

"Attention is all you need".  
-My ex.

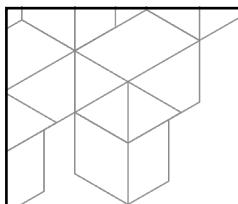


"The art of knowing is knowing what to ignore."  
- Rumi

All images from Wikimedia Commons unless specified.



1



# Natural Language Processing (NLP)

COMP6713 – 2025 Term 1



## Convener

Dr. Aditya Joshi

[aditya.joshi@unsw.edu.au](mailto:aditya.joshi@unsw.edu.au)



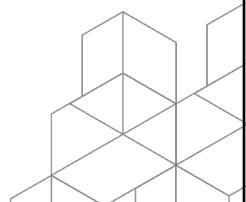
## Week 3

Attention! Transformer



## Schedule

2025 Term 1



2

**UNSW** SYDNEY | Australia's Global University

**Week 3**  
**Attention! Transformer**

**Attention**  
Intuition  
Attention in recurrent LM & limitations  
Attention as key, query, value

**Transformer**  
Inside-out view of Transformer  
Multi-head attention  
Encoder  
Decoder  
Byte-pair encoding  
Positional encoding  
Impact and legacy of Transformer

**And as usual...**  
code samples, pen-and-paper demos, linguistic examples and dad jokes.

3

**Topics this week**

Data → Features to look for → Ways to combine the features → Learn a model based on these combinations

Generation	1	2	3
Human Engineering	Icon: Human head with gear	Icon: Human head with gear	Icon: Human head with gear
Computational model	Icon: Laptop with gear	Icon: Laptop with gear	Icon: Laptop with gear
Neural model	Icon: Brain with gear	Icon: Brain with gear	Icon: Brain with gear

**Attention Transformer**

4

## Before we begin...

Congratulations on completing the first quiz!

### Final exam

No memorization of code syntax: there may be fill in the gaps or 'what does this part of the code do' style questions. Hints along the way in lectures.

Invigilated BYOD exam. Timetable to be communicated.

### Individual Assignment:

The individual assignment is now available

*Poetry generation ---- a popular use case for neural language models.... but using probabilistic language modeling*

Covers material from weeks 1 - 3

Due on Friday of Week 5

Details in your inbox. (Open the assignment document at this point).



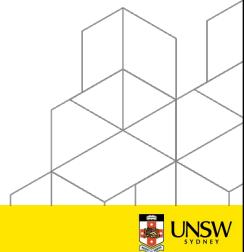
5

## ...and before we begin...

What do you know about Transformer?

Bullet points are okay.

.. and we will revisit them at the end of week 4.



6



7

**How many cyclists do you see in the picture?**

You paid attention to the parts of the **picture** that are important to the key term in the **question**.

--Cross-attention

The photograph shows a street corner with a rainbow-paved crosswalk. Three cyclists are visible: one on the left, one in the middle, and one on the right. Each cyclist is enclosed in a red square box. The background includes buildings and traffic lights.

Try this: How many piraxes do you see in the picture?

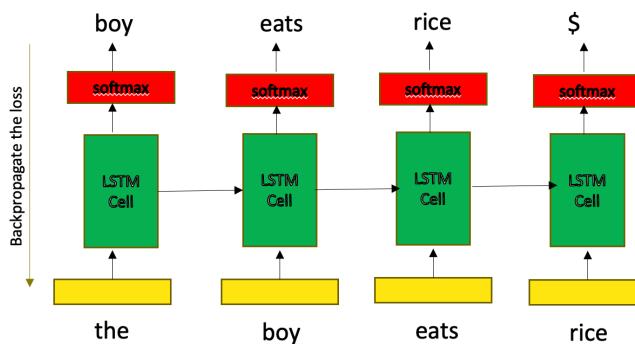
Just kidding.

The bottom right corner of the slide features a geometric pattern composed of overlapping triangles, creating a tessellated effect.

**UNSW** SYDNEY

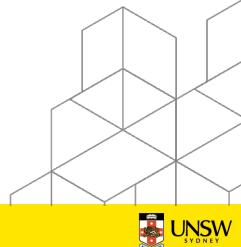
8

## Recap from week 2



Recurrent language models rely on sequential information passing.

**Fixed-length vectors become a bottleneck.**



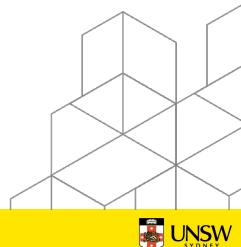
9

## Recap from week 2

“you can’t cram the meaning  
of a whole %&#&ing  
sentence into a single  
\$\*&@ing vector!”

— Ray Mooney (NLP professor at UT Austin)

**Fixed-length vectors become a bottleneck.**

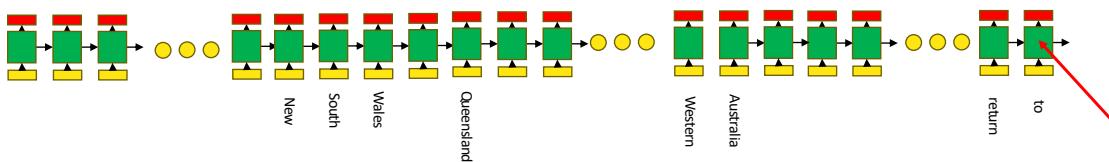


10

## Why a bottleneck?

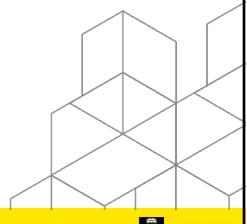
### - Long-range semantics

Thunderstorms are bringing intense deluges to parts of New South Wales and Queensland, Cyclone Lincoln is likely to re-emerge off the north coast of Western Australia, while elevated fire danger will return to \_\_\_\_\_



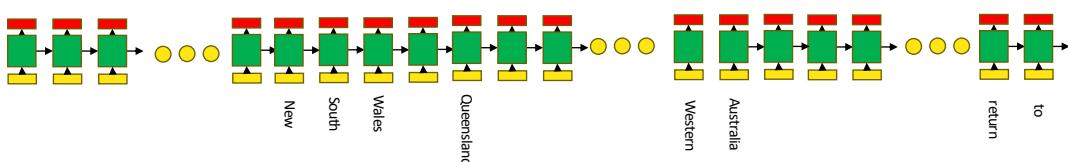
### - Coreferential pronoun selection

The restaurant first opened in 1997, and ran successfully until 2023 when \_\_\_\_\_



Thunderstorm example from: <https://www.abc.net.au/news/2024-02-21/thunderstorms-extreme-heat-cyclones-wild-weather-australia/103491052>

11

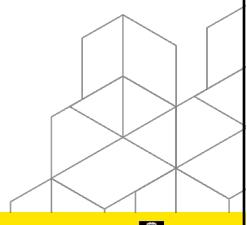


When generating the next word in the sequence,

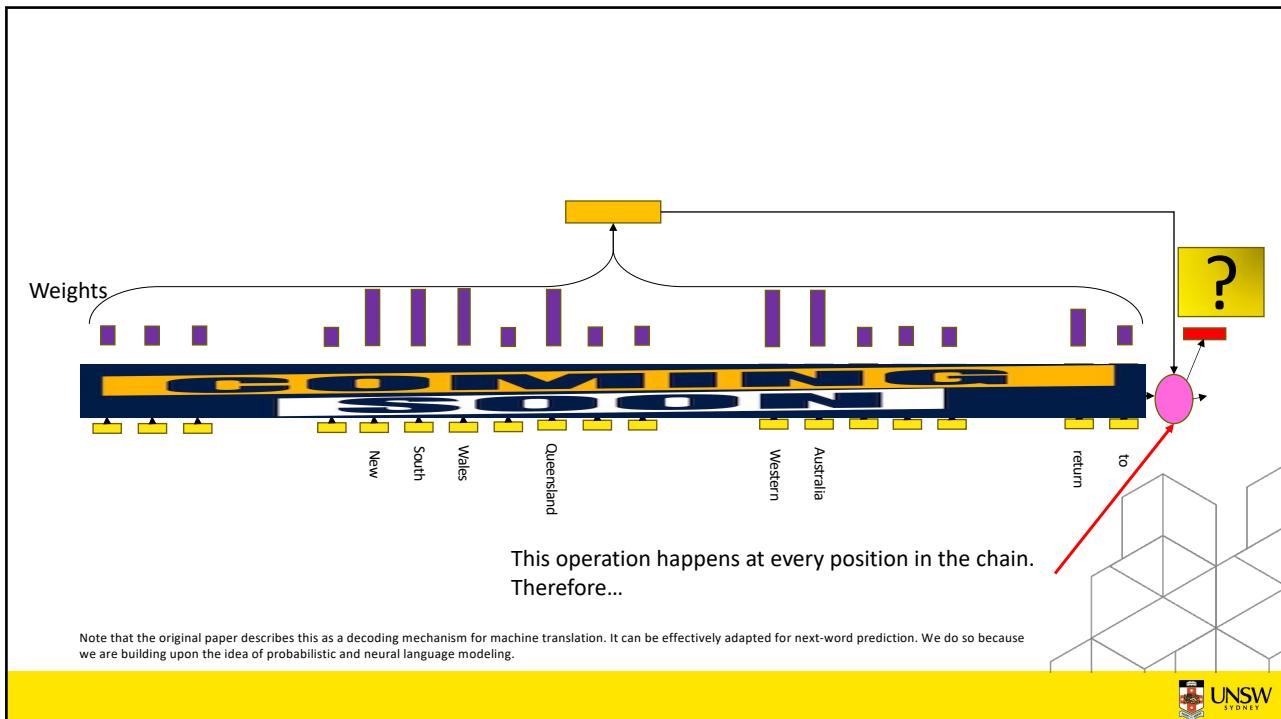
We want to pay 'attention' to some words.

i.e., some words are more important than others.

**Note:** Bidirectional LSTMs



12



13

## Attention

Attention mechanism is a technique that allows a neural layer to focus on specific parts of a sequence.

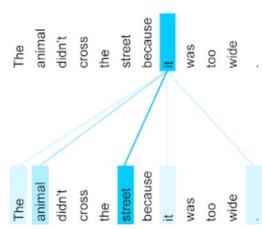
- I missed my bus because it was \_\_\_\_\_

*Which words are important when predicting the next word in the input?*

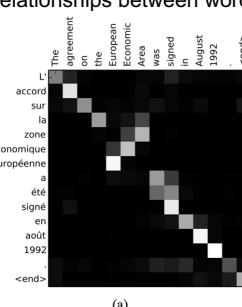
- Input: Who is the Prime Minister of Australia? Output: \_\_\_\_\_

*Which words are important when predicting the next word in the output?*

**Self-attention** models relationships between words in a sentence.



**Cross-attention** models relationships between words in pairs of sentences.



<https://blog.research.google/2017/08/transformer-novel-neural-network.html>  
Cross-attention from Bahdanau et al (2015)



14

## PyTorch

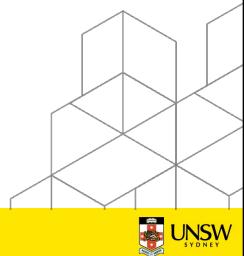
Originally developed by Meta; now a part of Linux Foundation

Tensor: Similar to numpy arrays, can be manipulated using GPUs

Alternative: TensorFlow



Demo time!



15



UNSW  
SYDNEY | Australia's  
Global University

Part 1  
**Attention in recurrent  
language modeling**

16

**Math board**

Advanced Reading (Optional): xLSTM  
<https://arxiv.org/pdf/2405.04517.pdf>

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate.". ICLR 2015.

17

**Math board**

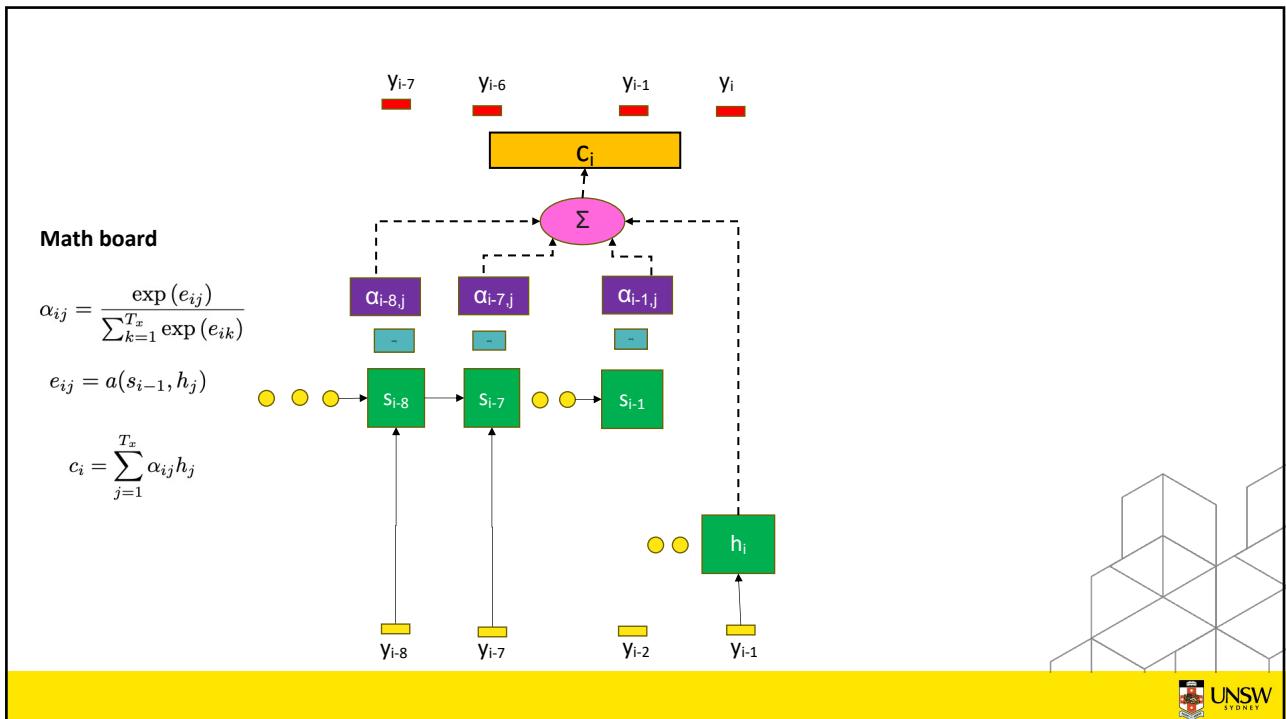
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$$e_{ij} = a(s_{i-1}, h_j)$$

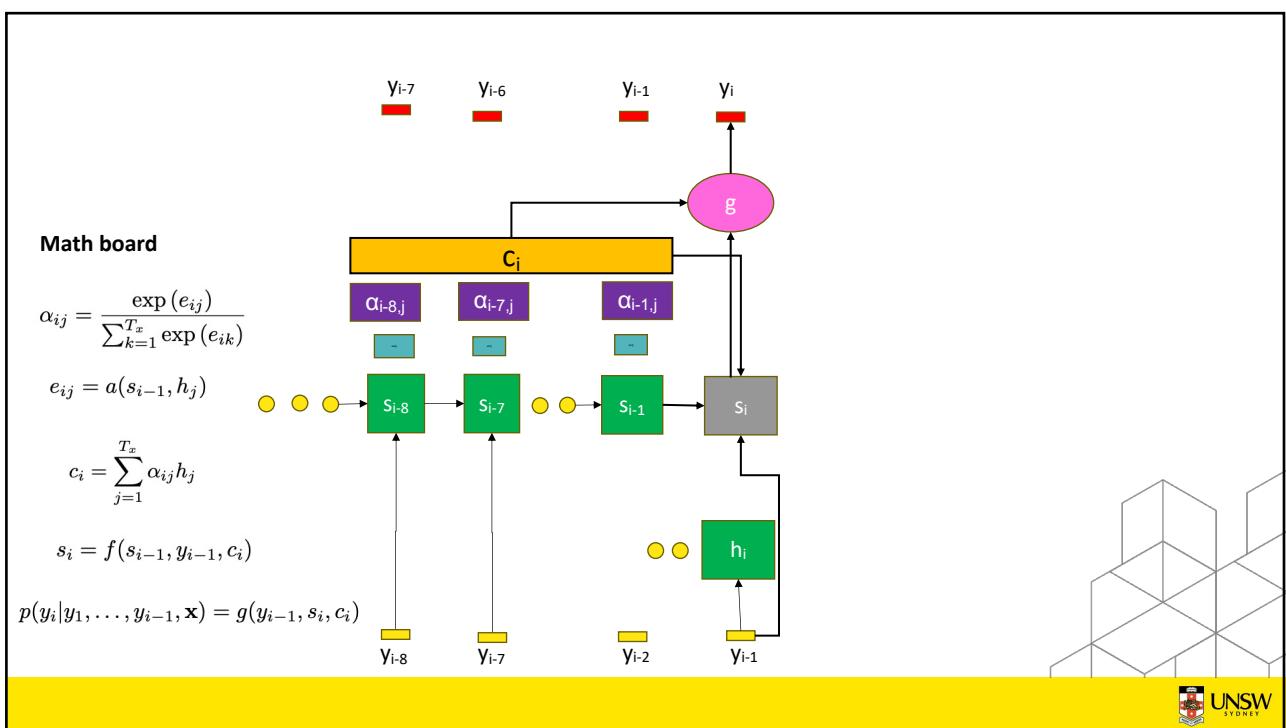
**attention between associated positions**

**energy between associated positions**  
*"how well do they vibe?"*

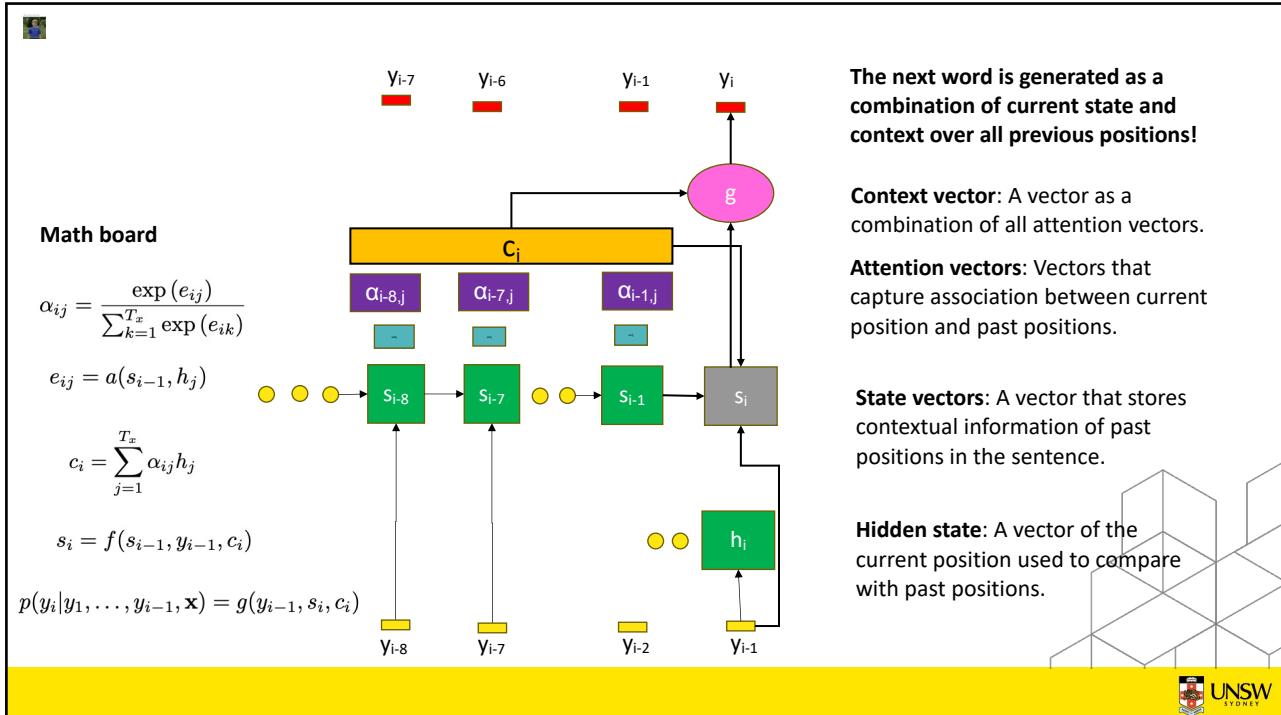
18



19



20



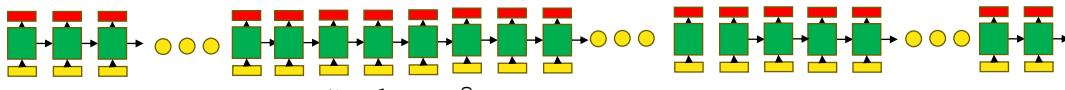
21



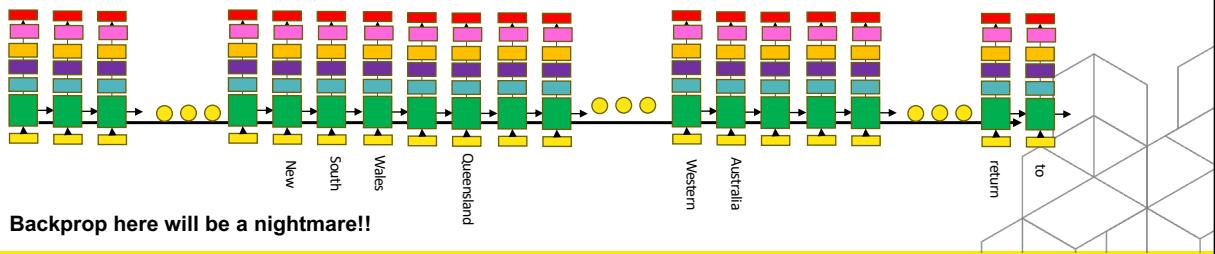
22

## Where did this land us?

We went from this...



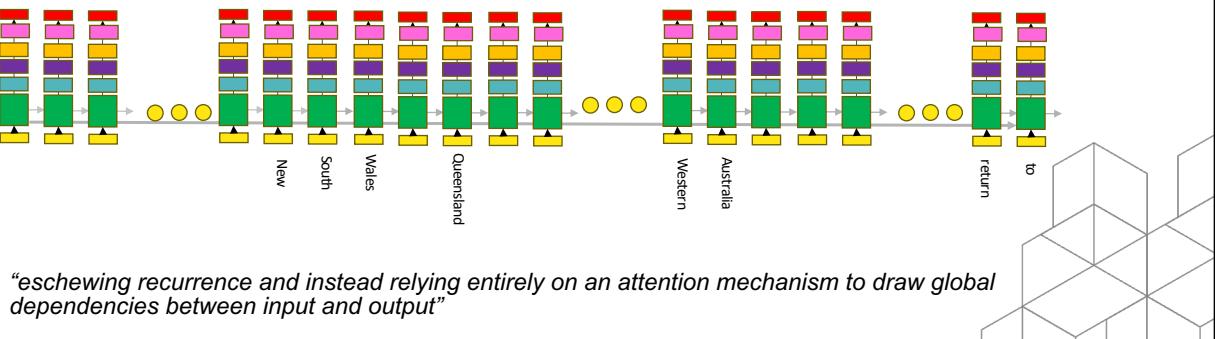
to this.....



23

## "Attention is all you need"

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).



24

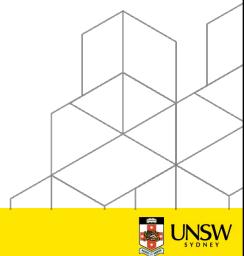
## A popular attention visualization

<https://colab.research.google.com/drive/1S-Cnm5KpH5rJXfDZH5flR0T58yT2dOF#scrollTo=rZgdi4xCK2X4>

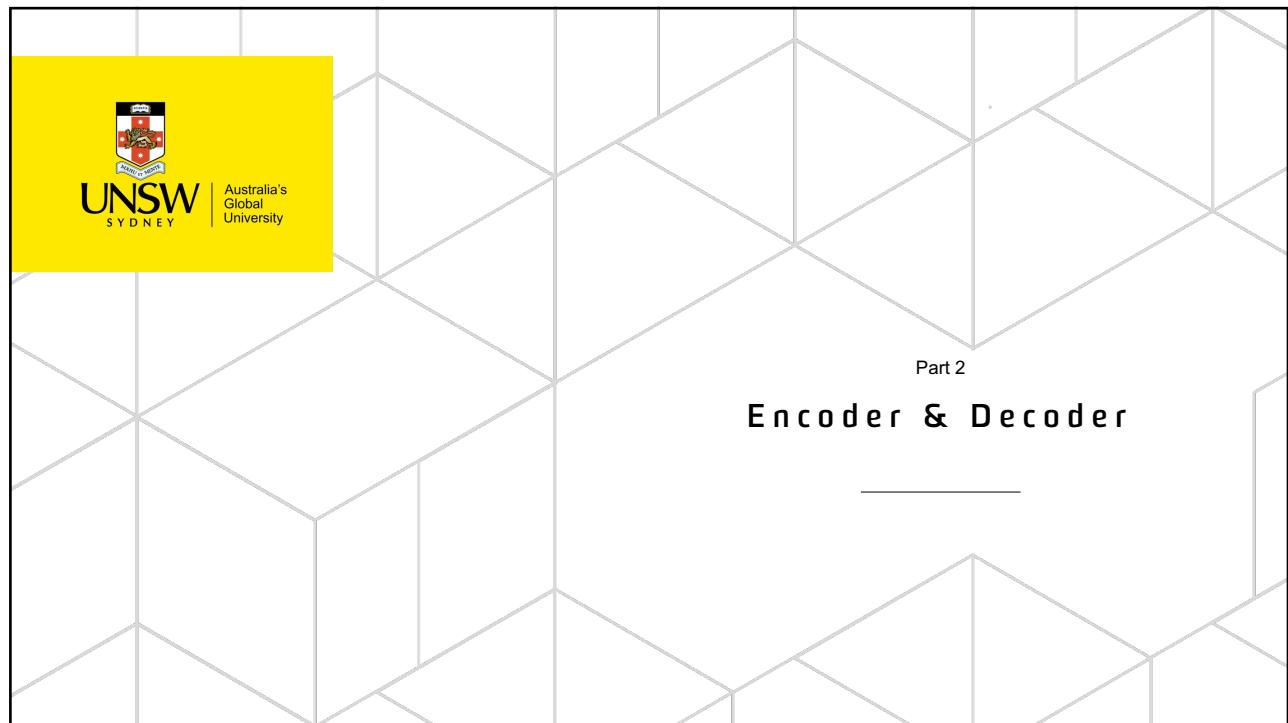
Source: <https://github.com/jessevig/bertviz>



Demo time!



25



26

## Transformer

Transformer is a sequence-to-sequence (seq2seq) model that uses the attention mechanism

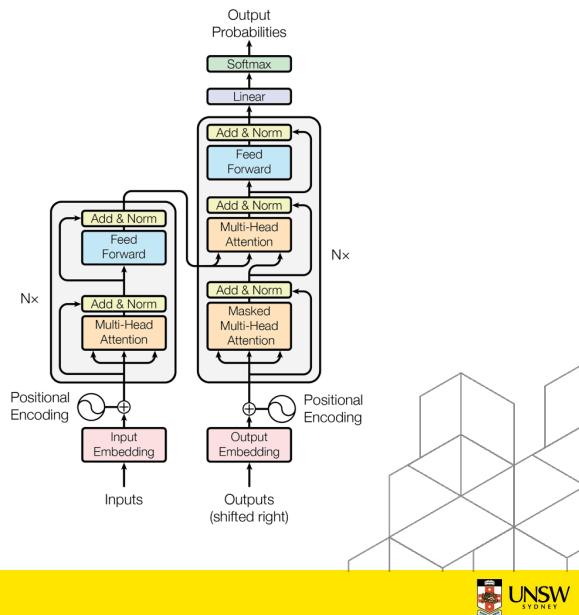
Transformer is a transducer that transforms a sentence into another

Originally evaluated on machine translation and other seq2seq tasks

Spawned several other kinds of NLP models (next module: encoder-only and decoder-only models)

Examples: BERT, PaLM, GPT, XLNet, OPT, BLOOM, etc.

... and non-NLP models



27

## Components of a Transformer

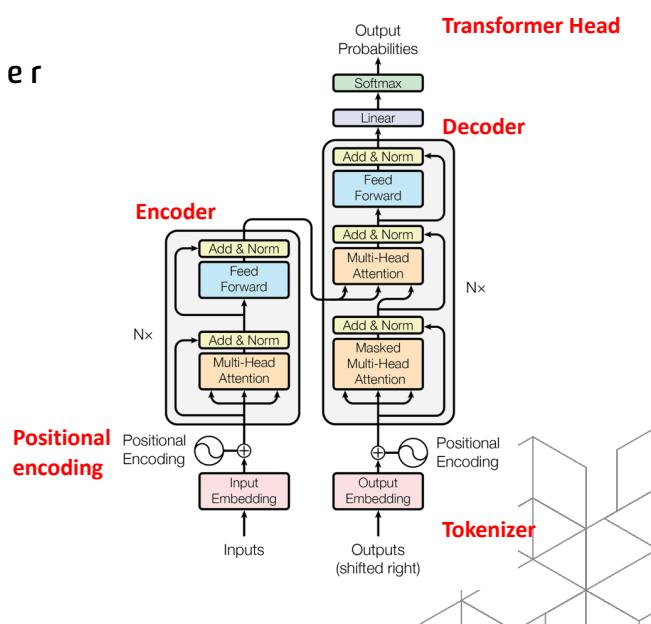
**Encoder:** A stack of encoders receive the input sentence and encode them into a hidden representation

**Decoder:** A stack of decoders receive the hidden representation from the encoder AND the text generated so far and produce the next token in the sequence.

**Tokenizer:** Represents a sentence as a set of tokens

**Positional Encoding:** Injects sequence information into the Transformer: a key reason why it can get rid of recurrence

**Transformer Head:** Attaches the prediction formulation used to optimize the model



[https://www.researchgate.net/figure/The-total-number-of-publications-in-the-field-of-NLP-The-statistics-are-from-the-Web-of\\_fig1\\_372092521](https://www.researchgate.net/figure/The-total-number-of-publications-in-the-field-of-NLP-The-statistics-are-from-the-Web-of_fig1_372092521)

28

## Pseudocode time

**Algorithm 8:**  $P \leftarrow \text{EDTransformer}(z, x|\theta)$

/\* Encoder-decoder transformer forward pass \*/

Input:  $z, x \in V^*$ , two sequences of token IDs.

Output:  $P \in (0, 1)^{N_v \times \text{length}(x)}$ , where the  $t$ -th column of  $P$  represents  $\hat{P}_\theta(x[t+1] | x[1:t], z)$ .

Hyperparameters:  $\ell_{\max}, L_{\text{enc}}, L_{\text{dec}}, H, d_e, d_{\text{mlp}} \in \mathbb{N}$

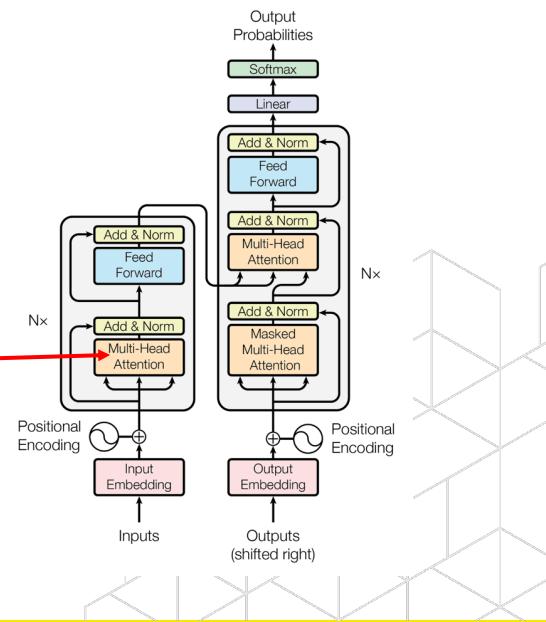
$z$ : Input;  $x$ : Output

Next word in the output sequence; conditional on input sequence and the output sequence so far.



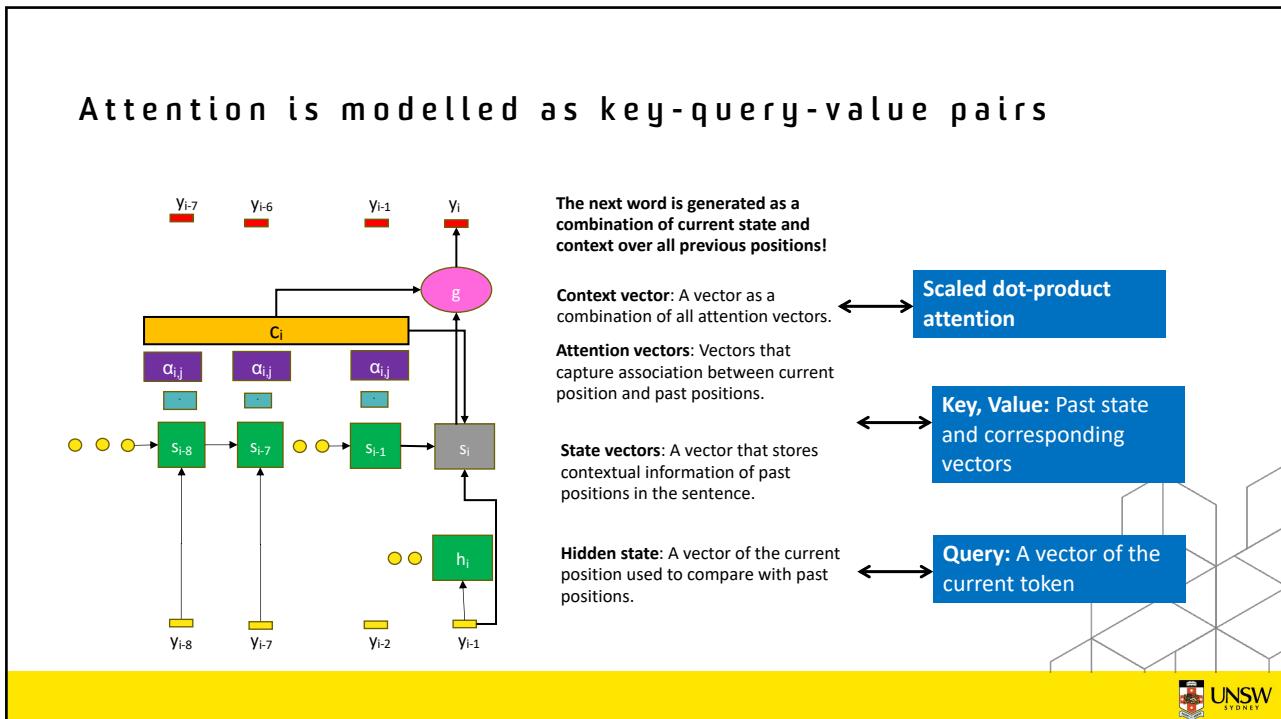
29

Let's go deep into the model



30

## Attention is modelled as key-query-value pairs

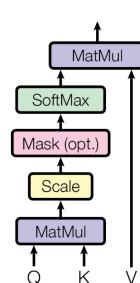


31

## Scaled-dot product attention

Query, keys, values, and output are all vectors

Scaled Dot-Product Attention



### Scaled dot-product attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



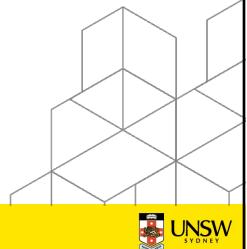
32

.. what? THREE vectors?

.... and the 'supermarket' example!



Key	Value
Key 1	Value 1
Key 2	Value 2
..	...



33

## A simple attention formulation

Many variants of attention

- Original formulation:  $a(\mathbf{q}, \mathbf{k}) = \mathbf{w}_2^T \tanh(W_1[\mathbf{q}; \mathbf{k}])$



Pen-and-paper time!

- Bilinear product:  $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T W \mathbf{k}$

Luong et al., 2015

- Dot product:  $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$

Luong et al., 2015

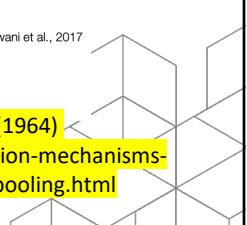
- Scaled dot product:  $a(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{|\mathbf{k}|}}$

Vaswani et al., 2017

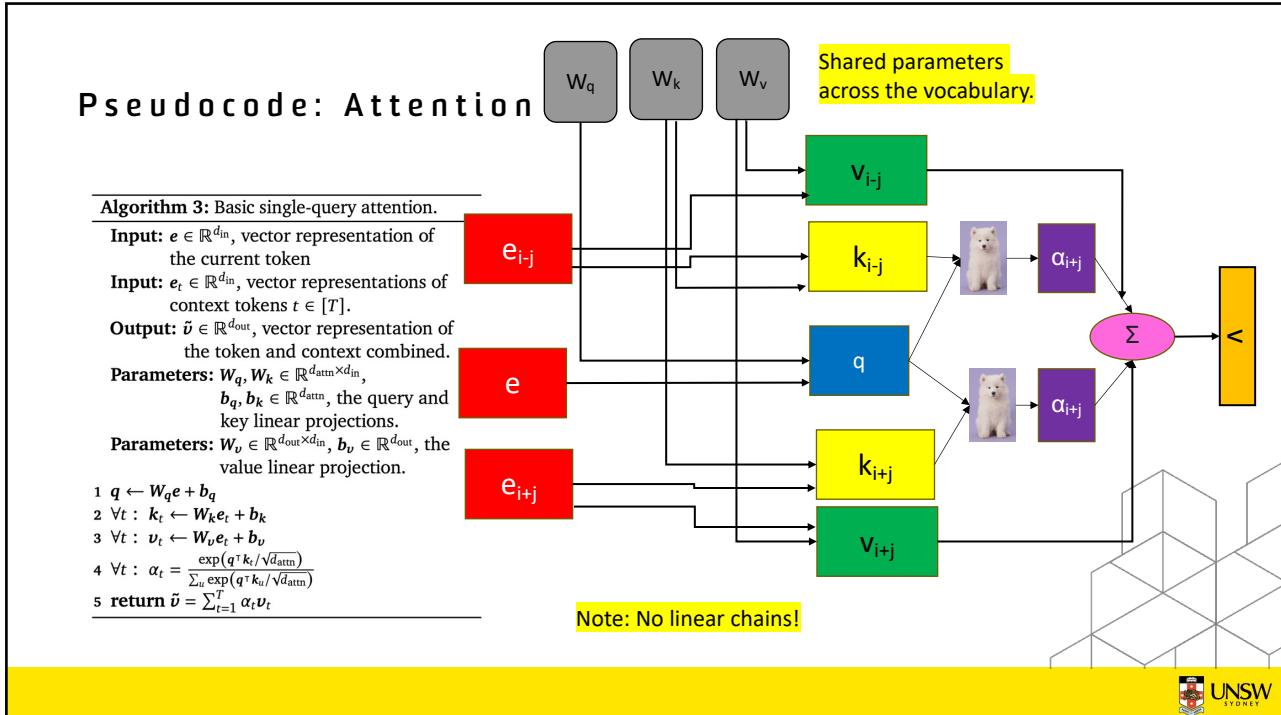
**Advanced Reading**  
**(Optional for the course)**

Nadaraya-Watson Estimator (1964)  
[https://d2l.ai/chapter\\_attention-mechanisms-and-transformers/attention-pooling.html](https://d2l.ai/chapter_attention-mechanisms-and-transformers/attention-pooling.html)

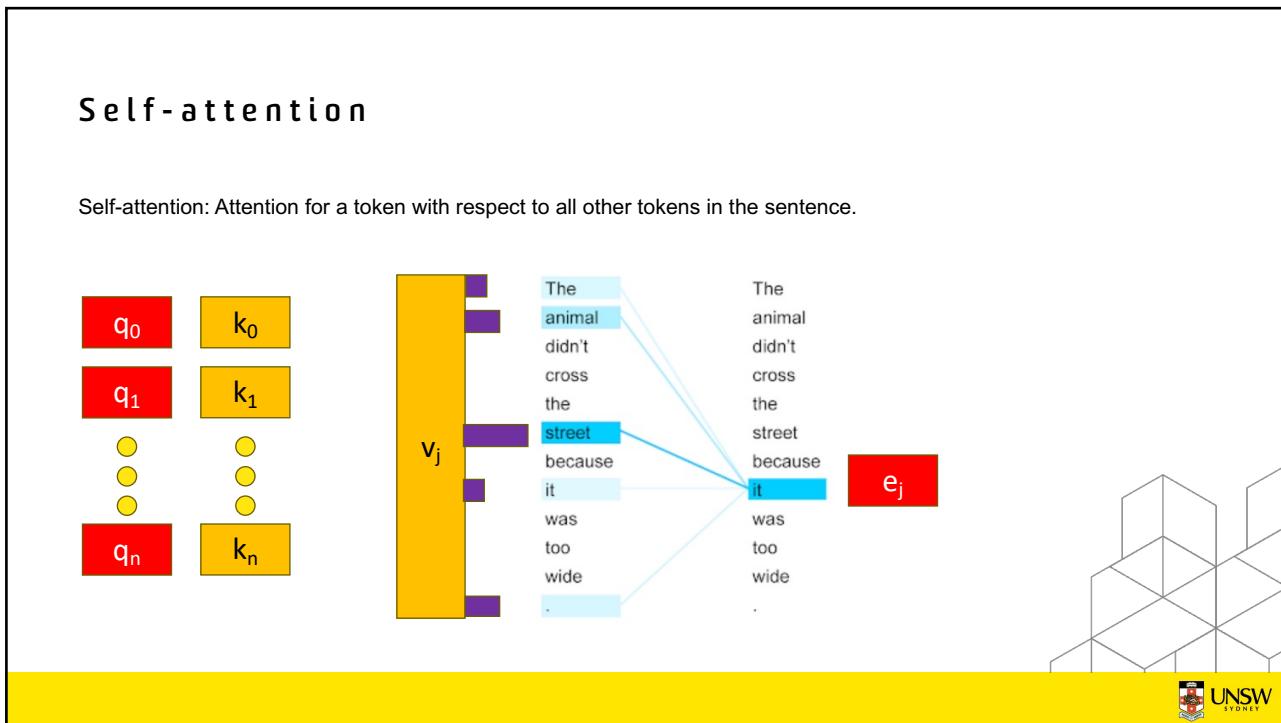
Credit: Mohit Iyyer, UMass



34



35

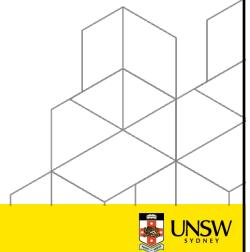


36

## Semantics is multi-dimensional

I love nature: green trees, pink flowers, red apples

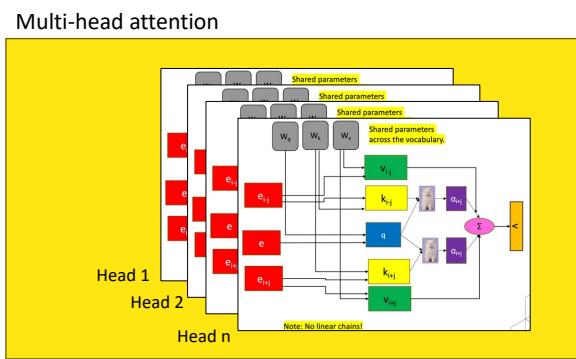
Remember syntagmatic and paradigmatic similarity?



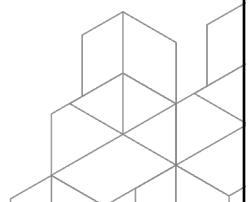
37

## Multi-head attention

Can we learn multiple projections to capture different aspects of semantics?  
They are independent. Therefore, can be parallelized.



The original Transformer paper uses 8 heads.



38

19

We need a way to combine these outputs....

---

**Algorithm 6:**  $\hat{e} \leftarrow \text{layer\_norm}(e|\gamma, \beta)$

---

```
/* Normalizes layer activations e. */
Input:  $e \in \mathbb{R}^{d_e}$ , neural network activations.
Output:  $\hat{e} \in \mathbb{R}^{d_e}$ , normalized activations.
Parameters:  $\gamma, \beta \in \mathbb{R}^{d_e}$ , element-wise scale and offset.
```

- 1  $m \leftarrow \sum_{i=1}^{d_e} e[i]/d_e$
- 2  $v \leftarrow \sum_{i=1}^{d_e} (e[i] - m)^2/d_e$
- 3 return  $\hat{e} = \frac{e-m}{\sqrt{v}} \odot \gamma + \beta$ , where  $\odot$  denotes element-wise multiplication.

---

Add & Normalize

Multi-head self-attention

Head 1 Head 2 Head n

Input Representation

UNSW SYDNEY

39

We need a way to combine these outputs....

---

...and some residual connections.

Add & Normalize

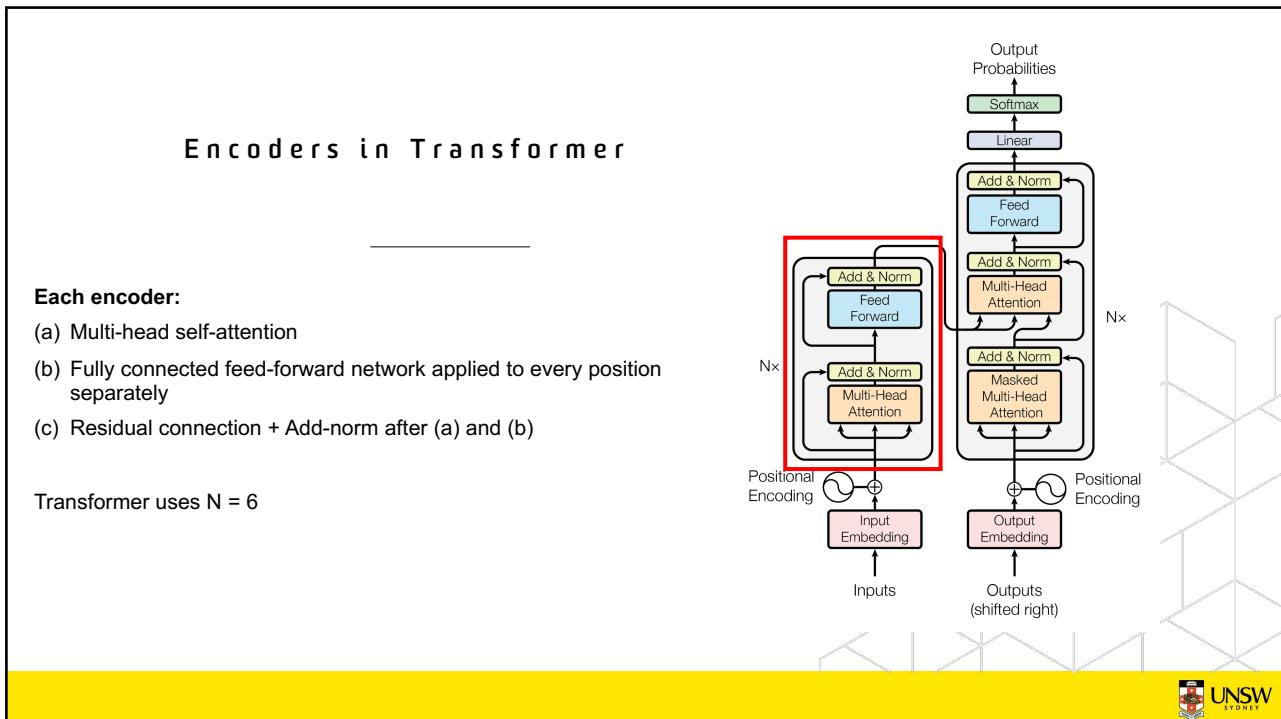
Multi-head self-attention

Head 1 Head 2 Head n

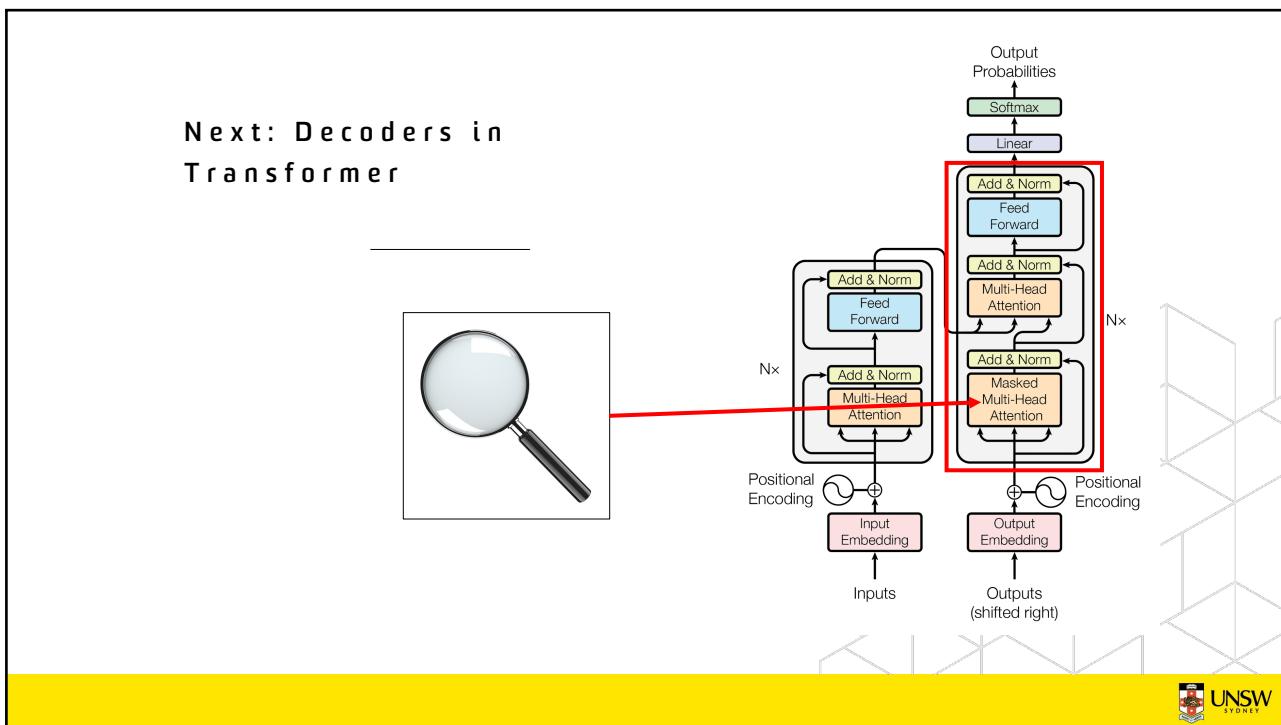
Input Representation

UNSW SYDNEY

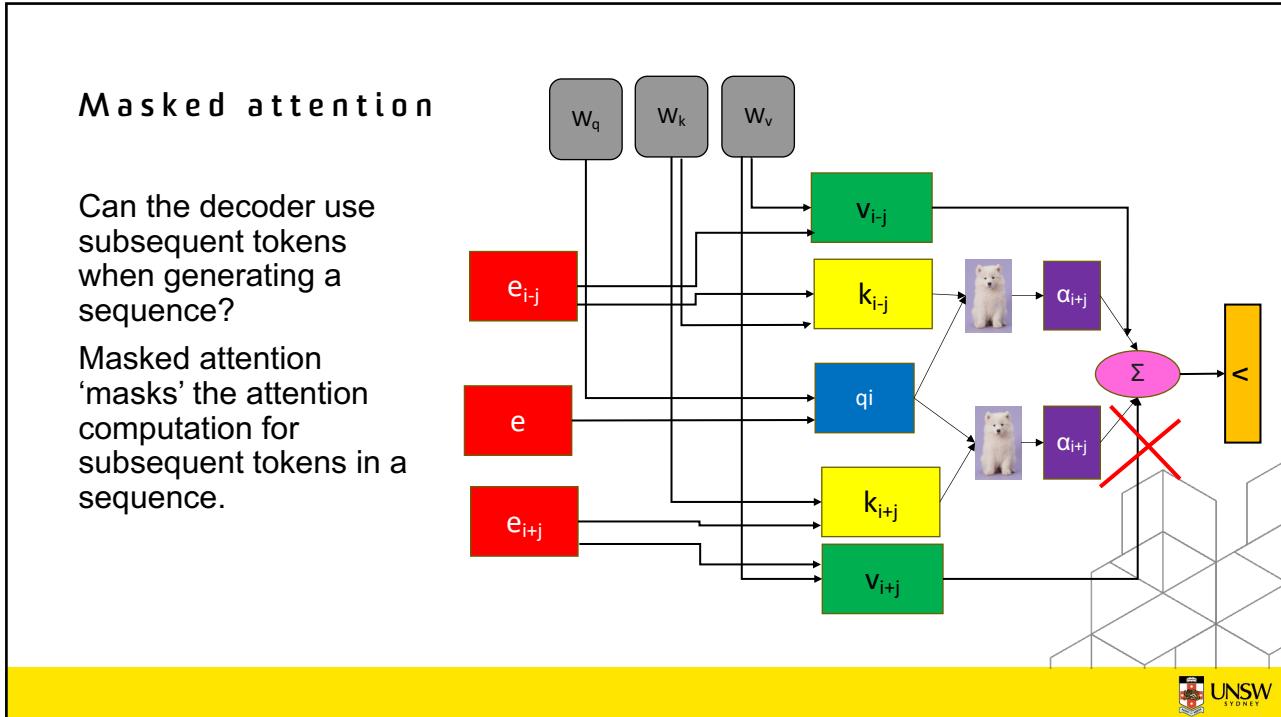
40



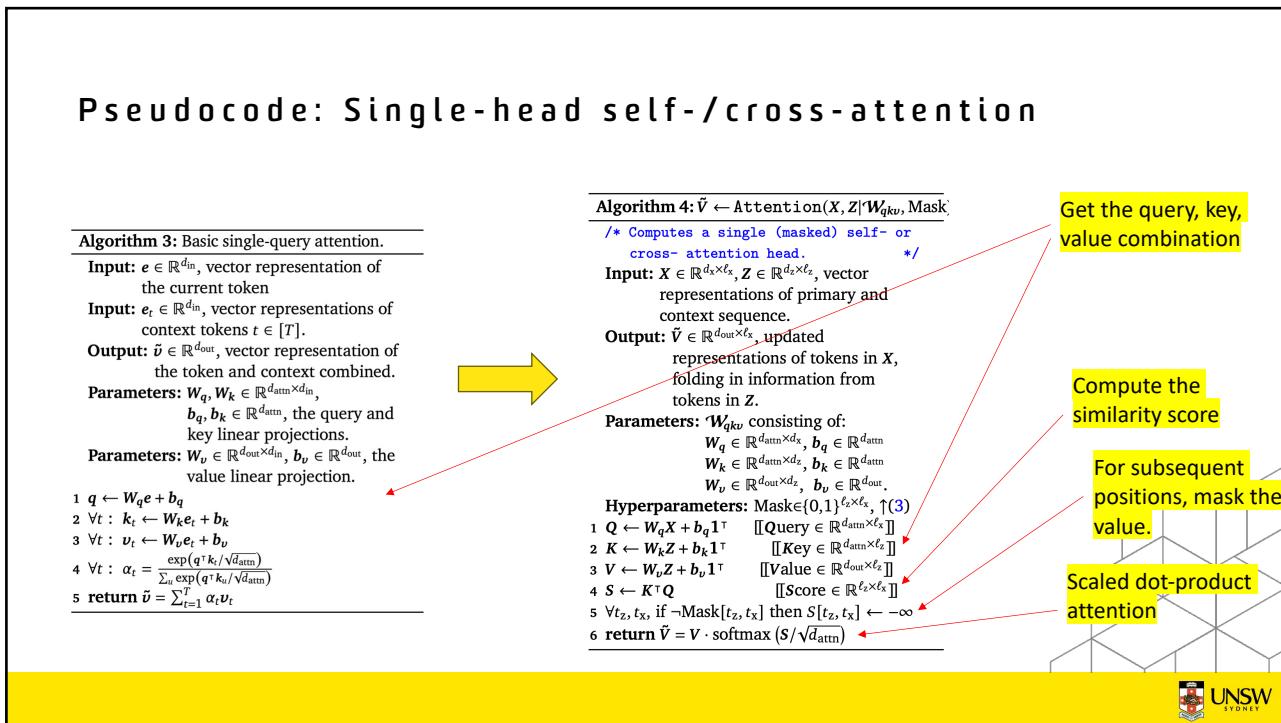
41



42



43



44

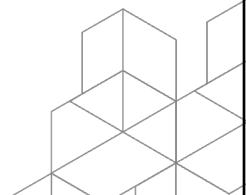
## ... and multi-head attention

```

Algorithm 5:  $\tilde{V} \leftarrow \text{MHAttention}(X, Z | \mathcal{W}, \text{Mask})$ 
/* Computes Multi-Head (masked) self-
   or cross- attention layer. */
Input:  $X \in \mathbb{R}^{d_x \times \ell_x}$ ,  $Z \in \mathbb{R}^{d_z \times \ell_z}$ , vector
   representations of primary and
   context sequence.
Output:  $\tilde{V} \in \mathbb{R}^{d_{\text{out}} \times \ell_x}$ , updated
   representations of tokens in  $X$ ,
   folding in information from
   tokens in  $Z$ .
Hyperparameters:  $H$ , number of
   attention heads
Hyperparameters: Mask  $\in \{0, 1\}^{\ell_z \times \ell_x}$ ,  $\uparrow(3)$ 
Parameters:  $\mathcal{W}$  consisting of
   For  $h \in [H]$ ,  $\mathcal{W}_{qkv}^h$  consisting of:
   |  $W_q^h \in \mathbb{R}^{d_{\text{attn}} \times d_x}$ ,  $b_q^h \in \mathbb{R}^{d_{\text{attn}}}$ ,
   |  $W_k^h \in \mathbb{R}^{d_{\text{attn}} \times d_x}$ ,  $b_k^h \in \mathbb{R}^{d_{\text{attn}}}$ ,
   |  $W_v^h \in \mathbb{R}^{d_{\text{mid}} \times d_x}$ ,  $b_v^h \in \mathbb{R}^{d_{\text{mid}}}$ .
    $W_o \in \mathbb{R}^{d_{\text{out}} \times H d_{\text{mid}}}$ ,  $b_o \in \mathbb{R}^{d_{\text{out}}}$ .
1 For  $h \in [H]$ :
2    $Y^h \leftarrow \text{Attention}(X, Z | \mathcal{W}_{qkv}^h, \text{Mask})$ 
3    $Y \leftarrow [Y^1; Y^2; \dots; Y^H]$ 
4 return  $\tilde{V} = W_o Y + b_o 1^\top$ 

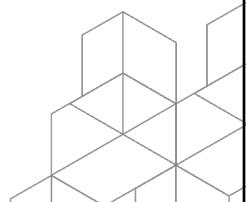
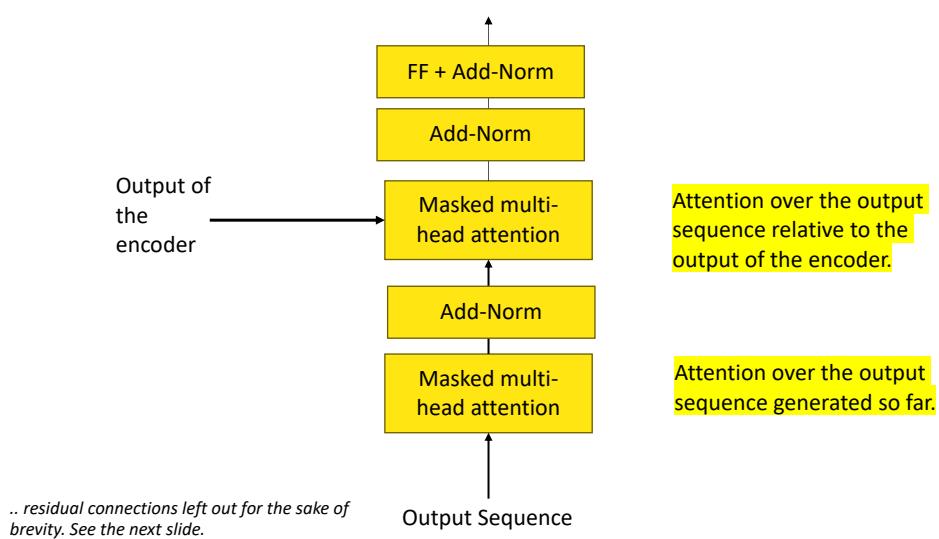
```

**h iterations that run in parallel!**

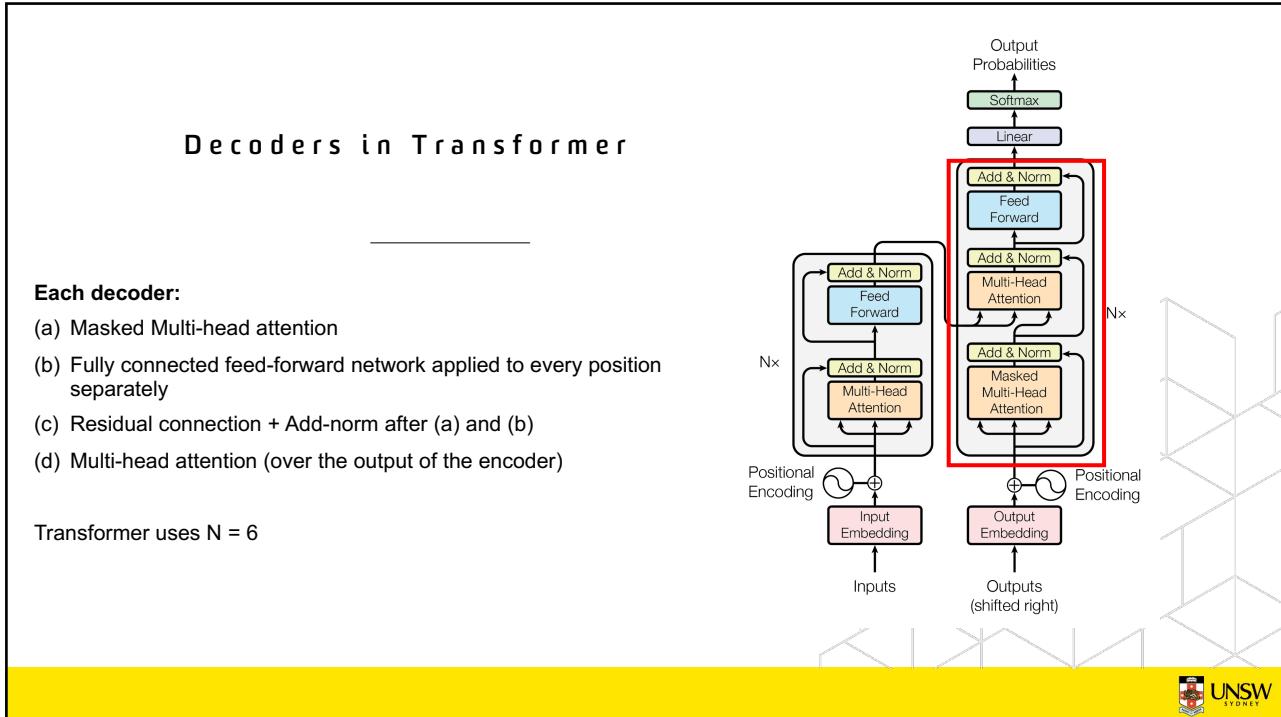


45

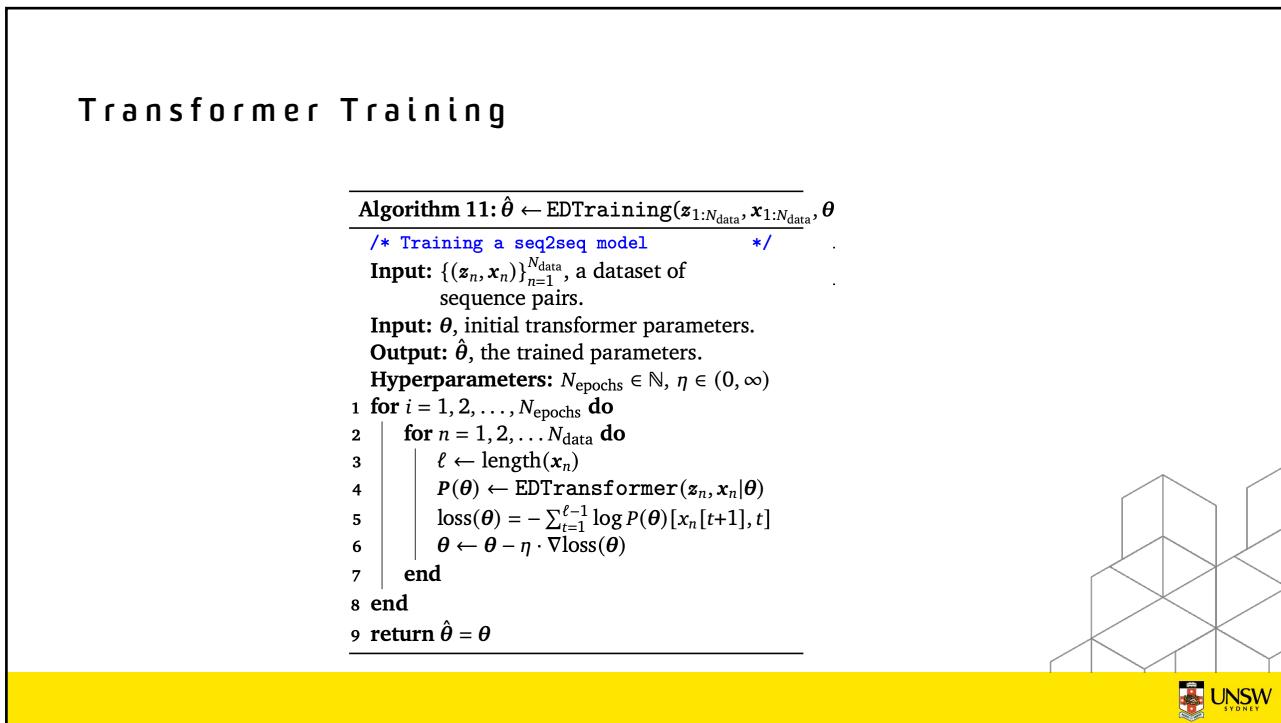
## Let's construct the decoder



46



47

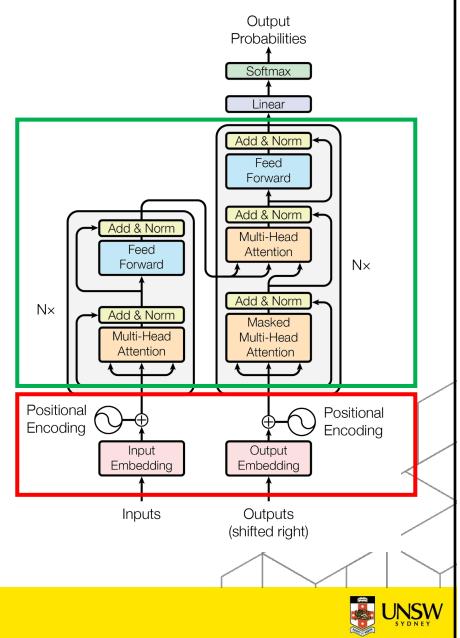


48



49

Let's now look at the input representations



50

## Tokenization

Tokenization: The process of splitting a sequence into tokens  
 Representations will be learned for tokens

Types of tokenization	Input	Output
Word-level tokenization	"Are you feeling sleepy?"	Are, you, feeling, sleepy
Character-level tokenization	"Are you feeling sleepy?"	A, r, e, , y, o.....
Sub-word tokenization	"Are you feeling sleepy?"	A, #r, #e, you, feel, #ing, sleep#, #y

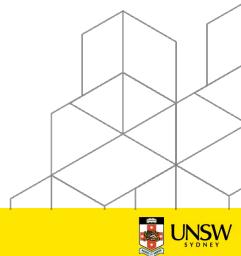
Remember stemming and lemmatization?

Sub-word tokenization

### Example: Byte-pair encoding

Originally developed as a text compression algorithm

Also, WordPiece Tokenizer: [https://www.youtube.com/watch?v=qpv6ms\\_t\\_1A](https://www.youtube.com/watch?v=qpv6ms_t_1A)



51

## Byte-pair encoding algorithm

### Intuition:

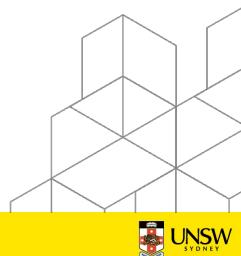
Retain the most common words as tokens

Break the less common words into sub-word tokens

### Why?

Each token will be represented by an ID. Common words will be represented wholly. See how this relates to text compression?

**Algorithm:** Corpus—based tokenizer



52

## Corpus-based tokenization

Step 1: Learn tokenization from a corpus (Learner)

Vocabulary = {Letters}

Find most common vocabulary pair

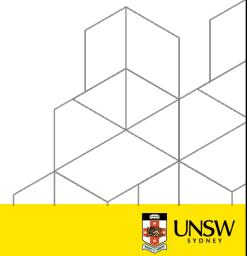
Add merged symbol to vocabulary

Replace all occurrences with the merged symbol

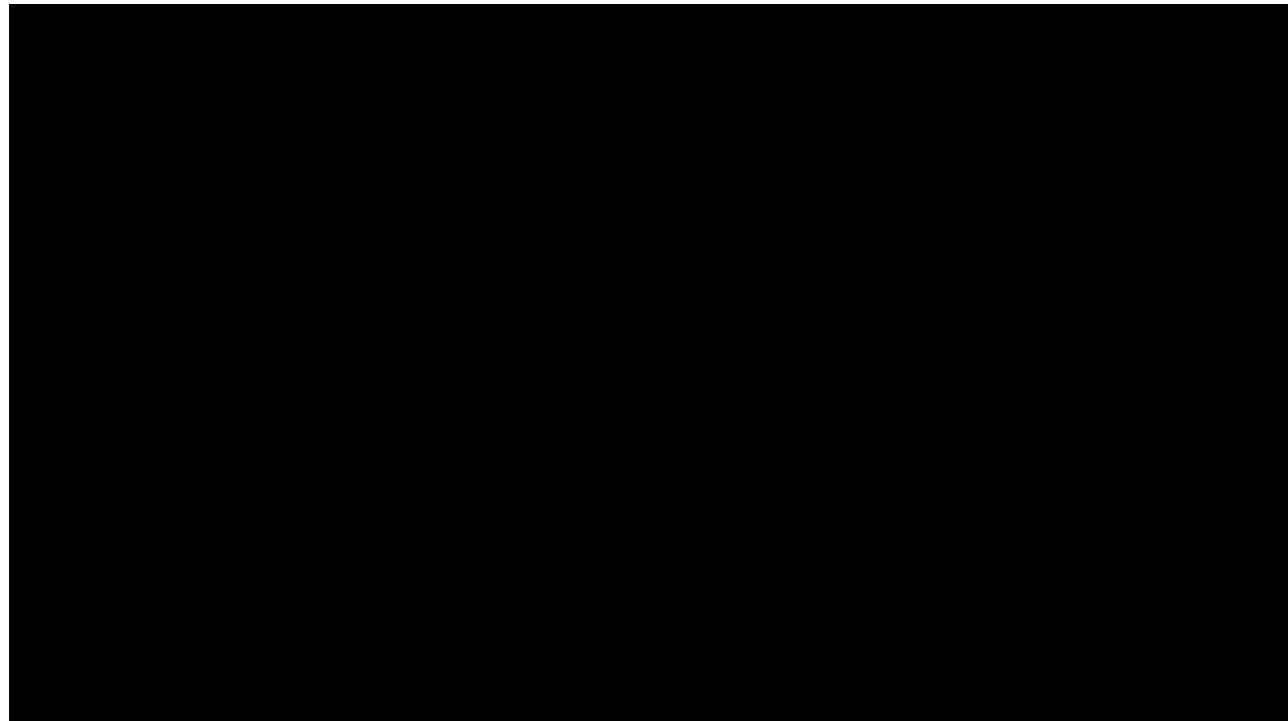
Store new symbols in the dictionary

Save merge rules

Step 2: Use the learned tokenizer on test sentences (Tokenizer)



53

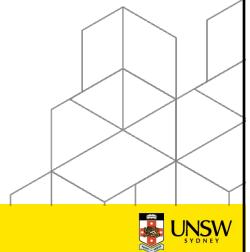


54

## Tokenizer



Demo time!



55

## Token embedding

Once input has been split into tokens, token embeddings are learned during training.  
Therefore, during inference:

---

### Algorithm 1: Token embedding.

---

**Input:**  $v \in V \cong [N_V]$ , a token ID.

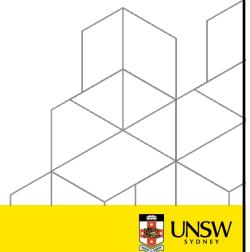
**Output:**  $e \in \mathbb{R}^{d_e}$ , the vector representation of the token.

**Parameters:**  $W_e \in \mathbb{R}^{d_e \times N_V}$ , the token embedding matrix.

---

1 **return**  $e = W_e[:, v]$

---



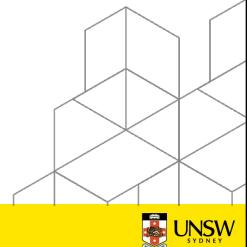
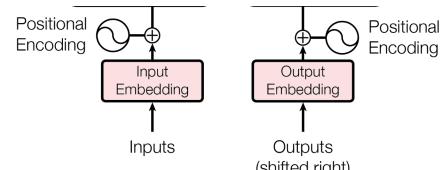
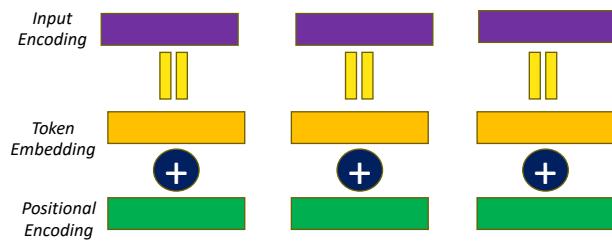
56

## What next...

Tokens have an embedding

....but remember that there's no recurrence.

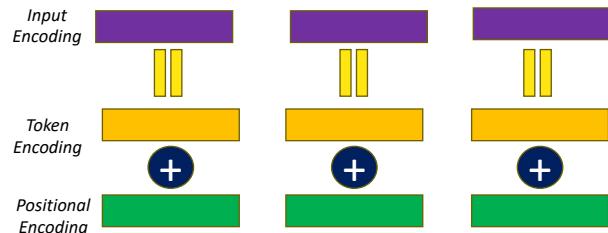
So, Transformer 'injects' a positional encoding.



57

## Positional encoding

Vectors corresponding to each position in a sequence.  
Assume token embeddings are of length 50.



If `max_length = 25`, we would like to have 25 position vectors of length 50, such that...

- Each position vector is unique
- Should generalize to longer lengths

Original Transformer paper uses fixed positional encoding. Learnable encodings used in newer adaptations.



58

29

## Fixed positional encoding

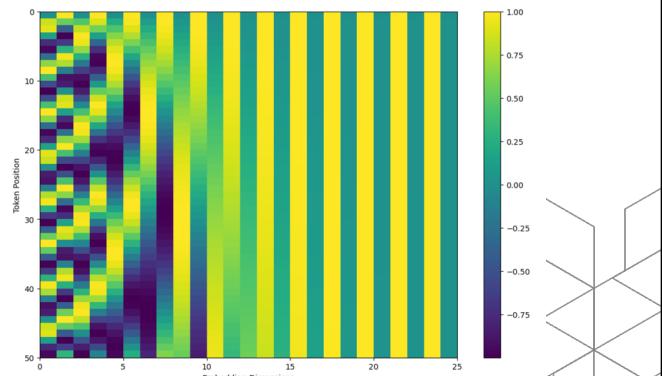
"we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k,  $P_{Epos+k}$  can be represented as a linear function of  $P_{Epos}$ ."

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



Demo time!



59



UNSW  
SYDNEY | Australia's  
Global University

Part 3  
**Transformer: Looking  
Back**

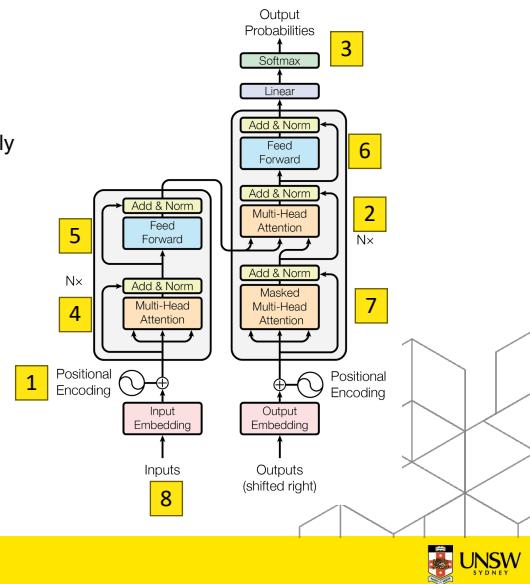
60

### Quiz: Match textual descriptors (A)-(D) with one or more of (1)-(8) in the diagram

Example: Textual sequence that is input to the Transformer: 8

- (A) Implemented as fixed embeddings
- (B) Attention module that looks at preceding tokens in the sequence only
- (C) Language model head that predicts the output sequence
- (D) Uses key, query and values

Write all that apply.



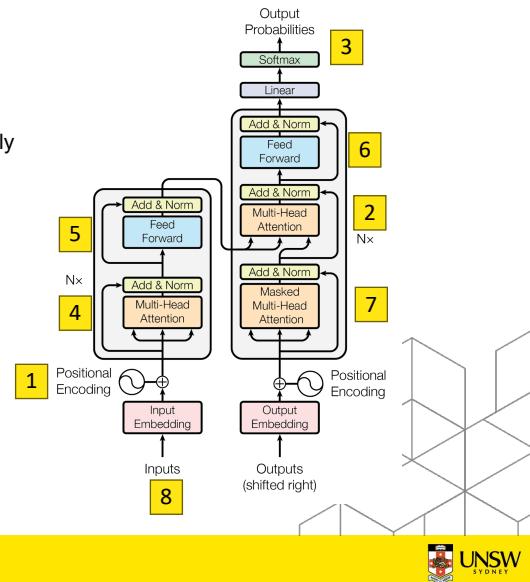
61

### Quiz: Match textual descriptors (A)-(D) with one or more of (1)-(8) in the diagram

Example: Textual sequence that is input to the Transformer: 8

- (A) Implemented as fixed embeddings: [1]
- (B) Attention module that looks at preceding tokens in the sequence only  
[7]
- (C) Language model head that predicts the output sequence: [3]
- (D) Uses key, query and values: [2, 4, 7]

Select all that apply.



62

Let's look at the pseudocode  
of the Transformer



63

## Specification: Reminder & the full thing

---

**Algorithm 8:**  $P \leftarrow \text{EDTransformer}(z, x|\theta)$ 


---

*/\* Encoder-decoder transformer forward pass*

*\*/*

Input:  $z, x \in V^*$ , two sequences of token IDs.

$z$ : Input;  $x$ : Output

Output:  $P \in (0, 1)^{N_v \times \text{length}(x)}$ , where the  $t$ -th column of  $P$  represents  $\hat{P}_\theta(x[t+1] | x[1:t], z)$ .

Hyperparameters:  $\ell_{\max}, L_{\text{enc}}, L_{\text{dec}}, H, d_e, d_{\text{mlp}} \in \mathbb{N}$

Parameters:  $\theta$  includes all of the following parameters:

$W_e \in \mathbb{R}^{d_e \times N_v}$ ,  $W_p \in \mathbb{R}^{d_e \times \ell_{\max}}$ , the token and positional embedding matrices.

For  $l \in [L_{\text{enc}}]$ :

- |  $W_l^{\text{enc}}$ , multi-head encoder attention parameters for layer  $l$ , see (4),
- |  $\gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2 \in \mathbb{R}^{d_e}$ , two sets of layer-norm parameters,
- |  $W_l^{\text{mlp}} \in \mathbb{R}^{d_{\text{mlp}} \times d_e}$ ,  $b_{\text{mlp}1}^l \in \mathbb{R}^{d_{\text{mlp}}}$ ,  $W_{\text{mlp}2}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}$ ,  $b_{\text{mlp}2}^l \in \mathbb{R}^{d_e}$ , MLP parameters.

For  $l \in [L_{\text{dec}}]$ :

- |  $W_l^{\text{dec}}$ , multi-head decoder attention parameters for layer  $l$ , see (4),
- |  $W_l^{\text{c/d}}$ , multi-head cross-attention parameters for layer  $l$ , see (4),
- |  $\gamma_l^3, \beta_l^3, \gamma_l^4, \beta_l^4, \gamma_l^5, \beta_l^5 \in \mathbb{R}^{d_e}$ , three sets of layer-norm parameters,
- |  $W_l^{\text{mlp3}} \in \mathbb{R}^{d_{\text{mlp}} \times d_e}$ ,  $b_{\text{mlp3}}^l \in \mathbb{R}^{d_{\text{mlp}}}$ ,  $W_{\text{mlp4}}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}$ ,  $b_{\text{mlp4}}^l \in \mathbb{R}^{d_e}$ , MLP parameters.

$W_u \in \mathbb{R}^{N_v \times d_e}$ , the unembedding matrix.

Next word in the output sequence; conditional on

Next word in the sequence output sequence;  $\exists$  output

conditional on uence so far.

input sequence

and the output

sequence so far.

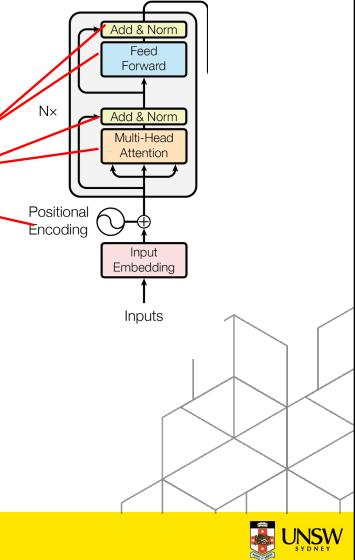
64

## Encoder

```

/* Encode the context sequence:
1  $\ell_z \leftarrow \text{length}(z)$ 
2 for  $t \in [\ell_z]$  :  $e_t \leftarrow W_e[:, z[t]] + W_p[:, t]$ 
3  $Z \leftarrow [e_1, e_2, \dots e_{\ell_z}]$ 
4 for  $l = 1, 2, \dots, L_{\text{enc}}$  do
5    $Z \leftarrow Z + \text{MHAttention}(Z | \mathcal{W}_l^{\text{enc}}, \text{Mask} \equiv 1)$ 
6   for  $t \in [\ell_z]$  :  $Z[:, t] \leftarrow \text{layer\_norm}(Z[:, t] | \gamma_l^1, \beta_l^1)$ 
7    $Z \leftarrow Z + W_{\text{mlp2}}^l \text{ReLU}(W_{\text{mlp1}}^l Z + b_{\text{mlp1}}^l \mathbf{1}^\top) + b_{\text{mlp2}}^l \mathbf{1}^\top$ 
8   for  $t \in [\ell_z]$  :  $Z[:, t] \leftarrow \text{layer\_norm}(Z[:, t] | \gamma_l^2, \beta_l^2)$ 
9 end

```



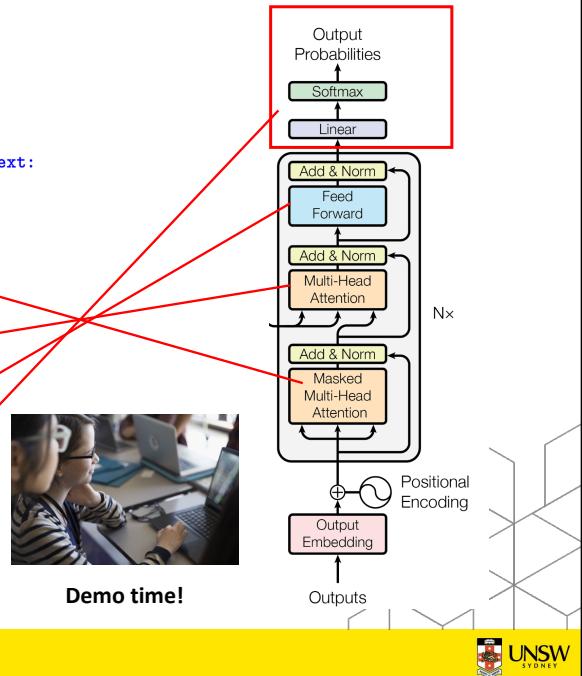
65

## Decoder

```

/* Decode the primary sequence, conditioning on the context:
10  $\ell_x \leftarrow \text{length}(x)$ 
11 for  $t \in [\ell_x]$  :  $e_t \leftarrow W_e[:, x[t]] + W_p[:, t]$ 
12  $X \leftarrow [e_1, e_2, \dots e_{\ell_x}]$ 
13 for  $i = 1, 2, \dots, L_{\text{dec}}$  do
14    $X \leftarrow X + \text{MHAttention}(X | \mathcal{W}_l^{\text{dec}}, \text{Mask}[t, t'] \equiv [[t \leq t']] )$ 
15   for  $t \in [\ell_x]$  :  $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] | \gamma_l^3, \beta_l^3)$ 
16    $X \leftarrow X + \text{MHAttention}(X, Z | \mathcal{W}_l^{e/d}, \text{Mask} \equiv 1)$ 
17   for  $t \in [\ell_x]$  :  $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] | \gamma_l^4, \beta_l^4)$ 
18    $X \leftarrow X + W_{\text{mlp4}}^l \text{ReLU}(W_{\text{mlp3}}^l X + b_{\text{mlp3}}^l \mathbf{1}^\top) + b_{\text{mlp4}}^l \mathbf{1}^\top$ 
19   for  $t \in [\ell_x]$  :  $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] | \gamma_l^5, \beta_l^5)$ 
20 end
/* Derive conditional probabilities and return:
21 return  $P = \text{softmax}(W_u X)$ 

```



Demo time!



66

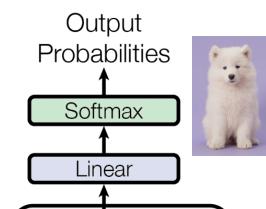
## Final note: The model head

Model head

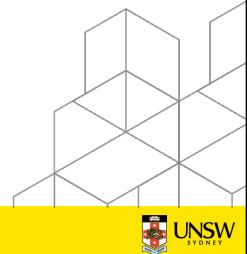
A linear layer that predicts probability over the vocabulary.

Loss gets calculated at this stage

"Language model head is the circuitry we do to do language model



You will see the model head in Week 4 when we discuss BERT and GPT.



67

## Configuration Parameters

Vocabulary: 37000 tokens

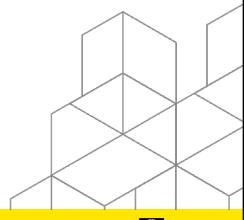
8 NVIDIA P100 GPUs

$d_{model} = 512$

$N = 6$

Trained on the English-to-German, English-to-French WMT benchmark.

It is customary to include configuration parameters in research papers to ensure reproducibility.



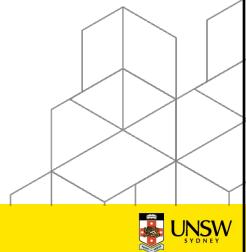
68

## Results

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.8</b>	$2.3 \cdot 10^{19}$	

BLEU: Higher the better. FLOPs: Lower the better



69

Let's look at another view  
of Transformer.

<https://jalammar.github.io/illustrated-transformer/>

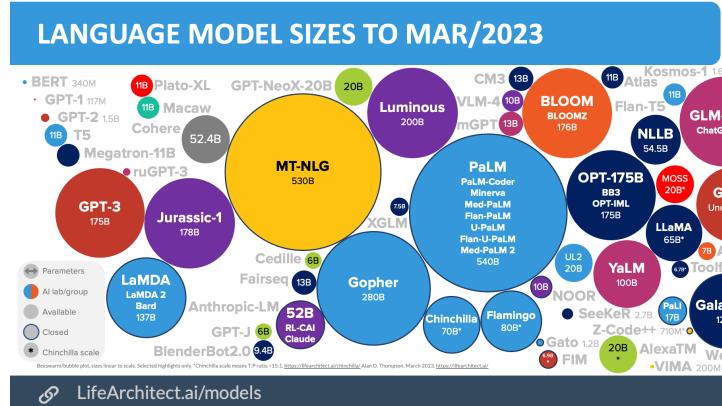


70

## The impact of Transformer on NLP

Transformer ushered a new revolution in NLP models

Advent of 'large' language models



<https://www.fullydistributed.co/p/language-models-size-matters>



71

## The impact of Transformer on the world

LLMs are not going to destroy the human race

It's just a chatbot, dude.

NOAH SMITH  
8 MAR 2023



"This is the sort of thing you lifeforms enjoy, is it?" — Marvin the Paranoid Android

Five years ago, if you told me that in 2023 I would be writing a post arguing that chatbots are not likely to wipe out humanity, I would have said "...Yeah, sounds about right." Not because I knew anything about progress in large language models five years ago, but because on this crazy internet of ours, it's always something.

ARTIFICIAL INTELLIGENCE / TECH / US & WORLD

The EU AI Act passed – now comes the waiting



/ Delays in implementing the AI Act means nothing changes for now.

By Emilia David, a reporter who covers AI. Prior to joining The Verge, she covered the intersection between technology, finance, and the economy.

Dec 15, 2023, 8:47 AM GMT+11



[1] Blog by Noah Smith  
[2] <https://www.theverge.com/2023/12/14/24001919/eu-ai-act-foundation-models-regulation-data>



72

## Many types of attentions....

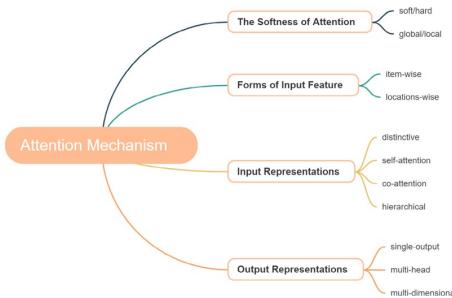
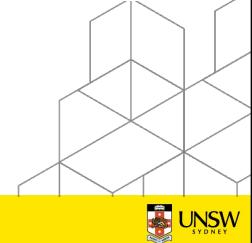


Fig. 1. Several typical approaches to attention mechanisms.

**Table 1**

Summary of score function  $f$ . Here,  $\mathbf{k}$  is an element of  $\mathbf{K}$ ,  $\mathbf{v}, \mathbf{b}, \mathbf{W}_1, \mathbf{W}_2$  are learnable parameters,  $d_k$  is the dimension of the input vector. The  $\text{act}$  is a nonlinear activation function, such as tanh and ReLU.

Name	Equation	Ref.
Additive	$f(q, k) = \mathbf{v}^\top \text{act}(\mathbf{W}_1 \mathbf{k} + \mathbf{W}_2 q + b)$	[15]
Multiplicative (dot-product)	$f(q, k) = \mathbf{q}^\top \mathbf{k}$	[15]
Scaled multiplicative	$f(q, k) = \frac{\mathbf{q}^\top \mathbf{k}}{\sqrt{d_k}}$	[16]
General	$f(q, k) = \mathbf{q}^\top \mathbf{W}\mathbf{k}$	[14]
Concat	$f(q, k) = \mathbf{v}^\top \text{act}(\mathbf{W}[\mathbf{k}; \mathbf{q}] + b)$	[14]
Location-based	$f(q, k) = f(q)$	[14]
Similarity	$f(q, k) = \frac{q^\top k}{\ \mathbf{q}\  \ \mathbf{k}\ }$	[60]

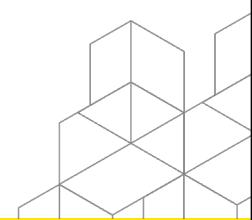


Niu, Zhaoyang, Guoqiang Zhong, and Hui Yu. "A review on the attention mechanism of deep learning." *Neurocomputing* 452 (2021): 48-62.

73

## Many types of position encodings....

Model	Data Structure	Relative Position	Invert. Mat.	Learnable	Recurrent	Unbounded	#Params
Transformer w/o emb. (Vaswani et al. 2017)							
BERT (Devlin et al. 2019)	A	APE	✓	✗	✗		$t_{max}d$
Reformer (Kitaev, Kaiser, and Levskaya 2020)							$(d - d_1)t_{max} + d_1t_1$
FLOATER (Lin et al. 2020)	A	APE	✓	✗	✓	✓	0 or more
Shortformer (Press, Smith, and Lewis 2021)	A	APE	✗	✓	✓	✓	0
Wang et al. (2020)	A	-	✓	✓	✓	✓	$2d$
Shazeer et al. and Vaswani (2018) (abs)	A	MAM	✓	✓	✓	✗	$2d^2 \  \mathbf{M} \ $
Shazeer, Uszkoreit, and Vaswani (2019) (rel)	R	MAM	✓	✓	✓	✗	$2(2t_{max} - 1)d$
T5 (Raffel et al. 2020)	R	MAM	✓	✓	✓	✗	$(2t_{max} - 1)b$
Huang et al. (2020)							$b(2t_{max} - 1)$
DeBERTa (He et al. 2021)	B	Both	✓	✓	✓	✗	$3t_{max}d$
Transformer XL (Dai et al. 2019)	R	MAM	✓	✓	✓	✓	$2d + d^2 b$
TENER (Yan et al. 2019)	R	-	✓	✗	✓	✓	$6d^2 + 3d$
DA-Transformer (Wu, Wu, and Huang 2021)	R	MAM	✓	✓	✓	✗	$3kd^2 + kd$
TUPE (Ke, He, and Liu 2021)	B	MAM	✓	✗	✗	✗	$2d^2 + k_{max}(d + 2)$
RNN-Transf. (Nishii and Yoshihaga 2019)	R	-	✓	✗	✓	✓	$6d^2 + 3d$
SPI (Lukosz et al. 2021)	R	MAM	✓	✓	✓	✗	$3kd^2 + kd$
Transformer w/ sin. (Vaswani et al. 2017)							
Liu et al. (2019)	A	APE	✗	✗	✓	✓	0
Takase and Okazaki (2019)							
Oka et al. (2020)							
Universal Transf. (Dehghani et al. 2019)	A	APE	✗	✓	✓	✓	0
DSAN (Shen et al. 2018)	R	MAM	✗	✓	✓	✓	0
RoBERTa (Liu et al. 2019)							
SPI-abs (Yang et al. 2019)	A	APE	✗	✗	✓	✓	0
SPI-rel (Wang et al. 2019)	R	MAM	✓	✓	✗	✗	$2(2t_{max} + 1)d$
TFE (Sho et Quirk 2019)	A	APE	✓	✓	✗	✗	$\frac{2d}{2t_{max}}$
Graph Transformer (Zhu et al. 2019)	R	MAM	✓	✓	✓	✓	$5d^2 + (d + 1)t_1$
Graph Transformer (Cai and Lam 2020)							$7d^2 + 3d$
Graphformer (Schmitz et al. 2021)	R	MAM	✓	✓	✓	✗	$2(2t_{max} + 1)d$
GRAPH-BERT (Zhang et al. 2020)	A	APE	✗	✗	✓	✓	0
	B						



Dufner, P., Schmitt, M., & Schütze, H. (2022). Position information in transformers: An overview. *Computational Linguistics*, 48(3), 733-763.

74

**Advanced Reading:**

... and beyond Transformer

(a) DenseSSM in autoregressive mode.

He, Wei, Kai Han, Yehui Tang, Chengcheng Wang, Yujie Yang, Tianyu Guo, and Yunhe Wang. "Densemamba: State space models with dense hidden connection for efficient large language models." *arXiv preprint arXiv:2403.00818* (2024).

UNSW SYDNEY

75

**Surge in NLP research (NOT all from Big Tech)...  
...and a 9-week course**

This ‘introductory’ course is a distilled version of NLP.

Why do you want to study NLP?

or

(Reminder: Lecture 1: “Course Philosophy”; “Why are we even learning pre-deep learning NLP?”)

Try: <https://saifmohammad.com/WebPages/nlp scholar.html>

Year	Number
2010	2315
2011	2134
2012	2528
2013	3169
2014	3776
2015	4577
2016	5024
2017	4669
2018	5621
2019	8755
2020	8835
2021	9407
2022	7381

[https://www.researchgate.net/figure/The-total-number-of-publications-in-the-field-of-NLP-The-statistics-are-from-the-Web-of\\_fig1\\_372092521](https://www.researchgate.net/figure/The-total-number-of-publications-in-the-field-of-NLP-The-statistics-are-from-the-Web-of_fig1_372092521)

UNSW SYDNEY

76

## Summary

Part	Key Idea	Demos
Why Attention?	Intuition behind attention	Pytorch tutorial
Attention in recurrent language modeling	Additive attention in recurrent LM and limitations	Attention Visualization
Encoders & decoders	Scaled dot-product attention -> Multi-head attention -> Self-attention, masked attention -> Encoders and decoders. Pseudocode + Animation.	Transformer tutorial
Tokenization and positional encoding	BPE algorithm, positional encoding, language model head	Transformer tutorial continued. BPE Tokenizer. Positional encoding Visualization
Transformer: Looking Back	Pseudocode; components; impact of Transformer.	



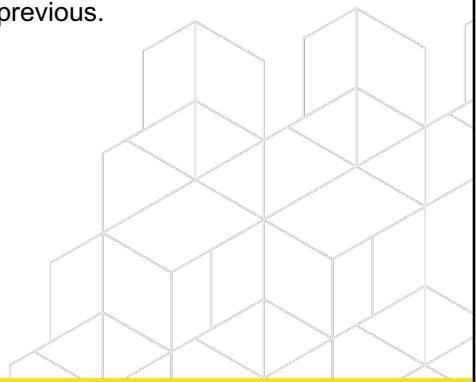
77

**Transformer is the basis of  
the state-of-the-art models  
in NLP**

New models continue to be released – one ‘larger’ than the previous.

Transformer spawned several derivative models...

‘Large language models’  
(Next week!)



78

## Suggested Reading

**All pseudocodes from:** Phuong et al., Formal Algorithms for Transformers, <https://arxiv.org/abs/2207.09238>.

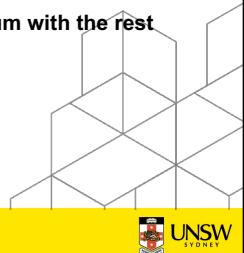
Annotated Transformer: <https://jalammar.github.io/illustrated-transformer/>

Byte-pair encoding: <https://huggingface.co/learn/nlp-course/en/chapter6/5>

Video: <https://www.youtube.com/watch?v=zduSFxRajkE>

Visualisation of Transformer (in Keras):  
<https://colab.research.google.com/github/tensorflow/text/blob/master/docs/tutorials/transformer.ipynb>

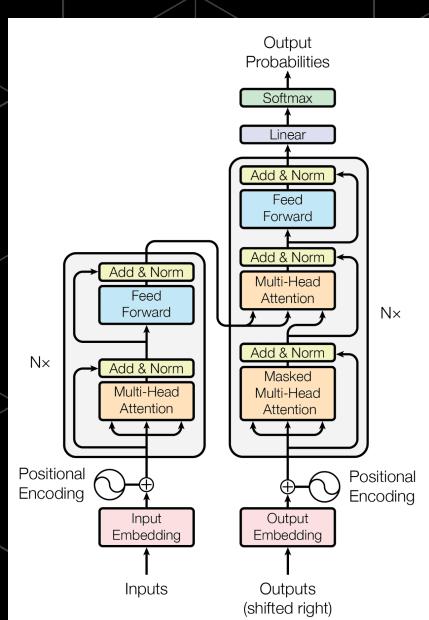
I hope the above links lead you into a rabbit-hole of interesting material. Do share it on the forum with the rest of your class.



79

So what kind of models did  
Transformer lead to?

→ Large Language Models  
(Next week!)



80