

# Discussion 7

EE599: Deep Learning

Olaoluwa Adigun

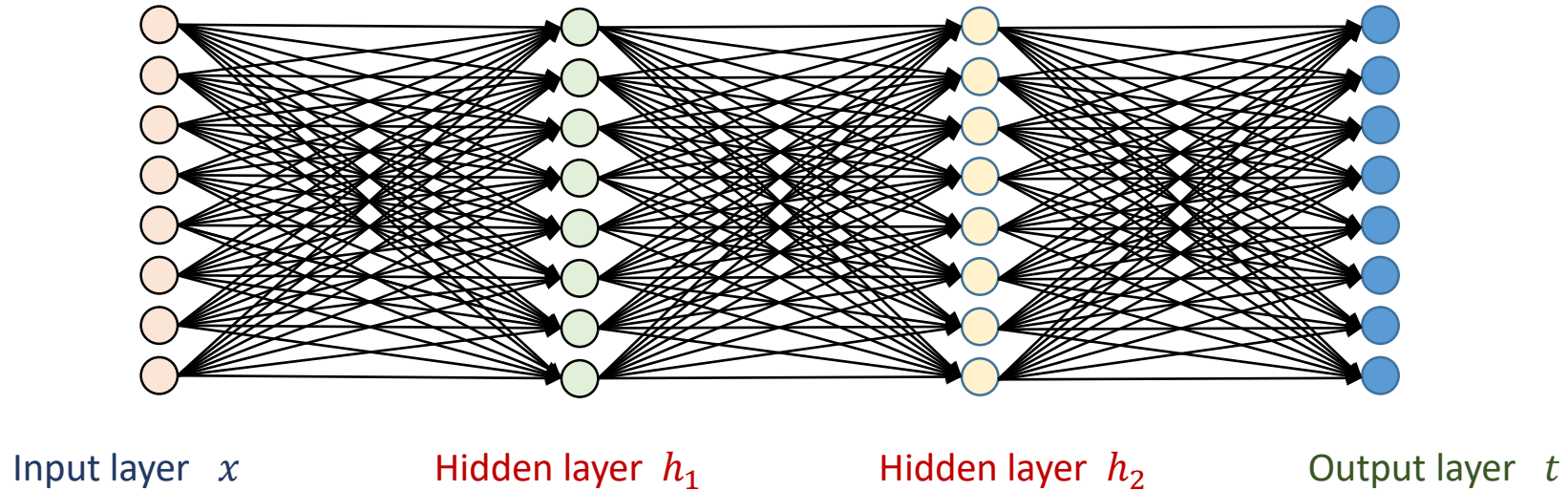
Spring 2020



**USC** University of  
Southern California

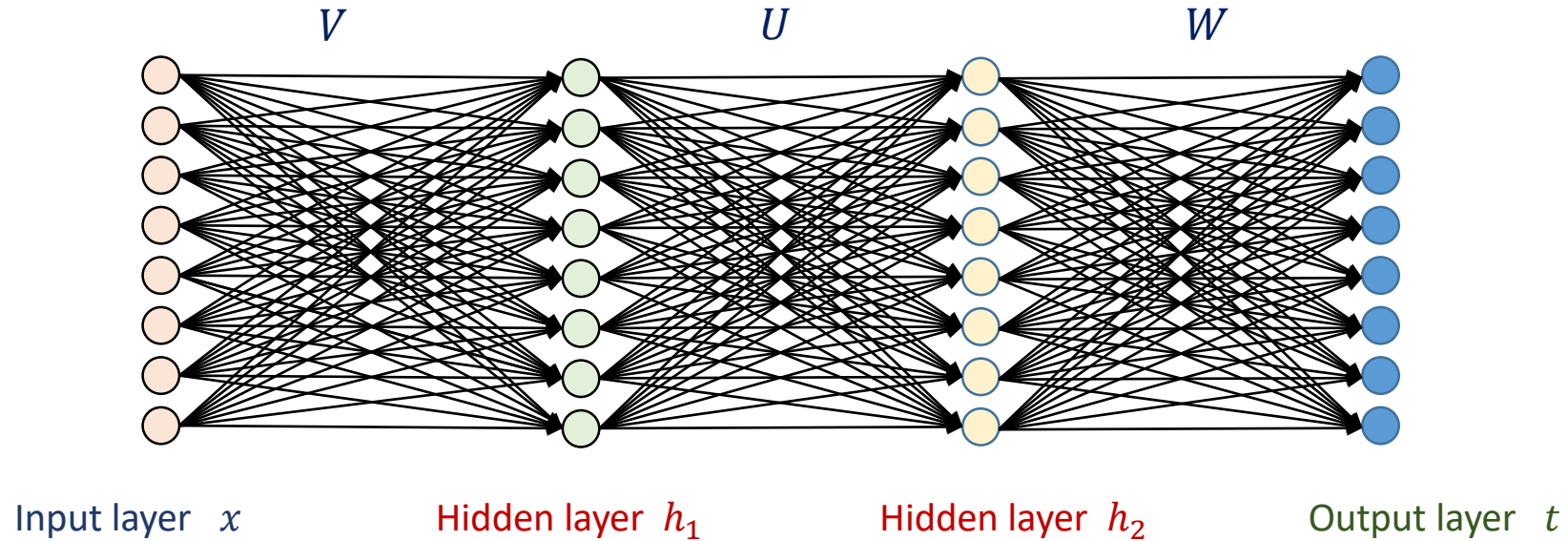
# Training MLP with Backpropagation

Find the best parameter  $\Theta^*$  that minimizes the cost function



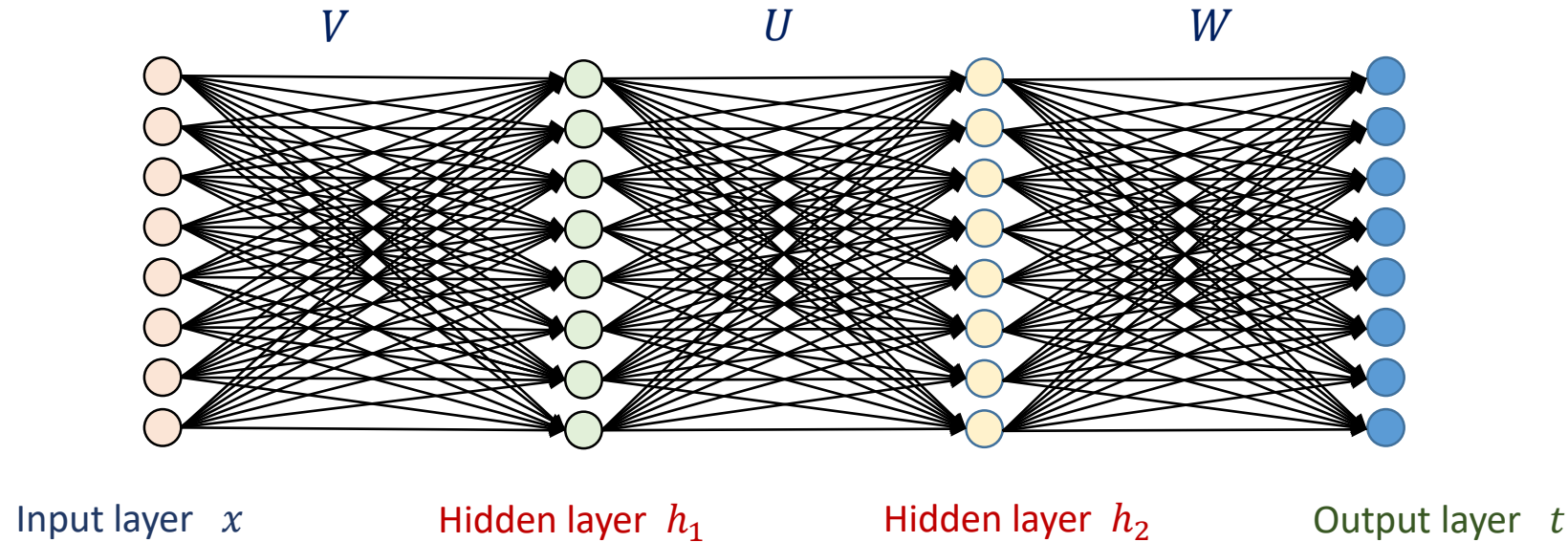
Backpropagation uses a variant of stochastic gradient descent to updates the weights

# Model Architecture



- Input Layer  $x$ :  $L$  neurons with *identity* activation
- Hidden Layer  $h_1$ :  $I$  neurons with *sigmoid* activation
- Hidden Layer  $h_2$ :  $J$  neurons with *sigmoid* activation
- Output Layer  $t$ :  $J$  neurons with *softmax* activation

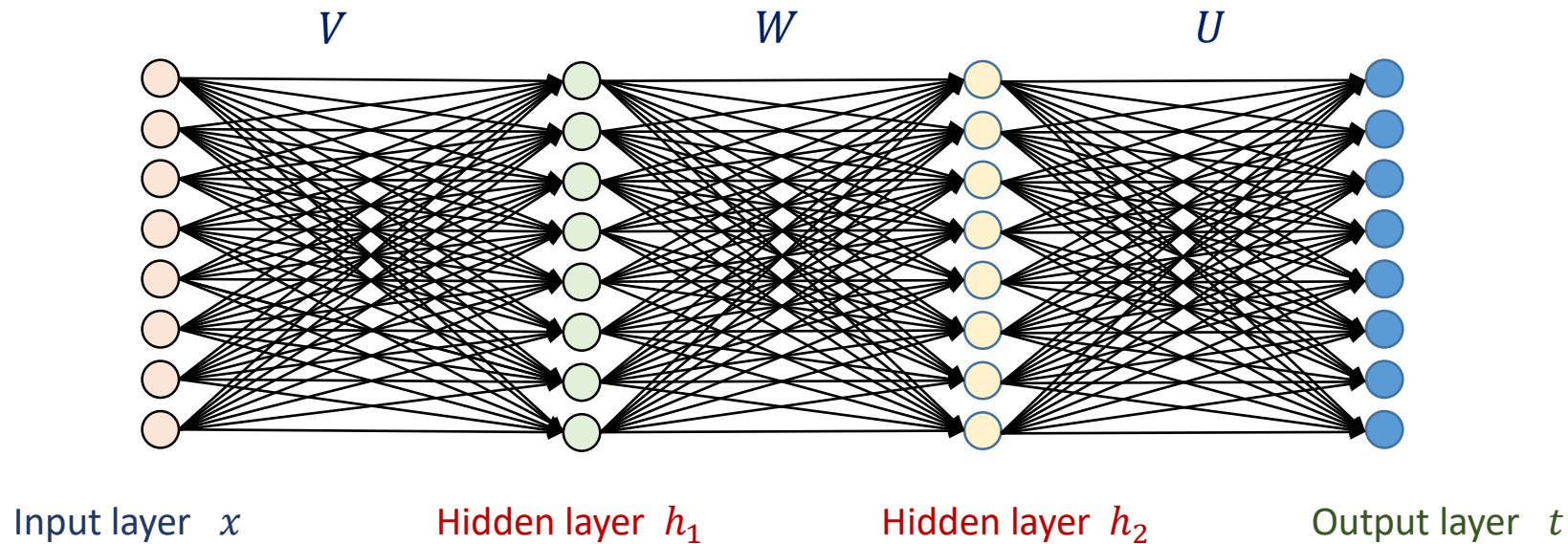
# Forward Pass over MLP Model (Inference)



- $b_l^{h_1}$  : Bias to the  $l^{th}$  neuron at the hidden layer  $h_1$
- $v_{li}$  : Weight connecting the  $l^{th}$  neuron of the hidden layer  $h_1$  to the  $i^{th}$  input neuron
- $b_j^{h_2}$  : Bias to the  $j^{th}$  neuron at the hidden layer  $h_2$
- $w_{jl}$  : Weight connecting the  $l^{th}$  neuron of layer  $h_1$  to the  $l^{th}$  neuron of layer  $h_2$
- $b_k^t$  : Bias to the  $k^{th}$  output neuron
- $u_{kj}$  : Weight connecting the  $j^{th}$  neuron of the hidden layer  $h_1$  to the  $k^{th}$  output neuron

$$\begin{aligned} o_l^{h_1} &= \sum_{i=1}^I v_{li} a_i^x + b_l^{h_1} \\ a_l^{h_1} &= \frac{1}{1 + \exp^{-o_l^{h_1}}} \\ o_j^{h_2} &= \sum_{l=1}^L w_{jl} a_l^{h_1} + b_j^{h_2} \\ a_j^{h_2} &= \frac{1}{1 + \exp^{-o_j^{h_2}}} \\ o_k^t &= \sum_{j=1}^J u_{kj} a_j^{h_2} + b_k^t \\ a_k^t &= \frac{\exp(o_k^t)}{\sum_{r=1}^K \exp(o_r^t)} \end{aligned}$$

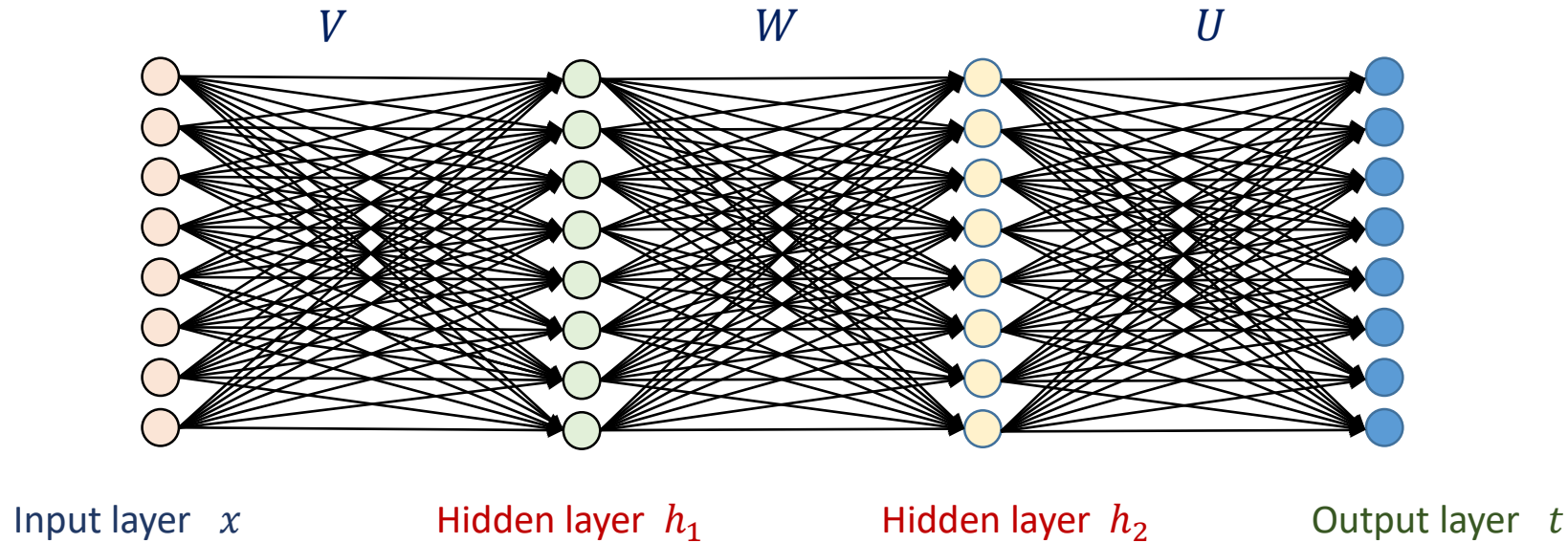
# Forward Pass over MLP Model (Inference)



- $\mathbf{b}^{h_1}$ : Bias vector to the hidden layer  $h_1$
- $\mathbf{V}$ : Weight connecting the input layer to hidden layer  $h_1$
- $\mathbf{b}^{h_2}$ : Bias vector to the hidden layer  $h_2$
- $\mathbf{W}$ : Weight connecting the hidden layer  $h_1$  to hidden layer  $h_2$
- $\mathbf{b}^t$ : Bias vector to the output layer
- $\mathbf{U}$ : Weight connecting the hidden layer  $h_2$  to the output layer

$$\begin{aligned}\mathbf{o}^{h_1} &= (\mathbf{V} \times \mathbf{a}^x) + \mathbf{b}^{h_1} \\ \mathbf{a}^{h_1} &= \sigma(\mathbf{o}^{h_1}) \\ \mathbf{o}^{h_2} &= (\mathbf{W} \times \mathbf{a}^{h_1}) + \mathbf{b}^{h_2} \\ \mathbf{a}^{h_2} &= \sigma(\mathbf{o}^{h_2}) \\ \mathbf{o}^t &= (\mathbf{U} \times \mathbf{a}^{h_2}) + \mathbf{b}^t \\ \mathbf{a}^t &= \text{Softmax}(\mathbf{o}^t)\end{aligned}$$

# Error Function and Backpropagation



$$\begin{aligned} \mathbf{o}^{h_1} &= (\mathbf{V} \times \mathbf{a}^x) + \mathbf{b}^{h_1} \\ \mathbf{a}^{h_1} &= \sigma(\mathbf{o}^{h_1}) \\ \mathbf{o}^{h_2} &= (\mathbf{W} \times \mathbf{a}^{h_1}) + \mathbf{b}^{h_2} \\ \mathbf{a}^{h_2} &= \sigma(\mathbf{o}^{h_2}) \\ \mathbf{o}^t &= (\mathbf{U} \times \mathbf{a}^{h_2}) + \mathbf{b}^t \\ \mathbf{a}^t &= \text{Softmax}(\mathbf{o}^t) \end{aligned}$$

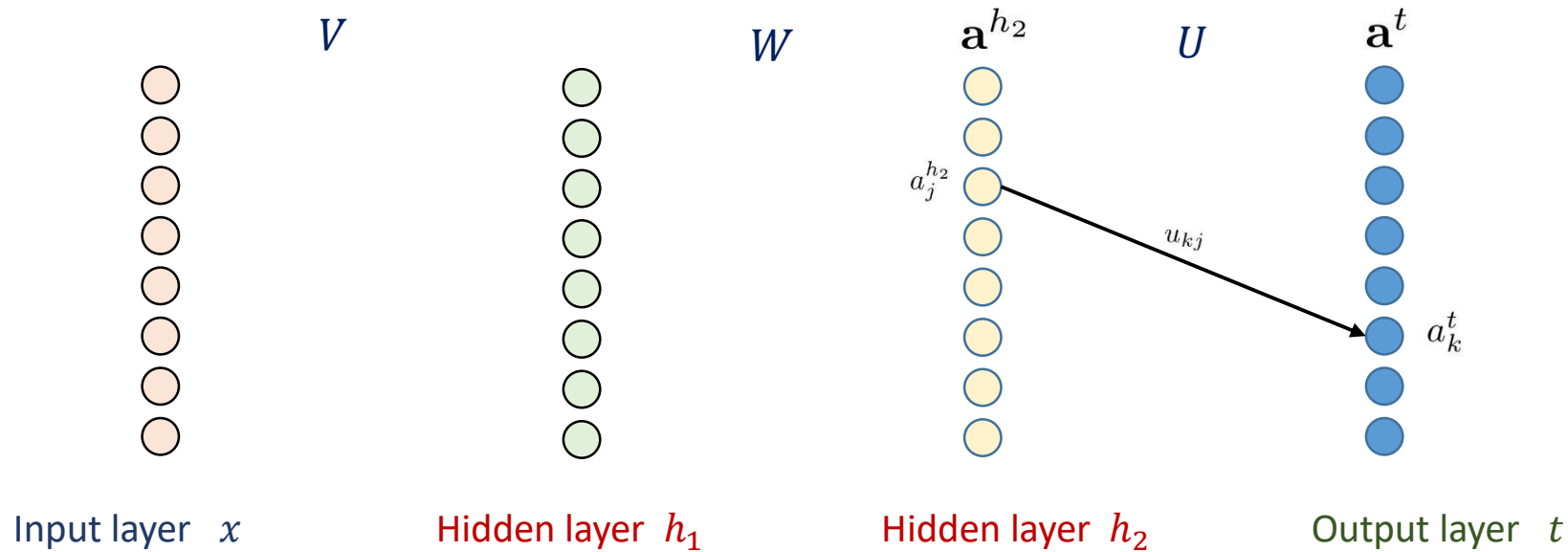
Cross entropy

Update rule

$$E(\Theta) = - \sum_{k=1}^K y_k \log a_k^t = -\mathbf{y}^T \log \mathbf{a}^t$$

$$\Theta^{(n+1)} = \Theta^{(n)} - \eta \nabla E(\Theta) \Big|_{\Theta=\Theta^{(n)}}$$

# Backpropagation Update Rule



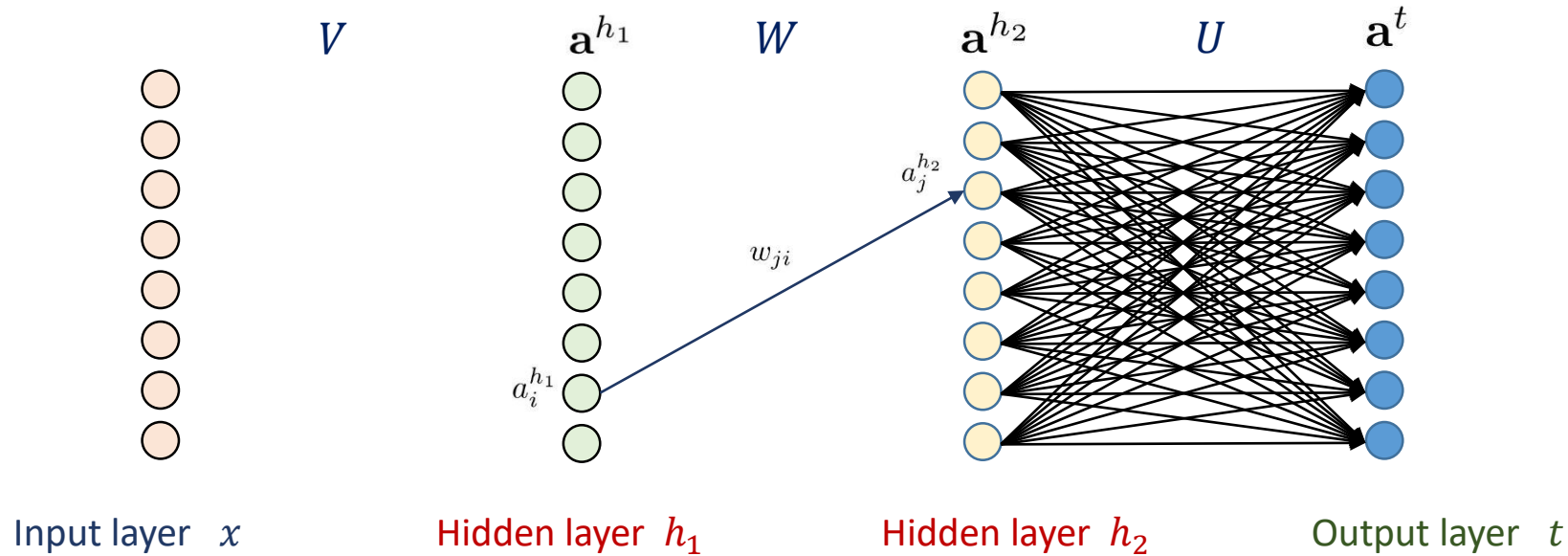
$$\begin{aligned}
 \mathbf{o}^{h_1} &= (\mathbf{V} \times \mathbf{a}^x) + \mathbf{b}^{h_1} \\
 \mathbf{a}^{h_1} &= \sigma(\mathbf{o}^{h_1}) \\
 \mathbf{o}^{h_2} &= (\mathbf{W} \times \mathbf{a}^{h_1}) + \mathbf{b}^{h_2} \\
 \mathbf{a}^{h_2} &= \sigma(\mathbf{o}^{h_2}) \\
 \mathbf{o}^t &= (\mathbf{U} \times \mathbf{a}^{h_2}) + \mathbf{b}^t \\
 \mathbf{a}^t &= \text{Softmax}(\mathbf{o}^t)
 \end{aligned}$$

$$\frac{\partial E}{\partial u_{kj}} = (y_k - a_k^t) a_j^{h_2}$$

$$\nabla_U E(\Theta) = (\mathbf{y} - \mathbf{a}^t)^T \mathbf{a}^{h_2}$$



# Backpropagation Update Rule



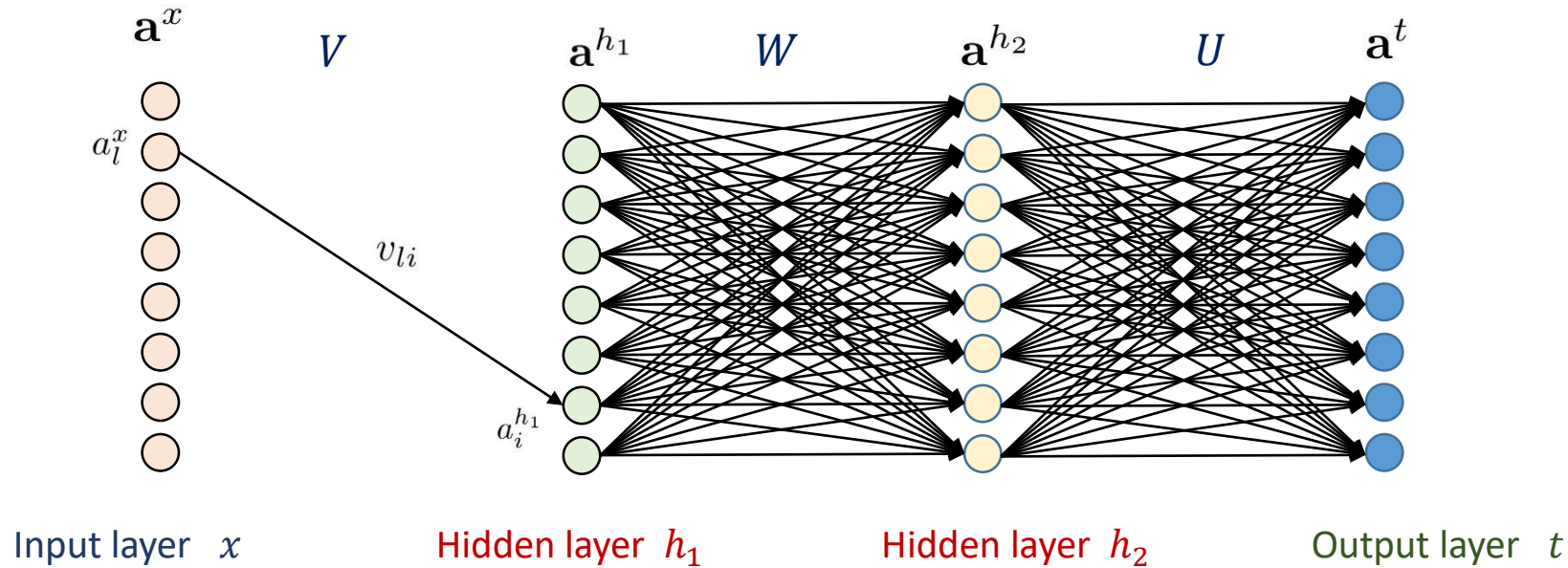
$$\begin{aligned} \mathbf{o}^{h_1} &= (\mathbf{V} \times \mathbf{a}^x) + \mathbf{b}^{h_1} \\ \mathbf{a}^{h_1} &= \sigma(\mathbf{o}^{h_1}) \\ \mathbf{o}^{h_2} &= (\mathbf{W} \times \mathbf{a}^{h_1}) + \mathbf{b}^{h_2} \\ \mathbf{a}^{h_2} &= \sigma(\mathbf{o}^{h_2}) \\ \mathbf{o}^t &= (\mathbf{U} \times \mathbf{a}^{h_2}) + \mathbf{b}^t \\ \mathbf{a}^t &= \text{Softmax}(\mathbf{o}^t) \end{aligned}$$

$$\frac{\partial E}{\partial w_{jl}} = \left( \sum_{k=1}^K (y_k - a_k^t) u_{kj} \right) a_j^{h_2} (1 - a_j^{h_2}) a_l^{h_1}$$

$$\nabla_W E(\Theta) = \left( \left( (\mathbf{y} - \mathbf{a}^t)^T \mathbf{U} \right) \odot \mathbf{a}^{h_2} \odot (1 - \mathbf{a}^{h_2}) \right)^T \mathbf{a}^{h_1}$$



# Backpropagation Update Rule



Update rule

$$\begin{aligned} \mathbf{o}^{h_1} &= (\mathbf{V} \times \mathbf{a}^x) + \mathbf{b}^{h_1} \\ \mathbf{a}^{h_1} &= \sigma(\mathbf{o}^{h_1}) \\ \mathbf{o}^{h_2} &= (\mathbf{W} \times \mathbf{a}^{h_1}) + \mathbf{b}^{h_2} \\ \mathbf{a}^{h_2} &= \sigma(\mathbf{o}^{h_2}) \\ \mathbf{o}^t &= (\mathbf{U} \times \mathbf{a}^{h_2}) + \mathbf{b}^t \\ \mathbf{a}^t &= \text{Softmax}(\mathbf{o}^t) \end{aligned}$$

$$\frac{\partial E}{\partial v_{li}} = \left( \sum_{j=1}^J \left( \sum_{k=1}^K (y_k - a_k^t) u_{kj} \right) a_j^{h_2} (1 - a_j^{h_2}) w_{jl} \right) a_l^{h_1} (1 - a_l^{h_1}) a_i^x$$

$$\nabla_V E(\Theta) = \left( \left( ((\mathbf{y} - \mathbf{a}^t)^T \mathbf{U}) \odot \mathbf{a}^{h_2} \odot (1 - \mathbf{a}^{h_2}) \right) \times \mathbf{W}^T \right) \odot \mathbf{a}^{h_1} \odot (1 - \mathbf{a}^{h_1}) \right)^T \mathbf{a}^x$$

