

# Forward $\chi^2$ divergence based variational importance sampling

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

# Contents

1. Recap: variational inference (VI)
2. Importance sampling (IS)
3. VIS as the best way of doing IS
4. Applications to three latent variable models
  - Mixture model: GMM-Bernoulli
  - VAE on MNIST
  - Partially observable GLM

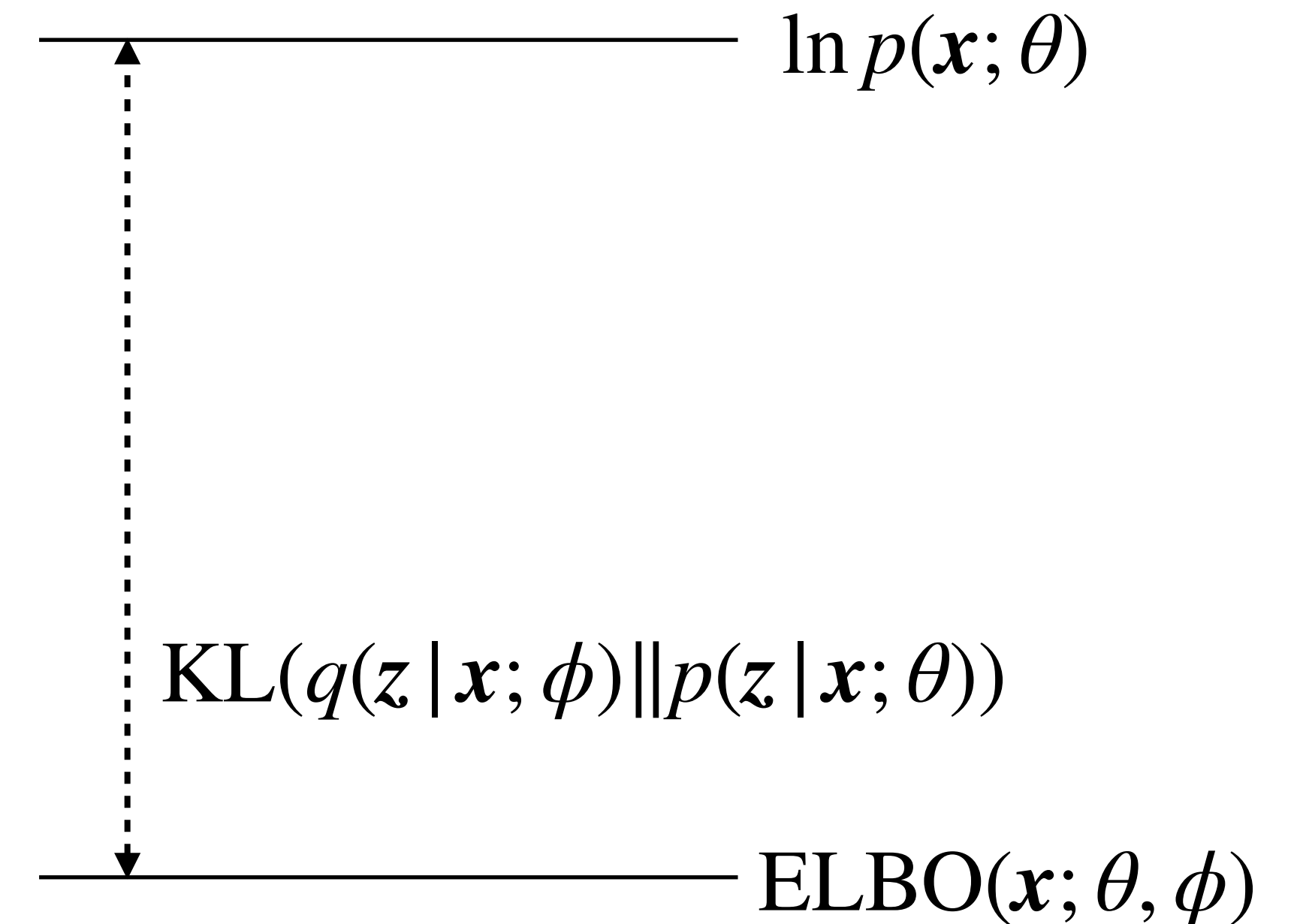
# 1 Recap: *variational* inference

# Latent variable model (LVM)

- Notation:
  - $\mathbf{x}$ : observed variables
  - $\mathbf{z}$ : latent variables
  - $\theta$ : parameter set
- Complete data likelihood:  $p(\mathbf{x}, \mathbf{z}; \theta)$
- Marginal likelihood:  $p(\mathbf{x}; \theta) = \int p(\mathbf{x}, \mathbf{z}; \theta) \, d\mathbf{z}$ , which is intractable when the problem is complicated.
- Maximum likelihood estimation:  $\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathbf{x}; \theta)$

# Variational inference

- Use a variational distribution  $q(\mathbf{z} | \mathbf{x}; \phi)$  to approximate  $p(\mathbf{z} | \mathbf{x}; \theta)$ .
- Consider the reverse KL divergence:
$$\begin{aligned}\text{KL}(q(\mathbf{z} | \mathbf{x}; \phi) || p(\mathbf{z} | \mathbf{x}; \theta)) &= \int q(\mathbf{z} | \mathbf{x}; \phi) \ln \frac{q(\mathbf{z} | \mathbf{x}; \phi)}{p(\mathbf{z} | \mathbf{x}; \theta)} d\mathbf{z} \\ &= \ln p(\mathbf{x}; \theta) - \mathbb{E}_q[\ln p(\mathbf{x}, \mathbf{z}; \theta) - \ln q(\mathbf{z} | \mathbf{x}; \phi)] \\ &= \ln p(\mathbf{x}; \theta) - \text{ELBO}(\mathbf{x}; \theta, \phi)\end{aligned}$$
- So, ELBO is a lower bound of  $\ln p(\mathbf{x}; \theta)$ .
- Maximize ELBO w.r.t.  $\theta$  increases  $\ln p(\mathbf{x}; \theta)$ , and maximize ELBO w.r.t.  $\phi$  reduces the reverse KL divergence  $\text{KL}(q(\mathbf{z} | \mathbf{x}; \phi) || p(\mathbf{z} | \mathbf{x}; \theta))$ .



# Numerical estimator of ELBO

- Numerical estimator of ELBO:

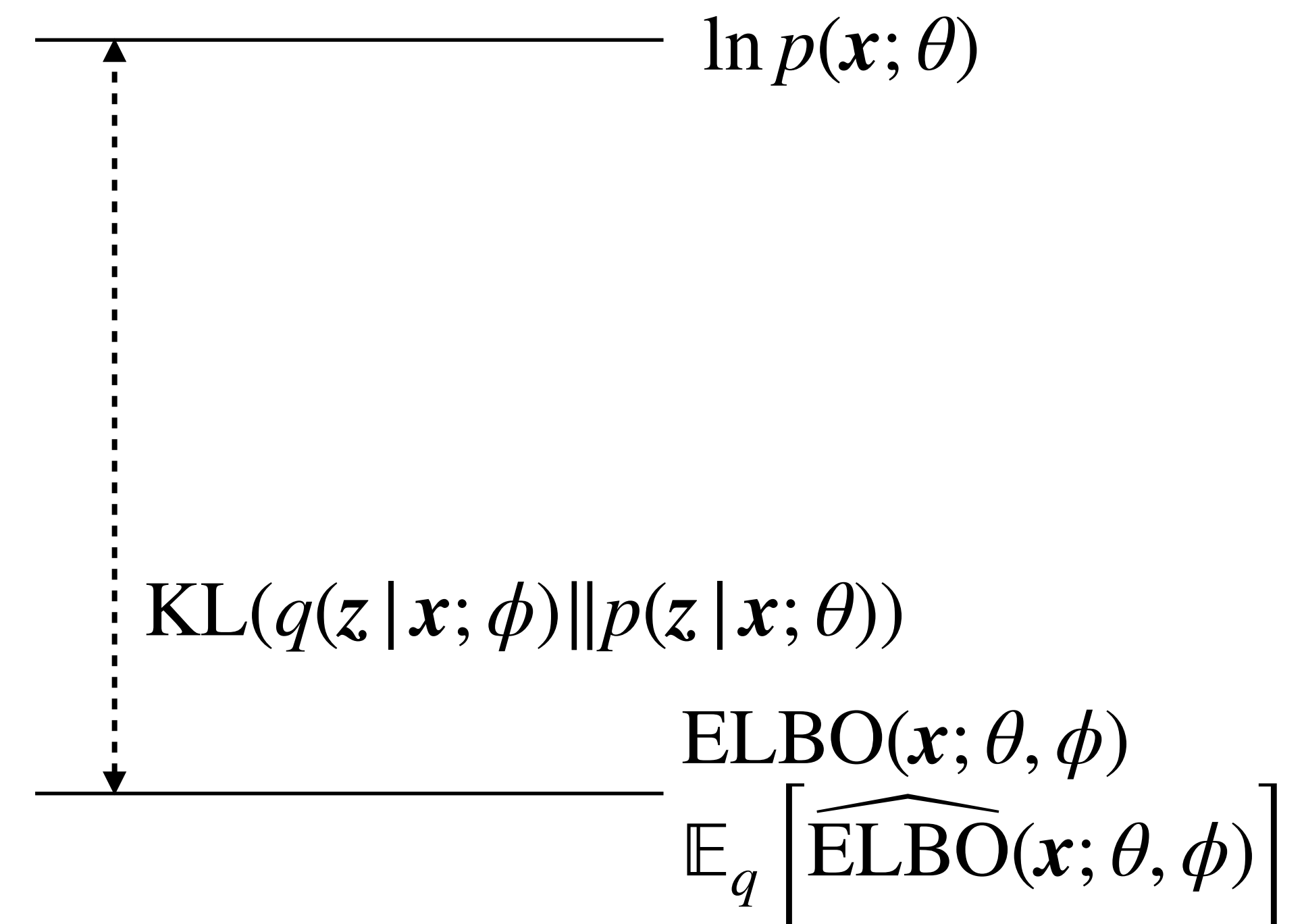
$$\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi) = \frac{1}{K} \sum_{k=1}^K \left[ \ln p(\mathbf{x}, \mathbf{z}^{(k)}; \theta) - \ln q(\mathbf{z}^{(k)} | \mathbf{x}; \phi) \right]$$

where  $\{\mathbf{z}^{(k)}\}_{k=1}^K \sim q(\mathbf{z} | \mathbf{x}; \phi)$ .

- The bias of the ELBO estimator:

$$\mathbb{E}_q \left[ \widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi) \right] = \text{ELBO}(\mathbf{x}; \theta, \phi)$$

- Is there any other estimator that approximates  $\ln p(\mathbf{x}; \theta)$  better?



# 2 Importance sampling

# Directly approximate the marginal

- Importance sampling uses a proposal distribution  $q(z | x; \phi)$  to estimate the integration:

$$p(\mathbf{x}; \theta) = \int p(\mathbf{x}, \mathbf{z}; \theta) \, d\mathbf{z} = \mathbb{E}_q[p(\mathbf{x}, \mathbf{z}; \theta) / q(\mathbf{z} | \mathbf{x}; \phi)] \approx \frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}^{(k)}; \theta)}{q(\mathbf{z}^{(k)} | \mathbf{x}; \phi)} =: \hat{p}(\mathbf{x}; \theta, \phi)$$

where  $\{\mathbf{z}^{(k)}\}_{k=1}^K \sim q(\mathbf{z} | \mathbf{x}; \phi)$ .

- Numerical stable IS estimator in log space

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

- And its gradient estimator

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



# Compare VI and IS

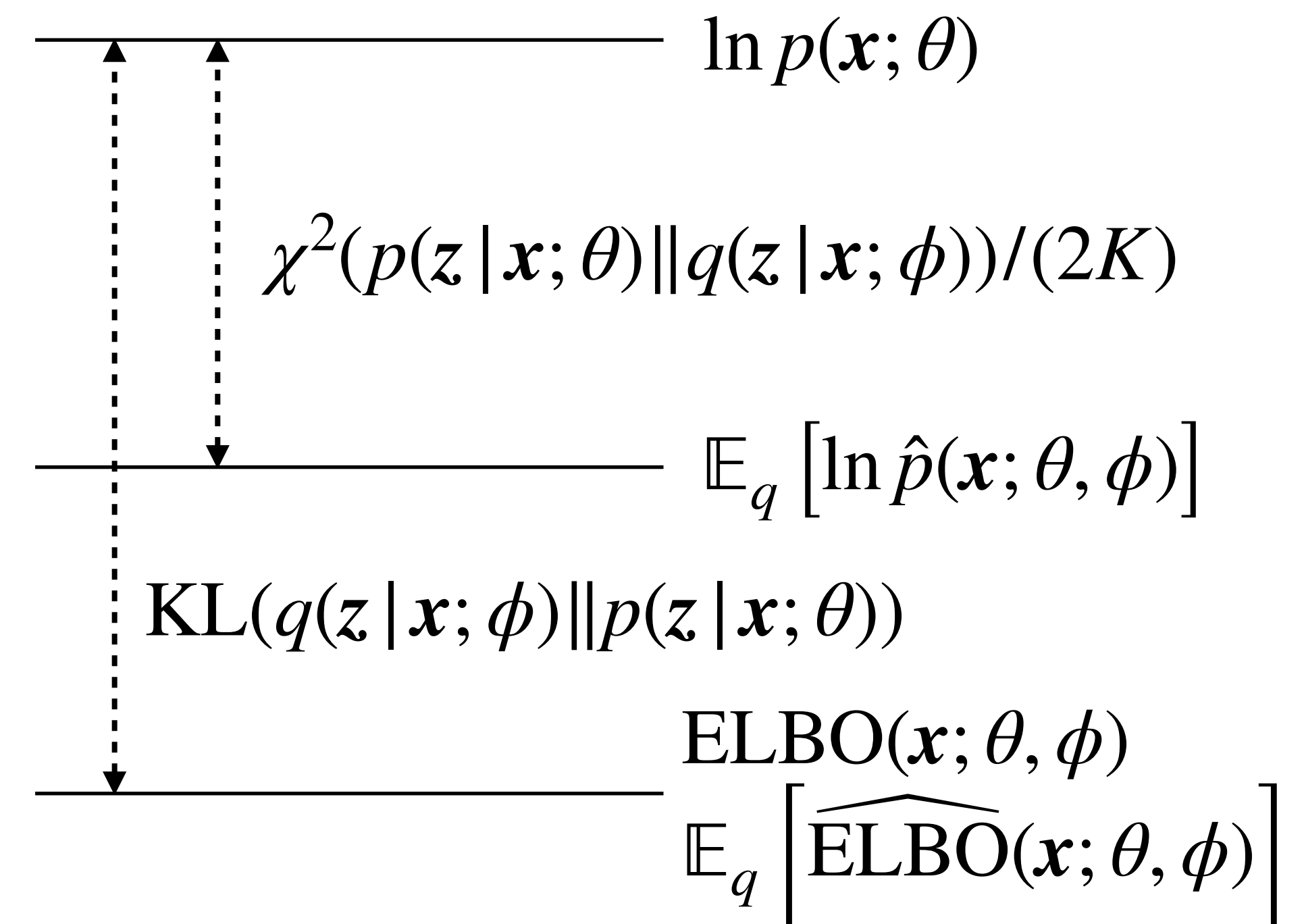
- The bias of the IS estimator:

$$\mathbb{E}_q [\ln \hat{p}(\mathbf{x}; \theta, \phi) - \ln p(\mathbf{x}; \theta)] \approx -\frac{1}{2K} \chi^2(p(\mathbf{z} | \mathbf{x}; \theta) \| q(\mathbf{z} | \mathbf{x}; \phi))$$

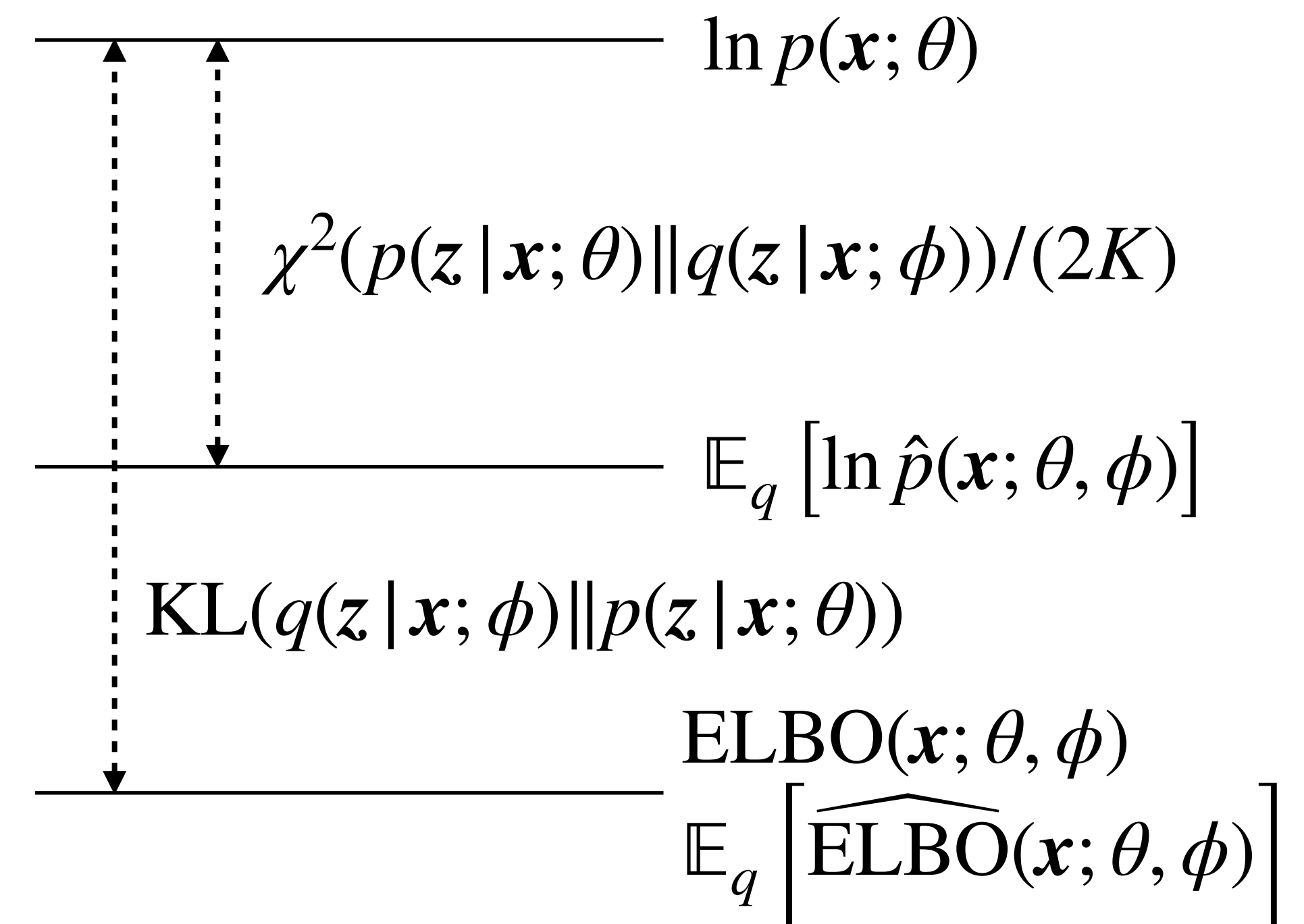
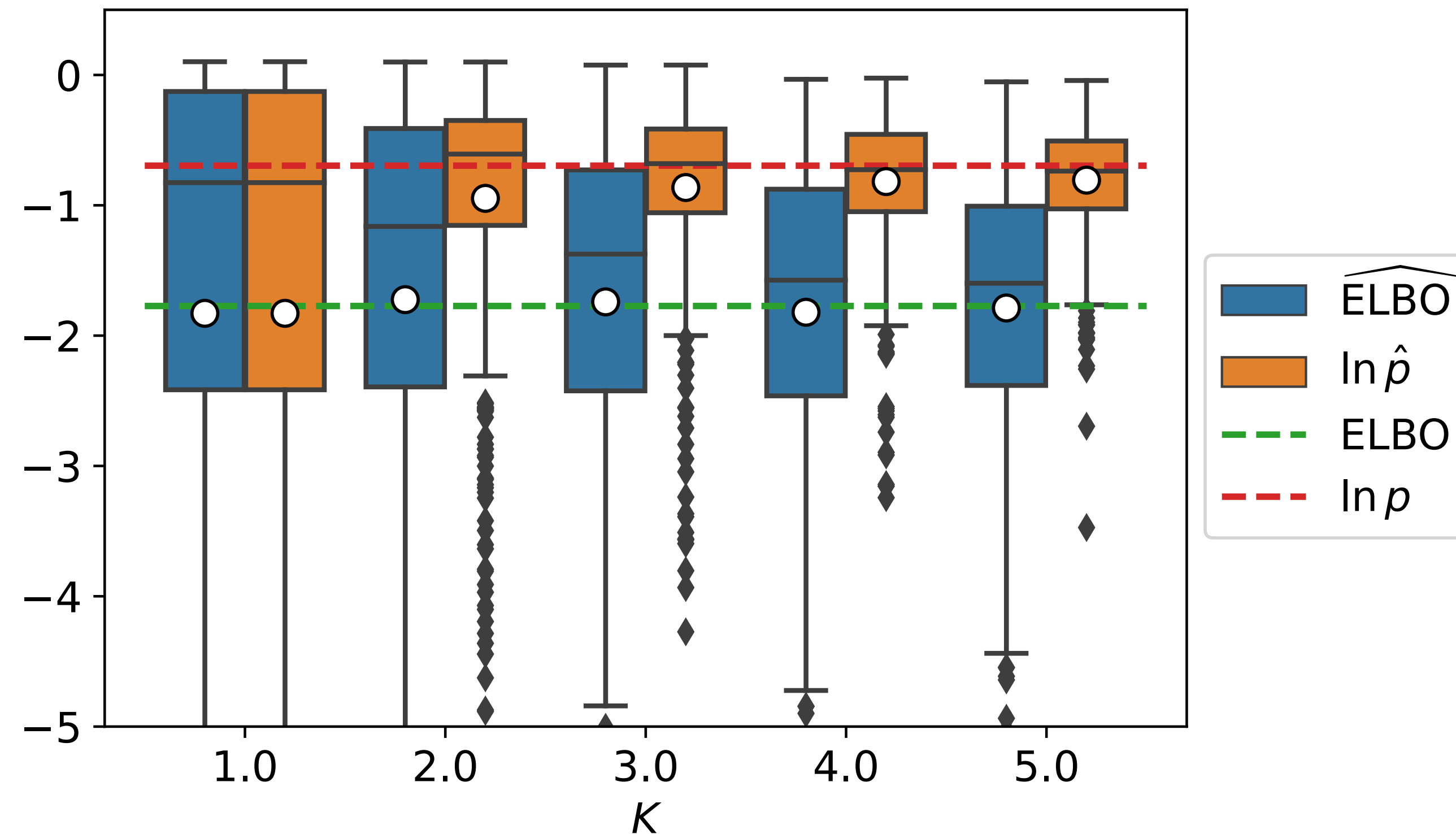
- The bias of the IS estimator  $\rightarrow 0$  as  $K \rightarrow \infty$ .

- Particularly, when  $K = 1$ ,  
 $\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi) = \ln \hat{p}(\mathbf{x}; \theta, \phi)$ .

- So, compared with  $\widehat{\text{ELBO}}(\mathbf{x}; \theta, \phi)$ ,  $\ln \hat{p}(\mathbf{x}; \theta, \phi)$  is an asymptotically tighter lower bound.



# Validate this by numerical simulations



Since IS estimates  $\ln p(\mathbf{x}; \theta)$  better than VI with a large  $K$ , we can use IS to learn  $\theta$ .

# 3 VIS as the best way of doing IS

# Choosing the optimal proposal distribution

- Remember the bias of the IS estimator is  $-\chi^2(p(\mathbf{z} | \mathbf{x}; \theta) || q(\mathbf{z} | \mathbf{x}; \phi)) / (2K)$ .
- The smaller the absolute bias, the better the estimator.
- In fact, the effectiveness of the estimator  $\text{Var}_q [\hat{p}(\mathbf{x}; \theta, \phi)] = \frac{p(\mathbf{x}; \theta)^2}{K} \chi^2(p(\mathbf{z} | \mathbf{x}; \theta) || q(\mathbf{z} | \mathbf{x}; \phi))$
- So, we should minimize this forward  $\chi^2$  divergence w.r.t.  $\phi$  to find the optimal choice of the proposal distribution  $q(\mathbf{z} | \mathbf{x}; \phi)$  for the current  $p(\mathbf{z} | \mathbf{x}; \theta)$ .
- Similar to the case we encountered in the reverse KL divergence, we don't know  $p(\mathbf{z} | \mathbf{x}; \theta)$ , so we derive the following equation.
- $$\chi^2(p(\mathbf{z} | \mathbf{x}; \theta) || q(\mathbf{z} | \mathbf{x}; \phi)) = \frac{1}{p(\mathbf{x}; \theta)^2} \int \frac{p(\mathbf{x}, \mathbf{z}; \theta)^2}{q(\mathbf{z} | \mathbf{x}; \phi)} d\mathbf{z} - 1 =: \frac{1}{p(\mathbf{x}; \theta)^2} V(\mathbf{x}; \theta, \phi) - 1$$
- We convert minimizing  $\chi^2(p(\mathbf{z} | \mathbf{x}; \theta) || q(\mathbf{z} | \mathbf{x}; \phi))$  to minimizing  $V(\mathbf{x}; \theta, \phi)$ .

# Estimate and minimize $V(\mathbf{x}; \theta, \phi)$ in log space

- For numerical stability,  $V(\mathbf{x}; \theta, \phi)$  should be minimized in log space.
- $\ln V(\mathbf{x}; \theta, \phi) \approx \text{logsumexp} \left[ 2 \ln p(\mathbf{x}, \mathbf{z}^{(k)}; \theta) - 2 \ln q(\mathbf{z}^{(k)} | \mathbf{x}; \phi) \right] - \ln K = \ln \hat{V}(\mathbf{x}; \theta, \phi)$
- The score function gradient estimator of  $\ln V(\mathbf{x}; \theta, \phi)$  is

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

# Comparison of VI and VIS

Only changes two lines in the code

	VI	VIS
Sample	Sample $\{z^{(k)}\}_{k=1}^K$ from $q(z   x; \phi)$	
E-step	Minimize $\text{KL}(q(z   x; \phi)    p(z   x; \theta))$ w.r.t. $\phi$ by maximizing $\widehat{\text{ELBO}}(x; \theta, \phi)$ w.r.t. $\phi$	Minimize $\chi^2(p(z   x; \phi)    q(z   x; \theta))$ w.r.t. $\phi$ by minimizing $\ln \hat{V}(x; \theta, \phi)$ w.r.t. $\phi$
M-step	Maximize $\ln p(x; \theta)$ w.r.t. $\theta$ by maximizing $\widehat{\text{ELBO}}(x; \theta, \phi)$ w.r.t. $\theta$	Maximize $\ln p(x; \theta)$ w.r.t. $\theta$ by maximizing $\ln \hat{p}(x; \theta, \phi)$ w.r.t. $\theta$

# 4 Applications to three latent variable models

# Mixture model: GMM-Bernoulli

- Model

- $p(z; \theta) = \sum_{i=1}^4 \mathcal{N}(z; \mu_i, 1^2)$  with  $\pi_1 = \pi_2 = \frac{1 - \pi}{2}$ ,  $\pi_3 = \pi_4 = \frac{\pi}{2}$

- $p(x | z; \theta) = \text{Bernoulli}(x; \text{sigmoid}(z))$

- $\theta = \{\pi, \mu_1, \mu_2, \mu_3, \mu_4\}$

- Variational/proposal distribution family

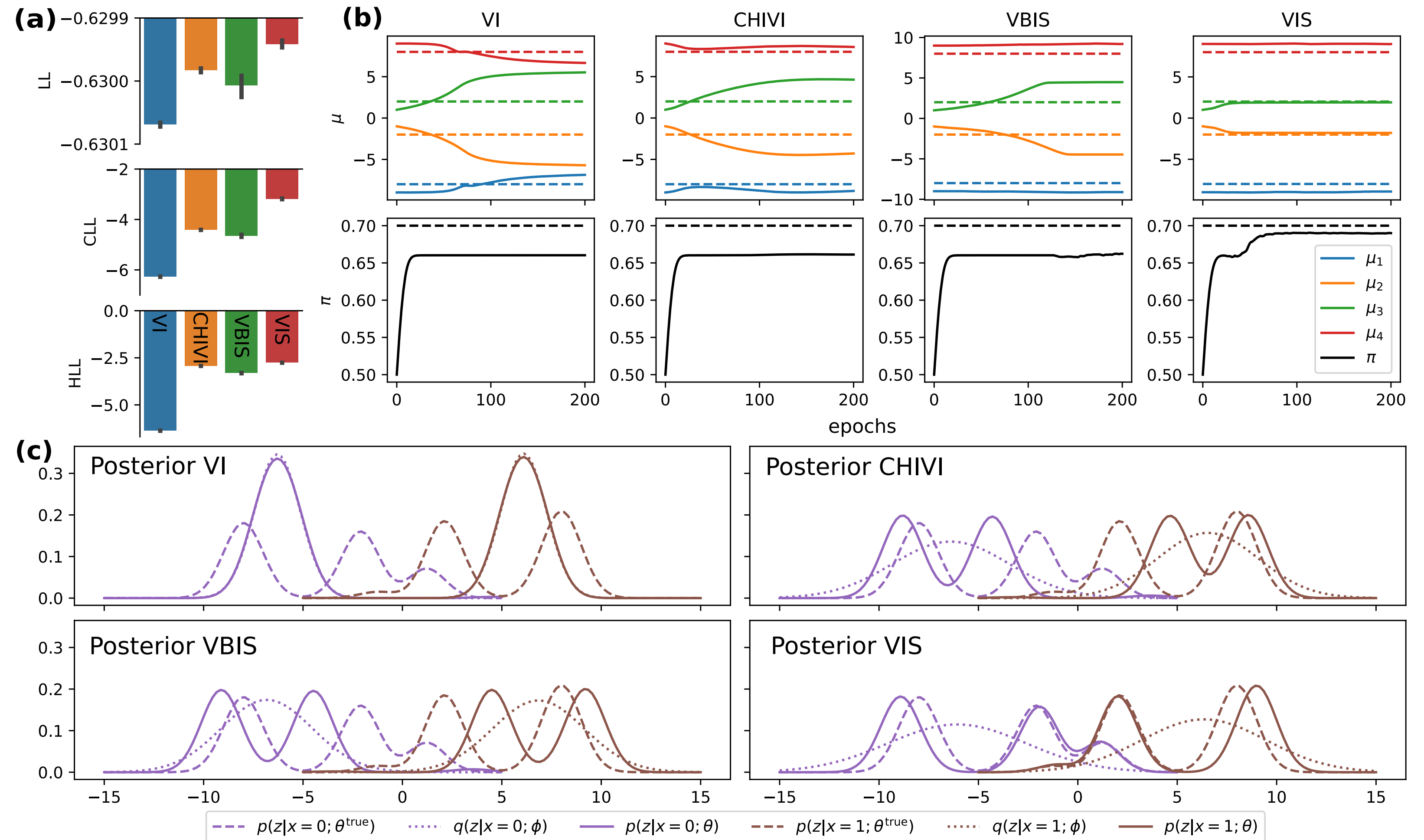
- $q(z | x; \phi) = \mathcal{N}(z; c_x, \sigma_x^2)$  for  $x \in \{0, 1\}$

- $\phi = \{c_0, c_1, \sigma_0, \sigma_1\}$



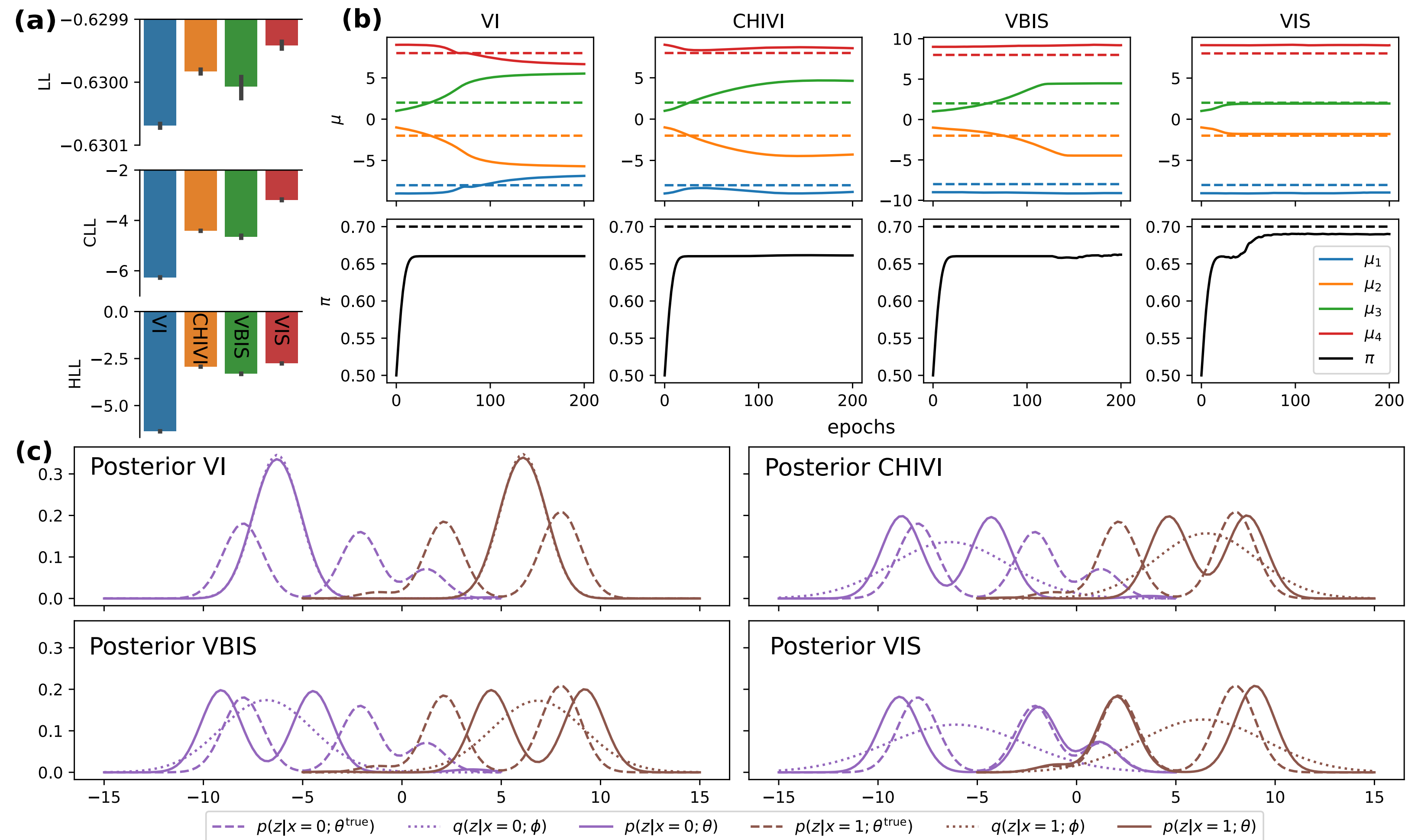
# Mixture model: GMM-Bernoulli

- (a) Quantitative comparison on the test dataset (e.g., test marginal log-likelihood).
- (b) parameter recovery.
- (c) Visualization of the true posterior (dashed), learned posterior (solid), and approximated posterior (dotted).



# 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  $p(z|x;\theta)$  and  $q(z|x;\phi)$  are far from  $p(z|x;\theta^{\text{true}})$ .
- VIS: mass-covering/mean-seeking behavior of minimizing the forward  $\chi^2$ . This enlarges the effective support range of  $q(z|x;\phi)$  for sampling.

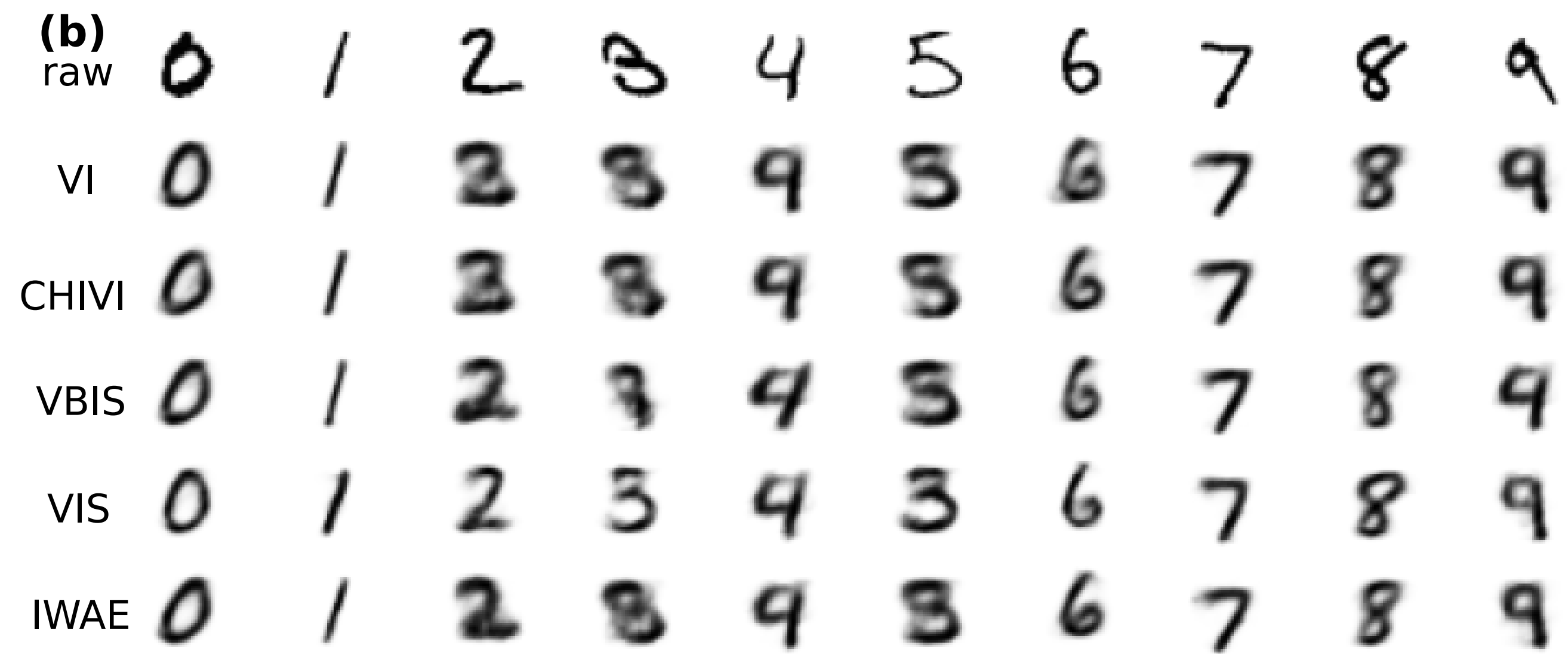
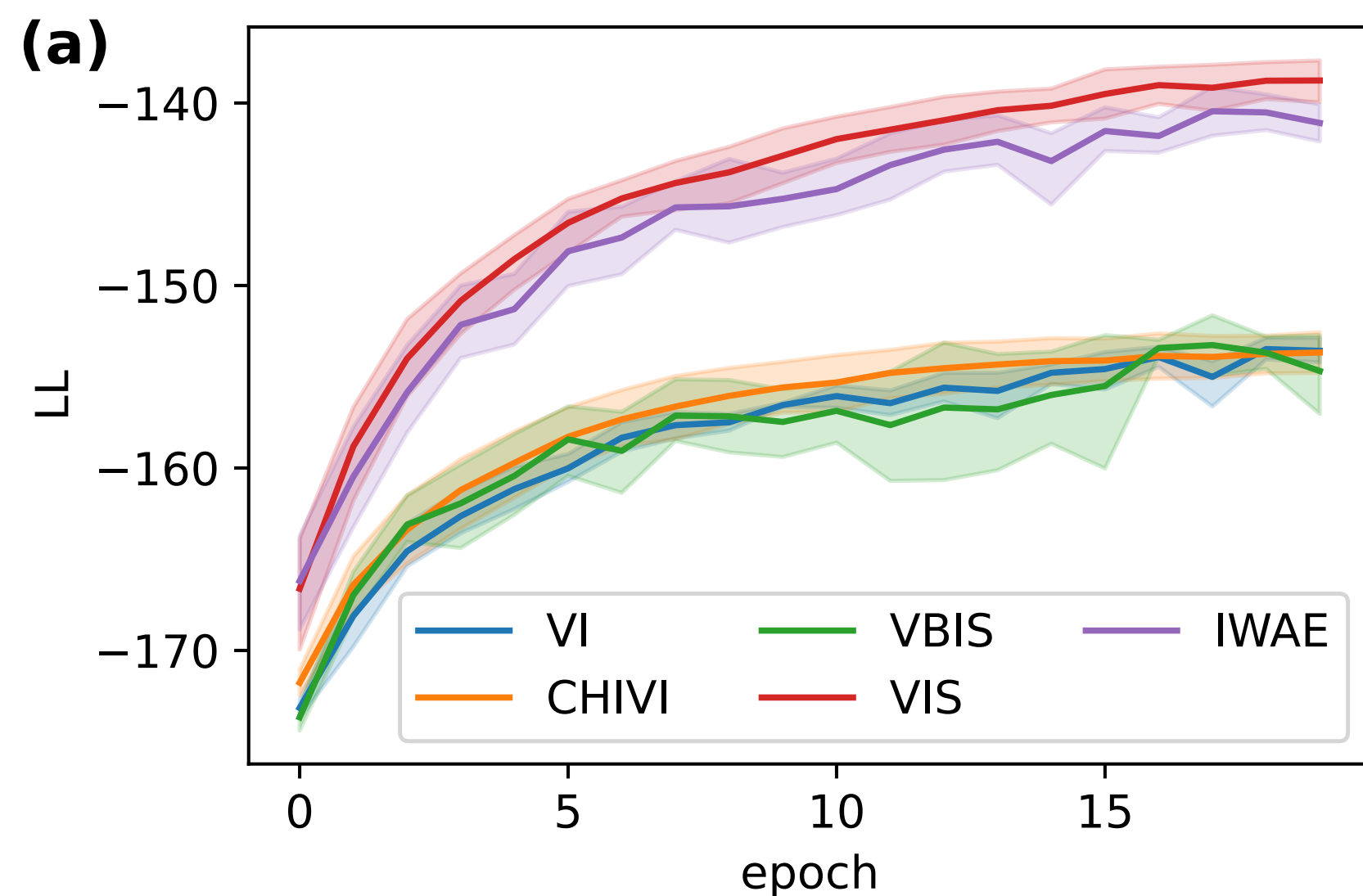


# VAE on MNIST

- Model
  - $p(\mathbf{z}; \theta) = \mathcal{N}(\mathbf{z}; \mathbf{0}, I)$
  - $p(\mathbf{x} | \mathbf{z}; \theta) = \text{Bernoulli}(\mathbf{x}; \text{sigmoid}(\text{decoder}(\mathbf{z})))$
  - $\theta$ : decoder's parameters
- Variational/proposal distribution family
  - $q(\mathbf{z} | \mathbf{x}; \phi) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}(\mathbf{x}), \text{diag } \boldsymbol{\sigma}^2(\mathbf{x}))$ , where  $\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\sigma}^2(\mathbf{x}) = \text{encoder}(\mathbf{x})$
  - $\phi$ : encoder's parameters

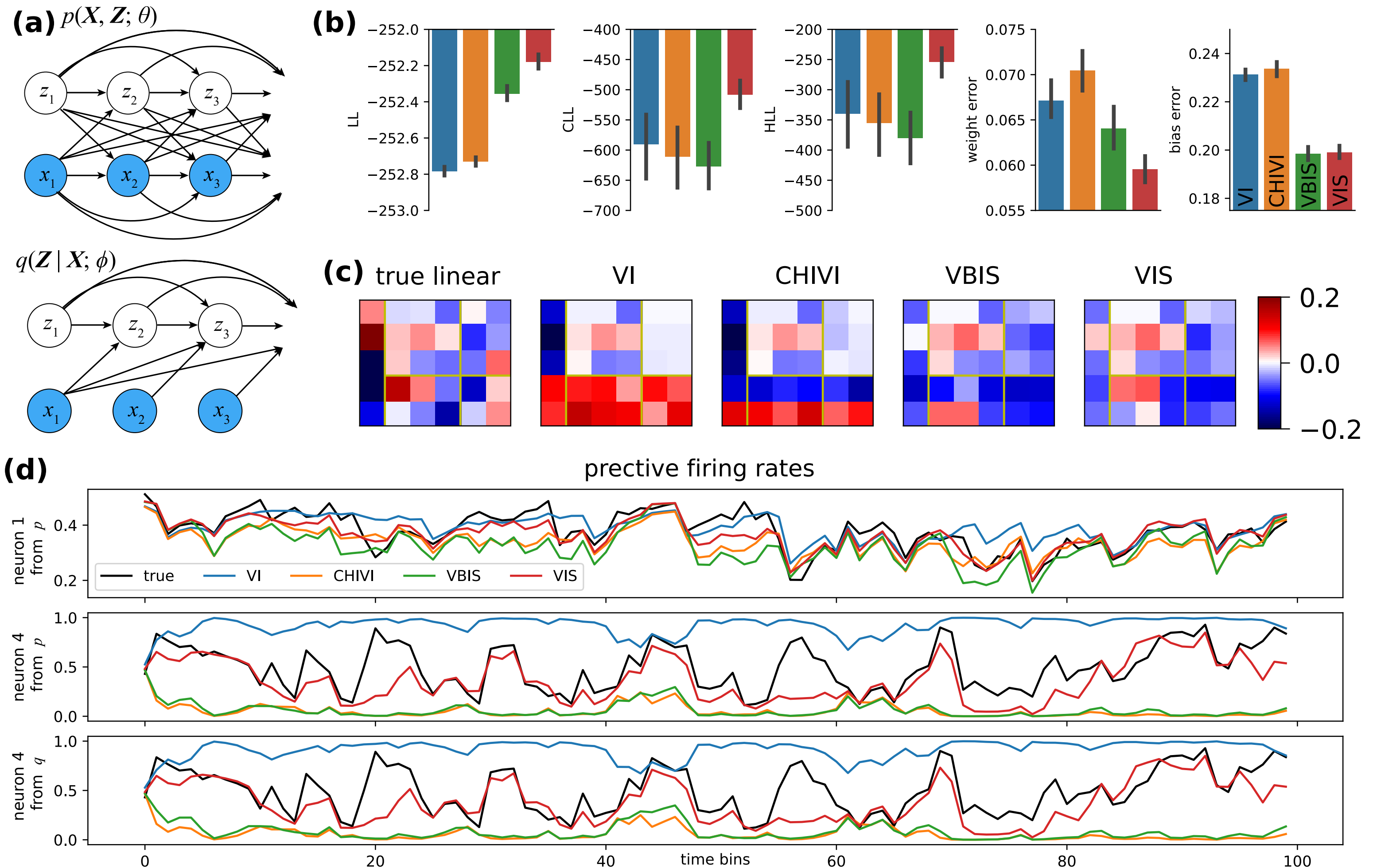
# VAE on MNIST

- (a) Convergence curve of the marginal log-likelihood on the test set.
- (b) Examples of raw images and the reconstructed images by different methods.



# Partially observable GLM (synthetic)

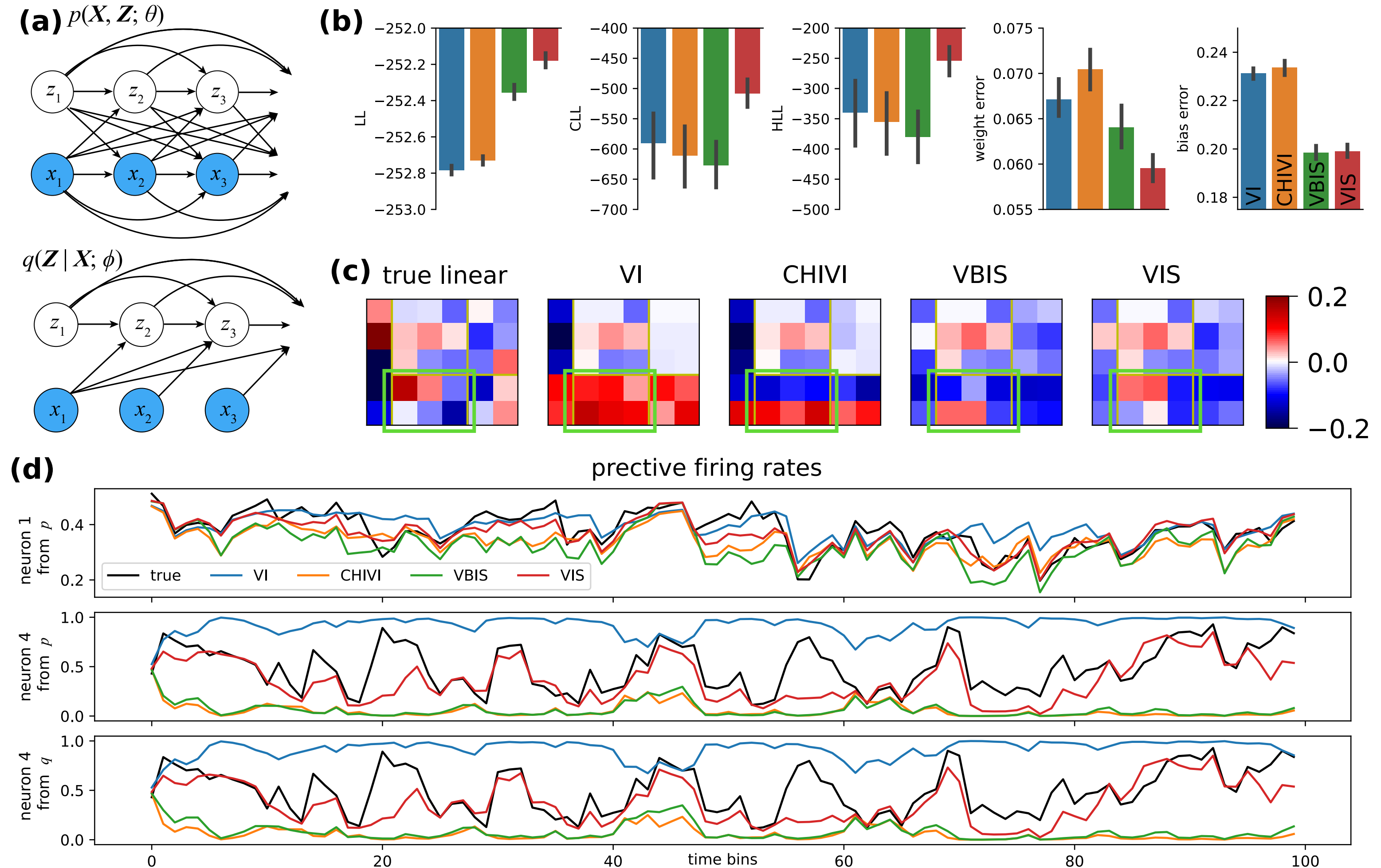
- This is a very hard problem, since  $p(\mathbf{x}, \mathbf{z}; \theta)$  cannot be explicitly factored as  $p(\mathbf{x} | \mathbf{z}; \theta)p(\mathbf{z}; \theta)$  (e.g., mixture model, HMM, PLDS, VAE, etc).
- $V$  of  $N$  neurons are visible and the remaining  $H$  neurons are hidden.
- $\mathbf{X}$  are spike trains from visible neurons, and  $\mathbf{Z}$  are spike trains from hidden neurons.
- $\theta = \{\mathbf{b} \in \mathbb{R}^N, \mathbf{W} \in \mathbb{R}^{N \times N}\}$ , where  $w_{n \leftarrow n'}$  represents the influence from neuron  $n'$  to neuron  $n$ .





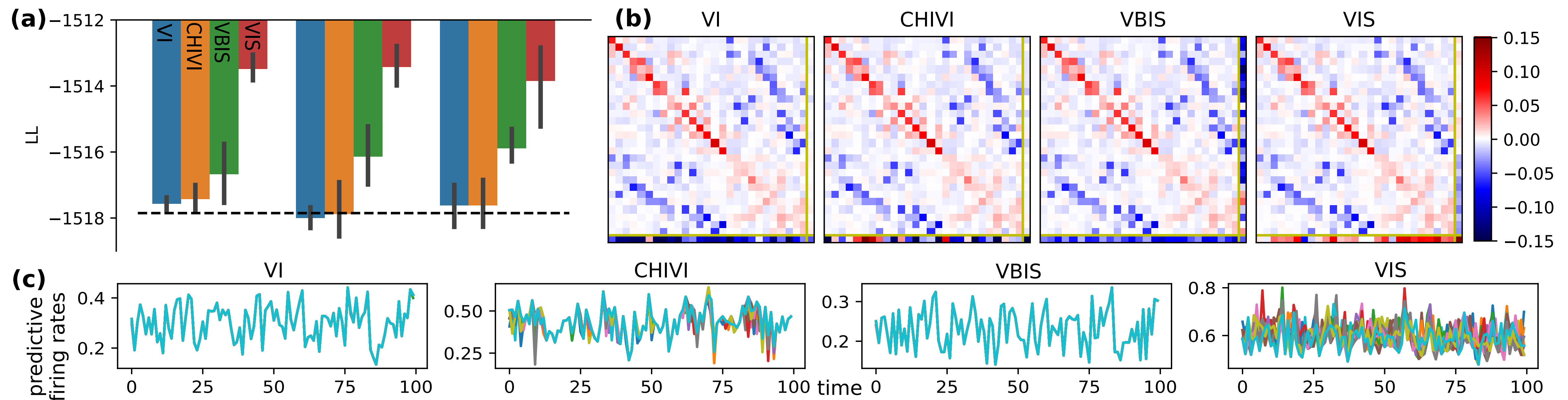
# Partially observable GLM (synthetic)

- (b) Quantitative comparison on the test dataset; weight and bias error.
- (c) The estimated weight  $W$  and the bias  $b$  by different methods compared with the true. The visible-to-hidden block learned by VIS is significantly better than others.
- (d) Predictive firing rates from different methods compared with the true. VIS matches the true the best.



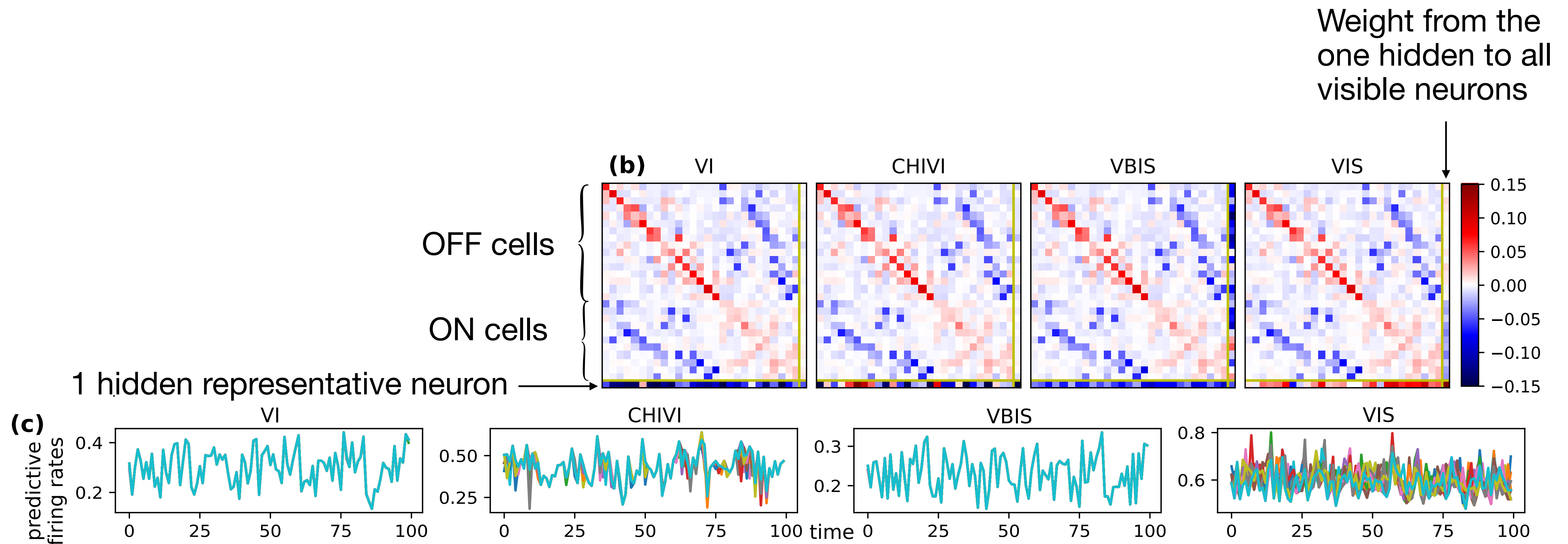
# Partially observable GLM (retinal ganglion cell)

- $V = 27$  basal ganglion neurons are recorded while a mouse is performing a visual test for about 20 mins (Pillow & Scott, 2012). Neuron 1-16 are OFF cells, neuron 17-27 are ON cells.
- (a) Quantitative comparison on the test dataset, with different numbers of hidden neurons  $H \in \{1, 2, 3\}$ .



# Partially observable GLM (retinal ganglion cell)

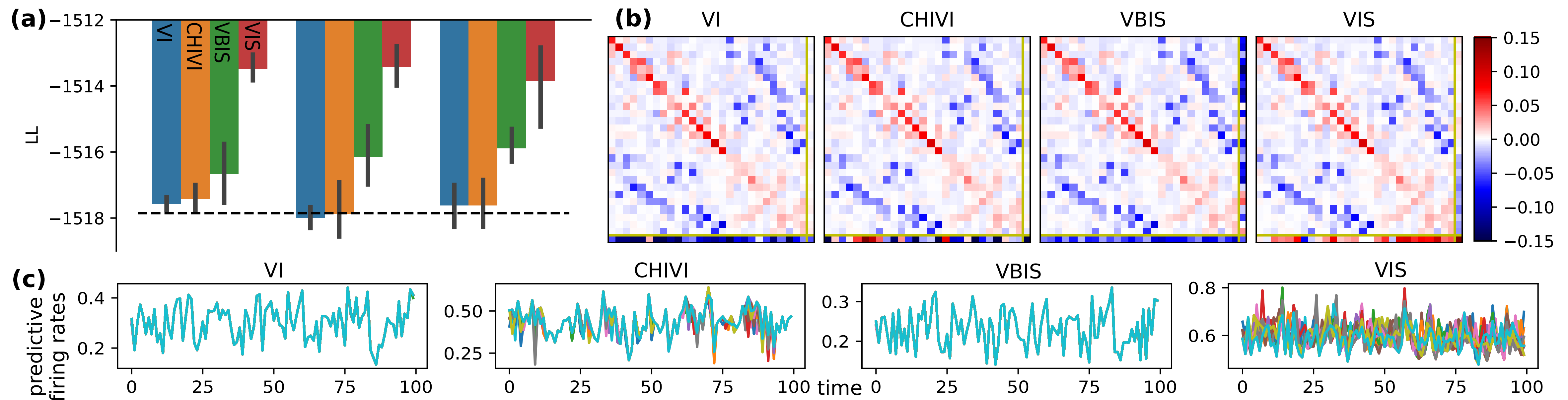
- (b) The hidden representative neuron learned by VIS behaves like an OFF cell. The sign of the weight from this hidden neuron to the visible neurons clearly tells us the type of those visible post-synaptic neurons.





# Partially observable GLM (retinal ganglion cell)

- (c) 20 randomly sampled predictive firing rates from  $q(\mathbf{Z} | \mathbf{X}; \phi)$ . The  $q(\mathbf{Z} | \mathbf{X}; \phi)$  learned by VIS provides the largest variability, which further improves the effectiveness of learning  $\theta$ .



# Summary

- VIS is the best way of doing importance sampling (IS), from the perspective of statistics.
- VIS has succinct and numerically stable gradient estimator derived in log space.
- By only changing two lines in the code, VIS learns significantly better model parameters, and achieves better test performance.

**Thanks for listening!**