# Generalized Doubly-Reparameterized Gradient Estimators

**Matthias Bauer**                                                    msbauer@google.com
**Andriy Mnih**                                                        amnih@google.com
*DeepMind, London, UK*

## 1. Introduction

In probabilistic machine learning, we are often interested in optimizing expectations of the form $\mathcal{L}_{\phi,\theta} = \mathbb{E}_{q_\phi(z)}[f_{\phi,\theta}(z)]$ w.r.t. to their parameters, where $f_{\phi,\theta}(z)$ is some objective function, and $\phi$ and $\theta$ denote the parameters of the sampling distribution $q_\phi(z)$ and other (e.g. model) parameters, respectively. In the case of widely used variational autoencoders (VAEs, Kingma and Welling (2014), Rezende et al. (2014)), $q_\phi(z)$ is the variational posterior and $\theta$ denotes the model parameters.

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 $z \sim q_\phi(z)$. 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 gradient variance; however, if $f$ depends on $\phi$ only through $z$, we can use reparameterization (Kingma and 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, $f_{\phi,\theta}(z)$ typically depends on $\phi$ not only through $z$ but also through the value of the log density $\log q_\phi(z)$. Then, the gradient estimators still involve the score function $\nabla_\phi \log q_\phi(z)$ 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 for the second time, targeting the score function terms. The DREGs estimator has, however, 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 and introduce GDREGs, a generalized doubly-reparameterized gradient estimator that applies to general score functions; we also extend the DREGs estimator to hierarchical models and show that both estimators improve performance on conditional and unconditional image modelling tasks.

## 2. Background

We are interested in computing gradients of variational objectives $\mathcal{L}_{\phi,\theta} = \mathbb{E}_{z \sim q_\phi(z)}[f_{\phi,\theta}(z)]$ w.r.t. the variational parameters $\phi$ of the sampling distribution $q_\phi(z)$, and parameters $\theta$ of a second distribution $p_\theta(z)$, such as a learnable prior. Here $f_{\phi,\theta}(z)$ is a general function that can depend on both $q_\phi(z)$ and $p_\theta(z)$ explicitly.

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

$$\mathcal{L}_{\phi,\theta}^{\text{IWAE}} = \mathbb{E}_{z_1,\ldots,z_K \sim q_\phi(z|x)}\left[\log\left(\frac{1}{K}\sum_{k=1}^{K} w_k\right)\right] \qquad w_k = \frac{p_\theta(z_k)p_\lambda(x|z_k)}{q_\phi(z_k|x)} \qquad (1)$$

**Gradient estimation.** In practise, we approximate the expectation in $\mathcal{L}_{\phi,\theta}$ by Monte Carlo sampling, so that our estimates of the expectation and of its gradients become random variables. We can distinguish between two types of gradient estimators in this setting: (i) *score function (SF) estimators* and (ii) *pathwise estimators*. Score functions are gradients of a log probability density w.r.t. its parameters, for example $\nabla_\phi \log q_\phi(z)$; SF estimators treat the function $f_{\phi,\theta}$ 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 $z$ of the function $f_{\phi,\theta}(z)$ and differentiate the computation path to often achieve lower variance gradients by using the knowledge of $\nabla_z f_{\phi,\theta}(z)$; 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, $q_\phi(z)$, as well as the function $f_{\phi,\theta}(z)$,

$$\nabla_\phi^{\text{TD}} \mathbb{E}_{q_\phi(z)}[f_{\phi,\theta}(z)] = \mathbb{E}_{q_\phi(z)}\left[\nabla_\phi f_{\phi,\theta}(z) + f_{\phi,\theta}(z)\nabla_\phi \log q_\phi(z)\right] \qquad (2)$$

$$\nabla_\theta^{\text{TD}} \mathbb{E}_{q_\phi(z)}[f_{\phi,\theta}(z)] = \mathbb{E}_{q_\phi(z)}\left[\nabla_\theta f_{\phi,\theta}(z)\right], \qquad (3)$$

and both can give rise to score functions (note $\nabla_\phi f_{\phi,\theta}(z) = \nabla_{\log q_\phi(z)} f_{\phi,\theta}(z)\nabla_\phi \log q_\phi(z)$).

In the following, we recapitulate how to address the score functions in Eq. (2) using the reparameterization trick and doubly-reparameterized gradients (DREGs, Tucker et al. (2019)), respectively. In Sec. 3 we introduce GDREGs, a generalization of DREGs, that allows us to eliminate the score function in Eq. (3).

**Reparameterization.** We can use the *reparameterization trick* (Kingma and Welling, 2014; Rezende et al., 2014) to turn the score function, $\nabla_\phi \log q_\phi(z)$, inside the expectation in Eq. (2) into a pathwise derivative of the function $f_{\phi,\theta}(z)$ as follows: we express the latent variables $z \sim q_\phi(z)$ through a bijection of new random variables $\epsilon \sim q(\epsilon)$, which are independent of $\phi$, $z = \mathcal{T}_q(\epsilon;\phi) \Leftrightarrow \epsilon = \mathcal{T}_q^{-1}(z;\phi)$. This allows us to rewrite expectations w.r.t. $q_\phi(z)$ as $\mathbb{E}_{q_\phi(z)}[f_{\phi,\theta}(z)] = \mathbb{E}_{q(\epsilon)}[f_{\phi,\theta}(\mathcal{T}_q(\epsilon;\phi))]$, which moves the parameter dependence into the argument of $f_{\phi,\theta}(z)$ and gives rise to a pathwise gradient:

$$\nabla_\phi^{\text{TD}} \mathbb{E}_{q_\phi(z)}[f_{\phi,\theta}(z)] = \mathbb{E}_{q(\epsilon)}\left[\nabla_\phi f_{\phi,\theta}(z) + \nabla_z f_{\phi,\theta}(z)\nabla_\phi \mathcal{T}_q(\epsilon;\phi)\right]_{z=\mathcal{T}_q(\epsilon;\phi)}. \qquad (4)$$

2

**Double reparameterization.** Tucker et al. (2019) reduce gradient variance by replacing the remaining score function in Eq. (4) with its reparameterized counterpart. Double reparameterization is based on the identity (Eq. (5) in Tucker et al. (2019))

$$\mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z})} \left[ g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z}) \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\boldsymbol{z}) \right] = \mathbb{E}_{\boldsymbol{\epsilon} \sim q(\boldsymbol{\epsilon})} \left[ \nabla_{\boldsymbol{z}}^{\text{TD}} g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z}) \big|_{\boldsymbol{z}=\mathcal{T}_q(\boldsymbol{\epsilon};\boldsymbol{\phi})} \nabla_{\boldsymbol{\phi}} \mathcal{T}_q(\boldsymbol{\epsilon};\boldsymbol{\phi}) \right] \qquad (5)$$

which follows from the fact that both the score function and the reparameterization estimators are unbiased and thus equal in expectation. For the IWAE objective Eq. (1), Tucker et al. (2019) derived the following doubly-reparametererized gradients (DREGs) estimator:

$$\widehat{\nabla}_{\boldsymbol{\phi}}^{\text{DREGs}} \mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\text{IWAE}} = \sum_{k=1}^{K} \widetilde{w}_k^2 \nabla_{\boldsymbol{z}_k}^{\text{TD}} \log w_k \nabla_{\boldsymbol{\phi}} \mathcal{T}_q(\boldsymbol{\epsilon}_k;\boldsymbol{\phi}); \quad \boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon}); \quad \widetilde{w}_k = \frac{w_k}{\sum_{k'=1}^{K} w_{k'}}. \qquad (6)$$

While the DREGs estimator reparameterizes the score function in Eq. (4), the previously proposed STL estimator (Roeder et al., 2017) simply drops it and is usually biased as a result. Crucially, because DREGs relies on reparameterization, it is limited to score functions of the sampling distribution $q_{\boldsymbol{\phi}}(\boldsymbol{z})$, making it inapplicable in the more general setting of arbitrary score functions (Eq. (3)). In the following Sec. 3 we introduce the GDREGs estimator that can be applied to these more general score function terms.

In App. C we discuss that the seemingly pathwise gradient in Eq. (4) can actually contain score functions for hierarchical models and explain how to extend DREGs to this case. Additional score functions arise because the distribution parameters (e.g. the mean and covariance) of one stochastic layer depend on the latent variables of previous layers.

## 3. Generalized DREGs

Here, we generalize DREGs to score function terms that involve distributions $p_{\boldsymbol{\theta}}(\boldsymbol{z})$ different from the sampling distribution $q_{\boldsymbol{\phi}}(\boldsymbol{z})$, such as $\mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z})}[g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z}) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{z})]$. Such terms appear, for example, when training a VAE with a trainable prior $p_{\boldsymbol{\theta}}(\boldsymbol{z})$ with the ELBO or IWAE objectives. We cannot use DREGs directly here, as it relies on reparameterization of the sampling distribution $q_{\boldsymbol{\phi}}(\boldsymbol{z})$, which means the path would then depend on the parameters $\boldsymbol{\phi}$ and double-reparameterization would only apply to its parameters $\boldsymbol{\phi}$, whereas the score function is with w.r.t. parameters $\boldsymbol{\theta}$ of a different distribution $p_{\boldsymbol{\theta}}(\boldsymbol{z})$.

To make progress *we need to make the path depend on the parameters $\boldsymbol{\theta}$* while still sampling $\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z})$ during training. Our solution is to: ⓐ *temporarily* change the path such that it depends on $\boldsymbol{\theta}$; ⓑ perform the reparameterized gradient computation; ⓒ change



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

the path back so we can use samples $\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z})$ to estimate the expectation. We change the path by first using an importance sampling reweighting to temporarily re-write the expectation, $\mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z})}[*] = \mathbb{E}_{p_{\boldsymbol{\theta}}(\boldsymbol{z})}\left[\frac{q_{\boldsymbol{\phi}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{z})}*\right]$, and then applying reparameterization to the new sampling distribution $p_{\boldsymbol{\theta}}(\boldsymbol{z})$: $\boldsymbol{z} = \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}; \boldsymbol{\theta})$ with $\widetilde{\boldsymbol{\epsilon}} \sim q(\widetilde{\boldsymbol{\epsilon}})$. We derive the gradient identity in Eq. (7) for a general $g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})$; we refer to it as the generalized DREGs (GDREGs) identity.

GDREGs identity

$$\mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{\phi}}}[g(\boldsymbol{z})\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{z})] = \mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{\phi}}}\left[\left(g(\boldsymbol{z})\nabla_{\boldsymbol{z}}^{\mathrm{TD}} \log \frac{q_{\boldsymbol{\phi}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{z})} + \nabla_{\boldsymbol{z}}^{\mathrm{TD}} g(\boldsymbol{z})\right)\nabla_{\boldsymbol{\theta}} \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}; \boldsymbol{\theta})|_{\widetilde{\boldsymbol{\epsilon}} = \mathcal{T}_p^{-1}(\boldsymbol{z},\boldsymbol{\theta})}\right] \quad (7)$$

$$\nabla_{\boldsymbol{\theta}}^{\mathrm{TD}}\mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z})}\left[g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})\right] \overset{\text{\textcircled{a}}}{=} \nabla_{\boldsymbol{\theta}}^{\mathrm{TD}}\mathbb{E}_{p_{\boldsymbol{\theta}}(\boldsymbol{z})}\left[\frac{q_{\boldsymbol{\phi}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{z})}g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})\right] \overset{\text{\textcircled{a}}}{=} \nabla_{\boldsymbol{\theta}}^{\mathrm{TD}}\mathbb{E}_{q(\widetilde{\boldsymbol{\epsilon}})}\left[\frac{q_{\boldsymbol{\phi}}(\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta}))}{p_{\boldsymbol{\theta}}(\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta}))}g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta}))\right] \quad (8)$$

$$\overset{\text{\textcircled{b}}}{=} \mathbb{E}_{q(\widetilde{\boldsymbol{\epsilon}})}\left[\nabla_{\boldsymbol{z}}^{\mathrm{TD}}\left(\frac{q_{\boldsymbol{\phi}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{z})}g(\boldsymbol{z})\right)\nabla_{\boldsymbol{\theta}}\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta}) + \frac{q_{\boldsymbol{\phi}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{z})}\left(\nabla_{\boldsymbol{\theta}}g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z}) - g(\boldsymbol{z})\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{z})\right)\right]_{\boldsymbol{z} = \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta})} \quad (9)$$

$$(10)$$

$$\overset{\text{\textcircled{c}}}{=} \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z})}\left[\left(g(\boldsymbol{z})\nabla_{\boldsymbol{z}}^{\mathrm{TD}} \log \frac{q_{\boldsymbol{\phi}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{z})} + \nabla_{\boldsymbol{z}}^{\mathrm{TD}} g(\boldsymbol{z})\right)\nabla_{\boldsymbol{\theta}}\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta})|_{\widetilde{\boldsymbol{\epsilon}} = \mathcal{T}_p^{-1}(\boldsymbol{z};\boldsymbol{\theta})} + \nabla_{\boldsymbol{\theta}}g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z}) - g(\boldsymbol{z})\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{z})\right]$$

In the derivation we have used the identity $x\nabla_* \log x = \nabla_* x$ repeatedly. By noting that $\nabla_{\boldsymbol{\theta}}^{\mathrm{TD}}\mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z})}[g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})] = \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z})}[\nabla_{\boldsymbol{\theta}}g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})]$, we can cancel these terms on the left hand side of Eq. (8) and right hand side of Eq. (10). By moving $-\mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z})}[g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{z})]$ to the other side we obtain the desired result.

Similar to DREGs (Eq. (5)), GDREGs allows us to transform score function gradients into pathwise gradients. But, unlike DREGs, GDREGs applies to general score functions and contains a correction term that vanishes when $p_{\boldsymbol{\theta}}(\boldsymbol{z})$ and $q_{\boldsymbol{\phi}}(\boldsymbol{z})$ are identical ($\log \frac{q_{\boldsymbol{\phi}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{z})}$ in Eq. (7)). Note that the pathwise derivative $\nabla_{\boldsymbol{\theta}}\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta})$ in Eq. (7) looks like a reparameterization of $p_{\boldsymbol{\theta}}(\boldsymbol{z})$ in terms of a noise variable $\widetilde{\boldsymbol{\epsilon}} = \mathcal{T}_p^{-1}(\boldsymbol{z};\boldsymbol{\theta})$ with $\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z})$. We can interpret this sequence of transformations as a normalizing flow (Rezende and Mohamed, 2015) $\boldsymbol{z} \to \widetilde{\boldsymbol{\epsilon}} \to \boldsymbol{z}$, such that $\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta}) = \mathcal{T}_p(\mathcal{T}_p^{-1}(\boldsymbol{z};\boldsymbol{\theta});\boldsymbol{\theta}) = \boldsymbol{z}$. We can think of this procedure as *re-expressing the sample $\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z})$ as if it came from $p_{\boldsymbol{\theta}}(\boldsymbol{z})$*: Its numerical value $\boldsymbol{z}$ remains unchanged and it is still distributed according to $q_{\boldsymbol{\phi}}(\boldsymbol{z})$, yet its pathwise gradient $\nabla_{\boldsymbol{\theta}}\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}};\boldsymbol{\theta})$ depends on $\boldsymbol{\theta}$. We illustrate the corresponding computational flow in Fig. 1. Note that to derive the GDREGs identity, we only require $p_{\boldsymbol{\theta}}(\boldsymbol{z})$ to be reparameterizable ( ☐ in Fig. 1). While $q_{\boldsymbol{\phi}}(\boldsymbol{z})$ may be reparameterizable as well ( ⌐ ⌐ in Fig. 1), this is not necessary; we only need to be able to evaluate its density in Eq. (7).

We use the GDREGs identity to address general score functions of the form Eq. (3) and derive the GDREGs estimator for the IWAE objective w.r.t. the prior parameters $\boldsymbol{\theta}$:

$$\widehat{\nabla}_{\boldsymbol{\theta}}^{\mathrm{GDREGs}}\mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\mathrm{IWAE}} = \sum_{k=1}^{K}\left(\widetilde{w}_k\nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_k) - \widetilde{w}_k^2\nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log w_k\right)\nabla_{\boldsymbol{\theta}}\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k;\boldsymbol{\theta})\Big|_{\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}(\boldsymbol{z}_k,\boldsymbol{\theta})} \quad (11)$$

with $\boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})$; full derivation is in App. D. The second term in Eq. (11) looks like the DREGs estimator for $\boldsymbol{\phi}$ in Eq. (6) except that the samples $\boldsymbol{z}_k$ are now re-expressed as if they came from $p_{\boldsymbol{\theta}}(\boldsymbol{z})$. In addition we obtain a term that involves the likelihood $p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z})$ and is linear in $\widetilde{w}_k$. We do not apply GDREGs to the likelihood parameters $\boldsymbol{\lambda}$ because $p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z})$ is a distribution over $\boldsymbol{x}$ rather than $\boldsymbol{z}$; in the following we therefore drop the subscript $\boldsymbol{\lambda}$.

Similarly to the DREGs estimator, we can also extend the GDREGs estimator to hierarchical objectives as we explain in App. C.
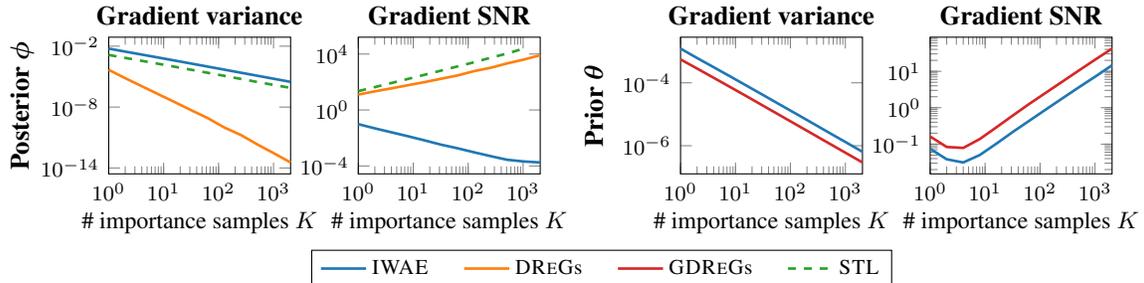
**Figure 2:** Average gradient variance and gradient signal-to-noise ratio (SNR) for the variational posterior parameters $\phi$ and the prior parameters $\theta$.

We learn all parameters by optimizing the same objective Eq. (1), 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. F for details. While separate objectives might seem computationally expensive, most terms are shared between them, and modern frameworks would avoid such duplicate computation. In practise, we found the runtime increase for training with DREGs and GDREGs estimators to be less than 10% without any optimization of the implementation.

## 4. Experiments

Here, we empirically evaluate the proposed GDREGs estimator as well as the proposed hierarchical extensions of DREGs and GDREGs, and compare them to the naive IWAE gradient estimator (labelled as IWAE) as well as STL (Roeder et al., 2017).

**Illustrative example.** We first consider an extended version of the illustrative example introduced 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 and variational posterior and find that (see Fig. 2): (i) for the parameters $\phi$ of the variational posterior, our extended version of the DREGs estimator also resolves the vanishing SNR problem of the naive IWAE estimator (Rainforth et al., 2018) in the hierarchical case by reducing the gradient variance at a faster rate with the number of importance samples; (ii) for the paramters $\theta$ of the prior, the GDREGs estimator has smaller gradient variance and better SNR than the naive IWAE estimator but scales at the same rate with the number of importance samples.

**Image modelling with VAEs.** We also consider conditional and unconditional image modelling tasks with single layer and hierarchical (multi-layer) VAEs on several standard benchmark datasets: MNIST (LeCun and Cortes, 2010), Omniglot (Lake et al., 2015), and FashionMNIST (Xiao et al., 2017). We use the dynamically binarized versions of the datasets to minimize overfitting. In the hierarchical case, the generative path (prior and likelihood) is top-down whereas the variational posterior is bottom-up; for conditional modelling we predict the bottom half of an image given its top half, as in Tucker et al. (2019); in this
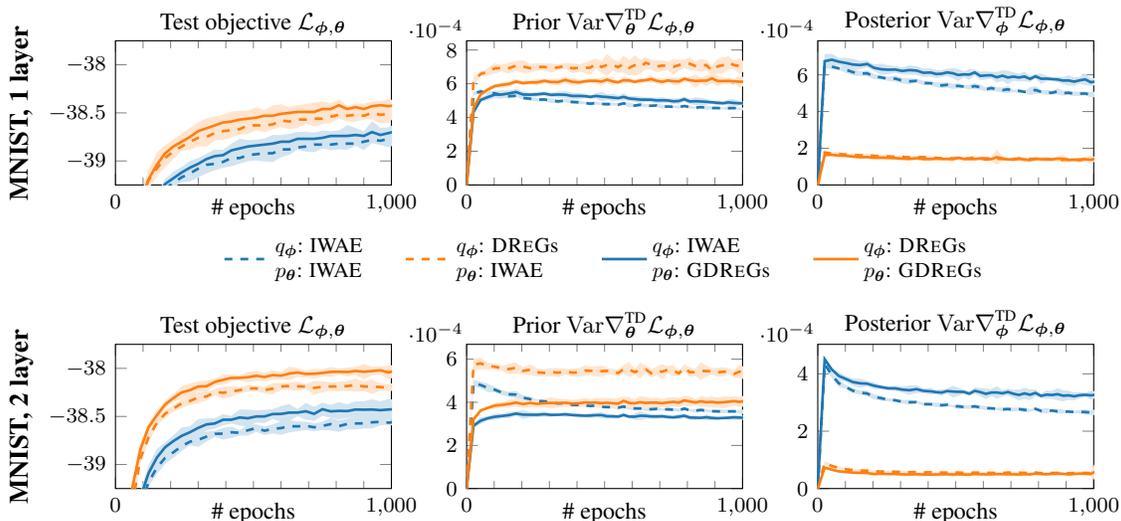
**Figure 3:** *Conditional* image modelling of MNIST with a VAE with 1 layer (top) and 2 layers (bottom). Shaded areas denote $\pm$ 1.96 standard deviations $\sigma$ over 5 reruns.

case, both the prior and variational posterior additionally depend on a context variable $\boldsymbol{c}$ ($q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}, \boldsymbol{c})$ and $p_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{c})$, respectively).

First, we evaluate the choice of estimator for the parameters of $q_{\boldsymbol{\phi}}(\boldsymbol{z})$. 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 $\boldsymbol{\phi}$ on all tasks (third column in Fig. 3), 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.

Second, we consider the estimators for the $\boldsymbol{\theta}$ parameters of the prior $p_{\boldsymbol{\theta}}(\boldsymbol{z})$. Using the GDReGs estimator instead of the naive IWAE estimator consistently improves the train and test performance when combined with *any* estimator for the variational posterior, especially for conditional image modelling with deeper models. For unconditional image modelling the improvements are marginal, though using GDReGs never hurts. In terms of gradient variance for the prior parameters $\boldsymbol{\theta}$, GDReGs consistently performs better in the beginning of training, when it always has lower variance. Later in training this is only consistently true when also using the DReGs estimator for the variational posterior parameters $\boldsymbol{\phi}$.

## 5. 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 to 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 lead to better train and test performance on conditional and unconditional image modelling tasks.

## Acknowledgments

## References

Yuri Burda, Roger B. Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. In *Proceedings of the 4th International Conference on Learning Representations*, 2016.

Tomas Geffner and Justin Domke. Approximation based variance reduction for reparameterization gradients. *arXiv preprint arXiv:2007.14634*, 2020.

Diederik Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations*, 2014.

Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 2015.

Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL http://yann.lecun.com/exdb/mnist/.

Andrew Miller, Nick Foti, Alexander D'Amour, and Ryan P Adams. Reducing reparameterization gradient variance. In *Advances in Neural Information Processing Systems*, 2017.

Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte Carlo gradient estimation in machine learning. *Journal of Machine Learning Research*, 2020.

Tom Rainforth, Adam Kosiorek, Tuan Anh Le, Chris Maddison, Maximilian Igl, Frank Wood, and Yee Whye Teh. Tighter variational bounds are not necessarily better. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.

Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning*, 2015.

Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.

Geoffrey Roeder, Yuhuai Wu, and David Duvenaud. Sticking the landing: Simple, lower-variance gradient estimators for variational inference. In *Advances in Neural Information Processing Systems 30*, 2017.

Francisco J. R. Ruiz, Michalis K. Titsias, and David M. Blei. Overdispersed black-box variational inference. In *Uncertainty in Artificial Intelligence*, 2016.

George Tucker, Dieterich Lawson, Shixiang Gu, and Chris J. Maddison. Doubly reparameterized gradient estimators for monte carlo objectives. In *Proceedings of the 7th International Conference on Learning Representations*, 2019.

Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 1992.

Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms, 2017.

## Appendix A. Additional results

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

**Table A.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.

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

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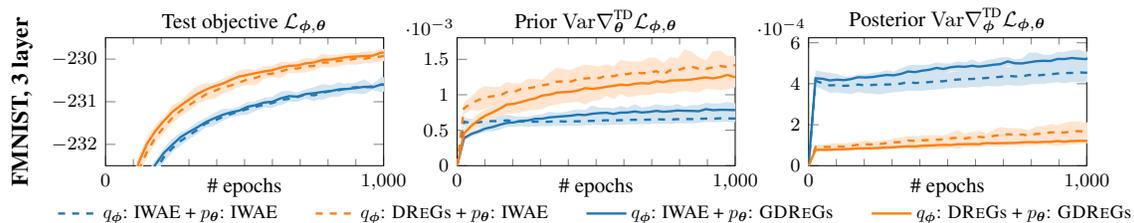


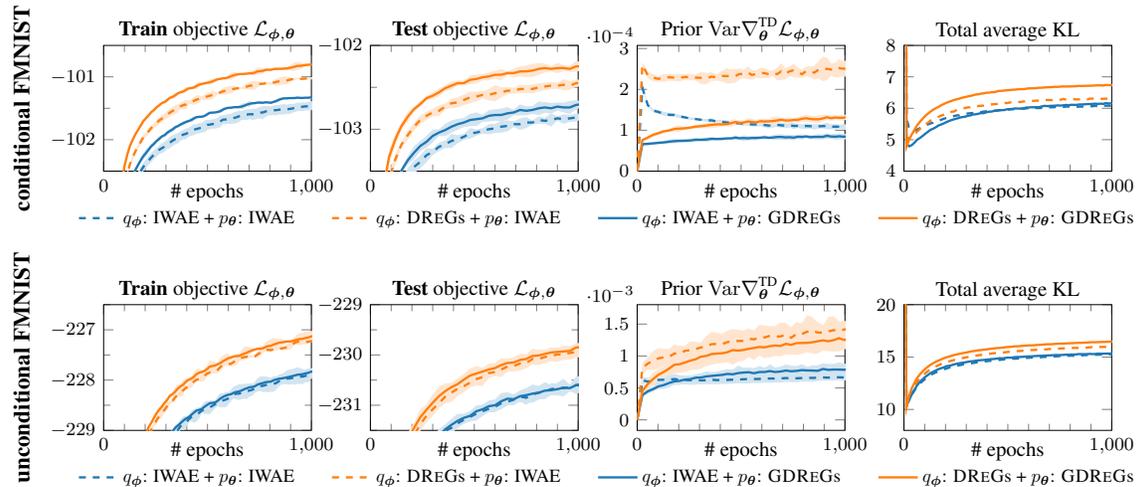**Figure A.1:** *Unconditional* image modelling on FashionMNIST; 3 layers.

**Figure A.2: Train** objective *(leftmost column)* and total average KL *(rightmost column)* in addition to the test objective and prior gradient variance for conditional and unconditional models with 3 stochastic layers on FashionMNIST. We find that the KL is lower for conditional models than unconditional models, which indicates that the variational posterior and the prior are closer to each other in this case.

## Appendix B. 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 obtained by dropping this score function to reduce the estimator variance. They applied the estimator to hierarchical ELBO models but did 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 bound. 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.

A number of classic techniques from the variance reduction literature have also been applied to variational inference and reparameterization. For example, Miller et al. (2017) and Geffner and Domke (2020) proposed control variates for reparameterization gradients; while 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.

## Appendix C. DREGs and GDREGs for hierarchical models

We now consider models with hierarchically structured latent variables and show that even terms that look like pathwise gradients, such as $\nabla_{\boldsymbol{z}}^{\mathrm{TD}} f_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})$ in Eq. (4), can give rise to score

functions in this case. *The score functions appear because the distribution parameters of one stochastic layer depend on the latent variables of another layer.*

To illustrate this, consider a hierarchical model with two layers where we first sample $z_2 \sim q_{\phi_2}(z_2)$ and then $z_1 \sim q_{\phi_1}(z_1|z_2)$.[1] 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 the subsequent level (we omit the true pathwise gradients $(\dots)$):

$$\nabla_{z_2}^{\mathrm{TD}} \log q_{\phi_1}(z_1|z_2) = \nabla_{z_2}^{\mathrm{TD}} \log q_{\alpha_{1|2}(z_2,\phi_1)}(z_1) \tag{C.1}$$

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

Similar additional score functions arise for seemingly pathwise gradients of hierarchical and/or autoregressive priors and variational posteriors.

### C.1. 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. (1)), 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 $\nabla_{\phi_2}^{\mathrm{TD}} \log q_{\phi_1,\phi_2}(z_1, z_2)$ 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 $z_1(\phi_1, \phi_2) \equiv \mathcal{T}_{q_1}(\epsilon_1; \alpha_{1|2}(z_2(\phi_2), \phi_1))$:

$$\nabla_{\phi_2}^{\mathrm{TD}} \log q_{\alpha_2(\phi_2)}(z_2(\phi_2)) \, q_{\alpha_{1|2}(z_2(\phi_2),\phi_1)}(z_1(\phi_1,\phi_2)) \tag{C.3}$$

The total derivative w.r.t. parameters of the upper layer, $\phi_2$, gives rise to three types of gradients: (true) pathwise gradients w.r.t. $z_1$ and $z_2$, a *direct* score function because the distribution parameters $\alpha_2(\phi_2)$ directly depend on $\phi_2$, and an *indirect* score function because $\alpha_{1|2}(z_2(\phi_2), \phi_1)$ indirectly depends on $\phi_2$ through $z_2$. Other terms in the gradient of the objective w.r.t. $\phi$ as well as gradients w.r.t. the $\theta$ parameters decompose similarly. We have three options to estimate each score function individually: (1) leave it—this naive estimator is unbiased but potentially has high variance; (2) drop it, similar to STL—this estimator is generally biased; (3) doubly-reparameterize it using DREGs—this estimator is unbiased, but can generate further score function terms.

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

---

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

variance for DREGs. Moreover, for most model structures it is impossible to consistently replace all successively arising score functions by doubly-reparameterized gradients. Thus, to extend DREGs to hierarchical models, we leave the indirect score functions unchanged and only doubly reparameterize the direct score functions. We provide detailed derivations and a general DREGs estimator for arbitrary hierarchical structures in App. E, and show how to implement the corresponding surrogate loss functions in App. F.

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 dropping the direct score functions while maintaining the indirect ones, similar to how we extend DREGs to hierarchical models; the STL estimator is biased for IWAE objectives (Tucker et al., 2019).
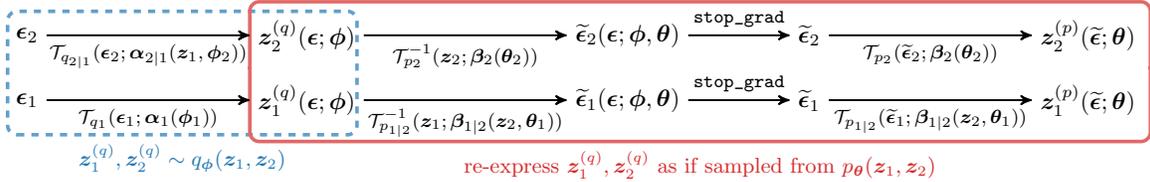
## C.2. Extending GDRᴇGs to hierarchical VAEs



**Figure C.1:** 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 $\widetilde{\epsilon}_i$ follows a different, usually more complex, distribution from $\epsilon_i$. $\alpha_i$ and $\beta_i$ denote the distribution parameters of the variatonal posterior and the prior, respectively.

When extending GDRᴇGs to hierarchical models, we again encounter direct and indirect score functions (see App. C), and we apply GDRᴇGs to the direct score functions but leave the indirect score functions. See App. E for derivations of a general GDRᴇGs estimator for arbitrary hierarchical structures and App. F for surrogate losses to implement them.

To apply GDRᴇGs we need to re-express samples from $q_\phi(z)$ as if they came from $p_\theta(z)$. We do this for the entire hierarchy jointly. In Fig. C.1 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. E 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 $(\widetilde{\epsilon}_1, \widetilde{\epsilon}_2)$ that follow a more complicated distribution than $(\epsilon_1, \epsilon_2)$. As in the single-layer case, only $p_\theta(z)$ needs to be reparameterizable.

## Appendix D. Derivation of the GDReGs estimator for the IWAE objective

In this section we apply the GDReGs identity derived above to derive the GDReGs estimator for the IWAE objective, Eq. (11) in the main paper.

### D.1. Preliminaries on the IWAE objective

The importance weighted autoencoder (IWAE) objective is given by

$$\mathcal{L}_{\phi,\theta}^{\text{IWAE}} = \mathbb{E}_{\boldsymbol{z}_{1:K} \sim q_{\phi}(\boldsymbol{z}_k|\boldsymbol{x})} \left[ \log \left( \frac{1}{K} \sum_{k=1}^{K} w_k \right) \right] \qquad w_k = \frac{p_{\theta}(\boldsymbol{z})p(\boldsymbol{x}|\boldsymbol{z}_k)}{q_{\phi}(\boldsymbol{z}_k|\boldsymbol{x})} \qquad (\text{D.1})$$

where $w_k$ are the importance weights (Burda et al., 2016).

Due to the structure of the IWAE objective, any gradient w.r.t. any of its parameters can be written as

$$\nabla_*^{\text{TD}} \mathcal{L}_{\phi,\theta}^{\text{IWAE}} = \mathbb{E}_{\boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon})} \left[ \sum_{k=1}^{K} \widetilde{w}_k \nabla_*^{\text{TD}} \log w_k \right]; \qquad \widetilde{w}_k = \frac{w_k}{\sum_j w_j} \qquad (\text{D.2})$$

using the chain rule and $\nabla_* w_k = w_k \nabla_* \log w_k$. $\widetilde{w}_k$ are the normalized importance weights, and we have reparameterized $\boldsymbol{z}_k$ as $\mathcal{T}_q(\boldsymbol{\epsilon}_k; \boldsymbol{\phi})$. Typically, the derivatives we are interested in are w.r.t. the parameters $\boldsymbol{\phi}$ and $\boldsymbol{\theta}$.

We also note the following identity that we use in the derivation of the doubly reparameterized estimators,

$$\nabla_{\boldsymbol{z}}^{\text{TD}} \widetilde{w}_k = \left( \widetilde{w}_k - \widetilde{w}_k^2 \right) \nabla_{\boldsymbol{z}}^{\text{TD}} \log w_k \qquad (\text{D.3})$$

which follows from applying the chain-rule and using $\nabla_* w_k = w_k \nabla_* \log w_k$.

Tucker et al. (2019) derive the DReGs identity (Eq. (5)) and use it to derive the following doubly-reparameterized gradient estimator (DReGs) w.r.t. the approximate posterior parameters $\boldsymbol{\phi}$ as:

$$\widehat{\nabla}_{\boldsymbol{\phi}}^{\text{DReGs}} \mathcal{L}_{\text{IWAE}} = \sum_{k=1}^{K} \widetilde{w}_k^2 \nabla_{\boldsymbol{z}_k}^{\text{TD}} \log w_k \nabla_{\boldsymbol{\phi}}^{\text{TD}} \mathcal{T}_q(\boldsymbol{\epsilon}_k; \boldsymbol{\phi}). \qquad \boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon}) \qquad (\text{D.4})$$

### D.2. Derivation of the GDReGs estimator

Similarly, we can derive a generalized doubly-reparameterized gradient (GDReGs) estimator w.r.t. the prior parameters $\boldsymbol{\theta}$. We use the GDReGs identity (Eq. (7)) derived above with $g_{\phi,\theta}(\boldsymbol{z}) = \widetilde{w}_k$ and note that the reweighting term $\log \frac{q_{\phi}(z)}{p_{\theta}(z)}$ looks like a log importance weight

except for the missing likelihood:

$$\nabla_{\boldsymbol{\theta}}^{\mathrm{TD}} \mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\mathrm{IWAE}} = \mathbb{E}_{\boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})} \left[ \sum_{k=1}^{K} \widetilde{w}_k \nabla_{\boldsymbol{\theta}}^{\mathrm{TD}} \log w_k \right] = \mathbb{E}_{\boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})} \left[ \sum_{k=1}^{K} \widetilde{w}_k \nabla_{\boldsymbol{\theta}}^{\mathrm{TD}} \log p_{\boldsymbol{\theta}}(\boldsymbol{z}) \right]$$

(D.5)

$$\stackrel{(7)}{=} \mathbb{E}_{\boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})} \left[ \sum_{k=1}^{K} \left( \widetilde{w}_k \nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log \frac{q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})}{p_{\boldsymbol{\theta}}(\boldsymbol{z}_k)} + \nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \widetilde{w}_k \right) \nabla_{\boldsymbol{\theta}} \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k; \boldsymbol{\theta})|_{\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}(\boldsymbol{z}_k;\boldsymbol{\theta})} \right]$$

(D.6)

$$\stackrel{(D.3)}{=} \mathbb{E}_{\boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})} \left[ \sum_{k=1}^{K} \left( \widetilde{w}_k \nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log \frac{q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})}{p_{\boldsymbol{\theta}}(\boldsymbol{z}_k)} + \left( \widetilde{w}_k - \widetilde{w}_k^2 \right) \nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log w_k \right) \nabla_{\boldsymbol{\theta}} \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k; \boldsymbol{\theta})|_{\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}(\boldsymbol{z}_k;\boldsymbol{\theta})} \right]$$

(D.7)

$$= \mathbb{E}_{\boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})} \left[ \sum_{k=1}^{K} \left( \widetilde{w}_k \nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log p(\boldsymbol{x}|\boldsymbol{z}_k) - \widetilde{w}_k^2 \nabla_{\boldsymbol{z}}^{\mathrm{TD}} \log w_k \right) \nabla_{\boldsymbol{\theta}} \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k; \boldsymbol{\theta})|_{\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}(\boldsymbol{z}_k;\boldsymbol{\theta})} \right].$$

(D.8)

Thus, the GDREGs estimator is given by:

$$\widehat{\nabla}_{\boldsymbol{\theta}}^{\mathrm{GDREGs}} \mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\mathrm{IWAE}} = \sum_{k=1}^{K} \left( \widetilde{w}_k \nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log p(\boldsymbol{x}|\boldsymbol{z}_k) - \widetilde{w}_k^2 \nabla_{\boldsymbol{z}}^{\mathrm{TD}} \log w_k \right) \nabla_{\boldsymbol{\theta}} \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k; \boldsymbol{\theta})|_{\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}(\boldsymbol{z}_k;\boldsymbol{\theta})} \quad \boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x}).$$

(11)

Note that the $\boldsymbol{z}_k$ are sampled from $q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})$ but re-rexpressed as if they came from $p_{\boldsymbol{\theta}}(\boldsymbol{z})$.

We can rewrite the importance weights as

$$w_k = \frac{p_{\boldsymbol{\theta}}(\boldsymbol{z}_k) p(\boldsymbol{x}|\boldsymbol{z}_k)}{q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})} = \frac{p_{\boldsymbol{\theta}}(\boldsymbol{z}_k|\boldsymbol{x}) p_{\boldsymbol{\theta}}(\boldsymbol{x})}{q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})}.$$

(D.9)

Thus, if the variational posterior $q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})$ is equal to the true posterior $p_{\boldsymbol{\theta}}(\boldsymbol{z}_k|\boldsymbol{x})$, all weights $w_k$ become equal to $p_{\boldsymbol{\theta}}(\boldsymbol{x})$ and thus constant w.r.t. $\boldsymbol{z}_k$. In that case the second term in the GDREGs estimator Eq. (11) vanishes and the overall expression simplifies to

$$\widehat{\nabla}_{\boldsymbol{\theta}}^{\mathrm{GDREGs}} \mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\mathrm{IWAE}} = \sum_{k=1}^{K} \widetilde{w}_k \nabla_{\boldsymbol{z}_k}^{\mathrm{TD}} \log p(\boldsymbol{x}|\boldsymbol{z}_k) \nabla_{\boldsymbol{\theta}} \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k; \boldsymbol{\theta})|_{\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}(\boldsymbol{z}_k;\boldsymbol{\theta})} \cdot \quad \boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})$$

(D.10)

In contrast, the usual IWAE gradient involves the score function for $p_{\boldsymbol{\theta}}(\boldsymbol{z}_k)$:

$$\widehat{\nabla}_{\boldsymbol{\theta}}^{\mathrm{naive}} \mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\mathrm{IWAE}} = \sum_{i=1}^{K} \widetilde{w}_k \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{z}_k), \qquad \boldsymbol{z}_{1:K} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x}).$$

(D.11)

14

## Appendix E. Derivation of the DREGs and GDREGs estimator for IWAE objectives of hierarchical VAEs

In this section we derivations of and further details on the extension of DREGs and GDREGs to hierarchical VAEs with the IWAE objective.

### E.1. Preliminaries and notation for the hierarchical IWAE objective

For a hierarchically structured model with $L$ stochastic layers the IWAE objective is still given by Eq. (D.1) but with importance weights $w_k$ given by

$$w_k = \frac{p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_{k1}, \ldots, \boldsymbol{z}_{kL})\, p_{\boldsymbol{\theta}}(\boldsymbol{z}_{k1}, \ldots, \boldsymbol{z}_{kL})}{q_{\boldsymbol{\phi}}(\boldsymbol{z}_{k1}, \ldots, \boldsymbol{z}_{kL}|\boldsymbol{x})}. \tag{E.1}$$

Here, $\boldsymbol{z}_{kl}$ denotes the $k$th importance sample ($k \in \{1, \ldots, K\}$) for the $l$th layer ($l \in \{1, \ldots, L\}$). Both the variational posterior and the prior distribution factorize according to their respective hierarchical structure. While the prior factorizes top-down in most cases, the variational posterior can have many different structures. In order for the distributions to be valid in the context of a VAE, we require the individual dependency graphs for the prior (generative path) and the variational posterior (inference path) to be directed acyclic graphs. Cycles would mean that a latent variable conditionally dependent on itself. To keep the dependency structure general, we write the factorization of the variational posterior and prior as follows:

$$q_{\boldsymbol{\phi}}(\boldsymbol{z}_{k1}, \ldots, \boldsymbol{z}_{kL}|\boldsymbol{x}) = \prod_{l=1}^{L} q_{\boldsymbol{\phi}_l}(\boldsymbol{z}_{kl}|\mathrm{pa}_{\boldsymbol{\alpha}}(l), \boldsymbol{x}) = \prod_{l=1}^{L} q_{\boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l); \boldsymbol{\phi}_l)}(\boldsymbol{z}_{kl}) \tag{E.2}$$

$$p_{\boldsymbol{\theta}}(\boldsymbol{z}_{k1}, \ldots, \boldsymbol{z}_{kL}) = \prod_{l=1}^{L} p_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{kl}|\mathrm{pa}_{\boldsymbol{\beta}}(l)) = \prod_{l=1}^{L} p_{\boldsymbol{\beta}_l(\mathrm{pa}_{\boldsymbol{\beta}}(l); \boldsymbol{\theta}_l)}(\boldsymbol{z}_{kl}) \tag{E.3}$$

Here, $\boldsymbol{\alpha}_l(\cdot; \boldsymbol{\phi}_l)$ and $\boldsymbol{\beta}_l(\cdot; \boldsymbol{\theta}_l)$ are the distribution parameters of the variational posterior and prior distribution in the $l$th layer, respectively, and we have made the dependencies of the conditional distributions explicit; $\mathrm{pa}_{\boldsymbol{\alpha}}(l)$ denotes the "parents" of the latent variable $\boldsymbol{z}_{kl}$ according to the dependency graph of the inference path (the factorization of the posterior); similarly, $\mathrm{pa}_{\boldsymbol{\beta}}(l)$ denotes the latent variables that $\boldsymbol{z}_{kl}$ directly depends on according to the factorization of the prior $p_{\boldsymbol{\theta}}$. Typically, the prior is assumed to factorize top-down, such that $\mathrm{pa}_{\boldsymbol{\beta}}(l) = \boldsymbol{z}_{k(l+1)}$ for all but the top-most layer.

The samples $\boldsymbol{z}_{kl}$ are drawn from the variational posterior and can be expressed through reparameterization as $\boldsymbol{z}_{kl} = \mathcal{T}_{q_l}(\boldsymbol{\epsilon}_{kl}; \boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l), \boldsymbol{\phi}_l))$, where $\boldsymbol{\epsilon}_{kl}$ is an independent noise variable per importance sample and layer.

We note that it is these dependencies of the distribution parameters $\boldsymbol{\alpha}_l$ and $\boldsymbol{\beta}_l$ on $\mathrm{pa}_{\boldsymbol{\alpha}}(l)$ and $\mathrm{pa}_{\boldsymbol{\beta}}(l)$, respectively, that give rise to the indirect score functions as discussed in App. C.

## E.2. Derivation of the hierarchical DReGs estimator for IWAE

With notation fully set up we consider the reparameterized gradients of the IWAE objective w.r.t. the variational parameters in a particular stochastic layer $\boldsymbol{\phi}_l$:

$$\nabla_{\boldsymbol{\phi}_l}^{\mathrm{TD}} \mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\mathrm{IWAE}} = \mathbb{E}_{\boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon})} \left[ \sum_{k=1}^{K} \widetilde{w}_k \nabla_{\boldsymbol{\phi}_l}^{\mathrm{TD}} \log w_k \right] \tag{E.4}$$

$$= \mathbb{E}_{\boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon})} \left[ \sum_{k=1}^{K} \widetilde{w}_k \left( \nabla_{\boldsymbol{z}_{kl}}^{\mathrm{TD}} \log w_k \nabla_{\boldsymbol{\phi}_l} \mathcal{T}_{q_l} \left( \boldsymbol{\epsilon}_{kl}; \boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l), \boldsymbol{\phi}_l) \right) + \nabla_{\boldsymbol{\phi}_l} \log w_k \right) \right] \tag{E.5}$$

where we have used the chain-rule to arrive at Eq. (E.5); the first term contains both the (true) pathwise gradients as well as the indirect score functions; the second term only contains a direct score function as we only take the partial derivative w.r.t. $\boldsymbol{\phi}_l$.

We can rewrite this direct score function gradient because only one term in the (log-)importance weight directly depends on $\boldsymbol{\phi}_l$,

$$\nabla_{\boldsymbol{\phi}_l} \log w_k = -\nabla_{\boldsymbol{\phi}_l} \log q_{\boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l); \boldsymbol{\phi}_l)}(\boldsymbol{z}_{kl}). \tag{E.6}$$

Applying the DReGs identity to this term and using Eq. (D.3) yields:

$$\mathbb{E}_{\boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon})} \left[ \sum_{k=1}^{K} \widetilde{w}_k \nabla_{\boldsymbol{\phi}_l} \log w_k \right] = -\mathbb{E}_{\boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon})} \left[ \sum_{k=1}^{K} (\widetilde{w}_k - \widetilde{w}_k^2) \nabla_{\boldsymbol{z}_{kl}}^{\mathrm{TD}} \log w_k \nabla_{\boldsymbol{\phi}_l} \mathcal{T}_{q_l} \left( \boldsymbol{\epsilon}_{kl}; \boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l), \boldsymbol{\phi}_l) \right) \right] \tag{E.7}$$

which agrees with the first term in Eq. (E.5) up to the prefactor. Thus, both the true pathwise gradients as well as the indirect score functions appear twice and the prefactors partly cancel to give rise to the DReGs estimator for hierarchical IWAE objectives:

DReGs estimator for hierarchical IWAE objectives

$$\widehat{\nabla}_{\boldsymbol{\phi}_l}^{\mathrm{DReGs}} \mathcal{L}_{\boldsymbol{\phi},\boldsymbol{\theta}}^{\mathrm{IWAE}} = \sum_{k=1}^{K} \widetilde{w}_k^2 \nabla_{\boldsymbol{z}_{kl}}^{\mathrm{TD}} \log w_k \nabla_{\boldsymbol{\phi}_l} \mathcal{T}_{q_l} \left( \boldsymbol{\epsilon}_{kl}; \boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l), \boldsymbol{\phi}_l) \right); \quad \boldsymbol{\epsilon}_{1:K} \sim q(\boldsymbol{\epsilon}) \tag{E.8}$$

where $\boldsymbol{z}_{kl} = \mathcal{T}_{q_l}(\boldsymbol{\epsilon}_{kl}; \boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l), \boldsymbol{\phi}_l)), \forall l \in \{1, \ldots, L\}, \forall k \in \{1, \ldots, K\}$ through reparameterization.

We emphasize that the total derivative w.r.t. $\boldsymbol{z}_{kl}$ contains pathwise gradients as well as indirect score functions for both the variational posterior as well as for the prior. The hierarhical DReGs estimator otherwise looks very similar to the DReGs estimator in the single layer case (Tucker et al., 2019).

In App. F.1 we explain how to implement this estimator effectively and in a structure-agnostic way. That is, we do *not* have to derive a new estimator for each new dependency graph of the variational posterior or the prior.

## E.3. Derivation of the hierarchical GDReGs estimator for IWAE

Next, we derive the expression for the GDReGs estimator for hierarchical VAEs with IWAE objective.

Applying the GDREGs identity entails re-expressing the samples $z_{kl}$ from the variational posterior as if they were sampled from the prior. Starting form a sample $(z_{k1}, \ldots, z_{kL}) \sim q_\phi(z_1, \ldots, z_L | x)$, we use the inverse flow of $p_\theta$ to obtain new noise variables for each layer, $(\widetilde{\epsilon}_{k1}, \ldots, \widetilde{\epsilon}_{kL})$. We then use the forward flow of $p_\theta$ to obtain back $(z_{k1}, \ldots, z_{kL})$ but with the gradient path now depending on $\theta$ as discussed in Sec. 3.

More precisely, we find that

$$z_{kl}^{(q)} = \mathcal{T}_{q_l}\left(\epsilon_{kl}; \alpha_l(\mathrm{pa}_\alpha(l), \phi_l)\right) \qquad \text{original sampling of } (z_{k1}, \ldots, z_{kL}) \sim q_\phi(z_1, \ldots, z_L | x) \tag{E.9}$$

$$\widetilde{\epsilon}_{kl} = \mathcal{T}_{p_l}^{-1}\left(z_{kl}^{(q)}; \beta_l(\mathrm{pa}_\beta(l), \theta_l)\right) \quad \text{inverse prior flow to obtain new "noise" variables} \tag{E.10}$$

$$z_{kl} = \mathcal{T}_{p_l}\left(\widetilde{\epsilon}_{kl}; \beta_l(\mathrm{pa}_\beta(l), \theta_l)\right) \qquad \text{forward prior flow to re-express the } z_{kl} \tag{E.11}$$

where $\epsilon_{kl} \sim q(\epsilon)$ follows a simple distribution that is different from the more complicated distribution of $\widetilde{\epsilon}_{kl}$. Note how the initial reparameterization of a sample $z_{kl}$ depends on the dependency structure of the variational posterior (through $\mathrm{pa}_\alpha(\cdot)$), while the other transformations depend on the dependency structure of the prior ($\mathrm{pa}_\beta(\cdot)$).

As for DREGs, we note that only one term in the log importance weight directly depends on the variable $\theta_l$,

$$\nabla_{\theta_l} \log w_k = \nabla_{\theta_l} \log p_{\beta_l(\mathrm{pa}_\beta(l); \theta_l)}(z_{kl}). \tag{E.12}$$

With these prerequesits, we can compute the GDREGs estimator for parameters $\theta_l$ of the $l$th stochastic layer.

$$\nabla_{\theta_l}^{\mathrm{TD}} \mathcal{L}_{\phi,\theta}^{\mathrm{IWAE}} = \mathbb{E}_{z_{1:K} \sim q_\phi(z|x)}\left[\sum_{k=1}^{K} \widetilde{w}_k \nabla_{\theta_l} \log w_k\right] \tag{E.13}$$

$$\stackrel{E.12}{=} \mathbb{E}_{z_{1:K} \sim q_\phi(z|x)}\left[\sum_{k=1}^{K} \widetilde{w}_k \nabla_{\theta_l} \log p_{\beta_l(\mathrm{pa}_\beta(l); \theta_l)}(z_{kl})\right] \tag{E.14}$$

$$\stackrel{7}{=} \mathbb{E}_{z_{1:K} \sim q_\phi(z|x)}\left[\sum_{k=1}^{K}\left(\widetilde{w}_k \nabla_{z_{kl}}^{\mathrm{TD}} \log \frac{q_\phi(z_{k1}, \ldots, z_{kL}|x)}{p_\theta(z_{k1}, \ldots, z_{kL})} + \right.\right.$$
$$\left.\left. + \left(\widetilde{w}_k - \widetilde{w}_k^2\right) \nabla_{z_{kl}}^{\mathrm{TD}} \log w_k\right) \nabla_{\theta_l} \mathcal{T}_{p_l}\left(\widetilde{\epsilon}_{kl}; \theta_l\right)\big|_{\widetilde{\epsilon}_{kl} = \mathcal{T}_{p_l}^{-1}(z_{kl}; \theta_l)}\right] \tag{E.15}$$

---
**GDREGs estimator for hierarchical IWAE objectives**

$$\widehat{\nabla}_{\theta_l}^{\mathrm{GDREGs}} \mathcal{L}_{\phi,\theta}^{\mathrm{IWAE}} = \tag{E.16}$$
$$= \sum_{k=1}^{K}\left(\widetilde{w}_k \nabla_{z_{kl}}^{\mathrm{TD}} \log p_\lambda(x|z_{k1}, \ldots, z_{kL}) - \widetilde{w}_k^2 \nabla_{z_{kl}}^{\mathrm{TD}} \log w_k\right) \nabla_{\theta_l} \mathcal{T}_{p_l}\left(\widetilde{\epsilon}_{kl}; \theta_l\right)\big|_{\widetilde{\epsilon}_{kl} = \mathcal{T}_{p_l}^{-1}(z_{kl}; \theta_l)}$$

---

with $z_{1:K} \sim q_\phi(z_k|x)$, and where we suppressed dependencies on $\mathrm{pa}_*(l)$ where they are not necessary to simplify notation.

The estimator looks very similar to the GDREGs estimator for a single layer IWAE model Eq. (11). Note that just like above for hierarchical DREGs, the total gradients

w.r.t. $\boldsymbol{z}_{kl}$ give rise to both (true) pathwise gradients as well as indirect score functions through the hierarchical dependencies of the variational posterior and prior.

In App. F.2 we show how to implement the hierarchical GDRECs estimator Eq. (E.16) effectively and in a way that is agnostic to the structure of the model. That is, we do not have to derive a separate estimator for every dependency graph of the variational posterior and prior.

### E.4. Double reparameterization and indirect score functions

In principle, we could apply double-reparameterization to the indirect score functions as well. However, as we explain now, we often cannot doubly-reparameterize *all* indirect score functions; moreover, even in cases where this is possible, it is still impractical, as the corresponding estimator depends on the exact model structure and would require adaptation to each dependency graph of the prior and variational posterior.

Double reparameterization of indirect score functions works in the same way as for the direct score functions except that $g_{\boldsymbol{\phi},\boldsymbol{\theta}}(\boldsymbol{z})$ is given by $\widetilde{w}_k^2$ instead of $\widetilde{w}_k$ in this case. The derivatives of $\widetilde{w}_k^2$ have a similar reproducing property as we observed in Eq. (D.3):

$$\nabla_{\boldsymbol{z}}^{\mathrm{TD}}\widetilde{w}_k^2 = 2(\widetilde{w}_k^2 - \widetilde{w}_k^3)\nabla_{\boldsymbol{z}}^{\mathrm{TD}}\log w_k. \tag{E.17}$$

Thus, double reparameterization of the indirect score functions similarly gives rise to further indirect score functions. We note that these indirect score functions only appear for the "children" of the current stochastic layer, that is, stochastic variables in those layers that depend on the current layer. In this context, "children" refers to *all* children w.r.t. the dependency structure of both, the variational posterior *and* the prior. For a particular layer $l$ we obtain indirect score functions from double reparameterization of all of its (direct or indirect) parent nodes. Following the dependency structure, we could collect all of these terms and reparameterize them to obtain pathwise gradients only.

However, a problem arises, because we need to account for dependencies of both the variational posterior *and* the prior. Reparameterization of a score function gives rise to indirect score functions in all its "children" layers for both the variational posterior and the prior. For general hierarchical structures, this leads to cycles, in that some of the children of one dependency tree (the variational posterior) are the parents in the other (the prior) and/or vice versa. In this case we are never able to collect all the terms and fully reparameterize all the score functions.

Moreover, even if the joint dependency graph of the variational posterior and the prior were acyclic, this derivation would be structure-specific and would need to be repeated for each hierarchical structure. We therefore do not doubly reparameterize the indirect score functions.

## Appendix F. Surrogate losses to implement the DREGs and GDREGs estimators for IWAE objectives

As we discussed in Sec. 3 and similar to Tucker et al. (2019), we use surrogate loss functions to compute the gradients w.r.t. the likelihood, proposal, and prior parameters. That is, we use different losses, such that backpropagation results in the respective gradient estimator.

While Tucker et al. (2019) use a single surrogate loss to compute the gradient estimators for all parts of the objective, we choose to use separate surrogate losses for each of the three parameter groups (likelihood, variational posterior, prior). In principle, we could combine them into a single loss, but in order to keep presentation simple we keep them separate. Computationally this does not make a difference as modern deep learning frameworks avoid duplicate computation.

For the likelihood parameters, we use the regular (negative) IWAE objective Eq. (1) as a loss. That is, the gradient estimator for the likelihood parameters is given by the gradient of the negative IWAE objective.

To construct the other surrogate losses we need to stop the gradients at various points in the computation graph. In the following, we use the shorthand notation ⎵⎵⎵ to indicate that we stop gradients into the underlined parts of an expression. Where it might be ambiguous, or to highlight where we do *not* stop gradients, we use the shorthand ⎵⎵⎵ to indicate that gradients flow. For example, $f(\boldsymbol{\phi}, \boldsymbol{\theta})$ means that we backpropagate gradients into $\boldsymbol{\phi}$ but not into $\boldsymbol{\theta}$.

### F.1. DREGs for variational posterior parameters $\boldsymbol{\phi}$

#### F.1.1. Single stochastic layer

Here we reproduce part of the surrogate loss for the variational parameters $\boldsymbol{\phi}$ by Tucker et al. (2019) for the single stochastic layer case:

$$
L_{\mathrm{DREGs}}(\boldsymbol{\phi}) = \sum_{k=1}^{K} \widetilde{w}_k^2 \Big( \log p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_k) + \log p_{\boldsymbol{\beta}(\boldsymbol{\theta})}(\boldsymbol{z}_k) - \log q_{\boldsymbol{\alpha}(\boldsymbol{\phi})}(\boldsymbol{z}_k) \Big) \tag{F.1}
$$
$$
\boldsymbol{z}_k = \mathcal{T}_q(\boldsymbol{\epsilon}_k; \boldsymbol{\phi}) \qquad \boldsymbol{\epsilon}_k \sim q(\boldsymbol{\epsilon}_k)
$$

That is, we sample $\boldsymbol{z}_k \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}_k|\boldsymbol{x})$ as usual (by reparameterizing independent noise variables $\boldsymbol{\epsilon}_k$) but stop the gradients of the parameters that parameterize the distributions when evaluating their densities, $\log q_{\boldsymbol{\alpha}(\boldsymbol{\phi})}(\boldsymbol{z}_k)$. In addition we stop the gradients around the normalized importance weights $\widetilde{w}_k$. Differentiating $L_{\mathrm{DREGs}}$ w.r.t. the proposal parameters $\boldsymbol{\phi}$ yields the DREGs estimator Eq. (6). Note that we do not explicitly stop gradients into $\boldsymbol{\lambda}$ or $\boldsymbol{\theta}$ because we use separate surrogate losses for those parameter groups. If we were to use a combined loss, we would potentially have to stop gradients into these parameters as well, depending on the estimator used.

To practically implement this surrogate loss, we use two copies of the variational posterior distribution. An unaltered one (no stopped gradients) to sample $\boldsymbol{z}$ and one with gradients into the proposal parameters stopped to evaluate the log densities. The stopped gradient makes sure that we do not obtain a direct score function as we have doubly-reparameterized it.

Note that for single-stochastic-layer models we could also stop the gradients of the distribution parameters $\boldsymbol{\alpha}$ instead as they only depend on $\boldsymbol{\phi}$. We emphasize that this is not possible for hierarchical models as this would eliminate the indirect score functions and thus produce potentially biased gradients.

### F.1.2. Multiple stochastic layers

For multiple layers, the surrogate loss for the DReGs estimator Eq. (E.8) is given by:

$$L_{\mathrm{DReGs}}(\boldsymbol{\phi}) = \sum_{k=1}^{K} \widetilde{w}_k^2 \log w_k$$

$$\log w_k = \log p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_{k1},\ldots,\boldsymbol{z}_{kL}) + \sum_{l=1}^{L} \log p_{\boldsymbol{\beta}_l(\mathrm{pa}_{\boldsymbol{\beta}}(l);\boldsymbol{\theta}_l)}(\boldsymbol{z}_{kl}) - \sum_{l=1}^{L} \log q_{\boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l);\boldsymbol{\phi}_l)}(\boldsymbol{z}_{kl})$$

$$\boldsymbol{z}_{kl} = \mathcal{T}_{q_l}\left(\boldsymbol{\epsilon}_{kl};\boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l),\boldsymbol{\phi}_l)\right) \qquad \boldsymbol{\epsilon}_{kl} \sim q(\boldsymbol{\epsilon}_{kl})$$

$$(\mathrm{F.2})$$

Again, we do not explicitly stop gradients into $\boldsymbol{\lambda}$ or $\boldsymbol{\theta}_l$ as we only take gradients w.r.t. $\boldsymbol{\phi}_l$.

The indirect score functions arise due to the indirect dependence of the distribution parameters $\boldsymbol{\alpha}_l(\mathrm{pa}_{\boldsymbol{\alpha}}(l);\boldsymbol{\phi}_l)$ and $\boldsymbol{\beta}_l(\mathrm{pa}_{\boldsymbol{\beta}}(l);\boldsymbol{\theta}_l)$ on the parent latent variables $\mathrm{pa}_{\boldsymbol{\alpha}}(l)$ and $\mathrm{pa}_{\boldsymbol{\beta}}(l)$, respectively. Note how the former depends on the hierarchical structure of the variational posterior, whereas the latter depends on the hierarchical structure of the prior.

To implement this surrogate loss effectively, we again use two copies of the variational posterior distribution. One un-altered one (without stopped gradiends) from which we sample the individual reparameterized $\boldsymbol{z}_{kl}$ and through which gradients can flow; we use these samples to evaluate densities at and to parameterize the distribution parameters at subsequent layers. Derivatives w.r.t. $\boldsymbol{\phi}_l$ will then give rise to pathwise gradients and indirect score functions. We use the second copy of the variational posterior, where we have stopped the parameters $\boldsymbol{\phi}_l$, to evaluate the density at for the log importance weights in the last summand of Eq. (F.2).

## F.2. GDReGs for prior parameter $\boldsymbol{\theta}$

### F.2.1. Single stochastic layer

$$L_{\mathrm{GDReGs}}(\boldsymbol{\theta}) = \sum_{k=1}^{K} \widetilde{w}_k \log p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_k) - \widetilde{w}_k^2 \Big( \log p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_k) + \log p_{\boldsymbol{\beta}(\boldsymbol{\theta})}(\boldsymbol{z}_k) - \log q_{\boldsymbol{\alpha}(\boldsymbol{\phi})}(\boldsymbol{z}_k) \Big)$$

$$\boldsymbol{z}_k = \mathcal{T}_p\left(\widetilde{\boldsymbol{\epsilon}}_k;\boldsymbol{\theta}\right)$$

$$\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}\left(\mathcal{T}_q(\boldsymbol{\epsilon}_k;\boldsymbol{\phi});\boldsymbol{\theta}\right) \qquad \boldsymbol{\epsilon}_k \sim q(\boldsymbol{\epsilon}_k)$$

$$(\mathrm{F.3})$$

Taking the derivative of Eq. (F.3) w.r.t. $\boldsymbol{\theta}$ gives rise to the GDReGs estimator for the single stochastic layer IWAE objective. As explained in Sec. 3, we need to re-express $\boldsymbol{z}_k$ such that its path depends on $\boldsymbol{\theta}$. In effect, we first sample $\boldsymbol{z}_k = \mathcal{T}_q(\boldsymbol{\epsilon}_k;\boldsymbol{\phi})$, then compute the new noise variable $\widetilde{\boldsymbol{\epsilon}}_k = \mathcal{T}_p^{-1}(\boldsymbol{z};\boldsymbol{\theta})$, and re-compute $\boldsymbol{z}_k = \mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k;\boldsymbol{\theta})$. Note that we have to stop gradients into the noise variables $\widetilde{\boldsymbol{\epsilon}}_k$ to obtain the correct gradient estimator. This explains the `stop_grad` in Fig. 1.

As above, we do not explicitly stop gradients into $\boldsymbol{\lambda}$ and $\boldsymbol{\phi}$ as we use separate losses for these parameter groups and only compute gradients of Eq. (F.3) w.r.t. $\boldsymbol{\theta}$.

To effectively implement this loss, we use two copies of the prior distribution. One that we implement as a normalizing flow and a second one with stopped gradients into the parameters. We then proceed as follows:

- Compute the new noise variables $\widetilde{\boldsymbol{\epsilon}}_k$ by using the inverse flow $\mathcal{T}_p^{-1}$ on the samples $\boldsymbol{z}_k$ from the variational posterior.

- Stop the gradients into $\widetilde{\boldsymbol{\epsilon}}_k$.

- Use the forward flow $\mathcal{T}_p(\widetilde{\boldsymbol{\epsilon}}_k; \boldsymbol{\theta})$ to re-compute $\boldsymbol{z}_k$ but with path dependent on $\boldsymbol{\theta}$. These samples when derived w.r.t. $\boldsymbol{\theta}$ will give rise to the pathwise gradients.

- Use the second copy of the prior (with stopped gradients into its parameters) to evaluate the log density at the samples $\boldsymbol{z}_k$. The stopped gradients make sure that we do not obtain the direct score function.

### F.2.2. Multiple stochastic layers

For multiple stochastic layers the surrogate loss that gives rise to the GDREGs estimator Eq. (E.16) is given by:

$$L_{\text{GDREGs}}(\boldsymbol{\theta}) = \sum_{k=1}^{K} \widetilde{w}_k \log p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_{k1}, \ldots, \boldsymbol{z}_{kL}) - \widetilde{w}_k^2 \log w_k$$

$$\log w_k = \log p_{\boldsymbol{\lambda}}(\boldsymbol{x}|\boldsymbol{z}_{k1}, \ldots, \boldsymbol{z}_{kL}) + \sum_{l=1}^{L} \log p_{\boldsymbol{\beta}_l(\text{pa}_{\boldsymbol{\beta}}(l);\boldsymbol{\theta}_l)}(\boldsymbol{z}_{kl}) - \sum_{l=1}^{L} \log q_{\boldsymbol{\alpha}_l(\text{pa}_{\boldsymbol{\alpha}}(l);\boldsymbol{\phi}_l)}(\boldsymbol{z}_{kl})$$

$$\boldsymbol{z}_{kl} = \mathcal{T}_{p_l}\left(\widetilde{\boldsymbol{\epsilon}}_{kl}; \boldsymbol{\beta}_l(\text{pa}_{\boldsymbol{\beta}}(l), \boldsymbol{\theta}_l)\right)$$

$$\widetilde{\boldsymbol{\epsilon}}_{kl} = \mathcal{T}_{p_l}^{-1}\left(\boldsymbol{z}_{kl}^{(q)}; \boldsymbol{\beta}_l(\text{pa}_{\boldsymbol{\beta}}(l), \boldsymbol{\theta}_l)\right)$$

$$\boldsymbol{z}_{kl}^{(q)} = \mathcal{T}_{q_l}\left(\boldsymbol{\epsilon}_{kl}; \boldsymbol{\alpha}_l(\text{pa}_{\boldsymbol{\alpha}}(l), \boldsymbol{\phi}_l)\right) \qquad \boldsymbol{\epsilon}_{kl} \sim q(\boldsymbol{\epsilon}_{kl})$$

$$\tag{F.4}$$

As for the single layer case, we need to re-express variational posterior samples $\boldsymbol{z}_{kl}$ as if they were sampled from the prior. To obtain the correct gradients, we again have to stop gradients into the new noise variables $\widetilde{\boldsymbol{\epsilon}}_{kl}$, also see Fig. C.1.

As for hierarchical DREGs, the indirect score functions stem from the second and third term of $\log w_k$ and arise because the distribution parameters $\boldsymbol{\alpha}_l$ and $\boldsymbol{\beta}_l$ depend on the "parent" stochastic layers.

As before we use two copies of the prior distribution, one with regular gradients that is set up as a flow, and a second with stopped gradients into the parameters. This allows us to implement the GDREGs estimator regardless of the model structure.