

Forward χ^2 Divergence Based Variational Importance Sampling

3 ICLR 2024 [Spotlight]



1 2 3 4 3 6 7 8 9

IWAE 0 / 2 8 9 8 9 8 9

Chengrui Li, Yule Wang, Weihan Li, Anqi Wu @ GaTech CSE



 $p(\boldsymbol{x}, \boldsymbol{z}; \theta) d\boldsymbol{z} \longrightarrow \text{latent variable}$ marginal likelihood $\leftarrow p(\boldsymbol{x}; \theta) = \boldsymbol{I}$ observed variable . intractable integral complete likelihood Maximum likelihood estimation (MLE): $\hat{\theta} = \arg \max_{\theta} p(\boldsymbol{x}; \theta)$

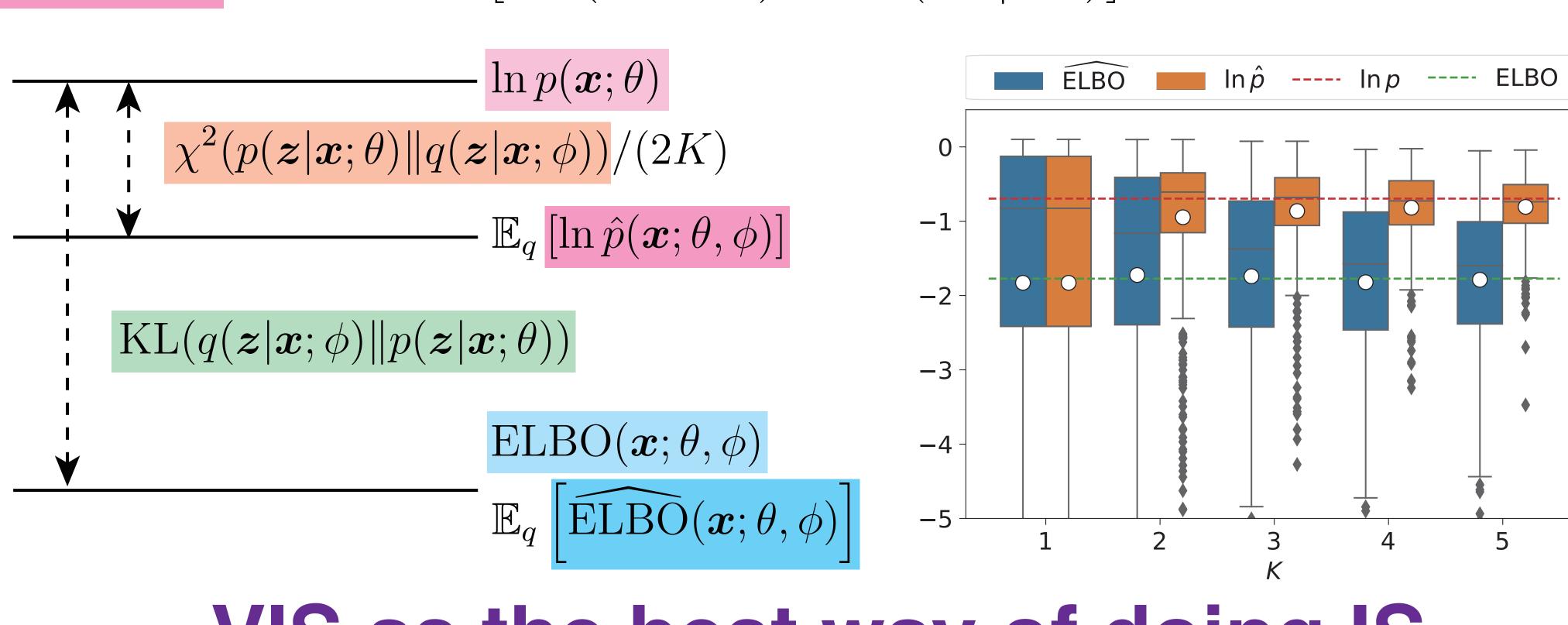
VI and IS, and their biases

variational inference (VI) importance sampling (IS) proposal distribution variational distribution $q(oldsymbol{z}|oldsymbol{x};\phi)$ $\mathrm{ELBO}(\boldsymbol{x}; \theta, \phi)$ $\ln p(\boldsymbol{x}; \theta)$ target function $\ln \hat{p}(\boldsymbol{x}; \theta, \phi)$ numerical estimator $\widehat{\mathrm{ELBO}}(oldsymbol{x}; heta,\phi)$

$$\underbrace{\text{ELBO}(\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\phi})}_{\text{ELBO}(\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\phi})} = \underbrace{\mathbb{E}_{q}[\ln p(\boldsymbol{x}, \boldsymbol{z}; \boldsymbol{\theta}) - \ln q(\boldsymbol{z}|\boldsymbol{x}; \boldsymbol{\phi})]}_{K} \quad \text{samples} \\
\underbrace{\text{ELBO}(\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\phi})}_{\text{ELBO}(\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\phi})} = \frac{1}{K} \sum_{k=1}^{K} \left[\ln p\left(\boldsymbol{x}, \boldsymbol{z}^{(k)}; \boldsymbol{\theta}\right) - \ln q\left(\boldsymbol{z}^{(k)}|\boldsymbol{x}; \boldsymbol{\phi}\right) \right] \quad \underbrace{\left\{\boldsymbol{z}^{(k)}\right\}_{k=1}^{K} \sim q(\boldsymbol{z}|\boldsymbol{x}; \boldsymbol{\phi})}_{k=1} \quad \text{samples}$$

 $\ln p(\boldsymbol{x}; \theta) = \ln \mathbb{E}_q[p(\boldsymbol{x}, \boldsymbol{z}; \theta) / q(\boldsymbol{z} | \boldsymbol{x}; \phi)]$

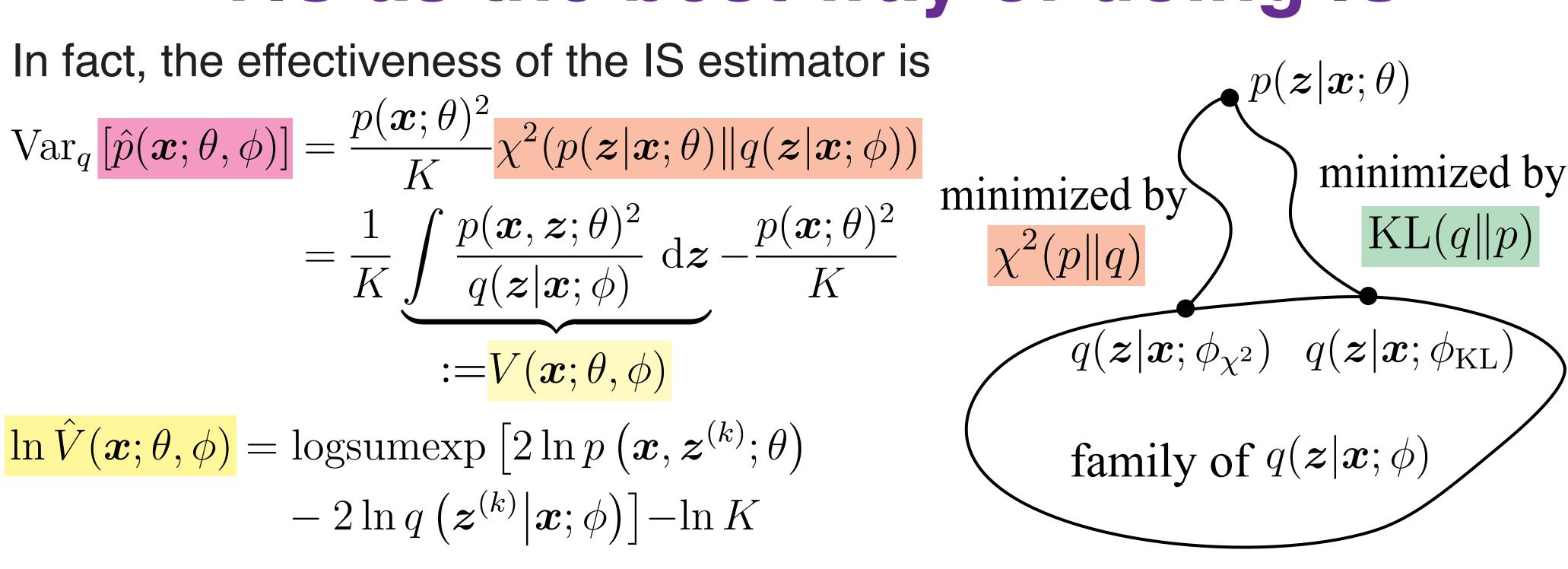
 $\ln \hat{p}(\boldsymbol{x}; \theta, \phi) = \text{logsumexp} \left[\ln p\left(\boldsymbol{x}, \boldsymbol{z}^{(k)}; \theta\right) - \ln q\left(\boldsymbol{z}^{(k)} \middle| \boldsymbol{x}; \phi\right) \right] - \ln K$



VIS as the best way of doing IS

 $:=V(\boldsymbol{x};\theta,\phi)$

 $\ln \hat{V}(\boldsymbol{x}; \theta, \phi) = \text{logsumexp} \left[2 \ln p \left(\boldsymbol{x}, \boldsymbol{z}^{(k)}; \theta \right) \right]$ $-2\ln q\left(\boldsymbol{z}^{(k)}|\boldsymbol{x};\phi\right)$ $-\ln K$



M-step

VIS algorithm

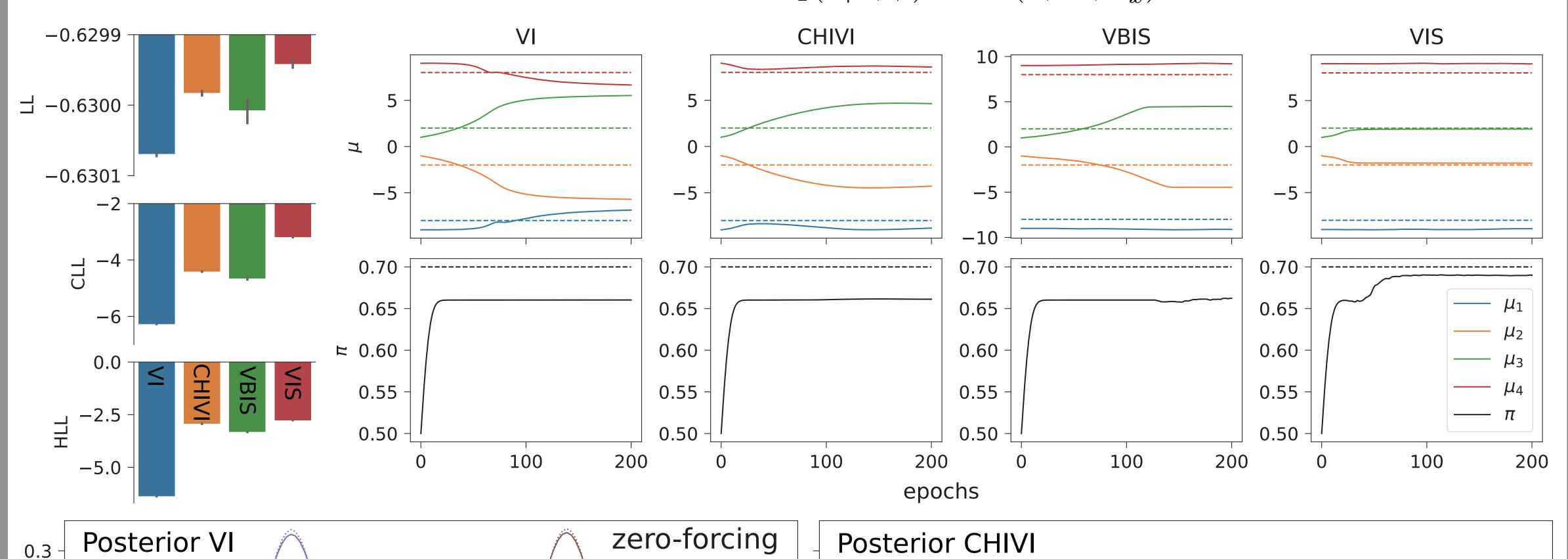
variational importance sampling (VIS)

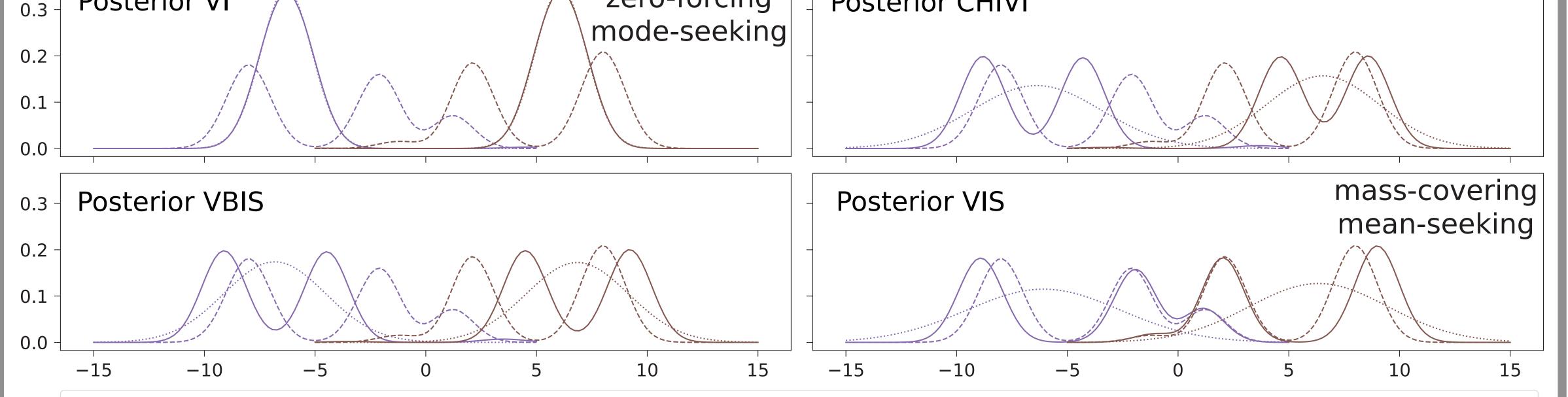
 $\left\{ oldsymbol{z}^{(k)}
ight\}_{k=1}^K \sim q(oldsymbol{z} | oldsymbol{x}; \phi)$ Sample $\min \chi^2(p(\boldsymbol{z}|\boldsymbol{x};\theta)||q(\boldsymbol{z}|\boldsymbol{x};\phi))$ $\min \mathrm{KL}(q(\boldsymbol{z}|\boldsymbol{x};\phi)||p(\boldsymbol{z}|\boldsymbol{x};\theta))$ $\rightarrow \min \ln V(\boldsymbol{x}; \theta, \phi)$ $\rightarrow \max \text{ELBO}(\boldsymbol{x}; \theta, \phi)$ E-step $\frac{\partial \operatorname{ELBO}(\boldsymbol{x}; \boldsymbol{\theta}, \phi)}{\partial \phi} \approx \frac{\partial}{\partial \phi} \frac{-1}{2K} \sum_{k=1}^{K}$ $\frac{\partial \ln V(\boldsymbol{x}; \boldsymbol{\theta}, \phi)}{\partial \phi} \approx \frac{\partial}{\partial \phi} \frac{1}{2} \ln \hat{V}(\boldsymbol{x}; \boldsymbol{\theta}, \phi)$ $\left[\ln p\left(\boldsymbol{x}, \boldsymbol{z}^{(k)}; \theta\right) - \ln q\left(\boldsymbol{z}^{(k)} | \boldsymbol{x}; \phi\right)\right]^{2}$

$\max_{\boldsymbol{\theta}} \ln p(\boldsymbol{x}; \boldsymbol{\theta})$ $\max \ln p(\boldsymbol{x}; \theta)$ $\rightarrow \max \text{ELBO}(\boldsymbol{x}; \theta, \phi)$ $\partial \operatorname{ELBO}(\boldsymbol{x}; \theta, \phi) \approx \frac{\partial}{\partial \theta} \widehat{\operatorname{ELBO}}(\boldsymbol{x}; \theta, \phi)$ $\frac{\partial \ln p(\boldsymbol{x}; \theta)}{\partial \theta} \approx \frac{\partial}{\partial \phi} \ln \hat{p}(\boldsymbol{x}; \theta, \phi)$

Experiments

Mixture model: $p(z;\theta) = \sum_{i=1}^{n} \pi_i \mathcal{N}(z;\mu_i, 1^2), \quad p(x|z;\theta) = \text{Bernoulli}(x; \text{logistic}(z))$ **GMM-Bernoulli** $q(z|x;\phi) = \mathcal{N}(z;c_x,\sigma_x^2)$



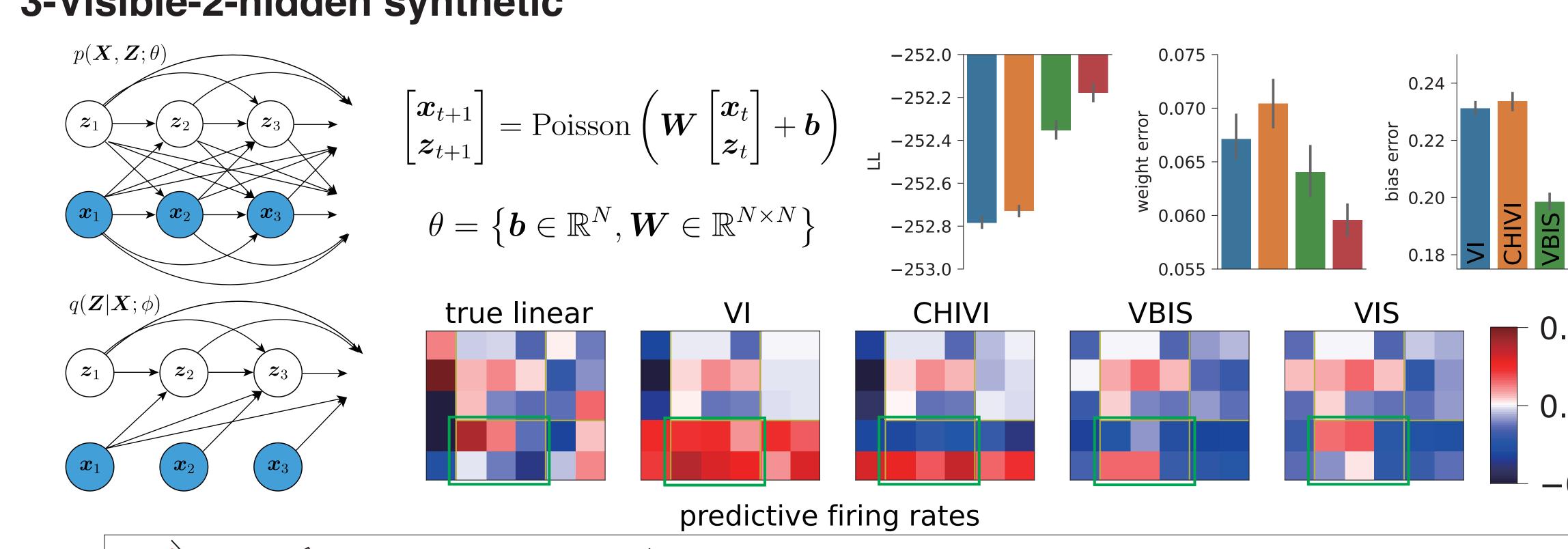


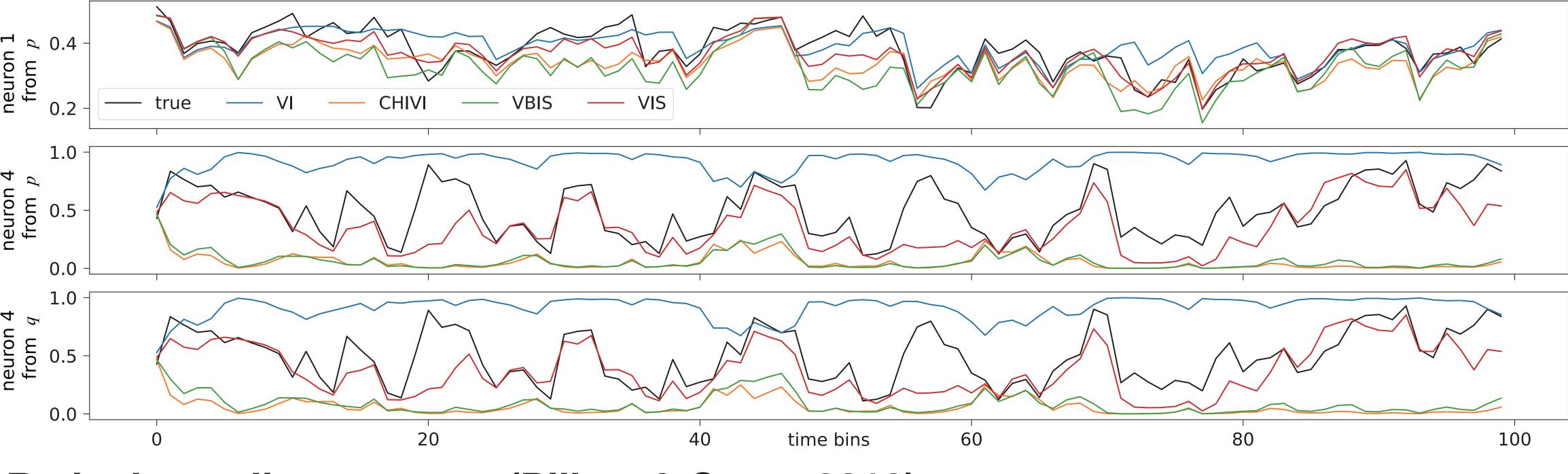
 $---- p(z|x=1;\theta)$

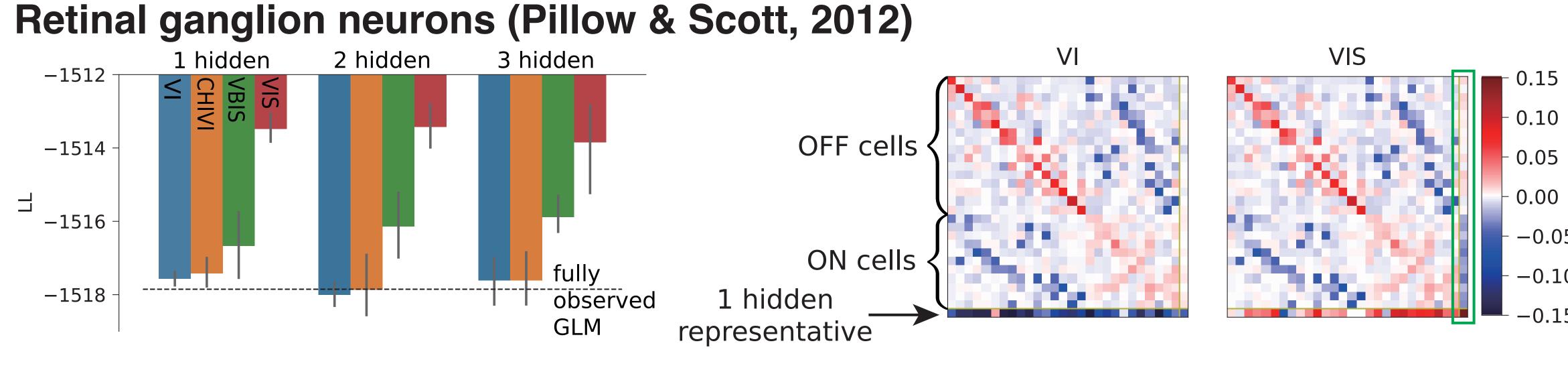
$p(z;\theta) = \mathcal{N}(z; \mathbf{0}, \mathbf{I}), \quad p(x|z;\theta) = \text{Bernoulli}(x; \text{logistic}(\text{decoder}(z)))$ VAE on MNIST $q(\boldsymbol{z}|\boldsymbol{x};\phi) = \mathcal{N}\left(\boldsymbol{x};\boldsymbol{\mu}(\boldsymbol{x}),\operatorname{diag}\boldsymbol{\sigma}^2(\boldsymbol{x})\right), \quad \boldsymbol{\mu}(\boldsymbol{x}),\boldsymbol{\sigma}^2(\boldsymbol{x}) = \operatorname{encoder}(\boldsymbol{x})$ raw **b** 1 2 3 4 5 6 7 8 9 -140 -V 0 / 3 5 9 5 9 9 CHIVI 0 / 3 3 9 5 6 7 8 9 VBIS 0 1 2 3 4 5 6 7 8 4

Partially observable GLM

A very hard problem since $p(\boldsymbol{x}, \boldsymbol{z}; \theta)$ cannot be explicitly factored as $p(\boldsymbol{x}|\boldsymbol{z}; \theta)p(\boldsymbol{z}; \theta)$. 3-Visible-2-hidden synthetic







References: [1] Burda et al., arXiv, 2015. [2] Dieng et al., NeurIPS, 2017. [3] Finke & Thiery, arXiv, 2019. [4] Jerfel et al., PMLR, 2018. [5] Domke & Sheldon, NeurIPS, 2021. [6] Su & Chen, Comp. Stat., 2021. [7] Akyildiz & Miguez, Stat. Comp., 2021.