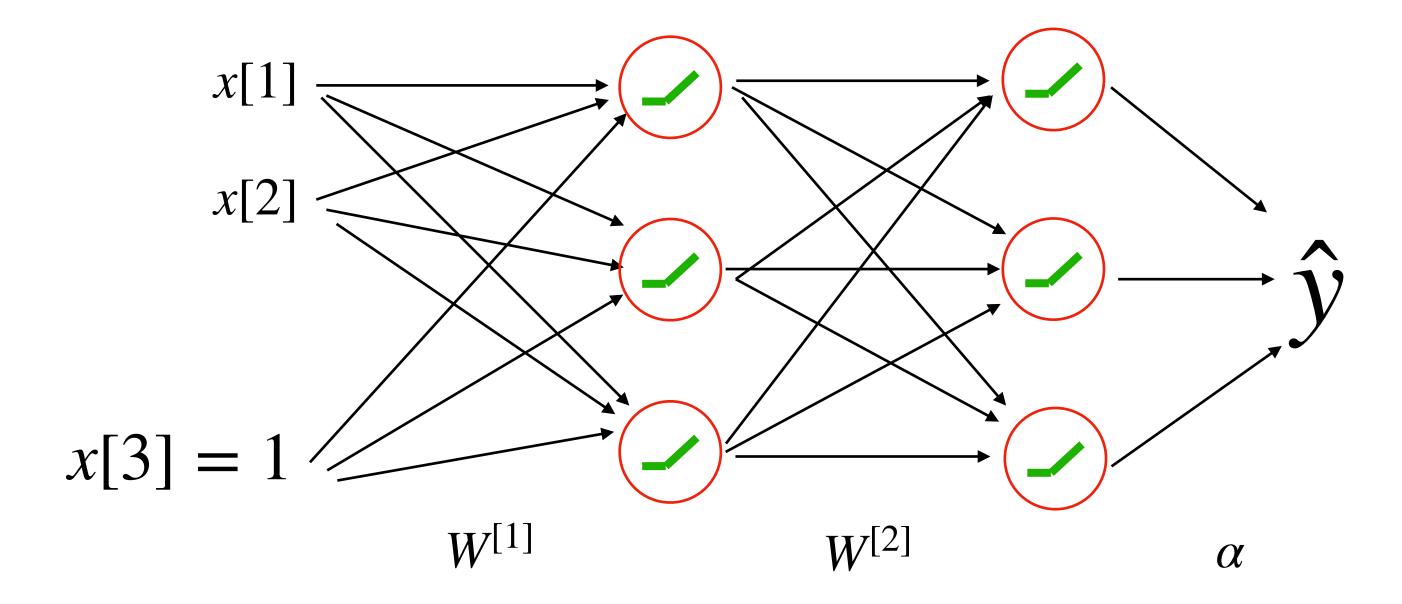
Neural Network: Training & Backpropagation

Announcements

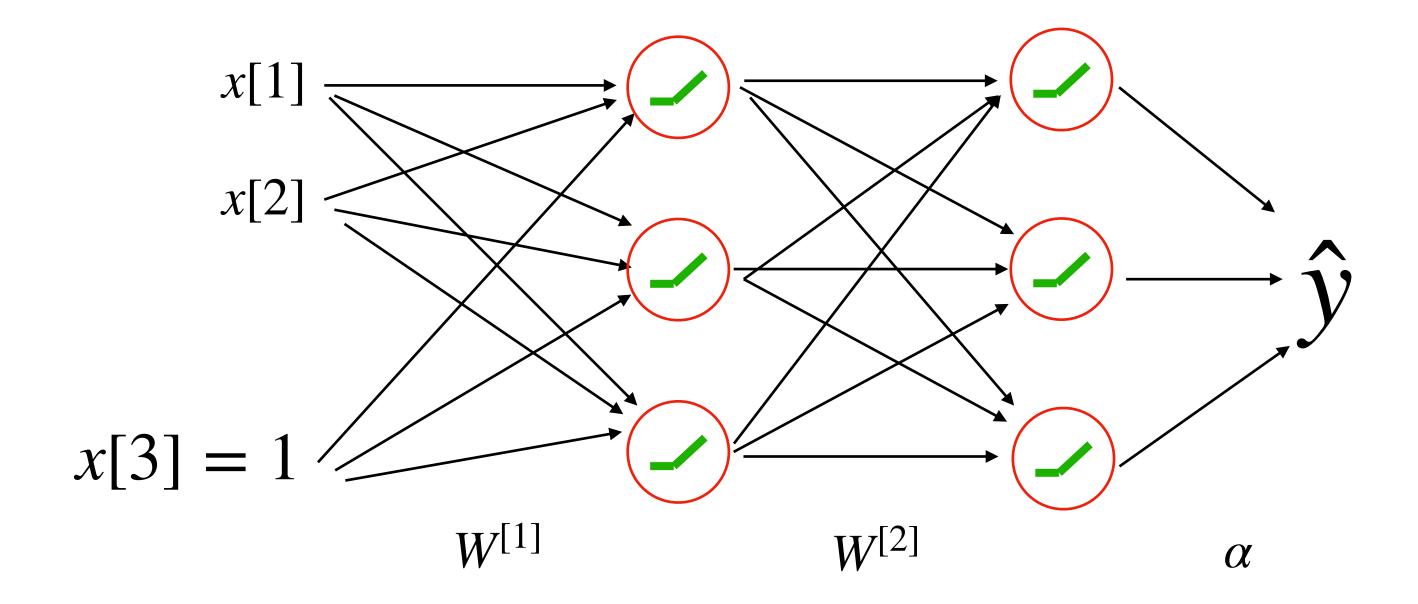
Recap

A two layer fully connected feedforward NN:



Recap

A two layer fully connected feedforward NN:



$$h(x) := \alpha^{\mathsf{T}} \mathsf{ReLU} \left(W^{[2]} \mathsf{ReLU} \left(W^{[1]} x \right) \right) + b$$

Outline of Today

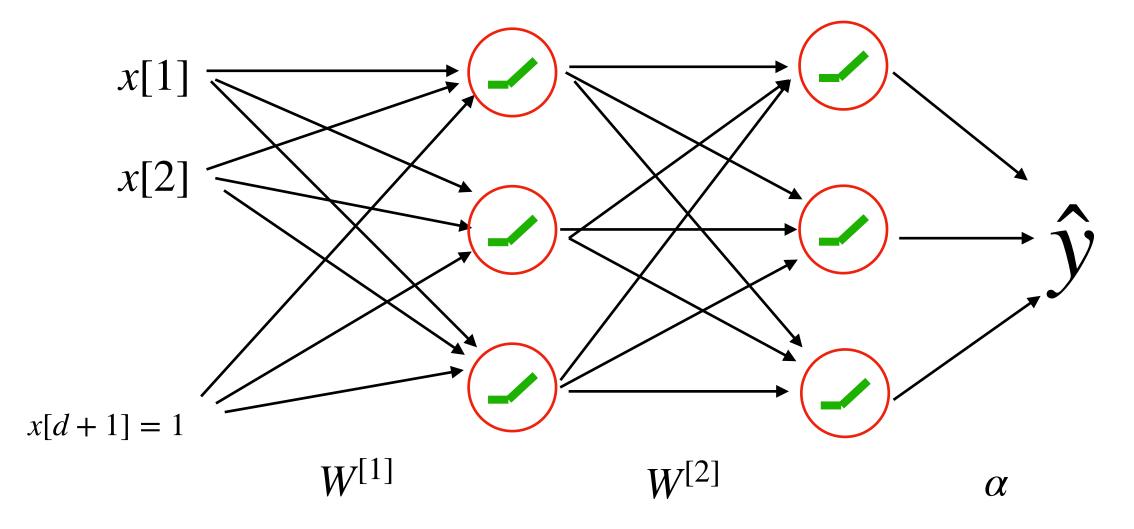
1. Training NNs via SGD

2. A Naive approach of computing gradients

3. Backpropagation: efficient way of computing gradients

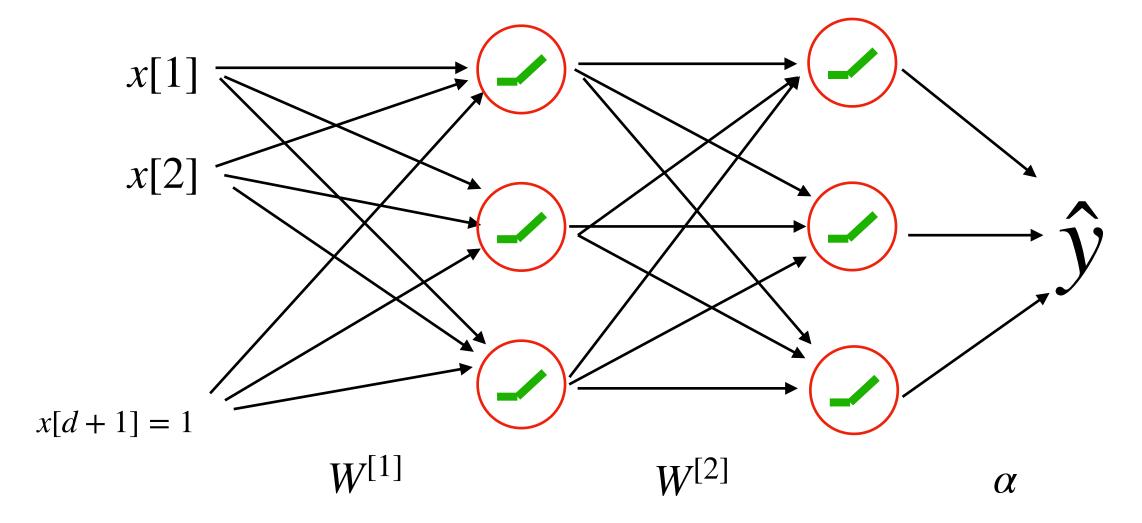
Square loss on training example (x, y)

$$h(x) := \alpha^{\mathsf{T}} \mathsf{ReLU}\left(W^{[2]} \mathsf{ReLU}\left(W^{[1]}x\right)\right) + b$$



Square loss on training example (x, y)

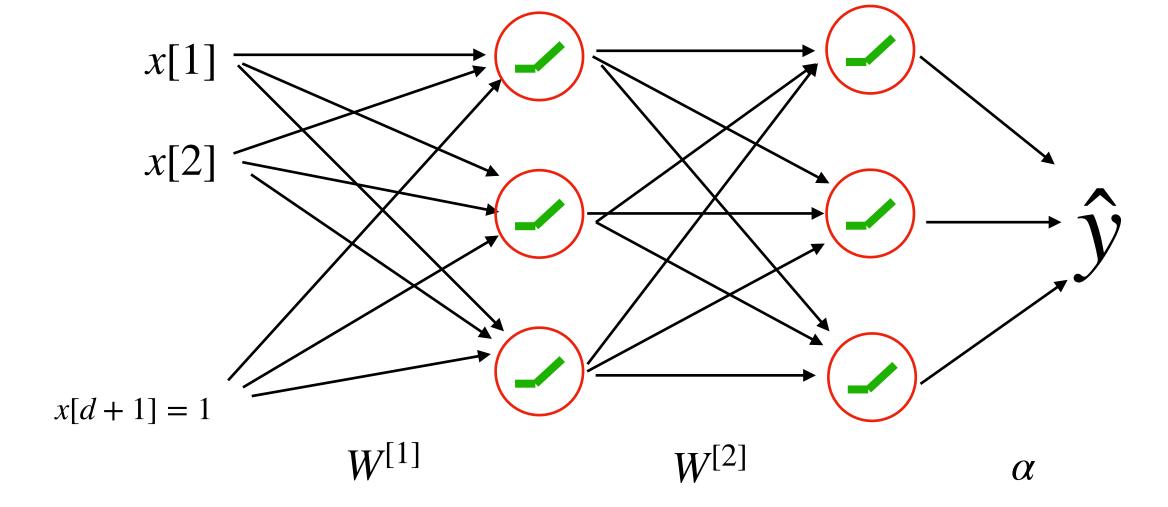
$$h(x) := \alpha^{\mathsf{T}} \mathsf{ReLU}\left(W^{[2]} \mathsf{ReLU}\left(W^{[1]}x\right)\right) + b$$



$$\mathcal{E}(h(x), y) = (\hat{y} - y)^2$$
, where $\hat{y} = h(x)$

Square loss on training example (x, y)

$$h(x) := \alpha^{\mathsf{T}} \mathsf{ReLU} \left(W^{[2]} \mathsf{ReLU} \left(W^{[1]} x \right) \right) + b$$

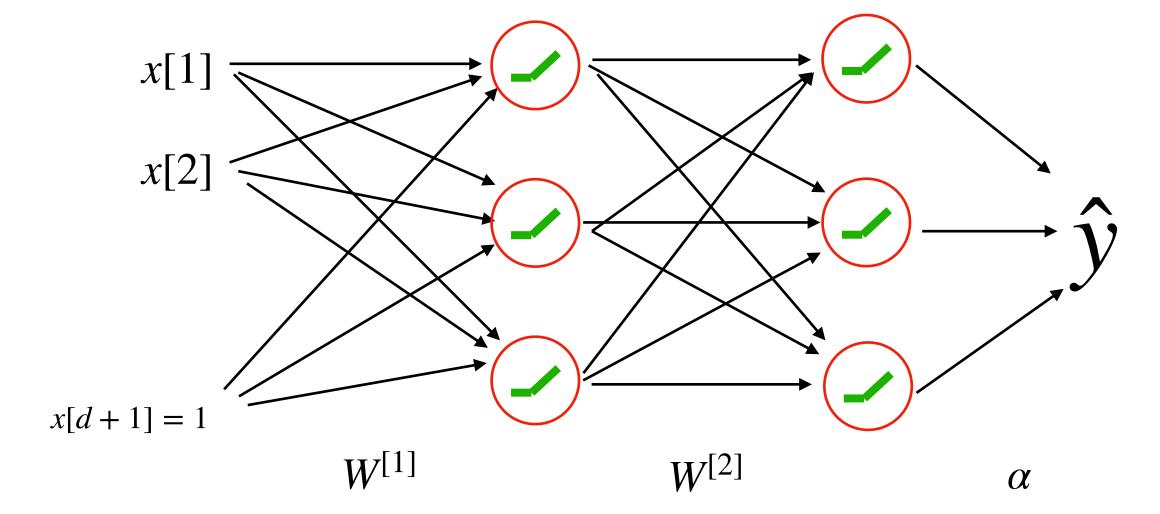


$$\mathcal{E}(h(x), y) = (\hat{y} - y)^2$$
, where $\hat{y} = h(x)$

Trainable parameters $W^{[1]}, W^{[2]}, \alpha, b$

Square loss on training example (x, y)

$$h(x) := \alpha^{\mathsf{T}} \mathsf{ReLU} \left(W^{[2]} \mathsf{ReLU} \left(W^{[1]} x \right) \right) + b$$



$$\mathcal{E}(h(x), y) = (\hat{y} - y)^2$$
, where $\hat{y} = h(x)$

Trainable parameters $W^{[1]}, W^{[2]}, \alpha, b$

Compute gradients:

$$\frac{\partial \mathcal{E}(h(x), y)}{\partial W^{[1]}} \qquad \frac{\partial \mathcal{E}(h(x), y)}{\partial W^{[2]}}$$

$$\frac{\partial \mathcal{E}(h(x), y)}{\partial \alpha} \qquad \frac{\partial \mathcal{E}(h(x), y)}{\partial b}$$

Mini-batch Stochastic gradient descent

$$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$$

For epoc t = 1 to T:

$$heta = [W^{[1]}, W^{[2]}, \alpha, b]$$
 // go through dataset multiple times For epoc $t=1$ to T :

Mini-batch Stochastic gradient descent

$$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$$
 // go through dataset multiple times For epoc $t=1$ to T :

Randomly shuffle the data

$$\theta = [W^{[1]}, W^{[2]}, \alpha, b] \text{// go through dataset multiple times}$$
 For epoc $t=1$ to T : // important (unbiased estimate of the true gradient)

$$heta = [W^{[1]}, W^{[2]}, \alpha, b]$$
 // go through dataset multiple times

For epoc $t=1$ to T : // important (unbiased estimate of the true gradient)

Split the data into n/B many batches \mathcal{D}_i , each w/ size B

$$\theta = [W^{[1]}, W^{[2]}, \alpha, b] \text{ // go through dataset multiple times}$$
 For epoc $t=1$ to T : // important (unbiased estimate of the true gradient) Split the data into n/B many batches \mathcal{D}_i , each w/ size B For i = 1 to n/B

$$\theta = [W^{[1]}, W^{[2]}, \alpha, b] \hspace{1cm} /\!/ \text{ go through dataset multiple times}$$
 For epoc $t=1$ to T : $\hspace{1cm} /\!/ \text{ important (unbiased estimate of the true gradient)}$ Split the data into n/B many batches \mathcal{D}_i , each w/ size B For $i=1$ to n/B

Mini-batch gradient
$$g = \sum_{x,y \in \mathcal{D}_i} \nabla_{\theta} \mathcal{E}(h_{\theta}(x),y)/B$$

$$heta = [W^{[1]}, W^{[2]}, \alpha, b]$$
 // go through dataset multiple times

For epoc $t=1$ to T : // important (unbiased estimate of the true gradient)

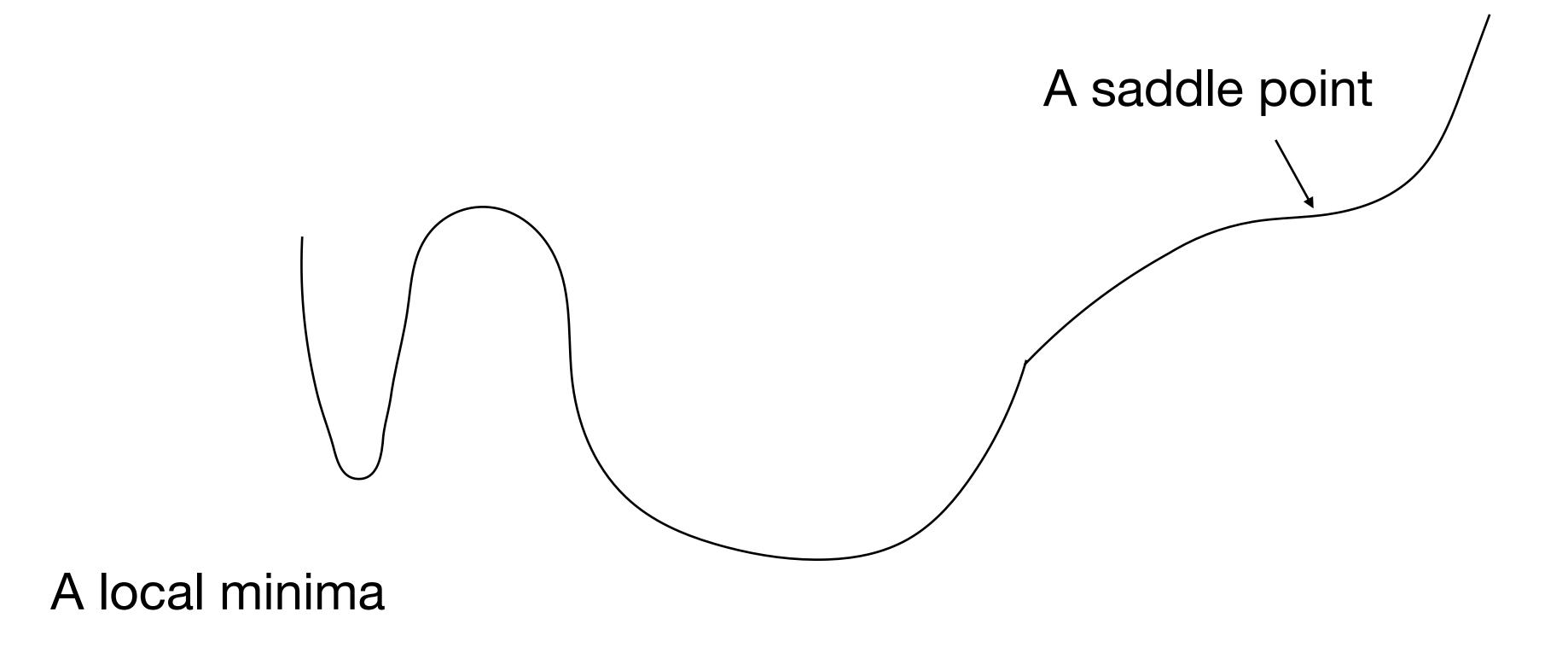
Split the data into n/B many batches \mathcal{D}_i , each w/ size B

For $i=1$ to n/B

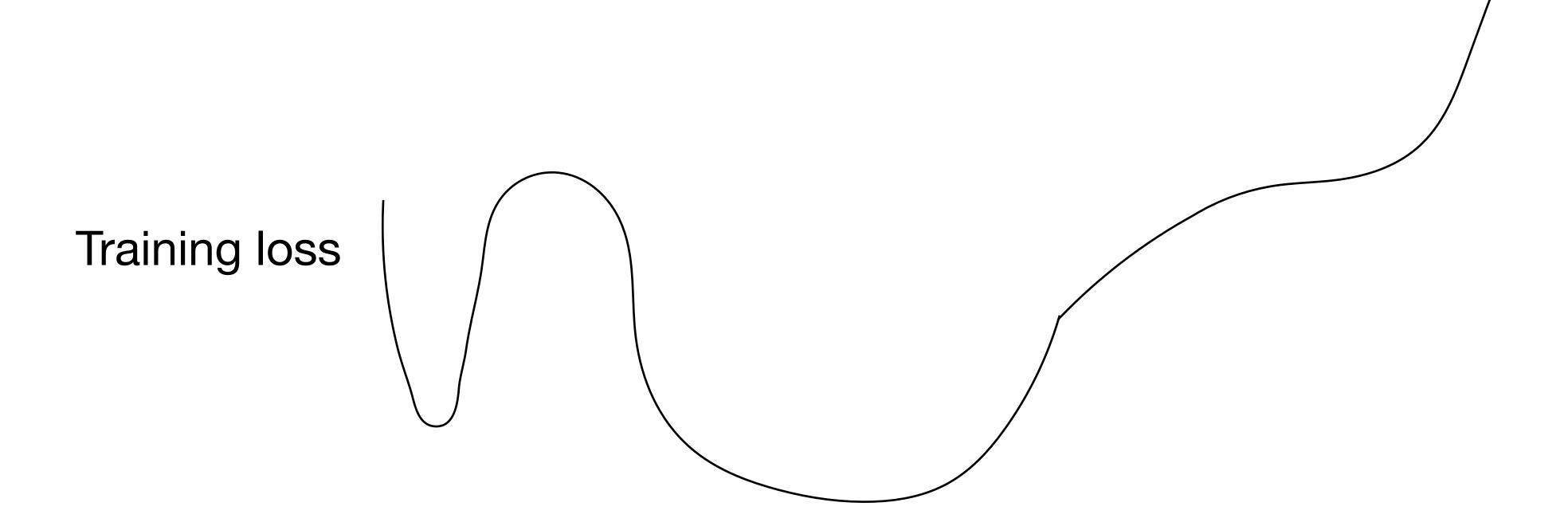
Mini-batch gradient
$$g=\sum_{x,y\in \mathcal{D}_i}\nabla_{\theta}\mathcal{E}(h_{\theta}(x),y)/B$$

$$\theta=\theta-\eta g$$

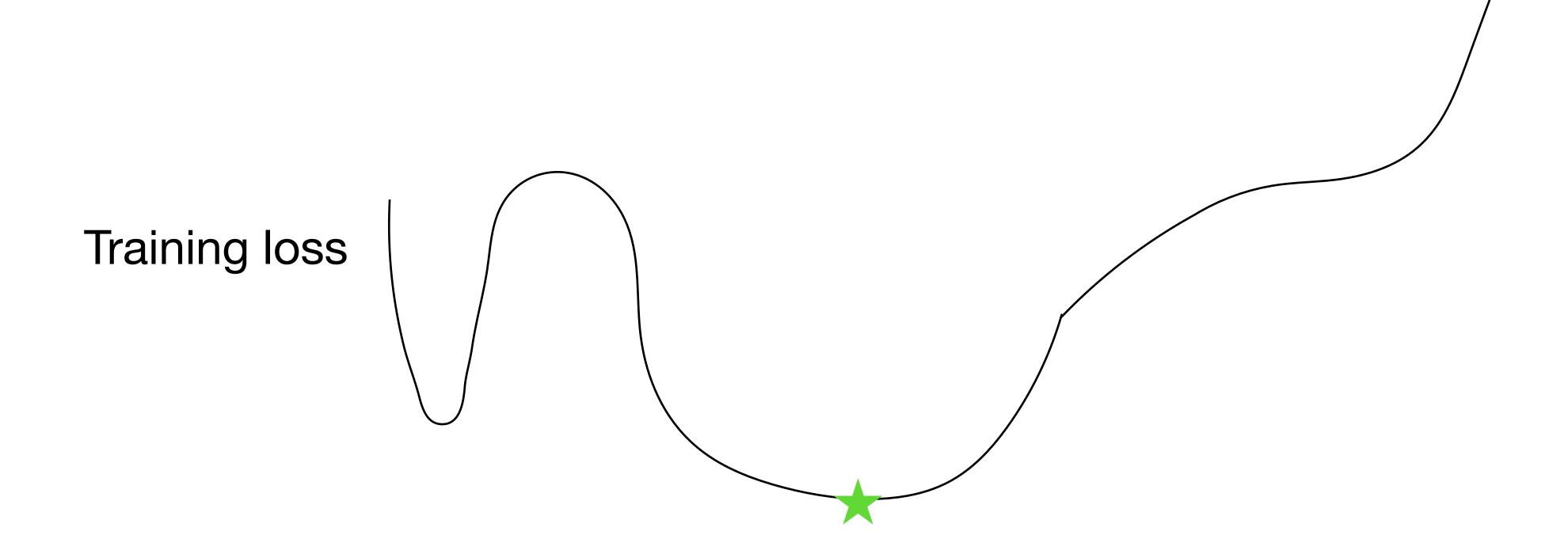
SGD helps avoiding local minima and saddle point



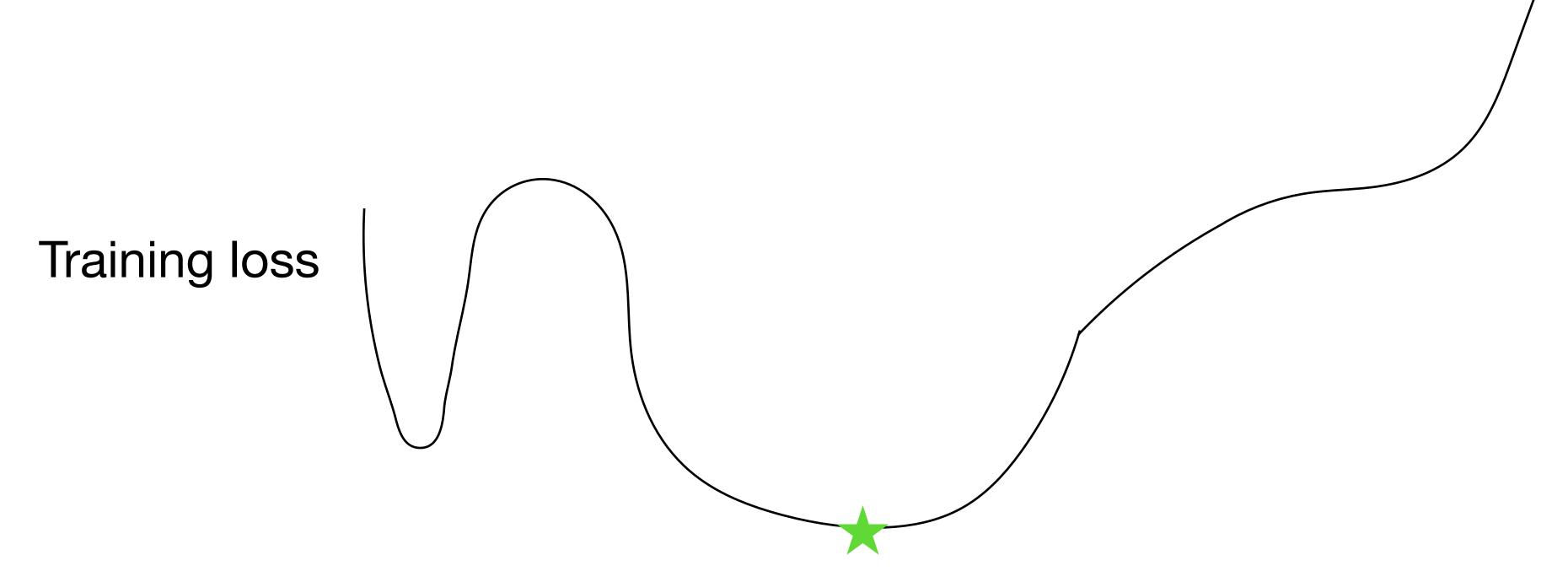
SGD tends to converge to a flat region



SGD tends to converge to a flat region

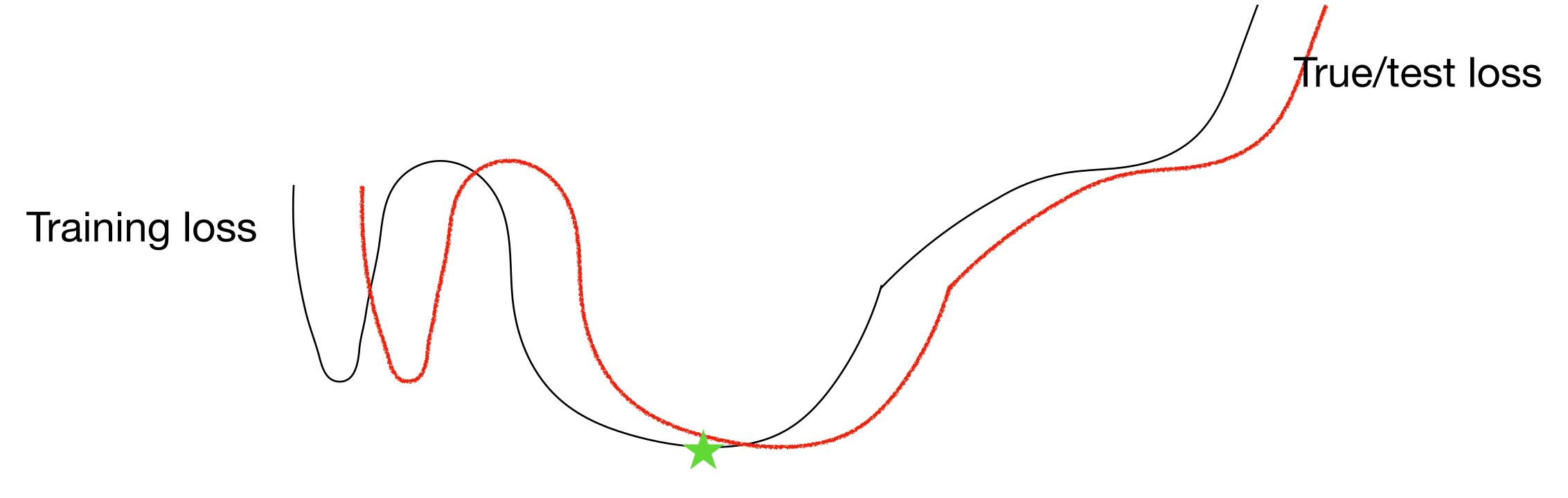


SGD tends to converge to a flat region



A flat local minima solution can help generalizes better to test data

SGD tends to converge to a flat region



A flat local minima solution can help generalizes better to test data

Outline of Today

1. Training NNs via SGD

2. A Naive approach of computing gradients

3. Backpropagation: efficient way of computing gradients

Consider the following one-dim case with identity transformation

$$x \xrightarrow{w_1} \bigcirc \xrightarrow{w_2} \bigcirc \xrightarrow{w_3} \dots \xrightarrow{w_T} \bigcirc \xrightarrow{a} \hat{y}$$

Consider the following one-dim case with identity transformation

$$x \xrightarrow{w_1} \bigcirc \xrightarrow{w_2} \bigcirc \xrightarrow{w_3} \dots \xrightarrow{w_T} \bigcirc \xrightarrow{a} \hat{y}$$

$$\hat{y} = aw_T \dots w_2 w_1 x$$

Consider the following one-dim case with identity transformation

$$x \xrightarrow{w_1} \bigcirc \xrightarrow{w_2} \bigcirc \xrightarrow{w_3} \dots \xrightarrow{w_T} \bigcirc \xrightarrow{a} \hat{y}$$

$$\hat{y} = aw_T \dots w_2 w_1 x$$

Let's compute derivatives $\partial \hat{y}/\partial w_i$, $\forall i = 1,...T$

Consider the following one-dim case with identity transformation

$$x \xrightarrow{w_1} \bigcirc \xrightarrow{w_2} \bigcirc \xrightarrow{w_3} \dots \xrightarrow{w_T} \bigcirc \xrightarrow{a} \hat{y}$$

$$\hat{y} = aw_T \dots w_2 w_1 x$$

Let's compute derivatives $\partial \hat{y}/\partial w_i$, $\forall i = 1,...T$

Via chain rule:
$$\frac{\partial \mathcal{E}}{\partial w_i} = \frac{\partial \mathcal{E}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_i}$$

$$x \xrightarrow{w_1} (z_1) \xrightarrow{w_2} (z_2) \xrightarrow{w_3} \dots \xrightarrow{w_T} (z_T) \xrightarrow{a} \hat{y}$$

$$z_0 = x$$

$$z_2 = w_2 z_1$$

$$x \xrightarrow{w_1} (z_1) \xrightarrow{w_2} (z_2) \xrightarrow{w_3} \dots \xrightarrow{w_T} (z_T) \xrightarrow{a} \hat{y}$$

$$z_0 = x$$

$$z_2 = w_2 z_1$$

$$x \xrightarrow{w_1} (z_1) \xrightarrow{w_2} (z_2) \xrightarrow{w_3} \dots \xrightarrow{w_T} (z_T) \xrightarrow{a} \hat{y}$$

$$z_0 = x$$

$$z_2 = w_2 z_1$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$x \xrightarrow{w_1} (z_1) \xrightarrow{w_2} (z_2) \xrightarrow{w_3} \dots \xrightarrow{w_T} (z_T) \xrightarrow{a} \hat{y}$$

$$z_0 = x$$

$$z_2 = w_2 z_1$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$
 // computation: T

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$
 // computation: T

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$
 // computation: T
$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$
 // computation: T-1

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$
 // computation: T
$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$
 // computation: T-1
$$\frac{\partial \hat{y}}{\partial w_T} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial w_T} \frac{\partial z_T}{\partial w_T}$$

$$x \xrightarrow{w_1} (z_1) \xrightarrow{w_2} (z_2) \xrightarrow{w_3} \dots \xrightarrow{w_T} (z_T) \xrightarrow{a} \hat{y}$$

$$z_0 = x$$

$$z_2 = w_2 z_1$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$
 // computation: T
$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$
 // computation: T-1
$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_2} \frac{\partial z_T}{\partial w_2}$$
 // computation: 1

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad \text{// computation: T}$$

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2} \quad \text{// computation: T-1}$$

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_2} \frac{\partial z_T}{\partial w_2} \quad \text{// computation: T-1}$$

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_2} \frac{\partial z_T}{\partial w_2} \quad \text{// computation: 1}$$

A naive algorithm

Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$
 // computation: T

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$
 // computation: T-1

$$\frac{\partial \hat{y}}{\partial w_T} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial w_T}$$
 // computation: 1

Total complexity:

$$1 + 2 + \dots + T = O(T^2)$$

Quadratic in size of the graph!

Summary so far

What we did:

for each edge weight w_i , apply chain rule to calculate $\partial \hat{y}/\partial w_i$

Summary so far

What we did:

for each edge weight w_i , apply chain rule to calculate $\partial \hat{y}/\partial w_i$

What we got:

Able to compute gradient in running time $O((size of graph)^2)$

Summary so far

What we did:

for each edge weight w_i , apply chain rule to calculate $\partial \hat{y}/\partial w_i$

What we got:

Able to compute gradient in running time $O\left((\text{size of graph})^2\right)$

Can we do better in running time?

Outline of Today

1. Training NNs via SGD

2. A Naive approach of computing gradients

3. Backpropagation: efficient way of computing gradients

...Hinton popularized what they termed a "backpropagation" algorithm ... in 1986.

...Hinton popularized what they termed a "backpropagation" algorithm ... in 1986.

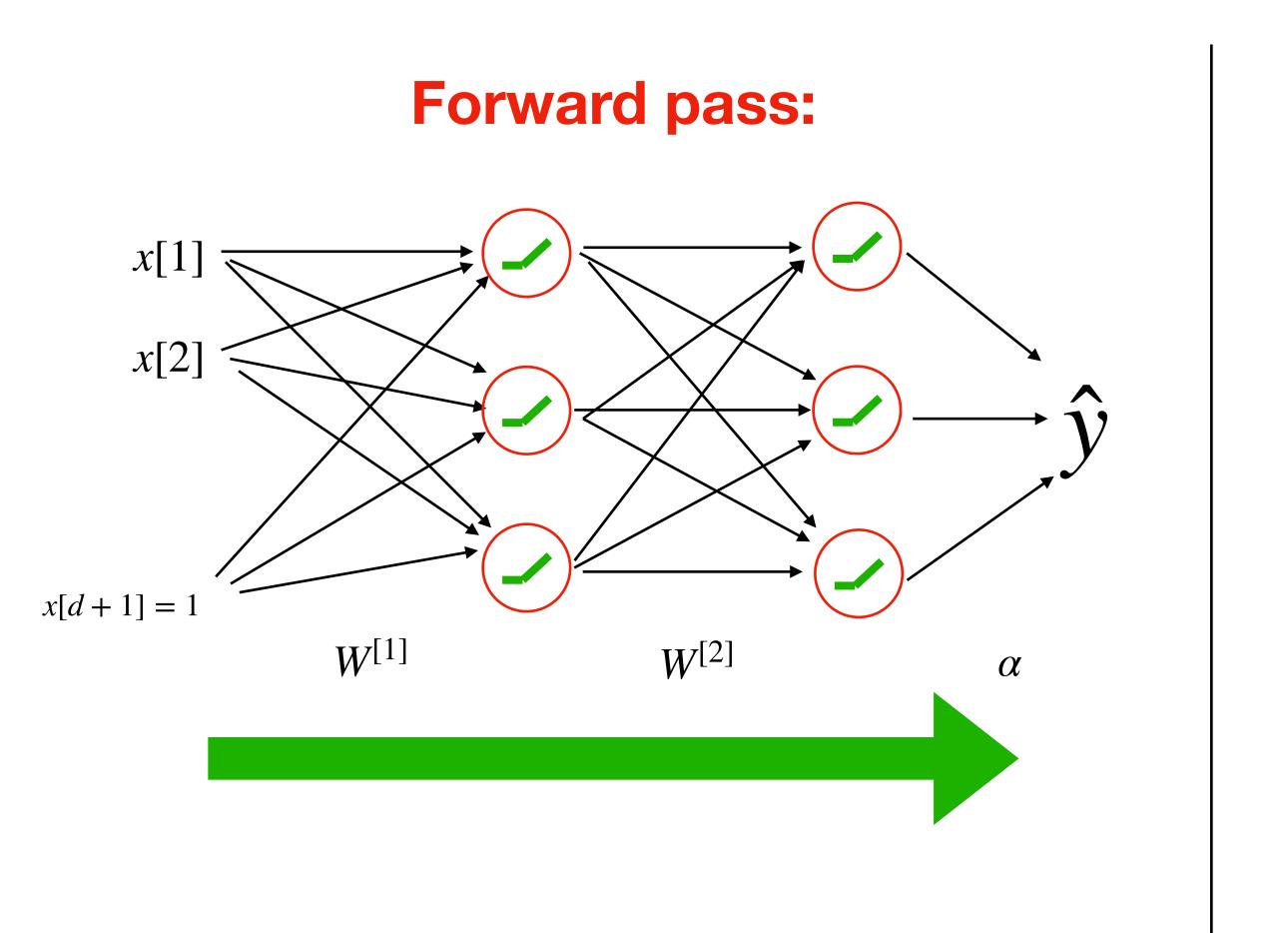
...the algorithm propagated measures of the errors produced by the network's guesses backwards through its neurons, starting with those directly connected to the outputs.

...Hinton popularized what they termed a "backpropagation" algorithm ... in 1986.

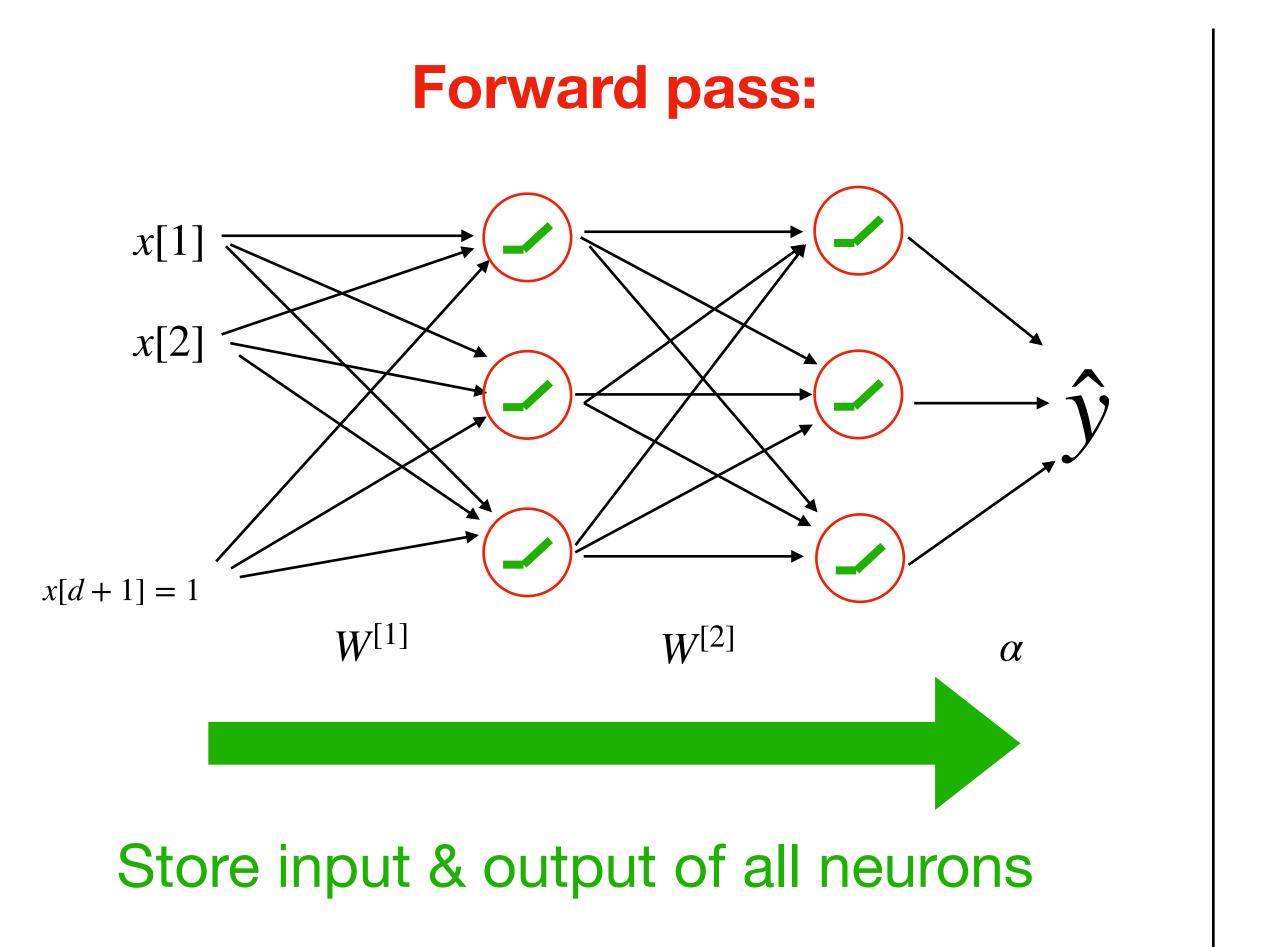
...the algorithm propagated measures of the errors produced by the network's guesses backwards through its neurons, starting with those directly connected to the outputs.

This allowed networks with intermediate "hidden" neurons between input and output layers to **learn efficiently**, overcoming the limitations noted by Minsky and Papert.

Forward pass followed by a backward pass

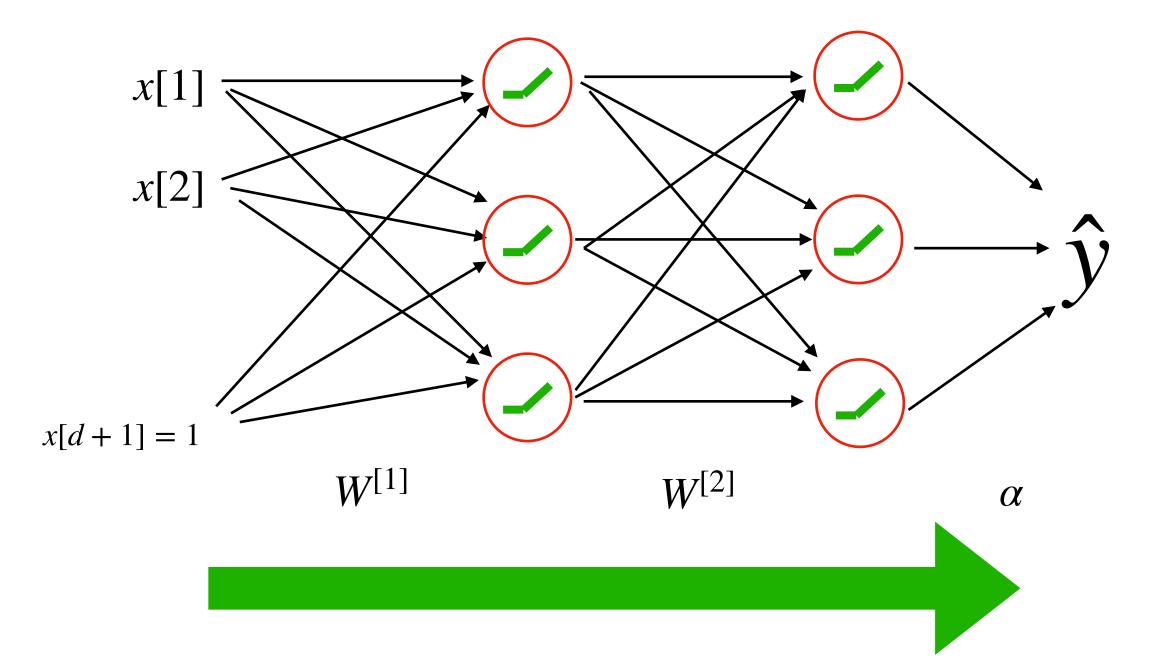


Forward pass followed by a backward pass



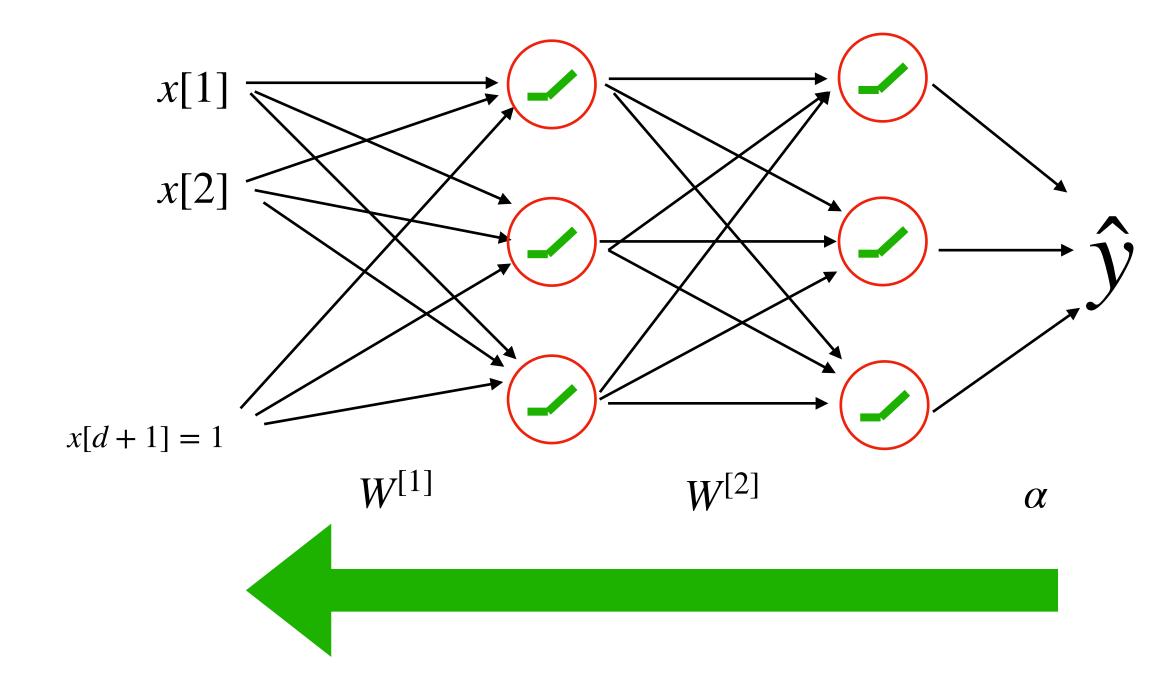
Forward pass followed by a backward pass

Forward pass:



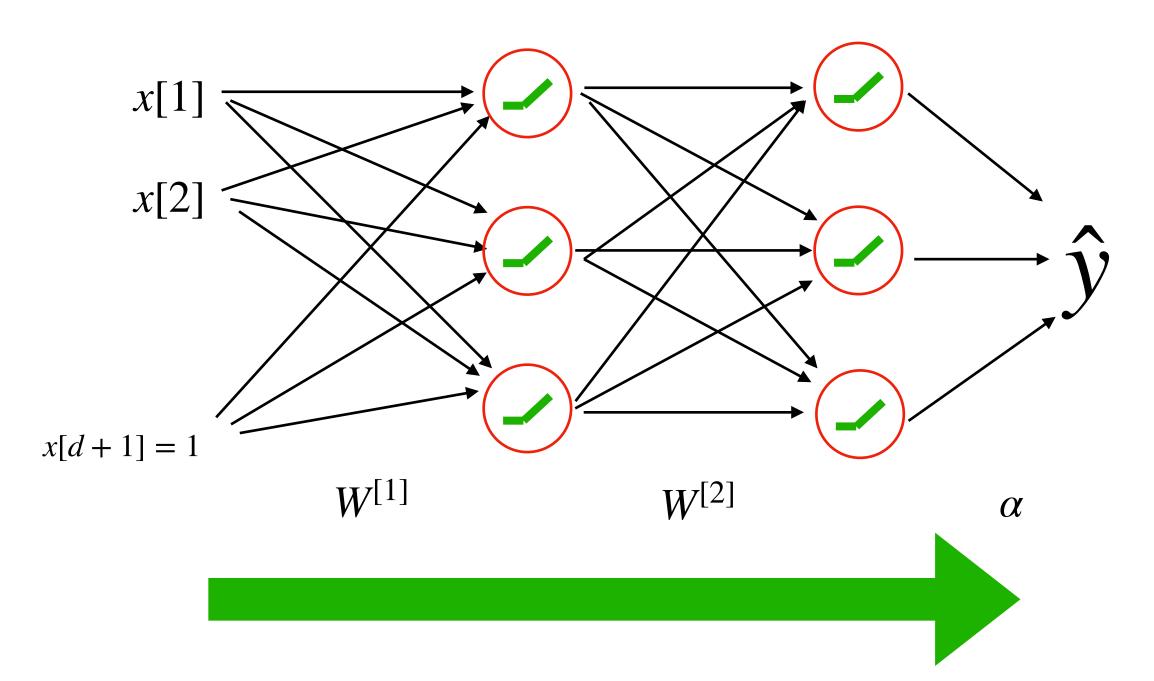
Store input & output of all neurons

backward pass:



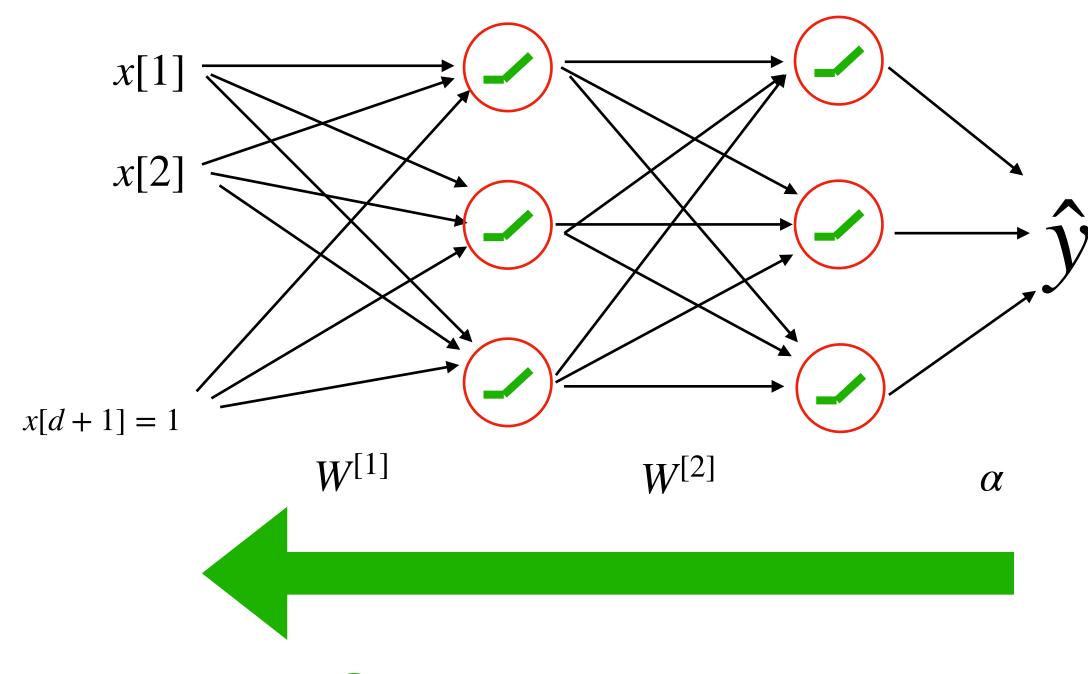
Forward pass followed by a backward pass

Forward pass:

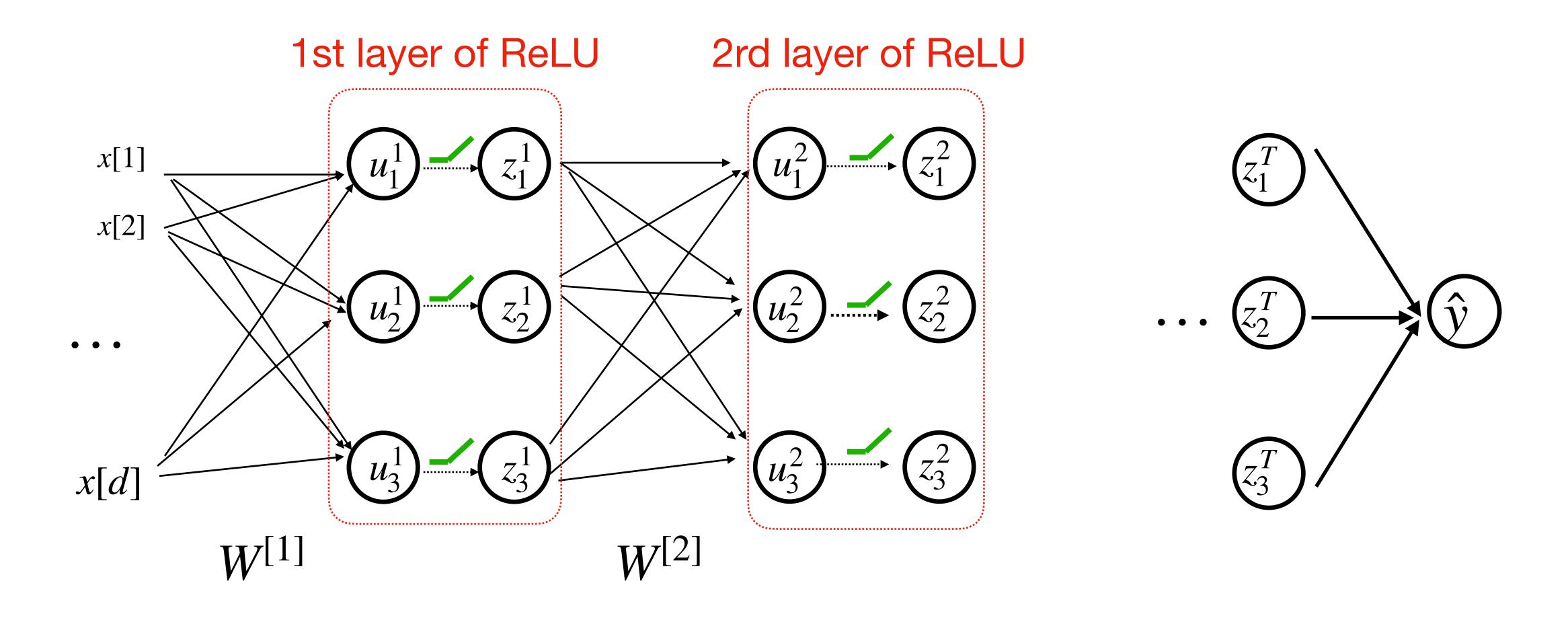


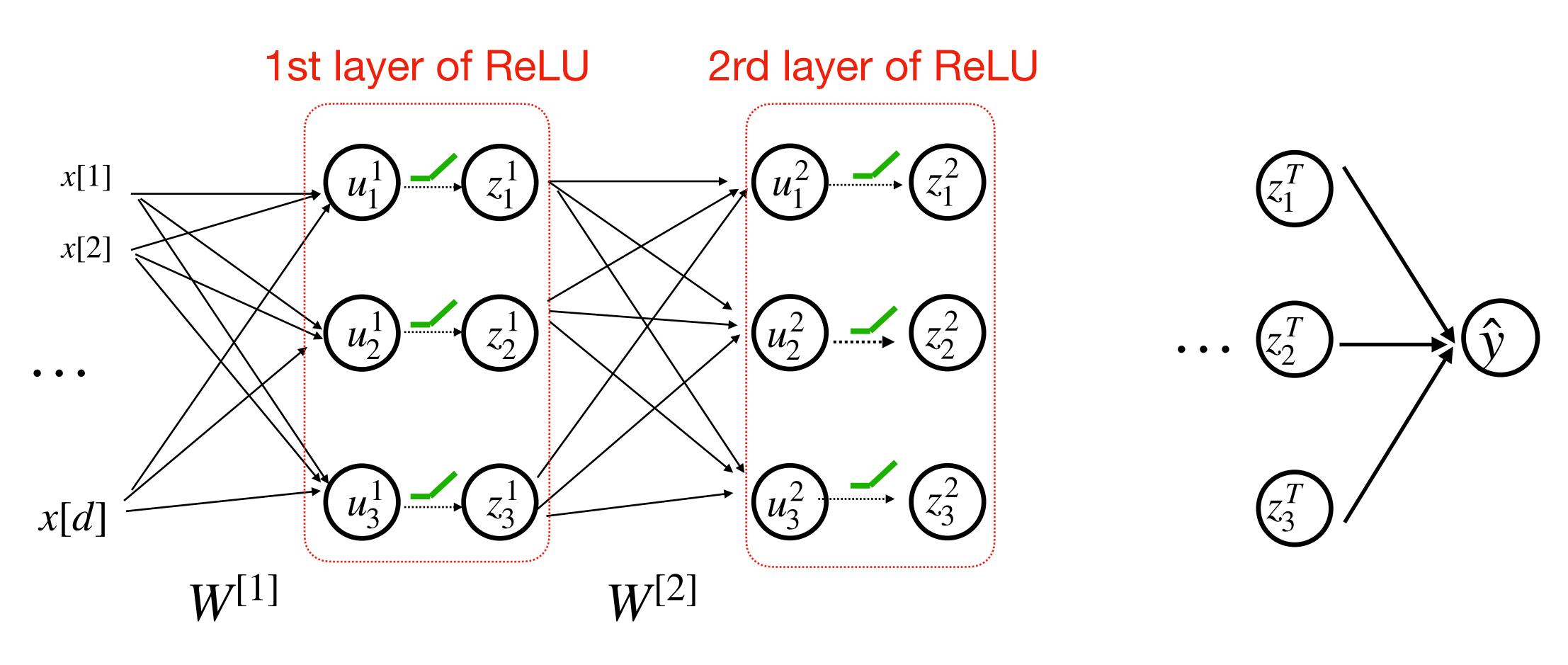
Store input & output of all neurons

backward pass:

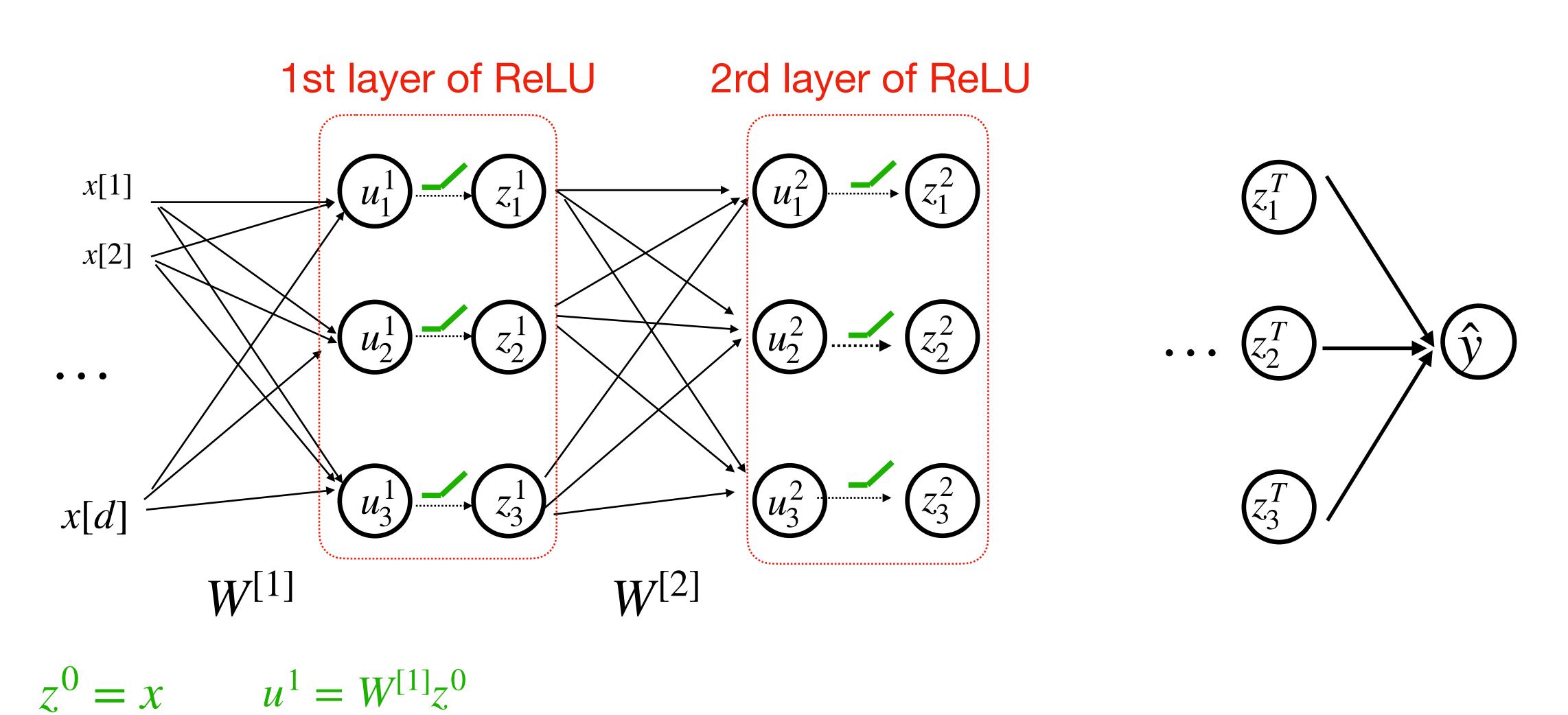


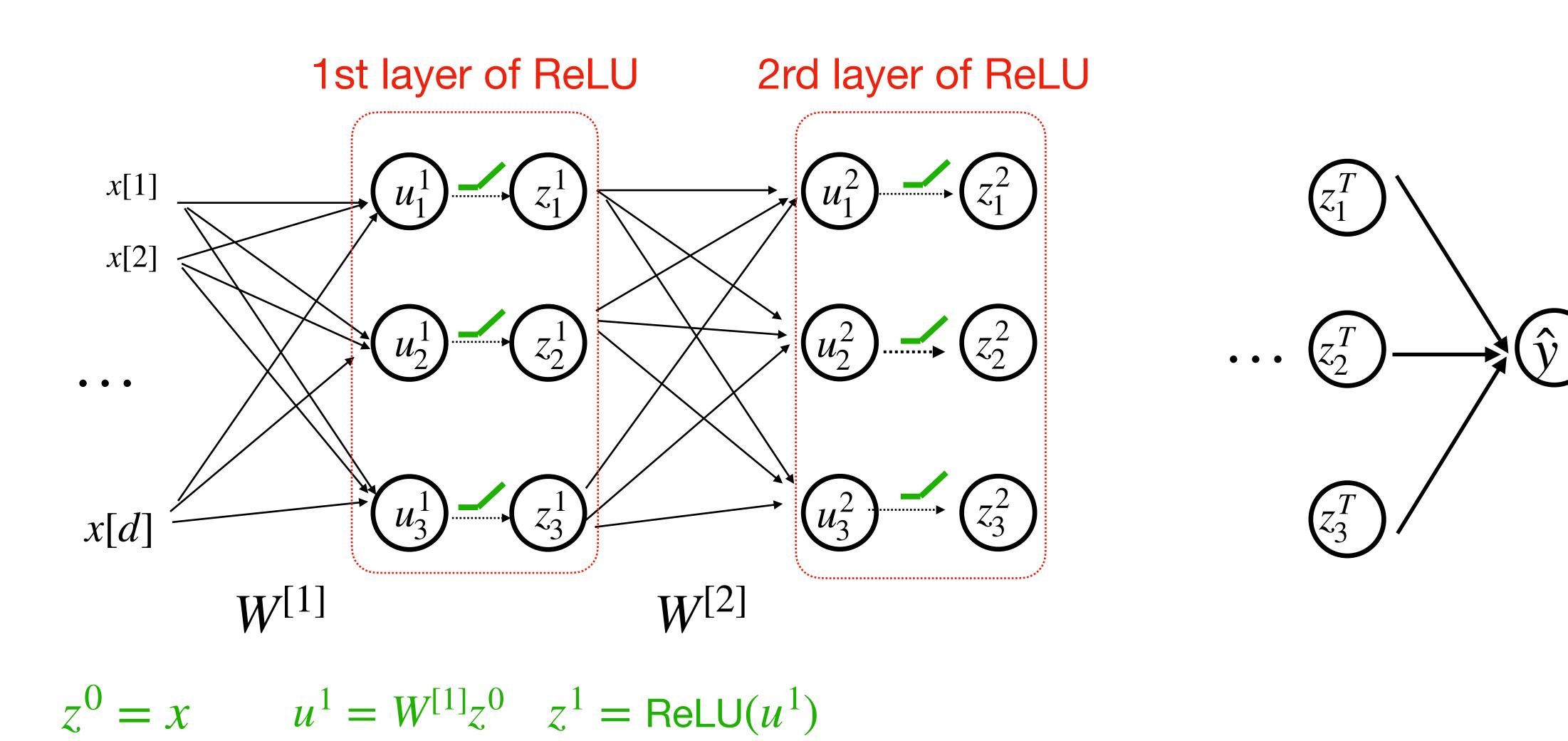
Compute derivatives

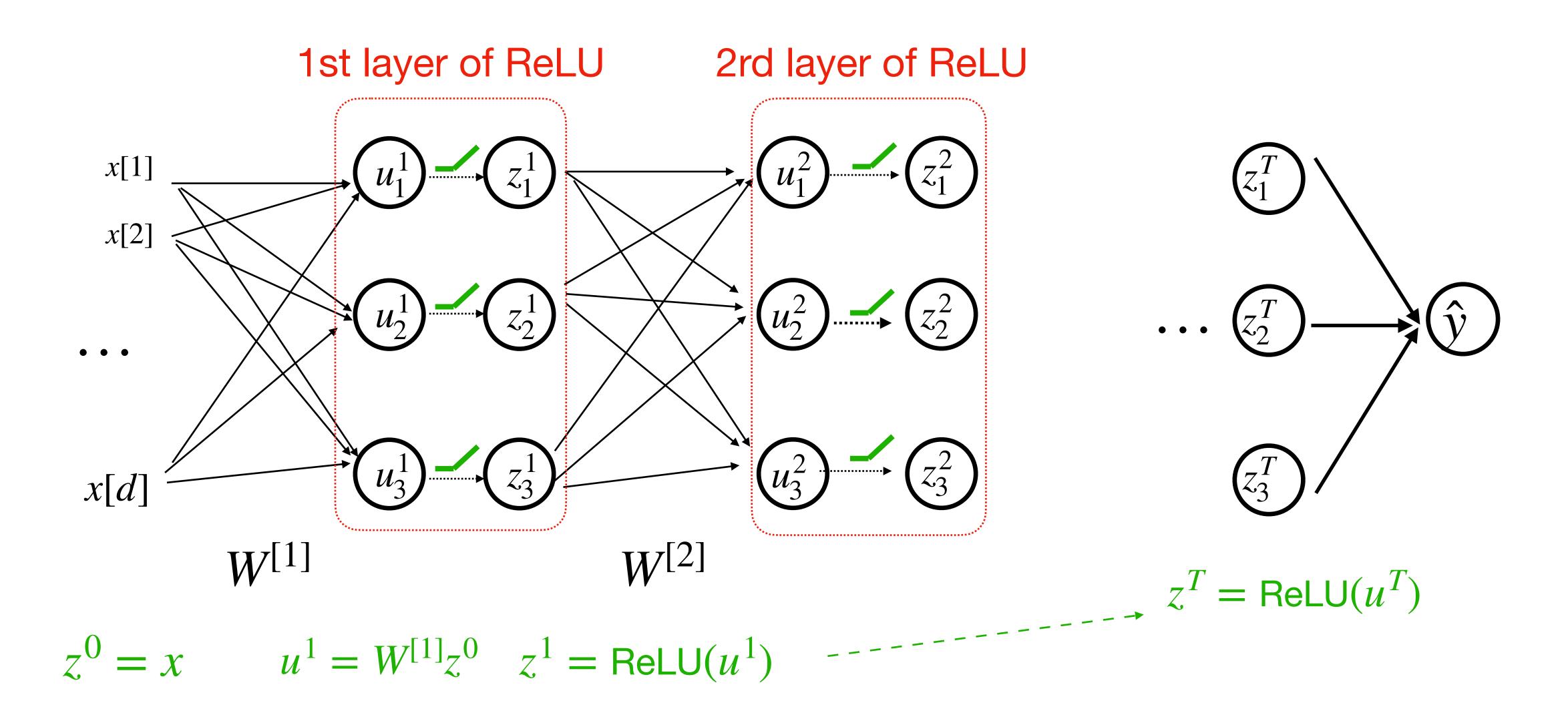


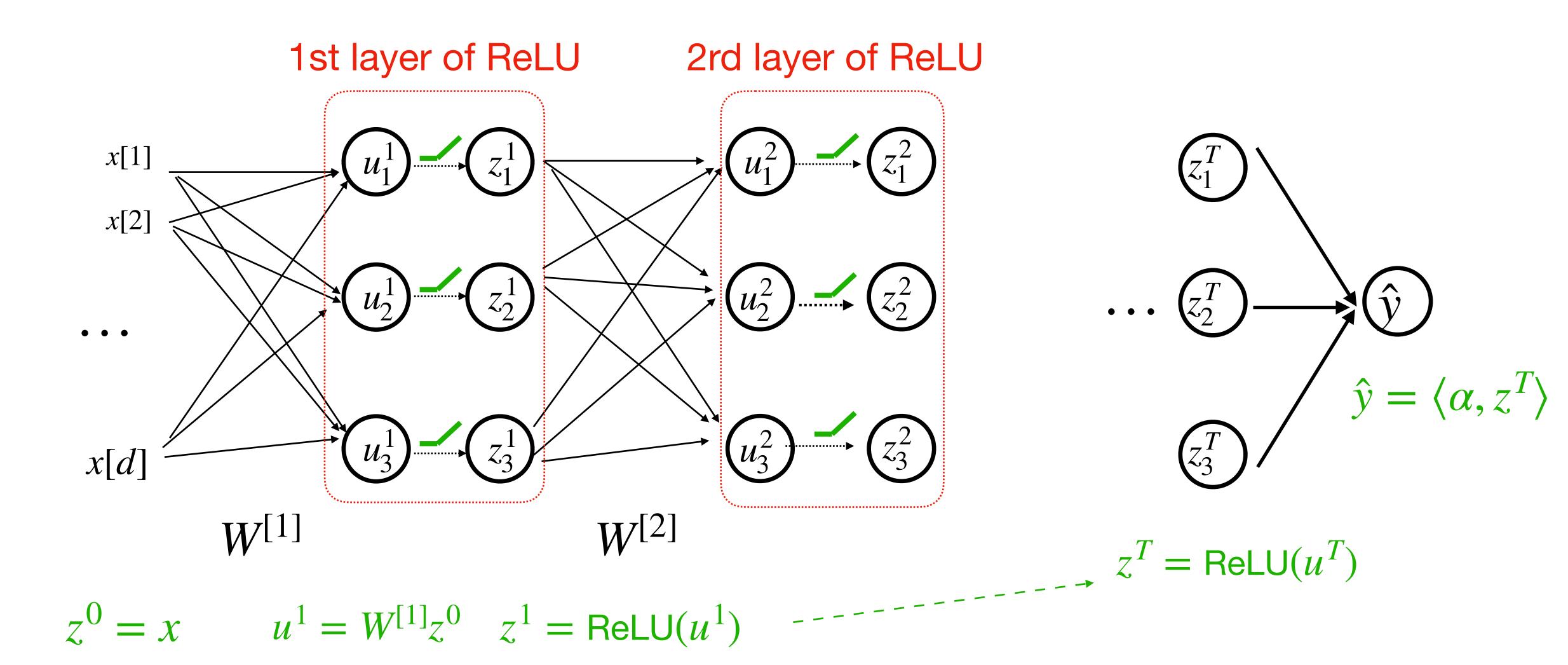


$$z^0 = x$$

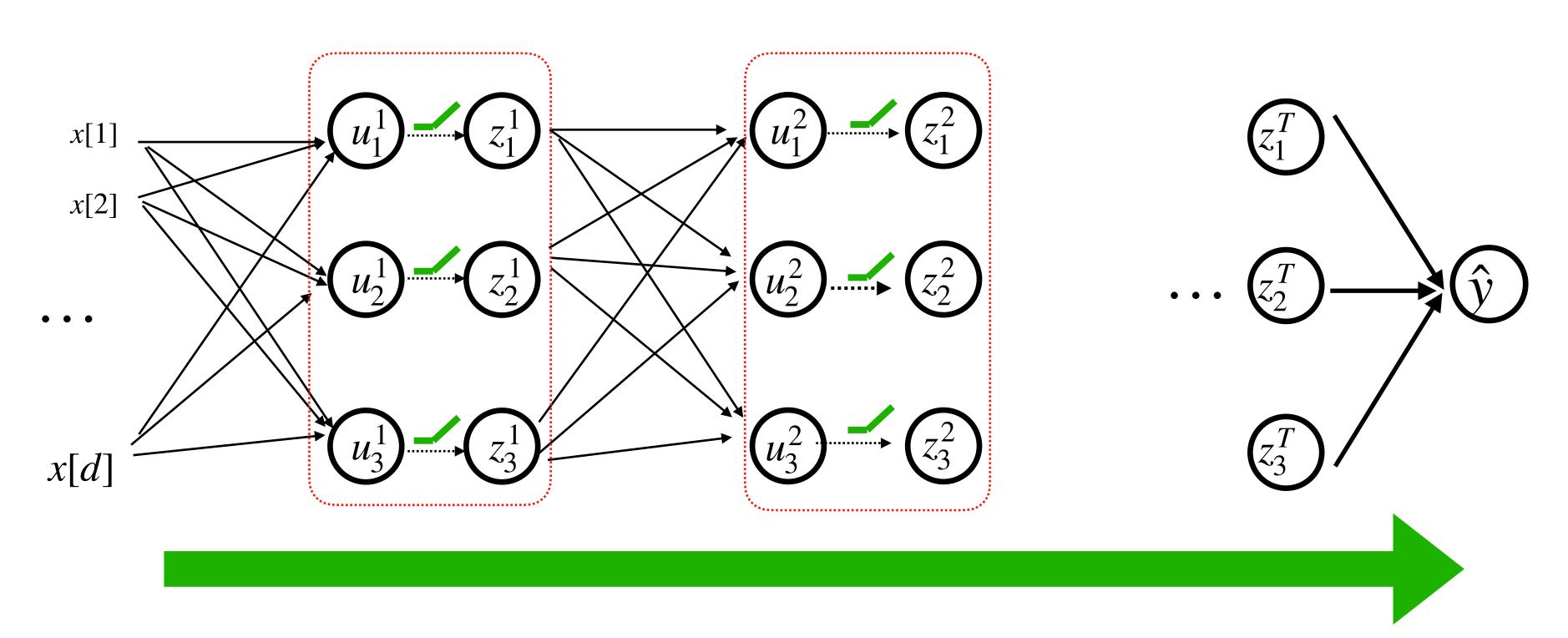






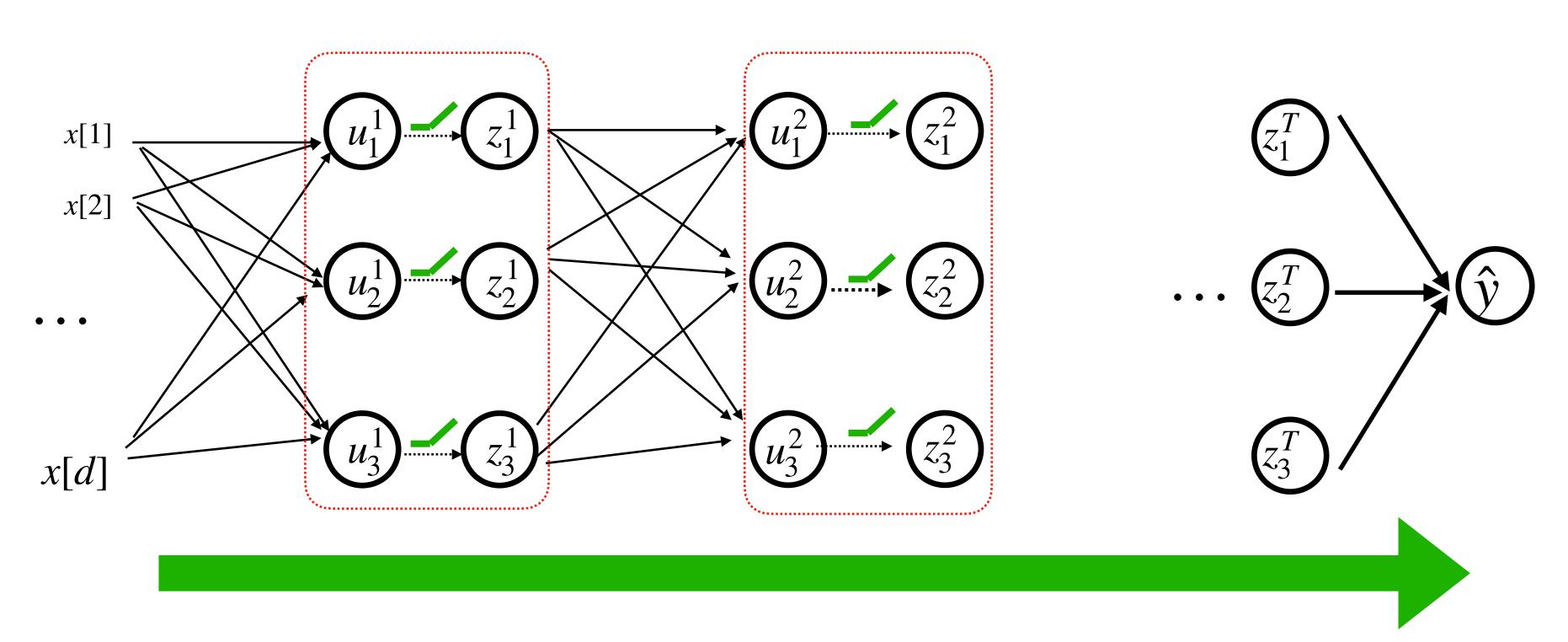


Summary of the forward pass



All nodes' values (i.e., z, u, \hat{y}) are computed and stored

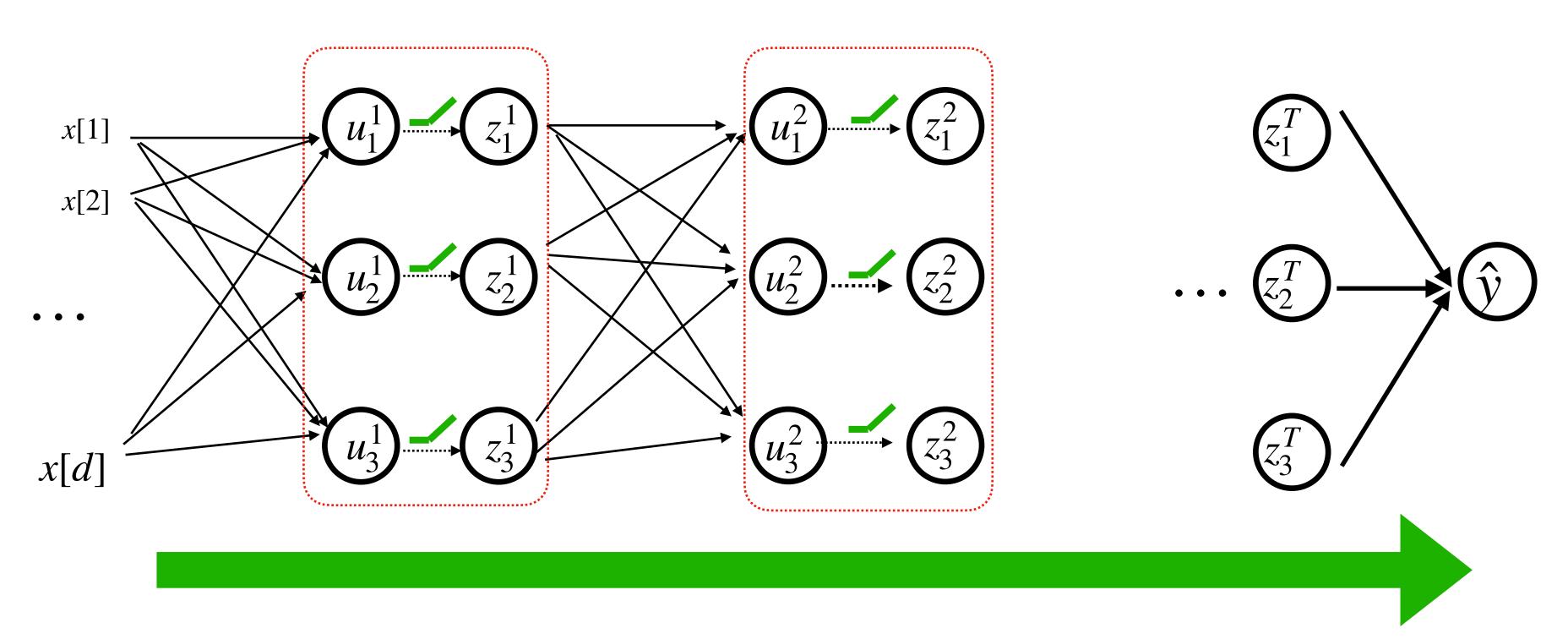
Summary of the forward pass



All nodes' values (i.e., z, u, \hat{y}) are computed and stored

Q: what is the computation complexity of the forward pass?

Summary of the forward pass



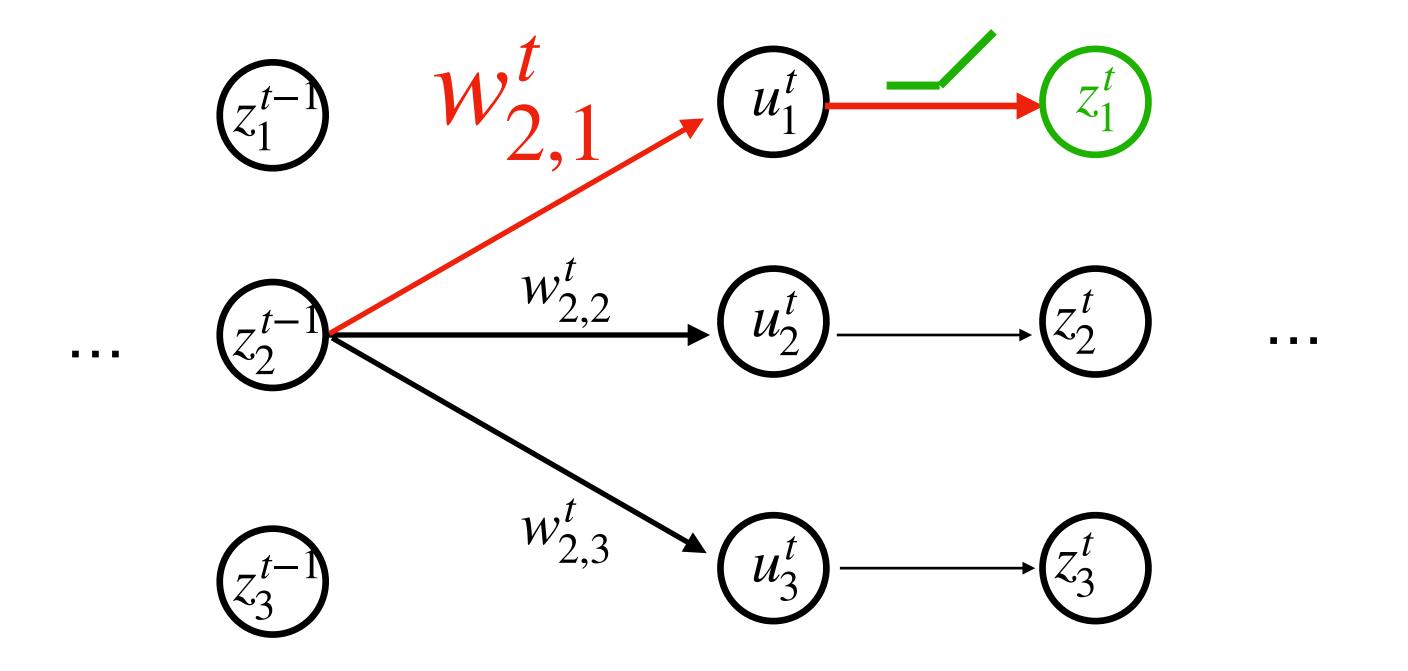
All nodes' values (i.e., z, u, \hat{y}) are computed and stored

Q: what is the computation complexity of the forward pass?

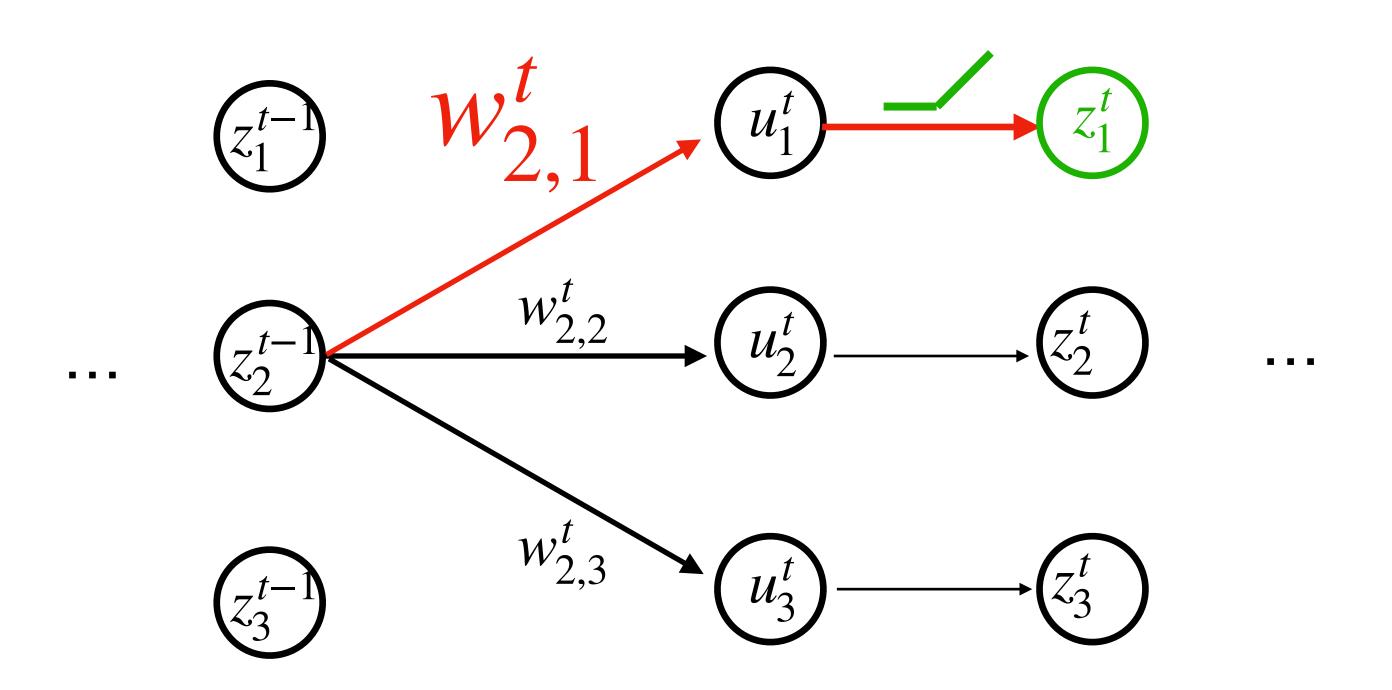
A: linear in # of Edges + # of nodes

Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.

Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.



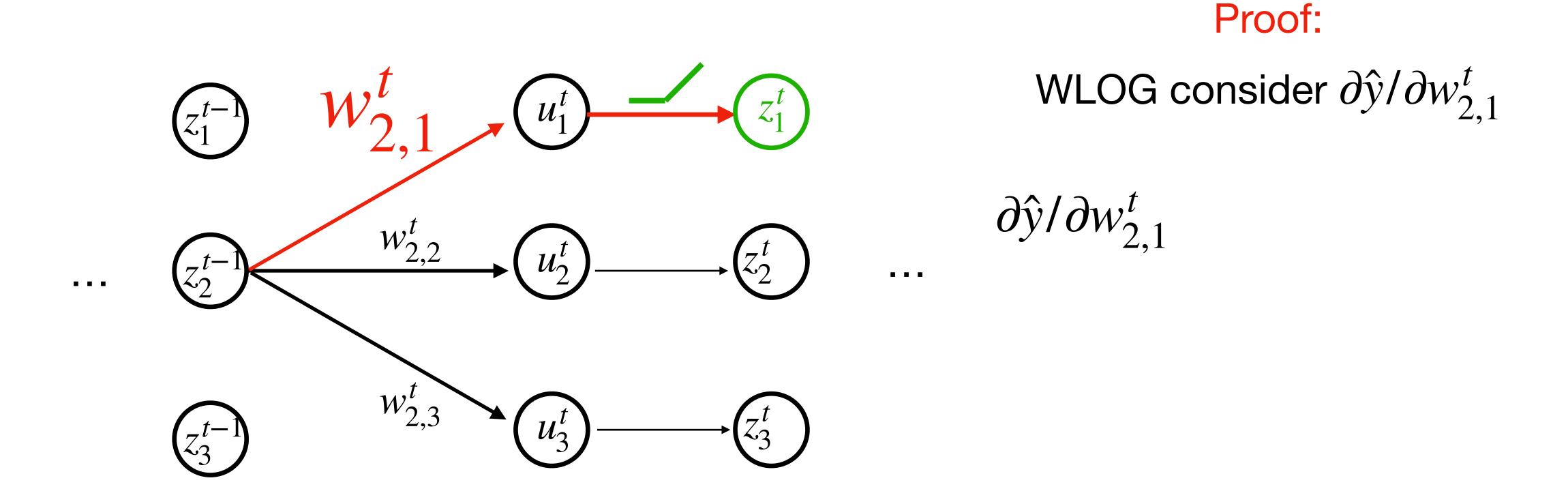
Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.



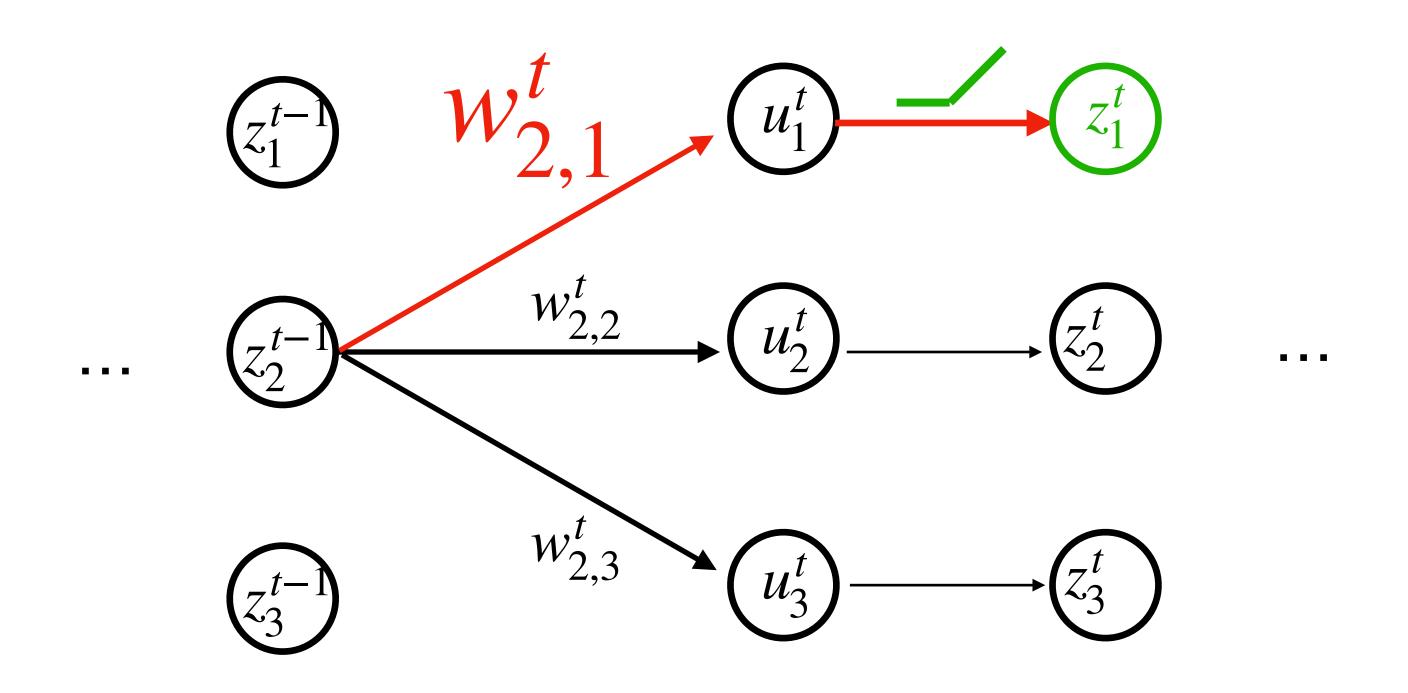
Proof:

WLOG consider $\partial \hat{y}/\partial w_{2,1}^t$

Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.



Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.

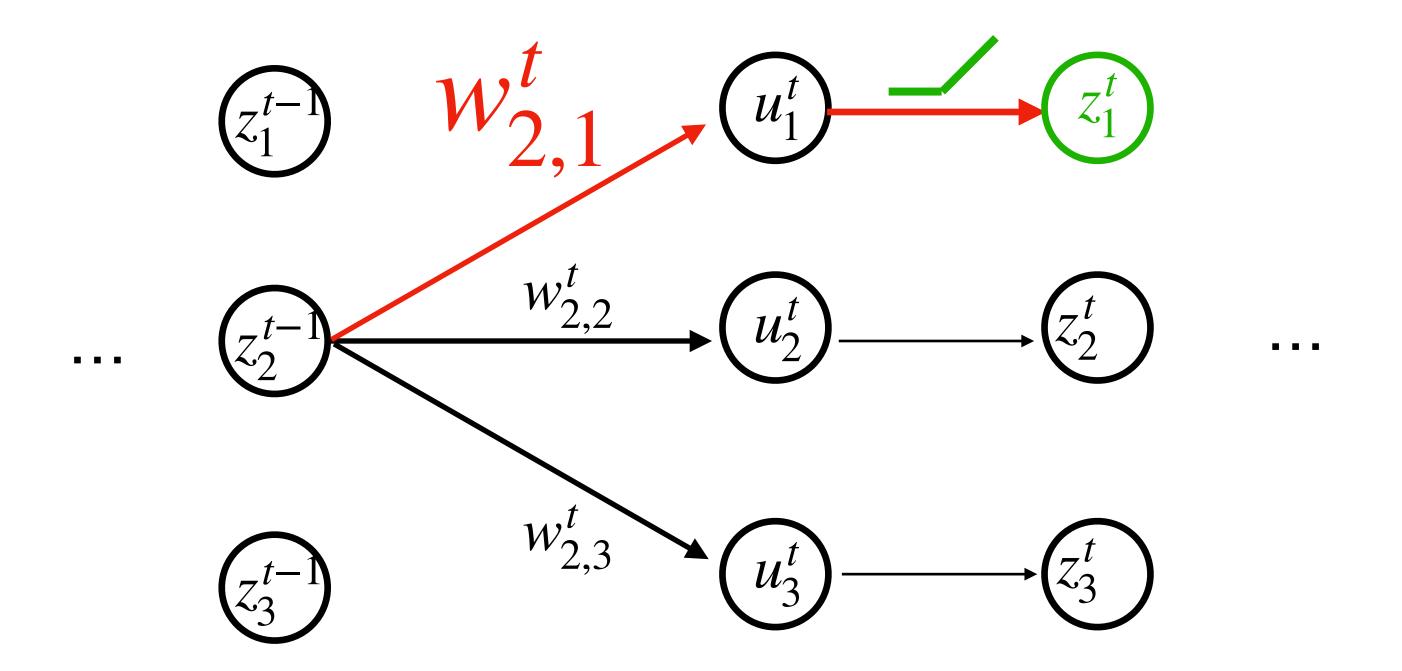


Proof:

WLOG consider $\partial \hat{y}/\partial w_{2,1}^t$

$$\partial \hat{y}/\partial w_{2,1}^t = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$

Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.

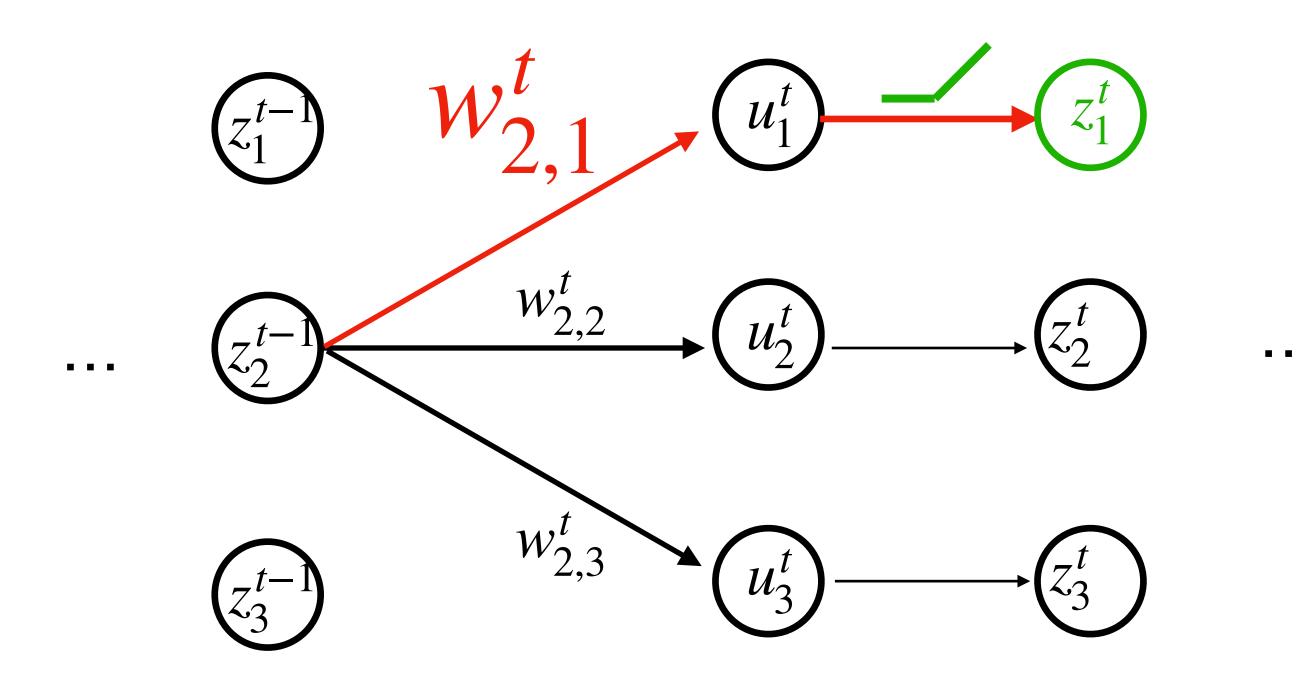


Proof:

WLOG consider $\partial \hat{y}/\partial w_{2,1}^t$

$$\frac{\partial \hat{y}}{\partial w_{2,1}^t} = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$
$$= \frac{\partial \hat{y}}{\partial z_1^t} \cdot \sigma'(u_1^t) \cdot z_2^{t-1}$$

Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.



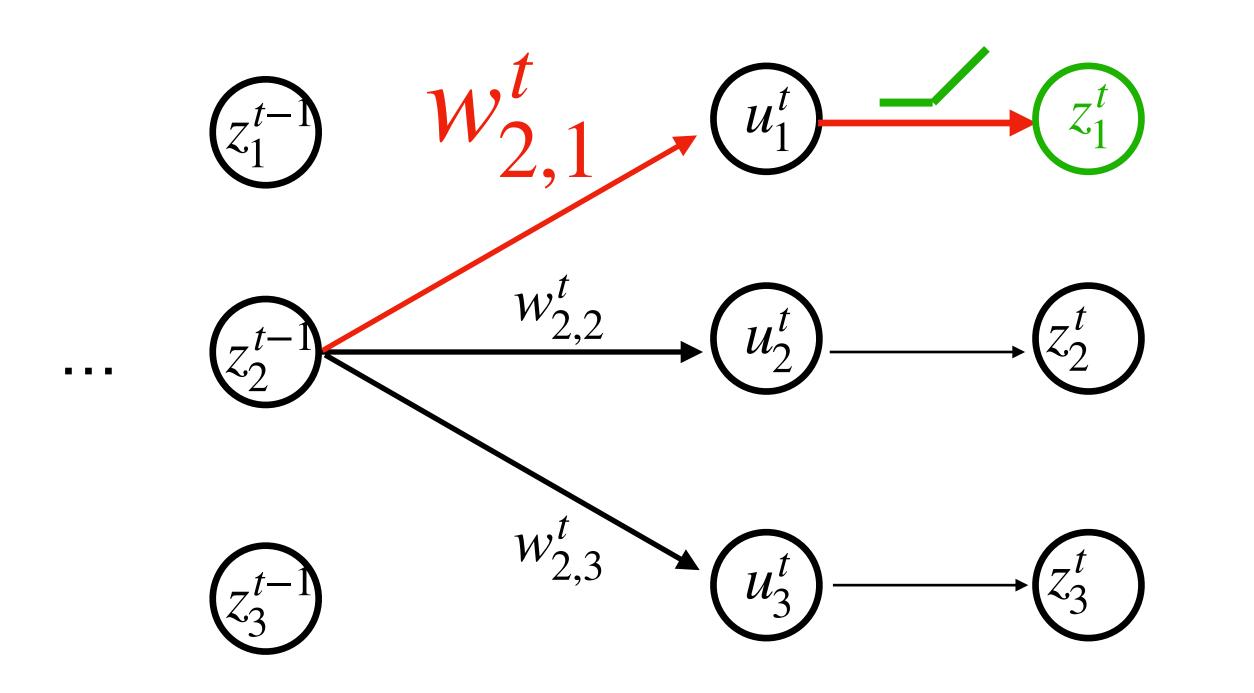
Proof:

WLOG consider $\partial \hat{y}/\partial w_{2,1}^t$

$$\frac{\partial \hat{y}}{\partial w_{2,1}^t} = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$
$$= \left(\frac{\partial \hat{y}}{\partial z_1^t}\right) \sigma'(u_1^t) \cdot z_2^{t-1}$$

Given by assumption

Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.



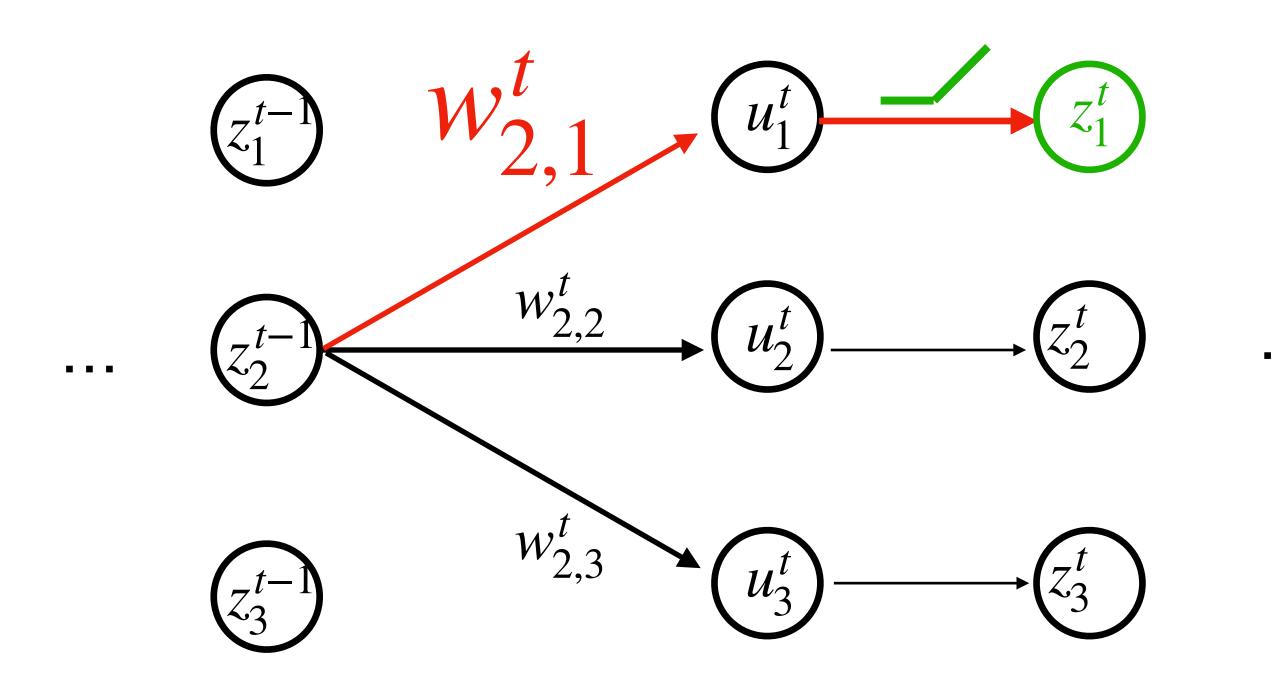
Proof:

WLOG consider $\partial \hat{y}/\partial w_{2,1}^t$

$$\frac{\partial \hat{y}}{\partial w_{2,1}^t} = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$
$$= \left(\frac{\partial \hat{y}}{\partial x_1^t}\right) \left(\sigma'(u_1^t) \cdot z_2^{t-1}\right)$$

Given by Derivative of assumption ReLU

Claim: to compute $\partial \hat{y}/\partial w$, \forall edge w, it suffices to compute $\partial \hat{y}/\partial z$, \forall node z.



Proof:

WLOG consider $\partial \hat{y}/\partial w_{2,1}^t$

$$\frac{\partial \hat{y}}{\partial w_{2,1}^t} = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$

$$= \left(\frac{\partial \hat{y}}{\partial z_1^t}\right) \left(\sigma'(u_1^t) \cdot z_2^{t-1}\right)$$

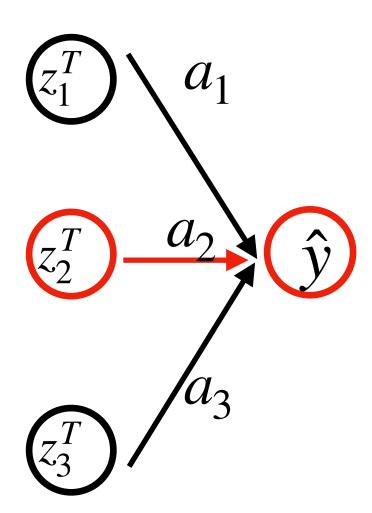
Known from forward pass

Given by Derivative of assumption ReLU

We compute $\partial \hat{y}/\partial z^t$ backwards in time from t=T to t=1:

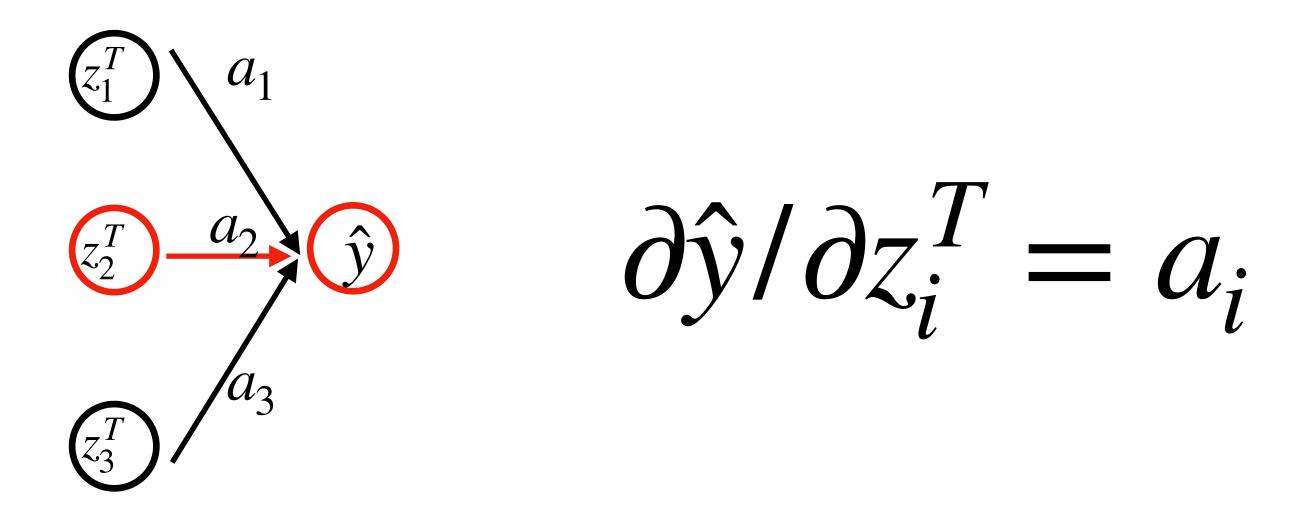
The backward Pass: base case

Base case: compute $\partial \hat{y}/\partial z^T$, for all node z at T-th Layer



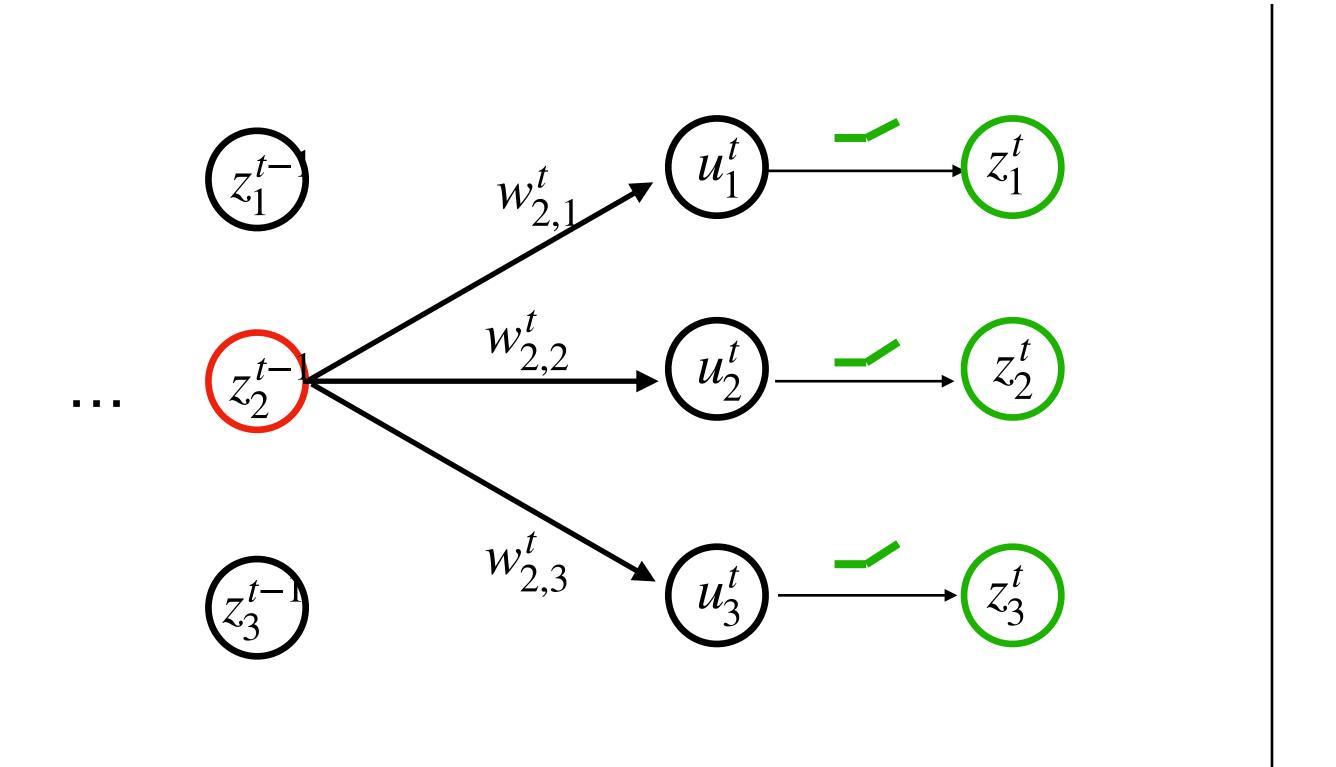
The backward Pass: base case

Base case: compute $\partial \hat{y}/\partial z^T$, for all node z at T-th Layer



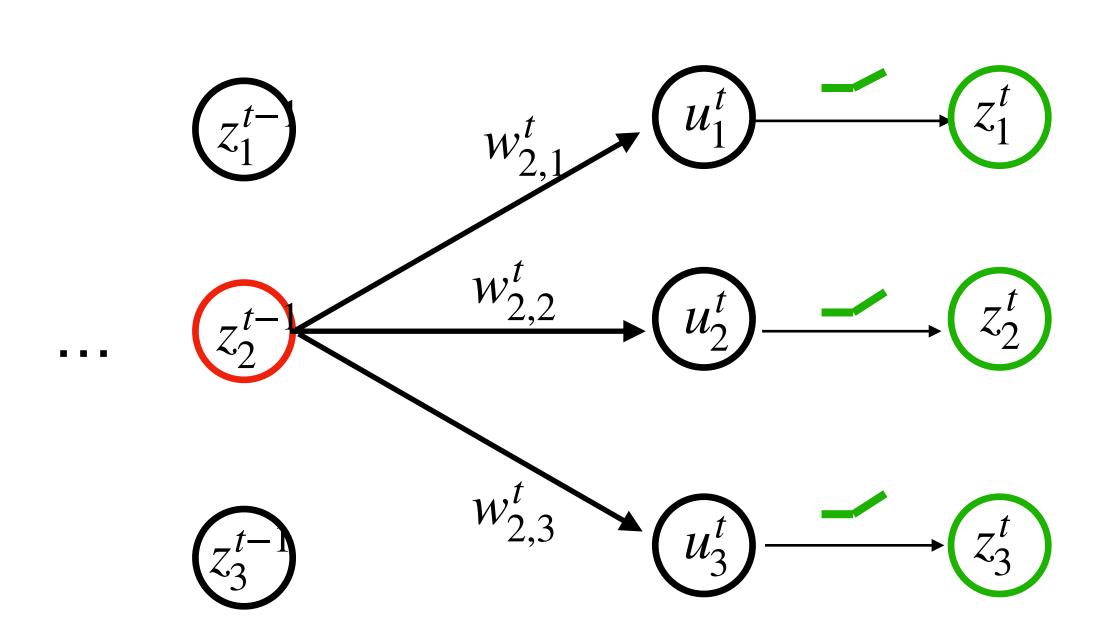
The backward Pass: induction step

Assume that we have computed $\partial \hat{y}/\partial z_i^t$, $\forall i$



The backward Pass: induction step

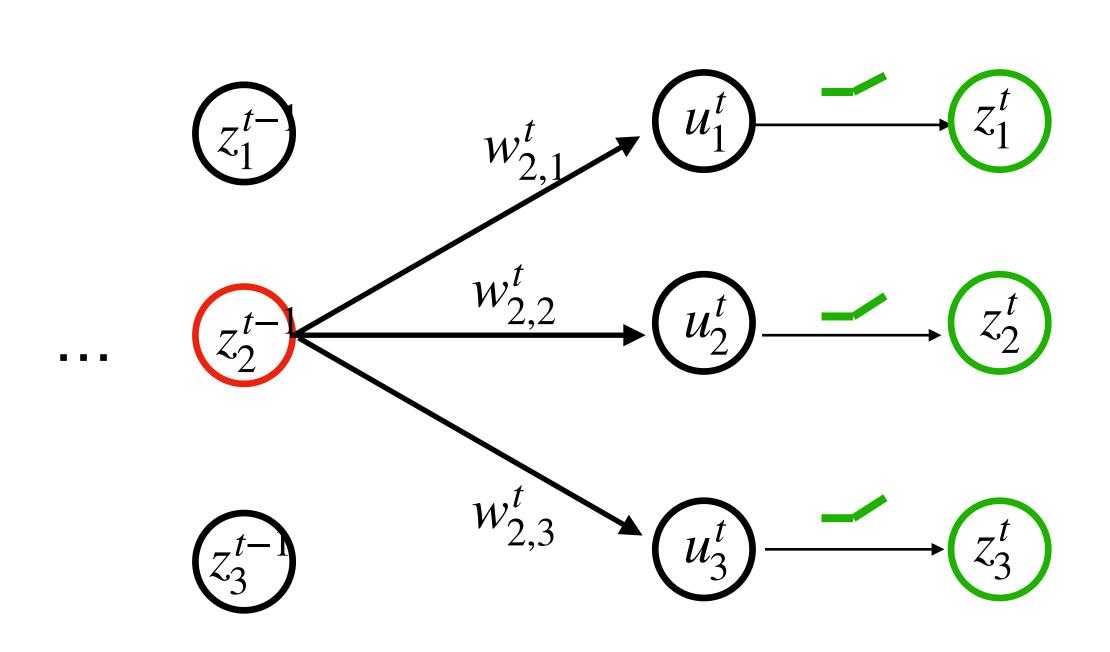
Assume that we have computed $\partial \hat{y}/\partial z_i^t$, $\forall i$



WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

The backward Pass: induction step

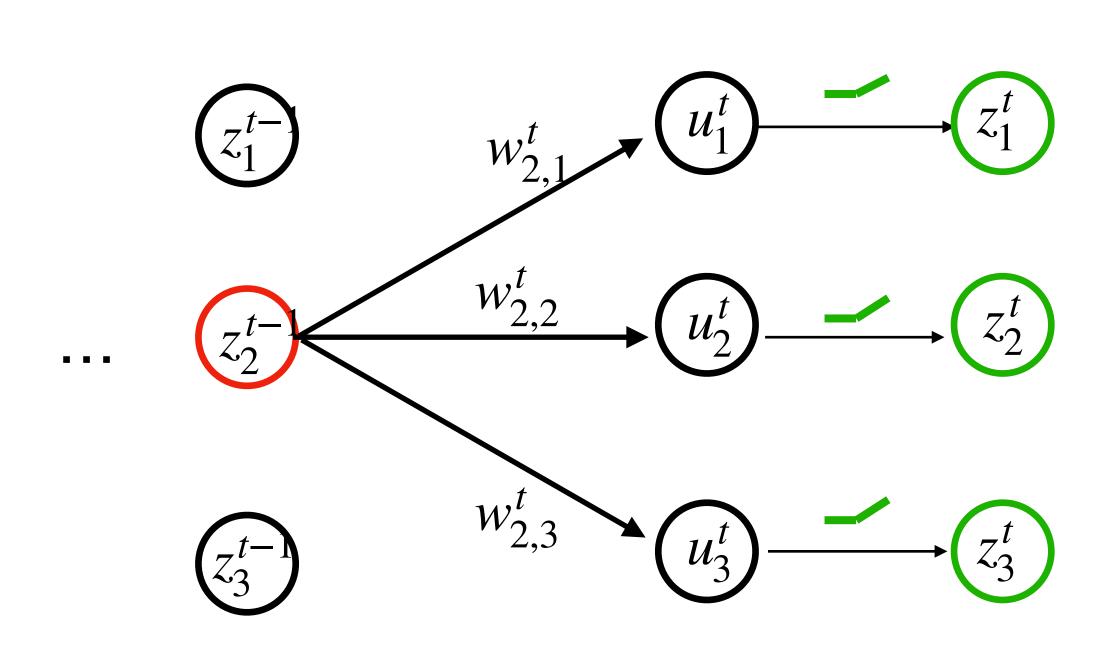
Assume that we have computed $\partial \hat{y}/\partial z_i^t$, $\forall i$



WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

Step 1: for all i,
$$\frac{\partial \hat{y}}{\partial u_i^t} = \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t}$$

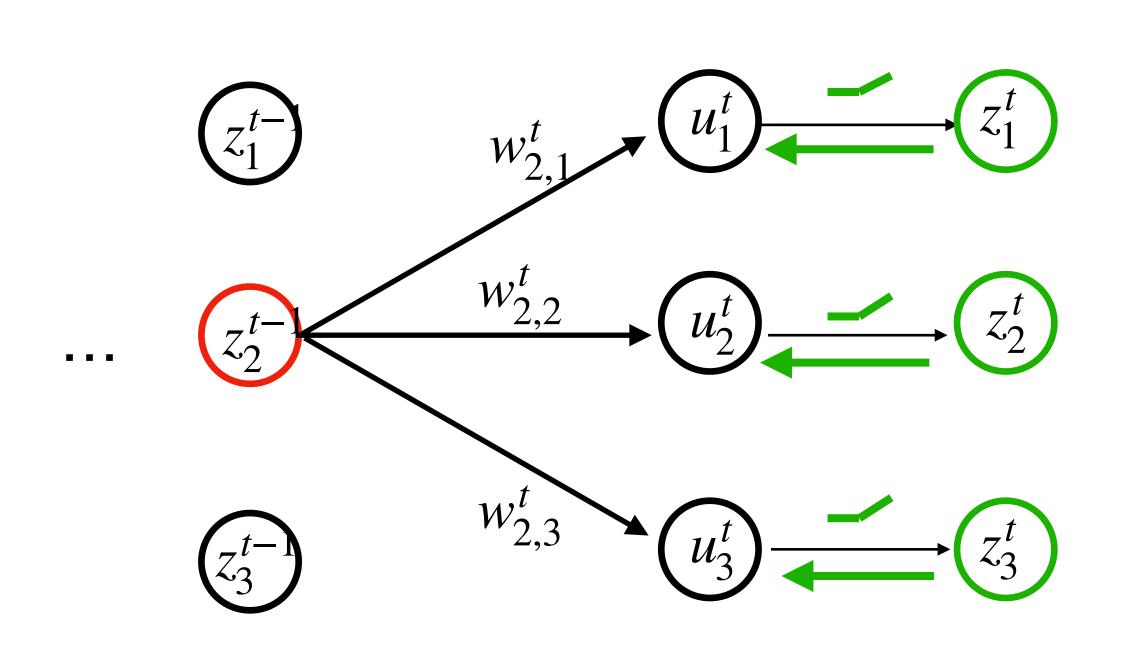
Assume that we have computed $\partial \hat{y}/\partial z_i^t$, $\forall i$



WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

Step 1: for all i,
$$\frac{\partial \hat{y}}{\partial u_i^t} = \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t}$$
$$= \frac{\partial \hat{y}}{\partial z_i^t} \cdot \sigma'(u_i^t)$$

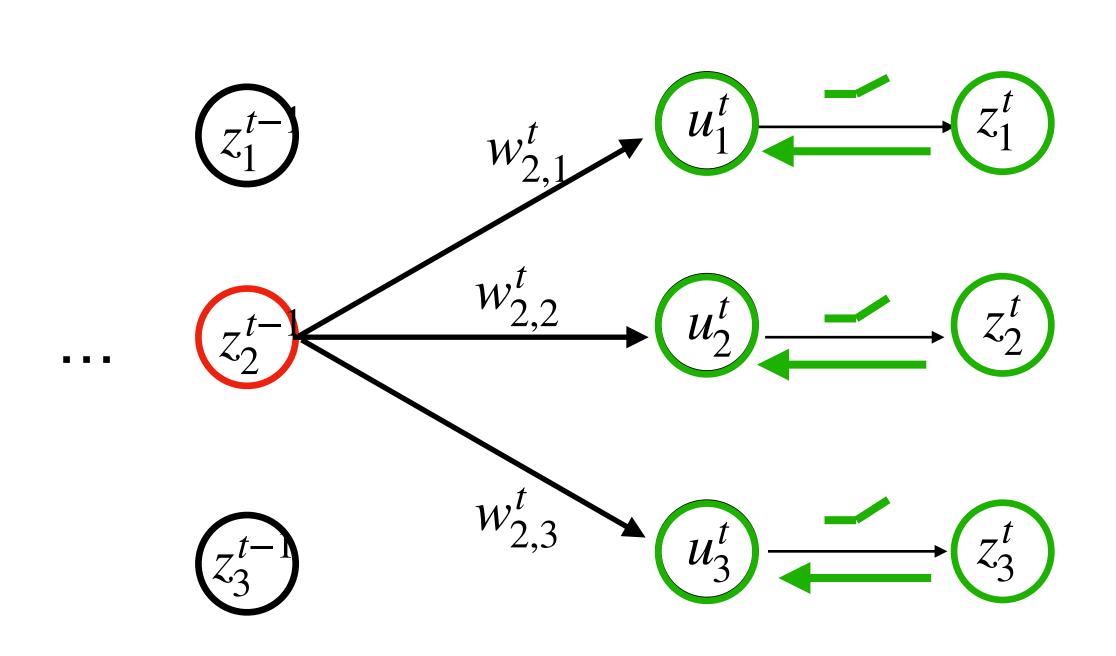
Assume that we have computed $\partial \hat{y}/\partial z_i^t$, $\forall i$



WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

Step 1: for all i,
$$\frac{\partial \hat{y}}{\partial u_i^t} = \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t}$$
$$= \frac{\partial \hat{y}}{\partial z_i^t} \cdot \sigma'(u_i^t)$$

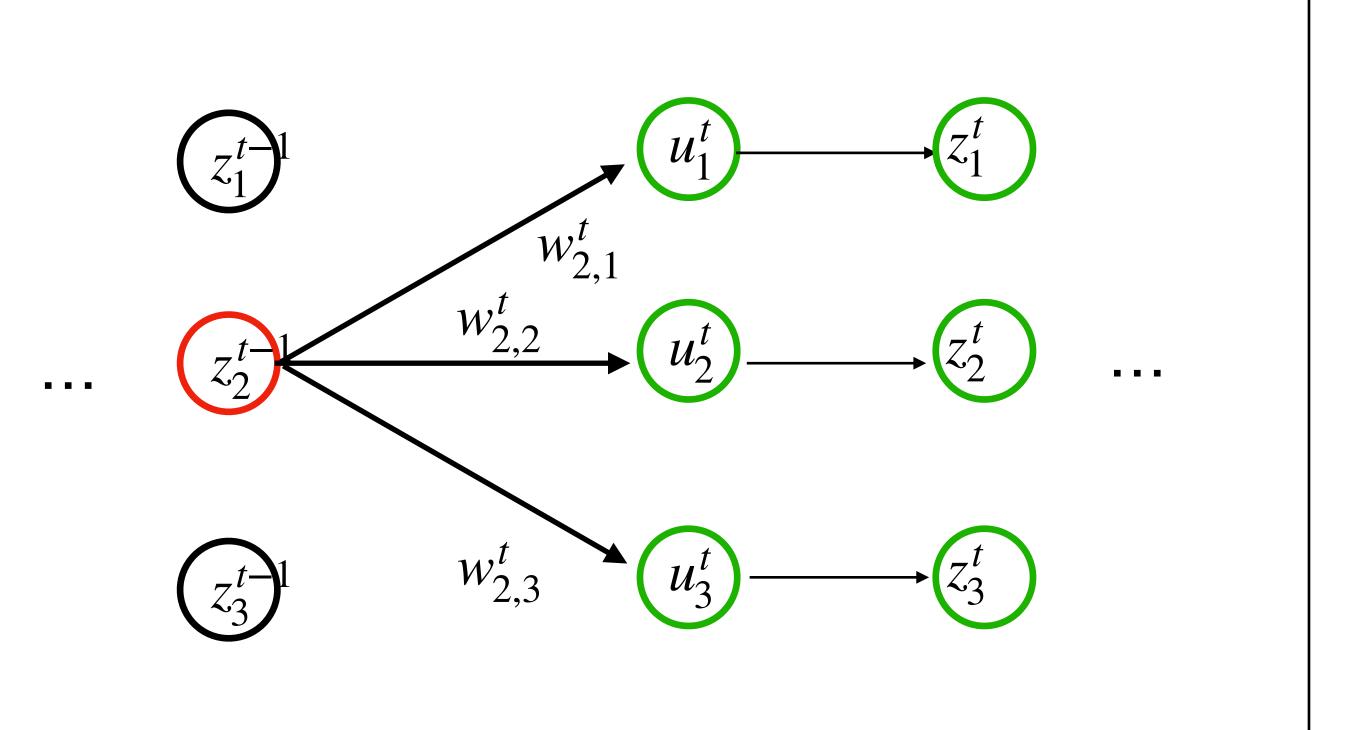
Assume that we have computed $\partial \hat{y}/\partial z_i^t$, $\forall i$



WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

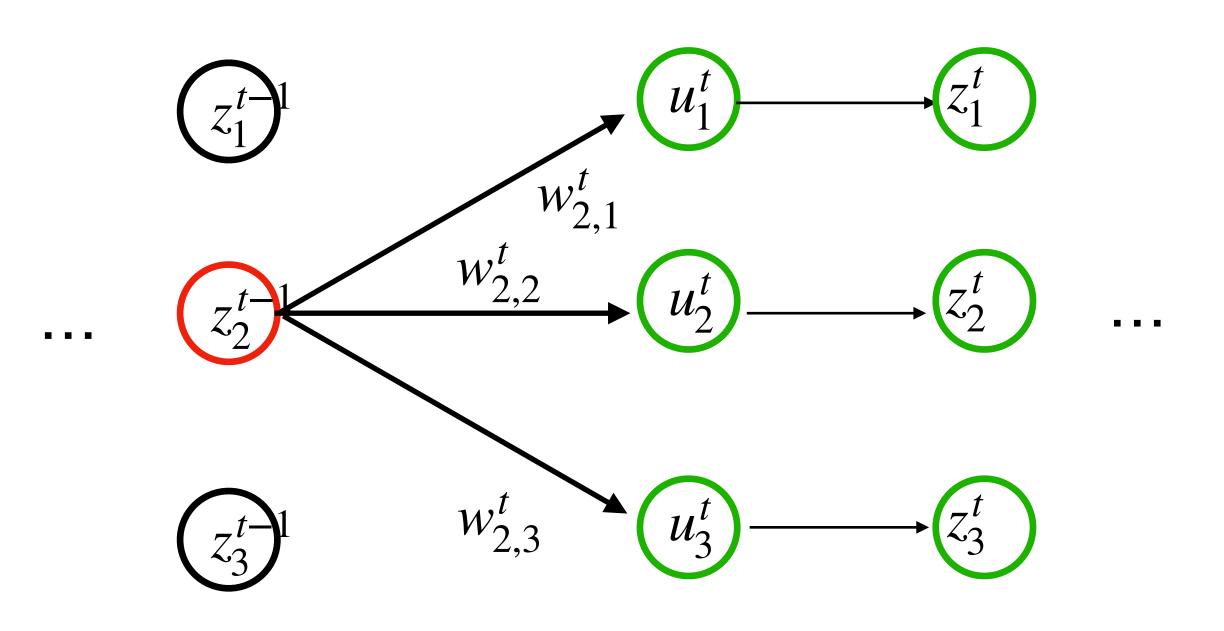
Step 1: for all i,
$$\frac{\partial \hat{y}}{\partial u_i^t} = \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t}$$
$$= \frac{\partial \hat{y}}{\partial z_i^t} \cdot \sigma'(u_i^t)$$

Assume that we have computed $\partial \mathcal{E}/\partial z_i^t$, $\forall i$



After step 1, we have $\partial \hat{y}/\partial u_i^t$, $\forall i$

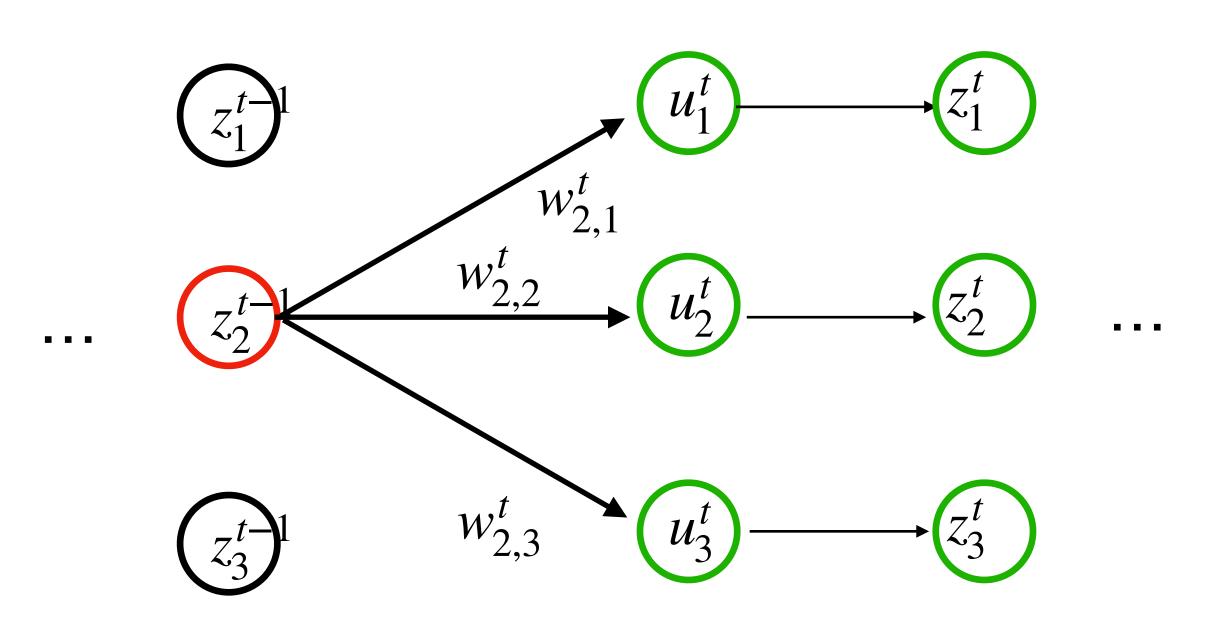
Assume that we have computed $\partial \mathcal{E}/\partial z_i^t$, $\forall i$



After step 1, we have $\partial \hat{y}/\partial u_i^t$, $\forall i$

Via multivariate chain rule:

Assume that we have computed $\partial \mathcal{C}/\partial z_i^t$, $\forall i$

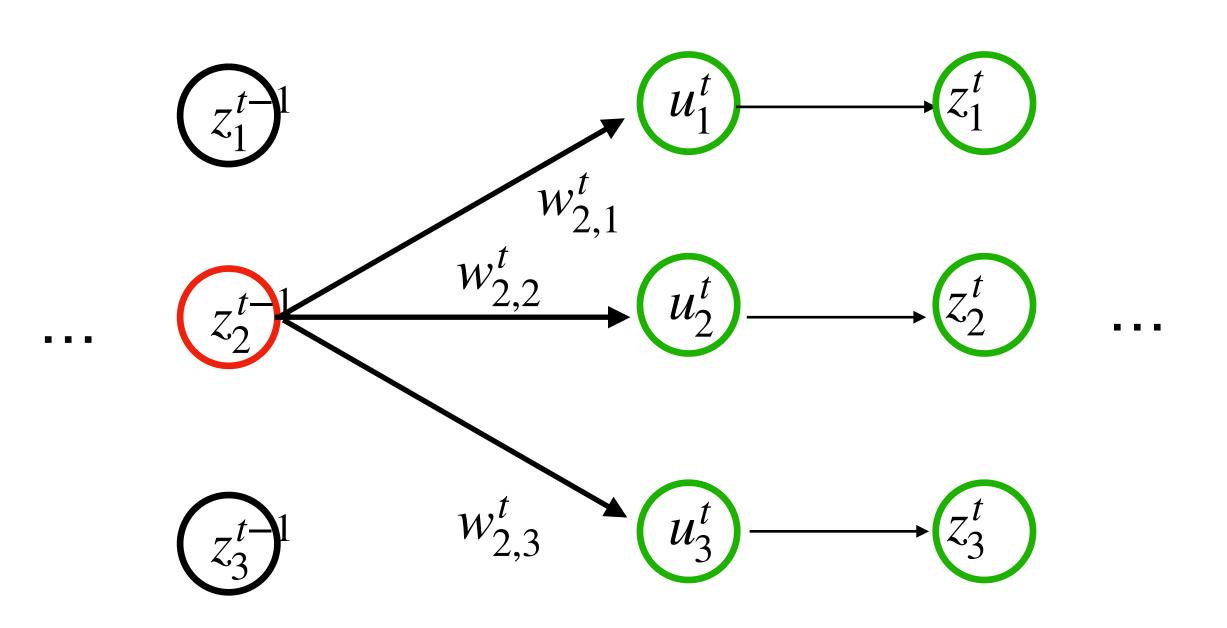


After step 1, we have $\partial \hat{y}/\partial u_i^t$, $\forall i$

Via multivariate chain rule:

Step 2:
$$\frac{\partial \hat{y}}{\partial z_2^{t-1}} = \sum_{i=1}^{K} \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}}$$

Assume that we have computed $\partial \mathcal{C}/\partial z_i^t$, $\forall i$

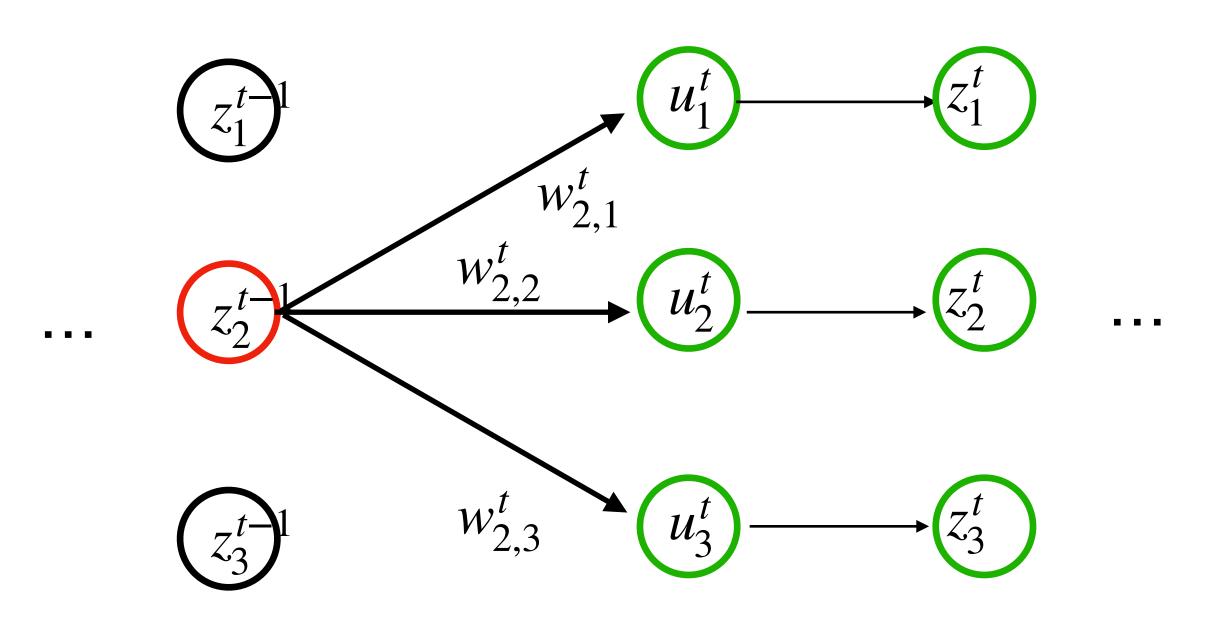


After step 1, we have $\partial \hat{y}/\partial u_i^t$, $\forall i$

Via multivariate chain rule:

Step 2:
$$\frac{\partial \hat{y}}{\partial z_2^{t-1}} = \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}}$$
$$= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t$$

Assume that we have computed $\partial \mathcal{C}/\partial z_i^t$, $\forall i$



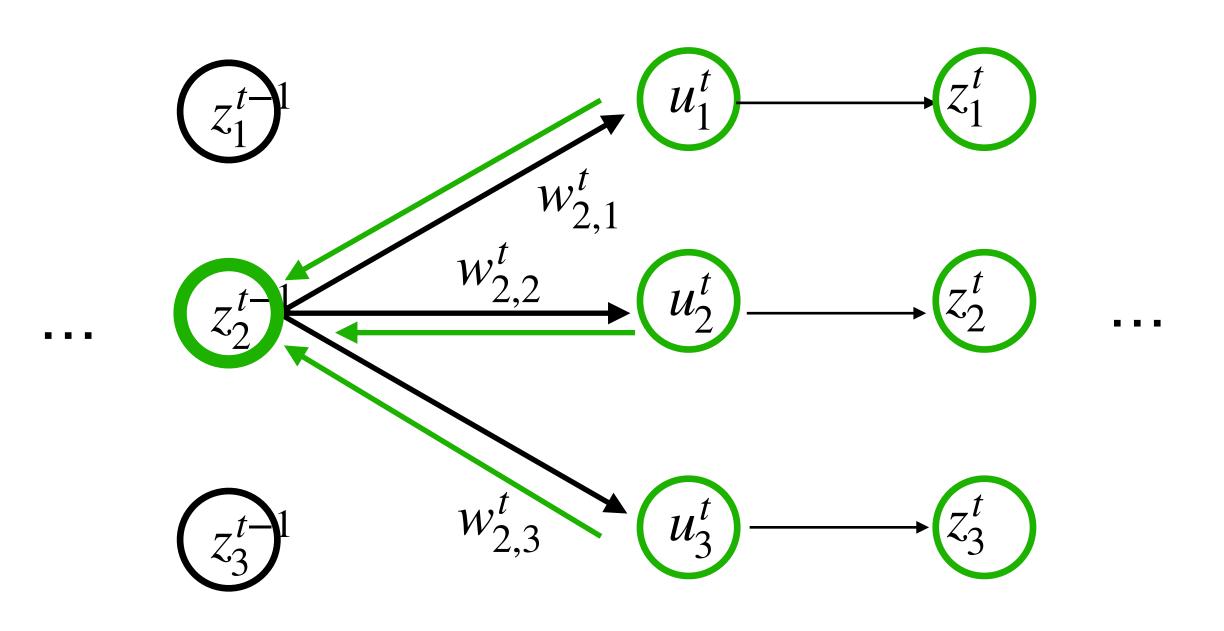
After step 1, we have $\partial \hat{y}/\partial u_i^t$, $\forall i$

Via multivariate chain rule:

Step 2:
$$\frac{\partial \hat{y}}{\partial z_2^{t-1}} = \sum_{i=1}^{K} \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}}$$
$$= \sum_{i=1}^{K} \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t$$

We are done at node z_2^{t-1} !

Assume that we have computed $\partial \mathcal{C}/\partial z_i^t$, $\forall i$



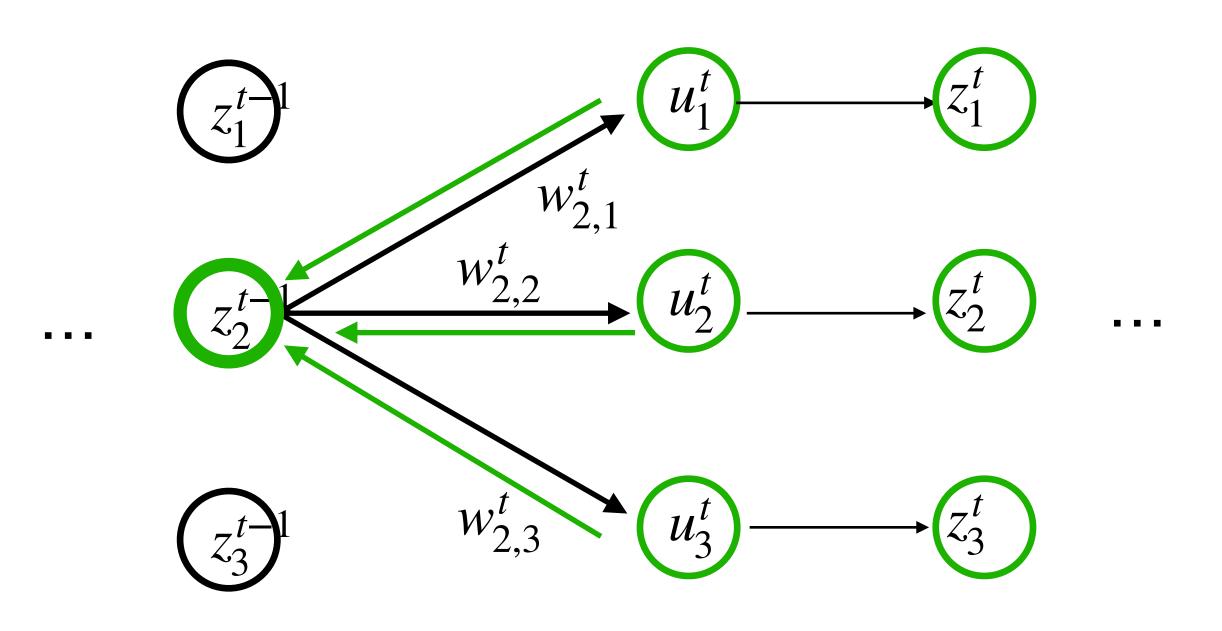
After step 1, we have $\partial \hat{y}/\partial u_i^t$, $\forall i$

Via multivariate chain rule:

Step 2:
$$\frac{\partial \hat{y}}{\partial z_2^{t-1}} = \sum_{i=1}^{K} \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}}$$
$$= \sum_{i=1}^{K} \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t$$

We are done at node z_2^{t-1} !

Assume that we have computed $\partial \mathcal{C}/\partial z_i^t$, $\forall i$



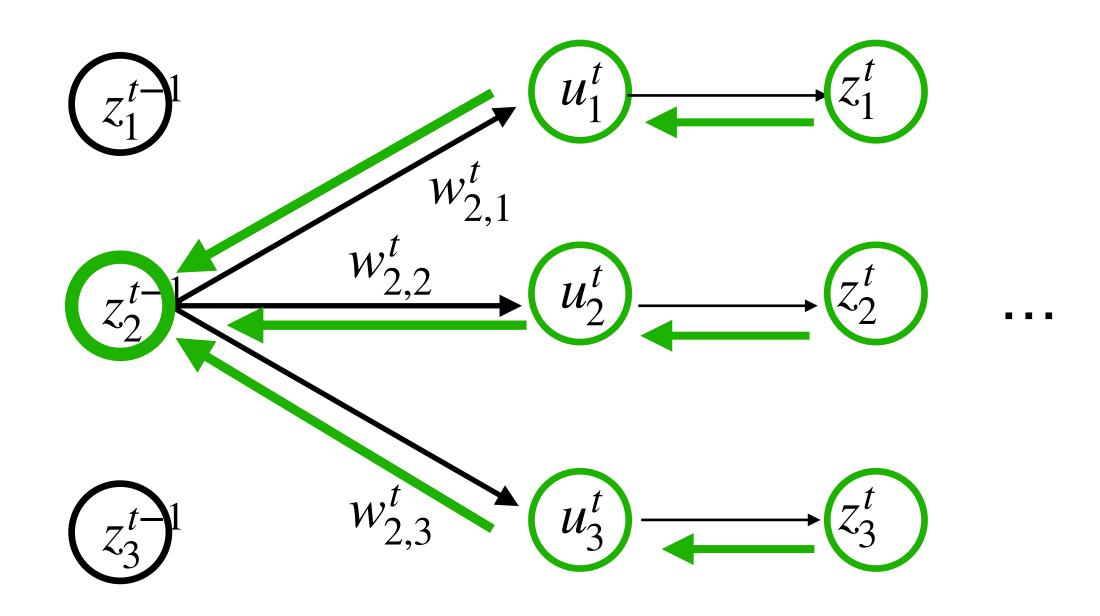
After step 1, we have $\partial \hat{y}/\partial u_i^t$, $\forall i$

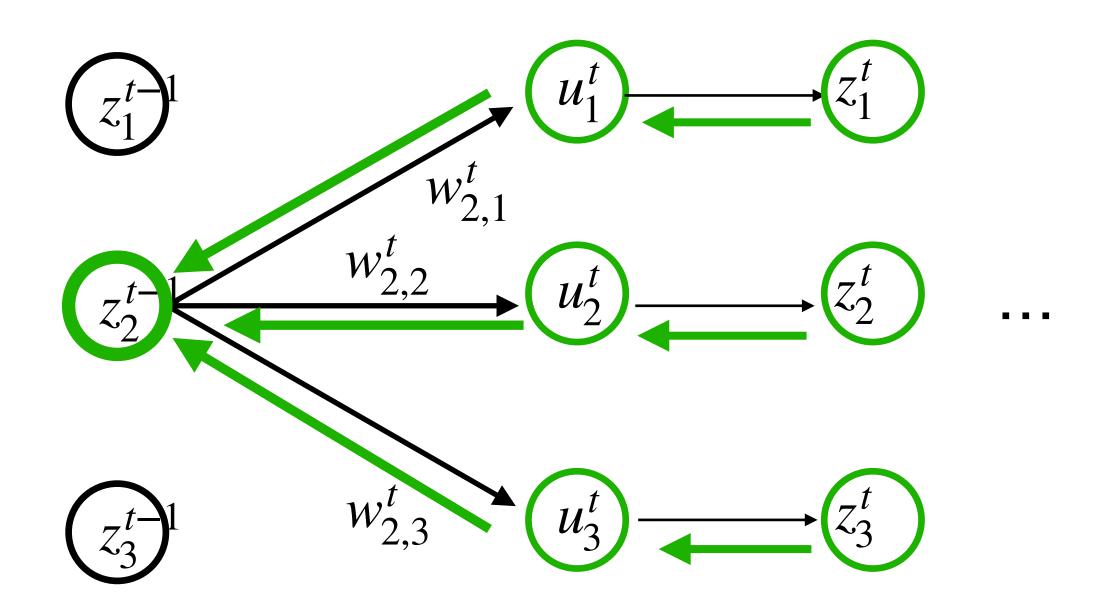
Via multivariate chain rule:

Step 2:
$$\frac{\partial \hat{y}}{\partial z_2^{t-1}} = \sum_{i=1}^{K} \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}}$$
$$= \sum_{i=1}^{K} \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t$$

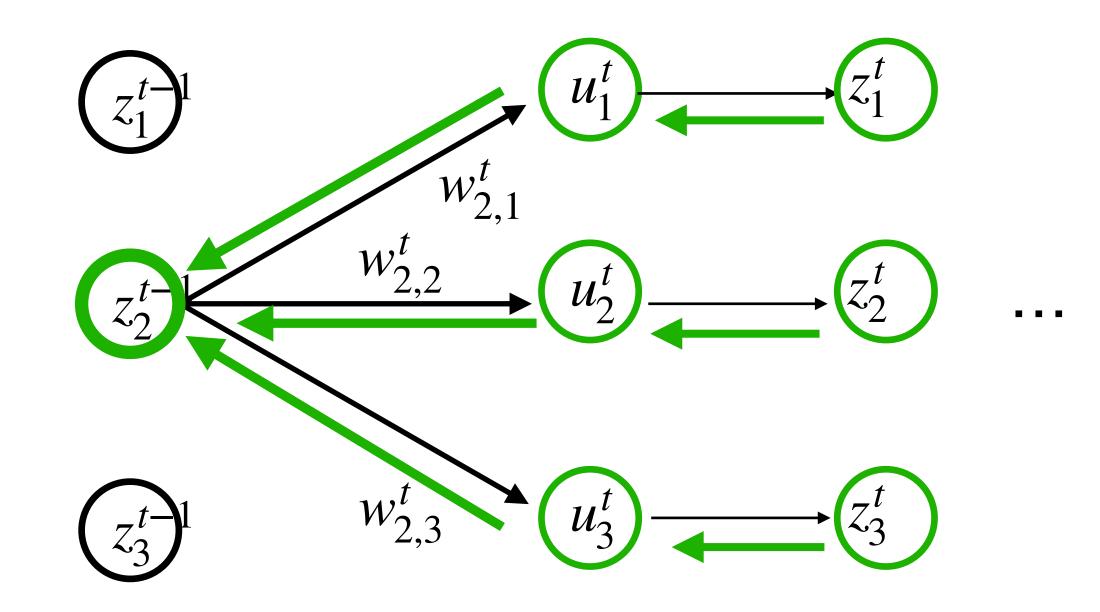
We are done at node z_2^{t-1} !

Repeat this for all z_i^{t-1} , $\forall i$



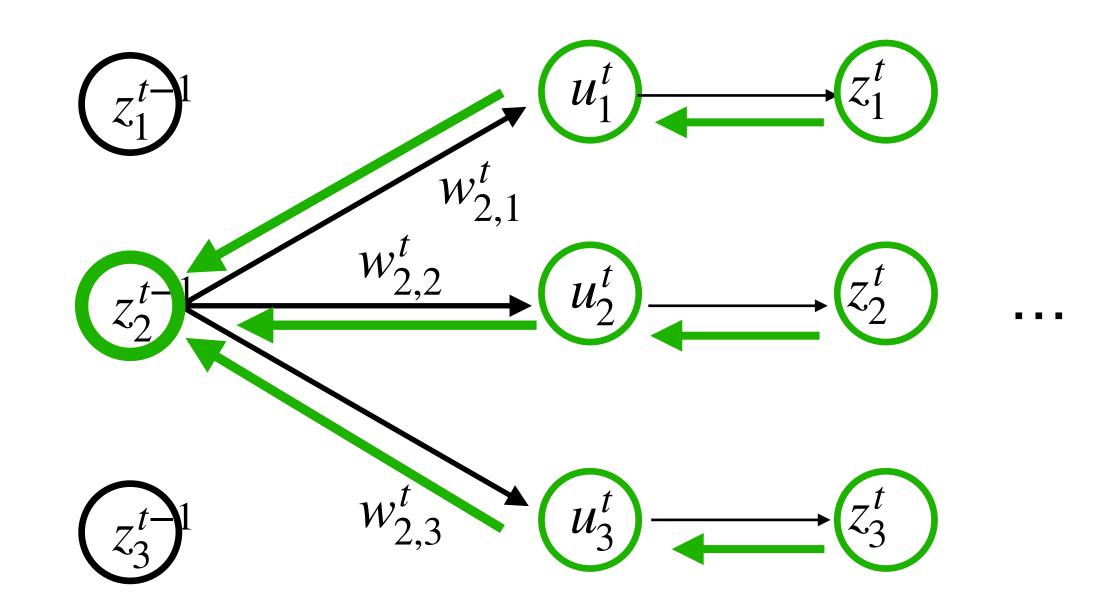


The computation from $\partial \hat{y}/z^t$ to $\partial \hat{y}/z^{t-1}$ is the # of all edges in the sub-graph



The computation from $\partial \hat{y}/z^t$ to $\partial \hat{y}/z^{t-1}$ is the # of all edges in the sub-graph

Total computation: # of edges + # of nodes!



The computation from $\partial \hat{y}/z^t$ to $\partial \hat{y}/z^{t-1}$ is the # of all edges in the sub-graph

Total computation: # of edges + # of nodes!

Exercise: can you express backward pass in matrix-vector format?

1. Naively compute all derivatives wrt edges using chain rule takes $(E+V)^2$ time

1. Naively compute all derivatives wrt edges using chain rule takes $(E+V)^2$ time

2. Backpropagation: forward pass & backward pass takes O(E+V) time

1. Naively compute all derivatives wrt edges using chain rule takes $(E+V)^2$ time

2. Backpropagation: forward pass & backward pass takes O(E+V) time

Forward pass: $x = z^0 \rightarrow u^1 \rightarrow z^1 \rightarrow \dots \rightarrow z^t \rightarrow u^{t+1} \rightarrow z^{t+1} \dots \rightarrow z^T \rightarrow \hat{y}$

1. Naively compute all derivatives wrt edges using chain rule takes $(E+V)^2$ time

2. Backpropagation: forward pass & backward pass takes O(E+V) time

Forward pass:
$$x = z^0 \rightarrow u^1 \rightarrow z^1 \rightarrow \dots \rightarrow z^t \rightarrow u^{t+1} \rightarrow z^{t+1} \dots \rightarrow z^T \rightarrow \hat{y}$$

Backward pass:
$$\frac{\partial \hat{y}}{\partial z^T} \rightarrow \frac{\partial \hat{y}}{\partial z^{T-1}} \rightarrow \dots \frac{\partial \hat{y}}{\partial z^1}$$