# VARIATIONAL AUTOENCODER
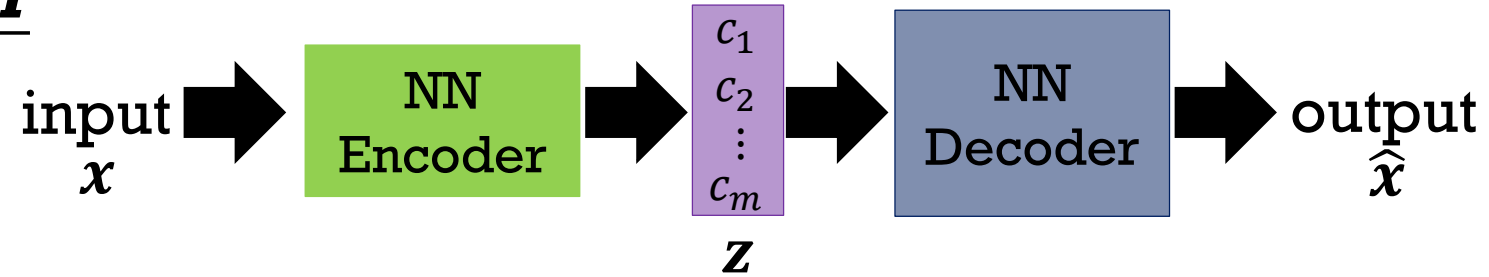
1

Pei-Yuan Wu

National Taiwan University
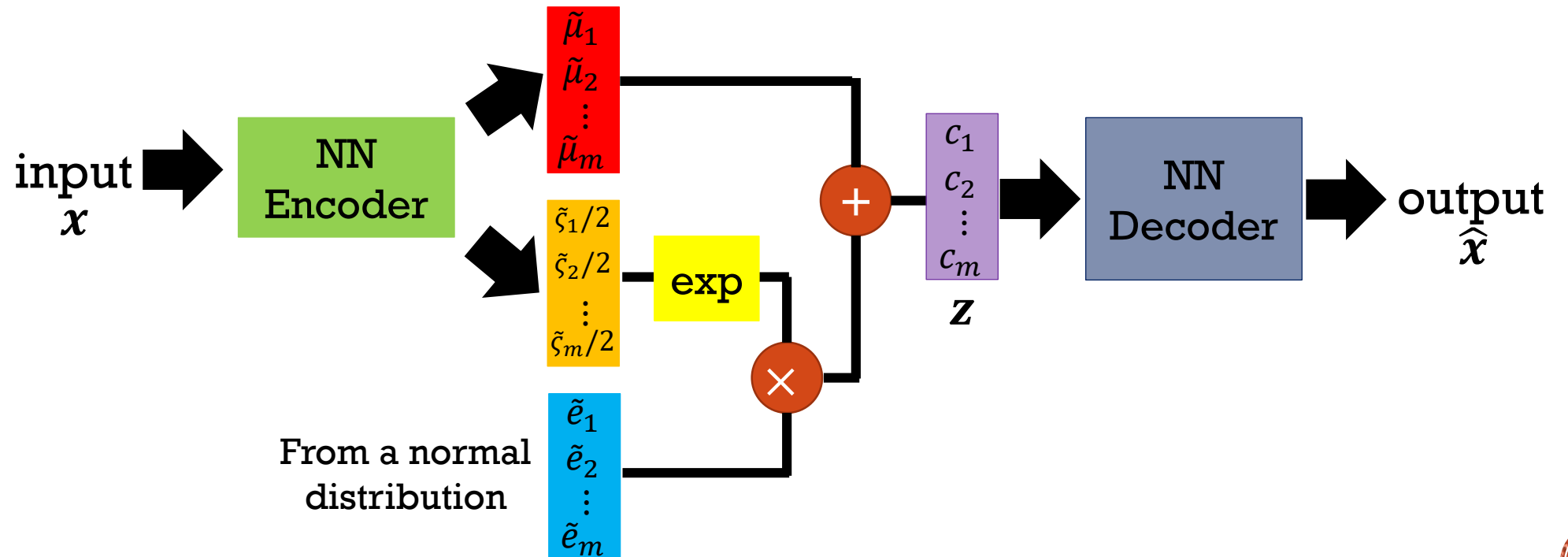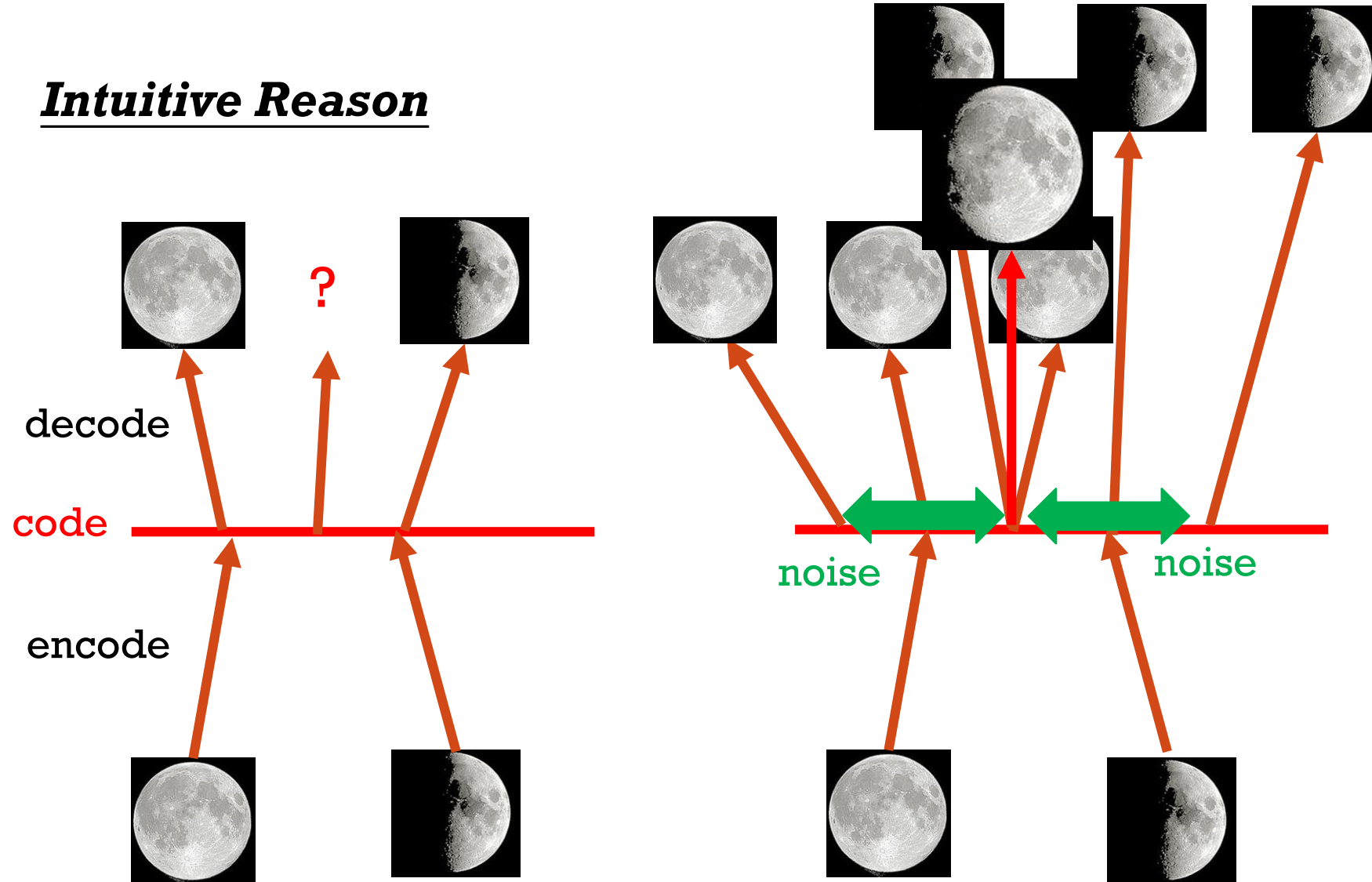
# AUTO ENCODER V.S. VAE

# WHY VAE?

**_Intuitive Reason_**
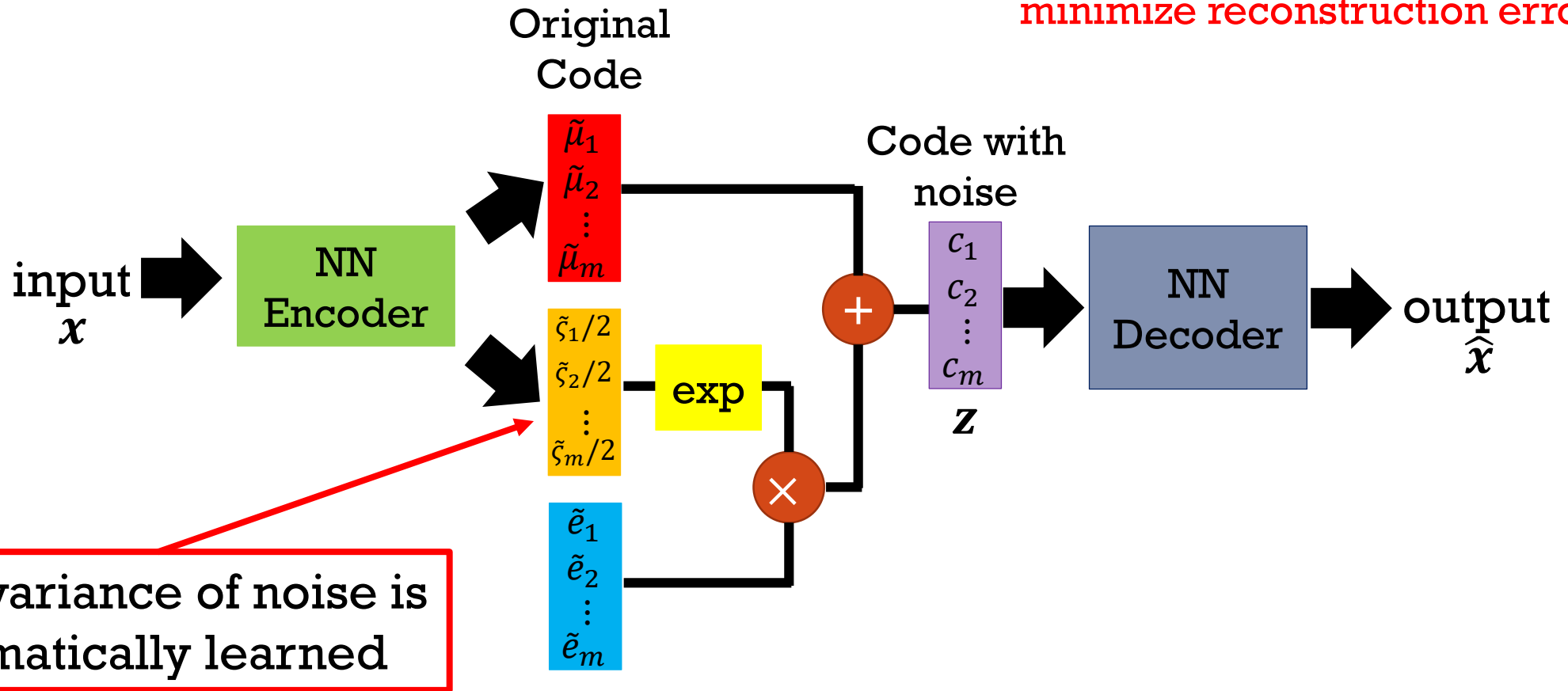
# VAE LOSS FUNCTION

What will happen if we only minimize reconstruction error?



Given data set $x_1, \dots, x_N$, minimize

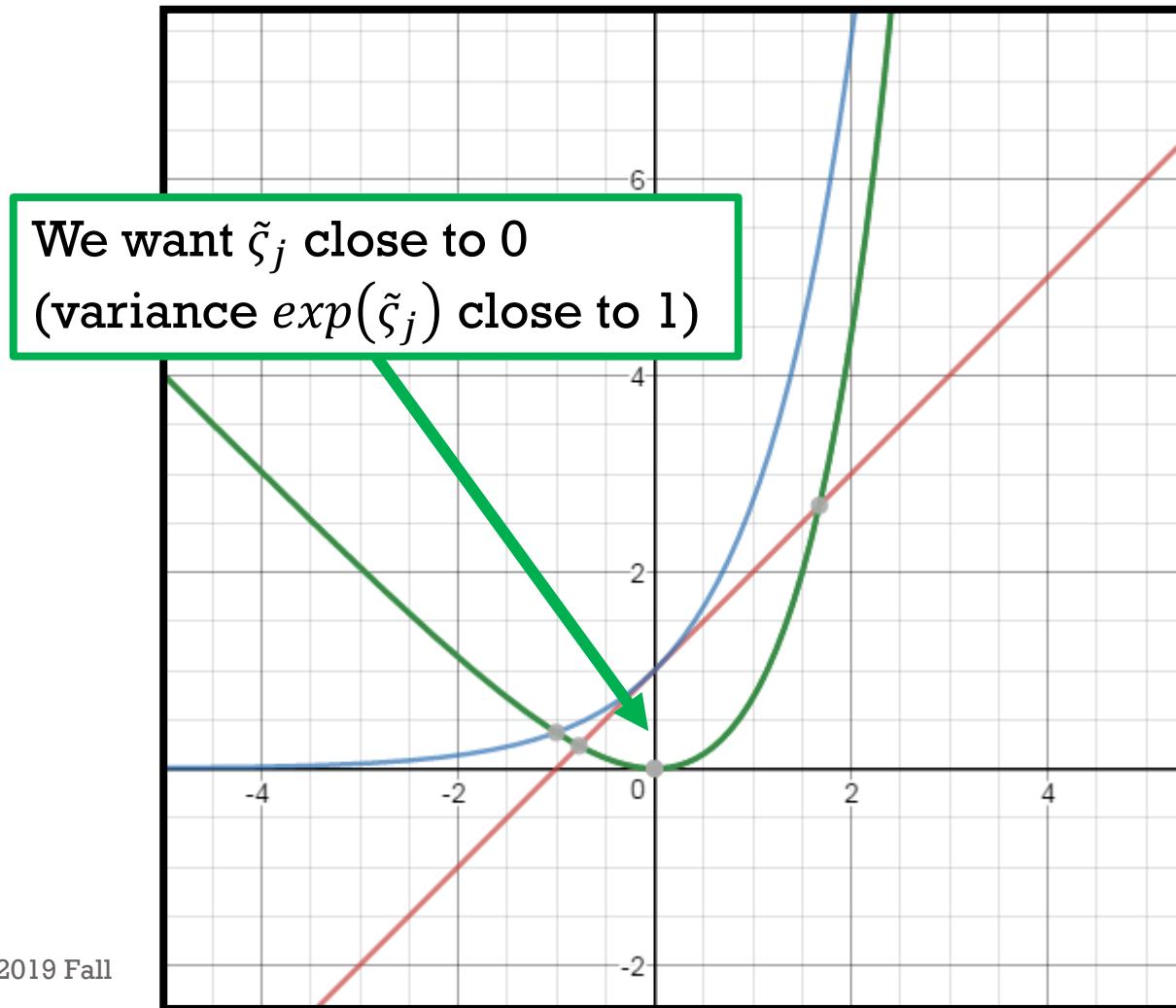$$L_{VAE} = \sum_{i=1}^{N} \left( \|x_i - \widehat{x}_i\|^2 + \sum_{j=1}^{m} \left( \|\tilde{\mu}_j(x_i)\|^2 + exp\left(\tilde{\varsigma}_j(x_i)\right) - \left(1 + \tilde{\varsigma}_j(x_i)\right) \right) \right)$$

Reconstruction loss          Regularization loss

# WHY VAE REGULARIZATION?

**_Intuitive Reason:_** _Want the code to have zero mean and unit variance_

We want $\tilde{\varsigma}_j$ close to 0
(variance $exp(\tilde{\varsigma}_j)$ close to 1)

Regularization loss
$$\left\|\tilde{\mu}_j(\boldsymbol{x}_i)\right\|^2 + exp\left(\tilde{\varsigma}_j(\boldsymbol{x}_i)\right) - \left(1 + \tilde{\varsigma}_j(\boldsymbol{x}_i)\right)$$

# VAE GENERATIVE MODEL AND MLE

- We have data set $\mathcal{X} = \{x_1, \ldots, x_N\}$, and assume each data point $x_i$ is generated according to generative model

$$p_\theta(x) = \int p_\theta(x|z)p_\theta(z)dz$$

Latent variable
(not directly observable) $z \in \mathbb{R}^m$

$\downarrow p_\theta(x|z)$

Data point
(directly observable) $x \in \mathbb{R}^M$

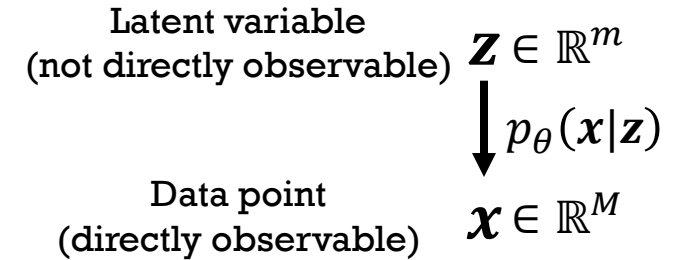- Goal: Maximum log-likelihood parameter estimation:

$$\theta^* = \max_{\theta \in \Theta} p_\theta(\mathcal{X}) = \max_{\theta \in \Theta} \log p_\theta(\mathcal{X})$$

  - If each data point $x_i$ is generated independently, then

$$\log p_\theta(\mathcal{X}) = \log\left(\prod_{i=1}^N p_\theta(x_i)\right) = \sum_{i=1}^N \log p_\theta(x_i)$$

$$p_\theta(x_i) = \int p_\theta(x_i|z_i)p_\theta(z_i)dz_i$$

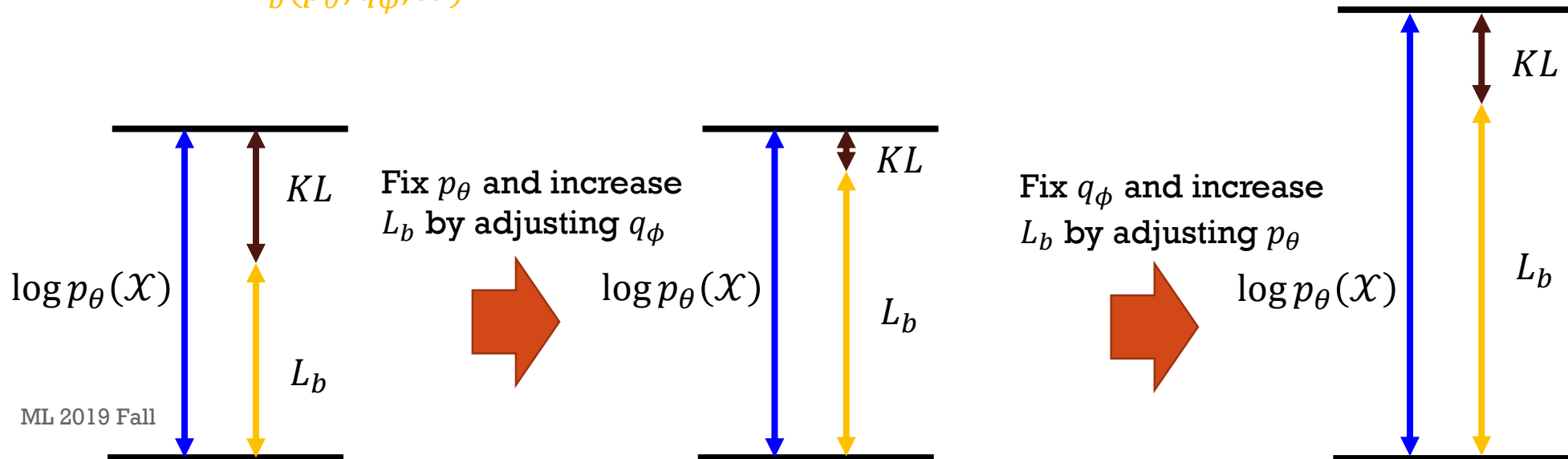- Introduce latent variables $\mathcal{Z} = \{z_1, \ldots, z_N\}$, where $z_i$ indicates the latent code of $x_i$.

# LOG-LIKELIHOOD LOWER BOUND

$$\log p_\theta(\mathcal{X}) = \int q_\phi(Z|\mathcal{X}) \log p_\theta(\mathcal{X}) \, dZ$$

$$= \int q_\phi(Z|\mathcal{X}) \left( \log \frac{p_\theta(Z,\mathcal{X})}{q_\phi(Z|\mathcal{X})} - \log \frac{p_\theta(Z|\mathcal{X})}{q_\phi(Z|\mathcal{X})} \right) dZ$$

$$= \underbrace{\int q_\phi(Z|\mathcal{X}) \log \frac{p_\theta(Z,\mathcal{X})}{q_\phi(Z|\mathcal{X})} \, dZ}_{\substack{\text{Tractable lower bound} \\ L_b(p_\theta, q_\phi, \mathcal{X})}} + \underbrace{KL\left( q_\phi(\cdot|\mathcal{X}), p_\theta(\cdot|\mathcal{X}) \right)}_{\text{Intractable but always} \geq 0}$$

$\log p_\theta(\mathcal{X})$    $KL$    $L_b$

**Fix $p_\theta$ and increase $L_b$ by adjusting $q_\phi$**

$\log p_\theta(\mathcal{X})$    $KL$    $L_b$

**Fix $q_\phi$ and increase $L_b$ by adjusting $p_\theta$**

$KL$    $\log p_\theta(\mathcal{X})$    $L_b$

# VAE V.S. EM ALGORITHM

- Randomly initialize parameters $\theta^{(1)}$.

- Iterate through step t=1,2,…

  ➤ Expectation Step (E-step): Compute
  $$Q\big(\theta\big|\theta^{(t)}\big) = \sum_{\mathcal{Z}} p\big(\mathcal{Z}\big|\mathcal{X};\theta^{(t)}\big) \log p(\mathcal{X}, \mathcal{Z};\theta)$$

  ➤ Maximization Step (M-step): Choose
  $$\theta^{(t+1)} = arg\max_{\theta \in \Theta} Q\big(\theta\big|\theta^{(t)}\big)$$

Fix $q_\phi$ and adjust $p_\theta$ to maximize
$$L_b\big(p_\theta, q_\phi, \mathcal{X}\big) = \mathbb{E}_{\mathcal{Z} \sim q_\phi(\cdot|\mathcal{X})}\left[\log \frac{p_\theta(\mathcal{X}, \mathcal{Z})}{q_\phi(\mathcal{Z}|\mathcal{X})}\right]$$

is equivalent to maximizing
$$\mathbb{E}_{\mathcal{Z} \sim \underset{\text{approx. } p_{\theta^{(t)}}(\cdot|\mathcal{X})}{q_\phi(\cdot|\mathcal{X})}}[\log p_\theta(\mathcal{X}, \mathcal{Z})] \approx Q\big(\theta\big|\theta^{(t)}\big)$$

**EM Algorithm**

Adjust $q_\phi$ to decrease $KL\left(q_\phi(\cdot|\mathcal{X}), p_{\theta^{(t)}}(\cdot|\mathcal{X})\right)$
will make $q_\phi(\cdot|\mathcal{X}) \approx p_{\theta^{(t)}}(\cdot|\mathcal{X})$

Fix $p_\theta$ and increase
$L_b$ by adjusting $q_\phi$

Fix $q_\phi$ and increase
$L_b$ by adjusting $p_\theta$

# VAE WITH INDEPENDENT SAMPLES

- Goal: Adjust $p_\theta$ and $q_\phi$ to maximize

$$L_b(p_\theta, q_\phi, \mathcal{X}) = \mathbb{E}_{Z \sim q_\phi(\cdot|\mathcal{X})}\left[\log \frac{p_\theta(\mathcal{X}, Z)}{q_\phi(Z|\mathcal{X})}\right]$$

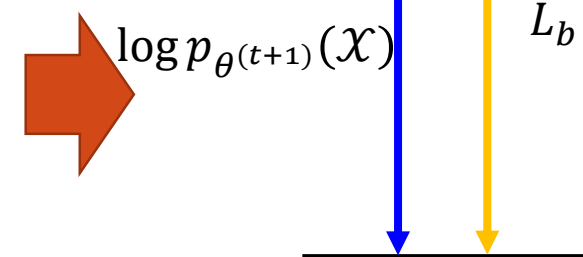$$= \mathbb{E}_{Z \sim q_\phi(\cdot|\mathcal{X})}\left[\log \frac{p_\theta(\mathcal{X}|Z)}{q_\phi(Z|\mathcal{X})}\right] + \mathbb{E}_{Z \sim q_\phi(\cdot|\mathcal{X})}[\log p_\theta(Z)]$$

- Assume $(\boldsymbol{x}_1, \boldsymbol{z}_1), \dots, (\boldsymbol{x}_N, \boldsymbol{z}_N)$ are independent w.r.t. $p_\theta$ and $q_\phi$, then

$$q_\phi(Z|\mathcal{X}) = \prod_{i=1}^{N} q_\phi(\boldsymbol{z}_i|\boldsymbol{x}_i), \ p_\theta(\mathcal{X}|Z) = \prod_{i=1}^{N} p_\theta(\boldsymbol{x}_i|\boldsymbol{z}_i), \ p_\theta(Z) = \prod_{i=1}^{N} p_\theta(\boldsymbol{z}_i)$$

$$\mathbb{E}_{Z \sim q_\phi(\cdot|\mathcal{X})}\left[\log \frac{p_\theta(\mathcal{X}|Z)}{q_\phi(Z|\mathcal{X})}\right] = \sum_{i=1}^{N} \mathbb{E}_{Z \sim q_\phi(\cdot|\mathcal{X})}\left[\log \frac{p_\theta(\boldsymbol{x}_i|\boldsymbol{z}_i)}{q_\phi(\boldsymbol{z}_i|\boldsymbol{x}_i)}\right] = \sum_{i=1}^{N} \mathbb{E}_{\boldsymbol{z}_i \sim q_\phi(\cdot|\boldsymbol{x}_i)}\left[\log \frac{p_\theta(\boldsymbol{x}_i|\boldsymbol{z}_i)}{q_\phi(\boldsymbol{z}_i|\boldsymbol{x}_i)}\right]$$

$$\mathbb{E}_{Z \sim q_\phi(\cdot|\mathcal{X})}[\log p_\theta(Z)] = \sum_{i=1}^{N} \mathbb{E}_{Z \sim q_\phi(\cdot|\mathcal{X})}[\log p_\theta(\boldsymbol{z}_i)] = \sum_{i=1}^{N} \mathbb{E}_{\boldsymbol{z}_i \sim q_\phi(\cdot|\boldsymbol{x}_i)}[\log p_\theta(\boldsymbol{z}_i)]$$

Hence

$$L_b(p_\theta, q_\phi, \mathcal{X}) = \sum_{i=1}^{N} \mathbb{E}_{\boldsymbol{z}_i \sim q_\phi(\cdot|\boldsymbol{x}_i)}\left[\log \frac{p_\theta(\boldsymbol{x}_i|\boldsymbol{z}_i)}{q_\phi(\boldsymbol{z}_i|\boldsymbol{x}_i)} + \log p_\theta(\boldsymbol{z}_i)\right]$$

# VAE AND GAUSSIAN DISTRIBUTION

- Assume latent code has Gaussian prior, and both encoder/decoder described by Gaussian distributions: $p_\theta(z) = \mathcal{N}(z; 0, I)$, $p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \Sigma_\theta(z))$, $q_\phi(z|x) = \mathcal{N}(z; \widetilde{\mu}_\phi(x), \widetilde{\Sigma}_\phi(x))$.

# CLOSE FORM OF $L_b$

- Assume latent code has Gaussian prior, and both encoder/decoder described by Gaussian distributions: $p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$, $p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_\theta(\mathbf{z}), \boldsymbol{\Sigma}_\theta(\mathbf{z}))$, $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}), \widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}))$. Then

$$\mathbb{E}_{\mathbf{z}_i \sim q_\phi(\cdot|\mathbf{x}_i)}\left[\log\frac{p_\theta(\mathbf{x}_i|\mathbf{z}_i)}{q_\phi(\mathbf{z}_i|\mathbf{x}_i)} + \log p_\theta(\mathbf{z}_i)\right] = \mathbb{E}_{\mathbf{z}_i \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i),\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i))}\left[\log\frac{\mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_\theta(\mathbf{z}_i), \boldsymbol{\Sigma}_\theta(\mathbf{z}_i))}{\mathcal{N}(\mathbf{z}_i; \widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i), \widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i))} + \log \mathcal{N}(\mathbf{z}_i; \mathbf{0}, \mathbf{I})\right]$$

$$= \mathbb{E}_{\mathbf{z}_i \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i),\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i))}[\log\mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_\theta(\mathbf{z}_i), \boldsymbol{\Sigma}_\theta(\mathbf{z}_i))] + \log\frac{\mathcal{N}(\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i); \mathbf{0}, \mathbf{I})}{\mathcal{N}(\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i); \widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i), \widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i))} - \frac{Tr(\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i) - \mathbf{I})}{2}$$

$$= \mathbb{E}_{\mathbf{z}_i \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i),\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i))}\left[\log\left(\frac{1}{\sqrt{(2\pi)^M|\boldsymbol{\Sigma}_\theta(\mathbf{z}_i)|}}\exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_\theta(\mathbf{z}_i))^T\boldsymbol{\Sigma}_\theta(\mathbf{z}_i)^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_\theta(\mathbf{z}_i))\right)\right)\right]$$

$$+ \log\left(\frac{\frac{1}{\sqrt{(2\pi)^m}}\exp\left(-\frac{1}{2}\|\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i)\|^2\right)}{1/\sqrt{(2\pi)^m|\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i)|}}\right) - \frac{Tr(\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i) - \mathbf{I})}{2}$$

$$= \mathbb{E}_{\mathbf{z}_i \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i),\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i))}\left[\log\frac{1}{\sqrt{(2\pi)^M|\boldsymbol{\Sigma}_\theta(\mathbf{z}_i)|}} - \frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_\theta(\mathbf{z}_i))^T\boldsymbol{\Sigma}_\theta(\mathbf{z}_i)^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_\theta(\mathbf{z}_i))\right] + \frac{1}{2}\log|\boldsymbol{\Sigma}_\theta(\mathbf{z}_i)|$$

$$- \frac{1}{2}\|\widetilde{\boldsymbol{\mu}}_\phi(\mathbf{x}_i)\|^2 - \frac{Tr(\widetilde{\boldsymbol{\Sigma}}_\phi(\mathbf{x}_i) - \mathbf{I})}{2}$$

**Lemma:** Let $\boldsymbol{\xi} \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \widetilde{\boldsymbol{\Sigma}})$ be a Gaussian-distributed r.v. in $\mathbb{R}^m$, then

$$\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \widetilde{\boldsymbol{\Sigma}})}[\log \mathcal{N}(\boldsymbol{\xi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})] = \log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}\left((\widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu}) + Tr(\widetilde{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1})\right)$$

$$= \log \mathcal{N}(\widetilde{\boldsymbol{\mu}}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) - Tr(\widetilde{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1})/2$$

In particular, if $\widetilde{\boldsymbol{\mu}} = \boldsymbol{\mu}$, $\widetilde{\boldsymbol{\Sigma}} = \boldsymbol{\Sigma}$, then
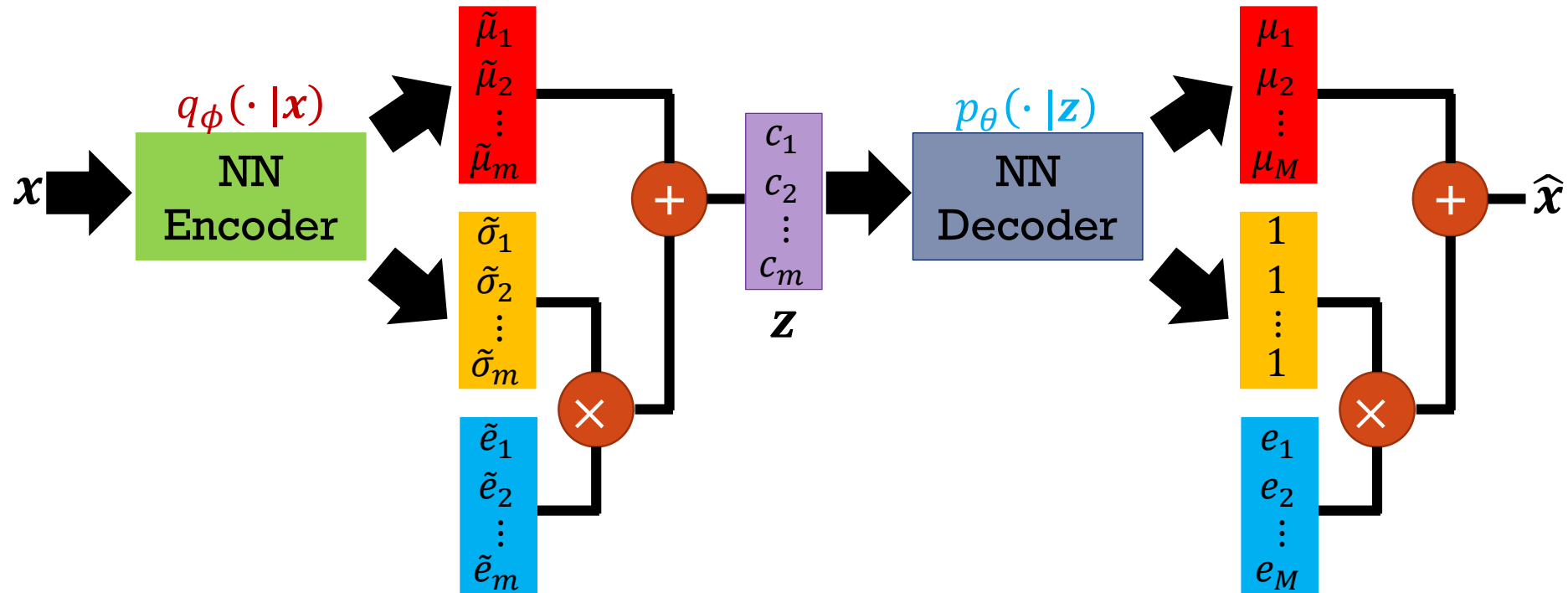
$$\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})}[\log \mathcal{N}(\boldsymbol{\xi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})] = -\frac{m}{2}\log\left(2\pi e |\boldsymbol{\Sigma}|^{1/m}\right)$$

Proof:

$$\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \widetilde{\boldsymbol{\Sigma}})}[\log \mathcal{N}(\boldsymbol{\xi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})] = \mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \widetilde{\boldsymbol{\Sigma}})}\left[\log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}(\boldsymbol{\xi} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\xi} - \boldsymbol{\mu})\right]$$

$$= \log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \widetilde{\boldsymbol{\Sigma}})}\left[Tr((\boldsymbol{\xi} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\xi} - \boldsymbol{\mu}))\right]$$

$$= \log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}Tr\left(\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \widetilde{\boldsymbol{\Sigma}})}[(\boldsymbol{\xi} - \boldsymbol{\mu})(\boldsymbol{\xi} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}]\right)$$

$$= \log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}Tr\left(\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \widetilde{\boldsymbol{\Sigma}})}[(\boldsymbol{\xi} - \boldsymbol{\mu})(\boldsymbol{\xi} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}]\right)$$

$$= \log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}Tr\left(\mathbb{E}_{\boldsymbol{\xi}' \sim \mathcal{N}(\mathbf{0}, \widetilde{\boldsymbol{\Sigma}})}[(\boldsymbol{\xi}' + \widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu})(\boldsymbol{\xi}' + \widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}]\right)$$

$$= \log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}Tr\left((\widetilde{\boldsymbol{\Sigma}} + (\widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu})(\widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu})^T)\boldsymbol{\Sigma}^{-1}\right) = \log \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} - \frac{1}{2}\left((\widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\widetilde{\boldsymbol{\mu}} - \boldsymbol{\mu}) + Tr(\widetilde{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1})\right)$$

# SIMPLER CLOSE FORM OF $L_b$

- For simplicity, assume $\mathbf{\Sigma}_\theta(\mathbf{z}_i) = \mathbf{I}, \widetilde{\mathbf{\Sigma}}_\phi(\mathbf{x}) = diag(\tilde{\sigma}_{\phi,1}^2(\mathbf{x}), \dots, \tilde{\sigma}_{\phi,m}^2(\mathbf{x}))$.

# VAE LOSS FUNCTION

- For simplicity, assume $\Sigma_\theta(z_i) = I$, $\widetilde{\Sigma}_\phi(x) = diag(\tilde{\sigma}_{\phi,1}^2(x), \ldots, \tilde{\sigma}_{\phi,m}^2(x))$, then

$$\mathbb{E}_{z_i \sim q_\phi(\cdot|x_i)} \left[ \log \frac{p_\theta(x_i|z_i)}{q_\phi(z_i|x_i)} + \log p_\theta(z_i) \right]$$

$$= \mathbb{E}_{z_i \sim \mathcal{N}(\widetilde{\mu}_\phi(x_i), \widetilde{\Sigma}_\phi(x_i))} \left[ \log \frac{1}{\sqrt{(2\pi)^M |\Sigma_\theta(z_i)|}} - \frac{1}{2}(x_i - \mu_\theta(z_i))^T \Sigma_\theta(z_i)^{-1}(x_i - \mu_\theta(z_i)) \right] + \frac{1}{2}\log|\Sigma_\theta(z_i)| - \frac{1}{2}\|\widetilde{\mu}_\phi(x_i)\|^2$$

$$- \frac{1}{2}Tr(\widetilde{\Sigma}_\phi(x_i) - I) = \log \frac{1}{\sqrt{(2\pi)^M}} - \frac{1}{2}\left( \mathbb{E}_{z_i \sim \mathcal{N}(\widetilde{\mu}_\phi(x_i), \widetilde{\Sigma}_\phi(x_i))}[\|x_i - \mu_\theta(z_i)\|^2] + \|\widetilde{\mu}_\phi(x_i)\|^2 + \sum_{j=1}^{m}(\tilde{\sigma}_{\phi,j}^2(x_i) - \log \tilde{\sigma}_{\phi,j}^2(x_i) - 1) \right)$$

Hence maximizing

$$L_b(p_\theta, q_\phi, \mathcal{X}) = N \log \frac{1}{\sqrt{(2\pi)^M}} - \frac{1}{2}\sum_{i=1}^{N}\left( \mathbb{E}_{z_i \sim \mathcal{N}(\widetilde{\mu}_\phi(x_i), \widetilde{\Sigma}_\phi(x_i))}[\|x_i - \mu_\theta(z_i)\|^2] + \|\widetilde{\mu}_\phi(x_i)\|^2 + \sum_{j=1}^{m}(\tilde{\sigma}_{\phi,j}^2(x_i) - \log \tilde{\sigma}_{\phi,j}^2(x_i) - 1) \right)$$

Is equivalent to minimizing

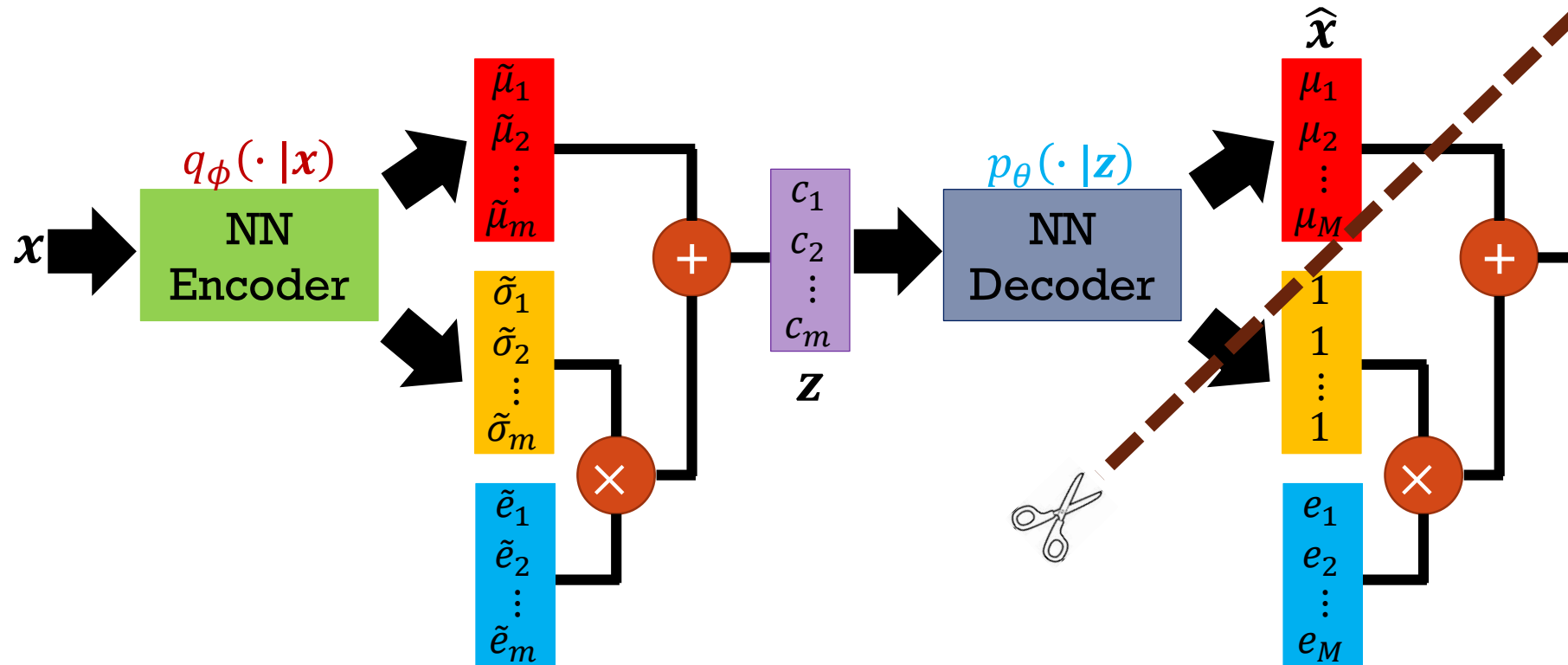$$\sum_{i=1}^{N}\left( \mathbb{E}_{z_i \sim \mathcal{N}(\widetilde{\mu}_\phi(x_i), \widetilde{\Sigma}_\phi(x_i))}[\|x_i - \mu_\theta(z_i)\|^2] + \|\widetilde{\mu}_\phi(x_i)\|^2 + \sum_{j=1}^{m}\left( \tilde{\sigma}_{\phi,j}^2(x_i) - (1 + \log \tilde{\sigma}_{\phi,j}^2(x_i)) \right) \right)$$

Approx. expectation by sample mean

# VAE LOSS FUNCTION

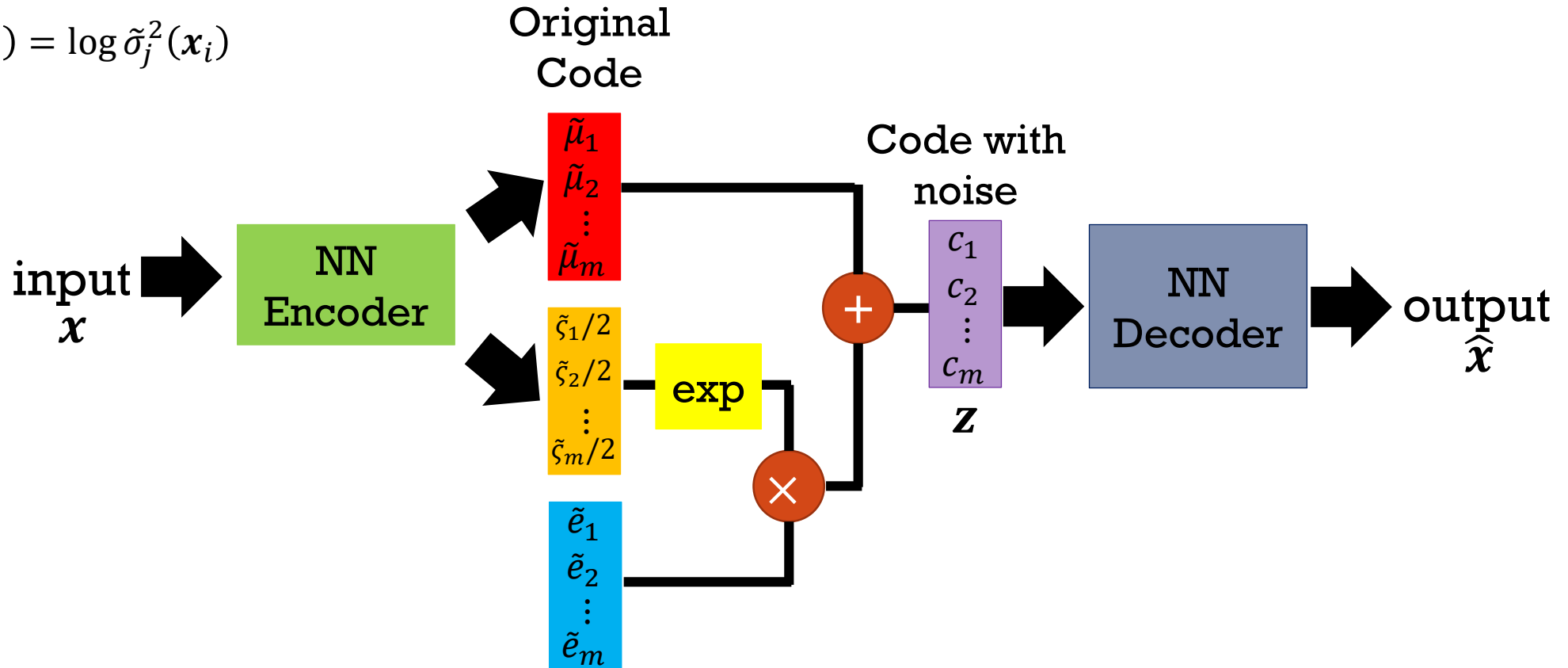- For simplicity, assume $\Sigma_\theta(z_i) = I, \tilde{\Sigma}_\phi(x) = diag(\tilde{\sigma}^2_{\phi,1}(x), \ldots, \tilde{\sigma}^2_{\phi,m}(x))$.



$$L_{VAE} = \sum_{i=1}^{N} \left( \|x_i - \hat{x}_i\|^2 + \sum_{j=1}^{m} \left( \|\tilde{\mu}_j(x_i)\|^2 + \tilde{\sigma}^2_j(x_i) - \left(1 + \log \tilde{\sigma}^2_j(x_i)\right) \right) \right)$$

# VAE LOSS FUNCTION

Take $\tilde{\varsigma}_j(\boldsymbol{x}_i) = \log \tilde{\sigma}_j^2(\boldsymbol{x}_i)$



$$L_{VAE} = \sum_{i=1}^{N}\left(\|\boldsymbol{x}_i - \widehat{\boldsymbol{x}}_i\|^2 + \sum_{j=1}^{m}\left(\|\tilde{\mu}_j(\boldsymbol{x}_i)\|^2 + exp\left(\tilde{\varsigma}_j(\boldsymbol{x}_i)\right) - \left(1 + \tilde{\varsigma}_j(\boldsymbol{x}_i)\right)\right)\right)$$