



Latent Variable Model

marginal likelihood $\leftarrow p(\mathbf{x}; \theta) = \int p(\mathbf{x}, \mathbf{z}; \theta) d\mathbf{z} \rightarrow$ latent variable

observed variable \leftarrow intractable integral \leftarrow complete likelihood \leftarrow parameter set

Maximum likelihood estimation (MLE): $\hat{\theta} = \arg \max_{\theta} p(\mathbf{x}; \theta)$

VI and IS, and their biases

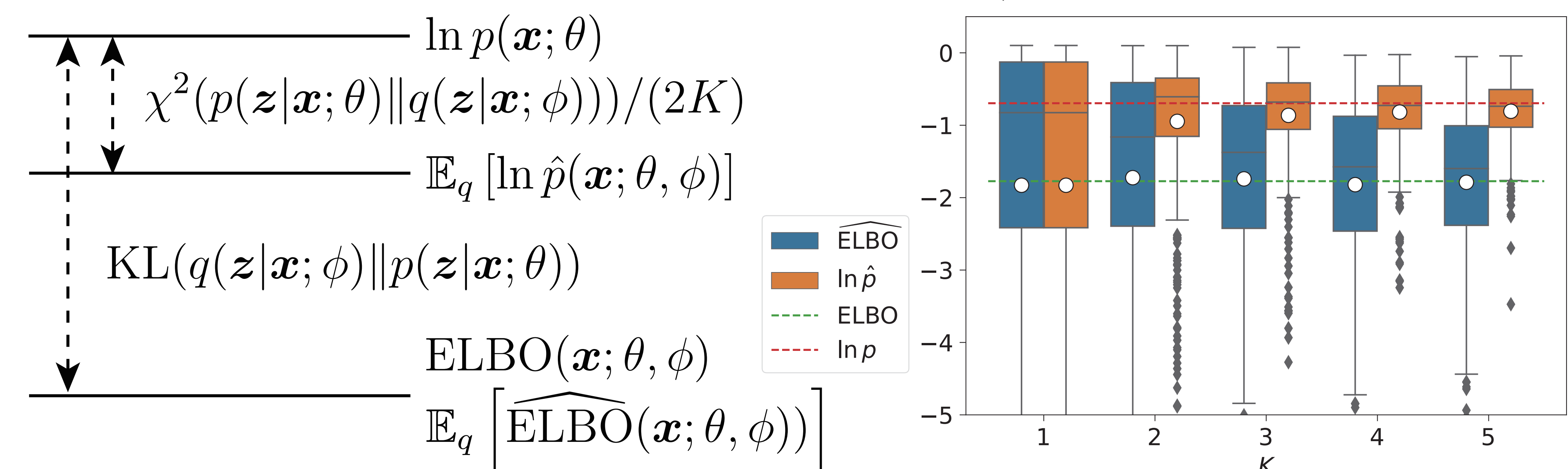
	variational inference (VI)	importance sampling (IS)
target function	$q(\mathbf{z} \mathbf{x}; \phi)$	variational distribution
numerical estimator	$\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)$	porposal distribution
	$\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)$	$\ln p(\mathbf{x}; \theta)$
		$\ln \hat{p}(\mathbf{x}; \theta, \phi)$

$$\text{ELBO}(\mathbf{x}; \theta, \phi) = \mathbb{E}_q[\ln p(\mathbf{x}, \mathbf{z}; \theta) - \ln q(\mathbf{z}|\mathbf{x}; \phi)]$$

$$\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi) = \frac{1}{K} \sum_{k=1}^K [\ln p(\mathbf{x}, \mathbf{z}^{(k)}; \theta) - \ln q(\mathbf{z}^{(k)}|\mathbf{x}; \phi)]$$

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

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



- Compared with $\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)$, $\ln \hat{p}(\mathbf{x}; \theta, \phi)$ is an asymptotically tighter lower bound of $\ln p(\mathbf{x}; \theta)$.
- When $K = 1$, $\text{ELBO}(\mathbf{x}; \theta, \phi) = \ln \hat{p}(\mathbf{x}; \theta, \phi)$.
- If $K \geq 2$, IS estimates $\ln p(\mathbf{x}; \theta)$ better than VI, so we can use IS to learn θ .

VIS as the best way of doing IS

In fact, the effectiveness of the IS estimator is

$$\text{Var}_q[\hat{p}(\mathbf{x}; \theta, \phi)] = \frac{p(\mathbf{x}; \theta)^2}{K} \chi^2(p(\mathbf{z}|\mathbf{x}; \theta) || q(\mathbf{z}|\mathbf{x}; \phi))$$

So, we should minimize this forward χ^2 divergence w.r.t. ϕ to find the optimal choice of the proposal distribution $q(\mathbf{z}|\mathbf{x}; \phi)$ for the current $p(\mathbf{z}|\mathbf{x}; \theta)$.

$$\chi^2(p(\mathbf{z}|\mathbf{x}; \theta) || q(\mathbf{z}|\mathbf{x}; \phi)) = \frac{1}{p(\mathbf{x}; \theta)^2} \int \frac{p(\mathbf{x}, \mathbf{z}; \theta)^2}{q(\mathbf{z}|\mathbf{x}; \phi)} d\mathbf{z} - 1 =: \frac{1}{p(\mathbf{x}; \theta)^2} V(\mathbf{x}; \theta, \phi) - 1$$

It is converted to minimizing $V(\mathbf{x}; \theta, \phi)$, which should be estimated and minimized in log space for numerical stability.

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

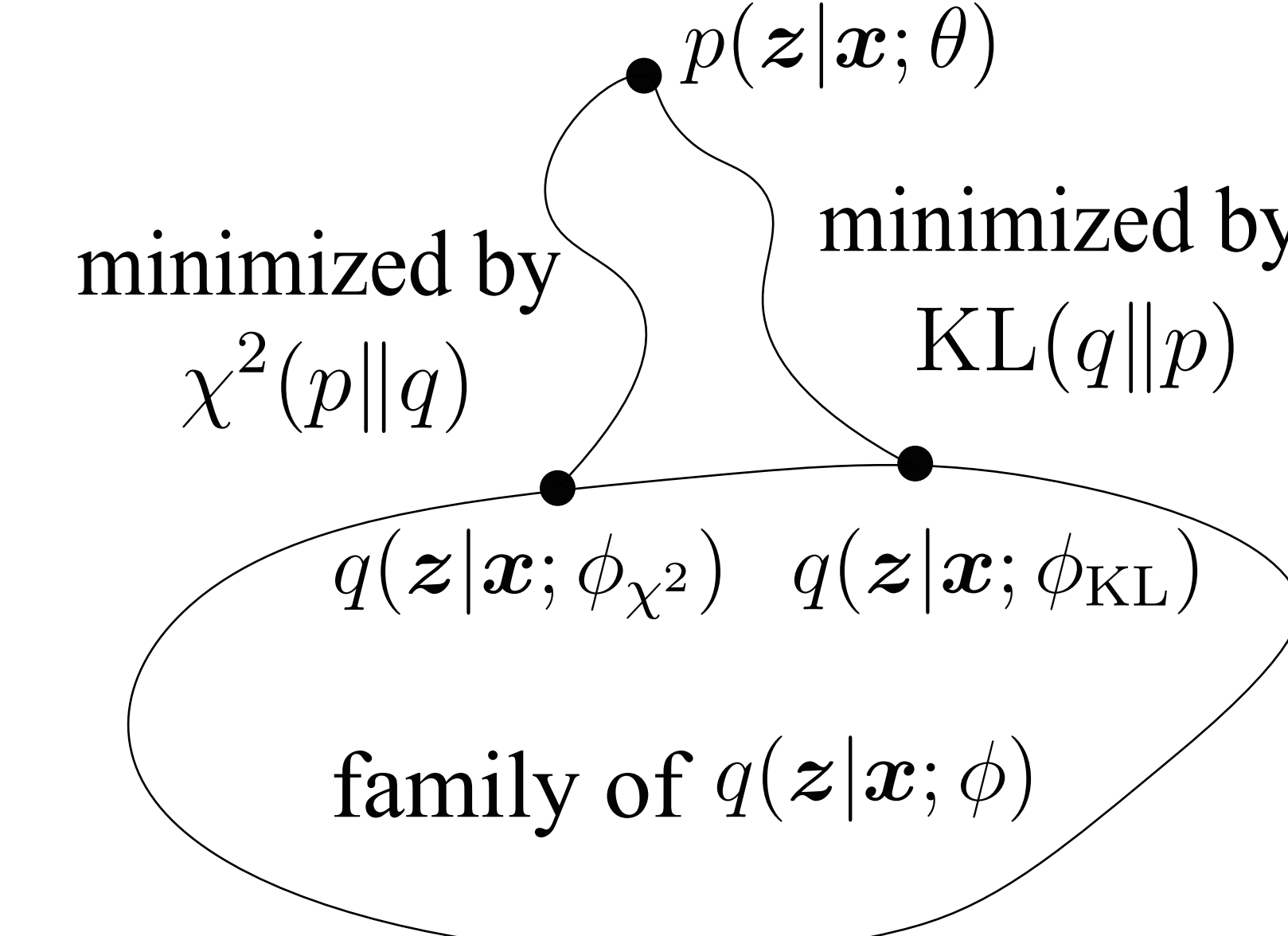
	VI	variational importance sampling (VIS)
Sample	$\{\mathbf{z}^{(k)}\}_{k=1}^K \sim q(\mathbf{z} \mathbf{x}; \phi)$	
E-step	Minimize $\text{KL}(q(\mathbf{z} \mathbf{x}; \phi) p(\mathbf{z} \mathbf{x}; \theta))$ w.r.t. ϕ by maximizing $\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)$ w.r.t. ϕ	Minimize $\text{KL}(q(\mathbf{z} \mathbf{x}; \phi) p(\mathbf{z} \mathbf{x}; \theta))$ w.r.t. ϕ by maximizing $\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)$ w.r.t. ϕ
M-step	Maximize $\ln p(\mathbf{x}; \theta)$ w.r.t. θ by maximizing $\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)$ w.r.t. θ	Maximize $\ln p(\mathbf{x}; \theta)$ w.r.t. θ by maximizing $\ln \hat{p}(\mathbf{x}; \theta, \phi)$ w.r.t. θ

VIS only changes two lines of the code.

The score function gradient estimators in the E-step:

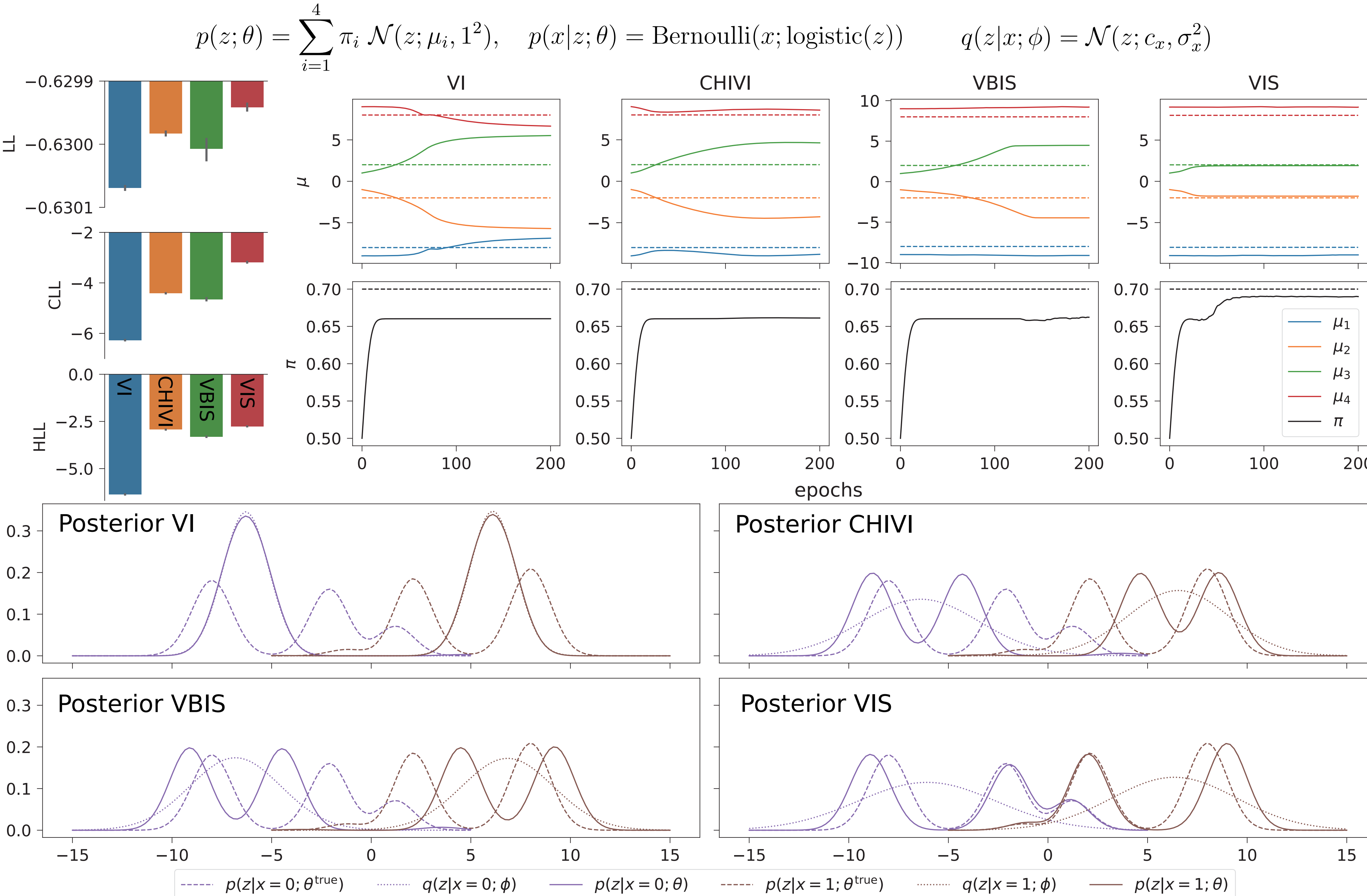
$$\frac{\partial \widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)}{\partial \phi} \approx \frac{1}{K} \sum_{k=1}^K \left\{ [\ln p(\mathbf{x}, \mathbf{z}^{(k)}; \theta) - \ln q(\mathbf{z}^{(k)}|\mathbf{x}; \phi)] \times \frac{\partial \ln q(\mathbf{z}^{(k)}|\mathbf{x}; \phi)}{\partial \phi} \right\}$$

$$\frac{\partial \ln V(\mathbf{x}; \theta, \phi)}{\partial \phi} \approx \frac{\partial}{\partial \phi} \frac{1}{2} \ln \hat{V}(\mathbf{x}; \theta, \phi)$$



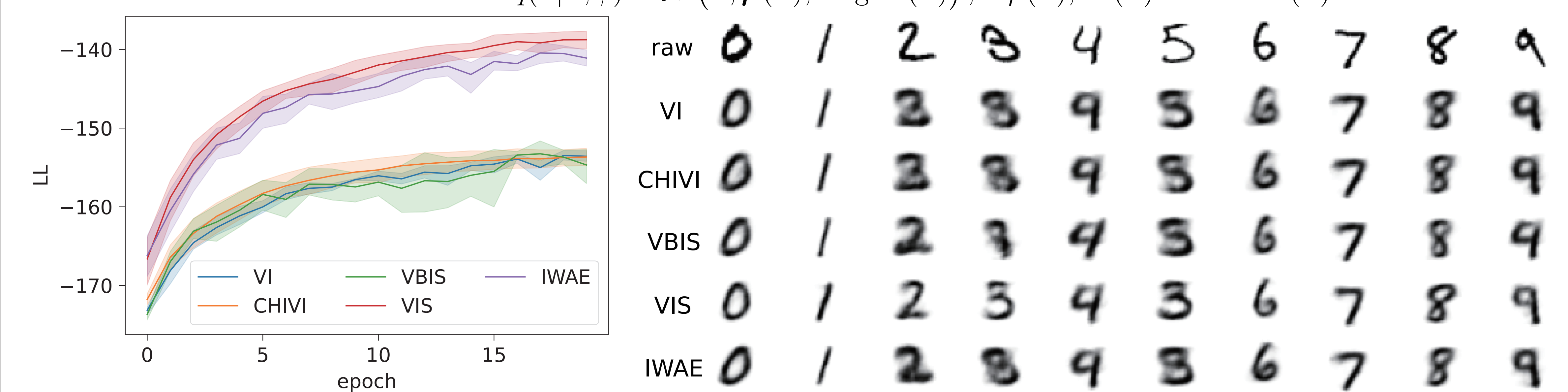
Experiments

Mixture model: GMM-Bernoulli



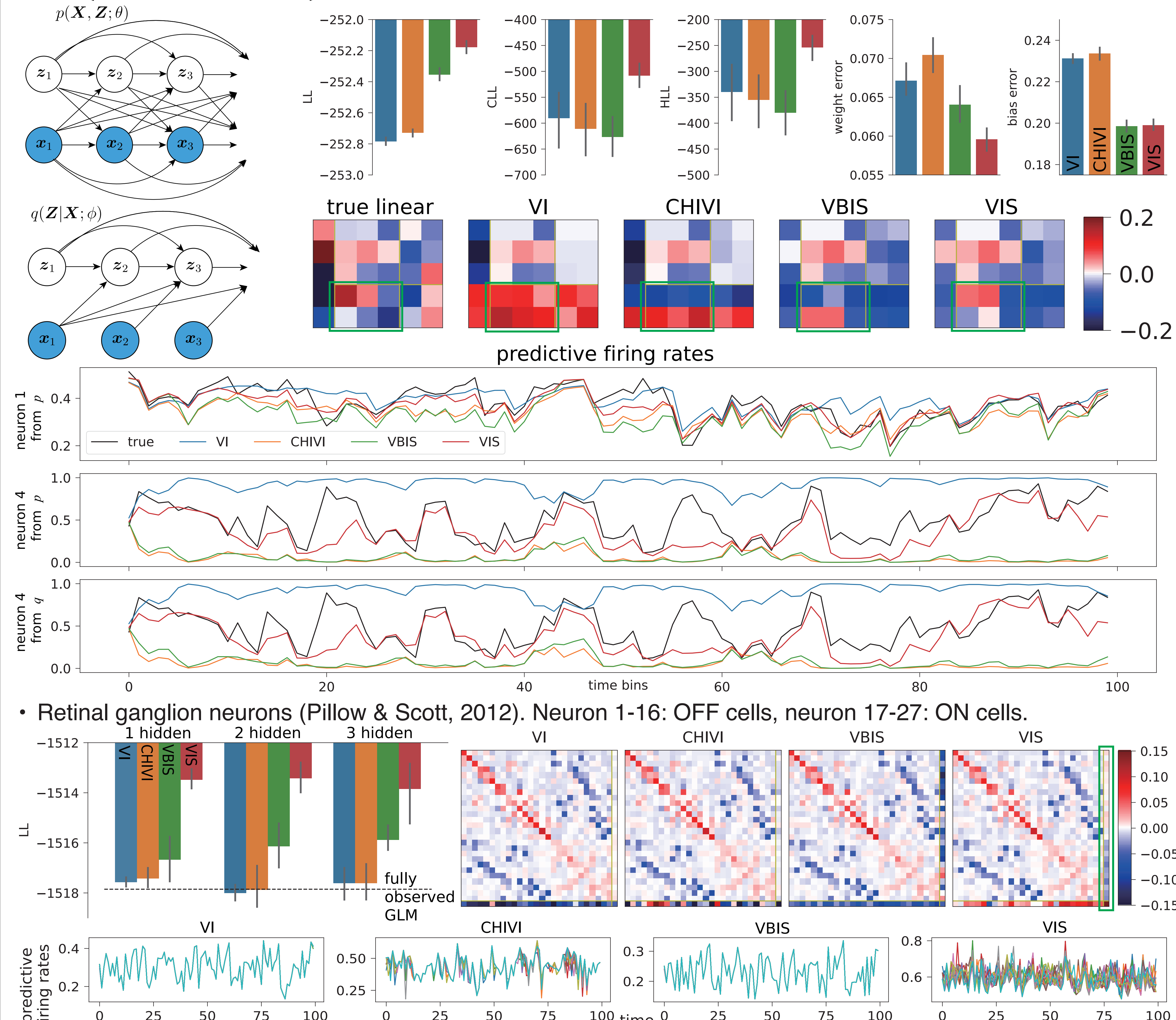
- VI: zero-forcing/mode-seeking behavior of minimizing the reverse KL. Good ELBO, reverse KL is nearly 0, but in fact both $p(\mathbf{z}|\mathbf{x}; \theta)$ and $q(\mathbf{z}|\mathbf{x}; \phi)$ are far from $p(\mathbf{z}|\mathbf{x}; \theta^{\text{true}})$.
- VIS: mass-covering/mean-seeking behavior of minimizing the forward χ^2 . This enlarges the effective support range of $q(\mathbf{z}|\mathbf{x}; \phi)$ for sampling.

VAE on MNIST



Partially observable GLM

- A very hard problem since $p(\mathbf{x}; \mathbf{z}; \theta)$ cannot be explicitly factored as $p(\mathbf{x}|\mathbf{z}; \theta)p(\mathbf{z}; \theta)$.
- \mathbf{X} are spike trains from V visible neurons and \mathbf{Z} are spike trains from $H = N - V$ hidden neurons.
- $\theta = \{\mathbf{b} \in \mathbb{R}^N, \mathbf{W} \in \mathbb{R}^{N \times N}\}$. $w_{n \leftarrow n'}$ is the weight from neuron n to neuron n' .



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.