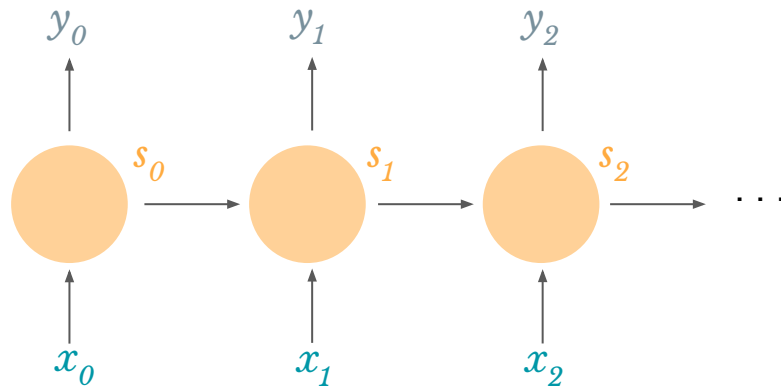


Sequence Modeling with Neural Networks

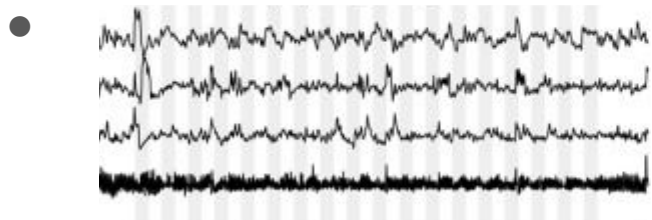
Harini Suresh



What is a sequence?

- “This morning I took the dog for a walk.”

sentence



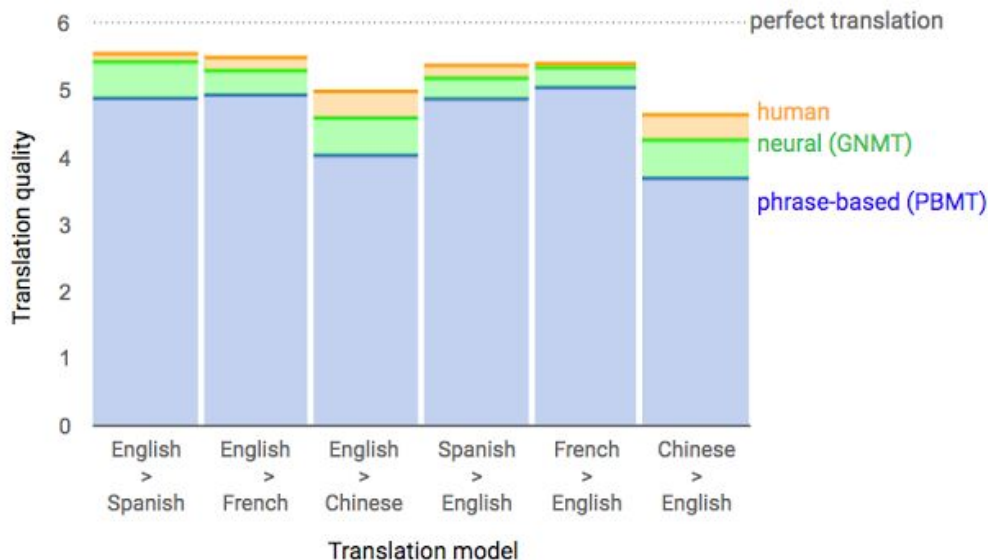
medical signals



speech waveform

Successes of deep models

Machine translation



Question Answering

Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion Denver Broncos defeated the National Football Conference (NFC) champion Carolina Panthers 24–10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California. As this was the 50th Super Bowl, the league emphasized the "golden anniversary" with various gold-themed initiatives, as well as temporarily suspending the tradition of naming each Super Bowl game with Roman numerals (under which the game would have been known as "Super Bowl L"), so that the logo could prominently feature the Arabic numerals 50.

Super Bowl 50 decided the NFL champion for what season?

Ground Truth Answers: 2015 the 2015 season 2015

Prediction: 2015

Left:

<https://research.googleblog.com/2016/09/a-neural-network-for-machine.html>

Right:

<https://rajpurkar.github.io/SQuAD-explorer/>

Successes of deep models

Audi Self-Parking
Demonstration

a sequence modeling problem:
predict the next word

a sequence modeling problem

“This morning I took the dog for a walk.”

a sequence modeling problem

“This morning I took the dog for a walk.”

given these words

a sequence modeling problem

“This morning I took the dog for a walk.”

given these words

*predict what
comes next?*

idea: use a fixed window

“This morning I took the dog **for a** **walk.**”

*given these 2
words, predict
the next word*

idea: use a fixed window

“This morning I took the dog **for a** walk.”

given these 2 words, predict the next word

[1 0 0 0 0 0 1 0 0 0]

for

a

One hot feature vector indicates what each word is

prediction

problem: we can't model long-term dependencies

“In **France**, I had a great time and I learnt some of the _____
language.”

We need information from the far past and future to accurately guess the correct word.

idea: use entire sequence, as a set of counts

This morning I took the dog for a

[0 1 0 0 1 0 0 ... 0 0 1 1 0 0 0 1]

“bag of words”

prediction

problem: counts don't preserve order

problem: counts don't preserve order

“The food was good, not bad at all.”

VS

“The food was bad, not good at all.”

idea: use a really big fixed window

“This morning I took the dog for a walk.”
given these 7 words, predict the next word

[1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 ...]

morning I took the dog ...



prediction

problem: no parameter sharing

 this morning
[1 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 ...]

each of these inputs has a *separate parameter*

problem: no parameter sharing

 this morning
[1 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 ...]

each of these inputs has a *separate parameter*

 this morning
[0 0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 ...]

to model sequences, we need:

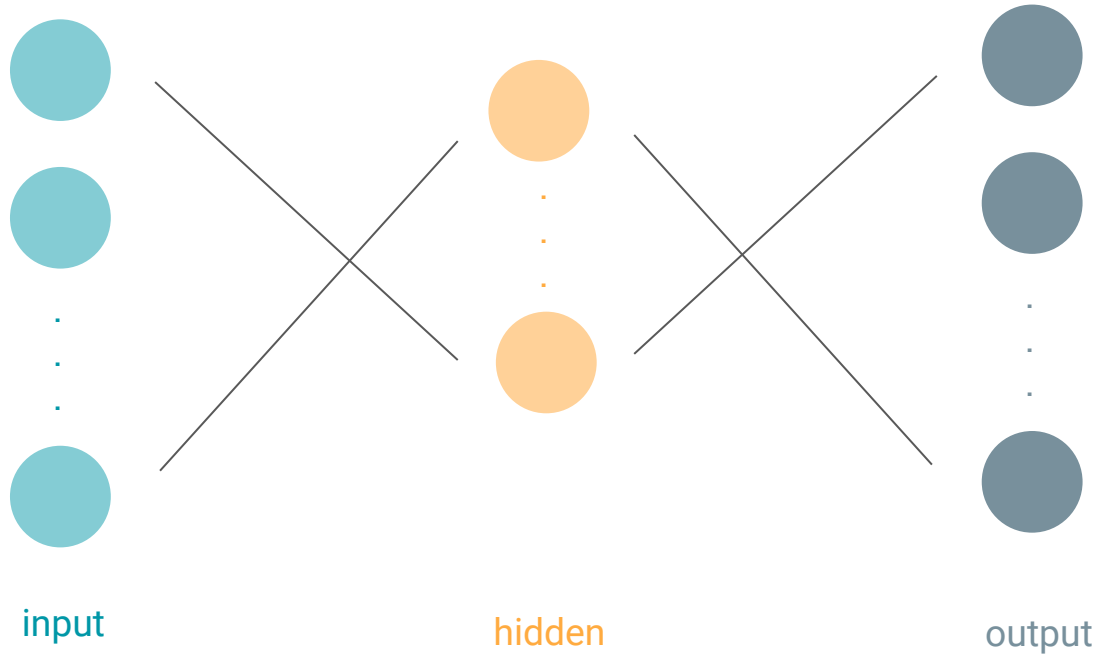
1. to deal with **variable-length** sequences
2. to maintain **sequence order**
3. to keep track of **long-term dependencies**
4. to **share parameters** across the sequence

to model sequences, we need:

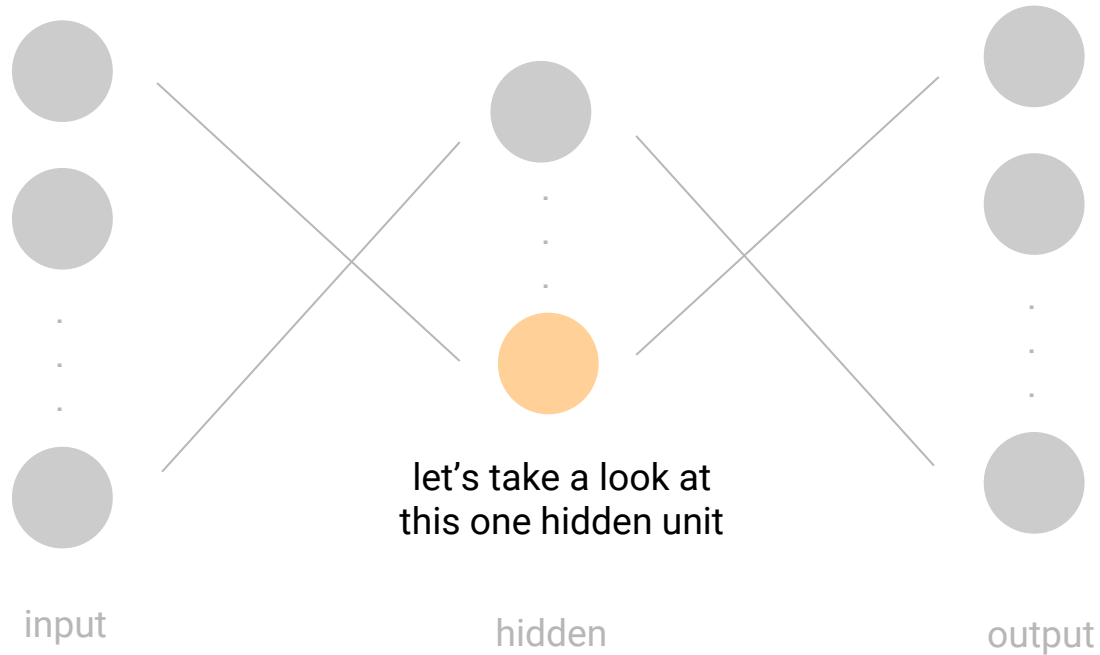
1. to deal with **variable-length** sequences
2. to maintain **sequence order**
3. to keep track of **long-term dependencies**
4. to **share parameters** across the sequence

let's turn to **recurrent neural networks**.

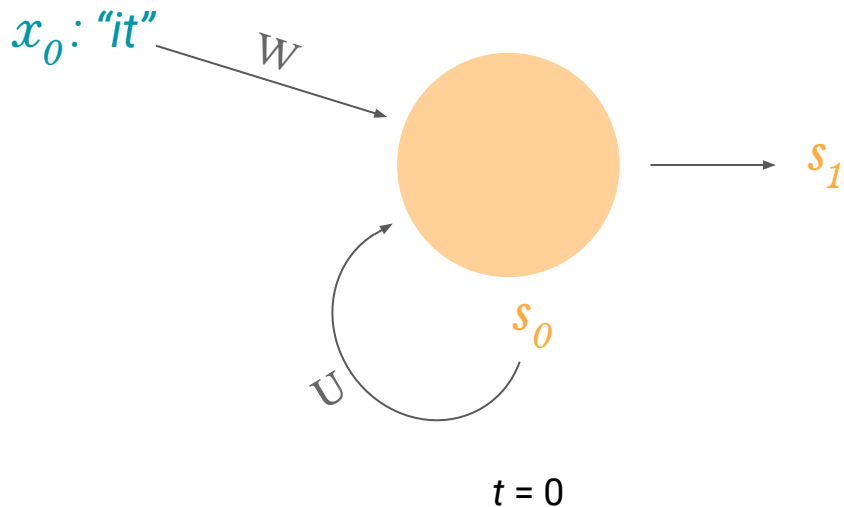
example network:



example network:



RNNS remember their previous state:

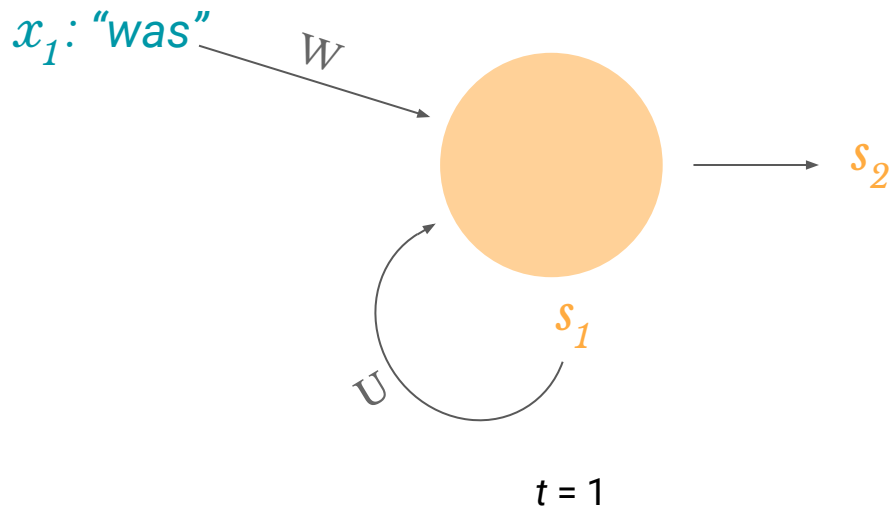


x_0 : vector representing first word
 s_0 : cell state at $t = 0$ (some initialization)
 s_1 : cell state at $t = 1$

$$s_1 = \tanh(Wx_0 + Us_0)$$

W, U : weight matrices

RNNS remember their previous state:



x_1 : vector representing second word

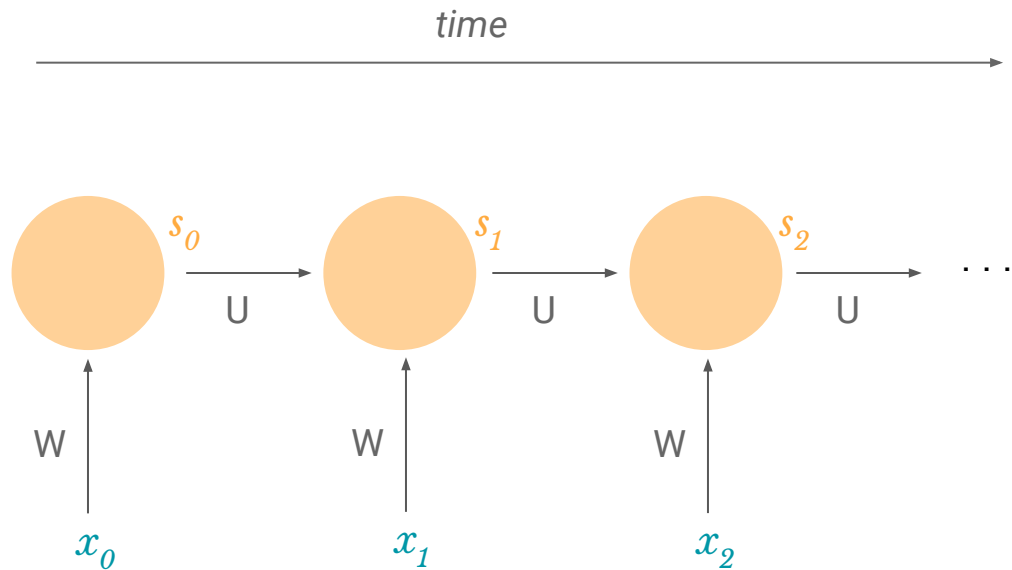
s_1 : cell state at $t = 1$

s_2 : cell state at $t = 2$

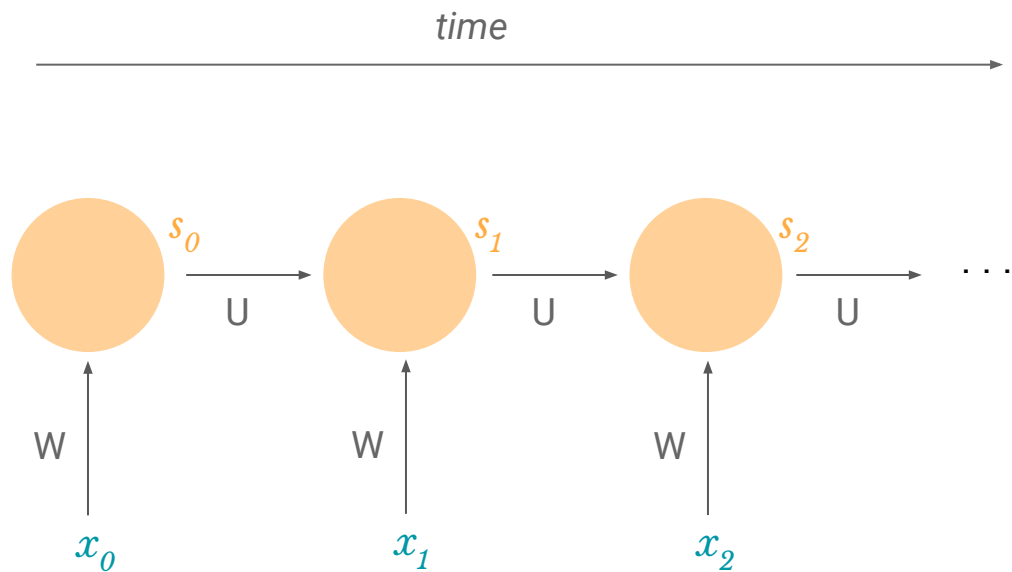
$$s_2 = \tanh(Wx_1 + Us_1)$$

W, U : weight matrices

“unfolding” the RNN across time:

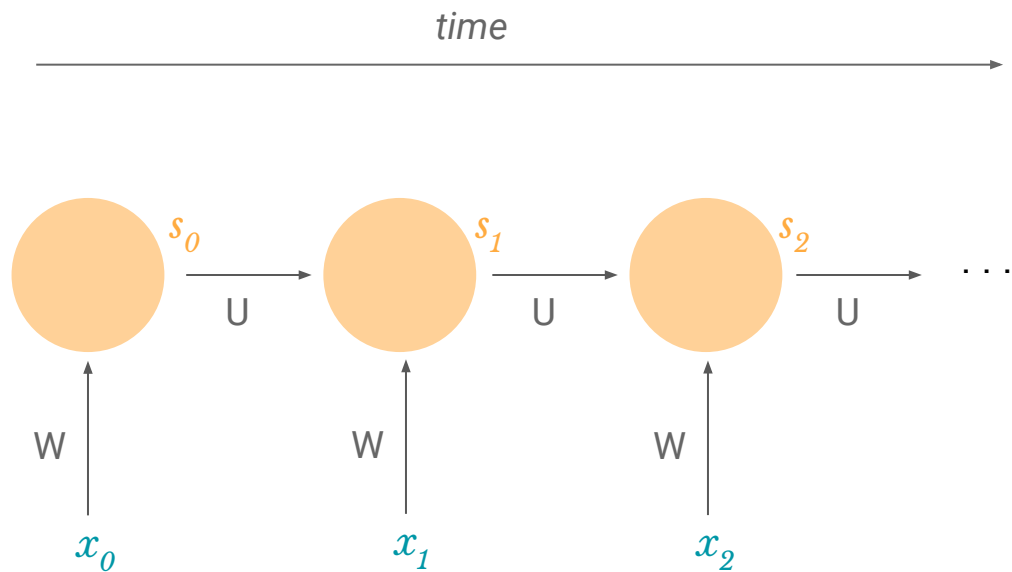


“unfolding” the RNN across time:



notice that we use the same parameters, W and U

“unfolding” the RNN across time:



s_n can contain information from all past timesteps

how do we **train** an RNN?

how do we **train** an RNN?

backpropagation!

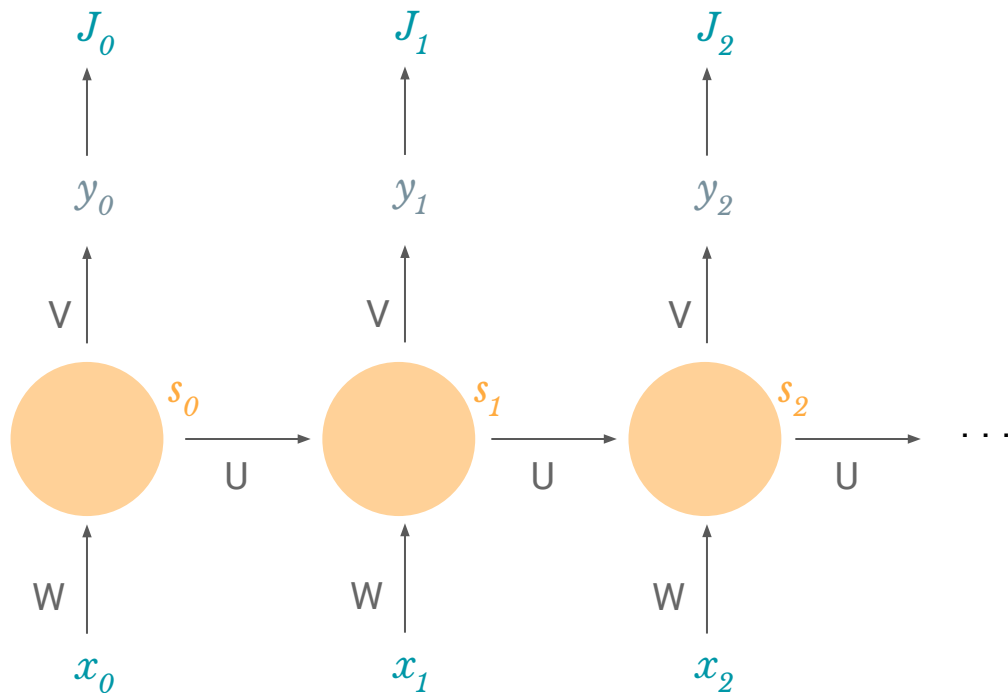
(through time)

remember: **backpropagation**

1. **take the derivative** (gradient) of the loss with respect to each parameter
2. **shift parameters in the opposite direction** in order to minimize loss

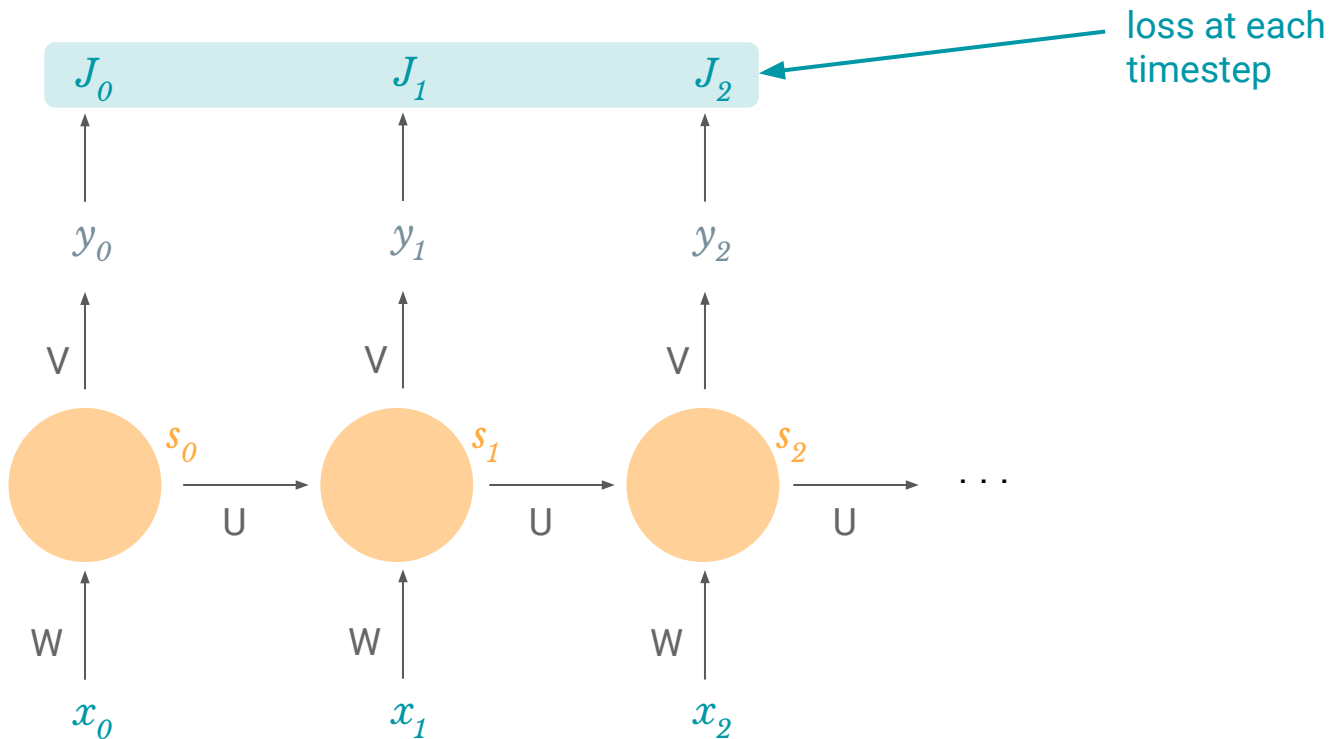
we have a **loss at each timestep**:

(since we're making a prediction at each timestep)



we have a **loss at each timestep**:

(since we're making a prediction at each timestep)



we **sum the losses** across time:

$$\text{loss at time } t = J_t(\Theta)$$

Θ = our
parameters, like
weights



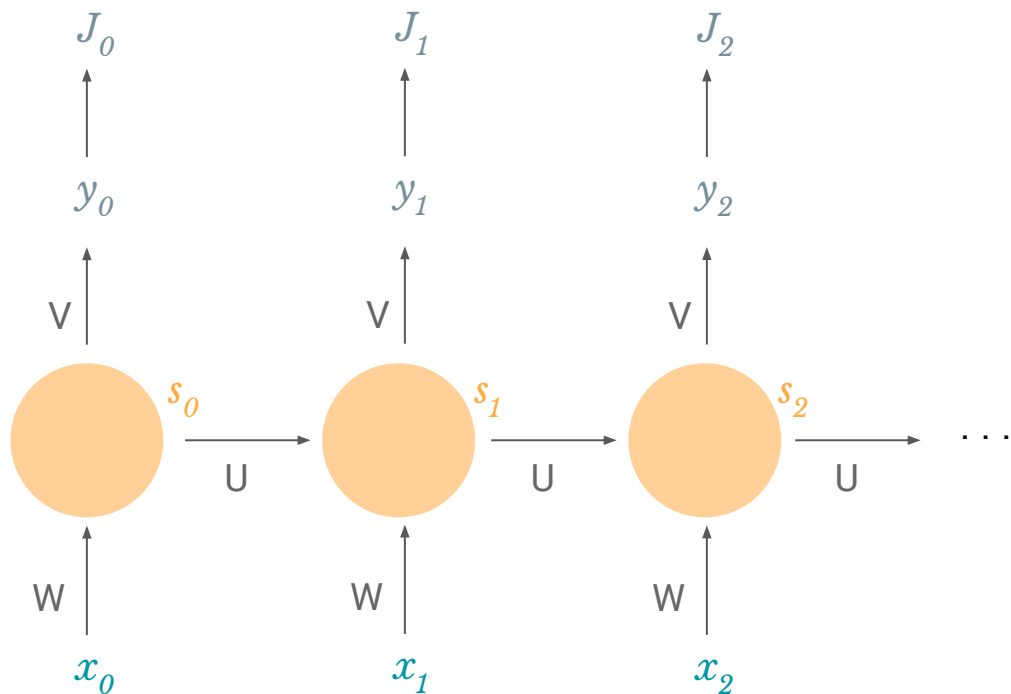
$$\text{total loss} = J(\Theta) = \sum_t J_t(\Theta)$$

what are our **gradients**?

we sum gradients across time for each parameter P :

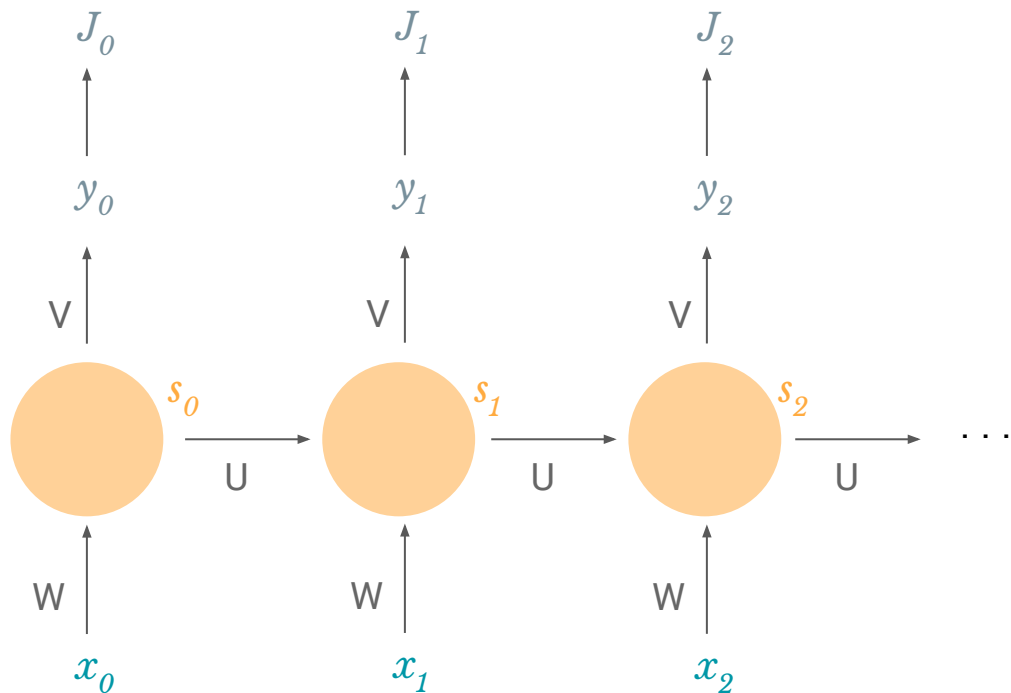
$$\frac{\partial J}{\partial P} = \sum_t \frac{\partial J_t}{\partial P}$$

let's try it out for W with the **chain rule**:



$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

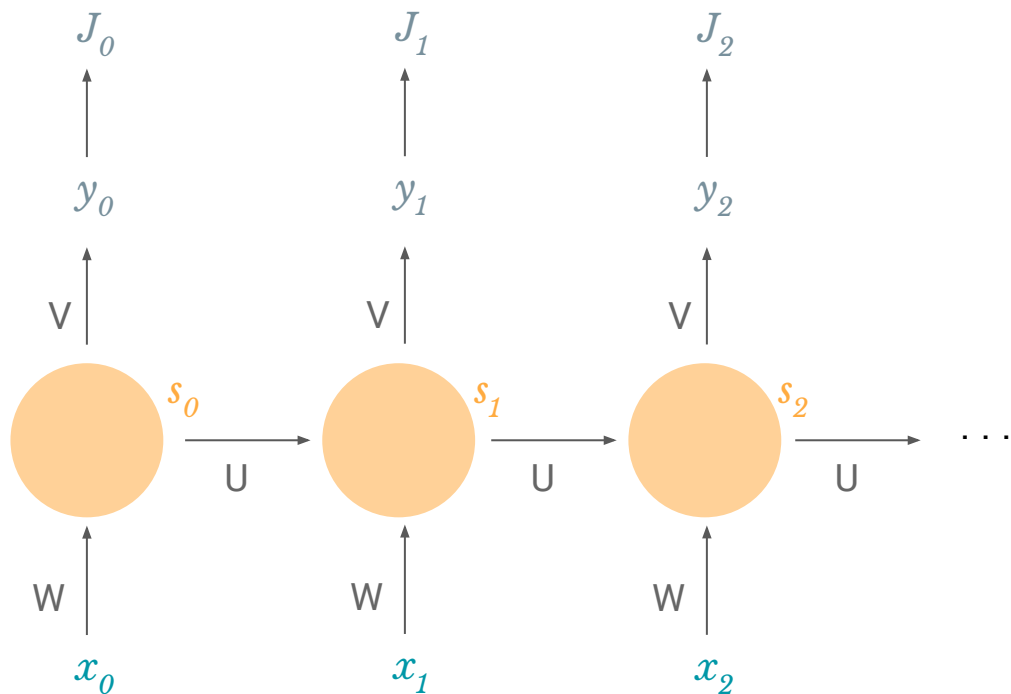
let's try it out for W with the **chain rule**:



$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

so let's take a single timestep t :

let's try it out for W with the **chain rule**:

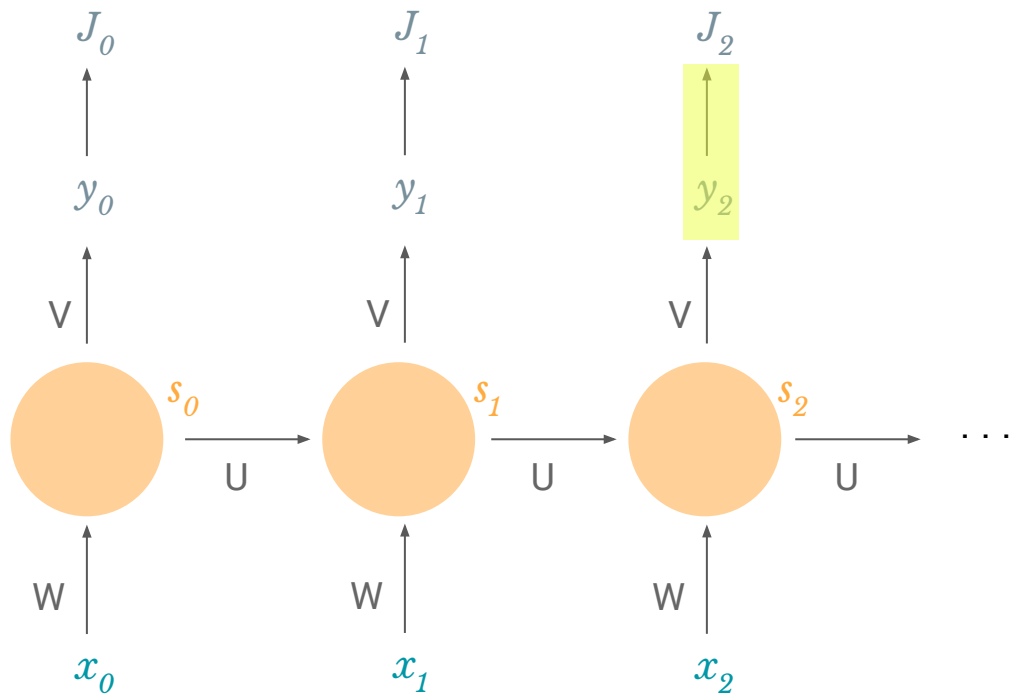


$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

so let's take a single timestep t :

$$\frac{\partial J_2}{\partial W}$$

let's try it out for W with the **chain rule**:

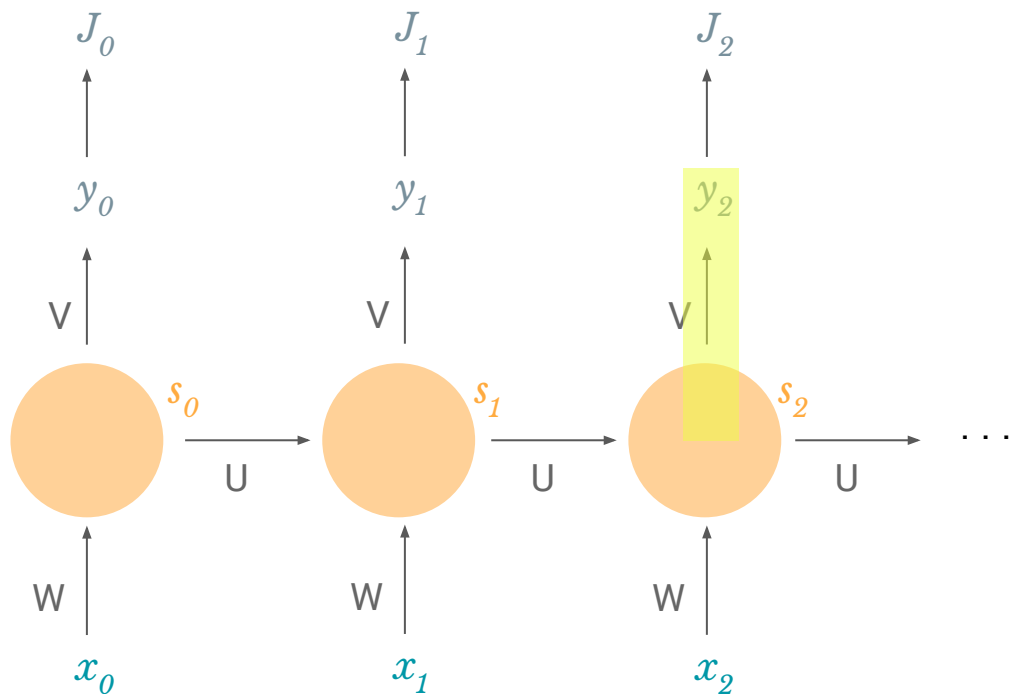


$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

so let's take a single timestep t :

$$\frac{\partial J_2}{\partial W} = \frac{\partial J_2}{\partial y_2}$$

let's try it out for W with the **chain rule**:

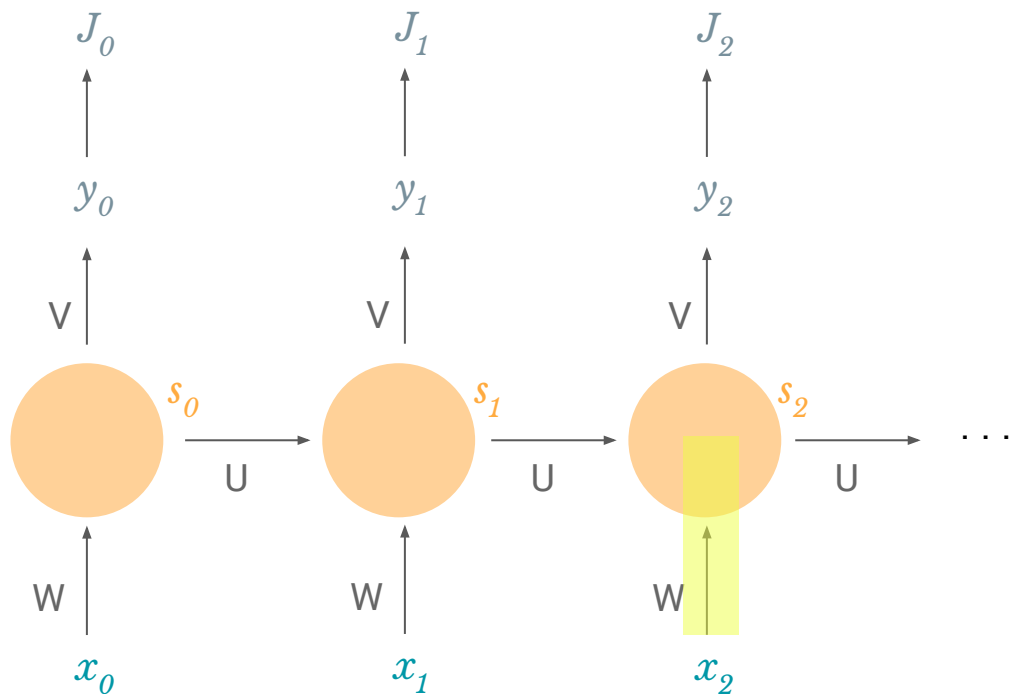


$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

so let's take a single timestep t :

$$\frac{\partial J_2}{\partial W} = \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2}$$

let's try it out for W with the **chain rule**:

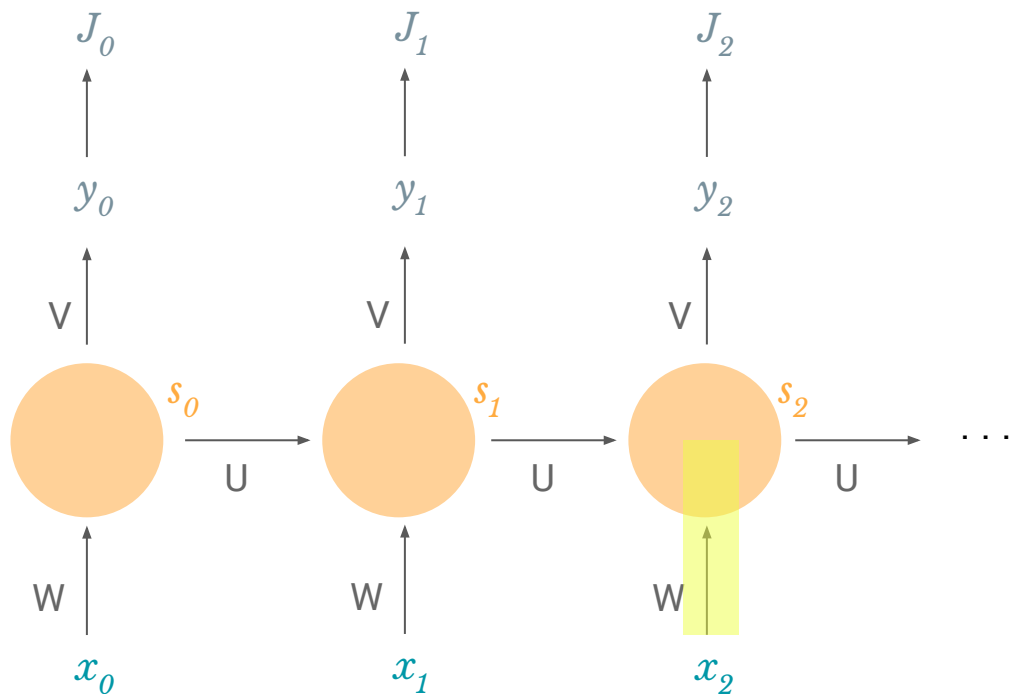


$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

so let's take a single timestep t :

$$\frac{\partial J_2}{\partial W} = \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \frac{\partial s_2}{\partial W}$$

let's try it out for W with the **chain rule**:



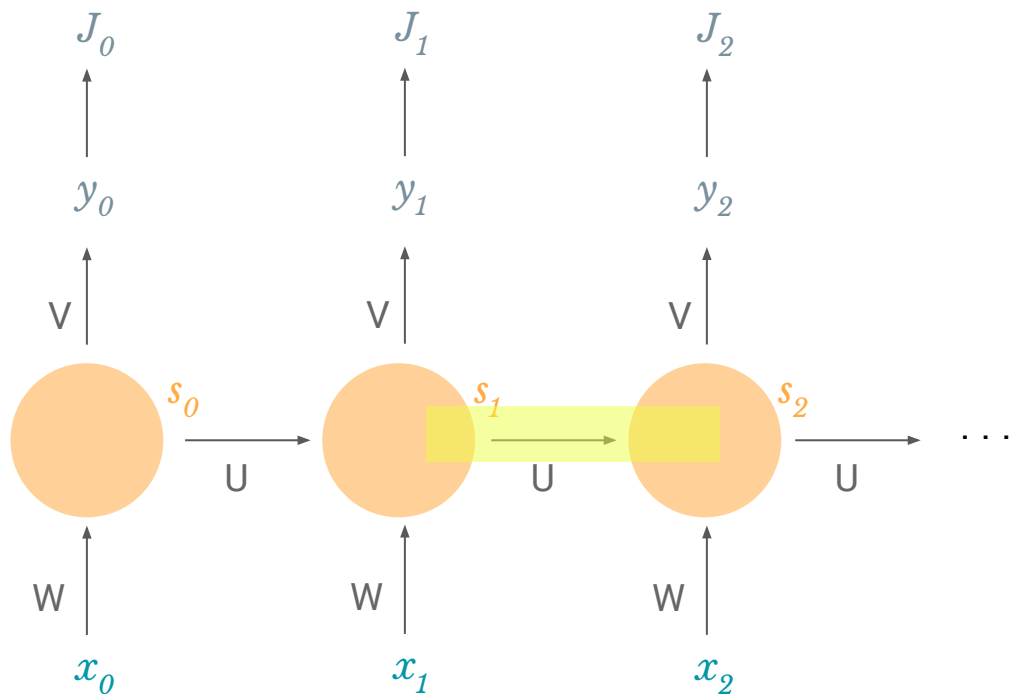
$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

so let's take a single timestep t :

$$\frac{\partial J_2}{\partial W} = \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \frac{\partial s_2}{\partial W}$$

but wait...

let's try it out for W with the **chain rule**:



$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

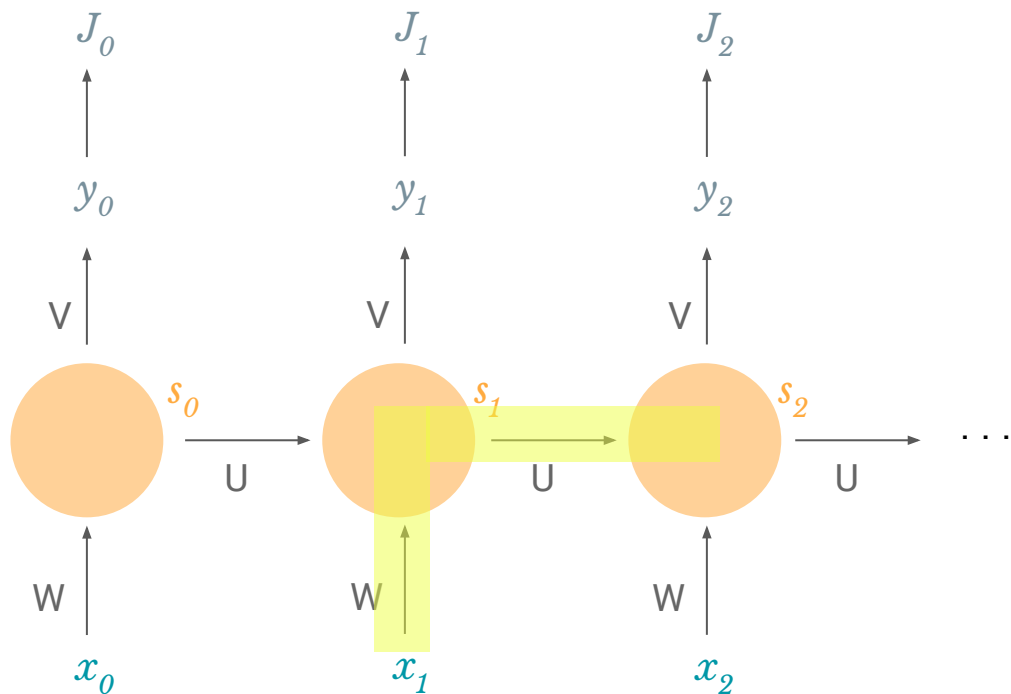
so let's take a single timestep t :

$$\frac{\partial J_2}{\partial W} = \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \frac{\partial s_2}{\partial W}$$

but wait...

$$s_2 = \tanh(U s_1 + W x_2)$$

let's try it out for W with the **chain rule**:



$$\frac{\partial J}{\partial W} = \sum_t \frac{\partial J_t}{\partial W}$$

so let's take a single timestep t :

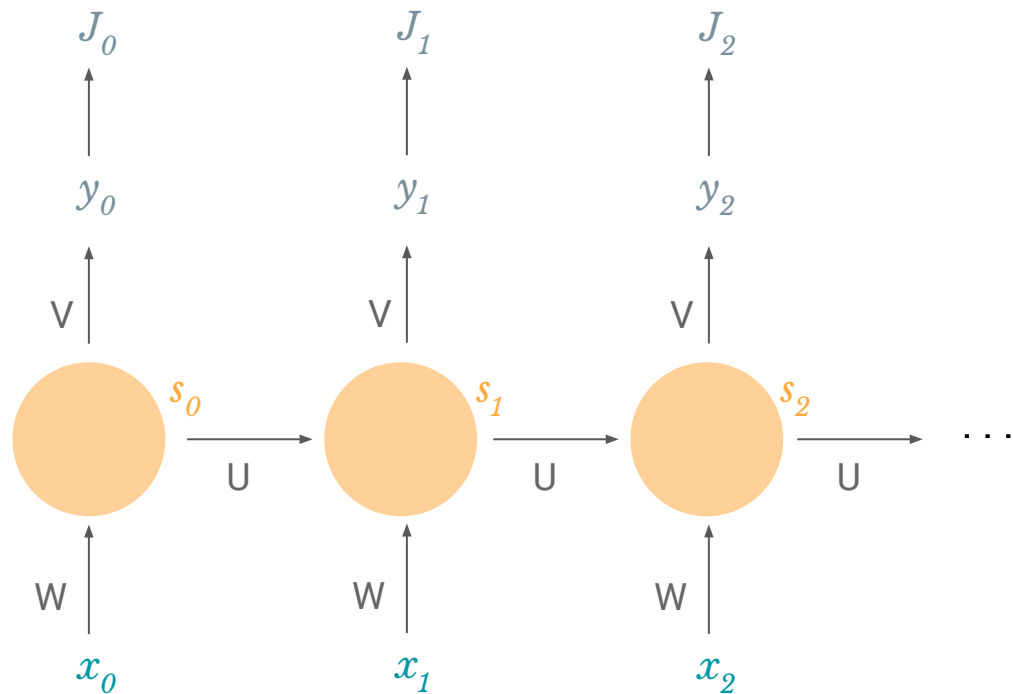
$$\frac{\partial J_2}{\partial W} = \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \frac{\partial s_2}{\partial W}$$

but wait...

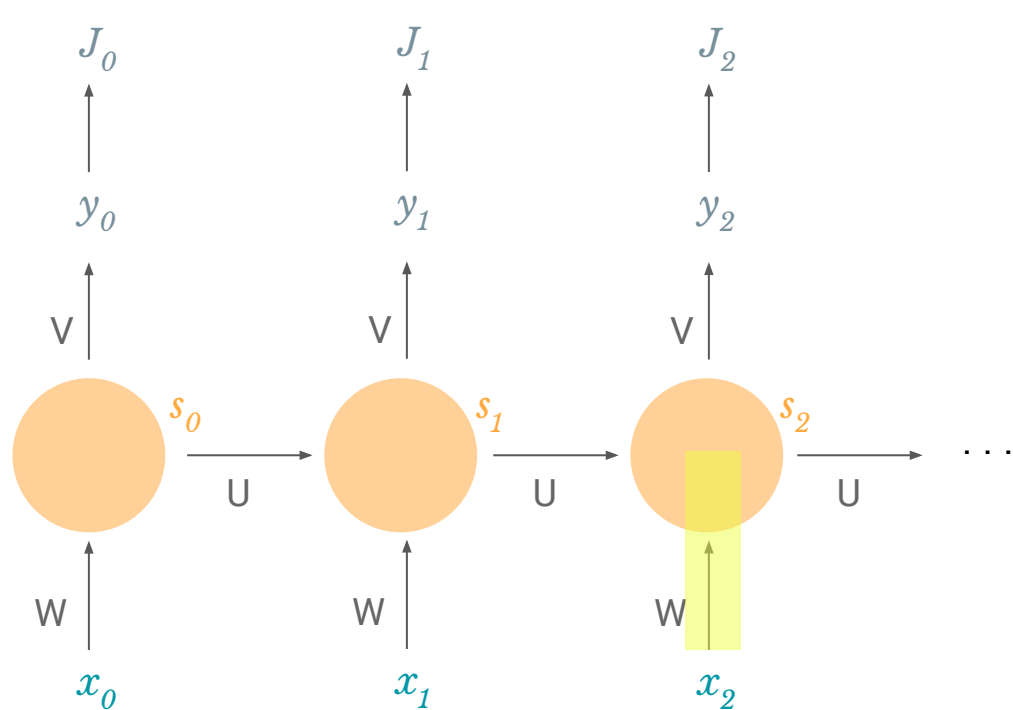
$$s_2 = \tanh(U s_1 + W x_2)$$

s_1 also depends on W so we can't just treat $\frac{\partial s_2}{\partial W}$ as a constant!

how does s_2 depend on W ?

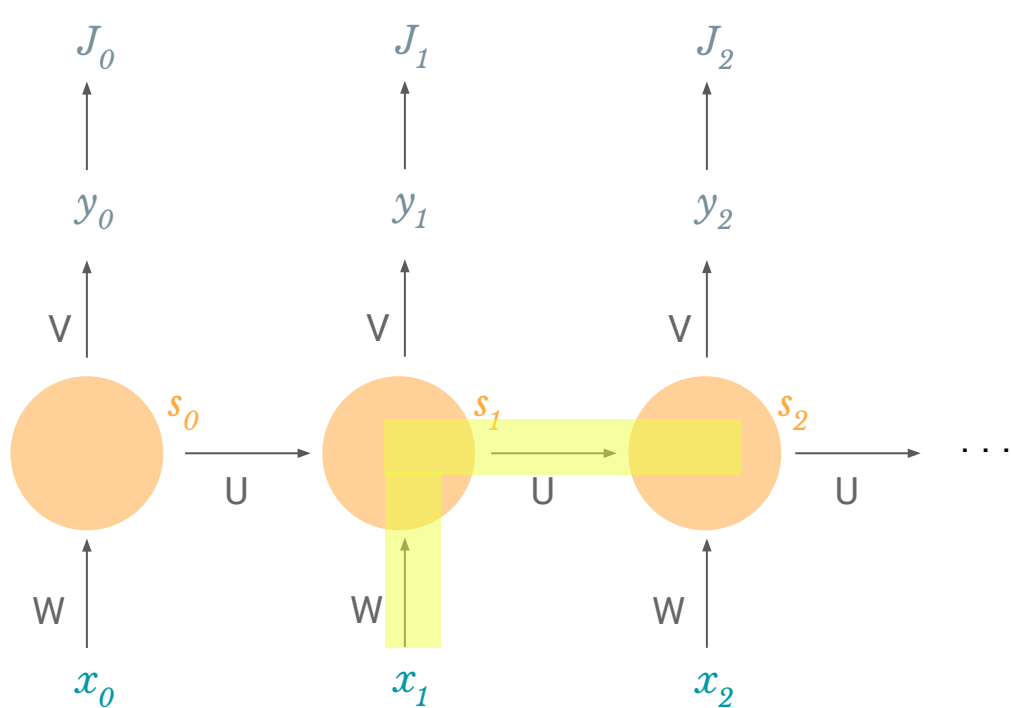


how does s_2 depend on W ?



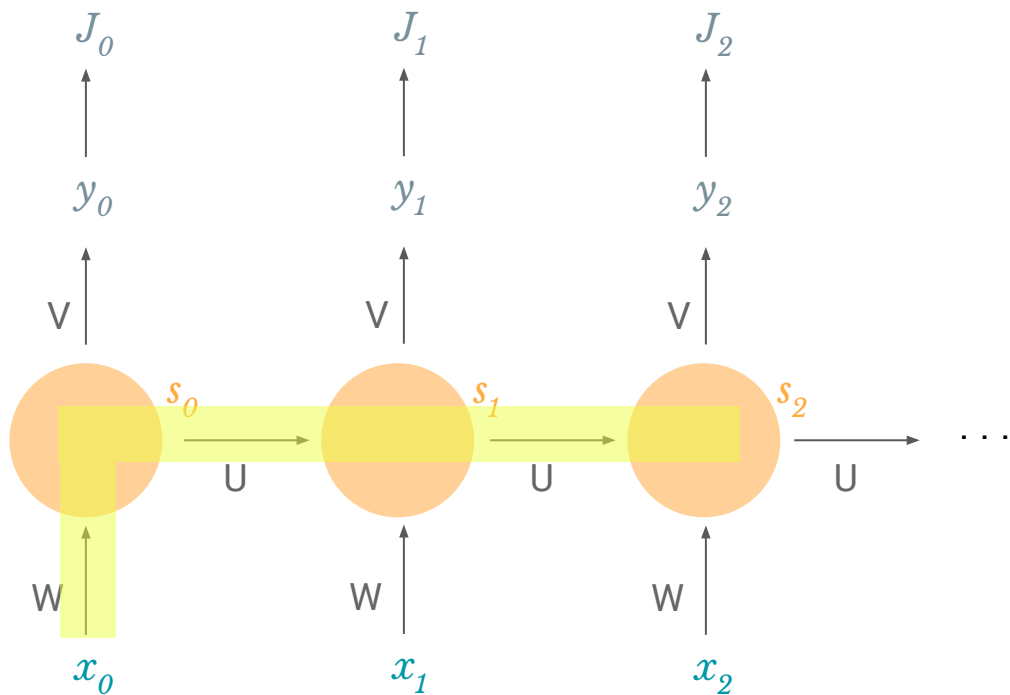
$$\frac{\partial s_2}{\partial W}$$

how does s_2 depend on W ?



$$\frac{\partial s_2}{\partial W} + \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W}$$

how does s_2 depend on W ?



$$\begin{aligned} & \frac{\partial s_2}{\partial W} \\ & + \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W} \\ & + \frac{\partial s_2}{\partial s_0} \frac{\partial s_0}{\partial W} \end{aligned}$$

backpropagation through time:

$$\frac{\partial J_2}{\partial W} = \sum_{k=0}^2 \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \underbrace{\frac{\partial s_2}{\partial s_k} \frac{\partial s_k}{\partial W}}_{\text{Contributions of } W \text{ in previous timesteps to the error at timestep } t}$$

Contributions of W in previous timesteps to the error at timestep t

backpropagation through time:

$$\frac{\partial J_t}{\partial W} = \sum_{k=0}^t \frac{\partial J_t}{\partial y_t} \frac{\partial y_t}{\partial s_t} \underbrace{\frac{\partial s_t}{\partial s_k} \frac{\partial s_k}{\partial W}}_{\text{Contributions of } W \text{ in previous timesteps to the error at timestep } t}$$

Contributions of W in previous timesteps to the error at timestep t

why are RNNs **hard to train**?

problem: vanishing gradient

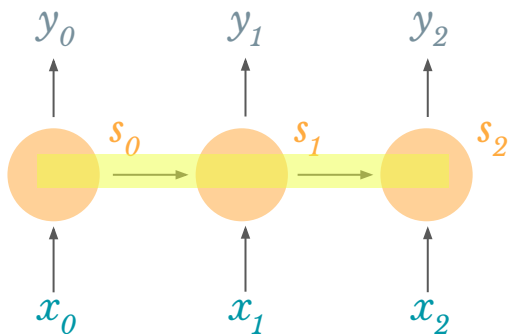
$$\frac{\partial J_2}{\partial W} = \sum_{k=0}^2 \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \frac{\partial s_2}{\partial s_k} \frac{\partial s_k}{\partial W}$$

problem: vanishing gradient

$$\frac{\partial J_2}{\partial W} = \sum_{k=0}^2 \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \frac{\partial s_2}{\partial s_k} \frac{\partial s_k}{\partial W}$$

problem: vanishing gradient

$$\frac{\partial J_2}{\partial W} = \sum_{k=0}^2 \frac{\partial J_2}{\partial y_2} \frac{\partial y_2}{\partial s_2} \frac{\partial s_2}{\partial s_k} \frac{\partial s_k}{\partial W}$$



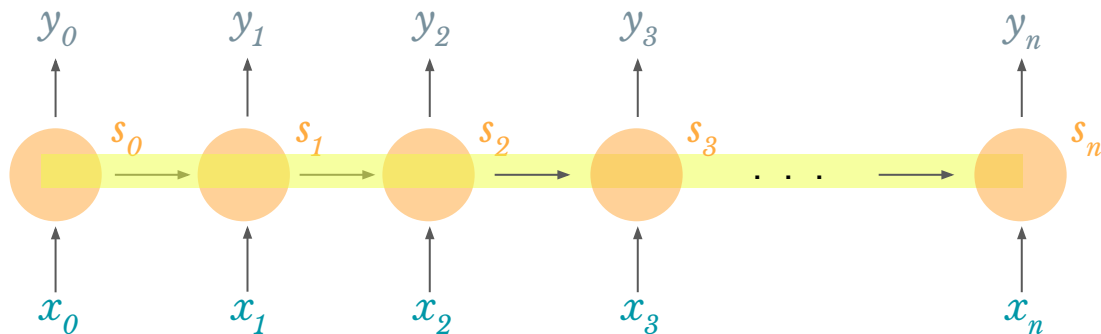
at $k = 0$:

$$\frac{\partial s_2}{\partial s_0} = \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial s_0}$$

problem: vanishing gradient

$$\frac{\partial J_n}{\partial W} = \sum_{k=0}^n \frac{\partial J_n}{\partial y_n} \frac{\partial y_n}{\partial s_n} \frac{\partial s_n}{\partial s_k} \frac{\partial s_k}{\partial W}$$

$\frac{\partial s_n}{\partial s_{n-1}} \frac{\partial s_{n-1}}{\partial s_{n-2}} \cdots \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial s_0}$

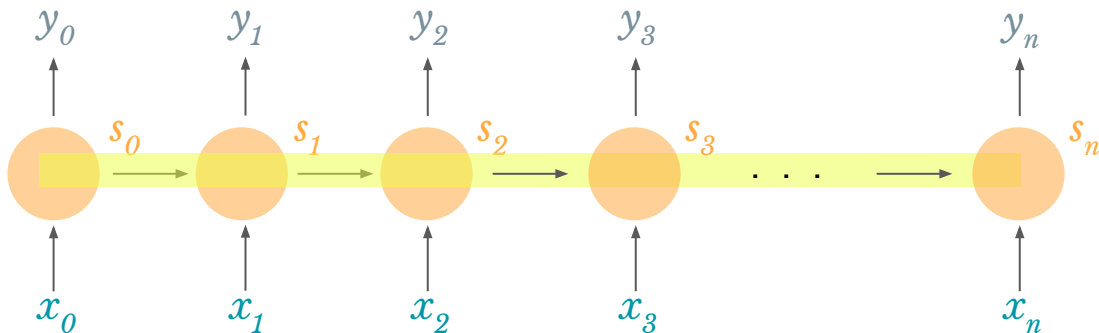


problem: vanishing gradient

$$\frac{\partial J_n}{\partial W} = \sum_{k=0}^n \frac{\partial J_n}{\partial y_n} \frac{\partial y_n}{\partial s_n} \frac{\partial s_n}{\partial s_k} \frac{\partial s_k}{\partial W}$$

$$\frac{\partial s_n}{\partial s_{n-1}} \frac{\partial s_{n-1}}{\partial s_{n-2}} \cdots \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial s_0}$$

as the gap between timesteps gets bigger, this product gets longer and longer!



problem: vanishing gradient

$$\frac{\partial s_n}{\partial s_{n-1}} \frac{\partial s_{n-1}}{\partial s_{n-2}} \cdot \dots \cdot \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial s_0}$$

problem: vanishing gradient

what are each of these terms?



$$\frac{\partial s_n}{\partial s_{n-1}} \frac{\partial s_{n-1}}{\partial s_{n-2}} \cdot \dots \cdot \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial s_0}$$

problem: vanishing gradient

what are each of these terms? →

$$\frac{\partial s_n}{\partial s_{n-1}} = W^T \text{diag}[f'(W s_{j-1} + U x_j)]$$

$$\frac{\partial s_n}{\partial s_{n-1}} \frac{\partial s_{n-1}}{\partial s_{n-2}} \cdots \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial s_0}$$

W = sampled from
standard normal
distribution = mostly < 1

f = tanh or sigmoid so $f' < 1$

problem: vanishing gradient

what are each of these terms? →

$$\frac{\partial s_n}{\partial s_{n-1}} = W^T \text{diag}[f'(W s_{j-1} + U x_j)]$$

$$\frac{\partial s_n}{\partial s_{n-1}} \frac{\partial s_{n-1}}{\partial s_{n-2}} \cdots \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial s_0}$$

W = sampled from standard normal distribution = mostly < 1

f = tanh or sigmoid so $f' < 1$

we're multiplying a lot of **small numbers** together.

we're multiplying a lot of **small numbers** together.

so what?

errors due to further back timesteps have increasingly **smaller gradients**.

so what?

parameters become biased to **capture shorter-term dependencies**.

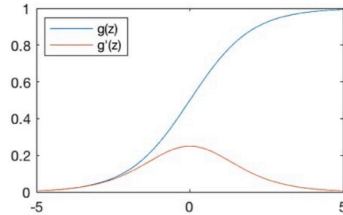
“In France, I had a great time and I learnt some of the _____ language.”



our parameters are not trained to capture long-term dependencies, so the word we predict will mostly depend on the previous few words, not much earlier ones

Common Activation Functions

Sigmoid Function

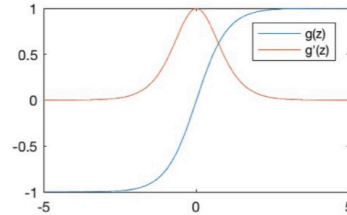


$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

 `tf.nn.sigmoid(z)`

Hyperbolic Tangent

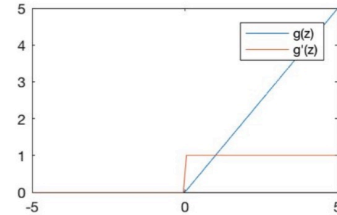


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

 `tf.nn.tanh(z)`

Rectified Linear Unit (ReLU)



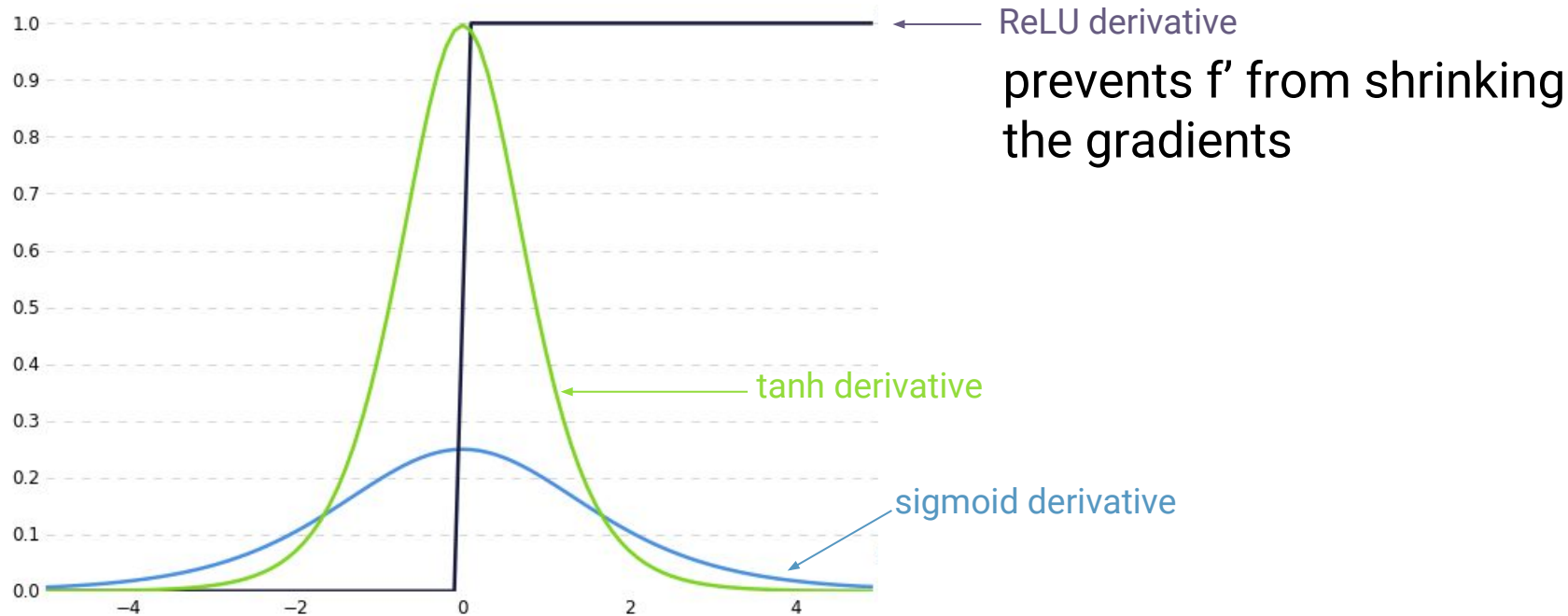
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

 `tf.nn.relu(z)`

NOTE: All activation functions are **non-linear**

solution #1: activation functions



solution #2: initialization

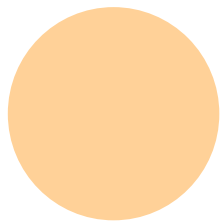
weights initialized to identity matrix $\longrightarrow I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$
biases initialized to zeros

prevents W from shrinking the gradients

a different type of solution:
more complex cells

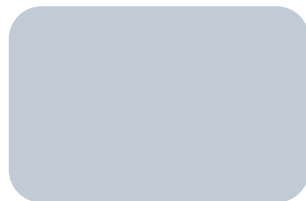
solution #3: **gated cells**

rather each node being just a simple RNN cell, make each node a more **complex unit with gates** controlling what information is passed through.



RNN

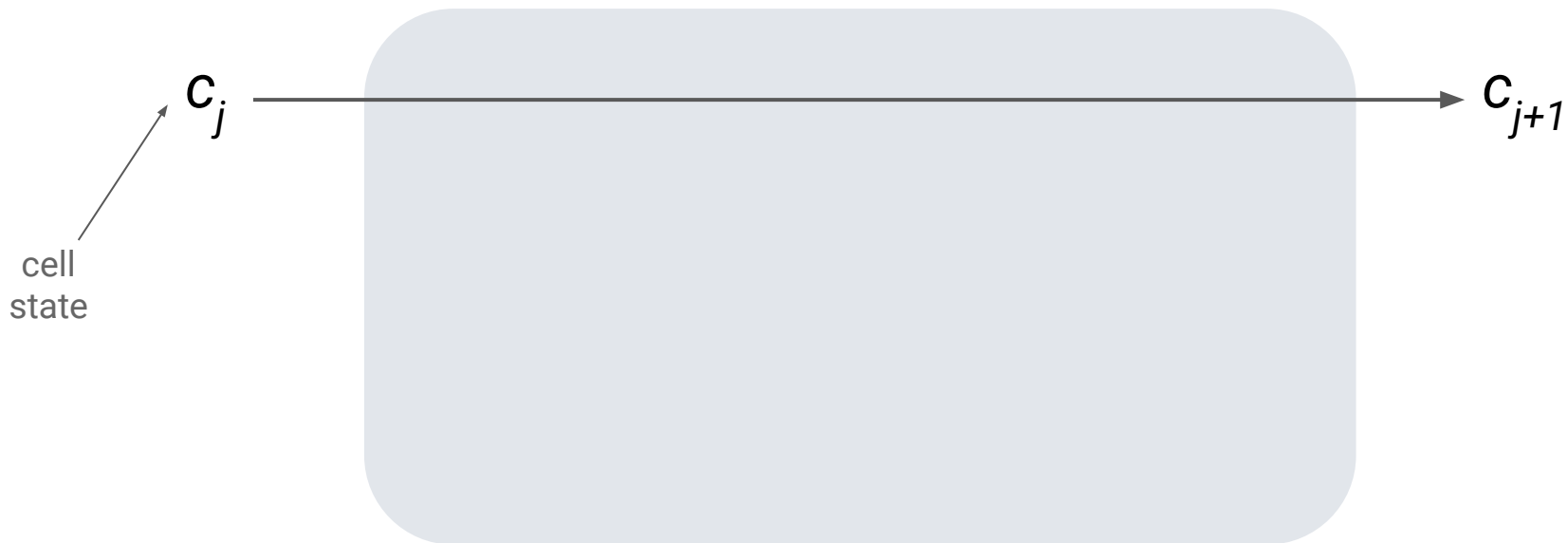
vs



LSTM, GRU, etc

Long short term memory cells are able to keep track of information throughout many timesteps.

solution #3: more on **LSTMs**



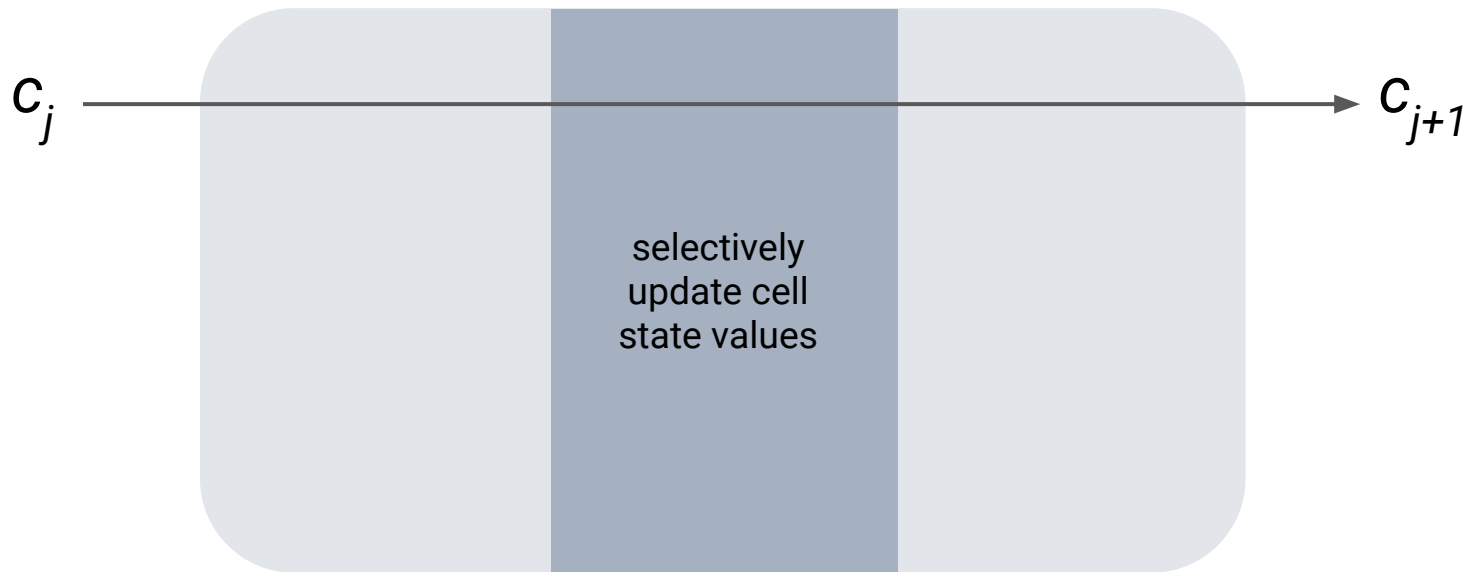
solution #3: more on **LSTMs**

e.g. forget an old subject when seeing a new one



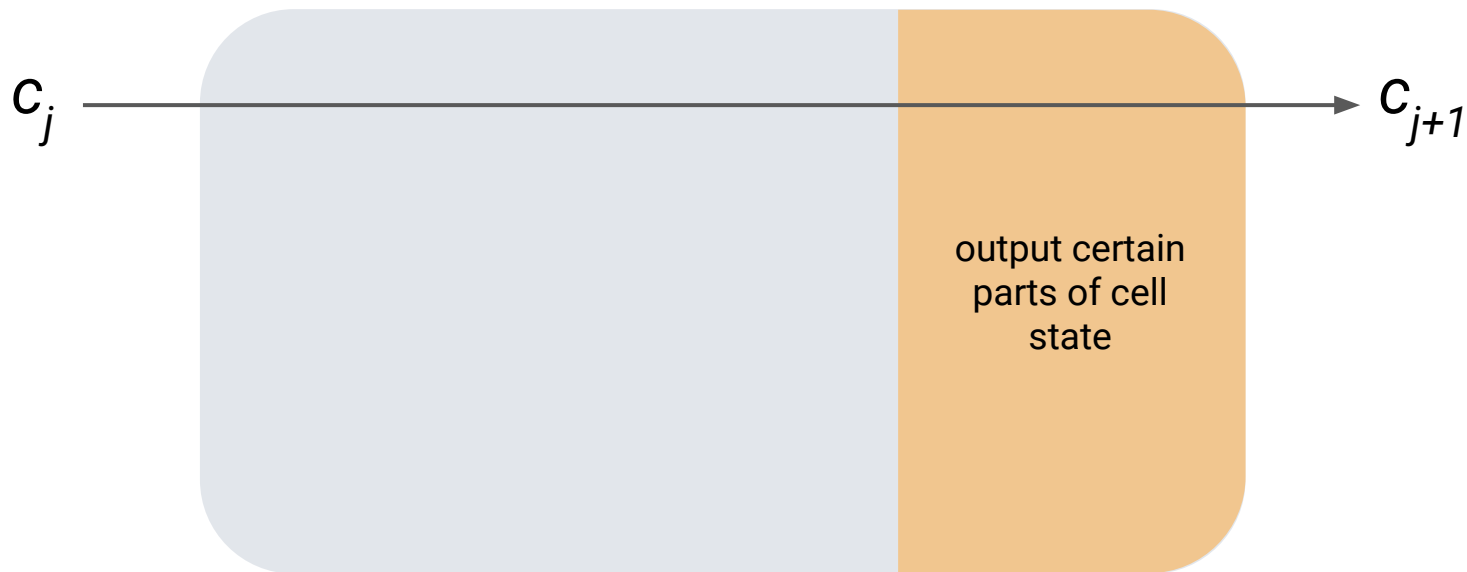
solution #3: more on **LSTMs**

E.g., based on the new subject, update gender, singular/plural



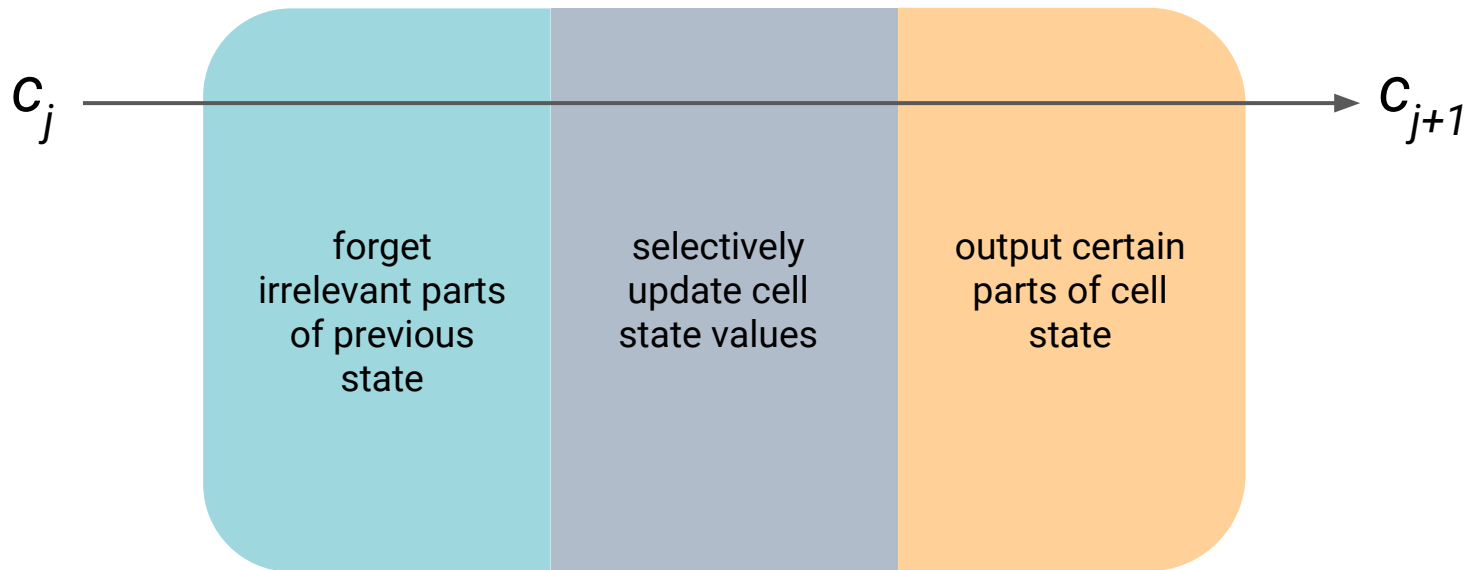
solution #3: more on **LSTMs**

E.g., output a verb after just seeing a subject



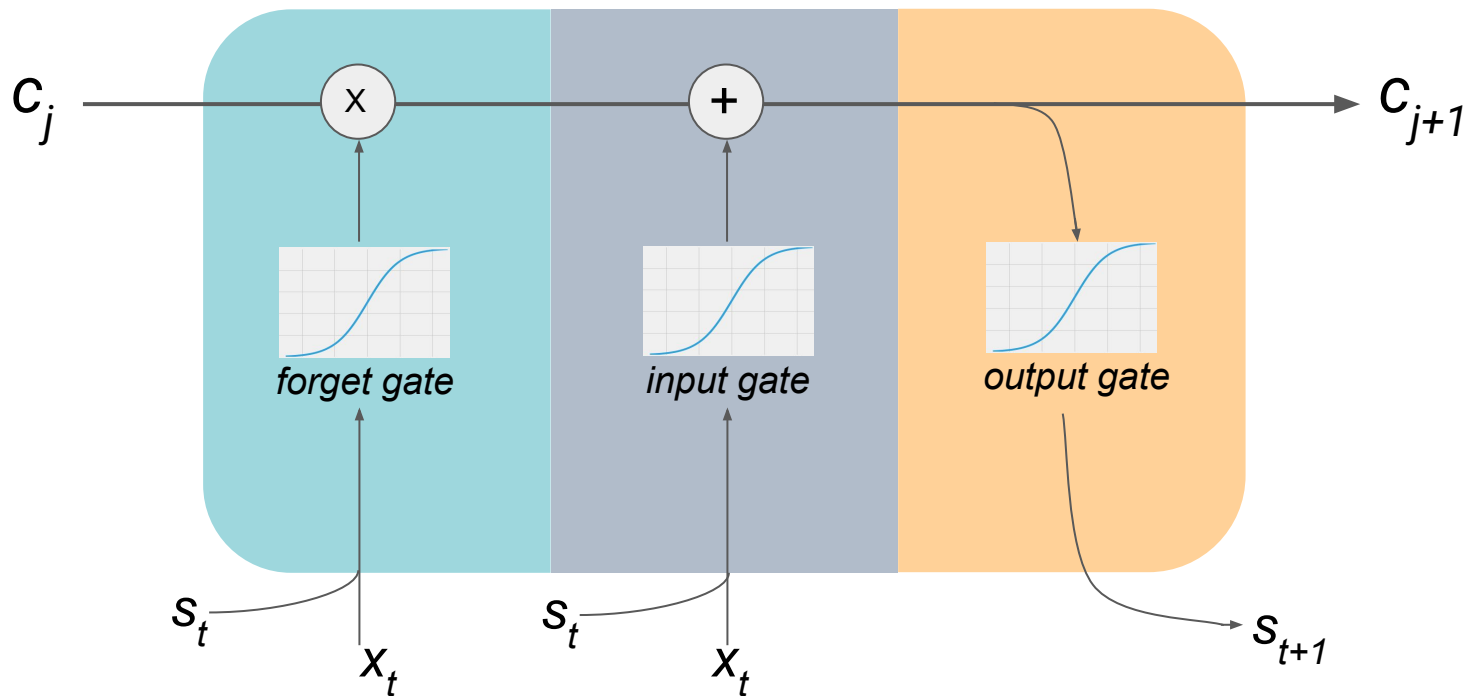
solution #3: more on **LSTMs**

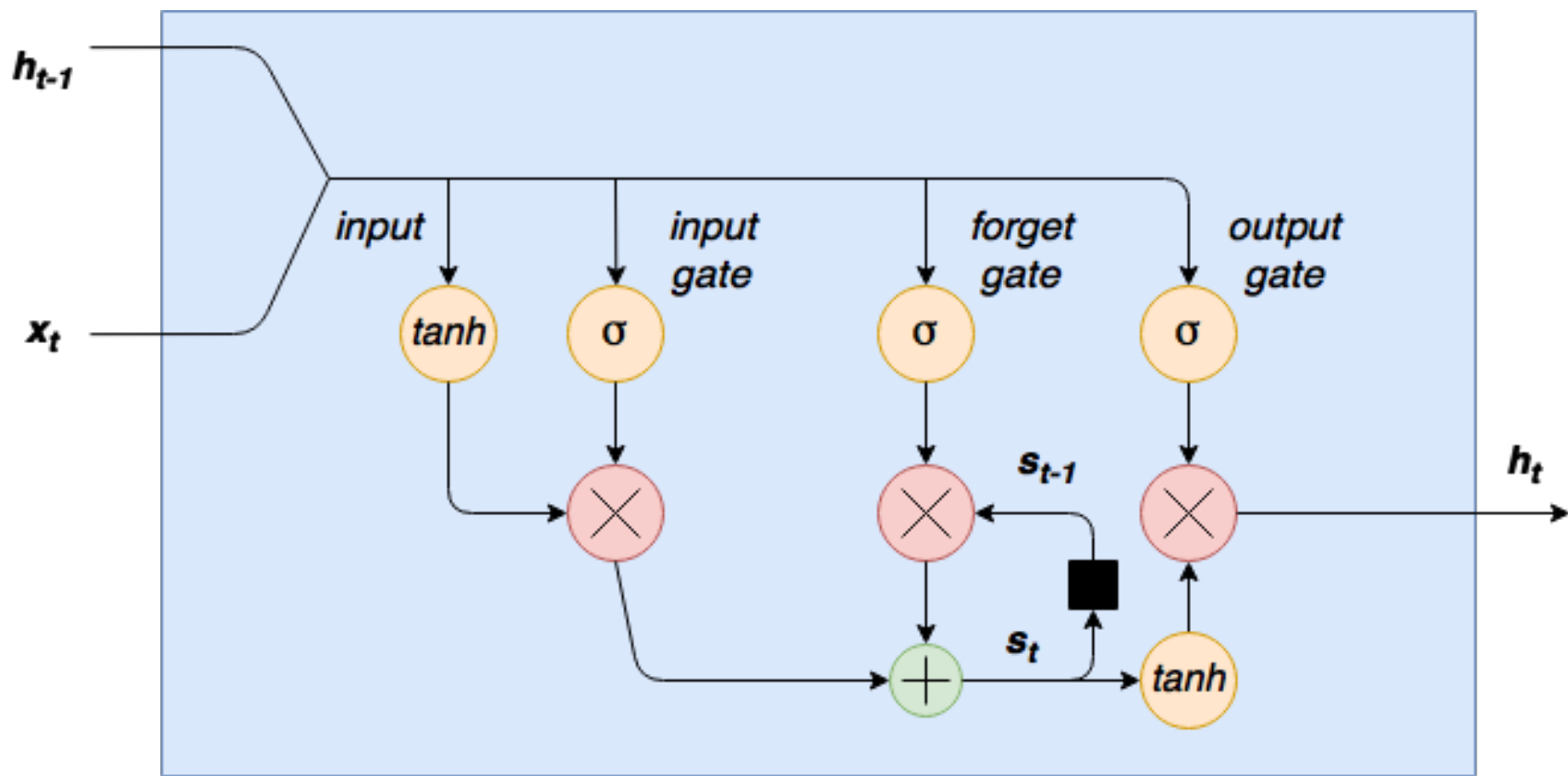
Implemented using logic gates with Sigmoid functions. See the end of slides for examples



solution #3: more on **LSTMs**

Cell states are now separate from the output states





why do LSTMs help?

1. forget gate allows information to **pass through unchanged**
2. **cell state is separate** from what's outputted
3. s_j depends on s_{j-1} through **addition!**
→ derivatives don't expand into a long product!

possible task: classification (i.e. sentiment)

Predict the end result (sentiment) after seeing the sequence, as opposed predicting the next word



@HVSVN



Don't fly with @British_Airways.
They can't keep track of your
luggage.



:(



Kim Kardashian ✓
@KimKardashian



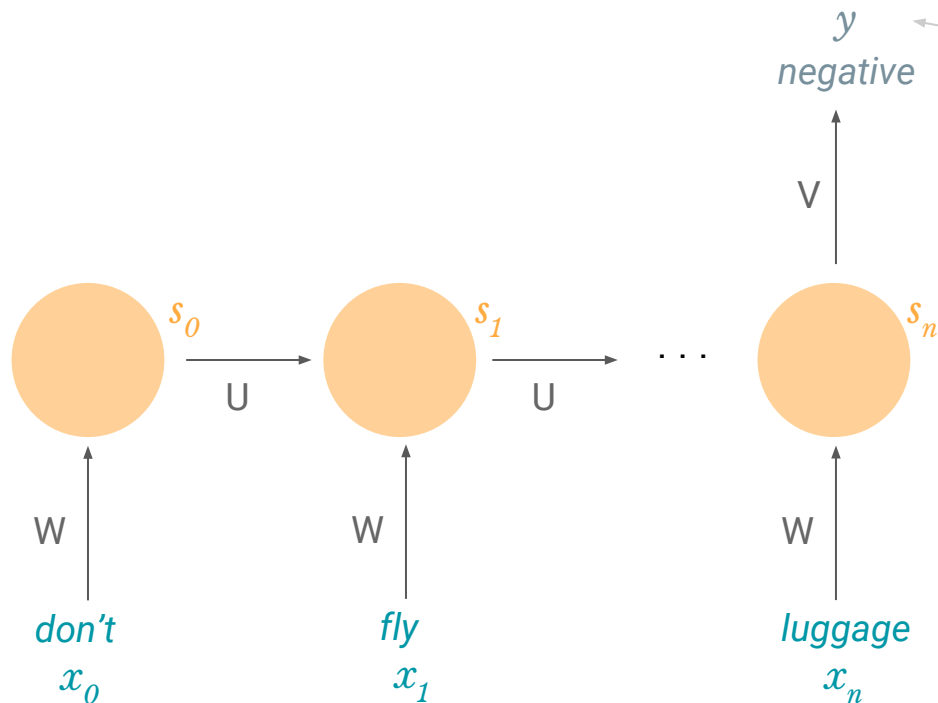
Following

Happy Birthday to my best friend, the ♥ of
my life, my soul!!!! I love you beyond words!
[instagram.com/p/aTgfl-OS-a/](https://www.instagram.com/p/aTgfl-OS-a/)



:)

possible task: classification (i.e. sentiment)



y is a probability distribution over possible classes (like positive, negative, neutral), aka a *softmax*

The final cell state s_n summarizes all the information in the sequence.

We only use s_n to do the prediction

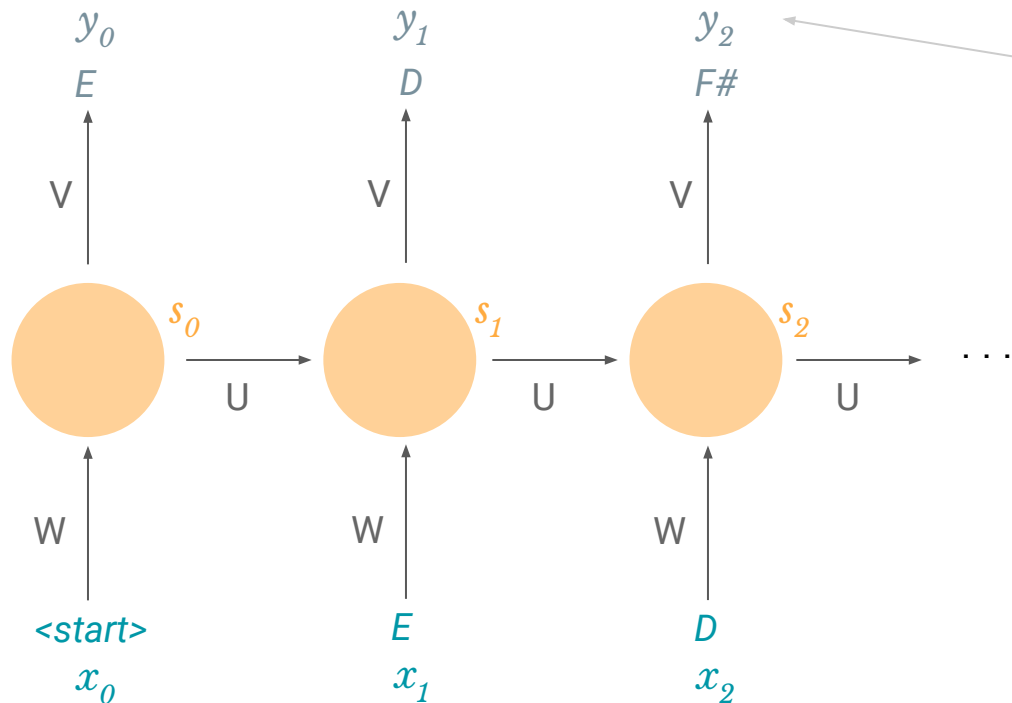
possible task: music generation



[Play video at 23:38](#)

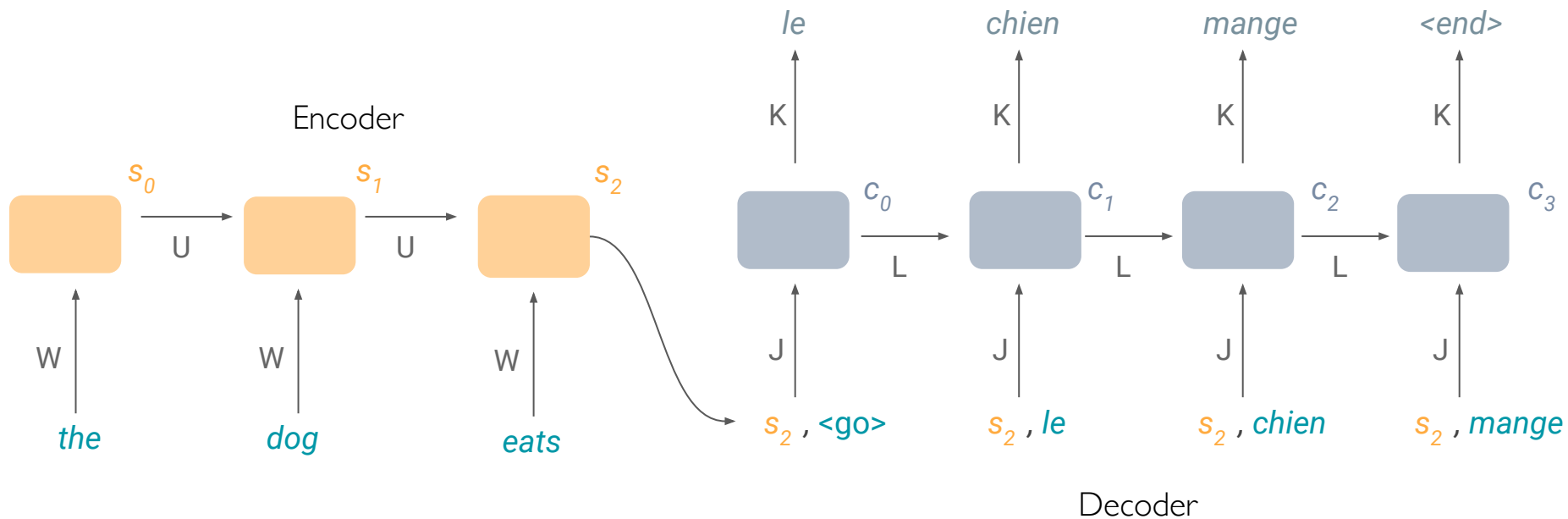
*Music by: Francesco Marchesani,
Computer Science Engineer, PoliMi*

possible task: music generation

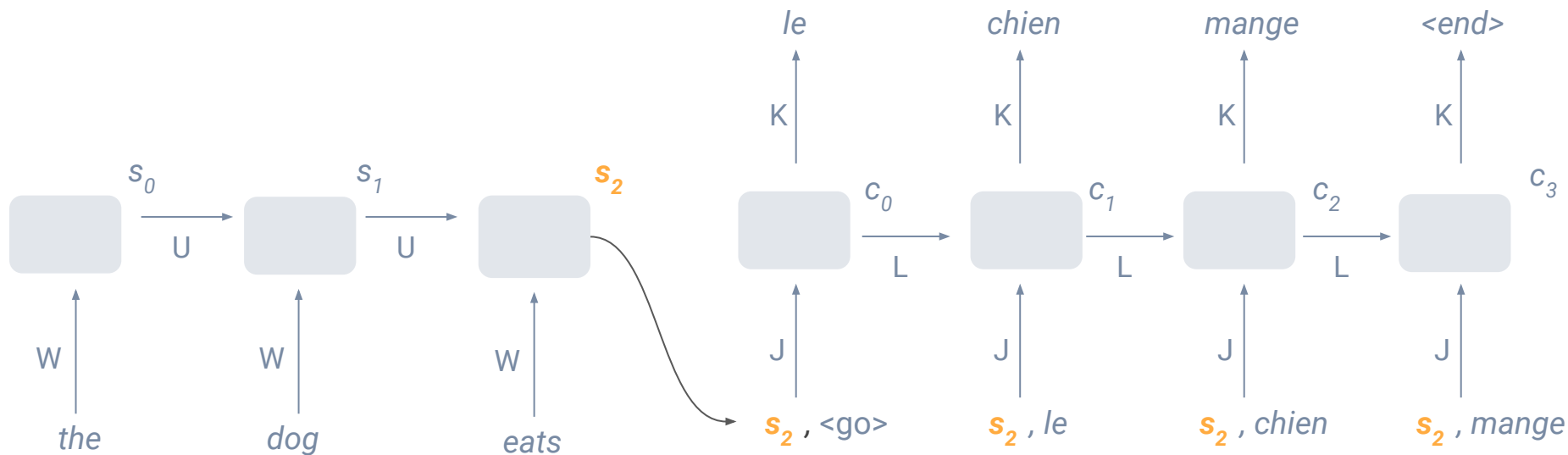


y_i is actually a probability distribution over possible next notes, aka a *softmax*

possible task: machine translation

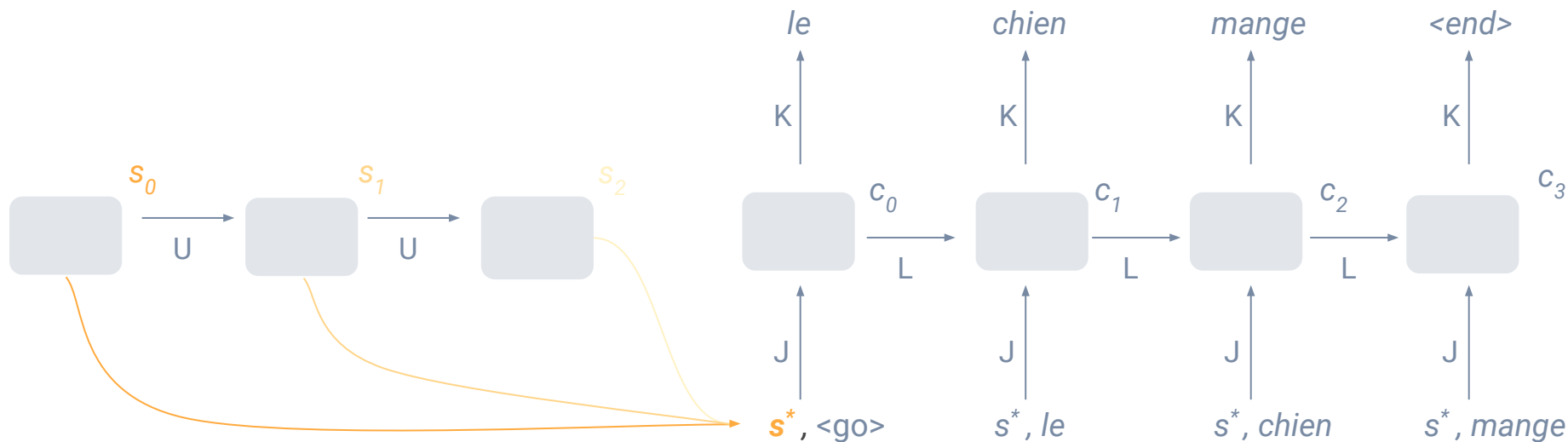


problem: a single encoding is limiting



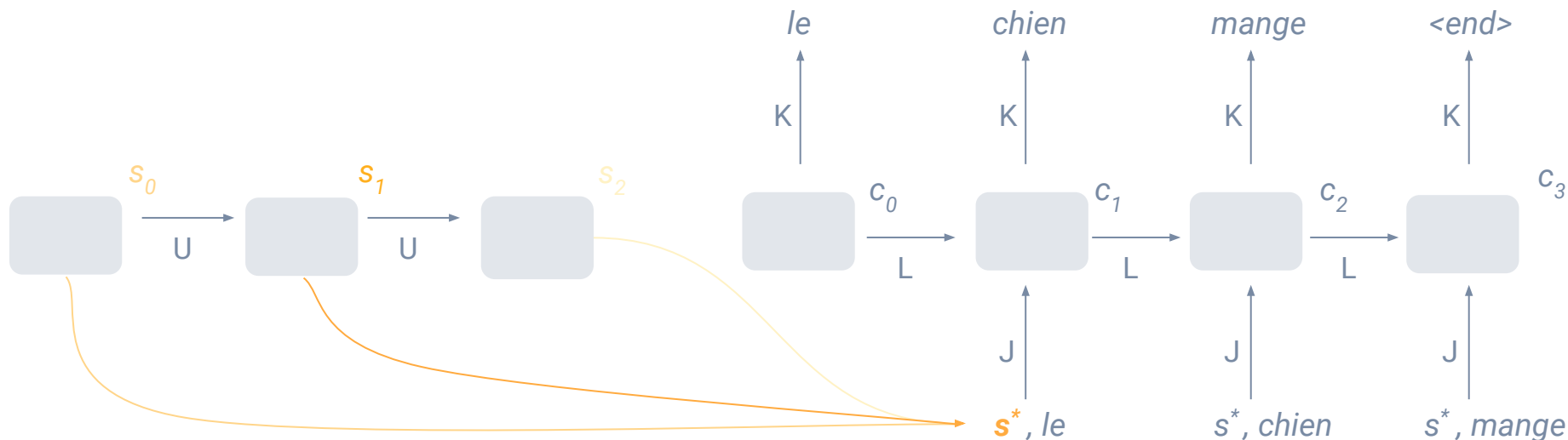
all the decoder knows about the input sentence is in one fixed length vector, s_2

solution: attend over all encoder states



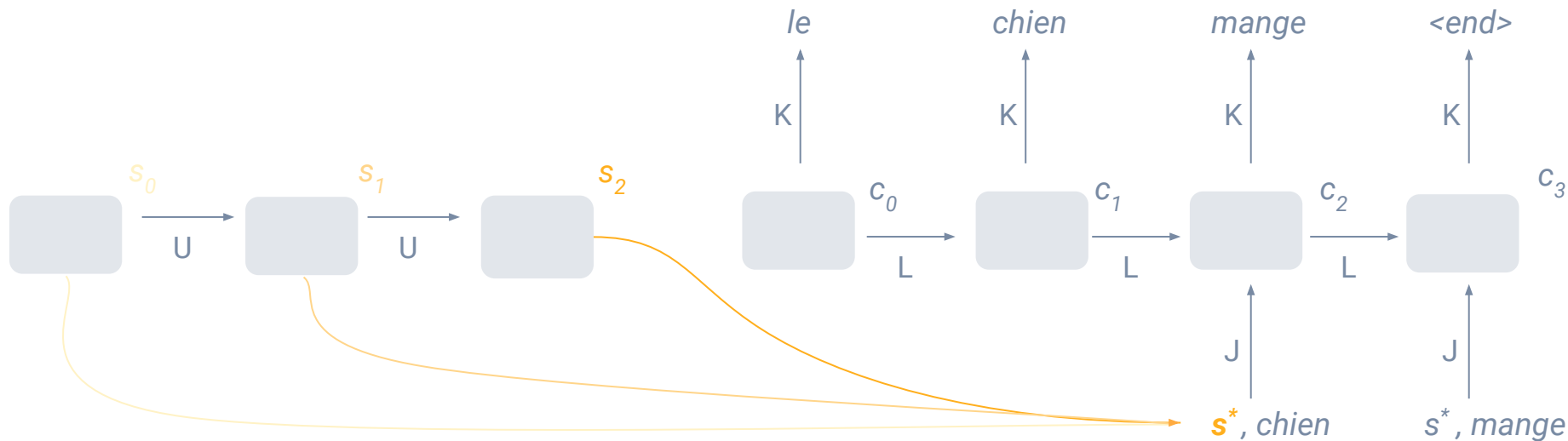
Attention: Instead of taking only the final state, we take a weighted sum of all states in the encoder.
To predict the first word in the target language, we probably want to give a higher weight to the first cell state s_0 .

solution: attend over all encoder states



similarly, to predict the second word, we will give a higher weight to the second cell state s_1 , but we will also consider s_0 and s_2 to obtain information about gender and tense.

solution: attend over all encoder states



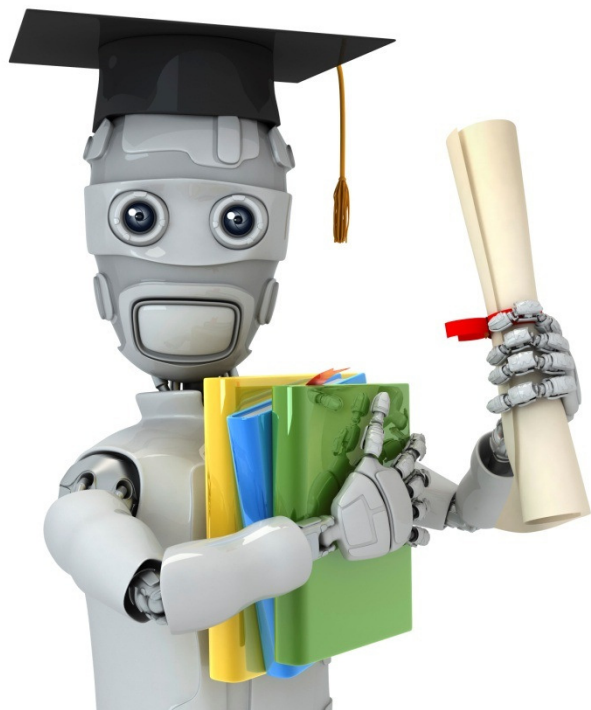
Training: include all the weights in the weight sum in the back propagation.

now we can model **sequences!**

- why recurrent neural networks?
- training them with backpropagation through time
- solving the vanishing gradient problem with activation functions, initialization, and gated cells (like LSTMs)
- building models for classification, music generation and machine translation
- using attention mechanisms

and there's lots more to do!

- extending our models to timeseries + waveforms
- complex language models to generate long text or books
- language models to generate code
- controlling cars + robots
- predicting stock market trends
- summarizing books + articles
- handwriting generation
- multilingual translation models
- ... many more!



Machine Learning

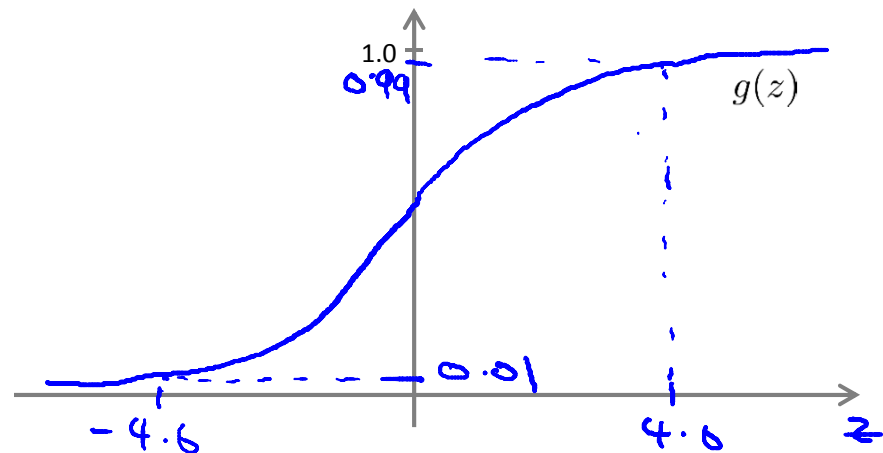
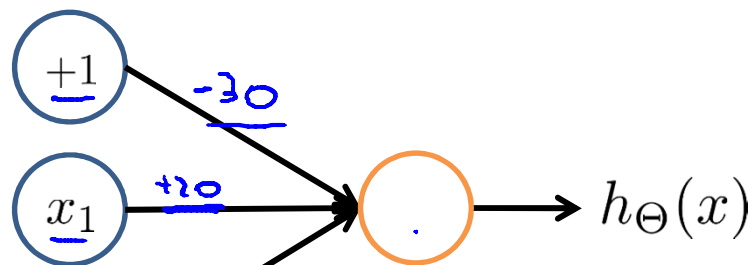
Neural Networks: Representation

Examples and intuitions I

Simple example: AND

→ $x_1, x_2 \in \{0, 1\}$

→ $y = x_1 \text{ AND } x_2$

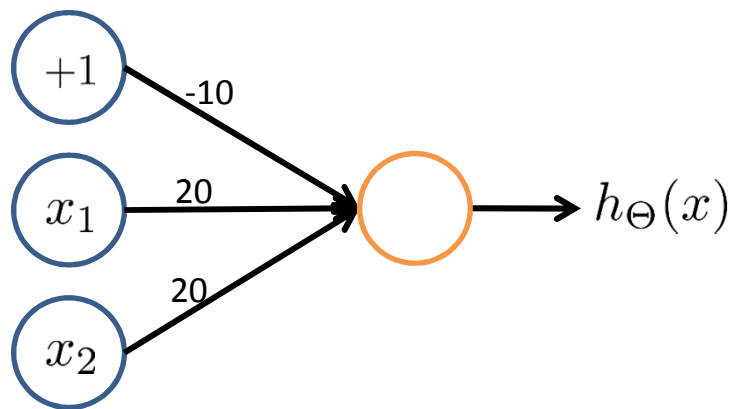


x_1	x_2	$h_{\Theta}(x)$
0	0	$g(-30) \approx 0$
→ 0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
→ 1	1	$g(10) \approx 1$

→ $h_{\Theta}(x) = g\left(\underbrace{-30}_{\textcircled{10}} + \underbrace{20}_{\textcircled{11}}x_1 + \underbrace{20}_{\textcircled{12}}x_2\right)$

$h_{\Theta}(x) \approx x_1 \text{ AND } x_2$

Example: OR function



$$g(-10 + 20x_1 + 20x_2)$$

x_1	x_2	$h_{\Theta}(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	≈ 1
1	1	≈ 1

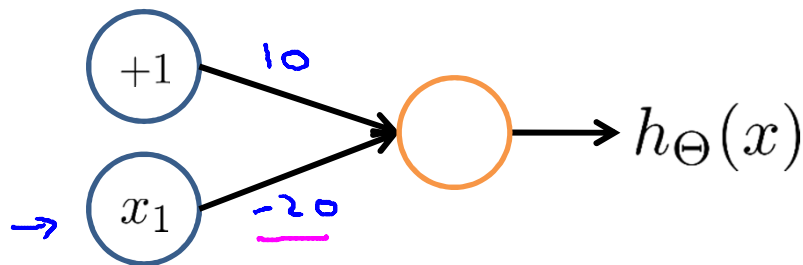
→ x_1 AND x_2

→ x_1 OR x_2

$\{0,1\}$.

Negation:

NOT x_1

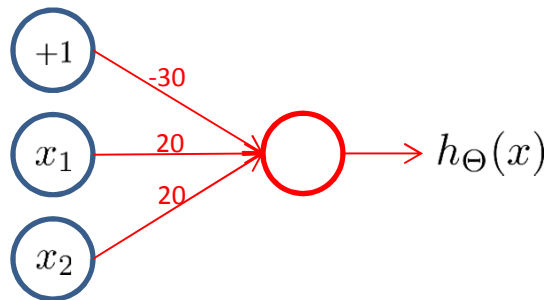
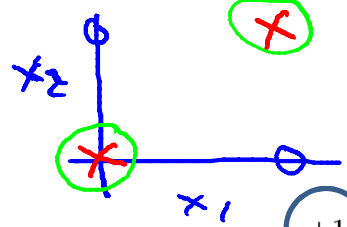


x_1	$h_{\Theta}(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

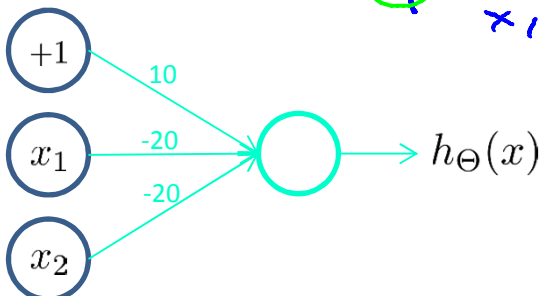
$$h_{\Theta}(x) = g(10 - 20x_1)$$

→ (NOT x_1) AND (NOT x_2)
(= 1 if and only if
→ $x_1 = x_2 = 0$)

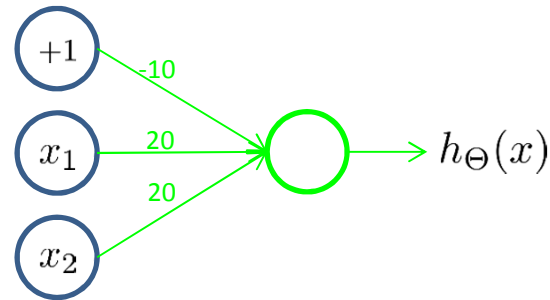
Putting it together: x_1 XNOR x_2



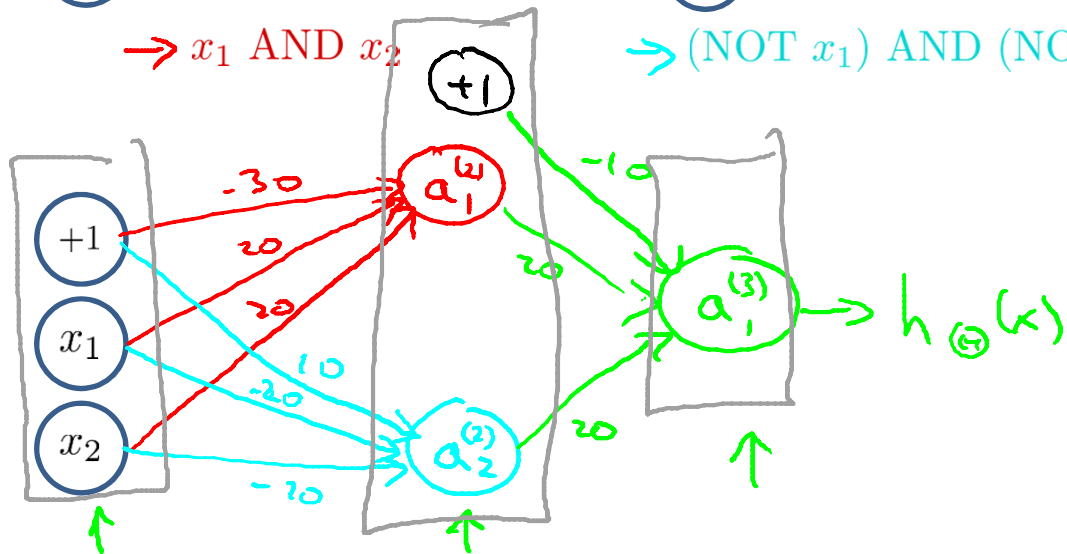
$\rightarrow x_1$ AND x_2



\rightarrow (NOT x_1) AND (NOT x_2)



$\rightarrow x_1$ OR x_2



x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1