

COMP 4446 / 5046

Lecture 3: Models – Non-linear

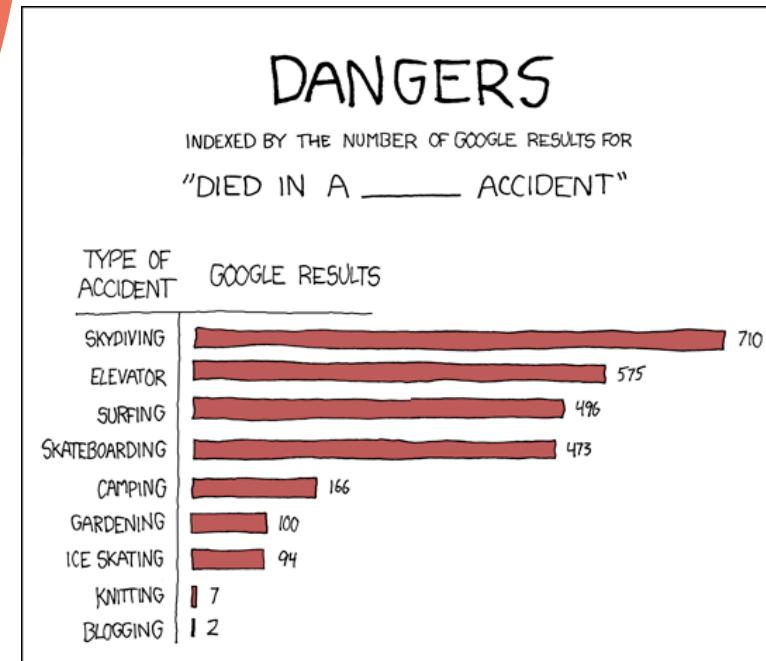
Jonathan K. Kummerfeld

Semester 1, 2025



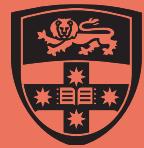
THE UNIVERSITY OF
SYDNEY

Dangers



[Zero results: 'snake charming' and 'haberdashery'.
(Things like 'car' and 'boating' and such are of course the highest, by a huge margin.)]

Source: <https://xkcd.com/369/>



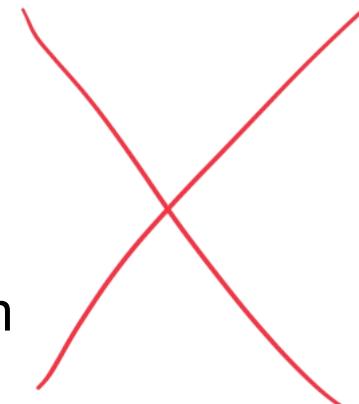
Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

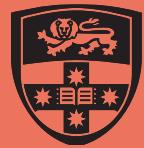
Deadlines

Assignment 1 - Due TOMORROW 11:59pm



Assignment 2 - Available on Thursday

Slip-day policy reminder:
Max 2 per assignment
5 over the semester



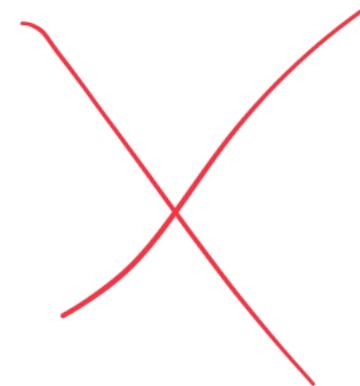
Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

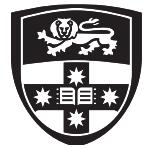
Quiz

10 minutes long at the start of the lecture



On paper - bring a dark pen

No devices permitted



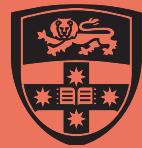
COMP 4446 / 5046
Lecture 3, 2025

Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

Neural Networks



Core idea: Add a non-linear function in the model

Model

Linear model

$$\text{scores(doc)} = \text{features(doc)} \cdot \text{weights}$$



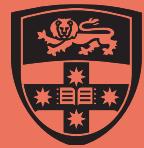
$$f(x) = x \cdot W$$

Non-Linear model

$$\begin{aligned} \text{scores(doc)} &= g(\text{features(doc)} \cdot \text{weights}_1) \cdot \text{weights}_2 \\ &= g(\begin{array}{ccccccccc} \text{---} & \text{---} \\ | & | & | & | & | & | & | & | & | \\ \text{---} & \text{---} \end{array} \cdot \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \end{array}) \cdot \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \end{array} \end{aligned}$$

Sigmoid,
etc.

$$f(x) = g(x \cdot W_1) \cdot W_2 + \text{bias}$$

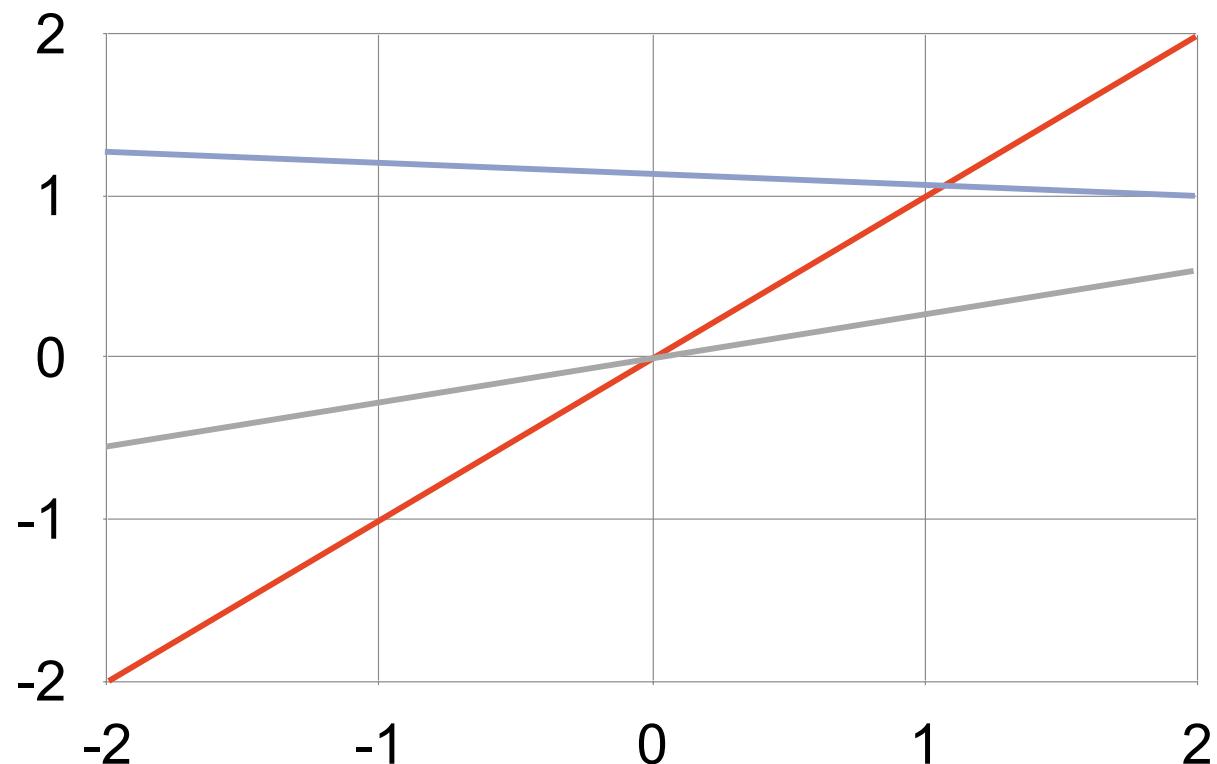


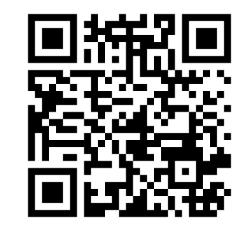
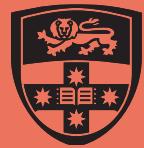
What is a non-linear function?

Linear functions

Model

$$f(x) = ax + b$$





What is a non-linear function?

Sigmoid

Model

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

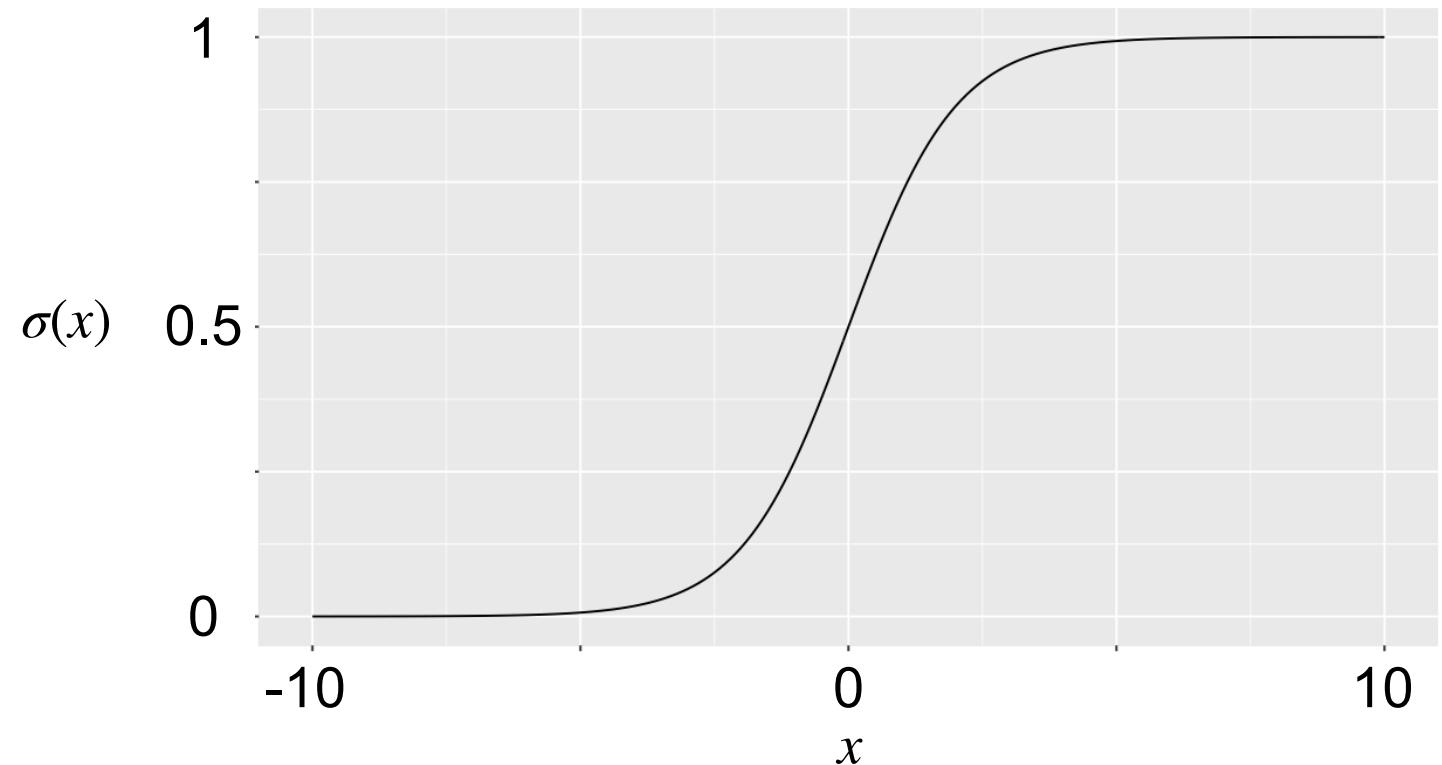
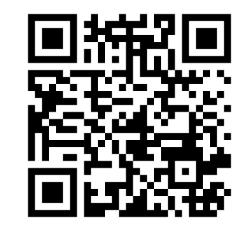
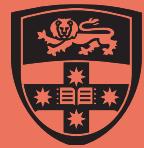


Figure from David Bamman's Lecture Slides for Info 159/259 at UC Berkeley



What is a non-linear function?

Hyperbolic Tangent

Model

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

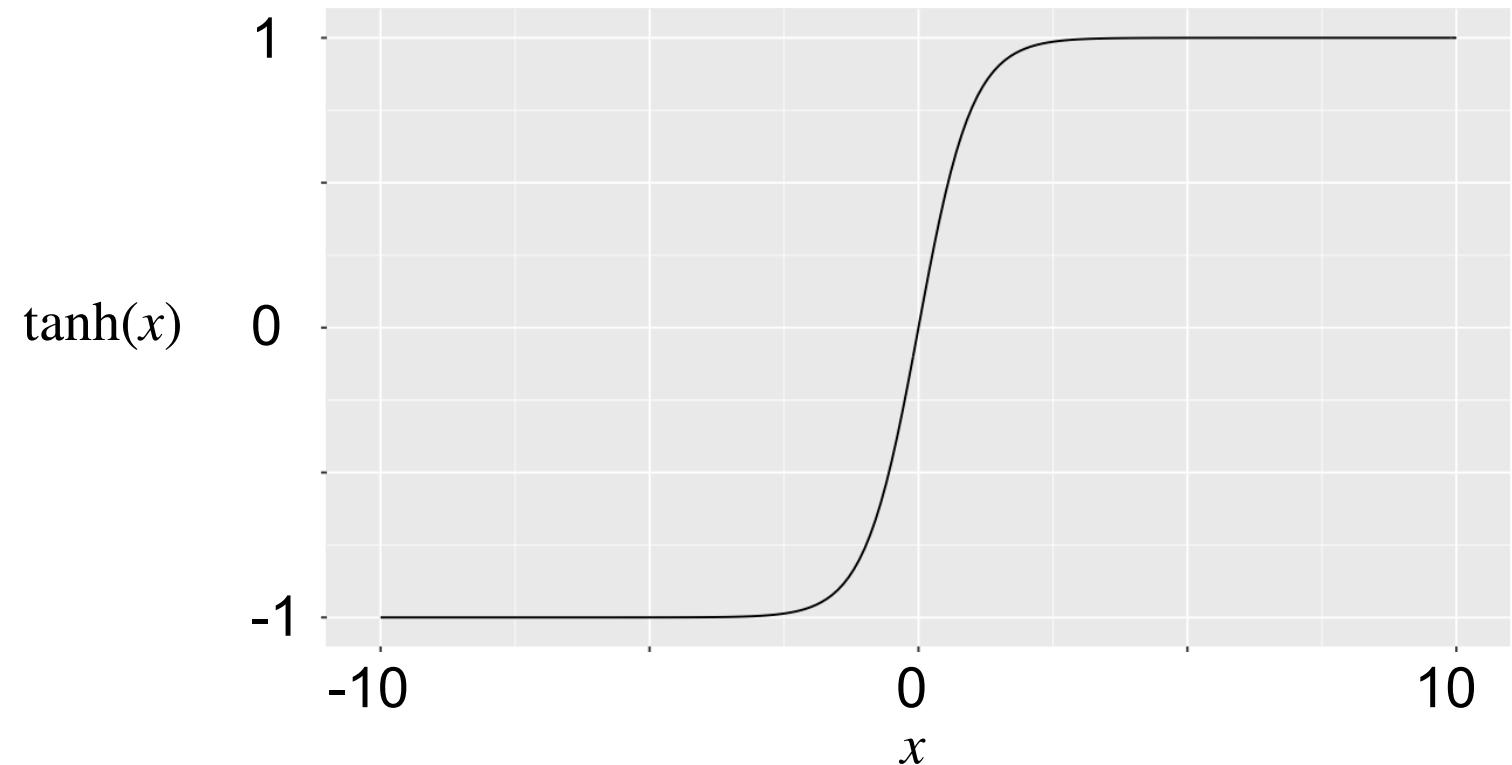
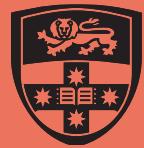


Figure from David Bamman's Lecture Slides for Info 159/259 at UC Berkeley



What is a non-linear function?

Rectified Linear Unit (ReLU)

Model

$$\text{ReLU}(x) = \max(0, x)$$

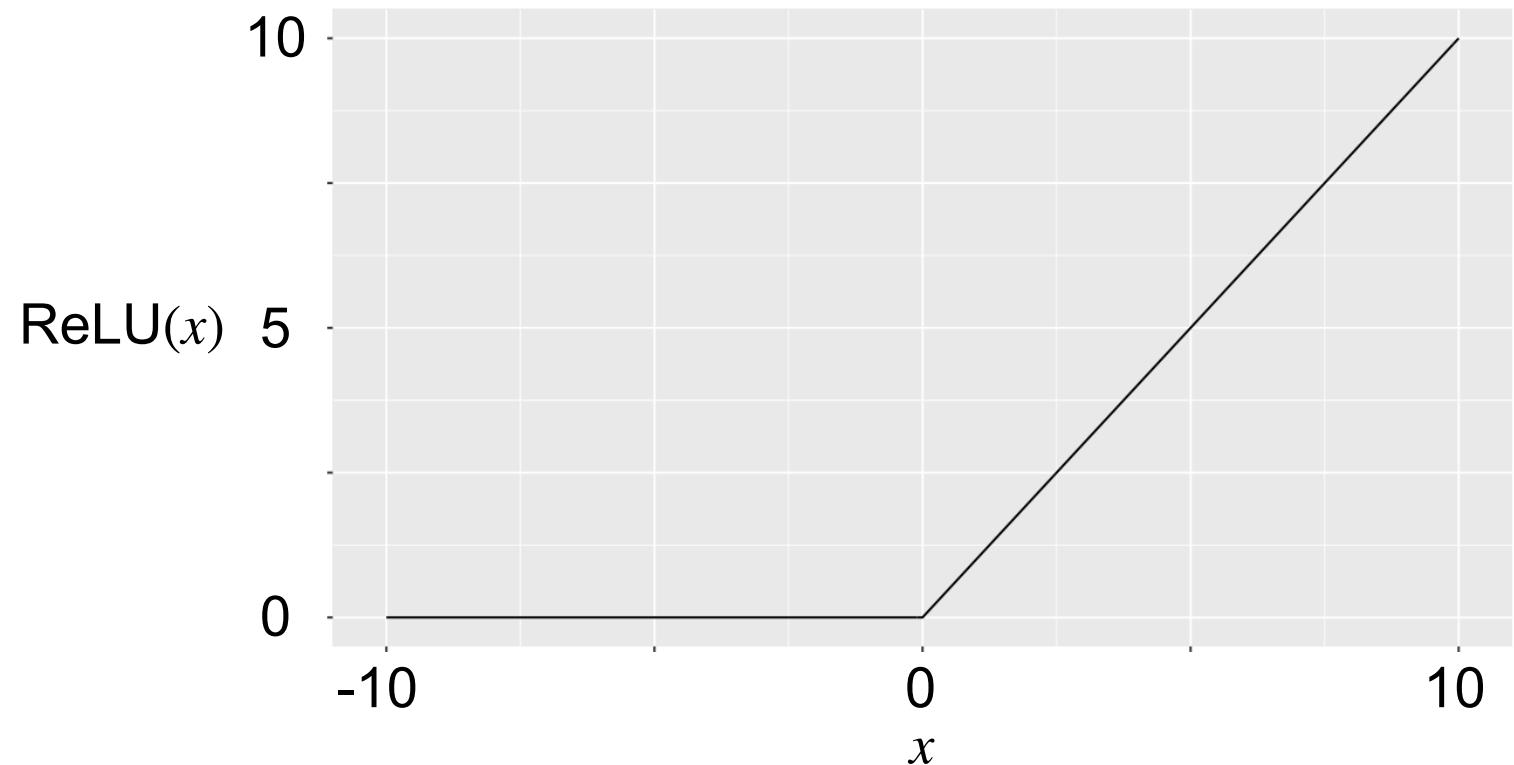
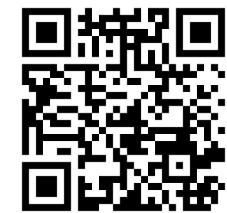
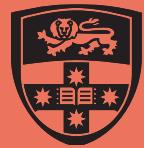


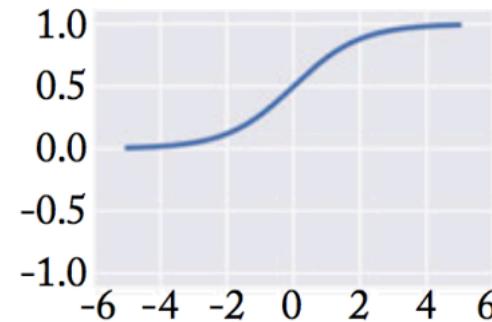
Figure from David Bamman's Lecture Slides for Info 159/259 at UC Berkeley



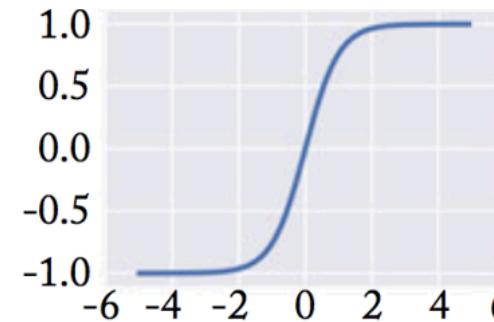
What is a non-linear function?

Model

① $\sigma(x) = \frac{1}{1 + e^{-x}}$



② $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



③ $\text{ReLU}(x) = \max(0, x)$

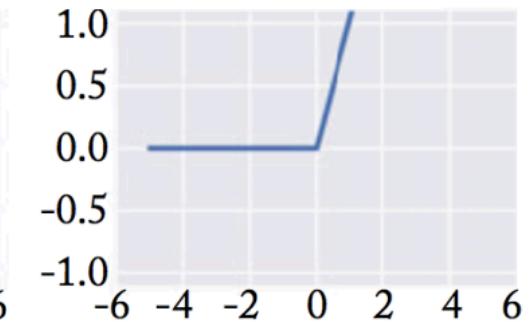
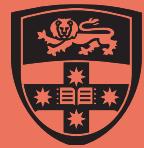
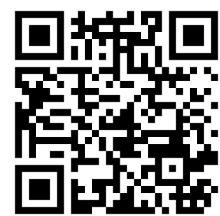


Figure from Goldberg, "Neural Network Methods for Natural Language Processing"



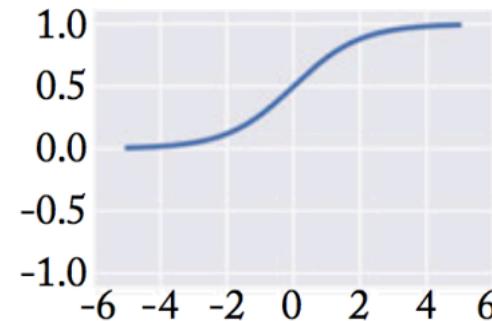
Neural Networks
Recurrent Models
Analysis
Workshop Preview



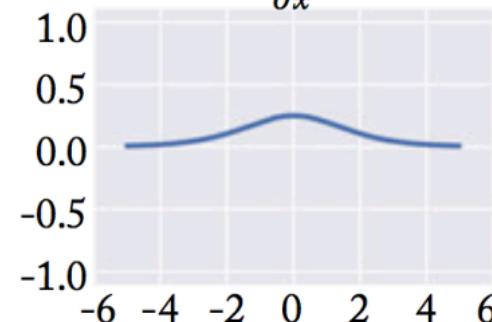
[menti.com 1750 7815](https://menti.com/17507815)

What is a non-linear function?

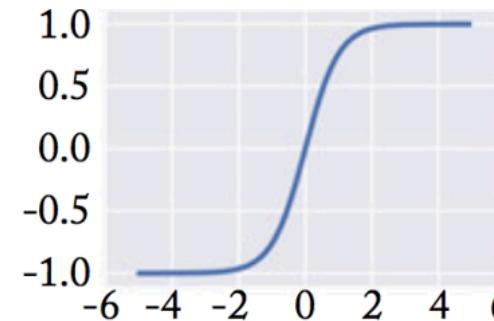
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



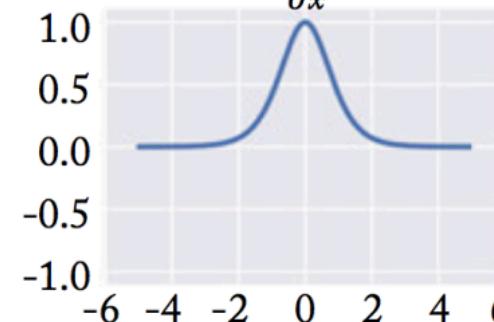
$$\frac{\partial f}{\partial x}$$



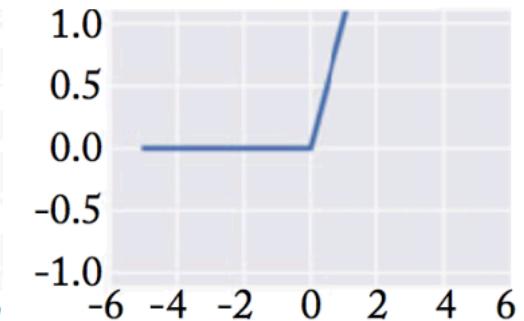
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



$$\frac{\partial f}{\partial x}$$



$$\text{ReLU}(x) = \max(0, x)$$



$$\frac{\partial f}{\partial x}$$

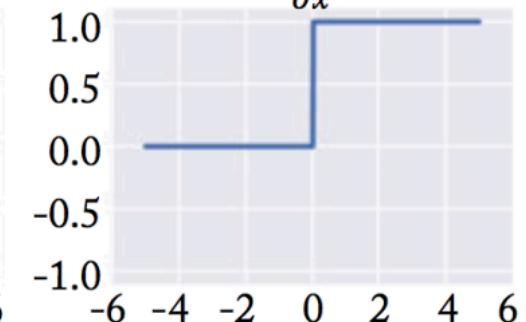
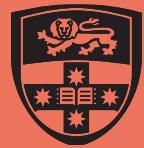


Figure from Goldberg, "Neural Network Methods for Natural Language Processing"



What is a non-linear function?

each layer has its
own weight

Model

$$\text{scores(doc)} = g(\text{features(doc)} \cdot \text{weights}_1) \cdot \text{weights}_2$$

$$= g(\begin{array}{cccccc} \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \end{array} \cdot \begin{array}{|c|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}) \cdot \begin{array}{|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}$$

$$= g(\begin{array}{cccc} \text{---} & \text{---} & \text{---} & \text{---} \end{array}) \cdot \begin{array}{|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}$$

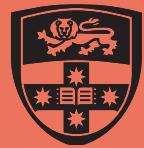
If $g(x) = \text{ReLU}(x) = \max(0, x)$

$$= \begin{array}{cccc} \text{---} & \text{---} & \text{---} & \text{---} \end{array} \cdot \begin{array}{|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}$$

If $g(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

$$= \begin{array}{cccc} \text{---} & \text{---} & \text{---} & \text{---} \end{array} \cdot \begin{array}{|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}$$

Let's see non linear layers



What is a non-linear function?

Model

$$f(x) = g(x \cdot W_1) \cdot W_2 \quad \leftarrow 1 \text{ layers}$$

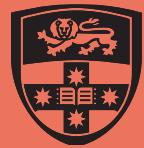
$$f(x) = h(g(x \cdot W_1) \cdot W_2) \cdot W_3 \quad \leftarrow 2 \text{ layers}$$

Each stage of non-linearity + weight multiplication is a layer

We typically also have a bias (fixed value that is added on):

$$f(x) = h(g(x \cdot W_1 + b_1) \cdot W_2 + b_2) \cdot W_3 + b_3$$

Can also think of this as an extra input that is always 1



[menti.com 1750 7815](https://menti.com/17507815)

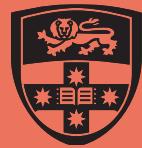
What is a non-linear function?

Model

So far: all feed-forward neural networks

Also called multi-layer perceptrons (MLP)

We'll see other types in the weeks ahead



How do we update the models? First consider linear

score(doc) = features(doc) · weights

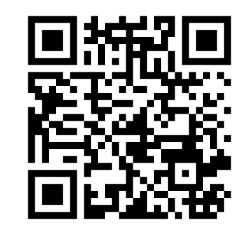
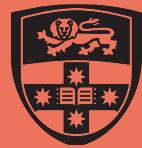
$$= \begin{array}{ccccccccc} \text{---} & \text{---} \end{array} \cdot$$



$$f(x) = x \cdot W \leftarrow \text{linear model}$$

Week 2: Perceptron


$$W_i = \begin{cases} W_i, & \text{if correct} \\ W_i - \text{sign}(x_i), & \text{if incorrect} \end{cases}$$



通过 Loss Function iteration

How do we update the models? First consider linear

score(doc) = features(doc) · weights

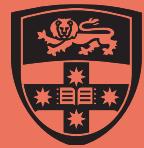
=  ·

$f(x) = x \cdot W$

Learning Method

Week 3: Move based on the derivative of loss / cost

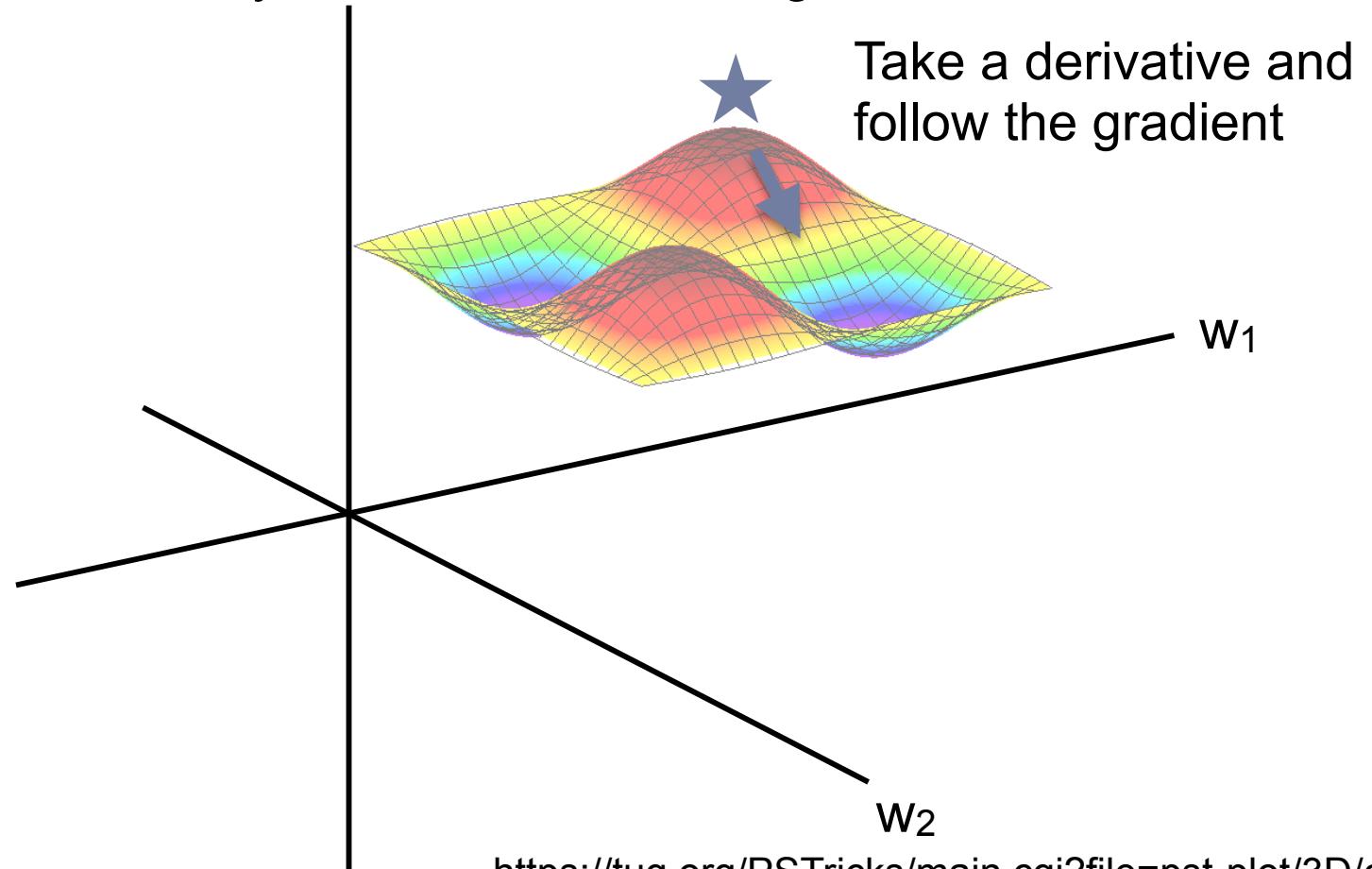
- Loss is a mathematical function measuring how wrong we are
- Derivative is the direction we should move to increase the function
- Move in the negative derivative to get better!

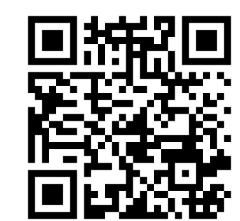


How do we update the models? First consider linear

Learning
Method

$y = \text{loss} / \text{cost}$, ie., how
many errors we are making





How do we update the models? First consider linear

score(doc) = features(doc) · weights

$$= \begin{array}{c} \text{[color-coded feature vector]} \\ \cdot \end{array}$$

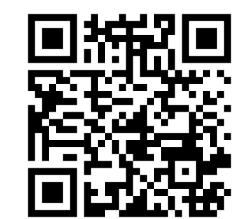


$$f(x) = x \cdot W$$

Week 3: Move based on the derivative of loss

$$\begin{aligned} \boxed{\text{loss}} &= -\log(p(y|x)) \\ &= -[y \log \hat{y} + (1-y) \log(1-\hat{y})] \\ &= -[y \log \sigma(x \cdot W) + (1-y) \log(1 - \sigma(x \cdot W))] \\ &= \begin{cases} -\log \sigma(x \cdot W), & \text{if } y = 1 \\ \log(1 - \sigma(x \cdot W)), & \text{if } y = 0 \end{cases} \end{aligned}$$

Sigmoid function



How do we update the models? First consider linear

score(doc) = features(doc) · weights

$$= \begin{array}{c} \text{[color-coded feature vector]} \\ \cdot \\ \text{[color-coded weight vector]} \end{array}$$

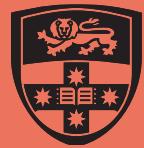
$$f(x) = x \cdot W$$

Week 3: Move b function from earlier, which ss

$$\begin{aligned} loss &= -\log(f) \\ &= -[y \log f + (1-y) \log(1-f)] \\ &= -[y \log \sigma(x \cdot W) + (1-y) \log(1 - \sigma(x \cdot W))] \end{aligned}$$

This is the sigmoid
function from earlier, which ss
conveniently maps values
to between 0 and 1

$$\frac{d \text{loss}}{dW_j} = (\sigma(x \cdot W) - y)x_j$$



How do we update the models? Now non-linear

→ backpropagation

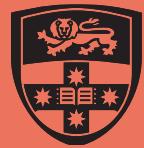
$$\text{scores(doc)} = g(\text{features(doc)} \cdot \text{weights}_1) \cdot \text{weights}_2$$

$$= g(\begin{array}{cccccc} \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \end{array} \cdot \begin{array}{|c|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}) \cdot \begin{array}{|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}$$

$$f(x) = g(x \cdot W_1) \cdot W_2$$

Just another function!

Learning Method



How can we avoid doing so much differentiation?

$$\text{scores(doc)} = g(\text{features(doc)} \cdot \text{weights}_1) \cdot \text{weights}_2$$

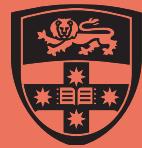
$$= g(\begin{array}{cccccc} \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \end{array} \cdot \begin{array}{|c|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}) \cdot \begin{array}{|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \end{array}$$

$$f(x) = g(x \cdot W_1) \cdot W_2$$

Need derivatives for W_1 and W_2 and more as we add layers

$$f(x) = h(g(x \cdot W_1) \cdot W_2) \cdot W_3$$

Chain rule!

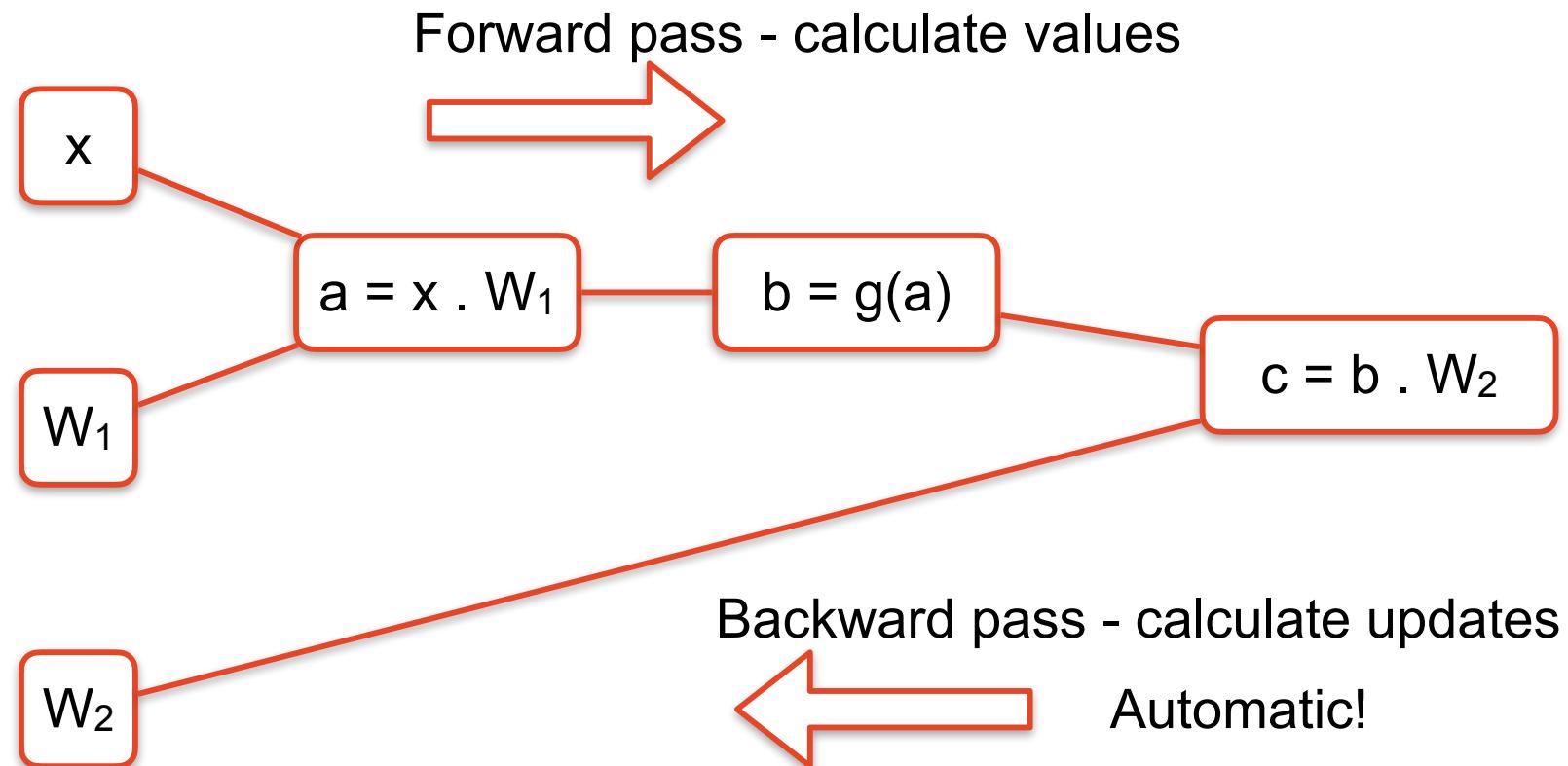


Computation graphs and backpropagation

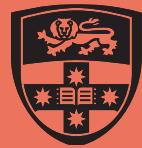
Learning
Method

$$f(x) = g(x \cdot W_1) \cdot W_2$$

The equation above can be represented as a graph



This is one loop



Neural Networks
Recurrent Models
Analysis
Workshop Preview



Computation graphs and backpropagation

Chain rule

$$f(x) = h(g(x))$$

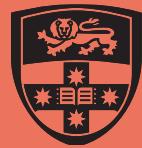
$$\frac{df}{dx} = \frac{dh}{dg} \cdot \frac{dg}{dx}$$

$$f(x) = g_3(g_2(g_1(x)))$$

$$\frac{df}{dx} = \frac{dg_3}{dg_2} \cdot \frac{dg_2}{dg_1} \cdot \frac{dg_1}{dx}$$

$$\frac{d}{dx} [f(g(x))] = f'(g(x)) g'(x) x'$$

[menti.com 1750 7815](https://menti.com/17507815)



Neural Networks
Recurrent Models
Analysis
Workshop Preview

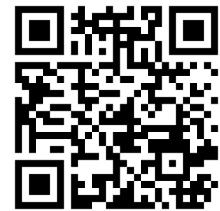
Computation graphs and backpropagation

There is some randomness to it because it depends on the order we see data in

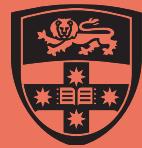
We go downwards in loss space

Stochastic gradient descent

We differentiate to get the gradient



[menti.com 1750 7815](https://menti.com/17507815)



COMP 4446 / 5046
Lecture 3, 2025

Neural Networks
Recurrent Models
Analysis
Workshop Preview



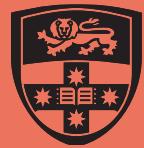
[menti.com 1750 7815](https://menti.com/17507815)

Computation graphs and backpropagation

Learning
Method

Libraries do the calculations for us:



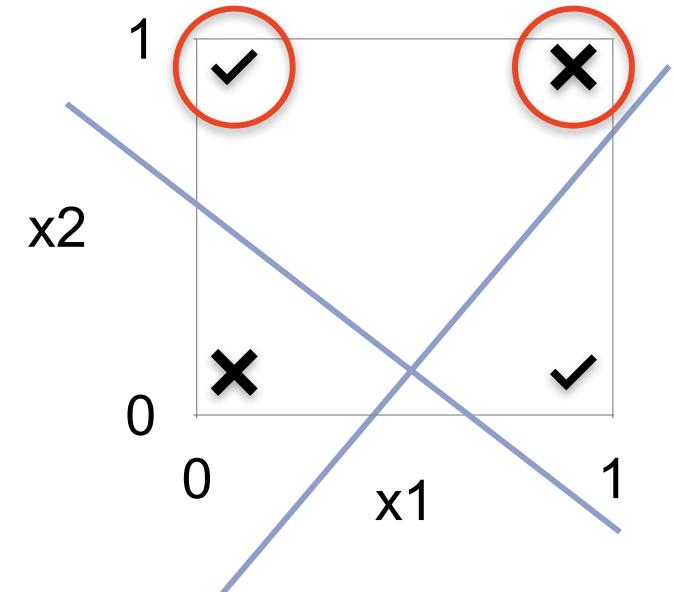


What can we learn with this?

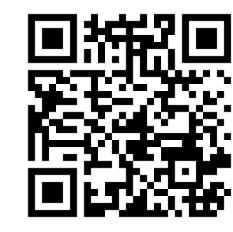
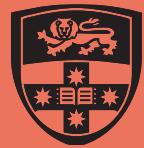
Model

XOR: Two inputs, x_1 and x_2

		x_1	
		0 1	
x_2	0	✗	✓
	1	✓	✗



cannot use linear to represent it

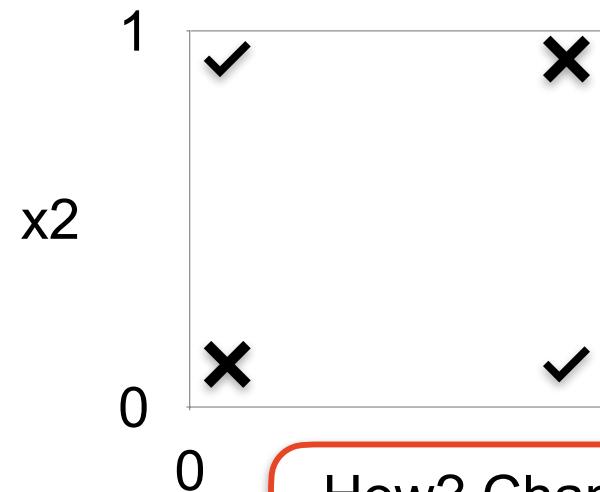


What can we learn with this?

Model

XOR: Two inputs, x_1 and x_2

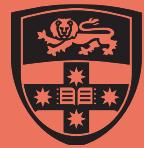
		x_1	
		0 1	
x_2	0	\times	\checkmark
	1	\checkmark	\times



How? Change representation / space

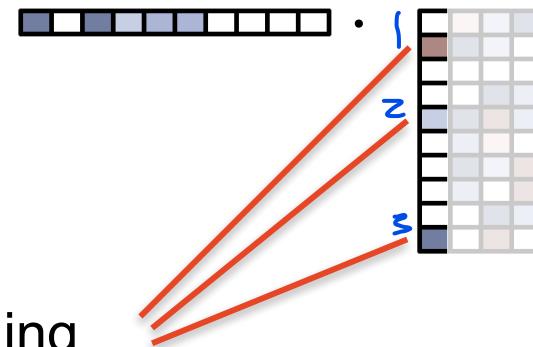
$$\text{ReLU}(x_1 + x_2) - 2 \quad (\text{ReLU}(x_1 + x_2 - 1))$$

$$\text{ReLU} \left[\begin{array}{cc} x & W_1 \\ \square & \blacksquare \end{array}, \begin{array}{cc} x & W_2 \\ (\square & \blacksquare) - 1 \end{array} \right] W_3 \quad \blacksquare$$



What can we learn with this?

Model

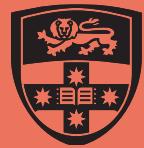


Combining
these 3 inputs

e.g., are 'very', and 'good' present, but not 'bad'?

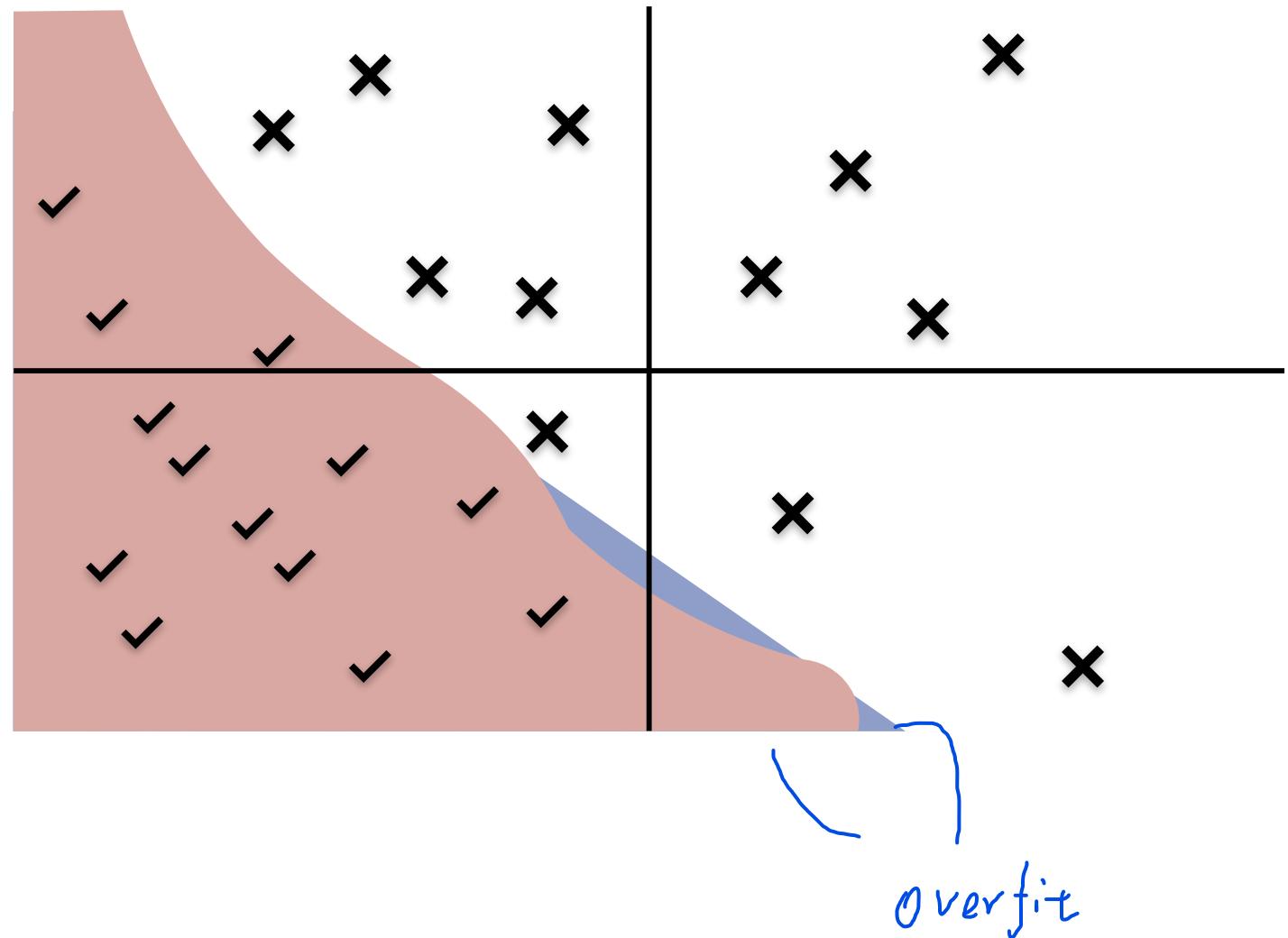
positive

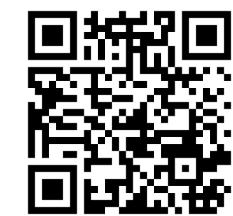
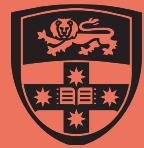
prediction



Model

Overfitting risk





Overfitting risk

Model

Original loss

$$-\log(p(y|x))$$

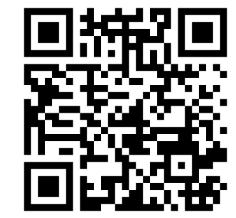
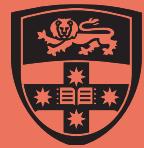
Loss with regularization

$$-\log(p(y|x)) + \text{regularizer}(W)$$

An extra cost that depends on the weights

Effect - push learning to favour smaller weights or more zero weights

越少 weight
更容易被 penalize

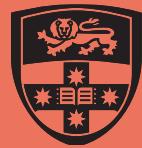


[menti.com 1750 7815](https://menti.com/17507815)

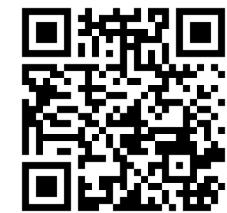
Methods of regularization

Learning
Method

- (1) Modify loss with cost, e.g, ℓ_1 , ℓ_2
- (2) Dropout, during training, randomly set some weights to zero
- (3) Early stopping, don't backpropagate error further once it gets small

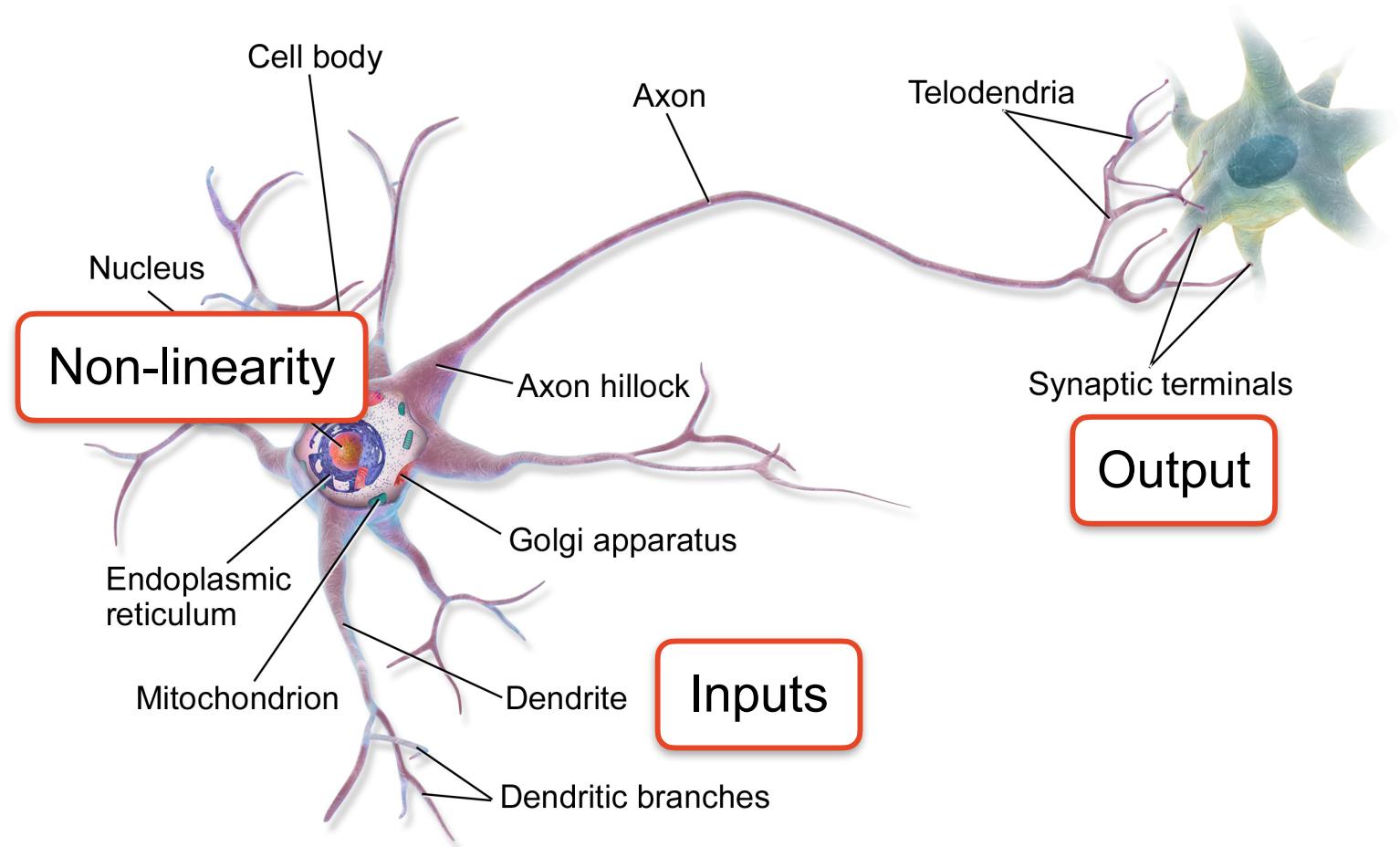


Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

Why are they called ‘neural networks’?



<https://commons.wikimedia.org/w/index.php?curid=28761830>



SUBSTITUTIONS

THAT MAKE READING THE NEWS MORE FUN:

WITNESSES	→ THESE DUDES I KNOW
ALLEGEDLY	→ KINDA PROBABLY
NEW STUDY	→ TUMBLR POST
REBUILD	→ AVENGE
SPACE	→ SPAACE
GOOGLE GLASS	→ VIRTUAL BOY
SMARTPHONE	→ POKÉDEX
ELECTRIC	→ ATOMIC
SENATOR	→ ELF-LORD
CAR	→ CAT
ELECTION	→ EATING CONTEST
CONGRESSIONAL LEADERS	→ RIVER SPIRITS
HOMELAND SECURITY	→ HOMESTAR RUNNER
COULD NOT BE REACHED FOR COMMENT	→ IS GUILTY AND EVERYONE KNOWS IT

Substitutions

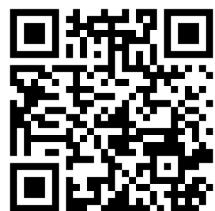
[INSIDE ELON MUSK'S NEW
ATOMIC CAT]

Source: <https://xkcd.com/1288/>



COMP 4446 / 5046
Lecture 3, 2025

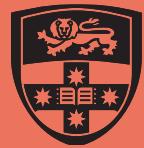
Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

Recurrent Models

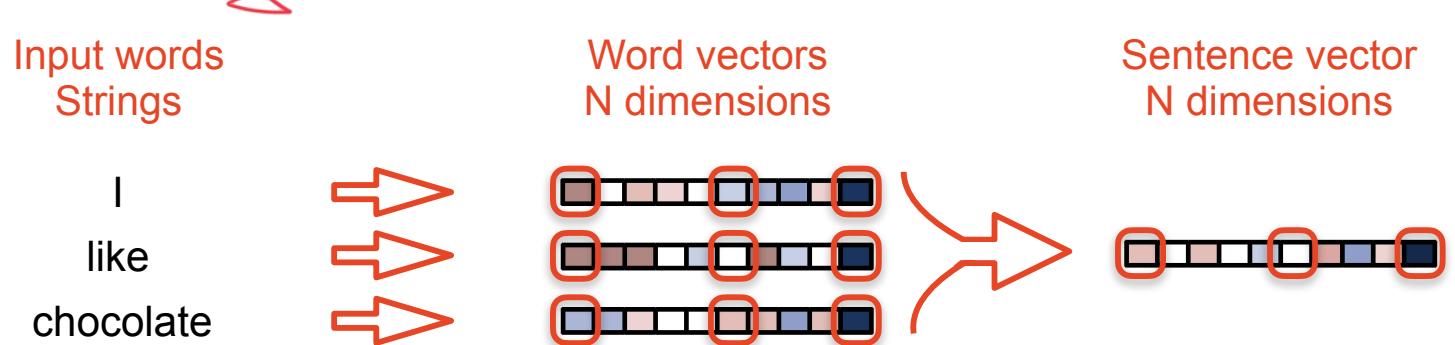
Structure and concepts drawn from CS288 at Berkeley,
which cites UT Austin and Stanford courses as sources



How to handle variable length input

Option 1

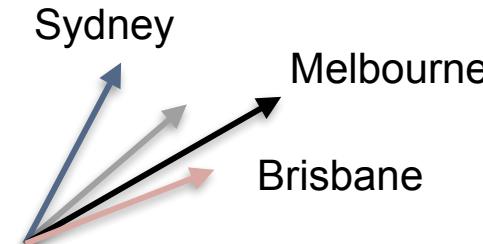
- Look up word vectors
- Take the mean



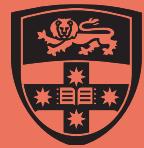
Problem: word properties are washed out / diluted

More broadly - why would this work?

Averaging city vectors will give something city-like in the middle of them



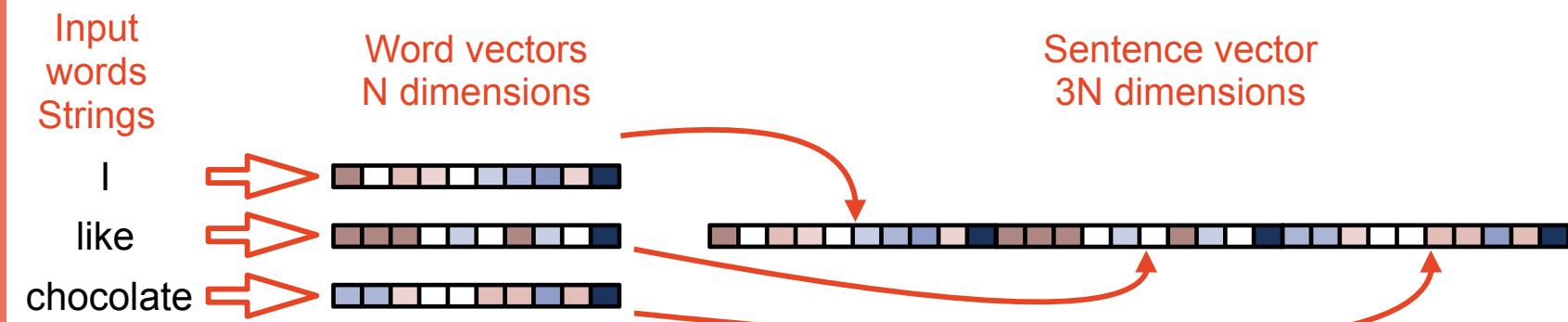
What does averaging all different words do?
I
chocolate
???
like

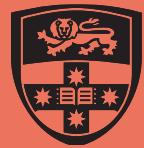


How to handle variable length input

Option 2

- Look up word vectors
- Concatenate
- Pad and/or Truncate

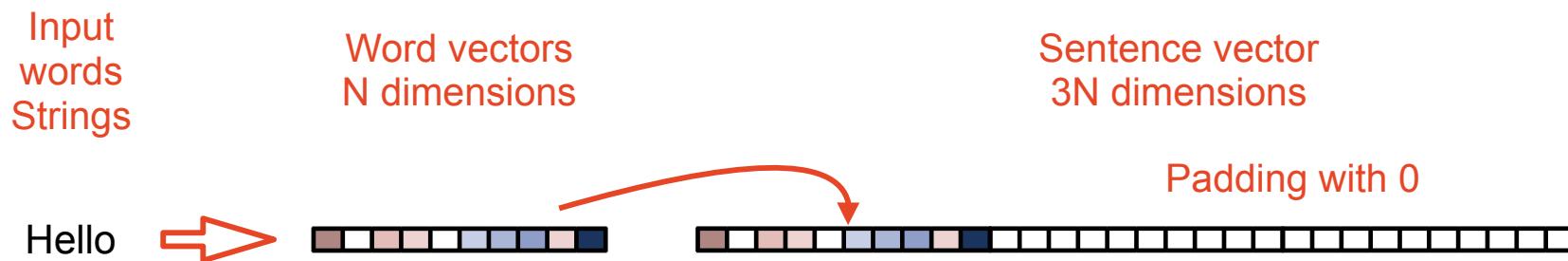


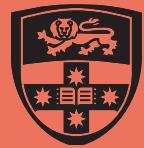


How to handle variable length input

Option 2

- Look up word vectors
- Concatenate
- **Pad and/or Truncate**

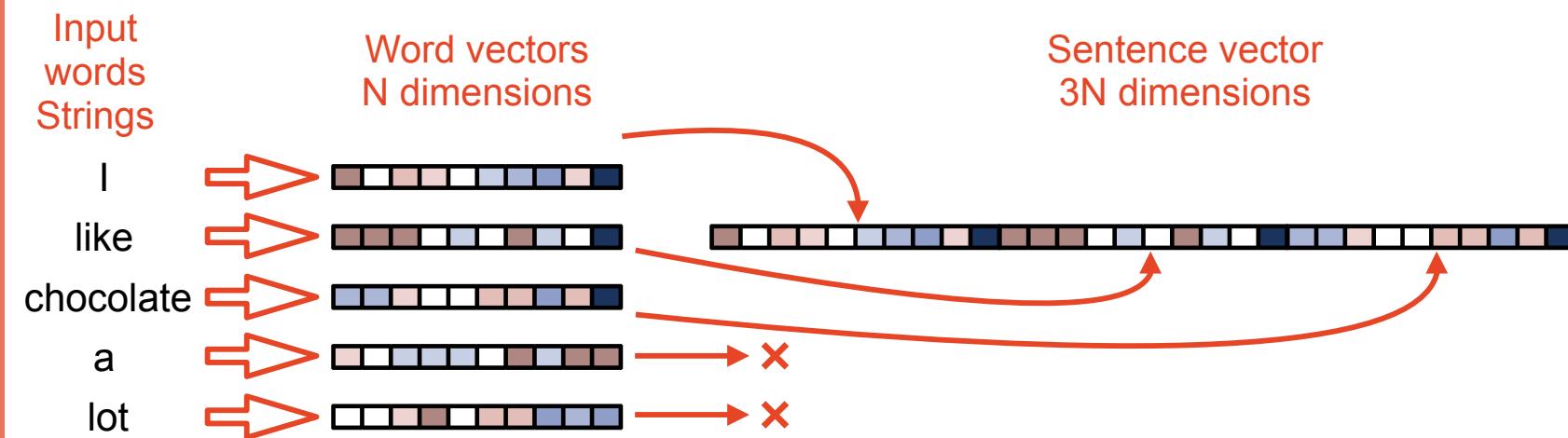


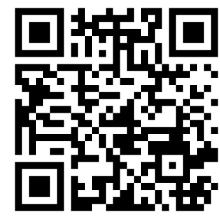
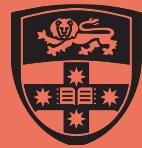


How to handle variable length input

Option 2

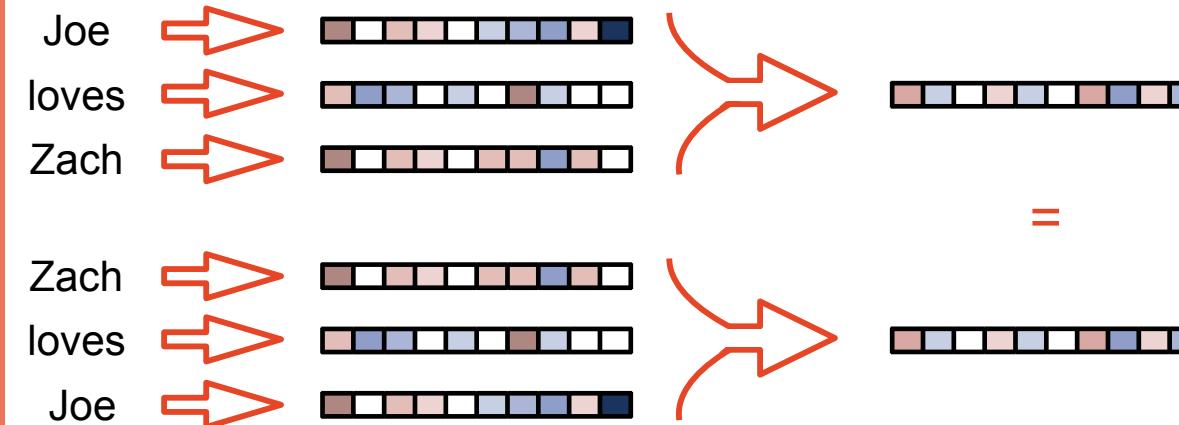
- Look up word vectors
- Concatenate
- Pad and/or **Truncate**





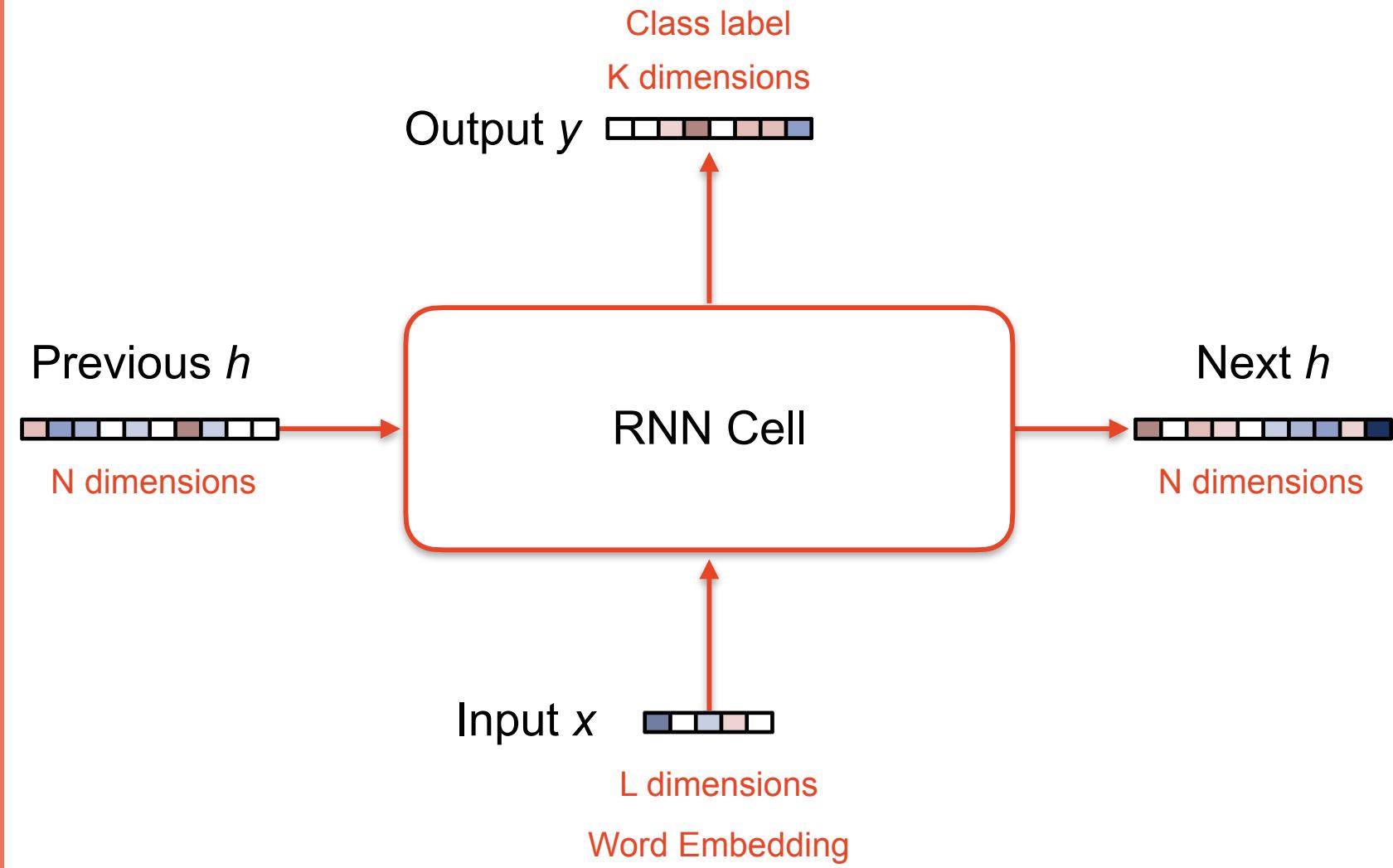
Issues with these two approaches

→ RNN 及其解决

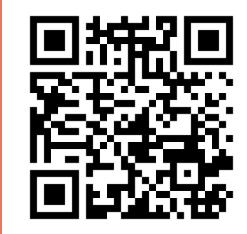




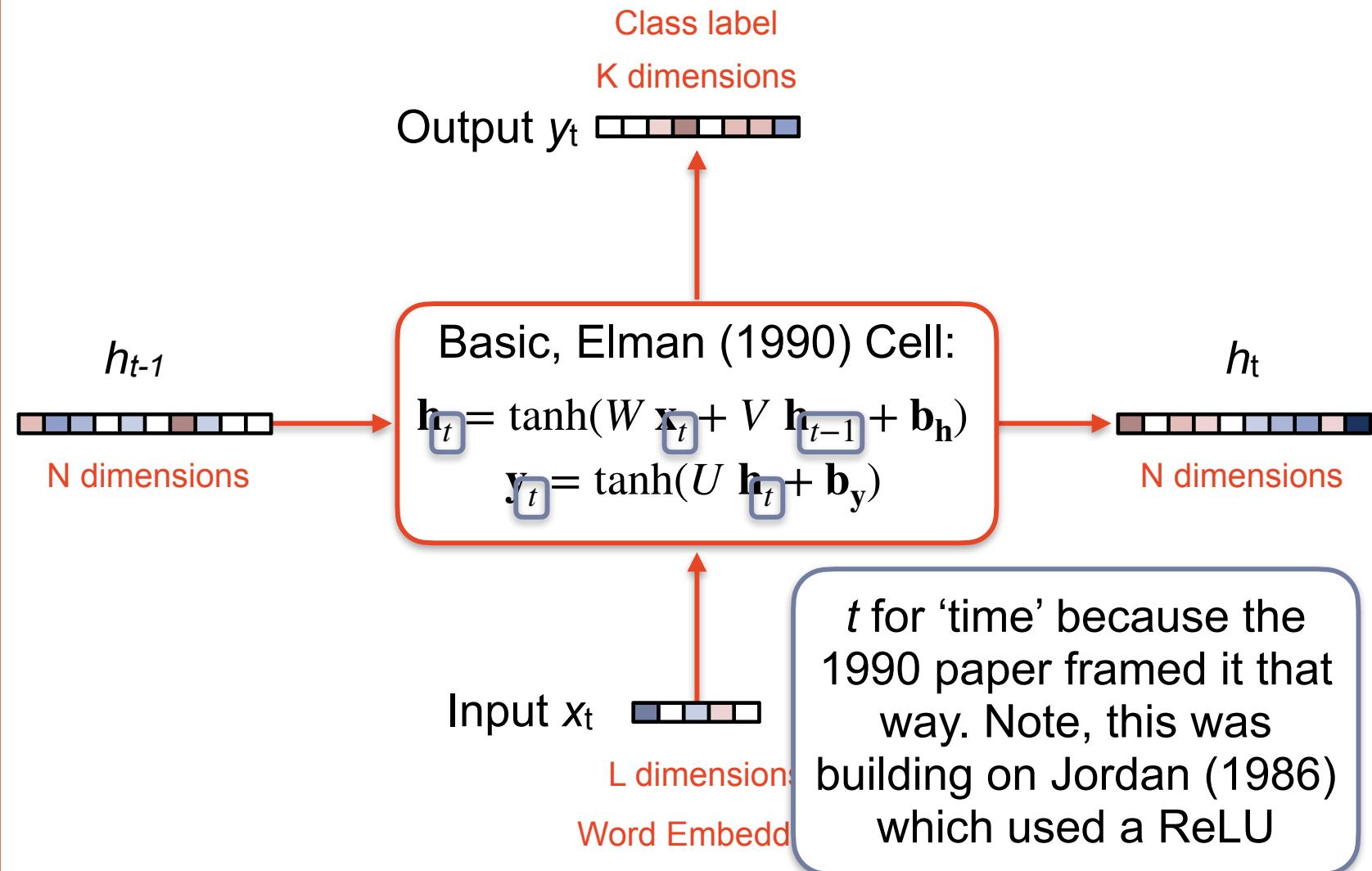
Recurrent Neural Networks (RNNs) address this issue with context-dependent representations



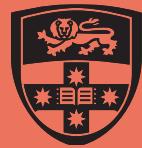
Simple RNN



Recurrent Neural Networks (RNNs) address this issue with context-dependent representations



Simple RNN

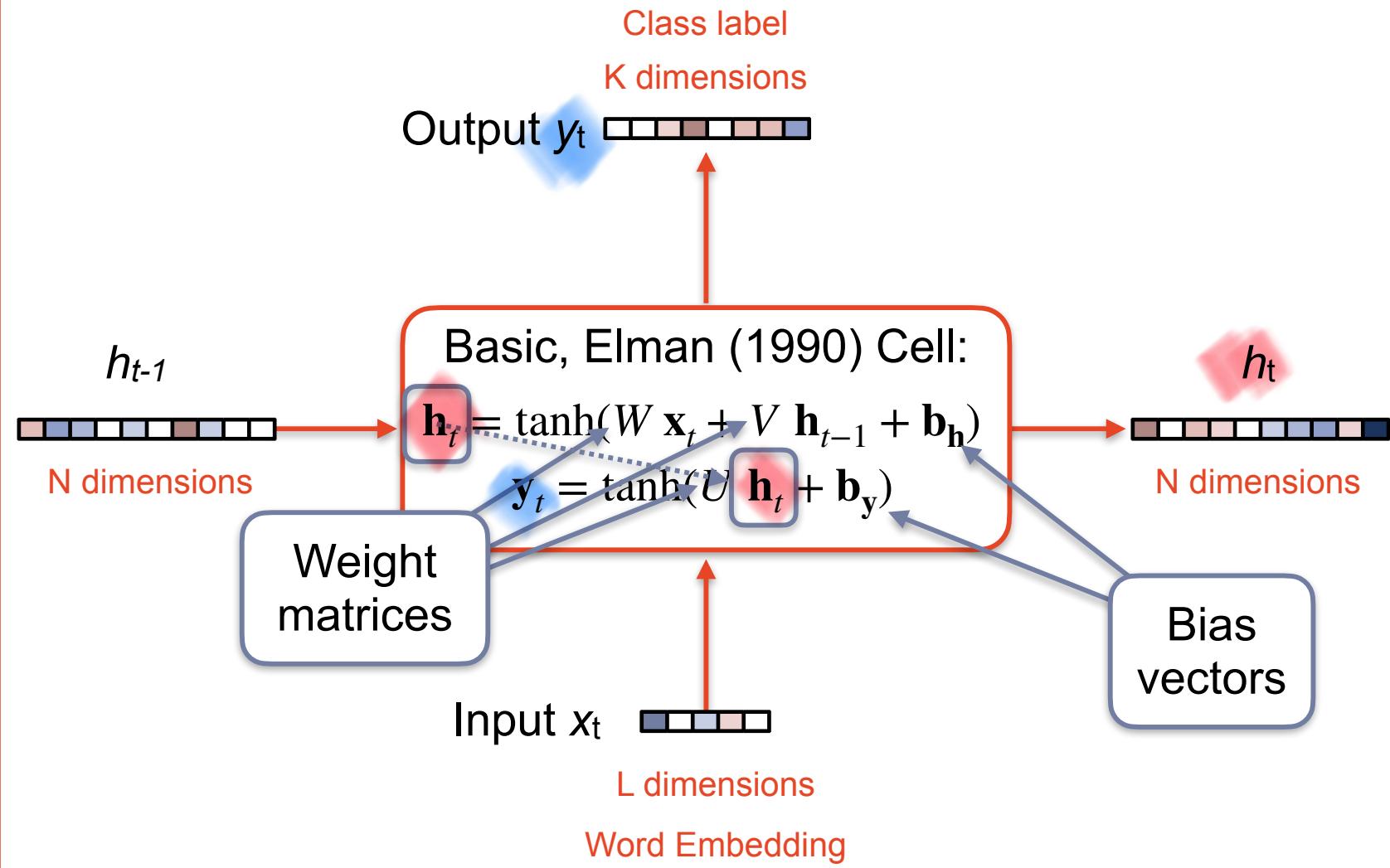


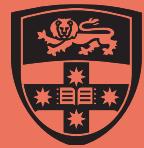
Neural Networks
Recurrent Models
Analysis
Workshop Preview



menti.com 1750 7815

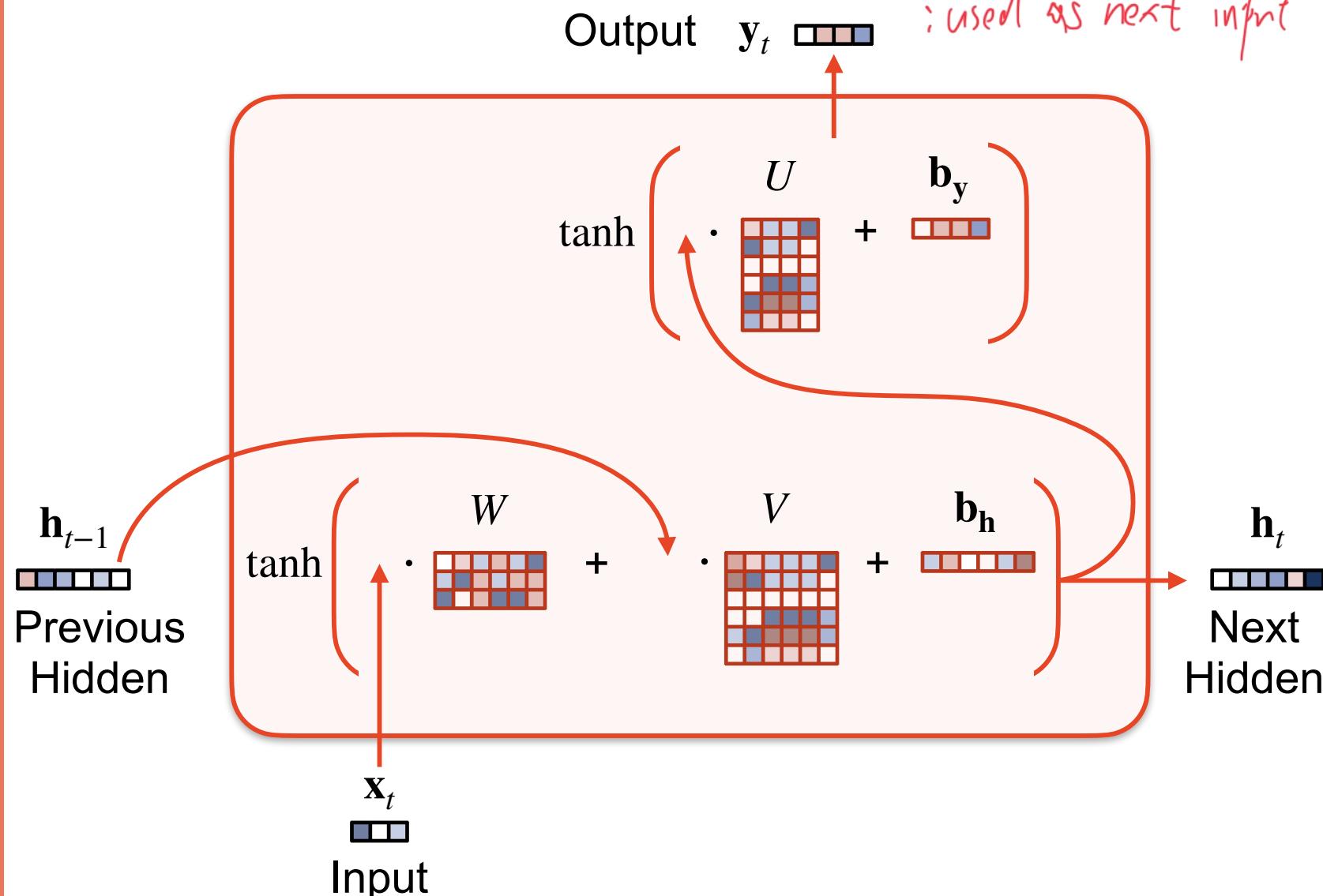
Recurrent Neural Networks (RNNs) address this issue with context-dependent representations

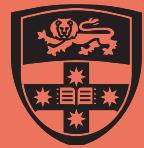




Simple RNN

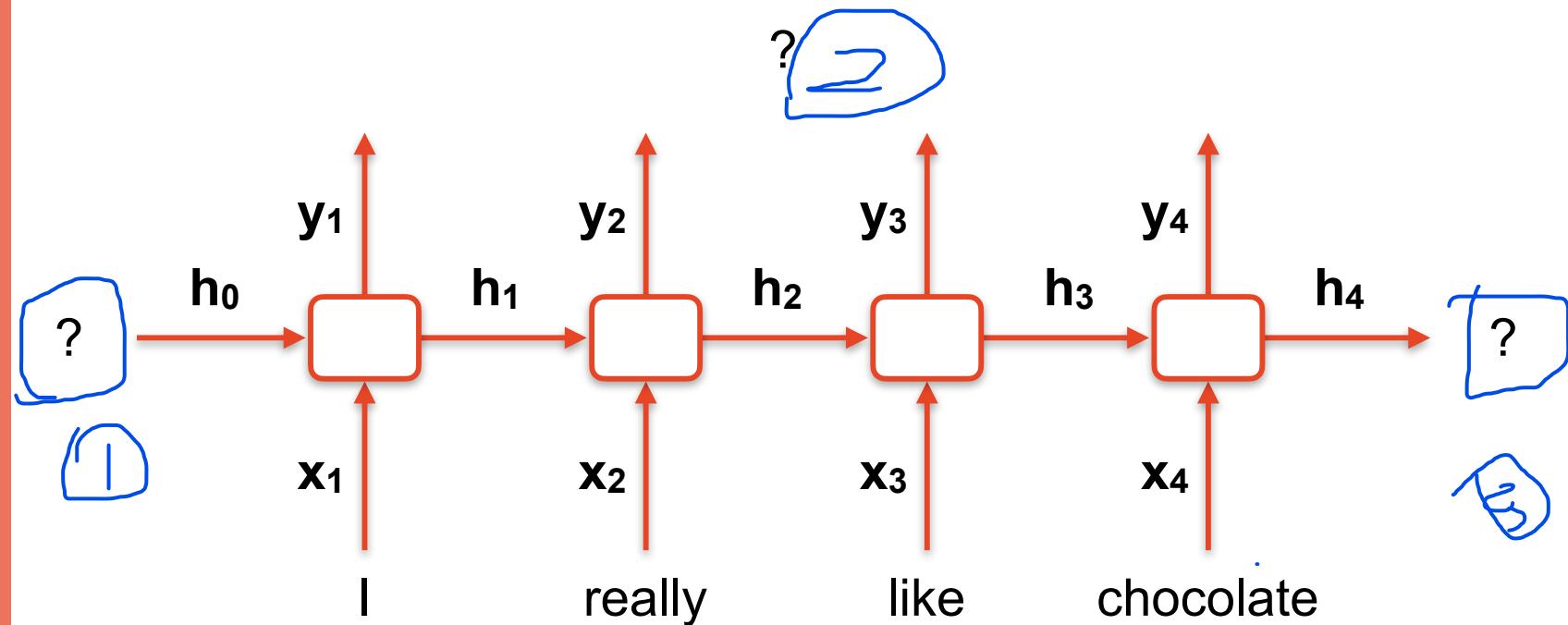
We can visualise these cells as a series of calculations

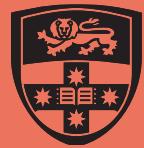




By putting a sequence of cells together, we process a multi-word input

— 一个输入词 Predict



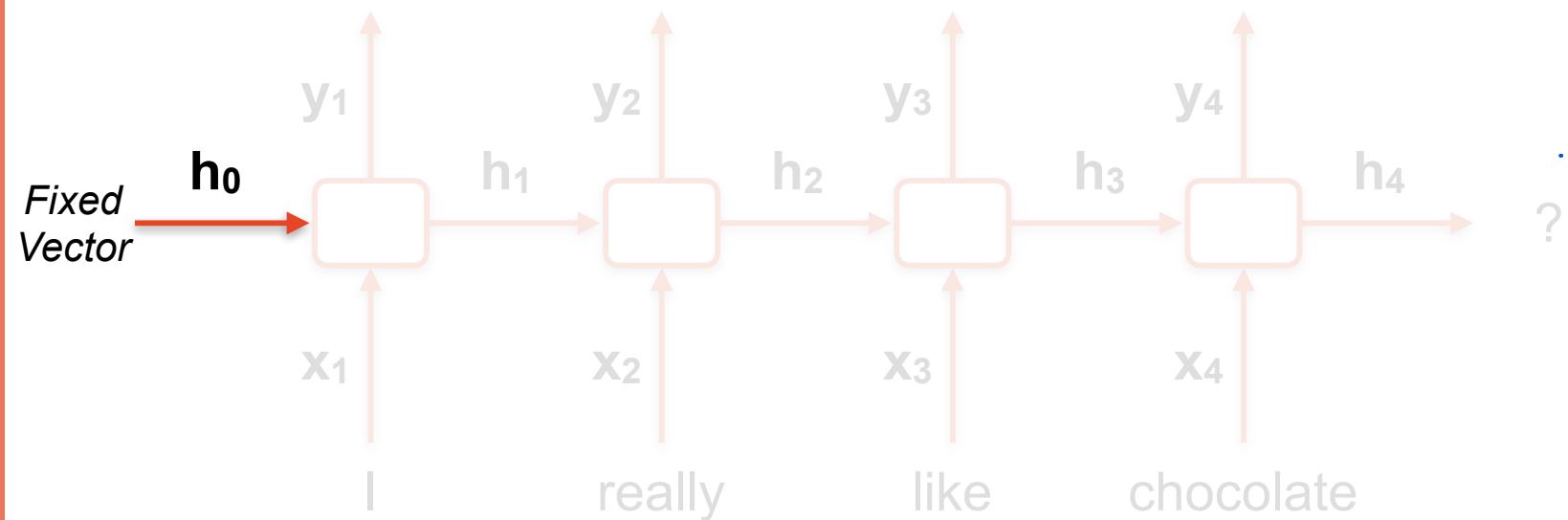


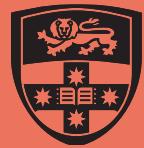
To initialise the sequence, use a fixed vector

Option 1: All zero



Option 2: Learn a set of weights



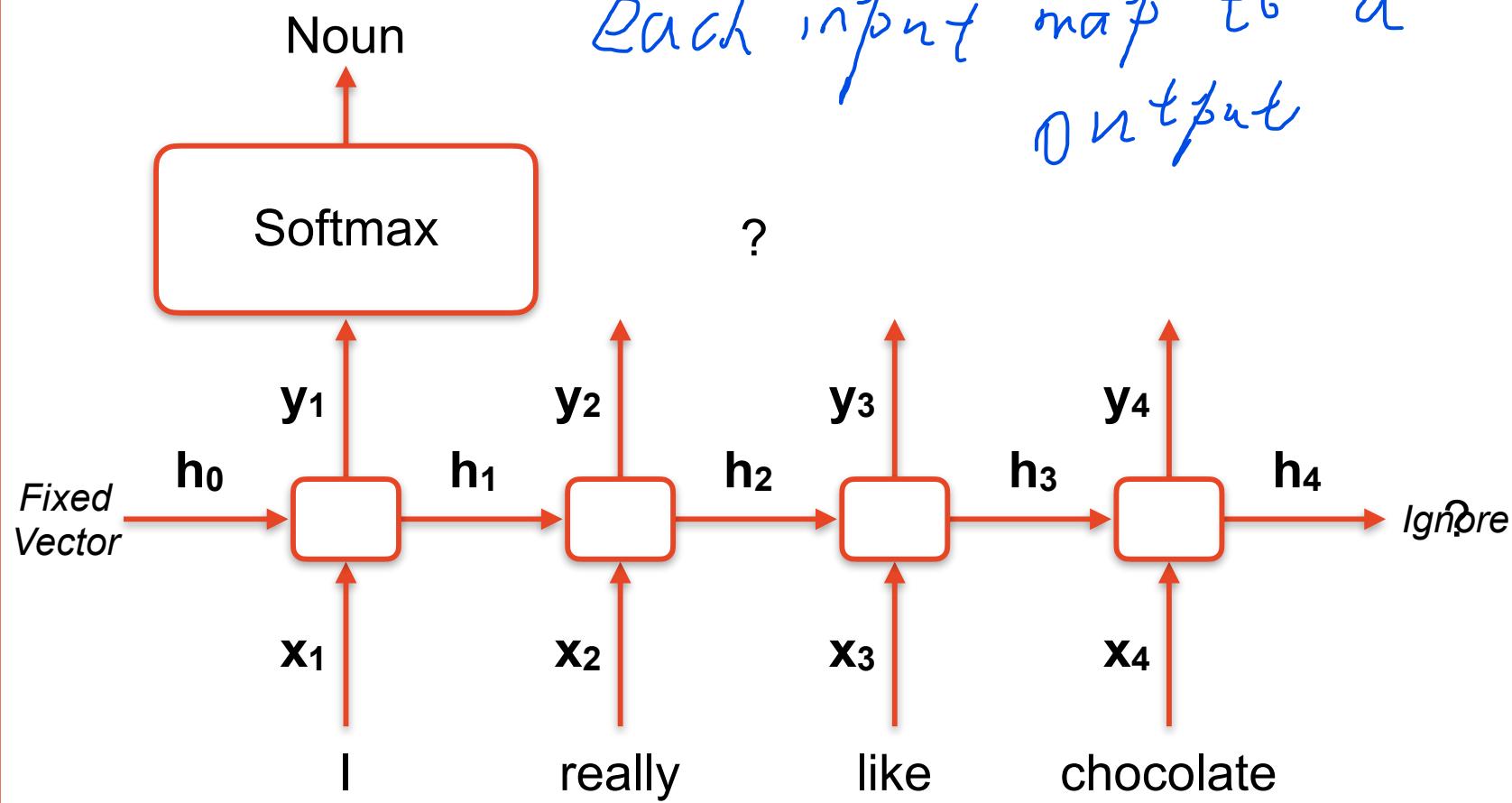


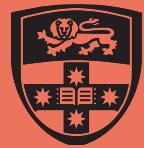
2

multiple output

One way to use them is as a **transducer**

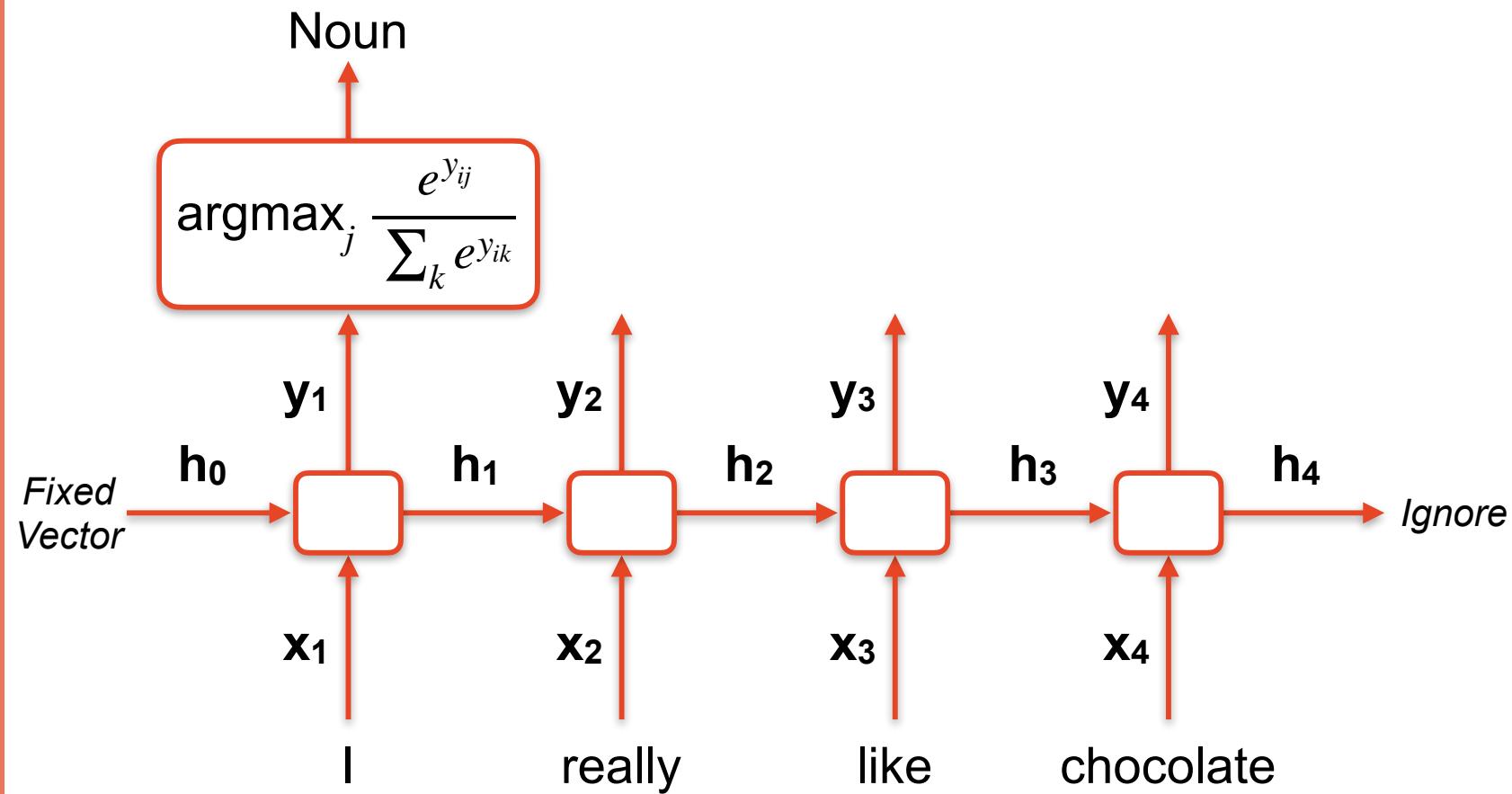
Each input map to a output

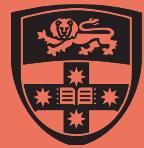




2

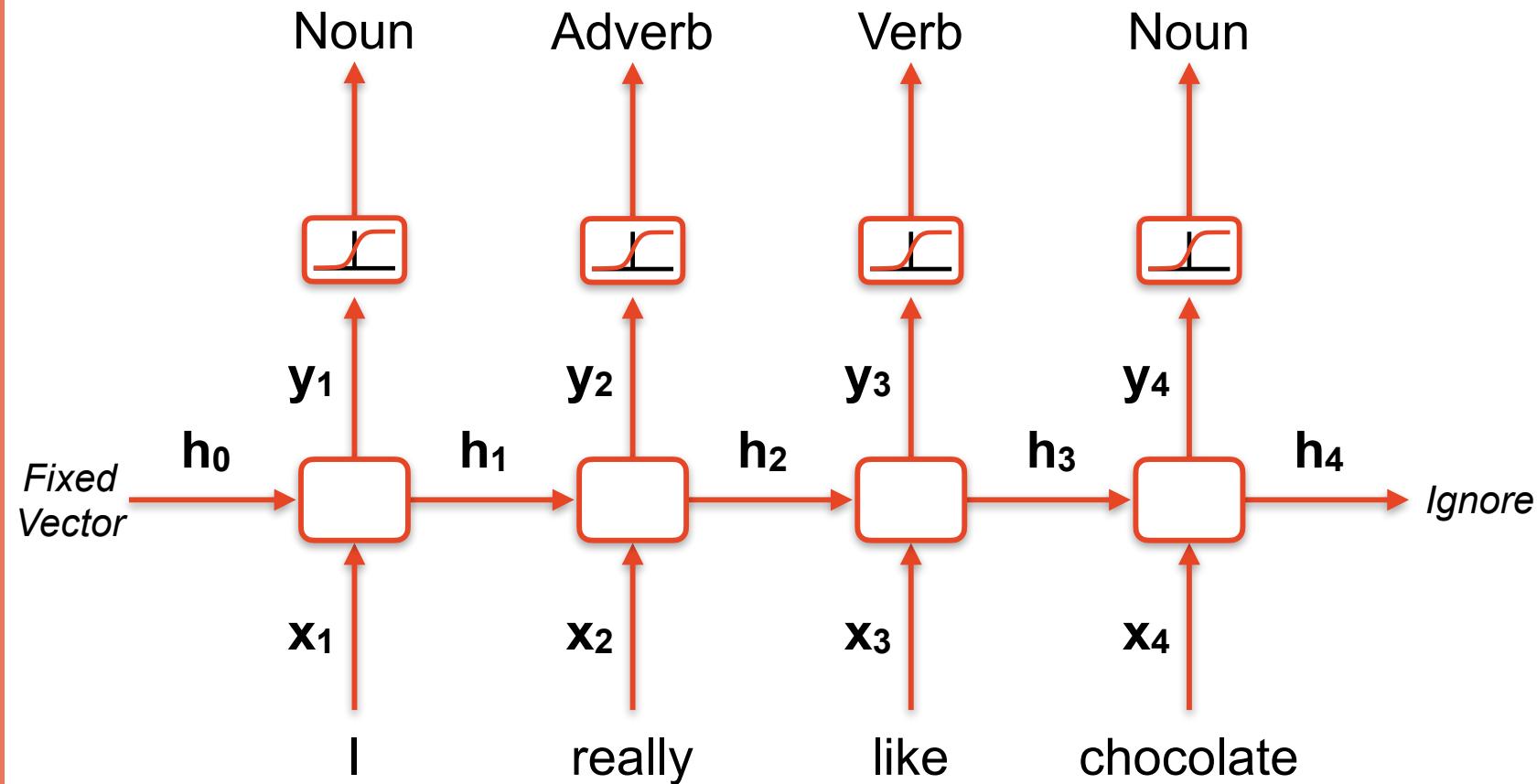
One way to use them is as a **transducer**

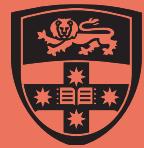




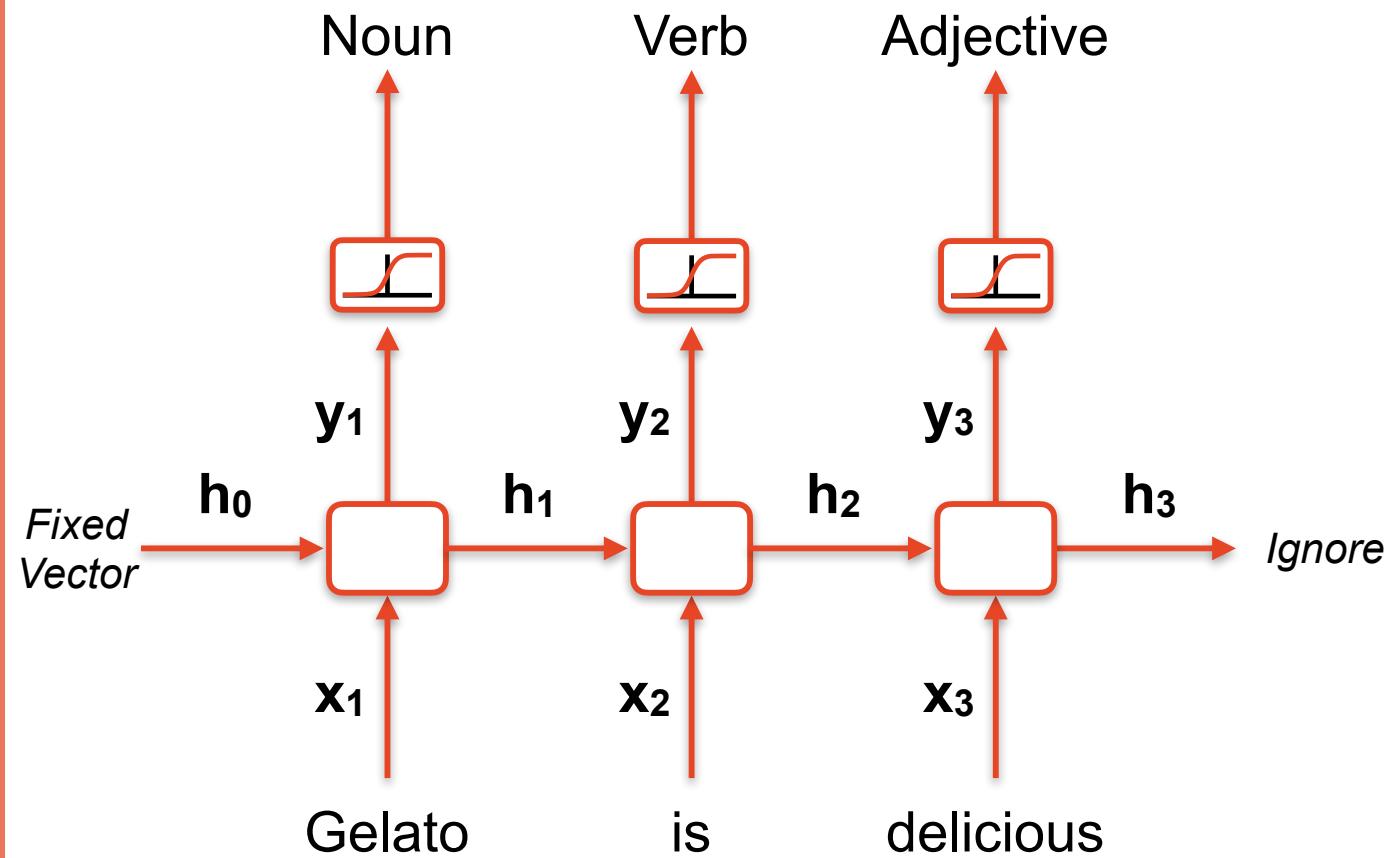
2

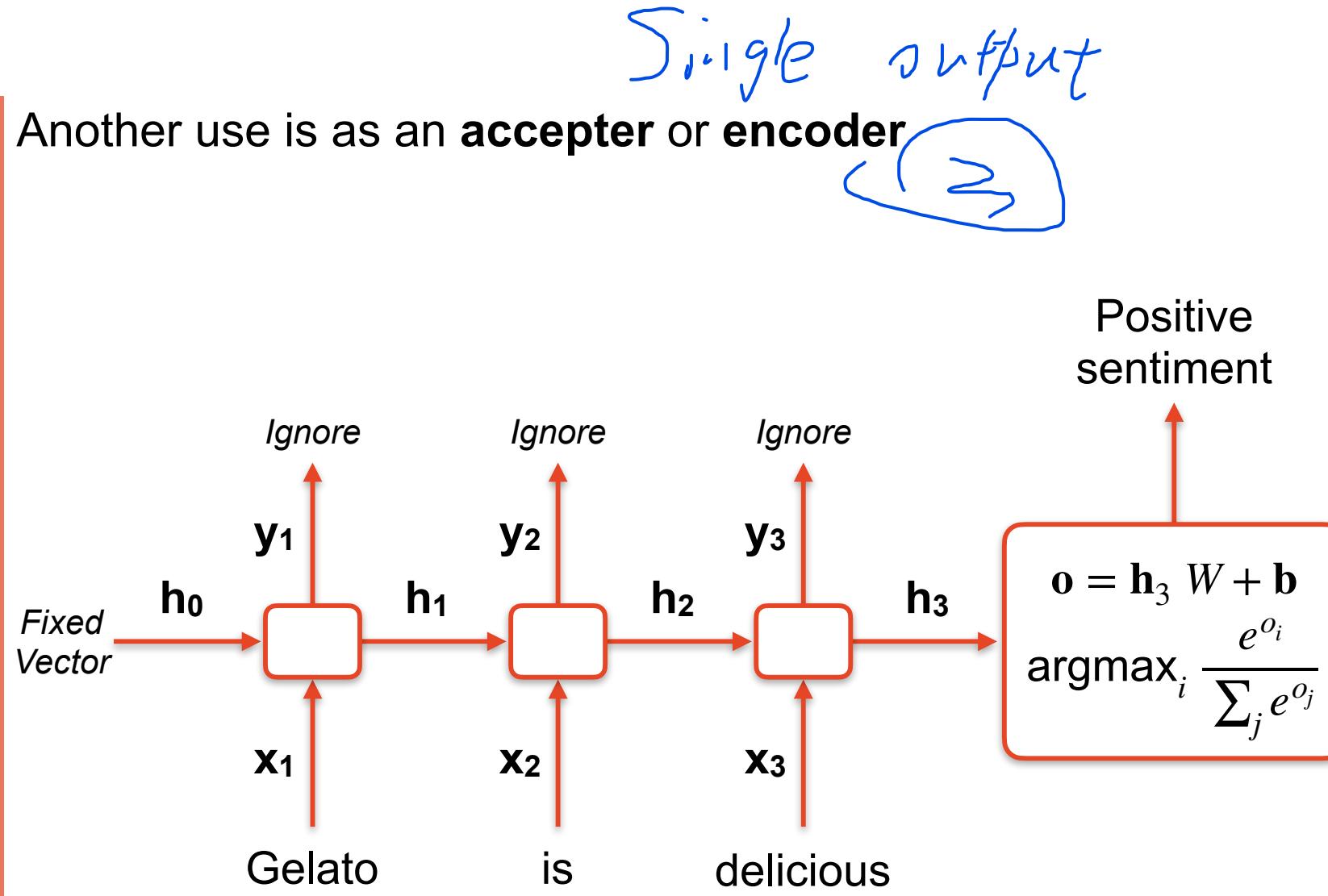
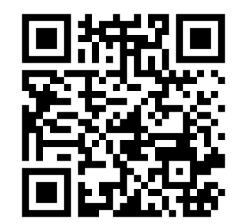
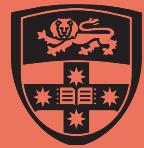
One way to use them is as a **transducer**

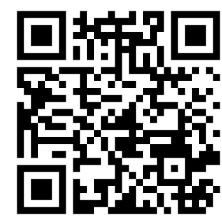
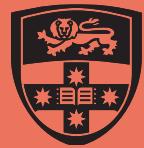




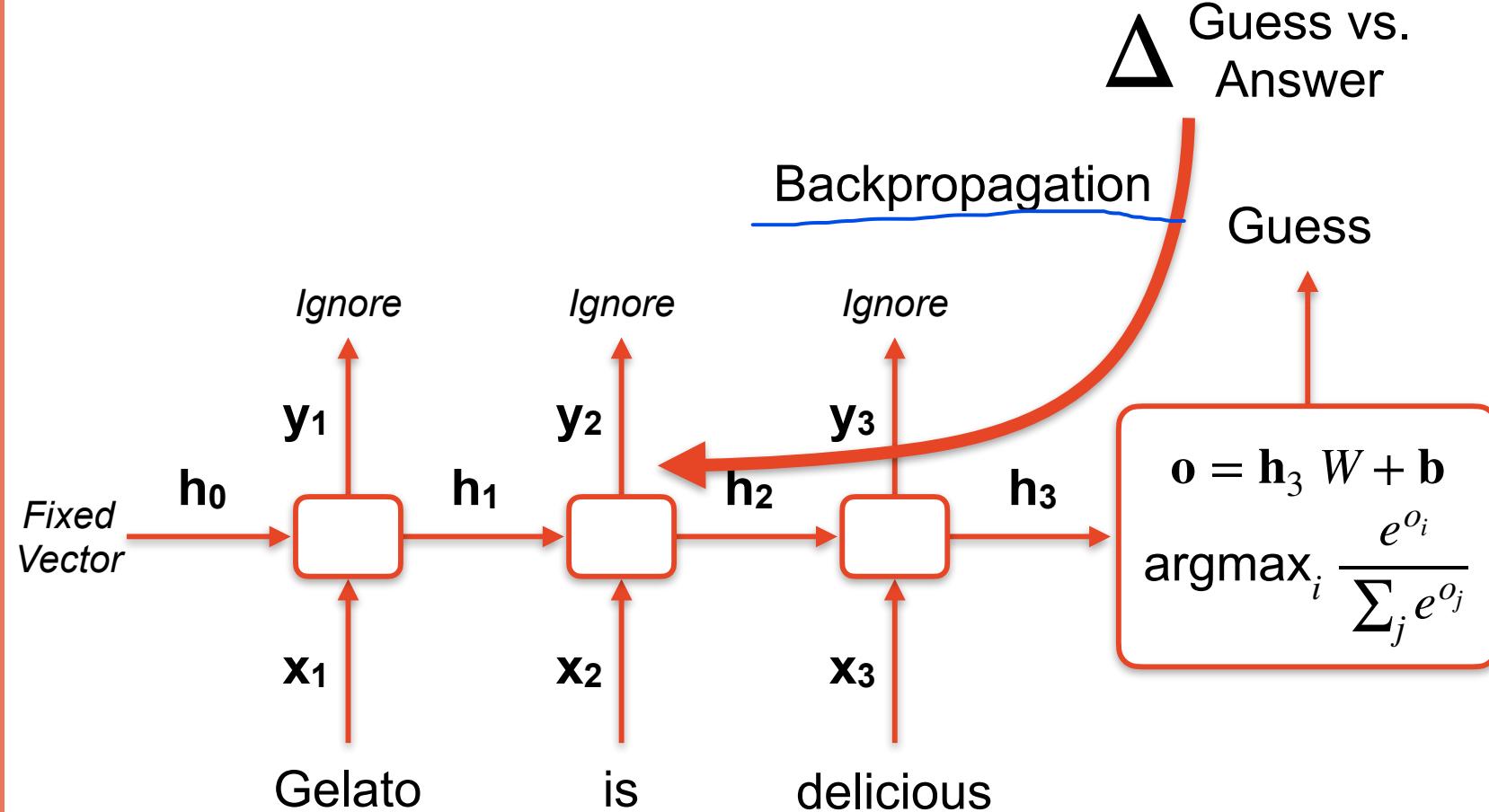
One way to use them is as a **transducer**

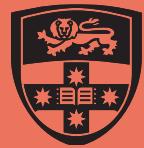




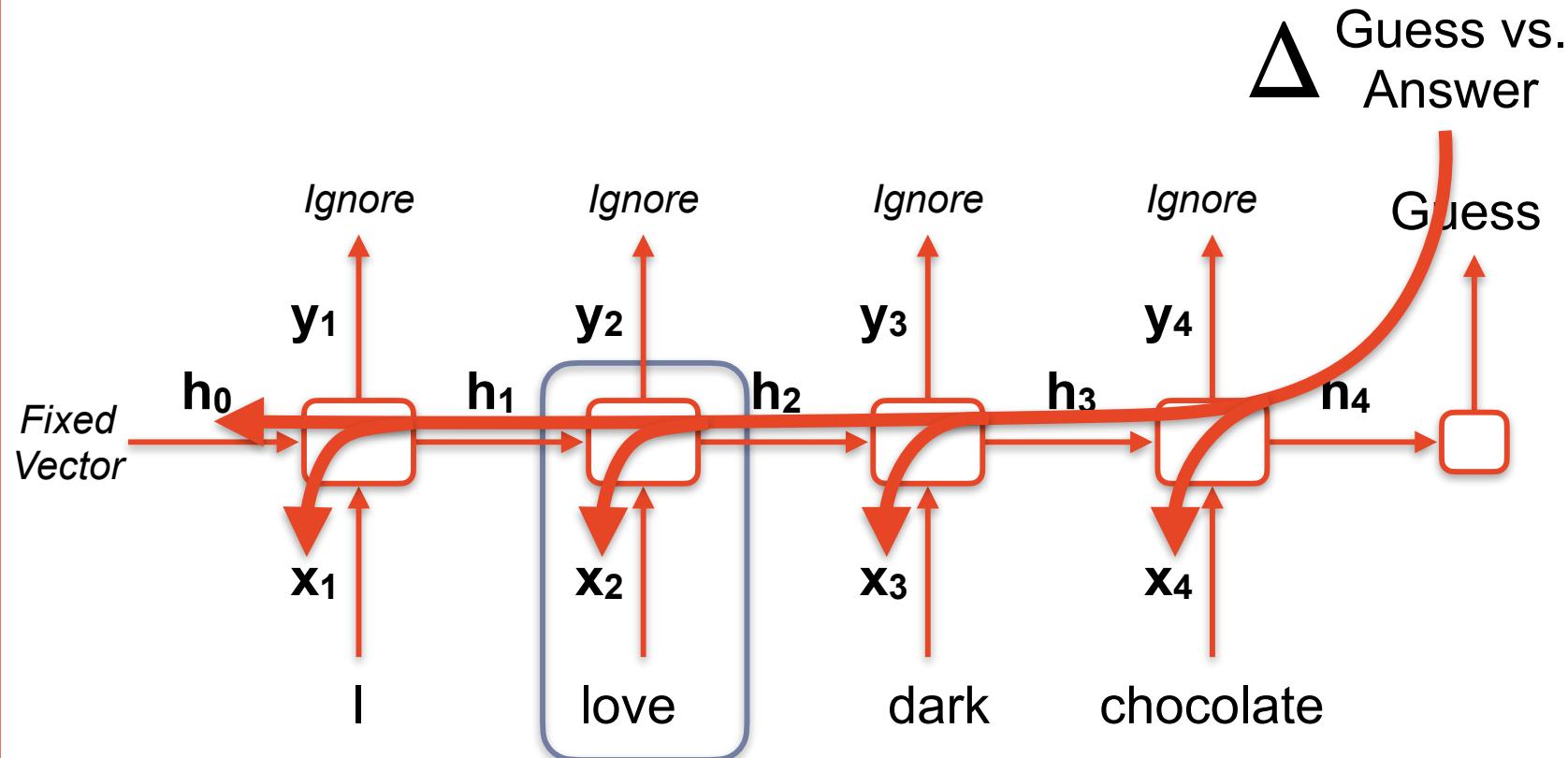


How does learning work?

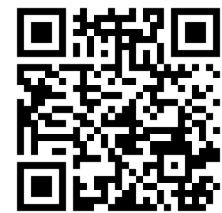
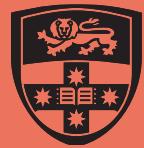




How does learning work?

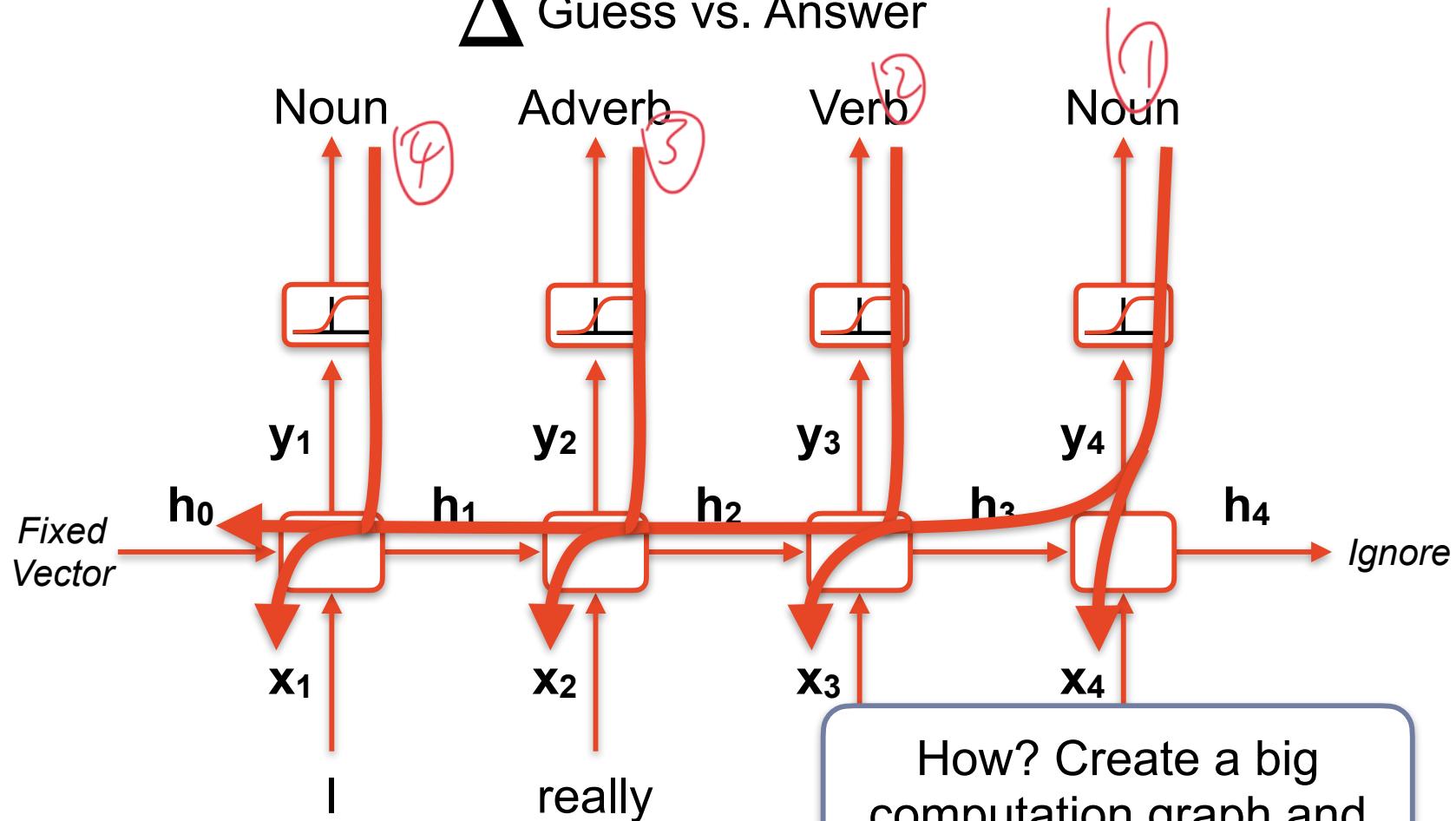


To fix the error, the derivative needs to update all the way back to x_2

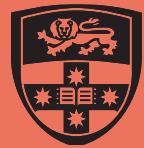


How does learning work?

△ Guess vs. Answer

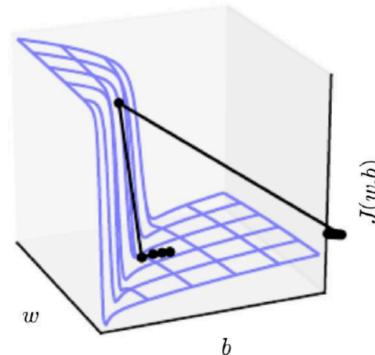


How? Create a big computation graph and calculate the derivatives



Two key challenges in learning: (1) Exploding gradients

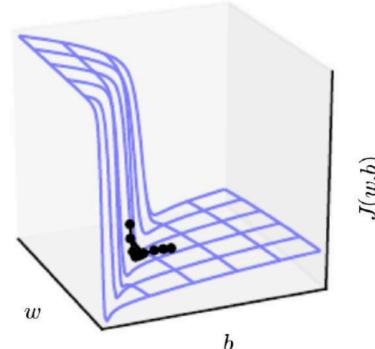
What happens if we have a very steep gradient?

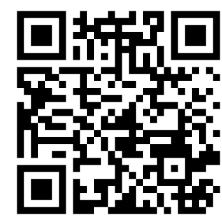
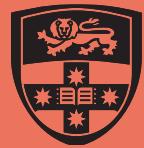


- Might jump far away to a bad weight
- Might overflow the floating point number storage

Soln

Clipping the gradient, i.e., using $\text{min}(\text{gradient}, k)$, solves this:





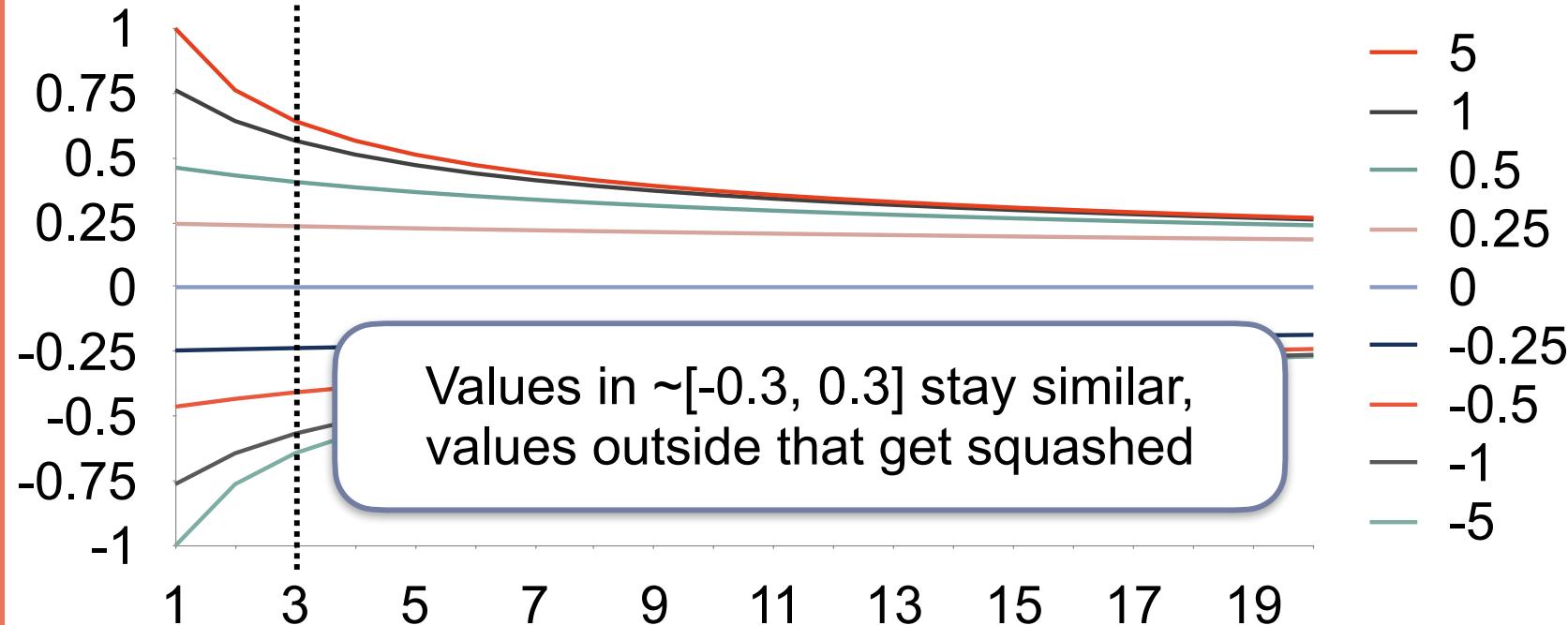
Two key challenges in learning: (2) Vanishing gradients

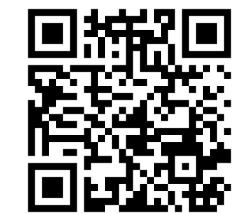
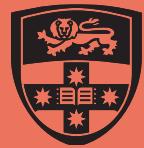
$$\mathbf{h}_t = \tanh(\mathbf{h}\mathbf{b} + W \mathbf{x}_t + V \mathbf{h}_{t-1})$$

$$= \tanh(\mathbf{h}\mathbf{b} + W \mathbf{x}_t + V \tanh(\mathbf{h}\mathbf{b} + W \mathbf{x}_{t-1} + V \mathbf{h}_{t-2}))$$

$$= \tanh(\mathbf{h}\mathbf{b} + W \mathbf{x}_t + V \tanh(\mathbf{h}\mathbf{b} + W \mathbf{x}_{t-1} + V \tanh(\mathbf{h}\mathbf{b} + W \mathbf{x}_{t-2} - V \mathbf{h}_{t-3})))$$

The values in this input go through tanh three times





Two key challenges in learning: (2) Vanishing gradients

The bias is added every time

$$x_t + V h_{t-1}$$

$$= \tanh(\mathbf{h}b + W x_t + V \tanh(\mathbf{h}b + W x_{t-1} + V h_{t-2}))$$

$$\tanh(\mathbf{h}b + W x_t)$$

$$= \tanh(\mathbf{h}b + W x_t + V \tanh(\mathbf{h}b + W x_{t-1} + V h_{t-2}) + V \tanh(\mathbf{h}b + W x_{t-2} + V h_{t-3}))$$

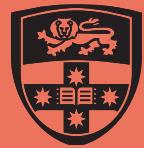
The values in this input is multiplied by W three times

New inputs are being added

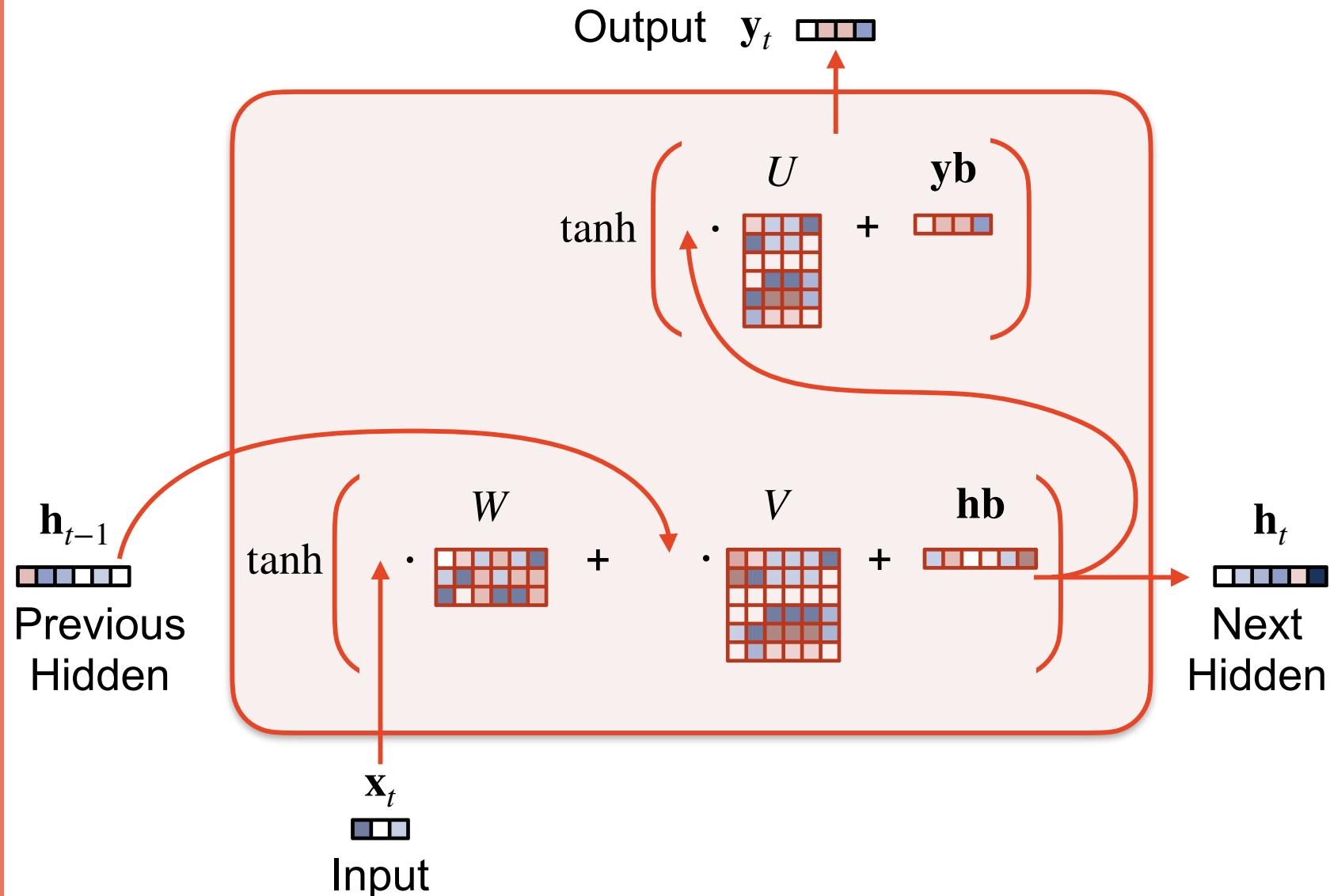
Result? model forgets quickly / mainly captures local information

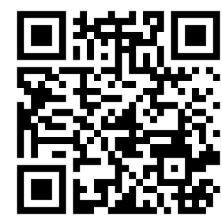
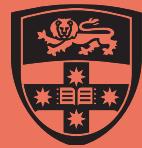
Small effect => Small gradient => Learning is hard

Vanishing Gradient description

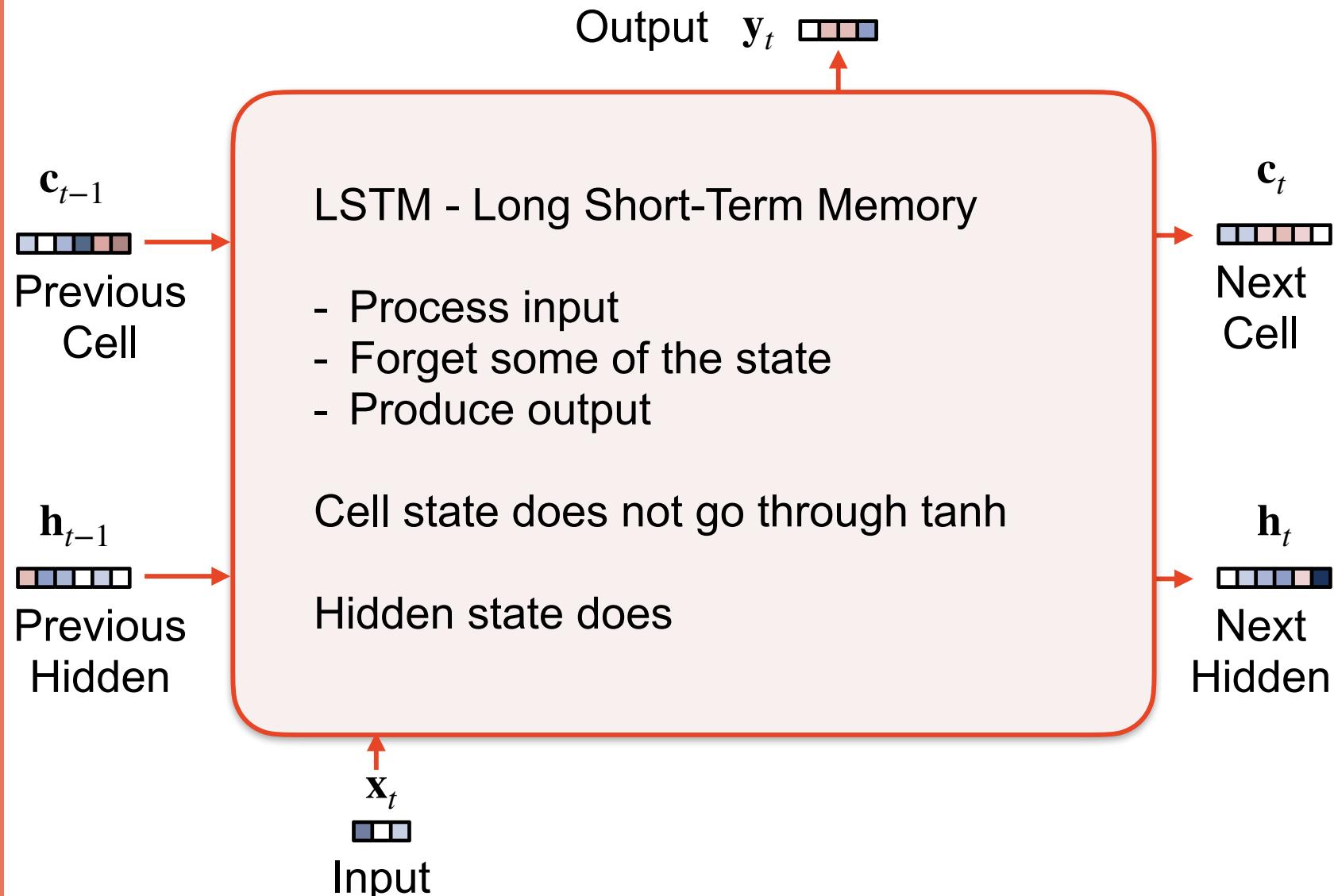


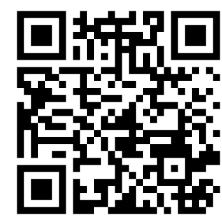
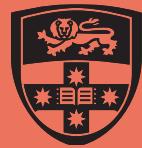
Key idea: maintain state in a non-decaying way





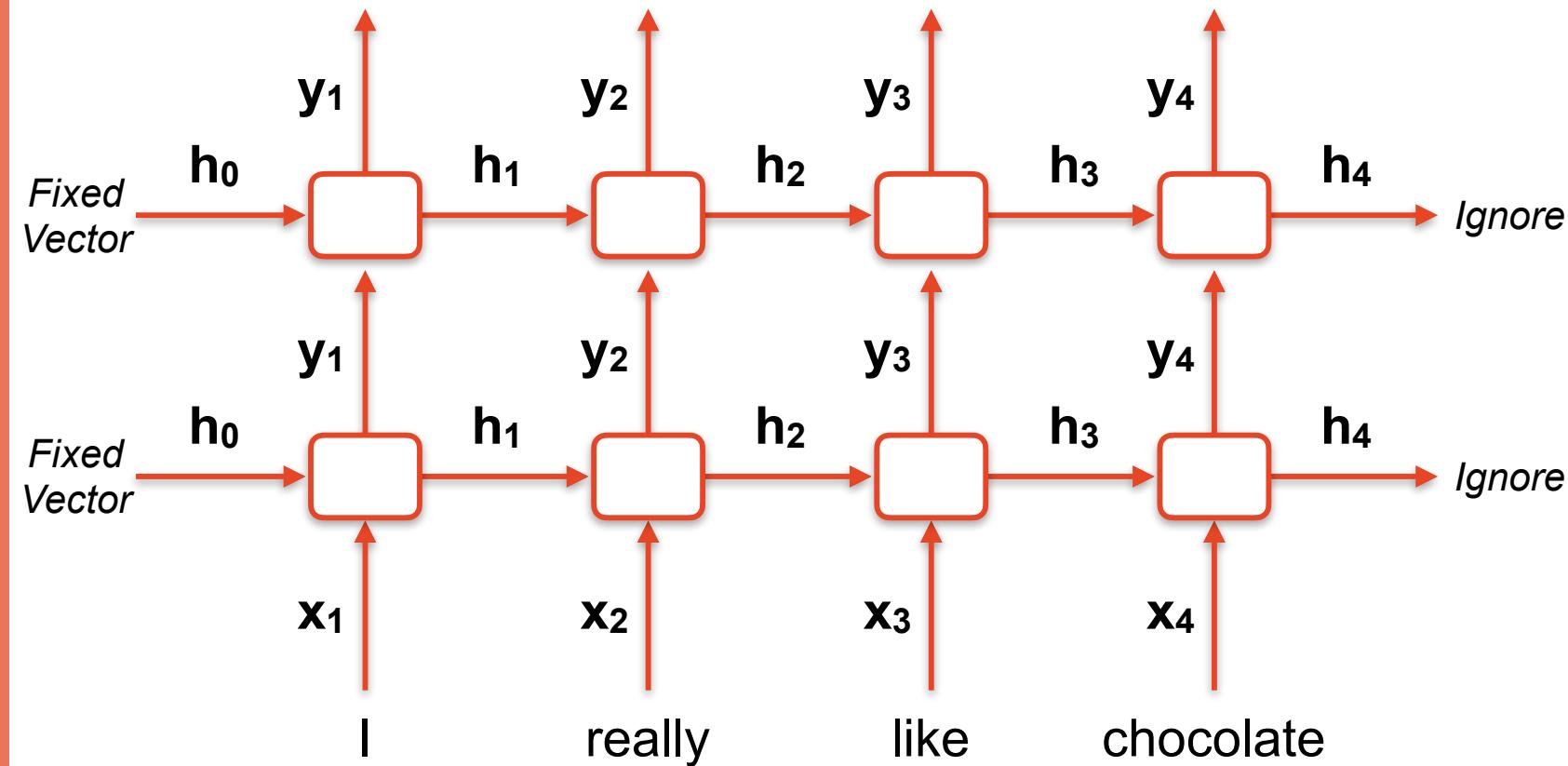
Key idea: maintain state in a non-decaying way

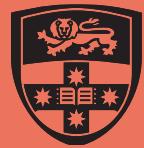




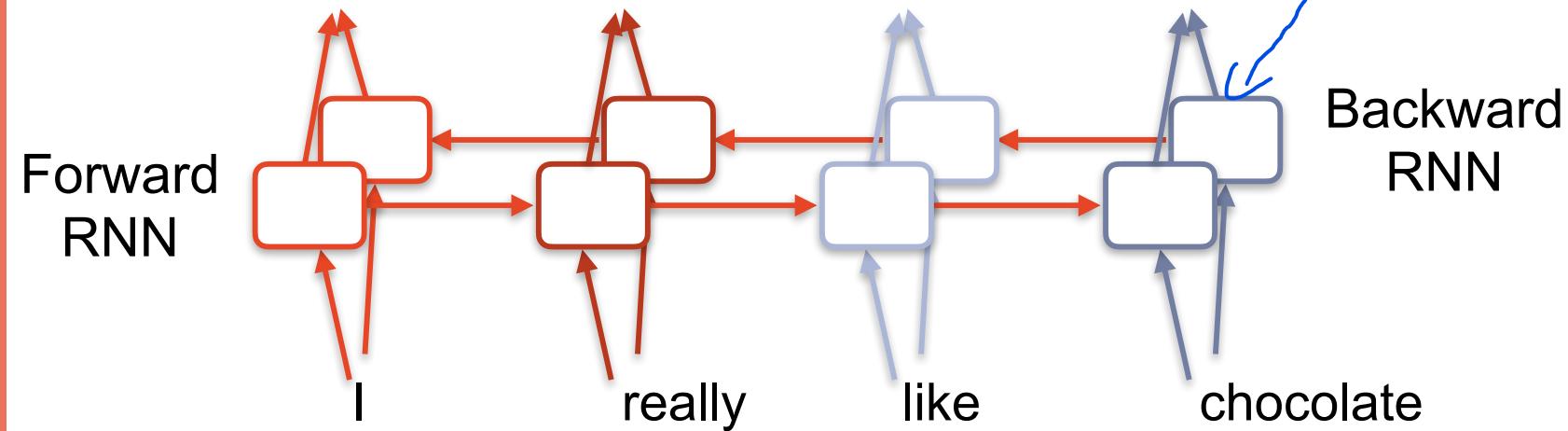
We can use the output of one as the input to another

Stack RNN

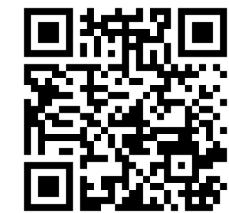




We can have two processing the input in reversed directions



Critically - this
does not
introduce loops



Part of Speech tagging

Noun	Adverb	Verb	Noun
I	really	like	chocolate

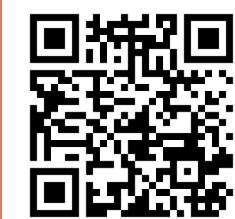
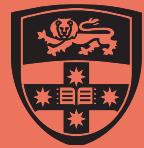
Categories of words, defined by:

Morphological context - combine with similar affixes
(ie., prefixes and suffixes)

Verbs

	-s	-ed	-ing
walk	walks	walked	walking
slice	slices	sliced	slicing
believe	believes	believed	believing

David Bamman's Info 159/259 lectures at Berkeley



Part of Speech tagging

Noun Adverb Verb Noun
I really like chocolate

Categories of words, defined by:

Morphological context - combine with similar affixes
(ie., prefixes and suffixes)

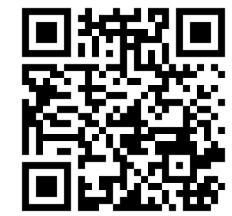
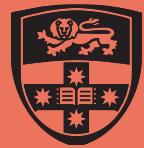
Verbs

Not Verbs

	-s	-ed	-ing
Verbs	walk	walks	walked
	slice	slices	sliced
	believe	believes	believed
Not Verbs	of	*ofs	*ofed
	red	*reds	*redded

Exceptions:
Sleep
Eat
Give

Problem
David Bamman's Info 159/259 lectures at Berkeley

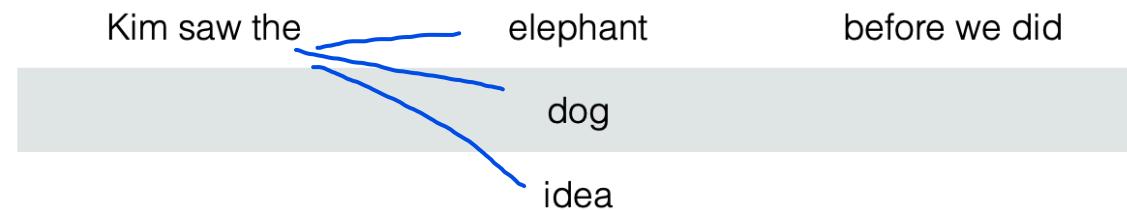


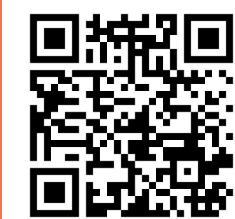
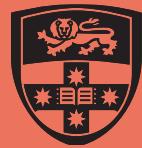
Part of Speech tagging

Noun	Adverb	Verb	Noun
I	really	like	chocolate

Categories of words, defined by:

Syntactic distribution - can be replaced with another word and remain grammatical





Part of Speech tagging

Noun	Adverb	Verb	Noun
I	really	like	chocolate

Categories of words, defined by:

Syntactic distribution - can be replaced with another word and remain grammatical

Kim saw the elephant before we did

dog

idea

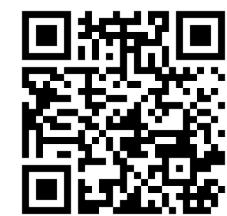
*of

*goes

-P r o f e m

May be too strict, e.g., “Jonathan” cannot go here

David Bamman's Info 159/259 lectures at Berkeley



不同词性导致不同意义

Part of Speech tagging can help distinguish meanings



NN



VB

DT

NN



NN

VBZ

IN

DT

NN

Time

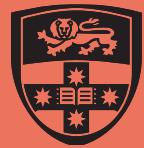
flies

like

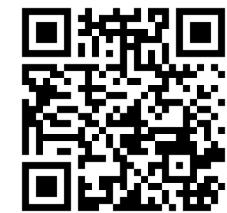
an

arrow



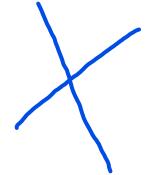


Neural Networks
Recurrent Models
Analysis
Workshop Preview



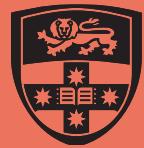
[menti.com 1750 7815](https://menti.com/17507815)

Part of Speech tagging can help indicate pronunciation



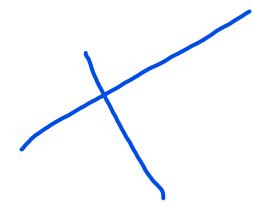
Noun	Verb
My conduct is great	I conduct myself well
She won the contest	I contest the ticket
He is my escort	He escorted me
That is an insult	Don't insult me
Rebel without a cause	He likes to rebel
He is a suspect	I suspect him

David Bamman's Info 159/259 lectures at Berkeley



[menti.com 1750 7815](https://menti.com/17507815)

How well do current models do?



English News

Rule based (guess most frequent label): 92%

Machine learning: 97%

But!

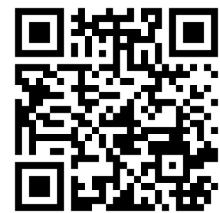
There are ~20 words per sentence...
sentence accuracy is 55%

Drops on other text types (out of domain)

David Bamman's Info 159/259 lectures at Berkeley



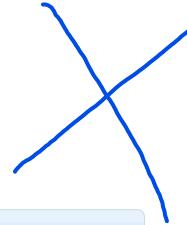
Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

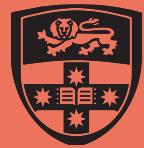
Where can I find data?

Universal Dependencies

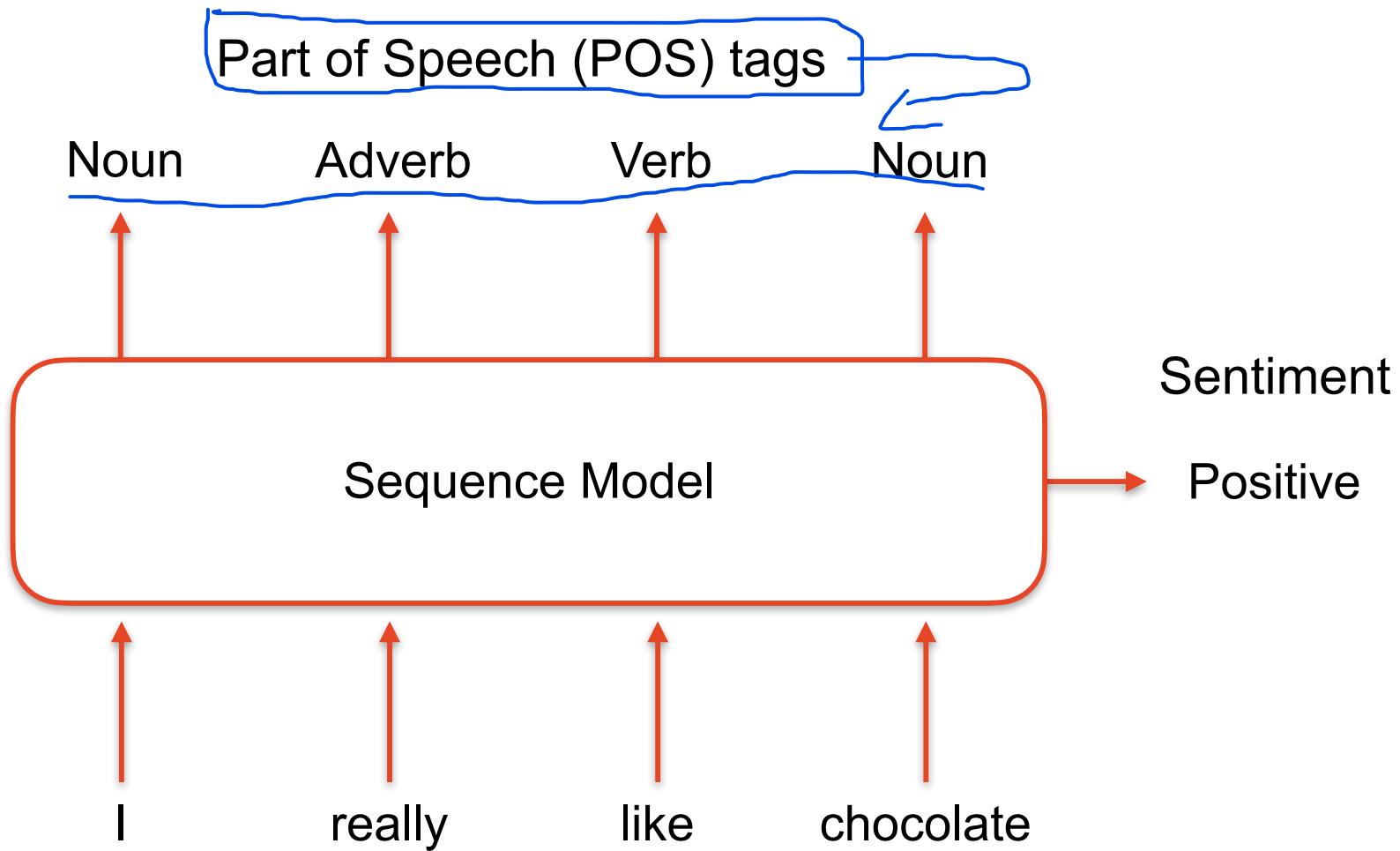


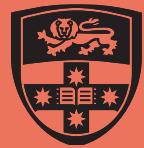
▶	Abaza	1	<1K	💬		Northwest Caucasian
▶	Afrikaans	1	49K	✍️ℹ️		IE, Germanic
▶	Akkadian	2	25K	📖ℹ️		Afro-Asiatic, Semitic
▶	Akuntsu	1	1K	📖ℹ️		Tupian, Tupari
▶	Albanian	1	<1K	W		IE, Albanian
▶	Amharic	1	10K	☁️✍️💡ℹ️		Afro-Asiatic, Semitic
▶	Ancient Greek	3	456K	☁️💡ℹ️		IE, Greek
▶	Ancient Hebrew	1	39K	☁️		Afro-Asiatic, Semitic
▶	Apurina	1	<1K	📖ℹ️		Arawakan
▶	Arabic	3	1,042K	📖W		Afro-Asiatic, Semitic
▶	Armenian	2	94K	✍️✍️✍️💡ℹ️W		IE, Armenian
▶	Assyrian	1	<1K	📖ℹ️		Afro-Asiatic, Semitic
▶	Bambara	1	13K	📖ℹ️		Mande
▶	Basque	1	121K	📖		Basque
▶	Beja	1	1K	💬		Afro-Asiatic, Cushitic
▶	Belarusian	1	305K	✍️✍️💡ℹ️🎵W		IE, Slavic
▶	Bengali	1	<1K	✍️		IE, Indic
▶	Bhojpuri	1	6K	📖ℹ️		IE, Indic
▶	Bororo	1	1K	✍️		Bororoan
▶	Breton	1	10K	✍️💡ℹ️🎵W		IE, Celtic
▶	Bulgarian	1	156K	✍️✍️		IE, Slavic
▶	Buryat	1	10K	✍️💡		Mongolic
▶	Cantonese	1	13K	💬		Sino-Tibetan
▶	Catalan	1	553K	📖		IE, Romance
▶	Cebuano	1	1K	✍️		Austronesian, Central Philippine
▶	Chinese	7	309K	✍️✍️💡ℹ️W		Sino-Tibetan
▶	Chukchi	1	6K	💬		Chukotko-Kamchatkan
▶	Classical Armenian	1	13K	☁️		IE, Armenian
▶	Classical Chinese	1	433K	ℹ️🎵		Sino-Tibetan
▶	Coptic	1	57K	☁️✍️		Afro-Asiatic, Egyptian
▶	Croatian	1	199K	📖💡W		IE, Slavic
▶	Czech	6	2,253K	✍️✍️💡ℹ️🎵W		IE, Slavic
▶	Danish	1	100K	✍️💡ℹ️		IE, Germanic
▶	Dutch	2	306K	📖W		IE, Germanic
▶	English	10	726K	✍️✍️✍️💡ℹ️🎵W		IE, Germanic
▶	Erzya	1	20K	✍️		Uralic, Mordvin
▶	Estonian	2	520K	✍️✍️💡ℹ️		Uralic, Finnic

<https://universaldependencies.org/>



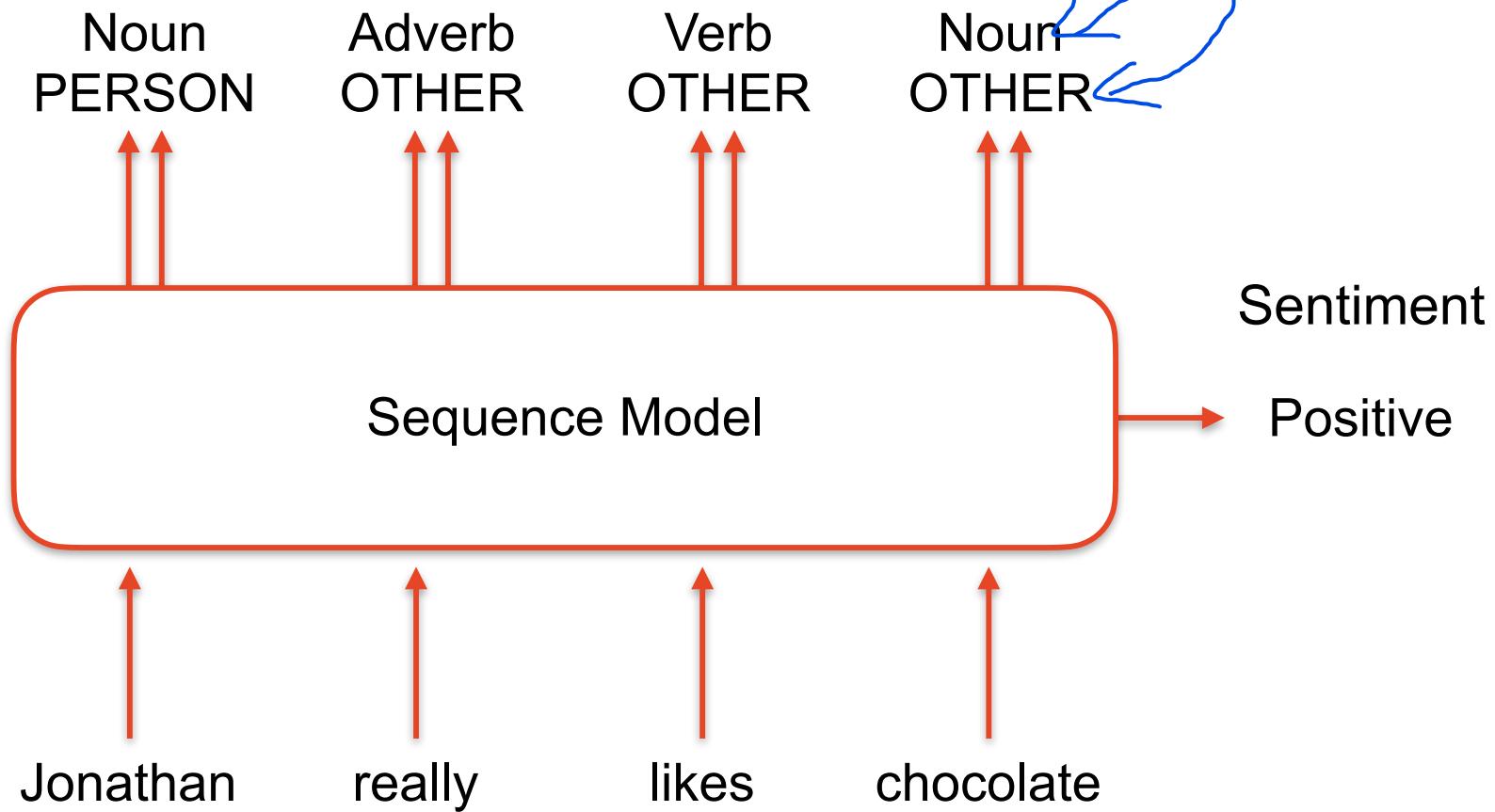
We can do multiple tasks simultaneously with a single model

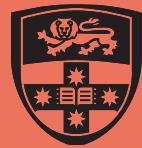




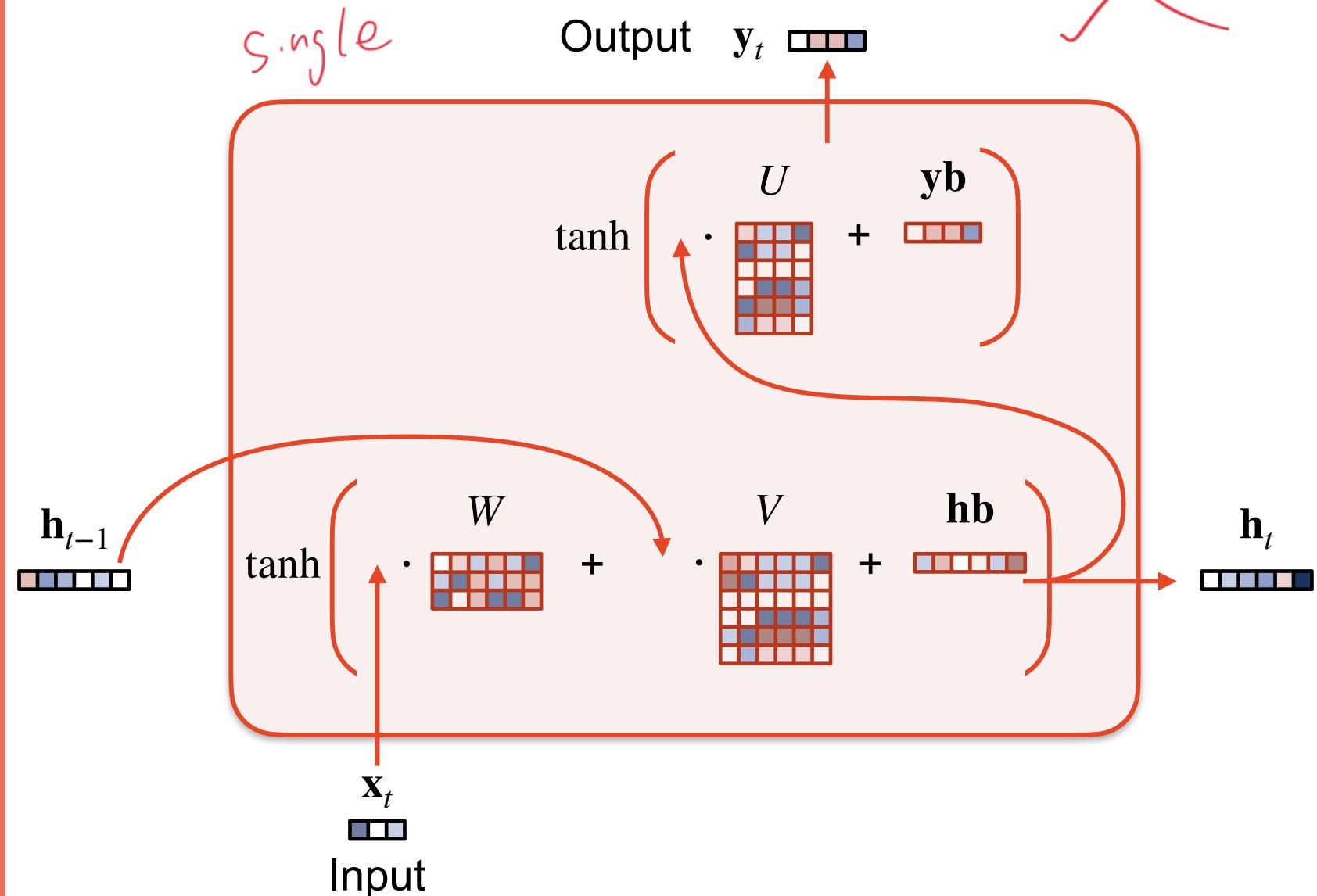
We can train on multiple related tasks

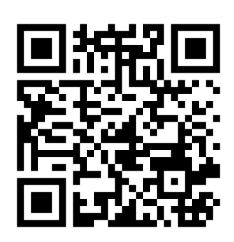
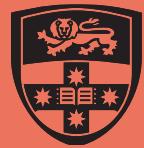
Part of Speech (POS) tags and
Named Entity Recognition (NER)





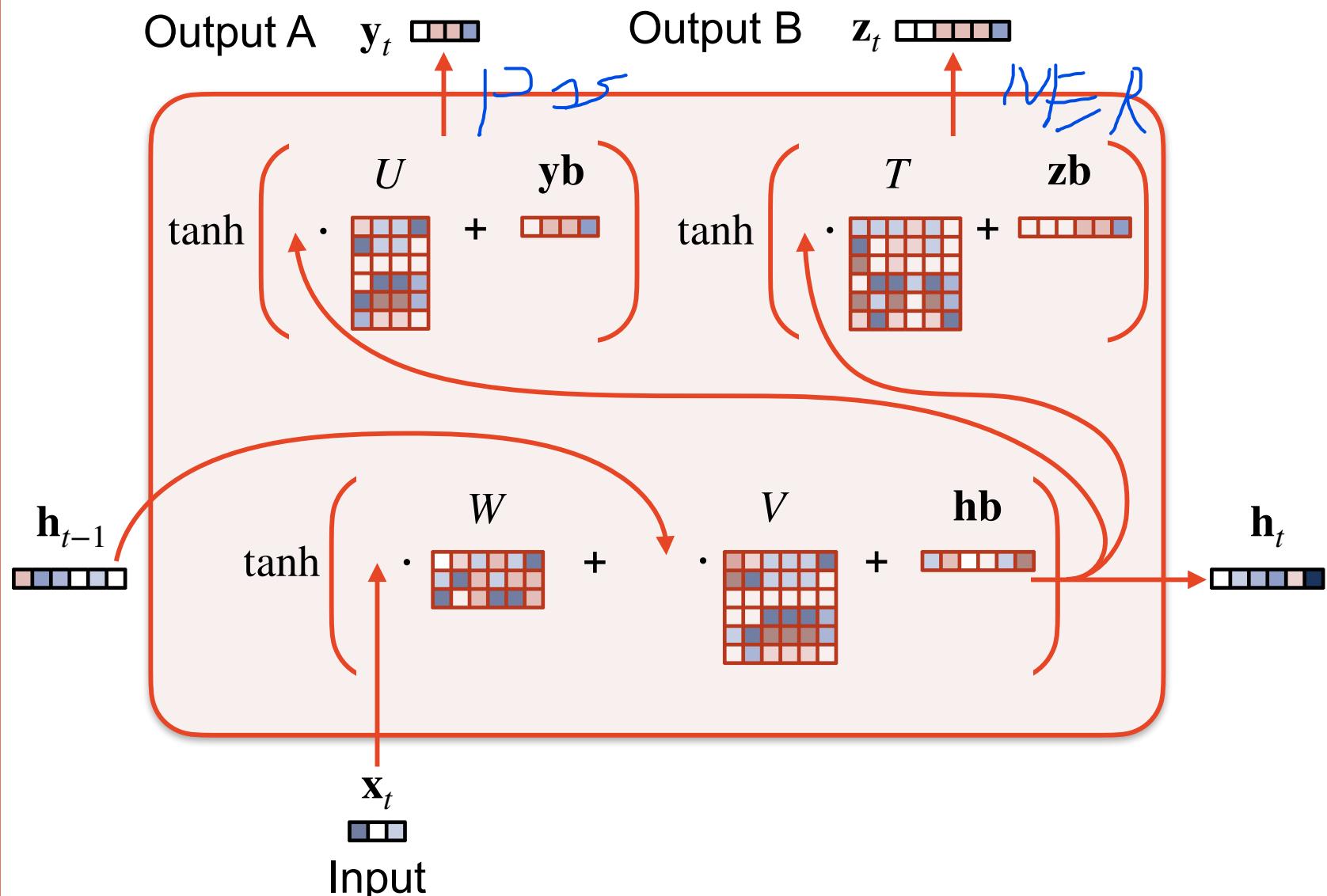
We can train on multiple related tasks





We can train on multiple related tasks

multiple





Named Entity Recognition

NER

Person	Other	Other	Other
Jonathan	really	likes	chocolate

Identifying proper nouns:

Coarse-grained:

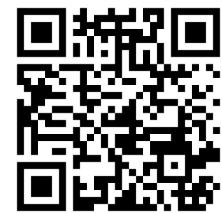
- Locations
- Geo-Political Entities
- People
- Organisations

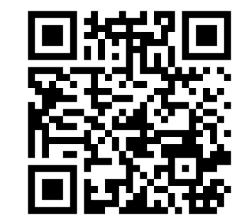
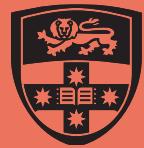
Example

- Blackwattle Bay
- Sydney
- Jonathan
- University of Sydney

Fine-grained:

112 categories in Ling and Weld (2012)!





Named Entity Recognition

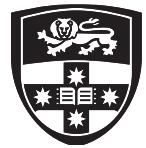
Person Other Other Other
Jonathan really begin likes chocolate
 other 是是? End
Identifying proper nouns:
Inside

Words	IO Label	BIO Label	BIOES Label
Jane	I-PER	B-PER	B-PER
Villanueva	I-PER	I-PER	E-PER
of	O	O	O
United	I-ORG	B-ORG	B-ORG
Airlines	I-ORG	I-ORG	I-ORG
Holding	I-ORG	I-ORG	E-ORG
discussed	O	O	O
the	O	O	O
Chicago	I-LOC	B-LOC	S-LOC
route	O	O	O
.	O	O	O

Why does it matter?

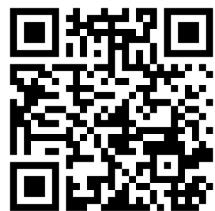
I gave Joe Smith a book
I gave Joe Zach's book

Jurafsky and Martin, chapter 8



COMP 4446 / 5046
Lecture 3, 2025

Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

Analysis



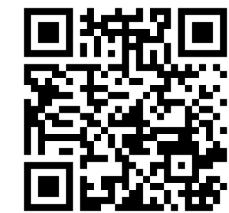
In lecture 2, we saw we could form a confusion matrix

	Answer: Spam	Answer: Not Spam
Guess: Spam		
Guess: Not Spam		

- Spam and predicted Spam
- Not spam and predicted Not spam
- Spam, predicted Not spam
- Not spam, predicted Spam



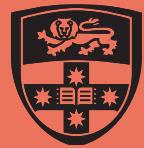
Neural Networks
Recurrent Models
Analysis
Workshop Preview



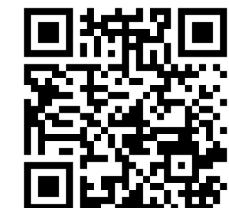
[menti.com 1750 7815](https://menti.com/17507815)

In lecture 2, we saw we could form a confusion matrix

		True Answer				
		A	B	C	...	None
Guess	A	TP				FP
	B		TP			FP
	C			TP		FP
	...				TP	FP
	None	FN	FN	FN	FN	TN



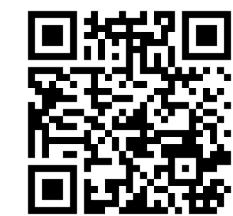
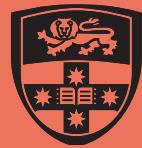
Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

In lecture 2, we saw we could form a confusion matrix

		True Answer				
		Person	Location	Org	...	Other
Guess	Person	TP				FP
	Location		TP			FP
	Org			TP		FP
	...				TP	FP
	Other	FN	FN	FN	FN	TN



For NER, we consider typed spans / chunks

Convert tags into labeled spans:

	Words	BIO Label	
0	Jane	B-PER	
1	Villanueva	I-PER	
2	of	O	
3	United	B-ORG	
4	Airlines	I-ORG	
5	Holding	I-ORG	
6	discussed	O	
7	the	O	
8	Chicago	B-LOC	
9	route	O	
10	.	O	

Annotations: A red box highlights "For NER, we consider typed spans / chunks". A red box highlights "Convert tags into labeled spans:". Blue arrows point from the BIO labels to the corresponding spans: [0, 1, PER], [3, 5, ORG], and [8, 8, LOC].

Compare true answer and guess:

Match? TP

Item in guess, but not answer? FP

Item in answer, but not guess? FN

D - I I - O
D - G I - G

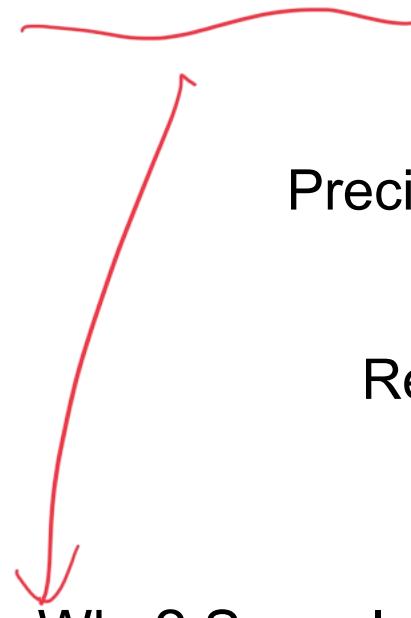
Note: Possible outputs = | Labels | * Length * Length
TN will be huge!

label

Jurafsky and Martin, chapter 8



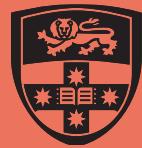
NER is a good example of where micro-P, micro-R, and micro-F are used



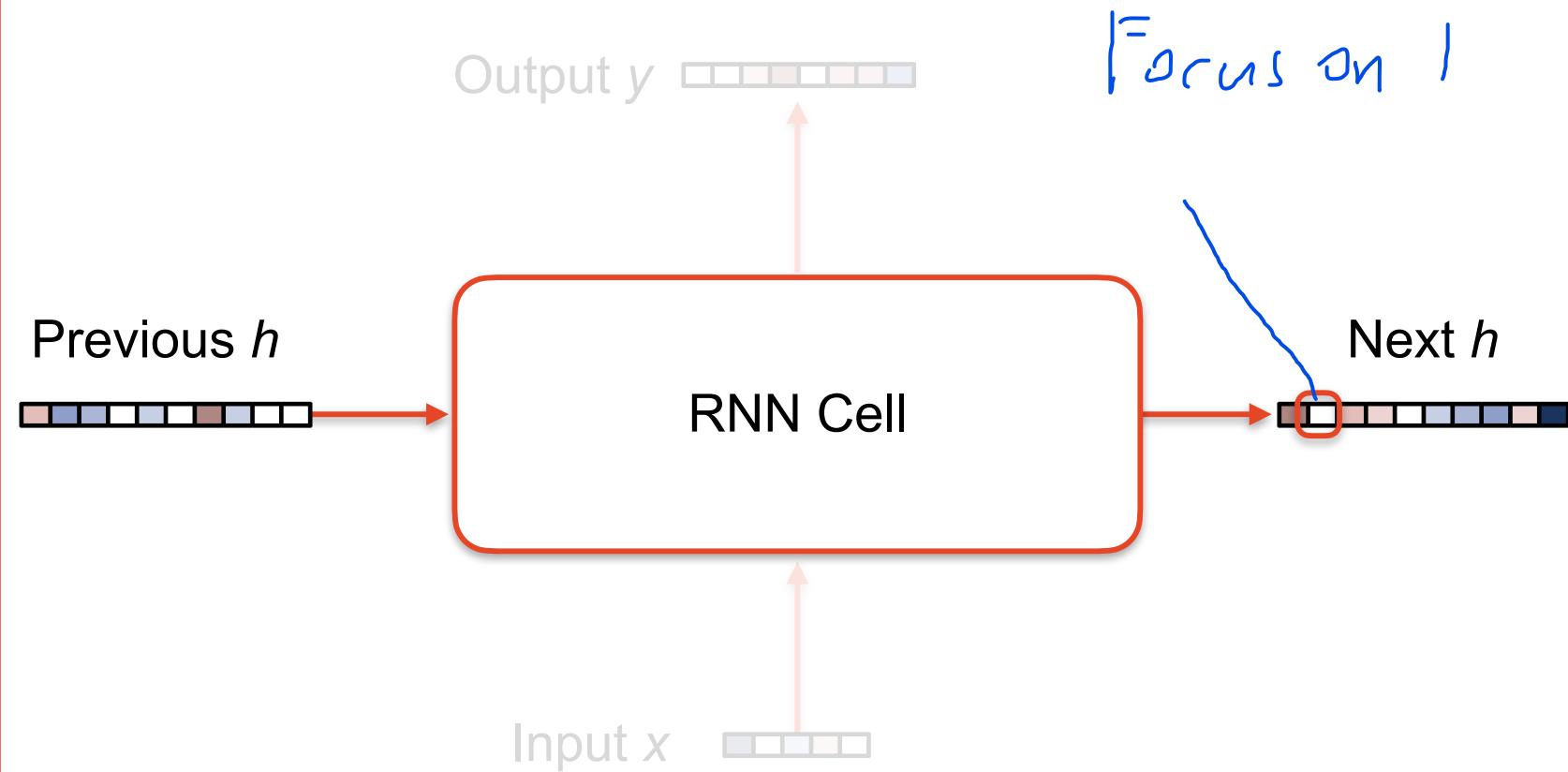
$$\text{Precision}_{\text{micro}} = \frac{\sum_{l \in \text{labels}} \text{TP}_l}{\sum_{l \in \text{labels}} \text{TP}_l + \text{FP}_l}$$

$$\text{Recall}_{\text{micro}} = \frac{\sum_{l \in \text{labels}} \text{TP}_l}{\sum_{l \in \text{labels}} \text{TP}_l + \text{FN}_l}$$

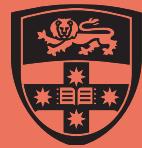
Why? Some labels are less common and we care about them in proportion to their frequency



We can visualise (1) values in the hidden state



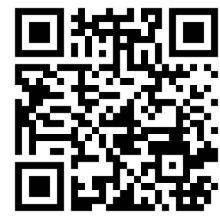
“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



Neural Networks
Recurrent Models
Analysis
Workshop Preview

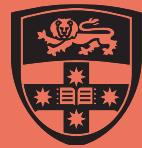
We can visualise (1) values in the hidden state

http://www.ynetnews.com/ English-language website of Israel's lar
gest newspaper '' [[Yedioth Ahronoth]] ''' Hebrew-language period
icals:''' *' [[Globes]]''' http://www.globes.co.il/ business da
ily *' [[Haaretz|Ha'aretz]]''' http://www.haaretz.co.il/ Relativ

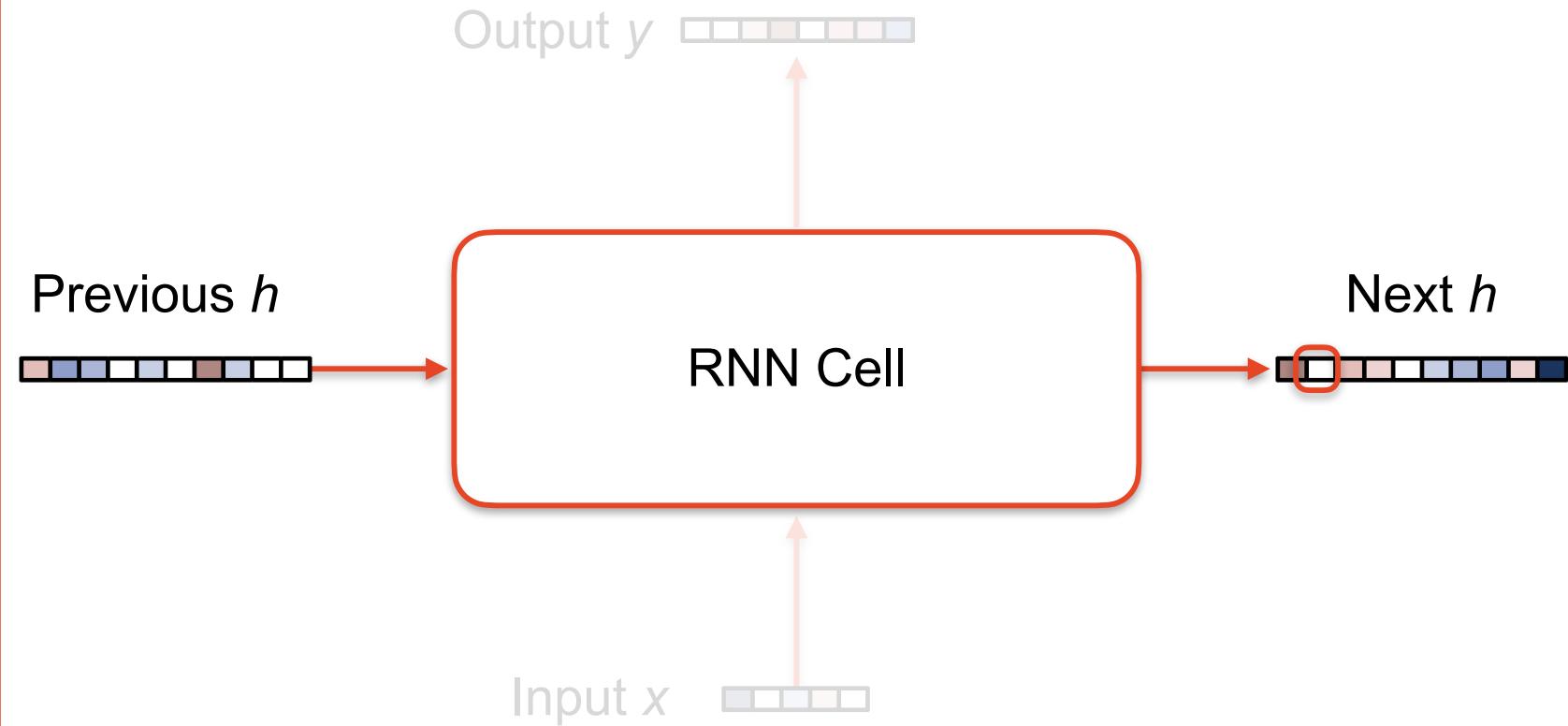


menti.com 1750 7815

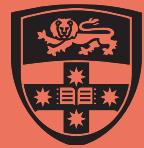
“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



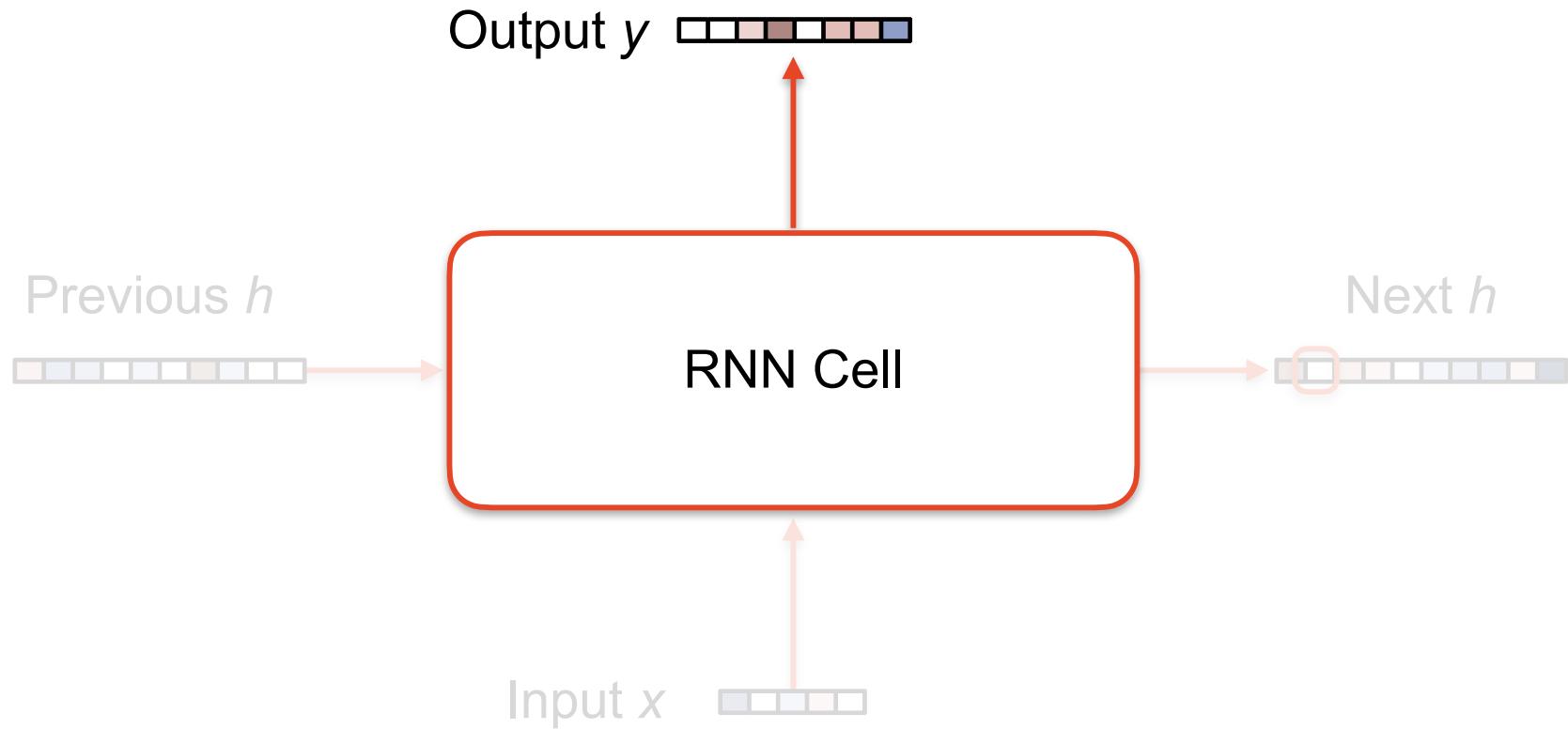
We can visualise (1) values in the hidden state



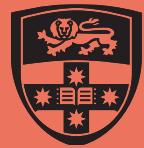
“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



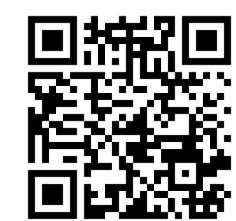
We can visualise (1) values in the hidden state



“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



Neural Networks Recurrent Models **Analysis** Workshop Preview



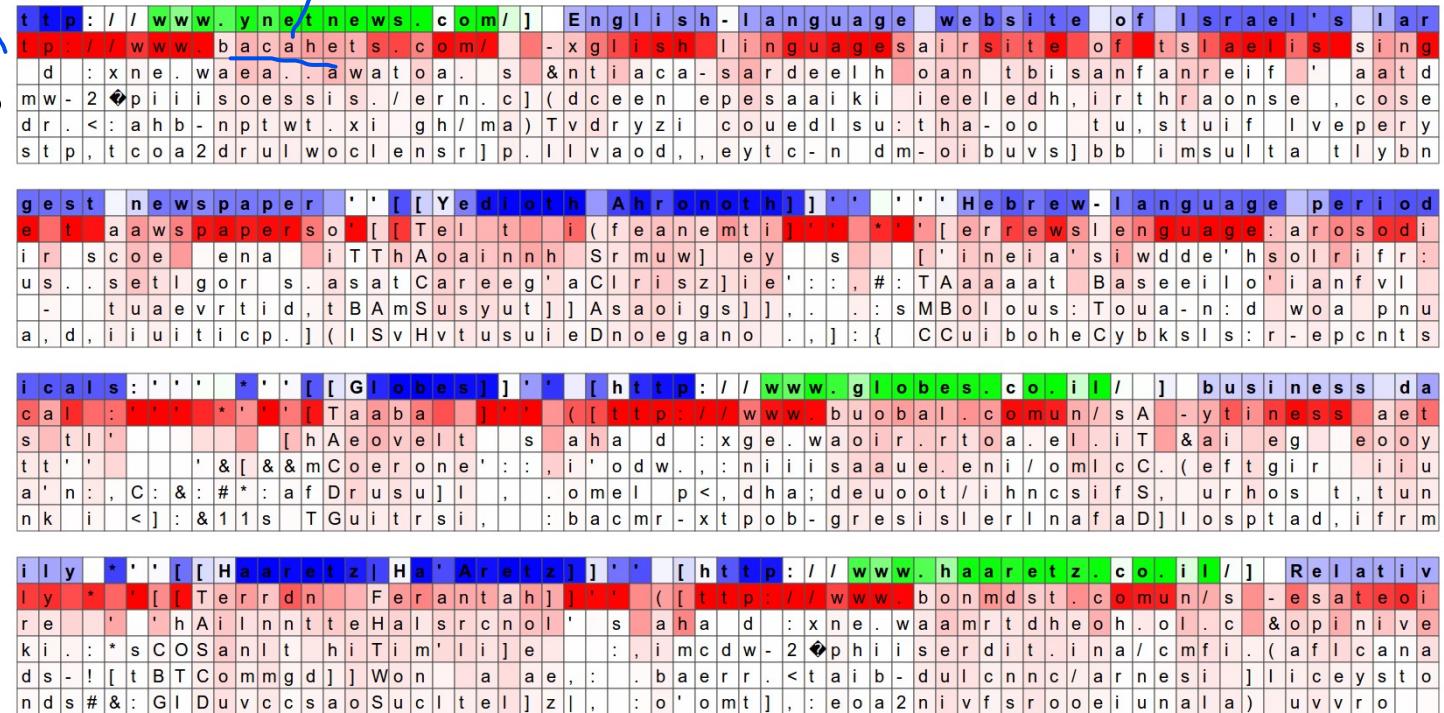
menti.com 1750 7815

We can visualise (1) values in the hidden state, and (2) the output distribution

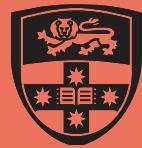
confident predict not confident

Input

Guesses



“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



Neural Networks Recurrent Models **Analysis** Workshop Preview

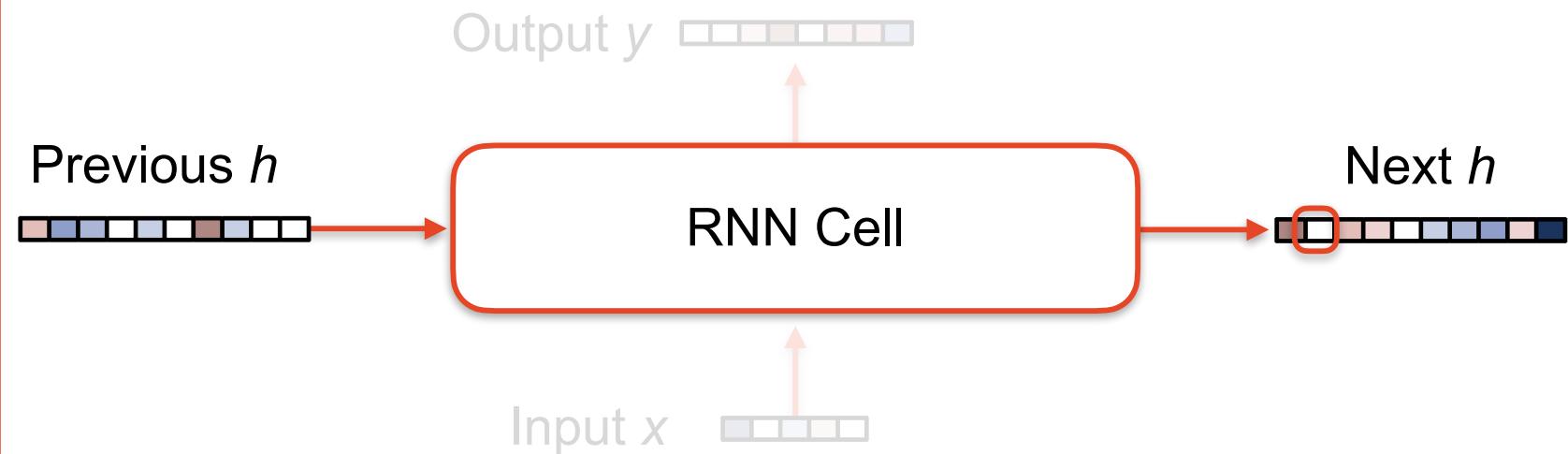


[menti.com 1750 7815](https://menti.com/17507815)

We can visualise (1) values in the hidden state, and (2) the output distribution

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.



“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



Neural Networks
Recurrent Models
Analysis
Workshop Preview



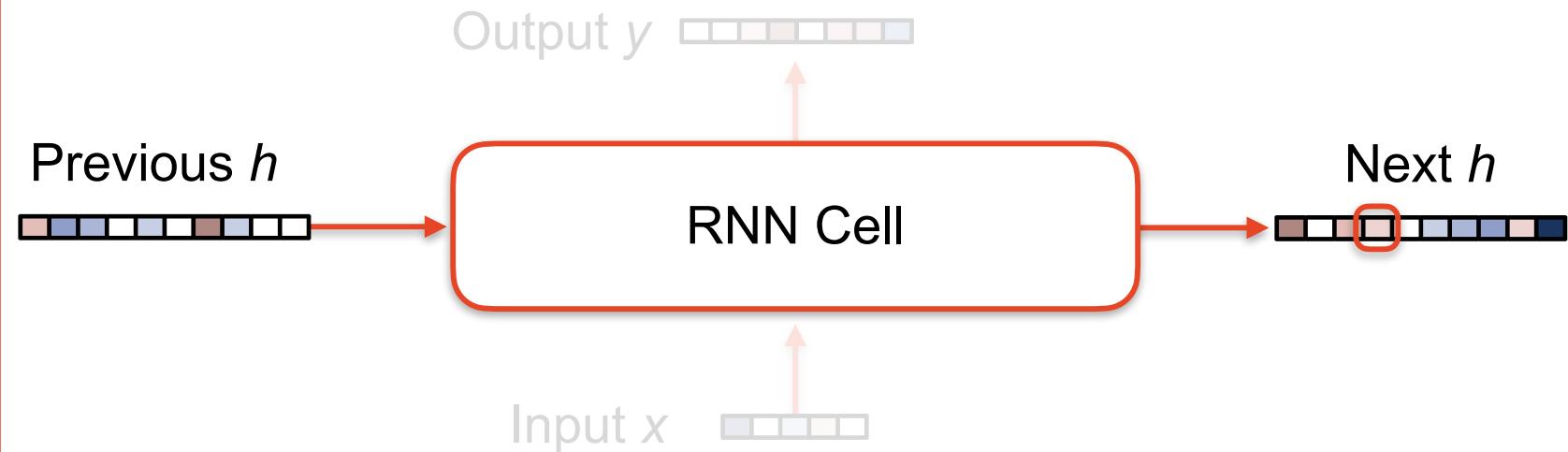
[menti.com 1750 7815](https://menti.com/17507815)

We can visualise (1) values in the hidden state, and (2) the output distribution

Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."



"The Unreasonable Effectiveness of Recurrent Neural Networks", Andrej Karpathy, 2015



Neural Networks
Recurrent Models
Analysis
Workshop Preview

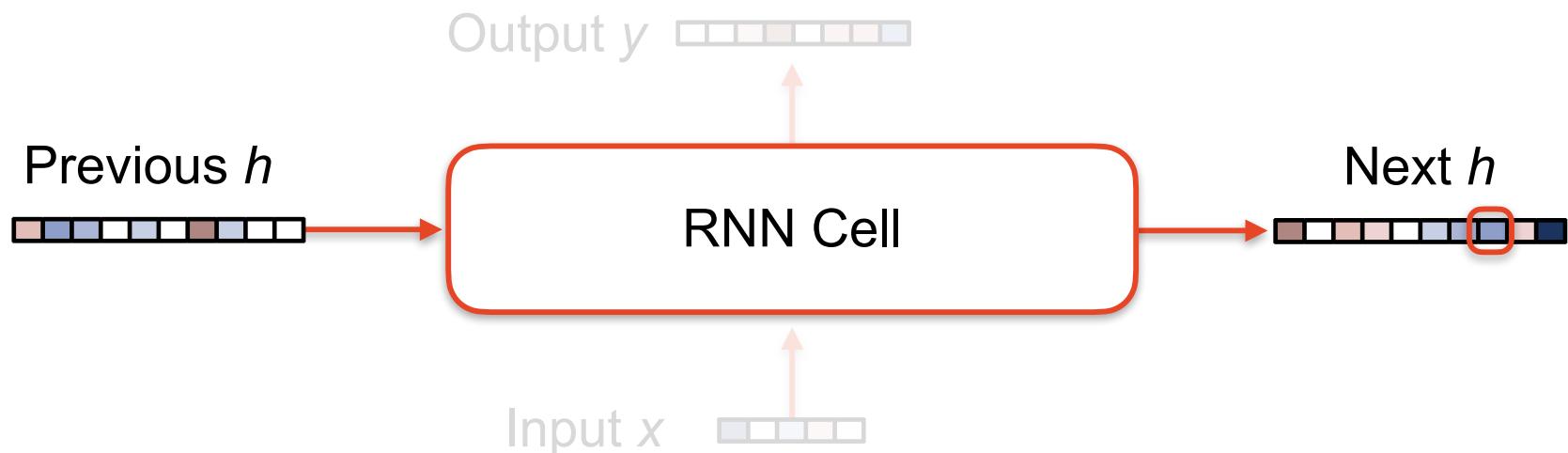


[menti.com 1750 7815](https://menti.com/17507815)

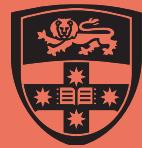
We can visualise (1) values in the hidden state, and (2) the output distribution

Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,  
    siginfo_t *info)  
{  
    int sig = next_signal(pending, mask);  
    if (sig) {  
        if (current->notifier) {  
            if (sigismember(current->notifier_mask, sig)) {  
                if (!!(current->notifier)(current->notifier_data)) {  
                    clear_thread_flag(TIF_SIGPENDING);  
                    return 0;  
                }  
            }  
        }  
        collect_signal(sig, pending, info);  
    }  
    return sig;  
}
```



“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



Neural Networks
Recurrent Models
Analysis
Workshop Preview

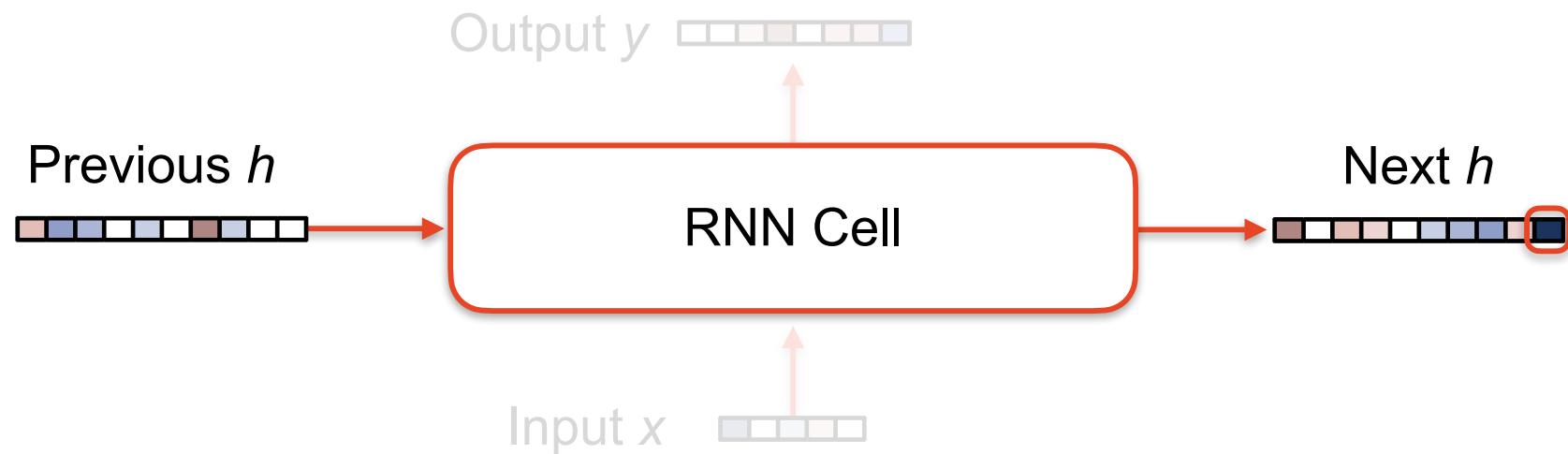


[menti.com 1750 7815](https://menti.com/17507815)

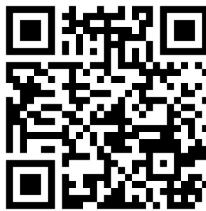
We can visualise (1) values in the hidden state, and (2) the output distribution

A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
    */
```



“The Unreasonable Effectiveness of Recurrent Neural Networks”, Andrej Karpathy, 2015



Boathouses and Houseboats

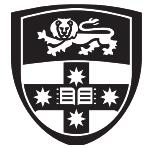
A THIS →
THAT HOLDS THIS
CAR

	CAR	HOUSE	BOAT
CAR	TOW TRUCK CARCAR	GARAGE CARHOUSE	CAR FERRY CARBOAT
HOUSE	MOBILE HOME HOUSECAR	APARTMENT HOUSEHOUSE	HOUSEBOAT
BOAT	BOAT TRAILER BOATCAR	BOATHOUSE	LIFEBOAT BOATBOAT

I REALLY LIKE THE WORDS FOR "BOATHOUSE" AND "HOUSEBOAT" AND THINK WE SHOULD APPLY THAT SCHEME MORE CONSISTENTLY.

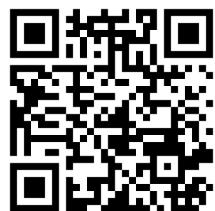
[The <x> that is held by <y> is also a <y><x>, so if you go to a food truck, the stuff you buy is truck food. A phone that's in your car is a carphone, and a car equipped with a phone is a phonecar. When you play a mobile racing game, you're in your phonecar using your carphone to drive a different phonecar. I'm still not sure about bananaphones.]

Source: <https://xkcd.com/2043/>



COMP 4446 / 5046
Lecture 3, 2025

Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 1750 7815](https://menti.com/17507815)

Workshop Preview



COMP 4446 / 5046
Lecture 3, 2025

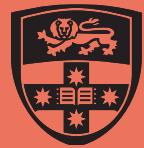
Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 7927 6661](https://menti.com/79276661)

Workshop:

 PyTorch



Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 7927 6661](https://menti.com/79276661)

Muddy Card

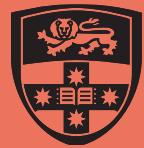
Week 2 Muddy Card (TEST1001)

Your Muddy Card Response: How does photosynthesis use water?

Would an answer to any of the following also answer your question? You can choose **zero or more** options. When finished, press "submit" to continue.

- Where does the water come from in photosynthesis?
- How does water split in photosynthesis?
- What exactly is the role of water in photosynthesis?
- Why is water necessary for photosynthesis?

Submit



Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 7927 6661](https://menti.com/79276661)

Muddy Card

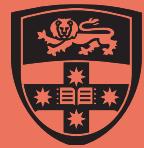
Week 2 Muddy Card (TEST1001)

Your Muddy Card Response: What is photosynthesis for?

Would an answer to any of the following also answer your question? You can choose **zero or more** options. When finished, press "submit" to continue.

- Why is photosynthesis important for humans?
- Why is water necessary for photosynthesis?
- How does photosynthesis affect the environment?
- How do scientists study photosynthesis?

Submit



COMP 4446 / 5046
Lecture 3, 2025

Neural Networks
Recurrent Models
Analysis
Workshop Preview

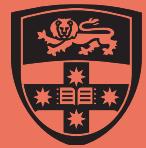


[menti.com 7927 6661](https://menti.com/79276661)

Muddy Card

A screenshot of a web browser window. The address bar shows the URL "saipll.shinyapps.io/student-interface/". The main content area displays the text "Week 2 Muddy Card (TEST1001)" at the top, followed by a large "Thank you!" message. Below this, there are three questions:

- Data collection consent status:** Does Consent
- SID:** 55555555
- Muddy card response:** What is photosynthesis for?



Neural Networks
Recurrent Models
Analysis
Workshop Preview



[menti.com 7927 6661](https://menti.com/79276661)

Muddy Card

Open in a moment, closes at 7:05pm

[https://saipll.shinyapps.io/
student-interface/](https://saipll.shinyapps.io/student-interface/)



If you do not wish to participate in the study, use
the Ed form instead

Go to Ed → Lessons → Muddy Card Lecture 3