

Training Restricted Boltzmann Machines

GRAHAM TAYLOR

VECTOR INSTITUTE

SCHOOL OF ENGINEERING
UNIVERSITY OF GUELPH

CANADIAN INSTITUTE
FOR ADVANCED RESEARCH

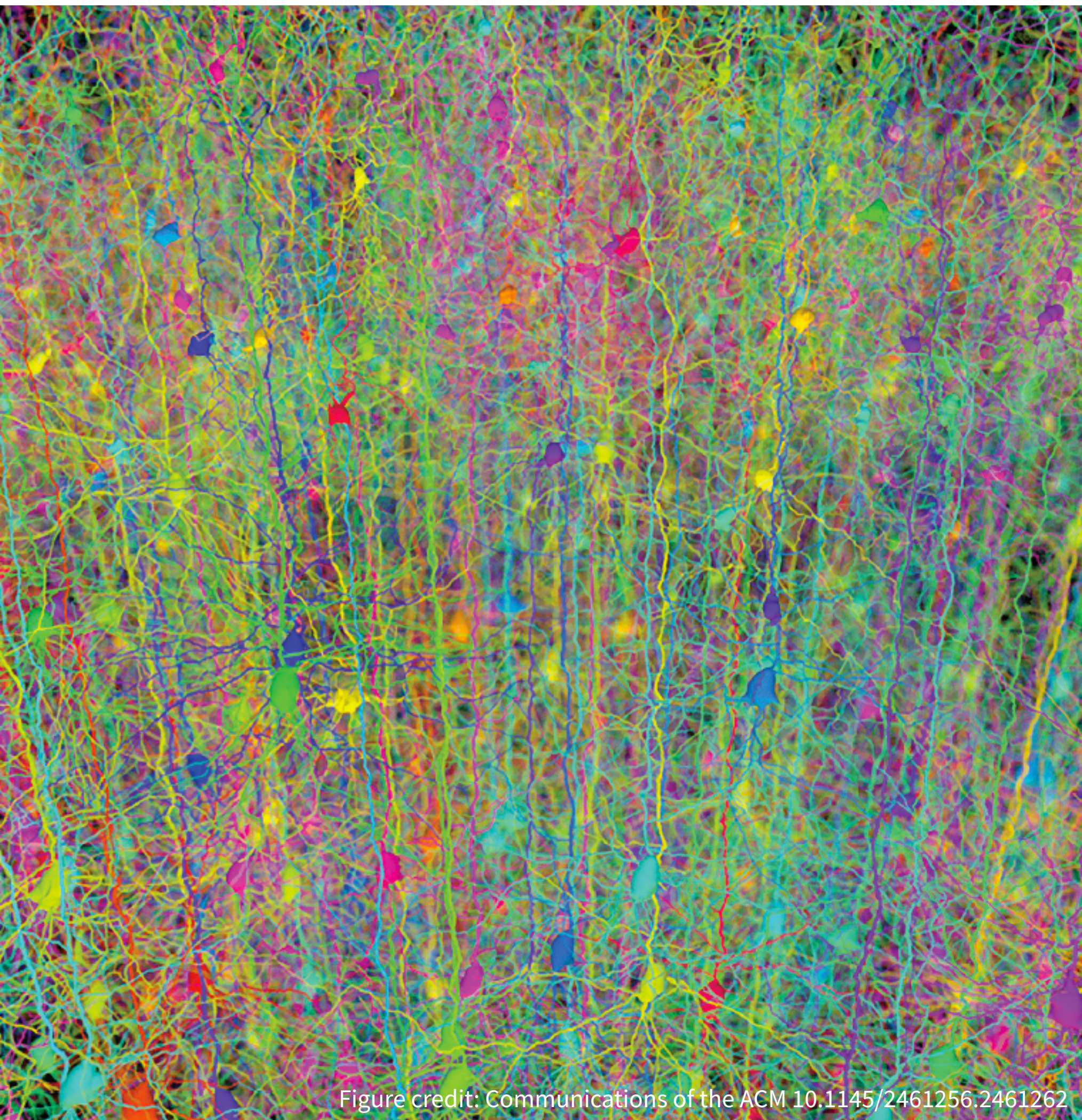


Figure credit: Communications of the ACM 10.1145/2461256.2461262

UNIVERSITY
of GUELPH

CHANGING LIVES
IMPROVING LIFE

CIFAR
CANADIAN
INSTITUTE
FOR
ADVANCED
RESEARCH

Training RBMs

- To train an RBM, we would like to minimize the average negative log likelihood:

$$\frac{1}{n} \sum_{i=1}^n -\log P(\mathbf{v}^{(i)})$$

- We'd then like to proceed by stochastic gradient descent:

$$\frac{\partial -\log P(\mathbf{v}^{(i)})}{\partial \boldsymbol{\theta}} = \underbrace{\mathbb{E}_{\mathbf{h}} \left[\frac{\partial E(\mathbf{v}^{(i)}, \mathbf{h})}{\partial \boldsymbol{\theta}} \mid \mathbf{v}^{(i)} \right]}_{\text{positive phase}} - \underbrace{\mathbb{E}_{\mathbf{v}, \mathbf{h}} \left[\frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right]}_{\text{negative phase}}$$

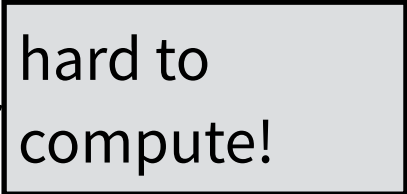
Training RBMs

- To train an RBM, we would like to minimize the average negative log likelihood:

$$\frac{1}{n} \sum_{i=1}^n -\log P(\mathbf{v}^{(i)})$$

- We'd then like to proceed by stochastic gradient descent:

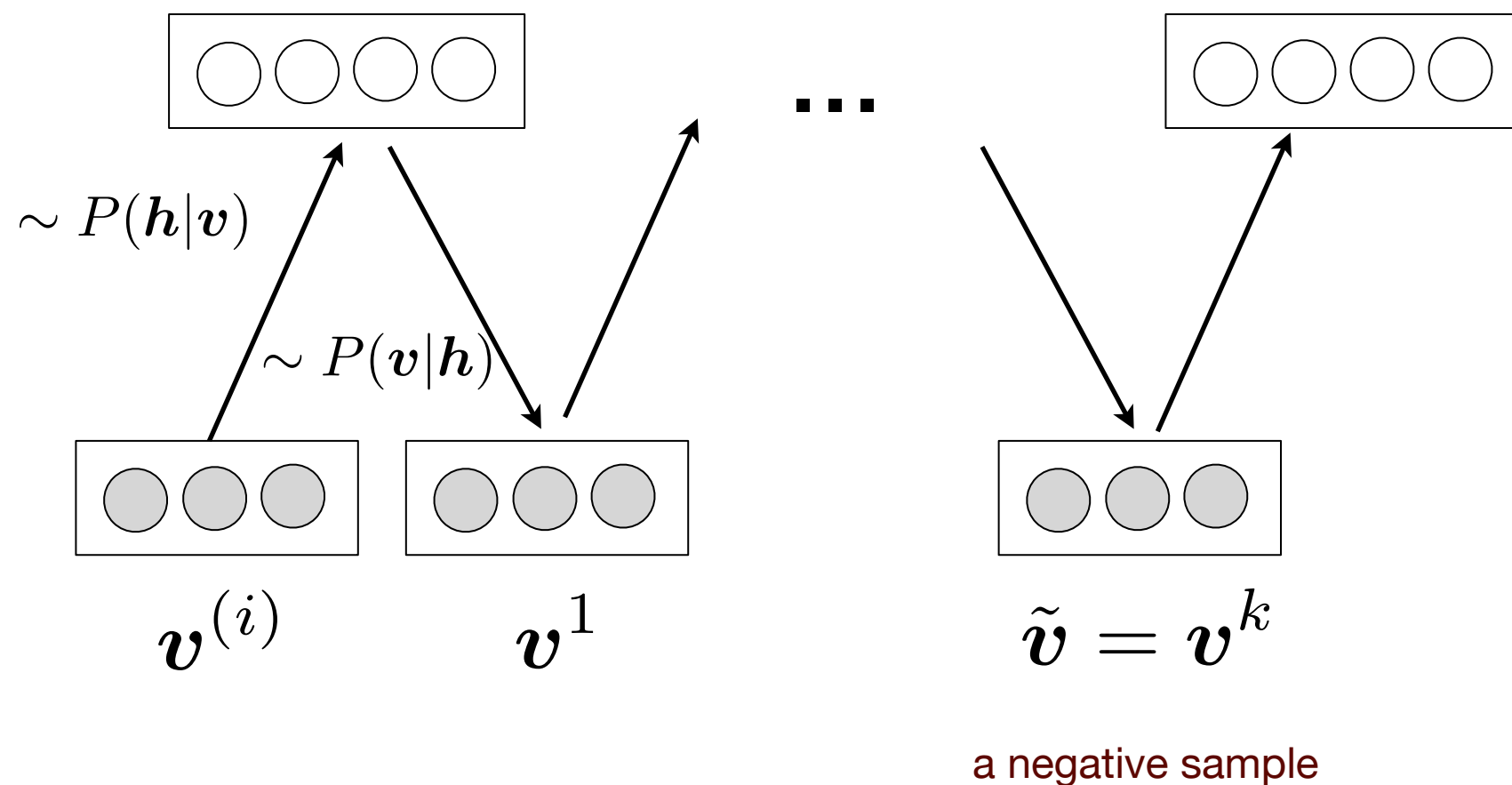
$$\frac{\partial -\log P(\mathbf{v}^{(i)})}{\partial \boldsymbol{\theta}} = \underbrace{\mathbb{E}_{\mathbf{h}} \left[\frac{\partial E(\mathbf{v}^{(i)}, \mathbf{h})}{\partial \boldsymbol{\theta}} \mid \mathbf{v}^{(i)} \right]}_{\text{positive phase}} - \underbrace{\mathbb{E}_{\mathbf{v}, \mathbf{h}} \left[\frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right]}_{\text{negative phase}}$$



Contrastive Divergence (CD)

Idea:

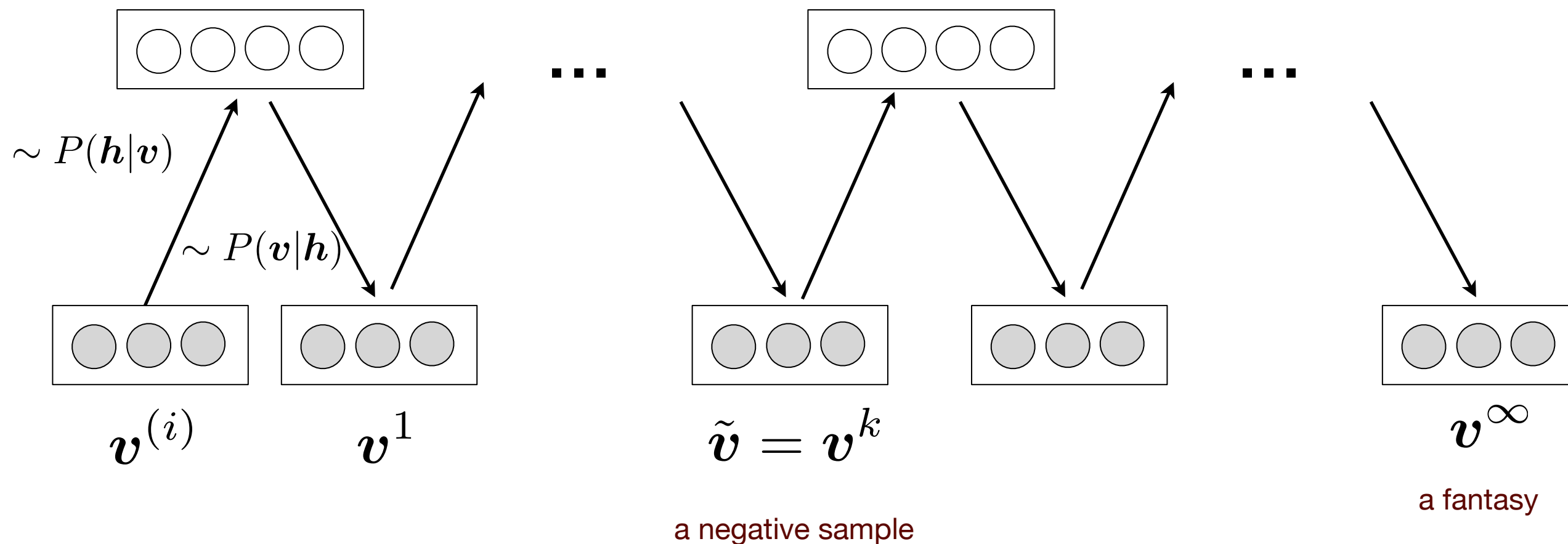
- replace the expectation by a point estimate at \tilde{v}
- obtain the point \tilde{v} by Gibbs sampling
- start sampling the chain at $\tilde{v}^{(i)}$



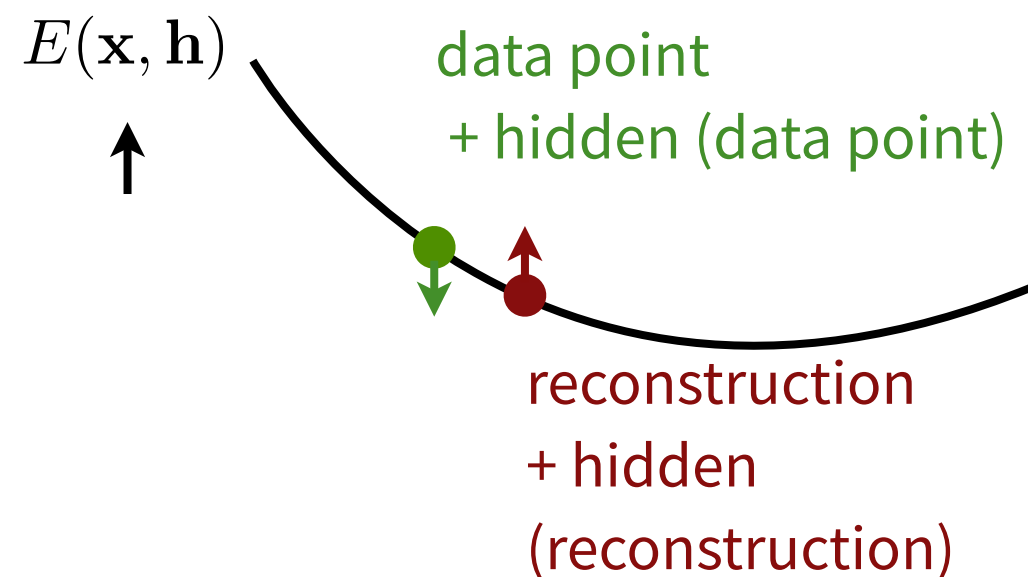
Contrastive Divergence (CD)

Idea:

- replace the expectation by a point estimate at \tilde{v}
- obtain the point \tilde{v} by Gibbs sampling
- start sampling the chain at $\tilde{v}^{(i)}$



Contrastive Divergence (A picture)



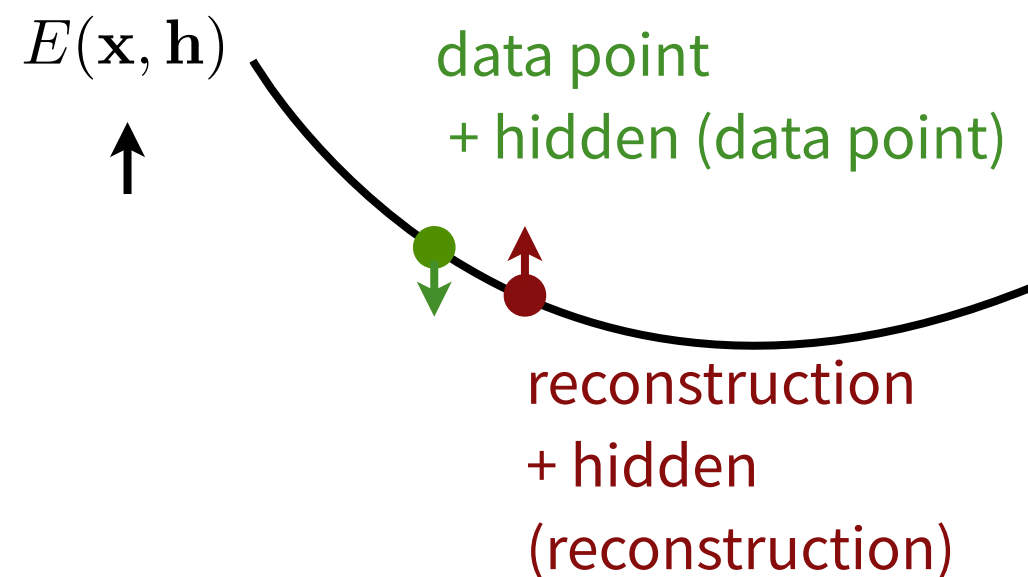
$$\mathbb{E}_{\mathbf{h}} \left[\frac{\partial E(\mathbf{v}^{(i)}, \mathbf{h})}{\partial \boldsymbol{\theta}} \mid \mathbf{v}^{(i)} \right] \approx \frac{\partial E(\mathbf{v}^{(i)}, \tilde{\mathbf{h}}^{(i)})}{\partial \boldsymbol{\theta}}$$

Change the weights to pull the energy down at the data point

$$\mathbb{E}_{\mathbf{v}, \mathbf{h}} \left[\frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \approx \frac{\partial E(\tilde{\mathbf{v}}, \tilde{\mathbf{h}})}{\partial \boldsymbol{\theta}}$$

Change the weights to pull the energy up at the reconstruction

Contrastive Divergence (A picture)

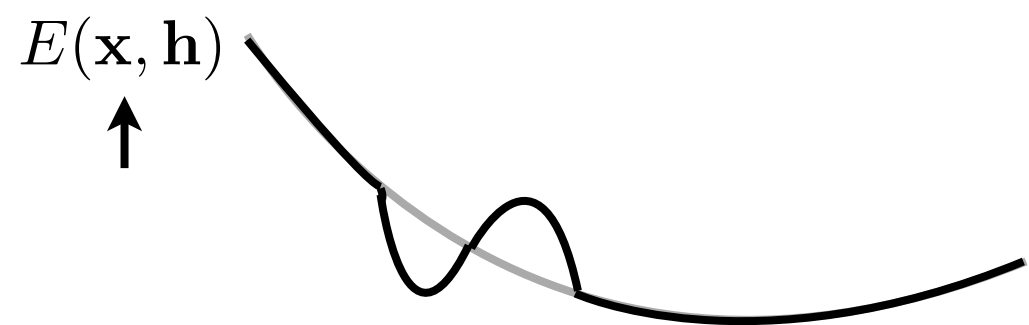


$$\mathbb{E}_{\mathbf{h}} \left[\frac{\partial E(\mathbf{v}^{(i)}, \mathbf{h})}{\partial \boldsymbol{\theta}} \mid \mathbf{v}^{(i)} \right] \approx \frac{\partial E(\mathbf{v}^{(i)}, \tilde{\mathbf{h}}^{(i)})}{\partial \boldsymbol{\theta}}$$

Change the weights to pull the energy down at the data point

$$\mathbb{E}_{\mathbf{v}, \mathbf{h}} \left[\frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \approx \frac{\partial E(\tilde{\mathbf{v}}, \tilde{\mathbf{h}})}{\partial \boldsymbol{\theta}}$$

Change the weights to pull the energy up at the reconstruction



Derivation of the Learning Rule

Given $\mathbf{v}^{(i)}$ and $\tilde{\mathbf{v}}$ the learning rule for $\theta = \mathbf{W}$ becomes:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \epsilon \left(\nabla_{\mathbf{W}} \left\{ -\log P(\mathbf{v}^{(i)}) \right\} \right) \\ &= \mathbf{W} - \epsilon \left(\mathbb{E}_{\mathbf{h}} \left[\nabla_{\mathbf{W}} E(\mathbf{v}^{(i)}, \mathbf{h}) \mid \mathbf{v}^{(i)} \right] - \mathbb{E}_{\mathbf{v}, \mathbf{h}} [\nabla_{\mathbf{W}} E(\mathbf{v}, \mathbf{h})] \right) \\ &\approx \mathbf{W} - \epsilon \left(\mathbb{E}_{\mathbf{h}} \left[\nabla_{\mathbf{W}} E(\mathbf{v}^{(i)}, \mathbf{h}) \mid \mathbf{v}^{(i)} \right] - \mathbb{E}_{\mathbf{h}} [\nabla_{\mathbf{W}} E(\tilde{\mathbf{v}}, \mathbf{h}) \mid \tilde{\mathbf{v}}] \right) \\ &= \mathbf{W} + \epsilon \left(\mathbf{h}(\mathbf{v}^{(i)}) \mathbf{v}^{(i)\top} - \mathbf{h}(\tilde{\mathbf{v}}) \tilde{\mathbf{v}}^\top \right)\end{aligned}$$

CD-k: Algorithm

1. For each training example $\mathbf{v}^{(i)}$
 - i. generate a negative sample $\tilde{\mathbf{v}}$ using k steps of Gibbs sampling, starting at $\mathbf{v}^{(i)}$

- ii. update parameters:

$$\mathbf{W} \leftarrow \mathbf{W} + \epsilon \left(\mathbf{h}(\mathbf{v}^{(i)}) \mathbf{v}^{(i)\top} - \mathbf{h}(\tilde{\mathbf{v}}) \tilde{\mathbf{v}}^\top \right)$$

$$\mathbf{c} \leftarrow \mathbf{c} + \epsilon \left(\mathbf{h}(\mathbf{v}^{(i)}) - \mathbf{h}(\tilde{\mathbf{v}}) \right)$$

$$\mathbf{b} \leftarrow \mathbf{b} + \epsilon \left(\mathbf{v}^{(i)} - \tilde{\mathbf{v}} \right)$$

2. Go back to 1 until stopping criteria

CD-k: Algorithm

1. For each training example $\mathbf{v}^{(i)}$
 - i. generate a negative sample $\tilde{\mathbf{v}}$ using k steps of Gibbs sampling, starting at $\mathbf{v}^{(i)}$

- ii. update parameters:

$$\mathbf{W} \leftarrow \mathbf{W} + \epsilon \left(\mathbf{h}(\mathbf{v}^{(i)}) \mathbf{v}^{(i)\top} - \mathbf{h}(\tilde{\mathbf{v}}) \tilde{\mathbf{v}}^\top \right)$$

$$\mathbf{c} \leftarrow \mathbf{c} + \epsilon \left(\mathbf{h}(\mathbf{v}^{(i)}) - \mathbf{h}(\tilde{\mathbf{v}}) \right)$$

$$\mathbf{b} \leftarrow \mathbf{b} + \epsilon \left(\mathbf{v}^{(i)} - \tilde{\mathbf{v}} \right)$$

2. Go back to 1 until stopping criteria

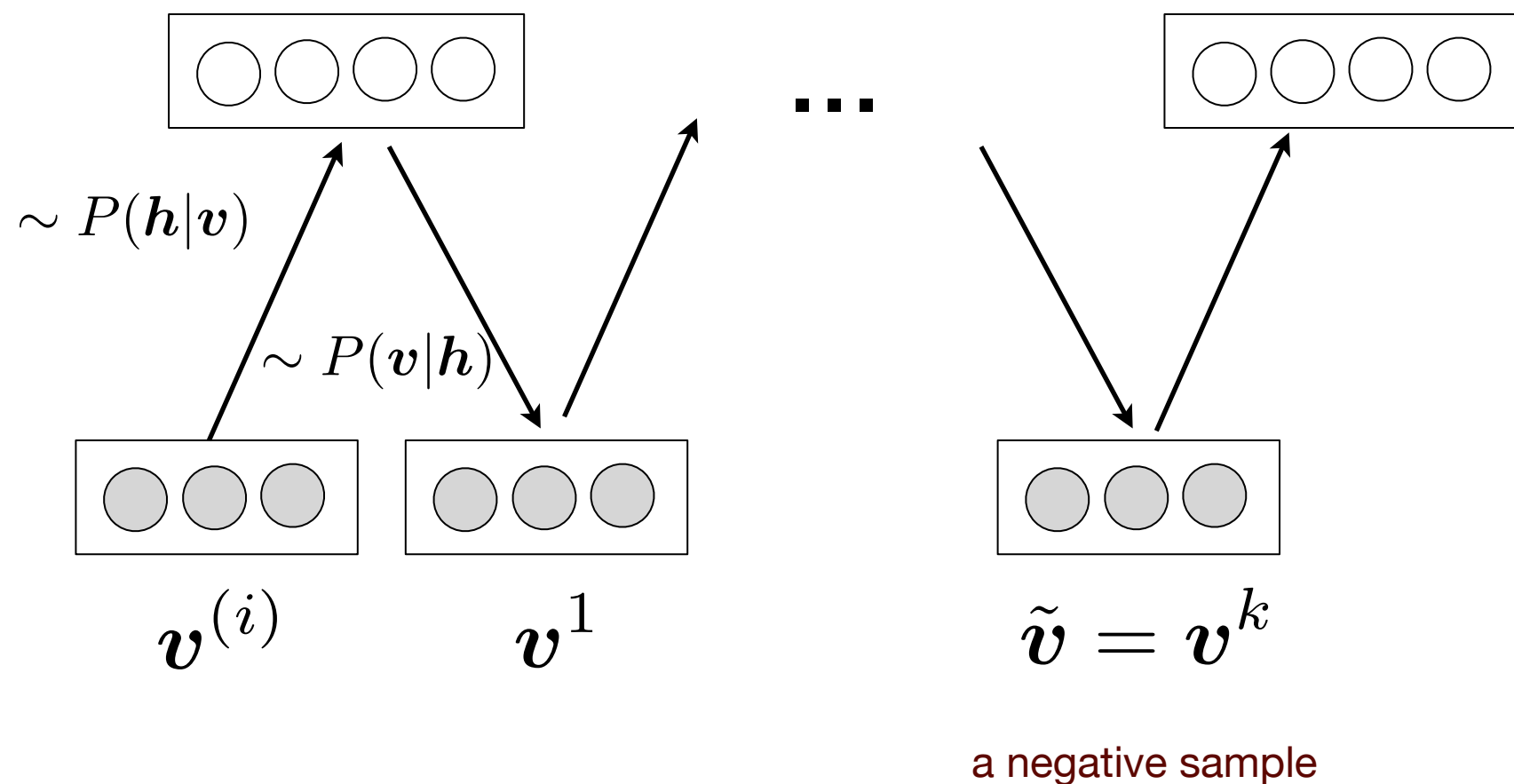
In general, the bigger k , the less **biased** the estimate of the gradient. In practice, $k = 1$ works well for pre-training.

Persistent CD (PCD)

(Tieleman 2008)

Idea:

- instead of initializing the chain to $\tilde{v}^{(i)}$, initialize the chain to the negative sample of the last iteration

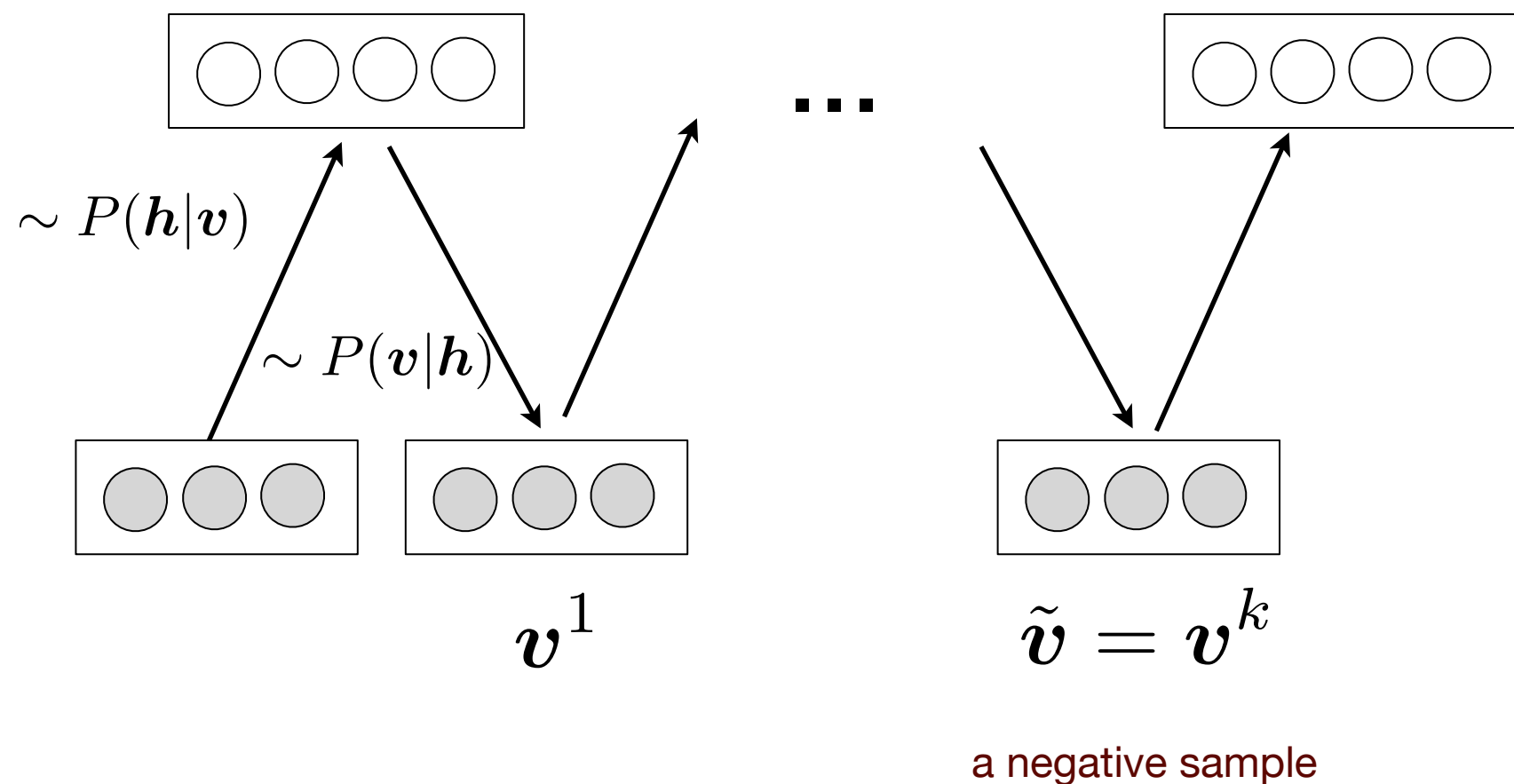


Persistent CD (PCD)

(Tieleman 2008)

Idea:

- instead of initializing the chain to $\tilde{v}^{(i)}$, initialize the chain to the negative sample of the last iteration

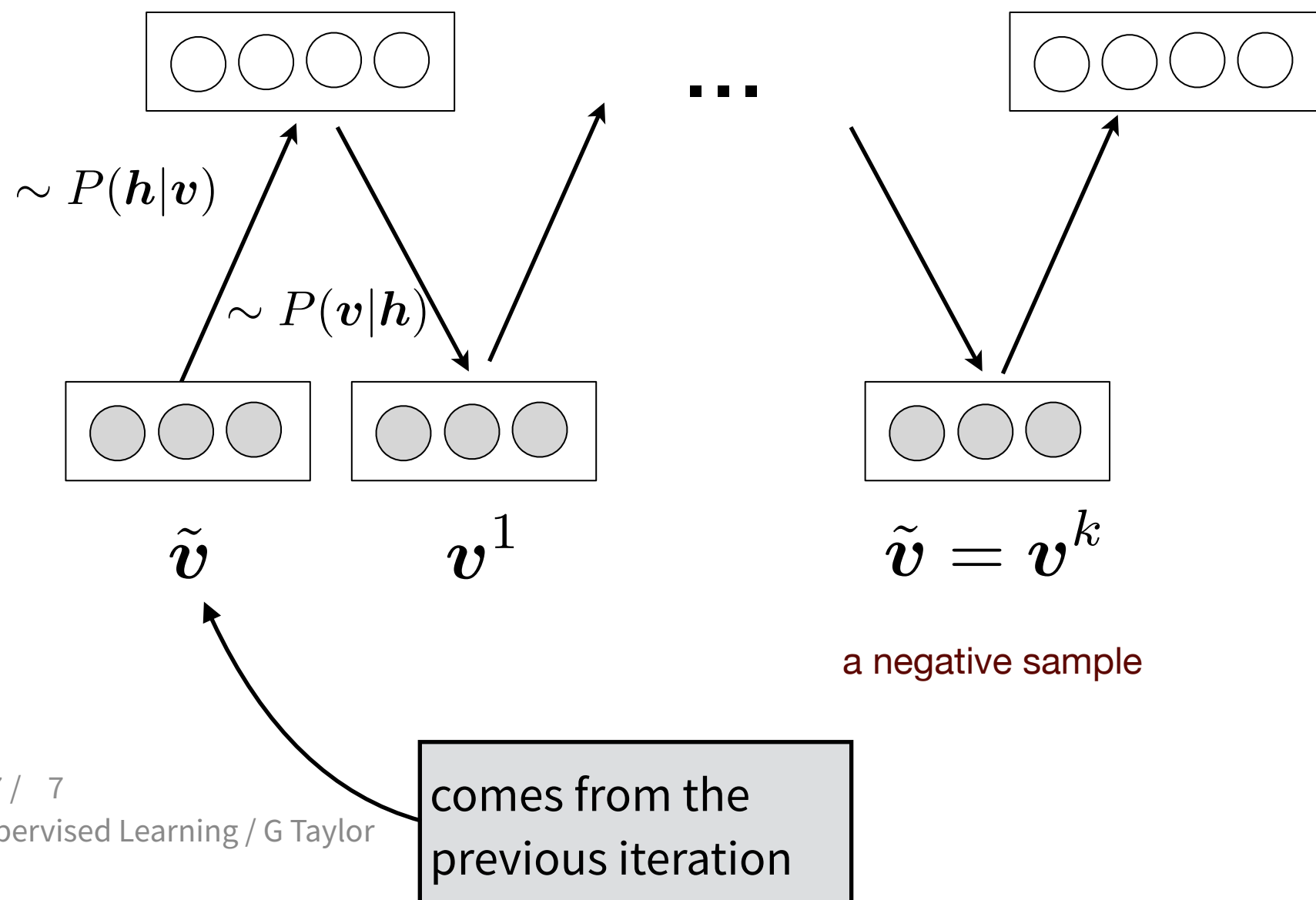


Persistent CD (PCD)

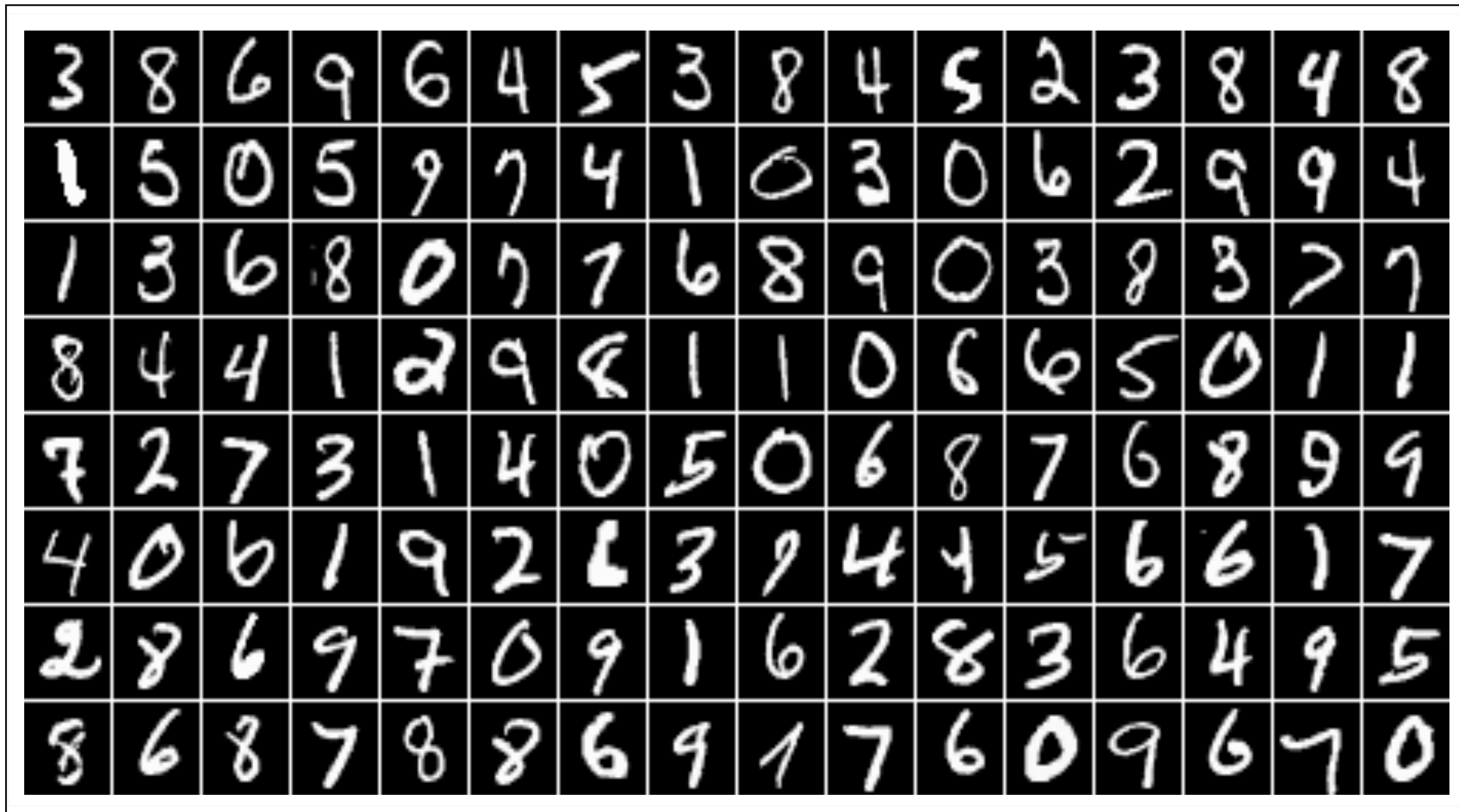
(Tieleman 2008)

Idea:

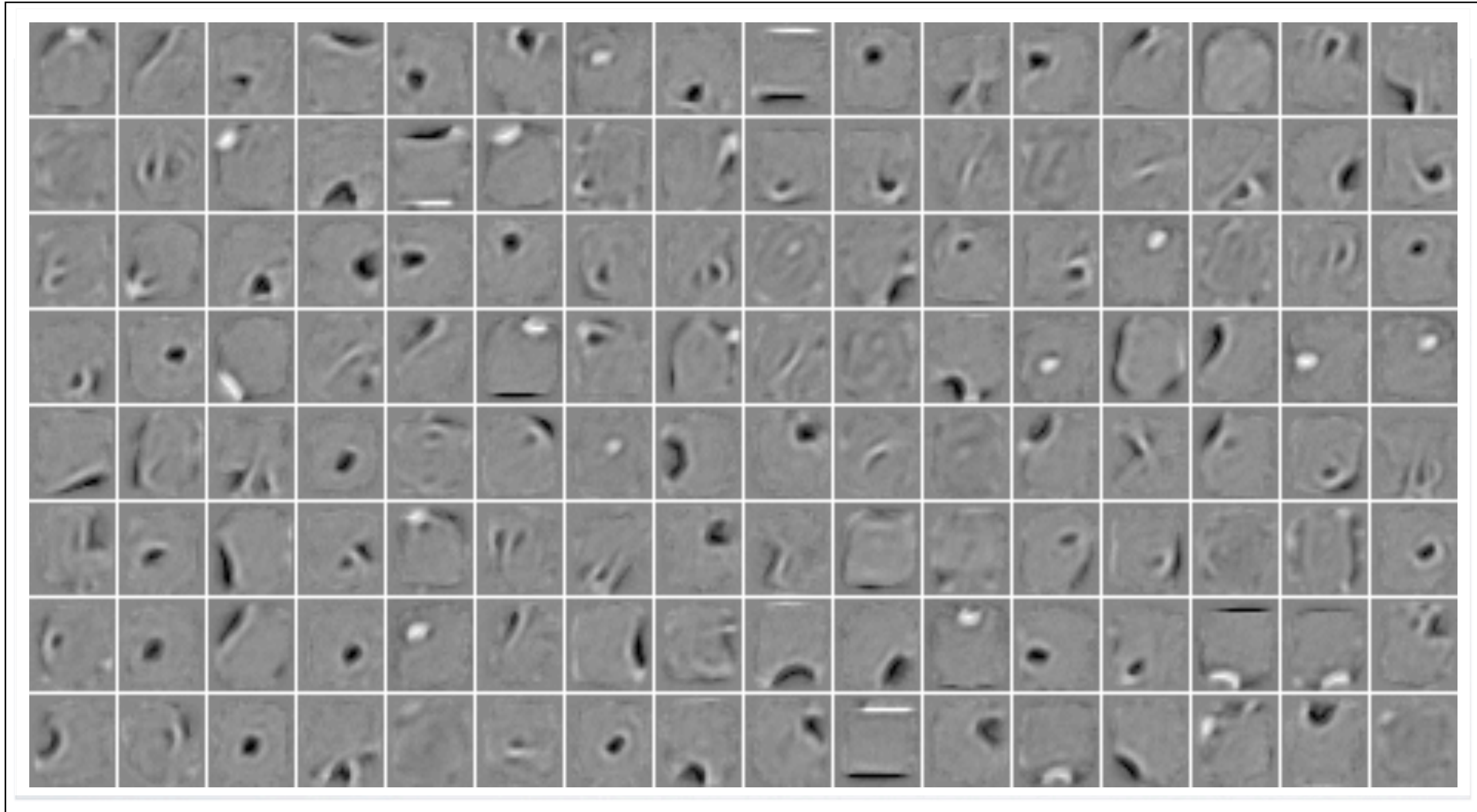
- instead of initializing the chain to $\tilde{v}^{(i)}$, initialize the chain to the negative sample of the last iteration



RBM: Example (MNIST)



RBM: Example (Filters)



RBM's: Debugging

- It's very difficult to **monitor** the training of RBMs because we do not have access to the log likelihood
- It's also not possible to check gradients with finite differences
- Can instead rely on some approximate “tricks”:
 - plot average stochastic reconstruction $\|v^{(i)} - \tilde{v}\|^2$ and see if it tends to decrease
 - for inputs that correspond to images, **visualize** the connection coming into each hidden unit (the filters)
 - can also try to approximate the partition function Z and see whether the (approximated) NLL decreases

Gaussian-Bernoulli RBM

What if inputs are **unbounded real values**?

- Add a quadratic term to the energy function

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h} - \frac{1}{2} \mathbf{v}^\top \mathbf{v}$$

- Only thing that changes is that $P(\mathbf{v}|\mathbf{h})$ is now a Gaussian distribution with mean $\mu = \mathbf{b} + \mathbf{W}^\top \mathbf{h}$ and identity covariance matrix
- Recommended to **normalize** the training set by
 - subtracting the mean of each input
 - dividing each input v_k by the training set standard deviation
- Should use a **smaller learning rate** than in a binary RBM