



# Quiz instructions

Get a blue or black pen out now, before we hand out the quiz.

**Do not talk or use electronic devices once we start handing out the quiz.**

Leave the quiz facing down until we say it is time to start.

There are two versions of the quiz. The people either side of you must have a different coloured quiz to you.

When we say it is time to stop, hand your quiz to the end of your row.

Make sure you write your name and SID.

This is a closed-book, closed-note quiz. No electronic devices may be used in any way.

Note your responses by completely filling in the relevant circle(s) and square(s): ●

If you make a mistake, put an X over the filled in circle / square: ✗

# COMP 4446 / 5046

## Lecture 6: Models – Encoder-Decoder

*Jonathan K. Kummerfeld*

Semester 1, 2025



[menti.com 4843 3031](https://menti.com/48433031)

Simple

DO YOU HAVE ANY THOUGHTS REGARDING THE PARTICLE ACCELERATOR'S TERTIARY F.E.L. GUIDANCE SYSTEM?

WE CAN'T PUT THE BROKEN PART IN THE MACHINE. IT WOULDN'T SMASH THE RIGHT TINY THINGS TOGETHER. THEN THE MACHINE MIGHT BREAK.

THAT WOULD BE VERY BAD.



I SPENT ALL NIGHT READING SIMPLE.WIKIPEDIA.ORG, AND NOW I CAN'T STOP TALKING LIKE THIS.

[Actually, I think if all higher math professors had to write for the Simple English Wikipedia for a year, we'd be in much better shape academically.]

Source: <https://xkcd.com/547/>

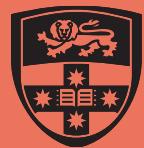


**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

# Contextual Representations

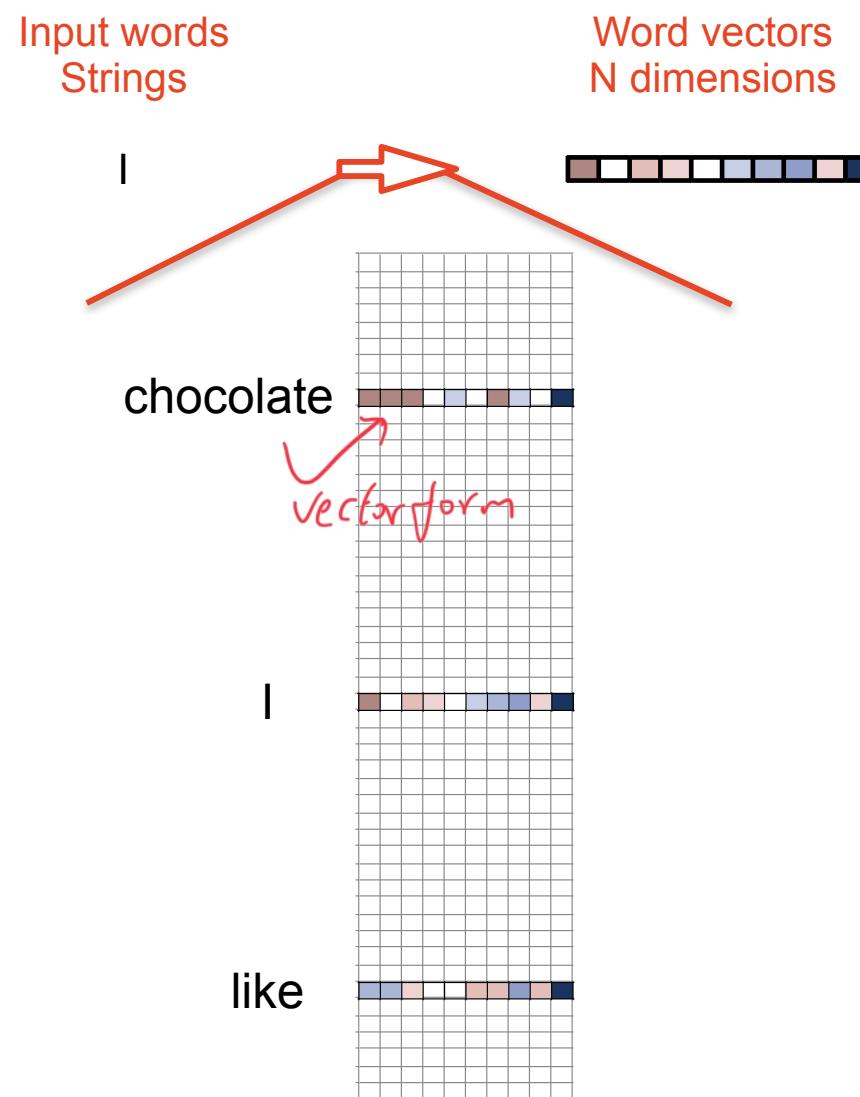


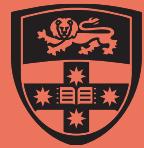
**Contextual Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



menti.com 4843 3031

So far, our word embeddings have been looked up in a table





**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



menti.com 4843 3031

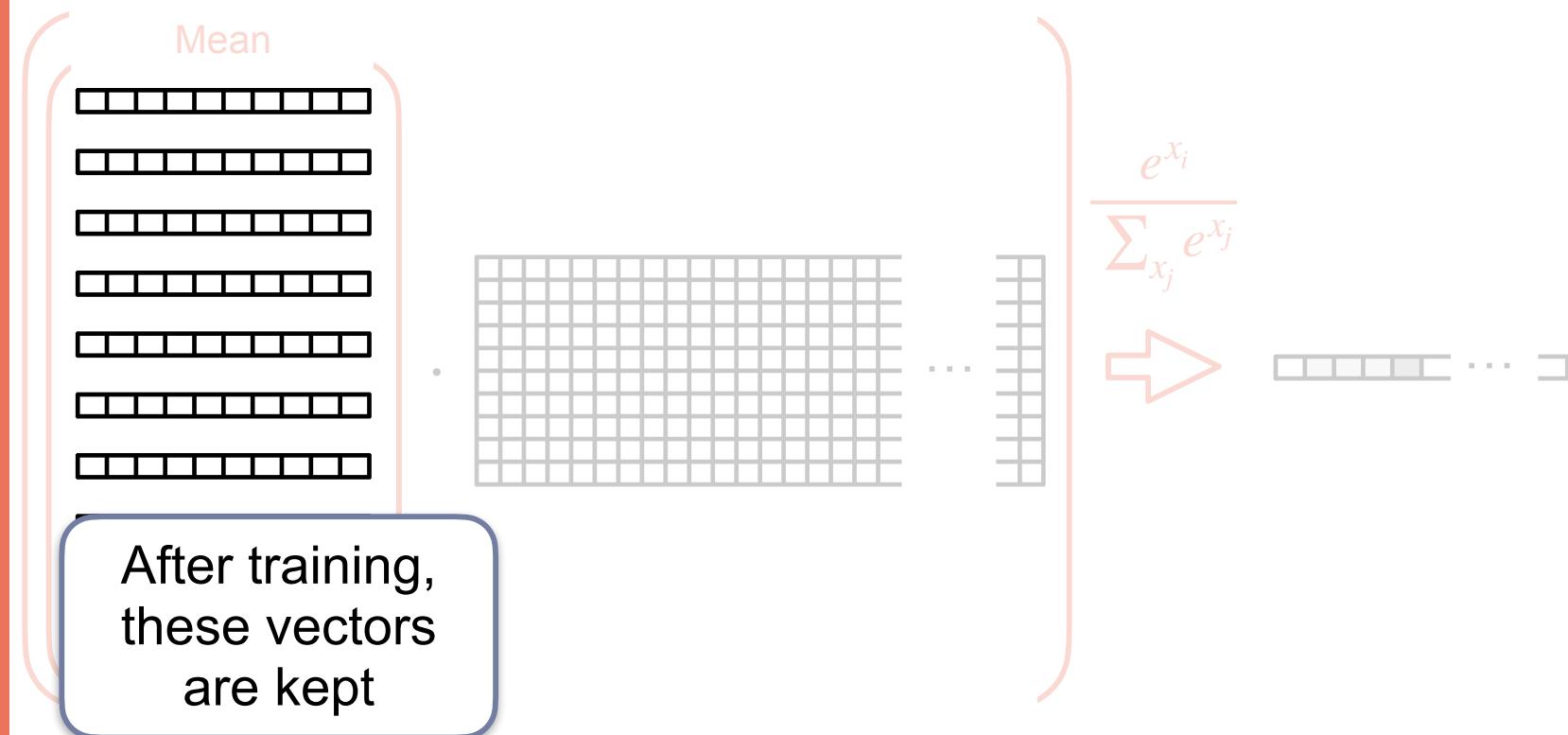
Where does the table come from?

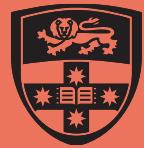
Review ✓

word2vec - Continuous Bag of Words

Input: Context words

Output: One word





**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

Where does the table come from?

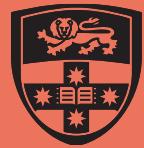
word2vec - SkipGram

Input: One word      After training, these vectors are kept

Output: Set of context words

?

!



**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview

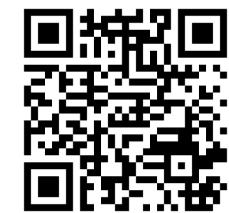


menti.com 4843 3031

Where does the table come from? *look for every word*

GloVe 

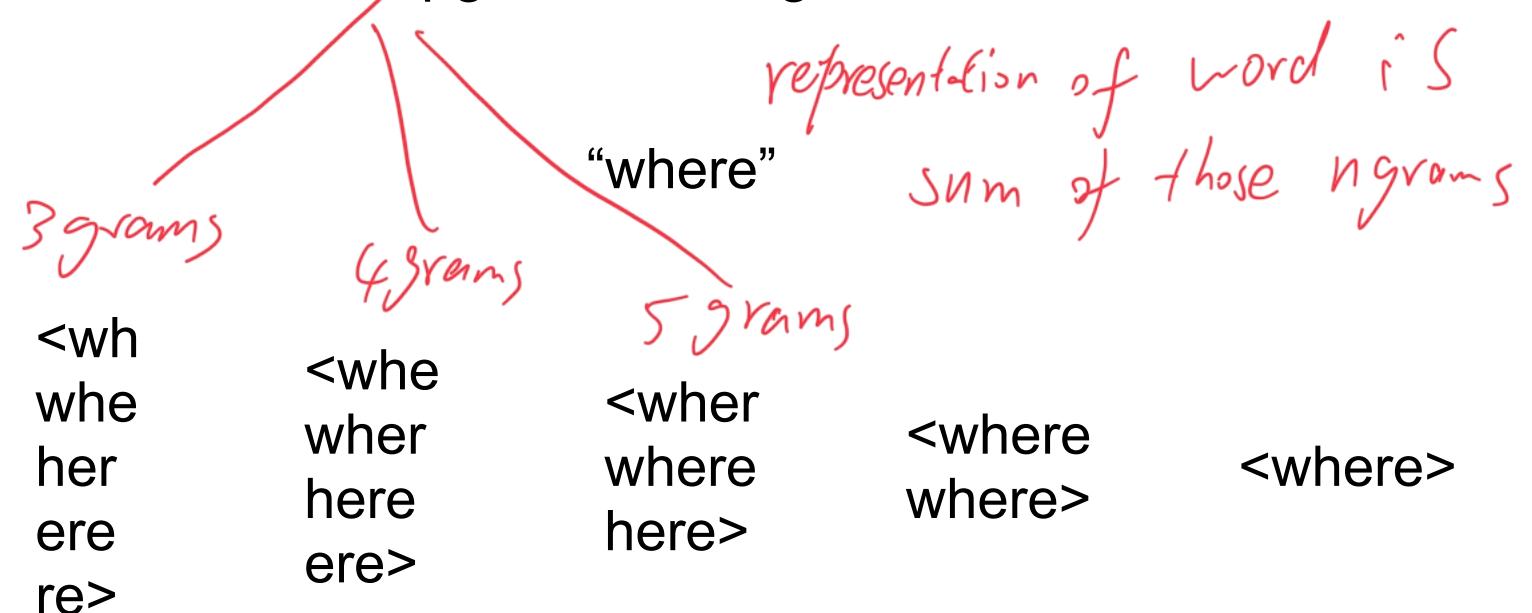
- Calculate co-occurrence statistics
- Form random vectors for each word
- Update vectors so the dot product of two vectors is approximately the co-occurrence value



Where do these come from?

### FastText

- Represent words as a set of character ngrams
- Learn vectors for ngrams
- Words are the sum of vectors for their ngrams
- Use word2vec skipgram learning



Ex. maybe never see a word before, but see a gram before.



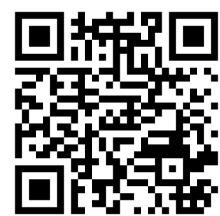
## Contextual Representations

Encoder-Decoder

Tokenisation

Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

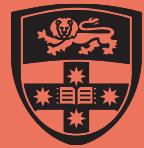
# What if my data is not the same as the data used for training?

Standard sources of data:

- Websites
- Books
- News
- Research literature

Not:

- Medical records
- Internal company documents
- Email
- Instant messaging
- Text messages



**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



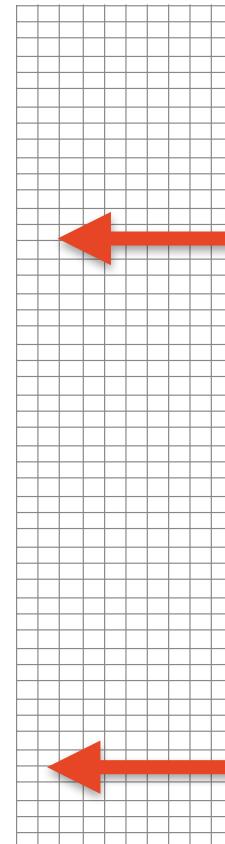
menti.com 4843 3031

训练和测试的数据不同

What if my data is not the same as the data used for training?

Fine-tuning - Update the embeddings for your task

?  
:



Method 1: retrain?

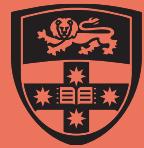
word2vec



Method 2: update model

Task specific  
model

Method 3: back prop



**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## What about word senses?



NN

NN



VB

DT

NN



NN

VBZ

IN

DT

NN

Time

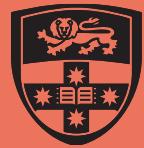
flies

like

an

arrow



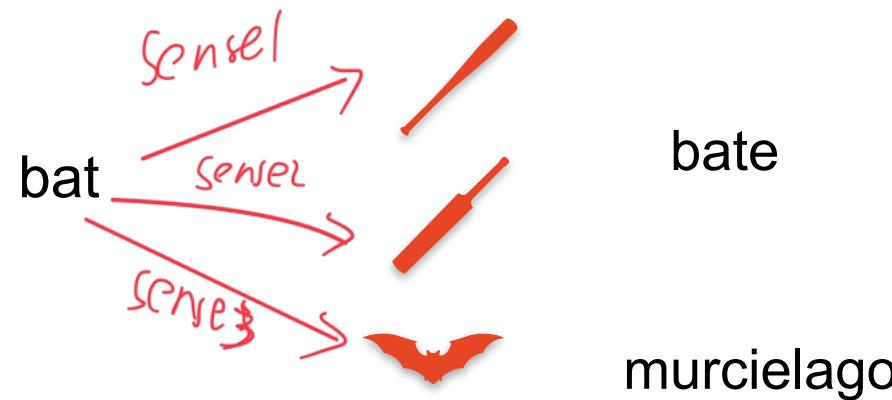


**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



menti.com 4843 3031

## What about word senses?



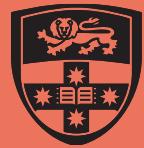
Which of those flights serve breakfast?

Does Air France serve Philadelphia? *flight*

?Does Air France serve breakfast and Philadelphia?

*weird*,

Jurafsky and Martin, Appendix G



**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview

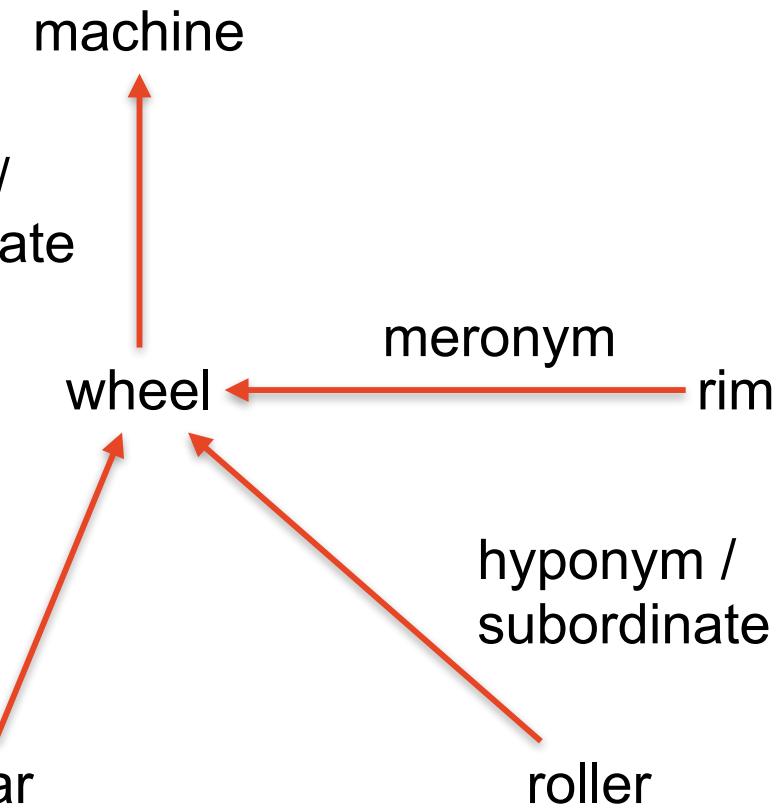


menti.com 4843 3031

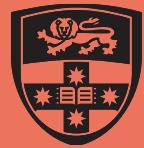
## What about word senses?

### WordNet

a simple machine  
consisting of a  
circular frame ...



<http://wordnetweb.princeton.edu/perl/webwn>



## Contextual Representations

Encoder-Decoder

Tokenisation

Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

Train a model with multiple word vectors, one per sense?

Challenge: data

SemCor - 226,036 words

Others in the 1,000 - 10,000 range



**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



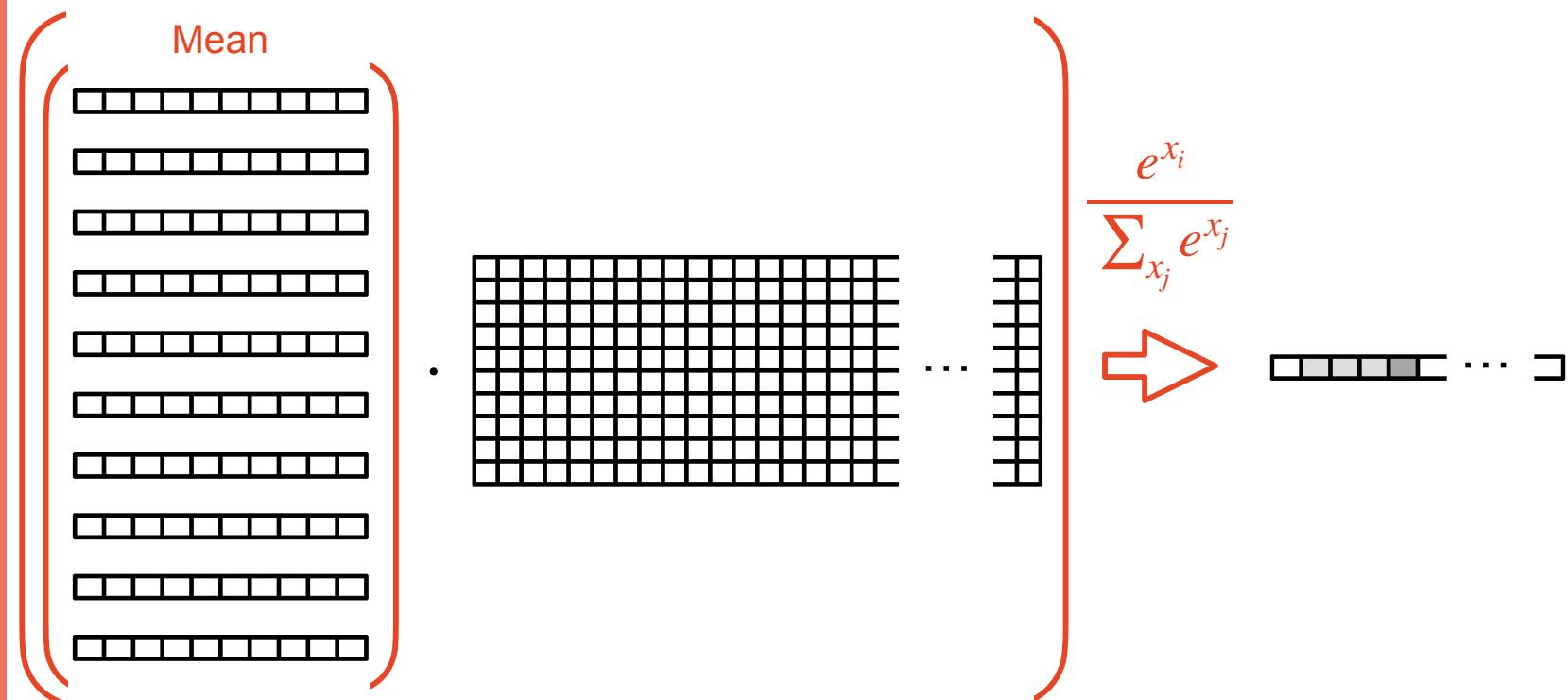
[menti.com 4843 3031](https://menti.com/48433031)

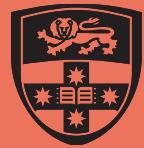
## Train a model with **contextual** representations

word2vec - Continuous Bag of Words

Input: Context words

Output: One word





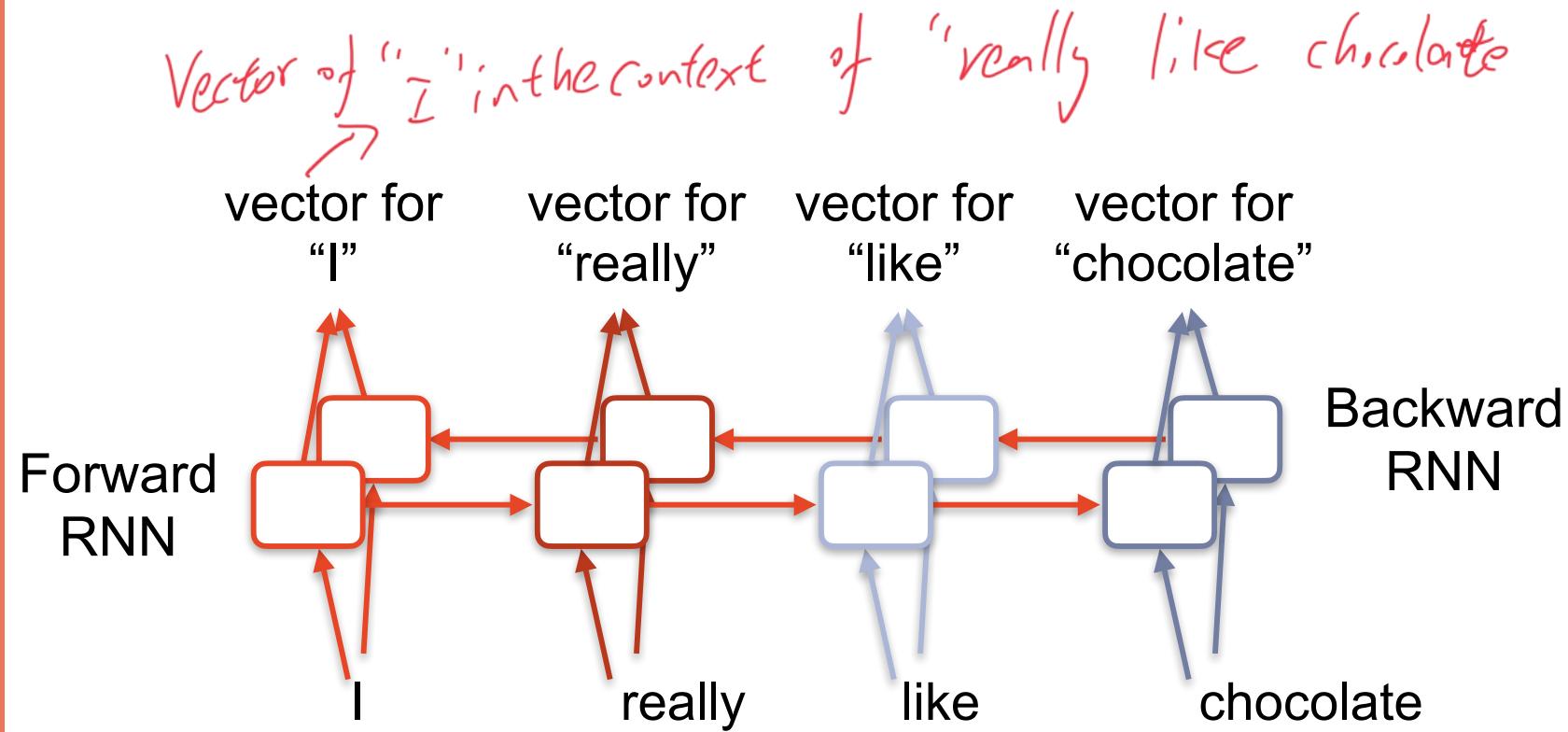
## Contextual Representations

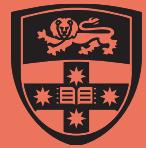
- Encoder-Decoder
- Tokenisation
- Attention
- Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Form contextual representations using an RNN





## Contextual Representations

Encoder-Decoder

Tokenisation

Attention

Workshop Preview

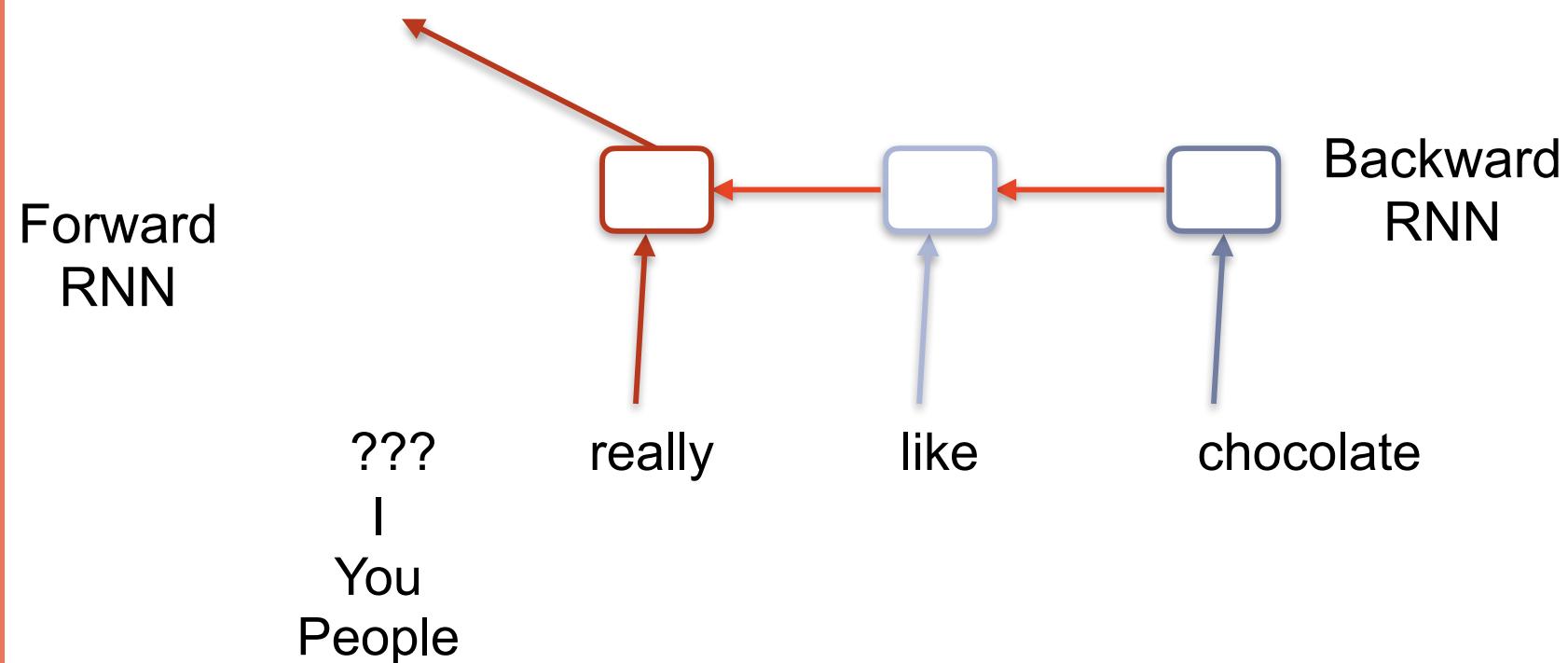


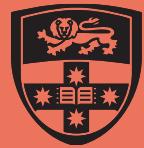
[menti.com 4843 3031](https://menti.com/48433031)

# How do we train the model?

Input: Context words

Output: One word





**Contextual Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview

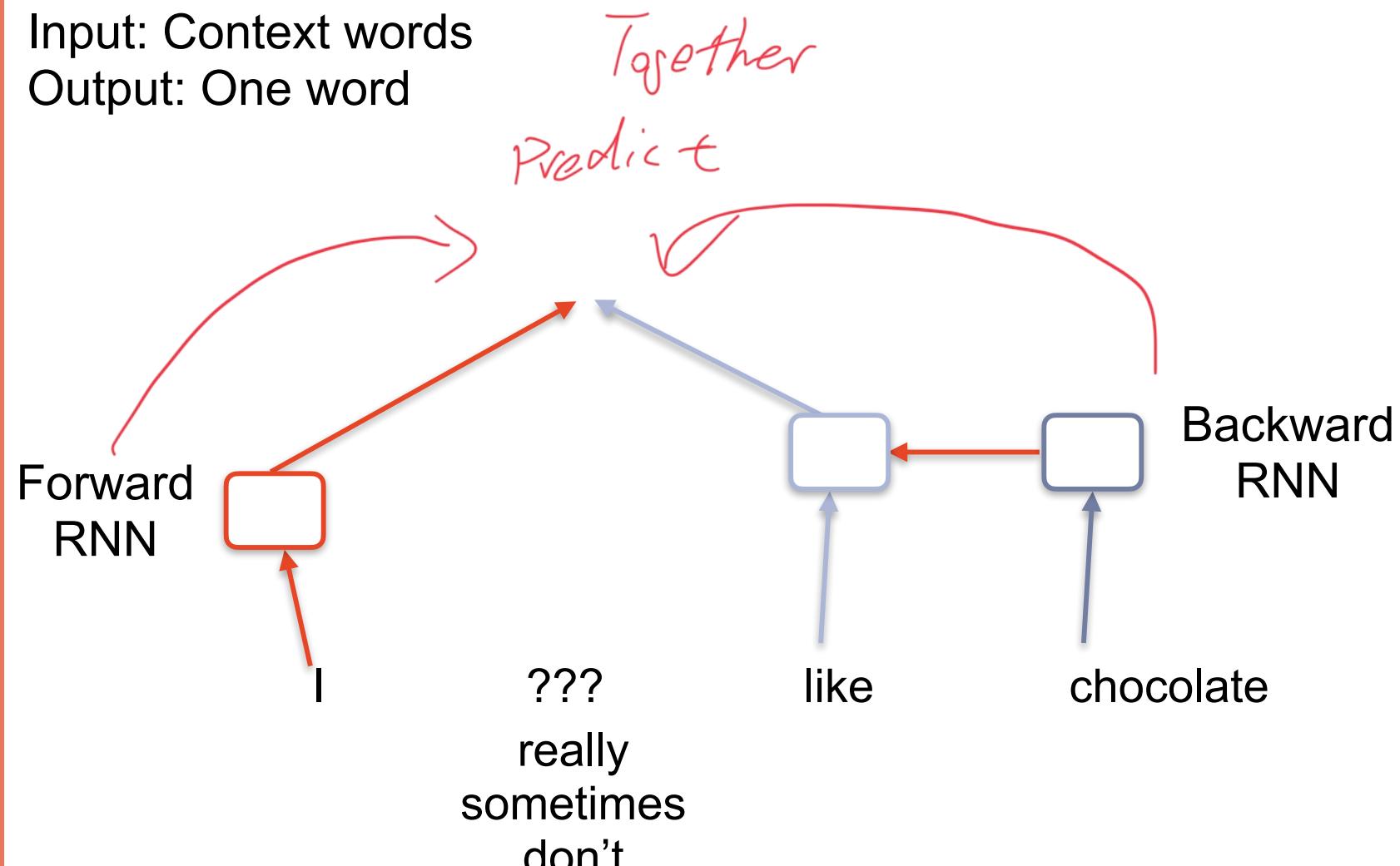


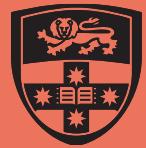
[menti.com 4843 3031](https://menti.com/48433031)

How do we train the model?

Input: Context words

Output: One word





**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview

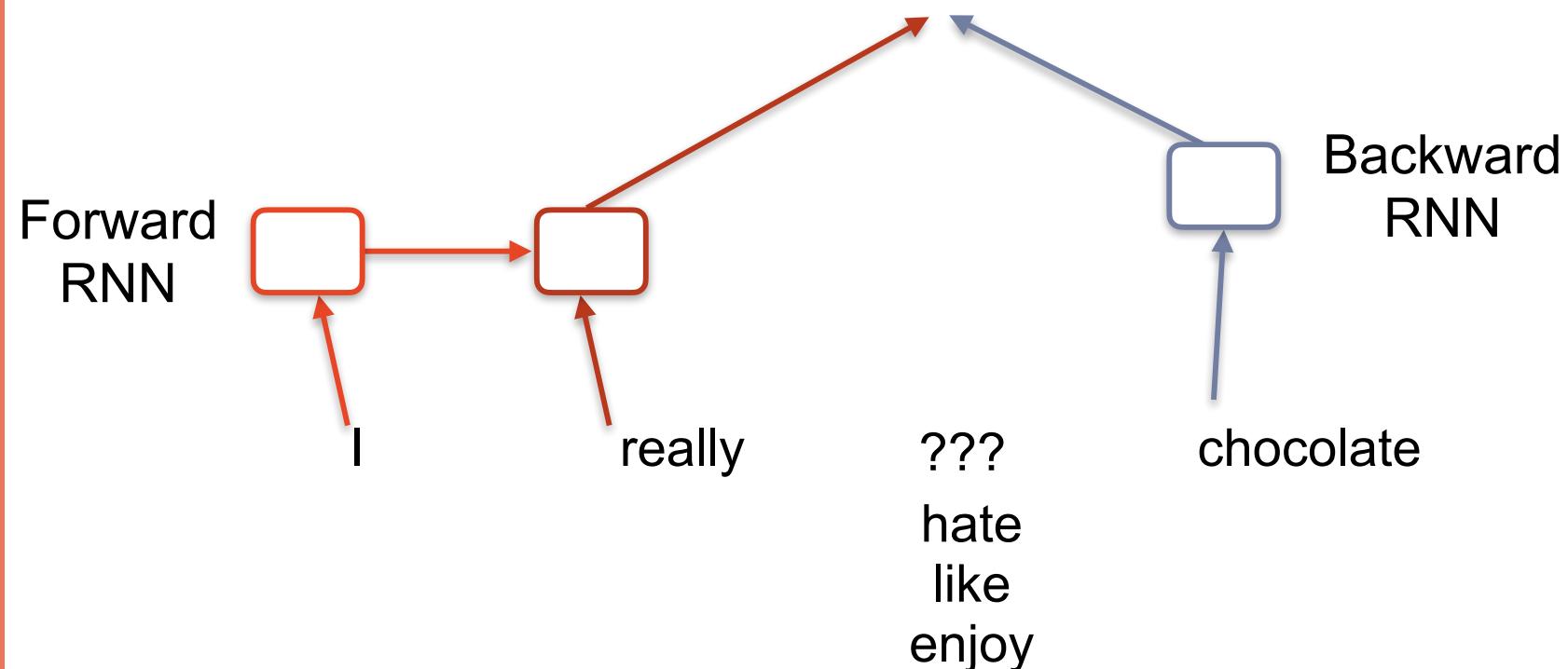


[menti.com 4843 3031](https://menti.com/48433031)

How do we train the model?

Input: Context words

Output: One word





## Contextual Representations

Encoder-Decoder

Tokenisation

Attention

Workshop Preview

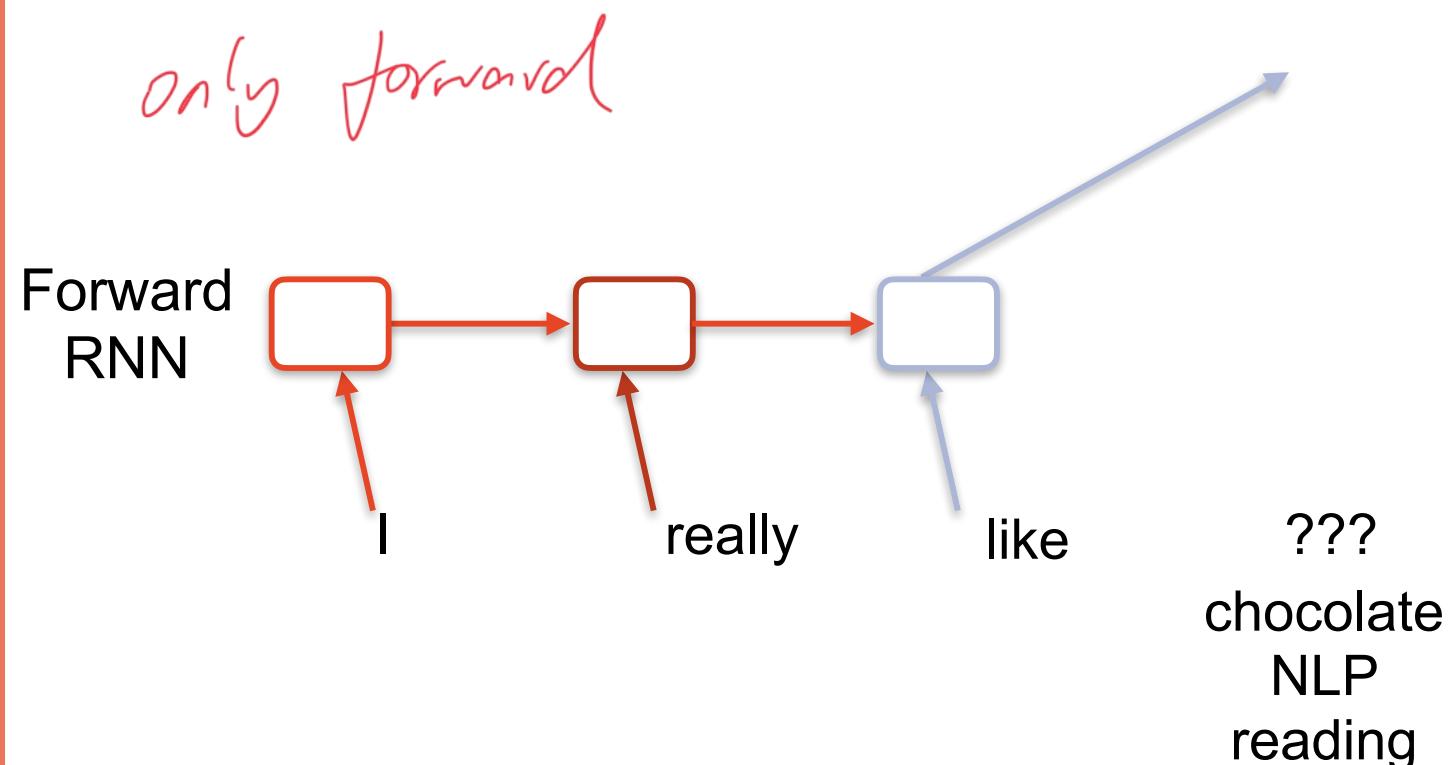


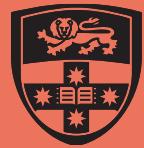
[menti.com 4843 3031](https://menti.com/48433031)

# How do we train the model?

Input: Context words

Output: One word





**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview

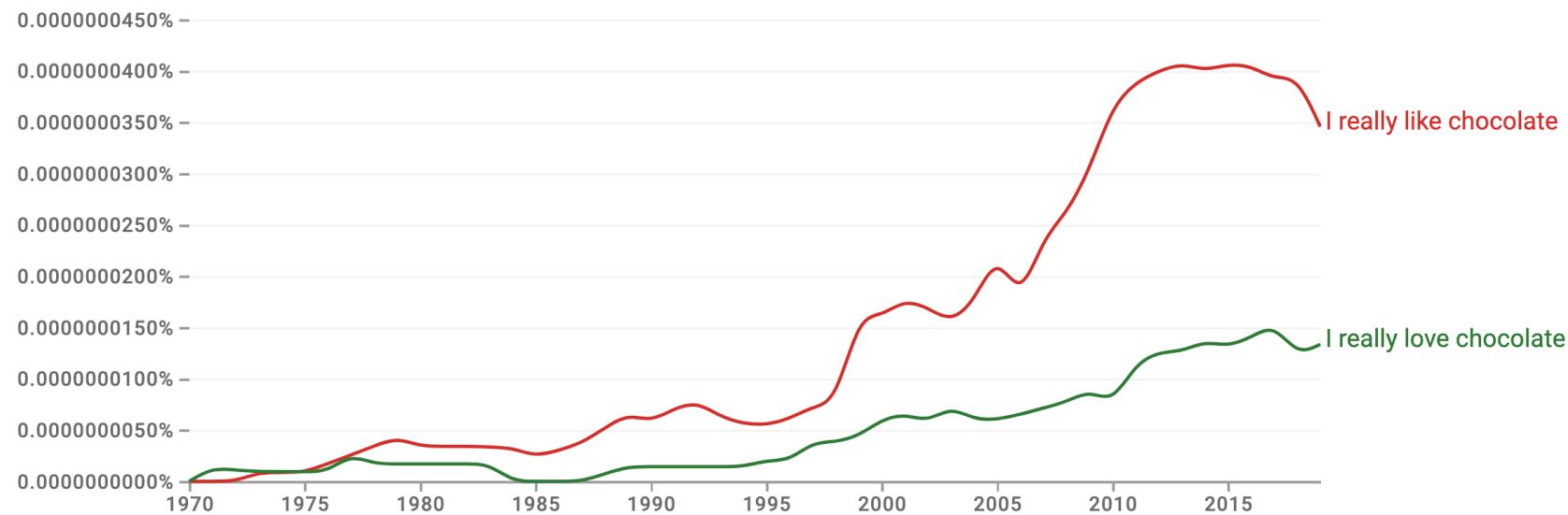


[menti.com 4843 3031](https://menti.com/48433031)

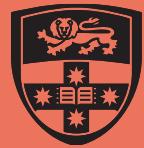
[aside - Google shows that books agree!]

Counts in books since 1970 of:

“I really \* chocolate”



<https://books.google.com/ngrams>



**Contextual Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview

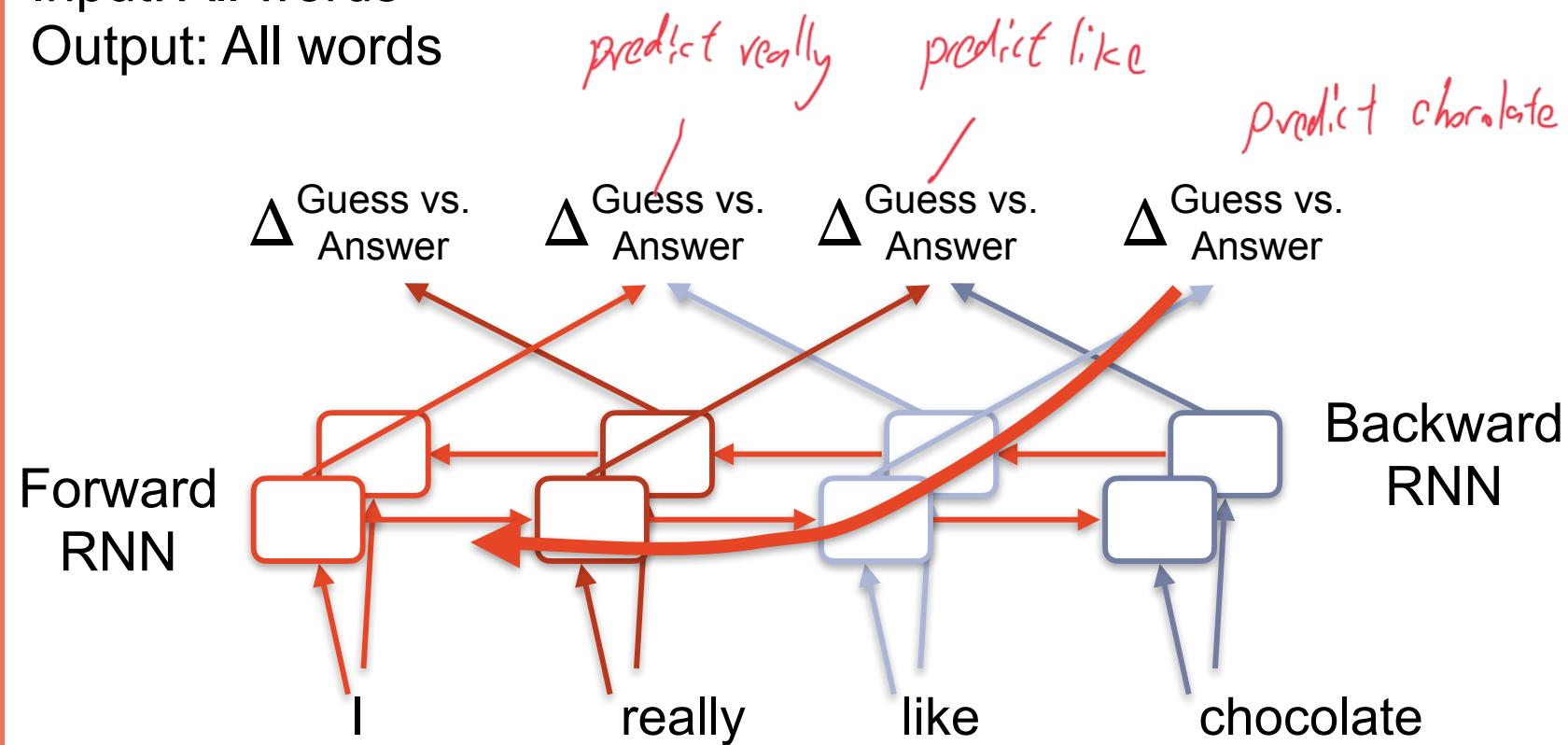


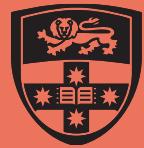
menti.com 4843 3031

We can train on multiple words at once

Input: All words

Output: All words





**Contextual Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



menti.com 4843 3031

## Major turning point in NLP

Earliest use I know of - 2015

No major awareness at first (only slight improvements)

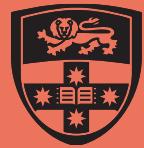
2018 - ELMo

“Deep contextualized word representations”

TASK	PREVIOUS SOTA	OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/RELATIVE)
SQuAD	Lai et al. (2015)	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.0	88.7 ± 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	84.6	3.2 / 17.2%
Coref	Lee et al. (2017)	67.0	70.1	3.2 / 9.8%
NER	Peters et al. (2017)	57.0	60.1	2.06 / 21%
SST-5	MG et al. (2017)	59.1	62.1	3.3 / 6.8%

[https://sesameworkshop.org/our-work/shows/  
sesame-street/sesame-street-characters/](https://sesameworkshop.org/our-work/shows/sesame-street/sesame-street-characters/)



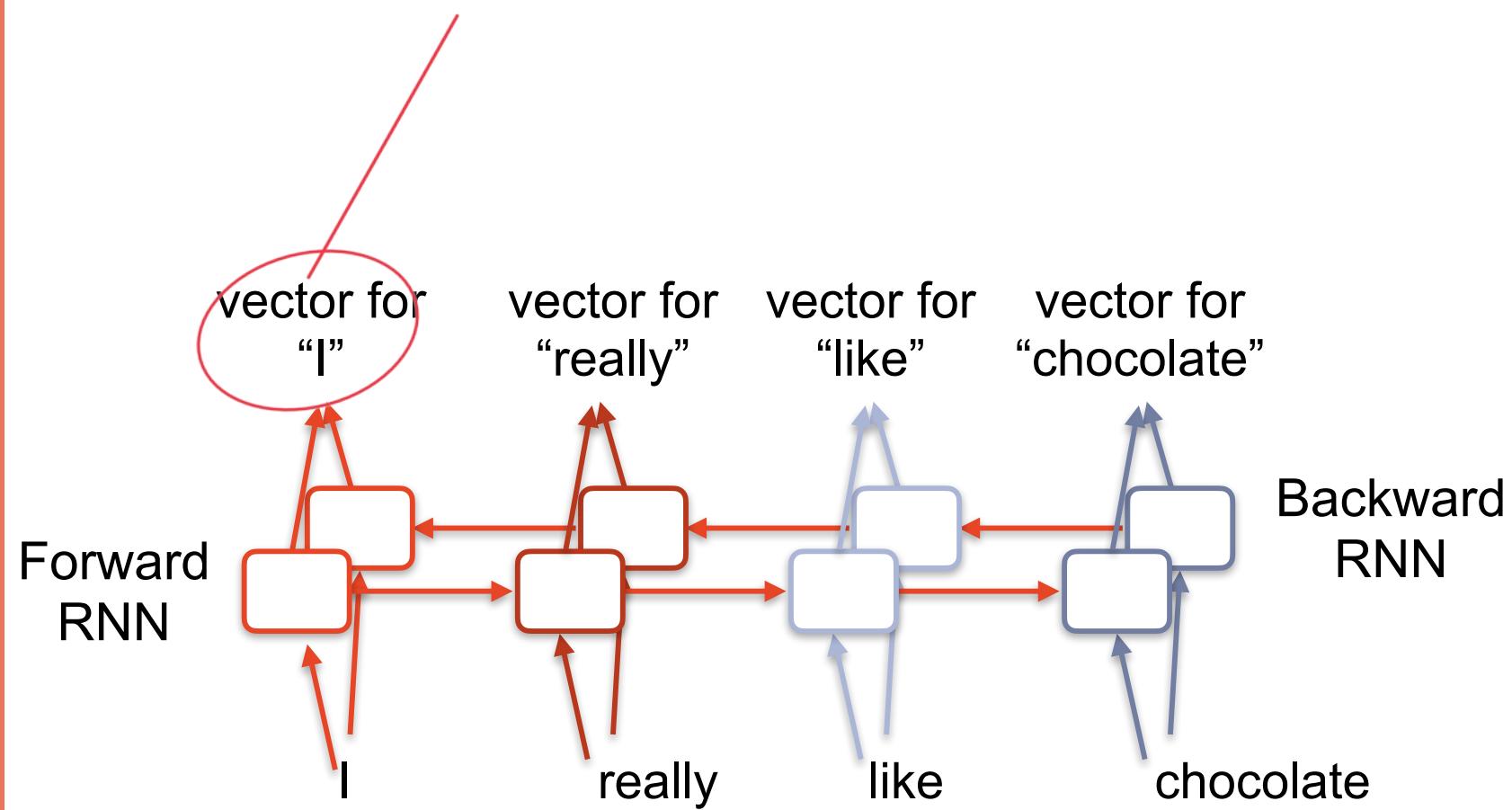


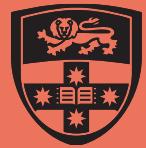
**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

How are these used?





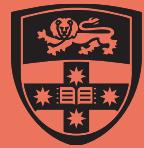
COMP 4446 / 5046  
Lecture 6, 2025

# What are they learning?

**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)



COMP 4446 / 5046  
Lecture 6, 2025

**Contextual  
Representations**  
Encoder-Decoder  
Tokenisation  
Attention  
Workshop Preview



[menti.com 6165 8383](https://menti.com/61658383)



26

The LSTM is just one possible model...



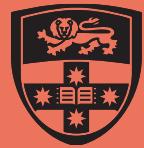


Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

# Encoder-Decoder



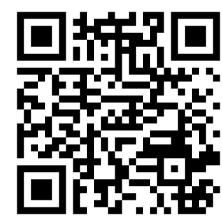
Contextual  
Representations

**Encoder-Decoder**

Tokenisation

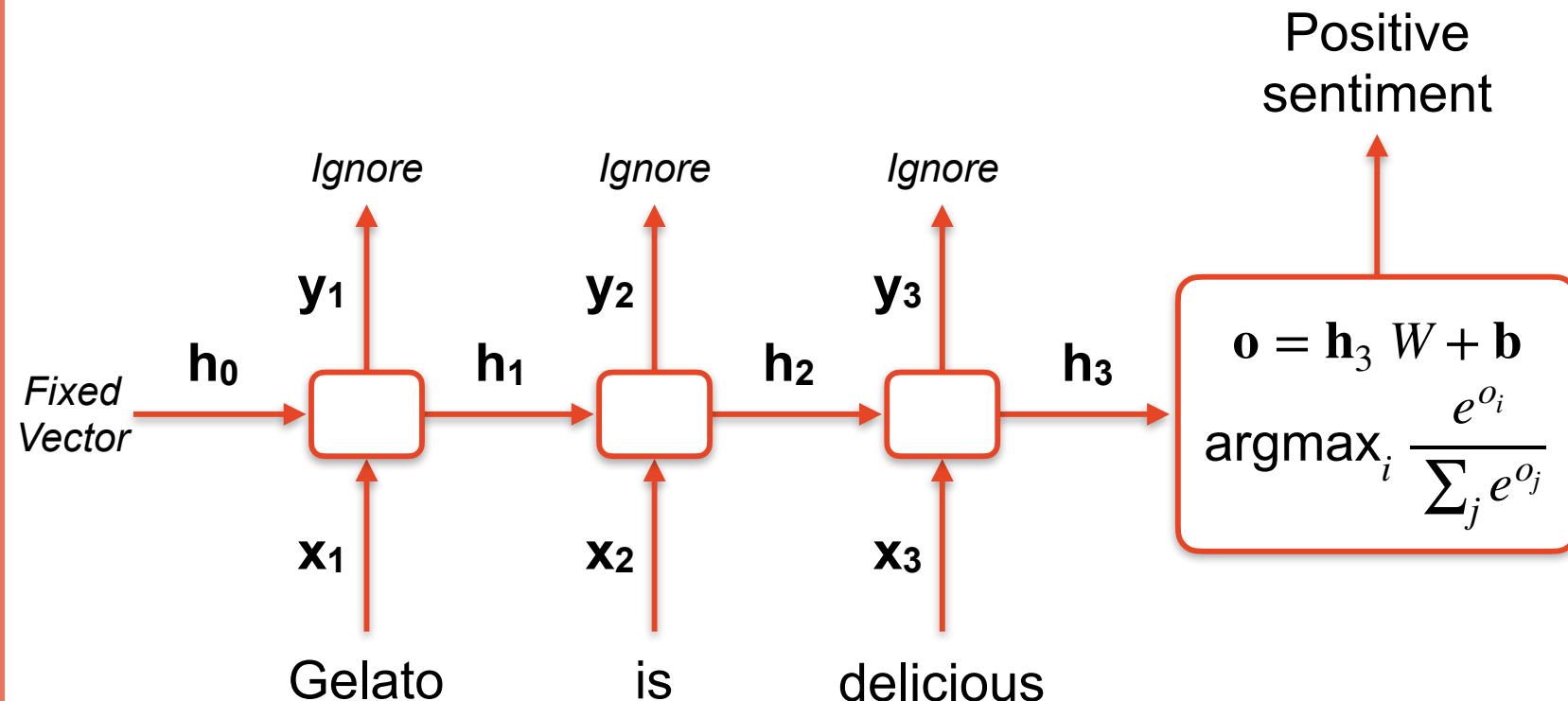
Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

RNNs can be used as an **accepter** or encoder





Contextual

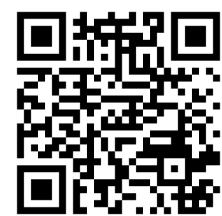
Representations

**Encoder-Decoder**

Tokenisation

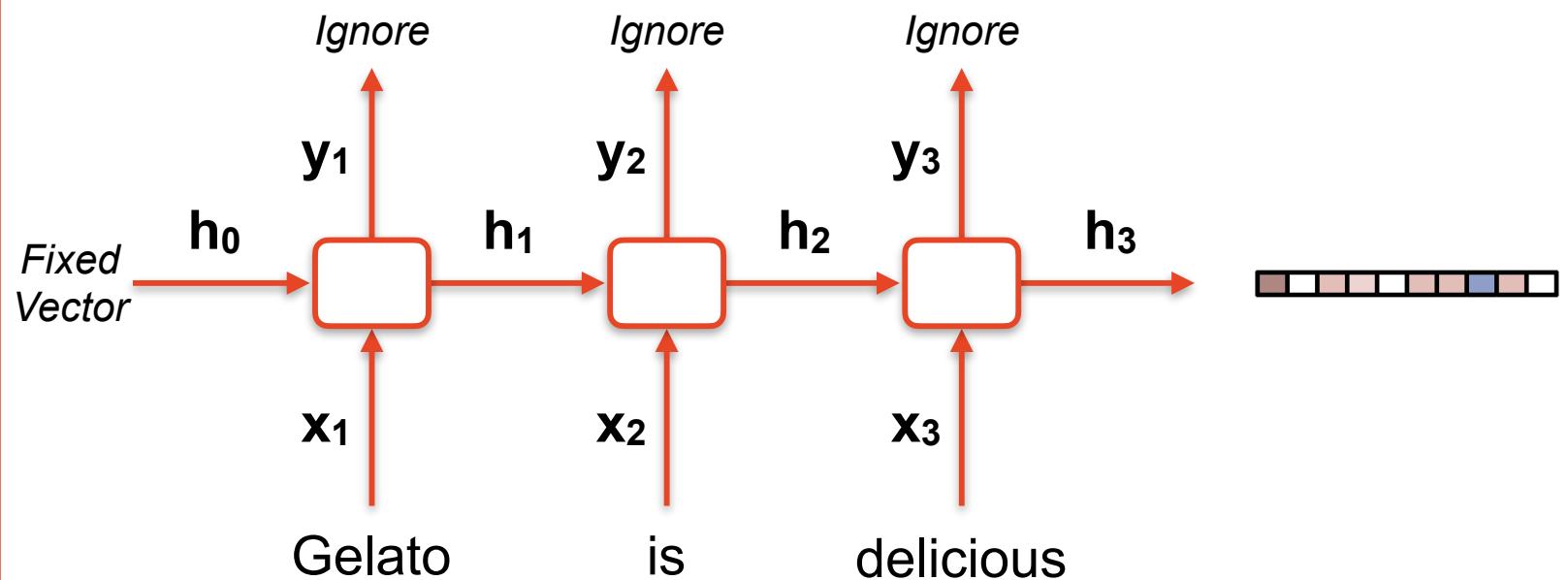
Attention

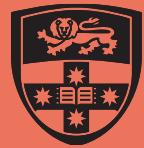
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

RNNs can be used as an accepter or **encoder**





Contextual

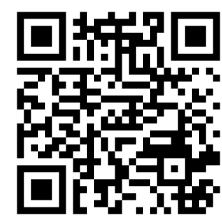
Representations

**Encoder-Decoder**

Tokenisation

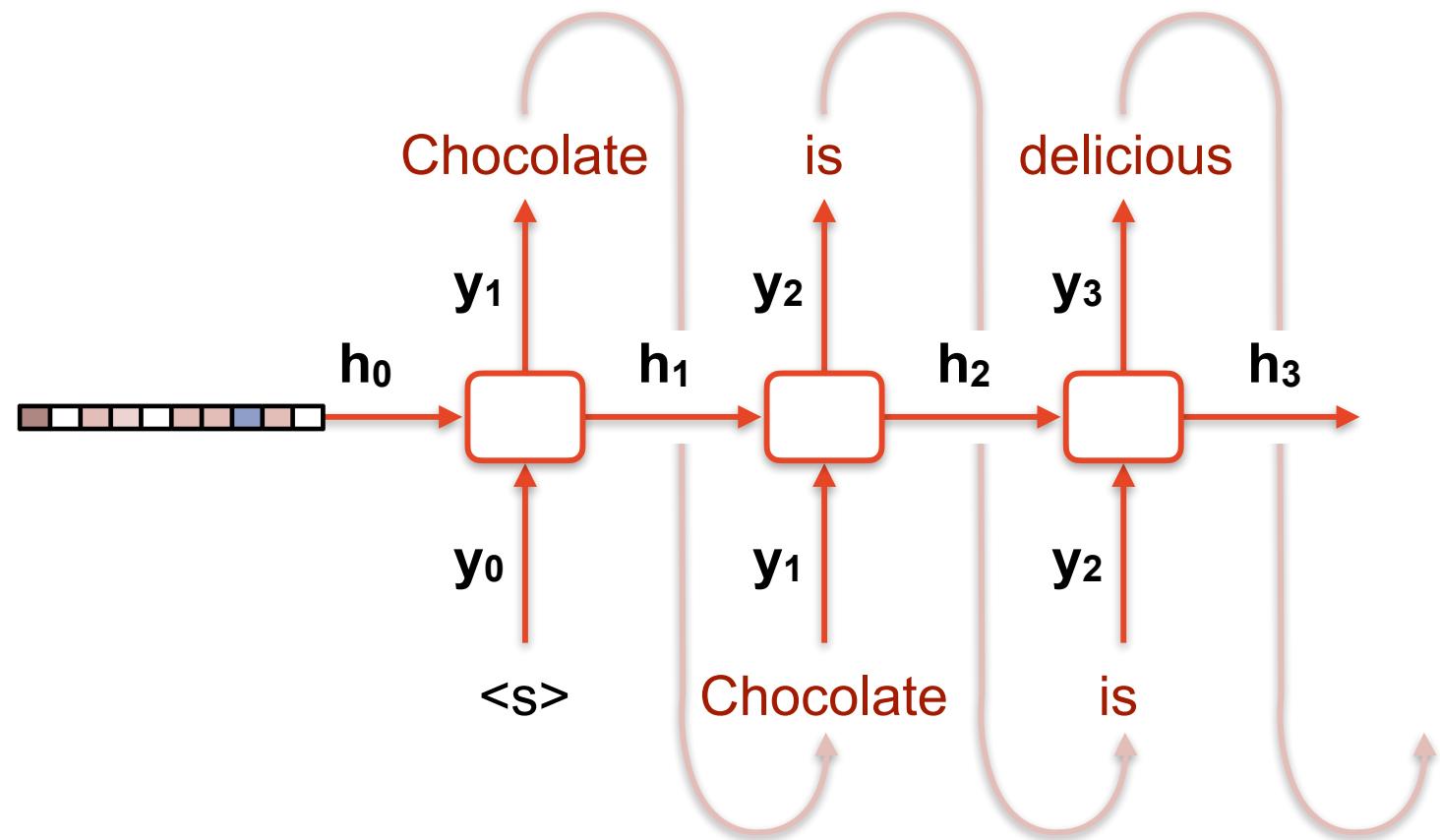
Attention

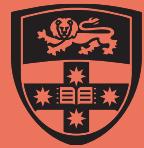
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

What is a decoder?





Contextual

Representations

**Encoder-Decoder**

Tokenisation

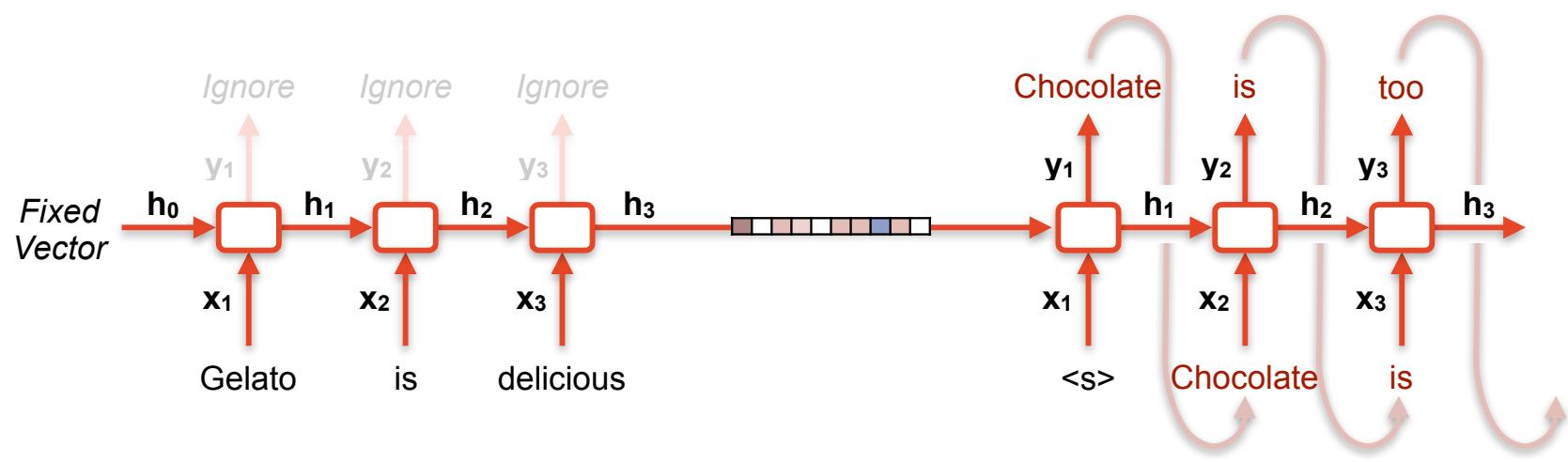
Attention

Workshop Preview

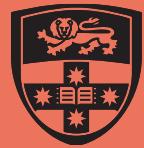


[menti.com 4843 3031](https://menti.com/48433031)

We can put these two pieces together



This is also called a  
'sequence to sequence'  
model



Contextual  
Representations

## Encoder-Decoder

Tokenisation

Attention

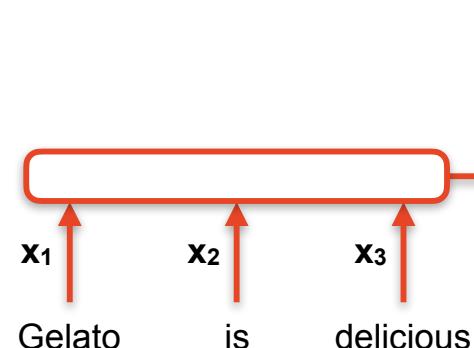
Workshop Preview



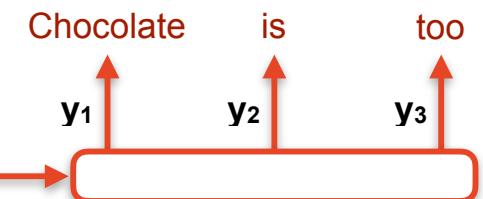
[menti.com 4843 3031](https://menti.com/48433031)

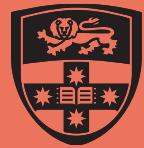
We can put these two pieces together

Encoder



Decoder





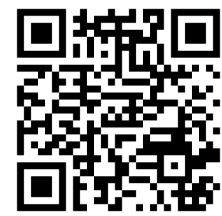
Contextual  
Representations

## Encoder-Decoder

Tokenisation

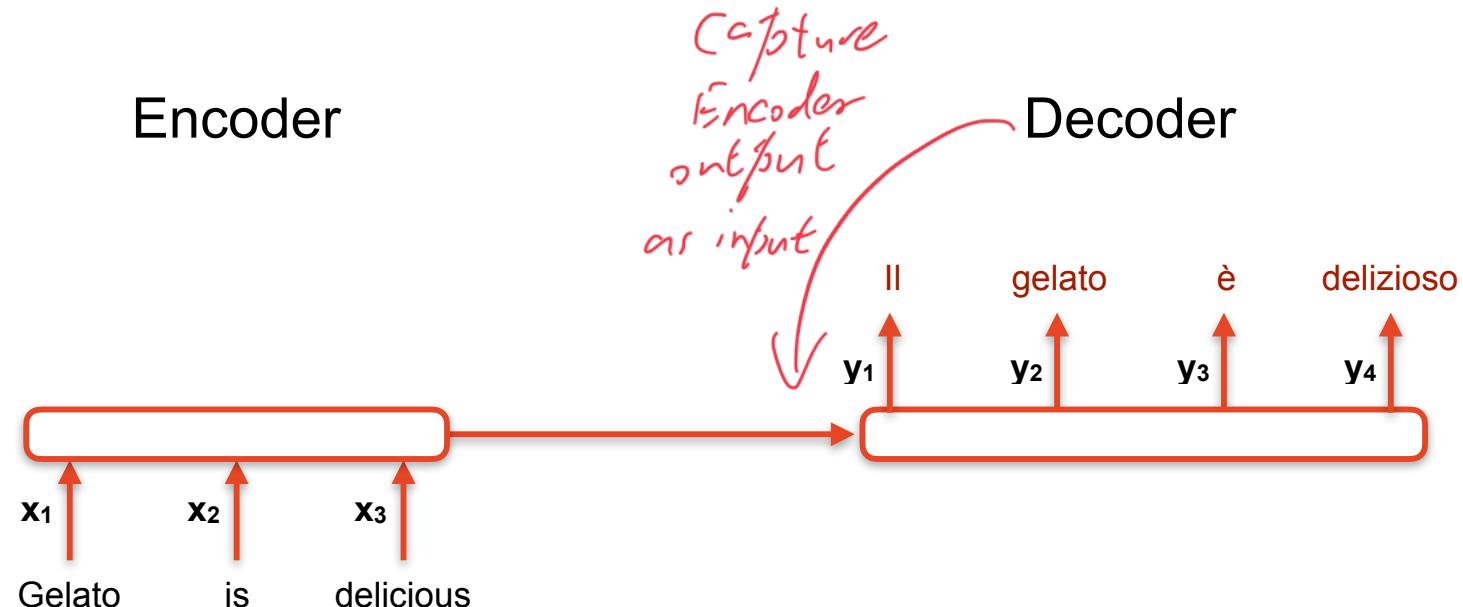
Attention

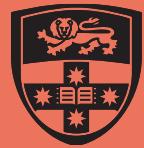
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

The first successful application was Machine Translation





Contextual  
Representations

**Encoder-Decoder**

Tokenisation

Attention

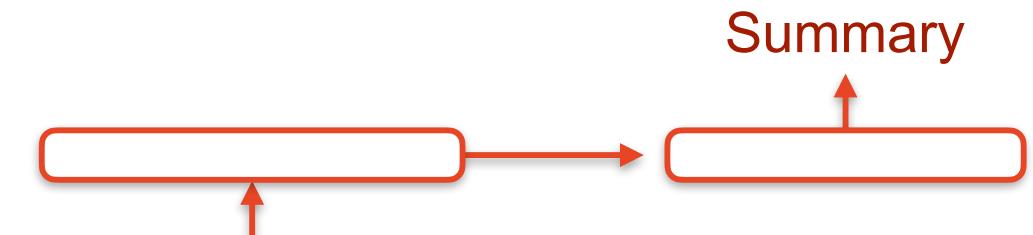
Workshop Preview



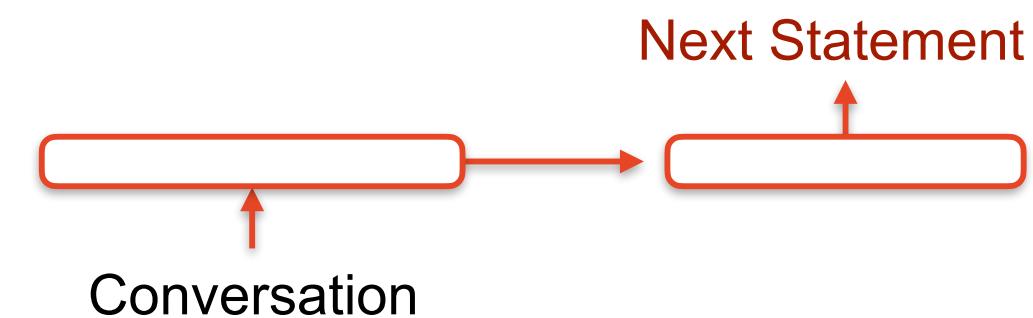
[menti.com 4843 3031](https://menti.com/48433031)

But it can be used for a wide range of applications

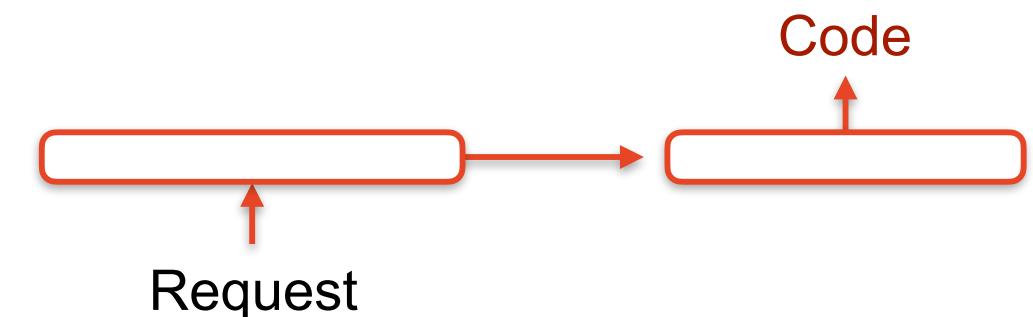
Summarisation



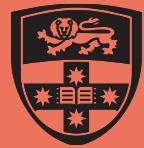
Dialogue



Code Generation



...



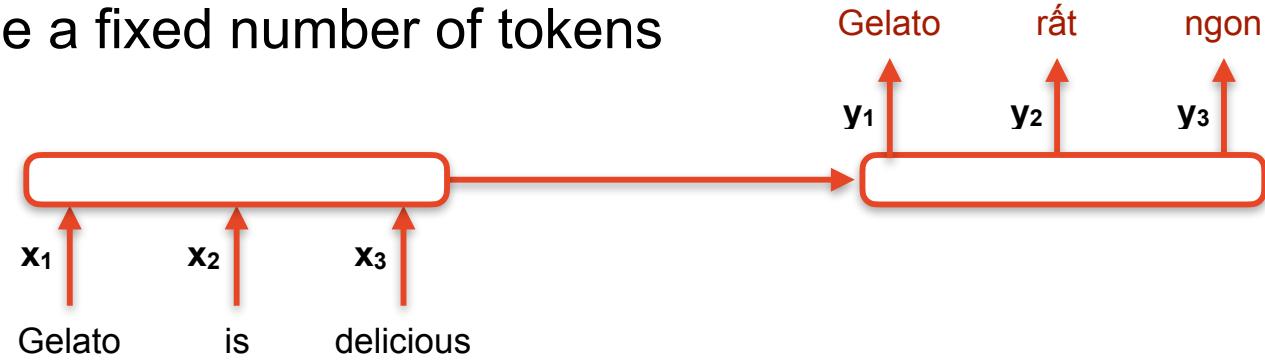
Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



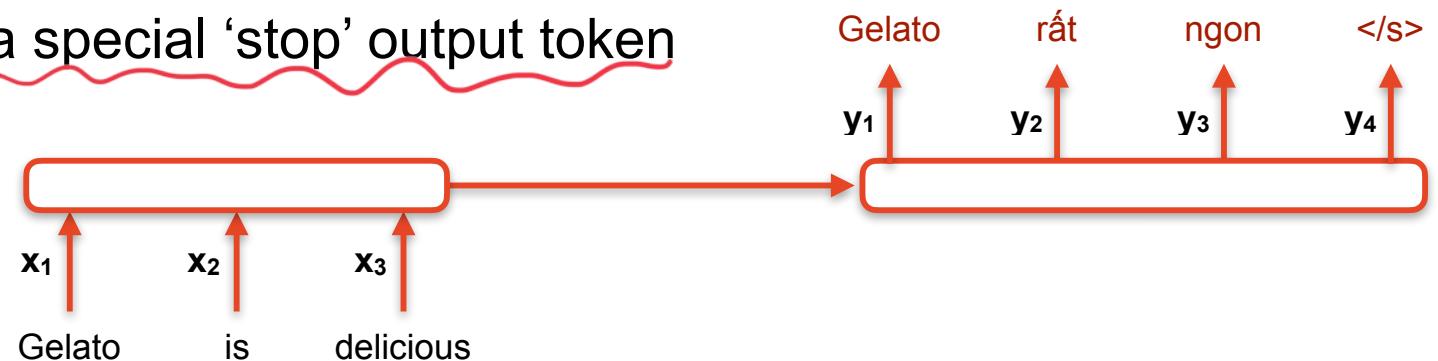
[menti.com 4843 3031](https://menti.com/48433031)

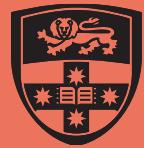
When do you stop decoding?

Option 1:  
Choose a fixed number of tokens



Option 2:  
Have a special 'stop' output token





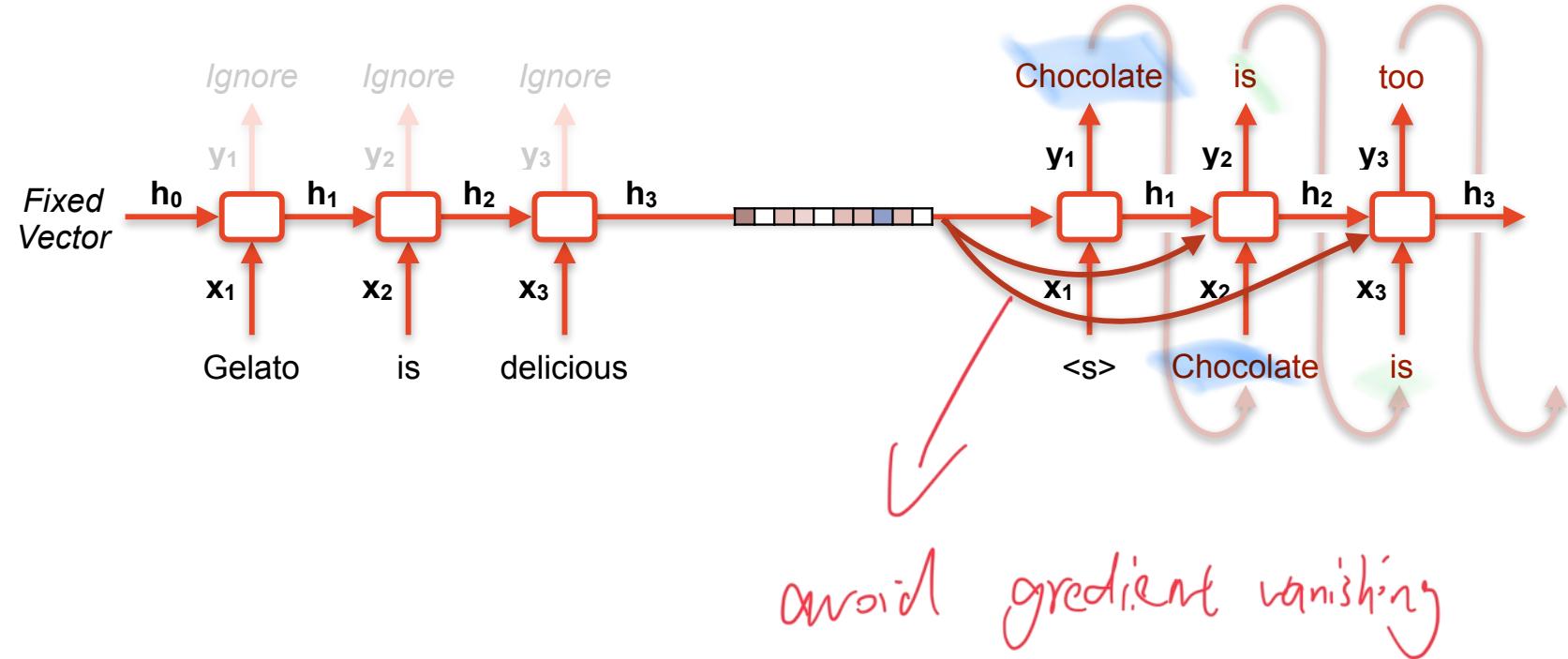
Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview

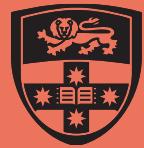


menti.com 4843 3031

In practise, the structure is slightly more complex

Encoding / Context is passed in at each step





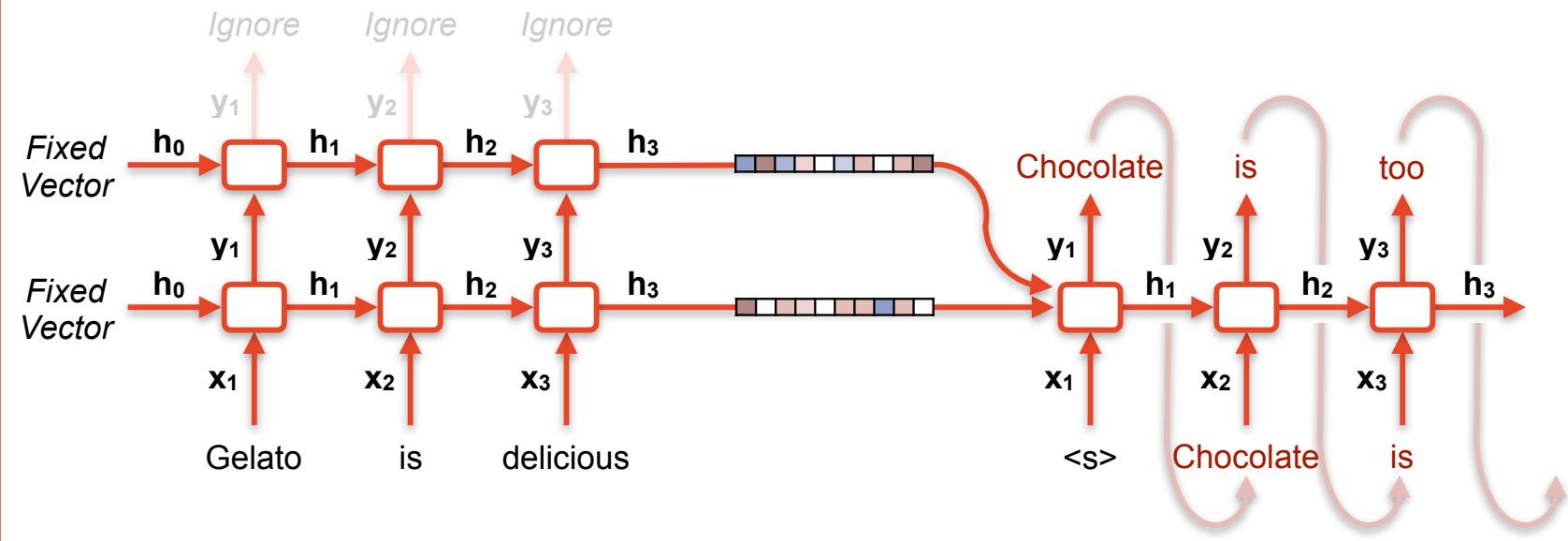
Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

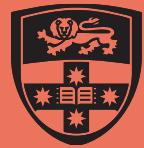
In practise, the structure is slightly more complex

## Multi-layer models



Can also use **bi-directional models**

- Encoder, easy to set up
- Decoder, trickier, but also possible at the cost of computation



Contextual  
Representations

**Encoder-Decoder**

Tokenisation

Attention

Workshop Preview

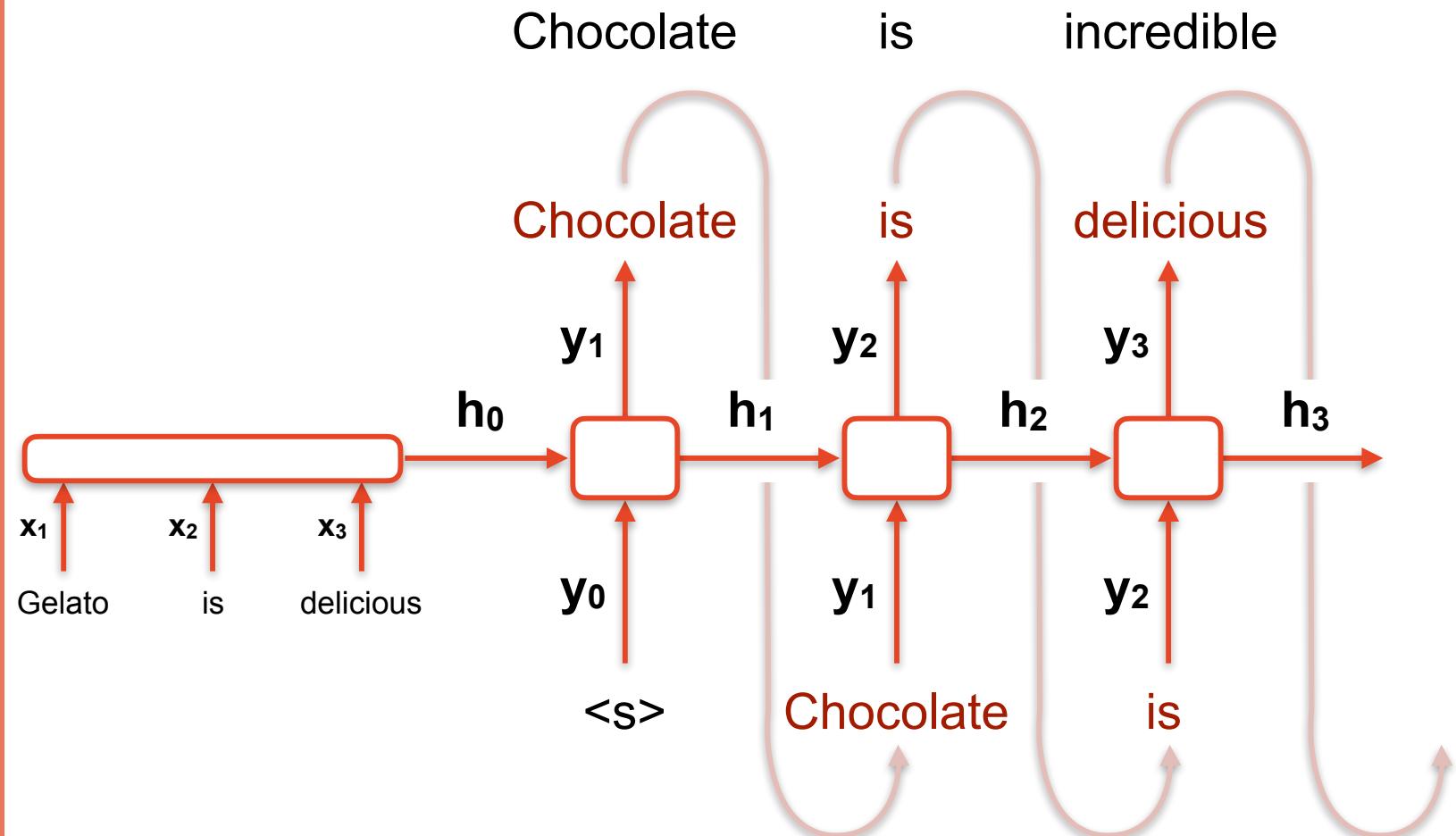


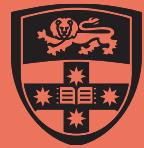
[menti.com 4843 3031](https://menti.com/48433031)

How do you train the decoder?

*back prop*

△ Guess vs.  
Answer





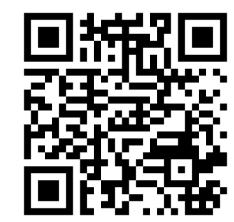
Contextual  
Representations

**Encoder-Decoder**

Tokenisation

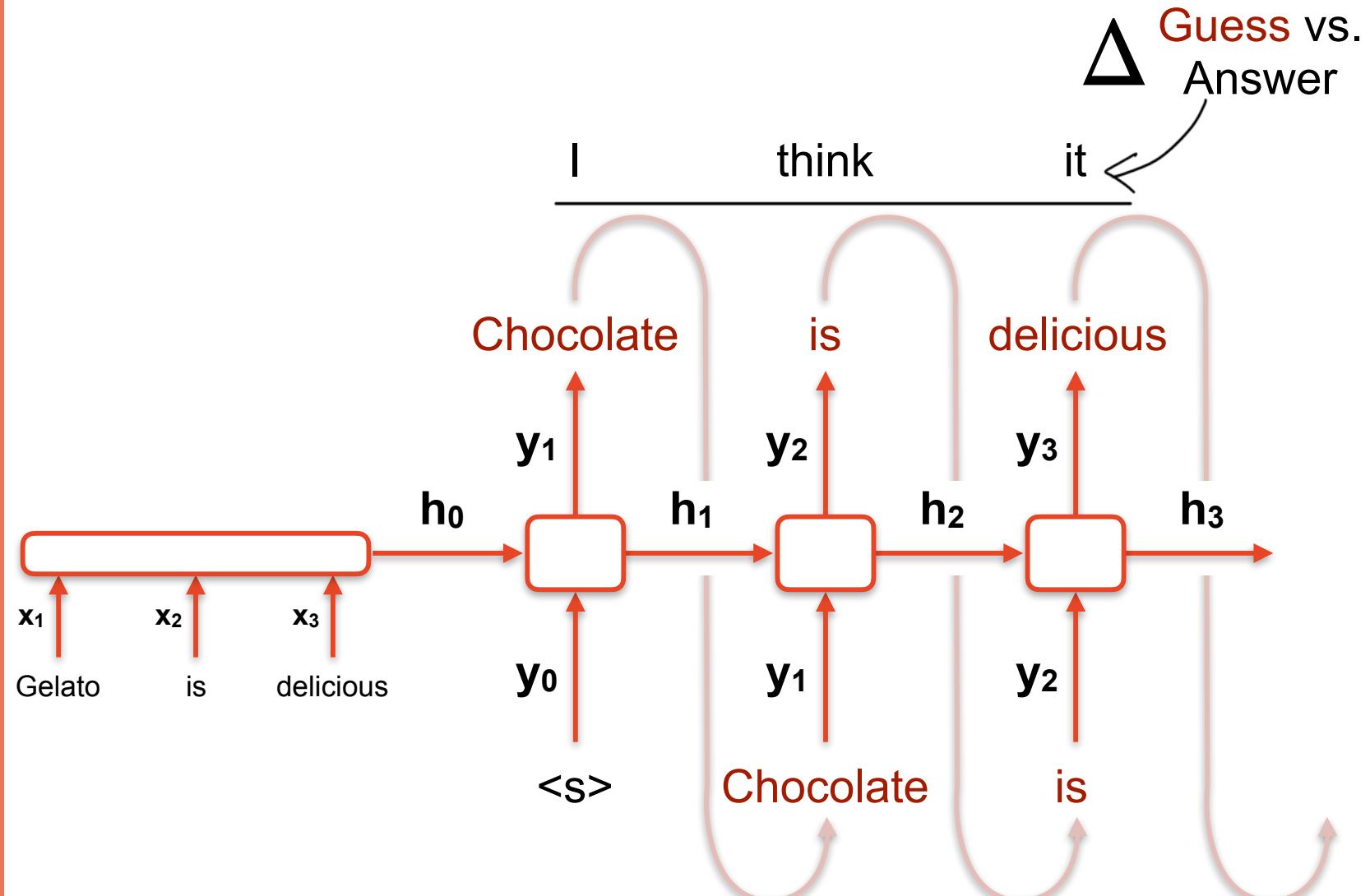
Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

How do you train the decoder?





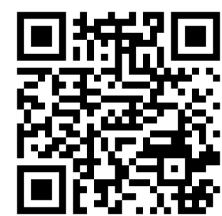
Contextual  
Representations

**Encoder-Decoder**

Tokenisation

Attention

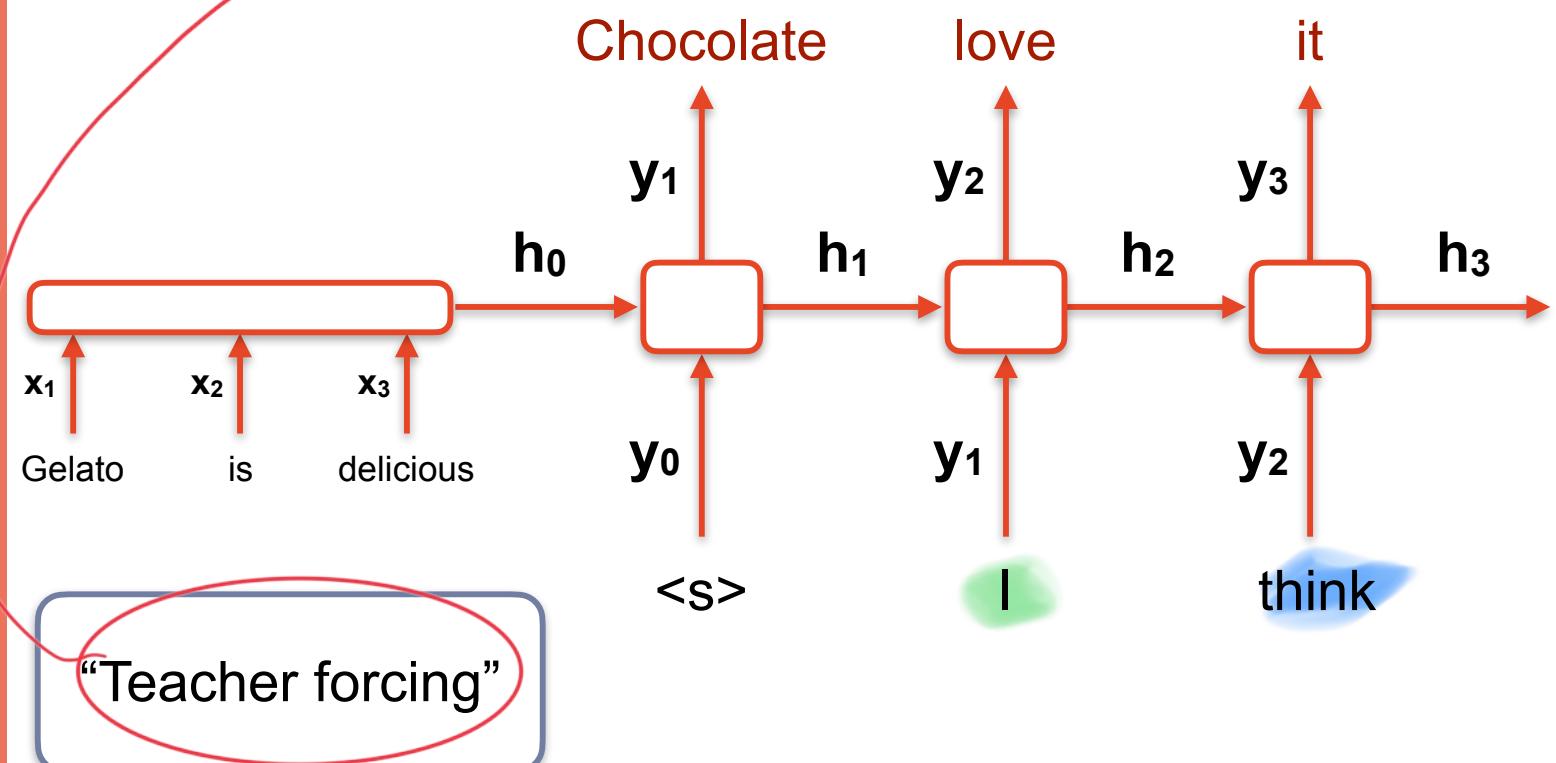
Workshop Preview



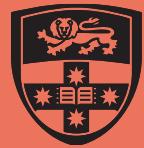
menti.com 4843 3031

How do you train the decoder?

If totally wrong,  
then just force teach  
it



△ Guess vs.  
Answer



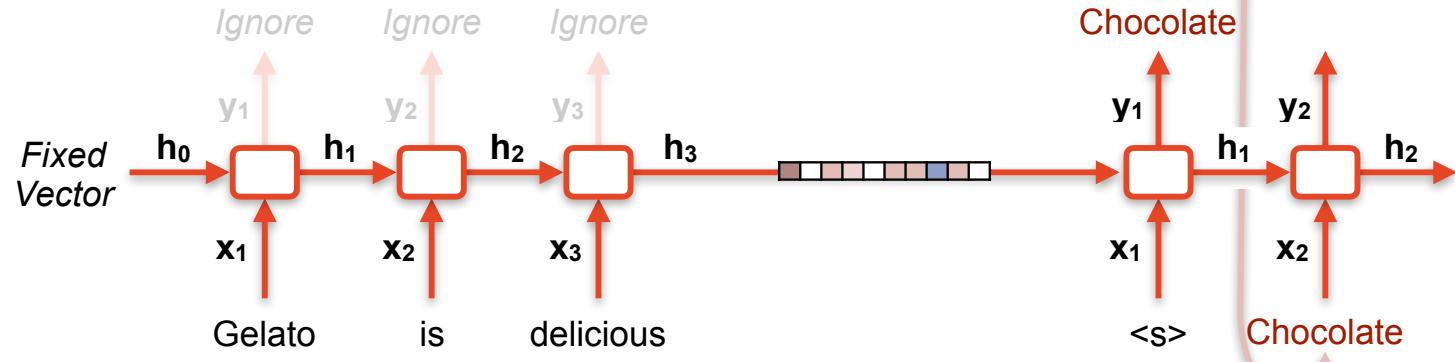
Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



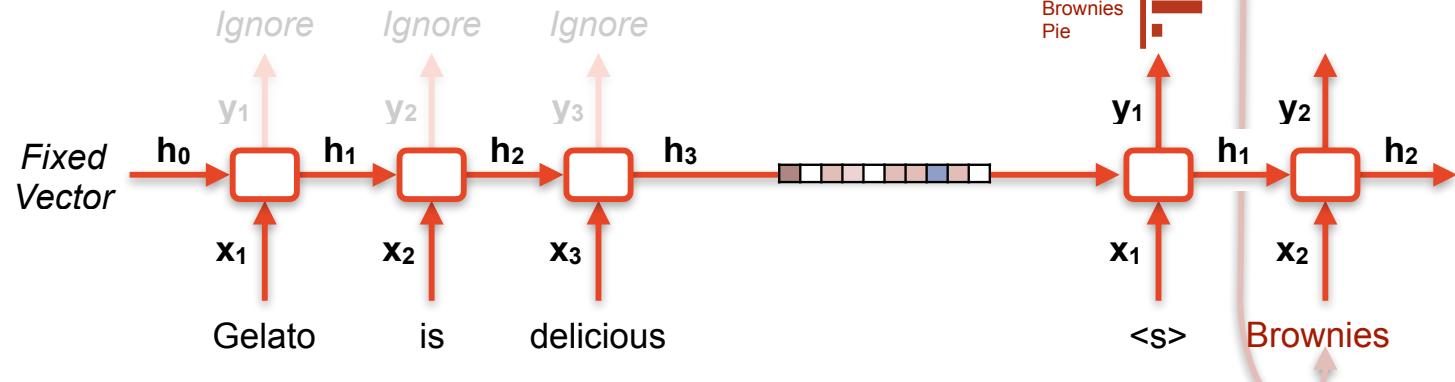
[menti.com 4843 3031](https://menti.com/48433031)

## Greedy inference with an encoder-decoder

Top-1



Sample





# Contextual Representations

## Encoder-Decoder

### Tokenisation

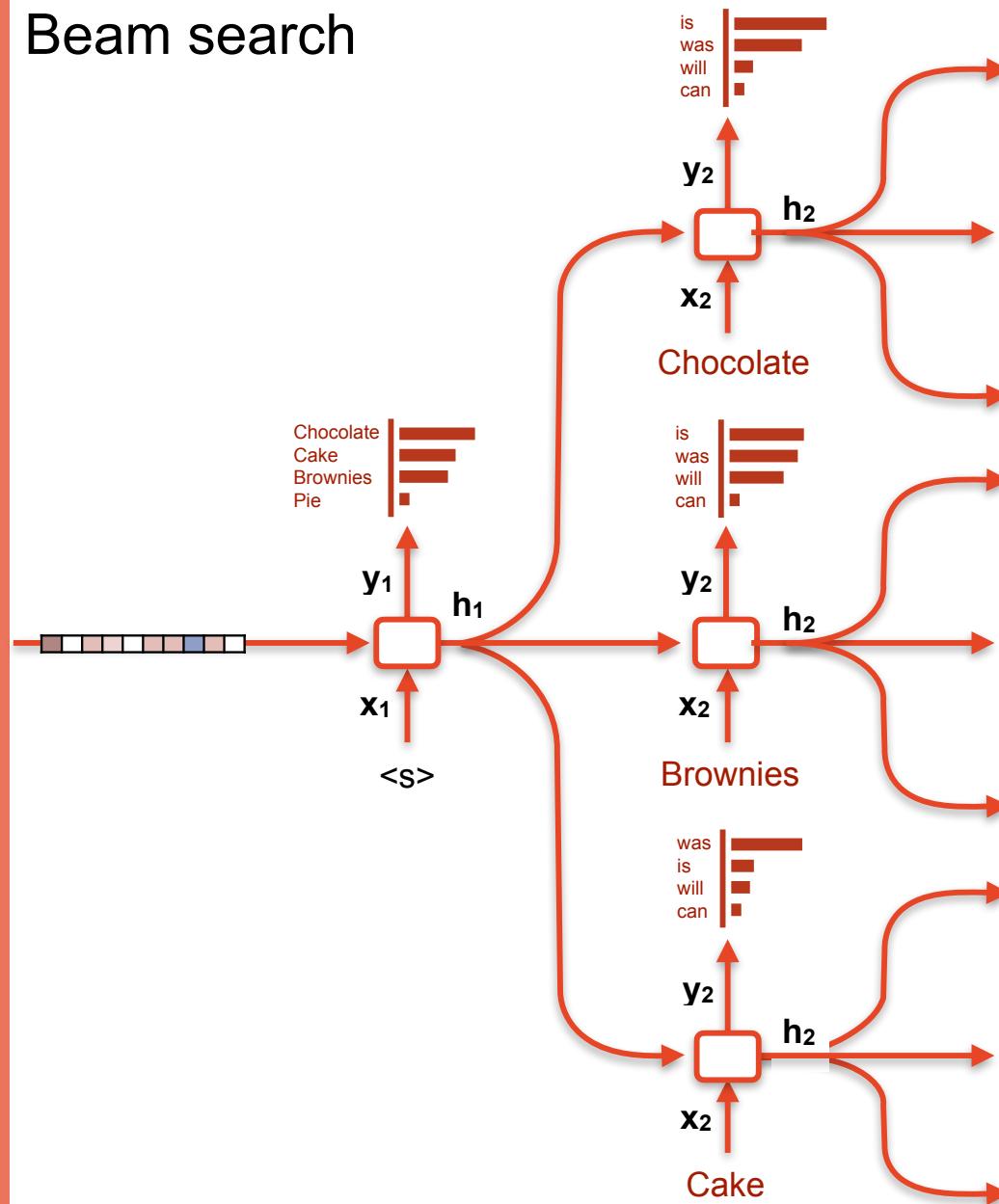
### Attention

## Workshop Preview



menti.com 4843 3031

## Beam search



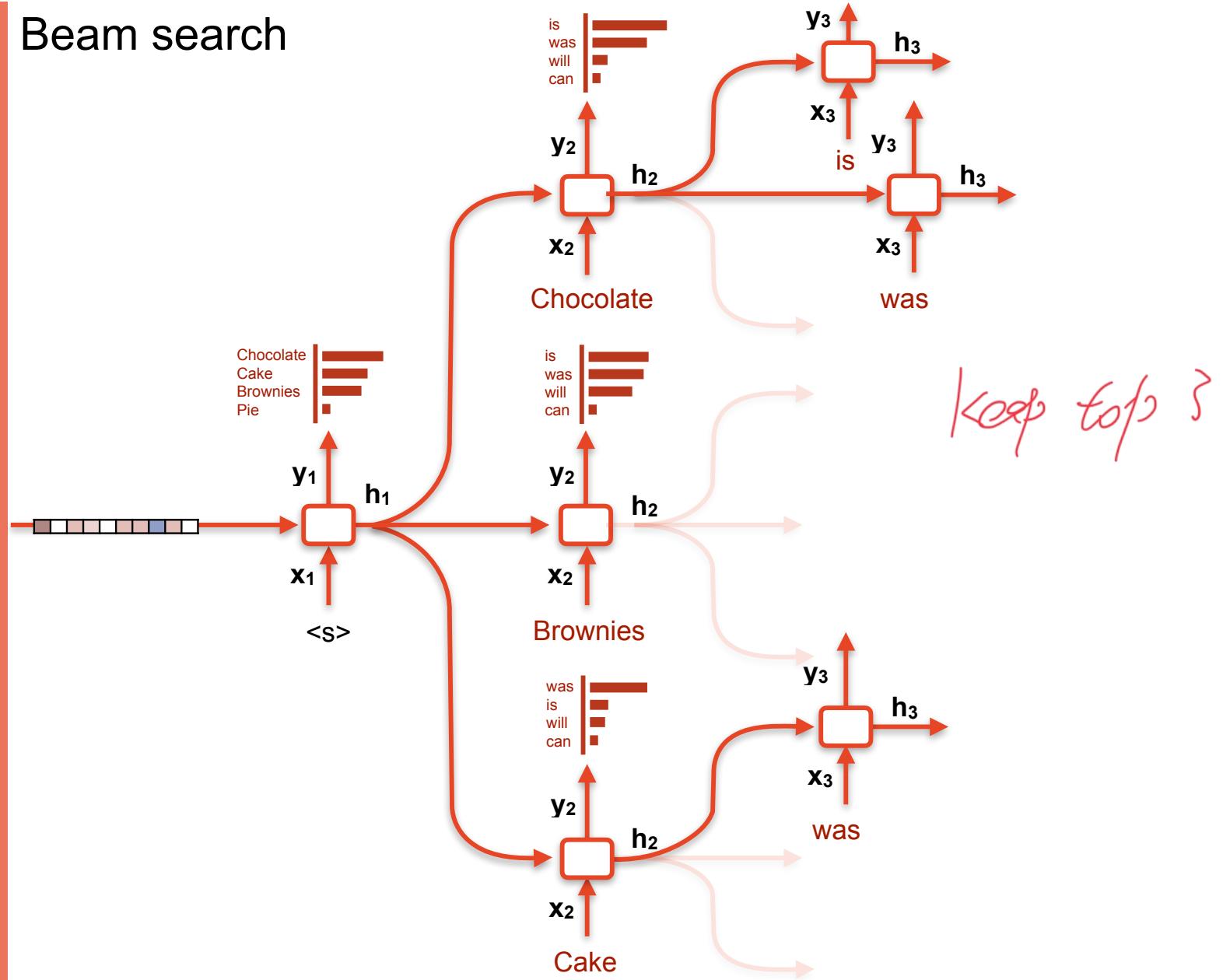


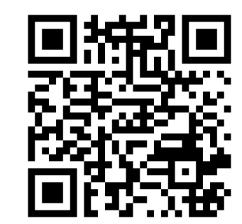
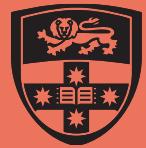
Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



menti.com 4843 3031

## Beam search





## Beam search and stopping (at different lengths)

Chocolate is delicious .

Cake is delicious .

Chocolate is the most delicious food.

(1)

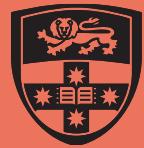
Keep going until we have N outputs.

(2)

Adjust the scores to get the average probability per word

$$\frac{1}{|\text{words}|} \sum \log P(\text{word} | \text{context})$$

deal with input with different length



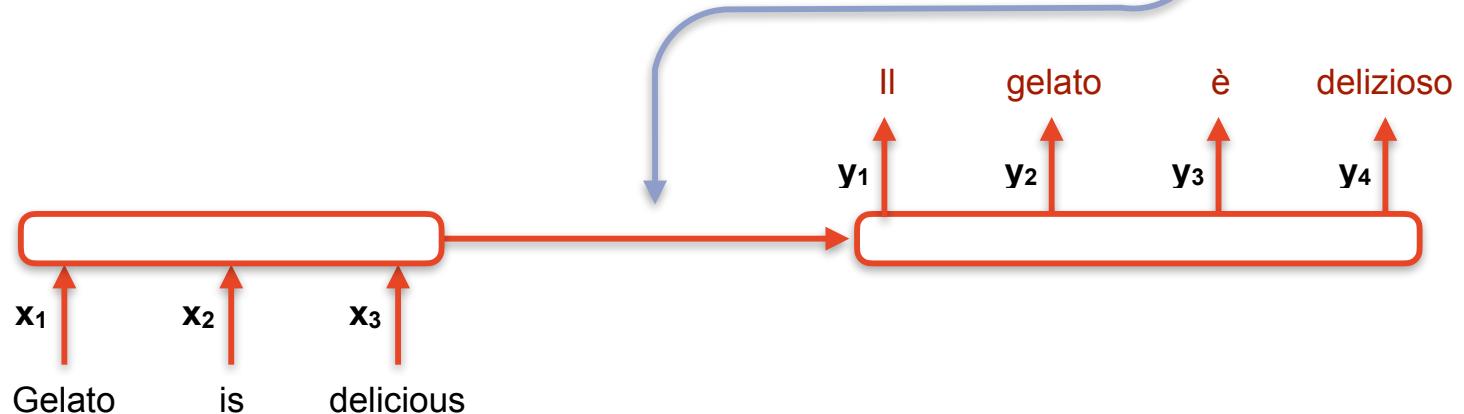
Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



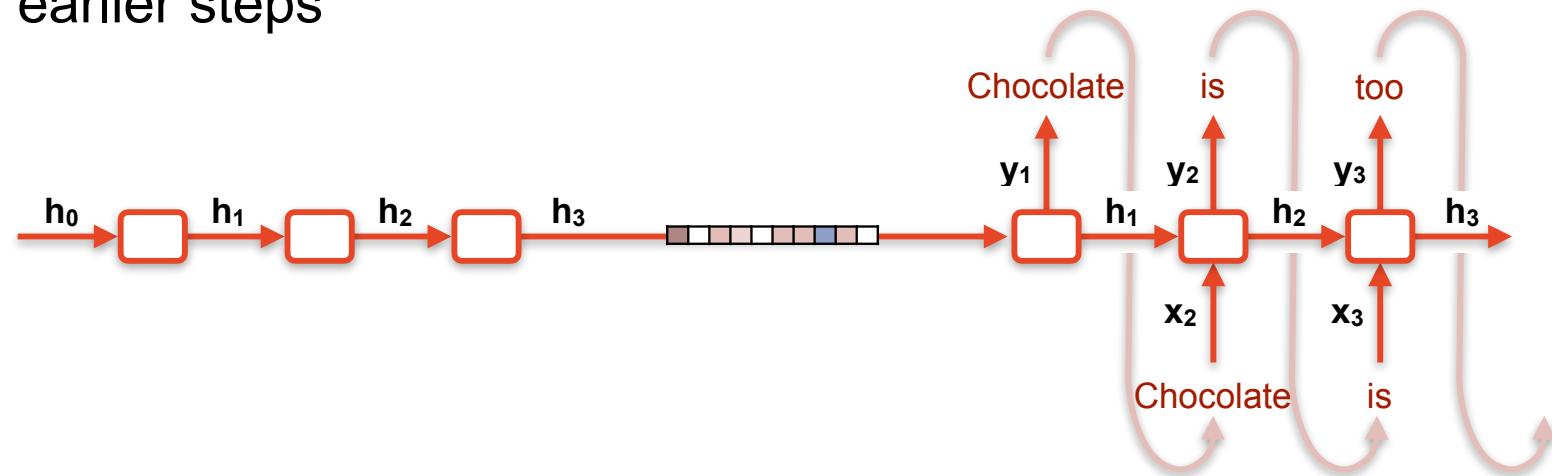
menti.com 4843 3031

## Issues with the encoder-decoder as described so far

① Bottleneck - all information has to go through one vector

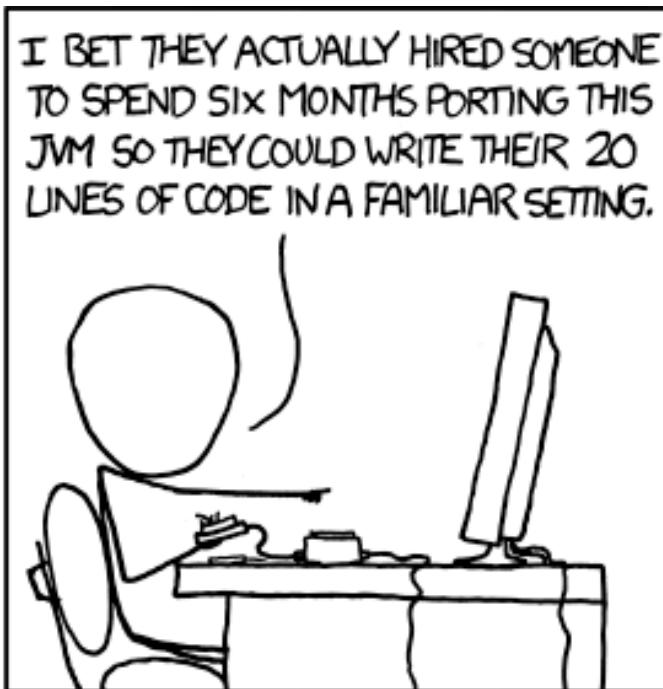
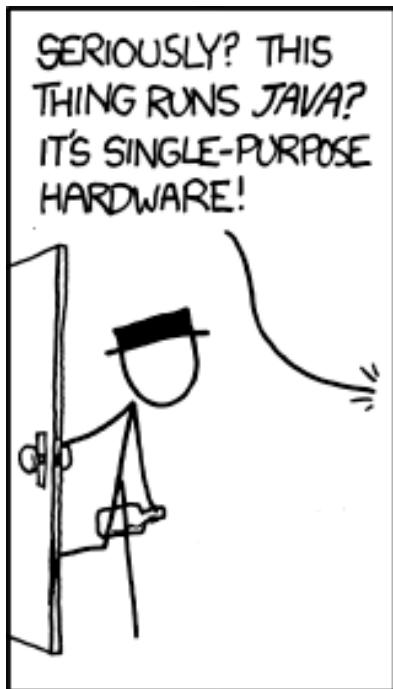


② Hard to parallelise - every step of the calculation depends on the earlier steps





## Golden Hammer



[Took me five tries to find the right one, but I managed to salvage our night out--if not the boat--in the end.]

Source: <https://xkcd.com/801/>

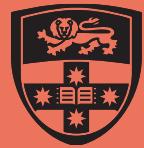


Contextual  
Representations  
Encoder-Decoder  
**Tokenisation**  
Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

# Tokenisation



Contextual  
Representations

**Encoder-Decoder**

Tokenisation

Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## How do we measure performance of translation systems?

Human

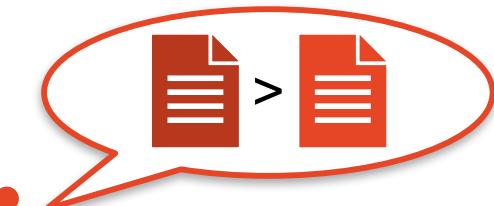
- Fluency
- Adequacy

Rate



4

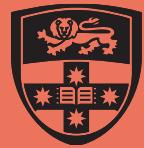
Rank



Edit &  
Compare



Bilingual? Not necessary If we have a sample translation



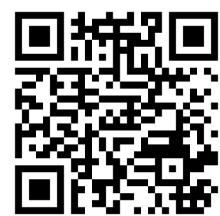
Contextual  
Representations

**Encoder-Decoder**

Tokenisation

Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

# How do we measure performance of translation systems?

Automatic

*Method*

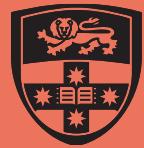
⑩

chrF - compare character ngrams

Human translation:  
I like to read too

Machine translation:  
Reading, I too like

I	II	III	IIlik	R	Re	Rea	Read
I	li	lik	like	e - 2	ea	ead	ead
i	ik	ike	iket	a	ad	adi	adin
k	ke	ket	keto	d	di	din	ding
e - 2	et	eto	etor	i - 2	in	ing	ing,
t - 2	to - 2	tor	tore	n	ng	ng,	ng,I
o - 3	or	ore	orea	g	g,	g,I	g,lt
r	re	rea	read	,	,I	,lt	,lto
a	ea	ead	eadt	l	lt	lto	ltoo
d	ad	adt	adto	t	to	too	tool
	dt	dto	dtoo	o - 2	oo	ool	ooli
	oo	too		l	ol	oli	olik
				k	li	lik	like
							ke



Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



menti.com 4843 3031

## How do we measure performance of translation systems?

Automatic

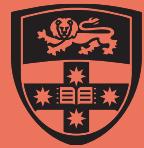
Human translation:  
I like to read too

Machine translation:  
Reading, I too like

chrF - compare character ngrams

$$\text{Precision} = \frac{TP}{TP + FP}$$
$$= \frac{23}{23 + 35}$$
$$= 0.40$$
$$\text{Recall} = \frac{TP}{TP + FN}$$
$$= \frac{23}{23 + 27}$$
$$= 0.46$$

	I	II	III	IIlik	R	Re	Rea	Read
I	li	lik	like	like	e - 2	ea	ead	ead
i	ik	ike	ket	ket	a	ad	adi	adin
k	ke	ket	keto	keto	d	di	din	ding
e - 2	et	eto	etor	etor	i - 1,1	in	ing	ing,
t - 1,1	to - 1,1	tor	tore	tore	n	ng	ng,	ng,I
o - 2,1	or	ore	orea	orea	g	g,	g,l	g,lt
r	re	rea	read	read	,	,l	,lt	,lto
a	ea	ead	eadt	eadt	l	lt	lto	ltoo
d	ad	adt	adto	adto	t	to	too	tool
dt	dto	dto	dtoo	dtoo	o - 2	oo	ool	ooli
oo	too	too			l	ol	oli	olik
2			3	4	k	li	lik	like
					ke			
					2	3	4	



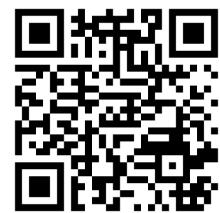
Contextual  
Representations

**Encoder-Decoder**

Tokenisation

Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## How do we measure performance of translation systems?

Automatic

Human translation:  
I like to read too

Machine translation:  
Reading, I too like

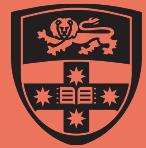
chrF - compare character ngrams

$$\text{Precision} = \frac{TP}{TP + FP}$$
$$= \frac{23}{23 + 35}$$
$$= 0.40$$

$$\text{Recall} = \frac{TP}{TP + FN}$$
$$= \frac{23}{23 + 27}$$
$$= 0.46$$

$$F_{\beta}\text{-score} = \frac{(1 + \beta^2)TP}{(1 + \beta^2)TP + \beta^2FN + FP}$$
$$F_2\text{-score} = \frac{5TP}{5TP + 4FN + FP}$$
$$= \frac{5 \cdot 23}{5 \cdot 23 + 4 \cdot 27 + 35}$$
$$= 0.45$$

The value of  $\beta$  leads  
the score to be closer  
to either P or R



## How do we measure performance of translation systems?

Automatic

Human translation  
I like to read too

chrF - compare character n-grams

$$\text{Precision} = \frac{TP}{TP + FP}$$
$$= \frac{23}{23 + 35}$$
$$= 0.40$$

$$\text{Recall} = \frac{TP}{TP + FN}$$
$$= \frac{23}{23 + 27}$$
$$= 0.46$$

$\beta = 0$

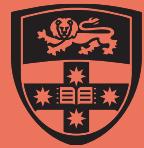
$$F_0\text{-score} = \frac{(1+0)TP}{(1+0)TP + 0FN + FP}$$
$$= \frac{TP}{TP + FP} = \text{Precision}$$

If N is big

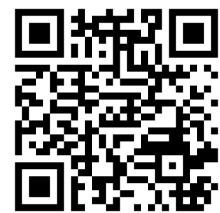
$\beta \approx \infty$

$$F_N\text{-score} = \frac{(1+N)TP}{(1+N)TP + N \cdot FN + FP}$$
$$\approx \frac{N \cdot TP}{N \cdot TP + N \cdot FN + FP}$$
$$= \frac{TP}{TP + FN + \frac{FP}{N}}$$
$$\approx \frac{TP}{TP + FN} = \text{Recall}$$

to either P or R



Contextual  
Representations  
**Encoder-Decoder**  
Tokenisation  
Attention  
Workshop Preview



menti.com 4843 3031

## How do we measure performance of translation systems?

### Automatic

Human  
I like to

Can also use  
word ngrams

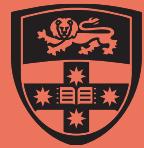
Machine translation:  
Reading, I too like

### chrF - compare character ngrams

I	II	III	IIlik	R	Re	Rea	Read
	li	lik	like	e - 2	ea	ead	ead
i	ik	ike	iket	a	ad	adi	adin
k	ke	ket	keto	d	di	din	ding
e - 2	et	eto	etor	i - 1,1	in	ing	ing,
t - 1,1	to - 2	tor	tore	n	ng	ng,	ng,I
o - 1,2	or	ore	orea	g	g,	g,I	g,It
		rea	read	,	,I	,It	,Ito
		ead	eadt	I	It	Ito	Itoo
		adt	adto	t	to	too	tool
		dto	dtoo	o - 2	oo	ool	ooli
		too		I	ol	oli	olik
				k	li	lik	ike
				ke			

Can move large chunks of  
the sentence and not impact  
the score. Maybe good?  
Maybe not?

move two part of sentence,  
only middle will be effect



Contextual  
Representations

**Encoder-Decoder**

Tokenisation

Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## How do we measure performance of translation systems?

Automatic

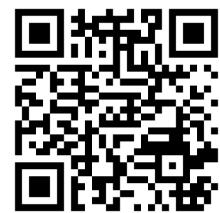
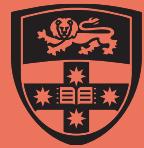
Method  
②

BLEU - compare word ngrams

Human translation:  
I like to read too

Machine translation:  
Reading, I too like

I	Reading
like	,
to	I
read	too
too	like
I like	Reading ,
like to	, I
to read	I too
read too	too like
I like to	Reading , I
like to read	, I too
to read too	I too like
I like to read	Reading , I too
like to read too	, I too like



## How do we measure performance of translation systems?

### Automatic

Human translation:  
I like to read too

BLEU - compare word ngrams

- Calculate precision for unigrams, bigrams, trigrams, 4grams separately
- Take the geometric mean
- Apply this to each sentence

Geometric mean only works well with multiple sentences - no 4grams matches means a score of 0!

I  
like  
to  
read  
too  
I like

to  
read  
too  
I like

to  
read  
too  
I like

to  
read  
too  
I like

to  
read  
too  
I like

to  
read  
too  
I like

to  
read  
too  
I like

Machine translation:  
Reading, I too like

Tokenisation?

Reading

No credit for  
read vs. Reading

Reading ,

, I

I too

too like

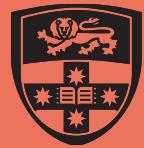
Reading , I

, I too

I too like

Reading , I too

, I too like



Contextual  
Representations  
Encoder-Decoder  
**Tokenisation**  
Attention  
Workshop Preview



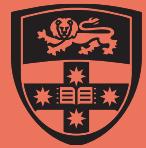
So far, we have always split on whitespace to get tokens

“Chocolate is delicious”

`text.split()`

“Chocolate”      “is”      “delicious”

[menti.com 4843 3031](https://menti.com/48433031)



Contextual  
Representations  
Encoder-Decoder  
**Tokenisation**  
Attention  
Workshop Preview

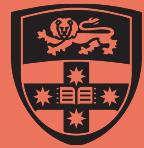


[menti.com 4843 3031](https://menti.com/48433031)

What about these cases?

Contractions:      won't

Punctuation:      Great!  
                          e.g.,



And can tokenisation help with rare words?

Variations      taaaaaaaasty

Misspellings      laern

New words      transformerify

Previously we used  
method's like Fasttext's  
representations to  
address this

Examples from Stanford's cs224n course



## Split into sub-tokens

Variations      taaaaaaaasty      →      ta## aaaa## asty

Missepllings      laern      →      la## ern

New words      transformerify      →      transformer## ify

How do we determine where to split?



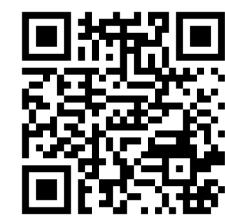
Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention

Workshop Preview

OCCUR  
5 times



menti.com 4843 3031

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowest\_  
6 newer\_  
3 wider\_  
2 new\_

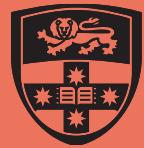
### Vocab:

\_, d, e, i, l, n, o, r, s, t, w

### Pairs and counts:

9 r _	2 t _
9 e r	2 s t
8 w e	2 e s
8 n e	
8 e w	
7 w _	
7 o w	
7 l o	
3 w i	
3 i d	
3 d e	

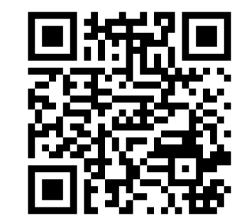
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowest\_  
6 newer\_ r  
3 wider\_ r  
2 new\_

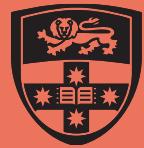
### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_

### Pairs and counts:

9 e r_	2 s t
8 w e	2 e s
8 n e	
8 e w	
7 w _	
7 o w	
7 l o	
3 w i	
3 i d	
3 d e	
2 t _	

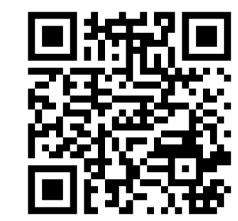
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



menti.com 4843 3031

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowest\_  
6 newer\_  
3 wider\_  
2 new\_

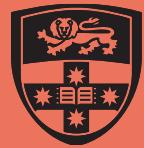
### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er  
*new*

### Pairs and counts:

8 ne 2 st  
8 ew 2 es  
7 w\_  
7 ow  
7 lo  
6 wr\_  
3 wi  
3 id  
3 der\_  
2 we  
2 t\_

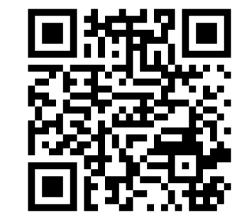
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



menti.com 4843 3031

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowest\_  
6 newer\_  
3 wider\_  
2 new\_

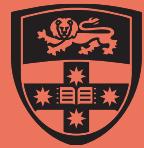
### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne

### Pairs and counts:

8 new 2 es  
7 w\_  
7 o w  
7 l o  
6 w er\_  
3 wi  
3 id  
3 der\_  
2 we  
2 t\_  
2 st

Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



menti.com 4843 3031

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowest\_  
6 newer\_  
3 wider\_  
2 new\_

### Vocab:

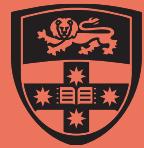
\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new

### Pairs and counts:

7 o w 2 e s  
7 l o  
6 new er\_  
5 w \_  
3 w i  
3 i d  
3 d er\_  
2 w e  
2 t \_  
2 s t  
2 new \_

Examples from Jurafsky  
and Martin textbook

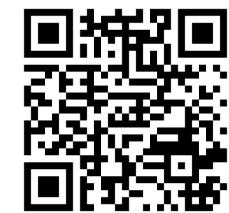
breaking ie in various way



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 I ow \_  
2 I ow e s t \_  
6 new er\_  
3 w i d er\_  
2 new \_

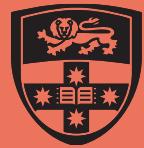
### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow

### Pairs and counts:

7 I ow  
6 new er\_  
5 ow \_  
3 w i  
3 i d  
3 d er\_  
2 t \_  
2 s t  
2 ow e  
2 new \_  
2 e s

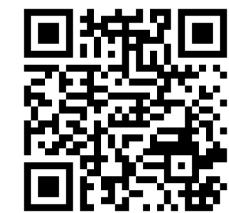
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low \_  
2 low e s t \_  
6 new er \_  
3 w i d er \_  
2 new \_

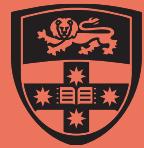
### Vocab:

\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low

### Pairs and counts:

6 new er \_  
5 low \_  
3 w i  
3 i d  
3 d er \_  
2 t \_  
2 s t  
2 new \_  
2 low e  
2 e s

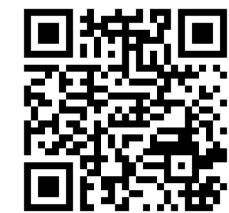
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low \_  
2 low e s t \_  
6 newer \_  
3 w i d er \_  
2 new \_

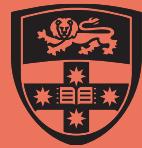
### Vocab:

\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_

### Pairs and counts:

5 low \_  
3 w i  
3 i d  
3 d er \_  
2 t \_  
2 s t  
2 new \_  
2 low e  
2 e s

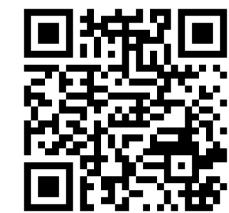
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 low e s t \_  
6 newer\_  
3 w i d er\_  
2 new \_

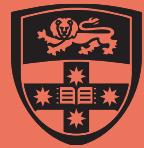
### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_

### Pairs and counts:

3 w i  
3 i d  
3 d er\_  
2 t \_  
2 s t  
2 new \_  
2 low e  
2 e s

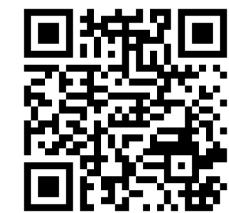
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 low e s t \_  
6 newer\_  
3 wi d er\_  
2 new \_

### Vocab:

\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi

### Pairs and counts:

3 wi d  
3 d er\_  
2 t \_  
2 s t  
2 new \_  
2 low e  
2 e s

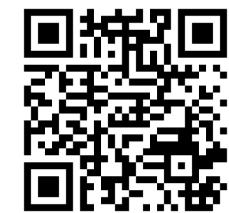
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 low e s t \_  
6 newer\_  
3 wid er\_  
2 new \_

### Pairs and counts:

3 wid er\_  
2 t \_  
2 s t  
2 new \_  
2 low e  
2 e s

### Vocab:

\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid

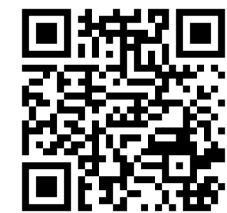
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 low e s t \_  
6 newer\_  
3 wider\_  
2 new \_

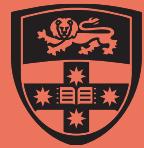
### Pairs and counts:

2 t \_  
2 s t  
2 new \_  
2 low e  
2 e s

### Vocab:

\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid, wider\_

Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 low e s t\_  
6 newer\_  
3 wider\_  
2 new \_

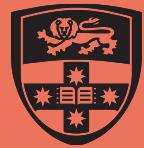
### Pairs and counts:

2 s t\_  
2 new \_  
2 low e  
2 e s

### Vocab:

\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid, wider\_, t\_

Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 low e st\_  
6 newer\_  
3 wider\_  
2 new \_

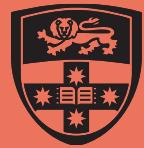
### Pairs and counts:

2 new \_  
2 low e  
2 e st\_

### Vocab:

\_ , d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid, wider\_, t\_, st\_

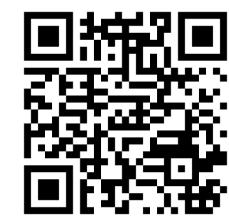
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 low e st\_  
6 newer\_  
3 wider\_  
2 new\_

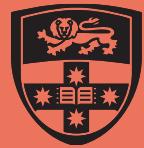
### Pairs and counts:

2 low e  
2 e st\_

### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid, wider\_, t\_, st\_,  
new\_

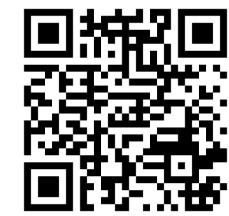
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowe st\_  
6 newer\_  
3 wider\_  
2 new\_

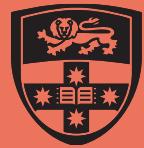
### Pairs and counts:

2 lowe st\_

### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid, wider\_, t\_, st\_,  
new\_, lowe

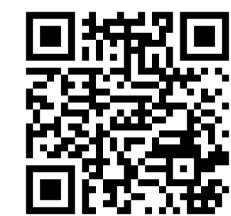
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowest\_  
6 newer\_  
3 wider\_  
2 new\_

### Pairs and counts:

### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid, wider\_, t\_, st\_,  
new\_, lowe, lowest\_

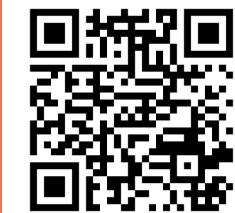
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



menti.com 4843 3031

## Byte-Pair Encoding (BPE)

### Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items that are adjacent most frequently
3. Create a new vocabulary item for that pair and update data
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

### Data:

5 low\_  
2 lowest\_  
6 newer\_  
3 wider\_  
2 new\_

### Pairs and counts:

### Vocab:

\_, d, e, i, l, n, o, r, s, t, w, r\_,  
er\_, ne, new, ow, low, newer\_,  
low\_, wi, wid, wider\_, t\_, st\_,  
new\_, lowe, lowest\_

*make new words*

⇒ low er\_  
new e st\_  
wid e st\_

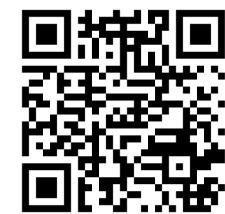
Examples from Jurafsky  
and Martin textbook



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

Other similar methods vary step 2

Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items using some method
3. Create a new vocabulary item for that pair
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

D BPE

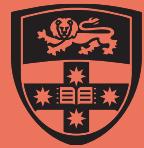
b WordPiece

Step 2 ‘some method’ =  
adjacent most frequently

Step 2 ‘some method’ =  

- Train an n-gram language model
- Consider all possible vocabulary pairs
- Choose the one that will decrease perplexity most if added to the model

Step 2 is different for those methods



Contextual  
Representations  
Encoder-Decoder

## Tokenisation

Attention  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

Other similar methods vary step 2

Algorithm:

1. Set vocabulary to be all characters
2. Find the two vocabulary items using some method
3. Create a new vocabulary item for that pair
4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.



HuggingFace  
WordPiece

Step 2 ‘some method’ =  
Pair that maximises

$$\frac{|\text{combined}|}{|\text{first symbol}| \cdot |\text{second symbol}|}$$

<https://huggingface.co/learn/nlp-course/en/chapter6/6>



## Unigram / SentencePiece takes a different approach

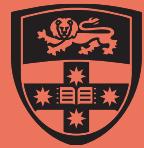
So far: start with small units, then combine units

Alternative: start with big and small units, then delete units

Algorithm:

1. Set vocabulary to be all characters and all frequent character sequences, up to full words.
2. Find vocabulary items to remove using a complex method.
3. Repeat until the vocabulary gets down to the desired size.

*How? I don't think it's downside*



## How do they compare?

**Original:** furiously

**BPE:** \_fur iously

**Uni. LM:** \_fur ious ly

**Original:** tricycles

**BPE:** \_t ric y cles

**Uni. LM:** \_tri cycle s

**Original:** nanotechnology

**BPE:** \_n an ote chn ology

**Uni. LM:** \_nano technology

**Original:** Completely preposterous suggestions

**BPE:** \_Comple t ely \_prep ost erous \_suggest ions

**Unigram LM:** \_Complete ly \_pre post erous \_suggestion s

**Original:** corrupted

**BPE:** \_cor rupted

**Unigram LM:** \_corrupt ed

**Original:** 1848 and 1852,

**BPE:** \_184 8 \_and \_185 2,

**Unigram LM:** \_1848 \_and \_1852 ,

**Original** 磁性は様々に分類がなされている。

**BPE** 磁 性は 様々 に 分類 がなされている。

**Unigram LM** 磁 性 は 様々 に 分類 がなされている。

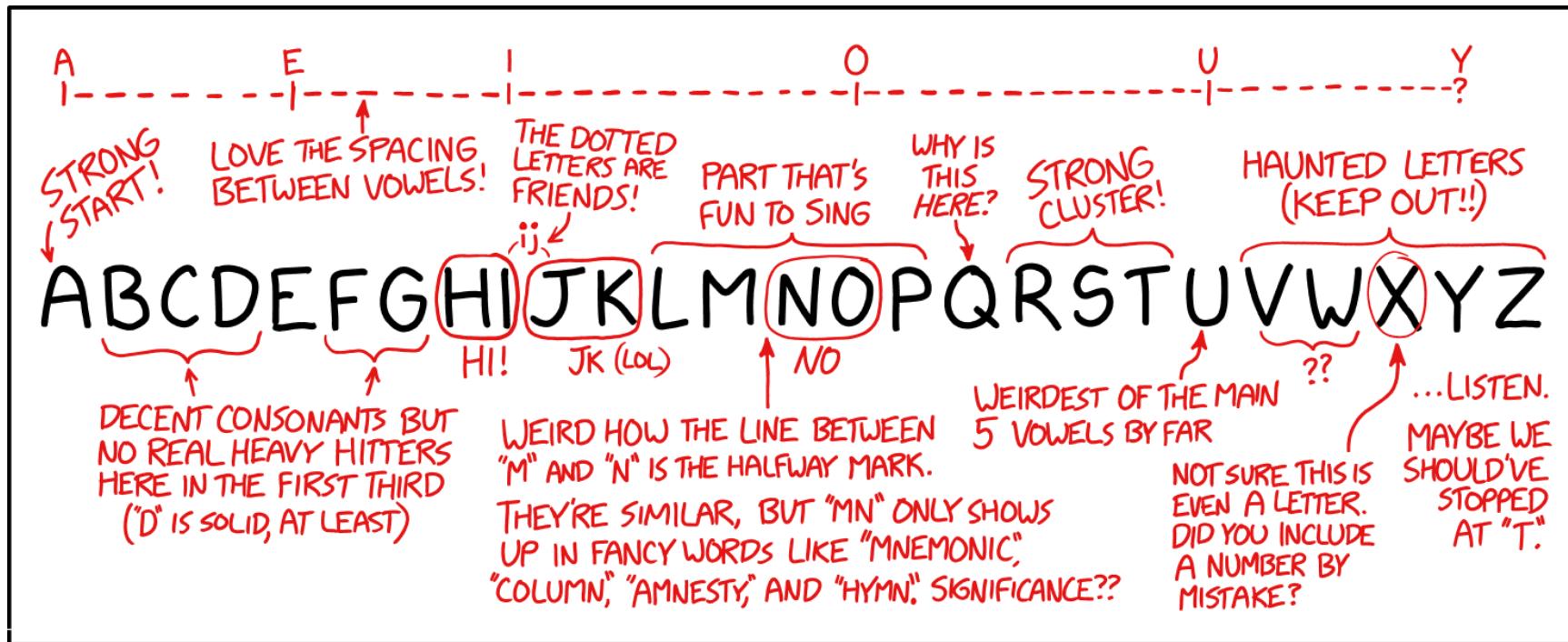
**Gloss** magnetism (top.) various ways in classification is done

**Translation** Magnetism is classified in various ways.

More frequent in									
BPE					Unigram LM				
H	L	M	T	B	s	.	,	ed	d
P	C	K	D	R	ing	e	ly	t	a



## Alphabet Notes



### DESIGN NOTES ON THE ALPHABET

[Listen, you're very cute, but if you rearrange the alphabet to put U and I together it will RUIN the spacing! ]

Source: <https://xkcd.com/2794/>



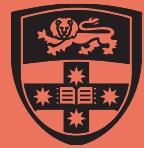
COMP 4446 / 5046  
Lecture 6, 2025

Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

# Attention



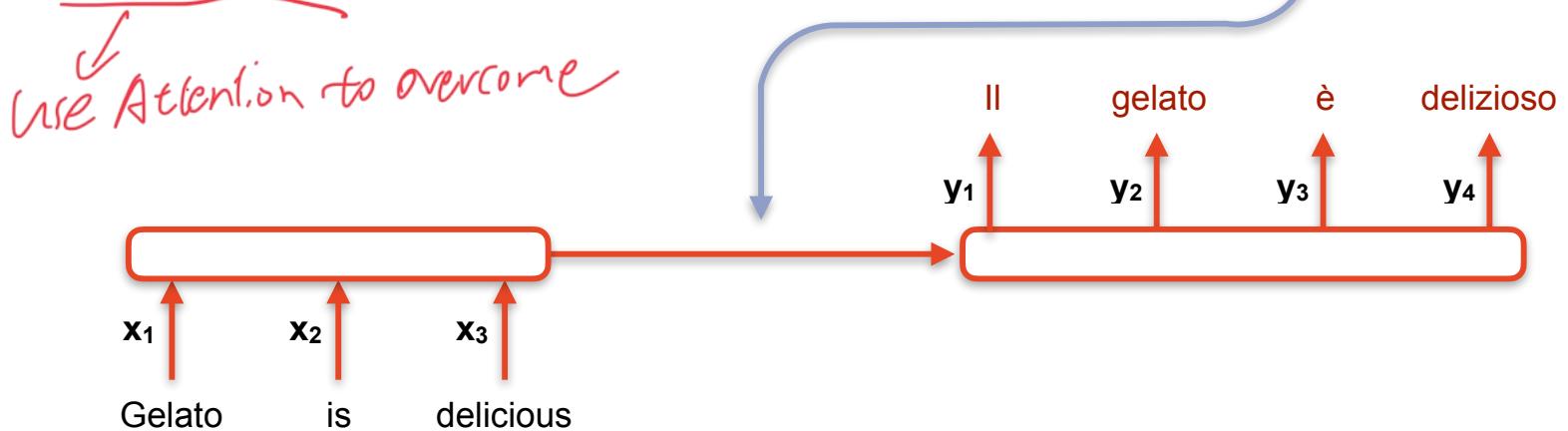
Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview



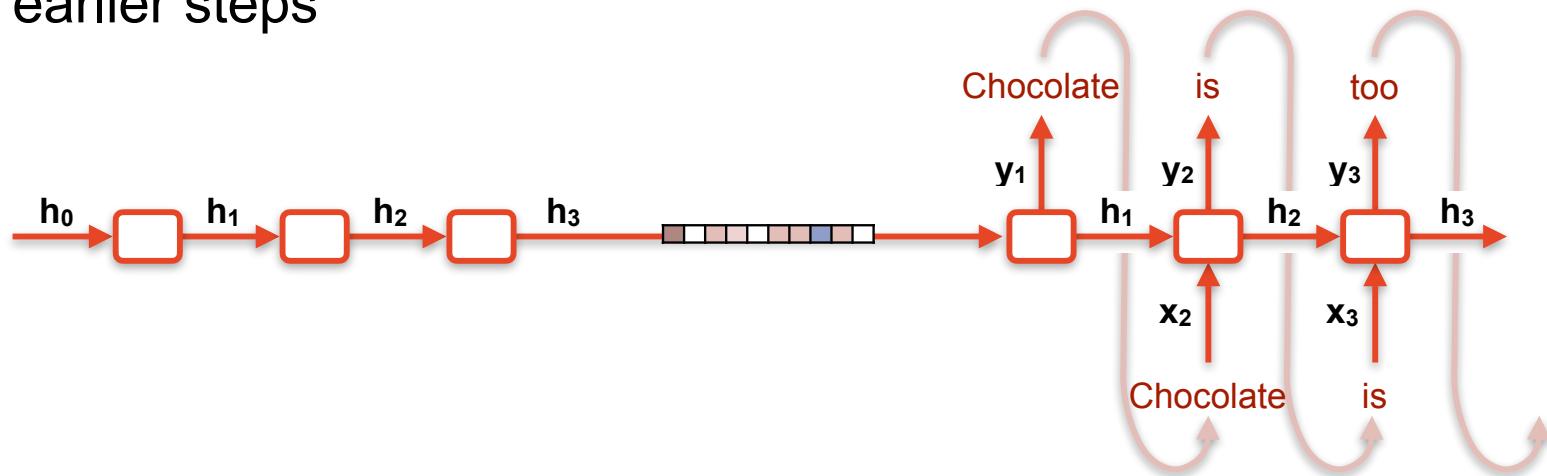
[menti.com 4843 3031](https://menti.com/48433031)

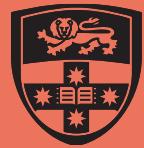
## Reminder: Issues with the encoder-decoder

Bottleneck - all information has to go through one vector



Hard to parallelise - every step of the calculation depends on the earlier steps





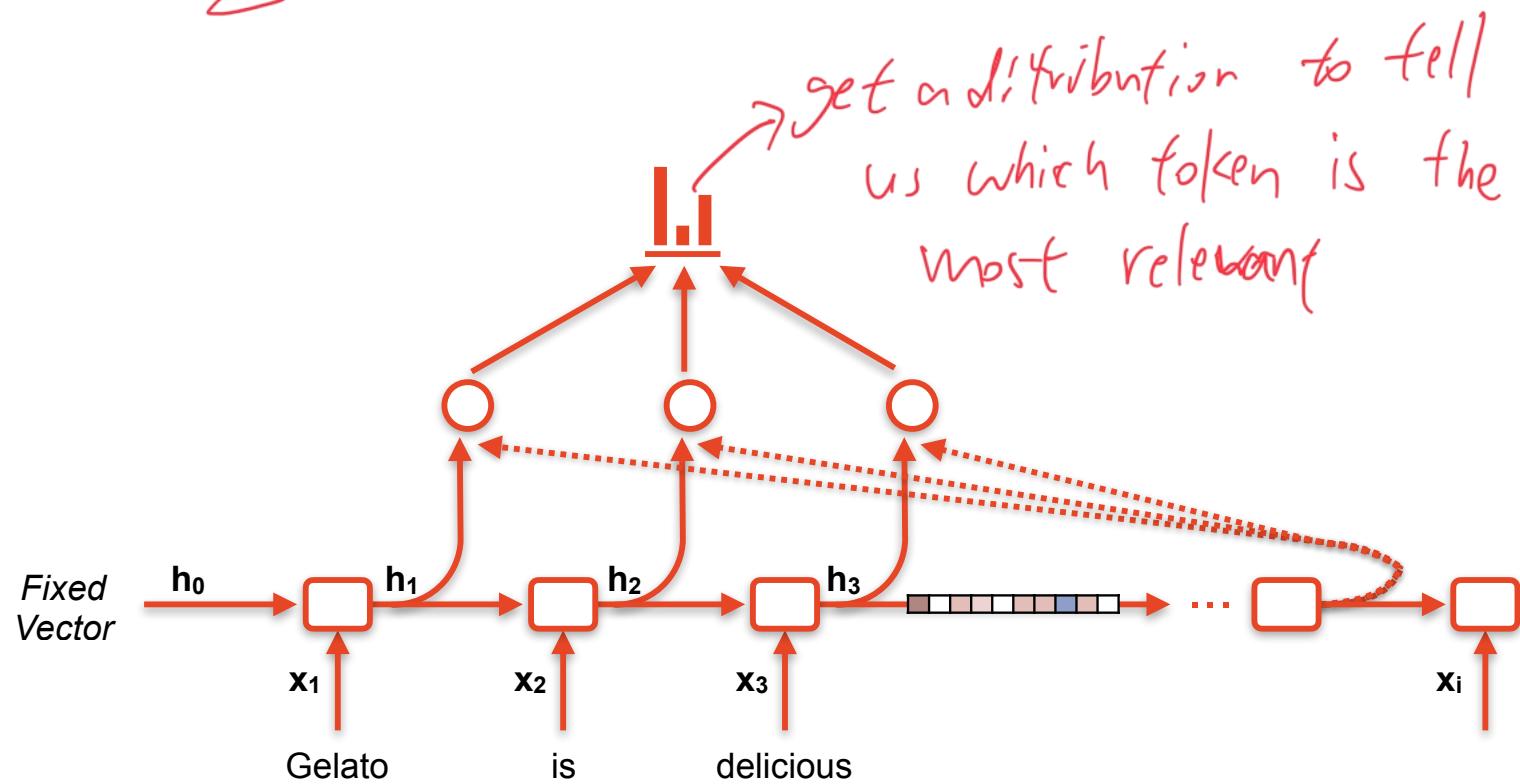
Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview



menti.com 4843 3031

We will solve this with ‘attention’

Give each output step  
a representation of the input  
that is most useful for that output decision





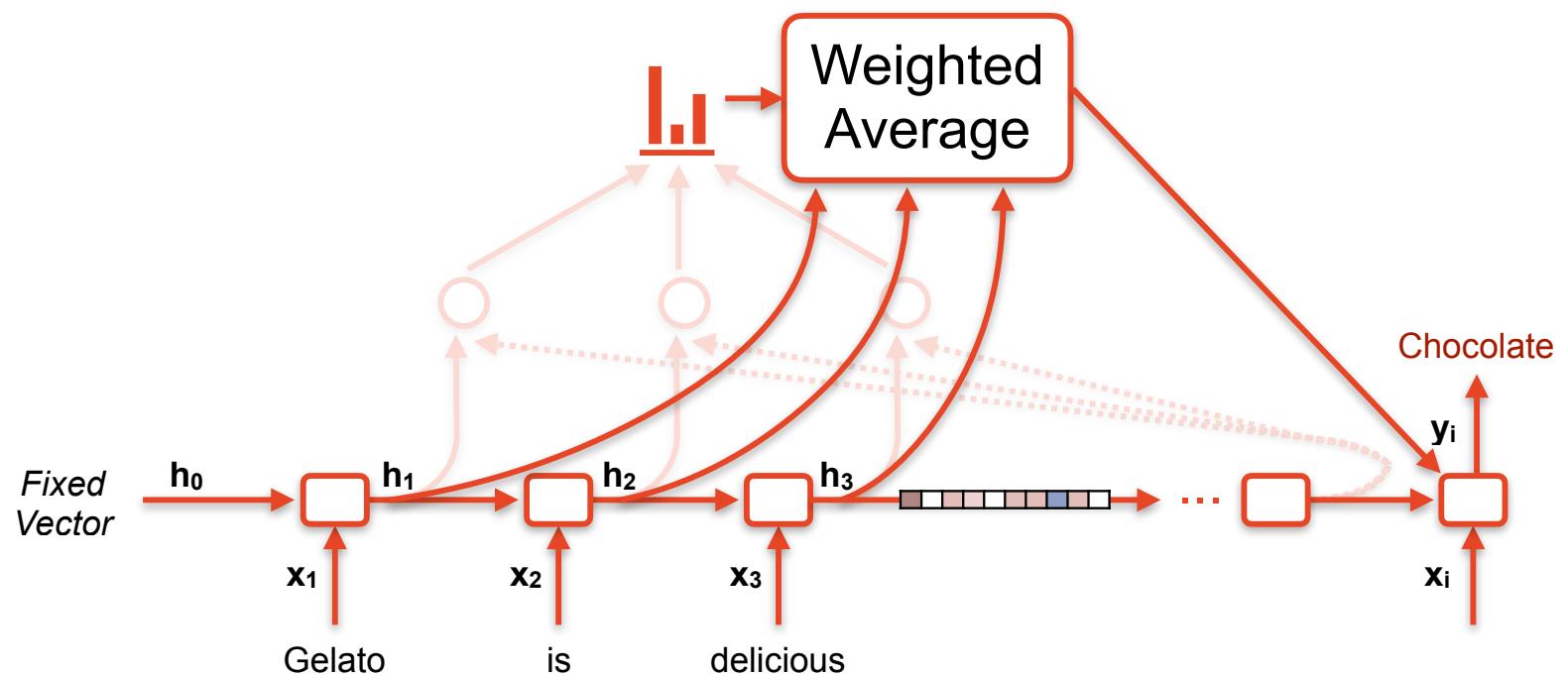
Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview

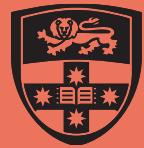


[menti.com 4843 3031](https://menti.com/48433031)

We will solve this with ‘attention’

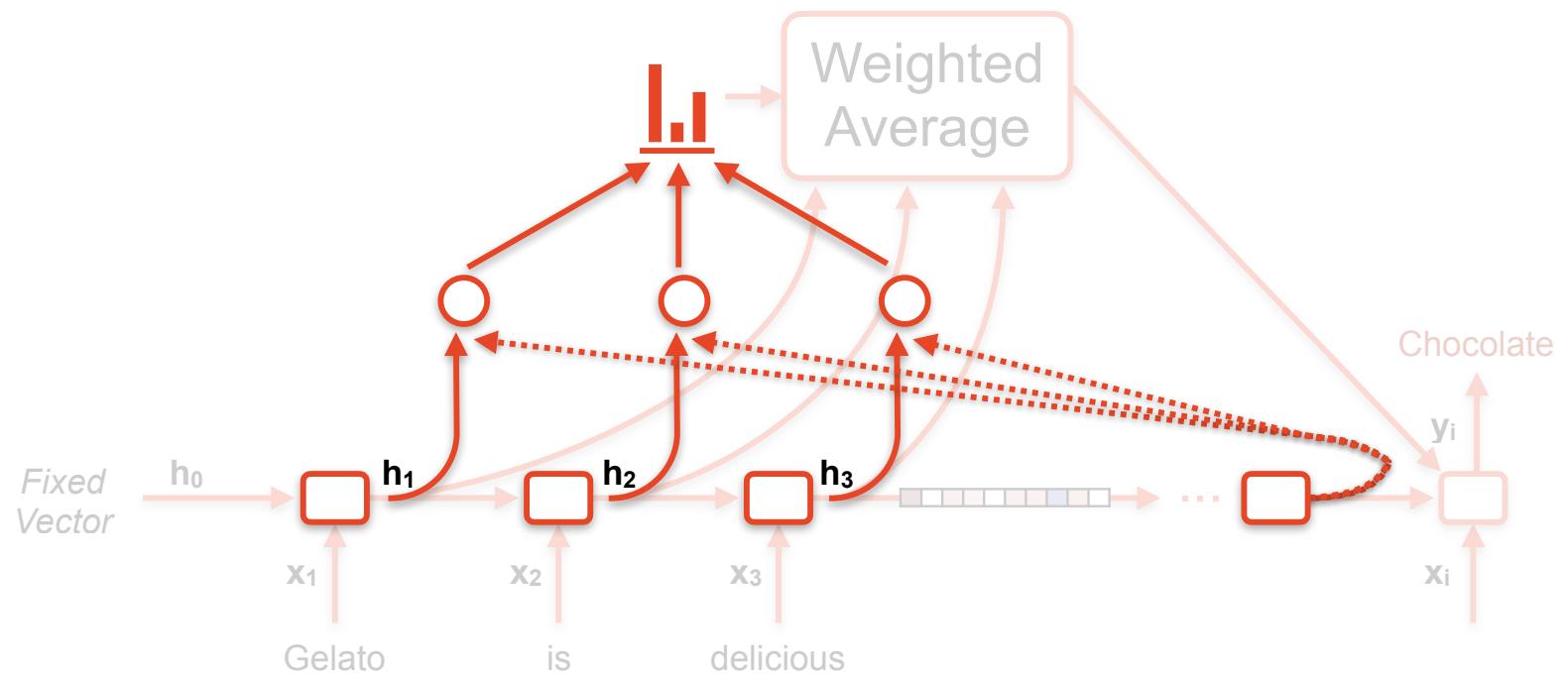
Give each output step  
a representation of the input  
that is most useful for that output decision

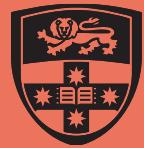




We will solve this with ‘attention’

Give each output step  
a representation of the input  
that is most useful for that output decision





We will solve this with ‘attention’

Give each output step  
a representation of the input  
that is most useful for that output decision

Hidden vectors from input

$$\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3 \dots \mathbf{h}_N \in \mathbb{R}^{d_h}$$

Hidden vector for current output step

$$\mathbf{s}_i \in \mathbb{R}^{d_h}$$

Calculate attention scores

$$\mathbf{e}_i = [\mathbf{s}_i^\top \mathbf{h}_1, \mathbf{s}_i^\top \mathbf{h}_2, \mathbf{s}_i^\top \mathbf{h}_3, \dots, \mathbf{s}_i^\top \mathbf{h}_N] \in \mathbb{R}^N$$

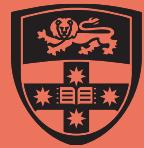
Normalise s

This is dot product attention

$$\alpha(\mathbf{e}_i) \in \mathbb{R}^N$$

Calculate weighted average

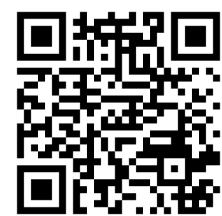
$$\mathbf{a}_i = \sum_{j=1}^N \alpha_{i,j} \mathbf{h}_j \in \mathbb{R}^{d_h}$$



Contextual  
Representations  
Encoder-Decoder  
Tokenisation

## Attention

Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

There are many forms of attention

Dot product attention

$$\mathbf{e} = \mathbf{s}^T \mathbf{h}$$

Scaled dot product  
attention

$$\mathbf{e} = \frac{\mathbf{s}^T \mathbf{h}}{\sqrt{d_h}}$$

Motivated by  
statistical properties  
of dot products

Multiplicative attention /  
Bilinear attention

$$\mathbf{e} = \mathbf{s}^T \mathbf{W} \mathbf{h}$$

s and h can be  
different sizes

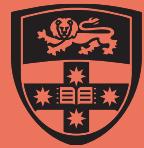
Reduced-rank  
multiplicative attention

$$\begin{aligned}\mathbf{e} &= \mathbf{s}^T (\mathbf{U}^T \mathbf{V}) \mathbf{h} \\ &= (\mathbf{s} \mathbf{U})^T (\mathbf{V} \mathbf{h})\end{aligned}$$

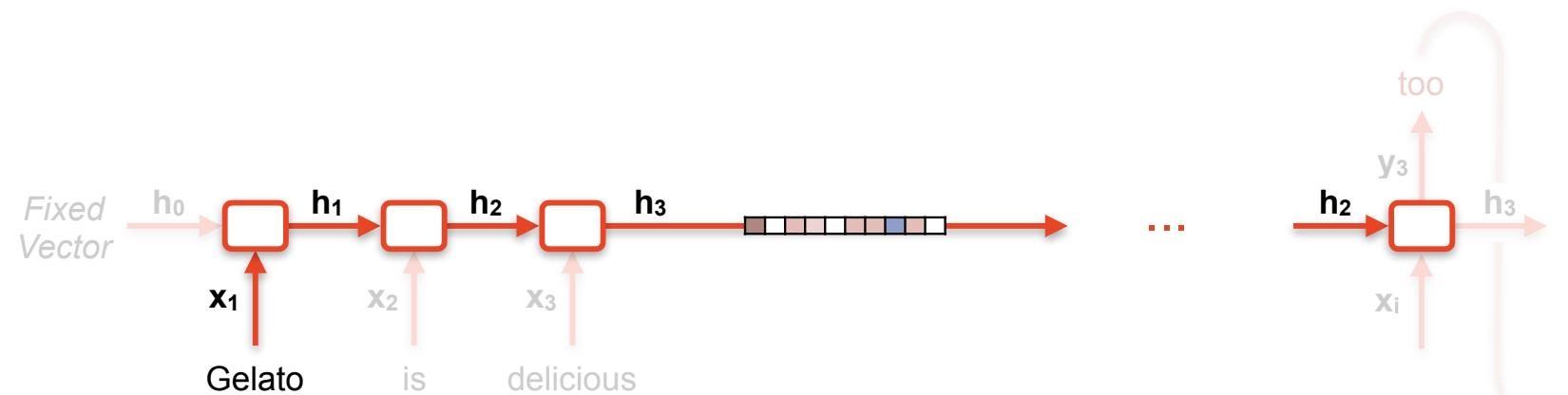
Can improve  
efficiency

Additive attention /  
Feedforward attention

$$\mathbf{e} = \mathbf{b} \tanh(\mathbf{W}_1 \mathbf{h} + \mathbf{W}_2 \mathbf{s})$$



This resolves the bottleneck problem without introducing many extra parameters



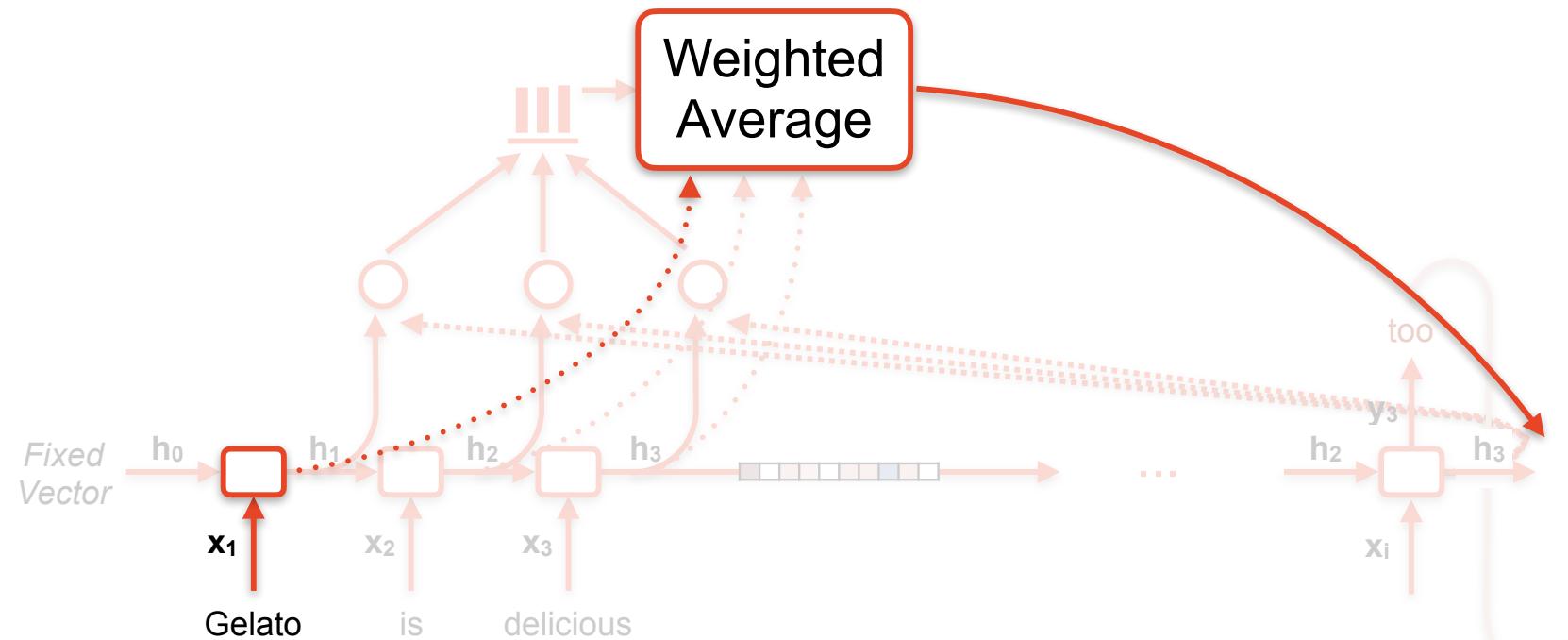


Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

This resolves the bottleneck problem without introducing many extra parameters





Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview

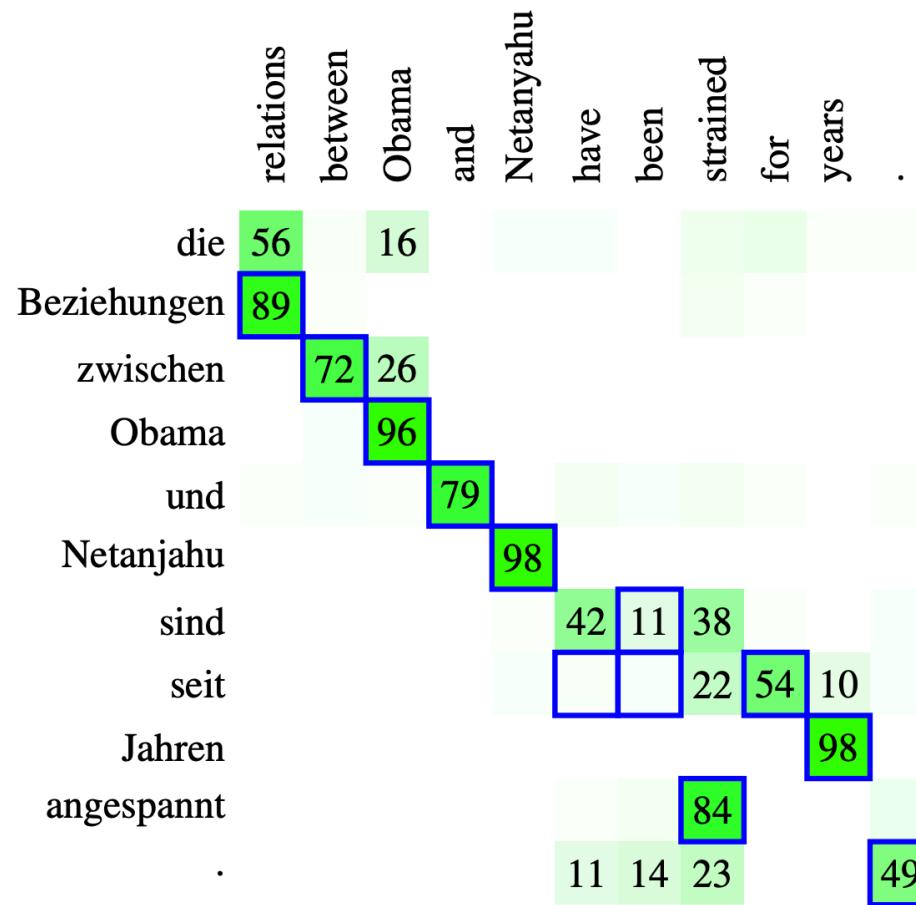


[menti.com 4843 3031](https://menti.com/48433031)

Attention in machine translation seems to capture word alignment!

Blue boxes:  
Alignment  
algorithm output

Numbers:  
Attention scores



Koehn and Knowles (2017)



Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview



[menti.com 4843 3031](https://menti.com/48433031)

We need to be careful not to read too much into these

## Understanding Neural Networks through Representation Erasure

### Attention is not Explanation

Sarthak Jain

Northeastern University

jain.sar@husky.neu.edu

### Attention is not not Explanation

Yuval Pinter\*

School of Interactive Computing

Georgia Institute of Technology

yvp@gatech.edu

## Is Attention Explanation? An Introduction to the Debate

Adrien Bibal, Rémi Cardon, David Alfter, Rodrigo Wilkens, Xiaoou Wang,  
Thomas François\* and Patrick Watrin\*

CENTAL, IL&C, University of Louvain, Belgium

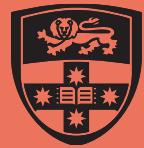
{adrien.bibal, remi.cardon, david.alfter, xiaoou.wang,  
thomas.francois, patrick.watrin}@uclouvain.be

### Abstract

The performance of deep learning models in NLP and other fields of machine learning has led to a rise in their popularity, and so the need for explanations of these models becomes paramount. Attention has been seen as a solution to increase performance, while providing some explanations. However, a debate has started to cast doubt on the explanation power of attention.

(global vs. local attention, according to Luong et al. (2015)) and where the query is generated (cross vs. self-attention as in the works of Bahdanau et al. (2015) and Vaswani et al. (2017)). In this paper, we focus on attention regardless of these technical differences. There are mainly two ways of computing the attention weights  $\hat{\alpha}$ : Bahdanau et al. (2015) introduced additive attention  $\hat{\alpha} = \text{softmax}(\mathbf{w}_3^T \tanh(\mathbf{W}_1 K + \mathbf{W}_2 Q))$ , where

ans for, e.g., model debugging or correction. A recent paper (Jain and Pinter, 2020) points to possible pitfalls that may arise from this approach. It also points to misapply attention scores as explanations. The authors argue that attention distributions should be interpreted as correlations between features and predictions, rather than as exclusive given a prediction.<sup>1</sup> Its main argument is that attention distributions exist that are very similar to those obtained by the original model's attention, but do not necessarily reflect the model's internal reasoning process. This is because the original model's attention is learned to optimize for a specific task, while the attention distribution used in the paper is learned to explain the model's behavior. The paper also shows that attention distributions can be manipulated to produce different explanations, which can lead to incorrect conclusions about the model's behavior. The paper concludes that attention distributions should be used with caution and that they should not be interpreted as exclusive given a prediction.



Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
**Attention**  
Workshop Preview

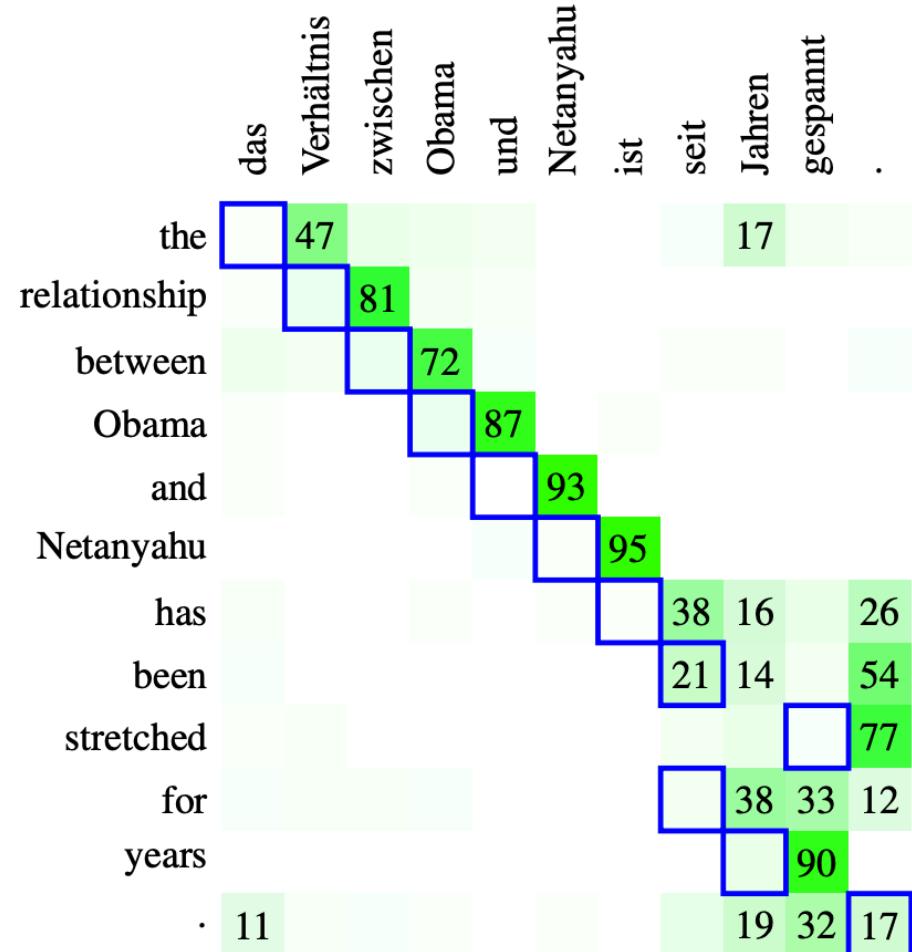


[menti.com 4843 3031](https://menti.com/48433031)

Attention in machine translation seems to capture word alignment! Or does it?

Blue boxes:  
Alignment  
algorithm output

Numbers:  
Attention scores



Koehn and Knowles (2017)



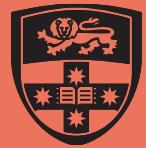
COMP 4446 / 5046  
Lecture 6, 2025

Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
Attention  
**Workshop Preview**



[menti.com 4843 3031](https://menti.com/48433031)

# Workshop Preview



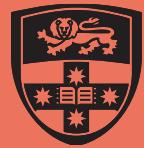
Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
Attention  
**Workshop Preview**



[menti.com 4843 3031](https://menti.com/48433031)

Pre-work: None this week

In-class: spaCy



Contextual  
Representations  
Encoder-Decoder  
Tokenisation  
Attention  
**Workshop Preview**



menti.com 4843 3031

## Muddy Card

Open shortly, 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 Cards Lecture 6