

---

# Divide-and-Conquer Predictive Coding: a Structured Bayesian Inference Algorithm

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

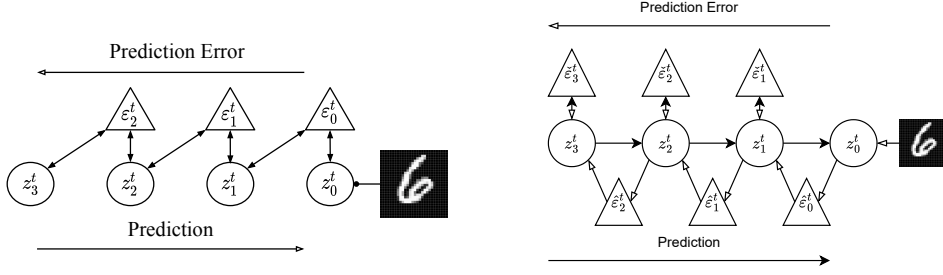
1 Unexpected stimuli induce “error” or “surprise” signals in the brain. The theory  
2 of predictive coding promises to explain these observations in terms of Bayesian  
3 inference by suggesting that the cortex implements variational inference in a proba-  
4 bilistic graphical model. However, when applied to machine learning tasks, this  
5 family of algorithms has yet to perform on par with other variational approaches in  
6 high-dimensional, structured inference problems. To address this, we introduce a  
7 novel predictive coding algorithm for structured generative models, that we call  
8 divide-and-conquer predictive coding (DCPC). DCPC differs from other formula-  
9 tions of predictive coding, as it respects the correlation structure of the generative  
10 model and provably performs maximum-likelihood updates of model parameters,  
11 all without sacrificing biological plausibility. Empirically, DCPC achieves better  
12 numerical performance than competing algorithms and provides accurate inference  
13 in a number of problems not previously addressed with predictive coding. We  
14 provide an open implementation of DCPC in Pyro on Github.

## 15 1 Introduction

16 In recent decades, the fields of cognitive science, machine learning, and theoretical neuroscience have  
17 borne witness to a flowering of successes in modeling intelligent behavior via statistical learning.  
18 Each of these fields has taken a different approach: cognitive science has studied probabilistic  
19 *inverse inference* [Chater et al., 2006, Pouget et al., 2013, Lake et al., 2017] in models of each task  
20 and environment, machine learning has employed the backpropagation of errors [Rumelhart et al.,  
21 1986, Lecun et al., 2015, Schmidhuber, 2015], and neuroscience has hypothesized that *predictive*  
22 *coding* [Srinivasan et al., 1982, Rao and Ballard, 1999, Friston, 2005, Bastos et al., 2012, Spratling,  
23 2017, Hutchinson and Barrett, 2019, Millidge et al., 2021] (PC) may explain neural activity in  
24 perceptual tasks. These approaches share in common a commitment to “deep” models, in which task  
25 processing emerges from the composition of elementary units.

26 At the computational level, probabilistic theories of perception suggest that the brain is an hypothesis  
27 testing machine, where the world is perceived via Bayesian inference [Doya, 2007]. In the PC  
28 framework, hypotheses correspond to prediction signals that flow down the cortical hierarchy to inhibit  
29 the bottom-up processing of predictable (or irrelevant) stimuli. Combining these top-down predictions  
30 with a bottom-up stimulus signal generates prediction errors, defined as the (weighted) difference  
31 between predicted and actual signals [Hoemann et al., 2017, Barrett, 2022]. Algorithmically, PC  
32 implements variational inference [Friston et al., 2006]: under some specific assumptions, a prediction  
33 error  $\varepsilon$  is the gradient of a variational free energy defined over a hierarchical Gaussian generative  
34 model, i.e.,  $\varepsilon := \nabla_{\mu} \log \mathcal{N}(\mu, \tau) = \tau(x - \mu)$ , with respect to the location parameter  $\mu$ , of a Gaussian  
35  $x \sim \mathcal{N}(\mu, \tau)$  log-density parameterized by mean  $\mu$  and precision  $\tau$ .

36 In machine learning, predictive coding algorithms have recently gained popularity for their theoretical  
37 potential to provide a more biologically plausible alternative to backpropagation for training neural



(a) Classical PC learns a mean-field approximate (b) Divide-and-conquer PC approximates the joint posterior with prediction error layers. posterior with bottom-up and recurrent errors.

Figure 1: Where classical predictive coding has layers communicate through shared error units, divide-and-conquer predictive coding separates recurrent from “bottom-up” error pathways to target complete conditional distributions rather than posterior marginal distributions.

networks [Salvatori et al., 2023, Song et al., 2024]. However, PC does not perform comparably in these tasks to backpropagation due to limitations in current formulations. First, predictive coding for gradient calculation typically models every node in the computation graph with a Gaussian, and hence fails to express many common generative models. Recent work on PC has addressed this by allowing approximating non-Gaussian energy functions with samples [Pinchetti et al., 2022]. Second, the Laplace approximation to the posterior infers only a maximum-a-posteriori (MAP) estimate and Gaussian covariance for each latent variable, keeping PC from capturing multimodal or correlated distributions. Third, this loose approximation to the posterior distribution results in inaccurate, high-variance updates to the generative model’s parameters.

In this work we propose a new algorithm, *divide-and-conquer predictive coding* (DCPC), for approximating structured target distributions with populations of Monte Carlo samples. DCPC goes beyond Gaussian assumptions, and decomposes the problem of sampling from structured targets into local coordinate updates to individual random variables. These local updates are informed by unadjusted Langevin proposals parameterized in terms of biologically plausible prediction errors. Nesting the local updates within divide-and-conquer Sequential Monte Carlo [Lindsten et al., 2017, Kuntz et al., 2024] ensures that DCPC can target any statically structured graphical model, while Theorem 2 provides a locally factorized way to learn model parameters by maximum marginal likelihood.

DCPC also provides a computational perspective on the canonical cortical microcircuit [Bastos et al., 2012, 2020, Campagnola et al., 2022] hypothesis in neuroscience. Experiments have suggested that deep laminar layers in the cortical microcircuit represent sensory imagery, while superficial laminar represent raw stimulus information [Bergmann et al., 2024]; experiments in a predictive coding paradigm specifically suggested that the deep layers represent “predictions” while the shallow layers represent “prediction errors”. This circuitry could provide the brain with its fast, scalable, generic Bayesian inference capabilities. Figure 1 compares the computational structure of DCPC with that of previous PC models. The following sections detail this work’s contributions:

- Section 3 defines the divide-and-conquer predictive coding algorithm and shows how to use it as a variational inference algorithm;
- Section 4 examines under what assumptions the cortex could plausibly implement DCPC, proving two theorems that contribute to biological plausibility;
- Section 5 demonstrates DCPC experimentally in head-to-head comparisons against recent generative models and inference algorithms from the predictive coding literature.

Section 2 will review the background for Section 3’s algorithm: the problem predictive coding aims to solve and a line of recent work addressing that problem from which this paper draws.

## 2 Background

This section reviews the background necessary to construct the divide-and-conquer predictive coding algorithm in Section 3. Let us assume we have a directed, acyclic graphical model with a joint density

split into observations  $x \in \mathbf{x}$  and latents  $z \in \mathbf{z}$ , parameterized by some  $\theta$  at each conditional density

$$p_\theta(\mathbf{x}, \mathbf{z}) := \prod_{x \in \mathbf{x}} p_\theta(x \mid \text{Pa}(x)) \prod_{z \in \mathbf{z}} p_\theta(z \mid \text{Pa}(z)), \quad (1)$$

where  $\text{Pa}(z)$  denotes the parents of the random variable  $z \in \mathbf{z}$ , while  $\text{Ch}(z)$  denotes its children.

**Empirical Bayes** *Empirical Bayes* consists of jointly estimating, in light of the data, both the parameters  $\theta^*$  and the Bayesian posterior over the latent variables  $\mathbf{z}$ , that is:

$$\theta^* = \arg \max_{\theta} p_\theta(\mathbf{x}) = \arg \max_{\theta} \int_{\mathbf{z} \in \mathbf{Z}} p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z}, \quad p_{\theta^*}(\mathbf{z} \mid \mathbf{x}) := \frac{p_{\theta^*}(\mathbf{x}, \mathbf{z})}{p_{\theta^*}(\mathbf{x})}.$$

Typically the marginal and posterior densities have no closed form, so learning and inference algorithms treat the joint distribution as a closed-form *unnormalized* density over the latent variables; its integral then gives the normalizing constant for approximation

$$\gamma_\theta(\mathbf{z}) := p_\theta(\mathbf{x}, \mathbf{z}), \quad Z_\theta := \int_{\mathbf{z} \in \mathbf{Z}} \gamma_\theta(\mathbf{z}) d\mathbf{z} = p_\theta(\mathbf{x}), \quad \pi_\theta(\mathbf{z}) := \frac{\gamma_\theta(\mathbf{z})}{Z_\theta}.$$

Neal and Hinton [1998] reduced empirical Bayes to minimization of the *variational free energy*:

$$\mathcal{F}(\theta, q) := \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[ -\log \frac{\gamma_\theta(\mathbf{z})}{q(\mathbf{z})} \right] \geq -\log Z(\theta). \quad (2)$$

The ratio of densities in Equation 2 is an example of a *weight* used to approximate a distribution known only up to its normalizing constant. The *proposal* distribution  $q(\mathbf{z})$  admits tractable sampling, while the unnormalized *target* density  $\gamma_\theta(\mathbf{z})$  admits tractable, pointwise density evaluation.

**Predictive Coding** Computational neuroscientists now often hypothesize that *predictive coding* (PC) can optimize the above family of objective functionals in a local, neurally plausible way [Millidge et al., 2021, 2023]. More in detail, it is possible to define this class of algorithms as follows:

**Definition 1** (Predictive Coding Algorithm). *Consider approximate inference in a model  $p_\theta(\mathbf{x}, \mathbf{z})$  using an algorithm  $\mathcal{A}$ . Salvatori et al. [2023] calls  $\mathcal{A}$  a predictive coding algorithm if and only if:*

1. *It maximizes the model evidence  $\log p_\theta(\mathbf{x})$  by minimizing a variational free energy;*
2. *The proposal  $q(\mathbf{z}) = \prod_{z \in \mathbf{z}} q(z)$  factorizes via a mean-field approximation; and*
3. *Each proposal factor is a Laplace approximation (i.e.  $q_\mu(z) := \mathcal{N}(\mu, \Sigma(\mu))$ ).*

**Particle Algorithms** In contrast to predictive coding, particle algorithms approach empirical Bayes problems by setting the proposal to a collection of weighted particles ( $w^k, \mathbf{z}^k$ ) drawn from a sampling algorithm meeting certain conditions (see Definition 2 in Appendix B). Any proposal meeting these conditions (see Proposition 1 in Appendix B and Naesseth et al. [2015], Stites et al. [2021]) defines a free energy functional, analogous to Equation 2 in upper-bounding the model surprisal:

$$\mathcal{F}(\theta, q) := \mathbb{E}_{w, \mathbf{z} \sim q(w, \mathbf{z})} [-\log w] \implies \mathcal{F}(\theta, q) \geq -\log Z(\theta).$$

This paper builds on the particle gradient descent (PGD) algorithm of Kuntz et al. [2023], that works as follows: At each iteration  $t$ , PGD diffuses the particle cloud  $q_K(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{z}^k}(\mathbf{z})$  across the target log-density with a learning rate  $\eta$  and independent Gaussian noise; it then updates the parameters  $\theta$  by ascending the gradient of the log-likelihood, estimated by averaging over the particles. The update rules are then the following:

$$\mathbf{z}^{t+1, k} := \mathbf{z}^{t, k} + \eta \nabla_{\mathbf{z}} \log \gamma_{\theta^t}(\mathbf{z}^{t, k}) + \sqrt{2\eta} \xi^k, \quad (3)$$

$$\theta^{t+1} := \theta^t + \eta \left( \frac{1}{K} \sum_{k=1}^K \nabla_{\theta} \log \gamma_{\theta^t}(\mathbf{z}^{t+1, k}) \right). \quad (4)$$

The above equations target the joint density of an entire graphical model<sup>1</sup>. When the prior  $p_\theta(\mathbf{z})$  factorizes into many separate conditional densities, achieving high inference performance often

<sup>1</sup>Kuntz et al. [2023] also interpreted Equation 3 as an update step along the Wasserstein gradient in the space of probability measures. Appendix C extends this perspective to predictive coding of discrete random variables.

	PC	LPC	MCPC	DCPC (ours)
Generative density	Gaussian	Differentiable	Gaussian	Differentiable
Inference approximation	Laplace	Gaussian	Empirical	Empirical
Posterior conditional structure	<b>X</b>	<b>X</b>	<b>X</b>	✓

Table 1: Comparison of divide-and-conquer predictive coding (DCPC) against other predictive coding algorithms. DCPC provides the greatest flexibility: arbitrary differentiable generative models, an empirical approximation to the posterior, and sampling according to the target’s conditional structure.

requires factorizing the inference network or algorithm into conditionals as well [Webb et al., 2018]. Estimating the gradient of the entire log-joint, as in PGD and amortized inference [Dasgupta et al., 2020, Peters et al., 2024], also requires nonlocal backpropagation. To provide a generic inference algorithm for high-dimensional, structured models using only local computations, Section 3 will apply Equation 3 to sample individual random variables in a joint density, combine the coordinate updates via sequential Monte Carlo, and locally estimate gradients for model parameters via Equation 4.

### 3 Divide-and-Conquer Predictive Coding

The previous section provided a mathematical toolbox for constructing Monte Carlo algorithms based on gradient updates and a working definition of predictive coding. This section will combine those tools to generalize the above notion of predictive coding, yielding the novel *divide-and-conquer predictive coding* (DCPC) algorithm. Given a causal graphical model, DCPC will approximate the posterior with a population  $q(\mathbf{z})$  of  $K$  samples, while also learning  $\theta$  explaining the data. This will require deriving local coordinate updates and then parameterizing them in terms of prediction errors.

Let us assume we again have a causal graphical model  $p_\theta(\mathbf{x}, \mathbf{z})$  locally parameterized by  $\theta$  and factorized (as in Equation 1) into conditional densities for each  $x \in \mathbf{x}$  and  $z \in \mathbf{z}$ . DCPC then requires two hyperparameters: a learning rate  $\eta \in \mathbb{R}^+$ , and particle count  $K \in \mathbb{N}^+$ , and is initialized (at  $t = 0$ ) via a population of predictions by ancestor sampling defined as  $\mathbf{z}^0 \sim \prod_{z \in \mathbf{z}} p_\theta(z^0 \mid \text{Pa}(z^0))$ .

DCPC aims to minimize the variational free energy (Equation 2). The optimal proposal  $q_*$  for each random variable would equal, if it had closed form, the *complete conditional* density for that variable, containing all information from other random variables

$$q_*(\mathbf{z}^t \mid \mathbf{z}^{t-1}) \propto \gamma_\theta(\mathbf{z}; \mathbf{z}_{\setminus z}) = p_\theta(z \mid \text{Pa}(z)) \prod_{v \in \text{Ch}(z)} p_\theta(v \mid \text{Pa}(v)). \quad (5)$$

We observe that the prediction errors  $\varepsilon_z$  in classical predictive coding, usually defined as the precision weighted difference between predicted and actual value of a variable, can be seen as the *score function* of a Gaussian, where the score is the gradient with respect to the parameter  $z$  of the log-likelihood:

$$\varepsilon_z := \nabla_z \log \mathcal{N}(z, \tau) = \tau(x - z);$$

When given the ground-truth parameter  $z$ , the *expected* score function  $\mathbb{E}_{x \sim p(x|z)} [\nabla_z \log p(x|z)] = 0$  under the likelihood becomes zero, making score functions a good candidate for implementing predictive coding. We therefore define  $\varepsilon_z$  in DCPC as the complete conditional’s score function

$$\varepsilon_z := \nabla_z \log \gamma_\theta(\mathbf{z}; \mathbf{z}_{\setminus z}) = \nabla_z \log p_\theta(z \mid \text{Pa}(z)) + \sum_{v \in \text{Ch}(z)} \nabla_z \log p_\theta(v \mid \text{Pa}(v)). \quad (6)$$

This gradient consists of a sum of local prediction-error terms: one for the local “prior” on  $z$  and one for each local “likelihood” of a child variable. By defining the prediction error as a sum of local score functions, we write Equation 3 in terms of  $\varepsilon_z$  (Equation 6):

$$q_\eta(z^t \mid \varepsilon_z^t, z^{t-1}) := \mathcal{N}(z^{t-1} + \eta \varepsilon_z^t, 2\eta I_z).$$

The resulting proposal now targets the complete conditional density (Equation 5), simultaneously meeting the informal requirement of Definition 1 for purely local proposal computations while also “dividing and conquering” the sampling problem into lower-dimensional coordinate updates.

Since the proposal from which we can sample by predictive coding is not the optimal coordinate update, we importance weight for the true complete conditional distribution that is optimal

$$z^t \sim q_\eta(z^t \mid z^{t-1}, \varepsilon_z^t) \quad u_z^t = \frac{\gamma_{\theta^{t-1}}(z^t; \mathbf{z}_{\setminus z})}{q_\eta(z^t \mid z^{t-1}, \varepsilon_z^t)}; \quad (7)$$

---

**Algorithm 1** Divide-and-Conquer Predictive Coding for empirical Bayes
 

---

**Require:** learning rate  $\eta \in \mathbb{R}^+$ , particle count  $K \in \mathbb{N}$ , number of sweeps  $S \in \mathbb{N}$

**Require:** initial particle vector  $\mathbf{z}^0$ , initial parameters  $\theta^0$ , observations  $\mathbf{x} \in \mathcal{X}$

```

1: for  $t \in [1 \dots T]$  do                                ▷ Loop through predictive coding steps
2:   for  $s \in [1 \dots S]$  do                                ▷ Loop through Gibbs sweeps over graphical model
3:     for  $z \in \mathbf{z}$  do                                    ▷ Loop through latent variables in graphical model
4:        $\varepsilon_z \leftarrow \nabla_{\mathbf{z}} \log p_{\theta^{t-1}}(z \mid \text{Pa}(z))$                                 ▷ Local prediction error
5:        $\varepsilon_z \leftarrow \varepsilon_z + \sum_{v \in \text{Ch}(z)} \nabla_{\mathbf{z}} \log p_{\theta^{t-1}}(v \mid \text{Pa}(v))$                                 ▷ Children's prediction errors
6:        $z^t \sim q_{\eta}(z^t \mid \varepsilon_z, z^{t-1})$                                 ▷ Sample coordinate update
7:        $u_z^t \leftarrow \frac{\gamma_{\theta^{t-1}}(z^t; \mathbf{z}_{\setminus z})}{q_{\eta}(z^t \mid \varepsilon_z, z^{t-1})}$                                 ▷ Correct coordinate update by weighing
8:        $z^t \leftarrow \text{RESAMPLE}(z^t, u_z^t)$                                 ▷ Resample from true coordinate update
9:        $\hat{Z}_{\theta^{t-1}}(\mathbf{z}_{\setminus z})^t \leftarrow \frac{1}{K} \sum_{k=1}^K u_z^{t,k}$                                 ▷ Estimate coordinate update's normalizer
10:     $\mathcal{F}^t \leftarrow -\frac{1}{K} \sum_{k=1}^K \log \left( \frac{p_{\theta^{t-1}}(\mathbf{x}, \mathbf{z}^{t,k})}{\prod_{z \in \mathbf{z}} \gamma_{\theta^{t-1}}(z^t; \mathbf{z}_{\setminus z})} \prod_{z \in \mathbf{z}} \hat{Z}_{\theta^{t-1}}(\mathbf{z}_{\setminus z})^t \right)$                                 ▷ Update free energy
11:     $\theta^t \leftarrow \theta^{t-1} + \eta \frac{1}{K} \sum_{k=1}^K \nabla_{\theta^{t-1}} \log p_{\theta^{t-1}}(\mathbf{x}, \mathbf{z}^{t,k})$                                 ▷ Update parameters
12: return  $\mathbf{z}^T, \theta^T, \mathcal{F}^T$                                 ▷ Output: updated particles, parameters, free energy

```

---

139 resampling with respect to these weights corrects for discretization error, yields particles distributed  
 140 according to the true complete conditional, and estimates the complete conditional's normalizer

$$\text{RESAMPLE}(z^t, u_z^t) \sim \pi_{\theta^{t-1}}(z^t \mid \mathbf{z}_{\setminus z}), \quad \hat{Z}_{\theta^{t-1}}(\mathbf{z}_{\setminus z})^t := \frac{1}{K} \sum_{k=1}^K u_z^{t,k}.$$

141 The recursive step of “Divide and Conquer” Sequential Monte Carlo [Lindsten et al., 2017, Kuntz  
 142 et al., 2024] exploits the estimates  $\hat{Z}_{\theta^{t-1}}(\mathbf{z}_{\setminus z})^t$  to weigh the samples for the complete target density

$$w_{\theta^{t-1}}^t = \frac{p_{\theta^{t-1}}(\mathbf{x}, \mathbf{z}^t)}{\prod_{z \in \mathbf{z}} \gamma_{\theta^{t-1}}(z^t; \mathbf{z}_{\setminus z})} \prod_{z \in \mathbf{z}} \hat{Z}_{\theta^{t-1}}(\mathbf{z}_{\setminus z})^t. \quad (8)$$

143 By Proposition 1, log-transforming these weights estimates the free energy (Equation 2):

$$\mathcal{F}^t(\mathbf{z}^{t-1}, \theta^{t-1}) := \mathbb{E}_{q_*(\mathbf{z}^t \mid \mathbf{z}^{t-1})} \left[ -\log \frac{p_{\theta^{t-1}}(\mathbf{x}, \mathbf{z}^t)}{q_*(\mathbf{z}^t \mid \mathbf{z}^{t-1})} \right] \approx \mathbb{E}_q [-\log w_{\theta^{t-1}}^t].$$

144 Theorem 3 in Appendix B shows that the gradient  $\nabla_{\theta^{t-1}} \mathcal{F}^t = \mathbb{E}_q [-\nabla_{\theta^{t-1}} \log p_{\theta^{t-1}}(\mathbf{x}, \mathbf{z}^t)]$  of the  
 145 above estimator equals the expected gradient of the log-joint distribution. Descending this gradient  
 146  $\theta^t := \theta^{t-1} - \eta \nabla_{\theta^{t-1}} \mathcal{F}^t$  enables DCPC to learn model parameters  $\theta$ .

147 The above steps describe a single pass of divide-and-conquer predictive coding over a causal graphical  
 148 model. Algorithm 1 shows the complete algorithm, consisting of nested iterations over latent variables  
 149  $z \in \mathbf{z}$  (inner loop) and iterations  $t \in T$  (outer loop). DCPC satisfies criteria (1) and (2) of Definition 1,  
 150 and relaxes criterion (3) to allow gradient-based proposals beyond the Laplace assumption. As with  
 151 Pinchetti et al. [2022] and Oliviers et al. [2024], relaxing the Laplace assumption enables much  
 152 greater flexibility in approximating the model's true posterior distribution.

## 153 4 Biological plausibility

154 Different works in the literature consider different criteria for biological plausibility. This paper  
 155 follows the non-spiking predictive coding literature and considers an algorithm biologically plausible  
 156 if it performs only spatially local computations in a probabilistic graphical model [Whittington and  
 157 Bogacz, 2017], without requiring a global control of computation. However, while in the standard  
 158 literature locality is either directly defined in the objective function [Rao and Ballard, 1999], or  
 159 derived from a mean-field approximation to the joint density [Friston, 2005], showing that the updates  
 160 of the parameters of DCPC require only local information is not as trivial. To this end, in this section  
 161 we first formally show that DCPC achieves decentralized inference of latent variables  $\mathbf{z}$  (Theorem 1),  
 162 and then that also the parameters  $\theta$  are updated via local information (Theorem 2).

163 Gibbs sampling provides the most widely-used algorithm for sampling from a high-dimensional  
 164 probability distribution by local signaling. It consists of successively sampling coordinate updates

to individual nodes in the graphical model by targeting their complete conditional densities  $\pi_\theta(z \mid \mathbf{x}, \mathbf{z}_{\setminus z})$ . Theorem 1 demonstrates that DCPC’s coordinate updates approximate Gibbs sampling.

**Theorem 1** (DCPC coordinate updates sample from the true complete conditionals). *Each DCPC coordinate update (Equation 7) for a latent  $z \in \mathbf{z}$  samples from  $z$ ’s complete conditional (the normalization of Equation 5). Formally, for every measurable  $h : \mathcal{Z} \rightarrow \mathbb{R}$ , resampled expectations with respect to the DCPC coordinate update equal those with respect to the complete conditional*

$$\mathbb{E}_{z \sim q_\eta(z \mid z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u), z' \sim \text{RESAMPLE}(z, u_z)} [h(z)]] = \int_{z \in \mathcal{Z}} h(z) \pi_\theta(z \mid \mathbf{z}_{\setminus z}) dz.$$

*Proof.* See Corollary 4.1 in Appendix B. □

We follow the canonical cortical microcircuit hypothesis of predictive coding [Bastos et al., 2012, Gillon et al., 2023] or predictive routing [Bastos et al., 2020]. Consider a cortical column representing  $z \in \mathbf{z}$ ; the  $\theta$ ,  $\alpha/\beta$ , and  $\gamma$  frequency bands of neuronal oscillations [Buzsáki and Draguhn, 2004] could synchronize parallelizations (known to exist for simple Gibbs sampling in a causal graphical model [Gonzalez et al., 2011]) of the loops in Algorithm 1. From the innermost to the outermost and following the neurophysiological findings of Bastos et al. [2015], Fries [2015],  $\gamma$ -band oscillations could synchronize the bottom-up conveyance of prediction errors (lines 4-6) from L2/3 of lower cortical columns to L4 of higher columns,  $\beta$ -band oscillations could synchronize the top-down conveyance of fresh predictions (implied in passing from  $s$  to  $s + 1$  in the loop of lines 2-9) from L5/6 of higher columns to L1+L6 of lower columns, and  $\theta$ -band oscillations could synchronize complete attention-directed sampling of stimulus representations (lines 1-11). Figure 5 in Appendix A visualizes these hypotheses for how neuronal areas and connections could implement DCPC.

Biological neurons often spike to represent *changes* in their membrane voltage [Mainen and Sejnowski, 1995, Lundstrom et al., 2008, Forkosh, 2022], and some have even been tested and found to signal the temporal derivative of the logarithm of an underlying signal [Adler and Alon, 2018, Borba et al., 2021]. Theorists have also proposed models [Chavlis and Poirazi, 2021, Moldwin et al., 2021] under which single neurons could calculate gradients internally. In short, if neurons can represent probability densities, as many theoretical proposals and experiments suggest they can, then they can likely also calculate the prediction errors used in DCPC. Theorem 2 will demonstrate that given the “factorization” above, DCPC’s model learning requires only local prediction errors.

**Theorem 2** (DCPC parameter learning requires only local gradients in a factorized generative model). *Consider a graphical model factorized according to Equation 1, with the additional assumption that the model parameters  $\theta \in \Theta = \prod_{x \in \mathbf{x}} \Theta_x \times \prod_{z \in \mathbf{z}} \Theta_z$  factorize disjointly. Then the gradient  $\nabla_\theta \mathcal{F}(\theta, q)$  of DCPC’s free energy similarly factorizes into a sum of local particle averages*

$$\nabla_\theta \mathcal{F} = \mathbb{E}_q [-\nabla_\theta \log p_\theta(\mathbf{x}, \mathbf{z})] \approx - \sum_{v \in (\mathbf{x}, \mathbf{z})} \frac{1}{K} \sum_{k=1}^K \nabla_{\theta_v} \log p_{\theta_v}(v^k \mid \text{Pa}(v)^k). \quad (9)$$

*Proof.* See Proposition 5 in Appendix B. □

Our practical implementation of DCPC, evaluated in the experiments above, takes advantage of Theorem 2 to save memory by detaching samples from the automatic differentiation graph in the forward ancestor-sampling pass through the generative model.

Finally, DCPC passes from local coordinate updates to the joint target density via an importance resampling operation, requiring that implementations synchronously transmit numerical densities or log-densities for the freshly proposed particle population. While phase-locking to a cortical oscillation may make this biologically feasible, resampling then requires normalizing the weights. Thankfully, divisive normalization appears ubiquitously throughout the brain [Carandini and Heeger, 2012], as well as just the type of “winner-take-all” circuit that implements a softmax function (e.g. for normalizing and resampling importance weights) being ubiquitous in crosstalk between superficial and deep layers of the cortical column [Liu, 1999, Douglas and Martin, 2004].

## 5 Experiments

Divide-and-conquer predictive coding is not the first predictive coding algorithm to incorporate sampling into the inference process, and certainly not the first variational inference algorithm for structured graphical models. This section therefore evaluates DCPC’s performance against both

Inference algorithm	Dataset	NLL ↓	Mean Squared Error ↓
MCPC	MNIST	$144.6 \pm 0.7$	$(8.29 \pm 0.05) \times 10^{-2}$
DCPC	MNIST	<b><math>102.5 \pm 0.01</math></b>	<b><math>0.01 \pm 7.2 \times 10^{-6}</math></b>
DCPC	EMNIST	$160.8 \pm 0.05$	$3.3 \times 10^{-6} \pm 3.5 \times 10^{-9}$
DCPC	Fashion MNIST	$284.1 \pm 0.05$	$0.03 \pm 2.7 \times 10^{-5}$

Table 2: Negative log-likelihood and mean squared error for MCPC against DCPC on held-out images from the MNISTs. Means and standard deviations are taken across five random seeds.

models from the predictive coding literature and against a standard deep generative model. Each experiment holds the generative model, dataset, and hyperparameters constant except where noted.

We have implemented DCPC as a variational proposal or “guide” program in the deep probabilistic programming language Pyro [Bingham et al., 2019]; doing so enables us to compute free energy and prediction errors efficiently in graphical models involving neural networks. Since the experiments below involve minibatched subsampling of observations  $\mathbf{x} \sim \mathcal{B}$  from a dataset  $\mathcal{D} \sim p(\mathcal{D})$  of unknown distribution, we replace Equation 9 with a subsampled form (see Welling and Teh [2011] for derivation) of the variational Sequential Monte Carlo gradient estimator [Naesseth et al., 2018]

$$\nabla_{\theta} \mathcal{F} \approx |\mathcal{D}| \mathbb{E}_{\mathcal{B} \sim p(\mathcal{D})} \left[ \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}^b \in \mathcal{B}} \mathbb{E}_{(\mathbf{z}, w)^{1:K} \sim q} \left[ \log \left( \frac{1}{K} \sum_{k=1}^K w^k \right) \mid \mathbf{x}^b \right] \right]. \quad (10)$$

We optimized the free energy in all experiments using Adam [Kingma and Ba, 2014], making sure to call `detach()` after every Pyro `sample()` operation to implement the purely local gradient calculations of Theorem 2 and Equation 10. The first experiment below considers a hierarchical Gaussian model on three simple datasets. The model consists of two latent codes above an observation.

**Deep latent Gaussian models with predictive coding** Oliviers et al. [2024] brought together predictive coding with neural sampling hypotheses in a single model: Monte Carlo predictive coding (MCPC). Their inference algorithm functionally backpropagated the score function of a log-likelihood, applying Langevin proposals to sample latent variables from the posterior joint density along the way. They evaluated MCPC’s performance on MNIST with a deep latent Gaussian model [Rezende et al., 2014] (DLGM). Their model’s conditional densities consisted of nonlinearities followed by linear transformations to parameterize the mean of each Gaussian conditional, with learned covariances. Figure 2 shows that the DLGM structure already requires DCPC to respect hierarchical dependencies.

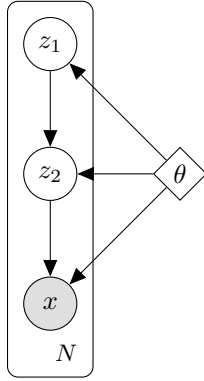


Figure 2: Hierarchical graphical model for DLGM’s.

We tested DCPC’s performance on elementary reconstruction and generation tasks by using it to train this exact generative model, changing only the likelihood from a discrete Bernoulli to a continuous Bernoulli [Loaiza-Ganem and Cunningham, 2019]. After training we evaluated with a discrete Bernoulli likelihood. Table 2 shows that in terms of both surprise (negative log evidence, with the discrete Bernoulli likelihood) and mean squared reconstruction error, DCPC enjoys better average performance with a lower standard deviation of performance, the latter by an order of magnitude. All experiments used a learning rate  $\eta = 0.1$  and  $K = 4$  particles.

Figure 3 shows an extension of this experiment to EMNIST [Cohen et al., 2017] and Fashion MNIST [Xiao et al., 2017] as well as the original MNIST, with ground-truth images in the top row and their reconstructions from DCPC-inferred latent codes below. The ground-truth images come from a 10% validation split of each data-set, on which DCPC only infers particles  $q_{K=4}(\mathbf{z})$ .

The above datasets do not typically challenge a new inference algorithm. The next experiment will thus attempt to learn representations of color images, as in the widely-used variational autoencoder [Kingma and Welling, 2013] framework, without an encoder network or amortized inference.

**Image generation with representation learning** Zahid et al. [2024] have also recently designed and evaluated Langevin predictive coding (LPC), with differences from both MCPC and DCPC. While MCPC sends prediction errors up through a hierarchical model, LPC computed as its prediction error the log-joint gradient for all latent variables in the generative model. This meant that biological plausibility, and their goal of amortizing predictive coding inference, restricted them to single-level

Algorithm	Likelihood	Resolution $\uparrow$	$\nabla_{\theta}$ -evaluations $\times$ Epochs $\downarrow$	FID $\downarrow$
PGD	$\mathcal{N}$	$32 \times 32$	$1 \times 100$	$100 \pm 2.7$
DCPC	$\mathcal{N}$	$32 \times 32$	$1 \times 100$	<b><math>89.6 \pm 0.6</math></b>
LPC	$\mathcal{DN}$	$64 \times 64$	$300 \times 15 = 4500$	120 (approximate)
DCPC	$\mathcal{DN}$	$64 \times 64$	$10 \times 450 = 4500$	<b><math>96.0 \pm 0.3</math></b>

Table 3: FID score comparisons on the CelebA dataset [Liu et al., 2015]. The score for LPC comes from Figure 2 in Zahid et al. [2024], where they ablated warm-starts and initialized from the prior.

255 decoder adapted from Higgins et al. [2017]. We evaluated with their reported discretized Gaussian  
256 likelihood, taken from Cheng et al. [2020], Ho et al. [2020], fixing the variance at  $\mathcal{DN}(\mu_{\theta}(\mathbf{z}), 0.01^2)$ .

257 We compare DCPC to LPC using the Frechet Inception Distance (FID) [Seitzer, 2020] featured in  
258 Zahid et al. [2024], holding constant the prior, neural network architecture, learning rate on  $\theta$ , and  
259 number of gradient evaluations used to train the parameters  $\theta$  and latents  $\mathbf{z}$ . Zahid et al. [2024]  
260 evaluated a variety of scenarios and reported that their training could converge quite quickly when  
261 counted in epochs, but they accumulated gradients of  $\theta$  over inference steps. We compare to the  
262 results they report after 15 epochs with 300 inference steps applied to latents initialized from the  
263 prior, equivalent to  $15 \times 300 = 4500$  gradient steps on  $\theta$  per batch, replicating their batch size of 64.  
264 Since Algorithm 1 updates  $\theta$  only in its outer loop, we set  $S = 10$  and ran DCPC for 450 epochs.  
265 Table 3 shows that DCPC outperforms LPC in apples-to-apples generative quality, though not to the  
266 point of matching state-of-the-art generative model architectures with just a decoder network.

267 Figure 4 shows reconstructed images from the validation set (left) and samples from the posterior  
268 predictive generative model (right). There is blurriness in the reconstructions, as often occurs with  
269 variational autoencoders, but DCPC training allows the network to capture background color, hair  
270 color, direction in which a face is looking, and other visual properties. Figure 4a shows reconstructions  
271 over the validation set, while Figure 4b shows samples from the predictive distribution.

272 Kuntz et al. [2023] also reported an experiment on CelebA in terms of FID score, at the lower  $32 \times 32$   
273 resolution. Since they provided both source code and an exact mathematical description, we were  
274 able to run an exact, head-to-head comparison with PGD. The line in Table 3 evaluating DCPC with  
275 PGD’s example neural architecture at the  $32 \times 32$  resolution demonstrates a significant improvement  
276 in FID for DCPC, alongside a reduction in variance across random samples.

## 277 6 Related Work

278 Pinchetti et al. [2022] expanded predictive coding beyond Gaussian generative models for the first  
279 time, applying the resulting algorithm to train variational autoencoders by variational inference and  
280 transformer architectures by maximum likelihood. DCPC, in turn, broadens predictive coding to target  
281 arbitrary probabilistic graphical models, following the broadening in Salvatori et al. [2022] to arbitrary  
282 deterministic computation graphs. DCPC follows on incremental predictive coding [Salvatori et al.,  
283 2024] in quickly alternating between updates to random variables and model parameters, giving  
284 an incremental EM algorithm [Neal and Hinton, 1998]. Finally, Zahid et al. [2024] and Oliviers  
285 et al. [2024] also recognized the analogy between predictive coding’s prediction errors and the score  
286 functions used in Langevin dynamics for continuous random variables.

287 There exists a large body of work on how neurobiologically plausible circuits could implement  
288 probabilistic inference. Classic work by Shi and Griffiths [2009] provided a biologically plausible  
289 implementation of hierarchical inference via importance sampling; DCPC proceeds from importance  
290 sampling as a foundation, while parameterizing the proposal distribution via prediction errors. Recent  
291 work by Fang et al. [2022] studied neurally plausible algorithms for sampling-based inference with

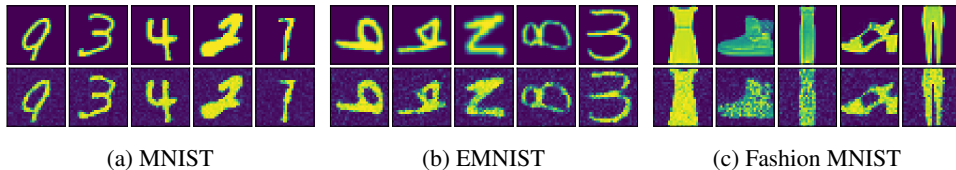
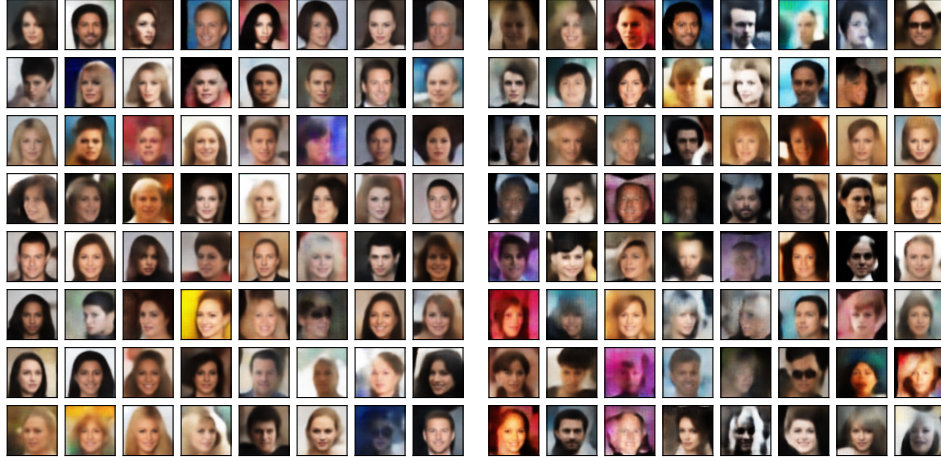


Figure 3: **Top**: images from validation sets of MNIST (left), EMNIST (middle), and Fashion MNIST (right). **Bottom**: reconstructions by deep latent Gaussian models trained with DCPC for MNIST (left), EMNIST (middle), and Fashion MNIST (right), averaging over  $K = 4$  particles. DCPC achieves quality reconstructions by inference over  $\mathbf{z}$  without training an inference network.





(a) Reconstructions of the CelebA validation set (b) Samples drawn *de novo* from the posterior predictive distribution of the trained network.

Figure 4: **Left:** reconstructions from the CelebA validation set. **Right:** samples from the generative model. DCPC achieves quality reconstructions by inference over  $\mathbf{z}$  with  $K = 16$  particles and no inference network, while the learned generative model captures variation in the data.

292 Langevin dynamics, though only for a Gaussian generative model of sparse coding. Golkar et al.  
 293 [2022] imposed a whitening constraint on a Gaussian generative model for biological plausibility.  
 294 Finally, Dong and Wu [2023] and Zahid et al. [2024] both suggest mechanisms for employing  
 295 momentum to reduce gradient noise in a biologically plausible sampling algorithm; the former  
 296 intriguingly analogize their momentum term to neuronal adaptation.

## 297 7 Conclusion

298 This paper proposed divide-and-conquer predictive coding (DCPC), an algorithm that efficiently  
 299 and scalably approximates Gibbs samplers by importance sampling; DCPC parameterizes efficient  
 300 proposals for a model’s complete conditional densities using local prediction errors. Section 4 showed  
 301 how Monte Carlo sampling can implement a form of “prospective configuration” [Song et al., 2024],  
 302 first inferring a sample from the joint posterior density (Theorem 1) and then updating the generative  
 303 model without a global backpropagation pass (Theorem 2). Experiments in Section 5 showed  
 304 that DCPC outperforms the state of the art Monte Carlo Predictive Coding from computational  
 305 neuroscience, head-to-head, on the simple generative models typically considered in theoretical  
 306 neuroscience; DCPC also outperforms the particle gradient descent algorithm of Kuntz et al. [2023]  
 307 while under the constraint of purely local computation. DCPC’s Langevin proposals admit the same  
 308 extension to constrained sample spaces as applied in Hamiltonian Monte Carlo [Brubaker et al.,  
 309 2012]; our Pyro implementation includes this extension via Pyro’s preexisting support for HMC.

310 DCPC offers a number of ways forward. Particularly, this paper employed naive Langevin proposals,  
 311 while Dong and Wu [2023], Zahid et al. [2024] applied momentum-based preconditioning to take  
 312 advantage of the target’s geometry. Yin and Ao [2006] demonstrated that gradient flows of this  
 313 general kind can also provide more efficient samplers by breaking the detailed-balance condition  
 314 necessary for the Metropolis-Hastings algorithm, motivating the choice of SMC over MCMC to  
 315 correct proposal bias. Appendix C derives a mathematical background for an extension of DCPC to  
 316 discrete random variables. Future work could follow Marino et al. [2018], Taniguchi et al. [2022] in  
 317 using a neural network to iteratively map from particles and prediction errors to proposal parameters.

### 318 7.1 Limitations

319 DCPC’s main limitations are its longer training time, and greater sensitivity to learning rates, than  
 320 state-of-the-art amortized variational inference trained end-to-end. Such limitations occur frequently  
 321 in the literature on neuroscience-inspired learning algorithms, as well as in the literature on particle-  
 322 based algorithms with no parametric form. Scaling up neuroscience-inspired algorithms is an  
 323 active area of research, and successes in this direction will naturally apply to DCPC, enabling the  
 324 training of larger models on more complex datasets by predictive coding. This work has no singular

325 ethical concerns specific only to DCPC, rather than the broader implications and responsibilities  
 326 accompanying advancements in biologically plausible learning and Bayesian inference.

## 327 **References**

- 328 Miri Adler and Uri Alon. Fold-change detection in biological systems. *Current Opinion in Systems*  
 329 *Biology*, 8:81–89, April 2018. ISSN 2452-3100. doi: 10.1016/j.coisb.2017.12.005.
- 330 Lisa Feldman Barrett. Context reconsidered: Complex signal ensembles, relational meaning, and  
 331 population thinking in psychological science. *American Psychologist*, 77(8):894, 2022.
- 332 Andre M Bastos, Julien Vezoli, and Pascal Fries. Communication through coherence with inter-  
 333 areal delays. *Current Opinion in Neurobiology*, 31:173–180, April 2015. ISSN 09594388. doi:  
 334 10.1016/j.conb.2014.11.001.
- 335 André M. Bastos, W Martin Usrey, Rick A Adams, George R Mangun, Pascal Fries, and Karl J  
 336 Friston. Canonical microcircuits for predictive coding. *Neuron*, 76(4):695–711, 2012.
- 337 André M. Bastos, Mikael Lundqvist, Ayan S. Waite, Nancy Kopell, and Earl K. Miller. Layer and  
 338 rhythm specificity for predictive routing. *Proceedings of the National Academy of Sciences*, 117  
 339 (49):31459–31469, December 2020. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.2014868117.
- 340 Johanna Bergmann, Lucy S Petro, Clement Abbatecola, Min S Li, A Tyler Morgan, and Lars Muckli.  
 341 Cortical depth profiles in primary visual cortex for illusory and imaginary experiences. *Nature*  
 342 *Communications*, 15(1):1002, 2024.
- 343 Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis  
 344 Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal  
 345 probabilistic programming. *Journal of Machine Learning Research*, 20(28):1–6, 2019. URL  
 346 <http://jmlr.org/papers/v20/18-403.html>.
- 347 Cezar Borba, Matthew J Kourakis, Shea Schwennicke, Lorena Brasnic, and William C Smith. Fold  
 348 change detection in visual processing. *Frontiers in Neural Circuits*, 15:705161, 2021.
- 349 Marcus Brubaker, Mathieu Salzmann, and Raquel Urtasun. A family of mcmc methods on implicitly  
 350 defined manifolds. In *Artificial intelligence and statistics*, pages 161–172. PMLR, 2012.
- 351 György Buzsáki and Andreas Draguhn. Neuronal oscillations in cortical networks. *Science*, 304  
 352 (5679):1926–1929, June 2004. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.1099745.
- 353 Luke Campagnola, Stephanie C. Seeman, Thomas Chartrand, Lisa Kim, Alex Hoggarth, Clare  
 354 Gamlin, Shinya Ito, Jessica Trinh, Pasha Davoudian, Cristina Radaelli, Mean-Hwan Kim, Travis  
 355 Hage, Thomas Braun, Lauren Alfiler, Julia Andrade, Phillip Bohn, Rachel Dalley, Alex Henry,  
 356 Sara Kebede, Alice Mukora, David Sandman, Grace Williams, Rachael Larsen, Corinne Teeter,  
 357 Tanya L. Daigle, Kyla Berry, Nadia Dotson, Rachel Enstrom, Melissa Gorham, Madie Hupp,  
 358 Samuel Dingman Lee, Kiet Ngo, Philip R. Nicovich, Lydia Potekhina, Shea Ransford, Amanda  
 359 Gary, Jeff Goldy, Delissa McMillen, Trangthanh Pham, Michael Tieu, La’Akea Siverts, Miranda  
 360 Walker, Colin Farrell, Martin Schroedter, Cliff Slaughterbeck, Charles Cobb, Richard Ellenbogen,  
 361 Ryder P. Gwinn, C. Dirk Keene, Andrew L. Ko, Jeffrey G. Ojemann, Daniel L. Silbergeld, Daniel  
 362 Carey, Tamara Casper, Kirsten Crichton, Michael Clark, Nick Dee, Lauren Ellingwood, Jessica  
 363 Gloe, Matthew Kroll, Josef Sulc, Herman Tung, Katherine Wadhwani, Krissy Brouner, Tom Egdorf,  
 364 Michelle Maxwell, Medea McGraw, Christina Alice Pom, Augustin Ruiz, Jasmine Bomben, David  
 365 Feng, Nika Hejazinia, Shu Shi, Aaron Szafer, Wayne Wakeman, John Phillips, Amy Bernard, Luke  
 366 Esposito, Florence D. D’Orazi, Susan Sunkin, Kimberly Smith, Bosiljka Tasic, Anton Arkhipov,  
 367 Staci Sorensen, Ed Lein, Christof Koch, Gabe Murphy, Hongkui Zeng, and Tim Jarsky. Local  
 368 connectivity and synaptic dynamics in mouse and human neocortex. *Science*, 375(6585):eabj5861,  
 369 2022. doi: 10.1126/science.abj5861. URL [https://www.science.org/doi/abs/10.1126/](https://www.science.org/doi/abs/10.1126/science.abj5861)  
 370 [science.abj5861](https://www.science.org/doi/abs/10.1126/science.abj5861).
- 371 Matteo Carandini and David J Heeger. Normalization as a canonical neural computation. *Nature*  
 372 *reviews neuroscience*, 13(1):51–62, 2012.

373 Nick Chater, Joshua B Tenenbaum, and Alan Yuille. Probabilistic models of cognition: Conceptual  
374 foundations. *Trends in cognitive sciences*, 10(7):287–291, 2006.

375 Spyridon Chavlis and Panayiota Poirazi. Drawing inspiration from biological dendrites to empower  
376 artificial neural networks. *Current Opinion in Neurobiology*, 70:1–10, October 2021. ISSN  
377 0959-4388. doi: 10.1016/j.conb.2021.04.007.

378 Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned image compression  
379 with discretized gaussian mixture likelihoods and attention modules. In *Proceedings of the 2020*  
380 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

381 Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to  
382 handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pages  
383 2921–2926. IEEE, 2017.

384 Ishita Dasgupta, Eric Schulz, Joshua B Tenenbaum, and Samuel J Gershman. A theory of learning to  
385 infer. *Psychological review*, 127(3):412, 2020.

386 Xingsi Dong and Si Wu. Neural Sampling in Hierarchical Exponential-family Energy-based Models.  
387 In *Advances in Neural Information Processing Systems*, New Orleans, LA, 2023. Curran Associates  
388 Inc.

389 Rodney J Douglas and Kevan AC Martin. Neuronal circuits of the neocortex. *Annu. Rev. Neurosci.*,  
390 27(1):419–451, 2004.

391 Kenji Doya. *Bayesian brain: Probabilistic approaches to neural coding*. MIT press, 2007.

392 Michael Y-S Fang, Mayur Mudigonda, Ryan Zarcone, Amir Khosrowshahi, and Bruno A Olshausen.  
393 Learning and inference in sparse coding models with langevin dynamics. *Neural Computation*, 34  
394 (8):1676–1700, 2022.

395 Oren Forkosh. Memoryless optimality: Neurons do not need adaptation to optimally encode stimuli  
396 with arbitrarily complex statistics. *Neural Computation*, 34(12):2374–2387, November 2022.  
397 ISSN 0899-7667, 1530-888X. doi: 10.1162/neco\_a\_01543.

398 Pascal Fries. Rhythms for cognition: Communication through coherence. *Neuron*, 88(1):220–235,  
399 October 2015. ISSN 0896-6273. doi: 10.1016/j.neuron.2015.09.034.

400 Karl Friston. A theory of cortical responses. *Philosophical transactions of the Royal Society B:*  
401 *Biological sciences*, 360(1456):815–836, 2005.

402 Karl Friston, James Kilner, and Lee Harrison. A free energy principle for the brain. *Journal*  
403 *of Physiology-Paris*, 100(1):70–87, 2006. ISSN 0928-4257. doi: [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.jphysparis.2006.10.001)  
404 [jphysparis.2006.10.001](https://www.sciencedirect.com/science/article/pii/S092842570600060X). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S092842570600060X)  
405 [S092842570600060X](https://www.sciencedirect.com/science/article/pii/S092842570600060X). Theoretical and Computational Neuroscience: Understanding Brain Func-  
406 tions.

407 Colleen J. Gillon, Jason E. Pina, Jérôme A. Lecoq, Ruweida Ahmed, Yazan N. Billeh, Shiella  
408 Caldejon, Peter Groblewski, Timothy M. Henley, India Kato, Eric Lee, Jennifer Luviano, Kyla  
409 Mace, Chelsea Nayan, Thuyanh V. Nguyen, Kat North, Jed Perkins, Sam Seid, Matthew T. Valley,  
410 Ali Williford, Yoshua Bengio, Timothy P. Lillicrap, Blake A. Richards, and Joel Zylberberg.  
411 Learning from unexpected events in the neocortical microcircuit. *bioRxiv*, 2023. doi: 10.1101/  
412 2021.01.15.426915. URL [https://www.biorxiv.org/content/early/2023/04/06/2021.](https://www.biorxiv.org/content/early/2023/04/06/2021.01.15.426915)  
413 [01.15.426915](https://www.biorxiv.org/content/early/2023/04/06/2021.01.15.426915).

414 Siavash Golkar, Tiberiu Tesileanu, Yanis Bahroun, Anirvan Sengupta, and Dmitri Chklovskii.  
415 Constrained predictive coding as a biologically plausible model of the cortical hierarchy. In  
416 S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in*  
417 *Neural Information Processing Systems*, volume 35, pages 14155–14169. Curran Associates,  
418 Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/5b5de8526aac159e37ff9547713677ed-Paper-Conference.pdf)  
419 [5b5de8526aac159e37ff9547713677ed-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/5b5de8526aac159e37ff9547713677ed-Paper-Conference.pdf).

- Joseph Gonzalez, Yucheng Low, Arthur Gretton, and Carlos Guestrin. Parallel gibbs sampling: From colored fields to thin junction trees. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 324–332, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR. URL <https://proceedings.mlr.press/v15/gonzalez11a.html>.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy2fzU9gl>.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, page 6840–6851. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html>.
- Katie Hoemann, Maria Gendron, and Lisa Feldman Barrett. Mixed emotions in the predictive brain. *Current Opinion in Behavioral Sciences*, 15:51–57, 2017. ISSN 2352-1546. doi: <https://doi.org/10.1016/j.cobeha.2017.05.013>. URL <https://www.sciencedirect.com/science/article/pii/S2352154616302686>. Mixed emotions.
- J. Benjamin Hutchinson and Lisa Feldman Barrett. The Power of Predictions: An Emerging Paradigm for Psychological Research. *Current Directions in Psychological Science*, 2019. ISSN 14678721. doi: 10.1177/0963721419831992.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Juan Kuntz, Jen Ning Lim, and Adam M Johansen. Particle algorithms for maximum likelihood training of latent variable models. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, volume 206, Valencia, Spain, April 2023. Proceedings of Machine Learning Research.
- Juan Kuntz, Francesca R. Crucinio, and Adam M. Johansen. The divide-and-conquer sequential Monte Carlo algorithm: Theoretical properties and limit theorems. *The Annals of Applied Probability*, 34(1B):1469 – 1523, 2024. doi: 10.1214/23-AAP1996. URL <https://doi.org/10.1214/23-AAP1996>.
- Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40:e253, 2017. doi: 10.1017/S0140525X16001837.
- Yann Lecun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. ISSN 14764687. doi: 10.1038/nature14539. Citation Key: Lecun2015.
- F. Lindsten, A. M. Johansen, C. A. Naesseth, B. Kirkpatrick, T. B. Schön, J. A.D. Aston, and A. Bouchard-Côté. Divide-and-conquer with sequential monte carlo. *Journal of Computational and Graphical Statistics*, 26(2):445–458, 2017. ISSN 15372715. doi: 10.1080/10618600.2016.1237363. arXiv: 1406.4993 Citation Key: Lindsten2017.
- Shih-Chii Liu. A winner-take-all circuit with controllable soft max property. In S. Solla, T. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL [https://proceedings.neurips.cc/paper\\_files/paper/1999/file/3e7e0224018ab3cf51abb96464d518cd-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1999/file/3e7e0224018ab3cf51abb96464d518cd-Paper.pdf).
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

- 468 Gabriel Loaiza-Ganem and John P Cunningham. The continuous bernoulli: fixing a pervasive error  
469 in variational autoencoders. In *Advances in Neural Information Processing Systems*, volume 32.  
470 Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/  
471 f82798ec8909d23e55679ee26bb26437-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/f82798ec8909d23e55679ee26bb26437-Abstract.html).
- 472 Brian N. Lundstrom, Matthew H. Higgs, William J. Spain, and Adrienne L. Fairhall. Fractional differ-  
473 entiation by neocortical pyramidal neurons. *Nature Neuroscience*, 11(11):1335–1342, November  
474 2008. ISSN 1546-1726. doi: 10.1038/nn.2212.
- 475 Zachary F Mainen and Terrence J Sejnowski. Reliability of spike timing in neocortical neurons.  
476 *Science*, 268(5216):1503–1506, 1995.
- 477 Joseph Marino, Yisong Yue, and Stephan Mandt. Iterative amortized inference. In *35th International  
478 Conference on Machine Learning, ICML 2018*, volume 8, page 5444–5462, 2018. ISBN 978-1-  
479 5108-6796-3. arXiv: 1807.09356 Citation Key: Marino2018a.
- 480 Beren Millidge, Anil Seth, and Christopher L Buckley. Predictive coding: a theoretical and experi-  
481 mental review. *arXiv preprint arXiv:2107.12979*, 2021.
- 482 Beren Millidge, Yuhang Song, Tommaso Salvatori, Thomas Lukasiewicz, and Rafal Bogacz. A  
483 theoretical framework for inference and learning in predictive coding networks. In *The Eleventh  
484 International Conference on Learning Representations*, 2023. URL [https://openreview.net/  
485 forum?id=ZCTvSF\\_uVM4](https://openreview.net/forum?id=ZCTvSF_uVM4).
- 486 Toviah Moldwin, Menachem Kalmenson, and Idan Segev. The gradient clusteron: A model neuron  
487 that learns to solve classification tasks via dendritic nonlinearities, structural plasticity, and gradient  
488 descent. *PLOS Computational Biology*, 17(5):e1009015, May 2021. ISSN 1553-7358. doi:  
489 10.1371/journal.pcbi.1009015.
- 490 Christian Naesseth, Fredrik Lindsten, and Thomas Schon. Nested sequential monte carlo methods.  
491 In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference  
492 on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1292–  
493 1301, Lille, France, 07–09 Jul 2015. PMLR. URL [https://proceedings.mlr.press/v37/  
494 naesseth15.html](https://proceedings.mlr.press/v37/naesseth15.html).
- 495 Christian Naesseth, Scott Linderman, Rajesh Ranganath, and David Blei. Variational sequen-  
496 tial monte carlo. In Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the  
497 Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of  
498 *Proceedings of Machine Learning Research*, pages 968–977. PMLR, 09–11 Apr 2018. URL  
499 <https://proceedings.mlr.press/v84/naesseth18a.html>.
- 500 Radford M. Neal and Geoffrey E. Hinton. *A View of the EM Algorithm that Justifies Incremental,  
501 Sparse, and other Variants*, page 355–368. NATO ASI Series. Springer Netherlands, Dordrecht,  
502 1998. ISBN 978-94-011-5014-9. doi: 10.1007/978-94-011-5014-9\_12. URL [https://doi.org/  
503 10.1007/978-94-011-5014-9\\_12](https://doi.org/10.1007/978-94-011-5014-9_12).
- 504 Gaspard Oliviers, Rafal Bogacz, and Alexander Meulemans. Learning probability distributions of  
505 sensory inputs with monte carlo predictive coding. *bioRxiv*, 2024.
- 506 Benjamin Peters, James J DiCarlo, Todd Gureckis, Ralf Haefner, Leyla Isik, Joshua Tenenbaum, Talia  
507 Konkle, Thomas Naselaris, Kimberly Stachenfeld, Zenna Tavares, et al. How does the primate  
508 brain combine generative and discriminative computations in vision? *ArXiv*, 2024.
- 509 Luca Pinchetti, Tommaso Salvatori, Yordan Yordanov, Beren Millidge, Yuhang Song, and  
510 Thomas Lukasiewicz. Predictive coding beyond gaussian distributions. In S. Koyejo,  
511 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neu-  
512 ral Information Processing Systems*, volume 35, pages 1280–1293. Curran Associates,  
513 Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/  
514 08f9de0232c0b485110237f6e6cf88f1-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/08f9de0232c0b485110237f6e6cf88f1-Paper-Conference.pdf).
- 515 Alexandre Pouget, Jeffrey M Beck, Wei Ji Ma, and Peter E Latham. Probabilistic brains: knowns and  
516 unknowns. *Nature neuroscience*, 16(9):1170–1178, 2013.

517 Rajesh P N Rao and Dana H Ballard. Predictive coding in the visual cortex: a functional interpretation  
518 of some extra-classical receptive-field effects. *Nature neuroscience*, 2(1):79–87, 1999. ISSN  
519 1097-6256. doi: 10.1038/4580. URL 10.1038/4580%5Cnhttp://www.nature.com/neuro/  
520 journal/v2/n1/abs/nn0199\_79.html.

521 Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and  
522 approximate inference in deep generative models. In *Proceedings of the 31st International  
523 Conference on Machine Learning*, volume 4, page 3057–3070, Beijing, China, 2014. arXiv:  
524 1401.4082 Citation Key: Rezende2014 ISBN: 9781634393973.

525 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by  
526 back-propagating errors. *Nature*, 323(6088):533–536, 1986.

527 Tommaso Salvatori, Luca Pinchetti, Beren Millidge, Yuhang Song, Tianyi Bao, Rafal Bogacz, and  
528 Thomas Lukasiewicz. Learning on arbitrary graph topologies via predictive coding. *Advances in  
529 neural information processing systems*, 35:38232–38244, 2022.

530 Tommaso Salvatori, Ankur Mali, Christopher L Buckley, Thomas Lukasiewicz, Rajesh PN Rao, Karl  
531 Friston, and Alexander Ororbia. Brain-inspired computational intelligence via predictive coding.  
532 *arXiv preprint arXiv:2308.07870*, 2023.

533 Tommaso Salvatori, Yuhang Song, Yordan Yordanov, Beren Millidge, Cornelius Emde, Zhenghua  
534 Xu, Lei Sha, Rafal Bogacz, and Thomas Lukasiewicz. A stable, fast, and fully automatic learning  
535 algorithm for predictive coding networks. In *International Conference on Learning Representations*,  
536 2024.

537 Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117,  
538 2015. ISSN 18792782. doi: 10.1016/j.neunet.2014.09.003. Citation Key: Schmidhuber2015.

539 Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. [https://github.com/mseitzer/  
540 pytorch-fid](https://github.com/mseitzer/pytorch-fid), August 2020. Version 0.3.0.

541 Lei Shi and Thomas L. Griffiths. Neural implementation of hierarchical bayesian inference by  
542 importance sampling. In *Advances in Neural Information Processing Systems*, page 1669–1677,  
543 2009. ISBN 978-1-61567-911-9. Citation Key: Shi2009.

544 Yuhang Song, Beren Millidge, Tommaso Salvatori, Thomas Lukasiewicz, Zhenghua Xu, and  
545 Rafal Bogacz. Inferring neural activity before plasticity as a foundation for learning beyond  
546 backpropagation. *Nature Neuroscience*, page 1–11, January 2024. ISSN 1546-1726. doi:  
547 10.1038/s41593-023-01514-1.

548 M. W. Spratling. A review of predictive coding algorithms. *Brain and Cognition*, 112:92–97, 2017.  
549 ISSN 10902147. doi: 10.1016/j.bandc.2015.11.003.

550 M. V. Srinivasan, S. B. Laughlin, and A. Dubs. Predictive coding: A fresh view of inhibition in the  
551 retina. *Proceedings of the Royal Society of London - Biological Sciences*, 216(1205):427–459,  
552 1982. ISSN 09628452. doi: 10.1098/rspb.1982.0085.

553 Sam Stites, Heiko Zimmermann, Hao Wu, Eli Sennesh, and Jan-Willem Van de Meent. Learning  
554 proposals for probabilistic programs with inference combinators. *37th Conference on Uncertainty  
555 in Artificial Intelligence (UAI 2021)*, 2021.

556 Shohei Taniguchi, Yusuke Iwasawa, Wataru Kumagai, and Yutaka Matsuo. Langevin autoencoders  
557 for learning deep latent variable models. *Advances in Neural Information Processing Systems*, 35:  
558 13277–13289, 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/  
559 hash/565f995643da6329cec701f26f8579f5-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/565f995643da6329cec701f26f8579f5-Abstract-Conference.html).

560 Stefan Webb, Adam Golinski, Rob Zinkov, Siddharth N, Tom Rainforth, Yee Whye Teh, and  
561 Frank Wood. Faithful inversion of generative models for effective amortized inference. In  
562 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, ed-  
563 itors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates,  
564 Inc., 2018. URL [https://proceedings.neurips.cc/paper\\_files/paper/2018/file/  
565 894b77f805bd94d292574c38c5d628d5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/894b77f805bd94d292574c38c5d628d5-Paper.pdf).

- 566 Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient langevin dynamics. In  
567 *Proceedings of the 28th International Conference on Machine Learning, ICML 2011*, Bellevue,  
568 WA, USA, 2011. Proceedings of Machine Learning Research.
- 569 James CR Whittington and Rafal Bogacz. An approximation of the error backpropagation algorithm  
570 in a predictive coding network with local hebbian synaptic plasticity. *Neural computation*, 29(5):  
571 1229–1262, 2017.
- 572 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking  
573 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 574 L. Yin and P. Ao. Existence and construction of dynamical potential in nonequilibrium processes  
575 without detailed balance. *Journal of Physics A: Mathematical and General*, 39(27):8593, June  
576 2006. ISSN 0305-4470. doi: 10.1088/0305-4470/39/27/003. URL [https://dx.doi.org/10.](https://dx.doi.org/10.1088/0305-4470/39/27/003)  
577 1088/0305-4470/39/27/003.
- 578 Umair Zahid, Qinghai Guo, and Zafeirios Fountas. Sample as you infer: Predictive coding with  
579 langevin dynamics. In *Proceedings of the 41st International Conference on Machine Learning*,  
580 volume 235, Vienna, Austria, 2024. Proceedings of Machine Learning Research.

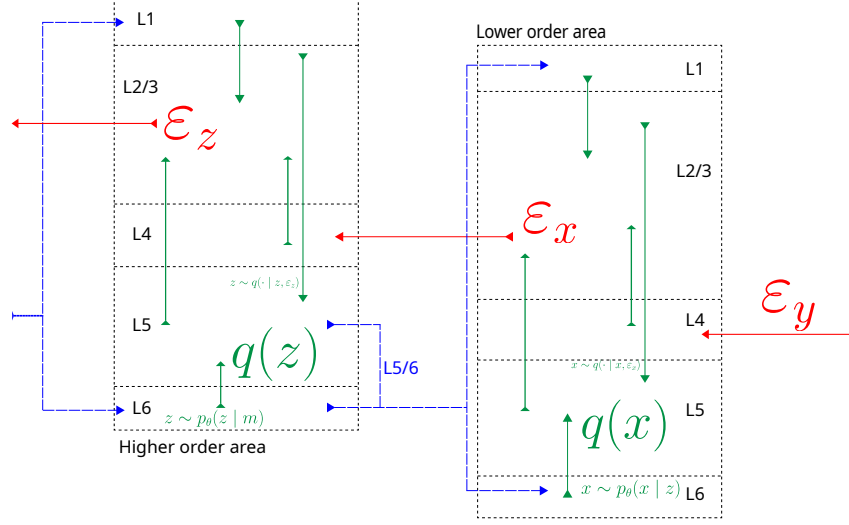


Figure 5: Divide-and-conquer predictive coding provides an algorithmic interpretation for some of the connections mapped in the canonical neocortical microcircuit [Bastos et al., 2012, 2020, Campagnola et al., 2022]: prediction errors (red) arrive through ascending pathways into the central laminar layer 4, which transmits them up to layers 2/3 (green). These layers combine the incoming errors with a present posterior estimate (green L5→L2/3 connection) to generate prediction errors for the next cortical area. Eventually the updated predictions flow back down the cortical hierarchy (blue).

## 581 A Further experiments and results

582 **Alternate image generation/ representation learning** As indicated in Section 2, this paper builds  
 583 upon the particle gradient descent (PGD) algorithm; Kuntz et al. [2023] demonstrated the algorithm’s  
 584 performance by training a generator network on CelebA. Their network employed a Gaussian  
 585 likelihood with a fixed standard deviation of 0.01, and evaluated a log-joint objective over 100 epochs  
 586 on exactly 10,000 subsampled data points. The paper then evaluated mean squared error on an  
 587 inpainting task and the Frechet Inception Distance over data images.

588 When applied to the resulting target density, DCPC amounts to PGD with a resampling step. Table 4  
 589 shows the results of training and evaluating the same model described above with DCPC. Since PGD  
 590 trained for 100 epochs with a batch size of 128, albeit on a 10,000-image subsample of CelebA, we  
 591 trained with the entire dataset for 100 epochs with batch-size 128.

Inference type	Log-joint	FID ↓
PGD ( $K = 10$ )	$-3.8 \times 10^5$	$100 \pm 2.7$
DCPC (ours, $K = 10$ )	$-3.0 \times 10^6$	$89.6 \pm 0.6$

Table 4: Log-joint probabilities and FID metrics show how DCPC performs against the original PGD.

592 We suspect that the supplied code for log-joint calculation averages over either particles or batch items  
 593 differently from how we have evaluated DCPC (e.g. we call `mean()` without dividing by any further  
 594 shape dimensions), accounting for the apparent order-of-magnitude difference between log-joints.

595 At the request of reviewers, we have substituted a simplified Figure 1 in the main text for Figure 5  
 596 showing how to map DCPC onto laminar microcircuit structure.

## 597 B Importance sampling and gradient estimation proofs

598 **Definition 2** (Strict proper weighting for a density). *Given an unnormalized density  $\gamma_\theta(\mathbf{z})$  with*  
 599 *corresponding normalizing constant  $Z(\theta)$  and normalized density  $\pi_\theta(\mathbf{z})$*

$$Z(\theta) := \int_{\mathbf{z} \in \mathcal{Z}} \gamma_\theta(\mathbf{z}) d\mathbf{z} \quad \pi_\theta(\mathbf{z}) := \frac{\gamma_\theta(\mathbf{z})}{Z(\theta)},$$



the random variables  $w, \mathbf{z} \sim q(w, \mathbf{z})$  are strictly properly weighted [Naesseth et al., 2015] with respect to  $\gamma_\theta(\mathbf{z})$  if and only if for any measurable test function  $h : \mathcal{Z} \rightarrow \mathbb{R}$ , the weighted expectation over the proposal  $q(w, \mathbf{z})$  equals the expectation under the target  $\gamma_\theta(\mathbf{z})$

$$\mathbb{E}_{w, \mathbf{z} \sim q(w, \mathbf{z})} [wh(\mathbf{z})] = \int_{\mathbf{z} \in \mathcal{Z}} h(\mathbf{z}) \gamma_\theta(\mathbf{z}) d\mathbf{z}. \quad (11)$$

The first two propositions come from the previous work by Wu et al. [2020], Stites et al. [2021] and Zimmermann et al. [2021]. The reader looking for foundations can see Naesseth et al. [2015] or Chopin and Papaspiliopoulos [2020].

**Proposition 1** (The free energy upper-bounds the surprisal). *Given a proposal  $q_\phi(w, \mathbf{z})$  strictly properly weighted (Definition 2) for the target  $\gamma_\theta(\mathbf{z})$ , the variational free energy provides an upper bound to the target’s surprisal*

$$\mathcal{F}(\theta, q) \geq -\log Z(\theta). \quad (12)$$

*Proof.* I begin by writing out the free energy (Equation 2) as an expectation of a negative logarithm

$$\mathcal{F}(\theta, q) = \mathbb{E}_{z, w \sim q(z, w)} [-\log w].$$

Jensen’s Inequality allows moving the expectation into the negative logarithm by relaxing the definition of the variational free energy from an equality to an upper bound

$$\mathcal{F}(\theta, q) \geq -\log \mathbb{E}_{z, w \sim q(z, w)} [w].$$

Setting  $h(z) = 1$ , strict proper weighting for an unnormalized density (Definition 2) says the expected weight will be the normalizing constant

$$\mathbb{E}_{z, w \sim q(z, w)} [w] = Z(\theta)$$

which I substitute back in to obtain the desired inequality  $\mathcal{F}(\theta, q) \geq -\log Z(\theta)$ .  $\square$

**Proposition 2** (Weighted expectations approximate the normalized target up to a constant). *Given a proposal  $q_\phi(w, \mathbf{z})$  strictly properly weighted (Definition 2) for the target  $\gamma_\theta(\mathbf{z})$  and a measurable test function  $h : \mathcal{Z} \rightarrow \mathbb{R}$ , weighted expectations under the proposal equal the target’s normalizing constant times the test function’s expectation under the normalized target*

$$\mathbb{E}_{(w, \mathbf{z}) \sim q_\phi(w, \mathbf{z})} [wh(\mathbf{z})] = Z(\theta) \mathbb{E}_{\mathbf{z} \sim \pi_\theta(\cdot)} [h(\mathbf{z})].$$

*Proof.* Strict proper weighting (Equation 11) states that weighted expectations under the proposal equal integrals over the unnormalized target, and by definition the normalized target equals the unnormalized density over its normalizing constant

$$\mathbb{E}_{w, \mathbf{z} \sim q(w, \mathbf{z})} [wh(\mathbf{z})] = \int_{\mathbf{z} \in \mathcal{Z}} h(\mathbf{z}) \gamma_\theta(\mathbf{z}) d\mathbf{z}, \quad \pi_\theta(\mathbf{z}) := \frac{\gamma_\theta(\mathbf{z})}{Z(\theta)}.$$

The second equation expresses the unnormalized target in terms of the normalized one

$$Z(\theta) \pi_\theta(\mathbf{z}) = \gamma_\theta(\mathbf{z}),$$

and substituting this expression into the definition of strict proper weighting leads to the desired result

$$\begin{aligned} \int_{\mathbf{z} \in \mathcal{Z}} h(\mathbf{z}) \gamma_\theta(\mathbf{z}) d\mathbf{z} &= \int_{\mathbf{z} \in \mathcal{Z}} h(\mathbf{z}) Z(\theta) \pi_\theta(\mathbf{z}) d\mathbf{z}, \\ &= Z(\theta) \int_{\mathbf{z} \in \mathcal{Z}} h(\mathbf{z}) \pi_\theta(\mathbf{z}) d\mathbf{z} \\ \mathbb{E}_{w, \mathbf{z} \sim q(w, \mathbf{z})} [wh(\mathbf{z})] &= Z(\theta) \mathbb{E}_{\pi_\theta(\mathbf{z})} [h(\mathbf{z})]. \end{aligned} \quad \square$$

**Proposition 3** (DCPC’s free energy has a pathwise derivative). *The free energy  $\mathcal{F}^{t+1} = \mathbb{E}_q [-\log w_{\theta^t}^{t+1}]$  constructed by the population predictive coding algorithm (Algorithm 1) has a pathwise derivative as the expectation of the negative gradient of the log-joint density*

$$\nabla_{\theta^t} \mathcal{F}^{t+1} = \mathbb{E}_q [-\nabla_{\theta^t} \log p_{\theta^t}(\mathbf{x}, \mathbf{z}^{t+1})].$$

627 *Proof.* The free energy has an expression in terms of Equation 8

$$\begin{aligned}\mathcal{F}^{t+1} &= \mathbb{E}_q [-\log w_{\theta^t}^{t+1}] & w_{\theta^t}^{t+1} &= \frac{p_{\theta^t}(\mathbf{x}, \mathbf{z})}{\prod_{z \in \mathbf{z}} \gamma_{\theta}(z_b^{t+1}; \mathbf{z}_{\setminus z})} \prod_{z \in \mathbf{z}} \hat{Z}_{\theta^t}(\mathbf{z}_{\setminus z})^{t+1}, \\ \hat{Z}_{\theta^t}(\mathbf{z}_{\setminus z})^{t+1} &= \frac{1}{K} \sum_{k=1}^K u_b^{t+1,k} & u_z^{t+1} &= \frac{\gamma_{\theta}(z^{t+1}; \mathbf{z}_{\setminus z})}{q(z^{t+1} \mid \varepsilon_z(z^t))},\end{aligned}$$

628 and writing out the free energy itself in full shows that many terms cancel

$$\begin{aligned}q(\mathbf{z}^{t+1} \mid \mathbf{z}^t) &= \prod_{z_b^{t+1} \in \mathbf{z}^{t+1}} q(z^{t+1} \mid z^t), \\ \mathcal{F}^{t+1} &= \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ -\log \frac{p_{\theta^t}(\mathbf{x}, \mathbf{z})}{\prod_{z \in \mathbf{z}} \gamma_{\theta}(z_b^{t+1}; \mathbf{z}_{\setminus z})} \prod_{z \in \mathbf{z}} \frac{1}{K} \sum_{k=1}^K \frac{\gamma_{\theta}(z_b^{t+1}; \mathbf{z}_{\setminus z})}{q(z^{t+1} \mid \varepsilon_z(z^t))} \right] \\ &= \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ -\log \frac{p_{\theta^t}(\mathbf{x}, \mathbf{z})}{\prod_{z \in \mathbf{z}} \gamma_{\theta}(z_b^{t+1}; \mathbf{z}_{\setminus z})} \frac{\prod_{z \in \mathbf{z}} \gamma_{\theta}(z_b^{t+1}; \mathbf{z}_{\setminus z})}{\prod_{z \in \mathbf{z}} q(z^{t+1} \mid \varepsilon_z(z^t))} \right] \\ &= \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ -\log \frac{p_{\theta^t}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \right].\end{aligned}$$

629 The proposal distribution  $q$  is a function of the random variable values themselves through the  
630 prediction errors, not of the parameters  $\theta$ . The above expression therefore admits a pathwise  
631 derivative [Schulman et al., 2015], moving the gradient operator into the expectation

$$\begin{aligned}\nabla_{\theta^t} \mathcal{F}^{t+1} &= \nabla_{\theta^t} \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ -\log \frac{p_{\theta^t}(\mathbf{x}, \mathbf{z}^{t+1})}{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \right] \\ &= \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ \nabla_{\theta^t} -\log \frac{p_{\theta^t}(\mathbf{x}, \mathbf{z}^{t+1})}{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \right] \\ &= \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ \nabla_{\theta^t} - [\log p_{\theta^t}(\mathbf{x}, \mathbf{z}^{t+1}) - \log q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)] \right] \\ &= \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ -[\nabla_{\theta^t} \log p_{\theta^t}(\mathbf{x}, \mathbf{z}^{t+1}) - \nabla_{\theta^t} \log q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)] \right] \\ \nabla_{\theta^t} \mathcal{F}^{t+1} &= \mathbb{E}_{q(\mathbf{z}^{t+1} \mid \mathbf{z}^t)} \left[ -\nabla_{\theta^t} \log p_{\theta^t}(\mathbf{x}, \mathbf{z}^{t+1}) \right].\end{aligned}\quad \square$$

632 **Proposition 4** (DCPC coordinate updates are strictly properly weighted for the complete conditionals).  
633 Each DCPC coordinate update (Equation 7) for a latent variable  $z \in \mathbf{z}$  is strictly properly weighted  
634 (Definition 2) for  $z$ 's unnormalized complete conditional. For every measurable  $h : \mathcal{Z} \rightarrow \mathbb{R}$

$$\mathbb{E}_{z \sim q_{\eta}(z^t \mid z^{t-1}, \varepsilon_z^t)} \left[ \mathbb{E}_{u \sim \delta(u), z' \sim \text{RESAMPLE}(z, u_z)} [h(z)] \right] = \int_{z \in \mathcal{Z}} h(z) \gamma_{\theta}(z; \mathbf{z}_{\setminus z}) dz. \quad (13)$$

635 *Proof.* Expanding the outer expectation into an integral and replacing the Dirac delta with the  
636 expression for the local weights transforms Equation 13 into

$$\begin{aligned}\int_{z \in \mathcal{Z}} \frac{\gamma_{\theta}(z; \mathbf{z}_{\setminus z})}{q_{\eta}(z \mid z^{t-1}, \varepsilon_z^t)} \mathbb{E}_{z' \sim \text{RESAMPLE}(z, u_z)} [h(z')] q_{\eta}(z \mid z^{t-1}, \varepsilon_z^t) dz &= \\ &= \int_{z \in \mathcal{Z}} h(z) \gamma_{\theta}(z; \mathbf{z}_{\setminus z}) dz;\end{aligned}$$

637 importance resampling also preserves strict proper weighting (see Naesseth et al. [2015], Stites et al.  
638 [2021] and Chopin and Papaspiliopoulos [2020] for proofs), and so this yields

$$\begin{aligned}\int_{z \in \mathcal{Z}} \mathbb{E}_{z' \sim \text{RESAMPLE}(z, u_z)} [h(z')] \gamma_{\theta}(z; \mathbf{z}_{\setminus z}) dz &= \int_{z \in \mathcal{Z}} h(z) \gamma_{\theta}(z; \mathbf{z}_{\setminus z}) dz \\ \int_{z' \in \mathcal{Z}} h(z') \gamma_{\theta}(z'; \mathbf{z}_{\setminus z}) dz' &= \int_{z \in \mathcal{Z}} h(z) \gamma_{\theta}(z; \mathbf{z}_{\setminus z}) dz.\end{aligned}$$

639

□

640 **Corollary 4.1** (DCPC coordinate updates sample from the true complete conditionals). *Each DCPC*  
641 *coordinate update (Equation 7) for a latent  $z \in \mathbf{z}$  samples from  $z$ 's complete conditional (the*  
642 *normalization of Equation 5). Formally, for every measurable  $h : \mathcal{Z} \rightarrow \mathbb{R}$ , resampled expectations*  
643 *with respect to the DCPC coordinate update equal those with respect to the complete conditional*

$$\mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u), z' \sim \text{RESAMPLE}(z, u_z)} [h(z')]] = \int_{z \in \mathcal{Z}} h(z) \pi_\theta(z | \mathbf{z}_{\setminus z}) dz.$$

644 *Proof.* Proposition 4 in Appendix B provides a lemma

$$\mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u), \hat{z} \sim \text{RESAMPLE}(z, u_z)} [h(z')]] = \int_{z \in \mathcal{Z}} h(z) \gamma_\theta(z; \mathbf{z}_{\setminus z}) dz,$$

645 which we can apply by observing that resampling sums over self-normalized weights

$$\begin{aligned} \mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u), z' \sim \text{RESAMPLE}(z, u_z)} [h(z)]] &= \\ \mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u)} [\mathbb{E}_{z' \sim \frac{u \delta_z(\cdot)}{\sum u'} [h(z')]]]] &, \end{aligned}$$

646 which is just a weighted sum that by Definition 2 is itself properly weighted

$$\begin{aligned} \mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u)} [\mathbb{E}_{z' \sim \frac{u \delta_z(\cdot)}{\sum u'} [h(z')]]]] &= \mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u)} [\frac{u}{\sum u} h(z)]] \\ &= \mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u)} [\frac{1}{\sum u} \int_{z \in \mathcal{Z}} h(z) \gamma_\theta(z; \mathbf{x}, \mathbf{z}_{\setminus z}) dz]] \\ &= \mathbb{E}_{z \sim q_\eta(z|z^{t-1}, \varepsilon_z^t)} [\mathbb{E}_{u \sim \delta(u)} [\frac{1}{\cancel{\hat{Z}_\theta(\mathbf{x}, \mathbf{z}_{\setminus z})} \cancel{Z_\theta(\mathbf{x}, \mathbf{z}_{\setminus z})}} \int_{z \in \mathcal{Z}} h(z) \pi_\theta(z | \mathbf{x}, \mathbf{z}_{\setminus z}) dz]] \\ &= \int_{z \in \mathcal{Z}} h(z) \pi_\theta(z | \mathbf{x}, \mathbf{z}_{\setminus z}) dz. \quad \square \end{aligned}$$

648 **Proposition 5** (DCPC parameter learning requires only local gradients in a factorized generative  
649 model). *Consider a graphical model factorized according to Equation 1, with the additional assump-*  
650 *tion that the model parameters  $\theta \in \Theta = \prod_{x \in \mathbf{x}} \Theta_x \times \prod_{z \in \mathbf{z}} \Theta_z$  share that factorization. Then the*  
651 *gradient  $\nabla_\theta \mathcal{F}(\theta, q)$  of DCPC's free energy similarly factorizes into a sum of local particle averages*

$$\begin{aligned} \nabla_\theta \mathcal{F} &= \mathbb{E}_q [-\nabla_\theta \log p_\theta(\mathbf{x}, \mathbf{z})] \\ &= \sum_{v \in (\mathbf{x}, \mathbf{z})} \mathbb{E}_{q(v, \text{Pa}(v) | \varepsilon_v, \varepsilon_{\text{Pa}(v)})} [-\nabla_{\theta_v} \log p_{\theta_v}(v | \text{Pa}(v))] \\ &= - \sum_{v \in (\mathbf{x}, \mathbf{z})} \frac{1}{K} \sum_{k=1}^K \nabla_{\theta_v} \log p_{\theta_v}(v^k | \text{Pa}(v)^k). \end{aligned}$$

652 *Proof.* Proposition 3 provides the lemma that  $\nabla_\theta \mathcal{F} = \mathbb{E}_q [-\nabla_\theta \log p_\theta(\mathbf{x}, \mathbf{z})]$ , and applying the  
653 factorization of the generative model demonstrates that

$$\nabla_\theta \mathcal{F} = \mathbb{E}_q \left[ -\nabla_\theta \sum_{v \in (\mathbf{x}, \mathbf{z})} \log p_\theta(v | \text{Pa}(v)) \right].$$

654 Since the proposal  $q$  does not depend on any  $\theta$  and consists of a particle cloud, we can rewrite it as a  
655 mixture over the particles (after sampling is performed)

$$\nabla_\theta \mathcal{F} \approx \frac{1}{K} \sum_{k=1}^K -\nabla_\theta \sum_{v \in (\mathbf{x}, \mathbf{z})} \log p_\theta(v^k | \text{Pa}(v)^k),$$

656 and then finally apply the assumption of this theorem that  $\theta \in \Theta = \prod_{x \in \mathbf{x}} \Theta_x \times \prod_{z \in \mathbf{z}} \Theta_z$ , moving  
657 the gradient operation into the sum over individual random variables

$$\approx \frac{1}{K} \sum_{k=1}^K \sum_{v \in (\mathbf{x}, \mathbf{z})} -\nabla_{\theta_v} \log p_{\theta_v}(v^k | \text{Pa}(v)^k). \quad \square$$

## C Extension to discrete sample spaces

Contemporaneously to the work of Kuntz et al. [2023] on particle gradient descent, Sun et al. [2023] derived a novel Wasserstein gradient flow and corresponding descent algorithm for discrete distributions. In their setting, each Wasserstein gradient step constructs a  $D$ -dimensional, finitely supported distribution over the  $C$ -Hamming ball of the starting sample, such that the distribution has  $DC$  possible states in total. Let  $z^{t+h} \in N_C(z^t)$  denote the resulting discrete random variable in the  $C$ -neighborhood around  $z^t$  with respect to the Hamming distance. The update rule relies on simulating the gradient flow for time  $h$ , sampling from a Markov jump process at time  $t + h$

$$z^{t+h} \sim \prod_{d \in [1 \dots D]} q(z_d^{t+h} | z_d^t).$$

A rate matrix  $Q_d(z^t)$  defined by the entire discrete variable  $z^t$  parameterizes the proposal distribution

$$q_h(z_d^{t+h} | z^t) = \exp(Q_d(z^t)h). \quad (14)$$

the rate matrix will have nondiagonal entries at indices  $i \neq j \in [1 \dots C]$  in the neighborhood  $N_C(z^t)$ ,

$$Q_d(z^t)_{i,j} = w_{i,j} g \left( \frac{\pi_\theta(z_{\setminus d}^t, z'_{d,j})}{\pi_\theta(z_{\setminus d}^t, z'_{d,i})} \right).$$

The above equation requires that  $\forall i, j \in [1 \dots C], w_{i,j} = w_{j,i} \in \mathbb{R}$  and  $g(a) = ag(\frac{1}{a})$ . The ratio of normalized target densities  $\pi$  will equal the ratio of unnormalized densities  $\gamma$

$$\frac{\pi_\theta(z_{\setminus d}^t, z'_{d,j})}{\pi_\theta(z_{\setminus d}^t, z'_{d,i})} = \frac{\gamma_\theta(z'_{d,j}; z_{\setminus d}^t) Z_{z_{\setminus d}^t}(z_{\setminus d}^t, \theta)}{Z_{z_{\setminus d}^t}(z_{\setminus d}^t) \gamma_\theta(z'_{d,i}; z_{\setminus d}^t)}$$

$$g \left( \frac{\pi_\theta(z_{\setminus d}^t, z'_{d,j})}{\pi_\theta(z_{\setminus d}^t, z'_{d,i})} \right) = g \left( \frac{\gamma_\theta(z'_{d,j}; z_{\setminus d}^t)}{\gamma_\theta(z'_{d,i}; z_{\setminus d}^t)} \right).$$

Based on the experimental recommendations of Sun et al. [2023], let  $w_{i,j} = w_{j,i} = 1$  and  $g(a) = \sqrt{a}$ . The rate matrix then simplifies to nondiagonal and diagonal terms

$$Q_d(z^t)_{i,j} = \sqrt{\frac{\gamma_\theta(z'_{d,j}; z_{\setminus d}^t)}{\gamma_\theta(z'_{d,i}; z_{\setminus d}^t)}}, \quad Q_d(z^t)_{i,i} = - \sum_{j \neq i} Q_d(z^t)_{i,j}. \quad (15)$$

Equations 14 and 15 give a distribution descending the Wasserstein gradient of the free energy with respect to a particle cloud in a discrete sample space. Applying Equation 15 to  $\gamma_\theta(z; \mathbf{z}_{\setminus z}^t)$  yields a factorization in log space

$$Q(z^t)_{i,j} = \sqrt{\frac{\gamma_\theta(z^t + i; \mathbf{z}_{\setminus z}^t)}{\gamma_\theta(z^t + j; \mathbf{z}_{\setminus z}^t)}} \log Q(z^t)_{i,j} = \frac{1}{2} \left( \log \gamma_\theta(z^t + i; \mathbf{z}_{\setminus z}^t) - \log \gamma_\theta(z^t + j; \mathbf{z}_{\setminus z}^t) \right).$$

This difference can be written as a difference of differences

$$\log \gamma_\theta(z^t + i; \mathbf{z}_{\setminus z}^t) - \log \gamma_\theta(z^t + j; \mathbf{z}_{\setminus z}^t) =$$

$$\left( \log \gamma_\theta(z^t + i; \mathbf{z}_{\setminus z}^t) - \log \gamma_\theta(z^t; \mathbf{z}_{\setminus z}^t) \right) - \left( \log \gamma_\theta(z^t + j; \mathbf{z}_{\setminus z}^t) - \log \gamma_\theta(z^t; \mathbf{z}_{\setminus z}^t) \right). \quad (16)$$

Recent work on efficient sampling for discrete distributions has focused on approximating density ratios, such as the one in Equation 15, with series expansions parameterized by error vectors. When the underlying discrete densities consist of exponentiating a differentiable energy function, as in Grathwohl et al. [2021], these error vectors have taken the form of gradients and the finite-series expansions have been Taylor series. When they do not, Xiang et al. [2023] showed how they take the form of finite differences and Newton's series

$$\log \gamma(z') - \log \gamma(z) \approx \Delta_1 (\log \gamma(z))^\top \cdot (z' - z). \quad (17)$$

682 Discrete DCPC would therefore use finite differences as discrete prediction errors, breaking each  
683 discrete  $z \in \mathbf{z}$  into dimensions and incrementing each dimension separately to construct a vector

$$\Delta_1 f(z) := (f(z_1 + 1, z_{2:D}), \dots, f(z_{1:i}, z_i + 1, z_{i+1:D}), \dots, f(z_{1:D-1}, z_D + 1)) \ominus f(z), \quad (18)$$

684 where  $\ominus$  subtracts the scalar  $f(z)$  from the vector elements and  $f : \mathbb{Z}^D \rightarrow \mathbb{R}$  is the target function.  
685 This would lead to defining the discrete prediction error as the finite difference

$$\varepsilon_z := \Delta_1 \log \gamma_\theta(z^t; \mathbf{z}_{\setminus z}^t). \quad (19)$$

686 Applying Equation 17 to the two terms of Equation 16, we obtain the approximations

$$\begin{aligned} \log \gamma_\theta(z^t + i; \mathbf{z}_{\setminus z}^t) - \log \gamma_\theta(z^t; \mathbf{z}_{\setminus z}^t) &\approx \Delta_1 \left( \log \gamma_\theta(z^t; \mathbf{z}_{\setminus z}^t) \right)^\top \cdot ((z^t + i) - z^t) \\ &\approx \varepsilon_z(z^t)^\top \cdot i \\ \log \gamma_\theta(z^t + j; \mathbf{z}_{\setminus z}^t) - \log \gamma_\theta(z^t; \mathbf{z}_{\setminus z}^t) &\approx \Delta_1 \left( \log \gamma_\theta(z^t; \mathbf{z}_{\setminus z}^t) \right)^\top \cdot ((z^t + j) - z^t) \\ &\approx \varepsilon_z(z^t)^\top \cdot j, \\ \log Q(z^t)_{i,j} &\approx \frac{1}{2} \varepsilon_z(z^t)^\top (i - j). \end{aligned}$$

687 Discrete DCPC would thus parameterize its discrete proposal (Equation 14) in terms of  $\varepsilon_z$  (Equa-  
688 tion 19), so that Equation 15 comes out to the (matrix) exponential of the (elementwise) exponential

$$q_h(z^{t+h} \mid \varepsilon_z) = \exp(Q(\varepsilon_z)h) \quad Q_d(\varepsilon_z)_{i,j} = \exp\left(\frac{(\varepsilon_z)_d^\top (i_d - j_d)}{2}\right).$$

## Supplementary References

- Nicolas Chopin and Omiros Papaspiliopoulos. *An Introduction to Sequential Monte Carlo*. Springer, 2020. ISBN 978-3-030-47844-5. doi: 10.1007/978-3-030-47845-2. Citation Key: Chopin2020 ISSN: 2197-568X.
- Will Grathwohl, Kevin Swersky, Milad Hashemi, David Duvenaud, and Chris Maddison. Oops I took a gradient: Scalable sampling for discrete distributions. In *Proceedings of the 38th International Conference on Machine Learning*, page 3831–3841. PMLR, July 2021. URL <https://proceedings.mlr.press/v139/grathwohl21a.html>.
- John Schulman, Nicolas Heess, Theophane Weber, and Pieter Abbeel. Gradient estimation using stochastic computation graphs. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL [https://proceedings.neurips.cc/paper\\_files/paper/2015/file/de03beffeed9da5f3639a621bcab5dd4-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/de03beffeed9da5f3639a621bcab5dd4-Paper.pdf).
- Haoran Sun, Hanjun Dai, Bo Dai, Haomin Zhou, and Dale Schuurmans. Discrete Langevin Samplers via Wasserstein Gradient Flow. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, Valencia, Spain, April 2023. Proceedings of Machine Learning Research.
- Hao Wu, Heiko Zimmermann, Eli Sennesh, Tuan Anh Le, and Jan Willem van de Meent. Amortized population Gibbs samplers with neural sufficient statistics. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- Yue Xiang, Dongyao Zhu, Bowen Lei, Dongkuan Xu, and Ruqi Zhang. Efficient Informed Proposals for Discrete Distributions via Newton’s Series Approximation. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, pages 7288–7310. PMLR, April 2023. URL <https://proceedings.mlr.press/v206/xiang23a.html>. ISSN: 2640-3498.
- Heiko Zimmermann, Hao Wu, Babak Esmaeili, and Jan-Willem van de Meent. Nested variational inference. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 20423–20435. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/ab49b208848abe14418090d95df0d590-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/ab49b208848abe14418090d95df0d590-Paper.pdf).