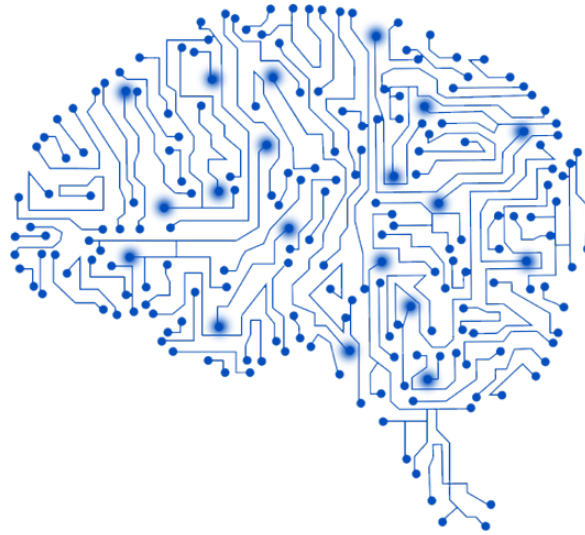# Improving Variational Inference for Complex Probabilistic Modeling
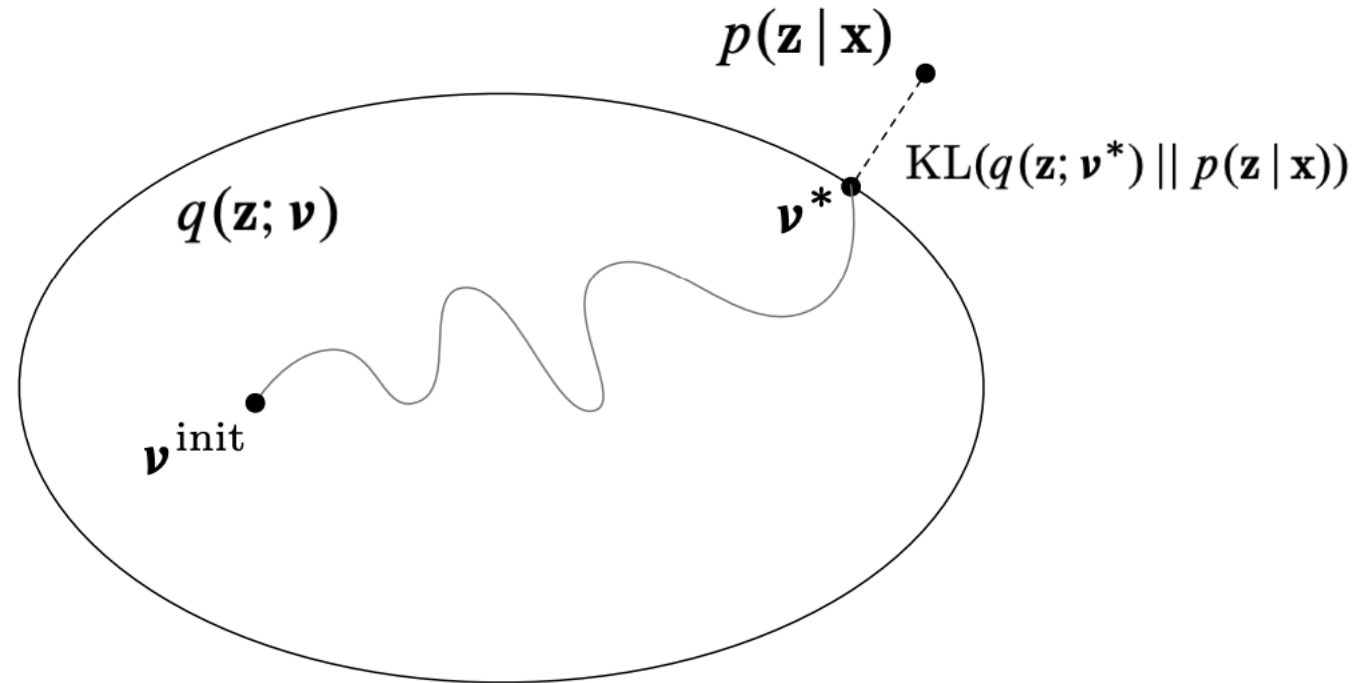
## Anqi Wu

**School of Computational Science and Engineering**

**Georgia Tech**

# Probabilistic Machine Learning

- A probabilistic model is a joint distribution of hidden variables z and observed variables **x, p(z, x)**.

- Inference about the unknowns is through the **posterior**, the conditional distribution of the hidden variables given the observations **p(z|x) = p(z, x)/p(x)**.

- For most interesting models, the denominator is not tractable. We appeal to **approximate posterior inference**.

# Variational Inference



- Variational inference turns **inference into optimization**.

- Posit a **variational family** of distributions over the latent variables, **q(z; v)**

- Fit the **variational parameters v** to be close (in KL) to the exact posterior. (There are alternative divergences, which connect to algorithms like EP, BP, and others.)

# Variational Inference

- Assume q(z; v) is an approximate posterior distribution (mean-field Gaussian)

$$v* = \text{argmin}_v \ D_{KL}[q(z;v)||p(z|x)]$$

Where

$$D_{KL}[q(z;v)||p(z|x)] = \mathbb{E}_q[\log \frac{q(z;v)}{p(z|x)}] = -\left(\underbrace{\mathbb{E}_{z \sim q}[\log p(x|z)] - D_{KL}(q(z;v)||p(z))}_{\substack{\text{ELBO} \\ \text{(evidence lower bound)}}}\right) + \underbrace{\log p(x)}_{\text{constant}}$$

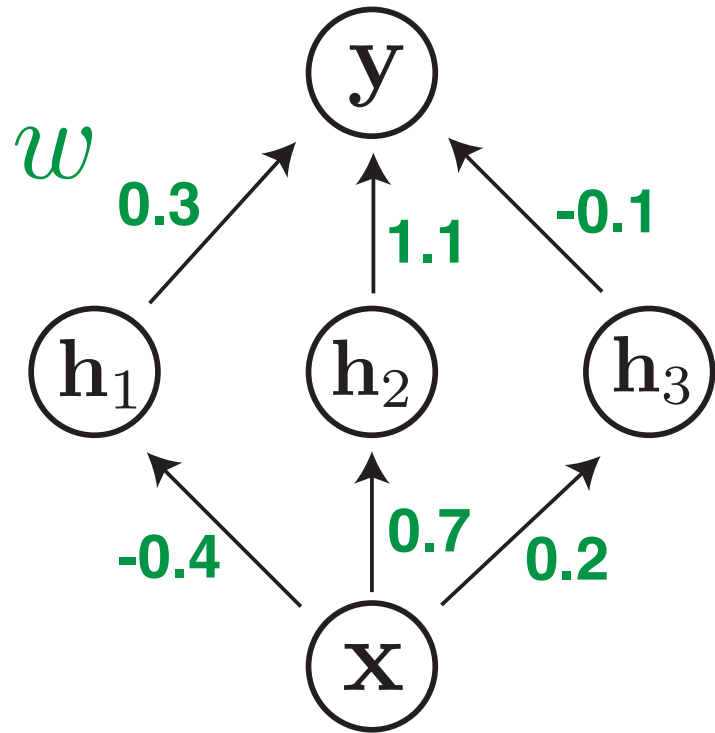- Thus, minimizing the KL is equivalent to maximizing the ELBO.

$$v* = \text{argmax}_v \ \mathbb{E}_{z \sim q}[\log p(x|z)] - D_{KL}(q(z;v)||p(z))$$

Still intractable!

# Outline

- **Determinstic variational inference** for Bayesian neural networks
  - Eliminate gradient variance in evaluating the expectation term
  - Empirical Bayes to avoid the prior tuning (*general approach*)

- **Variational importance sampling** for partially observed multivariate Hawkes process
  - VIS provides a tighter bound than ELBO (*general approach*)
  - Novel forward-backward approximate distribution

# Standard Neural Network



☺ Flexible class of models
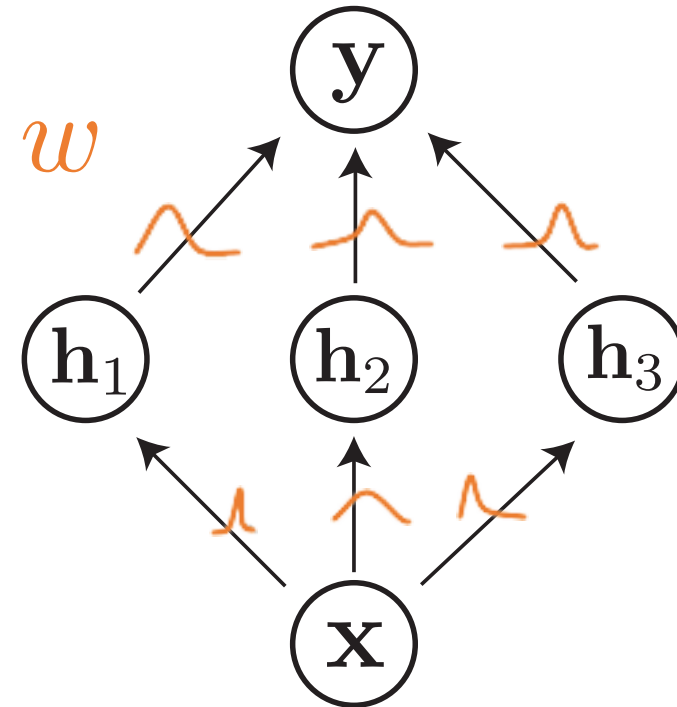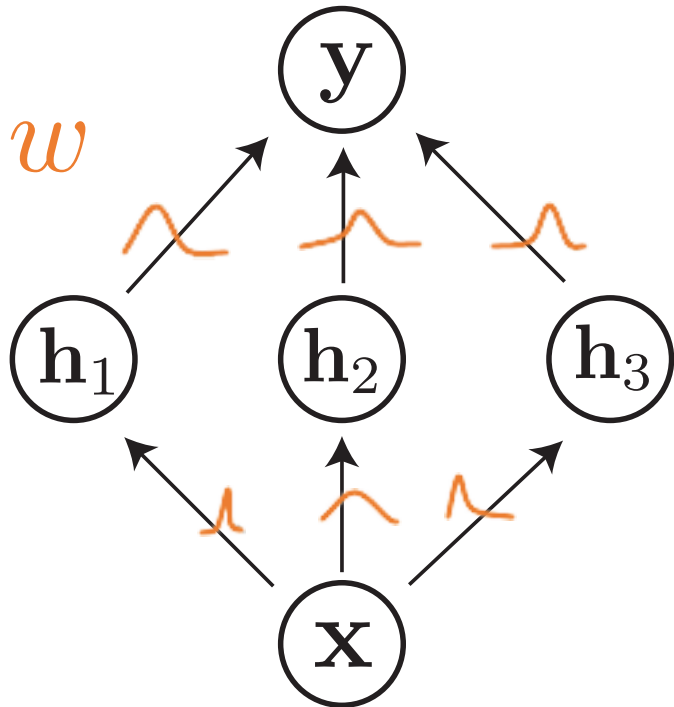
# Bayesian Neural Network



Image credit: Blundell et al., 2015

☺ Flexible class of models

☺ Principled handling of uncertainty

☺ Principled handling of regularization

# Bayesian Neural Network



**Goal**  The posterior distribution of $w$ is $p(w|x,y)$.

**Solution**  Variational Inference

variational approximate posterior    $q_\theta(w) \sim p(w|x,y)$

**ELBO (evidence lower bound)**

$$\max_\theta \; \mathbb{E}_{q_\theta(w)} \left[ \log p(y|x,w) \right] - D_{KL} \left[ q_\theta(w) \| p(w) \right]$$

Fit the data

Don't stray far
from the prior

challenge I

challenge II

# Challenge I: Gradient Variance



$$\mathbb{E}_{q_\theta(w)} \left[ \log p(y|x, w) \right] - D_{KL} \left[ q_\theta(w) \| p(w) \right]$$

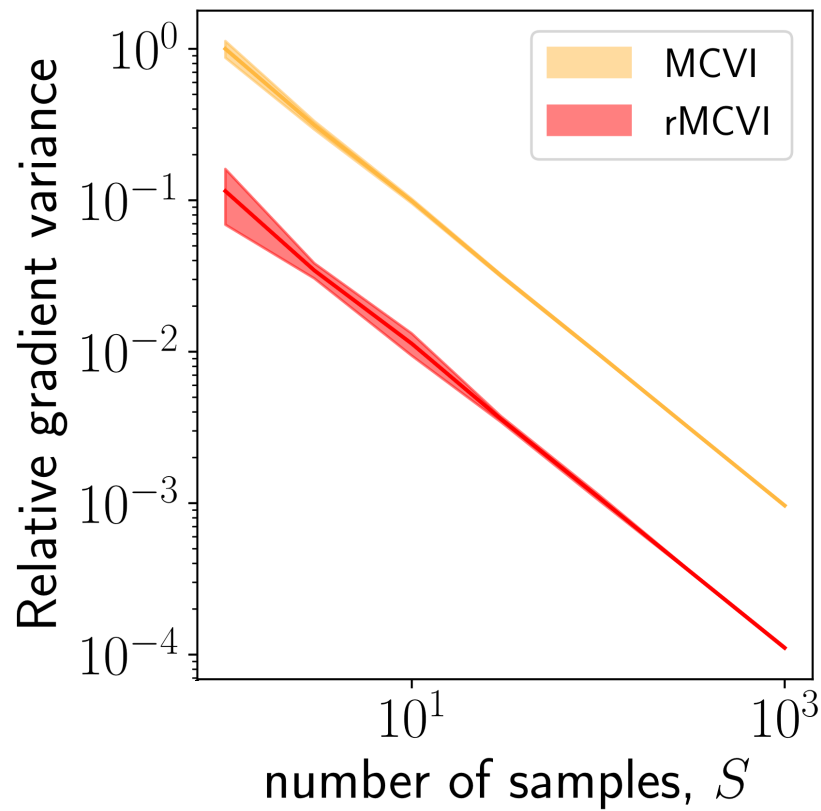Fit the data

**Monte Carlo sampling**

gradient variance

# Challenge I: Gradient Variance

**MCVI: Monte Carlo Variational Inference**

**DVI: Deterministic Variational Inference**

**reduce gradient variance**

**local reparameterization trick**



*nov et al., 2017*

**ient estimators**

**ximation instead**

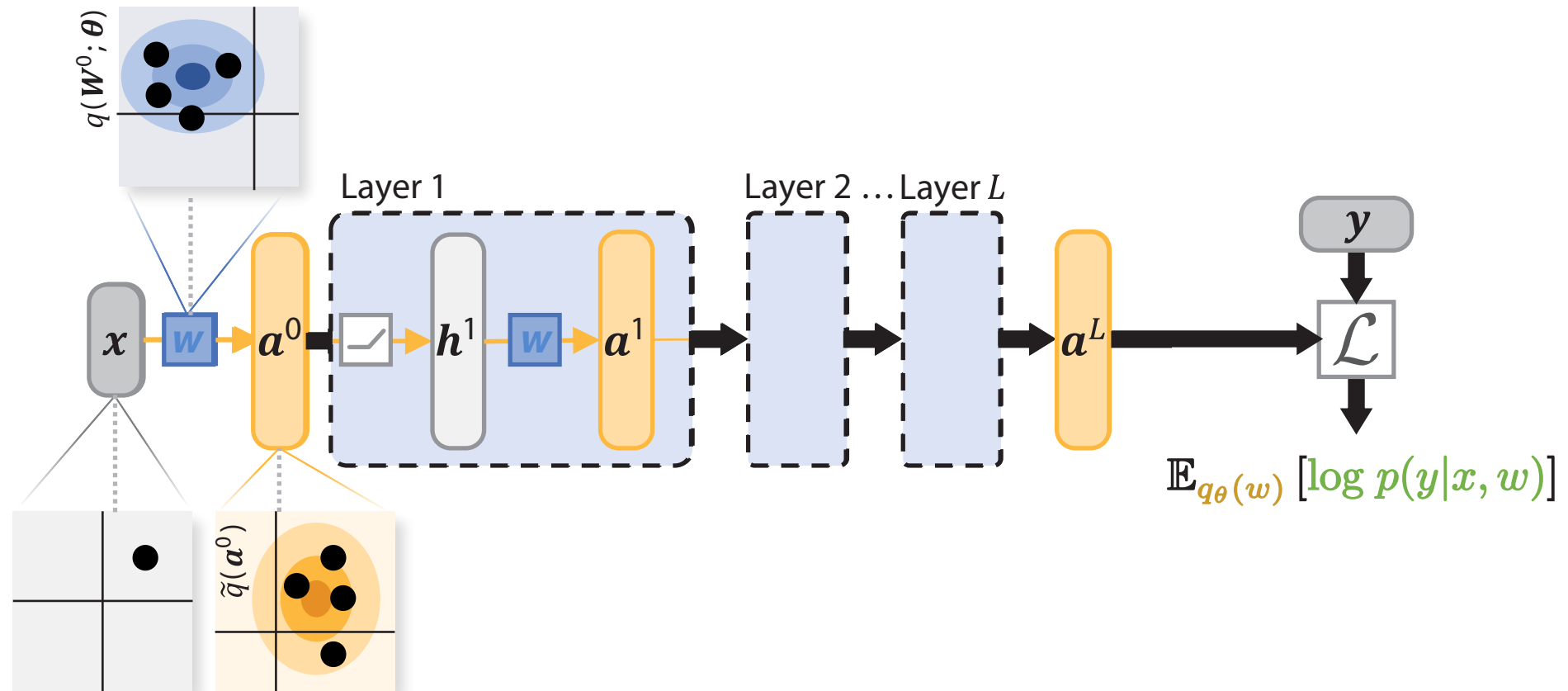**of MC, thus no gradient variance**

# Monte Carlo Approximation for ELBO

$$\mathbb{E}_{q_\theta(w)}\left[\log p(y|x,w)\right] \approx \frac{1}{S}\sum_{s=1}^{S} \log p(y|w^{(s)},x), \quad w^{(s)} \sim q_\theta(w)$$
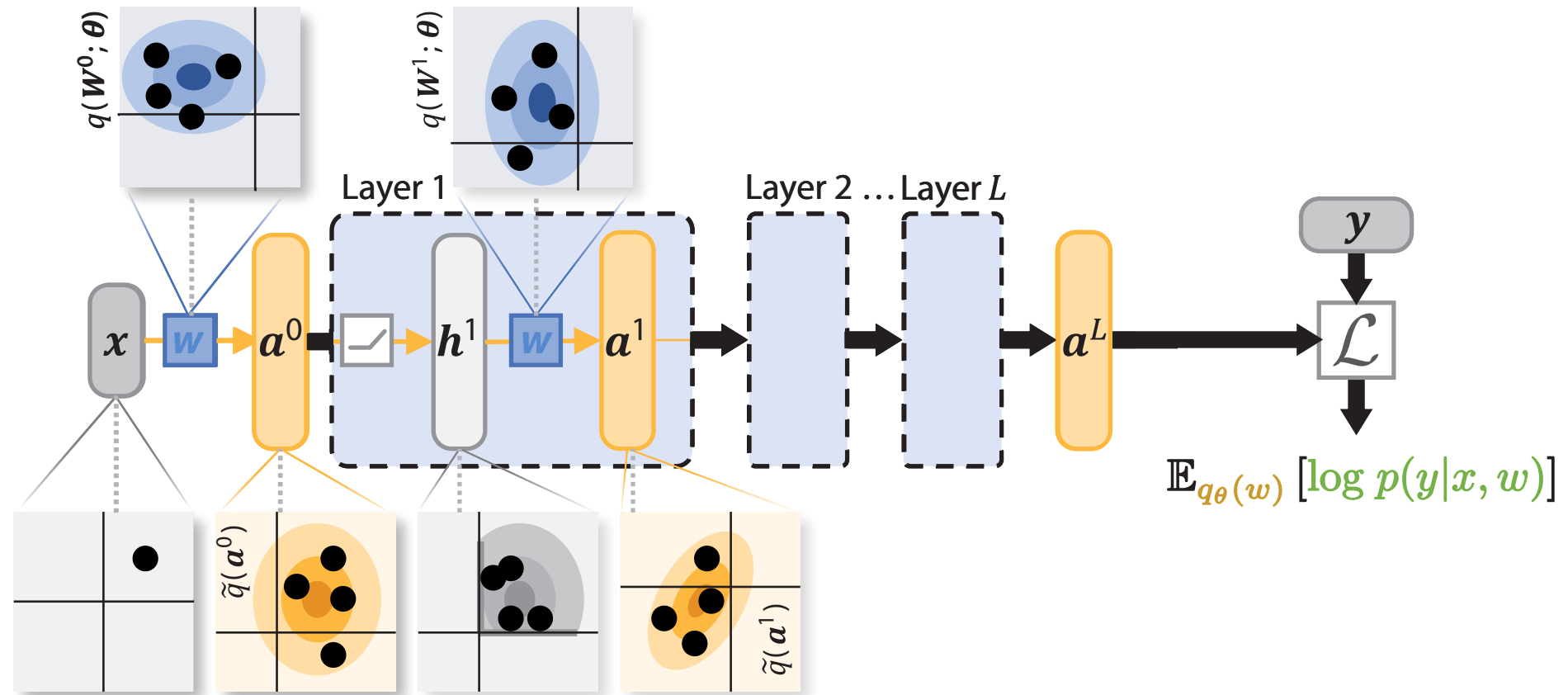
# Monte Carlo Approximation for ELBO

$$\mathbb{E}_{q_\theta(w)}\left[\log p(y|x,w)\right] \approx \frac{1}{S}\sum_{s=1}^{S}\log p(y|w^{(s)},x), \quad w^{(s)} \sim q_\theta(w)$$

# Monte Carlo Approximation for ELBO

$$\mathbb{E}_{q_{\theta}(w)}\left[\log p(y|x,w)\right] \approx \frac{1}{S}\sum_{s=1}^{S}\log p(y|w^{(s)},x), \quad w^{(s)} \sim q_{\theta}(w)$$
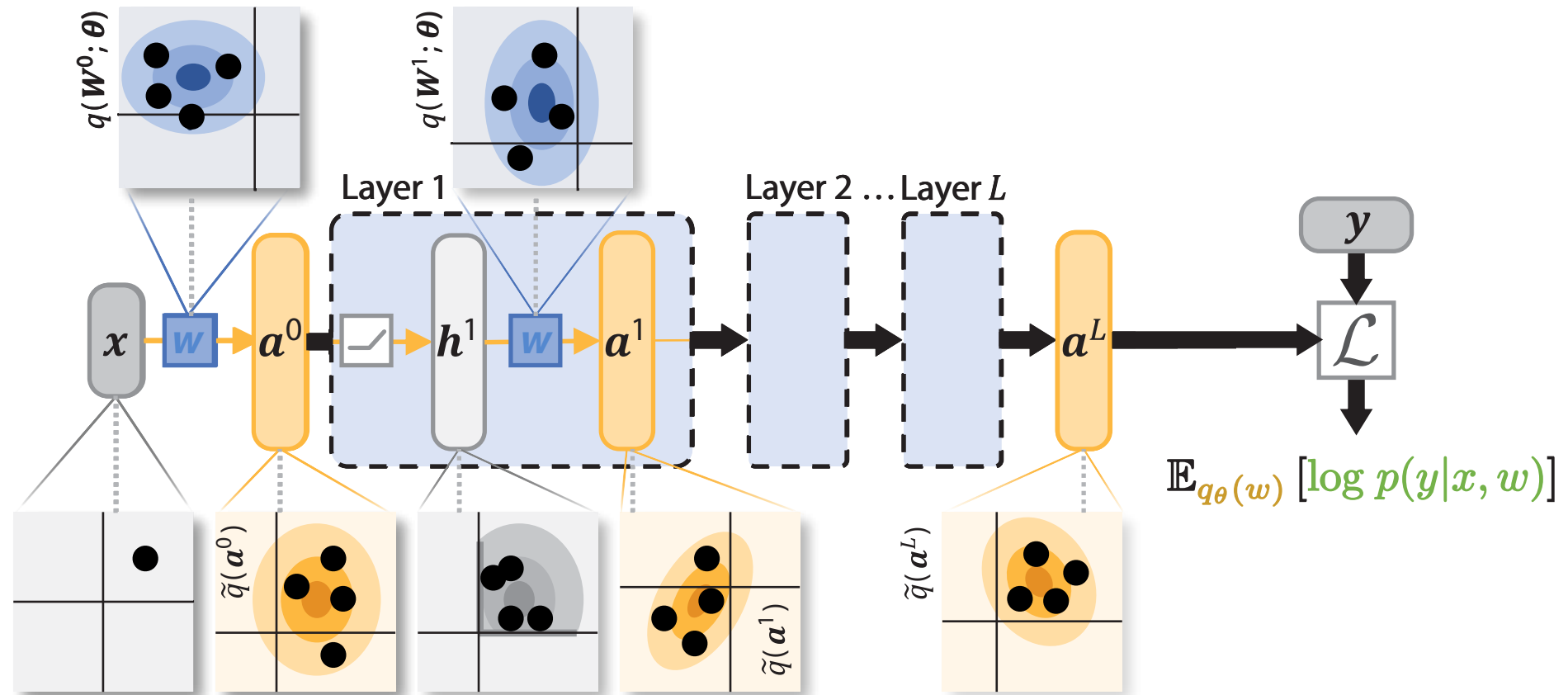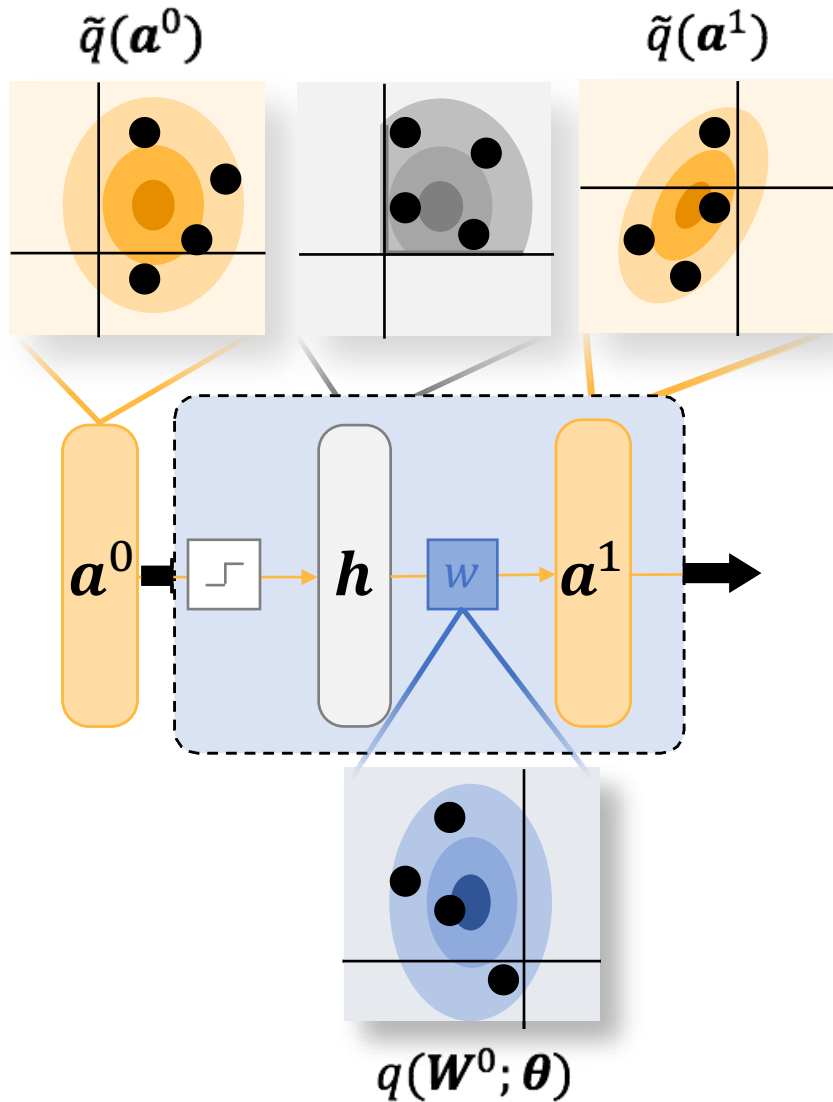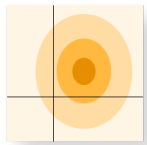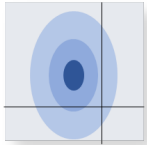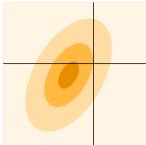
# Monte Carlo Approximation for ELBO

$$\mathbb{E}_{q_{\theta}(w)}\left[\log p(y|x,w)\right] \approx \frac{1}{S}\sum_{s=1}^{S}\log p(y|w^{(s)},x), \quad w^{(s)} \sim q_{\theta}(w)$$

# Challenge I: Deterministic Propagation of Uncertainties



$\tilde{q}(\boldsymbol{a}^0)$     $\tilde{q}(\boldsymbol{a}^1)$

$q(\boldsymbol{W}^0; \boldsymbol{\theta})$

Instead of propagating uncertainties via **samples**.

We can **deterministically** propagate **distributions**.

$$\tilde{q}(\boldsymbol{a}^0) \quad q(\boldsymbol{W}^0; \boldsymbol{\theta}) \quad \tilde{q}(\boldsymbol{a}^1)$$

bnn.activation_layer(   ) = 

$$a_i^1 = \sum_{i=1}^{d} w_{ij} h_j \ \sim \mathcal{N}(\boldsymbol{\mu^1}, \boldsymbol{\Sigma^1})$$

Gaussian

Central Limit Theorem:
1. $\boldsymbol{W}$ and $\boldsymbol{h}$ are i.i.d. samples
2. Large number of hidden nodes in $\boldsymbol{h}$

# Challenge I: Deterministic Propagation of Uncertainties



$\tilde{q}(\boldsymbol{a}^0)$

$\tilde{q}(\boldsymbol{a}^1)$

$\boldsymbol{a}^0$ ⟶ ⟶ $\boldsymbol{h}$ ⟶ $w$ ⟶ $\boldsymbol{a}^1$ ⟶

$q(\boldsymbol{W}^0; \boldsymbol{\theta})$

Instead of propagating uncertainties via **samples**.

We can **deterministically** propagate **distributions**.

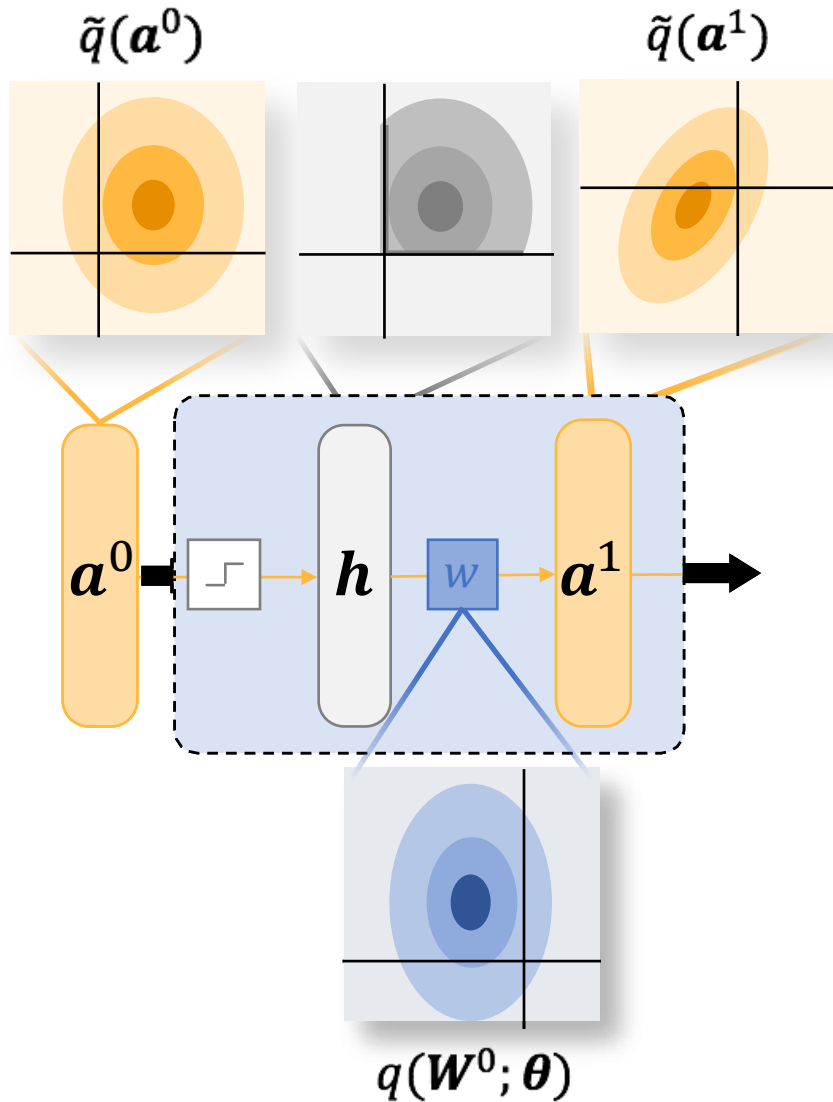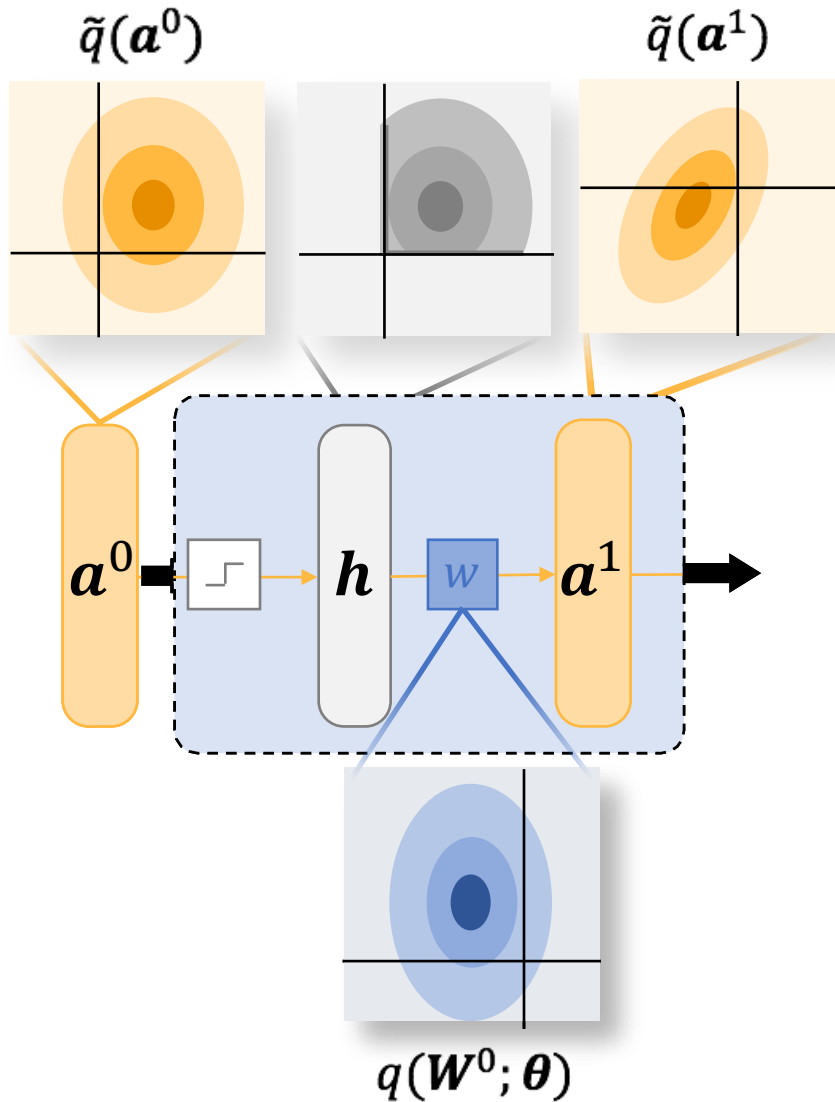$\tilde{q}(\boldsymbol{a}^0)$  $q(\boldsymbol{W}^0; \boldsymbol{\theta})$  $\tilde{q}(\boldsymbol{a}^1)$

bnn.heaviside_layer(   ) = 

Example: 2-dimensional $\boldsymbol{a}^1$

$$w_{ij} \sim \mathcal{N}^{2 \times d}$$

$$\boldsymbol{h} \sim \text{truncated}\mathcal{N}^d$$
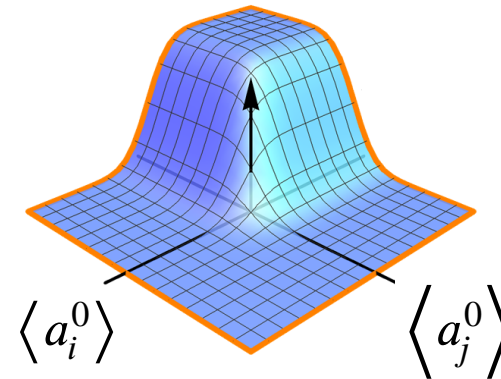
# Challenge I: Deterministic Propagation of Uncertainties



$$a^1 \sim \mathcal{N}(\boldsymbol{\mu}^1, \boldsymbol{\Sigma}^1)$$

Just need moments: $\langle h_i \rangle, \langle h_i h_j \rangle$

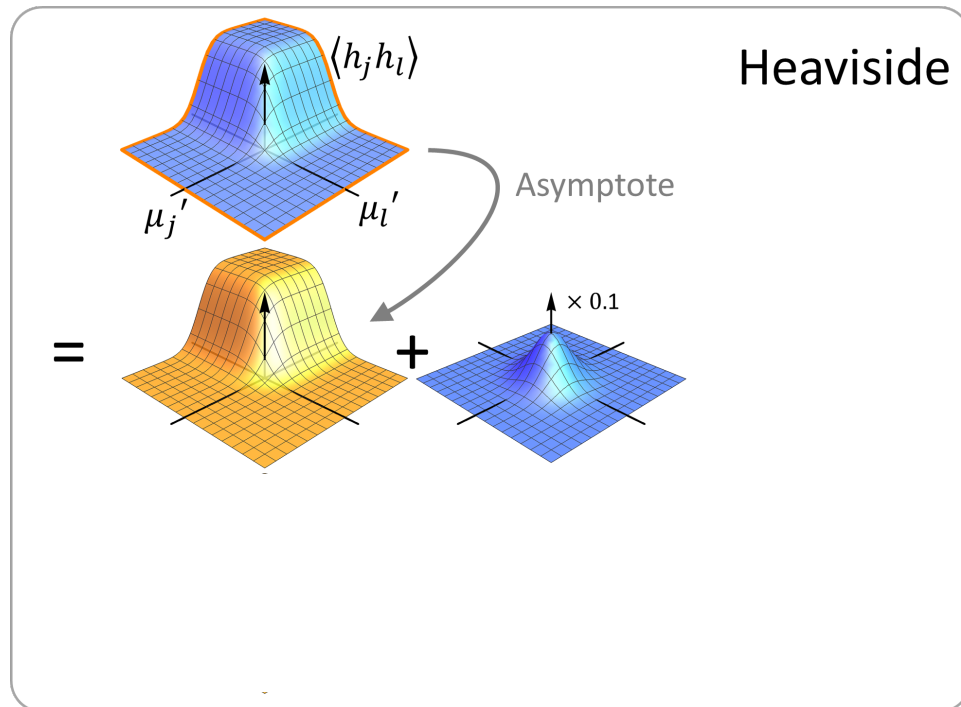$$\langle h_i \rangle = \mathbb{E}_{a^0 \sim \mathcal{N}(\boldsymbol{\mu}^0, \boldsymbol{\Sigma}^0)}\Big[f(a_i^0)\Big] = \int f(\alpha)\phi\left(\frac{\alpha - \langle a_i^0 \rangle}{\Sigma_{ii}^0}\right)d\alpha$$
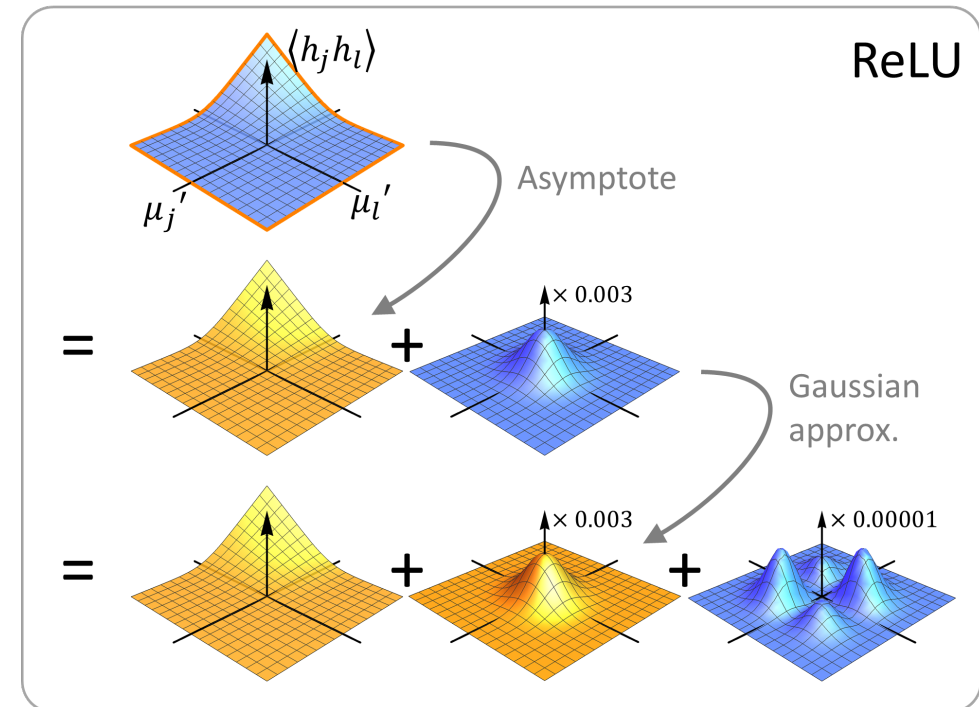
$$\langle h_i h_j \rangle =$$

# Challenge I: Deterministic Propagation of Uncertainties



```python
def heaviside_covariance(x):
    mu1 = tf.expand_dims(mu, 2)
    mu2 = tf.transpose(mu1, [0,2,1])

    s11s22 = tf.expand_dims(x_var_diag, axis=2) * tf.expand_dims(x_var_diag, axis=1)
    rho = x.var / (tf.sqrt(s11s22))# + EPSILON)
    rho = tf.clip_by_value(rho, -1/(1+EPSILON), 1/(1+EPSILON))

    return bu.heavy_g(rho, mu1, mu2)
```

```python
def relu_covariance(x):
    mu1 = tf.expand_dims(mu, 2)
    mu2 = tf.transpose(mu1, [0,2,1])

    s11s22 = tf.expand_dims(x_var_diag, axis=2) * tf.expand_dims(x_var_diag, axis=1)
    rho = x.var / (tf.sqrt(s11s22))# + EPSILON)
    rho = tf.clip_by_value(rho, -1/(1+EPSILON), 1/(1+EPSILON))

    return x.var * bu.delta(rho, mu1, mu2)
```
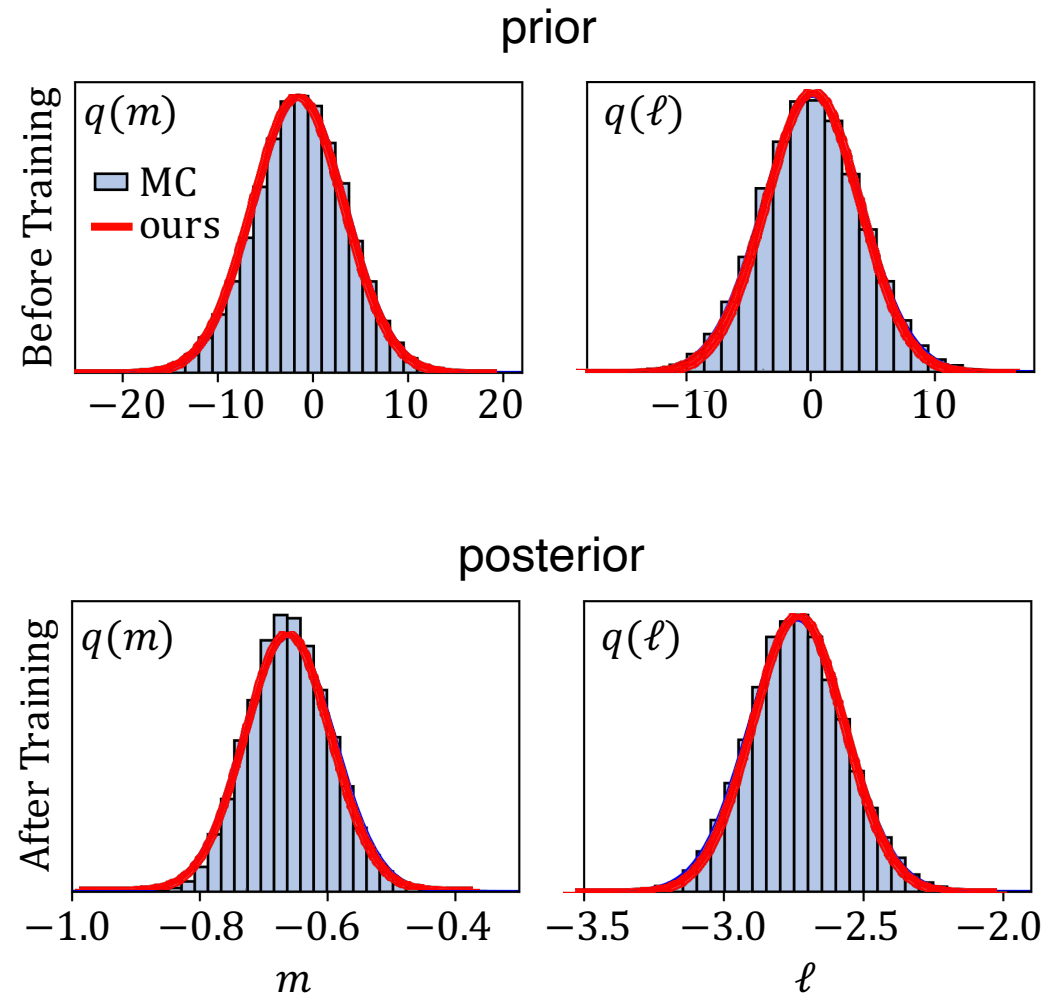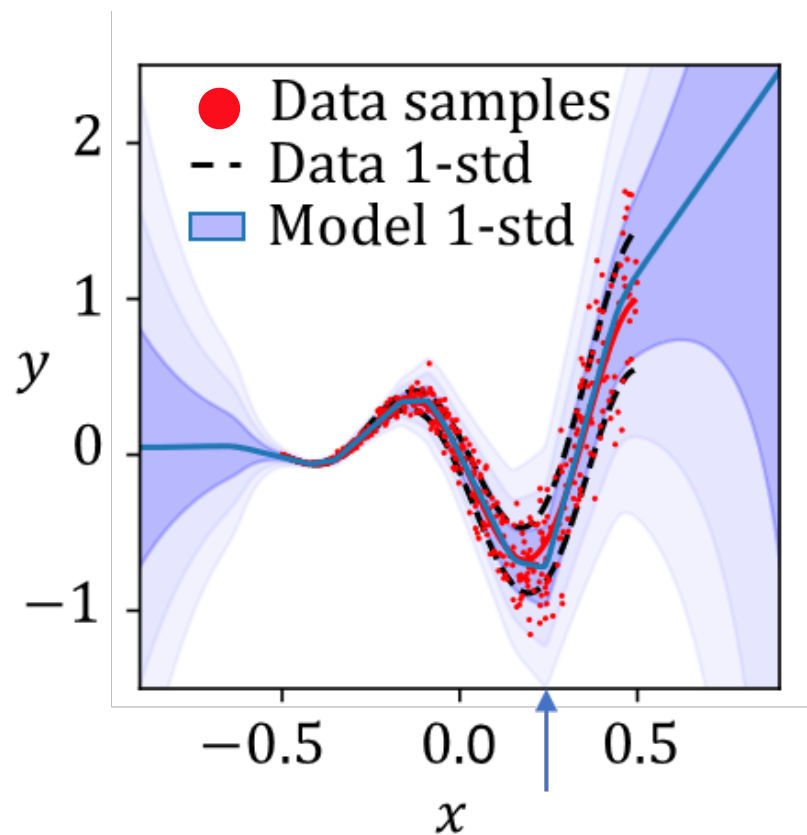
# Challenge II: Prior Tuning

$$\mathbb{E}_{q_\theta(w)}\left[\log p(y|x,w)\right] - D_{KL}\left[q_\theta(w)\|p(w)\right]$$

Fit the data

Don't stray far
from the prior

gradient
variance

prior tuning

solved

# Challenge II: Prior Tuning

UCI regression datasets

# Challenge II: Empirical Bayes for Prior Tuning

$$w \sim p(w|s) = \mathcal{N}(0, s)$$ prior variance

$$s \sim p(s) = \mathrm{InvGamma}(\alpha, \beta)$$ scale

shape

**Optimize ELBO** $$s^* = \underset{s}{\arg\max}\, ELBO(s, \theta) = \mathrm{function}(\theta)$$

variational parameter

**Empirical Bayes ELBO**

$$\mathbb{E}_{q_\theta(w)}\left[\log p(y|x, w)\right] - D_{KL}\left[q_\theta(w) \| p(w|s^*(\theta))\right]$$

# Challenge II: Empirical Verification

UCI regression datasets

# Deterministic VI + Empirical Bayes

**ELBO (evidence lower bound)**

$$\mathbb{E}_{q_\theta(w)} \left[ \log p(y|x, w) \right] - D_{KL} \left[ q_\theta(w) || p(w) \right]$$

Fit the data

Don't stray far from the prior

gradient variance

prior tuning

solved

solved

# Deterministic VI + Empirical Bayes

UCI regression datasets: test log likelihood

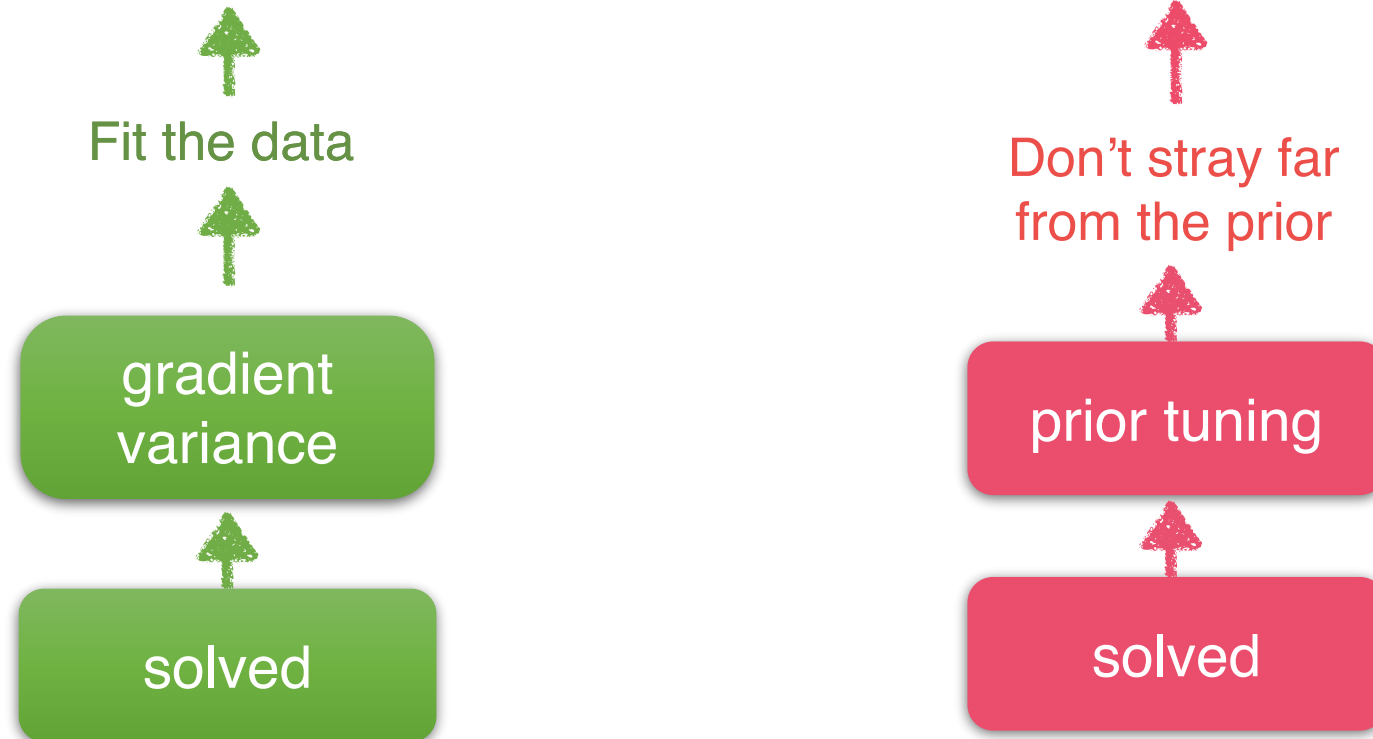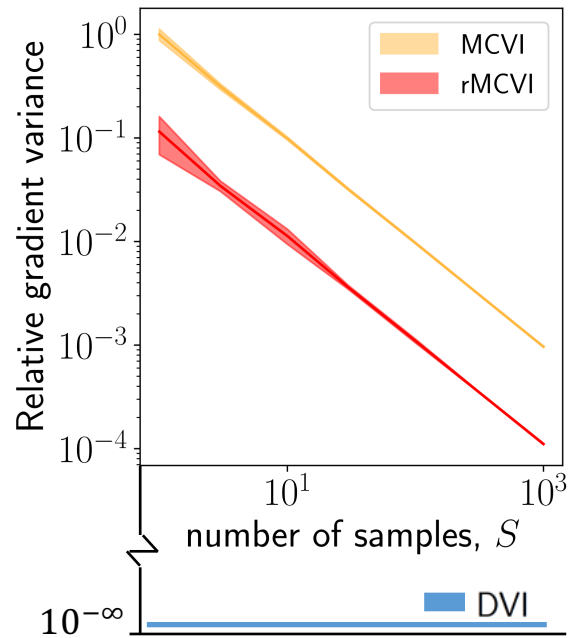| Dataset | $|\mathcal{D}|$ | $d_x$ | Deterministic+EB | Deterministic+fixed | MonteCarlo+EB | MonteCarlo+fixed |
|---------|------|------|------------------|---------------------|---------------|-------------------|
| bost | 506 | 13 | $\mathbf{-2.41 \pm 0.02}$ | $-2.46 \pm 0.02$ | $-2.46 \pm 0.02$ | $-2.48 \pm 0.02$ |
| conc | 1030 | 8 | $\mathbf{-3.06 \pm 0.01}$ | $-3.07 \pm 0.01$ | $-3.07 \pm 0.01$ | $-3.07 \pm 0.01$ |
| ener | 768 | 8 | $\mathbf{-1.01 \pm 0.06}$ | $-1.07 \pm 0.04$ | $-1.03 \pm 0.04$ | $-1.07 \pm 0.04$ |
| kin8 | 8192 | 8 | $1.13 \pm 0.00$ | $1.12 \pm 0.00$ | $\mathbf{1.14 \pm 0.00}$ | $1.13 \pm 0.00$ |
| nava | 11934 | 16 | $6.29 \pm 0.04$ | $\mathbf{6.32 \pm 0.04}$ | $5.94 \pm 0.05$ | $6.00 \pm 0.02$ |
| powe | 9568 | 4 | $\mathbf{-2.80 \pm 0.00}$ | $\mathbf{-2.80 \pm 0.01}$ | $\mathbf{-2.80 \pm 0.00}$ | $\mathbf{-2.80 \pm 0.00}$ |
| prot | 45730 | 9 | $-2.85 \pm 0.01$ | $\mathbf{-2.84 \pm 0.01}$ | $-2.87 \pm 0.01$ | $-2.89 \pm 0.01$ |
| wine | 1588 | 11 | $\mathbf{-0.90 \pm 0.01}$ | $-0.94 \pm 0.01$ | $-0.92 \pm 0.01$ | $-0.94 \pm 0.01$ |
| yach | 308 | 6 | $\mathbf{-0.47 \pm 0.03}$ | $-0.49 \pm 0.03$ | $-0.68 \pm 0.03$ | $-0.56 \pm 0.03$ |

# Deterministic:

**Eliminate Gradient variance**



# Robust:

**Less tuning required**



# Efficient:

**Just a few special function calls**

```python
def heaviside_covariance(x):
    mu1 = tf.expand_dims(mu, 2)
    mu2 = tf.transpose(mu1, [0,2,1])

    s11s22 = tf.expand_dims(x_var_diag, axis=2) * tf.expand_dims(x_var_diag, axis=1)
    rho = x.var / (tf.sqrt(s11s22))# + EPSILON)
    rho = tf.clip_by_value(rho, -1/(1+EPSILON), 1/(1+EPSILON))

    return bu.heavy_g(rho, mu1, mu2)
```
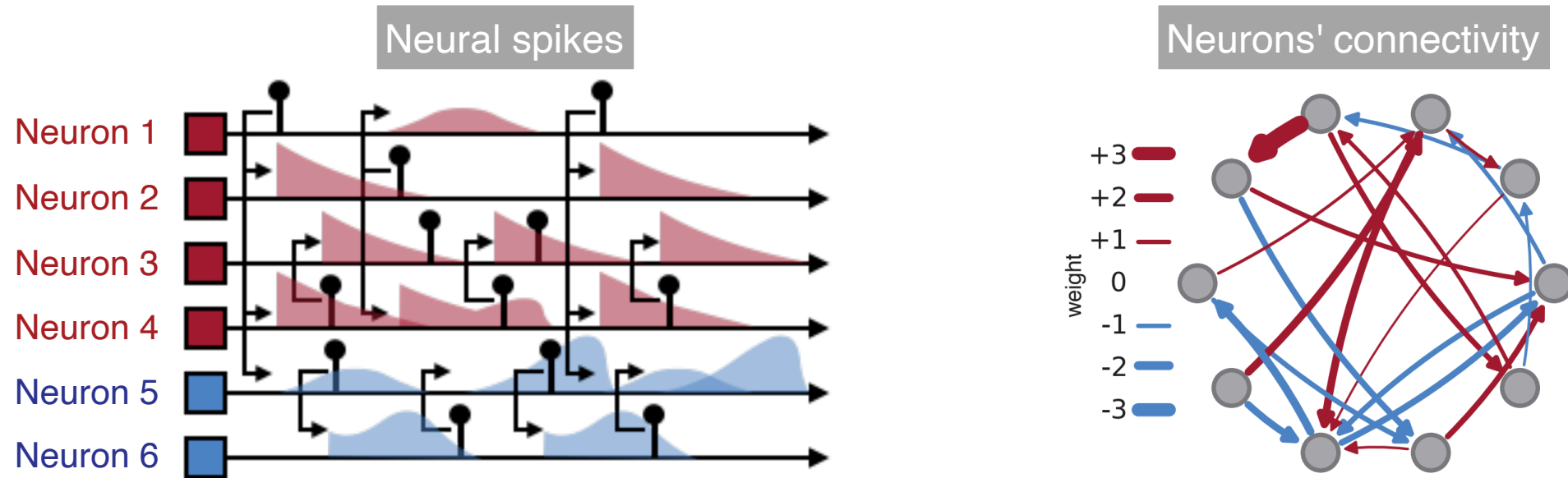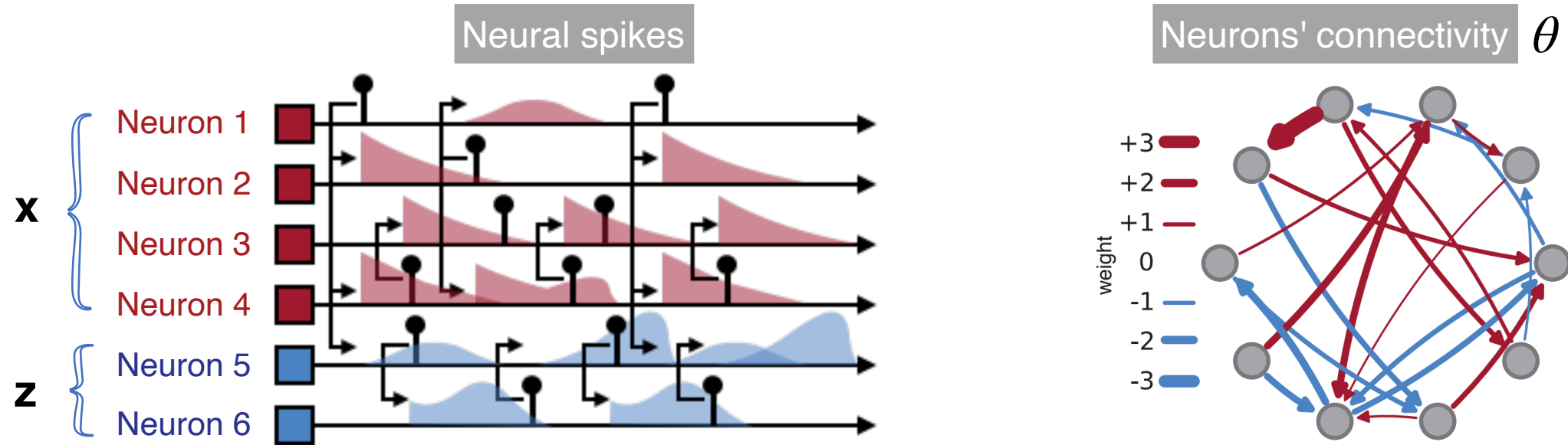
# Outline

- **Determinstic variational inference** for Bayesian neural networks
  - Eliminate gradient variance in evaluating the expectation term
  - Empirical Bayes to avoid the prior tuning (*general approach*)

- **Variational importance sampling** for partially observed multivariate Hawkes process
  - VIS provides a tighter bound than ELBO (*general approach*)
  - Novel forward-backward approximate distribution

# Partially observed multivariate Hawkes process (POMHP)



- Hawkes process is a self-exciting point process to describe neural spiking time.

- In a multivariate Hawkes process, each event can influence the occurrence of future events, **not just in the same dimension but also in other dimensions**.

- Partially observed means some events might be **hidden or unobserved**.

- Applications: finance, social networks, **neuroscience**, and so forth.

# Variational Inference



Neural spikes

Neurons' connectivity $\theta$

- Events from observed neurons, denoted as **x**.

- Events from unobserved neurons, denoted as **z**.

- Maximum likelihood estimation to maximize the marginal **p(x;** $\theta$ **)** with respect to the model parameters $\theta$ (such as **connectivity weights**).
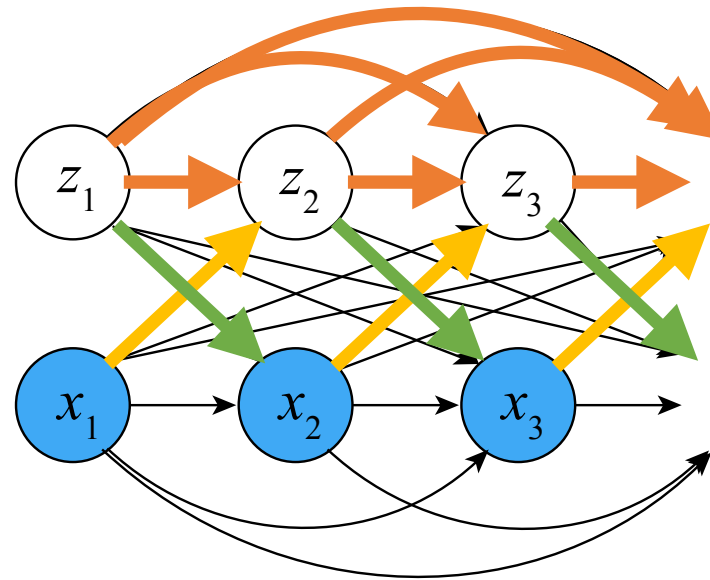
intractable!

# Variational Inference

- Maximize

$$\log p(x; \theta) = \log \int p(x, z; \theta) dz = \log \int q(z; \phi) \frac{p(x, z; \theta)}{q(z; \phi)} dz$$

*(Jensen's inequality)* $\quad \geq \mathbb{E}_{q(z;\phi)}[\log p(x, z; \theta) - \log q(z; \phi)]$

*(ELBO)* $\quad = \mathbb{E}_{z \sim q}[\log p(x \,|\, z, \theta)] - D_{KL}(q(z; \phi) \,||\, p(z))$
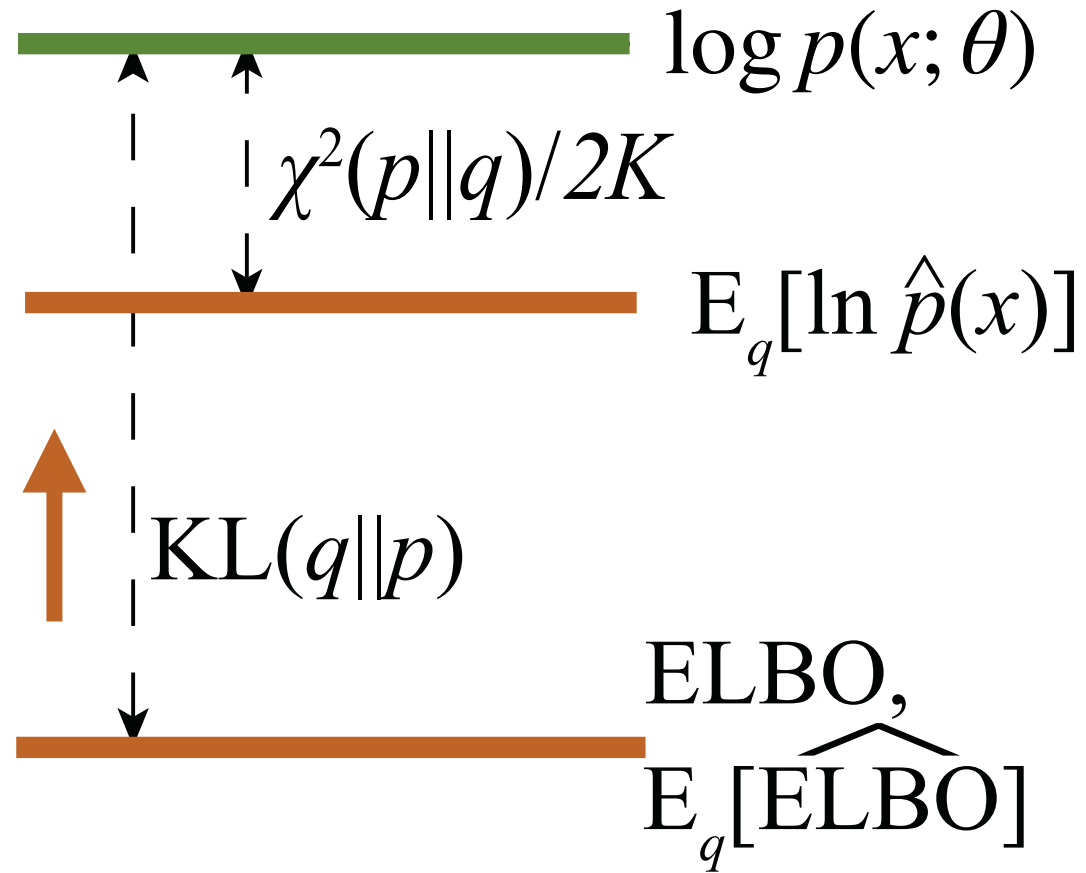
- **Two challenges**:

- ELBO doesn't always promise good parameter estimation or give the tightest lower bound, especially when a problem is very complicated like <u>POMHP</u>.

# Variational Inference



- **Two challenges**:

- ELBO doesn't always promise good parameter estimation or give the tightest lower bound, especially when a problem is very complicated like <u>POMHP</u>.

- The generally chosen $q(z; \phi)$ is an MHP considering only influence from <u>visible neurons</u> and <u>sampled history hidden neurons</u> to <u>future hidden neurons</u>.
  - inference is slow.
  - omits the influence from <u>hidden neurons</u> to <u>visible neurons</u>.

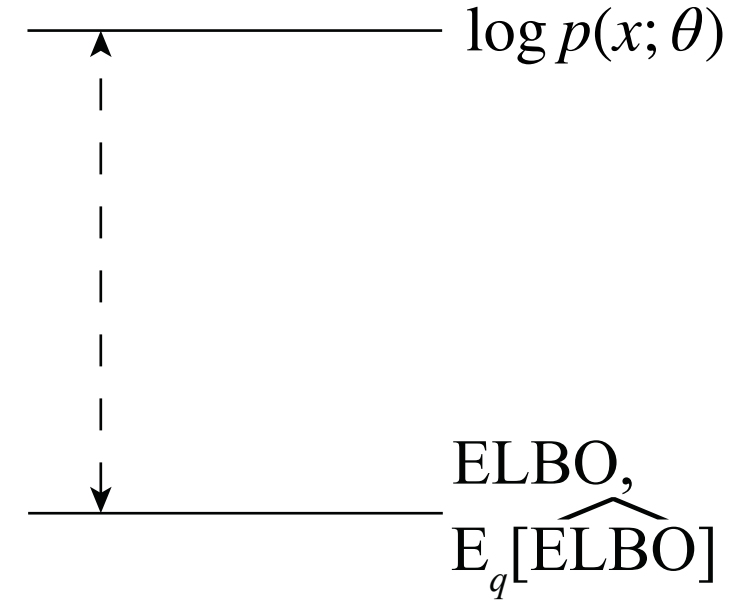# Challenge I: Tighter Lower Bound

# Variational Inference: ELBO

- ELBO $= \mathbb{E}_{z \sim q}[\log p(x \mid z, \theta)] - D_{KL}(q(z; \phi) \mid\mid p(z))$

- $\widehat{\text{ELBO}} = \dfrac{1}{K} \sum\limits_{k=1}^{K} [\log p(x, z^k; \theta) - \log q(z^k; \phi)]$

  where $\{z^k\}_{k=1}^{K}$ are K Monte Carlo samples from $q(z; \phi)$.

- $\widehat{\text{ELBO}}$ is an unbiased estimator of ELBO and **a down-biased estimator** of $\log p(x; \theta)$.

  i.e., $E_q[\widehat{\text{ELBO}}] = \text{ELBO} \leq \log p(x; \theta)$.

$\log p(x; \theta)$

$\text{ELBO,}$
$\text{E}_q[\widehat{\text{ELBO}}]$

# Importance Sampling

- Estimate the marginal with a proposal distribution $q(z; \phi)$

$$p(x; \theta) = \int p(x, z; \theta) dz = \int q(z; \phi) \frac{p(x, z; \theta)}{q(z; \phi)} dz$$

$$\approx \frac{1}{K} \sum_{k=1}^{K} \frac{p(x, z^k; \theta)}{q(z^k; \phi)} =: \hat{p}(x; \theta)$$
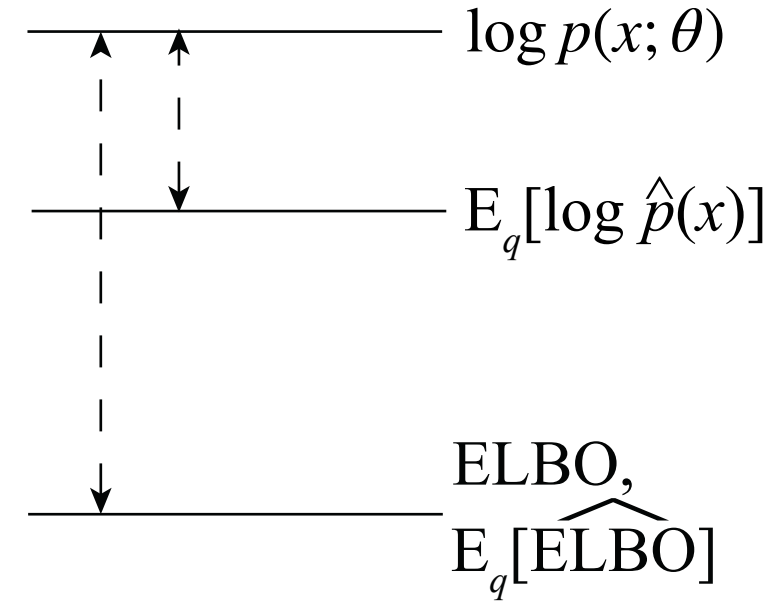
where $\{z^k\}_{k=1}^{K}$ are K Monte Carlo samples from $q(z; \phi)$.

- Since $E_q[\hat{p}(x; \theta)] = \frac{1}{K} K E_q[\frac{p(x, z; \theta)}{q(z; \phi)}] = \int p(x, z; \theta) dz = p(x; \theta)$

$\hat{p}(x; \theta)$ is an unbiased estimator of $p(x; \theta)$.

- Moreover, given Jensen's inequality $E_q[\log \hat{p}(x; \theta)] \leq \log E_q[\hat{p}(x; \theta)] = \log p(x; \theta)$

$\log \hat{p}(x; \theta)$ is **a down-biased estimator** of $\log p(x; \theta)$.

$\log p(x; \theta)$

$E_q[\log \hat{p}(x)]$

ELBO,
$E_q[\widehat{\text{ELBO}}]$

# VI vs IS

- The bias of $\widehat{\mathrm{ELBO}}$

$$E_q[\widehat{\mathrm{ELBO}} - \log p(x;\theta)] = \mathrm{ELBO} - \log p(x;\theta)$$

$$= - D_{KL}(q(z;\phi) \,||\, p(z\,|\,x,\theta))$$

- The bias of $\log \hat{p}(x;\theta)$   [Struski et al 2022]

$$E_q[\log \hat{p}(x;\theta) - \log p(x;\theta)] \approx - \frac{1}{2K}\chi^2(p(z\,|\,x,\theta) \,||\, q(z;\phi))$$

which converges to 0 when $K \to \infty$.

- When K=1, $\log \hat{p}(x;\theta) = \widehat{\mathrm{ELBO}}$ . Thus, $\log \hat{p}(x;\theta)$ is an <u>asymptotically tighter</u> lower bound compared with $\widehat{\mathrm{ELBO}}$.

$\log p(x;\theta)$

$\chi^2(p\|q)/2K$

$\mathrm{E}_q[\log \hat{p}(x)]$

$\mathrm{KL}(q\|p)$

$\mathrm{ELBO},$
$\mathrm{E}_q[\widehat{\mathrm{ELBO}}]$

# Variational Importance Sampling

---

**Algorithm 1** variational importance sampling

---

1: **function** $\text{VIS}(\boldsymbol{x}, p(\boldsymbol{x}, \boldsymbol{z}; \theta), q(\boldsymbol{z}|\boldsymbol{x}; \phi))$
2:     **for** $i = 0{:}N\text{-}1$ **do**
3:         Sample $\left\{\boldsymbol{z}^{(k)}\right\}_{k=1}^{K}$ from $q(\boldsymbol{z}|\boldsymbol{x}; \phi)$.
4:         Update $\theta$ by maximizing $\ln \hat{p}(\boldsymbol{x}; \theta)$.
5:         Update $\phi$ by minimizing $\chi^2(p(\boldsymbol{x}, \boldsymbol{z}; \theta)\|q(\boldsymbol{z}|\boldsymbol{x}; \phi))$.
6:     **end for**
7:     **return** $\theta, \phi$.
8: **end function**

---

- Inference with **importance sampling**.

- The **proposal distribution** is from minimizing $\chi^2(p(\boldsymbol{x}, \boldsymbol{z}; \theta)\|q(\boldsymbol{z}|\boldsymbol{x}; \phi))$.
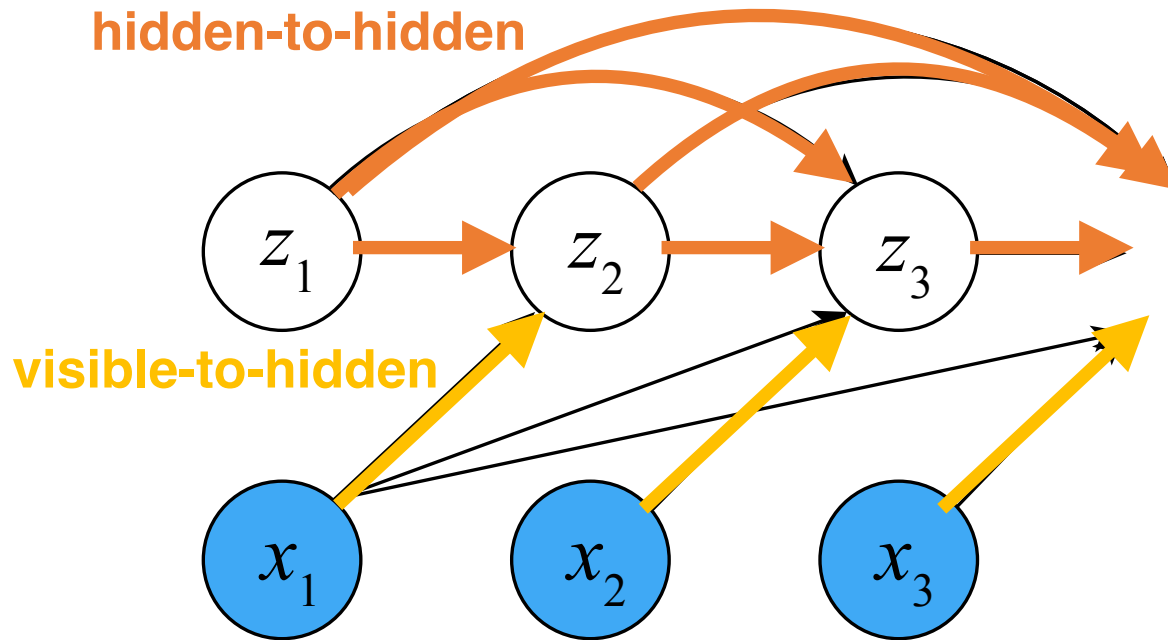
# Numerical Simulation

# Challenge II: VIS for POMHP



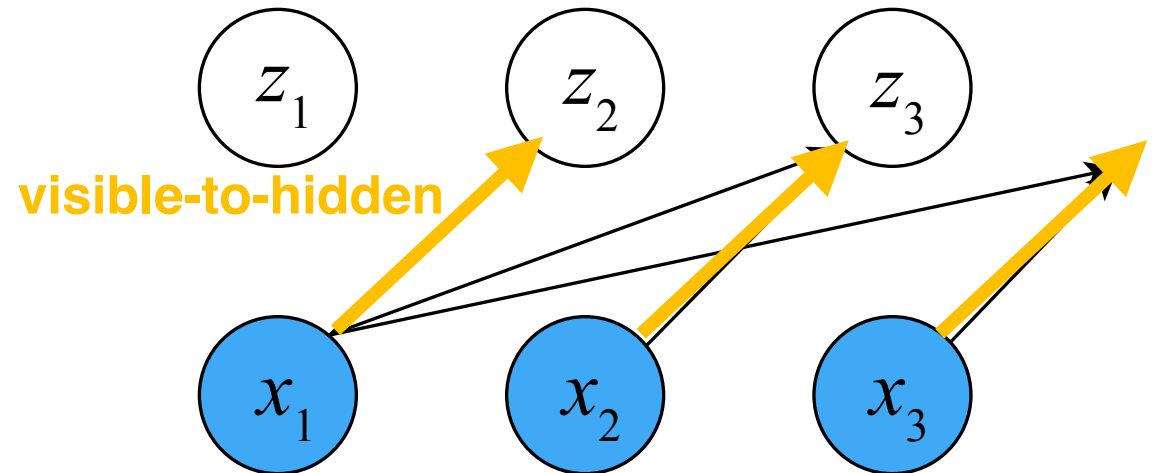The choice of the variational distribution family, $q(z \mid x; \phi)$, is important!

# Previous choices

**Forward-self** sampling
to formulate $q(z|x; \phi)$

**Forward** sampling
to formulate $q(z|x; \phi)$

# Our choice

**Forward-backward** sampling
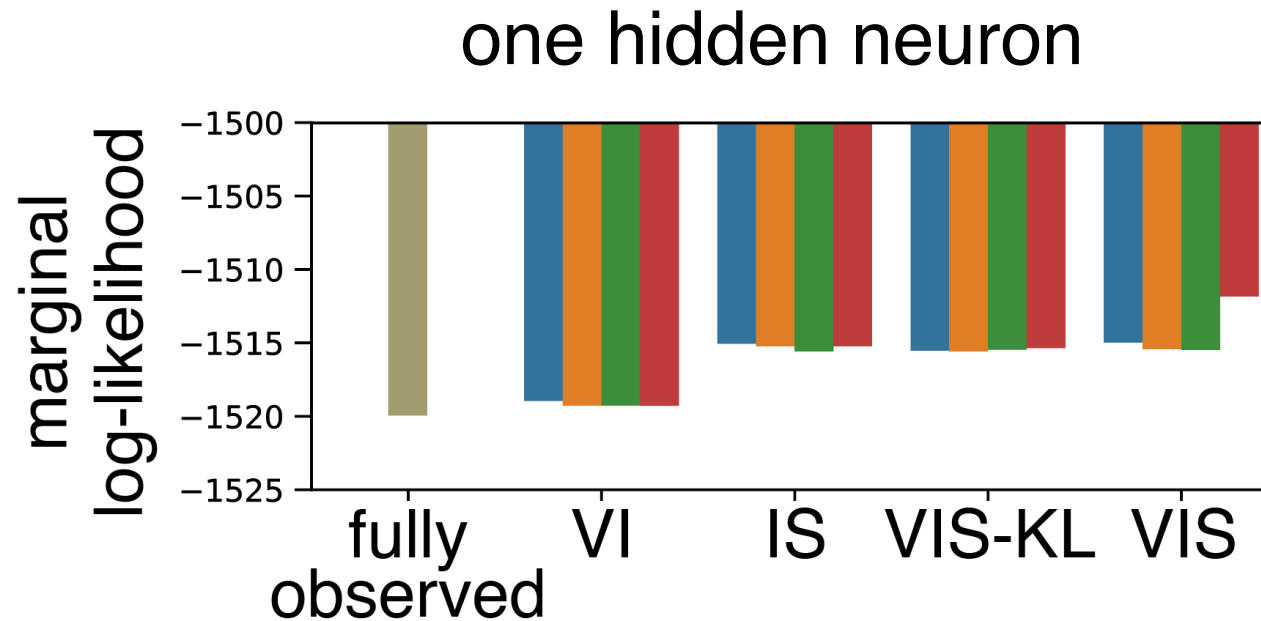to formulate $q(z \,|\, x; \phi)$



**visible-to-hidden**

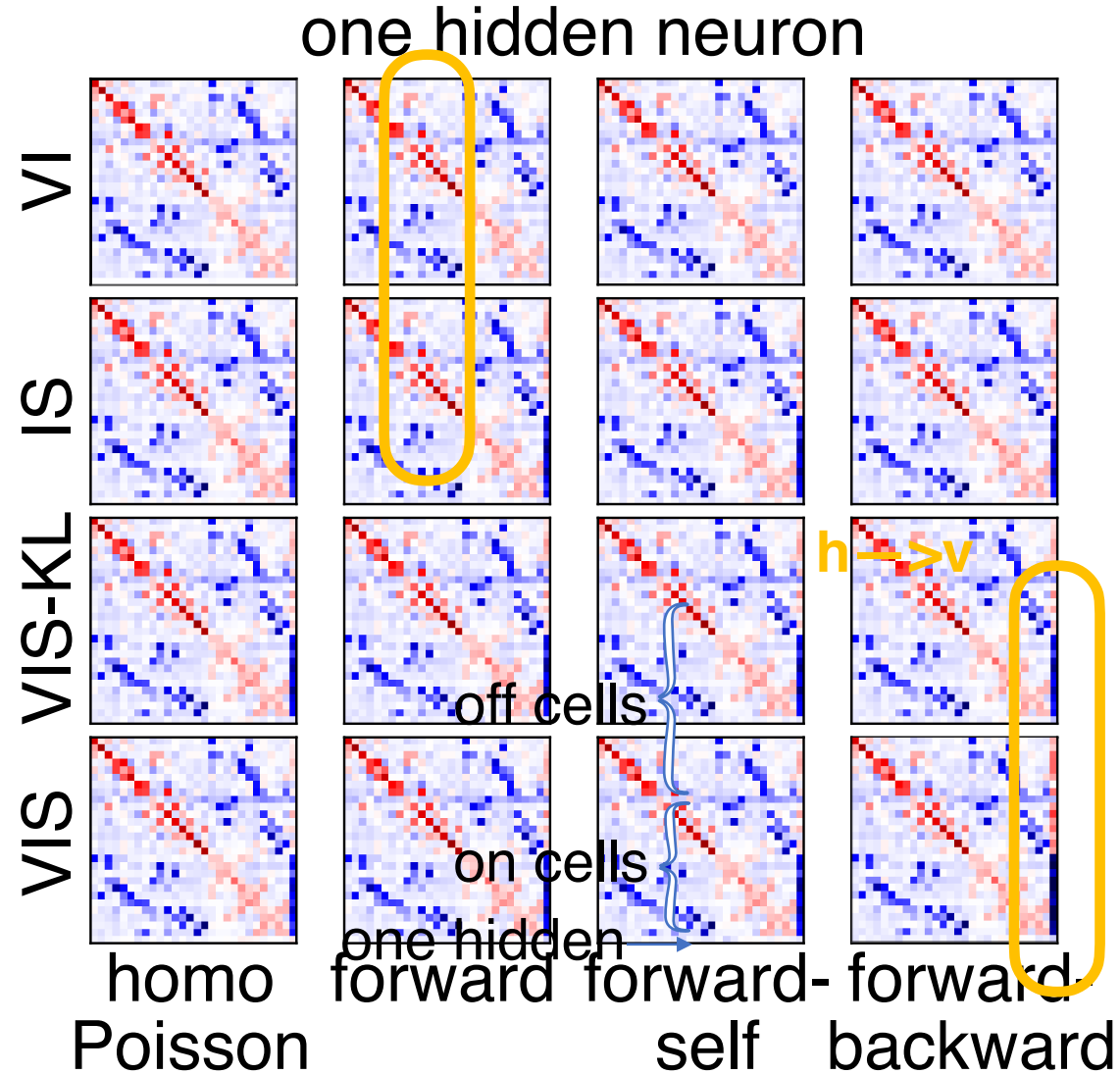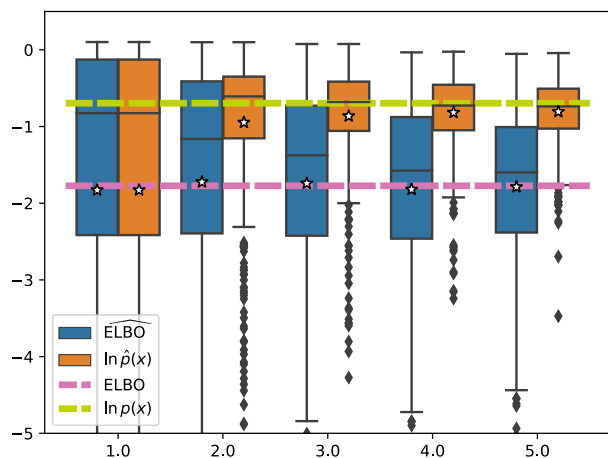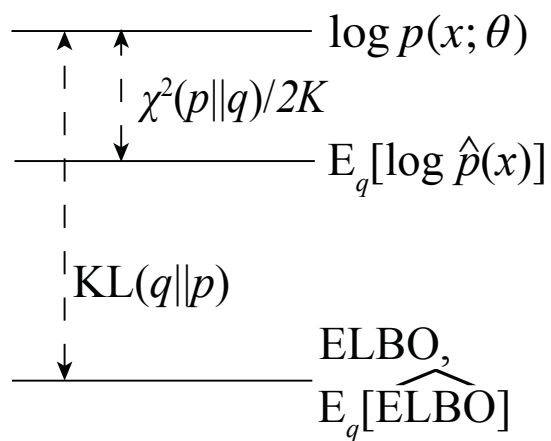**"hidden-to-visible"**

☺ efficient sampling

☺ improved accuracy

# A Toy Example
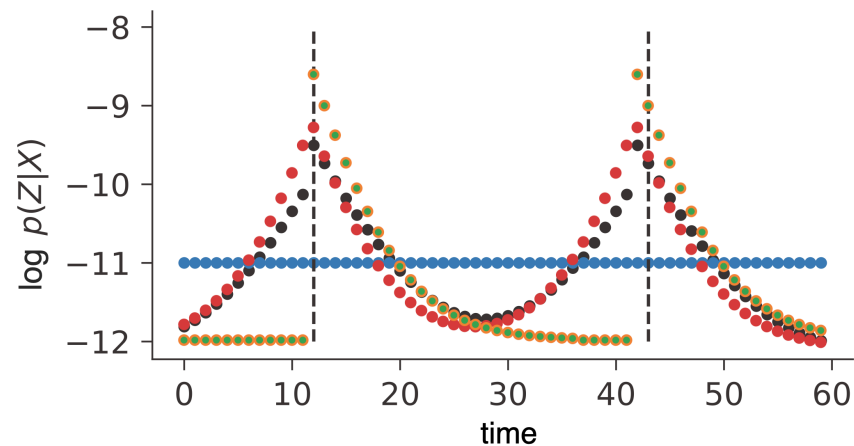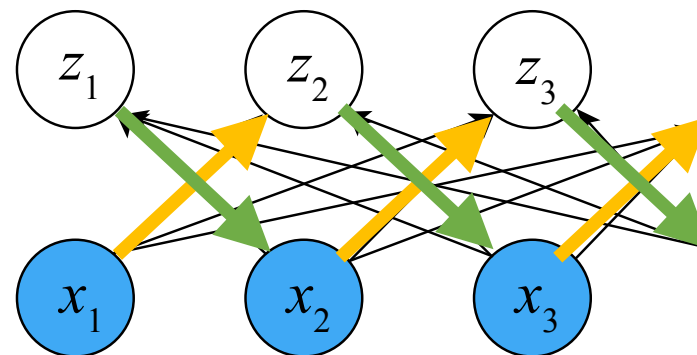
# Retinal ganglion cell (RGC) dataset

# Connectivity Weights

one hidden neuron



h→v

off cells

on cells

one hidden

VI    IS    VIS-KL    VIS

homo    forward    forward-    forward
Poisson                self    backward

# Asymptotically tighter lower bound

# Better variational distribution family
# Efficient parallel sampling



$$\log p(x; \theta)$$

$$\chi^2(p\|q)/2K$$

$$E_q[\log \hat{p}(x)]$$

$$KL(q\|p)$$

$$\text{ELBO,}$$
$$E_q[\widehat{\text{ELBO}}]$$

# Acknowledgement