

# Sticking the landing: A simple reduced-variance gradient for ADVI

Geoffrey Roeder, Yuhuai Wu, and David Duvenaud  
{roeder, ywu, duvenaud}@cs.toronto.edu



## Main Idea

- We give an estimator of the reparameterized ELBO gradient with lower variance when the variational approximation is close to truth
- Bigger improvement for flexible families like normalizing flows, IWAE, Hamiltonian variational inference
- Simple to implement

## Three forms of the ELBO:

$$\begin{aligned}\mathcal{L}(\phi) &= \mathbb{E}_{\mathbf{z} \sim q}[\log p(\mathbf{x}|\mathbf{z})] - KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) && \text{(exact KL)} \\ &= \mathbb{E}_{\mathbf{z} \sim q}[\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z})] + \mathbb{H}[q_\phi] && \text{(exact entropy)} \\ &= \mathbb{E}_{\mathbf{z} \sim q}[\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] && \text{(Monte Carlo)}\end{aligned}$$

- KL seems lowest variance, because it analytically integrates out some terms

## Monte Carlo variance goes to zero!

... as the approximate posterior gets close to the true posterior

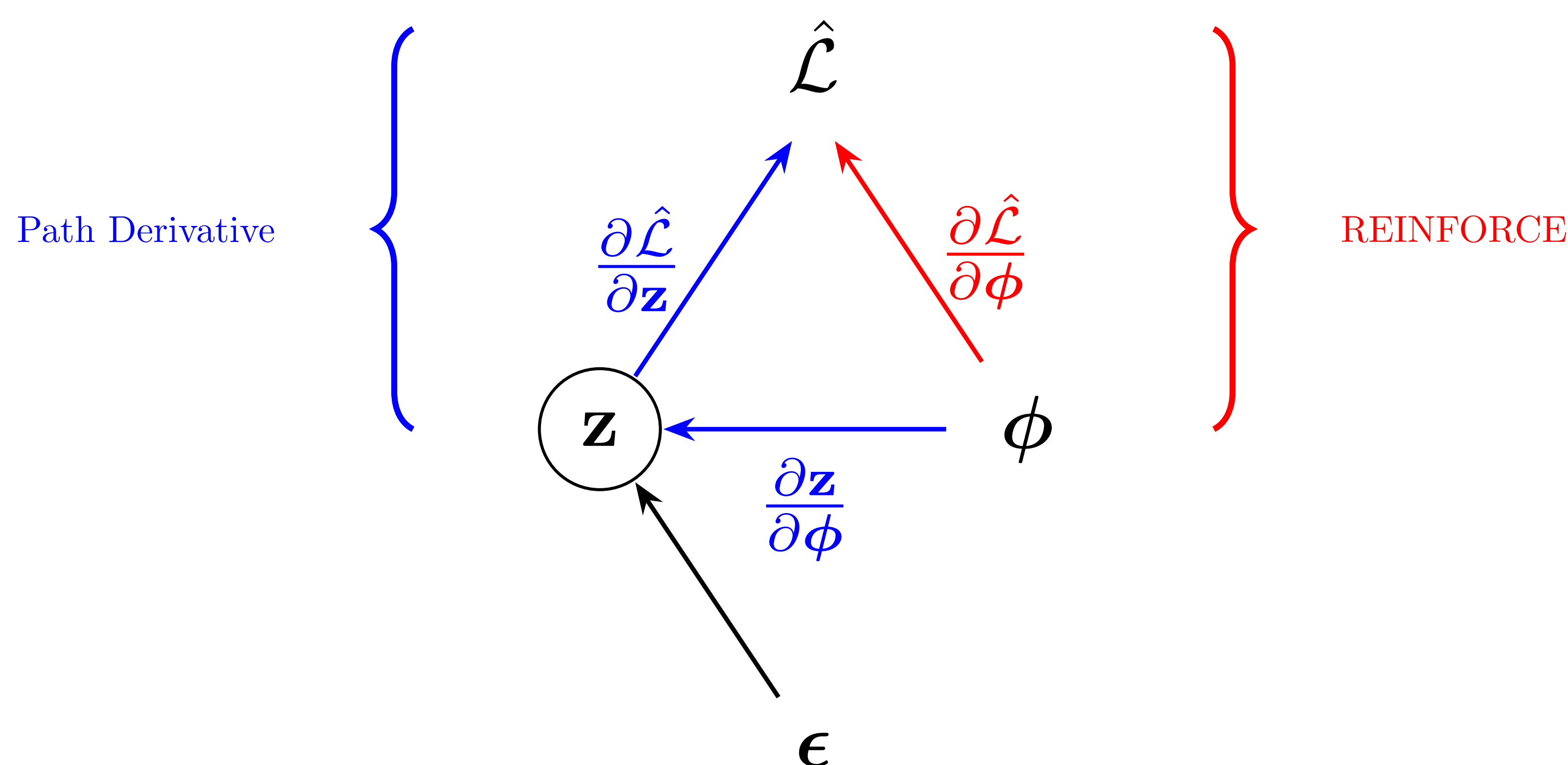
- If  $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$ , then a fully Monte Carlo estimator has 0 variance, since

$$\begin{aligned}\hat{\mathcal{L}}_{MC}(\phi) &= \log p(\mathbf{x}|\mathbf{z}_i) + \log p(\mathbf{z}_i) - \log q_\phi(\mathbf{z}_i|\mathbf{x}) && \mathbf{z}_i \sim_{\text{id}} q(\mathbf{z}) && (1) \\ &= \log p(\mathbf{z}_i|\mathbf{x}) + \log p(\mathbf{x}) - \log p(\mathbf{z}_i|\mathbf{x}) && \text{(using } q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}) \text{)} && (2) \\ &= \log p(\mathbf{x}) && && (3)\end{aligned}$$

a **constant** w.r.t.  $\mathbf{z}$ !

## But what about the gradient of the ELBO?

- It turns out that the naive fully Monte Carlo gradient estimator *doesn't* go to zero. Why not?



$$\hat{\nabla}_{MC} = \nabla_\phi [\log p(\mathbf{x}|\mathbf{z}_\phi) + \log p(\mathbf{z}_\phi) - \log q_\phi(\mathbf{z}_\phi|\mathbf{x})] \quad \epsilon \sim_{\text{id}} \mathcal{N}(0, I)$$

$$= \underbrace{\frac{\partial \log p(\mathbf{z}_\phi|\mathbf{x})}{\partial \mathbf{z}_\phi} \frac{\partial \mathbf{z}_\phi}{\partial \phi} + \frac{\partial \log p(\mathbf{z}_\phi)}{\partial \mathbf{z}_\phi} \frac{\partial \mathbf{z}_\phi}{\partial \phi}}_{\text{Path Derivative}} - \underbrace{\frac{\partial \log q(\mathbf{z}_\phi|\mathbf{x})}{\partial \mathbf{z}_\phi} \frac{\partial \mathbf{z}_\phi}{\partial \phi} - \frac{\partial \log q_\phi(\mathbf{z}|\mathbf{x})}{\partial \phi}}_{\text{REINFORCE}}$$

- The **Path Derivative** component is analytically 0, but the **REINFORCE** gradient has variance (equal to Fisher information) even when  $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$ .

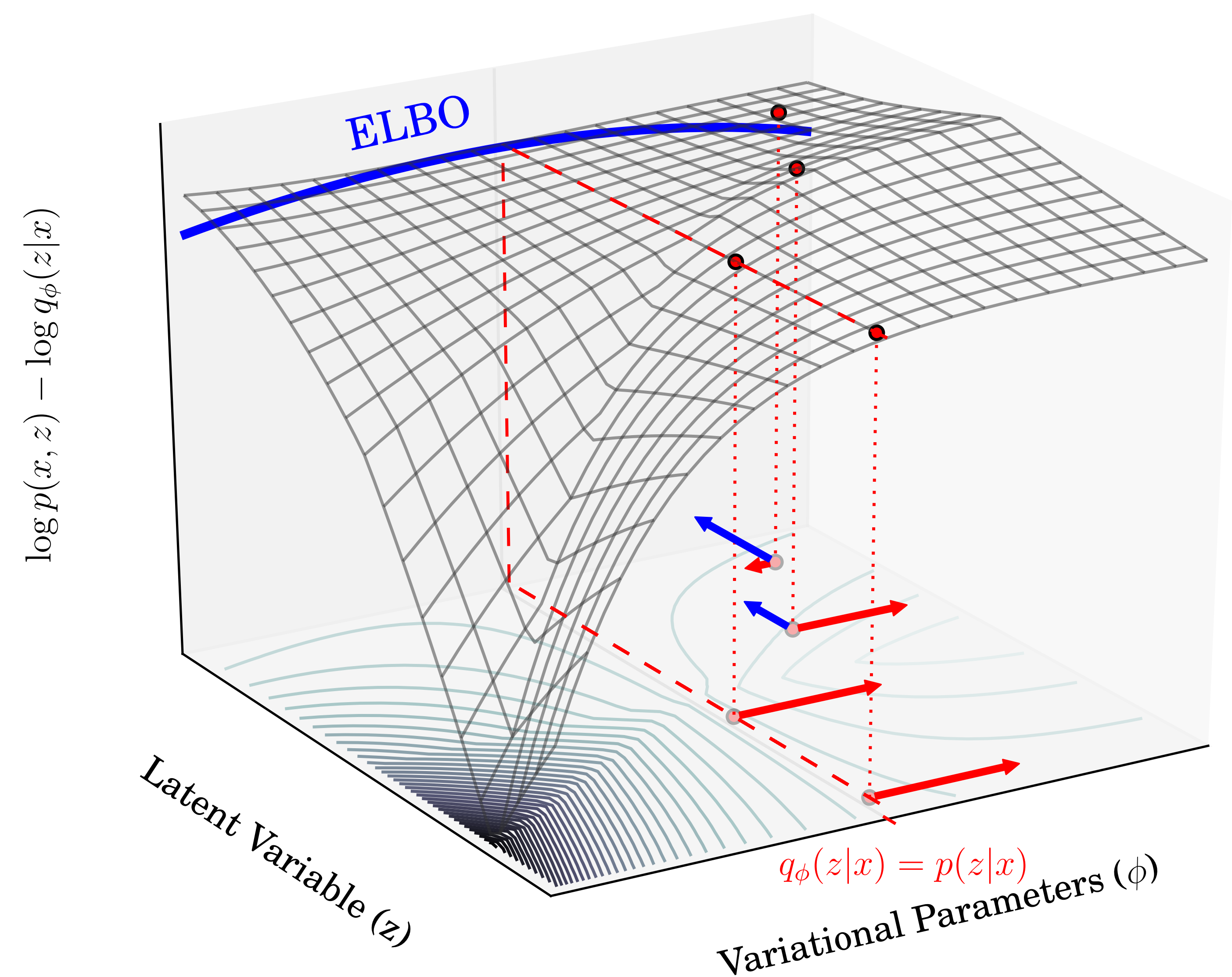
## Fix: remove the REINFORCE gradient!

- REINFORCE component (score function) has expectation 0
- **Still unbiased estimator of ELBO gradient**
- This can be interpreted as a control variate

## In other words:

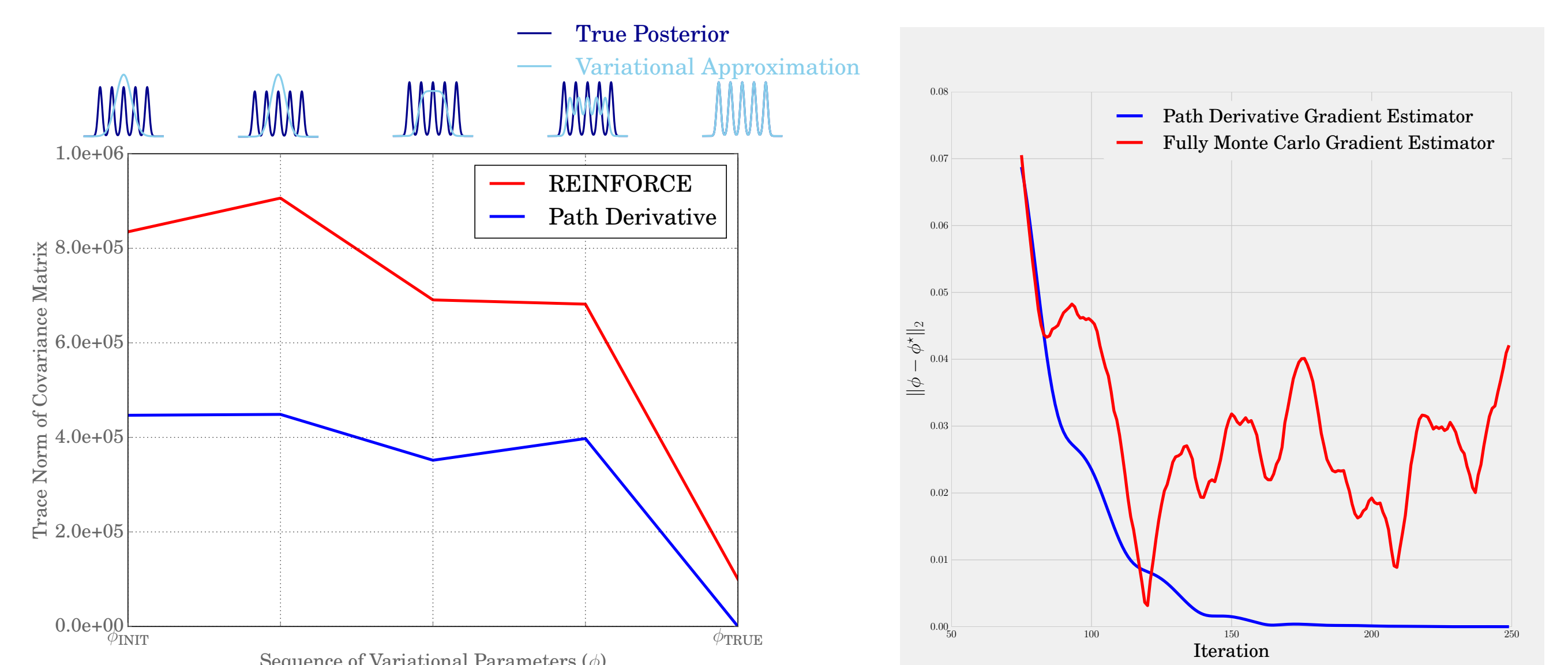
When the variational approximation is exact...

$\log p(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})$  Surface Along Trajectory through True  $\phi$



... the gradient w.r.t. the parameters is non-zero

## The new gradient estimator fixes this:



So the optimizer "sticks the landing"

## Implementing this is easy

- Block gradient through variational params. Ex. for Gaussian:

$$\log q = -\text{SUM}(\text{SQUARE}(\mathbf{z} - \mu_{\mathbf{z}})) / (2 * \text{EXP}(\log_{\text{sig\_z}})) - C - \text{SUM}(\log_{\text{sig\_z}})$$

becomes

$$\log q = -\text{SUM}(\text{SQUARE}(\mathbf{z} - \text{BLOCK}(\mu_{\mathbf{z}}))) / (2 * \text{EXP}(\text{BLOCK}(\log_{\text{sig\_z}}))) - C - 0.5 * \text{SUM}(\text{BLOCK}(\log_{\text{sig\_z}}))$$

- Use `gradient_disconnected`, `stop_gradient` in Theano, TensorFlow
- In autograd, define custom gradient that only evaluates path derivative

## Future work

- Unlikely that lower variance is maintained away from true posterior
- Control variates use an optimal scaling parameter: implement this!
- Empirical study of performance in multi-sample setting (IWAE)

## Related work

- BBSVI (no reparameterization) uses a similar control variate
- Tan et al. 2016 note phenomenon in sparse precision Gaussian VI
- Han et al. 2015 note phenomenon in Gaussian copula VI models
- Our goal is to present unified analysis with easy implementation