

E9: 309 ADL 23-11-2020



Housekeeping

* Mid-term exam

→ December 5th (Saturday) [Topics covered up to Dec 2nd]

→ Mode of exam

✓ Time to respond - 3 hours

○ Exam paper uploaded in Teams Channel and response (photo-scanned and uploaded in your folder).

○ Open book, open notes

★ Strictly no online communication or help sought.

★ Academic integrity and ethics strongly followed.



Recap of previous class

Restricted Boltzmann machines

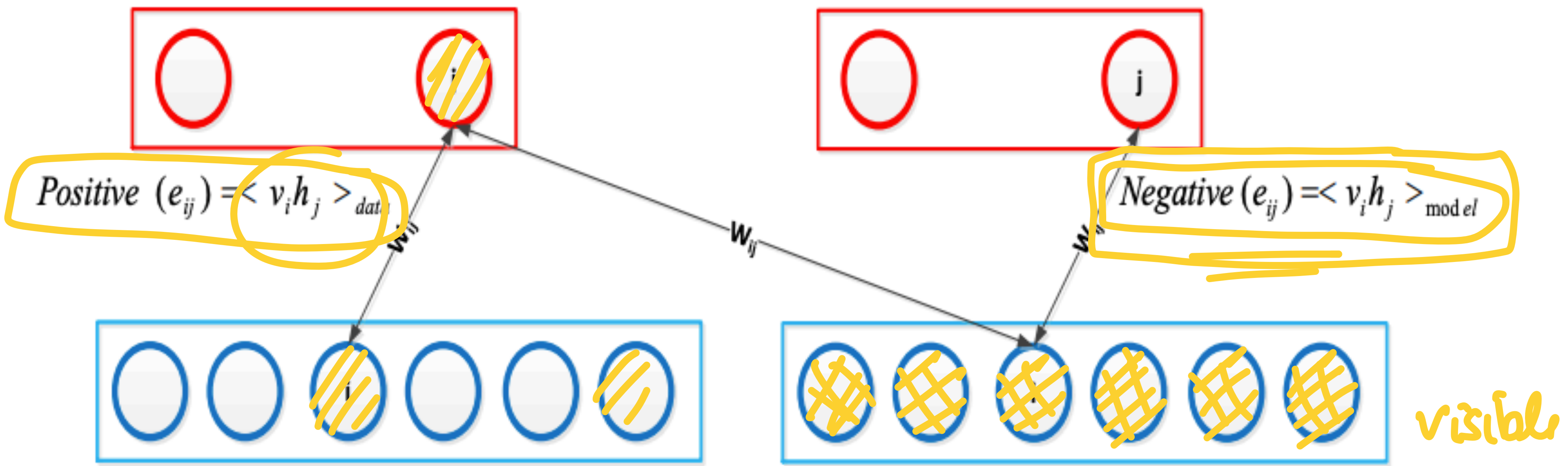
- Properties
- Training
 - Approximation



Contrastive Divergence

$$\sqrt{v_i} p(h_j = 1/v) = \sigma$$

$$p(h_j = 1/v) \neq j$$



One-step contrastive divergence

* Computing the gradient

$$\frac{\partial(p([\mathbf{v} \ \mathbf{h}], \Theta))}{\partial \mathbf{W}} \approx \underbrace{\frac{1}{N} \sum_{q=1}^N \mathbf{v}_q \mathbf{h}_q^T}_{\text{data}} - \underbrace{\frac{1}{N} \sum_{q=1}^N \tilde{\mathbf{v}}_q \tilde{\mathbf{h}}_q^T}_{\text{model}}$$

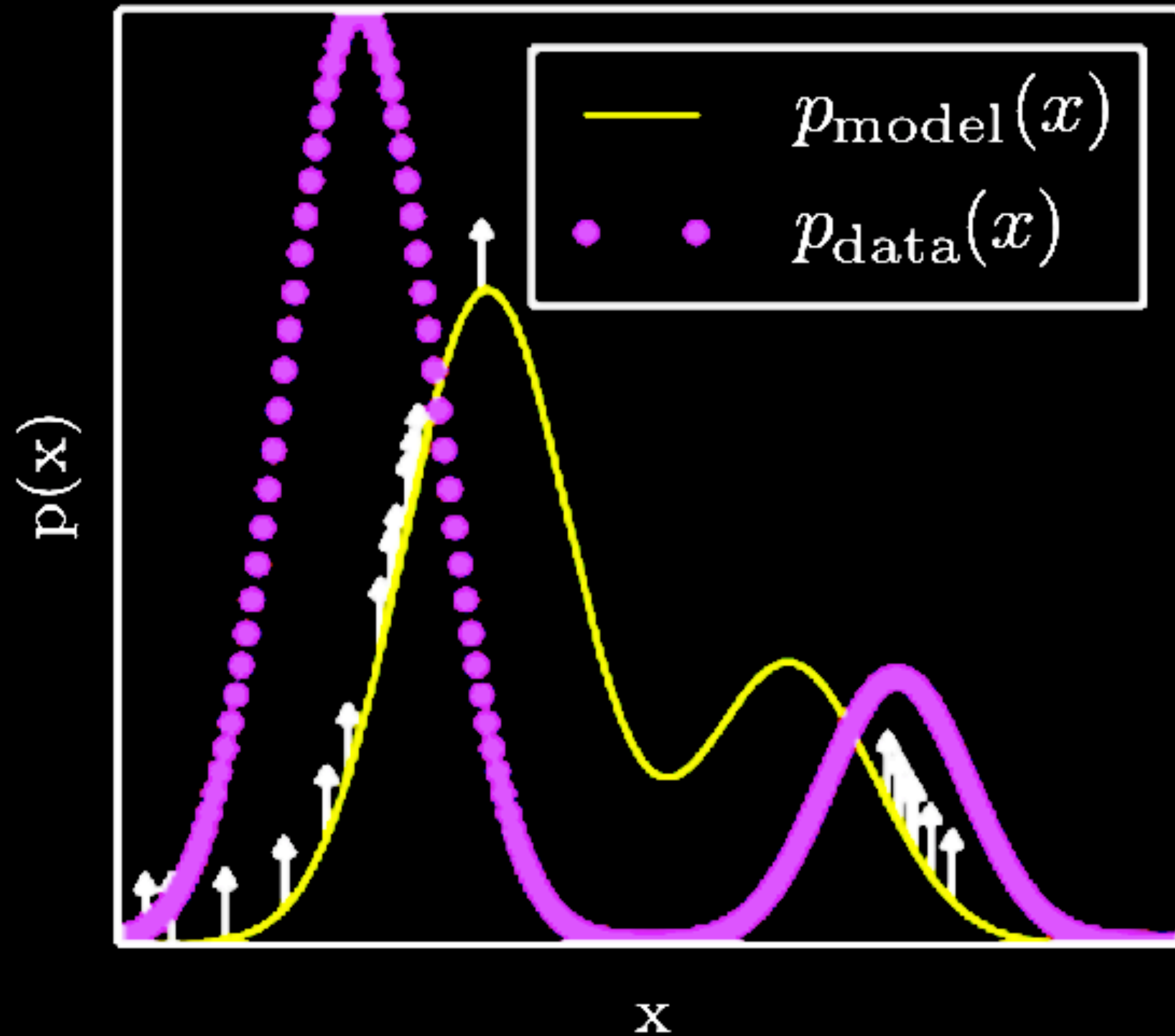
* Performing gradient ascent using the approximate gradient

$$\Theta^{k+1} = \Theta^k + \eta \left. \frac{\partial \log(p(\mathbf{x}, \Theta))}{\partial \Theta} \right|_{\Theta = \Theta^k}$$

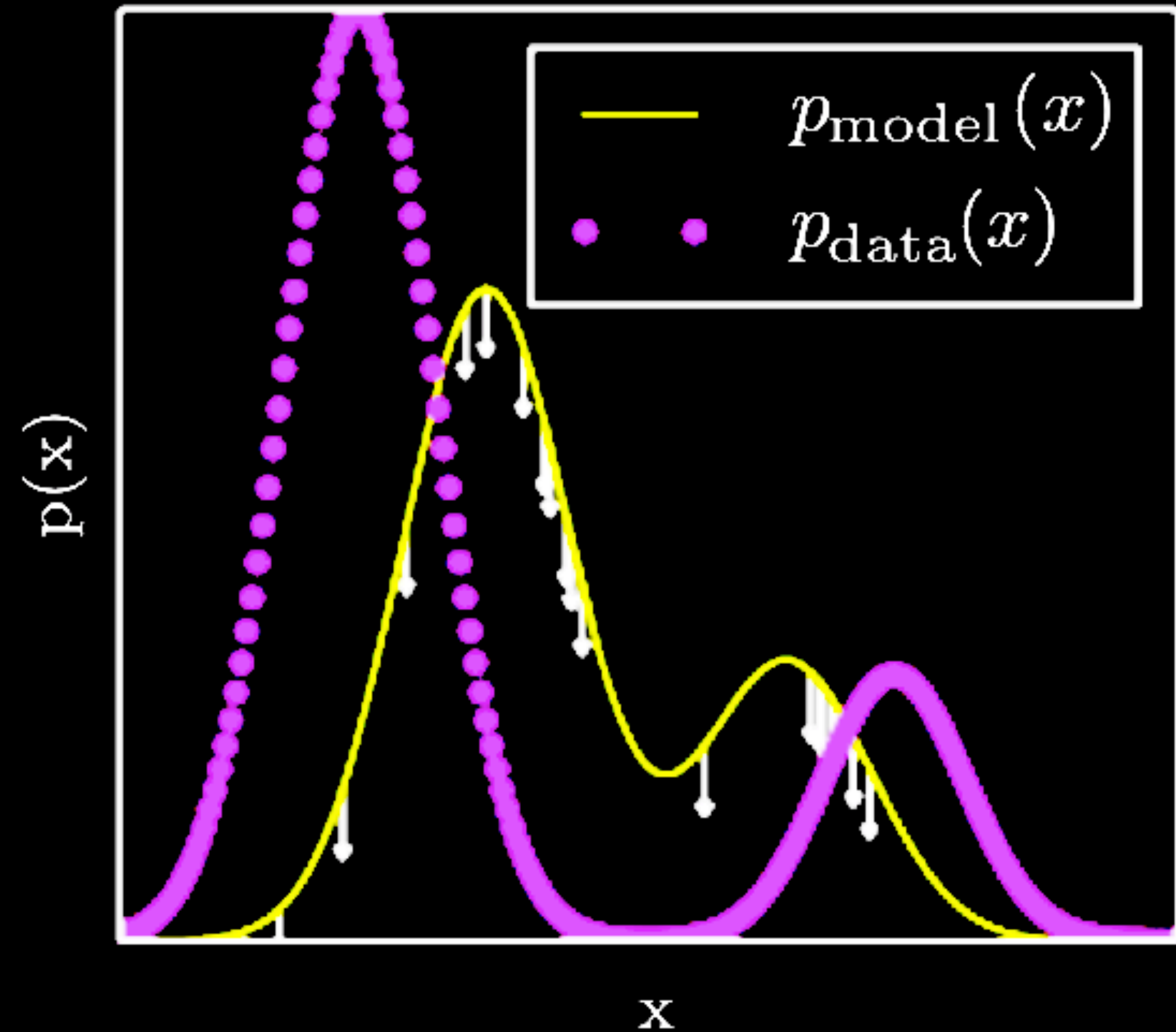


Positive phase and negative phase

The positive phase



The negative phase



Gaussian Bernoulli RBM

* For modeling real observations $\mathbf{v} \in \mathcal{R}^D$

* Define the energy function

$$E[\mathbf{v}, \mathbf{h}] = \frac{1}{2} \underbrace{(\mathbf{v} - \mathbf{a})^T (\mathbf{v} - \mathbf{a})} - \mathbf{v}^T \mathbf{W} \mathbf{h} - \mathbf{b}^T \mathbf{h}$$

$$p([\mathbf{v}, \mathbf{h}]) = \frac{e^{-E[\mathbf{v}, \mathbf{h}]}}{Z}$$

* The conditional distributions

$$p(\mathbf{v} | \mathbf{h}) = \mathcal{N}(\mathbf{W} \mathbf{h} + \mathbf{a}, \mathbf{I})$$

$$p(h_j = 1 | \mathbf{v}) = \sigma(\mathbf{v}^T \mathbf{W}_{:,j}^k + b_j)$$

✓
Perform these derivations



Properties of GRBM

* $d = 0$

$$E(\mathbf{v}) = \frac{1}{2}(\mathbf{v} - \mathbf{a})^T(\mathbf{v} - \mathbf{a})$$

$$p(\mathbf{v}) = \frac{e^{-E(\mathbf{v})}}{Z}$$

* The marginal distribution is a Gaussian.



Properties of GRBM

* $d = 1$ $E[\mathbf{v}, h] = \frac{1}{2}(\mathbf{v} - \mathbf{a})^T(\mathbf{v} - \mathbf{a}) - h\mathbf{v}^T\mathbf{w} - hb$

$$p([\mathbf{v}, h]) = \frac{e^{-E[\mathbf{v}, h]}}{Z}$$

$$p([\mathbf{v}, h = 0]) = \alpha \mathcal{N}(\mathbf{a}, \mathbf{I})$$

$$p([\mathbf{v}, h = 1]) = (1 - \alpha) \mathcal{N}(\mathbf{a} + \mathbf{w}, \mathbf{I})$$

* The marginal distribution is then

$$p(\mathbf{v}) = p([\mathbf{v}, h = 0]) + p([\mathbf{v}, h = 1])$$

✓ 2-mixture Gaussian



Properties of GRBM

$$\underline{h}_d = \begin{bmatrix} h_{d-1} \\ \vdots \\ h_d \end{bmatrix}$$

Handwritten diagram showing a vertical stack of terms h_{d-1}, \dots, h_d enclosed in large square brackets. An arrow labeled '0' points to the h_d term, and an arrow labeled '1' points to the h_{d-1} term.

* For any general d dimensions

$$E[\mathbf{v}, \mathbf{h}_d] = E[\mathbf{v}, \mathbf{h}_{d-1}] + h_d \mathbf{v}^T \mathbf{W}_{:,d} + b_d h_d$$

$$p([\mathbf{v}, [\mathbf{h}_{d-1}, h_d = 0]]) = \alpha p([\mathbf{v}, \mathbf{h}_{d-1}])$$

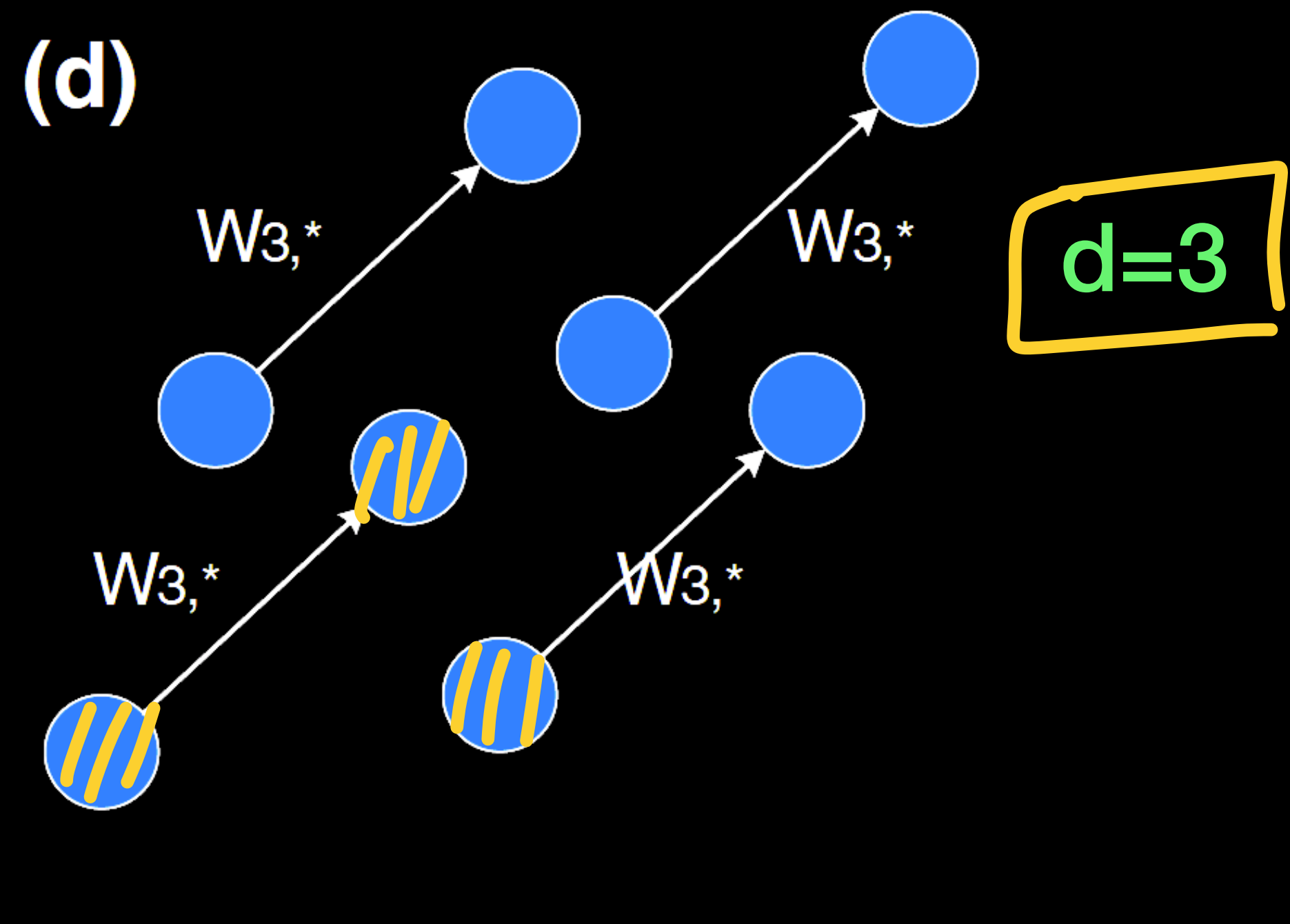
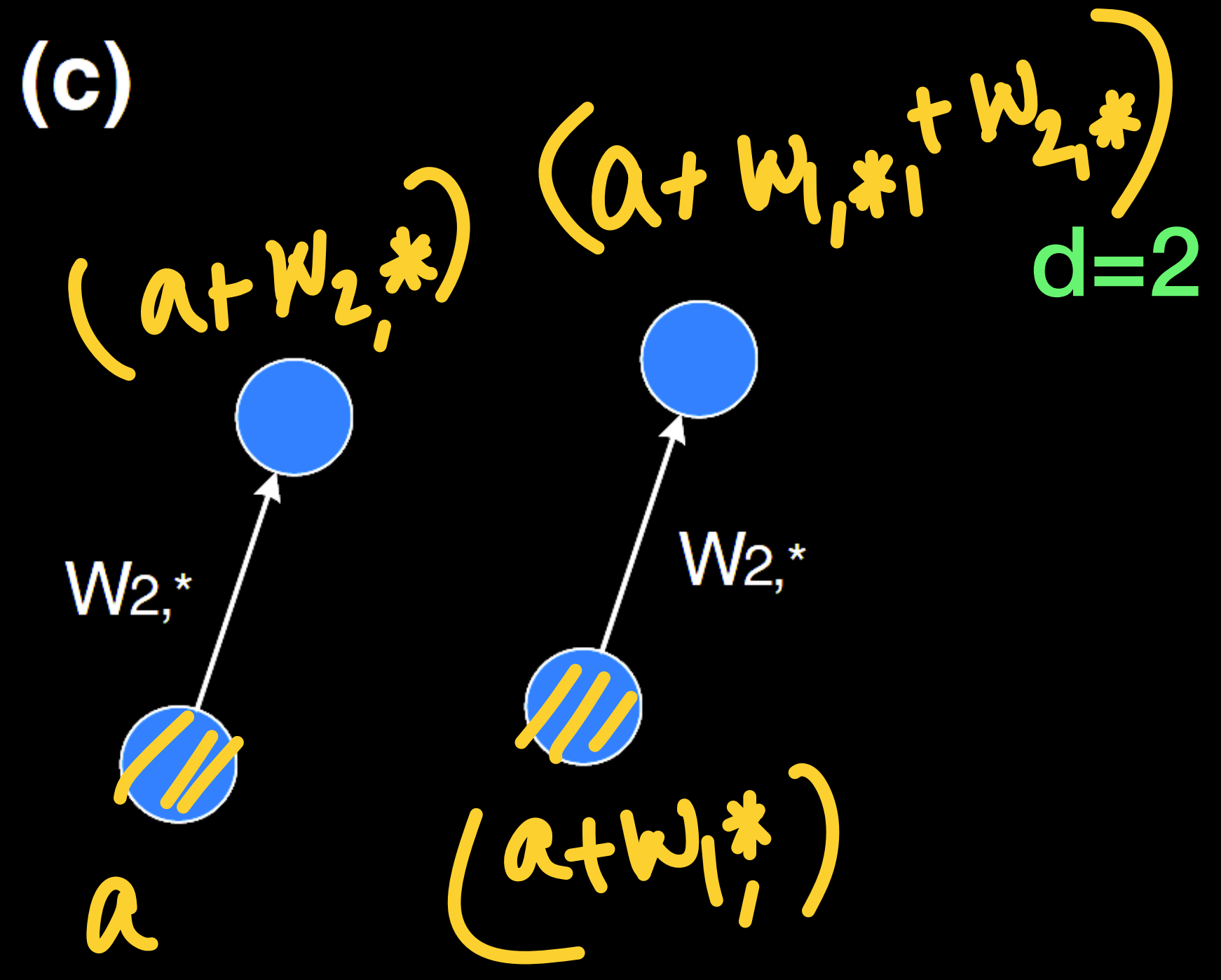
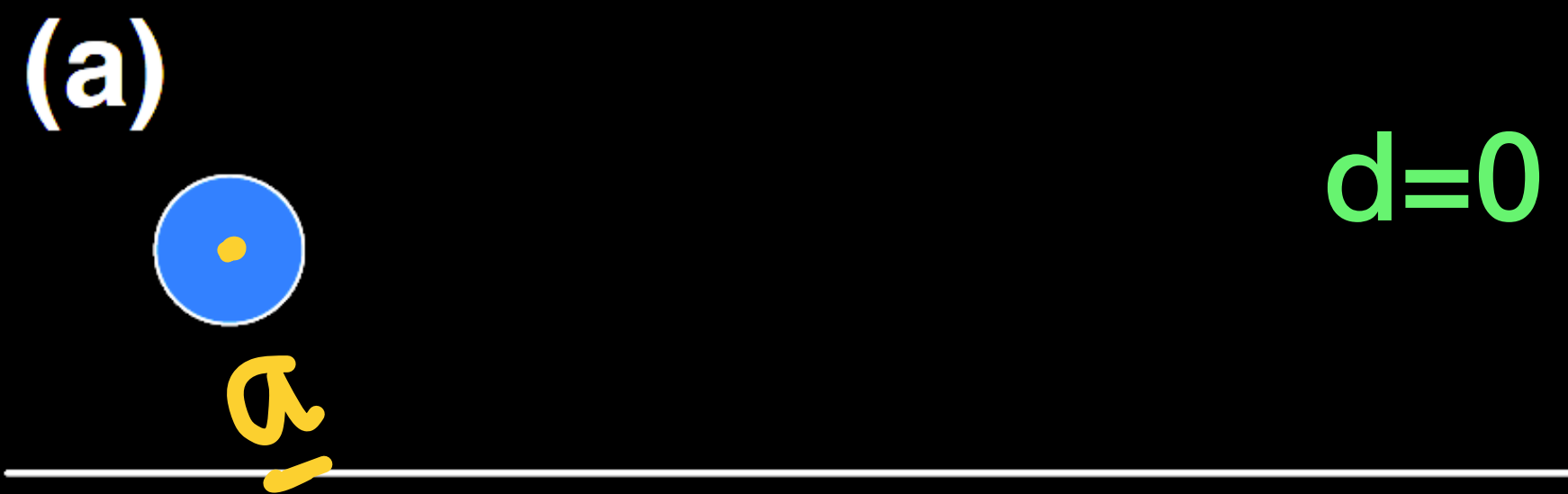
$$p([\mathbf{v}, [\mathbf{h}_{d-1}, h_d = 1]]) = (1 - \alpha) p([\mathbf{v} + \mathbf{W}_{:,d}, \mathbf{h}_{d-1}])$$

* For $d=0$, 1 Gaussian, $d=1$, 2-mix Gaussian, ...

→ 2^d mixture Gaussian for any arbitrary d .



GRBMs and GMMs



Deep Belief Networks (DBN)

* Stacking RBMs in a disjoint fashion

✓ Layer-wise training for each RBM.

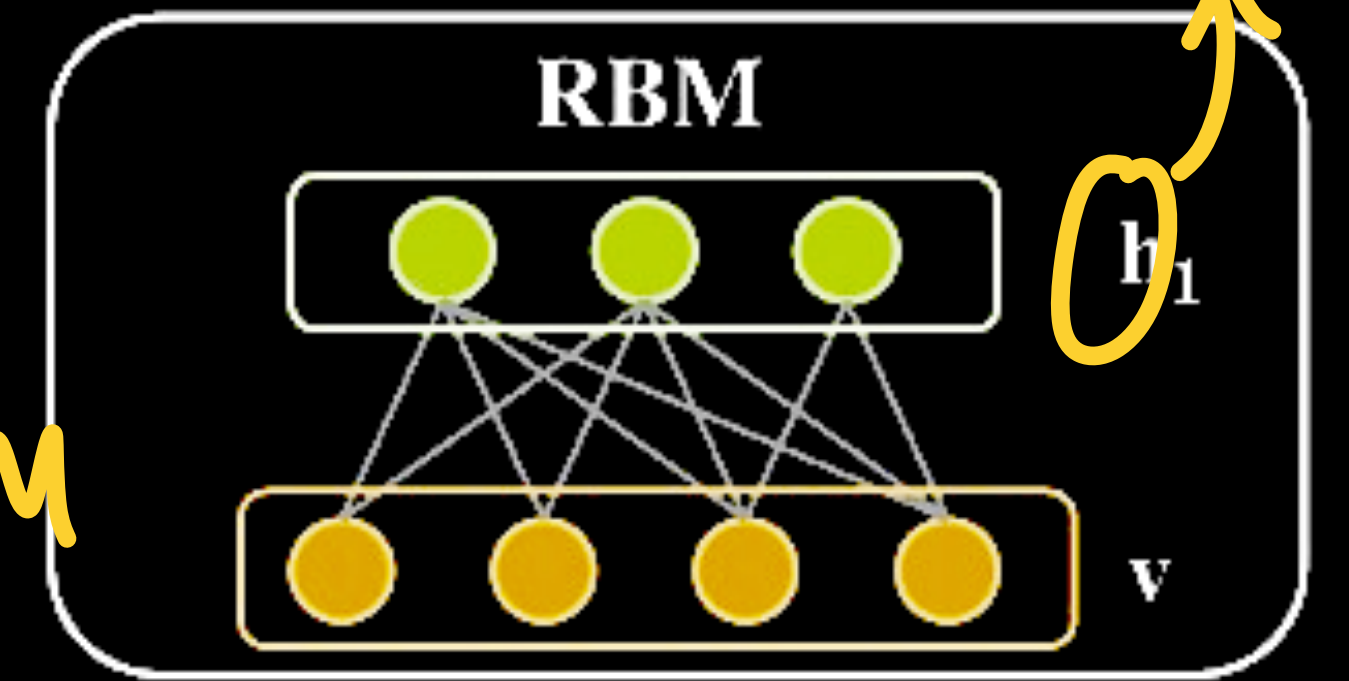
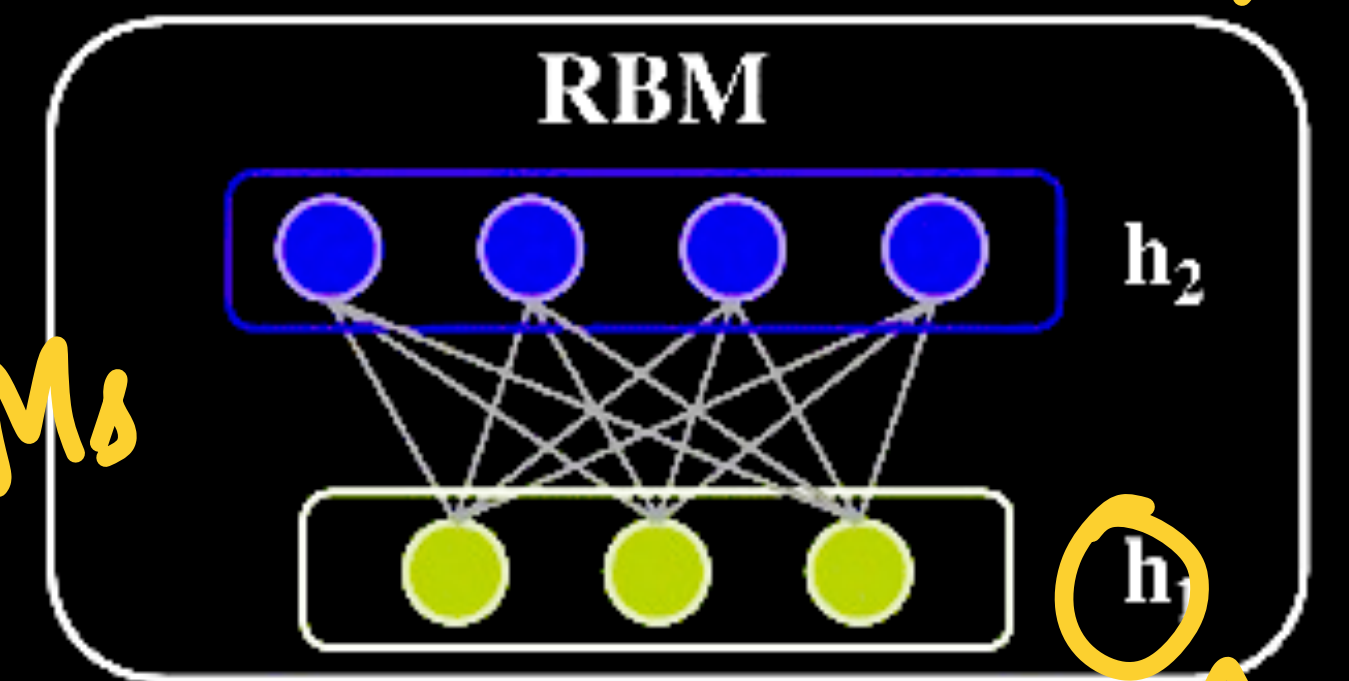
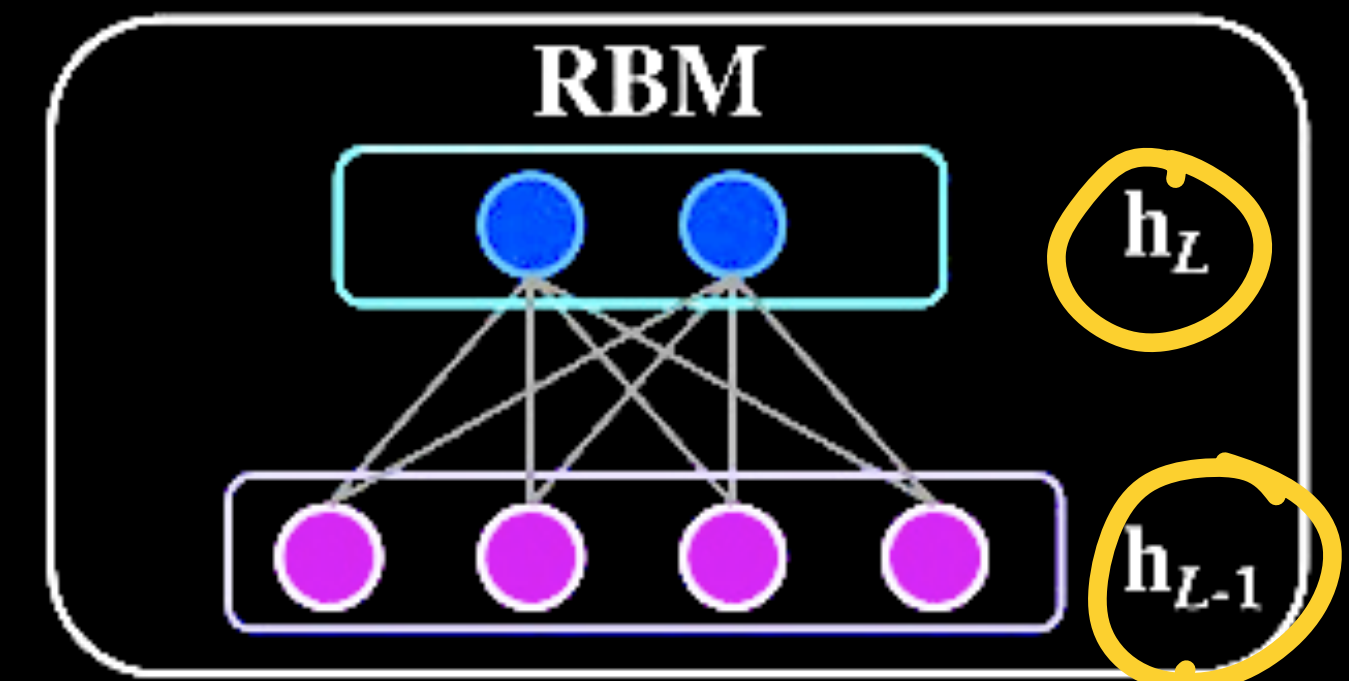
○ Weights are frozen each layer before training the next layer.

✓ Ancestral sampling can be performed for data generation

★ Lossy sample generation due to accumulation of errors.

* Most common use - pre-training of DNNs.

RBM

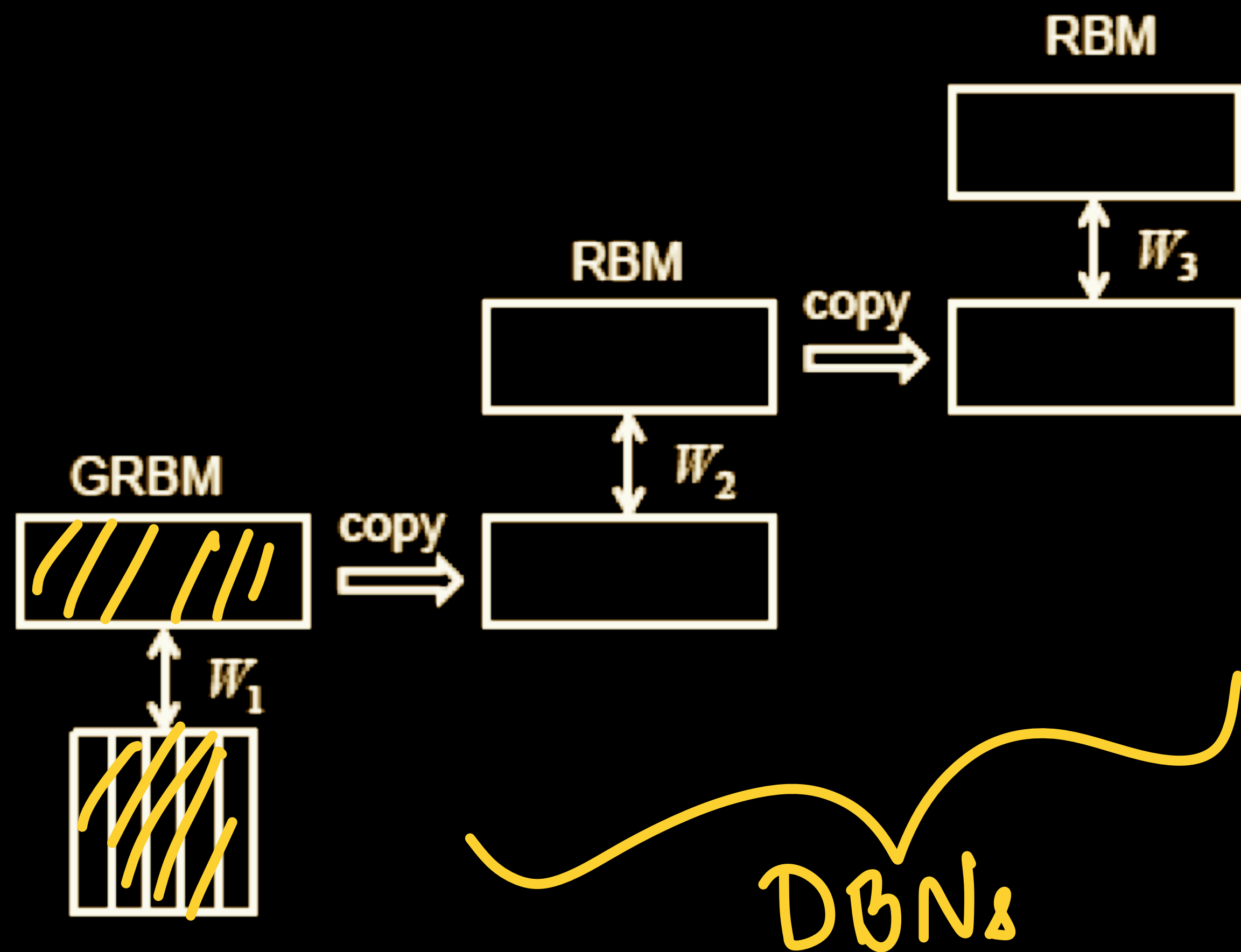


RBM's

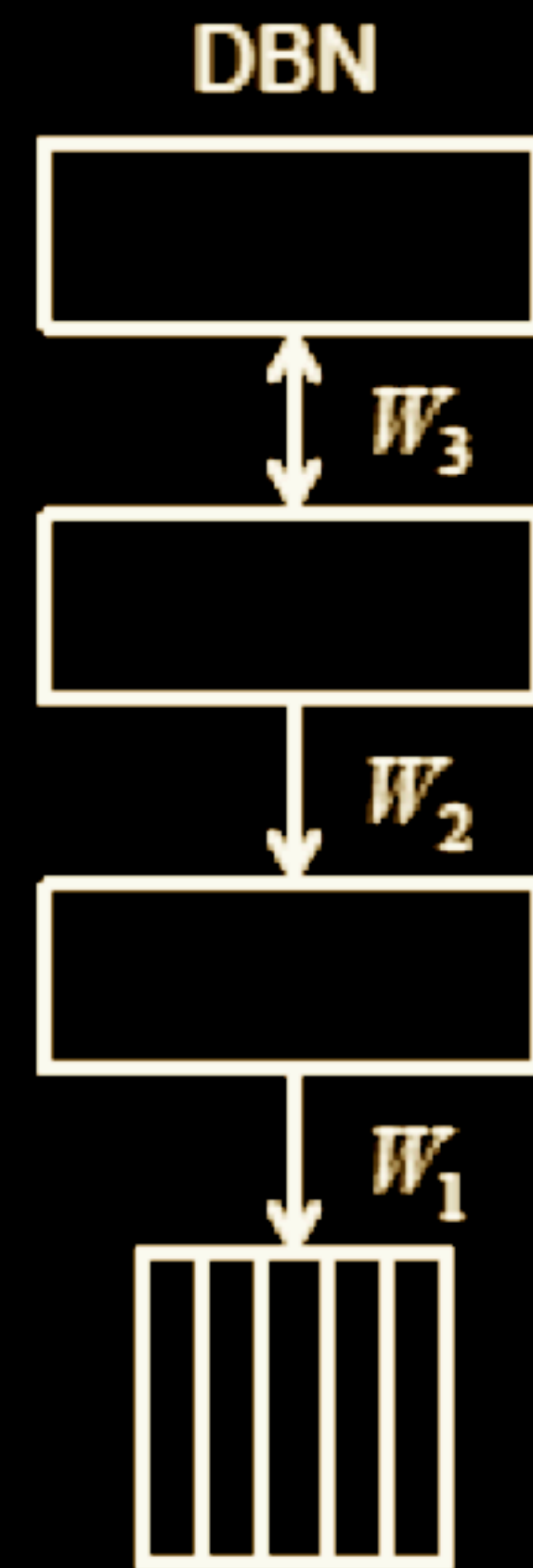
G RBM

DBNs for initialization

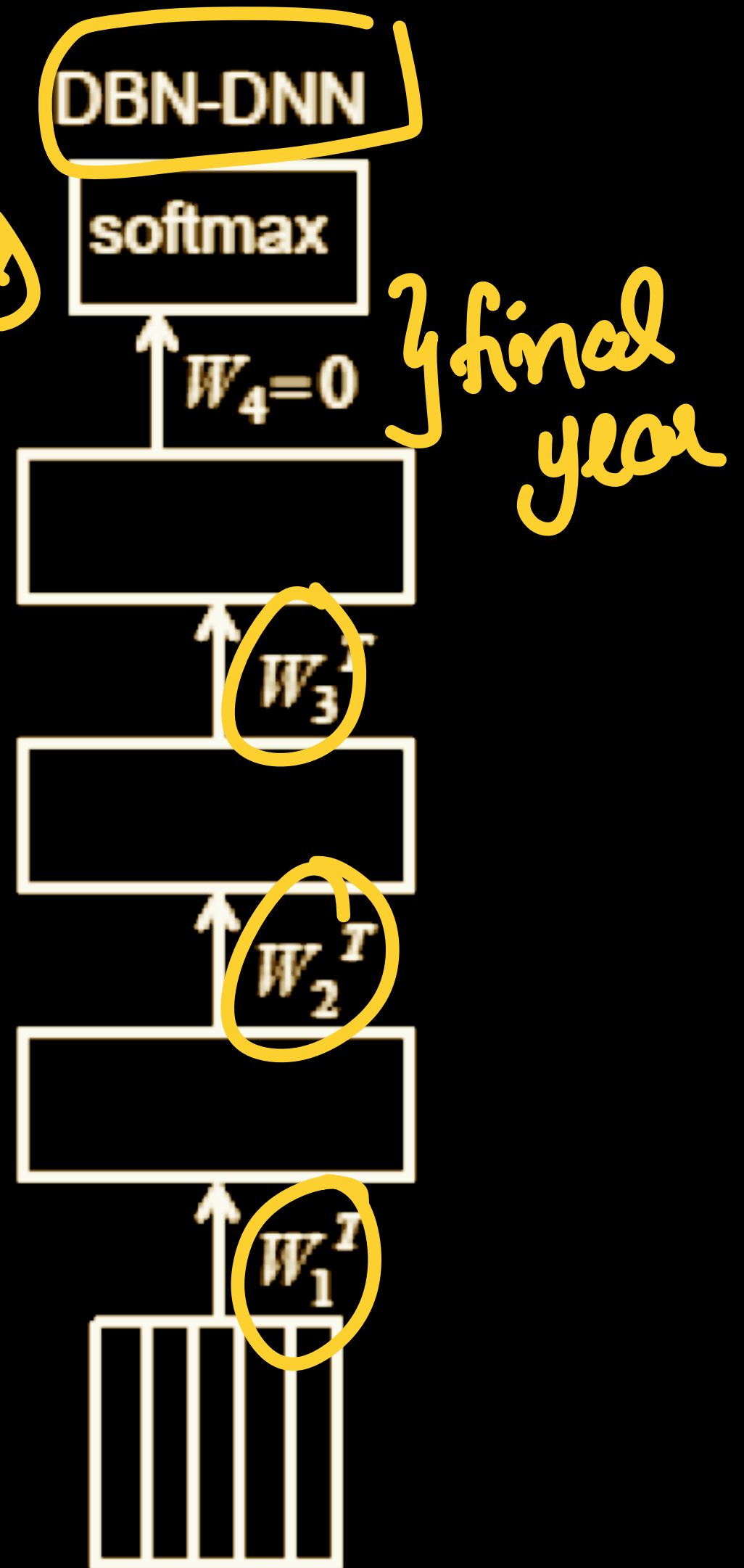
Unsupervised pre-training



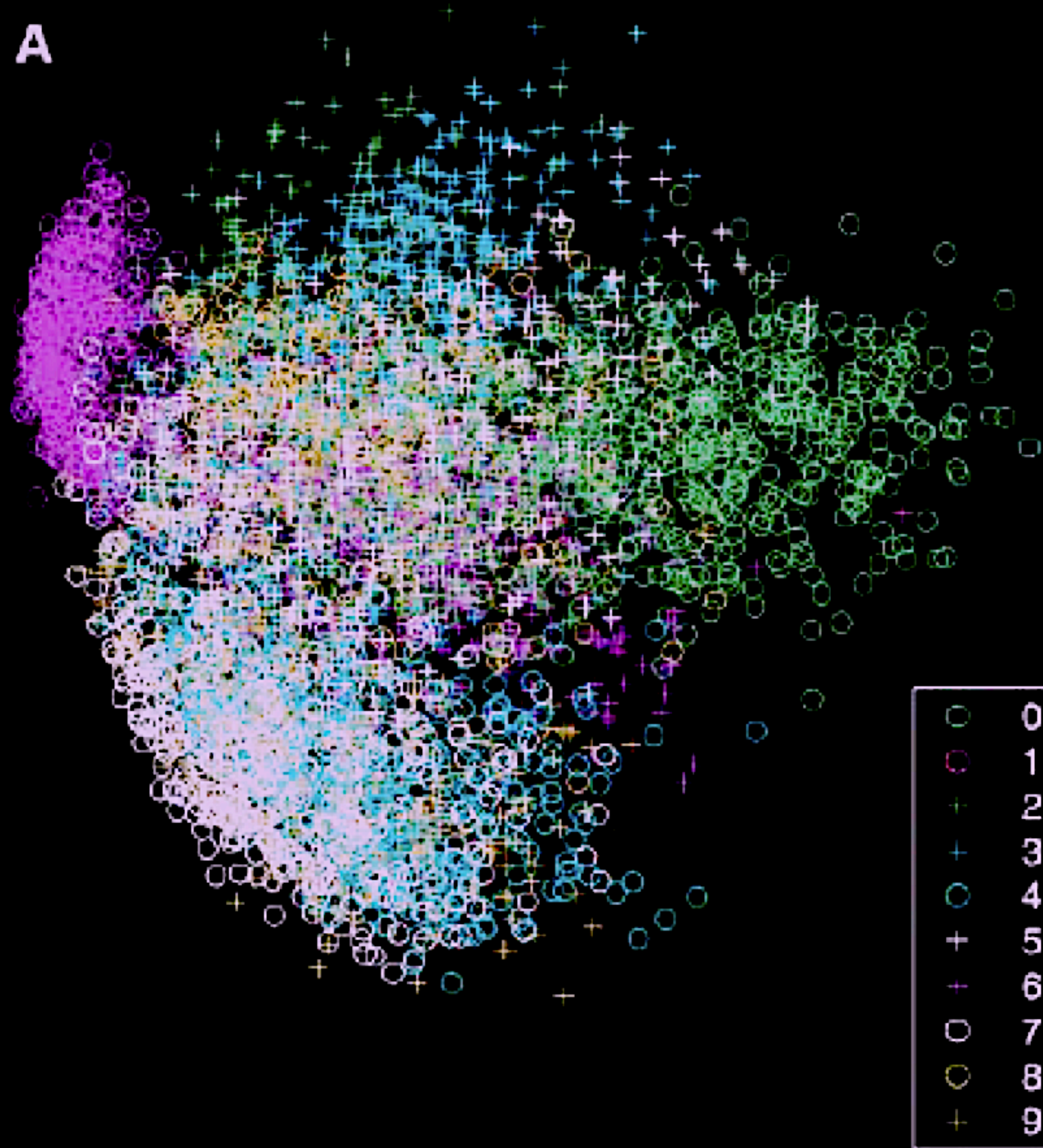
5k ←



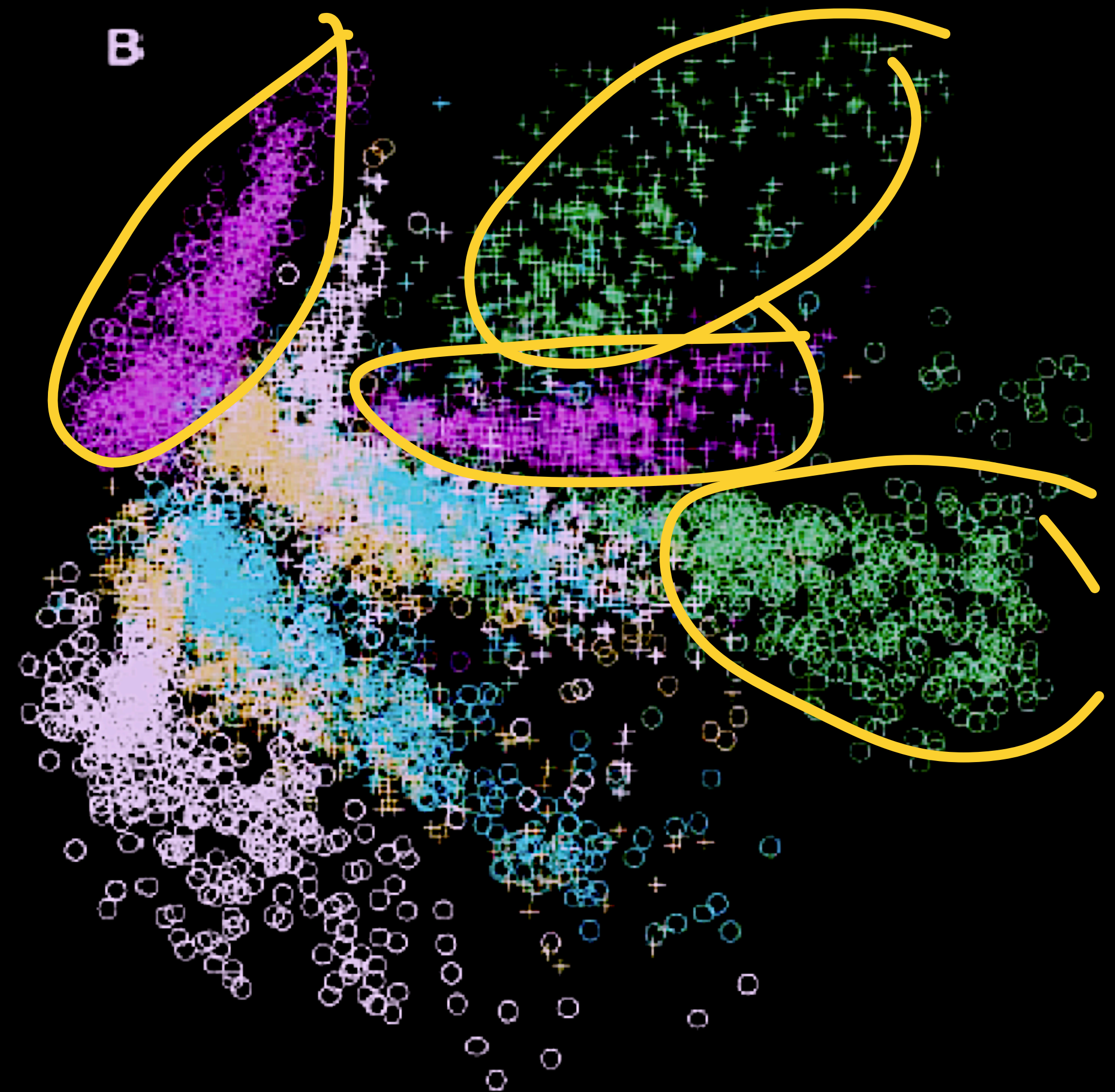
fine tuning



DBNs for visualization



PCA



RBM

t-SNE, PCA and DBNs

Unsupervised ✓

Unseen data
[PCA, DBN]

Neighborhood [t-SNE]

Hierarchy based [DBNs]
↑
binary

Distribution based [t-SNE, DBNs]

Data generation? [DBNs]

Initialization

* Generative modeling

Expected

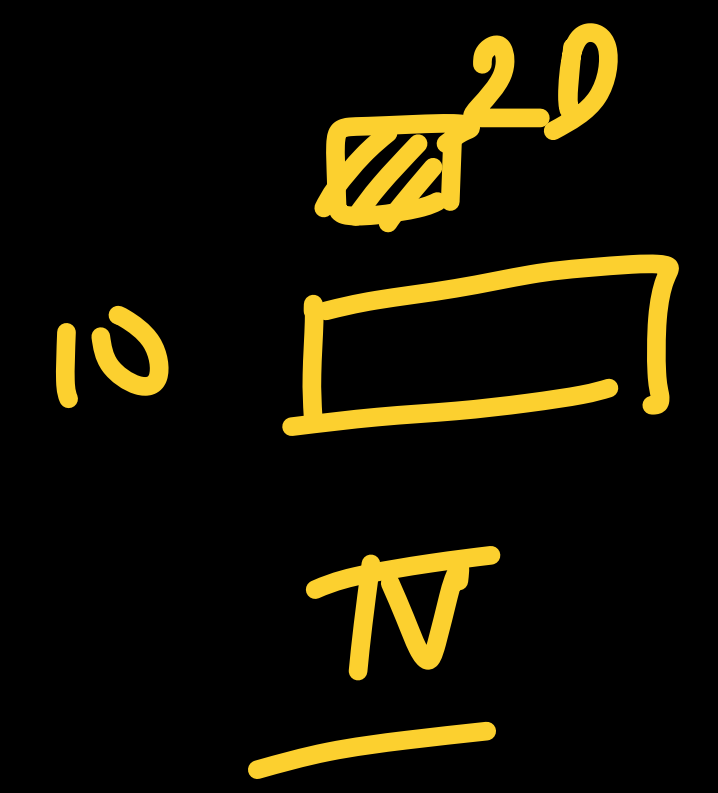
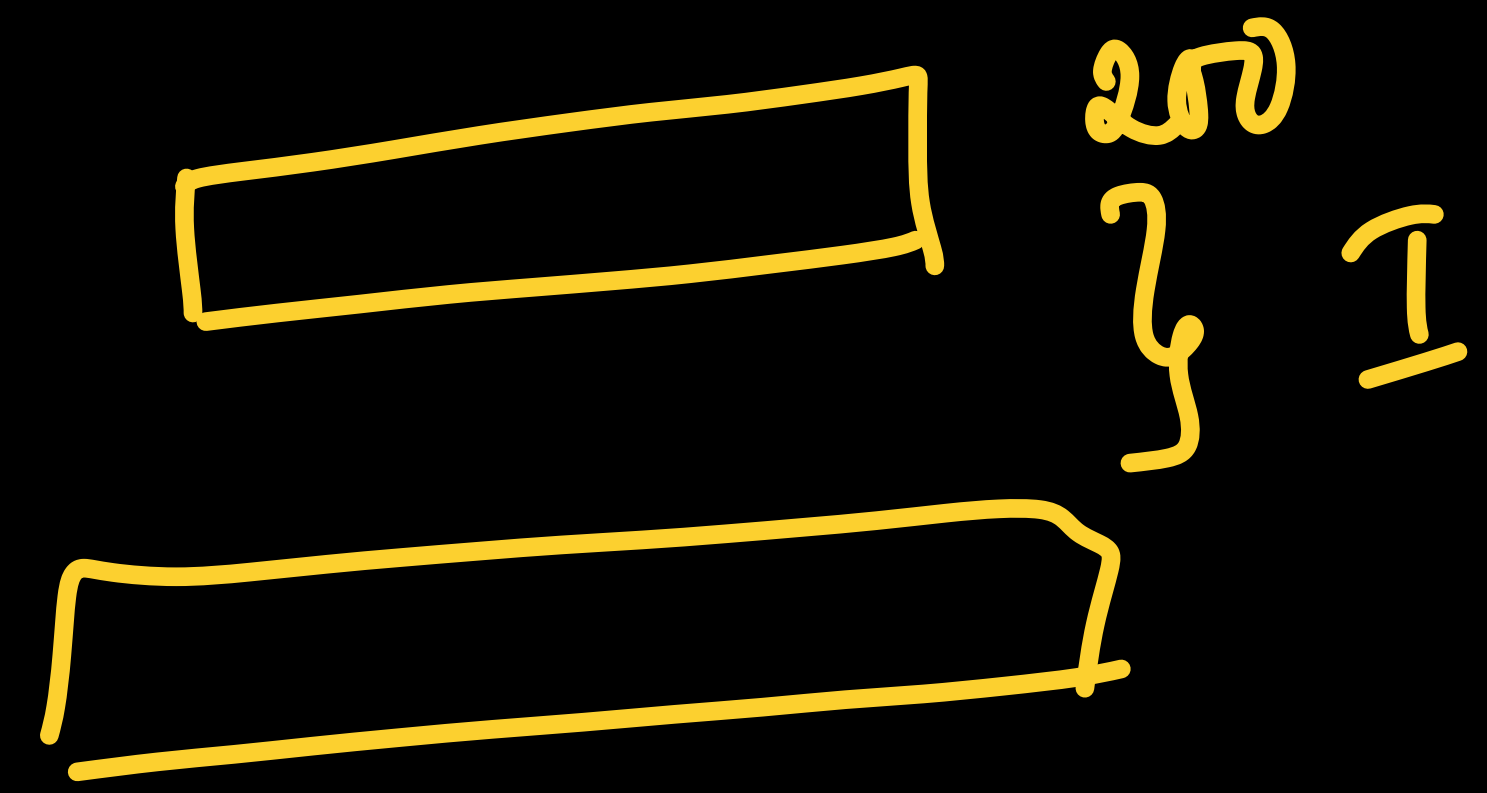
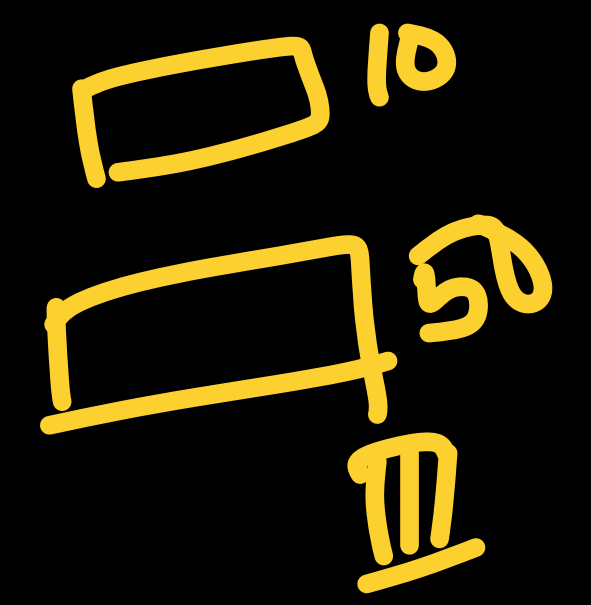
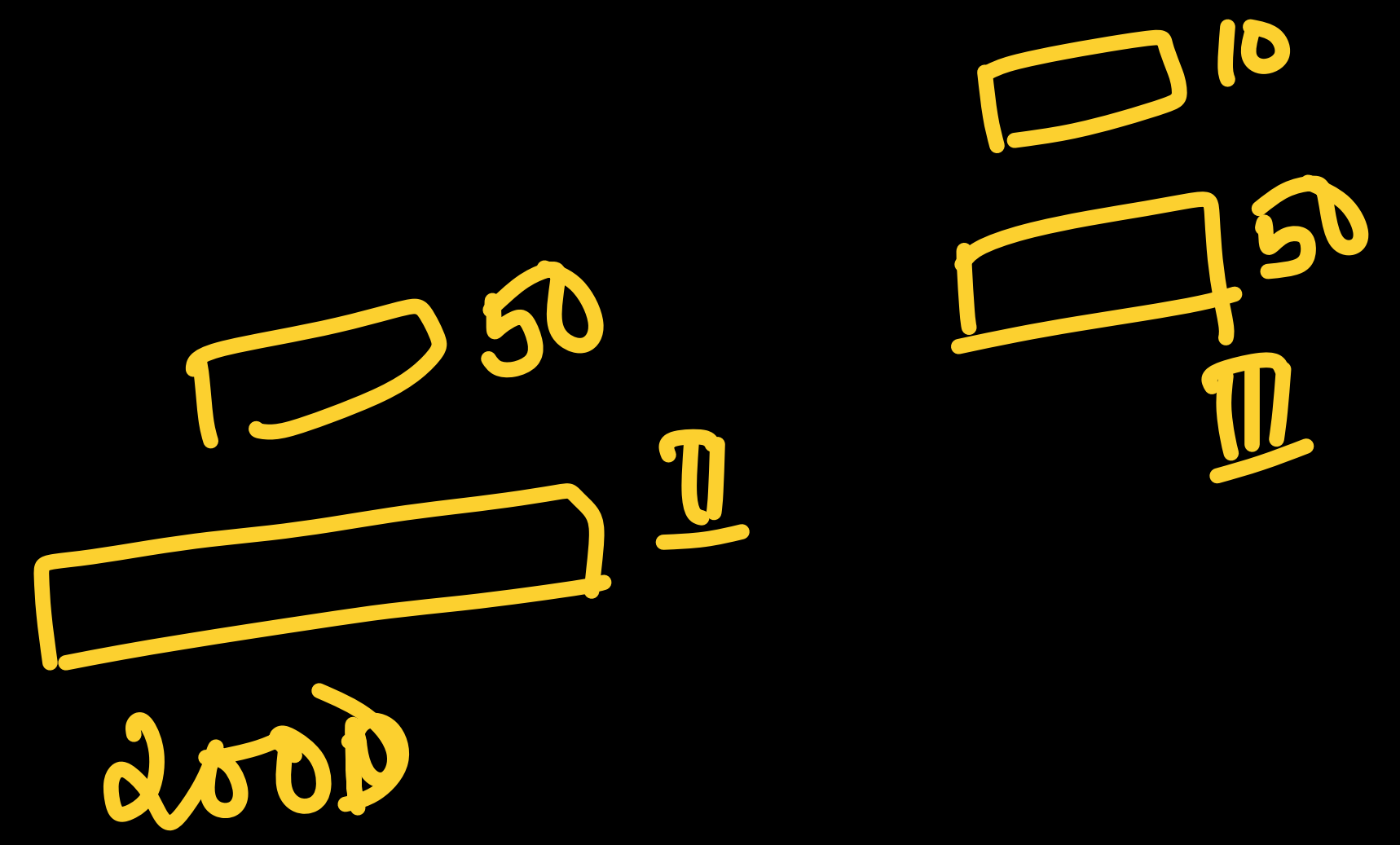
* \wedge Hidden units conditioned on visible



Affine + sigmoid [DNN]



[Hinton, 2006]



MNIST - 784

More reading

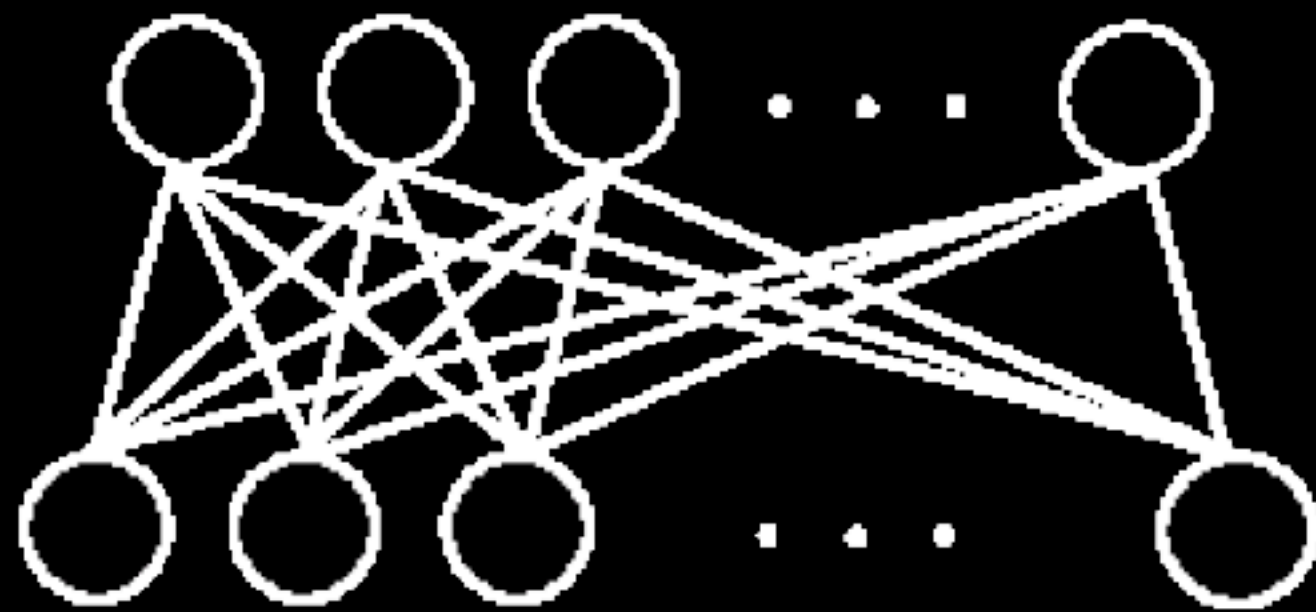
Salakhutdinov, Ruslan, Andriy Mnih, and Geoffrey Hinton. "Restricted Boltzmann machines for collaborative filtering." *Proceedings of the 24th international conference on Machine learning*. 2007.



Data generation using RBMs

learning

hidden units



visible units

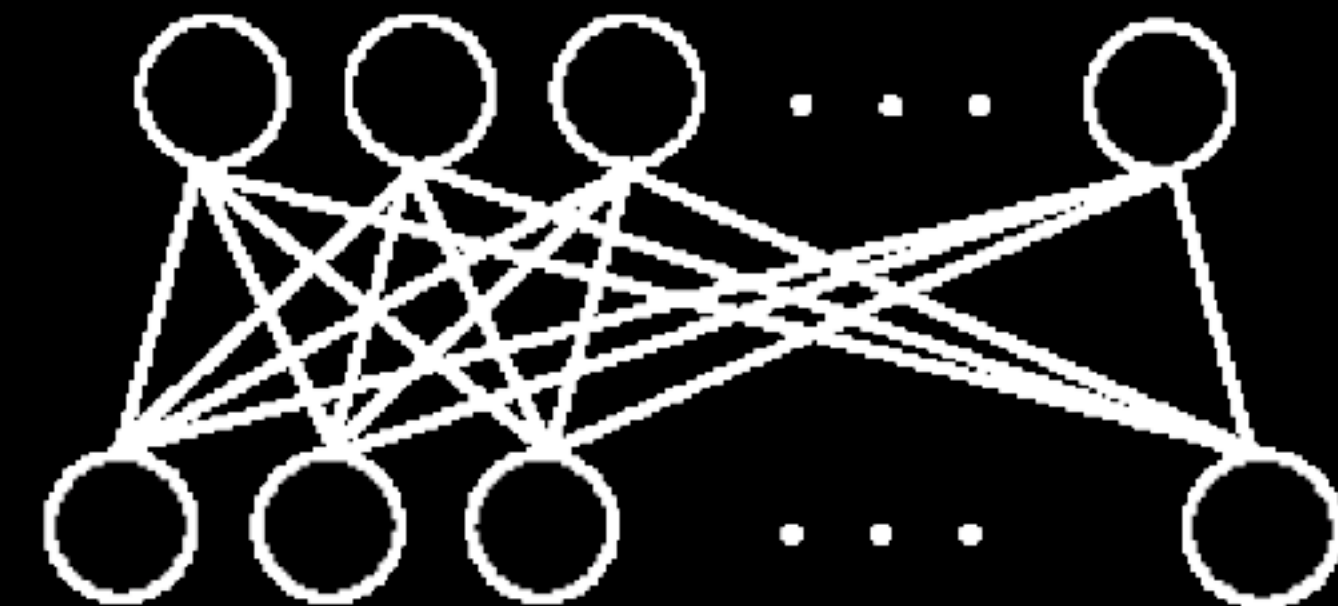
parameter fitting



training data

generating

hidden units



visible units

sampling



samples

Data generation using RBMs



Source — <https://lme.tf.fau.de/lecture-notes/lecture-notes-in-deep-learning-unsupervised-learning-part-1/>



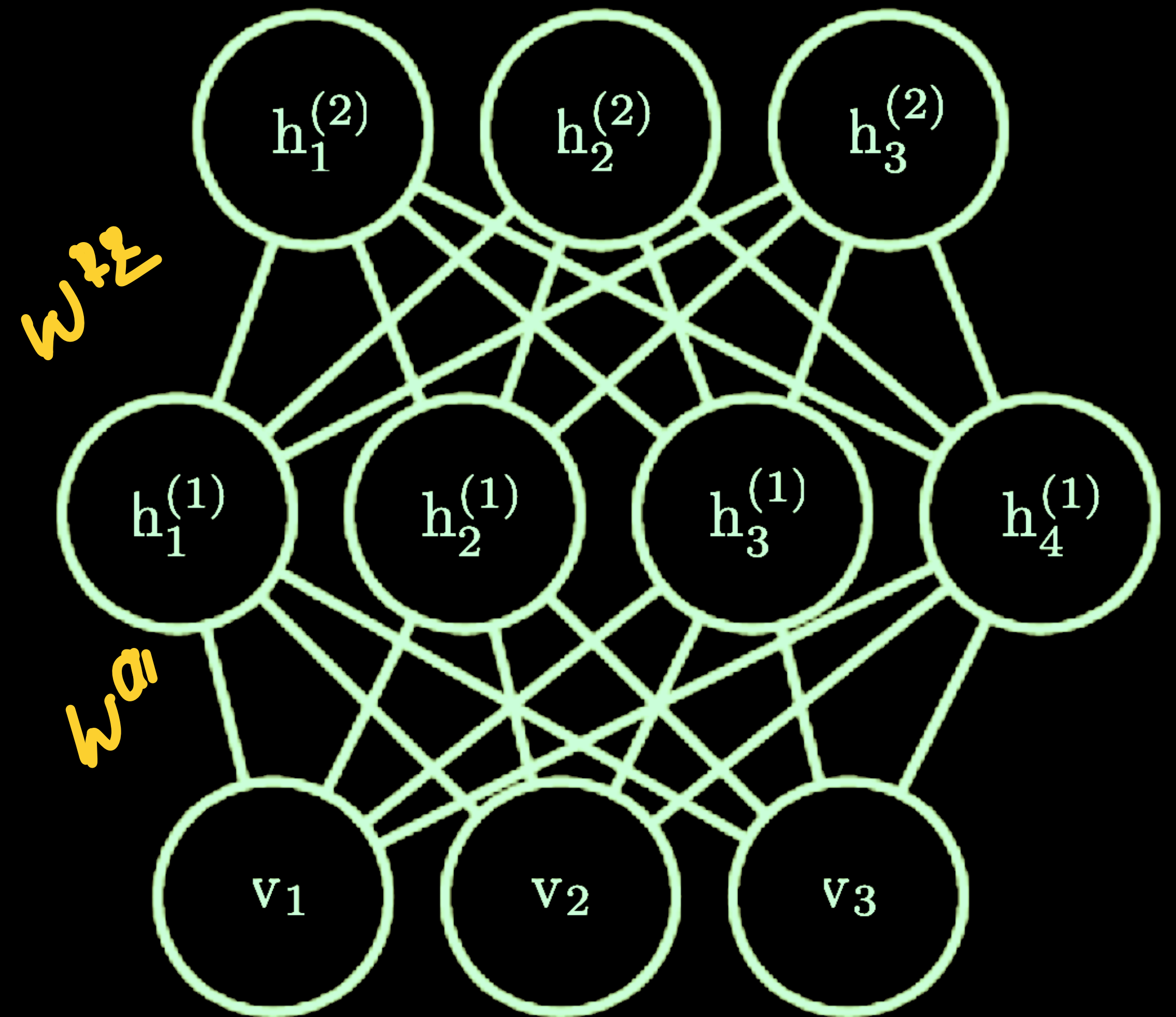
Deep Boltzmann machine (DBMs)

- * Deep layers of connections with RBM structure.

- Joint energy function.

- Undirected graph

- * Training and inference are more involved.



Autoencoders

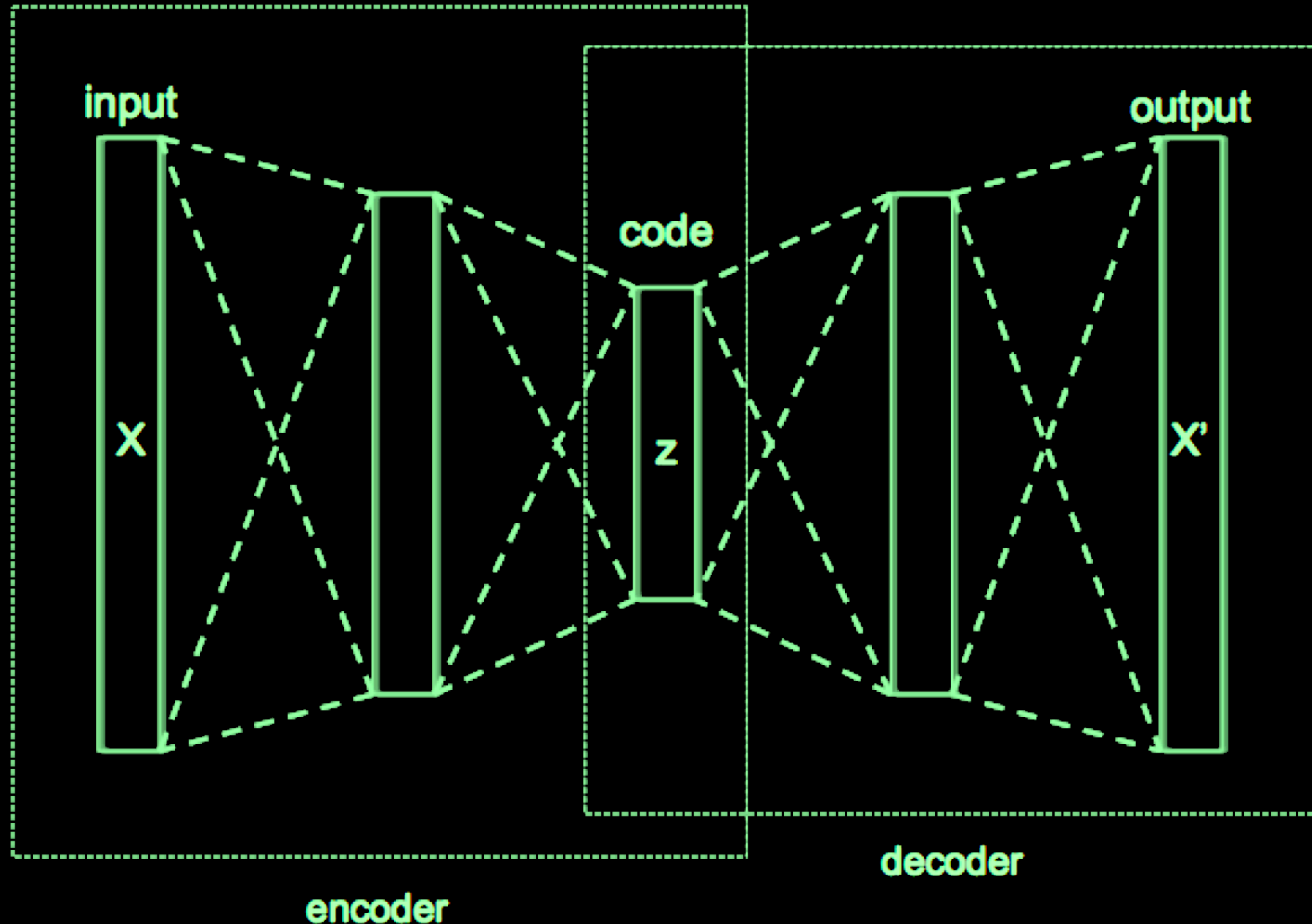
* Encoder decoder

→ Latent representations capture a lower dimensional embedding of the data.

✓ can be feedforward or convolutional layers.

* Model training

→ Using a reconstruction loss.



Variational auto encoders

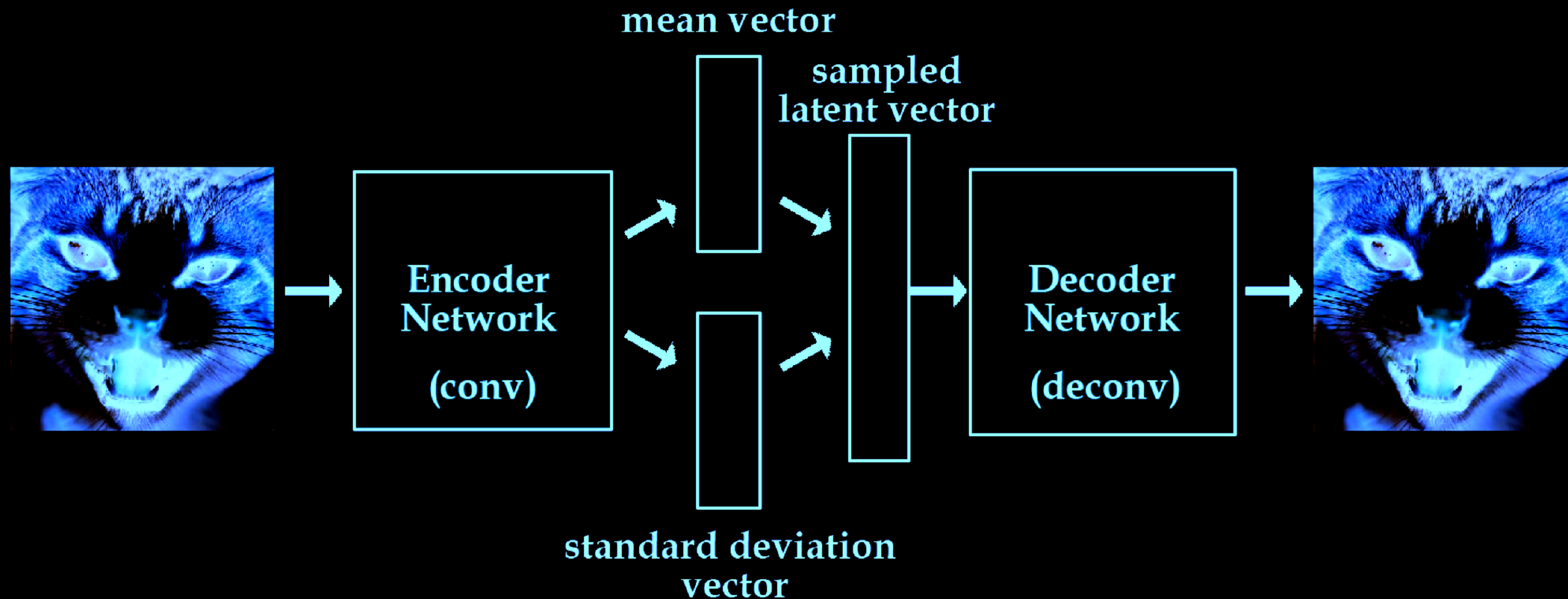
VAEs v/s RBMs

* Sampling from the latent space.

(VAE) latent variables \rightarrow real valued

\rightarrow Making Gaussian assumptions at the latent layer

(RBM) latent \rightarrow binary



Variational auto encoders

* The data \mathbf{x} and latent variable \mathbf{z}

* The forward model

- ✓ Sample the latent variable $p_{\theta}(\mathbf{z})$
- ✓ Sample the data given the latent $p_{\theta}(\mathbf{x}|\mathbf{z})$

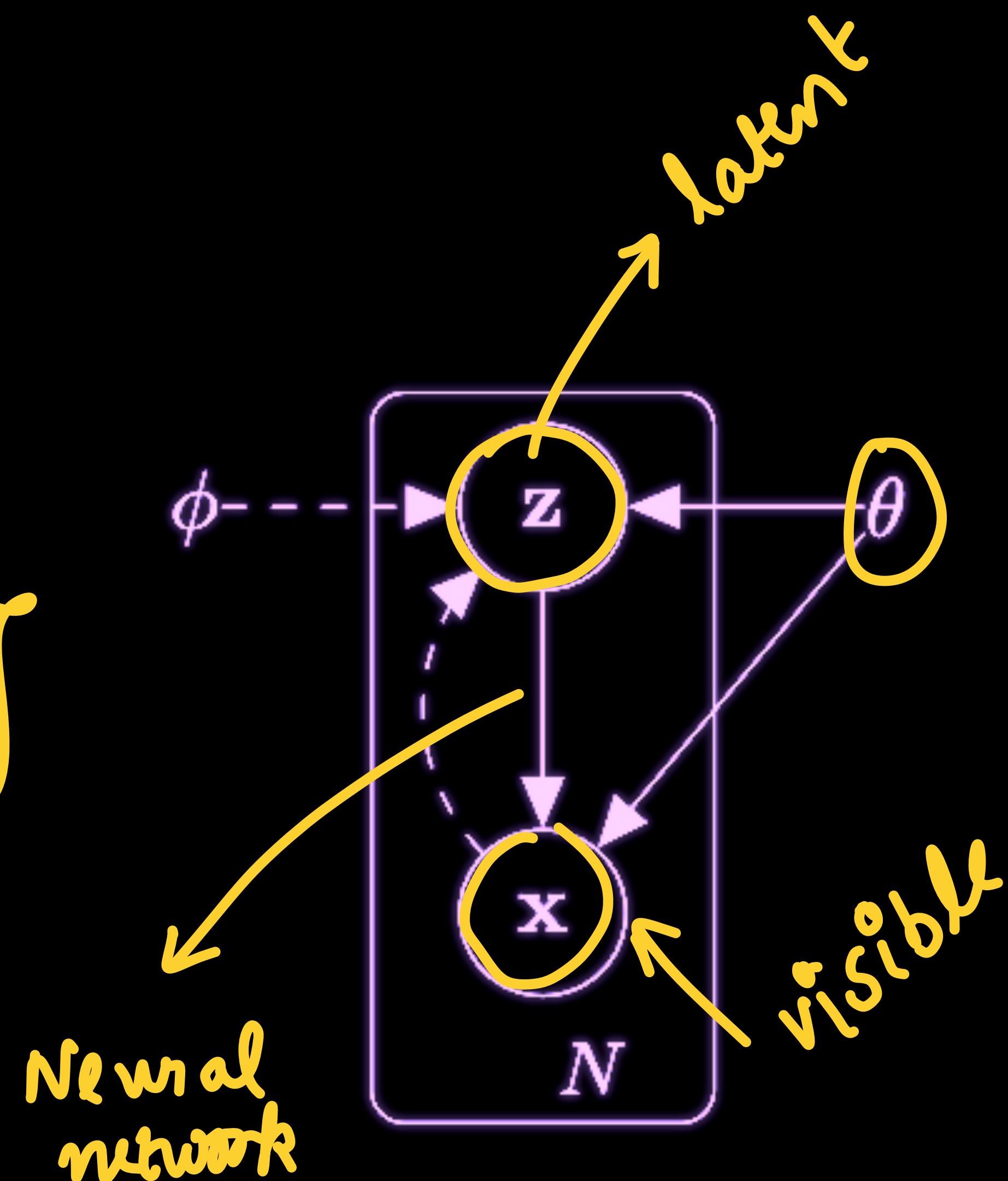
* The marginal distribution $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})d\mathbf{z}$

- ✓ maybe intractable

* The posterior distribution $p_{\theta}(\mathbf{z}|\mathbf{x})$

- may also be intractable

backward



Variational auto encoders

Variational lower bound

* Approximate the posterior

$$q_{\phi}(z|x) \sim p_{\theta}(z|x)$$

$$p_{\theta}(z|x)$$

approx.

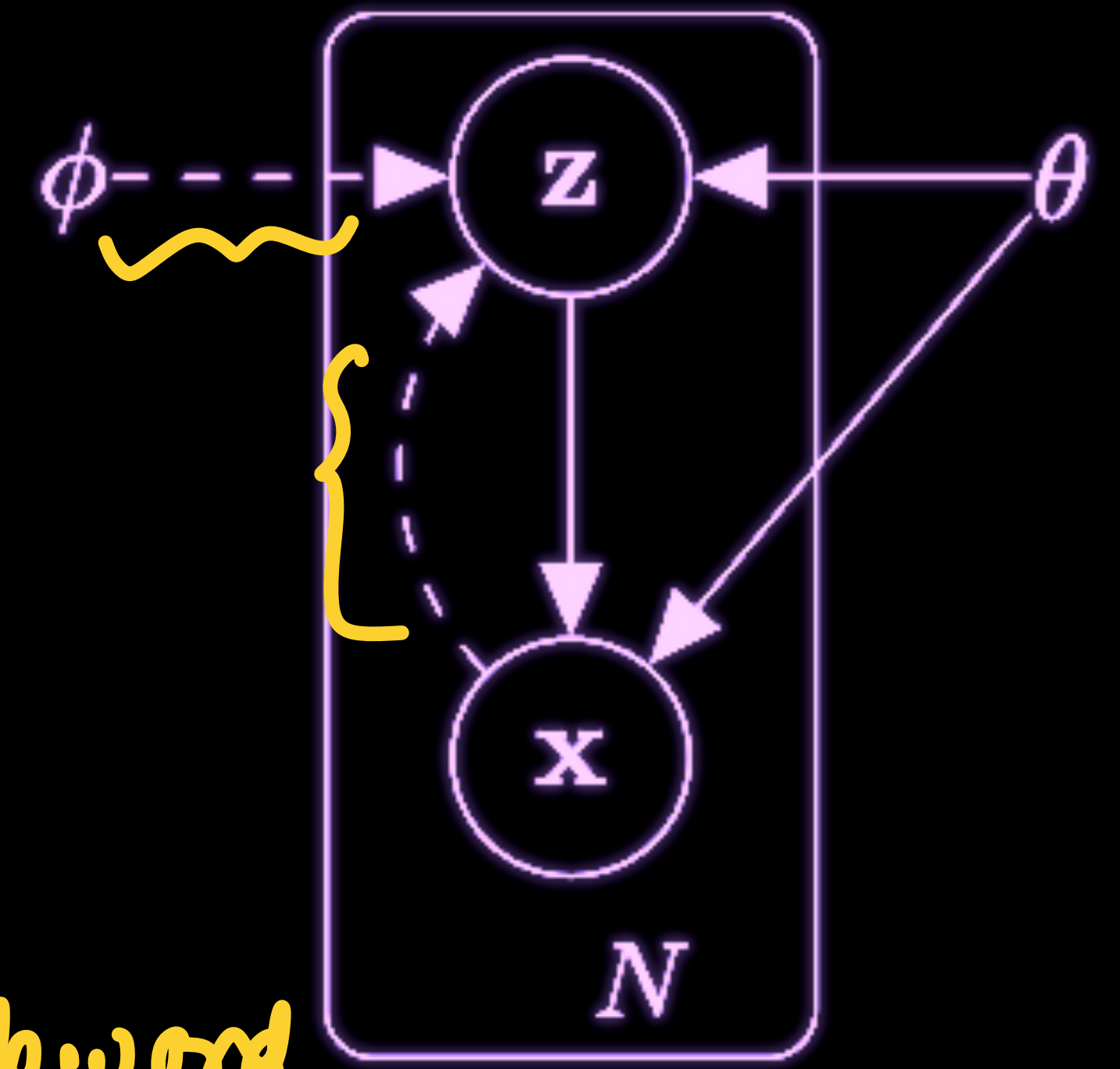
→ Using variational lower bound.

A neural network Θ → forward model

$$z \rightarrow x$$

A second neural network Φ → backward model

$$x \rightarrow z$$



Variational lower bound

$$\begin{aligned} \mathcal{L} &= \log p_{\theta}(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\ &= \log \int p_{\theta}(\mathbf{x}, \mathbf{z}) \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \log \left(\mathbb{E}_q \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right) \\ &\geq \mathbb{E}_q \left(\log \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right) \\ &= \mathbb{E}_q \left(\log \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right) \\ &= \mathbb{E}_q \left(\log \left[p_{\theta}(\mathbf{x}, \mathbf{z}) \right] \right) + H_q(\mathbf{z}|\mathbf{x}) \end{aligned}$$

$$\frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})}$$

Jensen's inequality

$$- \int_{\mathbf{z}} \log q_{\phi} \cdot q_{\phi} d\mathbf{z} = H_q$$



Variational lower bound

$$\begin{aligned}KL [q(Z) || p(Z|X)] &= \int_Z q(Z) \log \frac{q(Z)}{p(Z|X)} \\&= - \int_Z q(Z) \log \frac{p(Z|X)}{q(Z)} \\&= - \left(\int_Z q(Z) \log \frac{p(X, Z)}{q(Z)} - \int_Z q(Z) \log p(X) \right) \\&= - \int_Z q(Z) \log \frac{p(X, Z)}{q(Z)} + \log p(X) \int_Z q(Z) \\&= -L + \log p(X)\end{aligned}$$

