

---

# Generalized Doubly-Reparameterized Gradient Estimators

---

Matthias Bauer<sup>1</sup> Andriy Mnih<sup>1</sup>

## Abstract

Efficient low-variance gradient estimation enabled by the reparameterization trick (RT) has been essential to the success of variational autoencoders. Doubly-reparameterized gradients (DREGs) improve on the RT for multi-sample variational bounds by applying reparameterization a second time for an additional reduction in variance. Here, we develop two generalizations of the DREGs estimator and show that they can be used to train conditional and hierarchical VAEs on image modelling tasks more effectively. First, we extend the estimator to hierarchical models with several stochastic layers by showing how to treat additional score function terms due to the hierarchical variational posterior. We then generalize DREGs to score functions of arbitrary distributions instead of just those of the sampling distribution, which makes the estimator applicable to the parameters of the prior in addition to those of the posterior.

## 1. Introduction

In probabilistic machine learning we often optimize expectations of the form  $\mathcal{L}_{\phi, \theta} = \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\phi, \theta}(\mathbf{z})]$  w.r.t. to their parameters, where  $f_{\phi, \theta}(\mathbf{z})$  is some objective function, and  $\phi$  and  $\theta$  denote the parameters of the sampling distribution  $q_{\phi}(\mathbf{z})$  and other (e.g. model) parameters, respectively. In the case of the influential variational autoencoder (VAE, Kingma & Welling (2014), Rezende et al. (2014)),  $q_{\phi}(\mathbf{z})$  is the variational posterior,  $\theta$  denotes the model parameters, and  $\mathcal{L}_{\phi, \theta}$  is typically either the ELBO (Jordan et al., 1999; Blei et al., 2017) or IWAE (Burda et al., 2016) objective.

In most cases of interest, this expectation is intractable, and we estimate it and its gradients,  $\nabla_{\phi} \mathcal{L}$  and  $\nabla_{\theta} \mathcal{L}$ , using Monte Carlo samples  $\mathbf{z} \sim q_{\phi}(\mathbf{z})$ . The resulting gradient estimators are characterized by their *bias* and *variance*. We

usually prefer unbiased estimators as they tend to be better-behaved and are better understood. Lower variance is also preferable because it enables faster training by allowing using higher learning rates.

In this paper, we address gradient estimation for continuous variables in variational objectives. A naive implementation of  $\nabla_{\phi} \mathcal{L}$  results in a *score function*, or REINFORCE, estimator (Williams, 1992), which tends to have high variance; however, if  $f_{\phi, \theta}(\mathbf{z})$  depends on  $\phi$  only through  $\mathbf{z}$ , we can use reparameterization (Kingma & Welling, 2014; Rezende et al., 2014) to obtain an estimator with lower variance by replacing the score function estimator of the gradient with a *pathwise estimator*.

In variational inference,  $f_{\phi, \theta}(\mathbf{z})$  typically depends on  $\phi$  not only through  $\mathbf{z}$  but also through the value of the log density  $\log q_{\phi}(\mathbf{z})$ . Then, the gradient estimators still involve the score function  $\nabla_{\phi} \log q_{\phi}(\mathbf{z})$  despite using reparameterization. Roeder et al. (2017) propose the *sticking the landing* (STL) estimator, which simply drops these score function terms to reduce variance. Tucker et al. (2019) show that STL is biased in general, and introduce the *doubly-reparameterized gradient* (DREGs) estimator for IWAE objectives, which again yields unbiased lower-variance gradient estimates. This is achieved by applying reparameterization a second time, targeting the remaining score function terms.

However, the DREGs estimator has two major limitations: 1) it only applies to latent variable models with a single latent layer; 2) it only applies in cases where the score function depends on the same parameters as the sampling distribution. In this work we address both limitations. Moreover, we show that for hierarchical models with several stochastic layers, gradients that look like pathwise gradients can actually contain additional score function gradients that are not doubly reparameterizable. Despite this, we show that we can still obtain a simple estimator with a sizable reduction in gradient variance for hierarchical IWAE objectives.

Our main contributions are:

- We extend DREGs to hierarchical models;
- We introduce GDREGs, a generalization of DREGs to score functions that depend on a different distribution than the sampling distribution;

---

<sup>1</sup>DeepMind, London, UK. Correspondence to: Matthias Bauer <msbauer@deepmind.com>, Andriy Mnih <andriy@deepmind.com>.

- We show how to implement all proposed gradient estimators using automatic differentiation frameworks;
- We evaluate the proposed DREGs and GDREGs estimators on several conditional and unconditional unsupervised learning problems and find that they outperform the regular IWAE estimator.

## 2. Background

In this work we are interested in computing gradients of variational objectives of the form

$$\mathcal{L}_{\phi, \theta} = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z})} [f_{\phi, \theta}(\mathbf{z})] \quad (1)$$

w.r.t. the variational parameters  $\phi$  of the sampling distribution  $q_{\phi}(\mathbf{z})$ , and parameters  $\theta$  of a second distribution  $p_{\theta}(\mathbf{z})$ , such as a learnable prior. Here  $f_{\phi, \theta}(\mathbf{z})$  is a general function of  $\mathbf{z}$  that can also explicitly depend on both  $q_{\phi}(\mathbf{z})$  and  $p_{\theta}(\mathbf{z})$ . More precisely, we wish to compute

$$\nabla_{\phi}^{\text{TD}} \mathcal{L}_{\phi, \theta} \quad \text{and} \quad \nabla_{\theta}^{\text{TD}} \mathcal{L}_{\phi, \theta}, \quad (2)$$

where  $\nabla_{*}^{\text{TD}}$  denotes the total derivative, which we explicitly distinguish from the partial derivative  $\nabla_{*}$ .

Arguably the simplest objectives of this form are the negative entropy  $\mathcal{L}_{\phi, \theta}^{\text{ent}} = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z})} [\log q_{\phi}(\mathbf{z})]$  and negative cross-entropy  $\mathcal{L}_{\phi, \theta}^{\text{ce}} = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z})} [\log p_{\theta}(\mathbf{z})]$ .

**Importance weighted autoencoders.** Another such objective is the importance weighted autoencoder (IWAE) bound (Burda et al., 2016). For a VAE with likelihood  $p_{\lambda}(\mathbf{x}|\mathbf{z})$ , (learnable) prior  $p_{\theta}(\mathbf{z})$ , and variational posterior (or proposal)  $q_{\phi}(\mathbf{z}|\mathbf{x})$  the IWAE bound with  $K$  importance weights  $w_k = \frac{p_{\theta}(\mathbf{z}_k)p_{\lambda}(\mathbf{x}|\mathbf{z}_k)}{q_{\phi}(\mathbf{z}_k|\mathbf{x})}$  is given by

$$\mathcal{L}_{\phi, \theta}^{\text{IWAE}} = \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left( \frac{1}{K} \sum_{k=1}^K w_k \right) \right]. \quad (3)$$

Eq. (3) reduces to the regular ELBO objective for  $K = 1$  (Rezende et al., 2014; Kingma & Welling, 2014),

$$\mathcal{L}_{\phi, \theta}^{\text{ELBO}} = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{z})p_{\lambda}(\mathbf{x}|\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]. \quad (4)$$

Burda et al. (2016) showed that using multiple importance samples ( $K > 1$ ) provides the model with more flexibility to learn richer representations (fewer inactive units), and results in better log-likelihood estimates compared to VAEs trained with the single sample ELBO. The estimators discussed in this paper build on these results and lead to further improvements. While we focus on the IWAE objective, our proposed GDREGs estimator applies generally.

**Gradient estimation.** In practice, the expectation in Eq. (1) and its gradients are intractable, so we approximate them using Monte Carlo sampling, which makes the estimates of the objective and its gradients random variables.

The resulting gradient estimators will be unbiased but have non-zero variance. We prefer estimators with lower variance, as they enable fast training by allowing higher learning rates.

We can distinguish between two general types of gradient estimators in this setting: (i) *score function estimators* and (ii) *pathwise estimators*. Score functions are gradients of a log probability density w.r.t. its parameters, such as  $\nabla_{\phi} \log q_{\phi}(\mathbf{z})$ ; they treat the function  $f_{\phi, \theta}(\mathbf{z})$  as a black box and often yield high variance gradients. In contrast, pathwise estimators move the parameter-dependence from the probability density into the argument  $\mathbf{z}$  of the function  $f_{\phi, \theta}(\mathbf{z})$  and derive the computation path to often achieve lower variance gradients by using the knowledge of  $\nabla_{\mathbf{z}} f_{\phi, \theta}(\mathbf{z})$ ; see Mohamed et al. (2020) for a recent review.

When computing gradients of the objective  $\mathcal{L}_{\phi, \theta}$  we have to differentiate both the sampling distribution of the expectation,  $q_{\phi}(\mathbf{z})$ , as well as the function  $f_{\phi, \theta}(\mathbf{z})$ ,

$$\nabla_{\phi}^{\text{TD}} \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\phi, \theta}(\mathbf{z})] = \mathbb{E}_{q_{\phi}(\mathbf{z})} [\nabla_{\phi} f_{\phi, \theta}(\mathbf{z}) + f_{\phi, \theta}(\mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z})] \quad (5)$$

$$\nabla_{\theta}^{\text{TD}} \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\phi, \theta}(\mathbf{z})] = \mathbb{E}_{q_{\phi}(\mathbf{z})} [\nabla_{\theta} f_{\phi, \theta}(\mathbf{z})], \quad (6)$$

and both can give rise to score functions. To see that all of the underlined terms indeed contain score functions, note that we can rewrite  $\nabla_{\phi} f_{\phi, \theta}(\mathbf{z})$  as  $\nabla_{\phi} f_{\phi, \theta}(\mathbf{z}) = \nabla_{\log q_{\phi}(\mathbf{z})} f_{\phi, \theta}(\mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z})$  and similarly for  $\nabla_{\theta} f_{\phi, \theta}(\mathbf{z}) = \nabla_{\log p_{\theta}(\mathbf{z})} f_{\phi, \theta}(\mathbf{z}) \nabla_{\theta} \log p_{\theta}(\mathbf{z})$ .

In the following we recapitulate how to address the score functions w.r.t.  $\phi$  in Eq. (5) using the reparameterization trick and doubly-reparameterized gradients (DREGs, Tucker et al. (2019)), respectively. In Sec. 4 we introduce GDREGs, a generalization of DREGs, that eliminates the score function w.r.t.  $\theta$  in Eq. (6).

**Reparameterization.** We can use the *reparameterization trick* (Kingma & Welling, 2014; Rezende et al., 2014) to turn the score function  $\nabla_{\phi} \log q_{\phi}(\mathbf{z})$  inside the expectation in Eq. (5) into a pathwise derivative of the function  $f_{\phi, \theta}(\mathbf{z})$  as follows: we express the latent variables  $\mathbf{z} \sim q_{\phi}(\mathbf{z})$  through a bijection of new random variables  $\epsilon \sim q(\epsilon)$ , which are independent of  $\phi$ ,

$$\mathbf{z} = \mathcal{T}_q(\epsilon; \phi) \Leftrightarrow \epsilon = \mathcal{T}_q^{-1}(\mathbf{z}; \phi). \quad (7)$$

This allows us to rewrite expectations w.r.t.  $q_{\phi}(\mathbf{z})$  as  $\mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\phi, \theta}(\mathbf{z})] = \mathbb{E}_{q(\epsilon)} [f_{\phi, \theta}(\mathcal{T}_q(\epsilon; \phi))]$ , which moves the parameter dependence into the argument of  $f_{\phi, \theta}(\mathbf{z})$  and gives rise to a pathwise gradient:

$$\nabla_{\phi}^{\text{TD}} \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\phi, \theta}(\mathbf{z})] = \mathbb{E}_{q(\epsilon)} [\nabla_{\phi} f_{\phi, \theta}(\mathbf{z}) + \nabla_{\mathbf{z}} f_{\phi, \theta}(\mathbf{z}) \nabla_{\phi} \mathcal{T}_q(\epsilon; \phi)]_{\mathbf{z}=\mathcal{T}_q(\epsilon; \phi)}. \quad (8)$$

In Sec. 3 we discuss that this seemingly pathwise gradient in Eq. (8) can actually contain score functions for more structured or hierarchical models and explain how to extend DREGs to this case. For the remainder of this section we restrict ourselves to simple (single stochastic layer) models.

**Double reparameterization.** Tucker et al. (2019) further reduce gradient variance by replacing the remaining score function in the reparameterized gradient Eq. (8),  $\nabla_{\phi} f_{\phi, \theta}(z) = \nabla_{\log q_{\phi}(z)} f_{\phi, \theta}(z) \nabla_{\phi} \log q_{\phi}(z)$ , with its reparameterized counterpart. Double reparameterization is based on the identity Eq. (9) (Eq. 5 in Tucker et al. (2019)),

$$\begin{aligned} \mathbb{E}_{z \sim q_{\phi}(z)} [g_{\phi, \theta}(z) \nabla_{\phi} \log q_{\phi}(z)] &= \\ &= \mathbb{E}_{\epsilon \sim q(\epsilon)} \left[ \nabla_z^{\text{TD}} g_{\phi, \theta}(z) \Big|_{z=\mathcal{T}_q(\epsilon; \phi)} \nabla_{\phi} \mathcal{T}_q(\epsilon; \phi) \right] \quad (9) \\ &= \mathbb{E}_{z \sim q_{\phi}(z)} \left[ \nabla_z^{\text{TD}} g_{\phi, \theta}(z) \nabla_{\phi} \mathcal{T}_q(\epsilon; \phi) \Big|_{\epsilon=\mathcal{T}_q^{-1}(z; \phi)} \right] \quad (10) \end{aligned}$$

which follows from the fact that both the score function and the reparameterization estimators are unbiased and thus equal in expectation. This identity holds for arbitrary  $g_{\phi, \theta}(z)$ ; to match the score function in Eq. (8) with the LHS of Eq. (9), we have to choose  $g_{\phi, \theta}(z) = \nabla_{\log q_{\phi}(z)} f_{\phi, \theta}(z)$ .

In Eq. (10) we have rewritten the expectation over  $\epsilon \sim q(\epsilon)$  in terms of  $z \sim q_{\phi}(z)$ , as this will become useful for our later generalization. Note how, to compute the pathwise gradient, the sample  $z$  is mapped back to the noise variable  $\epsilon = \mathcal{T}_q^{-1}(z; \phi)$ .  $\nabla_{\phi} \mathcal{T}_q(\epsilon; \phi)$  is also sometimes written as  $\nabla_{\phi} z(\epsilon; \phi)$  (Tucker et al., 2019).

**Gradient estimation for the IWAE objective.** For the IWAE objective Eq. (3), Tucker et al. (2019) derived the following doubly-reparameterized gradients (DREGs) estimator, which supersedes the previously proposed STL estimator (Roeder et al., 2017):

$$\begin{aligned} \widehat{\nabla}_{\phi}^{\text{DREGs}} \mathcal{L}_{\phi, \theta}^{\text{IWAE}} &= \sum_{k=1}^K \tilde{w}_k^2 \nabla_{z_k}^{\text{TD}} \log w_k \nabla_{\phi} \mathcal{T}_q(\epsilon_k; \phi) \quad (11) \\ \widehat{\nabla}_{\phi}^{\text{STL}} \mathcal{L}_{\phi, \theta}^{\text{IWAE}} &= \sum_{k=1}^K \tilde{w}_k \nabla_{z_k}^{\text{TD}} \log w_k \nabla_{\phi} \mathcal{T}_q(\epsilon_k; \phi) \quad (12) \end{aligned}$$

with normalized importance weights  $\tilde{w}_k = \frac{w_k}{\sum_{k'=1}^K w_{k'}}$  and  $\epsilon_{1:K} \sim q(\epsilon)$ . While the DREGs estimator doubly-reparameterizes the score function in Eq. (8), the STL estimator simply drops it and is biased as a result. Crucially, because DREGs relies on reparameterization, it is limited to score functions of the sampling distribution  $q_{\phi}(z)$ , making it inapplicable in the more general setting of arbitrary score functions, such as  $\mathbb{E}_{q_{\phi}(z)} [\nabla_{\theta} f_{\phi, \theta}(z)]$  in Eq. (6).

### 3. DREGs for hierarchical models

We now show that for models with hierarchically structured latent variables even terms that look like pathwise gradients, such as  $\nabla_z^{\text{TD}} f_{\phi, \theta}(z)$  in Eq. (8) or  $\nabla_{z_k}^{\text{TD}} \log w_k$  in the DREGs or STL estimator for the IWAE objective Eqs. (11)

and (12), can give rise to additional score functions. *These additional score functions appear because the distribution parameters of one stochastic layer depend on the latent variables of another layer.* Their appearance is contrary to the intuition that doubly-reparameterized gradient estimators only contain pathwise gradients.

#### 3.1. An illustrative example

To illustrate this, consider a hierarchical model with two layers where we first sample  $z_2 \sim q_{\phi_2}(z_2)$  and then  $z_1 \sim q_{\phi_1}(z_1|z_2)$ .<sup>1</sup> Note that the conditioning on  $z_2$  is through the distribution parameters of  $q_{\phi_1}(z_1|z_2)$ ; to highlight this dependence of  $z_1$  on  $z_2$ , we rewrite  $q_{\phi_1}(z_1|z_2) = q_{\alpha_{1|2}(z_2, \phi_1)}(z_1)$ , where we explicitly distinguish between the *distribution* parameters  $\alpha_{1|2}$ , such as the mean and covariance of a Gaussian, and the *network* parameters  $\phi_1$  that parameterize them together with the previously sampled latent  $z_2$ . A derivative w.r.t.  $z_2$  that looks like a pathwise gradient actually gives rise to a score function term at a subsequent level:

$$\nabla_{z_2}^{\text{TD}} \log q_{\phi_1}(z_1|z_2) = \nabla_{z_2}^{\text{TD}} \log q_{\alpha_{1|2}(z_2, \phi_1)}(z_1) \quad (13)$$

$$= \nabla_{\alpha_{1|2}} \log q_{\alpha_{1|2}}(z_1) \nabla_{z_2} \alpha_{1|2}(z_2, \phi_1) + \dots \quad (14)$$

We omitted (true) pathwise gradients ( $\dots$ ), as the samples  $z_1$  also depend on  $z_2$  through reparameterization. Similar additional score functions arise for seemingly pathwise gradients of hierarchical or autoregressive priors and variational posteriors.

#### 3.2. Extending DREGs to hierarchical VAEs

Here we show how to extend DREGs to hierarchical VAEs to effectively reduce gradient variance for the variational posterior despite the results in the previous section. We still consider the IWAE objective (Eq. (3)), but now the latent space  $z$  is structured, and both  $p_{\theta}$  and  $q_{\phi}$  are hierarchically factorized distributions.

Let us consider a 2-layer VAE  $(z_2) \rightarrow (z_1) \rightarrow (x)$  and examine the term  $\nabla_{\phi_2}^{\text{TD}} \log q_{\phi_1, \phi_2}(z_1, z_2)$ , which appears in the total derivative of the IWAE objective, as a concrete example. We have sampled  $z_1$  and  $z_2$  hierarchically using reparameterization:  $z_2(\phi_2) \equiv \mathcal{T}_{q_2}(\epsilon_2; \alpha_2(\phi_2))$  and  $z_1(\phi_1, \phi_2) \equiv \mathcal{T}_{q_1}(\epsilon_1; \alpha_{1|2}(z_2(\phi_2), \phi_1))$ :

$$\nabla_{\phi_2}^{\text{TD}} \log \left[ q_{\alpha_2(\phi_2)}(z_2(\phi_2)) q_{\alpha_{1|2}(z_2(\phi_2), \phi_1)}(z_1(\phi_1, \phi_2)) \right] \quad (15)$$

When computing total derivatives w.r.t. parameters  $\phi_2$  of the upper layer, we distinguish between three types of gradients: the (true) pathwise gradients w.r.t.  $z_1$  and  $z_2$ , a direct score function because the distribution parameters  $\alpha_2(\phi_2)$

<sup>1</sup>The subscript indices refer to the latent layer indices and not to the importance samples in this case.

directly depend on  $\phi_2$ , and an *indirect* score function because the distribution parameter  $\alpha_{1|2}(z_2(\phi_2), \phi_1)$  indirectly depends on  $\phi_2$  through  $z_2(\phi_2)$ . Indirect score functions do not arise in single stochastic layer models considered by Tucker et al. (2019), and we have three options to estimate them: (1) leave them—this naive estimator is unbiased but potentially has high variance; (2) drop them, similar to STL—this estimator is generally biased; (3) doubly-reparameterize them using DREGs again—this estimator is unbiased, but can generate further indirect score functions.

Total derivatives of other terms in the objective similarly decompose into pathwise gradients as well as direct and indirect score functions. Notably, this includes indirect score functions of the prior  $\log p_{\theta_1, \theta_2}(z_1, z_2)$ , to which DREGs does not apply. In Sec. 4 we introduce the generalized DREGs (GDREGs) estimator that applies in this case.

### 3.3. DREGs for hierarchical IWAE objectives

For IWAE objectives we find that the indirect score functions come up twice: once when computing pathwise gradients of the initial reparameterization, and a second time (with a different prefactor) when computing pathwise gradients for the double-reparameterization of the direct score functions. The same happens for the (true) pathwise gradients, and it is this double-appearance and the resulting cancellation of prefactors that helps reduce gradient variance for DREGs. Moreover, for general model structures it is impossible to replace all successively arising indirect score functions with pathwise gradients, even by applying GDREGs. For example, when the prior and posterior do not factorize in the same way, double-reparameterization of one continuously creates indirect score functions of the other and vice versa, see Apps. C and C.4 for more details.

Thus, to extend DREGs to hierarchical models, we leave the indirect score functions unchanged and only doubly reparameterize the direct score functions. The extended DREGs estimator for IWAE models with arbitrary hierarchical structures is given by Eq. (20)

$$\widehat{\nabla}_{\phi_l}^{\text{DREGs}} \mathcal{L}_{\phi, \theta}^{\text{IWAE}} = \mathbb{E}_{\{\epsilon_{1:K}\}_l \sim q(\epsilon)} \left( \sum_{k=1}^K \tilde{w}_k^2 \nabla_{z_{kl}}^{\text{TD}} \log w_k \nabla_{\phi_l} \mathcal{T}_{q_l}(\epsilon_{kl}; \alpha_l(\text{pa}_{\alpha}(l), \phi_l)) \right) \quad (20)$$

where  $l$  denotes the stochastic layer,  $\text{pa}_{\alpha}(l)$  is the set of latent variables that  $z_{kl}$  depends on, and  $z_{kl} = \mathcal{T}_{q_l}(\epsilon_{kl}; \alpha_l(\text{pa}_{\alpha}(l), \phi_l))$  through reparameterization. We provide a detailed derivation in App. C and a worked example for a VAE with two stochastic layers in App. D. In Apps. E and E.1 we show how to implement this estimator using automatic differentiation by using a *surrogate loss function*, whose forward computation we discard, but whose backward computation exactly corresponds to the estimator in Eq. (20). Alternatively, one could implement a custom gradient for the objective that directly implements Eq. (20); however, we found our approach using a surrogate loss function to be simpler both conceptually and implementation-wise.

Roeder et al. (2017) apply the STL estimator to hierarchical ELBO objectives but do not discuss indirect score functions. Their experimental results are consistent with maintaining the indirect score functions, similar to how we extend DREGs to hierarchical models; the STL estimator is biased for IWAE objectives (Tucker et al., 2019).

## 4. Generalized DREGs

Here, we generalize DREGs to score functions that involve distributions  $p_{\theta}(z)$  different from the sampling distribution  $q_{\phi}(z)$ , as in Eq. (6). In other words, we would like to replace score function terms of the form  $\mathbb{E}_{q_{\phi}(z)} [g_{\phi, \theta}(z) \nabla_{\theta} \log p_{\theta}(z)]$  with doubly-reparameterized pathwise gradients. Such terms appear, for example, when training a VAE with a trainable prior  $p_{\theta}(z)$  with the ELBO or IWAE objectives. DREGs cannot be used directly in this case as it relies on reparameterization of the sampling distribution  $q_{\phi}(z)$ , so that the path depends on the parameters  $\phi$ ,

$$\mathbb{E}_{z \sim q_{\phi}(z)} [g_{\phi, \theta}(z) \nabla_{\theta} \log p_{\theta}(z)] = \mathbb{E}_{z \sim q_{\phi}(z)} \left[ \left( g_{\phi, \theta}(z) \nabla_z^{\text{TD}} \log \frac{q_{\phi}(z)}{p_{\theta}(z)} + \nabla_z^{\text{TD}} g_{\phi, \theta}(z) \right) \nabla_{\theta} \mathcal{T}_p(\tilde{\epsilon}; \theta) \Big|_{\tilde{\epsilon} = \mathcal{T}_p^{-1}(z, \theta)} \right] \quad (16)$$

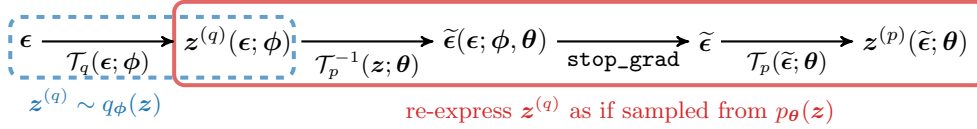
$$\stackrel{\textcircled{1}}{=} \nabla_{\theta}^{\text{TD}} \mathbb{E}_{q_{\phi}(z)} [g_{\phi, \theta}(z)] \stackrel{\textcircled{1}}{=} \nabla_{\theta}^{\text{TD}} \mathbb{E}_{p_{\theta}(z)} \left[ \frac{q_{\phi}(z)}{p_{\theta}(z)} g_{\phi, \theta}(z) \right] \stackrel{\textcircled{1}}{=} \nabla_{\theta}^{\text{TD}} \mathbb{E}_{q(\tilde{\epsilon})} \left[ \frac{q_{\phi}(\mathcal{T}_p(\tilde{\epsilon}; \theta))}{p_{\theta}(\mathcal{T}_p(\tilde{\epsilon}; \theta))} g_{\phi, \theta}(\mathcal{T}_p(\tilde{\epsilon}; \theta)) \right]; \quad q(\tilde{\epsilon}) = \mathcal{N}(0, \mathbb{I}) \quad (17)$$

$$\stackrel{\textcircled{2}}{=} \mathbb{E}_{q(\tilde{\epsilon})} \left[ \nabla_z^{\text{TD}} \left( \frac{q_{\phi}(z)}{p_{\theta}(z)} g(z) \right) \nabla_{\theta} \mathcal{T}_p(\tilde{\epsilon}; \theta) + \frac{q_{\phi}(z)}{p_{\theta}(z)} (\nabla_{\theta} g_{\phi, \theta}(z) - g(z) \nabla_{\theta} \log p_{\theta}(z)) \right]_{z = \mathcal{T}_p(\tilde{\epsilon}; \theta)} \quad (18)$$

$$\stackrel{\textcircled{3}}{=} \mathbb{E}_{q_{\phi}(z)} \left[ \left( g(z) \nabla_z^{\text{TD}} \log \frac{q_{\phi}(z)}{p_{\theta}(z)} + \nabla_z^{\text{TD}} g(z) \right) \nabla_{\theta} \mathcal{T}_p(\tilde{\epsilon}; \theta) \Big|_{\tilde{\epsilon} = \mathcal{T}_p^{-1}(z, \theta)} + \nabla_{\theta} g_{\phi, \theta}(z) - g(z) \nabla_{\theta} \log p_{\theta}(z) \right] \quad (19)$$

**Figure 1:** The GDREGs identity and a brief derivation in three steps:  $\textcircled{1}$  temporarily change the path so that it depends on  $\theta$ ;  $\textcircled{2}$  perform the reparameterized gradient computation;  $\textcircled{3}$  change the path back so we can use samples  $z \sim q_{\phi}(z)$  to estimate the expectation. See App. A for details and an alternative derivation using the change of density formula.





**Figure 2:** Computational flow to re-express a sample  $z$  from  $q_\phi(z)$  as if it were sampled from  $p_\theta(z)$ . Its numerical value and distribution remain unchanged but the pathwise gradient through it now depends on  $\theta$ :  $\nabla_\theta \mathcal{T}_p(\tilde{\epsilon}; \theta)|_{\tilde{\epsilon}=\mathcal{T}_p^{-1}(z, \theta)}$ .  $\tilde{\epsilon} = \mathcal{T}_p^{-1}(\mathcal{T}_q(\epsilon; \phi); \theta)$  follows a different, usually more complex, distribution from  $\epsilon \sim q(\epsilon)$ .

whereas the score function is with w.r.t. parameters  $\theta$  of a different distribution  $p_\theta(z)$ .

To make progress *we need to make the path depend on the parameters  $\theta$*  while still sampling  $z \sim q_\phi(z)$  during training. Our solution consists of three steps (also see Fig. 1):

- ① temporarily change the path so that it depends on  $\theta$ ;
- ② perform the reparameterized gradient computation;
- ③ change the path back so we can use samples  $z \sim q_\phi(z)$  to estimate the expectation.

We change the path by first using an importance sampling reweighting to temporarily re-write the expectation,  $\mathbb{E}_{q_\phi(z)}[*] = \mathbb{E}_{p_\theta(z)}\left[\frac{q_\phi(z)}{p_\theta(z)}*\right]$ , and then employing reparameterization on the new sampling distribution  $p_\theta(z)$ :  $z = \mathcal{T}_p(\tilde{\epsilon}; \theta)$  with  $\tilde{\epsilon} \sim q(\tilde{\epsilon})$ . Following this recipe, we derive the gradient identity in Eq. (16) for general  $g_{\phi, \theta}(z)$  that we refer to as generalized DREGs (or GDREGs for short) identity.

Like DREGs (Eq. (10)), GDREGs allows us to transform score functions into pathwise gradients. Yet, unlike DREGs, GDREGs applies to general score functions and contains a correction term that vanishes when  $p_\theta(z)$  and  $q_\phi(z)$  are identical ( $\log \frac{q_\phi(z)}{p_\theta(z)}$  term in Eq. (16)).

Note that the pathwise derivative  $\nabla_\theta \mathcal{T}_p(\tilde{\epsilon}; \theta)$  in Eq. (16) looks like we reparameterized an independent noise variable  $\tilde{\epsilon}$  using  $p_\theta(z)$ , where the numerical value of the noise variable is given by  $\tilde{\epsilon} = \mathcal{T}_p^{-1}(z; \theta)$  and  $z \sim q_\phi(z)$ . We can interpret this sequence of transformations as a normalizing flow (Rezende & Mohamed, 2015)  $z \rightarrow \tilde{\epsilon} \rightarrow z$ , such that  $\mathcal{T}_p(\tilde{\epsilon}; \theta) = \mathcal{T}_p(\mathcal{T}_p^{-1}(z; \theta); \theta) = z$ . We can think of this procedure as *re-expressing the sample  $z \sim q_\phi(z)$  as if it came from  $p_\theta(z)$* : Its numerical value  $z$  remains unchanged and it is still distributed according to  $q_\phi(z)$ , yet its pathwise gradient  $\nabla_\theta \mathcal{T}_p(\tilde{\epsilon}; \theta)$  depends on  $\theta$ . We illustrate the corresponding computational flow in Fig. 2 and provide an example implementation with code in App. F

Note that to derive the GDREGs identity, we only require  $p_\theta(z)$  to be reparameterizable (red box in Fig. 2). While  $q_\phi(z)$  may be reparameterizable as well (dashed blue box in Fig. 2), this is not necessary; we only need to be able to evaluate its density in Eq. (16).

In the simplest case of the cross-entropy objective  $\mathcal{L}^{\text{ce}} = \mathbb{E}_{q_\phi(z)}[\log p_\theta(z)]$  (as in the ELBO with a sample-based KL estimate),  $g_{\phi, \theta}(z) = 1$ , and the GDREGs identity Eq. (16) gives rise to the following GDREGs estimator:

$$\widehat{\nabla}_\theta^{\text{GDREGs}} \mathcal{L}^{\text{ce}} = \nabla_z \log \frac{q_\phi(z)}{p_\theta(z)} \nabla_\theta \mathcal{T}_p(\tilde{\epsilon}; \theta) \Big|_{\tilde{\epsilon}=\mathcal{T}_p^{-1}(z, \theta)} \quad (21)$$

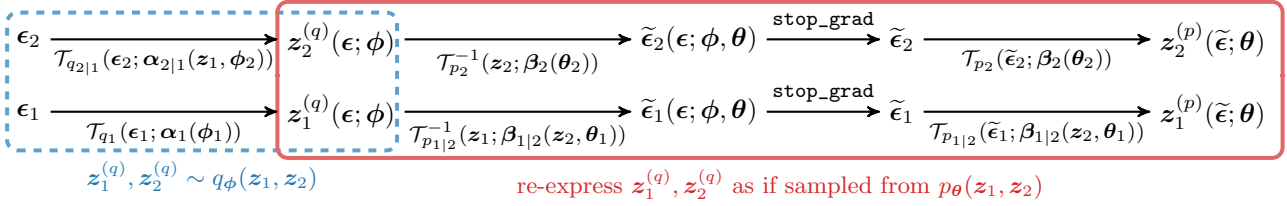
with  $z \sim q_\phi(z)$ . For Gaussian distributions  $q_\phi(z)$  and  $p_\theta(z)$ , the cross-entropy and its gradients can be computed in closed form, which we can think of as a perfect estimator with zero bias and zero variance. Moreover, the expectation and variance of both the naive score function as well as the GDREGs estimator in Eq. (21) can be computed in closed form. We provide full derivations and a discussion of this special case in App. H as well as an example implementation in terms of (pseudo-)code in App. F. The main results are: (i) GDREGs has lower variance gradients than the score function when  $q_\phi(z)$  and  $p_\theta(z)$  overlap substantially, which is typically the case at the beginning of training; (ii) we can derive a closed-form control variate that depends on a ratio of the means and variances of the two distributions and that is strictly superior to the naive score function estimator and the GDREGs estimator in terms of gradient variance. However, the analytic expression for the cross-entropy has even lower (zero) gradient variance in this case.

#### 4.1. GDREGs for VAE objectives

We can now use the GDREGs identity Eq. (16) to derive generalized doubly-reparameterized estimators for expectations of general score functions of the form  $\mathbb{E}_{q_\phi(z)}[g_{\phi, \theta}(z) \nabla_\theta \log p_\theta(z)]$ , also see Eq. (6). In App. B we derive the following GDREGs estimator of the IWAE objective w.r.t. the prior parameters  $\theta$ :

$$\widehat{\nabla}_\theta^{\text{GDREGs}} \mathcal{L}_{\phi, \theta}^{\text{IWAE}} = \sum_{k=1}^K (\tilde{w}_k \nabla_{z_k}^{\text{TD}} \log p_\lambda(\mathbf{x}|z_k) - \tilde{w}_k^2 \nabla_{z_k}^{\text{TD}} \log w_k) \nabla_\theta \mathcal{T}_p(\tilde{\epsilon}_k; \theta) \Big|_{\tilde{\epsilon}_k=\mathcal{T}_p^{-1}(z_k, \theta)} \quad (22)$$

with  $z_{1:K} \sim q_\phi(z|\mathbf{x})$ . The second term in Eq. (22) looks like the DREGs estimator for  $\phi$  in Eq. (11) except that the samples  $z_k$  are now re-expressed as if they came from  $p_\theta(z)$ . In addition we obtain a term that involves the likelihood  $p_\lambda(\mathbf{x}|z)$  and is linear in  $\tilde{w}_k$ . Note that we do not apply GDREGs to the likelihood parameters  $\lambda$  because  $p_\lambda(\mathbf{x}|z)$



**Figure 3:** Computational flow to re-express samples  $z_1, z_2$  from  $q_\phi(z_1, z_2) = q_{\phi_1}(z_1)q_{\phi_2}(z_2|z_1)$  as if they were sampled from  $p_\theta(z_1, z_2) = p_{\theta_2}(z_2)p_{\theta_1}(z_1|z_2)$ . Their numerical values and distribution remain unchanged but the gradient flow through them changes. Note that  $\tilde{\epsilon}_i$  follows a different, usually more complex, distribution from  $\epsilon_i$ .  $\alpha_i$  and  $\beta_i$  denote the distribution parameters of the variational posterior and the prior, respectively.

is a distribution over  $\mathbf{x}$  rather than  $\mathbf{z}$ ; in the following we therefore drop the subscript  $\lambda$ .

We learn all parameters by optimizing the same objective function Eq. (3), but employ different gradient estimators for different subsets of parameters. In practice, we implement these estimators using different *surrogate objectives* for the likelihood, proposal, and prior parameters, see App. E for details. While separate objectives seem computationally expensive, most terms are shared between them, and modern frameworks avoid such duplicate computation. In practice, we found the runtime increase for training with DREGs and GDREGs estimators to be smaller than 10% without any optimization of the implementation.

## 4.2. Extending GDREGs to hierarchical VAEs

When extending GDREGs to hierarchical models, we again encounter direct and indirect score functions (see Sec. 3). Like for the posterior parameter  $\phi$  we apply GDREGs to the direct score functions but leave the indirect score functions unchanged. The full GDREGs estimator for IWAE objectives with arbitrary hierarchical structure is given in App. C Eq. (C.16), see App. C for a derivation. In App. D we provide a worked example and in App. E we again show how to use surrogate losses to implement the estimator in practice.

To apply GDREGs we need to re-express samples from  $q_\phi(\mathbf{z})$  as if they came from  $p_\theta(\mathbf{z})$ . We do this for the entire hierarchy jointly. In Fig. 3 we illustrate the necessary computational flow for the example of a 2-layer VAE with the variational posterior factorized in the opposite direction from the generative process; see App. C for the general case. We draw samples  $z_1, z_2 \sim q_\phi(z_1, z_2) = q_{\phi_1}(z_1)q_{\phi_2}(z_2|z_1)$  (by transforming independent variables  $\epsilon_i$ ) and then re-express them as if they were sampled from the prior  $p_{\theta_2}(z_2)p_{\theta_1}(z_1|z_2)$ , which factorizes in the opposite direction. While the numerical values of  $z_1$  and  $z_2$  remain unchanged,  $z_1$  is now dependent on  $z_2$  and both depend on the respective  $\theta$  parameters when computing gradients; we can view  $(z_1, z_2)$  as samples that were obtained by transforming independent variables  $(\tilde{\epsilon}_1, \tilde{\epsilon}_2)$  that follow

a more complicated distribution than  $(\epsilon_1, \epsilon_2)$ . As in the single-layer case, only  $p_\theta(\mathbf{z})$  needs to be reparameterizable.

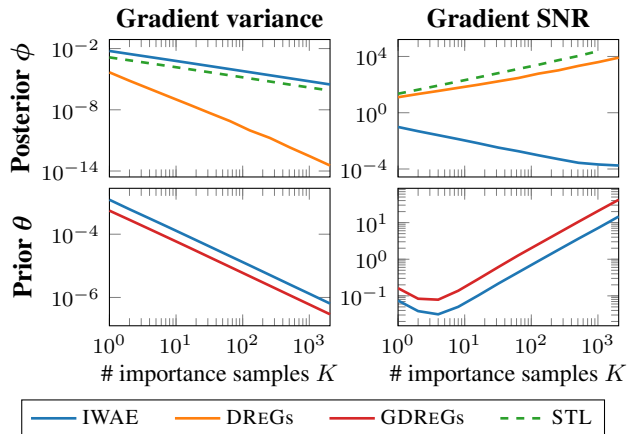
## 5. Experiments

In this section we empirically evaluate the hierarchical extension of DREGs and its generalization to GDREGs, and compare them to the naive IWAE gradient estimator (labelled as IWAE) as well as STL (Roeder et al., 2017). First, we illustrate that DREGs and GDREGs increase the gradient signal-to-noise ratio (SNR) and reduce gradient variance compared to the naive estimator on a simple hierarchical example (Sec. 5.1); second, we show that they also reduce gradient variance in practice and improve test performance on several generative modelling tasks with VAEs with one or more stochastic layers (Sec. 5.2). We highlight that both the extension of DREGs to more than one stochastic layer as well as training the prior with GDREGs are novel contributions of this work.

### 5.1. Illustrative example: linear VAE

We first consider an extended version of the illustrative example by Rainforth et al. (2018) and Tucker et al. (2019) to show that hierarchical DREGs and GDREGs increase the gradient signal-to-noise ratio (SNR) and reduce gradient variance compared to the naive IWAE gradient estimator.

We consider a 2-layer linear VAE with hierarchical prior  $z_2 \sim \mathcal{N}(0, \mathbb{I})$ ,  $z_1|z_2 \sim \mathcal{N}(\mu_\theta(z_2), \sigma_\theta^2(z_2))$ , Gaussian noise likelihood  $\mathbf{x}|z_1 \sim \mathcal{N}(z_1, \mathbb{I})$ , and bottom up variational posterior  $q_{\phi_1}(z_1|\mathbf{x}) = \mathcal{N}(\mu_{\phi_1}(\mathbf{x}), \sigma_{\phi_1}^2(\mathbf{x}))$  and  $q_{\phi_2}(z_2|z_1) = \mathcal{N}(\mu_{\phi_2}(z_1), \sigma_{\phi_2}^2(z_1))$ . All  $\mu_i$  and  $\sigma_i$  are linear functions, and  $z_1, z_2, \mathbf{x} \in \mathbb{R}^D$ . We sample 512 datapoints in  $D = 5$  dimensions from a model with  $\mu_\theta(z_2) = z_2$  and  $\sigma_\theta(z_2) = 1$ . We then train the parameters  $\phi$  and  $\theta$  using SGD and the IWAE objective til convergence and evaluate the gradient variance and signal-to-noise ratio for each estimator. For the proposal parameters  $\phi$  we compare DREGs to the naive score function (labelled as IWAE) and to STL; for the prior parameters  $\theta$  we compare GDREGs to IWAE.



**Figure 4:** Average gradient variance (*left*) and signal-to-noise ratio (SNR) (*right*) for the proposal parameters  $\phi$  (*top*) and the prior parameters  $\theta$  (*bottom*).

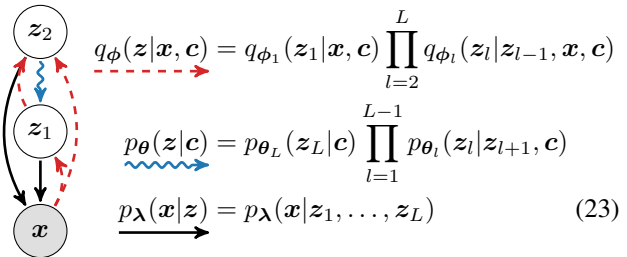
We find that our extension of DREGs to hierarchical models behaves qualitatively the same as in the single layer case considered by Tucker et al. (2019), see Fig. 4 (*top*): While the SNR for the naive estimator vanishes with increasing number of importance samples, the SNR increases for DREGs. This can be explained by the faster rate with which the gradient variance decreases for DREGs compared to IWAE and STL. While STL has an even better SNR, its gradients are biased.

When considering gradients w.r.t. the prior parameters we find that the SNR is higher and the gradient variance is lower for GDREGs compared to the naive estimator (IWAE), see Fig. 4 (*bottom*). However, they grow and shrink at the same rate for both estimators as the number of importance samples is increased.

## 5.2. Image modelling with VAEs

In the remainder of this paper we consider image modelling tasks with VAEs on several standard benchmark datasets: MNIST (LeCun & Cortes, 2010), Omniglot (Lake et al., 2015), and FashionMNIST (Xiao et al., 2017). We use dynamically binarized versions of all datasets to minimize overfitting.

We consider both single layer and hierarchical (multi-layer) VAEs and evaluate them on unconditional and conditional modelling tasks using the IWAE objective, Eq. (3). In the hierarchical case, the generative path (prior and likelihood) is top-down whereas the variational posterior is bottom-up, see Fig. 5 and Eq. (23) for a full description of the model and a 2-layer unconditional example. For conditional modelling we predict the bottom half of an image given its top half, as in Tucker et al. (2019); in this case, both the prior and variational posterior also depend on a context variable



**Figure 5:** Model specification and 2-layer example for conditional and unconditional image modelling.

$c$ ,  $q_\phi(z|x, c)$  and  $p_\theta(z|c)$ , respectively. We use a factorized Bernoulli likelihood along with factorized Gaussians for the variational posterior and prior. Every conditional distribution in Eq. (23) is parameterized by an MLP with two hidden layers of 300  $\tanh$  units each, and all latent spaces have 50 dimensions. Unless stated otherwise, we train all models for 1000 epochs using the Adam optimizer (Kingma & Ba, 2015) with default learning rate of  $3 \cdot 10^{-4}$ , a batch size of 64, and  $K = 64$  importance samples; see App. G for details.

As mentioned in Sec. 4.1, we use separate surrogate objectives to compute the gradient estimators for the likelihood, posterior, and prior parameters. While we always train the likelihood parameters  $\lambda$  on the naive IWAE objective, we consider the naive IWAE estimator (labelled as IWAE), STL, and DREGs for the variational posterior parameters  $\phi$ , and IWAE and GDREGs for the prior parameters  $\theta$ . See App. E for details on the implementation of the estimators. We present the results for conditional modelling in Tab. 1 and Fig. 6, and for unconditional modelling in Tab. 2 and Fig. 7; see App. G for more experimental results.

**Estimators for the variational parameters  $\phi$ .** First, we evaluate the choice of estimator for the parameters of  $q_\phi(z)$ . Like Tucker et al. (2019) for the single layer case, we find that our extension of DREGs to hierarchical models leads to a dramatic reduction in gradient variance for the variational posterior parameters  $\phi$  on all tasks (third column in Figs. 6 and 7), which translates to an improved test objective in all cases considered. DREGs is unbiased and typically outperforms the (biased) STL estimator. We also observed similar improvements on the training objective.

**Estimators for the prior parameters  $\theta$ .** Second, we consider the estimators for the  $\theta$  parameters of the prior  $p_\theta(z)$ . Using the GDREGs estimator instead of the naive IWAE estimator consistently improves the train and test objective when combined with *any* estimator for the variational posterior, especially for conditional image modelling with deeper models. For unconditional image modelling the improvements are only marginal, though using GDREGs never hurts. In terms of gradient variance for the prior parameters  $\theta$ , GDREGs consistently performs better in the beginning of

estimator $\nabla_{\phi}^{\text{TD}}$	estimator $\nabla_{\theta}^{\text{TD}}$	IWAE		STL		DREGs	
		IWAE	GDREGs	IWAE	IWAE	IWAE	GDREGs
MNIST	1 layer	-38.77±0.01	-38.71±0.02	-38.76±0.03	-38.68±0.03	-38.50±0.01	<b>-38.44±0.01</b>
	2 layer	-38.55±0.02	-38.42±0.03	-38.24±0.02	-38.14±0.02	-38.20±0.01	<b>-38.02±0.02</b>
	3 layer	-38.63±0.01	-38.44±0.02	-38.20±0.01	-38.10±0.02	-38.20±0.01	<b>-38.04±0.01</b>
Omniglot	1 layer	-55.84±0.02	-55.66±0.03	-55.80±0.05	-55.62±0.05	-55.34±0.02	<b>-55.24±0.02</b>
	2 layer	-55.27±0.03	-54.98±0.02	-54.66±0.03	<b>-54.28±0.02</b>	-54.73±0.02	-54.36±0.03
	3 layer	-55.35±0.02	-54.93±0.02	-54.64±0.03	<b>-54.21±0.03</b>	-54.72±0.02	-54.28±0.02
FMNIST	1 layer	-102.84±0.02	-102.80±0.02	-102.99±0.02	-102.88±0.02	-102.61±0.01	<b>-102.58±0.01</b>
	2 layer	-102.74±0.02	-102.68±0.01	-102.65±0.02	-102.48±0.03	-102.40±0.01	<b>-102.30±0.02</b>
	3 layer	-102.86±0.01	-102.71±0.01	-102.68±0.01	-102.42±0.02	-102.46±0.01	<b>-102.26±0.01</b>

Table 1: Test objective values (higher is better) on *conditional* image modelling with a VAE model trained with IWAE. Higher is better; errorbars denote  $\pm 1.96$  standard errors ( $\sigma/\sqrt{5}$ ) over 5 reruns.

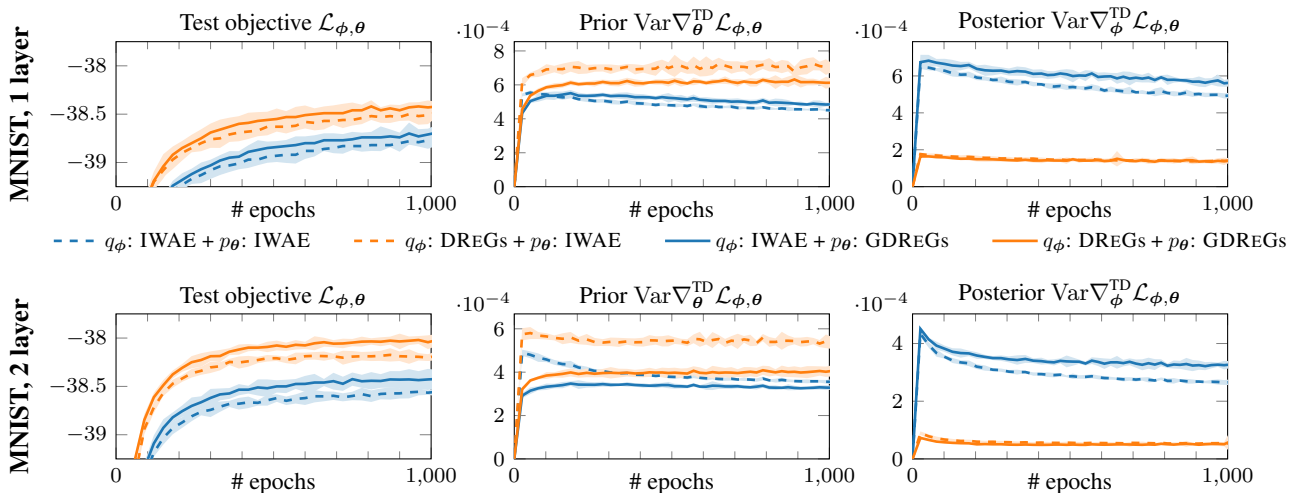


Figure 6: *Conditional* image modelling of MNIST with a VAE with 1 layer (top) and 2 layers (bottom). Shaded areas denote  $\pm 1.96$  standard deviations  $\sigma$  over 5 reruns.

training, when it always has lower variance. However, later in training this is only consistently true when also using the DREGs estimator for the variational posterior parameters  $\phi$ . We hypothesize that the GDREGs estimator yields larger improvements for conditional modelling because the prior and posterior distribution are closer to each other due to the conditioning, and we saw that GDREGs works particularly well in this case for Gaussian distributions, also see App. H. To quantify this “closeness” we compared the KL of the variational posterior to the prior on the same dataset and found it to be about twice as large for unconditional modelling than for conditional modelling, see App. G.

We also note that the gradient variance for the prior parameters  $\theta$  is higher when using the DREGs estimator for the variational posterior parameters  $\phi$ , compared to the naive IWAE estimator (compare orange and blue lines in the middle column of Figs. 6 and 7). This is an indirect effect of altered learning dynamics. We suspect that better posterior gradient estimates with DREGs lead to generative models

that fit the data better, which in turn results in larger gradient variance for the prior. This effect is absent in the illustrative example in Fig. 4 because we evaluate the gradient variance on the same fixed model for all estimators. In App. G.3 we compare the estimators *offline* for different combinations of estimators during training. The results are in line with our online results in this section: for the gradients of the variational posterior the DREGs estimator *always* has lower variance than the naive (IWAE) estimator; for the gradients of the prior the GDREGs estimator typically has lower variance, though in some cases only in the beginning of training.

## 6. Related work

Roeder et al. (2017) observed that the reparameterization gradient estimator for the ELBO contains a score function term and proposed the STL estimator that simply drops this term to reduce the estimator variance. They considered



estimator	$\nabla_{\phi}^{\text{TD}}$ estimator $\nabla_{\theta}^{\text{TD}}$	IWAE		STL		DREGs	
		IWAE	GDREGs	IWAE	GDREGs	IWAE	GDREGs
MNIST	2 layer	-86.07±0.02	-86.04±0.03	-85.29±0.02	<b>-85.23±0.03</b>	<b>-85.25±0.02</b>	-85.32±0.02
	3 layer	-85.69±0.02	-85.70±0.02	-85.01±0.03	-84.94±0.05	<b>-84.87±0.03</b>	<b>-84.90±0.04</b>
Omniglot	2 layer	-105.20±0.02	-105.11±0.02	-104.10±0.05	<b>-104.00±0.05</b>	-104.12±0.05	<b>-104.05±0.04</b>
	3 layer	-104.68±0.02	-104.71±0.03	-104.02±0.02	<b>-103.55±0.03</b>	-104.71±0.03	<b>-103.51±0.06</b>
FMNIST	2 layer	-230.65±0.03	-230.61±0.02	-230.14±0.02	<b>-229.98±0.02</b>	-230.04±0.03	<b>-229.98±0.03</b>
	3 layer	-230.60±0.03	-230.59±0.03	-230.26±0.04	-229.92±0.03	-229.92±0.02	<b>-229.87±0.03</b>

Table 2: Test objective values on *unconditional* image modelling with a VAE model trained with IWAE.

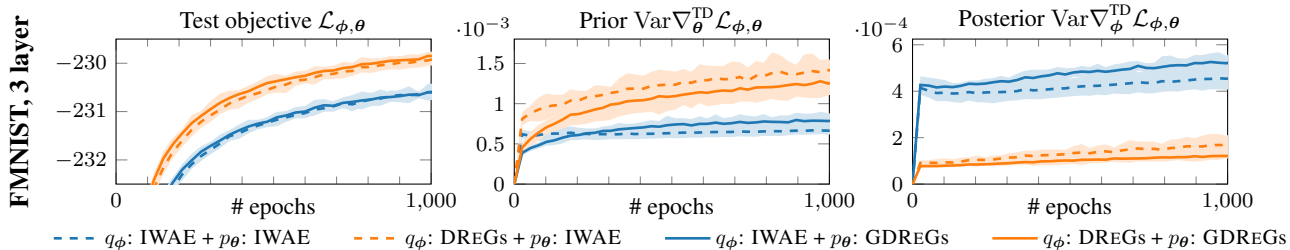


Figure 7: Unconditional image modelling on FashionMNIST; 3 layers.

hierarchical ELBO models but do not discuss how to treat indirect score functions. While the STL estimator is unbiased for the ELBO objective, Tucker et al. (2019) showed that it is biased for more general objectives such as the IWAE. They proposed the DREGs estimator that yields unbiased and low variance gradients for IWAE and resolves the diminishing signal-to-noise issue of the naive IWAE gradients first discussed by Rainforth et al. (2018). We extend DREGs to hierarchical models, discuss how to treat the indirect score functions, and generalize it to general score functions by introducing GDREGs.

Several classic techniques from the variance reduction literature have been applied to variational inference and reparameterization. For example, Miller et al. (2017) and Geffner & Domke (2020) proposed control variates for reparameterization gradients; Ruiz et al. (2016) used importance sampling with a proposal optimized to reduce variance. Such approaches are orthogonal to methods such as (G)DREGs and STL, and can be combined with them for greater variance reduction (Agrawal et al., 2020).

### 7. Conclusion

In this paper we generalized the recently proposed doubly-reparameterized gradients (DREGs, Tucker et al. (2019)) estimator for variational objectives in two ways. First, we showed that for hierarchical models such as VAEs seemingly pathwise gradients can actually contain score functions, and how to consistently and effectively extend DREGs in this case. Second, we introduced GDREGs, a doubly-reparameterized gradient estimator that applies to general

score functions, while DREGs is limited to score functions of the variational distribution. Finally, we demonstrated that both generalizations can improve performance on conditional and unconditional image modelling tasks.

While we present and discuss the GDREGs estimator in the context of deep probabilistic models, it applies generally to score function gradients of the form  $\mathbb{E}_{q_{\phi}(z)} [\nabla_{\theta} \log p_{\theta}(z)]$ . Applying it to other problem settings of this type such as normalizing flows is an exciting area of future research.

### Acknowledgements

We thank Chris Maddison as well as the anonymous reviewers for feedback on the manuscript.

### References

Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, 2016.

Agrawal, A., Sheldon, D. R., and Domke, J. Advances in black-box VI: normalizing flows, importance weighting, and optimization. In *Advances in Neural Information Processing Systems 33*, 2020.

Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 2017.

Bradbury, J., Frostig, R., Hawkins, P., Johnson, M. J., Leary,

- C., Maclaurin, D., and Wanderman-Milne, S. JAX: composable transformations of Python+NumPy programs, 2018.
- Burda, Y., Grosse, R. B., and Salakhutdinov, R. Importance weighted autoencoders. In *Proceedings of the 4th International Conference on Learning Representations*, 2016.
- Dillon, J. V., Langmore, I., Tran, D., Brevdo, E., Vasudevan, S., Moore, D., Patton, B., Alemi, A., Hoffman, M. D., and Saurous, R. A. Tensorflow distributions. *arXiv preprint arXiv:1711.10604*, 2017.
- Geffner, T. and Domke, J. Approximation based variance reduction for reparameterization gradients. In *Advances in Neural Information Processing Systems 33*, 2020.
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Rio, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. Array programming with numpy. *Nature*, Sep 2020.
- Hennigan, T., Cai, T., Norman, T., and Babuschkin, I. Haiku: Sonnet for JAX, 2020.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. *An Introduction to Variational Methods for Graphical Models*. MIT Press, Cambridge, MA, USA, 1999.
- Kingma, D. and Ba, J. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations*, 2015.
- Kingma, D. and Welling, M. Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations*, 2014.
- Lake, B. M., Salakhutdinov, R., and Tenenbaum, J. B. Human-level concept learning through probabilistic program induction. *Science*, 2015.
- LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Miller, A., Foti, N., D’Amour, A., and Adams, R. P. Reducing reparameterization gradient variance. In *Advances in Neural Information Processing Systems*, 2017.
- Mohamed, S., Rosca, M., Figurnov, M., and Mnih, A. Monte carlo gradient estimation in machine learning. *Journal of Machine Learning Research*, 2020.
- Rainforth, T., Kosiorek, A., Le, T. A., Maddison, C., Igl, M., Wood, F., and Teh, Y. W. Tighter variational bounds are not necessarily better. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- Rezende, D. J. and Mohamed, S. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning*, 2015.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, 2014.
- Roeder, G., Wu, Y., and Duvenaud, D. Sticking the landing: Simple, lower-variance gradient estimators for variational inference. In *Advances in Neural Information Processing Systems 30*, 2017.
- Ruiz, F. J. R., Titsias, M. K., and Blei, D. M. Overdispersed black-box variational inference. In *Uncertainty in Artificial Intelligence*, 2016.
- Tucker, G., Lawson, D., Gu, S., and Maddison, C. J. Doubly reparameterized gradient estimators for monte carlo objectives. In *Proceedings of the 7th International Conference on Learning Representations*, 2019.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 1992.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms, 2017.