

1 We thank the reviewers for their constructive feedback and address the common concerns across the four reviews.

2 **Difference with REINFORCE** We agree we didn't clearly explain the difference between our estimator and REIN-
 3 FORCE. We have made the following changes to the manuscript.

4 **1)** Rather than stating "the gradient estimator we derive ... does not require the high-variance REINFORCE..." in the
 5 introduction and throughout, we now say we derive a "score function estimator" to emphasize the fact that our estimator
 6 belongs in the family of estimators that use the log-derivative trick.

7 **2)** We have added a new appendix section clarifying the relationship as follows. Assuming $\pi_\beta(\mathbf{z}) = \hat{\pi}_\beta(\mathbf{z})/Z_\beta$ depends
 8 on parameters λ , to compute $\nabla_\lambda \mathbb{E}_{\pi_\beta(\mathbf{z})} [f(\mathbf{z})]$ one can use the:

$$\begin{aligned} \text{REINFORCE ESTIMATOR: } & \mathbb{E}_{\pi_\beta(\mathbf{z})} [f(\mathbf{z}) \nabla_\lambda \log \pi_\beta(\mathbf{z})] \\ \text{REINFORCE ESTIMATOR + BASELINE: } & \mathbb{E}_{\pi_\beta(\mathbf{z})} [(f(\mathbf{z}) - \mathbb{E}_{\pi_\beta(\mathbf{z})} [f(\mathbf{z})]) \nabla_\lambda \log \pi_\beta(\mathbf{z})] \\ \text{COVARIANCE ESTIMATOR (ours): } & \mathbb{E}_{\pi_\beta(\mathbf{z})} [(f(\mathbf{z}) - \mathbb{E}_{\pi_\beta(\mathbf{z})} [f(\mathbf{z})]) (\nabla_\lambda \log \hat{\pi}_\beta(\mathbf{z}) - \mathbb{E}_{\pi_\beta(\mathbf{z})} [\nabla_\lambda \log \hat{\pi}_\beta(\mathbf{z})])] \end{aligned}$$

9 We emphasize that our estimator applies in the general case of expectations over $\pi_\beta(\mathbf{z})$, where the REINFORCE
 10 estimator $\mathbb{E}_{\pi_\beta(\mathbf{z})} [f(\mathbf{z}) \nabla_\phi \log \pi_\beta(\mathbf{z})]$ would require differentiating through $\log \pi_\beta(\mathbf{z})$ which contains the intractable
 11 normalizing constant. Additionally unlike REINFORCE, where a baseline is typically added ad-hoc to reduce
 12 variance, the baseline $b = \mathbb{E}_{\pi_\beta(\mathbf{z})} [f(\mathbf{z})]$ naturally appears as a result of differentiating through $\pi_\beta(\mathbf{z})$ using the identity
 13 $\nabla_\lambda \log Z_{\lambda, \beta}(\mathbf{x}) = \mathbb{E}_{\pi_\beta(\mathbf{z})} [\nabla_\lambda \log \hat{\pi}_{\lambda, \beta}(\mathbf{z})]$ derived in appendix E.

14 **Low variance** To address concern 2.6 of reviewer 1 directly, in section 6.1 (Figures 3 and 4) we compare the TVO
 15 against VIMCO which uses REINFORCE updates, and in 6.2 we compare the TVO against VAEs and IWAE which
 16 use the reparameterization trick. In figure 4 we plot the ϕ gradient std of a discrete VAE and compare against VIMCO
 17 and reweighted-wake sleep. Our method has lower variance than VIMCO which uses REINFORCE updates. However
 18 we agree we need to discuss the low variance properties of our estimator further. We have made the following changes
 19 to the manuscript:

20 **1)** After clarifying the relationship between the covariance estimator and REINFORCE in the aforementioned addition
 21 to the appendix, we observe the $\mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z})]$ baseline that is implicit in our covariance estimator is equivalent to the
 22 "average baseline" commonly used in reinforcement learning to reduce variance. This provides a theoretical justification
 23 for using the average baseline which is typically chosen because of empirical success and intuitive appeal.

24 **2)** We include additional experimental results to report the effect of using the covariance estimator to train a model on
 25 the ELBO, and compare against the same model trained using REINFORCE and the reparameterization trick. We plot
 26 the std. of the θ, ϕ gradients across 10 trials. Preliminary results indicate the variance of our estimator is empirically
 27 equivalent to the variance of the reparameterization trick, while the variance of the reinforce estimator is unstable
 28 despite having 10x samples.

29 **3)** We include a table explicitly reporting the mean gradient std of the TVO compared to the REINFORCE-based
 30 VIMCO (reproduced below) and update the writing in the experimental section to make it clear that section 6.1 and 6.2
 31 are designed to compare our method against REINFORCE and the reparameterization trick respectively.

32 **4)** The second source of variance reduction comes from using the 'Common Random Numbers' (CRN) technique from
 33 Owen (chapter 8.6), which we now refer to by name in the manuscript. The terms in the TVO are highly correlated,
 34 thus we expect reusing a single batch of samples for each additional term will act to reduce variance according to
 35 equation 8.21 in Owen. However, because the covariance term breaks into both positive and negative terms, CRN could
 36 potentially increase variance. We therefore have included an additional experiment running the TVO with and w/o CRN
 37 and include the results in tabular form below.

38 **Connections to Wake-Sleep** The endpoints of the TVI, which the TVO approximates, corresponds to the two ob-
 39 jectives used in Wake-sleep to jointly learn a generative model and inference network. Therefore we can view the
 40 objectives as two approximations (a left and right Riemannian sum) of a single objective, the TVI. In section 4
 41 we discuss how the left endpoint (i.e. $\beta = 0$) corresponds to the first objective (the wake-phase θ update) which
 42 in turn corresponds to the ELBO discussed in detail in section 1. The right endpoint (i.e. $\beta = 1$) corresponds
 43 to the second objective, of which there are two variants (the wake-phase ϕ update and sleep-phase ϕ update) de-
 44 termined by whether one uses real or simulated data. Revisiting the text we feel we relied too heavily on the
 45 description of WS given by (Le et al, 2018b) and will revise the manuscript to make this connection more clear.
 46

Particles	2	5	10	50	Iterations	10	10m	20m	30m	40m
Reinforce (VIMCO)	4.48	8.10	5.72	5.57	TVO w/o CRN	52.33	2.88	2.57	2.39	2.47
TVO	5.48	4.31	2.39	1.34	TVO w/ CRN	8.19	1.38	1.17	1.05	1.03

Table 1: Mean gradient std across 10 seeds for TVO vs REINFORCE (Left) and TVO with and without common random numbers (CRN) (Right)