods $\log p(\vec{x}|\vec{z}) = \sum_{t=0}^T \log p(\mathbf{x}_t|\mathbf{z}_t)$. The KL-term of the ELBO (eq. (3)) also factorizes:

$$\begin{aligned} \mathcal{D}_{\text{KL}}\left(q(\vec{z}|\vec{x}) \parallel p(\vec{z})\right) &= \mathcal{D}_{\text{KL}}\left(q(\mathbf{z}_0|\mathbf{x}_0) \parallel p(\mathbf{z}_0)\right) + \sum_{t=1}^T \mathcal{D}_{\text{KL}}\left(q(\mathbf{z}_t|\mathbf{x}_t, \mathbf{x}_{t-1}) \parallel p(\mathbf{z}_t|\mathbf{x}_{t-1})\right) \\ &= \mathbf{r}_0 \cdot f(\delta\mathbf{r}(\mathbf{x}_0)) + \sum_{t=1}^T \mathbf{r}_t(\mathbf{x}_{t-1}) \cdot f(\delta\mathbf{r}(\mathbf{x}_t)), \end{aligned} \quad (6)$$

where \cdot represents a vector dot product, and $f(y) = 1 - y + y \log y$ is applied element-wise. Because rates are positive, the KL term penalizes large rates, acting like a sparsity penalty (Vafaii et al., 2024). The remaining sections describe how we specify the multiplicative gain, $\delta\mathbf{r}$, which results in adaptive Bayesian posterior updating in the dynamics of the model.

Bayesian posterior updates using membrane potential dynamics Because rates are positive and prior and posterior rates interact multiplicatively, it is difficult to implement dynamic updates directly on rates. A natural solution is to define updates on log rates, $\mathbf{u}(t) := \log \mathbf{r}(t)$, with \mathbb{R}^K as our state space for a K -dimensional latent space.

Dynamic updates on log-rates is both a mathematical convenience and biologically realistic. Because of internal noise, the spike threshold of real neurons is best modeled as an expansive nonlinearity like an exponential (Priebe et al., 2004; Fourcaud-Trocmé et al., 2003). Further, synapses have a compressive nonlinearity for incoming spikes because of synaptic depression (Abbott et al., 1997). Here, we take $\log(x)$ to be the synaptic nonlinearity and $\exp(x)$ to be the spiking nonlinearity. For the aforementioned reasons, $\mathbf{u}(t)$ can be interpreted quite literally as membrane potentials.

We define the model updates as $\mathbf{u}_{t+1} = \mathbf{u}_t + \delta\mathbf{u}_t$, with $\mathbf{r}_t = \exp(\mathbf{u}_t)$ acting as the corresponding prior rates at time t , and $\mathbf{r}_t \odot \delta\mathbf{r} = \exp(\mathbf{u}_{t+1})$, as the posterior rates at time t . When processing the next input in the sequence, we take the previous posterior and use it as our current prior. This works, because in the present paper, we restrict ourselves to stationary inputs comprised of the same image presented multiple times.

A natural choice for $\delta\mathbf{u}$ is the gradient of the loss with respect to \mathbf{u} , through the samples \mathbf{z} . However, the KL term results in high order terms, which for this implementation we approximate as the following dynamics (See appendix D for a detailed derivation):

$$\delta\mathbf{u}_t = \mathbf{J}_\theta \cdot \Delta_t = \left. \frac{\partial f_\theta(\mathbf{z})}{\partial \mathbf{z}} \right|_{\mathbf{z}=\mathbf{z}_t} \cdot (\mathbf{x}_t - f_\theta(\mathbf{z}_t)), \quad (7)$$

where \mathbf{J}_θ is the Jacobian of the decoder, f_θ , which is a function of sampled spike counts \mathbf{z} .

Importantly, this form aligns with real neuronal properties for several reasons. Since the comparison, $\mathbf{x}_t - f_\theta(\mathbf{z}_t)$, is based on spikes, each neuron’s update does not directly depend on the internal states of other neurons, which matches how real neurons function (Kandel et al., 2000). Additionally, because the comparison happens on membrane potential (log rates), feedback will appear as a modulatory signal on rate, which is also consistent with neuroscience literature (Gilbert & Li, 2013). Finally, this update (eq. (7)) resembles a generalization of Rao & Ballard (1999) for nonlinear generative models and avoids hacky solutions to keep rates positive, after subtracting them.

It is straightforward to see how this is an SNN for linear decoder networks. If $f_\theta(\mathbf{z}) = \Phi\mathbf{z}$, then

$$\begin{aligned} \delta\mathbf{u}_t &= \Phi^T(\mathbf{x}_t - \Phi\mathbf{z}) \\ &= \Phi^T\mathbf{x} - \Phi^T\Phi\mathbf{z} \\ &= \Phi^T\mathbf{x} - \mathbf{W}\mathbf{z}, \end{aligned} \quad (8)$$

where the first term is the feedforward receptive fields (the input current) and the second term, \mathbf{W} , are the recurrent weights between neurons, implementing lateral competition. Note that they only communicate with each other through spikes, \mathbf{z} . Thus for linear generative models, $i\mathcal{P}$ -VAE closely resembles the locally competitive algorithm for sparse coding (LCA; Rozell et al. (2008)), except that it is explicitly spiking and does not have a leak term (although this could be included by replacing the diagonal of the recurrent term with a leak rather than having neurons operate on their own spikes).

In this section, we showed how following some fairly general assumptions for optimizing ELBO with Poisson distribution, led us to a spiking neural network that implements Bayesian posterior updates via predictive coding in the membrane potential dynamics. In the next section, we evaluate $i\mathcal{P}$ -VAE and compare it to amortized \mathcal{P} -VAE, as well as iterative Gaussian VAEs.

4 EXPERIMENTS

We performed empirical analyses of $i\mathcal{P}$ -VAE and alternative iterative VAE models. In section 4.1, we test the general performance and stability of inference dynamics, including generalization to longer sequence lengths. Section 4.2 shows $i\mathcal{P}$ -VAE closes the gap with sparse coding. Section 4.3 demonstrates robustness to out-of-distribution (OOD) samples by evaluating models trained on MNIST (LeCun et al., 2010) with perturbed samples (e.g., rotated MNIST). We then evaluate OOD generalization from MNIST to other character-based datasets in section 4.3. Finally, in section 4.4, we visualize the learned weights of $i\mathcal{P}$ -VAE, revealing their compositional nature, which is consistent with $i\mathcal{P}$ -VAE’s strong generalization capabilities. We push the limits of MNIST-trained models by testing their performance on natural images.

Architecture notation. We experimented with both convolutional and multi-layer perceptron (MLP) architectures. We highlight the **encoder** and **decoder** networks using **red** and **blue**, respectively. We use the $\langle \text{enc}|\text{dec} \rangle$ convention to clearly specify which type was used. For example $\langle \text{mlp}|\text{mlp} \rangle$ means both encoder and decoder networks were mlp. We use the notation $\langle \text{jacob}|\text{mlp} \rangle$ to denote our fully iterative (non-amortized) $i\mathcal{P}$ -VAE. We chose symmetrical architectures, such that $\langle \text{mlp}|\text{mlp} \rangle$ has exactly twice as many parameters as $\langle \text{jacob}|\text{mlp} \rangle$.

Datasets. For the generalization results, we use MNIST, extended MNIST (EMNIST; Cohen et al. (2017)), Omniglot (Lake et al., 2015) and Imagenet32 (Chrabaszcz et al., 2017). We resize Omniglot and Imagenet32 to 28×28 for more straightforward comparisons. We also replicated the sparsity analysis in Fig. 3 of Vafaii et al. (2024) in our Table 1, using the van Hateren natural images dataset with whitened, contrast normalized 16×16 patches.

Alternative models. We compare our iterative \mathcal{P} -VAE ($i\mathcal{P}$ -VAE) to \mathcal{P} -VAE. The main difference between their two architectures is that the latter independently parameterizes an encoder, whereas the former constructs its encoder adaptively by inverting the decoder. We also compare to state-of-the-art methods that combine iterative with amortized inference. These include iterative amortized VAE (ia-VAE; Marino et al. (2018)), and semi-amortized VAE (sa-VAE; Kim et al. (2018)). Since ia-VAE comes with both hierarchical (h) and single-level (s) variants, we compare to each of these.

Number of iterations. For $i\mathcal{P}$ -VAE, we experimented with different numbers of training iterations, T_{train} . During training, we differentiate through the entire sequence of iterations, which can lead to qualitatively different dynamics. We report results for $T_{\text{train}} = 4, 16, 32, 64$. For generalization results, we use a model with $T_{\text{train}} = 64$. At test time, we report results using $T_{\text{test}} = 1,000$ iterations, unless stated otherwise. For semi-amortized models, we use their default number of train and test iterations found in their code, unless stated otherwise (sa-VAE: $T_{\text{train}} = T_{\text{test}} = 20$; ia-VAE: $T_{\text{train}} = T_{\text{test}} = 5$).

4.1 STABILITY BEYOND THE TRAINING REGIME AND CONVERGENCE.

An algorithm with strong generalization potential should learn how to perform inference that extends beyond the training regime. We evaluated this by training models on MNIST under different numbers of training iterations, $T_{\text{train}} = 4, 16, 32$, and 64 . We used both $\langle \text{jacob}|\text{mlp} \rangle$ and $\langle \text{jacob}|\text{conv} \rangle$ architectures and then tested each model on its ability to keep improving beyond the training number of iterations. In Fig. 1a, we show that $i\mathcal{P}$ -VAE converges. Even with as few as 4 iterations, $i\mathcal{P}$ -VAE learns to keep improving. We also observe that increasing the number of training iterations has an interesting effect: $i\mathcal{P}$ -VAE trained with a larger number of iterations starts from worse performance, but converge to better solutions (Fig. 1a). This suggests $i\mathcal{P}$ -VAE learns dynamics that depend on the training sequence length, but generalizes beyond the training set in all cases.

In contrast, the two hybrid models (sa-VAE and ia-VAE) start with strong amortized initial guesses, but plateau rapidly (Fig. 1a, right), and converge to a much higher MSE than $i\mathcal{P}$ -VAE models, which

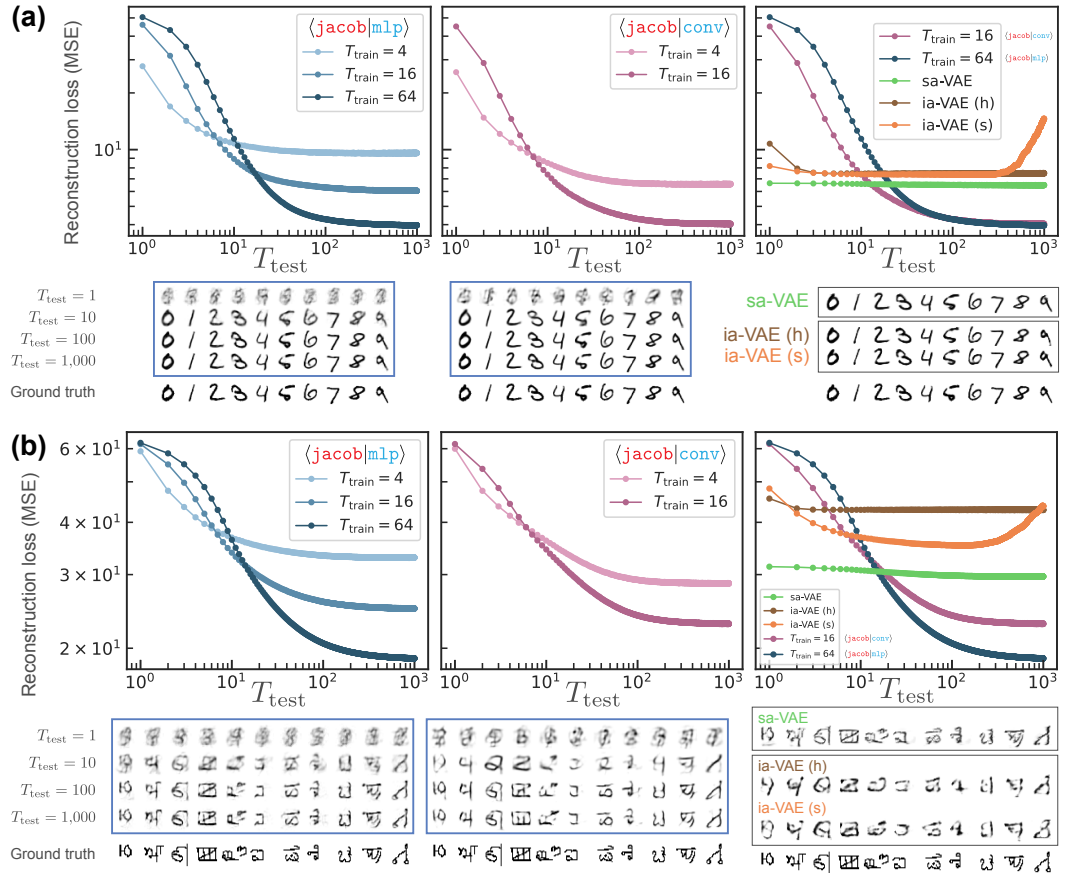
324 have a fraction of the parameters. The authors of sa-VAE were aware of issues regarding dominance
 325 of the iterative part of the algorithm for Omniglot, and reported using tricks like gradient clipping to
 326 mitigate it, which we suspect is the source for our observations on MNIST (see footnote 6 in Kim
 327 et al. (2018)). We also see that ia-VAE (single-level) starts to diverge outside its training regime. ¹

328 Overall, $i\mathcal{P}$ -VAE achieves the best reconstruction performance and continues to improve outside the
 329 training regime, unlike other models. This shows the first sign of OOD generalization in $i\mathcal{P}$ -VAE:
 330 temporal generalization. In later sections, we test whether $i\mathcal{P}$ -VAE can generalize OOD in vision
 331 tasks, but first, we evaluate the performance and sparsity on natural images as in Vafaii et al. (2024).
 332

333 4.2 IP-VAE CLOSES THE GAP WITH SPARSE CODING
 334

335 One of the limitations of previous work with \mathcal{P} -VAE, was that the authors identified a large perfor-
 336 mance gap between \mathcal{P} -VAE and LCA sparse coding (Vafaii et al., 2024). Here, we evaluated $i\mathcal{P}$ -VAE
 337 and compared models on their ability to reconstruct whitened natural image patches (table 1). Unlike
 338 \mathcal{P} -VAE, $i\mathcal{P}$ -VAE performs as well as LCA with similar sparsity levels. \mathcal{P} -VAE, and the two hybrid
 339 approaches, have many more parameters and achieve much worse performance. ²

340 ¹It's worth noting that in our hands, ia-VAE (s) often resulted in nans at test time upon going beyond T_{train} .
 341 ²The performance ia-VAE and sa-VAE might be modestly improved by tuning the tradeoff between recon-
 342 struction and the KL term.
 343



373 Figure 1: $i\mathcal{P}$ -VAE learns to learn. **(a)** Training $i\mathcal{P}$ -VAE on as few as $T_{\text{train}} = 4$ time steps allows
 374 it to generalize and keep improving its inference beyond the training domain. This holds true
 375 irrespective of the $i\mathcal{P}$ -VAE architecture; left, $\langle \text{jacob} | \text{mlp} \rangle$; middle, $\langle \text{jacob} | \text{conv} \rangle$. In contrast,
 376 hybrid amortized/iterative models do not improve, and either remain flat or diverge (right). **(b)**
 377 $i\mathcal{P}$ -VAE trained on MNIST generalizes to Omniglot at test time. All models in this figure were
 trained on MNIST, and tested either on MNIST (a), or Omniglot (b).

Table 1: Model performance and efficiency. We prefer lightweight models that achieve low reconstruction loss using sparse representations and fewer parameters. We reported results on natural image patches extracted from the van Hateren dataset (Van Hateren & van der Schaaf, 1998). All models have $K = 512$ dimensional latent space. For the $i\mathcal{P}$ -VAE models, we scaled the β parameter proportional to the number of training inference iterations. Specifically, we chose $\beta = 3/8 * T_{\text{train}}$. We found that $i\mathcal{P}$ -VAE results were robust to variations in β . Entries formatted as mean \pm std.

Model	β	Architecture	# params \downarrow	MSE \downarrow	Sparsity \uparrow		# iters	
					lifetime	%	train	test
$i\mathcal{P}$ -VAE	24.00	$\langle \text{jacob} \text{lin} \rangle$	0.13 M	12.0 \pm 2.6	0.79 \pm .03	60.0	64	1K
$i\mathcal{P}$ -VAE	3.00	$\langle \text{jacob} \text{lin} \rangle$	0.13 M	27.5 \pm 7.1	0.85 \pm .02	73.2	8	1K
$i\mathcal{P}$ -VAE	1.50	$\langle \text{jacob} \text{lin} \rangle$	0.13 M	50.4 \pm 15.5	0.90 \pm .03	83.3	4	1K
\mathcal{P} -VAE	0.50	$\langle \text{conv} \text{lin} \rangle$	3.44 M	101.9 \pm 25.3	0.76 \pm .16	65.9	1	1
\mathcal{P} -VAE	0.75	$\langle \text{conv} \text{lin} \rangle$	3.44 M	119.4 \pm 26.4	0.83 \pm .09	77.7	1	1
\mathcal{P} -VAE	1.00	$\langle \text{conv} \text{lin} \rangle$	3.44 M	131.8 \pm 31.2	0.90 \pm .08	84.1	1	1
LCA	0.28	-	0.13 M	16.1 \pm 8.1	0.79 \pm .02	65.6	1K	1K
LCA	0.44	-	0.13 M	28.5 \pm 14.1	0.86 \pm .02	73.9	1K	1K
LCA	0.70	-	0.13 M	50.1 \pm 25.2	0.92 \pm .01	83.4	1K	1K
ia-VAE (s)	1.00	$\langle \text{mlp} \text{mlp} \rangle$	39.55 M	80.08 \pm 21.06	0.36 \pm .00	\sim 0.0	5	10
sa-VAE	1.00	$\langle \text{conv} \text{conv} \rangle$	1.67 M	97.74 \pm 38.97	0.36 \pm .00	\sim 0.0	20	20

4.3 OUT-OF-DISTRIBUTION GENERALIZATION.

In this section, we evaluate whether MNIST-trained models generalize to OOD perturbations and dataset. First, we tested whether MNIST-trained models generalize to Omniglot (see Fig. 1b). We found that $i\mathcal{P}$ -VAE improves over iterations and outperforms alternative models in terms of reconstruction quality. In this section, we evaluate two levels of generalization tasks: (1) within-dataset perturbations; and, (2) across similar datasets (i.e., digits to characters).

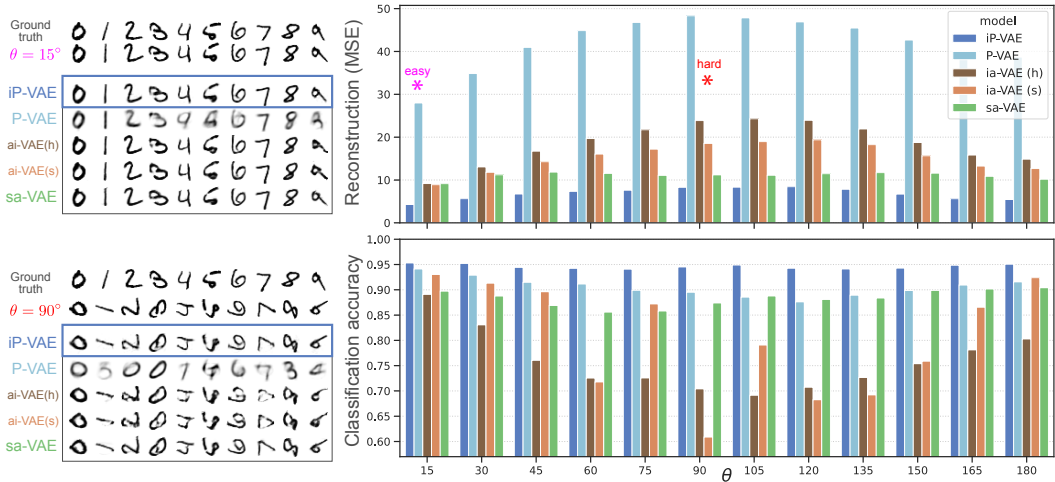


Figure 2: Robustness to training set perturbation. We rotated MNIST digits and evaluated model performance in both reconstruction of the perturbed inputs, and classification accuracy. On the left, we show reconstructed samples for easy ($\theta = 15^\circ$) and hard ($\theta = 90^\circ$) tasks across different models. On the right, we visualize the average reconstruction loss and classification accuracies over different rotations. Both visually and quantitatively, $i\mathcal{P}$ -VAE maintains a high performance regardless of the rotation, and outperforms alternative models.

OOD generalization to within-dataset perturbation. We tested whether models trained on standard MNIST generalized to rotated MNIST digits. We rotate MNIST between 0 and 180 degrees, with incremental steps of 15 degrees. We then test (a) whether models are capable of reconstructing the rotated digits, and (b) whether the representations of rotated digits can be used to classify them (Fig. 2). *iP*-VAE and *sa*-VAE demonstrated consistent performance across angles, both in terms of reconstruction loss and classification accuracy. Amortized *P*-VAE shows worse reconstruction performance than all iterative models, but its classification accuracy is remarkably consistent across angles, beating or matching all models except for *iP*-VAE. *ia*-VAE variants were greatly affected by the rotation, with significant falloff in both their classification score and reconstruction. Overall, *iP*-VAE maintains stable performance across rotations at levels above alternative models.

OOD generalization across similar datasets. If a model learns compositional features, and if it employs an effective inference algorithm that leverages those features, it should be able to represent datasets that are within the same distributional vicinity as the training set. To test this, we evaluated MNIST-trained models on EMNIST and Omniglot. We report both mean squared error (MSE) of reconstruction and classification accuracy ³.

Again, *iP*-VAE exhibited superior reconstruction performance over other models, both visually and MSE (Fig. 3). It also had substantially higher classification accuracy, suggesting it learns a compositional code and has strong generalization potential.

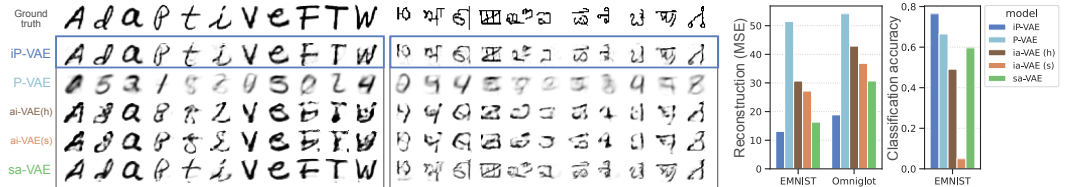


Figure 3: Evaluating generalization from models trained on MNIST digits to novel character datasets (EMNIST and Omniglot) at test time. The right panel shows the average classification performance on latent representations for EMNIST. The middle-right panel compares the reconstruction performance on EMNIST and Omniglot. The left two panels visualize the reconstructions on EMNIST and Omniglot, respectively. In both metrics, *iP*-VAE maintains high performance compared to alternative models.

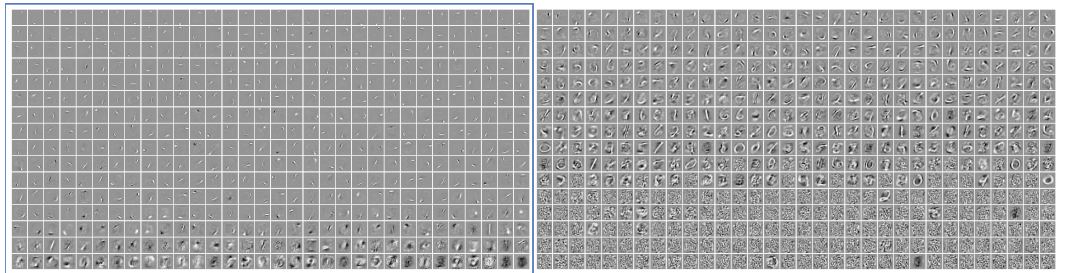


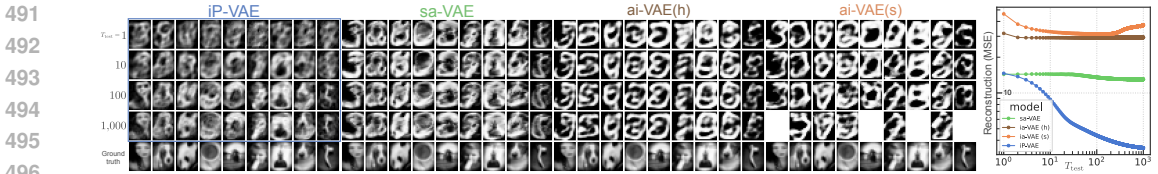
Figure 4: *iP*-VAE learns a compositional set of features for the last layer’s weights, enabling its generalization capacity. Left, *iP*-VAE with a *<jacob|mlp>* architecture; right, *P*-VAE with an *<mlp|mlp>* architecture. Both models were trained on MNIST, but only *iP*-VAE develops Gabor-like features. In contrast, the non-iterative, amortized *P*-VAE clearly overfits to MNIST. Features are ordered in ascending order of their weight distribution kurtosis to highlight the sparse nature of *iP*-VAE feature space. Best viewed when zoomed in.

4.4 A COMPOSITIONAL CODE THAT GENERALIZES ACROSS DOMAINS .

Using the *<jacob|mlp>* variant of *iP*-VAE, we visualized the 512 learned features of the last layer of the *mlp* decoder. In Fig. 4, we show the features learned by *iP*-VAE trained on MNIST and contrast

³We omit classification accuracy for Omniglot due to its large number of classes (over 1,000)

486 them to features learned by \mathcal{P} -VAE, also trained on MNIST. We see a stark contrast. $i\mathcal{P}$ -VAE features
 487 are Gabor-like, while \mathcal{P} -VAE features look like digits or strokes of the digits. While previous work
 488 highlighted strokes as the compositional subcomponents of digits (Lee et al., 2007), $i\mathcal{P}$ -VAE learns
 489 an even more general code that generalized to cropped, grey scaled natural images (Fig. 5).
 490



491
 492
 493
 494
 495
 496
 497 Figure 5: Evaluating generalization from models trained on MNIST digits to cropped, gray scaled
 498 natural images (ImageNet32) at test time. The right panel shows average reconstruction performance
 499 over inference iterations for the entire dataset. The left panels visualizes selected ground truth
 500 images compared with model reconstructions. The ai-VAE variants are unable to adapt to the new
 501 domain, whereas sa-VAE can capture more details. $i\mathcal{P}$ -VAE outperforms the alternatives, and its
 502 reconstructions are shown to maintain the semantic information of ground truth images.
 503

504 Since both $i\mathcal{P}$ -VAE and \mathcal{P} -VAE are spiking models, this result suggests that the difference lies in
 505 the inference algorithm: $i\mathcal{P}$ -VAE is iterative and adaptive; whereas, \mathcal{P} -VAE is one-shot amortized.
 506 Overall, our experiments provide strong evidence for the utility of iterative algorithms in practical
 507 settings.
 508

509 **5 DISCUSSION AND CONCLUSIONS**

510
 511 In this work, we introduce the $i\mathcal{P}$ -VAE, which is a spiking neural network that maximizes ELBO,
 512 while performing Bayesian posterior updates through membrane potential dynamics. Empirically,
 513 $i\mathcal{P}$ -VAE exhibits outstanding adaptability and robustness to OOD samples, while being able to
 514 dynamically trade off compute and performance. It outperforms amortized versions and recent
 515 iterative inference VAEs on every task we tested while using substantially fewer parameters.
 516

517 $i\mathcal{P}$ -VAE results directly from the choice of Poisson in the ELBO and it avoids many of the problems
 518 with predictive coding. First, there is no population-wide prediction signal, only a feedforward
 519 receptive field and recurrent terms. Second, neurons only communicate through spikes and all
 520 dynamics are private on the membrane potential. And finally, additive terms in the membrane
 521 potential appear as gains in the spike rate, which avoids negative rates, and is more consistent with
 522 real neurons Gilbert & Li (2013).

523 We believe $i\mathcal{P}$ -VAE is well positioned for a neuromorphic implementation. The recent rise of
 524 neuromorphic hardware as an avenue for performance improvements requires new algorithms that
 525 can make use of its architecture (Schuman et al., 2022). We found that $i\mathcal{P}$ -VAE with a linear decoder
 526 reduces to a spiking LCA, addressing the performance gap noted by Vafaii et al. (2024). Both
 527 algorithms share key features: sparsity, recurrence, and parameter efficiency. Since LCA has been
 528 implemented as an SNN (Zylberberg et al., 2011) and on neuromorphic hardware (Du et al., 2024),
 we expect the same for $i\mathcal{P}$ -VAE.

529 In summary, the choice of Poisson in the ELBO results in a spiking neural network, $i\mathcal{P}$ -VAE,
 530 that performs iterative Bayesian inference. This lays the groundwork for a prescriptive theoretical
 531 framework for building brain-like generative models that can leverage neuromorphic hardware.
 532

533 **Limitations and future work.** In our experiments, we tested the simplest version of $i\mathcal{P}$ -VAE,
 534 showing the practical benefits of the derived theory. There are a few avenues that we did not test,
 535 and we think are exciting for future work. The design of a hierarchical model is a natural extension
 536 for brain-like algorithm, especially given evidence that hierarchical VAE are more aligned to the
 537 brain (Vafaii et al., 2023). In addition, training and evaluating on nonstationary sequences like videos
 538 would be a straightforward extension, as we derived the theory with this in mind. When attempting
 539 to use such sequences, it may also be beneficial to explore more sophisticated forward-predictive
 models that “evolve” current posteriors to future priors.

REFERENCES

- 540
541
542 Larry F Abbott, JA Varela, Kamal Sen, and SB Nelson. Synaptic depression and cortical gain control.
543 *Science*, 275(5297):221–224, 1997.
- 544 Laurence Aitchison and Máté Lengyel. With or without you: predictive coding and bayesian inference
545 in the brain. *Current opinion in neurobiology*, 46:219–227, 2017.
- 546
547 Alhazen. *Book of optics (Kitab Al-Manazir)*. 1011–1021 AD.
- 548
549 Christina Allen and Charles F Stevens. An evaluation of causes for unreliability of synaptic transmis-
550 sion. *Proceedings of the National Academy of Sciences*, 91(22):10380–10383, 1994.
- 551
552 Mel Andrews. The math is not the territory: navigating the free energy principle. *Biology &*
553 *Philosophy*, 36(3):30, 2021.
- 554
555 Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul,
556 Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient
557 descent. *Advances in neural information processing systems*, 29, 2016.
- 558
559 Bruno B Averbeck, Peter E Latham, and Alexandre Pouget. Neural correlations, population coding
560 and computation. *Nature reviews neuroscience*, 7(5):358–366, 2006.
- 561
562 Jimmy Ba, Geoffrey E Hinton, Volodymyr Mnih, Joel Z Leibo, and Catalin Ionescu. Using fast
563 weights to attend to the recent past. *Advances in neural information processing systems*, 29, 2016.
- 564
565 Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. Deep equilibrium models. In H. Wal-
566 lach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Ad-
567 vances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,
568 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/
569 file/01386bd6d8e091c2ab4c7c7de644d37b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/01386bd6d8e091c2ab4c7c7de644d37b-Paper.pdf).
- 570
571 Shaojie Bai, Vladlen Koltun, and J. Zico Kolter. Multiscale deep equilibrium models. In
572 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-
573 ral Information Processing Systems*, volume 33, pp. 5238–5250. Curran Associates, Inc.,
574 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/
575 file/3812f9a59b634c2a9c574610eaba5bed-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/3812f9a59b634c2a9c574610eaba5bed-Paper.pdf).
- 576
577 Arpit Bansal, Eitan Borgnia, Hong-Min Chu, Jie S. Li, Hamid Kazemi, Furong Huang, Micah
578 Goldblum, Jonas Geiping, and Tom Goldstein. Cold diffusion: Inverting arbitrary image transforms
579 without noise, 2022.
- 580
581 Andre M Bastos, W Martin Usrey, Rick A Adams, George R Mangun, Pascal Fries, and Karl J
582 Friston. Canonical microcircuits for predictive coding. *Neuron*, 76(4):695–711, 2012. doi:
583 10.1016/j.neuron.2012.10.038.
- 584
585 Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4.
586 Springer, 2006.
- 587
588 David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians.
589 *Journal of the American statistical Association*, 112(518):859–877, 2017.
- 590
591 Martin Boerlin, Christian K Machens, and Sophie Denève. Predictive coding of dynamical variables
592 in balanced spiking networks. *PLoS computational biology*, 9(11):e1003258, 2013.
- 593
594 Victor Boutin, Aïmen Zerroug, Minju Jung, and Thomas Serre. Iterative vae as a predictive brain
595 model for out-of-distribution generalization. *arXiv preprint arXiv:2012.00557*, 2020.
- 596
597 Daniel A Butts, Yuwei Cui, and Alexander RR Casti. Nonlinear computations shaping temporal
598 processing of precortical vision. *Journal of Neurophysiology*, 116(3):1344–1357, 2016.
- 599
600 William H Calvin and CHARLES F Stevens. Synaptic noise and other sources of randomness in
601 motoneuron interspike intervals. *Journal of neurophysiology*, 31(4):574–587, 1968.

- 594 Matteo Carandini. Amplification of trial-to-trial response variability by neurons in visual cortex.
595 *PLoS biology*, 2(9):e264, 2004.
596
- 597 Stanley H. Chan. Tutorial on diffusion models for imaging and vision. 2024. URL <https://arxiv.org/abs/2403.18103>.
598
- 599 Michael Chang, Thomas L. Griffiths, and Sergey Levine. Object representations as fixed points:
600 Training iterative refinement algorithms with implicit differentiation. In Alice H. Oh, Alekh Agar-
601 wal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing*
602 *Systems*, 2022. URL <https://openreview.net/forum?id=-5rFUTO2NWe>.
603
- 604 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary
605 differential equations. *Advances in neural information processing systems*, 31, 2018.
- 606 Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an
607 alternative to the cifar datasets. *arXiv preprint arXiv:1707.08819*, 2017.
608
- 609 Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and
610 Yoshua Bengio. A recurrent latent variable model for sequential data. In C. Cortes,
611 N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural In-*
612 *formation Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL
613 [https://proceedings.neurips.cc/paper_files/paper/2015/file/](https://proceedings.neurips.cc/paper_files/paper/2015/file/b618c3210e934362ac261db280128c22-Paper.pdf)
614 [b618c3210e934362ac261db280128c22-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/b618c3210e934362ac261db280128c22-Paper.pdf).
- 615 Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. Emnist: an extension of
616 mnist to handwritten letters. *arXiv preprint arXiv:1702.05373*, 2017.
617
- 618 Chris Cremer, Xuechen Li, and David Duvenaud. Inference suboptimality in variational autoencoders.
619 In *International Conference on Machine Learning*, pp. 1078–1086. PMLR, 2018.
- 620 Peter Dayan and Laurence F Abbott. *Theoretical neuroscience: computational and mathematical*
621 *modeling of neural systems*. MIT press, 2005.
622
- 623 Peter Dayan, Geoffrey E Hinton, Radford M Neal, and Richard S Zemel. The helmholtz machine.
624 *Neural Computation*, 7(5):889–904, 1995. doi: 10.1162/neco.1995.7.5.889.
- 625 AF Dean. The variability of discharge of simple cells in the cat striate cortex. *Experimental Brain*
626 *Research*, 44(4):437–440, 1981.
627
- 628 Mauricio Delbracio and Peyman Milanfar. Inversion by direct iteration: An alternative to denoising
629 diffusion for image restoration, 2024. URL <https://arxiv.org/abs/2303.11435>.
630
- 631 Xuexing Du, Zhong-qi K Tian, Songting Li, and Douglas Zhou. A generalized spiking locally
632 competitive algorithm for multiple optimization problems. *arXiv preprint arXiv:2407.03930*, 2024.
- 633 Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *Advances*
634 *in Neural Information Processing Systems*, 32, 2019.
635
- 636 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of
637 deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- 638 Nicolas Fourcaud-Trocmé, David Hansel, Carl Van Vreeswijk, and Nicolas Brunel. How spike gener-
639 ation mechanisms determine the neuronal response to fluctuating inputs. *Journal of neuroscience*,
640 23(37):11628–11640, 2003.
- 641 Karl Friston. A theory of cortical responses. *Philosophical transactions of the Royal Society B:*
642 *Biological Sciences*, 360(1456):815–836, 2005. doi: 10.1098/rstb.2005.1622.
- 643 Karl Friston. The free-energy principle: a rough guide to the brain? *Trends in cognitive sciences*, 13
644 (7):293–301, 2009.
- 645 Karl Friston. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2):
646 127–138, 2010. doi: 10.1038/nrn2787.
647

- 648 Samuel Gershman and Noah Goodman. Amortized inference in probabilistic reasoning. In
649 *Proceedings of the annual meeting of the cognitive science society*, volume 36, 2014. URL
650 <https://escholarship.org/uc/item/34j1h7k5>.
- 651 Samuel J Gershman. What does the free energy principle tell us about the brain? *arXiv preprint*
652 *arXiv:1901.07945*, 2019.
- 653 Charles D Gilbert and Wu Li. Top-down influences on visual processing. *Nature Reviews Neuro-*
654 *science*, 14(5):350–363, 2013.
- 655 Robbe LT Goris, J Anthony Movshon, and Eero P Simoncelli. Partitioning neuronal variability.
656 *Nature neuroscience*, 17(6):858–865, 2014.
- 657 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
658 *preprint arXiv:2312.00752*, 2023.
- 659 Geoffrey E Hinton and David C Plaut. Using fast weights to deblur old memories. In *Proceedings of*
660 *the ninth annual conference of the Cognitive Science Society*, pp. 177–186, 1987.
- 661 Geoffrey E Hinton, Peter Dayan, Brendan J Frey, and Radford M Neal. The” wake-sleep” algorithm
662 for unsupervised neural networks. *Science*, 268(5214):1158–1161, 1995.
- 663 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
664 *neural information processing systems*, 33:6840–6851, 2020.
- 665 Arthur Hobson. A new theorem of information theory. *Journal of Statistical Physics*, 1:383–391,
666 1969.
- 667 Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference.
668 *Journal of Machine Learning Research*, 2013.
- 669 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are uni-
670 versal approximators. *Neural Networks*, 2(5):359–366, 1989. ISSN 0893-6080. doi: [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8). URL <https://www.sciencedirect.com/science/article/pii/0893608089900208>.
- 671 Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural
672 networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):
673 5149–5169, 2021.
- 674 Kazuki Irie, Imanol Schlag, Róbert Csordás, and Jürgen Schmidhuber. Going beyond linear trans-
675 formers with recurrent fast weight programmers. *Advances in neural information processing*
676 *systems*, 34:7703–7717, 2021.
- 677 Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization
678 in diffusion models arises from geometry-adaptive harmonic representation. In *The Twelfth*
679 *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ANvmVS2Yr0>.
- 680 R. E. Kalman. A New Approach to Linear Filtering and Prediction Problems. *Journal of Basic*
681 *Engineering*, 82(1):35–45, 03 1960. ISSN 0021-9223. doi: 10.1115/1.3662552. URL <https://doi.org/10.1115/1.3662552>.
- 682 Eric R Kandel, James H Schwartz, Thomas M Jessell, Steven Siegelbaum, A James Hudspeth, Sarah
683 Mack, et al. *Principles of neural science*, volume 4. McGraw-hill New York, 2000.
- 684 Georg B Keller and Thomas D Mrsic-Flogel. Predictive processing: a canonical cortical computation.
685 *Neuron*, 100(2):424–435, 2018.
- 686 Minyoung Kim and Vladimir Pavlovic. Reducing the amortization gap in variational autoencoders:
687 A bayesian random function approach. *arXiv preprint arXiv:2102.03151*, 2021.
- 688 Yoon Kim, Sam Wiseman, Andrew Miller, David Sontag, and Alexander Rush. Semi-amortized
689 variational autoencoders. In *International Conference on Machine Learning*, pp. 2678–2687.
690 PMLR, 2018.

- 702 Diederik P Kingma and Ruiqi Gao. Understanding diffusion objectives as the ELBO with simple
703 data augmentation. In *Thirty-seventh Conference on Neural Information Processing Systems, 2023*.
704 URL <https://openreview.net/forum?id=NnMEadcdyD>.
705
- 706 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. 2014.
- 707 Naoki Kogo and Chris Trengove. Is predictive coding theory articulated enough to be testable?, 2015.
708
- 709 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
710 convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Wein-
711 berger (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Asso-
712 ciates, Inc., 2012. URL [https://proceedings.neurips.cc/paper_files/paper/
713 2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
- 714 Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning
715 through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
716
- 717 Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, Fugie Huang, et al. A tutorial on energy-based
718 learning. *Predicting structured data*, 1(0), 2006.
- 719 Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*.
720 Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
721
- 722 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444,
723 2015. doi: 10.1038/nature14539.
- 724 Honglak Lee, Chaitanya Ekanadham, and Andrew Ng. Sparse deep belief net model for visual area
725 v2. *Advances in neural information processing systems*, 20, 2007.
726
- 727 Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold,
728 Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot
729 attention. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-
730 vances in Neural Information Processing Systems*, volume 33, pp. 11525–11538. Curran Asso-
731 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/
732 2020/file/8511df98c02ab60aealb2356c013bc0f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/8511df98c02ab60aealb2356c013bc0f-Paper.pdf).
- 733 William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video
734 prediction and unsupervised learning. In *International Conference on Learning Representations*,
735 2017. URL <https://openreview.net/forum?id=Blawdt9xe>.
- 736 Calvin Luo. Understanding diffusion models: A unified perspective. arxiv 2022. *arXiv preprint*
737 *arXiv:2208.11970*, 2022.
738
- 739 Laurin Luttmann and Paolo Mercorelli. Comparison of backpropagation and kalman filter-based
740 training for neural networks. In *2021 25th International Conference on System Theory, Control
741 and Computing (ICSTCC)*, pp. 234–241, 2021. doi: 10.1109/ICSTCC52150.2021.9607274.
- 742 Zachary F Mainen and Terrence J Sejnowski. Reliability of spike timing in neocortical neurons.
743 *Science*, 268(5216):1503–1506, 1995.
744
- 745 Joe Marino, Yisong Yue, and Stephan Mandt. Iterative amortized inference. In Jennifer Dy and
746 Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*,
747 volume 80 of *Proceedings of Machine Learning Research*, pp. 3403–3412. PMLR, 7 2018. URL
748 <https://proceedings.mlr.press/v80/marino18a.html>.
- 749 Joseph Marino. Predictive coding, variational autoencoders, and biological connections. *Neural
750 Computation*, 34(1):1–44, 2022. doi: 10.1162/neco.a.01458.
- 751 Joseph Marino, Alexandre Piché, Alessandro Davide Ialongo, and Yisong Yue. Iterative amortized
752 policy optimization. *Advances in Neural Information Processing Systems*, 34:15667–15681, 2021.
753
- 754 Fabian A Mikulasch, Lucas Rudelt, Michael Wibral, and Viola Priesemann. Where is the error?
755 hierarchical predictive coding through dendritic error computation. *Trends in Neurosciences*, 46
(1):45–59, 2023.

- 756 Beren Millidge, Anil K. Seth, and Christopher L. Buckley. Predictive coding: a theoretical and
757 experimental review. *CoRR*, abs/2107.12979, 2021a. URL [https://arxiv.org/abs/2107.](https://arxiv.org/abs/2107.12979)
758 [12979](https://arxiv.org/abs/2107.12979).
- 759 Beren Millidge, Alexander Tschantz, Anil Seth, and Christopher Buckley. Neural kalman filtering,
760 2021b. URL <https://arxiv.org/abs/2102.10021>.
- 761 Beren Millidge, Tommaso Salvatori, Yuhang Song, Rafał Bogacz, and Thomas Lukasiewicz. Predictive
762 coding: Towards a future of deep learning beyond backpropagation? In *International Joint*
763 *Conference on Artificial Intelligence*, 2022. doi: 10.24963/ijcai.2022/774.
- 764 Beren Millidge, Mufeng Tang, Mahyar Osanlouy, Nicol S Harper, and Rafal Bogacz. Predictive
765 coding networks for temporal prediction. *PLOS Computational Biology*, 20(4):e1011183, 2024.
- 766 Sreyas Mohan, Joshua L Vincent, Ramon Manzorro, Peter Crozier, Carlos Fernandez-Granda,
767 and Eero Simoncelli. Adaptive denoising via gaintuning. In M. Ranzato, A. Beygelz-
768 imer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural In-*
769 *formation Processing Systems*, volume 34, pp. 23727–23740. Curran Associates, Inc.,
770 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/](https://proceedings.neurips.cc/paper_files/paper/2021/file/c7558e9d1f956b016d1fdb7ea132378-Paper.pdf)
771 [file/c7558e9d1f956b016d1fdb7ea132378-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/c7558e9d1f956b016d1fdb7ea132378-Paper.pdf).
- 772 David Mumford. On the computational architecture of the neocortex: Ii the role of cortico-cortical
773 loops. *Biological Cybernetics*, 66(3):241–251, 1992. doi: 10.1007/BF00198477.
- 774 Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning
775 a sparse code for natural images. *Nature*, 381(6583):607–609, 1996. doi: 10.1038/381607a0.
- 776 Bruno A Olshausen and David J Field. Sparse coding of sensory inputs. *Current opinion in*
777 *neurobiology*, 14(4):481–487, 2004. doi: 10.1016/j.conb.2004.07.007.
- 778 Nicholas J Priebe, Ferenc Mechler, Matteo Carandini, and David Ferster. The contribution of spike
779 threshold to the dichotomy of cortical simple and complex cells. *Nature neuroscience*, 7(10):
780 1113–1122, 2004. doi: 10.1038/nn1310.
- 781 Yuhui Quan, Mingqin Chen, Tongyao Pang, and Hui Ji. Self2self with dropout: Learning self-
782 supervised denoising from single image. In *IEEE/CVF Conference on Computer Vision and*
783 *Pattern Recognition (CVPR)*, June 2020.
- 784 Rajesh PN Rao and Dana H Ballard. Predictive coding in the visual cortex: a functional interpretation
785 of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87, 1999. doi:
786 10.1038/4580.
- 787 Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and
788 approximate inference in deep generative models. In *International Conference on Machine*
789 *Learning*, pp. 1278–1286. PMLR, 2014. URL [https://proceedings.mlr.press/v32/](https://proceedings.mlr.press/v32/rezende14.html)
790 [rezende14.html](https://proceedings.mlr.press/v32/rezende14.html).
- 791 Fred Rieke, David Warland, Rob de Ruyter Van Steveninck, and William Bialek. *Spikes: exploring*
792 *the neural code*. MIT press, 1999.
- 793 Christopher J Rozell, Don H Johnson, Richard G Baraniuk, and Bruno A Olshausen. Sparse coding
794 via thresholding and local competition in neural circuits. *Neural Computation*, 20(10):2526–2563,
795 2008. doi: 10.1162/neco.2008.03-07-486.
- 796 Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Pearson, 2016.
- 797 Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent
798 networks. *Neural Computation*, 4(1):131–139, 1992.
- 799 C. D. Schuman, S. R. Kulkarni, M. Parsa, J. P. Mitchell, P. Date, and B. Kay. Opportunities for
800 neuromorphic computing algorithms and applications. *Nature Computational Science*, 2022. doi:
801 10.1038/s43588-022-00223-2.
- 802
- 803
- 804
- 805
- 806
- 807
- 808
- 809

- 810 Terrence J Sejnowski. The unreasonable effectiveness of deep learning in artificial intelligence.
811 *Proceedings of the National Academy of Sciences*, 117(48):30033–30038, 2020.
812
- 813 Michael N Shadlen and William T Newsome. The variable discharge of cortical neurons: implications
814 for connectivity, computation, and information coding. *Journal of neuroscience*, 18(10):3870–3896,
815 1998.
- 816 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
817 learning using nonequilibrium thermodynamics. In *International conference on machine learning*,
818 pp. 2256–2265. PMLR, 2015.
- 819 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
820 *Advances in neural information processing systems*, 32, 2019.
821
- 822 Mandyam Veerambudi Srinivasan, Simon Barry Laughlin, and Andreas Dubs. Predictive coding:
823 a fresh view of inhibition in the retina. *Proceedings of the Royal Society of London. Series B.*
824 *Biological Sciences*, 216(1205):427–459, 1982. doi: 10.1098/rspb.1982.0085.
825
- 826 Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training
827 with self-supervision for generalization under distribution shifts. In *International conference on*
828 *machine learning*, pp. 9229–9248. PMLR, 2020.
- 829 Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei
830 Chen, Xiaolong Wang, Sanmi Koyejo, Tatsunori Hashimoto, and Carlos Guestrin. Learning to
831 (learn at test time): Rnns with expressive hidden states, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2407.04620)
832 [abs/2407.04620](https://arxiv.org/abs/2407.04620).
- 833 Malvin C Teich. Fractal character of the auditory neural spike train. *IEEE Transactions on Biomedical*
834 *Engineering*, 36(1):150–160, 1989.
835
- 836 Michael Teti. Lca-pytorch. [Computer Software] [https://doi.org/10.11578/dc.](https://doi.org/10.11578/dc.20230728.4)
837 [20230728.4](https://doi.org/10.11578/dc.20230728.4), jun 2023. URL <https://doi.org/10.11578/dc.20230728.4>.
838
- 839 David J Tolhurst, J Anthony Movshon, and Andrew F Dean. The statistical reliability of signals in
840 single neurons in cat and monkey visual cortex. *Vision research*, 23(8):775–785, 1983.
- 841 Margaret Trautner, Gabriel Margolis, and Sai Ravela. Informative neural ensemble kalman learning,
842 2020. URL <https://arxiv.org/abs/2008.09915>.
843
- 844 Wilson Truccolo, Uri T Eden, Matthew R Fellows, John P Donoghue, and Emery N Brown. A point
845 process framework for relating neural spiking activity to spiking history, neural ensemble, and
846 extrinsic covariate effects. *Journal of neurophysiology*, 93(2):1074–1089, 2005.
- 847 Hadi Vafaii, Jacob L. Yates, and Daniel A. Butts. Hierarchical VAEs provide a normative account
848 of motion processing in the primate brain. In *Thirty-seventh Conference on Neural Information*
849 *Processing Systems*, 2023. URL <https://openreview.net/forum?id=1wOkHN9JK8>.
850
- 851 Hadi Vafaii, Dekel Galor, and Jacob L. Yates. Poisson variational autoencoder. 2024. URL
852 <https://arxiv.org/abs/2405.14473>.
- 853 J Hans Van Hateren and Arjen van der Schaaf. Independent component filters of natural images
854 compared with simple cells in primary visual cortex. *Proceedings of the Royal Society of London.*
855 *Series B: Biological Sciences*, 265(1394):359–366, 1998.
856
- 857 Nicolaas Godfried Van Kampen. *Stochastic processes in physics and chemistry*, volume 1. Elsevier,
858 1992.
- 859 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
860 Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon,
861 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett
862 (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
863 Inc., 2017. URL [https://papers.nips.cc/paper_files/paper/2017/hash/](https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html)
[3f5ee243547dee91fbd053c1c4a845aa-Abstract.html](https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html).

- 864 Hermann Von Helmholtz. *Handbuch der physiologischen Optik*, volume 9. Voss, 1867. URL
865 <https://archive.org/details/handbuchderphysi00helm>.
866
- 867 Kevin S Walsh, David P McGovern, Andy Clark, and Redmond G O’Connell. Evaluating the
868 neurophysiological evidence for predictive processing as a model of perception. *Annals of the new*
869 *York Academy of Sciences*, 1464(1):242–268, 2020.
- 870 Alison I Weber and Jonathan W Pillow. Capturing the dynamical repertoire of single neurons with
871 generalized linear models. *Neural computation*, 29(12):3260–3289, 2017.
872
- 873 Haiguang Wen, Kuan Han, Junxing Shi, Yizhen Zhang, Eugenio Culurciello, and Zhongming Liu.
874 Deep predictive coding network for object recognition. In *International conference on machine*
875 *learning*, pp. 5266–5275. PMLR, 2018.
- 876 B. Widrow. *Adaptive "adaline" Neuron Using Chemical "memistors."*. 1960. URL <https://books.google.com/books?id=Yc4EAAAAIAAJ>.
877
878
- 879 Bernard Widrow and Samuel D. Stearns. *Adaptive Signal Processing*. Prentice-Hall PTR, 1985.
- 880 Robert Wilson and Leif Finkel. A neural implementation of the kalman filter. In
881 Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta (eds.), *Ad-*
882 *vances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc.,
883 2009. URL [https://proceedings.neurips.cc/paper_files/paper/2009/](https://proceedings.neurips.cc/paper_files/paper/2009/file/6d0f846348a856321729a2f36734d1a7-Paper.pdf)
884 [file/6d0f846348a856321729a2f36734d1a7-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2009/file/6d0f846348a856321729a2f36734d1a7-Paper.pdf).
885
- 886 Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang,
887 Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and
888 applications, 2024. URL <https://arxiv.org/abs/2209.00796>.
- 889 Han Yu, Jiashuo Liu, Xingxuan Zhang, Jiayun Wu, and Peng Cui. A survey on evaluation of
890 out-of-distribution generalization. *arXiv preprint arXiv:2403.01874*, 2024.
891
- 892 Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A
893 survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4396–4415, 2022.
894
- 895 Joel Zylberberg, Jason Timothy Murphy, and Michael Robert DeWeese. A sparse coding model with
896 synaptically local plasticity and spiking neurons can account for the diverse shapes of v1 simple
897 cell receptive fields. *PLoS computational biology*, 7(10):e1002250, 2011.
898

899 A EXPERIMENT DETAILS

900
901 In our comparisons to previous work, we utilized the code accompanied with sa-VAE (Kim et al.
902 (2018)), ai-VAE (Marino et al. (2018)), and \mathcal{P} -VAE Vafaii et al. (2024). Across models where code
903 was provided, we trained using the same train/validation split, and without changing the parameters
904 in the code unless we specify otherwise. For the locally competitive algorithm (LCA) baseline, we
905 used the library lca-pytorch (Teti, 2023) to replicate the analysis from Vafaii et al. (2024).

906 Since the code for sa-VAE was limited to a Bernoulli observation model, we adapted it for compati-
907 bility to Gaussian by removing the sigmoid in the decoder and replacing its reconstruction loss with
908 MSE (for the van Hateren dataset). For sa-VAE, only Omniglot parameters were provided, with
909 default batch size of 50, and default number of epochs of 100. We trained it on Omniglot with default
910 parameters, on van Hateren for 100 epochs and batch size 200, on MNIST for 32 epochs and batch
911 size 50, and EMNIST for 16 epochs and batch size 50, adjusting for the size and complexity of
912 datasets.

913 The codebase for ai-VAE included parameters for both Bernoulli and Gaussian observation models,
914 and we use them accordingly. We used their MNIST configuration for MNIST, EMNIST, and
915 Omniglot. We used their CIFAR configuration for van Hateren, except for increasing batch size
916 to 200 (van Hateren is much smaller spatially). For training the ai-VAE single-level model on van
917 Hateren, we matched the latent dimension to all other van Hateren models (512 dims instead of 1024
from the CIFAR configuration). The number of epochs in the ai-VAE code base is hardcoded to 2000,

918 but we stopped the models between 780 and 2000 epochs when the loss converged. We found that
 919 the training code occasionally resulted in nans, requiring rerunning the training from the checkpoint.
 920 In one case, the hierarchical van Hateren model, the training was unable to proceed past 61 epochs
 921 without stopping due to nans.

922 We obtained the \mathcal{P} -VAE code upon request from the authors and used the default parameters as
 923 described in the appendix of Vafaii et al. (2024).
 924

925 B ARE REAL NEURONS TRULY POISSON?

926 In this section, we discuss empirical and theoretical observations from neuroscience that support our
 927 Poisson assumption.

928 “Poisson-like” noise in neuroscience has a long history. It begins with observations that neurons do not
 929 fire the same sequence of spikes to repeated presentations of the same input and that the variance is
 930 proportional to the mean (Tolhurst et al., 1983; Dean, 1981) and was followed by the observation that
 931 for short counting windows, that proportion is 1 (Teich, 1989; Shadlen & Newsome, 1998; Averbeck
 932 et al., 2006; Rieke et al., 1999; Dayan & Abbott, 2005). Larger windows and higher visual areas are
 933 notably super-Poisson, but that can be attributed to a modulation of the rate of an inhomogeneous
 934 Poisson process (Goris et al., 2014). In other words, neurons are conditionally Poisson, not marginally
 935 Poisson (Truccolo et al., 2005).
 936
 937

938 Spike-generation, it is argued, is not noisy (Mainen & Sejnowski, 1995; Calvin & Stevens, 1968), but
 939 synaptic noise (Allen & Stevens, 1994) or noise on the membrane potential can create a Poisson-like
 940 distributions of spikes (Carandini, 2004). An important caveat is that the most famous examples
 941 of precision in spike generation, Mainen & Sejnowski (1995), is well captured well by a Poisson-
 942 process Generalized linear model (Weber & Pillow, 2017), although that precision depends on the
 943 Bernoulli approximation to a Poisson process in the limit where only 0 or 1 spikes are possible. There
 944 is a widely-held misconception that precise timing cannot be produced by spike-rate models, but
 945 inhomogeneous rate models can operate at high time resolution and produce precise spiking (Butts
 946 et al., 2016).

947 Importantly, to maximize the ELBO, one has to choose an approximate posterior and prior. Because
 948 spike counts are integer and cannot be negative, Poisson is a more natural choice than Gaussian
 949 without knowing anything about neural firing statistics. Here, we found that Poisson assumption
 950 produced a prescriptive theory for neural coding. Future work might interpret this assumption at
 951 higher time resolution using inhomogeneous Poisson processes in the limit of binary spiking.
 952

953 C EXTENDED RELATED WORKS

954 C.1 DIFFUSION MODELS

955 Diffusion models have recently gained significant traction in various generative tasks, demonstrating
 956 impressive performance across applications (Yang et al., 2024; Chan, 2024). Originally introduced
 957 by Sohl-Dickstein et al. (2015), these models iteratively restore data structure by learning a reverse
 958 diffusion process. Despite the dominance of one-shot feedforward methods, the success of diffusion
 959 models highlights the ongoing relevance of iterative approaches. Several studies have sought to
 960 explain why these models perform so well in tasks like image generation. In this section, we highlight
 961 three key findings.
 962

963 First, Delbracio & Milanfar (2024) and Bansal et al. (2022) showed that fully deterministic iterative
 964 restoration methods, without diffusion theory, can match the performance of conditional diffusion
 965 models. This suggests that the strength of diffusion models lies, at least partially, in their iterative
 966 nature.
 967

968 Second, Kingma & Gao (2023) revealed that despite their distinct loss functions, diffusion models
 969 essentially optimized the ELBO objective (identical under certain conditions), particularly in noise-
 970 perturbed data settings. This adds further support to the idea that diffusion models succeed not
 971 because of their diffusion-specific properties, but because they are iterative, aligning them closely
 with $i\mathcal{P}$ -VAE, which also optimizes an ELBO-like objective through iterative processes.

972 Finally, Kadkhodaie et al. (2024) found that diffusion models operate by applying a shrinkage
973 operation on an adaptive basis, a fundamental concept in signal processing. In methods like sparse
974 coding, this is represented by an $L1$ regularization term. Similarly, an $L1$ -like term appears in
975 $i\mathcal{P}$ -VAE, which also uses integer representations to zero out small values. These similarities suggest a
976 strong connection between $i\mathcal{P}$ -VAE and diffusion models, presenting an exciting direction for future
977 research.

978 979 C.2 ADAPTIVE FILTERS 980

981 Adaptive filters are a widely used class of algorithms capable of modeling signals with varying statis-
982 tics (Widrow & Stearns (1985)). Their applications are highly diverse, including communications,
983 control and robotics, weather prediction, and inverse problems such as denoising. Two of the most
984 popular adaptive filter classes, the Kalman filter (Kalman (1960)) and the Least mean squares (LMS)
985 filter (Widrow & Stearns (1985)), have close connections to machine learning. The LMS filter was
986 originally based on research aiming to train neural networks (Widrow (1960)). Backpropagation can
987 be understood as a generalization of the LMS filter when applied to multi-layer networks. Although
988 the Kalman filter has not had much use as a learning algorithm, a recent line of work shows that there
989 is a lot of potential benefits in doing so (Trautner et al. (2020); Luttmann & Mercorelli (2021)). Both
990 algorithms, when used in dynamic settings, encode the prediction residual (like $i\mathcal{P}$ -VAE), and can be
991 interpreted from the framework of predictive coding. More concretely, Millidge et al. (2021a) showed
992 predictive coding in the linear case corresponds to Kalman filtering, and also showed the relationship
993 between backpropagation (extension of LMS) and predictive coding. Later, Millidge et al. (2021b)
994 showed that predictive coding and Kalman filtering, although not identical in general, optimize the
995 same objective. In addition, they show a neurally plausible implementation of the Kalman filter (see
996 Wilson & Finkel (2009) for an earlier paper in this line of work).

997 In future work, it would be interesting to incorporate additional ideas from the rich literature of
998 Kalman filters. Particularly, extensions of Kalman filtering, such as the ensemble Kalman filtering,
999 tend to be better suited for nonlinear and nongaussian applications (albeit with the loss of guarantees).

1000 1001 C.3 TEST-TIME OPTIMIZATION 1002

1003 There has been a recent surge of work showing that incorporating test-time optimization leads to
1004 improved performance. One notable line of work is known as Test-Time-Training (TTT), introduced
1005 by Sun et al. (2020). TTT is a general approach for updating model parameters in test time using
1006 self-supervised learning, demonstrating increased performance and robustness. Around the same
1007 time Quan et al. (2020) introduced Self2Self, a denoising method that is only trained during test
1008 time. A follow-up to Self2Self instead optimized a per-layer gain value of a trained model Mohan
1009 et al. (2021). In a recent paper, Sun et al. (2024) extended the TTT framework to language modeling,
1010 introducing an architecture that outperforms transformers (Vaswani et al., 2017) and Mamba (Gu &
1011 Dao, 2023). The authors also showed that theoretically, transformers can be understood as a special
1012 case of their TTT algorithm. In this work, we found that $i\mathcal{P}$ -VAE can also be understood within
1013 the TTT framework. Overall, our results reveal a novel grounding of TTT within well-established
1014 theoretical concepts in neuroscience.

1015 1016 C.4 FEEDFORWARD VERSUS ITERATIVE COMPUTATION 1017

1018 Deep learning is currently the dominant paradigm in artificial intelligence (AI) research, driven
1019 largely by the success of feedforward neural networks (LeCun et al., 2015; Sejnowski, 2020).
1020 The deep learning era invoked the universal approximation theorem (Hornik et al., 1989) and
1021 emphasized parallelization of training (Krizhevsky et al., 2012; Vaswani et al., 2017) leading to an
1022 over-reliance on models that perform one-shot inference. This “unrolling” of inference diverged from
1023 the classic AI literature, which recognized the importance of iterative algorithms (Russell & Norvig,
1024 2016). Although feedforward models initially achieved remarkable results, their limitations became
1025 increasingly apparent as they struggled to generalize beyond their training distributions (Zhou et al.,
2022; Yu et al., 2024). To counter this limitation, iterative computation at test time has recently
resurfaced as a promising direction (Sun et al., 2020; 2024).

1026 Unlike feedforward models, iterative algorithms refine their predictions over multiple steps, allowing
 1027 them to adapt dynamically to new inputs. Examples include iterative amortized inference techniques
 1028 Marino et al. (2018); Kim et al. (2018), diffusion models Sohl-Dickstein et al. (2015); Ho et al.
 1029 (2020); Song & Ermon (2019), energy based models (Du & Mordatch, 2019; LeCun et al., 2006),
 1030 test-time training Sun et al. (2020; 2024), meta-learning algorithms (Andrychowicz et al., 2016; Finn
 1031 et al., 2017; Hospedales et al., 2021), neural ordinary differential equations (Chen et al., 2018), deep
 1032 equilibrium models (Bai et al., 2019; 2020), object-centric models (Locatello et al., 2020; Chang
 1033 et al., 2022), and many more. These methods have demonstrated that a dynamic, multi-step inference
 1034 process can help overcome many of the challenges faced by static models.

1035 C.5 FAST WEIGHTS

1036
 1037
 1038 In the late 1980s and early 1990s, Hinton & Plaut (1987) and Schmidhuber (1992) introduced the
 1039 concept of "fast weights" as a way to enhance the adaptability of neural networks through dynamic
 1040 memory. These innovations laid the foundation for modern models like transformers and recurrent
 1041 neural networks, significantly influencing memory-augmented architectures and iterative inference
 1042 methods. Fast weights are particularly relevant in iterative inference, where dynamic updates align
 1043 with the goal of flexible, adaptive neural computation (Ba et al., 2016; Irie et al., 2021). In our work,
 1044 the adaptive Bayesian posterior updates in $\mathbf{u}(t)$ —the membrane potential state of $i\mathcal{P}$ -VAE—closely
 1045 parallel the concept of fast weights.

1046 D DYNAMICS

1047
 1048
 1049
 1050 In this section, we will go through the derivation of the dynamics of $i\mathcal{P}$ -VAE (eq. (7) in the main
 1051 paper). Our goal is to define membrane potential updates in a way that the resulting dynamics will
 1052 minimize the ELBO loss.

1053
 1054 We begin with the general definition of the ELBO, $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$, and consider its Monte
 1055 Carlo estimate using a single sample, \mathbf{z} , drawn from the approximate posterior $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\begin{aligned}
 1056 \ell(\mathbf{x}, \mathbf{z}) &:= \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \\
 1057 &= \log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \\
 1058 &= \log p(\mathbf{x}|\mathbf{z}) + \log \frac{p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \\
 1059 &= -\text{MSE}(\mathbf{x}, \mathbf{z}) + \mathbf{r} \odot (\exp(\delta \mathbf{u}) - 1) - \mathbf{z} \odot \delta \mathbf{u}.
 \end{aligned}
 \tag{9}$$

1060
 1061
 1062
 1063 In the last line of eq. (9), we inserted our specific choice of Gaussian conditional density, resulting in
 1064 $\log p(\mathbf{x}|\mathbf{z}) = -\text{MSE}(\mathbf{x}, \mathbf{z}) = -\|\mathbf{x} - \mathbf{f}_\theta(\mathbf{z})\|^2$. We also expressed the log ratio between the prior
 1065 and approximate posterior distributions, both modeled as Poisson, as in the case in $i\mathcal{P}$ -VAE.

1066
 1067
 1068
 1069 Next, we take the partial derivative of $\ell(\mathbf{x}, \mathbf{z})$ w.r.t the samples \mathbf{z} and keep only the first order terms.
 1070 This results in:

$$\frac{\partial}{\partial \mathbf{z}} \ell(\mathbf{x}, \mathbf{z}) \approx -\frac{\partial}{\partial \mathbf{z}} \text{MSE}(\mathbf{x}, \mathbf{z}) - \delta \mathbf{u}.
 \tag{10}$$

1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079 If we define our posterior updates, $\delta \mathbf{u}$, to be proportional to the gradient of $\ell(\mathbf{x}, \mathbf{z})$ w.r.t the state
 variable, \mathbf{u} , we get:

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

$$\begin{aligned}
 \delta \mathbf{u} &:= \alpha \nabla_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{z}) \\
 &= \alpha \frac{\partial \mathbf{z}}{\partial \mathbf{u}} \frac{\partial}{\partial \mathbf{z}} \ell(\mathbf{x}, \mathbf{z}) \\
 &\approx -\alpha \frac{\partial \mathbf{z}}{\partial \mathbf{u}} \left[\frac{\partial}{\partial \mathbf{z}} \text{MSE}(\mathbf{x}, \mathbf{z}) + \delta \mathbf{u} \right],
 \end{aligned} \tag{11}$$

where α is a proportionality constant. We rearrange some terms to get the following update rule:

$$\delta \mathbf{u} = - \left(\frac{\alpha \partial \mathbf{z} / \partial \mathbf{u}}{1 + \alpha \partial \mathbf{z} / \partial \mathbf{u}} \right) \frac{\partial}{\partial \mathbf{z}} \text{MSE}(\mathbf{x}, \mathbf{z}). \tag{12}$$

The stochastic samples, \mathbf{z} , depend to the state variable, \mathbf{u} , through firing rates, $\mathbf{r} = \exp(\mathbf{u})$. Therefore, we have $\partial \mathbf{z} / \partial \mathbf{u} = (\partial \mathbf{z} / \partial \mathbf{r}) (\partial \mathbf{r} / \partial \mathbf{u})$. But $\partial \mathbf{r} / \partial \mathbf{u}$ is just \mathbf{r} , and if we approximate $\partial \mathbf{z} / \partial \mathbf{r}$ using the straight-through estimator, we will have $\partial \mathbf{z} / \partial \mathbf{u} \approx \mathbf{r}$. Plug this back into eq. (12) to get:

$$\delta \mathbf{u} \approx - \left(\frac{\alpha \mathbf{r}}{1 + \alpha \mathbf{r}} \right) \frac{\partial}{\partial \mathbf{z}} \text{MSE}(\mathbf{x}, \mathbf{z}). \tag{13}$$

The proportionality coefficient, $\alpha \mathbf{r} / (1 + \alpha \mathbf{r})$, can be interpreted as an adaptive learning rate that depends on the instantaneous firing rate of neurons. While this result is intriguing, in the present work we simplified our update rule by removing the proportionality coefficient. Instead, we simply used the gradient of the MSE to compute $\delta \mathbf{u}$:

$$\begin{aligned}
 \delta \mathbf{u} &\propto - \frac{\partial}{\partial \mathbf{z}} \text{MSE}(\mathbf{x}, \mathbf{z}) \\
 &= - \frac{\partial}{\partial \mathbf{z}} \|\mathbf{x} - f_{\theta}(\mathbf{x})\|^2 \\
 &\propto \frac{\partial f_{\theta}(\mathbf{z})}{\partial \mathbf{z}} \cdot (\mathbf{x}_t - f_{\theta}(\mathbf{z}_t)) \\
 &= \mathbf{J}_{\theta} \cdot \Delta_t.
 \end{aligned} \tag{14}$$

This concludes our derivation of eq. (7).