# Recurrent Neural Nets

## Introduction to NLP

**Rahul Mishra**
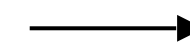IIIT-Hyderabad
March 05, 2024

# Sequence Data
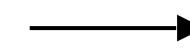
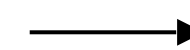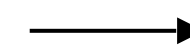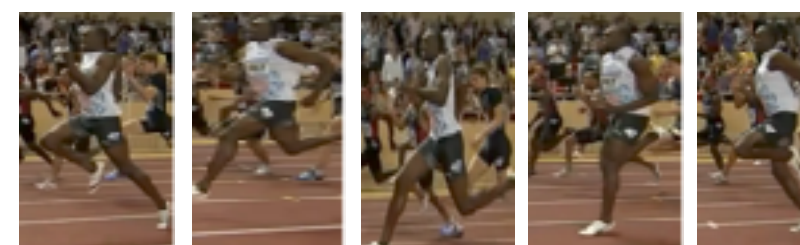| | | |
|---|---|---|
| Speech recognition |  $\varnothing$ | "The quick brown fox jumped over the lazy dog." |
| Music generation | |  |
| Sentiment classification | "There is nothing to like in this movie." | ★☆☆☆☆ |
| DNA sequence analysis | AGCCCCTGTGAGGAACTAG | AG**CCCCTGTGAGGAACT**AG |
| Machine translation | Voulez-vous chanter avec moi? | Do you want to sing with me? |
| Video activity recognition |  | Running |
| Name entity recognition | Yesterday, Harry Potter met Hermione Granger. | Yesterday, **Harry Potter** met **Hermione Granger**. |

# Motivating example

NLP

x:     Harry Potter and Hermione Granger invented a new spell.

$\rightarrow x^{\langle 1 \rangle} \quad x^{\langle 2 \rangle} \quad x^{\langle 3 \rangle} \quad \cdots \quad x^{\langle t \rangle} \quad \cdots \quad x^{\langle 9 \rangle}$

$T_x = 9$

$\rightarrow y:$     1     1     0     1     1     0   0   0   0

$y^{\langle 1 \rangle} \quad y^{\langle 2 \rangle} \quad y^{\langle 3 \rangle} \quad \cdots \quad y^{\langle 9 \rangle}$

$T_y = 9$

$x^{(i)\langle t \rangle}$       $T_x^{(i)} = 9$     15

$y^{(i)\langle t \rangle}$       $T_y^{(i)}$

# Representing Words

$x^{<t>}$       $(x, y)$

$x \rightarrow y$

x:      Harry Potter and Hermione Granger invented a new spell.

...

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$            $x^{<7>}$     $x^{<9>}$

Vocabulary

$$\begin{bmatrix} a \\ aaron \\ \vdots \\ and \\ \vdots \\ harry \\ \vdots \\ potter \\ \vdots \\ zulu \end{bmatrix} \begin{matrix} 1 \leftarrow \\ 2 \\ \vdots \\ 367 \leftarrow \\ \vdots \\ 4075 \\ \vdots \\ 6830 \\ \vdots \\ 10,000 \end{matrix}$$

<UNK>   10,000

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4075$$

One-hot

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 6830$$

$$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 367 \quad 10,000$$

$$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

# Why not a standard network?

$$x^{<1>}$$

$$x^{<2>}$$

$$\vdots$$

$$x^{<T_x>}$$

$$y^{<1>}$$

$$y^{<2>}$$

$$\vdots$$

$$y^{<T_y>}$$

Problems:

- Inputs, outputs can be different lengths in different examples.

- Doesn't share features learned across different positions of text.

# RNNs

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$

books

laptops

a          zoo

$U$

output distribution

$$\hat{y}^{(t)} = \mathrm{softmax}\left(Uh^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$$

$h^{(0)}$     $h^{(1)}$     $h^{(2)}$     $h^{(3)}$     $h^{(4)}$

**RNN Advantages:**
- Can process any length input
- Computation for step *t* can (in theory) use information from many steps back
- Model size doesn't increase for longer input context
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

hidden states

$$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$$

$h^{(0)}$ is the initial hidden state

$W_h$     $W_h$     $W_h$     $W_h$

$W_e$     $W_e$     $W_e$     $W_e$

**RNN Disadvantages:**
- Recurrent computation is slow
- In practice, difficult to access information from many steps back

word embeddings

$$e^{(t)} = Ex^{(t)}$$

$e^{(1)}$     $e^{(2)}$     $e^{(3)}$     $e^{(4)}$

$E$     $E$     $E$     $E$

words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

the          students          opened          their

$x^{(1)}$     $x^{(2)}$     $x^{(3)}$     $x^{(4)}$

# RNN

‣ Idea: We can
process a sequence
of vectors **x** by
applying a
recurrence formula
at every time step t

‣ Note: Same
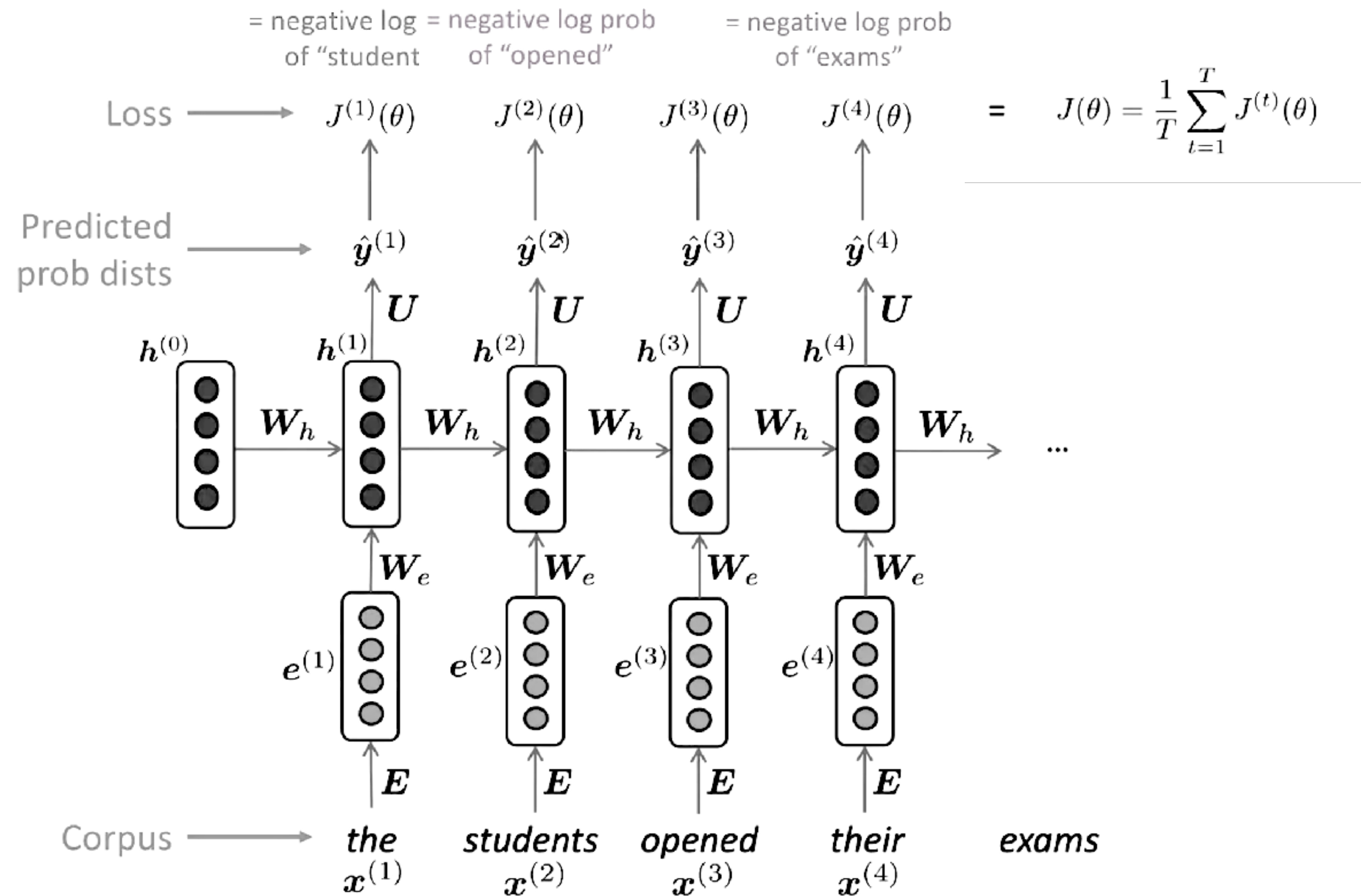parameters, weights
and function used at
each step.

# Training RNN

Loss function on step $t$ is cross-entropy between predicted probability distribution $\hat{y}^{(t)}$, and the true next word $y^{(t)}$ (one-hot for $x^{(t+1)}$):

$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_w^{(t)} \log \hat{\boldsymbol{y}}_w^{(t)} = -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$
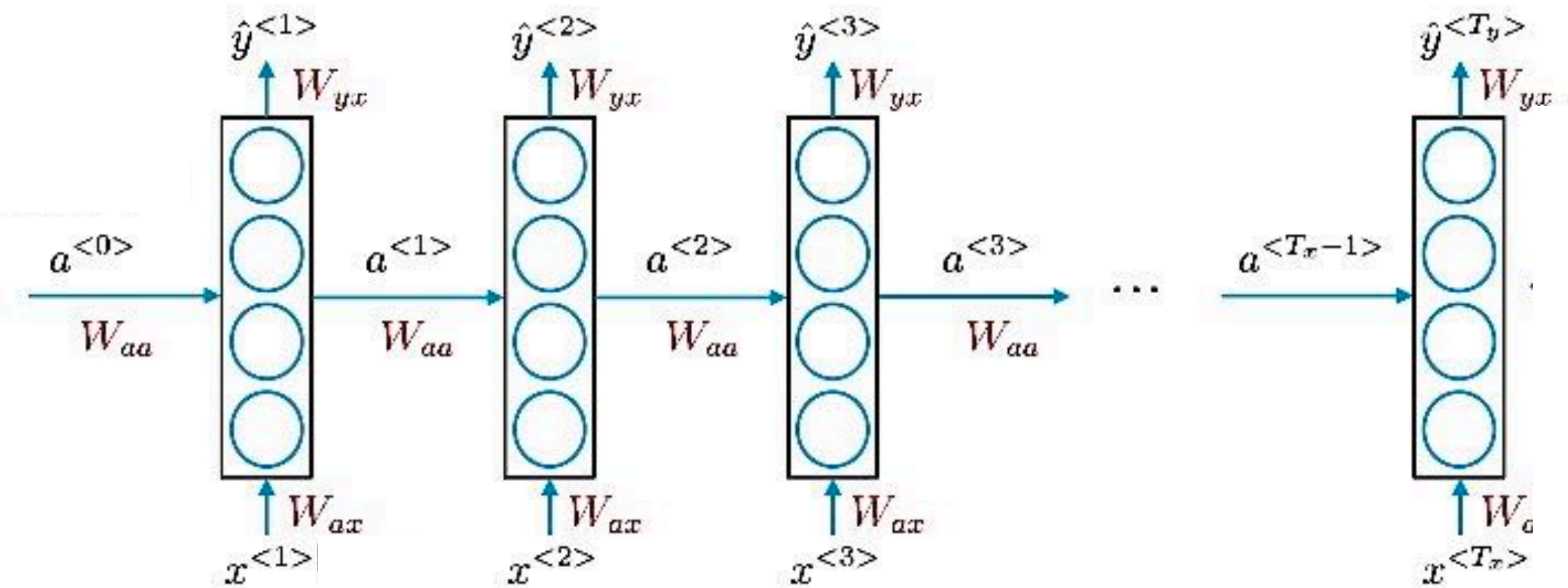
Average this to get overall loss

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

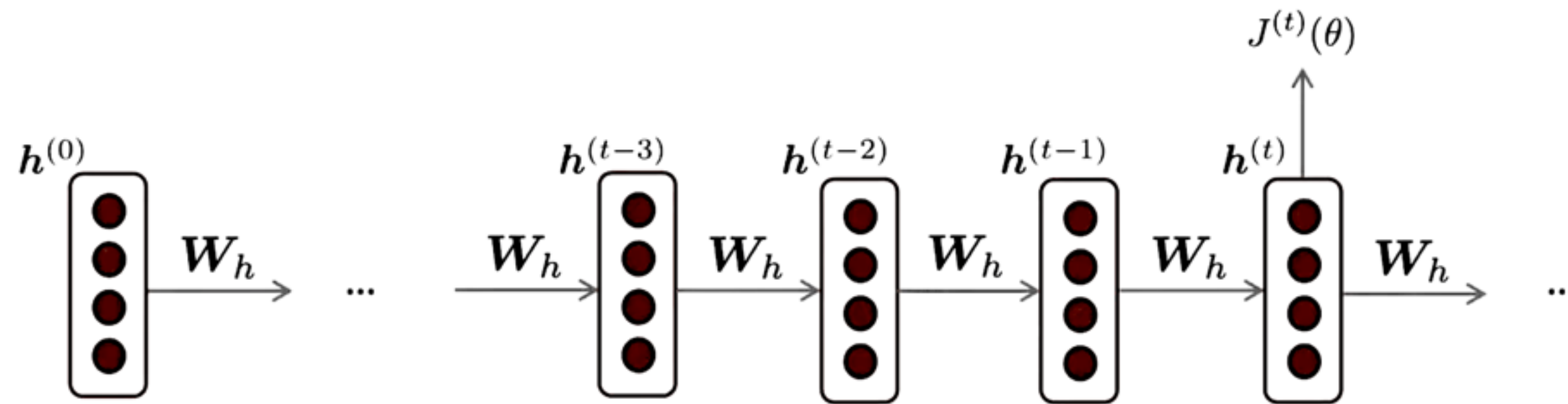# Training RNN

# Forward Propagation



$$a^{<0>} = 0$$

$$a^{<1>} = g_1\left(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a\right) \quad g_1() \leftarrow \tanh/\text{ReLU}$$

$$\hat{y}^{<1>} = g_2\left(W_{ya}a^{<1>} + b_y\right) \quad\quad g_2() \leftarrow \text{Sigmoid/Softmax/Linear}$$

$$a^{<t>} = g_1\left(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a\right)$$

$$\hat{y}^{<t>} = g_2\left(W_{ya}a^{<t>} + b_y\right)$$
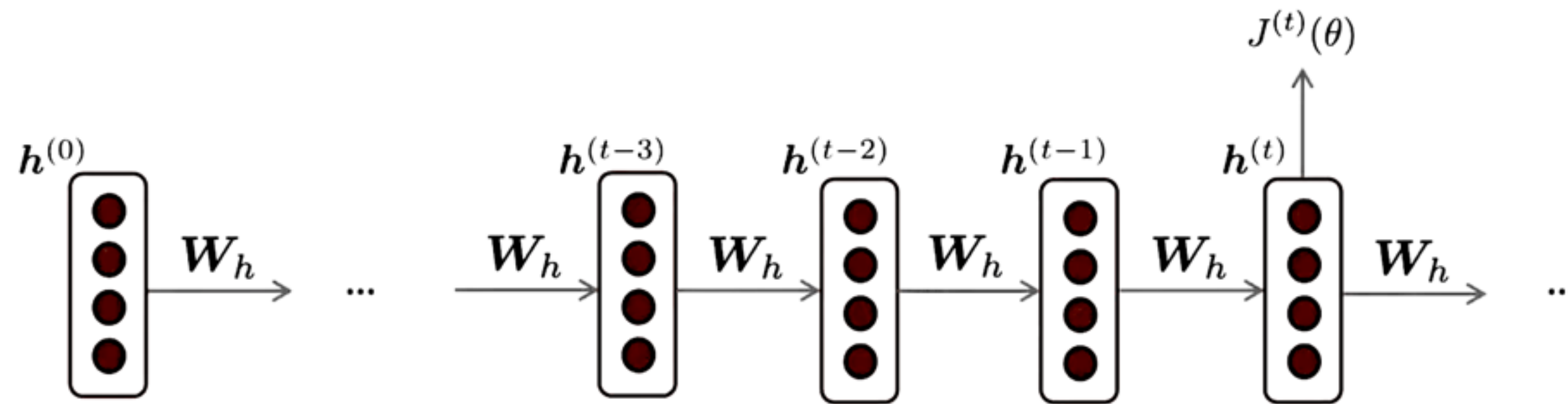
# Training RNN: Backdrop in RNN



**Question:** What's the derivative of $J^{(t)}(\theta)$ w.r.t. the repeated weight matrix $W_h$ ?
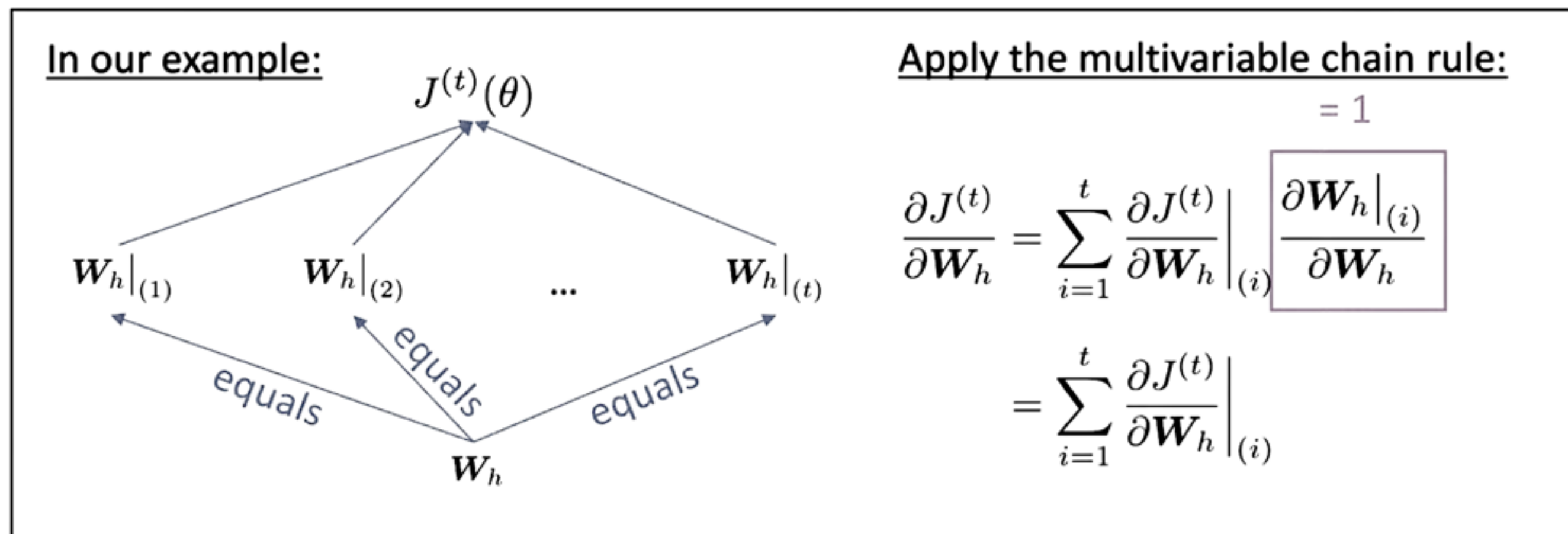
**Answer:** 
$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial W_h}\bigg|_{(i)}$$

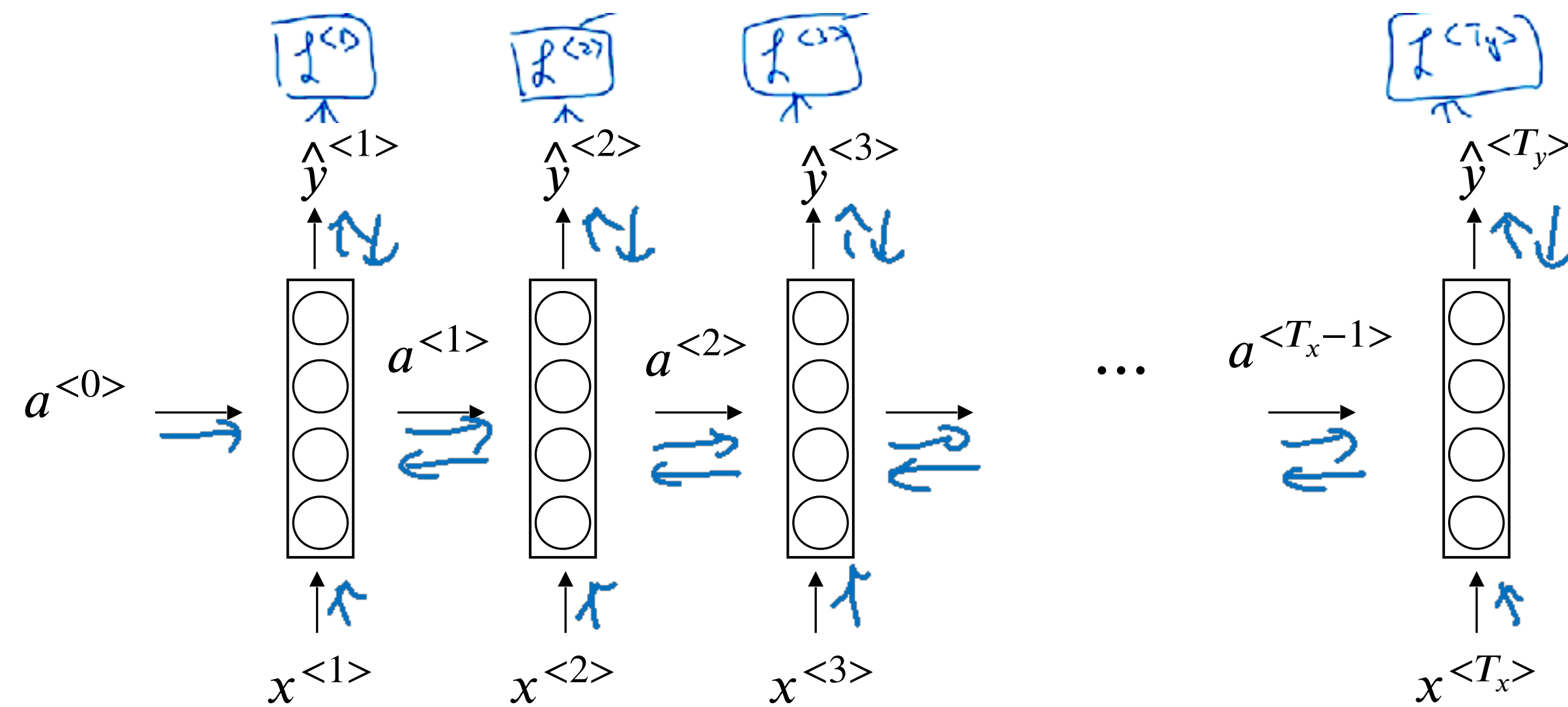"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

# Training RNN: Backdrop in RNN



**Question:** What's the derivative of $J^{(t)}(\theta)$ w.r.t. the repeated weight matrix $W_h$ ?

In our example:

Apply the multivariable chain rule:

$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^{t} \left.\frac{\partial J^{(t)}}{\partial W_h}\right|_{(i)} \boxed{\frac{\partial W_h|_{(i)}}{\partial W_h}}$$

$$= 1$$

$$= \sum_{i=1}^{t} \left.\frac{\partial J^{(t)}}{\partial W_h}\right|_{(i)}$$

# Backpropagation Through Time

# Backpropagation Through Time



$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}} = \boxed{\sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}}\bigg|_{(i)}}$$

**Question:** How do we calculate this?

**Answer:** Backpropagate over timesteps $i=t,...,0$, summing gradients as you go. This algorithm is called **"backpropagation through time"** [Werbos, P.G., 1988, *Neural Networks* **1**, and others]
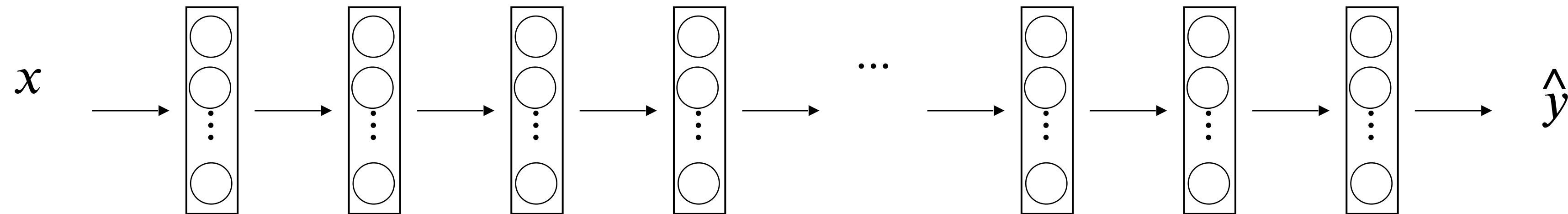
# Generating Text using RNNs



"Sorry," Harry shouted, panicking—"I'll leave those brooms in London, are they?"

"No idea," said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry's shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn't felt it seemed. He reached the teams too.

**Source:** https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803da6
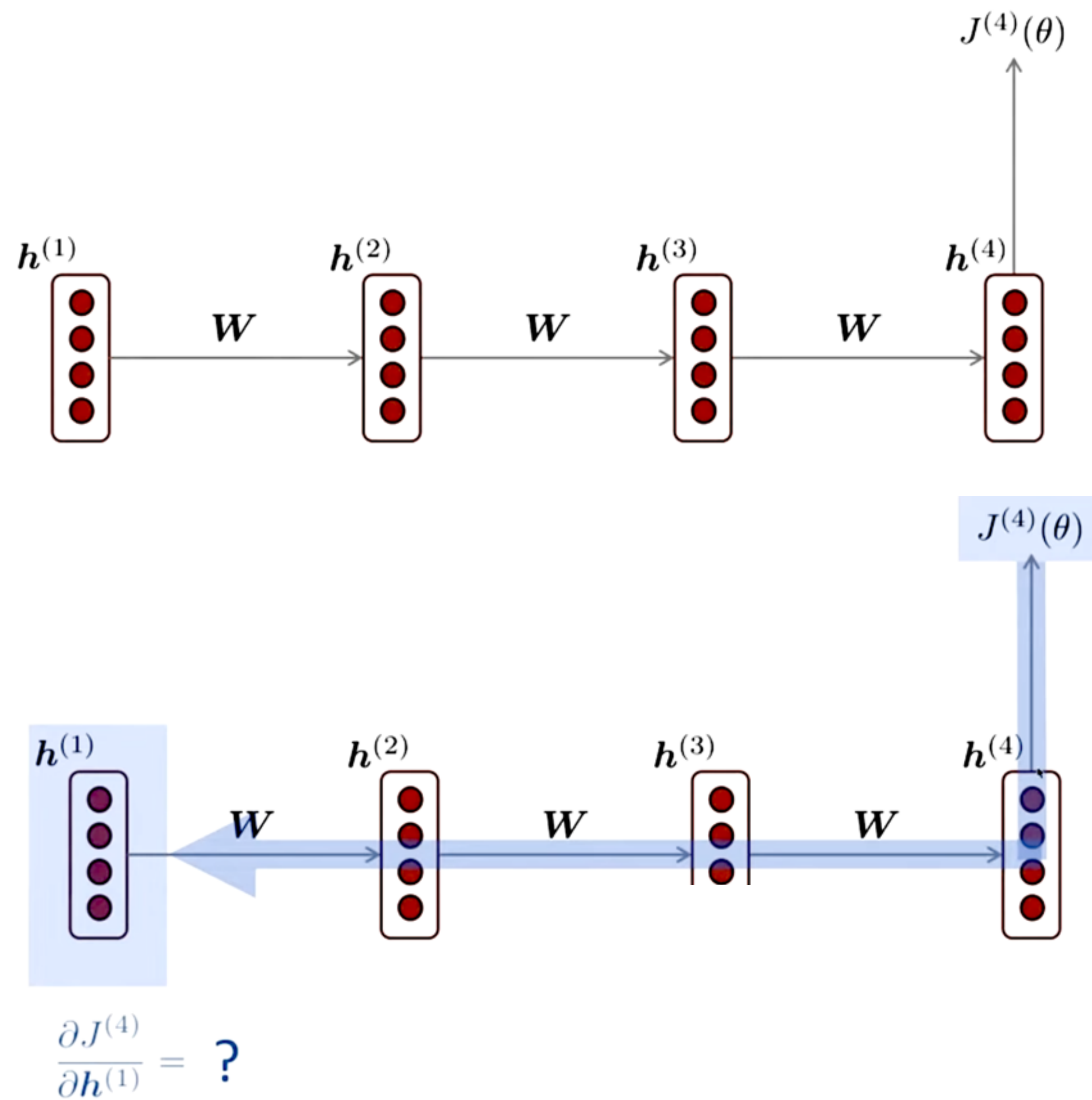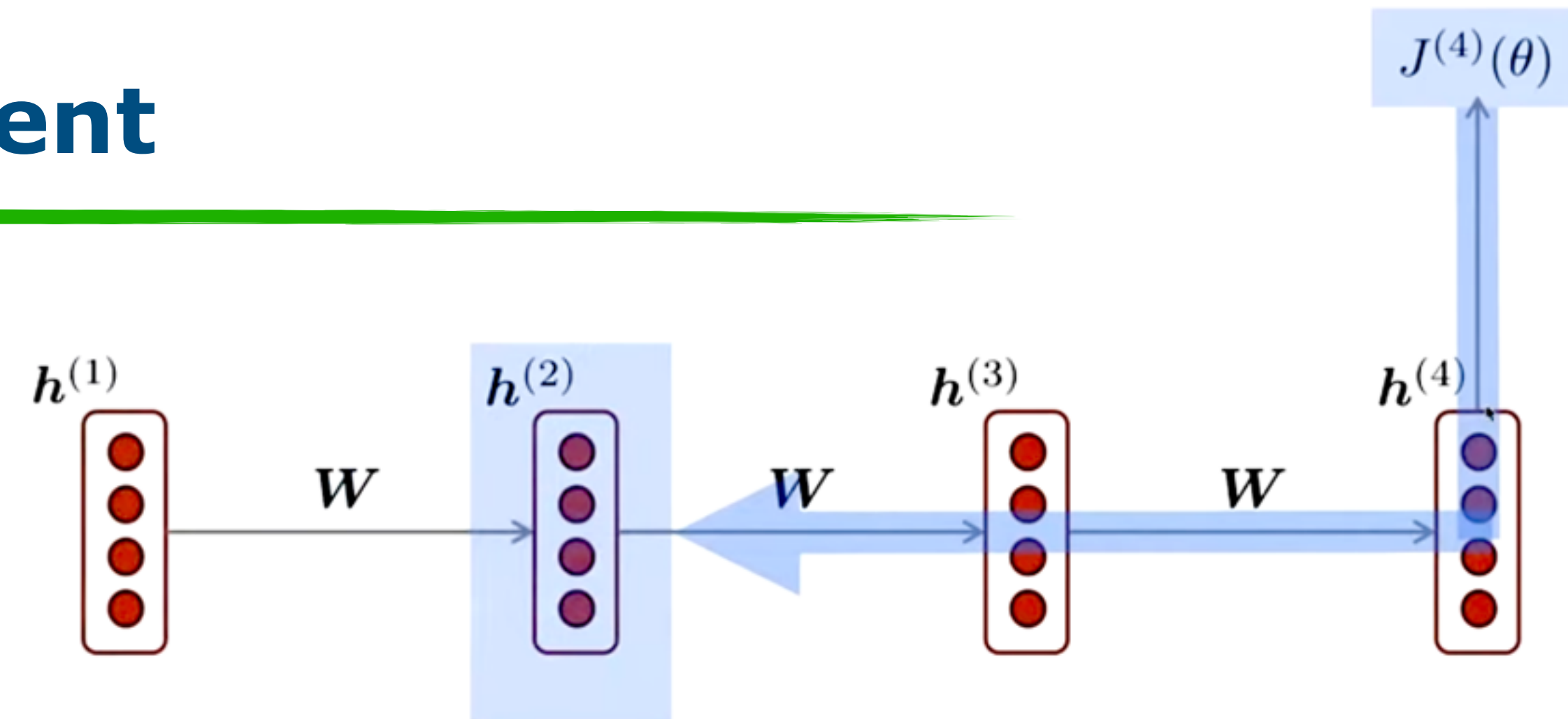
# Vanishing and Exploding Gradient



**The cat, that already ate..........., was full**

**The cats, that already ate.........., were full**

# Vanishing Gradient



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \, ?$$

# Vanishing Gradient



$$J^{(4)}(\theta)$$

$$h^{(1)} \quad h^{(2)} \quad h^{(3)} \quad h^{(4)}$$

$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial J^{(4)}}{\partial h^{(2)}}$$

chain rule!

$$h^{(1)} \quad h^{(2)} \quad h^{(3)} \quad h^{(4)}$$

$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \qquad \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \qquad \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$
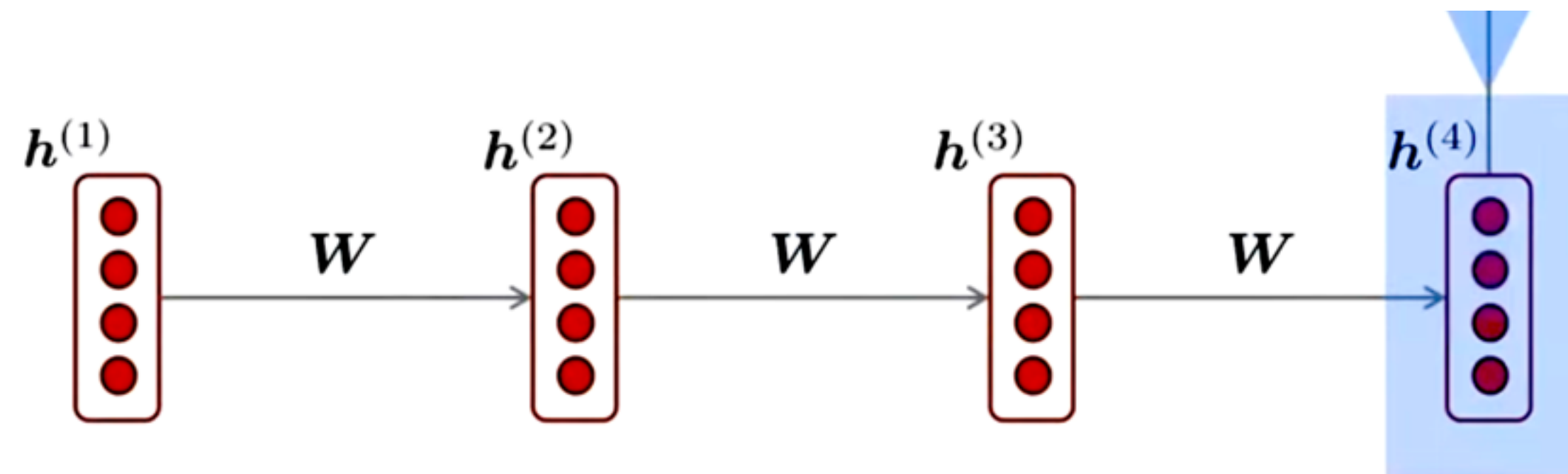
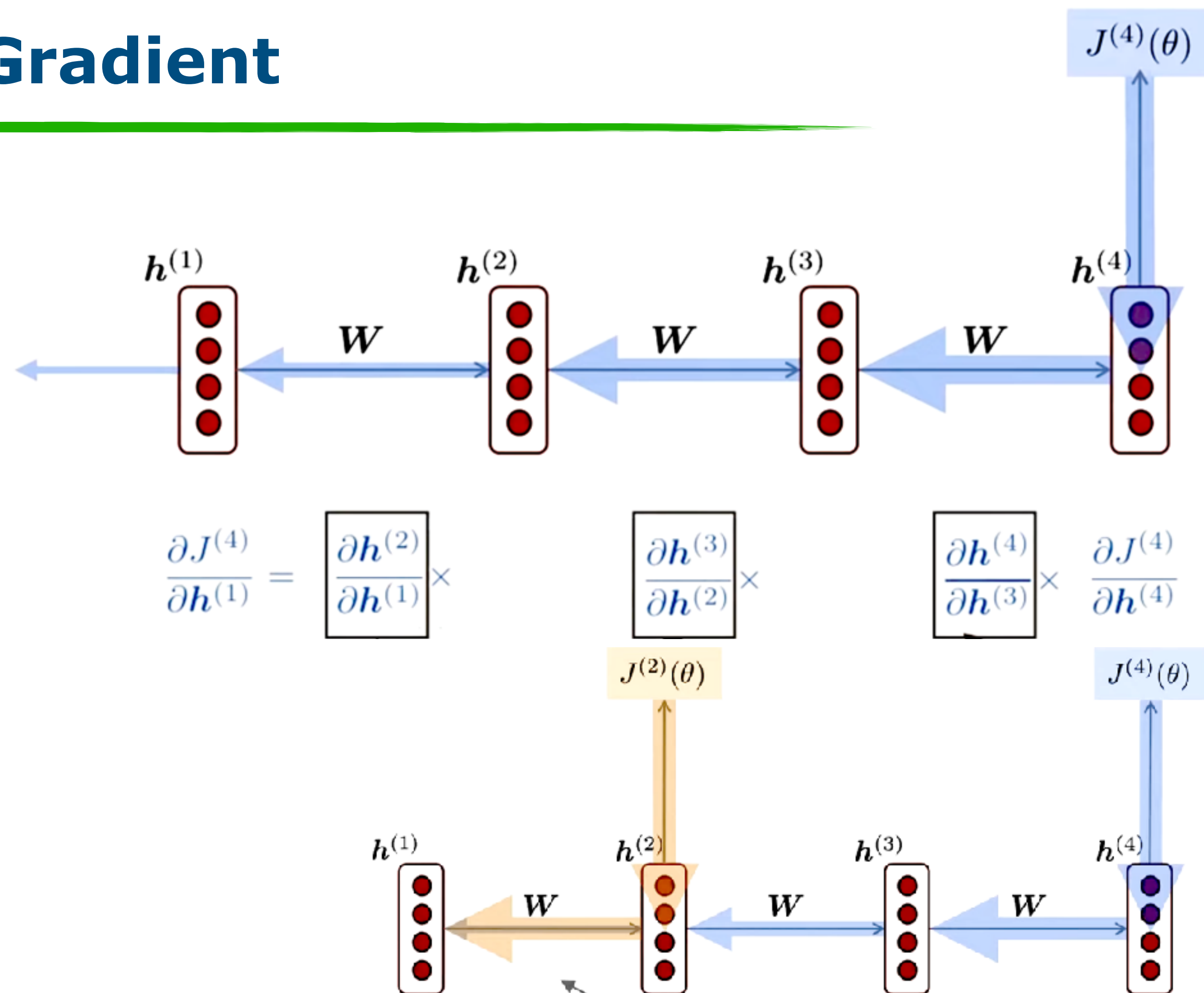chain rule!

# Vanishing Gradient



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \boxed{\frac{\partial h^{(2)}}{\partial h^{(1)}}} \times \boxed{\frac{\partial h^{(3)}}{\partial h^{(2)}}} \times \boxed{\frac{\partial h^{(4)}}{\partial h^{(3)}}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

Gradient signal from far away is lost because it's much smaller than gradient signal from close-by.

So, model weights are updated only with respect to near effects, not long-term effects.

# Exploding Gradient

If the gradient becomes too big, then the SGD update step becomes too big:

$$\theta^{new} = \theta^{old} - \overbrace{\alpha}^{\text{learning rate}} \underbrace{\nabla_\theta J(\theta)}_{\text{gradient}}$$

**Gradient clipping**: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update
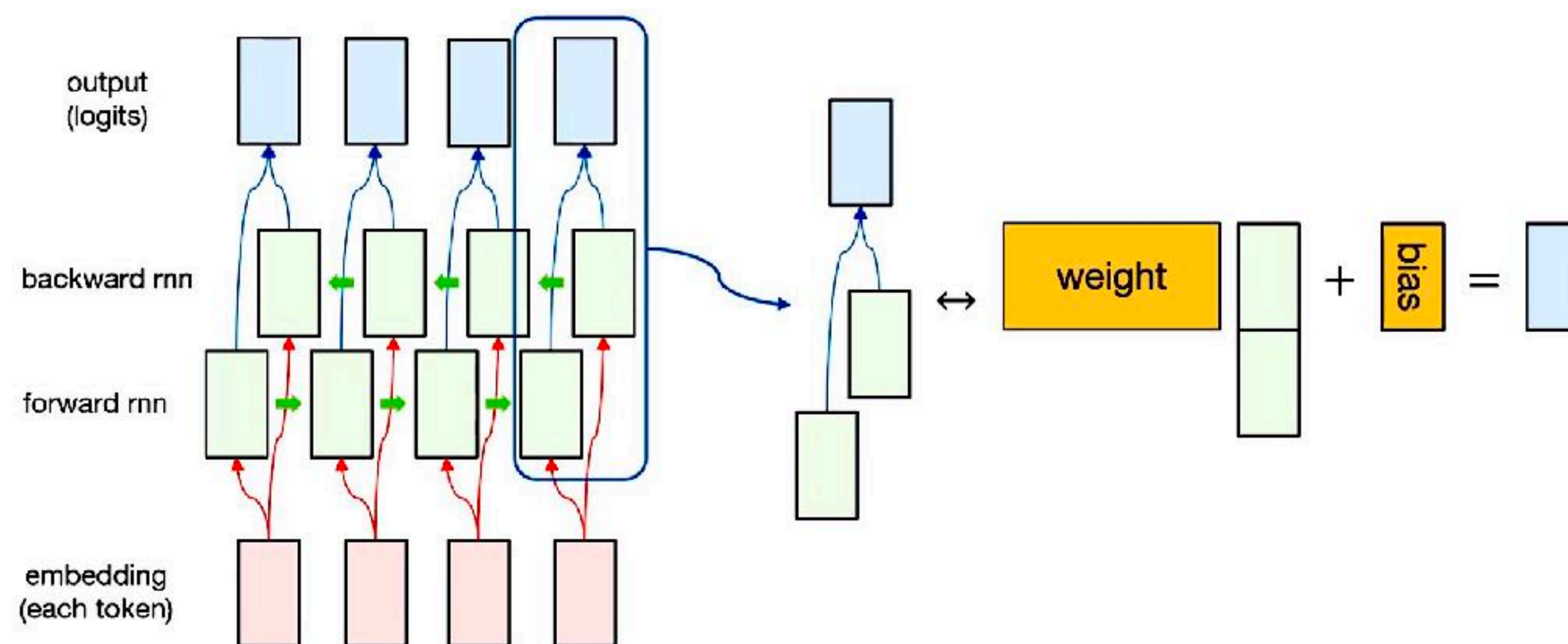
---
**Algorithm 1** Pseudo-code for norm clipping
---
$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$
**if** $\|\hat{\mathbf{g}}\| \geq threshold$ **then**
$\quad \hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$
**end if**

---

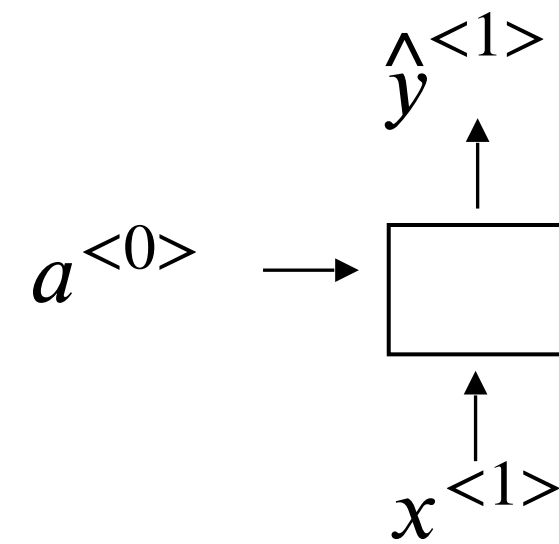**Intuition**: take a step in the same direction, but a smaller step

# Another Problem
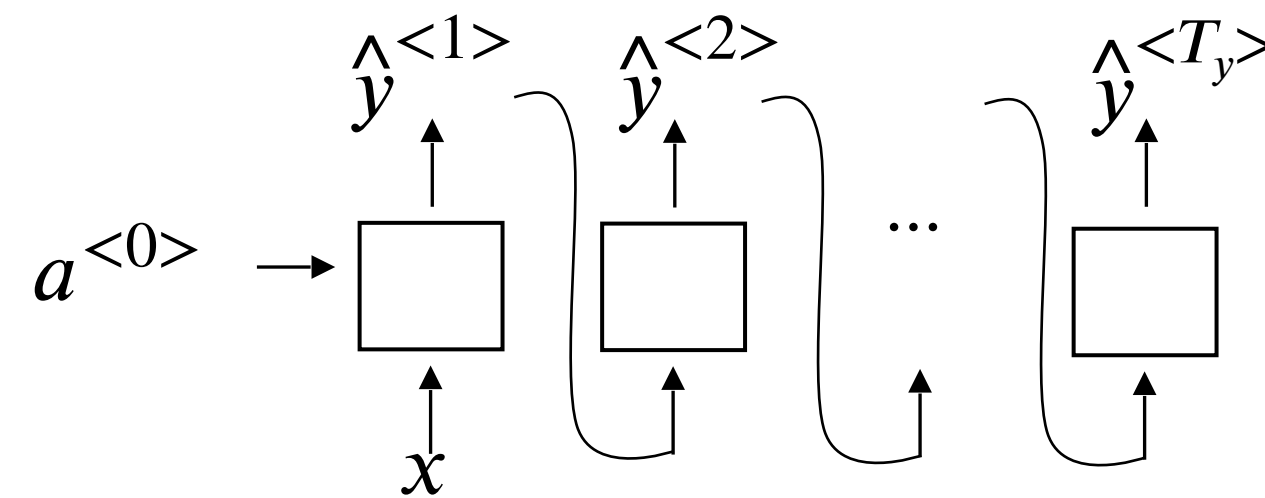
He said, "Teddy Roosevelt was a great President."
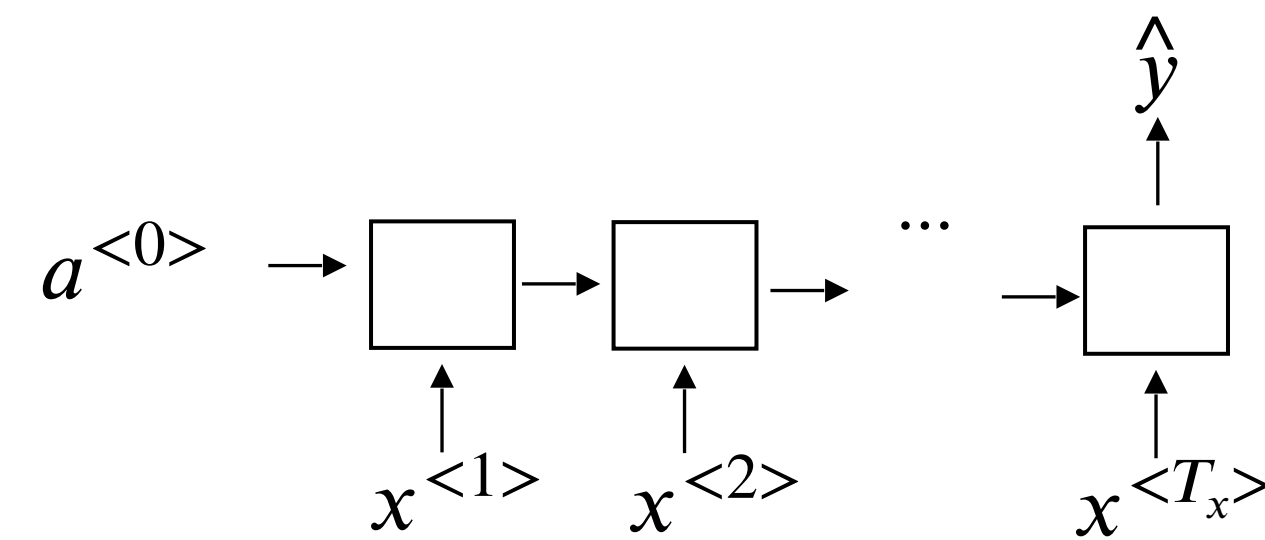
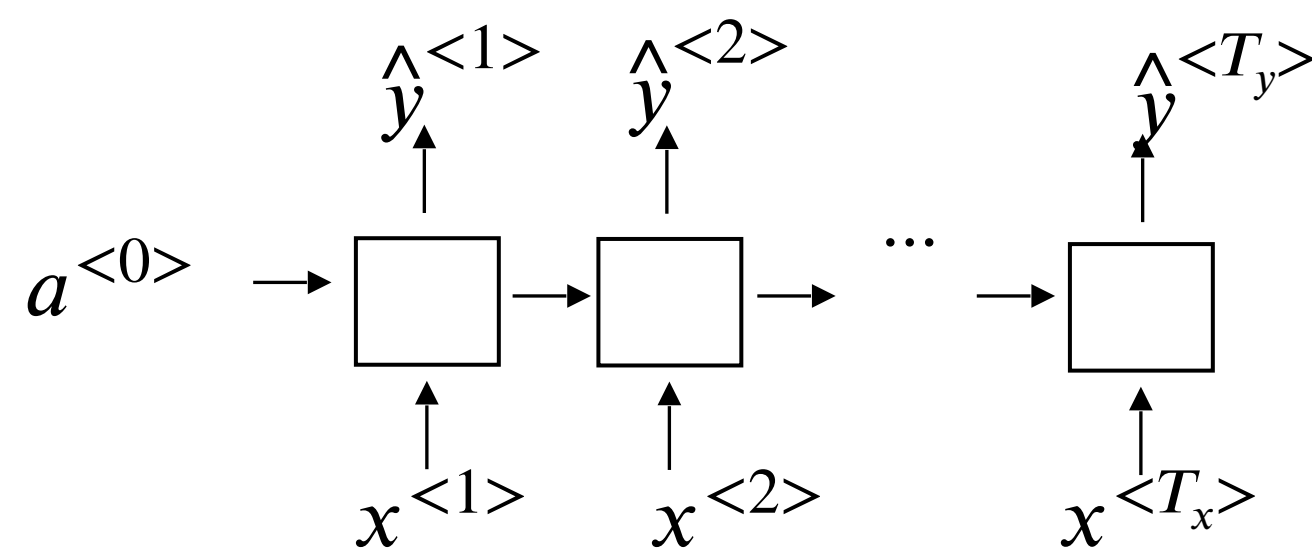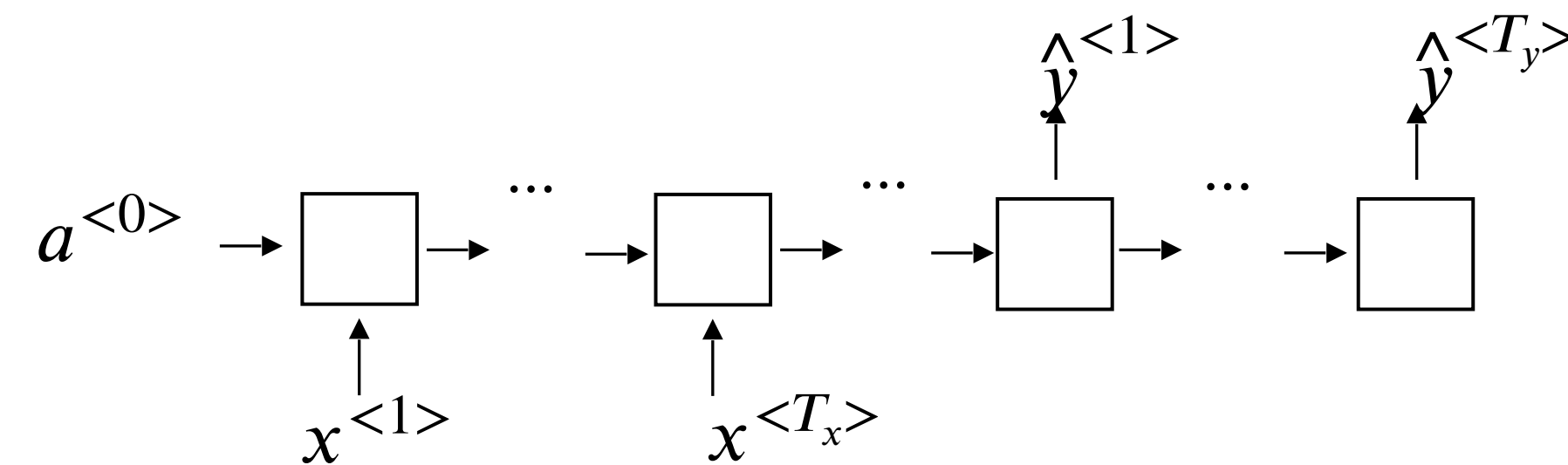He said, "Teddy bears are on sale!"

# RNN architecture types

$\hat{y}^{<1>}$

$a^{<0>} \rightarrow \boxed{\phantom{aa}}$

$x^{<1>}$

One to one

$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \cdots \quad \hat{y}^{<T_y>}$

$a^{<0>} \rightarrow \boxed{\phantom{a}} \quad \boxed{\phantom{a}} \quad \cdots \quad \boxed{\phantom{a}}$

$x$

One to many

$\hat{y}$

$a^{<0>} \rightarrow \boxed{\phantom{a}} \rightarrow \boxed{\phantom{a}} \rightarrow \cdots \rightarrow \boxed{\phantom{a}}$

$x^{<1>} \quad x^{<2>} \quad x^{<T_x>}$

Many to one

$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<T_y>}$

$a^{<0>} \rightarrow \boxed{\phantom{a}} \rightarrow \boxed{\phantom{a}} \rightarrow \cdots \rightarrow \boxed{\phantom{a}}$

$x^{<1>} \quad x^{<2>} \quad x^{<T_x>}$

Many to many

$\hat{y}^{<1>} \quad \hat{y}^{<T_y>}$

$a^{<0>} \rightarrow \boxed{\phantom{a}} \rightarrow \cdots \rightarrow \boxed{\phantom{a}} \rightarrow \cdots \rightarrow \boxed{\phantom{a}} \rightarrow \cdots \rightarrow \boxed{\phantom{a}}$

$x^{<1>} \quad x^{<T_x>}$

Many to many

# GRU: Gated Recurrent Unit

GRU unit

$= \Gamma_u * \check{c}^{<t>} + (1-\Gamma_u) * c^{<t-1>}$

Softmax

$\hat{y}^{<t>}$

$c^{<t-1>} = a^{<t-1>}$

$c^{<t>} = a^{<t>}$

$\check{c}^{<t>}$

$\Gamma_u$

tanh

$\sigma$

$x^{<t>}$

$X^{<t>}$

$C = \text{memory cell}$

$$\rightarrow C^{<t>} = a^{<t>}$$

$$\rightarrow \tilde{c}^{<t>} = \tanh \left( W_c \left[ c^{<t-1>}, x^{<t>} \right] + b_c \right)$$

$$\rightarrow \Gamma_u = \sigma \left( W_u \left[ c^{<t-1>}, x^{<t>} \right] + b_u \right)$$

"update"

$$C^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1-\Gamma_u) * c^{<t-1>}$$

$=1$

The cat, which already ate …, was full.

[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches]

[Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling]

# Full GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) + c^{<t-1>}$$

The cat, which ate already, was full.