

COMP 4446 / 5046

Lecture 7: Models – Transformer

Jonathan K. Kummerfeld

Semester 1, 2025



THE UNIVERSITY OF
SYDNEY

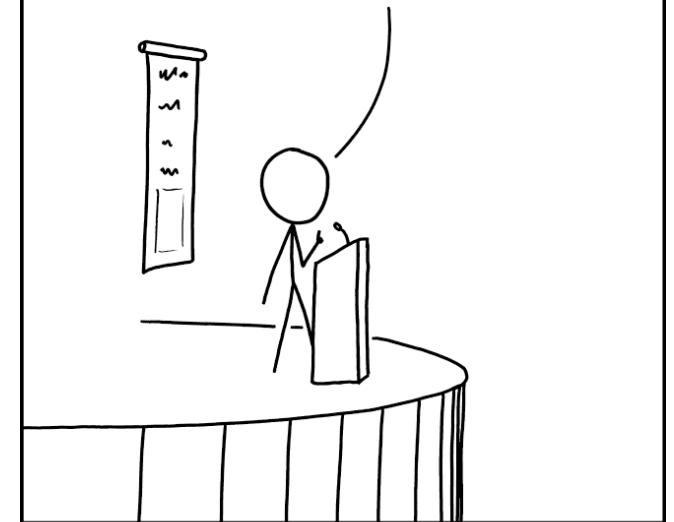
[menti.com 2520 1730](https://menti.com/25201730)

[The mainstream dogma sparked a wave of dogmatic revisionism, and this revisionist mainstream dogmatism has now given way to a more rematic mainvisionist dogstream.]

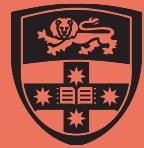
Source: <https://xkcd.com/2857/>

Rebuttals

IT'S BECOME CONVENTIONAL WISDOM THAT THE BACKLASH AGAINST THE PREVAILING CONSENSUS LED RESEARCHERS TO IGNORE INCONVENIENT NEW EVIDENCE. HOWEVER...



IN A FIELD THAT'S BEEN AROUND FOR A WHILE, IT CAN BE HARD TO FIGURE OUT HOW MANY LEVELS OF REBUTTAL DEEP YOU ARE.

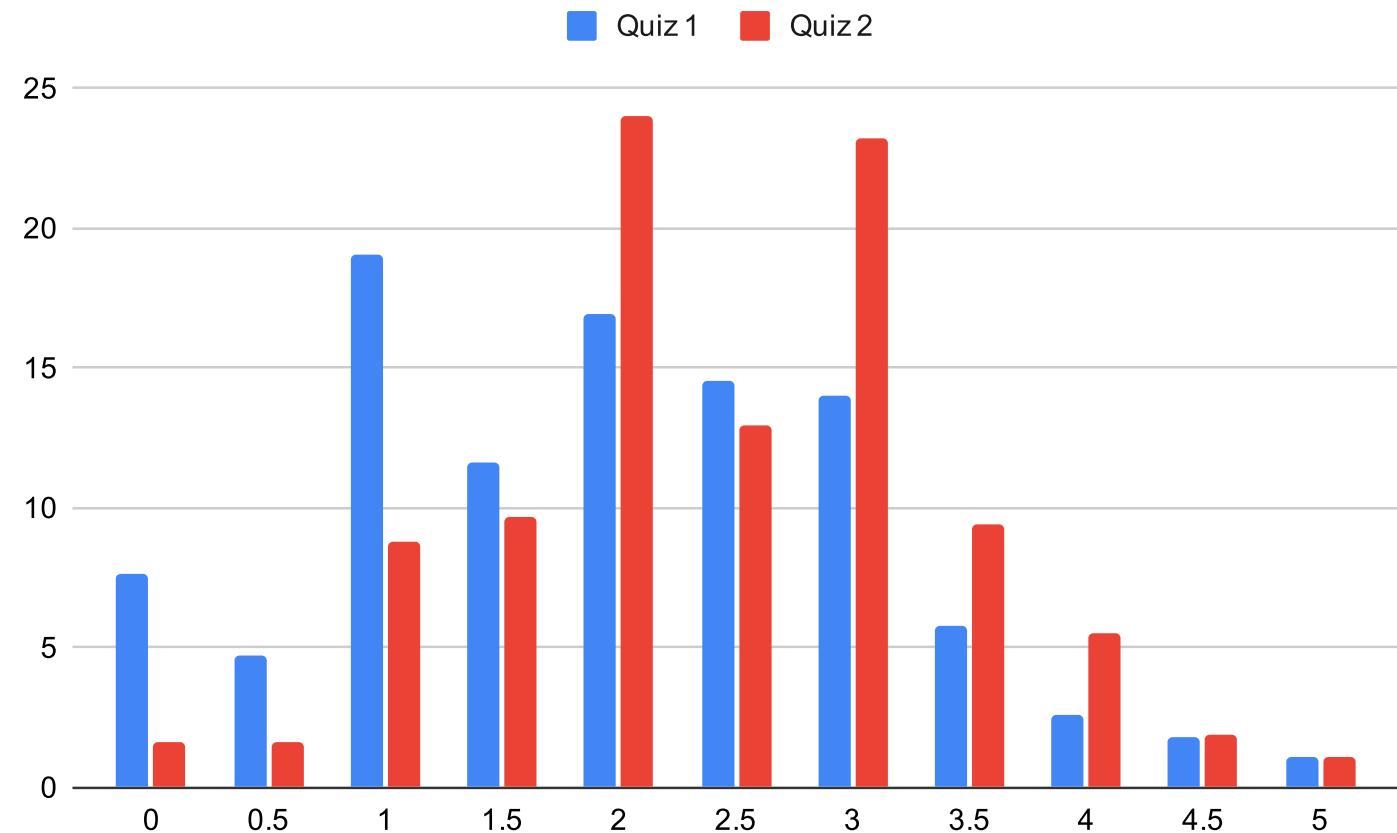


Self-Attention
Transformer
Training and Use
Lab Preview



menti.com 2520 1730

Quiz results - an improvement! But still 22% below 2



Note: results are preliminary and subject to change. I will be giving credit for the best result between Quiz 1 and Quiz 2.

To avoid the bugs from this quiz, I will be getting TAs to check them in future (just as exams are reviewed by another academic)

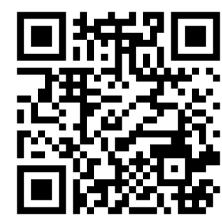


Self-Attention

Transformer

Training and Use

Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Special Consideration Reminder

Short release assignments - No simple extensions

Regular Special Consideration is not processed fast enough

Slip days instead!

For major issues (e.g., broke your arm and can't type for a month), go through Special Consideration process and a more significant change will be determined



COMP 4446 / 5046
Lecture 7, 2025

Self-Attention

Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Office Hours shift

This week Office Hours will be 2-3pm on Tuesday



COMP 4446 / 5046
Lecture 7, 2025

Self-Attention

Transformer

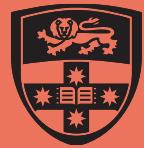
Training and Use

Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Self-Attention



Self-Attention

Transformer

Training and Use

Lab Preview

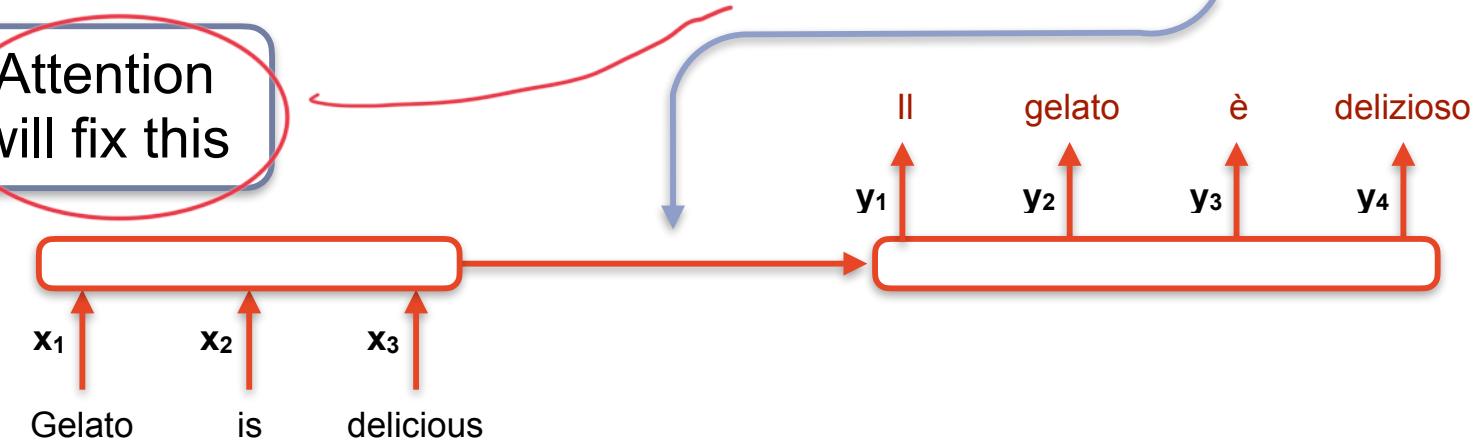


menti.com 2520 1730

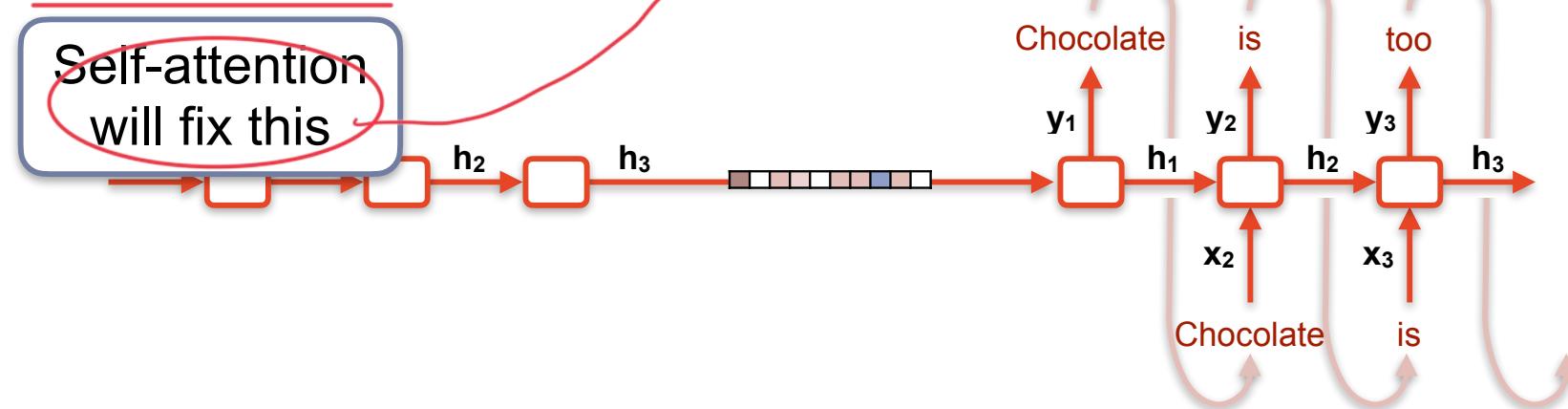
Reminder: Issues with the encoder-decoder

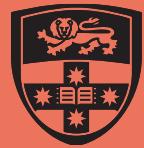
Review
Data loss

Bottleneck - all information has to go through one vector



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





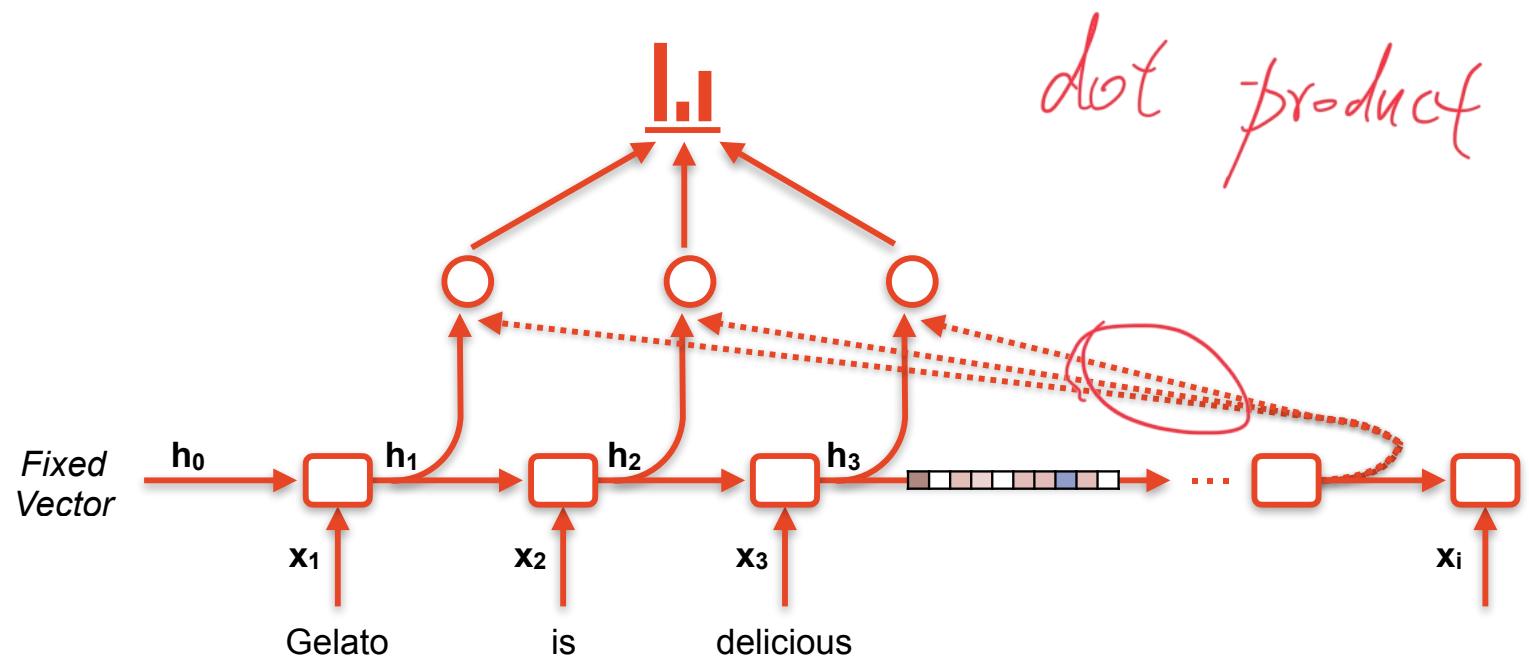
Self-Attention Transformer Training and Use Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

We will solve this with ‘attention’

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





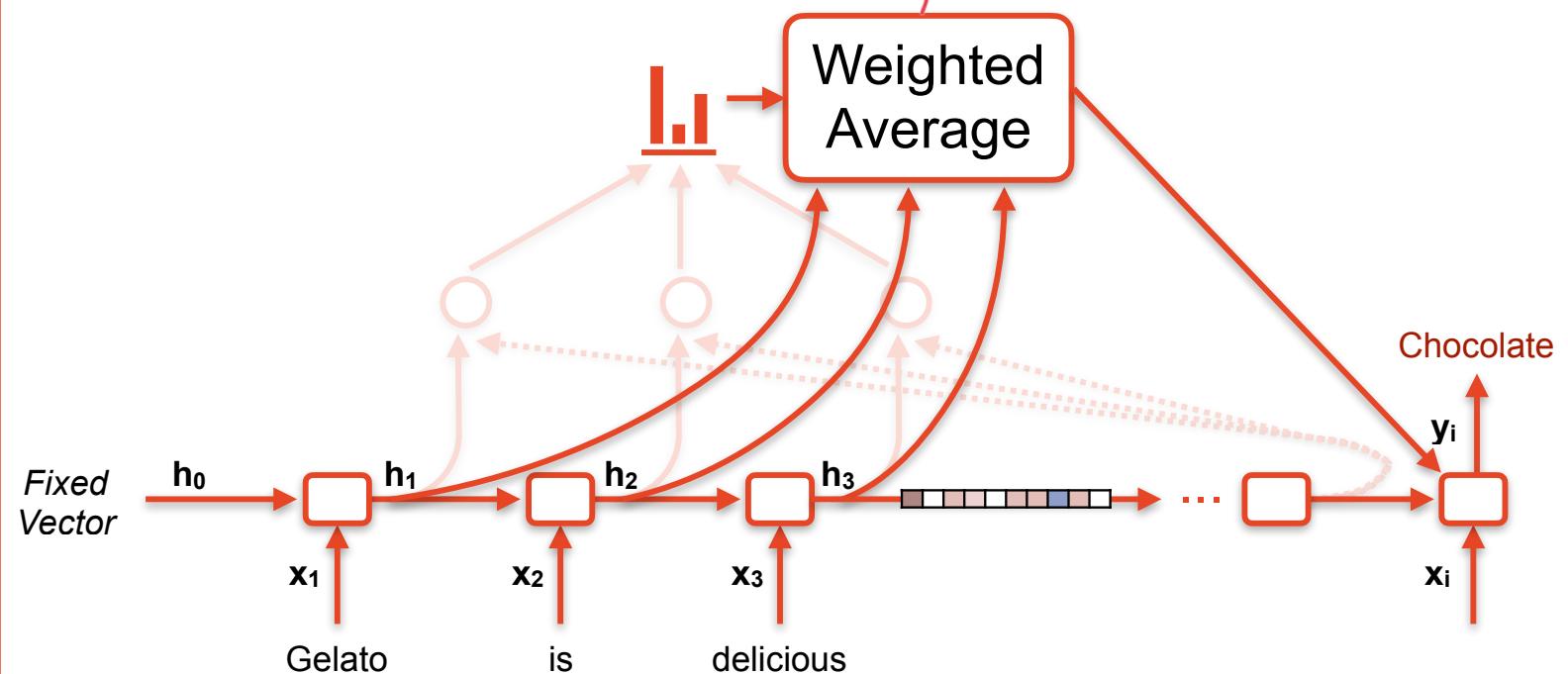
Self-Attention Transformer Training and Use Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

We will solve this with ‘attention’

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





Self-Attention

Transformer

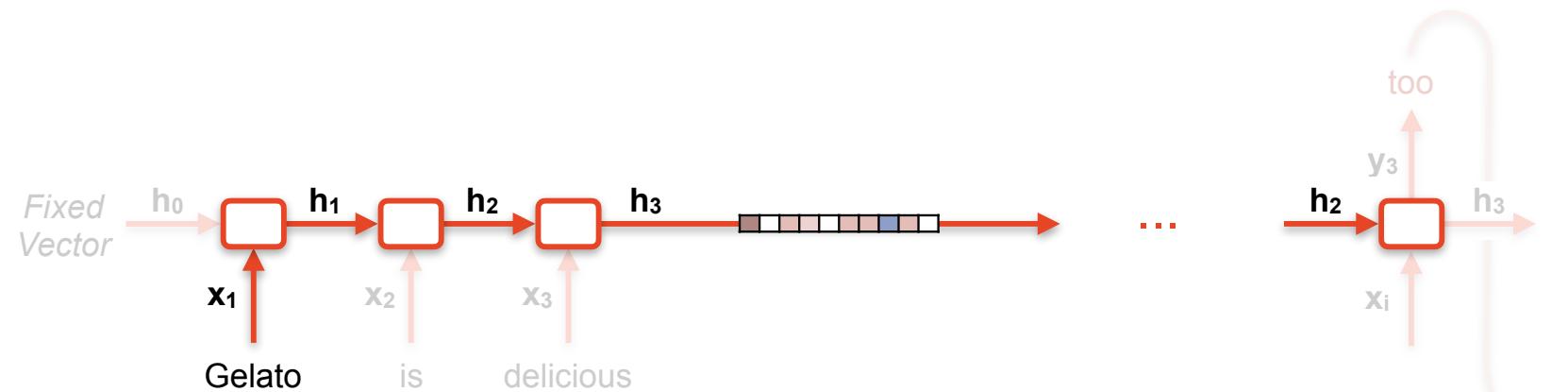
Training and Use

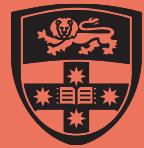
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

This resolves the bottleneck problem without introducing many extra parameters



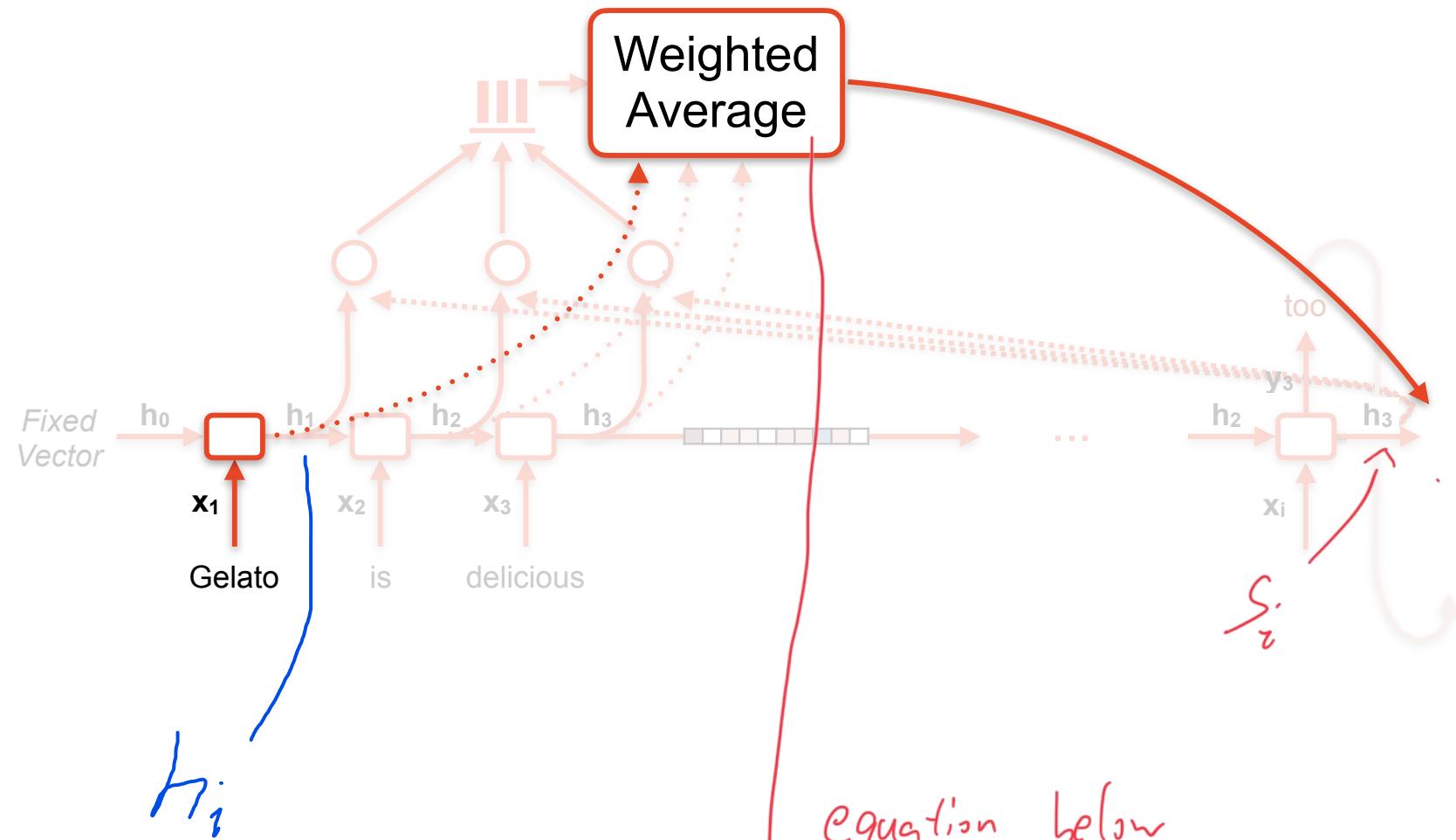


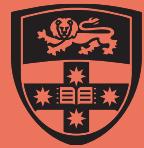
Self-Attention Transformer Training and Use Lab Preview



menti.com 2520 1730

This resolves the bottleneck problem without introducing many extra parameters





Self-Attention Transformer Training and Use Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

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

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

Normalise s

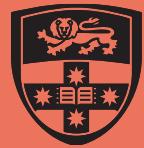
This is dot product attention

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

$$\alpha_i := \text{softmax}(\mathbf{e}_i) \in \mathbb{R}^N$$

Calculate weighted average

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

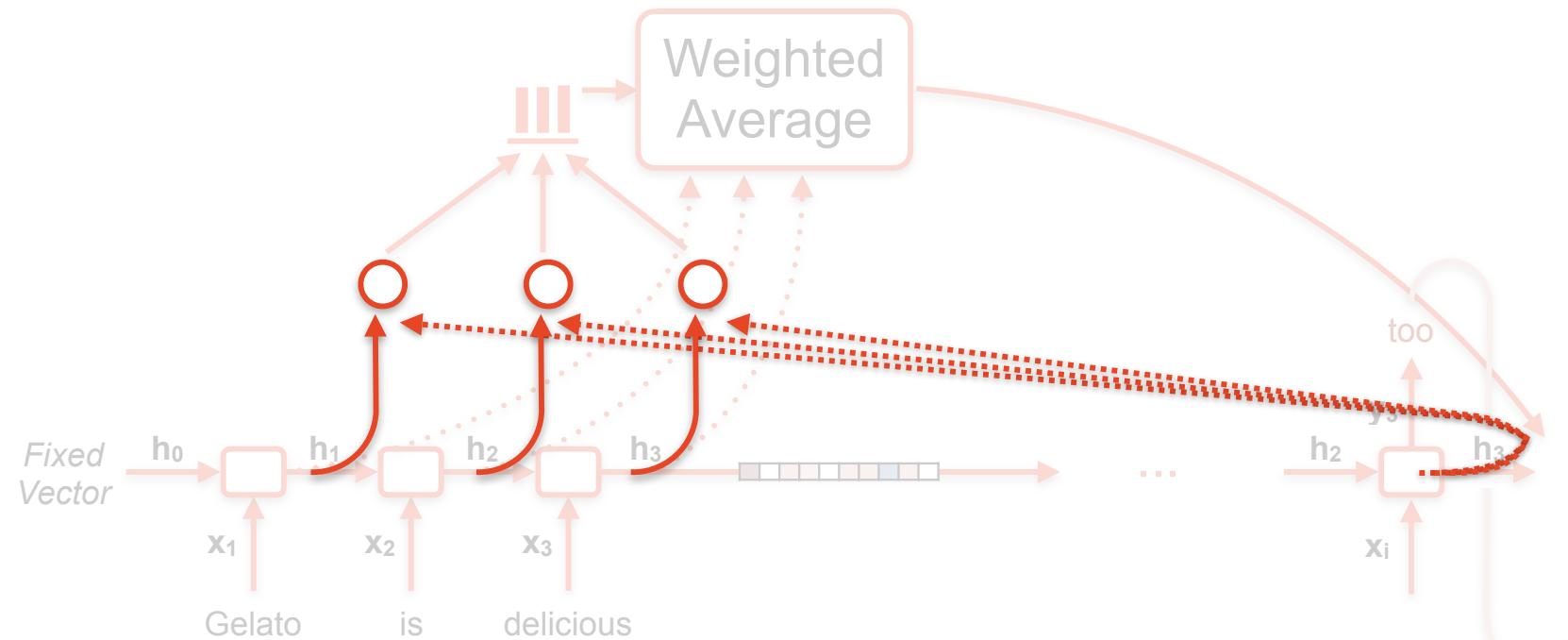


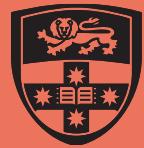
Self-Attention Transformer Training and Use Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

This resolves the bottleneck problem without introducing many extra parameters





Self-Attention

Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

There are many forms of attention

Dot product attention

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

s, h must have same size

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

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

Powerful /
Flexible



Self-Attention

Transformer

Training and Use

Lab Preview



menti.com 2520 1730

Another way of thinking about attention - looking up values

Consider a lookup table:

One Example

Input

Output

The diagram illustrates a search operation on a hash table. A 'Query' is shown on the left, with four red arrows pointing to the keys K0, K1, K2, and K3 in the 'Keys' column. The 'Values' column contains V0, V1, V2, and V3 respectively. A red box highlights the row for K1. A red arrow points from the 'the most similar value' label to V1. Another red arrow points from the 'Return' label to V1.

Keys	Values
K0	V0
K1	V1
K2	V2
K3	V3

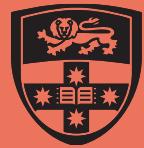
Compare

Query

the most similar value

V1

Return



Self-Attention

Transformer

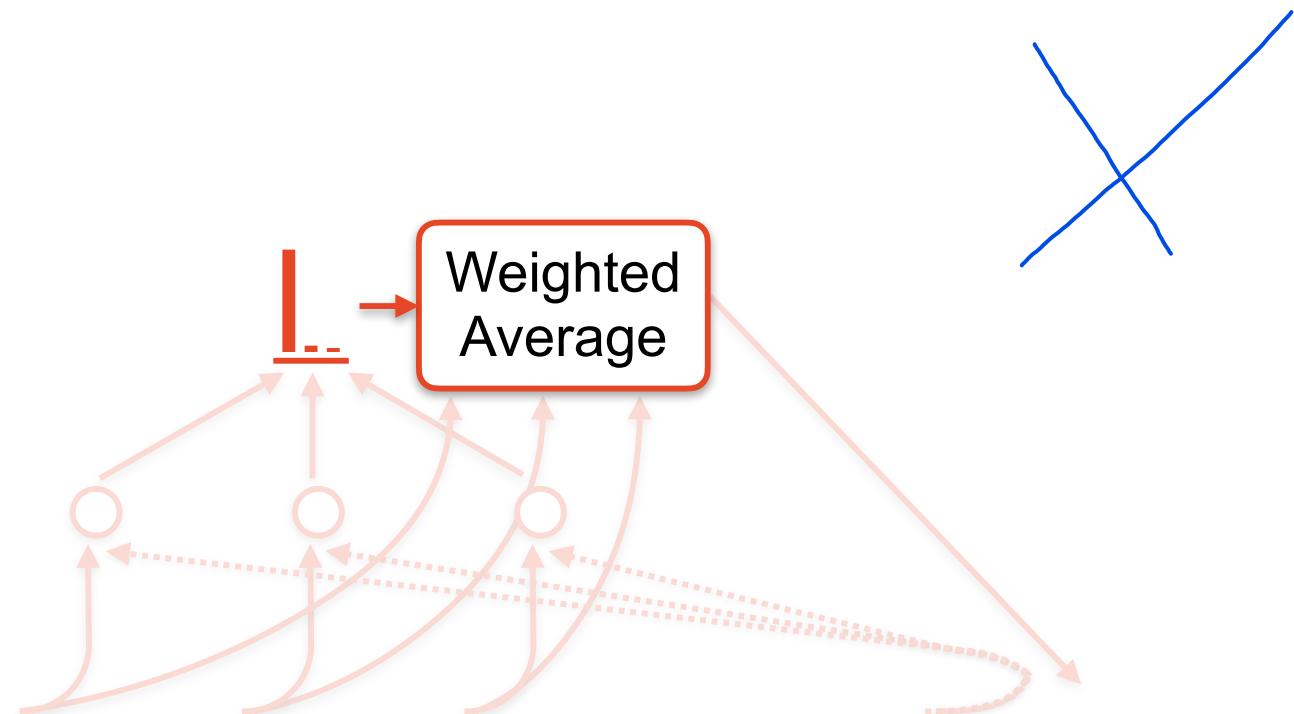
Training and Use

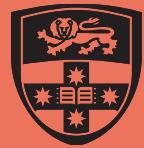
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Another way of thinking about attention - looking up values





Self-Attention

Transformer

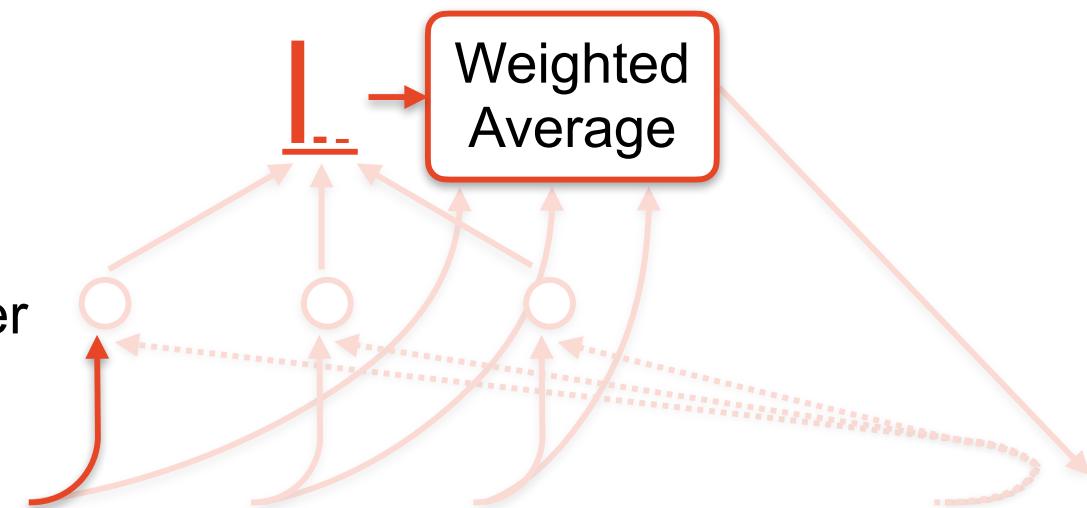
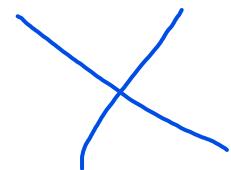
Training and Use

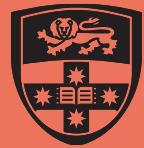
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Another way of thinking about attention - looking up values





Self-Attention

Transformer

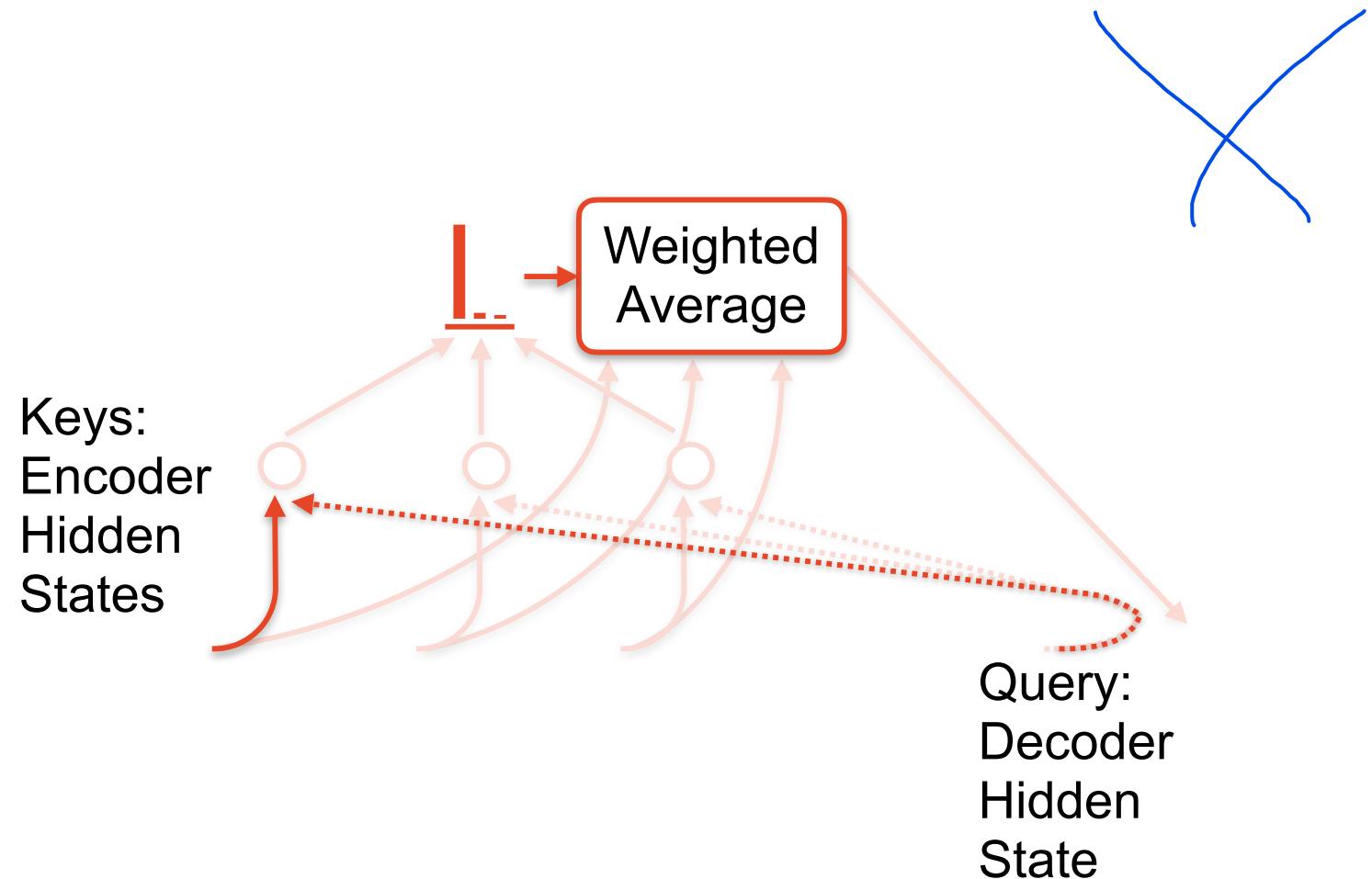
Training and Use

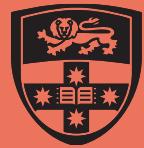
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Another way of thinking about attention - looking up values





Self-Attention

Transformer

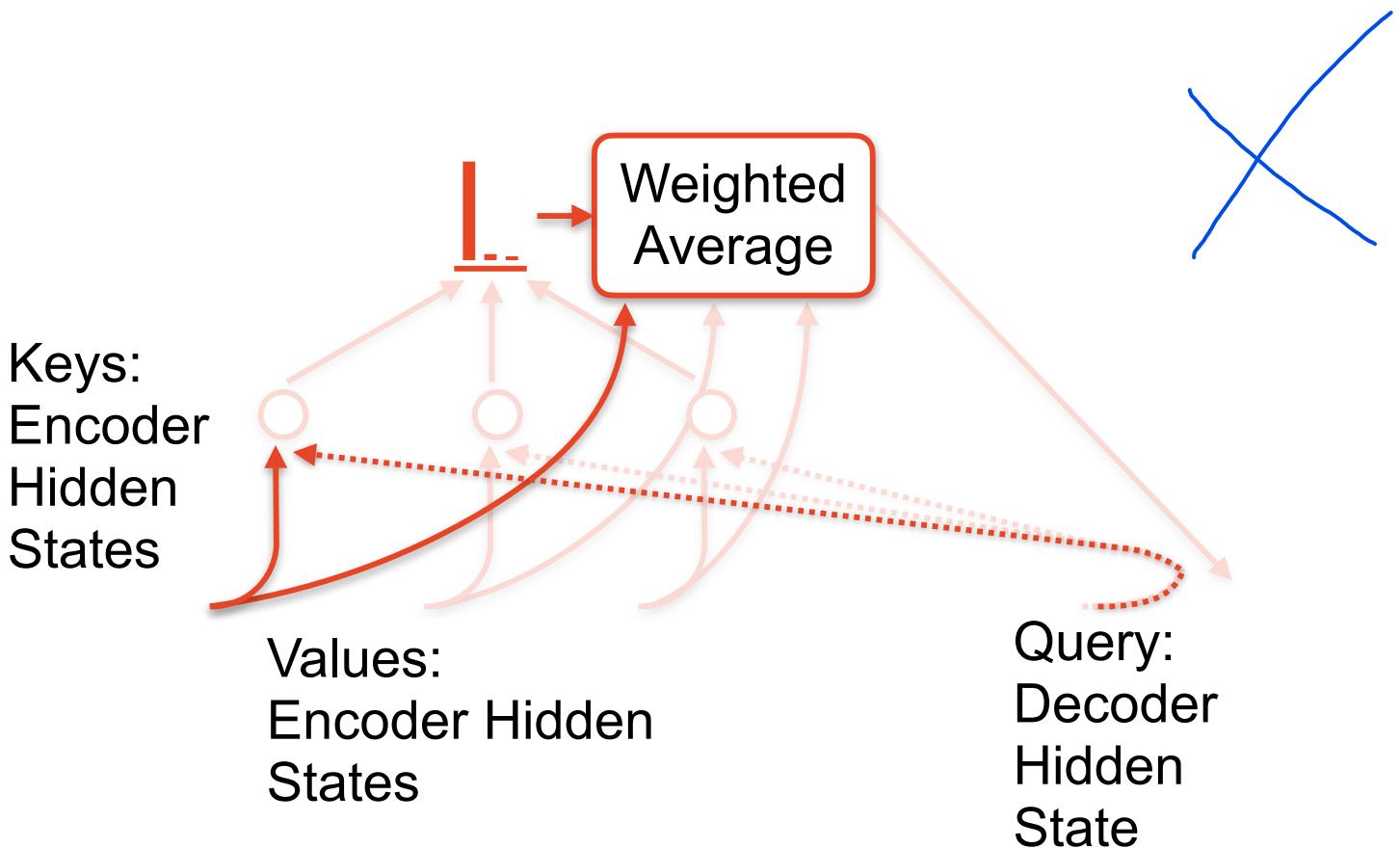
Training and Use

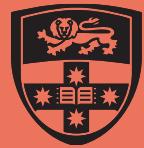
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Another way of thinking about attention - looking up values





Self-Attention

Transformer

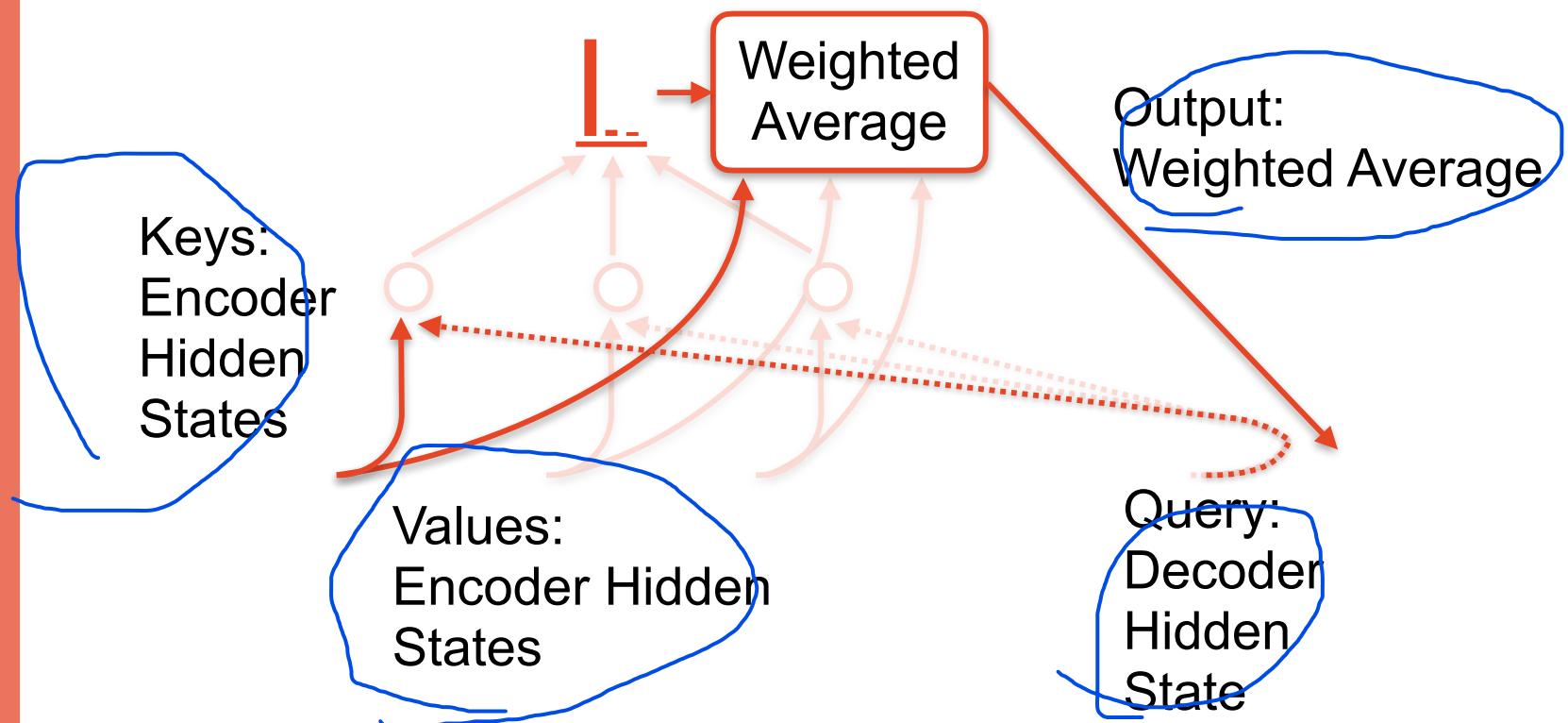
Training and Use

Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Another way of thinking about attention - looking up values





Self-Attention

Transformer

Training and Use

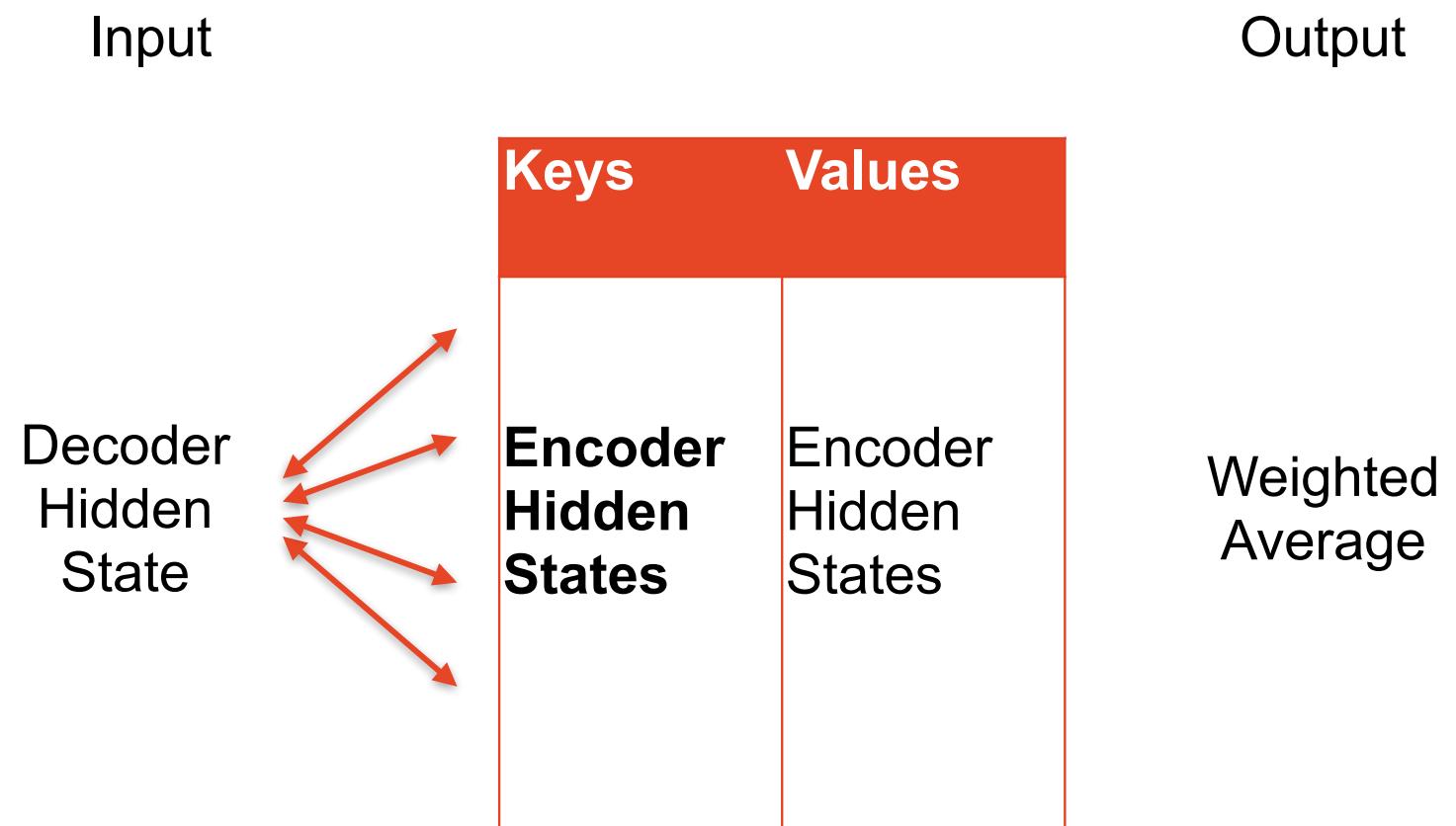
Lab Preview

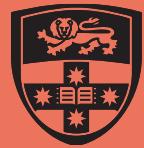


[menti.com 2520 1730](https://menti.com/25201730)

Another way of thinking about attention - looking up values

Consider a lookup table:





Do we need RNNs?

In lecture 3:

"How do we handle variable length input?"

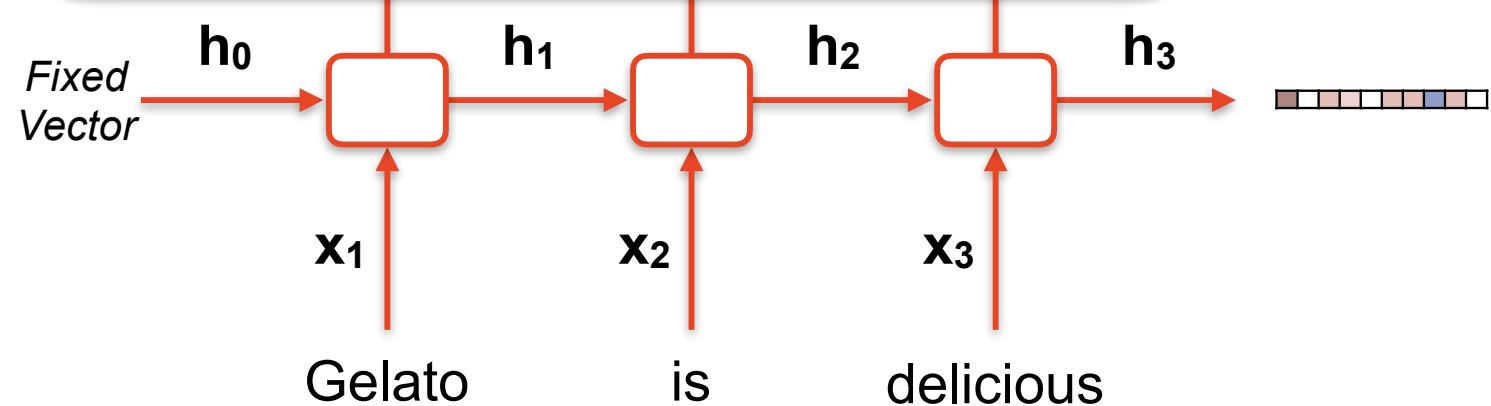
"Recurrent Neural Networks (RNNs) address this issue with **context-dependent representations**"

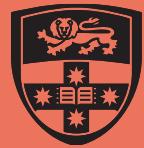
hidden layer

In lecture 6, attention:

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

average weight





Self-Attention

Transformer

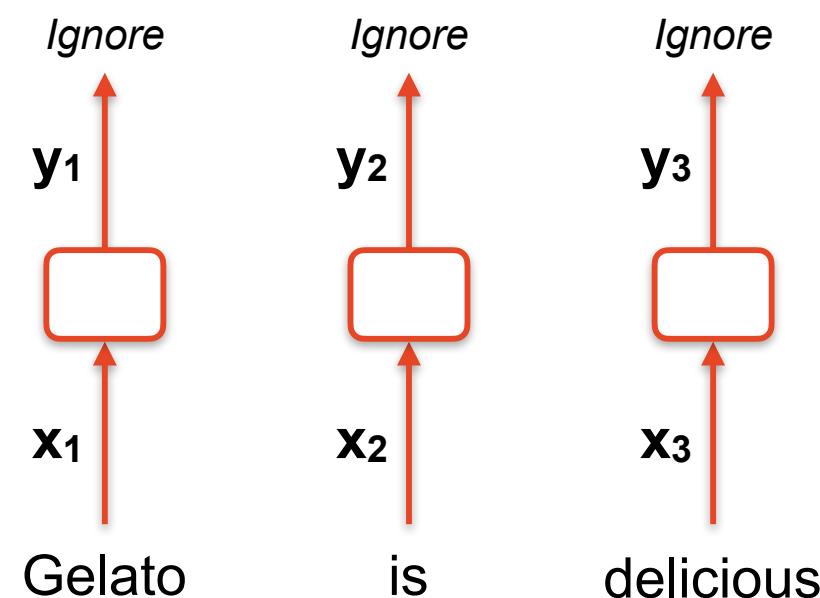
Training and Use

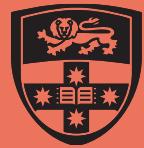
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Cut the recurrent edges and use attention instead





Self-Attention

Transformer

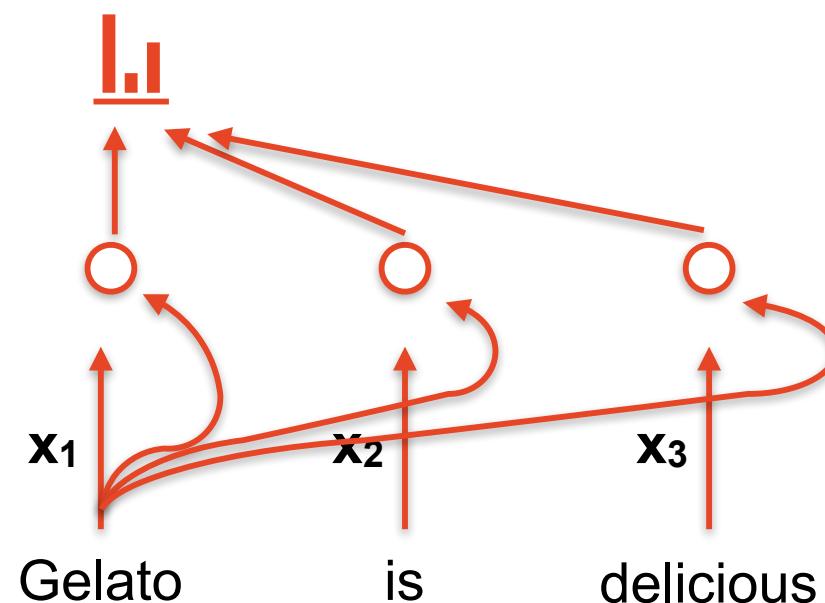
Training and Use

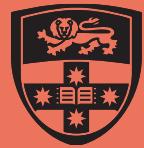
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Cut the recurrent edges and use attention instead





Self-Attention Transformer Training and Use Lab Preview

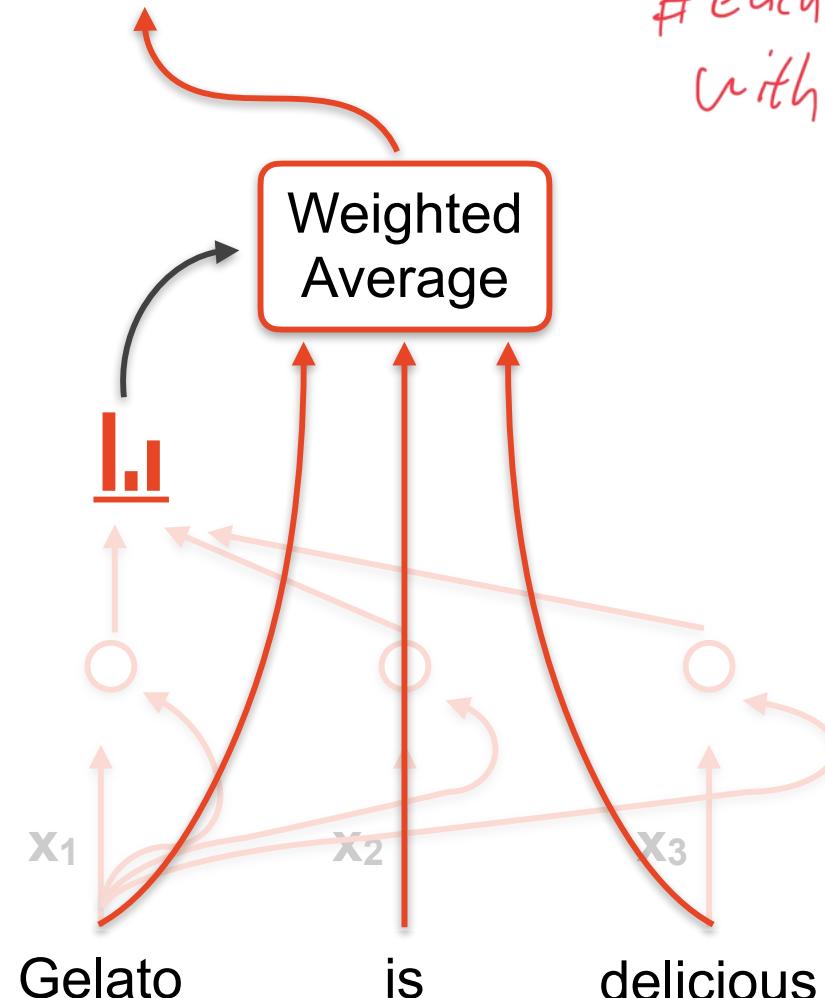


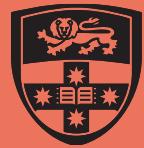
[menti.com 2520 1730](https://menti.com/25201730)

Cut the recurrent edges and use attention instead

'Gelato' contextual vector

each word compare
with 'Gelato'





Self-Attention

Transformer

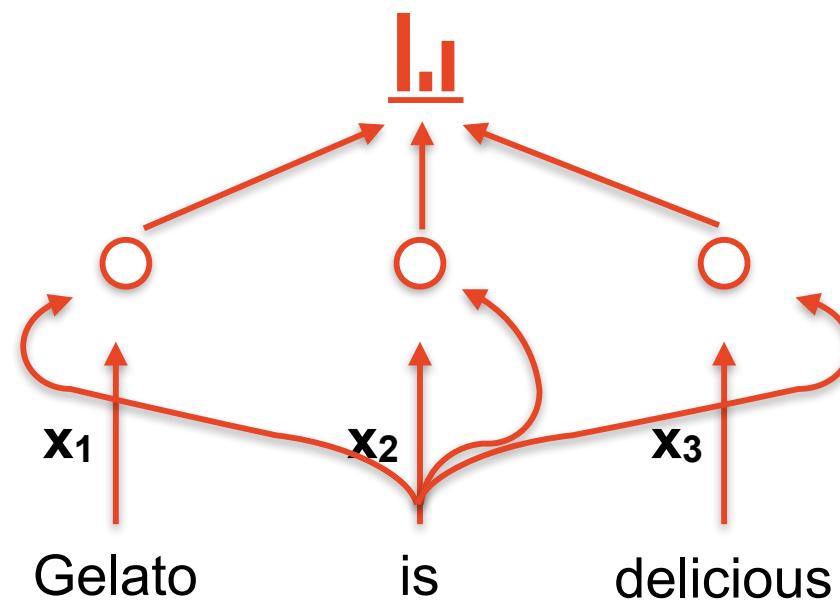
Training and Use

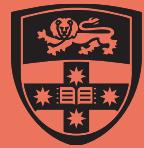
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Cut the recurrent edges and use attention instead





Self-Attention

Transformer

Training and Use

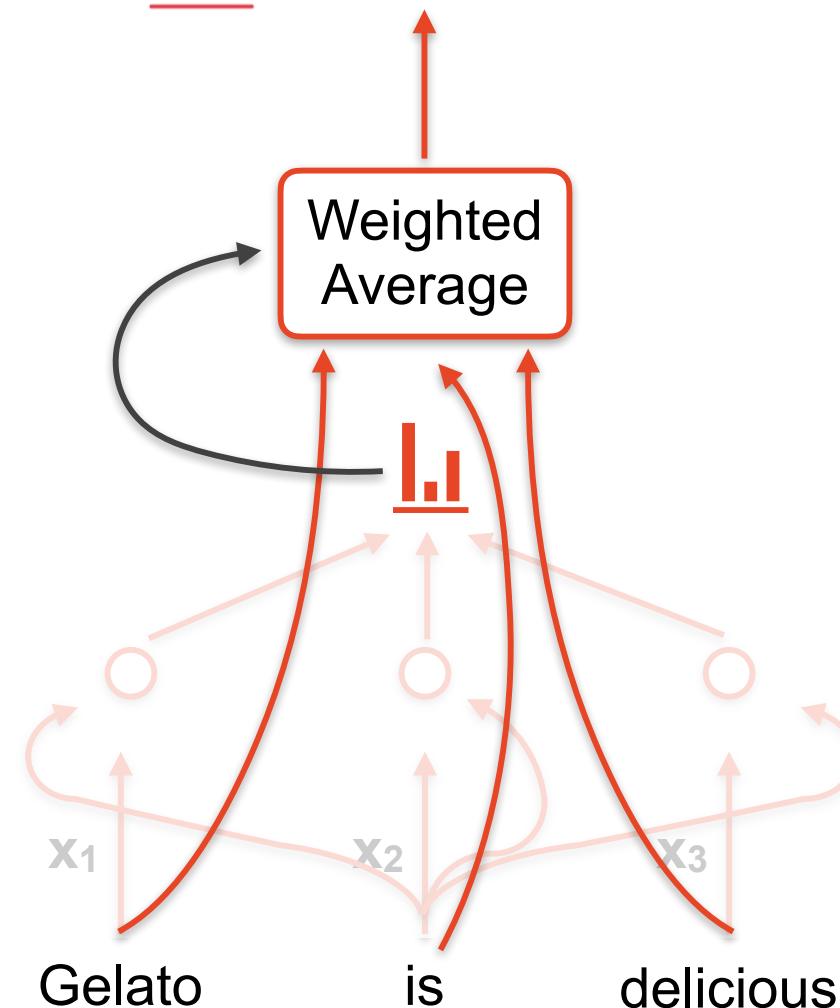
Lab Preview

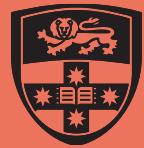


[menti.com 2520 1730](https://menti.com/25201730)

Cut the recurrent edges and use attention instead

'is' contextual vector





Self-Attention

Transformer

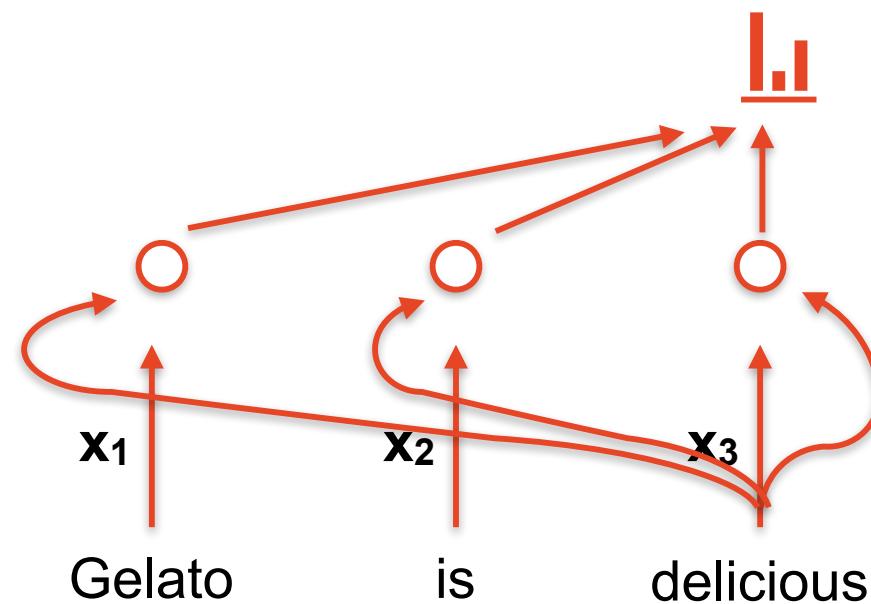
Training and Use

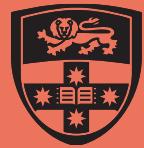
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Cut the recurrent edges and use attention instead





Self-Attention

Transformer

Training and Use

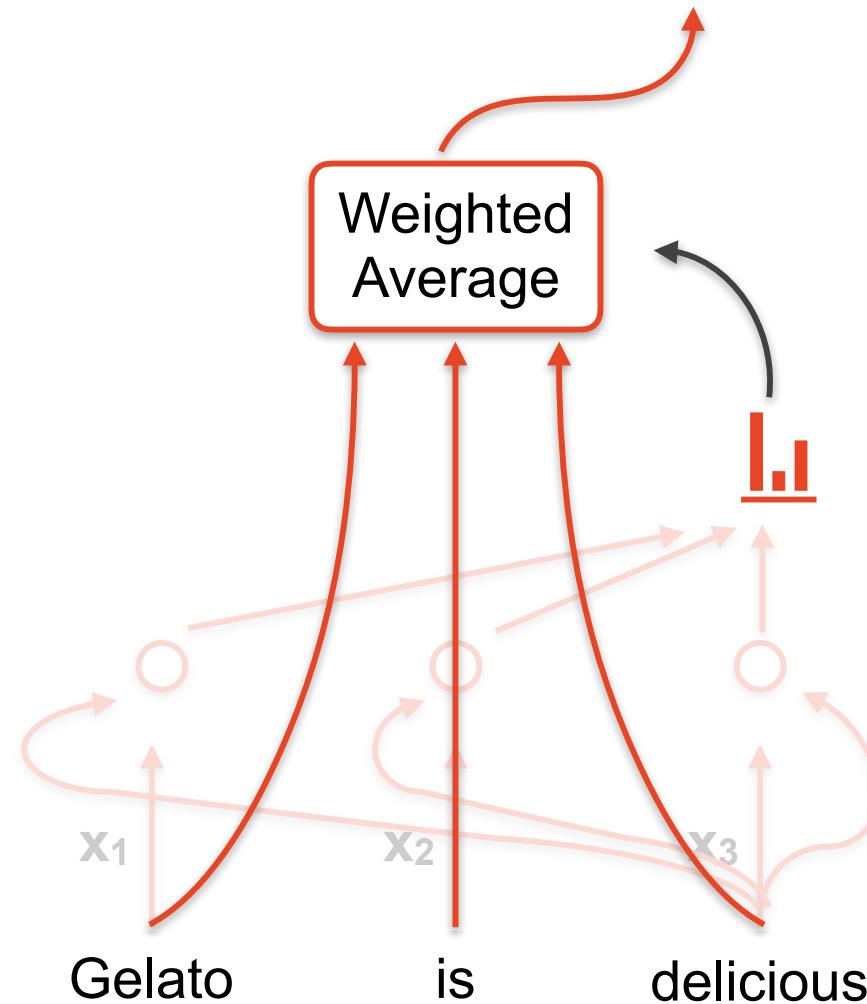
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Cut the recurrent edges and use attention instead

'delicious' contextual vector



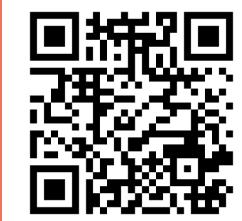


Self-Attention

Transformer

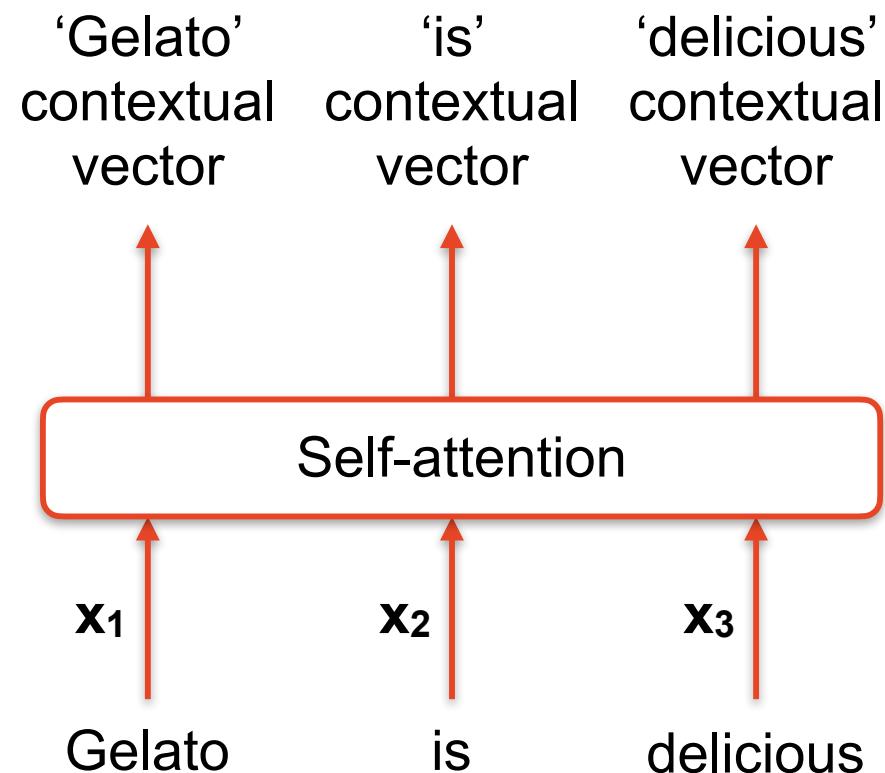
Training and Use

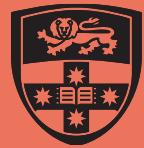
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Self-attention gives contextual representations of input



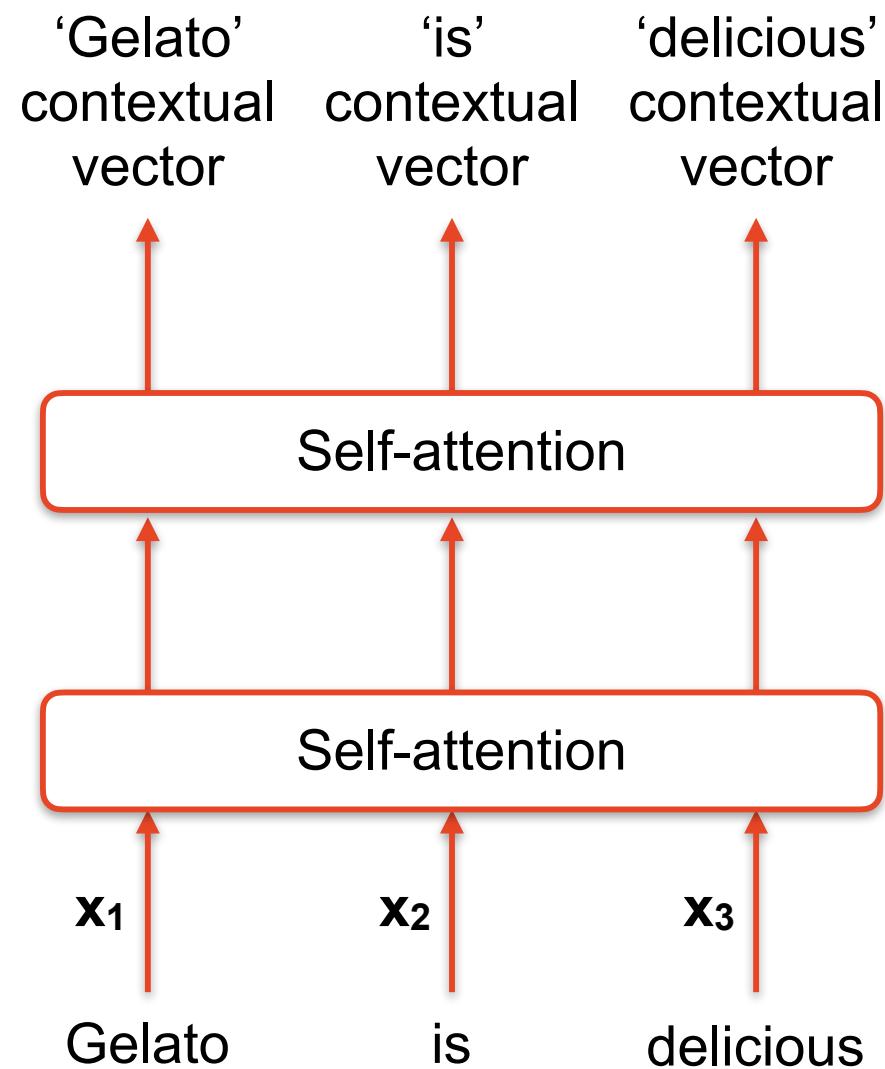


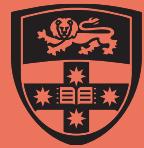
Self-Attention Transformer Training and Use Lab Preview



menti.com 2520 1730

Self-attention gives contextual representations of input





Self-Attention

Transformer

Training and Use

Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Using equations, here is what we have done

Self attention

Initial word vectors

$$\mathbf{x}_i \in \mathbb{R}^d$$

current word

Do dot product to get
attention scores

$$e_{ij} = \mathbf{x}_i^\top \mathbf{x}_j$$

\mathbf{x}_i is the query

Softmax to get a
distribution

$$\alpha_{ij} = \text{softmax}(e_{ij})$$

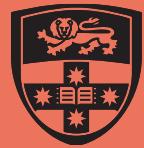
\mathbf{x}_j is the key

Weighted average

$$\mathbf{o}_i = \sum_j \alpha_{ij} \mathbf{x}_j$$

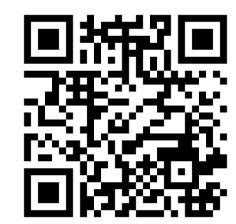
\mathbf{x}_j is the value

all words
include self



Self-Attention

Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Using equations, here is what we have done

Initial word vectors

$$\mathbf{x}_i \in \mathbb{R}^d$$

Transform into query,
keys, and values

$$\mathbf{q}_i = Q\mathbf{x}_i$$

$$\mathbf{k}_i = K\mathbf{x}_i$$

$$\mathbf{v}_i = V\mathbf{x}_i$$

**Q, K, and V
are learned
matrices**

Do dot product to get
attention scores

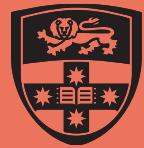
$$e_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

Softmax to get a
distribution

$$\alpha_{ij} = \text{softmax}(e_{ij})$$

Weighted average

$$\mathbf{o}_i = \sum_j \alpha_{ij} \mathbf{v}_j$$



(above is linear)

To recover all benefits of the RNN we need nonlinearities and position representation

Self-Attention

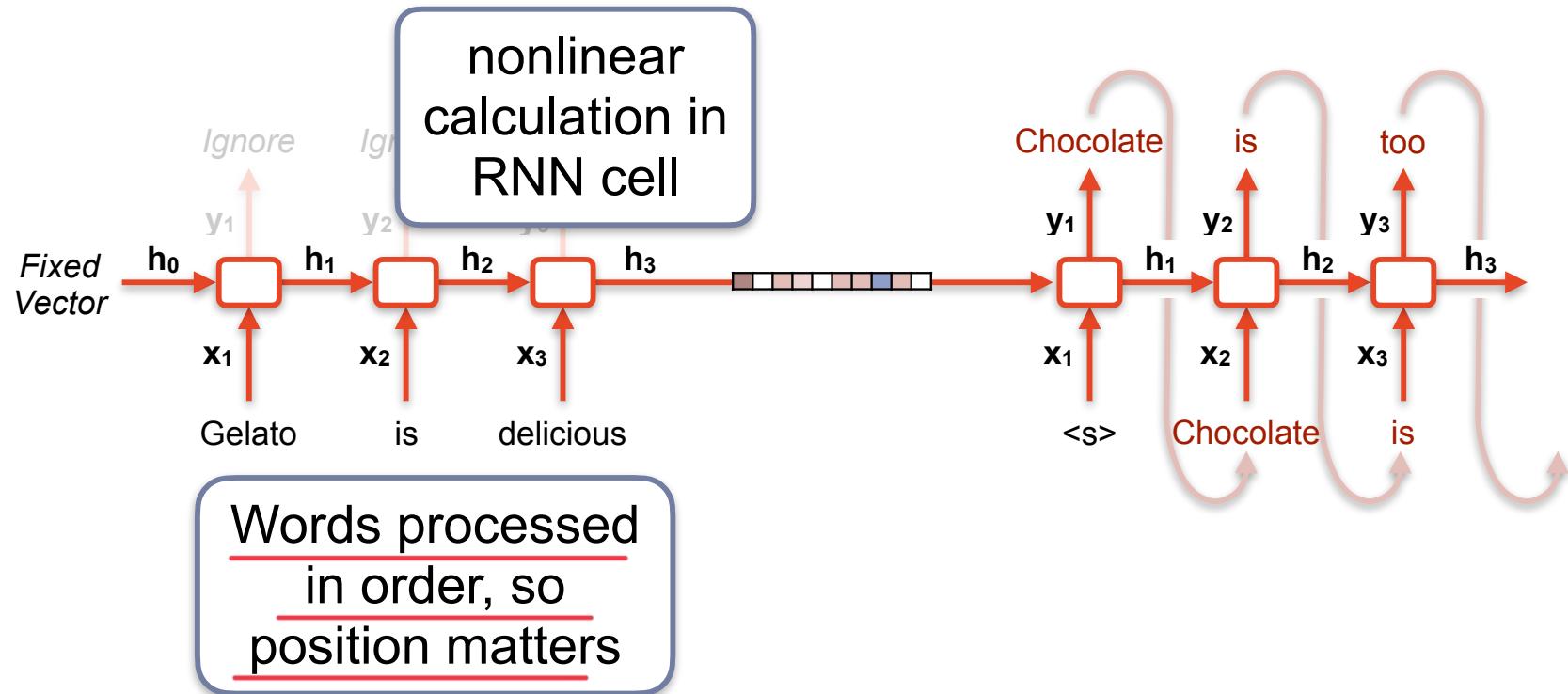
Transformer

Training and Use

Lab Preview



menti.com 2520 1730





Self-Attention

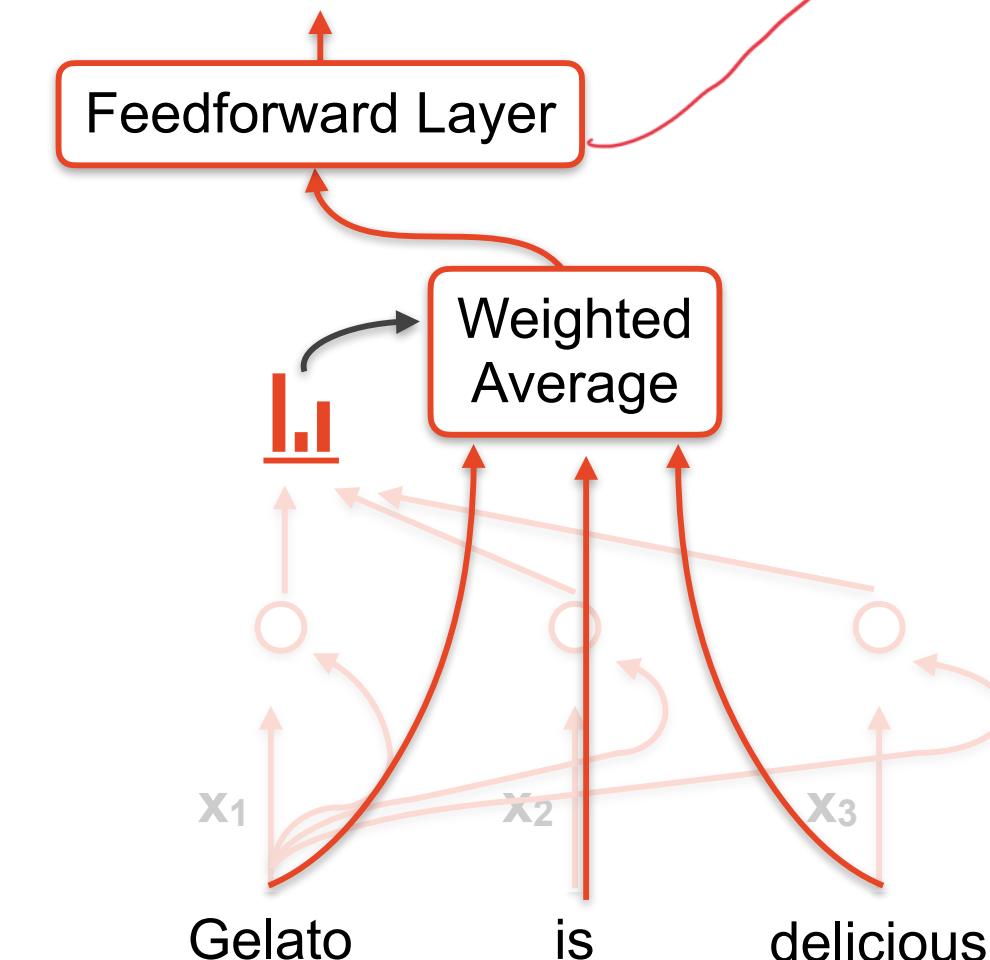
Transformer
Training and Use
Lab Preview

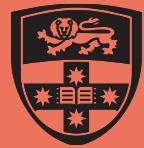


[menti.com 2520 1730](https://menti.com/25201730)

To recover all benefits of the RNN we need nonlinearities and position representation

'Gelato' contextual vector





Self-Attention

Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

To recover all benefits of the RNN we need **nonlinearities** and position representation

Initial word vectors

$$\mathbf{x}_i \in \mathbb{R}^d$$

Transform into query,
keys, and values

$$\mathbf{q}_i = Q\mathbf{x}_i$$

$$\mathbf{k}_i = K\mathbf{x}_i$$

$$\mathbf{v}_i = V\mathbf{x}_i$$

Do dot product to get
attention scores

$$e_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

Softmax to get a
distribution

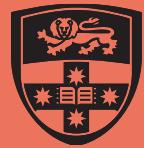
$$\alpha_{ij} = \text{softmax}(e_{ij})$$

Weighted average

$$\mathbf{t}_i = \sum_j \alpha_{ij} \mathbf{v}_j$$

Feedforward layer

$$\mathbf{o}_i = W_2 \text{ReLU}(W_1 \mathbf{t}_i + \mathbf{b}_1) + \mathbf{b}_2$$



Self-Attention

Transformer
Training and Use
Lab Preview



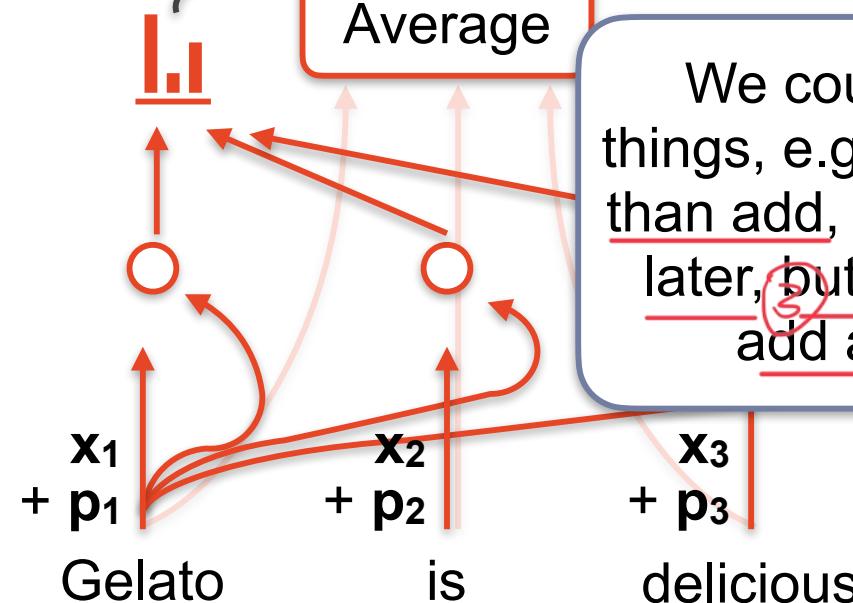
[menti.com 2520 1730](https://menti.com/25201730)

To recover all benefits of the RNN we need nonlinearities and **position representation**

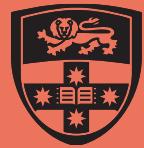
'Gelato' contextual vector

Feedforward Layer

Weighted Average



We could do many other things, e.g. ① concatenate rather than add, or add them at every later, but people tend to just ③ add at the first layer.



Self-Attention

Transformer
Training and Use
Lab Preview



menti.com 2520 1730

To recover all benefits of the RNN we need nonlinearities and **position representation**

Initial word vectors

$$\mathbf{x}_i \in \mathbb{R}^d$$

$$\mathbf{q}_i = Q\tilde{\mathbf{x}}_i$$

$$\mathbf{k}_i = K\tilde{\mathbf{x}}_i$$

$$\mathbf{v}_i = V\tilde{\mathbf{x}}_i$$

Transform into query, keys, and values

$$e_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

Do dot product to get attention scores

$$\alpha_{ij} = \text{softmax}(e_{ij})$$

Softmax to get a distribution

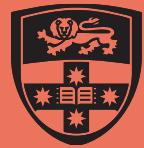
$$\mathbf{t}_i = \sum_j \alpha_{ij} \mathbf{v}_j$$

Weighted average

$$\mathbf{o}_i = W_2 \text{ReLU}(W_1 \mathbf{t}_i + \mathbf{b}_1) + \mathbf{b}_2$$

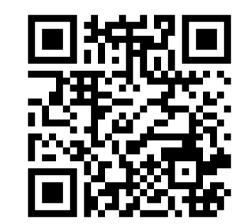
Feedforward layer

Original without position



Self-Attention

Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

To recover all benefits of the RNN we need nonlinearities and **position representation**

Initial word vectors with positions

Transform into query, keys, and values

Do dot product to get attention scores

Softmax to get a distribution

Weighted average

Feedforward layer

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \mathbf{p}_i$$

$$\mathbf{q}_i = Q\tilde{\mathbf{x}}_i$$

$$\mathbf{k}_i = K\tilde{\mathbf{x}}_i$$

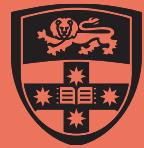
$$\mathbf{v}_i = V\tilde{\mathbf{x}}_i$$

$$e_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

$$\alpha_{ij} = \text{softmax}(e_{ij})$$

$$\mathbf{t}_i = \sum_j \alpha_{ij} \mathbf{v}_j$$

$$\mathbf{o}_i = W_2 \text{ReLU}(W_1 \mathbf{t}_i + \mathbf{b}_1) + \mathbf{b}_2$$

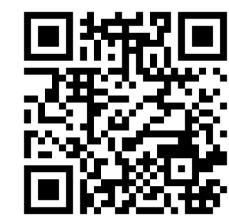


Self-Attention

Transformer

Training and Use

Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

To recover all benefits of the RNN we need nonlinearities and **position representation**

What is p_i ?

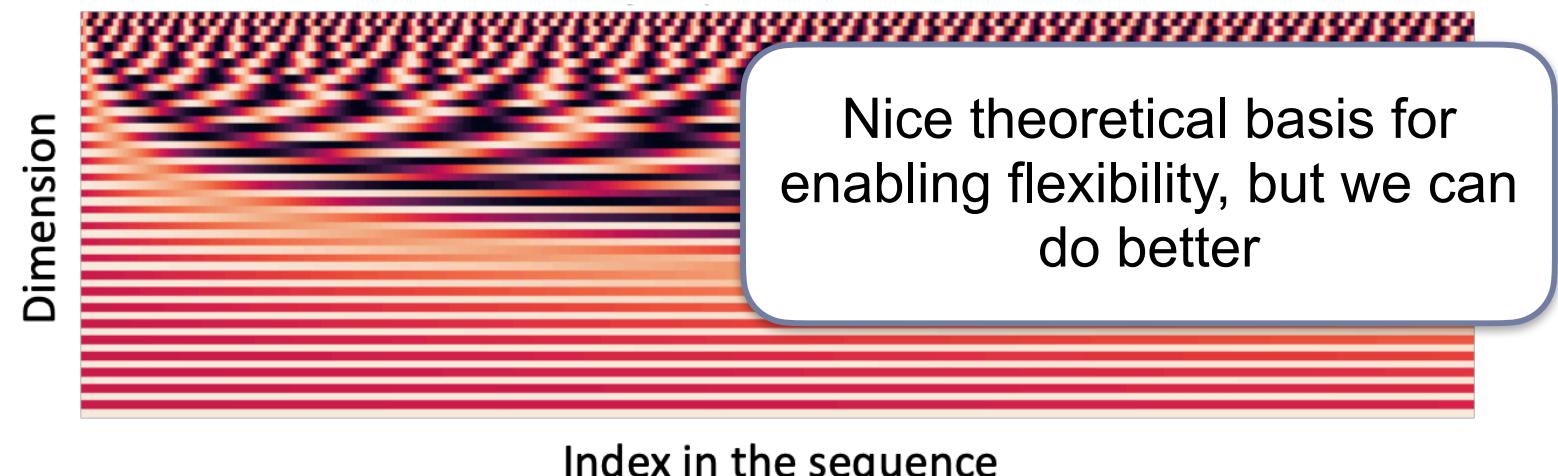
First idea:

$$p_i =$$

$$\begin{aligned} & \sin(i/10000^{2*1/d}) \\ & \cos(i/10000^{2*1/d}) \\ & \vdots \\ & \sin(i/10000^{2*\frac{d}{2}/d}) \\ & \cos(i/10000^{2*\frac{d}{2}/d}) \end{aligned}$$

at top change quickly

at bottom change slowly



Stanford CS224 lectures



Self-Attention

Transformer

Training and Use

Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

To recover all benefits of the RNN we need nonlinearities and **position representation**

What is p_i ?

Next idea:

① Learned values - for each position, learn a vector

problem

Note: can't generalise to sequences longer than those seen in training!

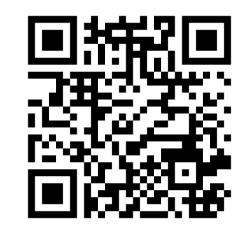


Self-Attention

Transformer

Training and Use

Lab Preview



menti.com 2520 1730

To recover all benefits of the RNN we need nonlinearities and **position dependent representation**

Rotary Positional Embeddings

This chocolate is delicious

I think chocolate is delicious

Most people know chocolate is delicious

(2)

可以看
到 chocolate

位置
嵌入

到
达



Self-Attention

Transformer

Training and Use

Lab Preview



menti.com 2520 1730

To recover all benefits of the RNN we need nonlinearities and **position dependent representation**



Rotary Positional Embeddings

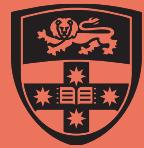
This chocolate is delicious

angel is same
因为两者之间
的间距是一致
没变

I think chocolate is delicious

Most people know chocolate is delicious

Cosine similarity remains the same for
two words the same distance apart!



Self-Attention

Transformer

Training and Use

Lab Preview

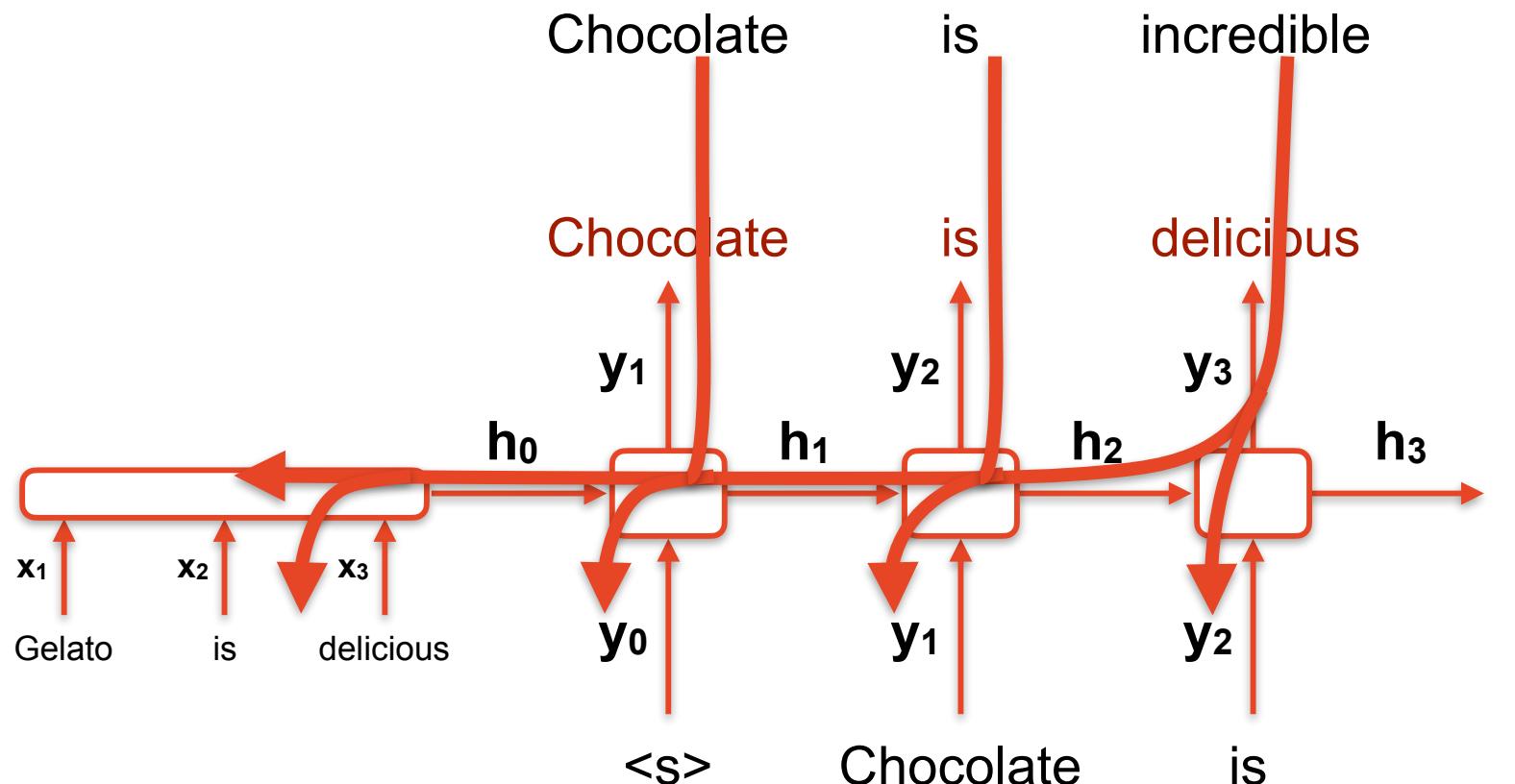


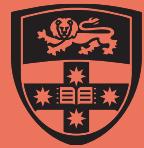
menti.com 2520 1730

Training needs to account for what is known at each step

backprop

△ Guess vs.
Answer





Self-Attention

Transformer
Training and Use
Lab Preview



menti.com 2520 1730

Training needs to account for what is known at each step

' $<\text{s}>$ ' contextual vector

Feedforward Layer

Weighted Average

$x_1 + p_1$
 $<\text{s}>$

What if we try to train
on all at once?

不能看前面
的 Train data

'is' contextual vector

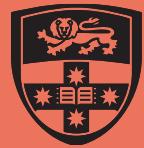
Feedforward Layer

Weighted Average

$x_1 + p_1$
 $<\text{s}>$

$x_2 + p_2$
Chocolate

$x_3 + p_3$
is



Self-Attention

Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Training needs to account for what is known at each step

否冗

'<s>' contextual vector

Feedforward Layer

'is' contextual vector

Feedforward Layer

Weighted Average

Cheating!

Weighted Average

Mask out
these values

$x_1 + p_1$
<math><math><math>

$x_2 + p_2$
<math><math><math>

$x_3 + p_3$
<math><math><math>

Chocolate

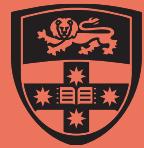
$x_1 + p_1$
<math><math><math>

<math><math><math>

$x_2 + p_2$
<math><math><math>

<math><math><math>

is



Self-Attention

Transformer
Training and Use
Lab Preview



menti.com 2520 1730

Training needs to account for what is known at each step

Initial word vectors
with positions

Transform into query,
keys, and values

Do dot product to get
attention scores

Softmax to get a
distribution

Weighted average

Feedforward layer

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \mathbf{p}_i$$

$$\mathbf{q}_i = Q\tilde{\mathbf{x}}_i$$

$$\mathbf{k}_i = K\tilde{\mathbf{x}}_i$$

$$\mathbf{v}_i = V\tilde{\mathbf{x}}_i$$

$$e_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

Training:

$$e_{ij} = \begin{cases} \mathbf{q}_i^\top \mathbf{k}_j, & j \leq i \\ -\infty, & j > i \end{cases}$$

*bav chart
in Graph*

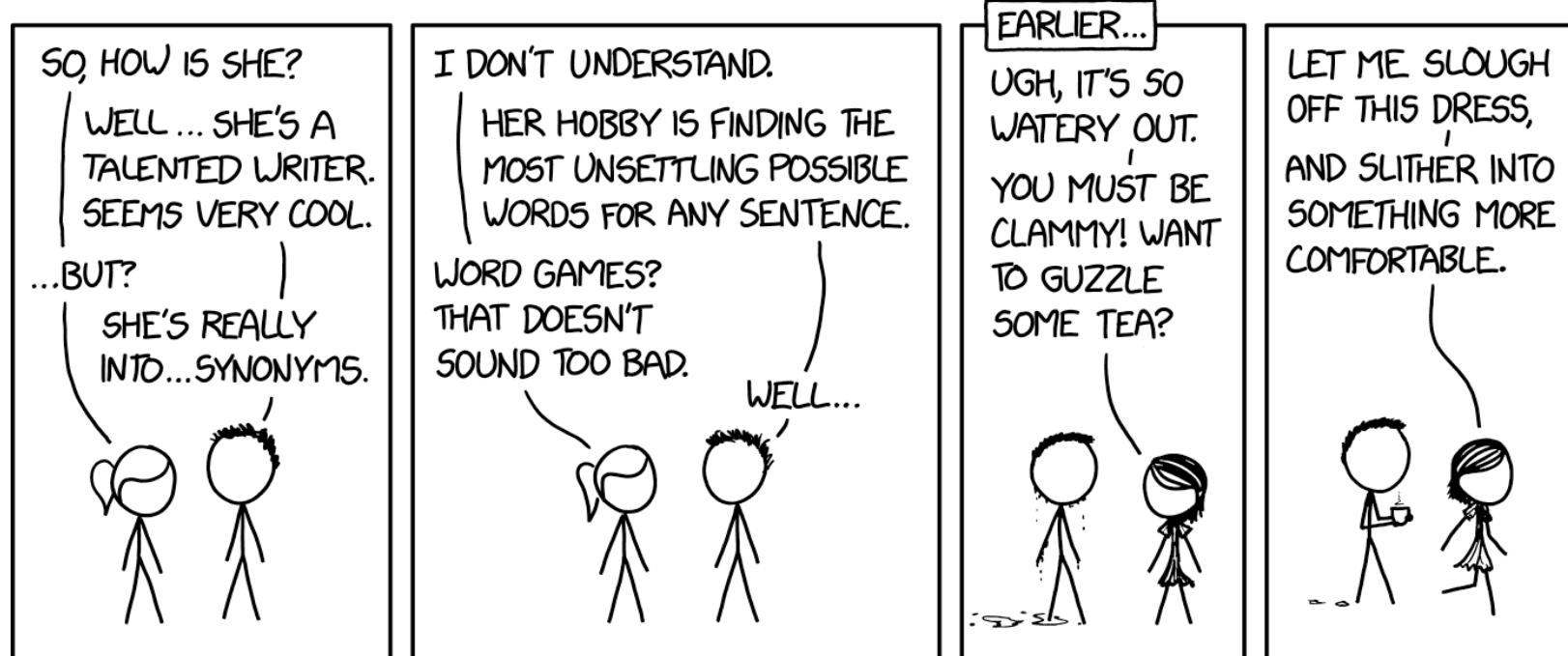
$$\alpha_{ij} = \text{softmax}(e_{ij})$$

$$\mathbf{t}_i = \sum_j \alpha_{ij} \mathbf{v}_j$$

$$\mathbf{o}_i = W_2 \text{ReLU}(W_1 \mathbf{t}_i + \mathbf{b}_1) + \mathbf{b}_2$$



Synonym Date



[We need some grub to munch—I'll go slouch over to the kitchen.]

Source: <https://xkcd.com/2352/>



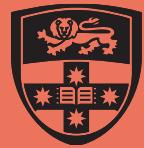
COMP 4446 / 5046
Lecture 7, 2025

Self-Attention
Transformer
Training and Use
Lab Preview



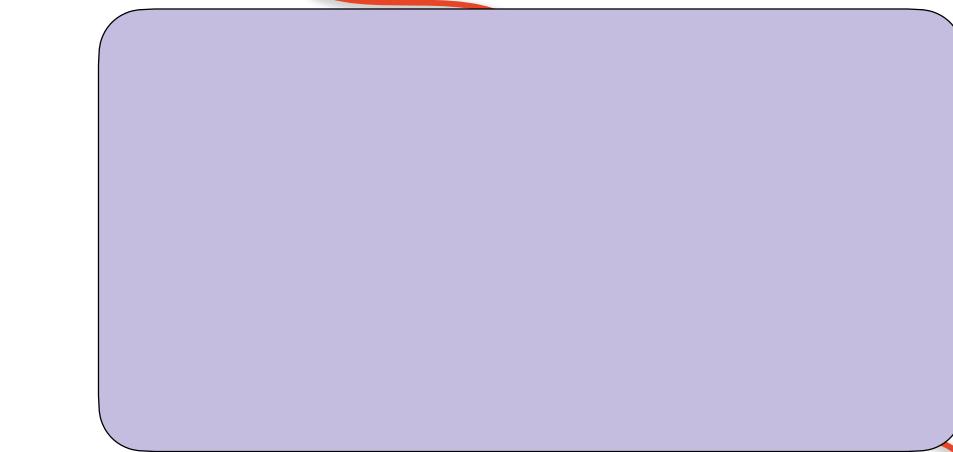
[menti.com 2520 1730](https://menti.com/25201730)

Transformer



Going from Self-Attention to the Transformer

'Gelato' contextual vector



用 box 表示之前在 self-attention 里做的

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

\oplus Element-wise sum

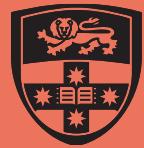
Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

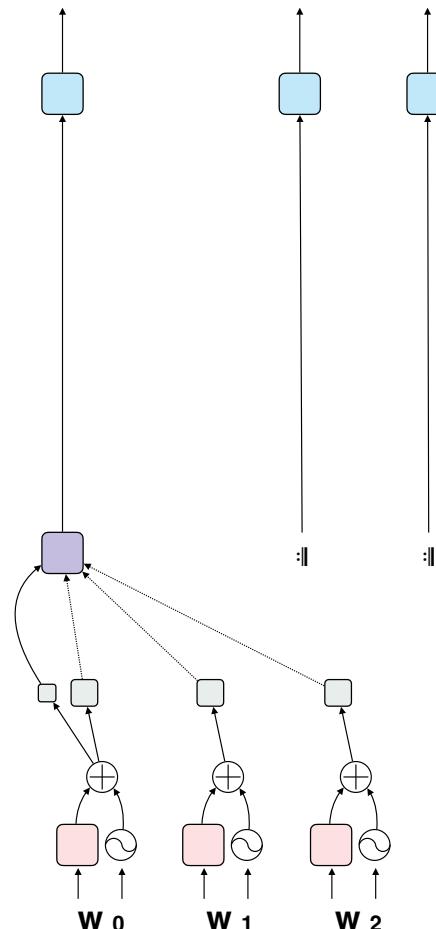
Linear transform

Look up vector

\odot Position as vector



Going from Self-Attention to the Transformer



Feed Forward

$$\max(\mathbf{0}, \mathbf{xW}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

\oplus Element-wise sum

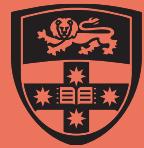
Attention:

$$\text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

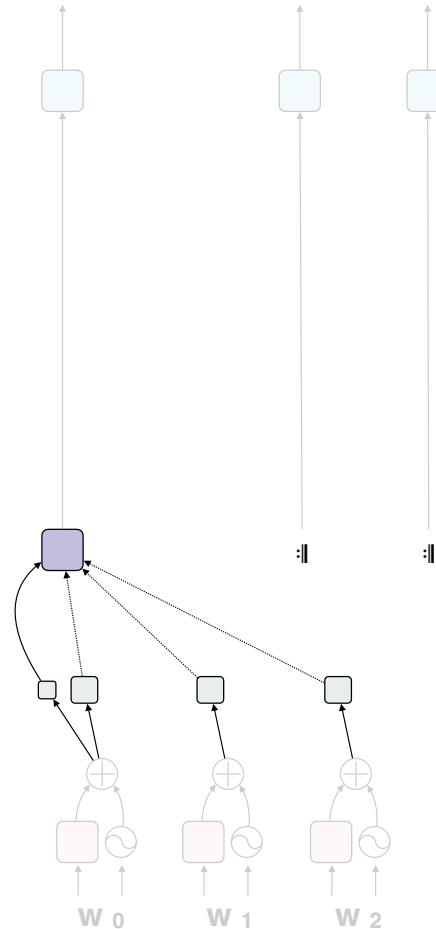
Linear transform

Look up vector

\circlearrowright Position as vector



Going from Self-Attention to the Transformer: Multi-Head Attention



Feed Forward

$$\max(\mathbf{0}, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

\oplus Element-wise sum

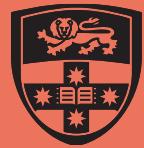
Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

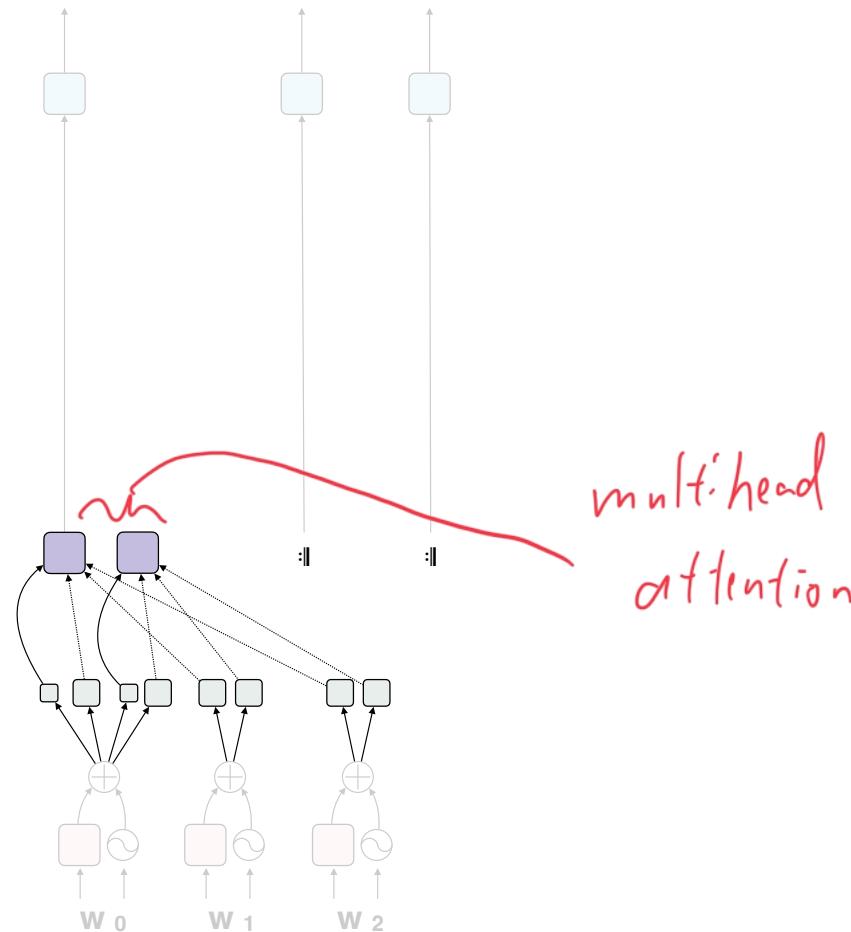
Linear transform

Look up vector

\circlearrowleft Position as vector



Going from Self-Attention to the Transformer: Multi-Head Attention



Feed Forward

$$\max(\mathbf{0}, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

\oplus Element-wise sum

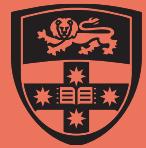
Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

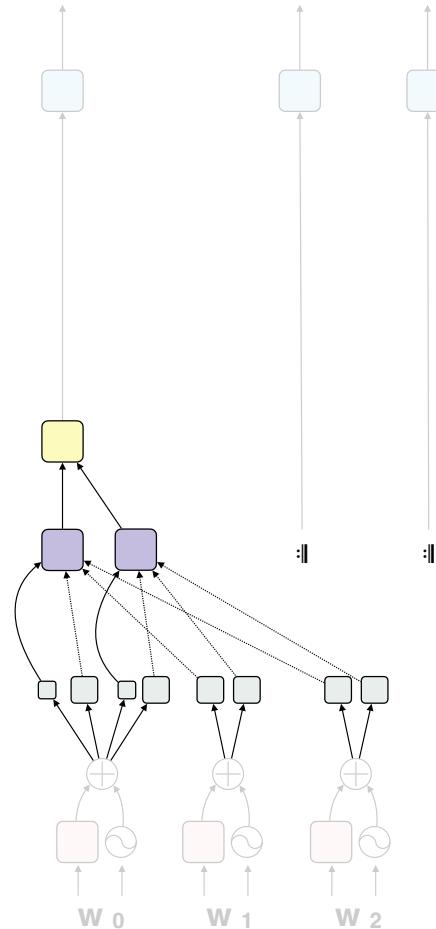
Linear transform

Look up vector

\odot Position as vector



Going from Self-Attention to the Transformer: Multi-Head Attention



Feed Forward

$$\max(\mathbf{0}, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

\oplus Element-wise sum

Concatenate

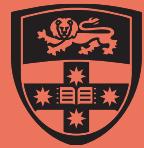
Attention:

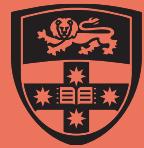
$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

Look up vector

\odot Position as vector





Why do scaled dot-product attention?

Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{V}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Can
Avoid

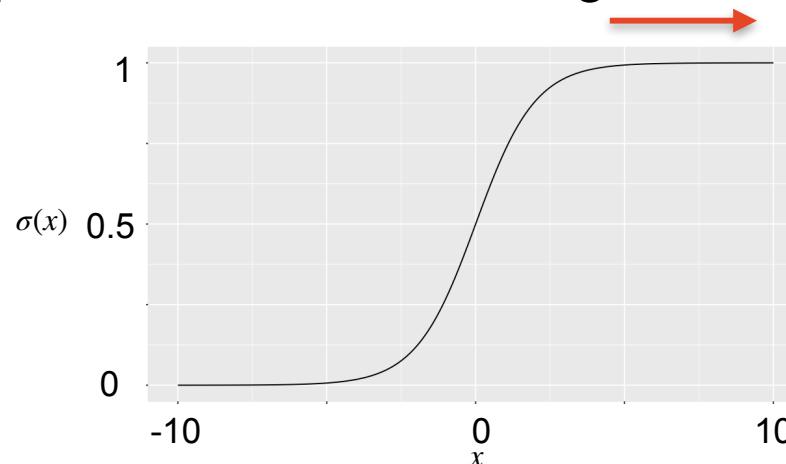
If we increase dimensionality, dot products become larger:

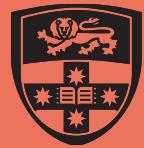
$$[1, 2] \cdot [2, 3] = 8$$

$$[1, 2, 1] \cdot [2, 3, 4] = 12$$

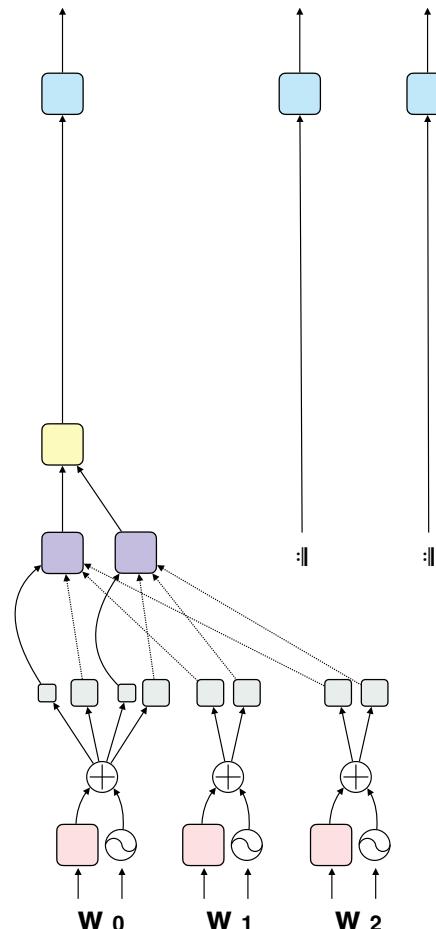
$$[1, 2, 1, 1] \cdot [2, 3, 4, 8] = 20$$

Larger dot products mean smaller gradients after softmax:





Going from Self-Attention to the Transformer: Multi-Head Attention, Residuals



Feed Forward

$$\max(\mathbf{0}, \mathbf{xW}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

\oplus Element-wise sum

Concatenate

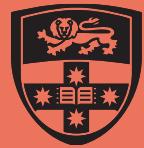
Attention:

$$\text{softmax} \left(\frac{\mathbf{QK}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

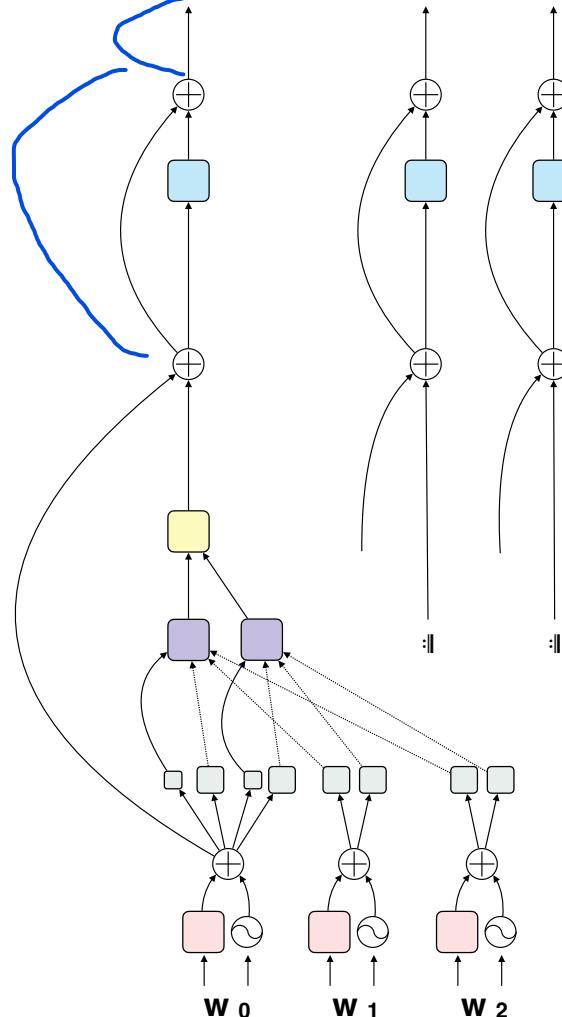
Linear transform

Look up vector

\odot Position as vector



Going from Self-Attention to the Transformer: Multi-Head Attention, Residuals



Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

⊕ Element-wise sum

Concatenate

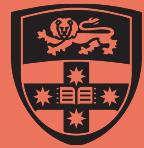
Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

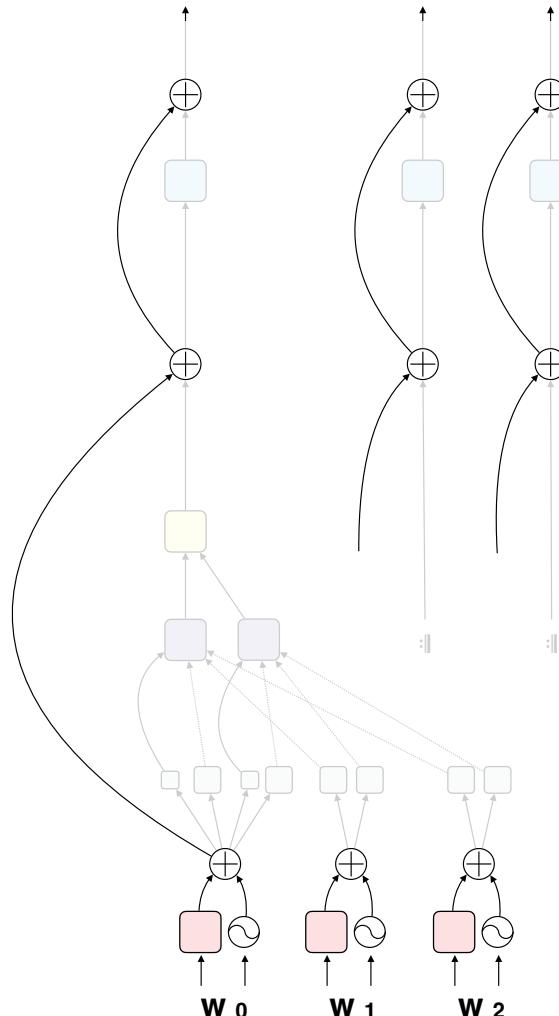
Linear transform

Look up vector

○ Position as vector



Going from Self-Attention to the Transformer: Multi-Head Attention, Residuals



Feed Forward

$$\max(\mathbf{0}, \mathbf{xW}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

⊕ Element-wise sum

Concatenate

Attention:

$$\text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

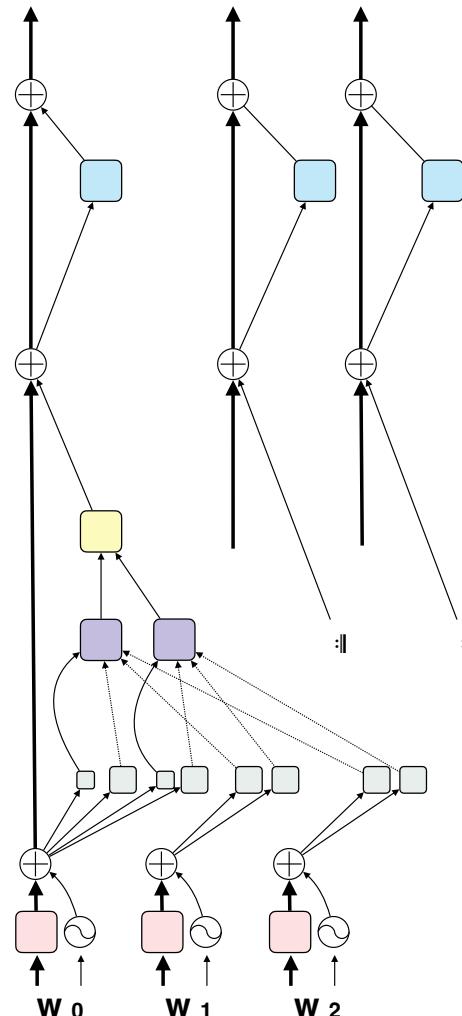
Linear transform

Look up vector

⊖ Position as vector



Going from Self-Attention to the Transformer: Multi-Head Attention, Residuals



Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

⊕ Element-wise sum

Concatenate

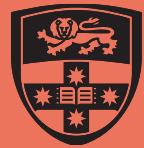
Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

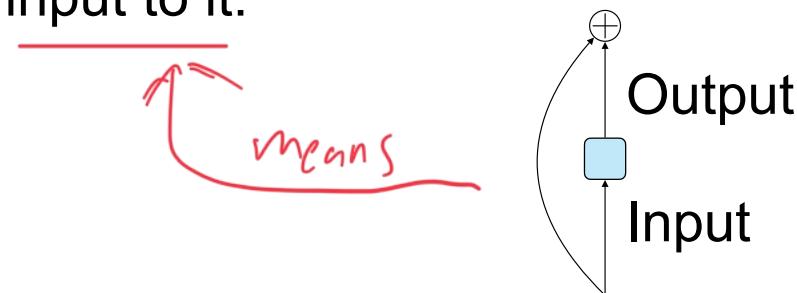
Look up vector

⊕ Position as vector



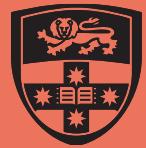
What is a residual?

Originally from Computer Vision.
Taking the output of some calculation and adding the input to it:



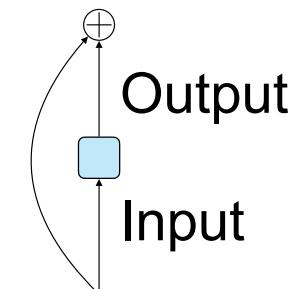
Why?

- Gradient passes back easily (no non-linearity on the edge we added!)
- Encourages our model to capture something close to the identity function
- Frees up the model to use the weights to capture stuff other than the identity function

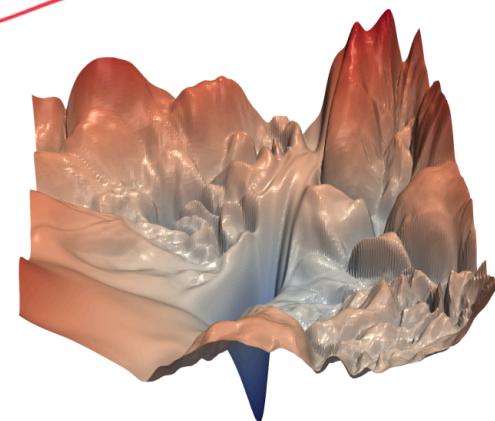


What is a residual?

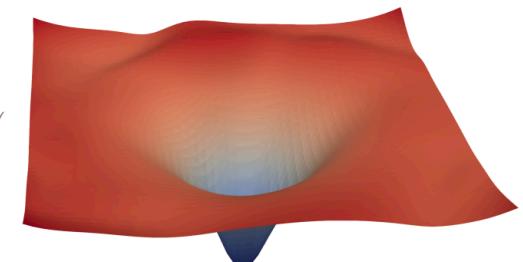
Originally from Computer Vision.
Taking the output of some calculation and adding the input to it:



Smooth gradients.

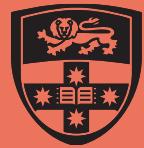


No residuals

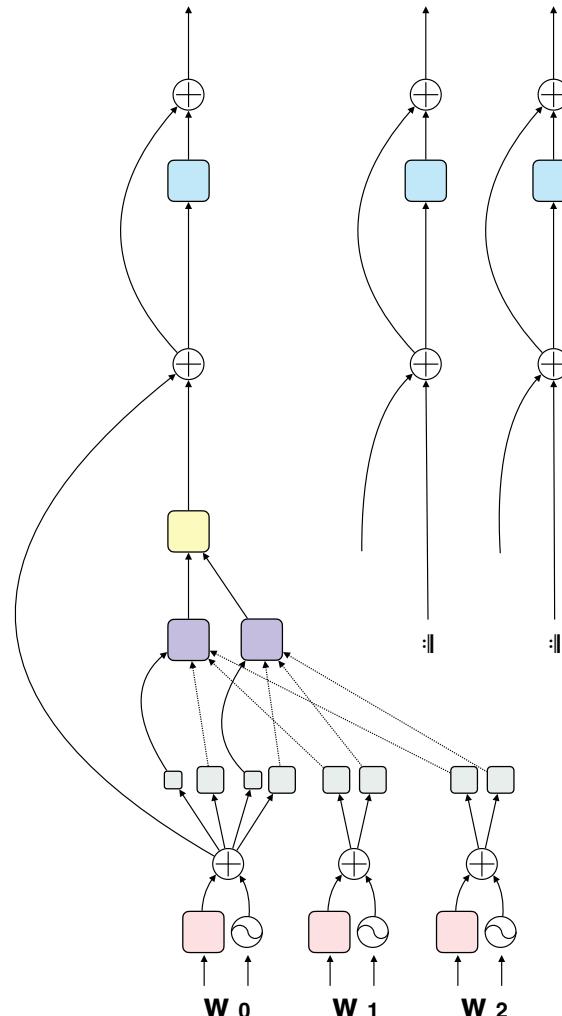


Residuals

Li, et. Al. (NeurIPS 2018)



Going from Self-Attention to the Transformer: Multi-Head Attention, Residuals, and Layer Normalization



Feed Forward

$$\max(\mathbf{0}, \mathbf{xW}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

⊕ Element-wise sum

Concatenate

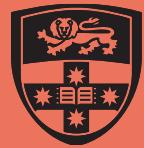
Attention:

$$\text{softmax} \left(\frac{\mathbf{QK}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

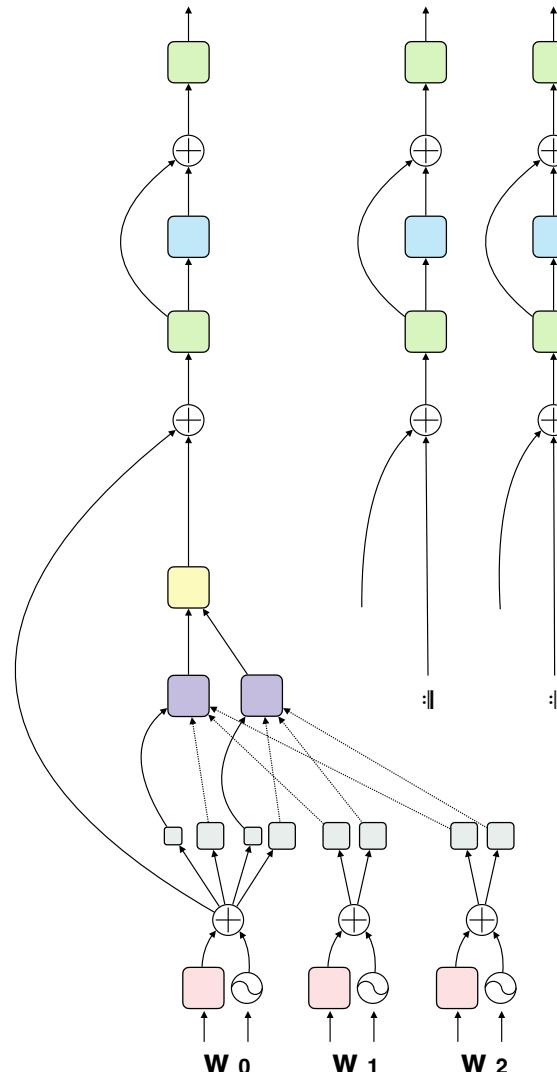
Linear transform

Look up vector

⊖ Position as vector



Going from Self-Attention to the Transformer: Multi-Head Attention, Residuals, and Layer Normalization



Feed Forward

$$\max(\mathbf{0}, \mathbf{xW}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

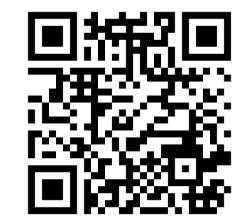
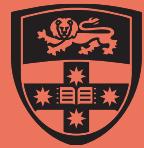
Attention:

$$\text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

Linear transform

Look up vector

⊖ Position as vector



Layer Normalization shifts values in a way that improves training

Take a vector

$$x \in \mathbb{R}^d$$

Calculate the mean and standard deviation of the vector's values

$$\mu = \sum_{j=1}^d x_j$$

$$\sigma = \sqrt{\frac{1}{d} \sum_{j=1}^d (x_j - \mu)^2}$$

Define two learned weight vectors (or set them to 1 and 0)

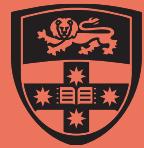
$$\gamma \in \mathbb{R}^d \quad \beta \in \mathbb{R}^d$$

Rescale the vector

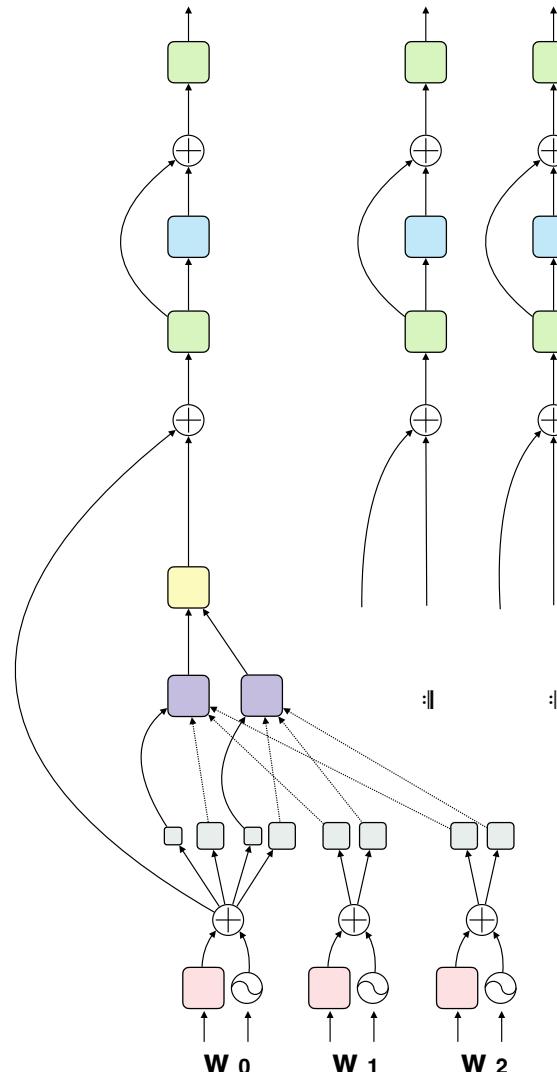
$$\text{output} = \frac{x - \mu}{\sqrt{\sigma} + \epsilon} * \gamma + \beta$$

has square root

A small value to prevent division by zero



Going from Self-Attention to the Transformer: Multi-Head Attention, Residuals, and Layer Normalization



Feed Forward

$$\max(\mathbf{0}, \mathbf{xW}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

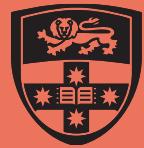
Attention:

$$\text{softmax} \left(\frac{\mathbf{QK}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Linear transform

Look up vector

⊖ Position as vector

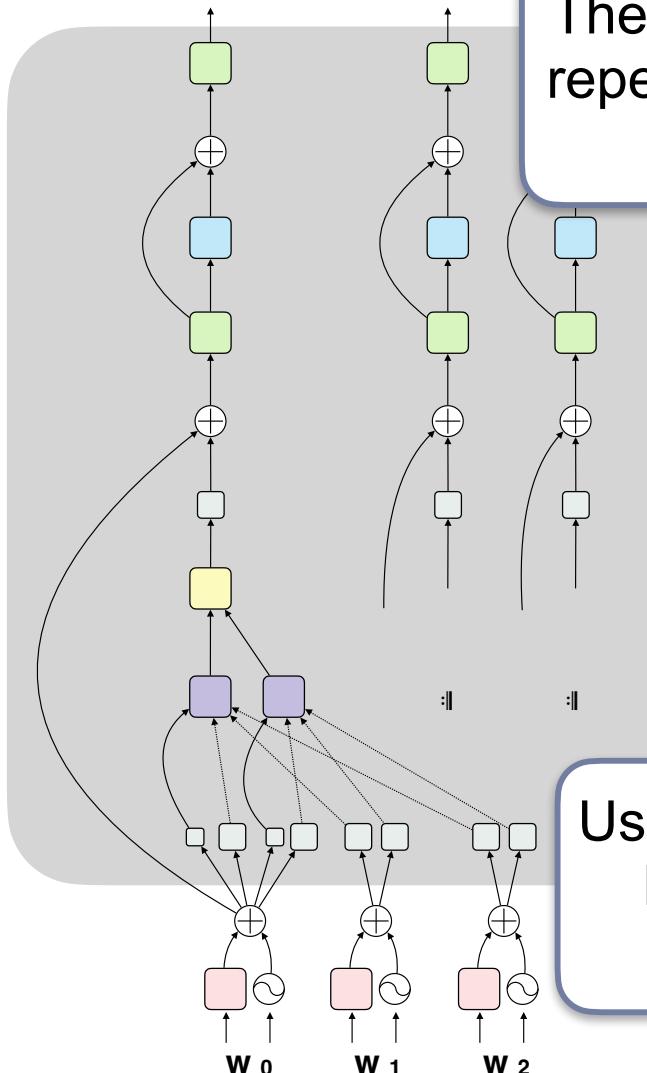


Self-Attention **Transformer** Training and Use Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

This is the Transformer encoder! (Almost)



The shaded area is repeated in a stack, 6 times

Usually more than 2 heads (original paper had 8)

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

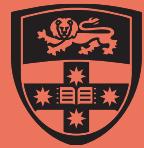
Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Linear transform

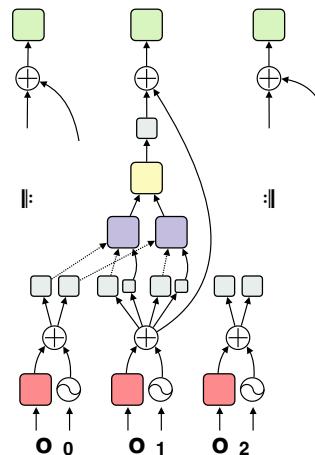
Look up vector

○ Position as vector



What about
the decoder?

Start the same as the encoder,
except attention does not apply to
the future



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

\oplus Element-wise sum

Concatenate

Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

Look up vector

\odot Position as vector



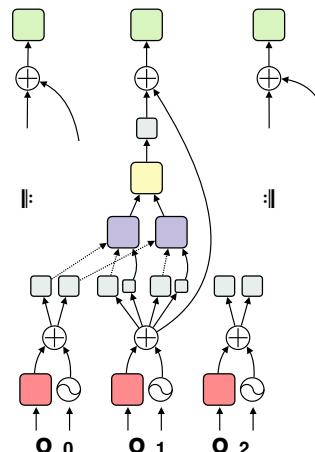
Self-Attention **Transformer** Training and Use Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

What about the decoder?

Provide encoder state
using attention, also
multi-head



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

\oplus Element-wise sum

Concatenate

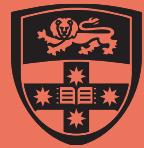
Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

Look up vector

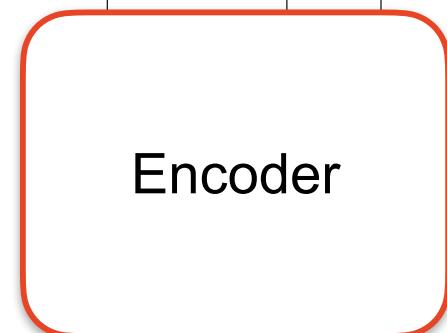
\odot Position as vector



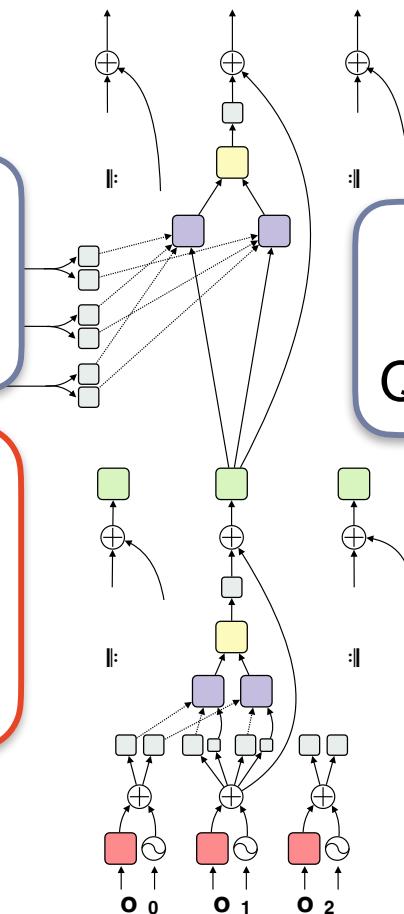
What about the decoder?

Provide encoder state using attention, also multi-head

Called 'cross attention' to distinguish it from 'self attention'



Input text



$T_{decoder}$

Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

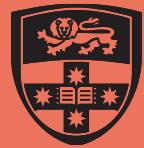
K and V come from encoder
Q comes from decoder

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

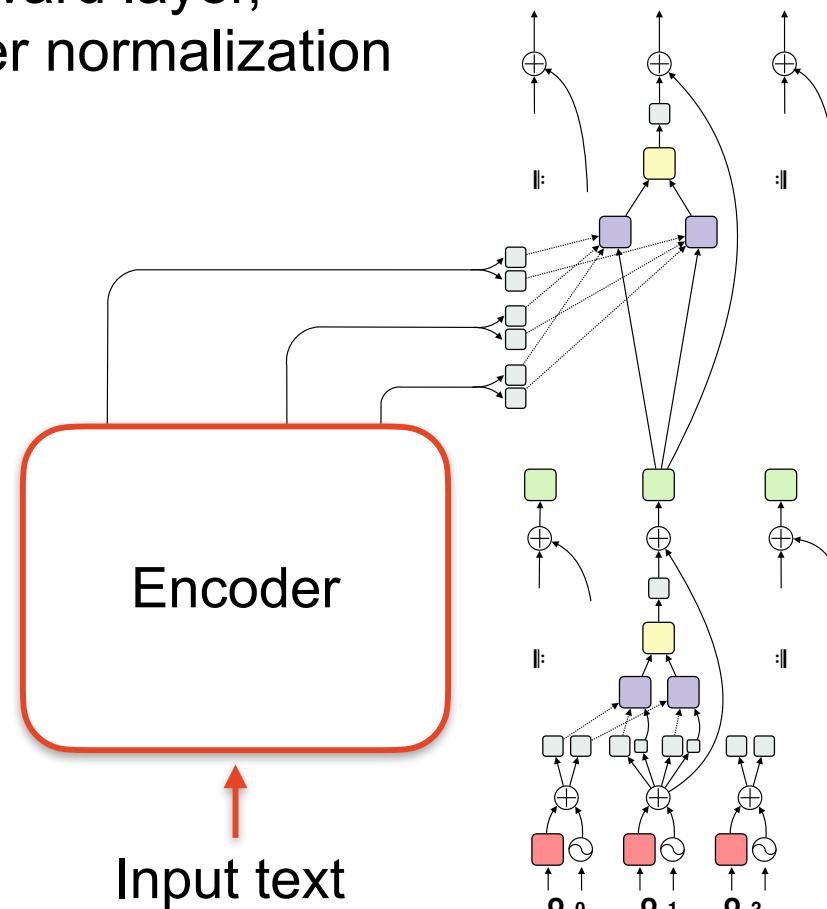
Look up vector

Position as vector



What about the decoder?

Add another
feed forward layer,
with layer normalization



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Linear transform

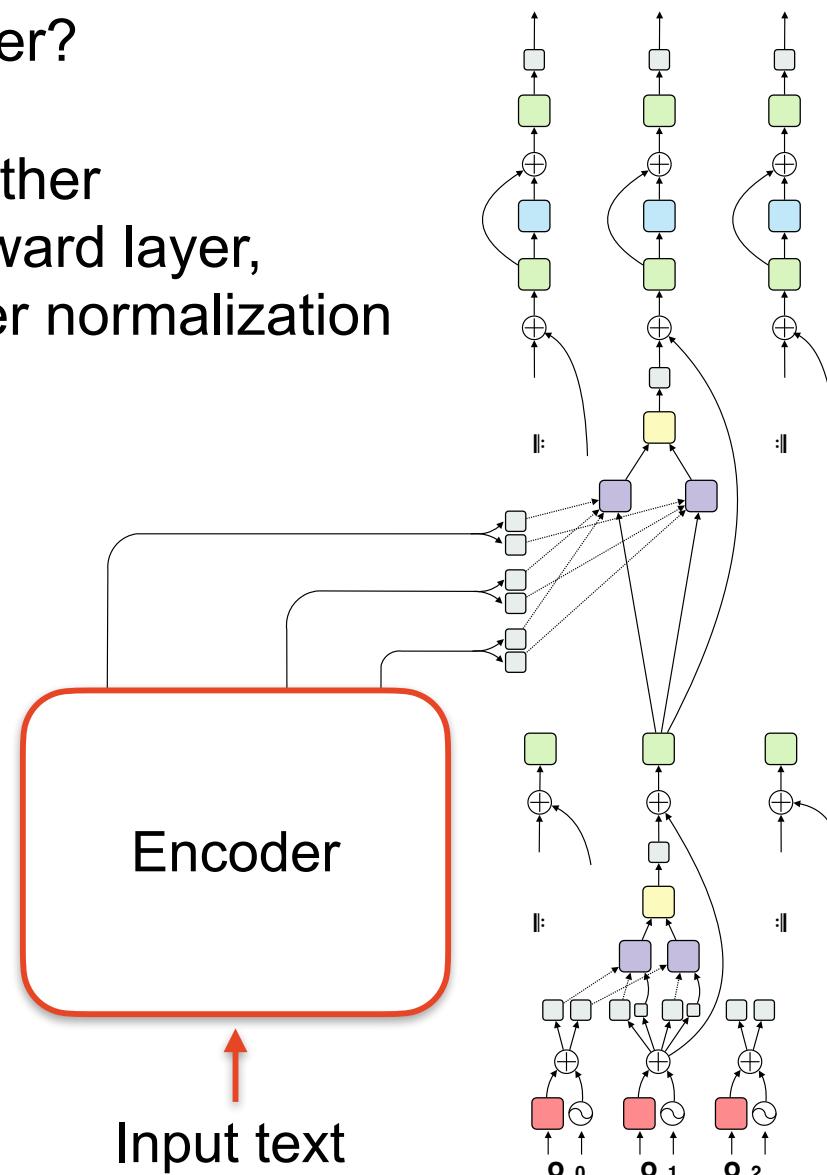
Look up vector

○ Position as vector



What about the decoder?

Add another
feed forward layer,
with layer normalization



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

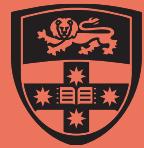
Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Linear transform

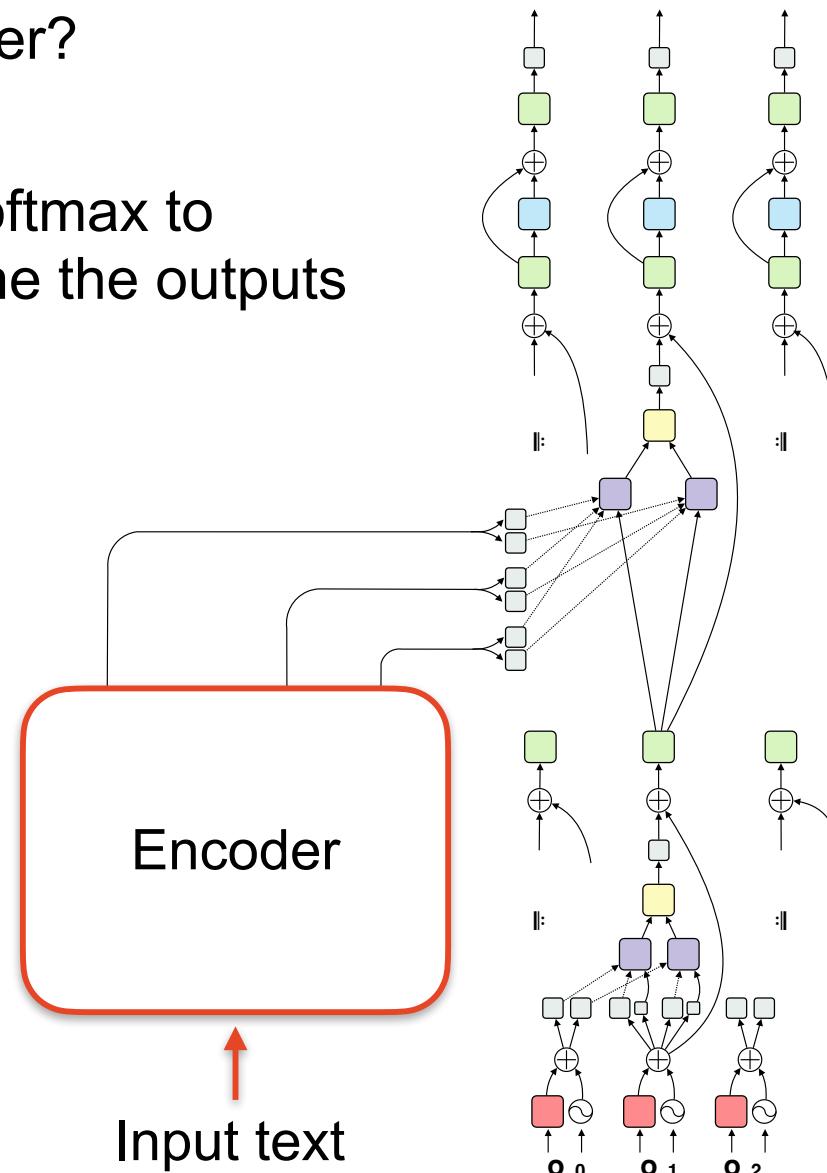
Look up vector

○ Position as vector



What about the decoder?

Apply softmax to
determine the outputs



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

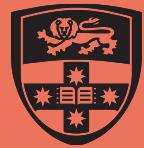
Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Linear transform

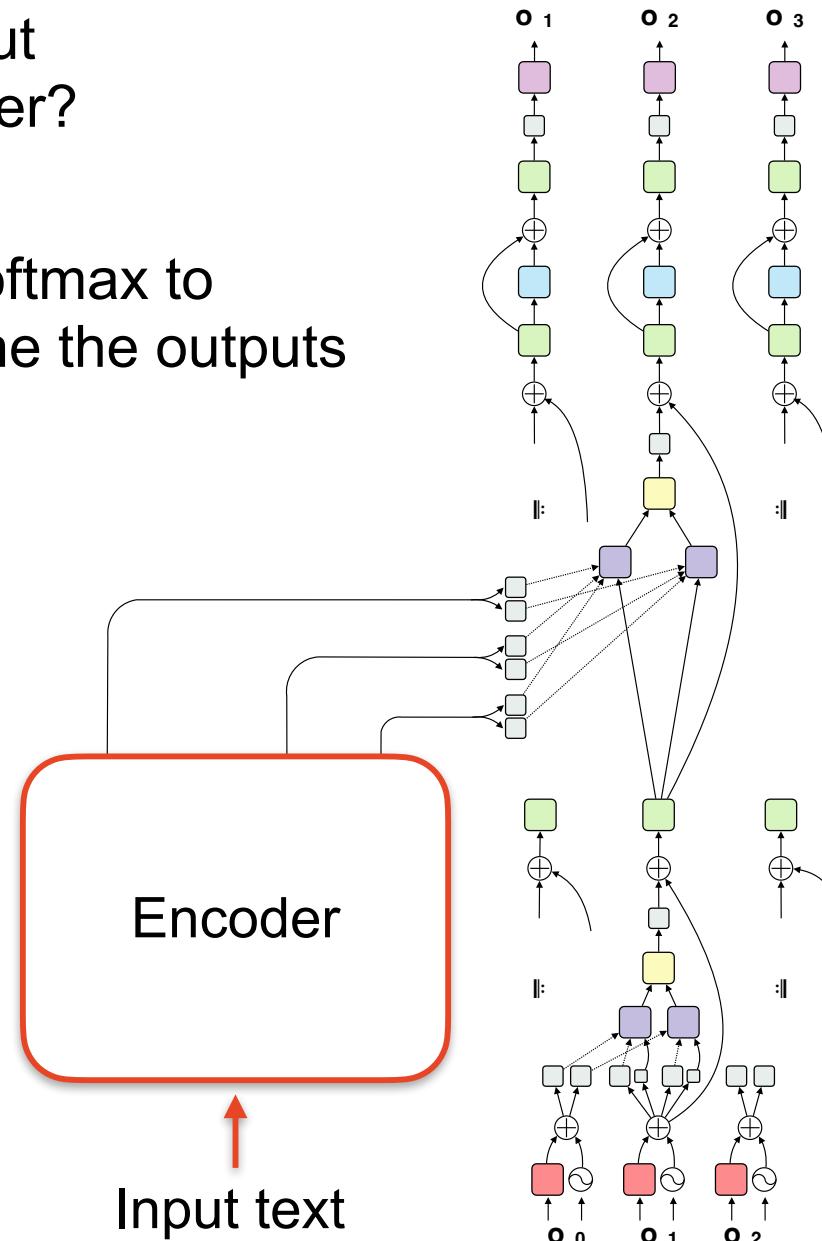
Look up vector

○ Position as vector



What about
the decoder?

Apply softmax to
determine the outputs



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

Element-wise sum

Concatenate

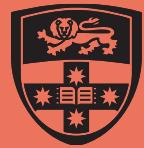
Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

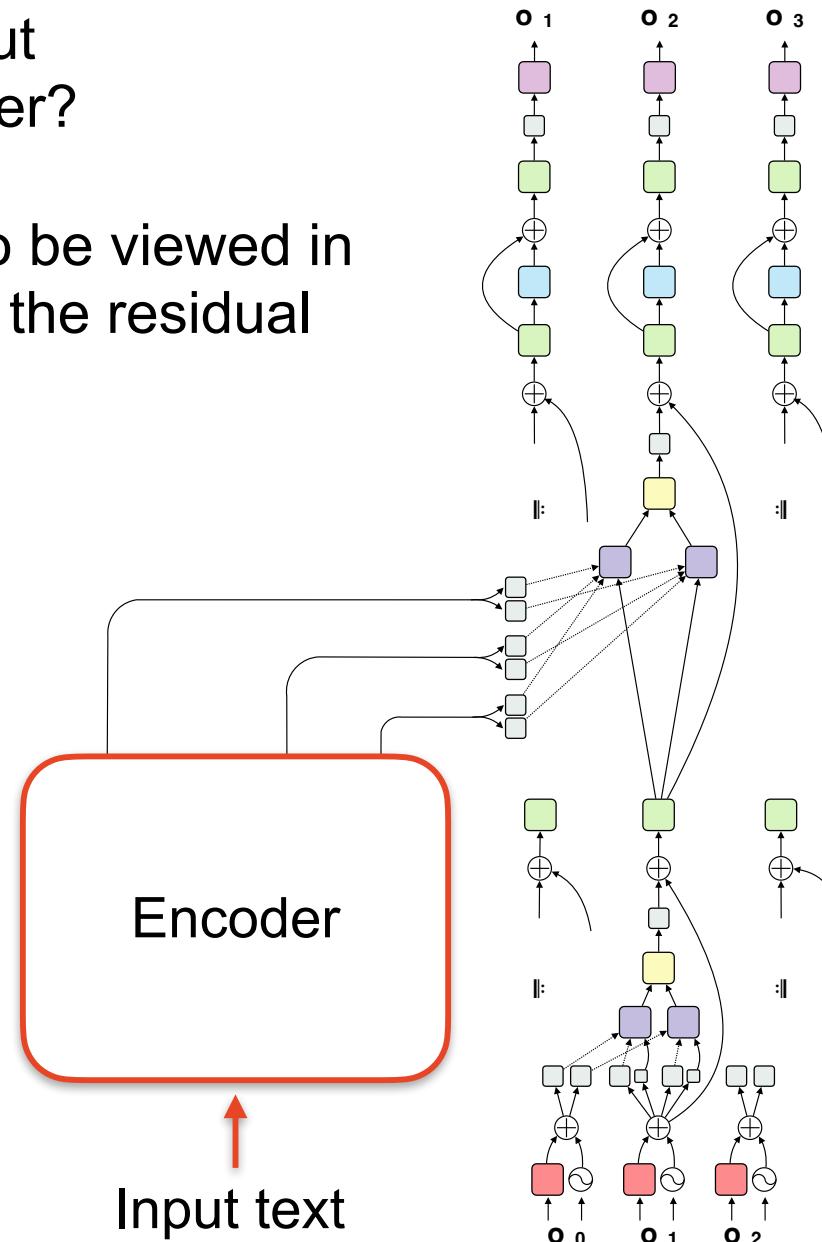
Look up vector

Position as vector



What about
the decoder?

Can also be viewed in
terms of the residual
stream



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

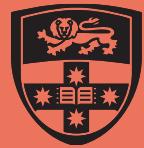
Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

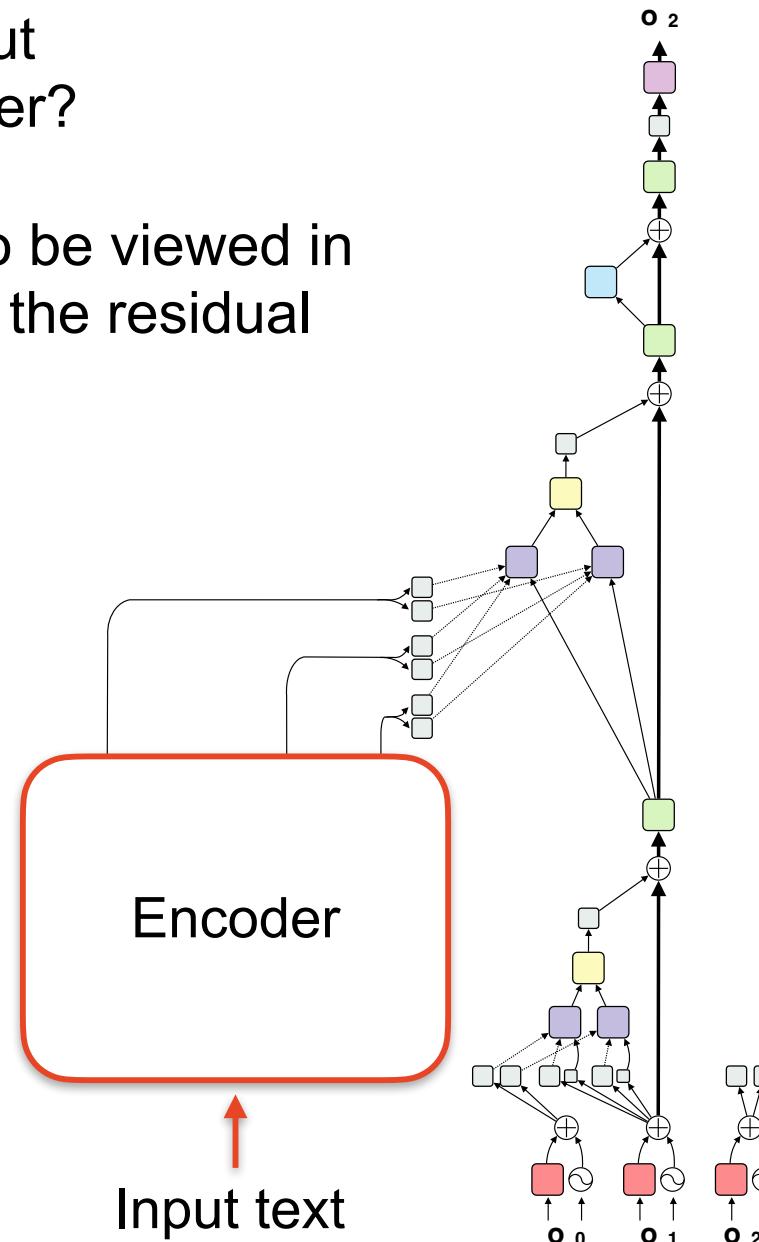
Look up vector

○ Position as vector



What about
the decoder?

Can also be viewed in
terms of the residual
stream



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

⊕ Element-wise sum

Concatenate

Attention:

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Linear transform

Look up vector

○ Position as vector



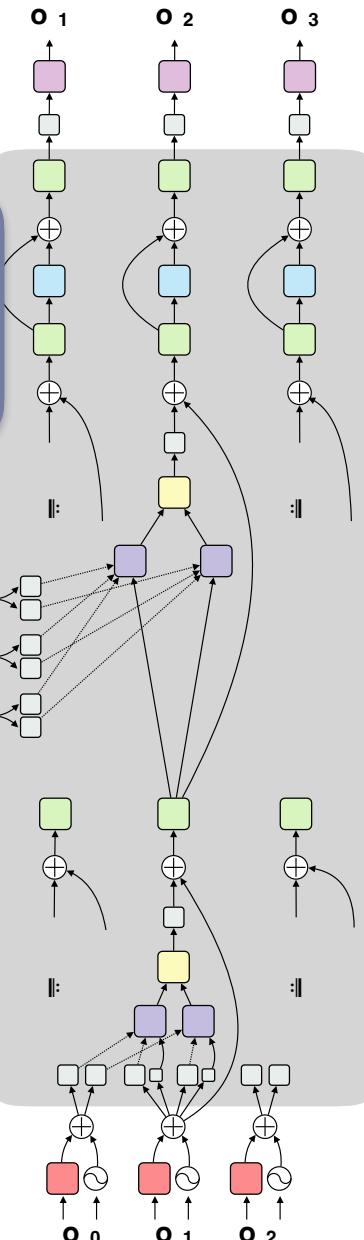
This is the Transformer decoder! (Almost)

The shaded area is repeated in a stack,
6 times

Usually more than 2
heads (original
paper had 8)

Encoder

Input text



Softmax

Feed Forward

$$\max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

Layer Norm

Element-wise sum

Concatenate

Attention:

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Linear transform

Look up vector

Position as vector



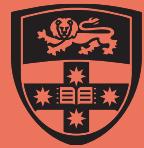
COMP 4446 / 5046
Lecture 7, 2025

Self-Attention
Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Training and Use

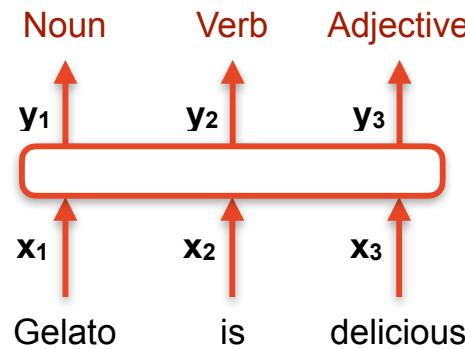


Self-Attention
Transformer
Training and Use
Lab Preview

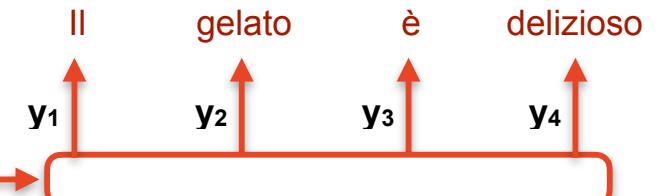
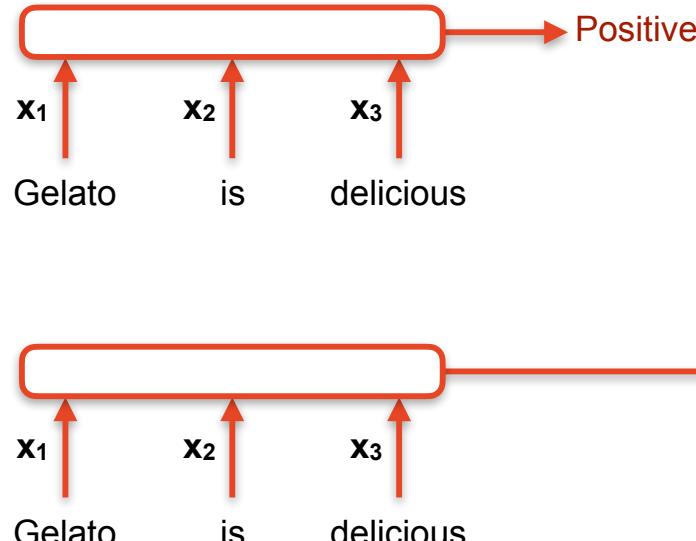


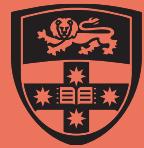
[menti.com 2520 1730](https://menti.com/25201730)

How do we use this? - Just like the RNN encoder-decoders

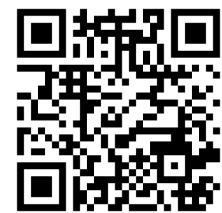


Or use it to create contextual word representations that go into another model





Self-Attention
Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Initial results showed slight accuracy gains, big efficiency gains

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	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		$3.3 \cdot 10^{18}$
Transformer (big)	28.4	41.8		$2.3 \cdot 10^{19}$

Vaswani, et al., (NeurIPS 2017)



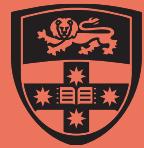
Over time, transformer models came to be dominant

Summarisation (Lin, et. al., ICLR 2018)

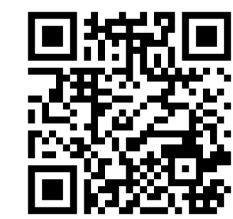
Model	Test perplexity	ROUGE-L
<i>seq2seq-attention, L = 500</i>	5.04952	12.7
<i>Transformer-ED, L = 500</i>	2.46645	34.2
<i>Transformer-D, L = 4000</i>	2.22216	33.6
<i>Transformer-DMCA, no MoE-layer, L = 11000</i>	2.05159	36.2
<i>Transformer-DMCA, MoE-128, L = 11000</i>	1.92871	37.9
<i>Transformer-DMCA, MoE-256, L = 7500</i>	1.90325	38.8

GLUE Benchmark - a range of tasks (Devlin, et. al., NAACL 2019)

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7
BiLSTM+ELMo	61.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8
OpenAI GPT	61.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0
BERT _{BASE}	64.0/63.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4
BERT _{LARGE}	65.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1



Self-Attention
Transformer
Training and Use
Lab Preview

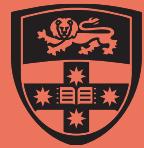


[menti.com 2520 1730](https://menti.com/25201730)

BERT
↓
include

WordPiece tokenisation (last lecture)
+
Transformer architecture (this lecture)
+
Data
=
High performance!

This was the moment when
training large models became
out of reach for most
university research labs



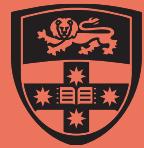
Faster training means they can scale

Original paper (2017)

	N	d_{model}	d_{ff}	h	d_k	d_v
base	6	512	2048	8	64	64

GPT-3 (2020)

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}
GPT-3 Small	125M	12	768	12	64
GPT-3 Medium	350M	24	1024	16	64
GPT-3 Large	760M	24	1536	16	96
GPT-3 XL	1.3B	24	2048	24	128
GPT-3 2.7B	2.7B	32	2560	32	80
GPT-3 6.7B	6.7B	32	4096	32	128
GPT-3 13B	13.0B	40	5140	40	128
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128

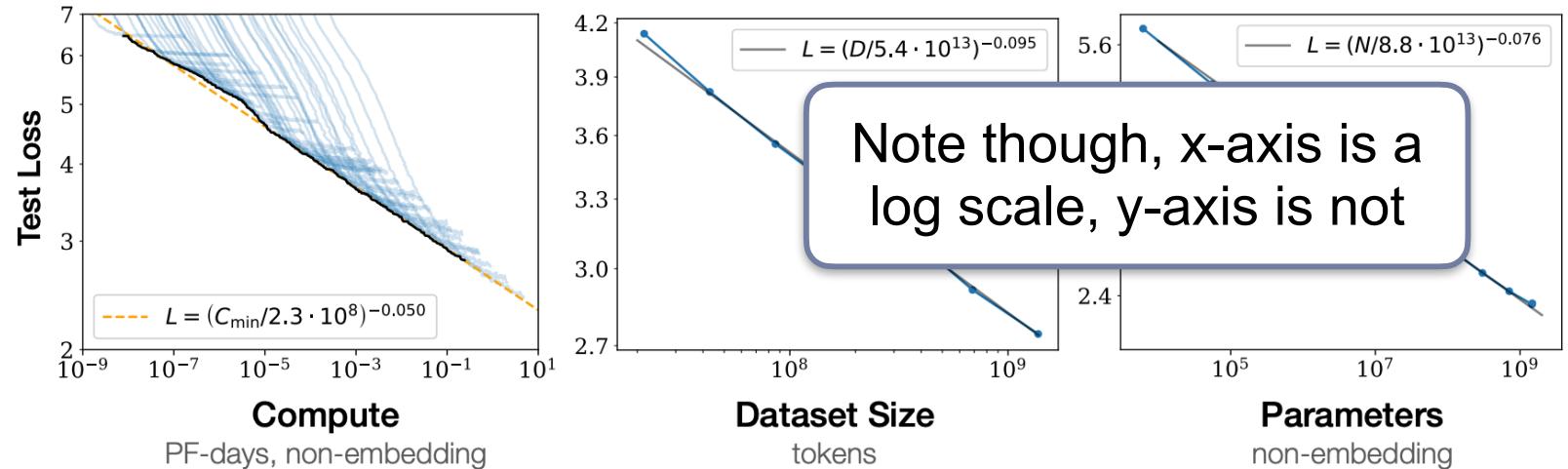


Self-Attention
Transformer
Training and Use
Lab Preview



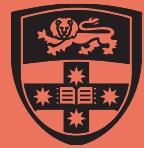
menti.com 2520 1730

They scale very well to more compute, data, and parameters

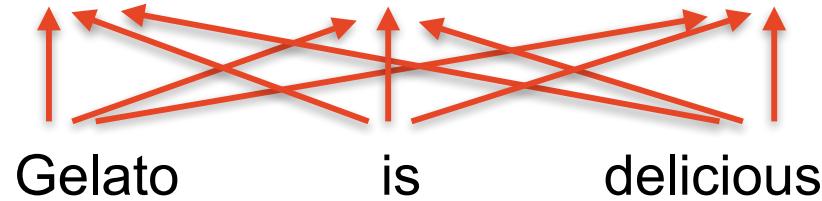


This is where 'scaling laws' come from - methods to estimate the best model size given a certain amount of data and a budget for training

Kaplan, et al., (arXiv 2020)



But, attention is quadratic, will that be an issue?



Lots of clever variants to avoid it!

- PerFormer
- LongFormer
- LinFormer

But... it turns out that the bigger issue is memory bandwidth

Most notable improvement - FlashAttention

Feedforward computation dominates as model size increases



Laser Products

ONLINE REVIEWS OF LASER PRODUCTS		
	...SURGERY	...REMOVAL
LASER	EYE...	★★★★★ "I DON'T NEED GLASSES ANYMORE!"
	JET...	★★★★★ "TOO NERVOUS TO TRY IT."
	HAIR...	★★★★★ "CONFUSING TERM FOR HAIRCUT. BURNING SMELL."
...PRINTER		
	★★★★★ "AAAAAAA! MISREAD THE DESCRIPTION! AAAAAAAAAAAA!!"	★★★★★ "EWU."
	★★★★★ "EFFECTIVE, BUT THE FAA GOT REALLY MAD"	★★★★★ "PRINTS GREAT!"
	★★★★★ "GREAT RESULTS!"	★★★★★ "DISGUSTING, WON'T TURN OFF, JAMS CONSTANTLY."

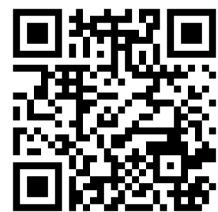
[ERRORS: HAIR JAM. COLOR-SAFE CONDITIONER CARTRIDGE RUNNING LOW. LEGAL-SIZE HAIR TRAY EMPTY, USING LETTER-SIZE HAIR ONLY.]

Source: <https://xkcd.com/1681/>



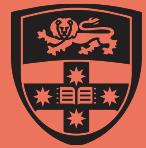
COMP 4446 / 5046
Lecture 7, 2025

Self-Attention
Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

Lab Preview



Self-Attention
Transformer
Training and Use
Lab Preview



[menti.com 2520 1730](https://menti.com/25201730)

1. multiple stack : better performance,
more scenario?

2. mask you pose

HuggingFace!

3. rotary positional embedding

Mid-semester survey

<https://forms.office.com/r/4aHbx4U4mY>



All questions optional

I will read (manually - no AI!)
and report back results + my
plans in response

Muddy Card

<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 7