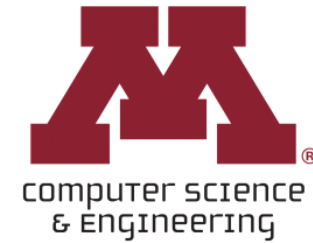


# CSCI 5541: Natural Language Processing

## Lecture 10: Deep Dive on Transformers



Using some slides borrowed from Anna Goldie (Google Brain) and John Hweitt (Stanford)

# Announcement (0304)

- ❑ Cancel Colab Pro Payment → Should be attached to your accounts now
- ❑ Proposal Report
  - Due Tonight
  - Lookout for our review of your report and suggestions for changes/improvements to what you submit
- ❑ HW4 Released (due 3 weeks from today)
- ❑ Small Due Date Changes
  - Homeworks moved from Tuesday to Thursday
  - Project Midterm pushed back one week



# ELMo

(Peters et al., 2018)



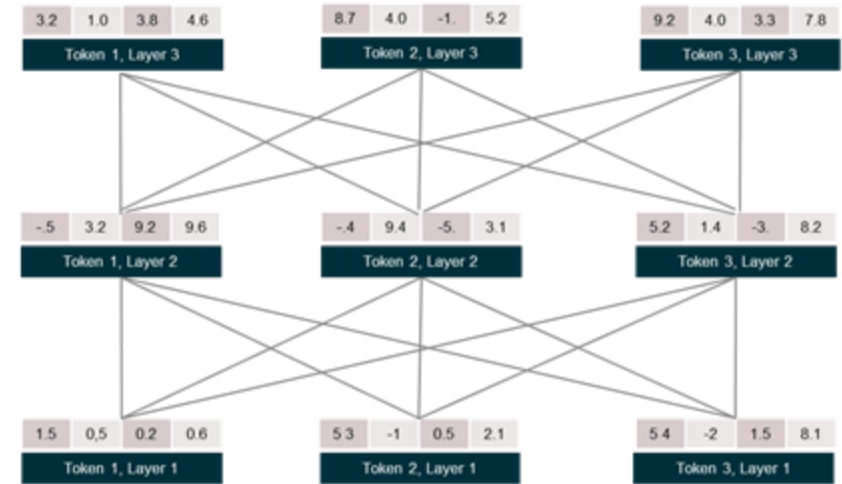
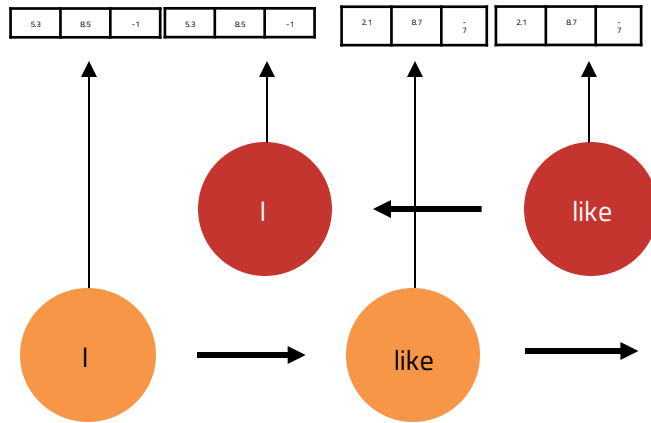
# BERT

(Devlin et al., 2019)



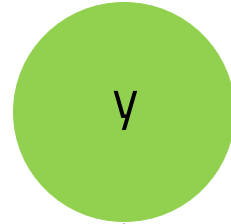
Stacked Bidirectional RNN trained to predict next word in language modeling task

Transformer-based model to predict masked word using bidirectional context and next sentence prediction





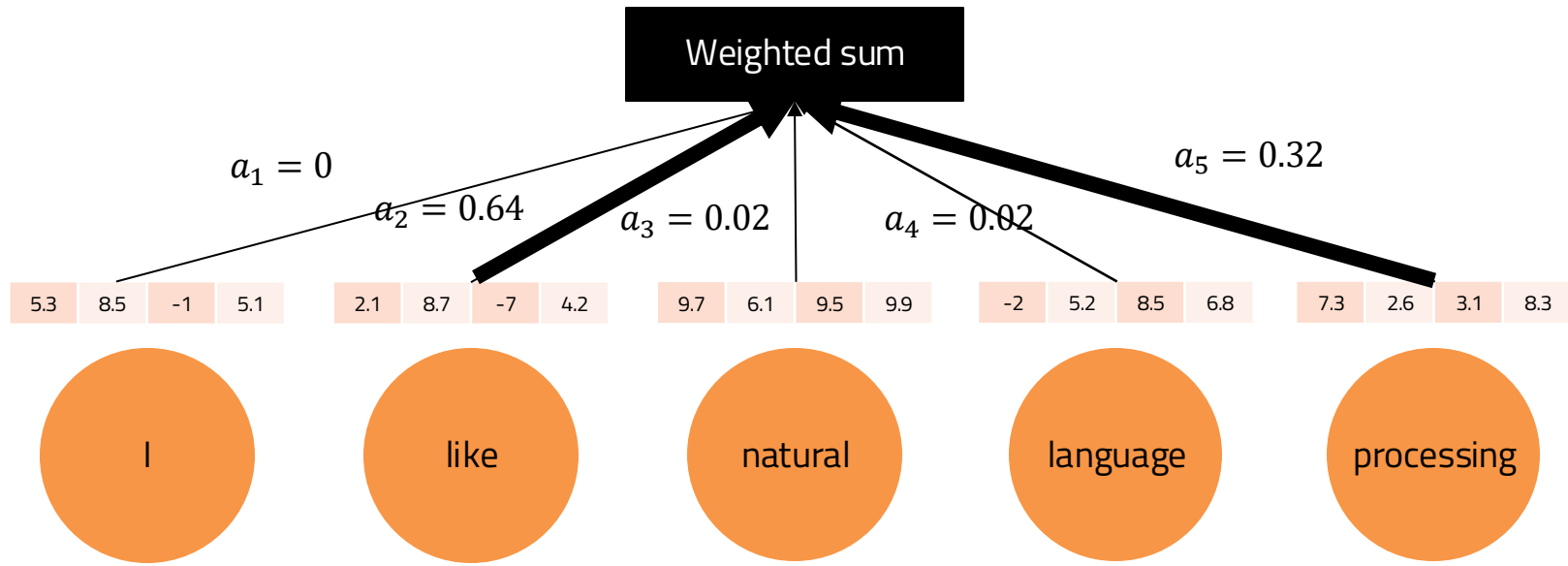
Positive / Negative



6.5 2.2 -3. 9.6

$$x_1a_1 + x_2a_2 + x_3a_3 + x_4a_4 + x_5a_5$$

Weighted sum

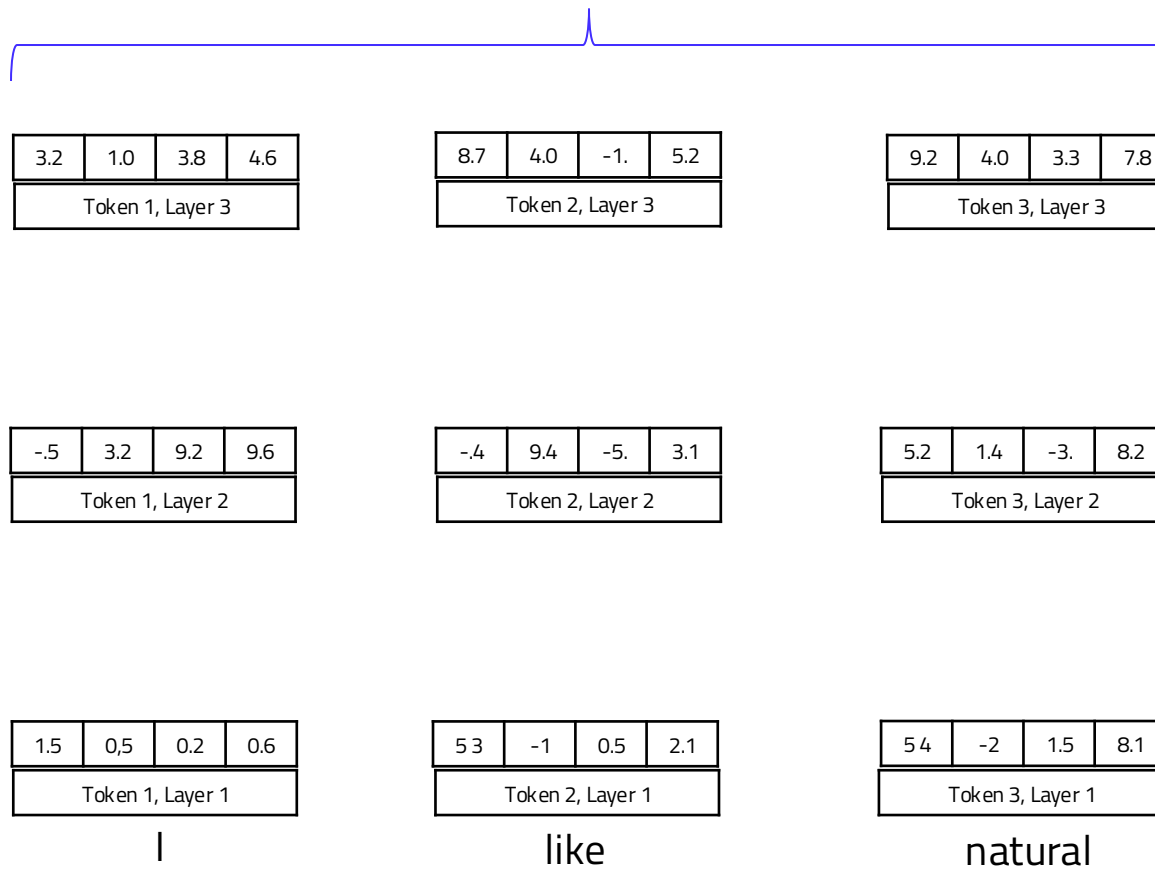




At the end, we have one representation  
for each layer for each token

## Input Length (T)

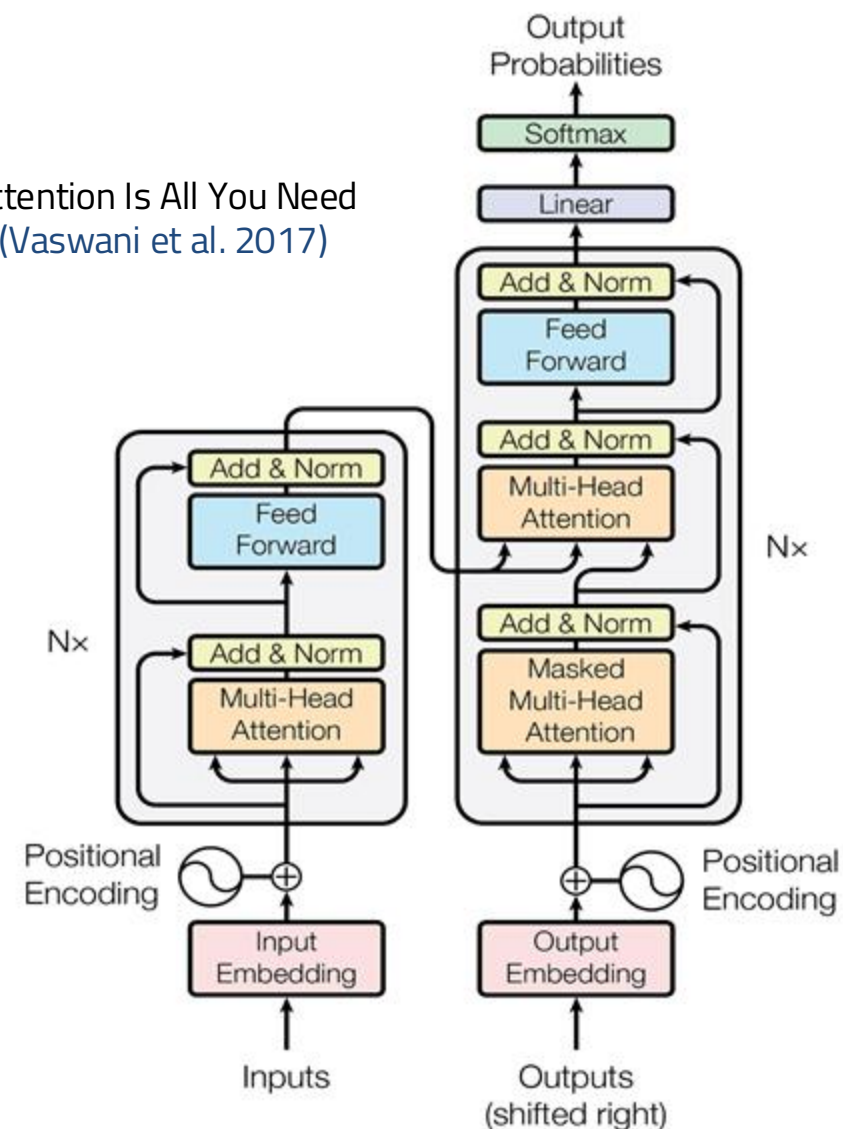
Layers (L)



# Summary of Transformers

- ❑ A sequence-to-sequence model based entirely on **attention**
- ❑ **Strong results** on translation and a wide variety of other tasks
- ❑ Faster: More easy to train in a parallel fashion
- ❑ (At right) Encoder-Decoder Transformer

Attention Is All You Need  
(Vaswani et al. 2017)



# Strong results/findings and applications of Transformers



# Strong results with Transformers on machine translation

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)	<b>27.3</b>	<b>38.1</b>	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.8</b>	<b><math>2.3 \cdot 10^{19}</math></b>	

[Test sets: WMT 2014 English-German and English-French]

(Vaswani et al. 2017)





# Strong results with Transformers on document summarization























<b>Model</b>	<b>Test perplexity</b>	<b>ROUGE-L</b>
<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

WikiSum dataset (Liu et al., 2018)



# Strong results with (pre-trained) Transformers on classification tasks

Sentiment classification  
on SST-2 dataset

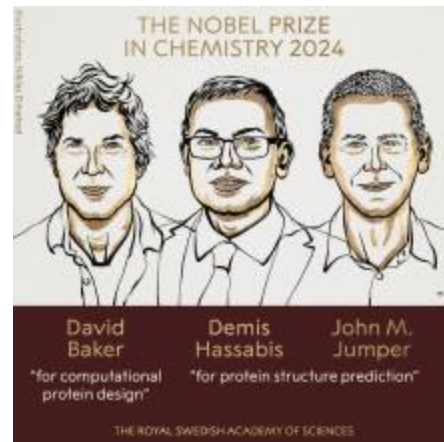
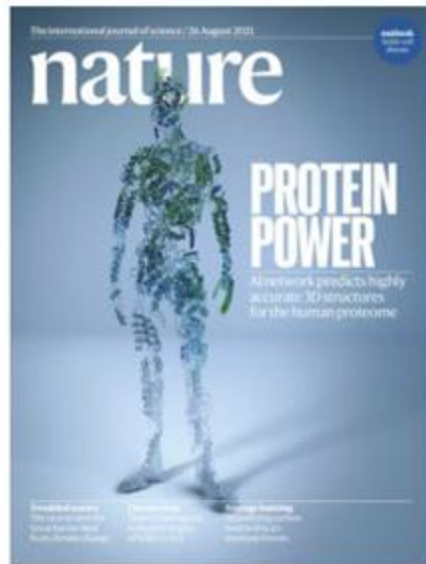
Rank	Model	Accuracy↑	Paper	Code	Result	Year	Tags 
1	SMART-RoBERTa Large	97.5	SMART: Robust and Efficient Fine-Tuning for Pre-trained Natural Language Models through Principled Regularized Optimization			2019	Transformer
2	T5-3B	97.4	Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer			2019	Transformer
3	MUPPET Roberta Large	97.4	Muppet: Massive Multi-task Representations with Pre-Finetuning			2021	
4	ALBERT	97.1	ALBERT: A Lite BERT for Self-supervised Learning of Language Representations			2019	Transformer
5	T5-11B	97.1	Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer			2019	Transformer
6	StructBERTRoBERTa ensemble	97.1	StructBERT: Incorporating Language Structures into Pre-training for Deep Language Understanding			2019	Transformer
7	XLNet (single model)	97	XLNet: Generalized Autoregressive Pretraining for Language Understanding			2019	Transformer
8	ELECTRA	96.9	ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators			2020	
9	EFL	96.9	Entailment as Few-Shot Learner			2021	Transformer
10	XLNet-Large (ensemble)	96.8	XLNet: Generalized Autoregressive Pretraining for Language Understanding			2019	Transformer
11	RoBERTa	96.7	RoBERTa: A Robustly Optimized BERT Pretraining Approach			2019	Transformer

<https://paperswithcode.com/>



# Transformers used outside of NLP

## Protein folding



AlphaFold2 (Jumper et al., 2021)

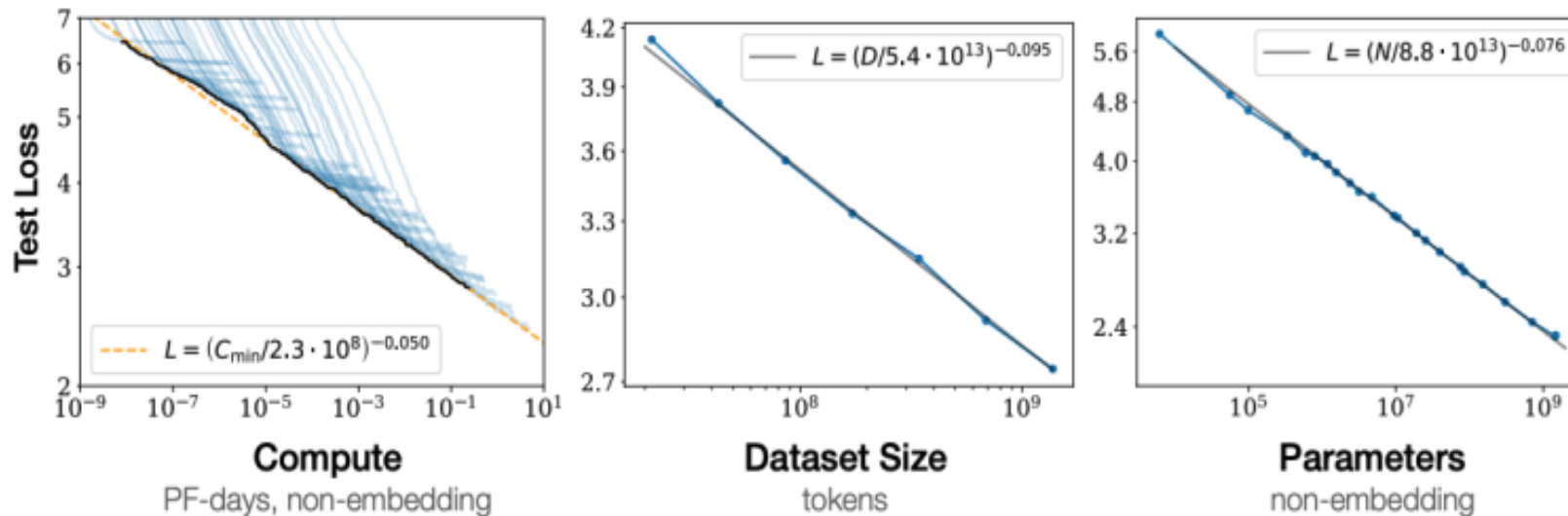
## Image Classification



Vision Transformer (ViT) outperforms ResNet-based baselines with substantially less compute (Dosovitskiy et al. 2020)

# Scaling laws

- With Transformers, language modeling performance improves smoothly as we increase **model size**, **training data**, and **computing resources**.
- This power-law relationship has been observed over multiple orders of magnitude with no sign of slowing down!



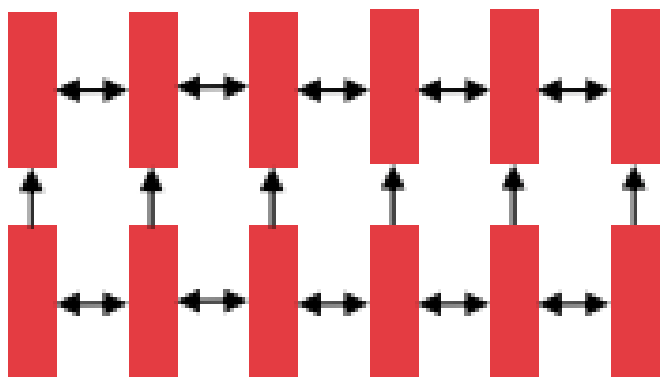
Kaplan et al., 2020, *Scaling Laws for Neural Language Models*



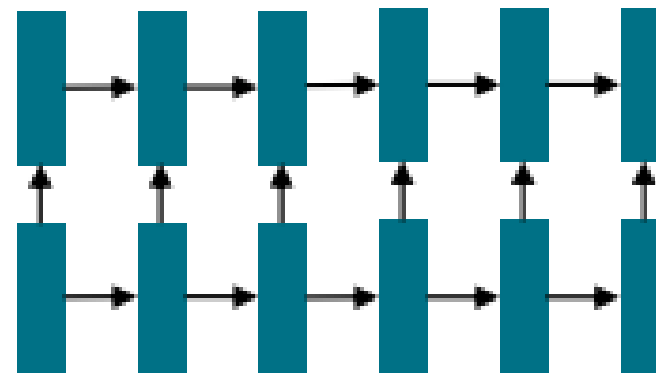
# Why self-attention?



# Recurrence in RNNs



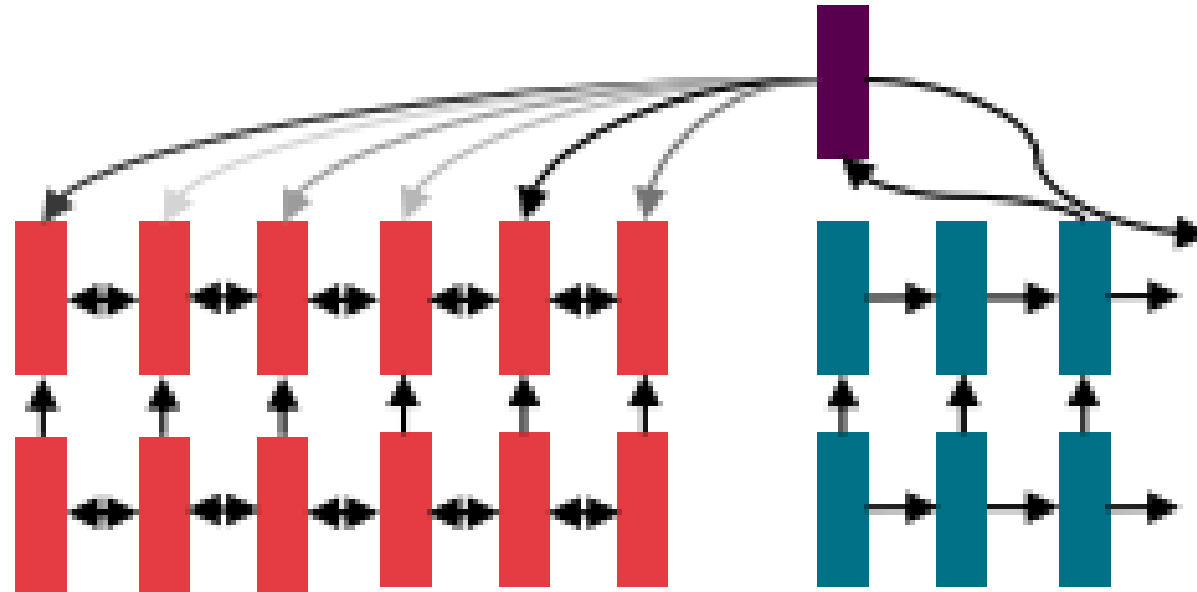
**Encoding:** Encode input sentences with bi-directional LSTM



**Decoding:** Define your outputs (parse, sentence, summary) as a sequence/label, and use LSTM to decode it.



# Sequence-to-sequence with attention



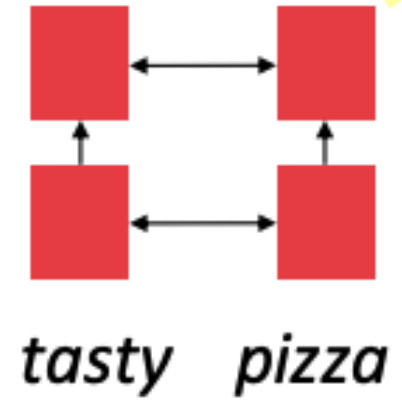
Use **attention** to allow **flexible access** to input memory



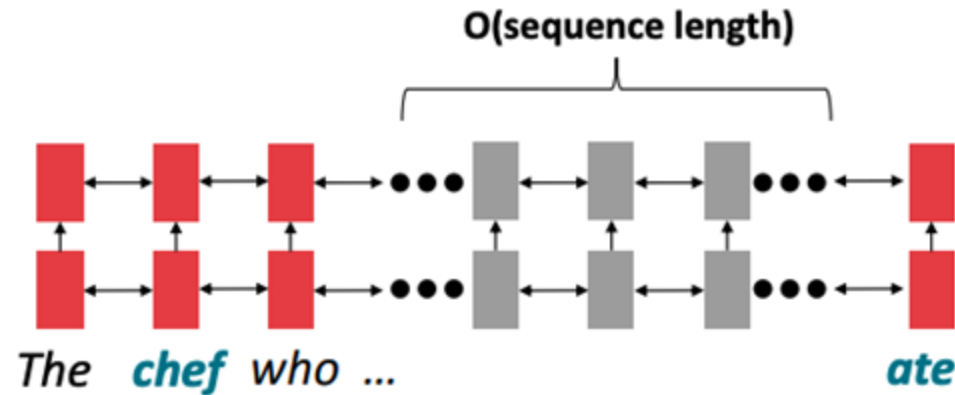
# Issues with recurrent models: **Linear interaction distance**



- ❑ Forward RNNs are unrolled “left-to-right”.
- ❑ It encodes linear locality:
  - Nearby words often affect each other’s meanings



- ❑ **Problem:** RNNs take  **$O(\text{sequence length})$  steps** for distant word pairs to interact



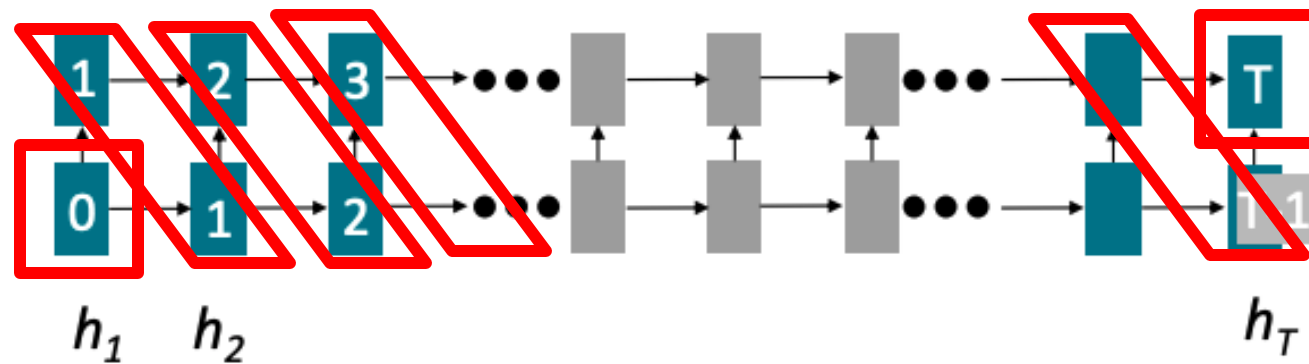
Info of *chef* has gone through  $O(\text{sequence length})$  many layers!





# Issues with recurrent models: **Lack of parallelizability**

- ❑ Forward and backward passes have  **$O(\text{seq length})$  un-parallelizable** operations
  - GPUs (and TPUs) can perform many independent computations at once! But future RNN hidden states can't be computed fully before past RNN hidden states have been computed
  - Particularly problematic as sequence length increases, as we can no longer batch many examples together due to memory limitations



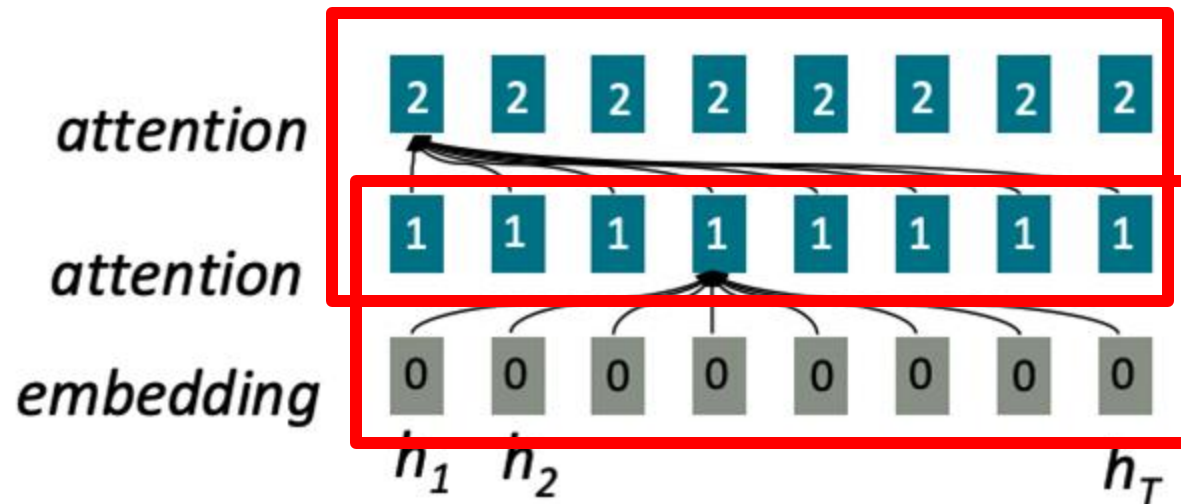
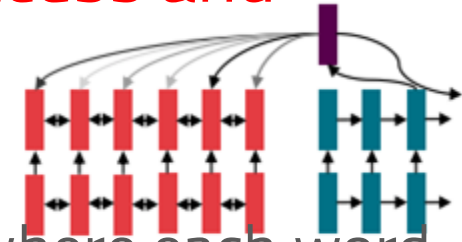
Numbers indicate min # of steps before a state can be computed



# If not recurrence, then what? How about (self) attention?

□ **Attention** treats each word's representation as a **query** to **access and incorporate information** from a **set of values**.

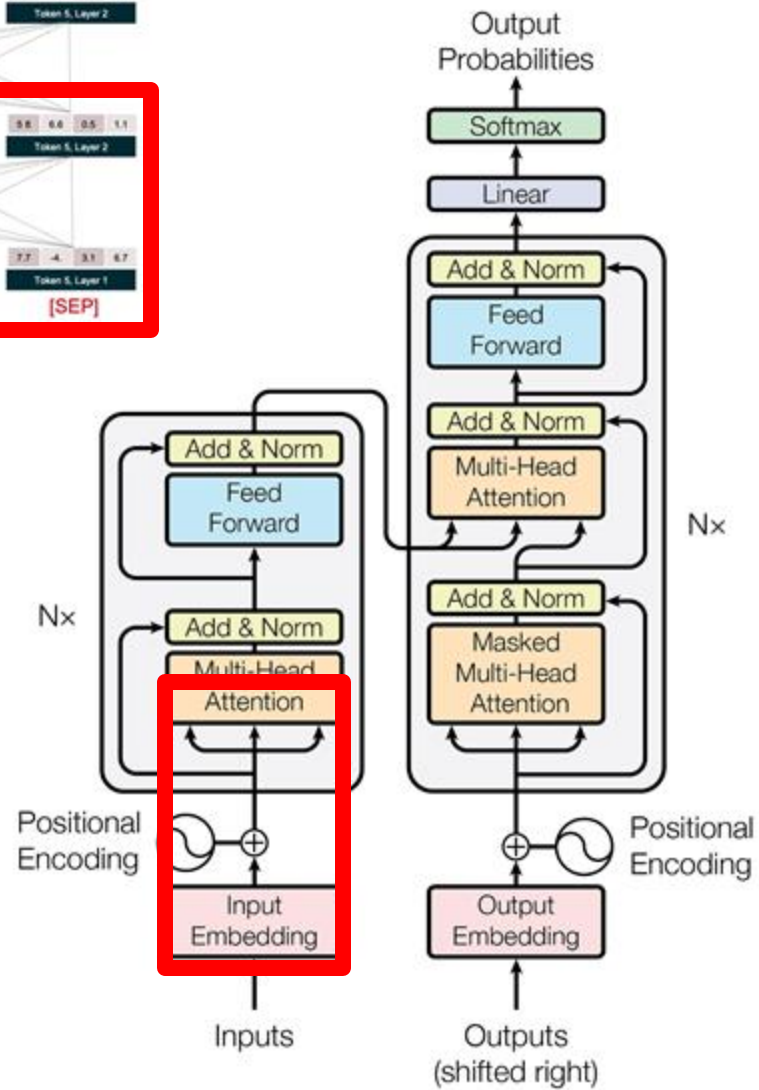
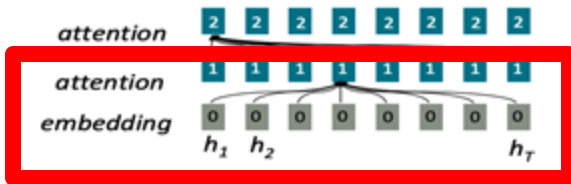
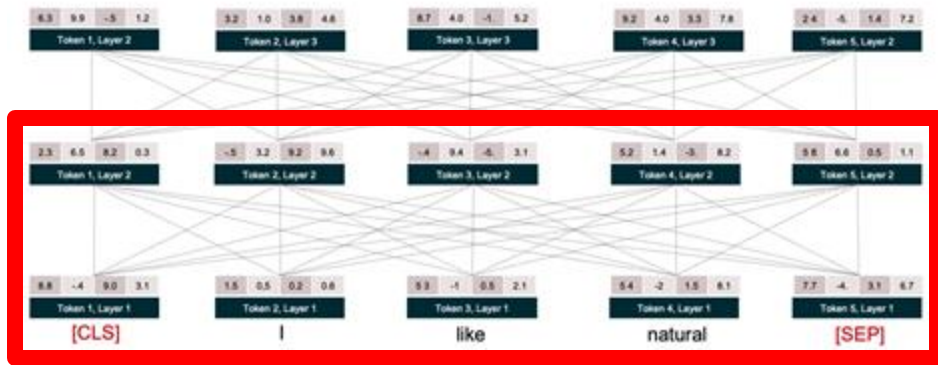
- We saw attention from the **decoder** to the **encoder**;
- Self-attention is **encoder-encoder** (or **decoder-decoder**) attention where each word attends to each other word within the input (or output).



All words attend to all words in previous layer; most arrows are omitted

**$O(\text{seq length}) O(\text{Layers})$**





Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Output Embedding

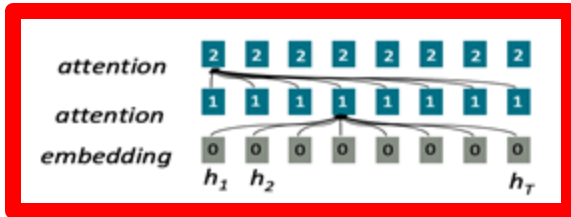
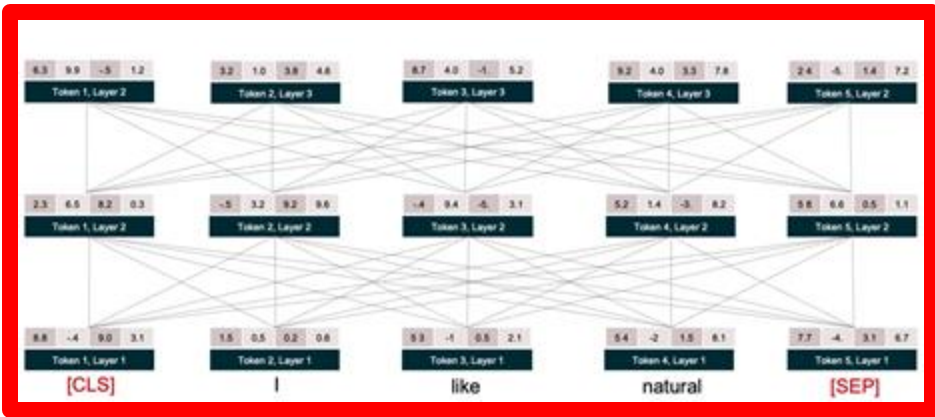
Positional Encoding

Positional Encoding

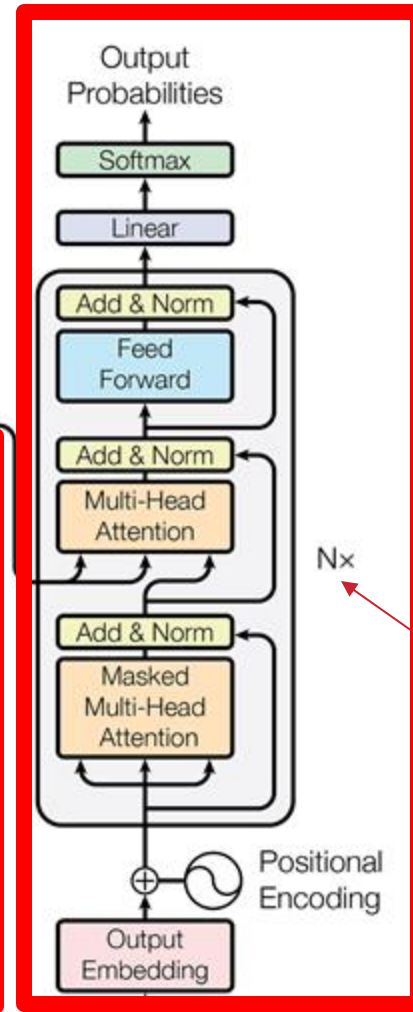
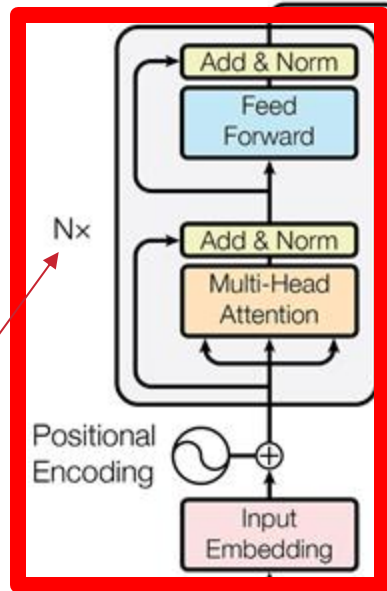
Inputs

Outputs (shifted right)





Repeat N times (number of layers)

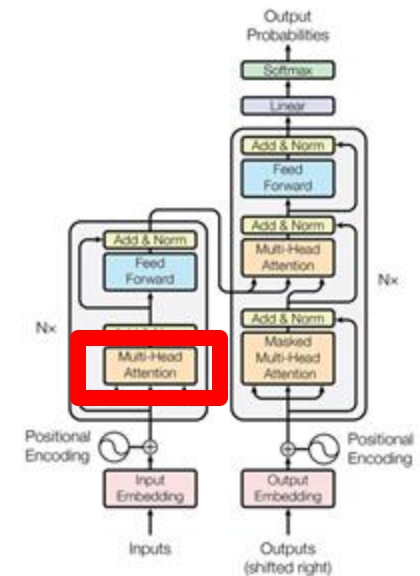
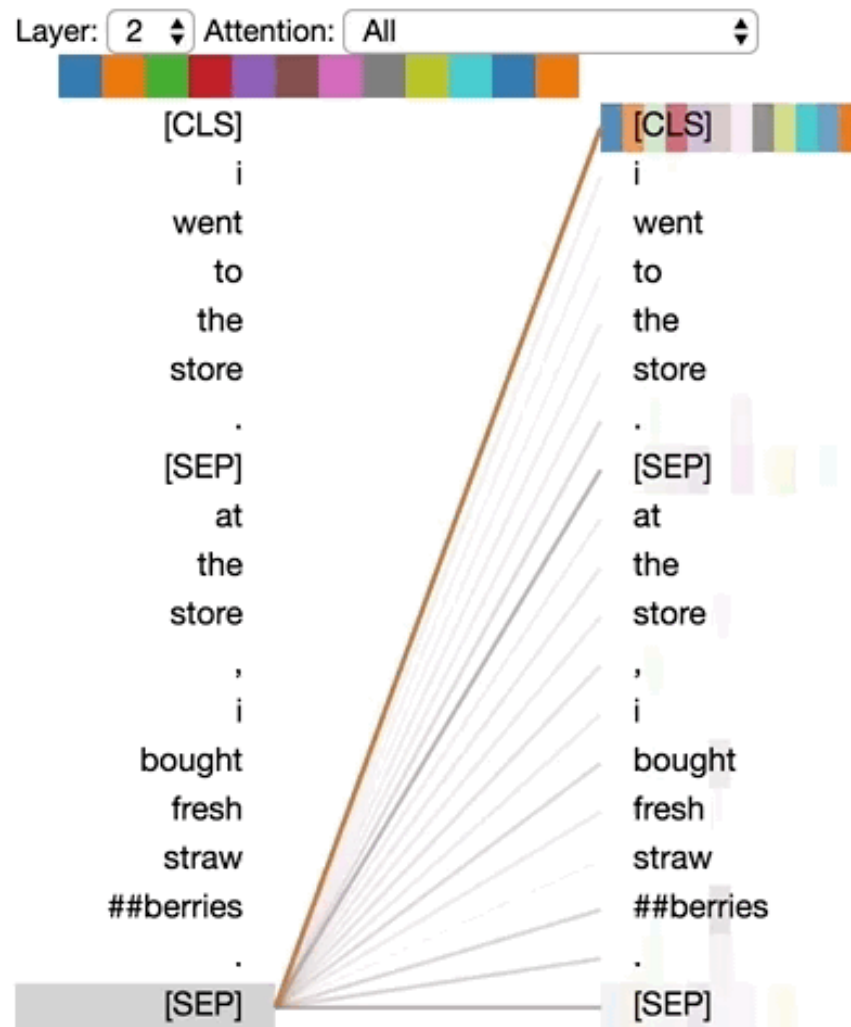
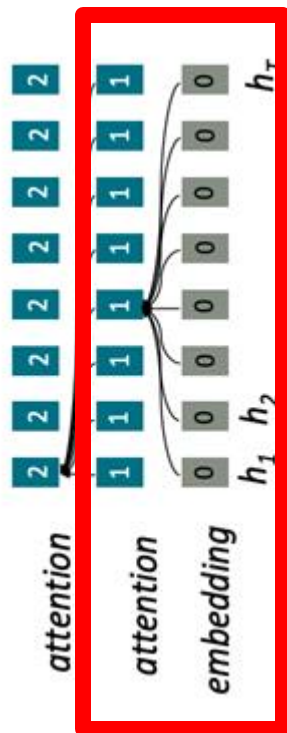


Decoder

Repeat N times (number of layers)



"I went to the store. At the store, I bought fresh strawberries."

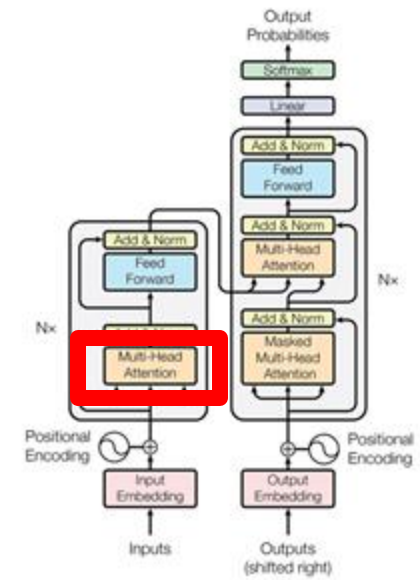
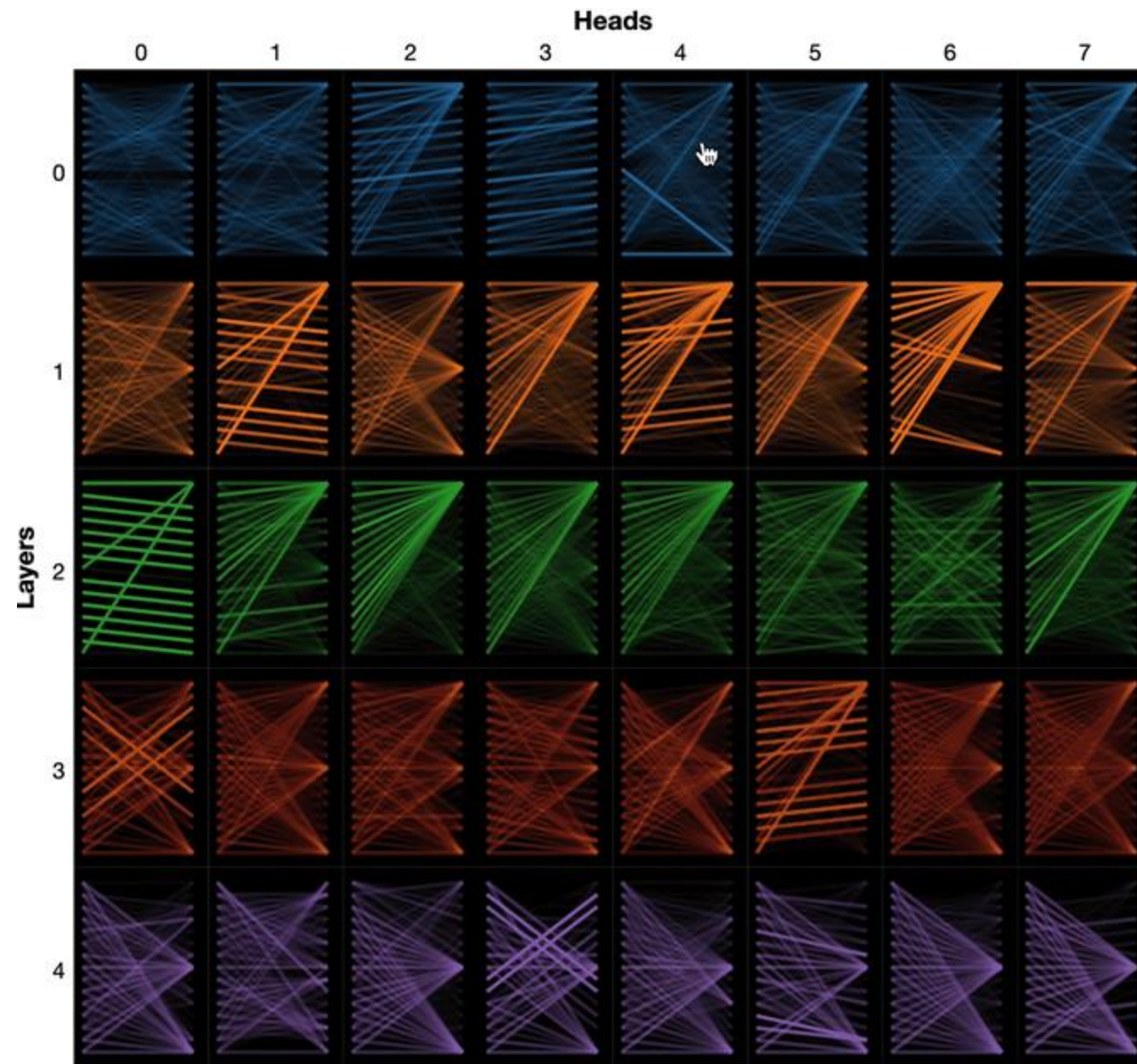


<https://github.com/jessevig/bertviz>

[https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\\_t2t.ipynb](https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb)







<https://github.com/jessevig/detwiz>

[https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\\_t2t.ipynb](https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb)

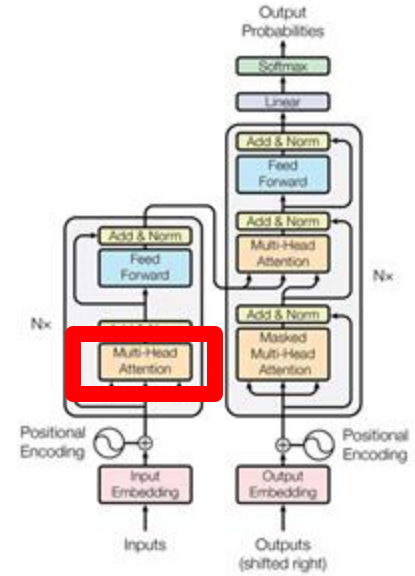
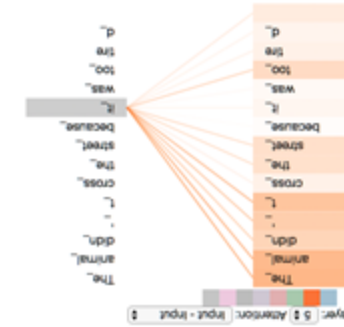


# Encoder: Self-Attention

Recap: Attention as a **query** to access and incorporate information from a set of **values**.

Let's think of attention as a "fuzzy" or approximate hashtable:

- To look up a **value**, we compare a **query** against **keys** in a table.
- In a hashtable
  - Each **query** (hash) maps to exactly one **key-value** pair.
- In (self-)attention:
  - Each **query** (token in current layer) matches each **key** to varying degrees.
  - We return a sum of **values** (token in previous layer) weighted by the query-key match (attention score).



# Encoder: Self-Attention

In (self-)attention: Each **query** (token in current layer) matches each **key** to varying degrees. We return a sum of **values** (token in previous layer) weighted by the query-key match (**attention score**).

query (token in current layer)

$$x_1 a_1 + x_2 a_2 + x_3 a_3$$

-5	3.2	9.2	9.6
Token 1, Layer 2			

$$a_1 = 0.4$$

$$a_2 = 0.64$$

$$a_3 = 0.02$$

query-key match (attention score)

1.5	0.5	0.2	0.6
Token 1, Layer 1			

I

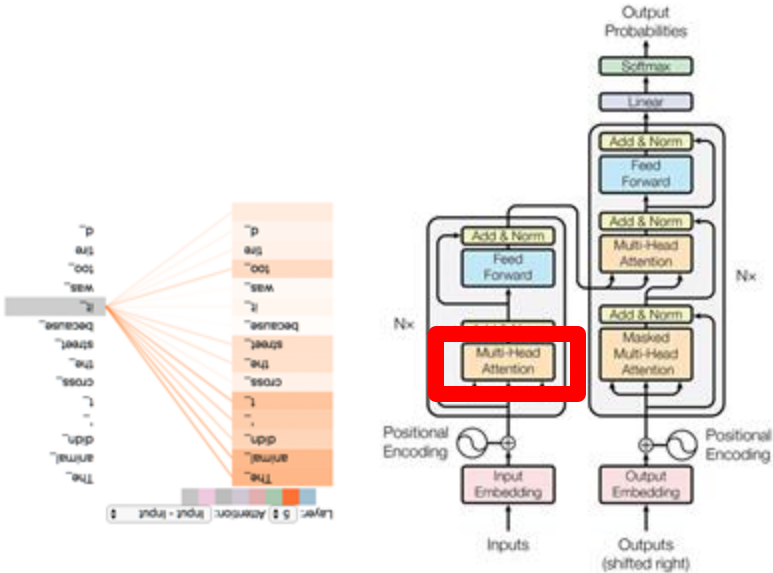
5.3	-1	0.5	2.1
Token 2, Layer 1			

like

5.4	-2	1.5	8.1
Token 3, Layer 1			

natural

values (token in previous layer)





# Recipe for Self-Attention in the Transformer Encoder

Model parameters to learn (randomly initialized)

- Step 1: For each word  $x_i$ , calculate its **query**, **key**, and **value**.

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$

- Step 2: Calculate attention score between **query** and **keys**.

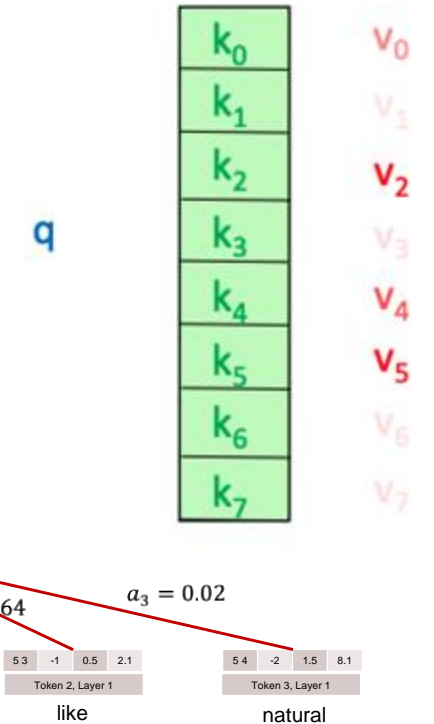
$$e_{ij} = q_i \cdot k_j$$

- Step 3: Take the softmax to normalize attention scores.

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

- Step 4: Take a weighted sum of **values**.

$$\text{Output}_i = \sum_j \alpha_{ij} v_j$$



# Recipe for (Vectorized) Self-Attention in the Transformer Encoder

- Step 1: For each word  $x$ , calculate its query, key, and value.

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$

- Step 2: Calculate attention score between query and keys.

$$E = QK^T$$

- Step 3: Take the softmax to normalize attention scores.

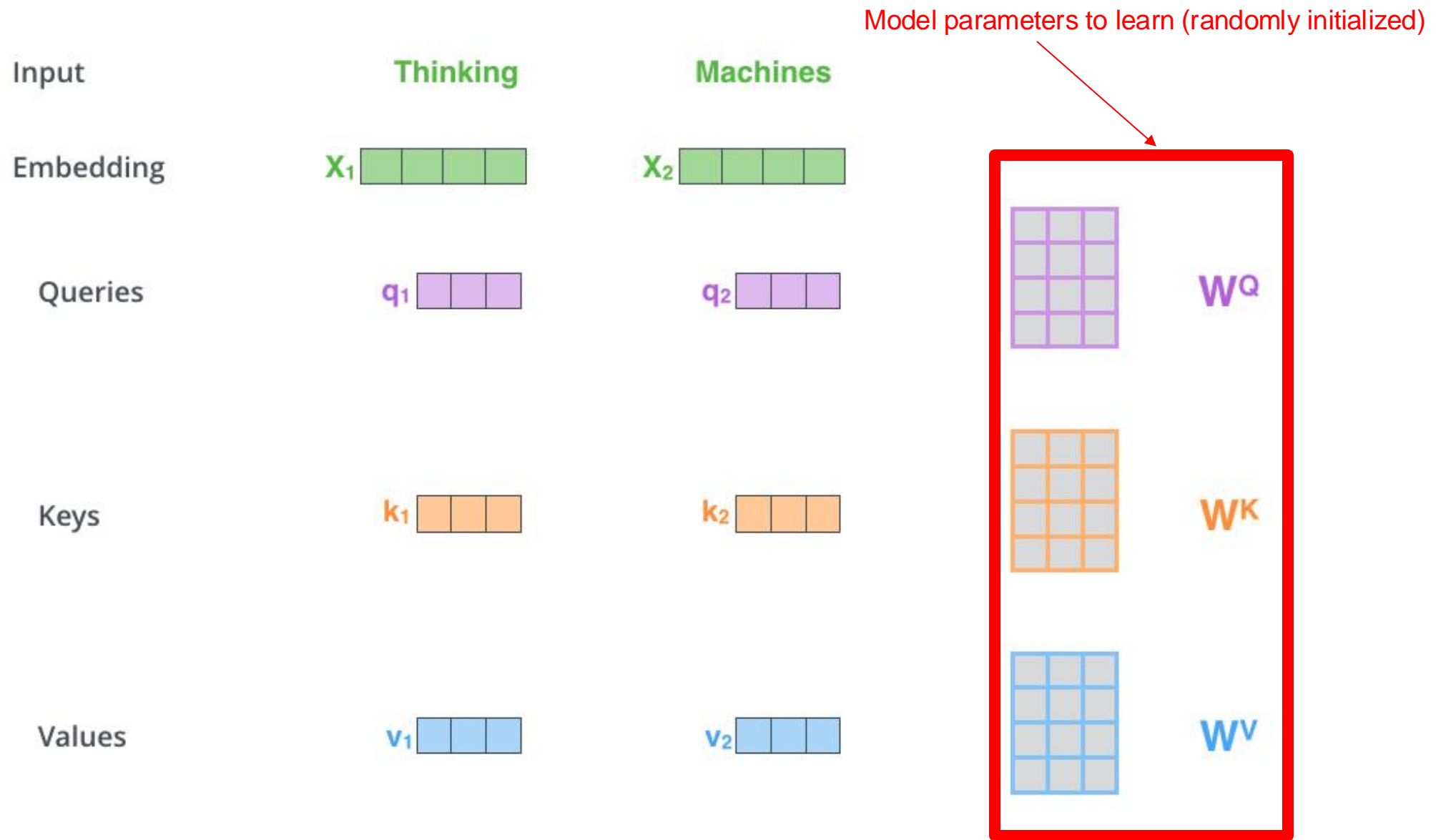
$$A = \text{softmax}(E)$$

- Step 4: Take a weighted sum of values.

$$\text{Output} = AV$$

$$\text{Output} = \text{softmax}(QK^T)V$$



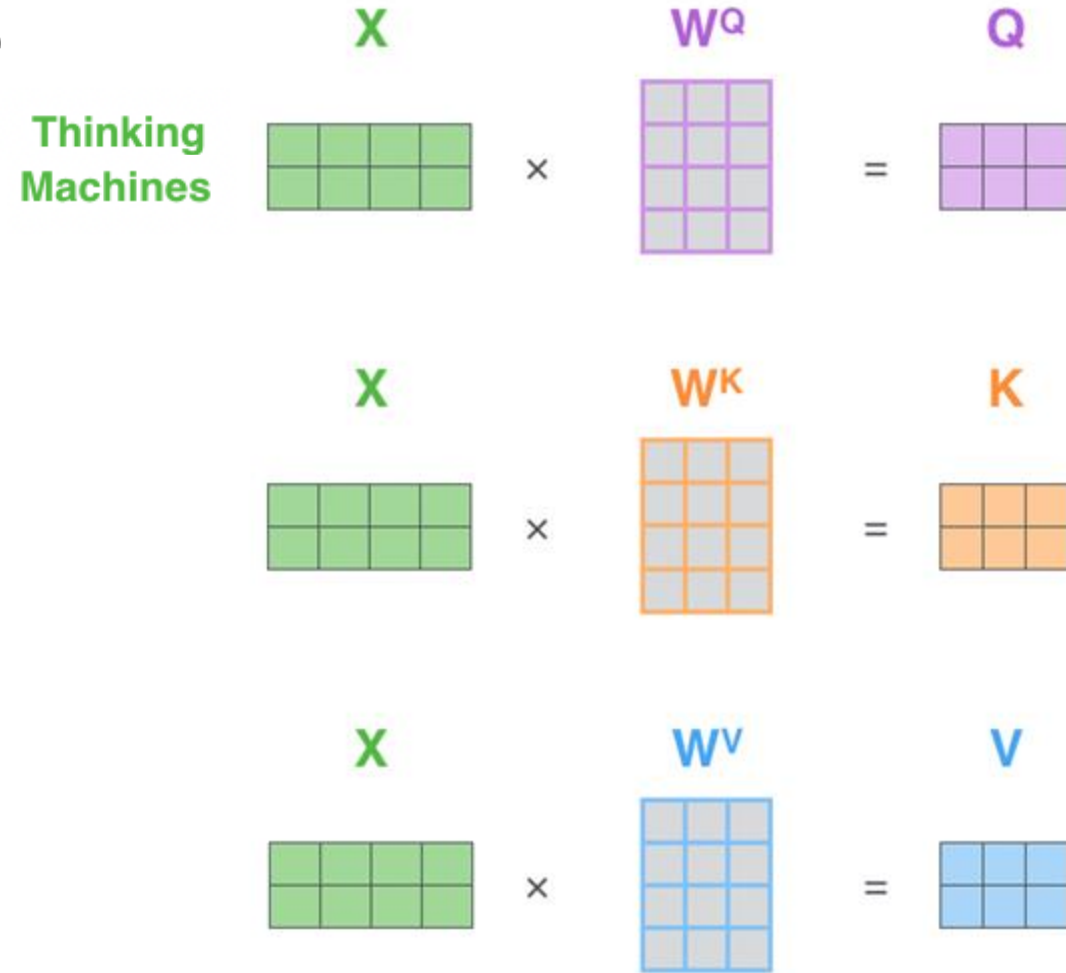


<https://jalammar.github.io/illustrated-transformer/>



□ Step 1: For each word, calculate its query, key, and value.

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$



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



- Step 2: Calculate attention score between query and keys.

$$E = QK^T$$

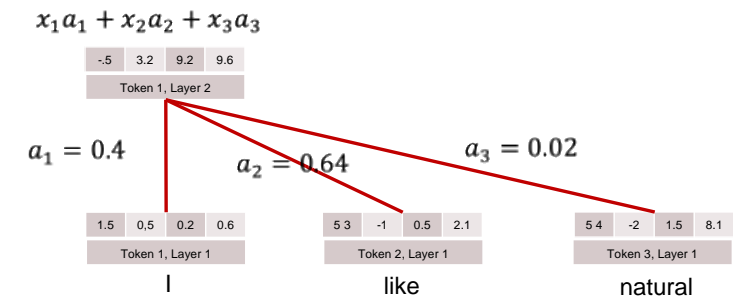
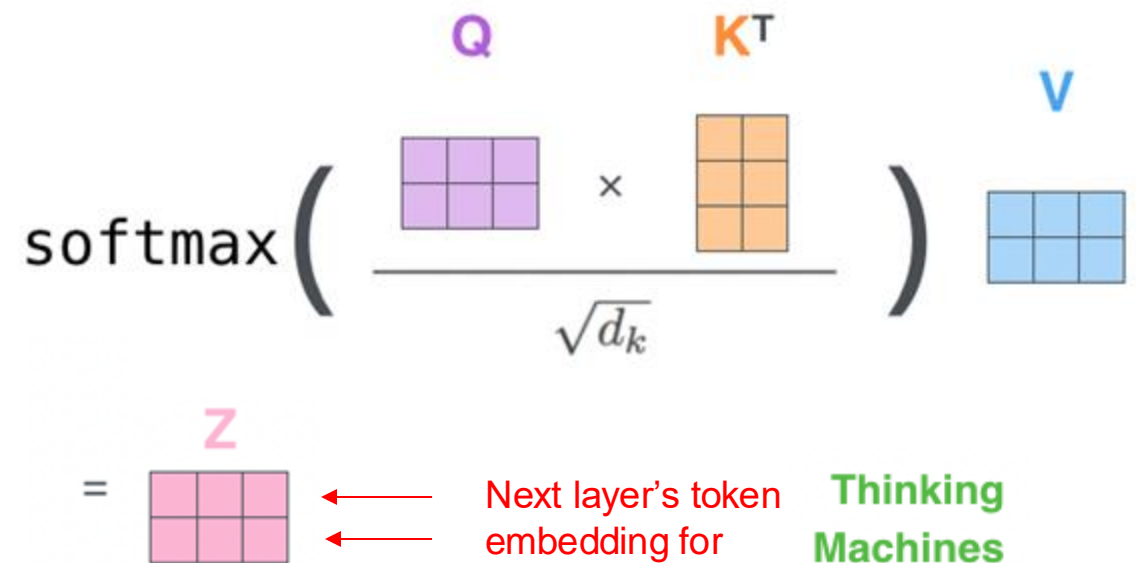
- Step 3: Take the softmax to normalize attention scores.

$$A = \text{softmax}(E)$$

- Step 4: Take a weighted sum of values.

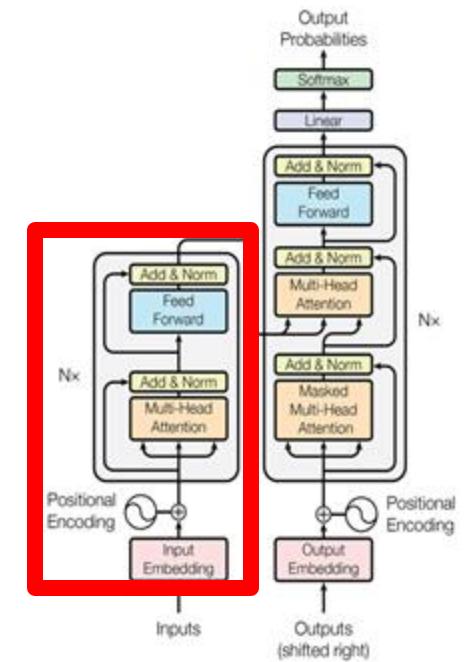
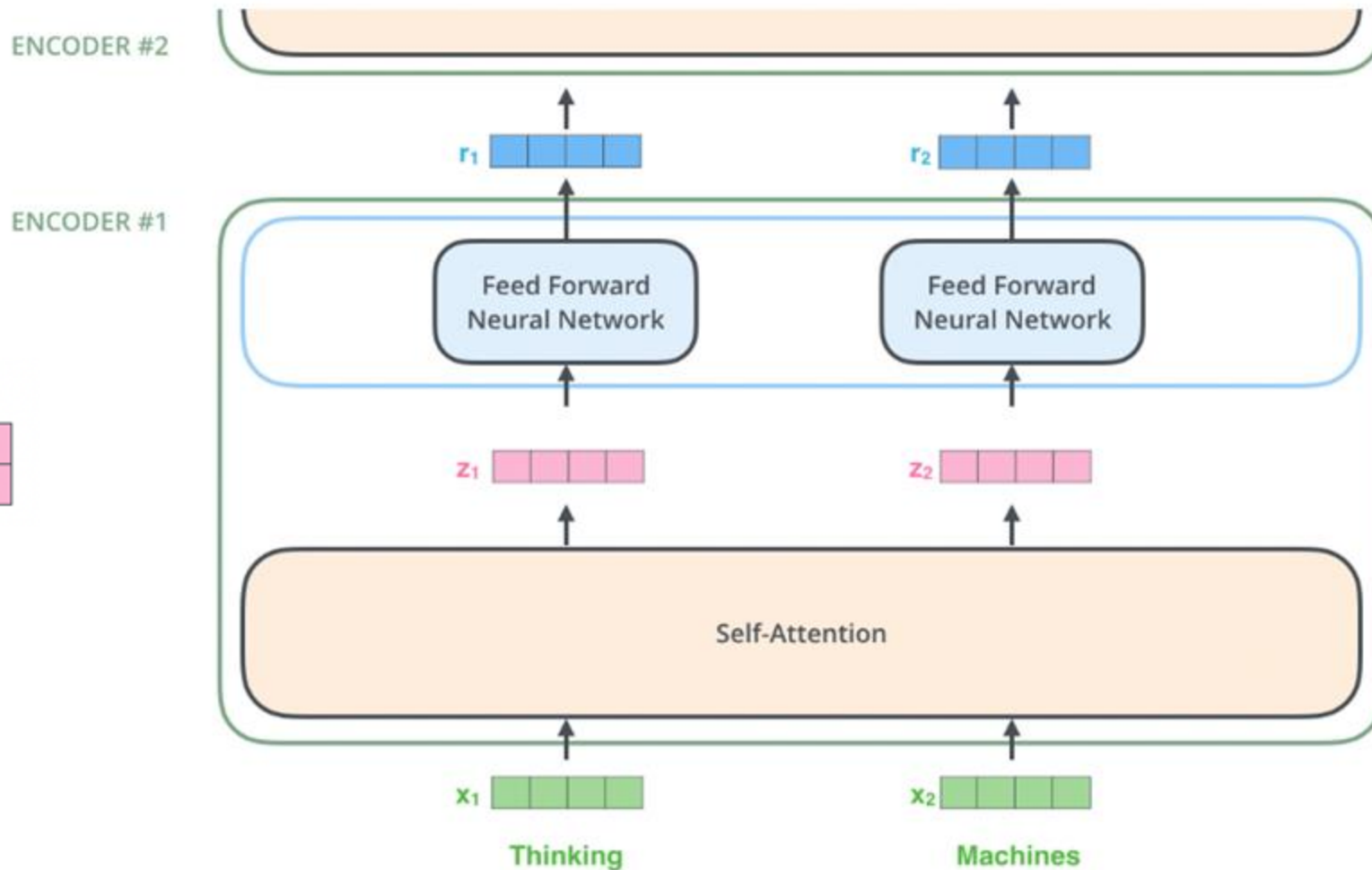
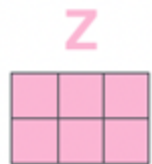
$$\text{Output} = AV$$

$$\text{Output} = \text{softmax}(QK^T)V$$



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





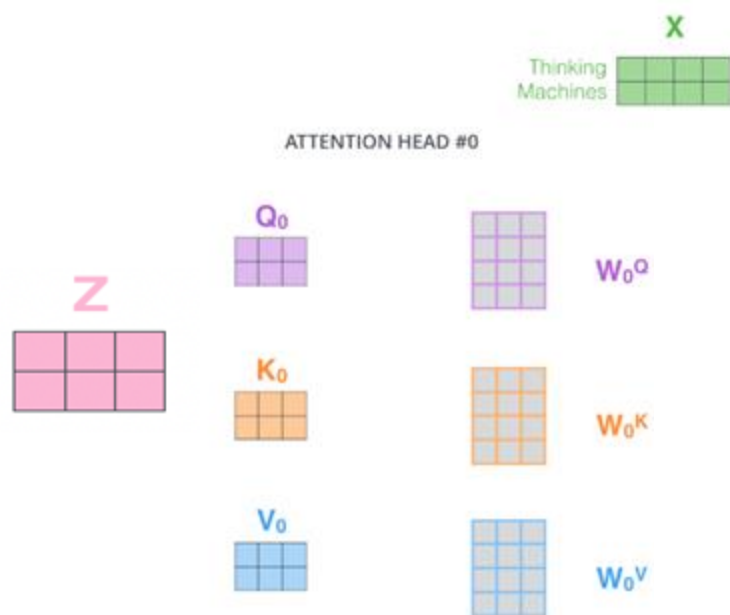
<https://jalammar.github.io/illustrated-transformer/>



# Multi-headed self-attention



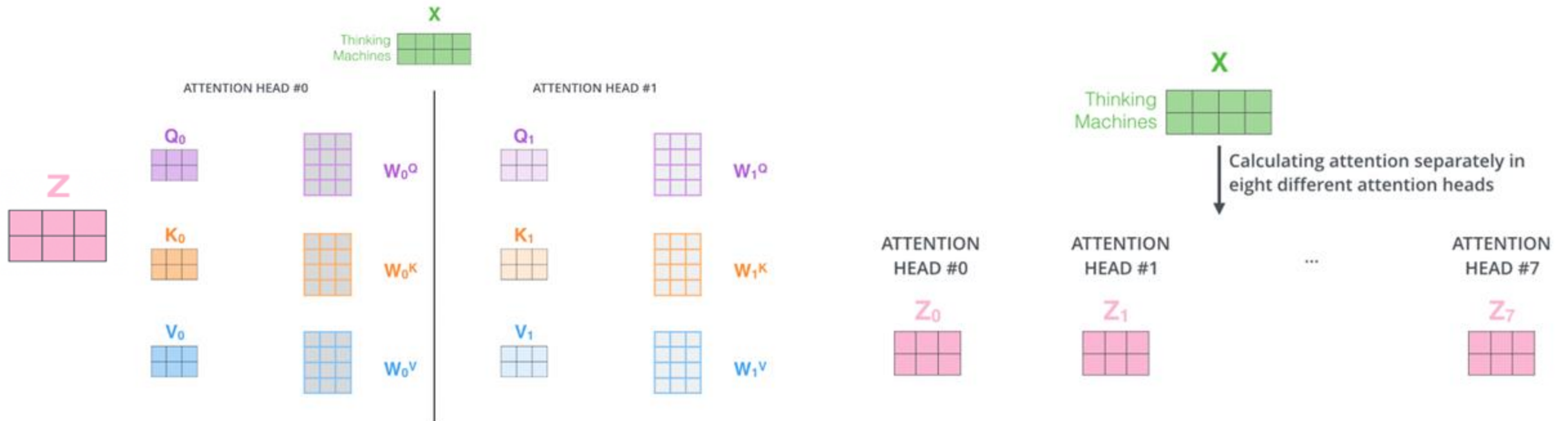
- It gives the attention layer multiple “representation subspaces”
- Multiple sets of Query/Key/Value weight matrices (Transformer uses eight attention heads, so we end up with eight sets for each encoder/decoder). Each of these sets is randomly initialized.



# Multi-headed self-attention



- It gives the attention layer multiple “representation subspaces”
- Multiple sets of Query/Key/Value weight matrices (Transformer uses eight attention heads, so we end up with eight sets for each encoder/decoder). Each of these sets is randomly initialized.





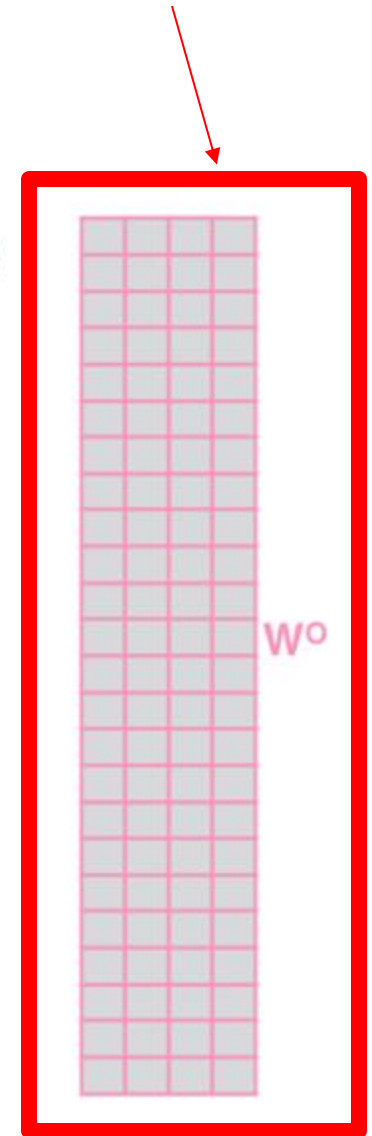
# Condensing multi-head attentions into a single matrix

1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^O$  that was trained jointly with the model

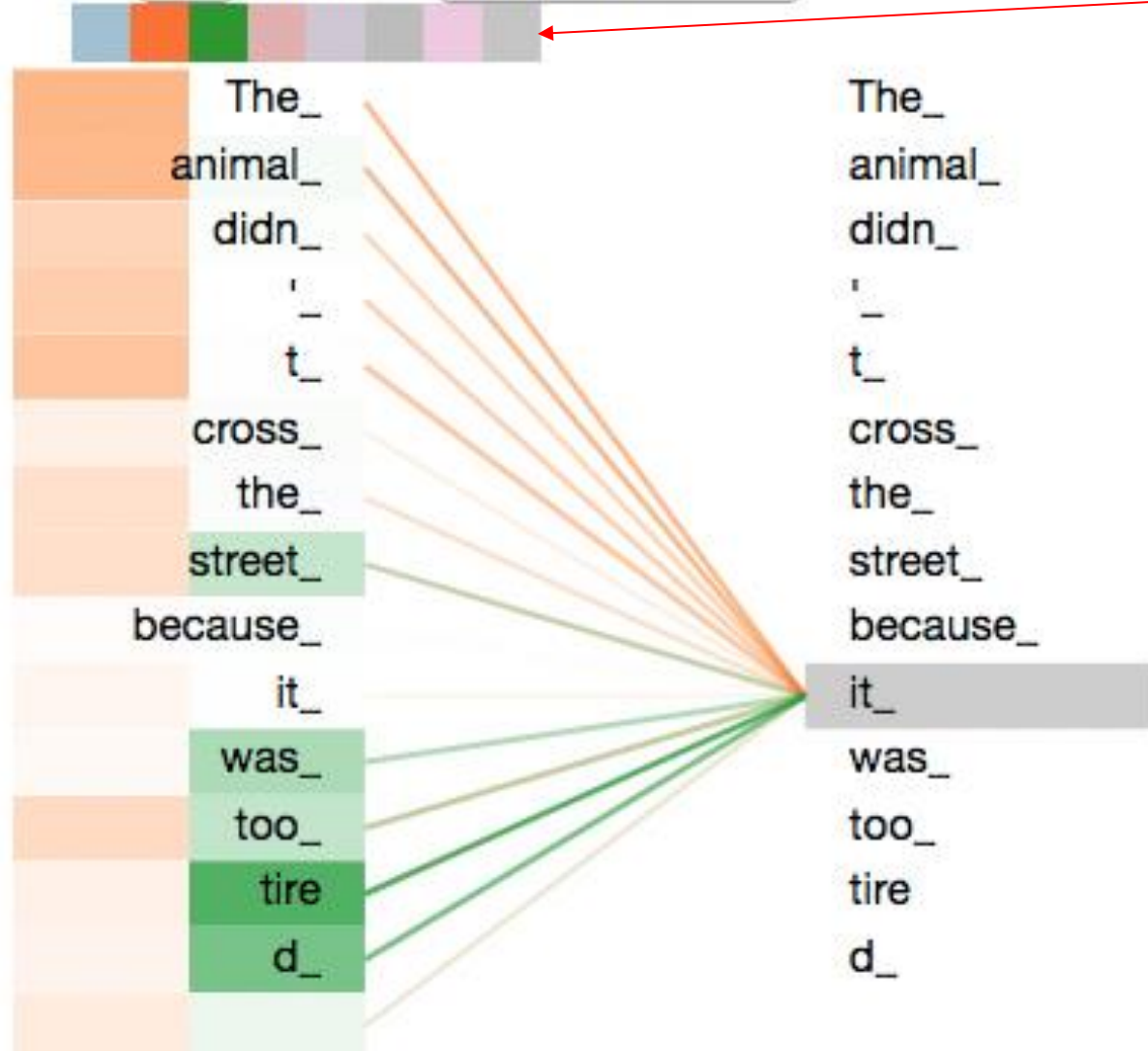
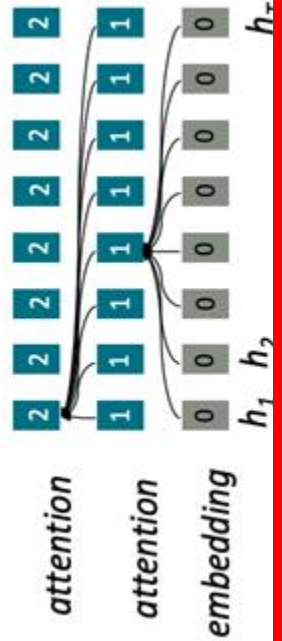
x



3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN



Layer: 5 Attention: Input - Input

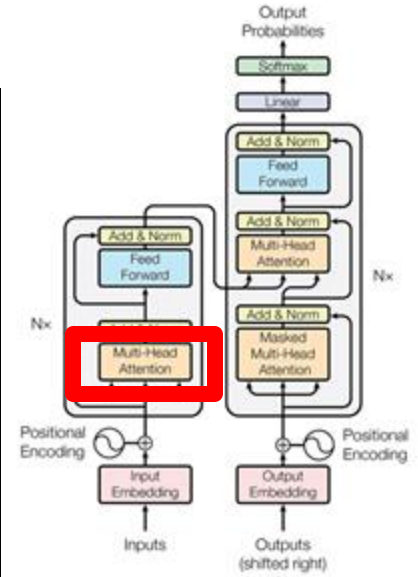


<https://github.com/jessevig/bertviz>

[https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\\_t2t.ipynb](https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb)



Layer: 4 Head: 3 Attention: All

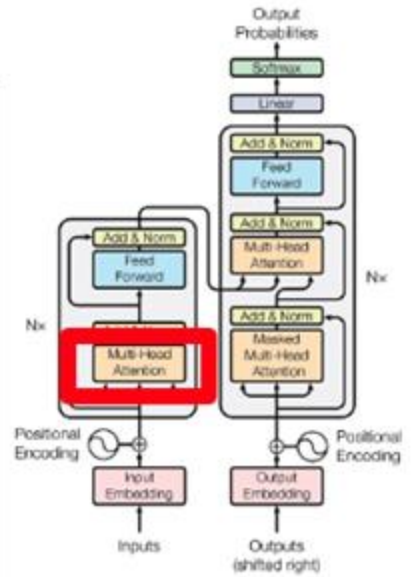
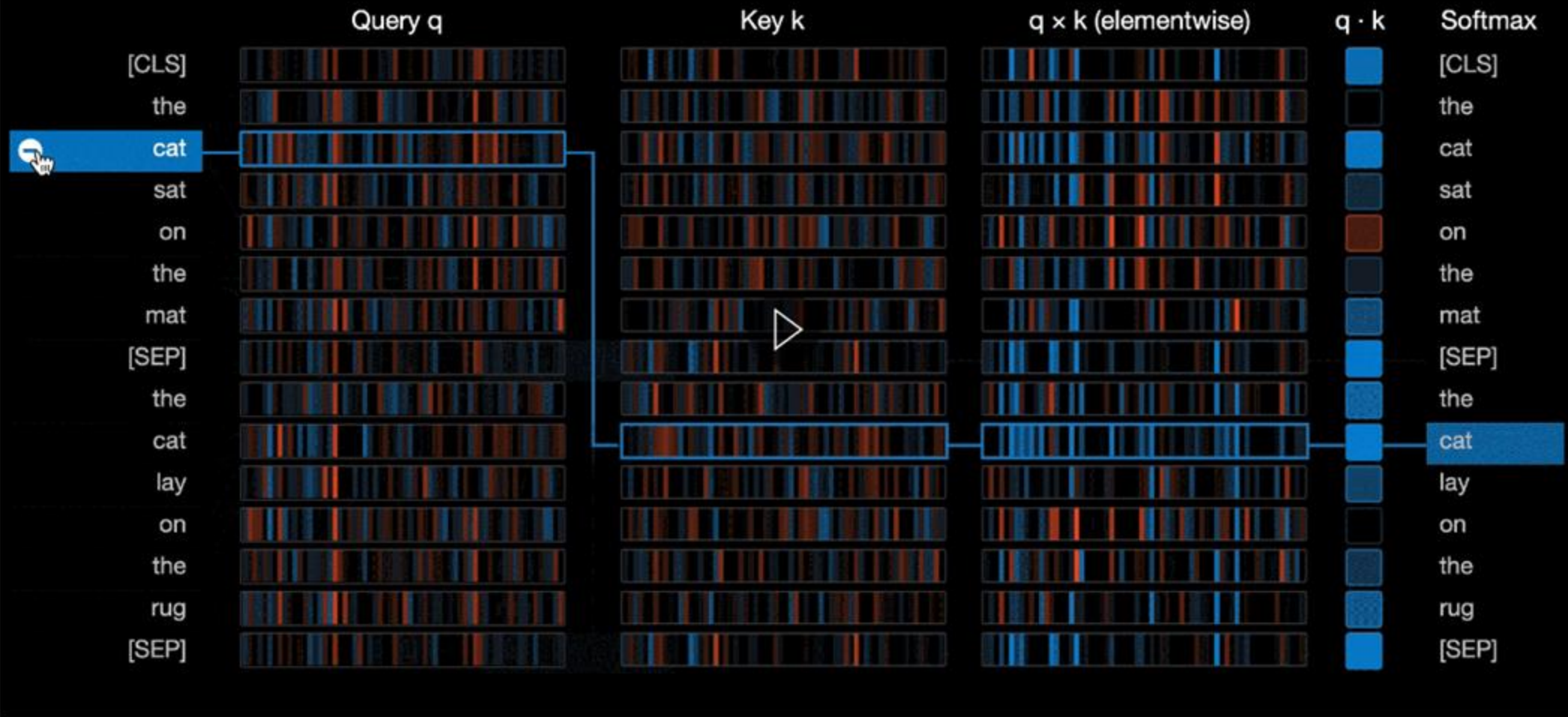


<https://github.com/jessevig/bertviz>

[https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\\_t2t.ipynb](https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb)



Layer: 4 Head: 3 Attention: All



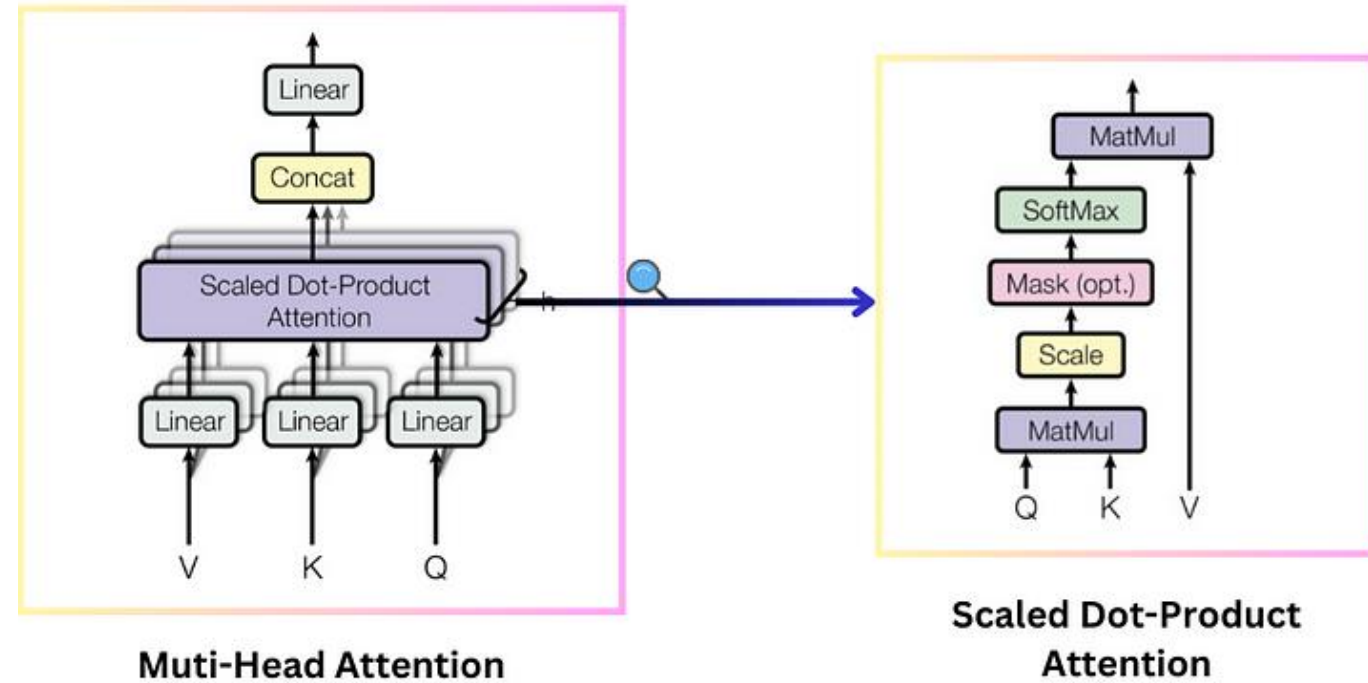
<https://github.com/lessevig/bertviz>

[https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\\_t2t.ipynb](https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb)

# Recap

- ❑ Pass our input through the  $W_v$ ,  $W_k$ ,  $W_q$  matrices for each head (corresponding to the 'Linear' boxes at right)
- ❑ Perform Scaled dot product attention for each head
- ❑ Concatenate the results for each head
- ❑ Use linear layer to project to original output dimension

Multi-Head Attention



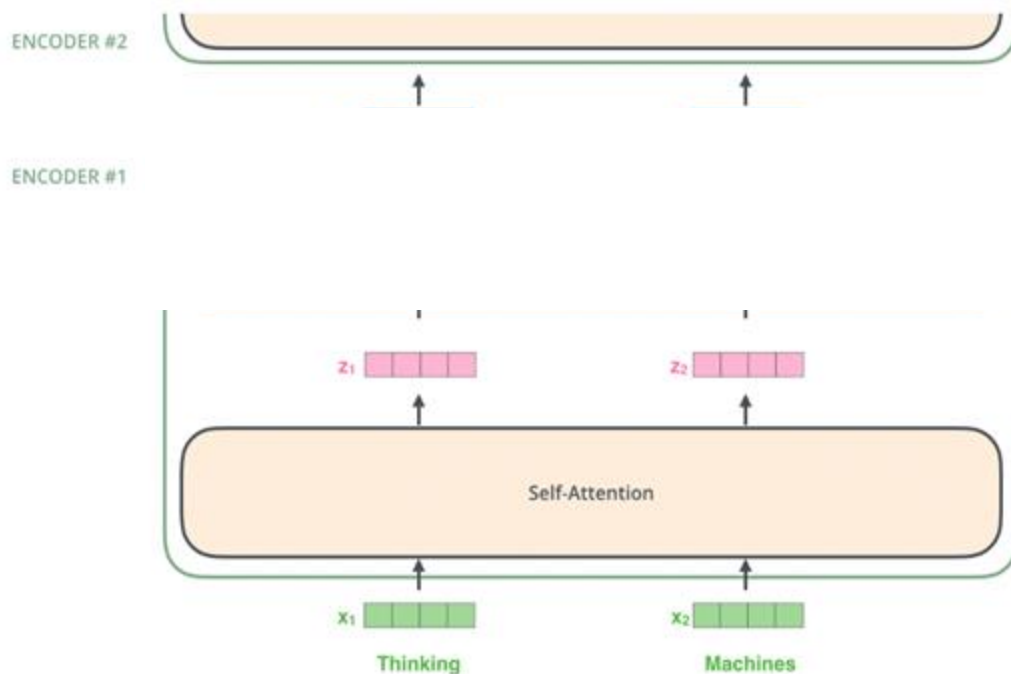
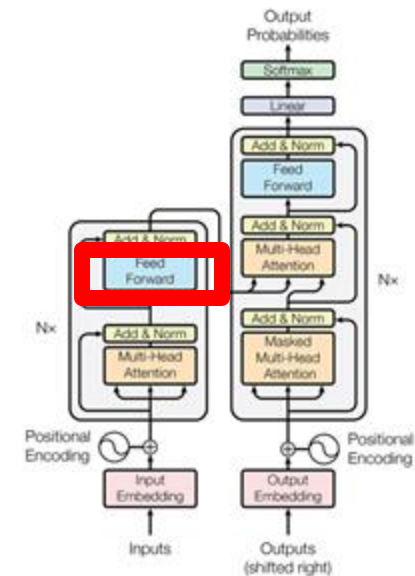
Other tricks than attention?





# But attention isn't quite all you need!

- ❑ **Problem:** Since there are no element-wise non-linearities, self-attention is simply performing a re-averaging of the value vectors.
- ❑ **Easy fix:** Apply a feedforward layer to the output of attention, providing non-linear activation (and additional expressive power).



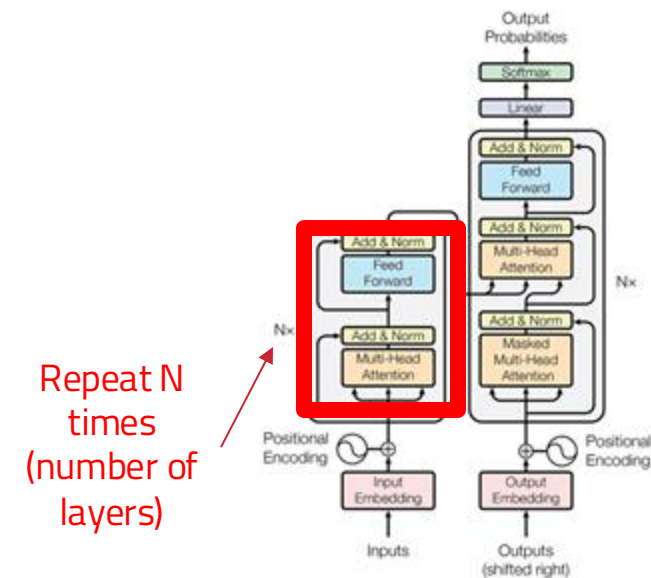
Equation for Feed-Forward layer

$$\begin{aligned} m_i &= MLP(\text{output}_i) \\ &= W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2 \end{aligned}$$



# Stacking deep neural nets

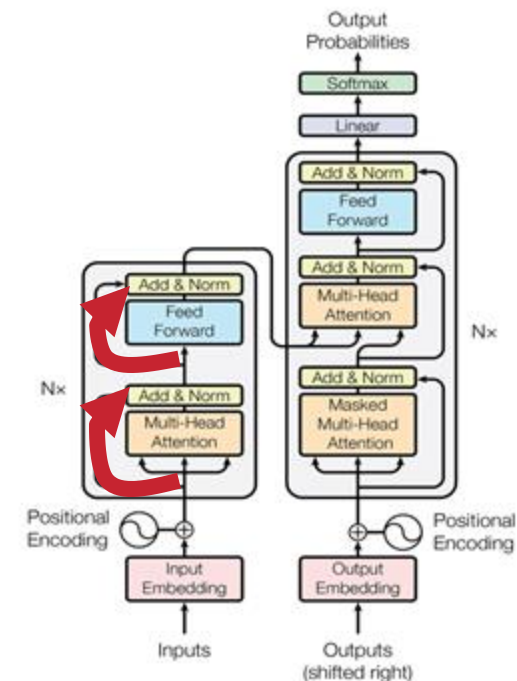
- ❑ Training trick #1: Residual Connections
- ❑ Training trick #2: LayerNorm
- ❑ Training trick #3: Scaled Dot Product Attention





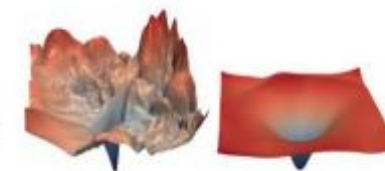
# Trick #1: Residual Connections [He et al., 2016]

- ❑ Residual connections are a simple but powerful technique from computer vision.
- ❑ Similar to additive connection in LSTM
- ❑ Directly passing "raw" embeddings to the next layer prevents the network from "forgetting" or distorting important information as it is processed by many layers.



$$x_{\ell} = F(x_{\ell-1}) + x_{\ell-1}$$

Residual connections are also thought to smooth the loss landscape and make training easier!



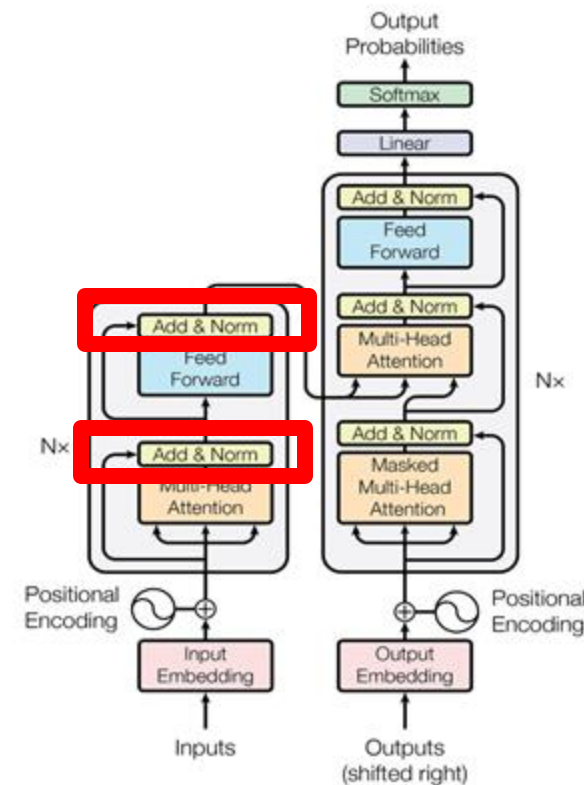
[no residuals] [residuals]

[Loss landscape visualization, Li et al., 2018, on a ResNet]



## Trick #2: Layer Normalization [Ba et al., 2016]

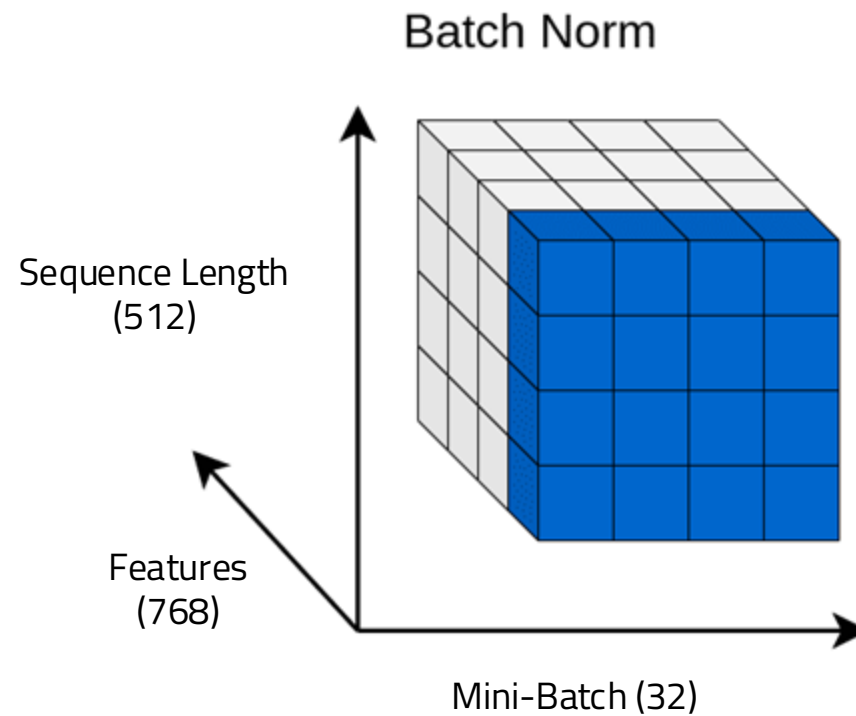
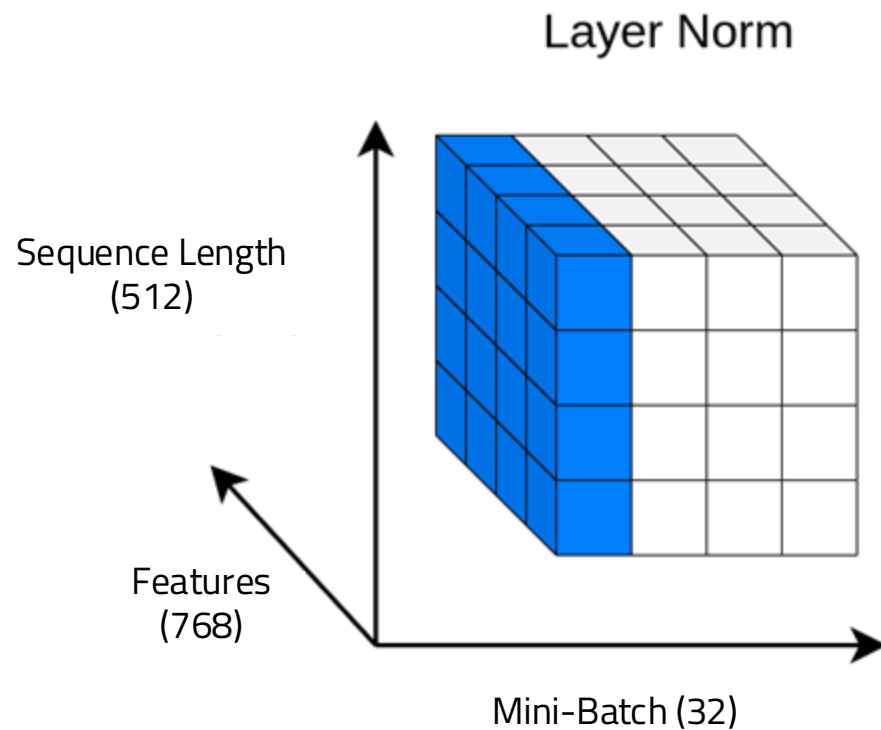
- ❑ **Problem:** Difficult to train the parameters of a given layer because its input from the layer beneath keeps shifting.
- ❑ **Solution:** Reduce uninformative variation by normalizing to zero mean and standard deviation of one within each layer.

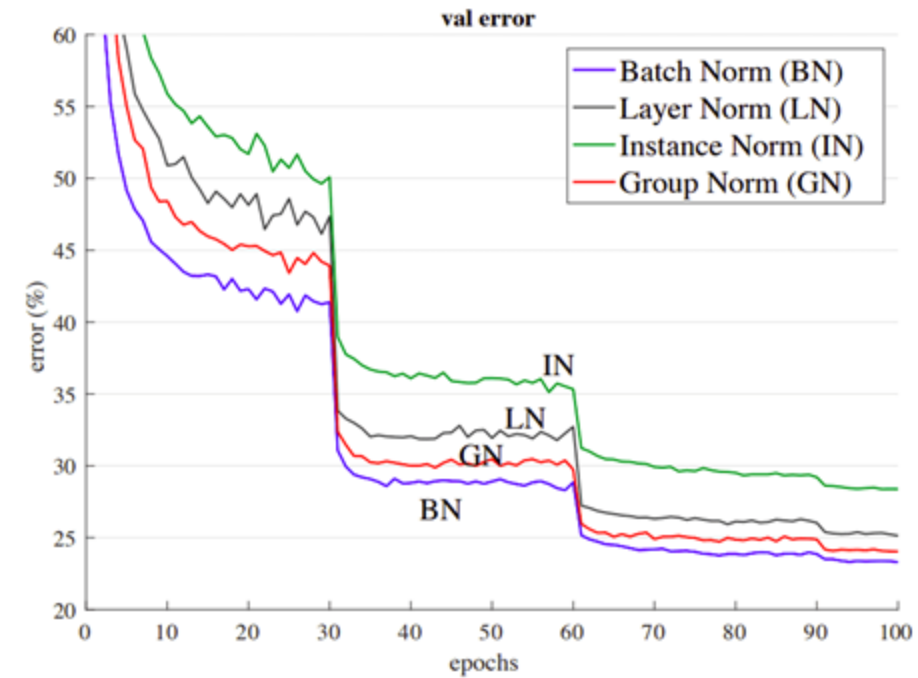
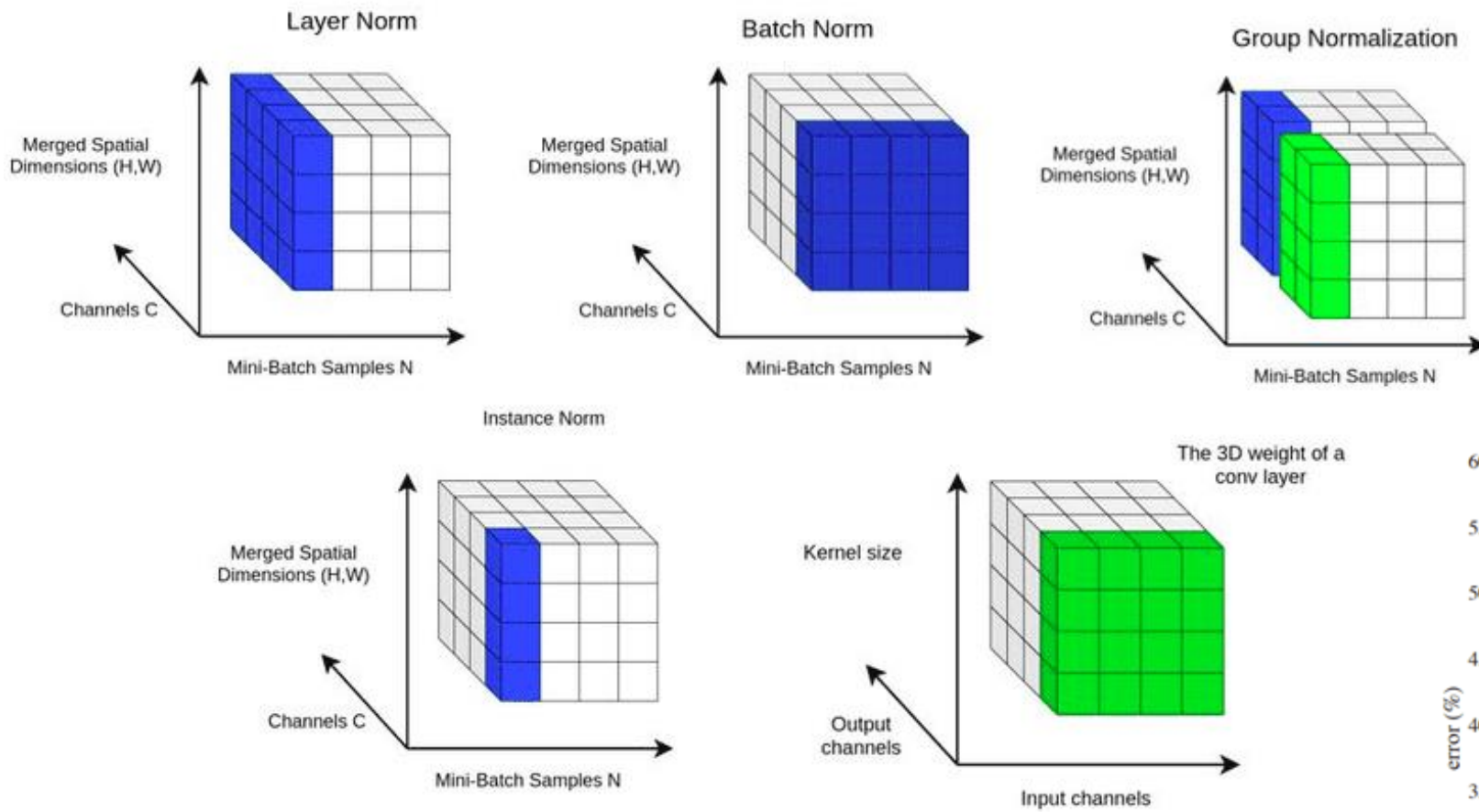


$$\text{Mean: } \mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \text{Standard Deviation: } \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$x^{l'} = \frac{x^l - \mu^l}{\sigma^l + \epsilon}$$

# Layer norm vs Batch norm



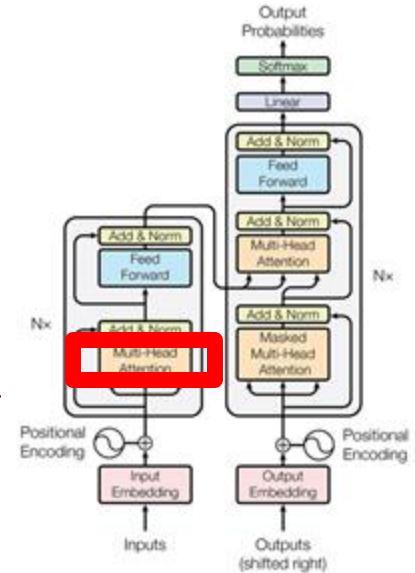
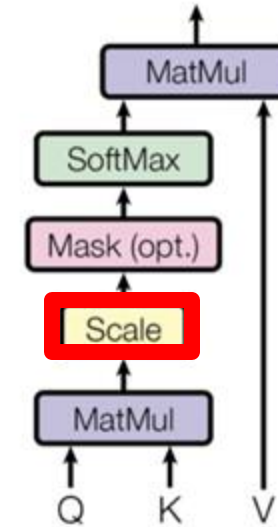


<https://theaisummer.com/normalization>

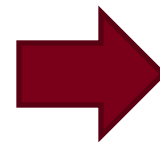


# Trick #3: Scaled Dot Product Attention

- ❑ After LayerNorm, the mean and var of vector elements is 0 and 1, respectively.
- ❑ But, the dot product still tends to take on extreme values, as its variance scales with dimensionality  $d_k$



$$Output = softmax(QK^T)V$$



Updated Self-Attention Equation

$$Output = softmax\left(QK^T / \sqrt{d_k}\right)V$$



# Representing The Order of The Sequence Using Positional Encoding

- Since self-attention doesn't build in order information, we need to encode the order of the sentence in our keys, queries, and values.
- Consider representing **each sequence index** as a **vector**

$p_i \in \mathbb{R}^d$ , for  $i \in \{1, 2, \dots, T\}$  are position vectors

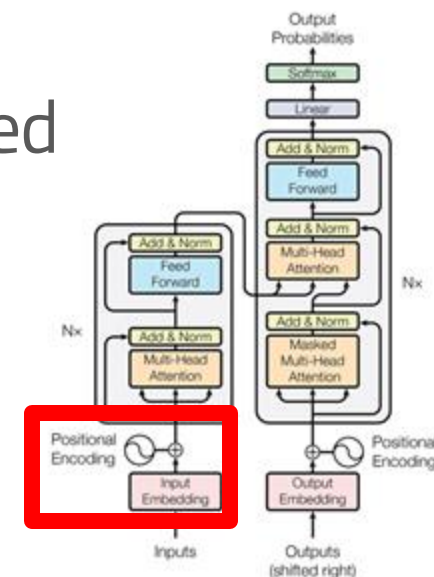
- Easy to incorporate this info into our self-attention block: just add the  $p_i$  to our inputs!

$$v_i = \tilde{v}_i + p_i$$

$$q_i = \tilde{q}_i + p_i$$

$$k_i = \tilde{k}_i + p_i$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add...



# Position representation vectors through sinusoids

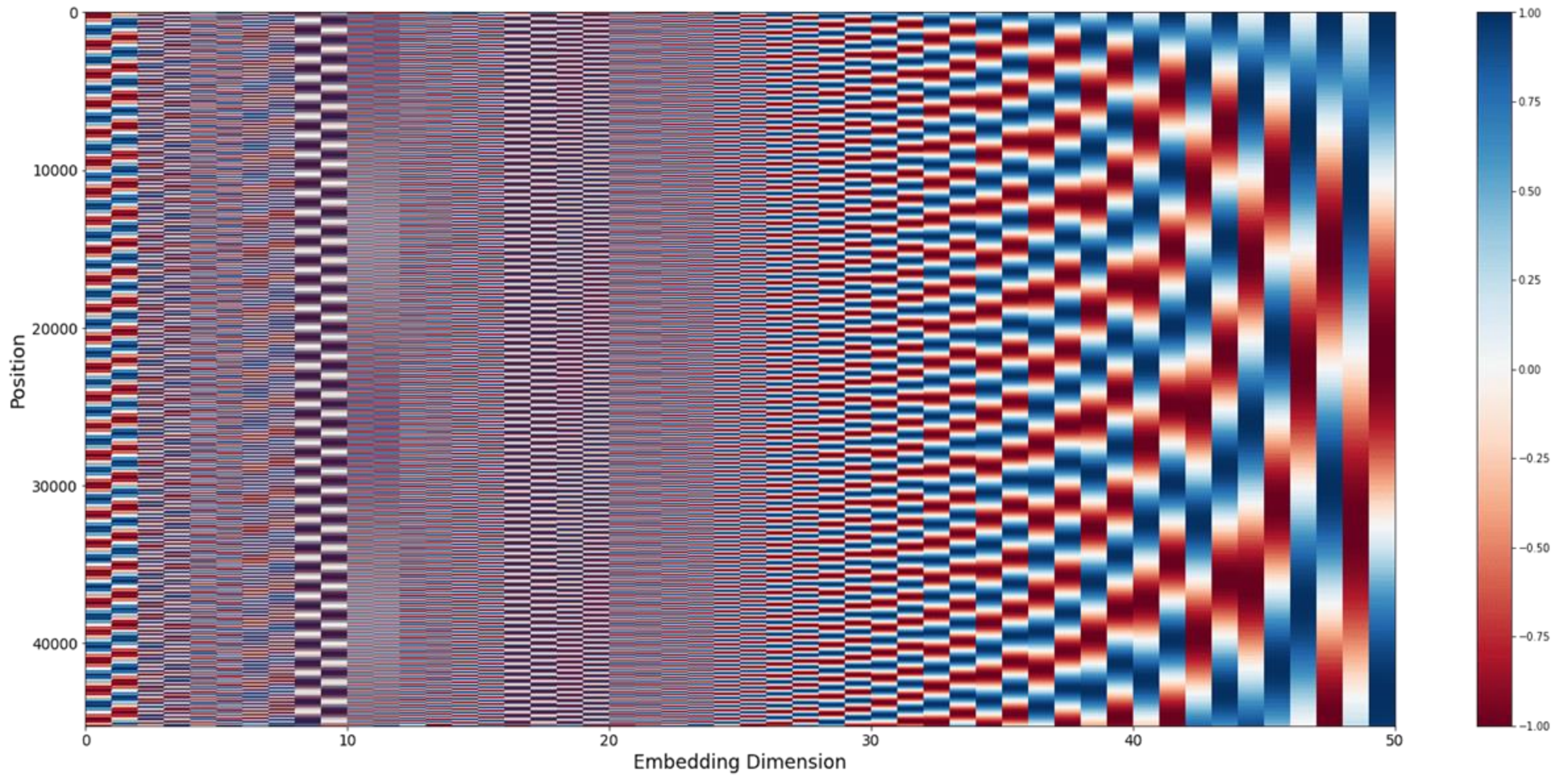
- Sinusoidal position representations: concatenate sinusoidal functions of varying periods:

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad p_i = \begin{pmatrix} \sin(i/10000^{2 \cdot 1/d}) \\ \cos(i/10000^{2 \cdot 1/d}) \\ \vdots \\ \sin(i/10000^{2 \cdot \frac{d}{2}/d}) \\ \cos(i/10000^{2 \cdot \frac{d}{2}/d}) \end{pmatrix}$$

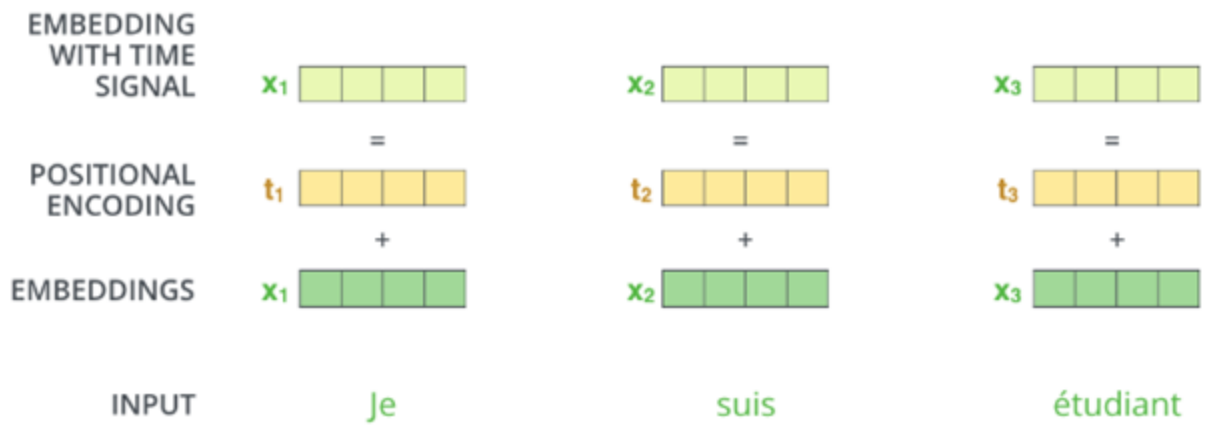
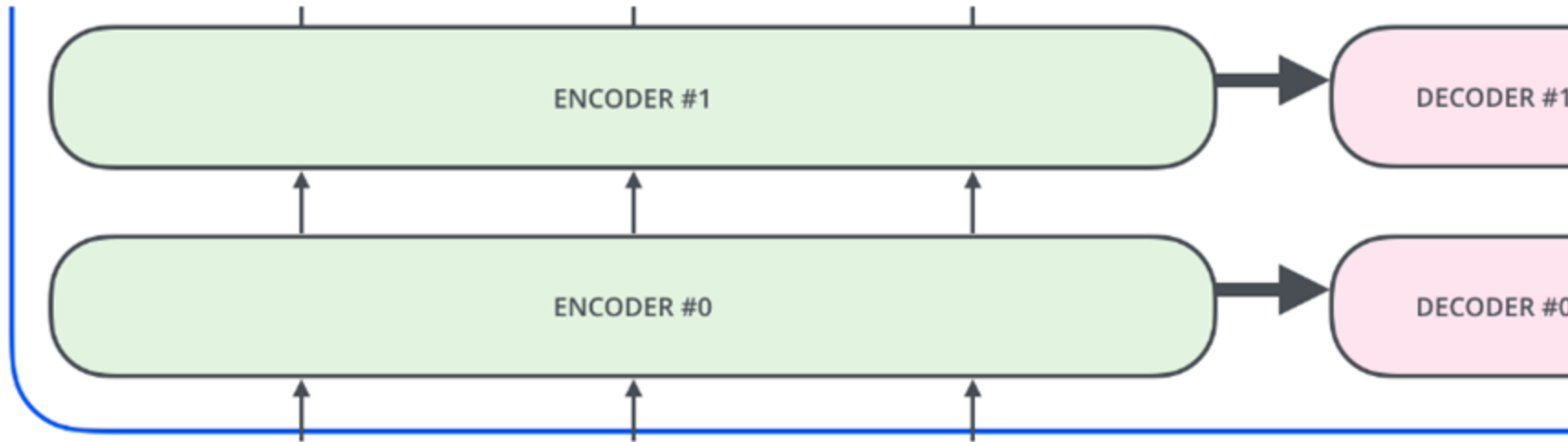
- Pros: Periodicity indicates that maybe “absolute position” isn’t as important
- Cons: Not learnable; also the extrapolation doesn’t really work

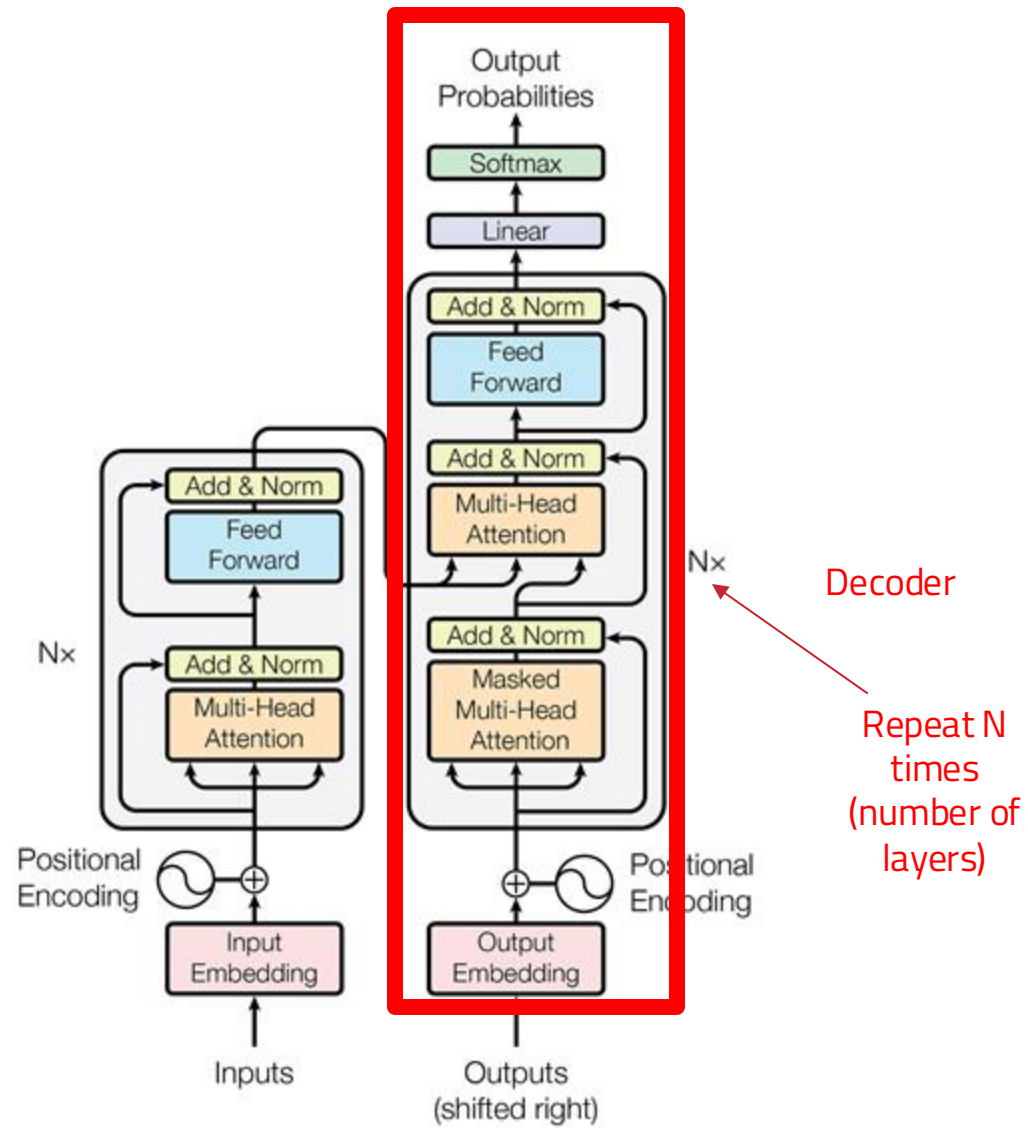










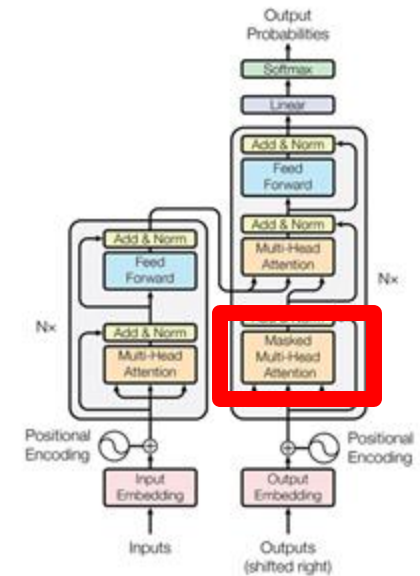


# Decoder: Masked Multi-Head Self-Attention

❑ **Problem:** How do we prevent the decoder from "cheating"? If we have a language modeling objective, can't the network just look ahead and "see" the answer?

❑ **Solution:** Masked Multi-Head Attention.

❑ At a high-level, we hide (mask) information about future tokens from the model.

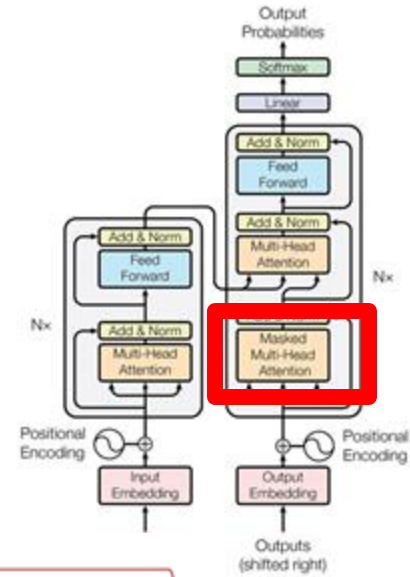


# Masking the future in self-attention

- To use self-attention in decoders, we need to ensure we can't peek at the future.
- At every timestep, we could change the set of keys and queries to include only past words. (Inefficient!)
- To enable parallelization, we mask out attention to future words by setting attention scores to  $-\infty$

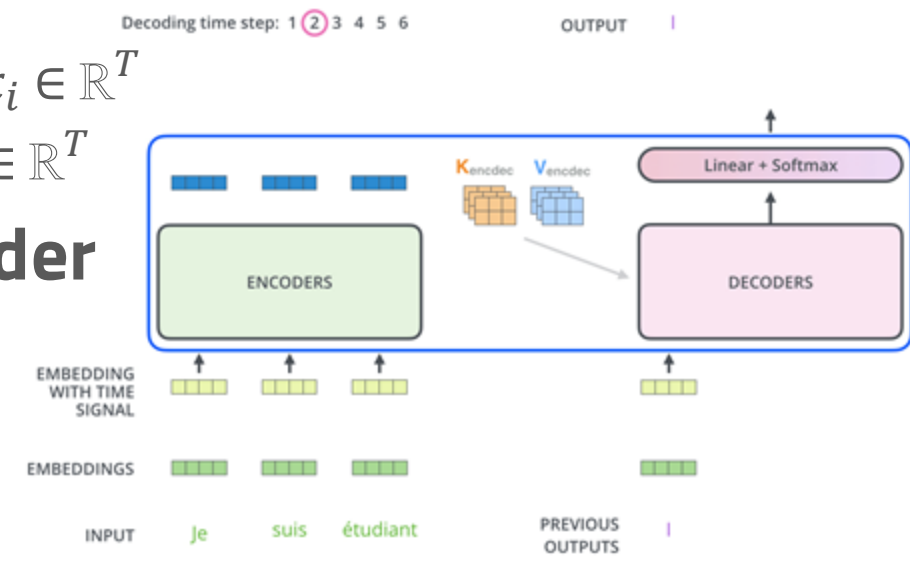
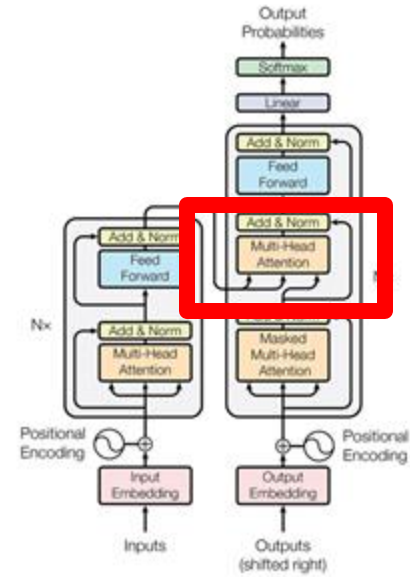
$$e_{ij} = \begin{cases} q_i^T k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

For encoding these words



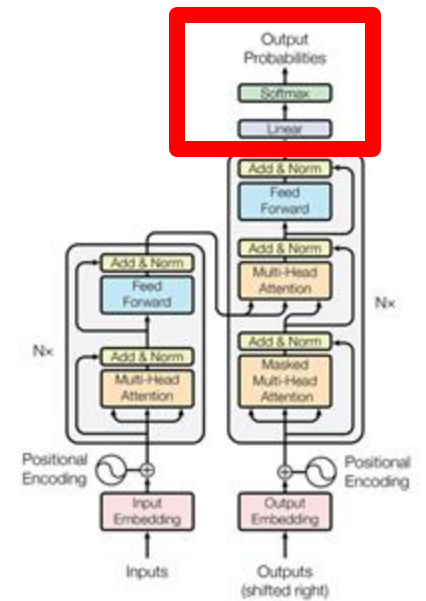
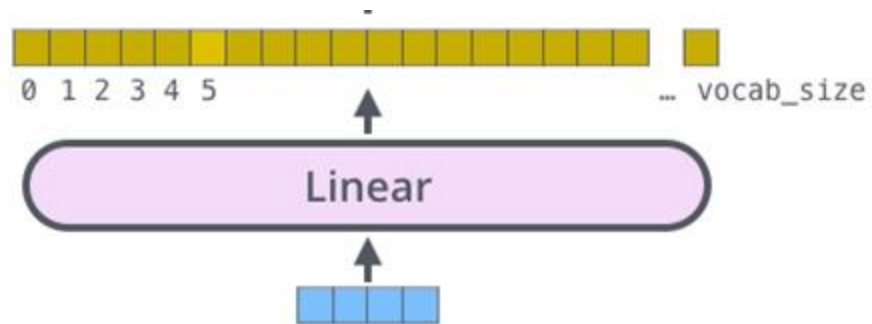
# Encoder-Decoder Attention

- We saw that **self-attention** is when keys, queries, and values come from the same source.
- In the decoder, we have attention that looks more like seq2seq with attention.
  - Let  $h_1.. h_T$  be output vectors from the Transformer encoder;  $x_i \in \mathbb{R}^T$
  - Let  $z_1.. z_T$  be input vectors from the Transformer decoder,  $z_i \in \mathbb{R}^T$
- Then **keys** and **values** are drawn from the **encoder** (like a **memory**):
  - $k_i = K h_i, v_i = V h_i.$
- And the **queries** are drawn from the **decoder**,
  - $q_i = Qz_i$



Decoder stack output

logits



Which word in our vocabulary  
is associated with this index?

am

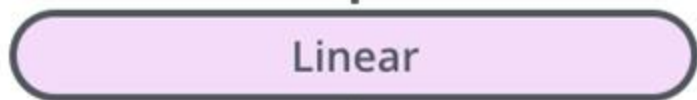
Get the index of the cell  
with the highest value  
(argmax)

5

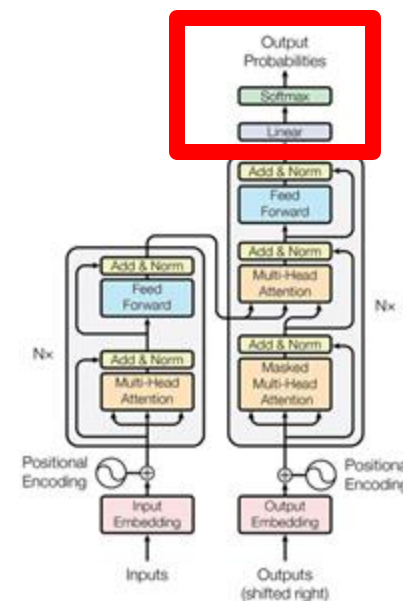
log\_probs



logits



Decoder stack output



# Drawback of Transformer





# Drawback of Transformer

## ❑ Static positional embedding representations:

- Are simple absolute indices the best we can do to represent position?
- Relative linear position attention [Shaw et al., 2018]

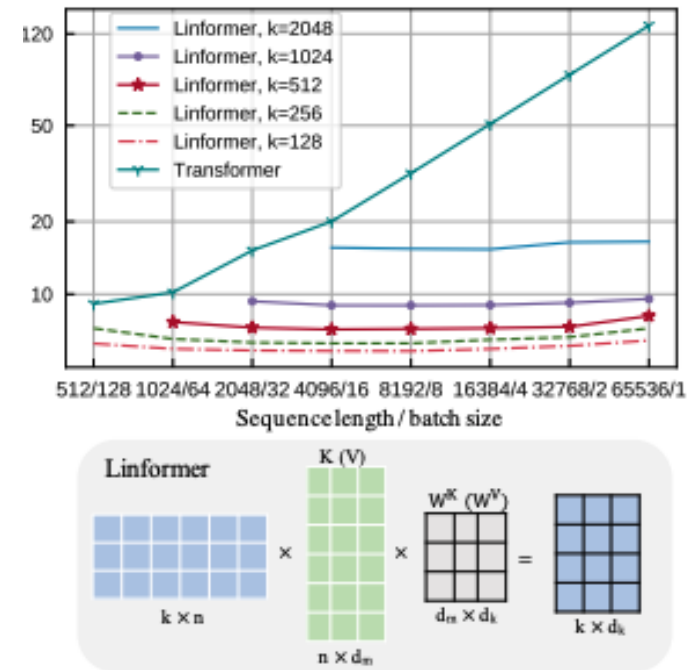
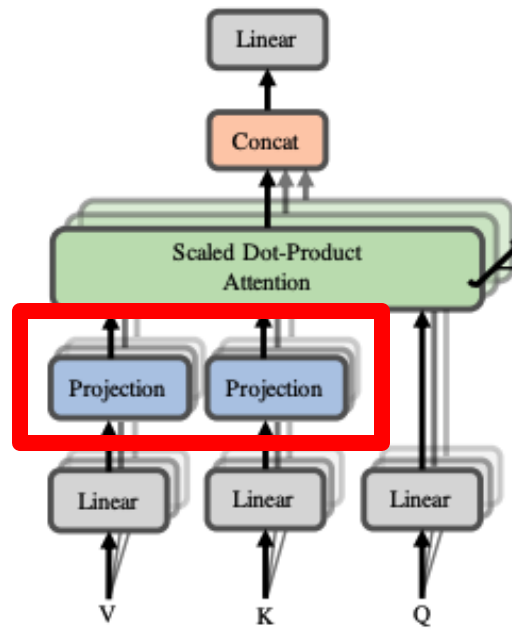
## ❑ Quadratic compute in self-attention:

- Computing all pairs of interactions ( $T^2$ ) means our computation grows quadratically with the sequence length! **For recurrent models, it only grew linearly!**
- Reduce  $O(T^2)$  all-pairs self-attention cost?



# Reduce $O(T^2)$ all-pairs self-attention cost?

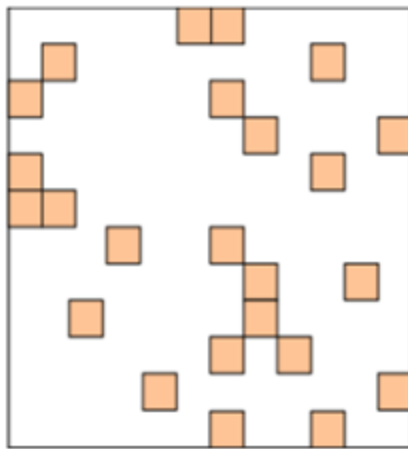
- LinFormer (Wang et al., 2020);  $O(T^2) \rightarrow O(T)$ 
  - Map the sequence length dimension to a lower-dimensional space for values, keys



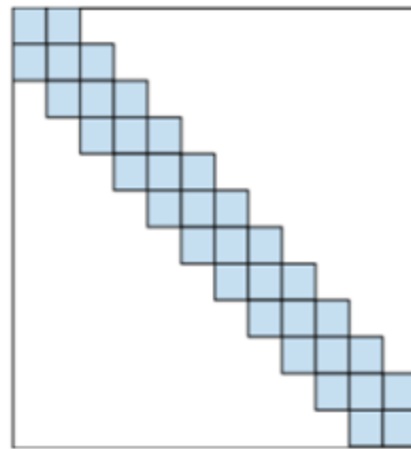
# Reduce $O(T^2)$ all-pairs self-attention cost?

## □ BigBird (Zaheer et al., 2021)

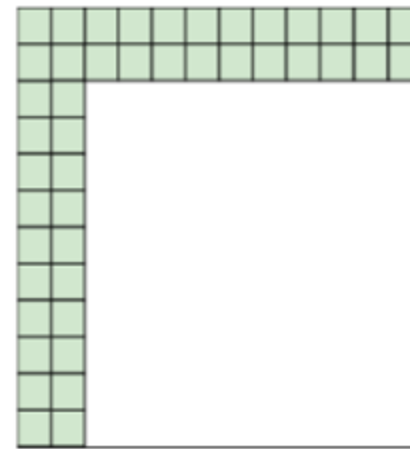
- Replace all-pairs interactions with a family of other interactions, like local windows, looking at everything, and random interactions.



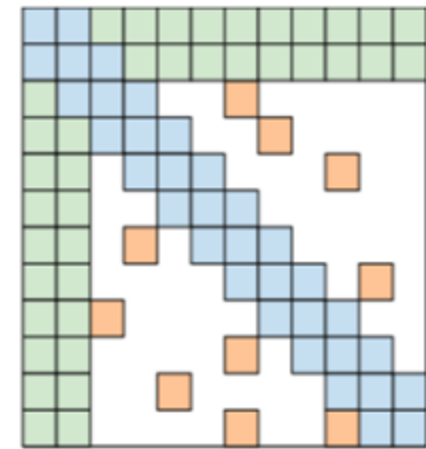
(a) Random attention



(b) Window attention

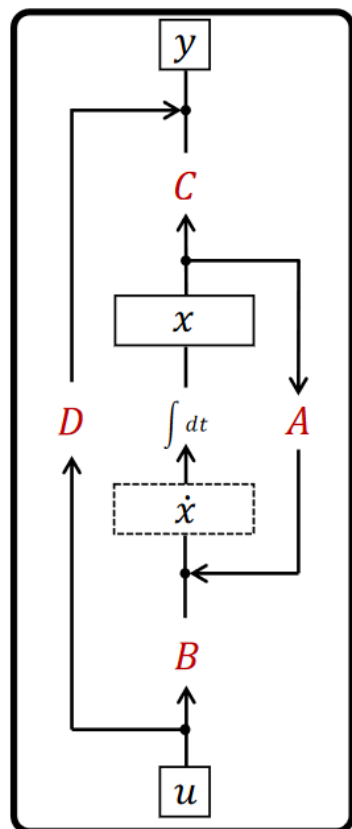


(c) Global Attention



(d) BIGBIRD

# State Space Models

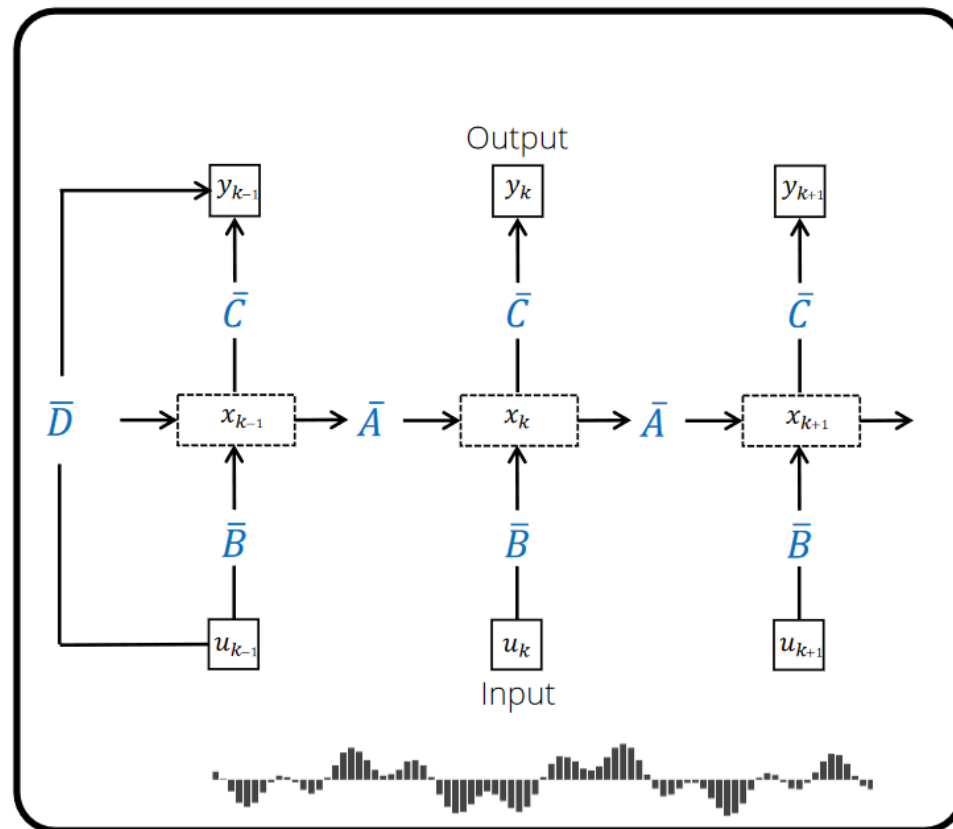


Continuous-time

- ✓ continuous data
- ✓ irregular sampling

Discretize

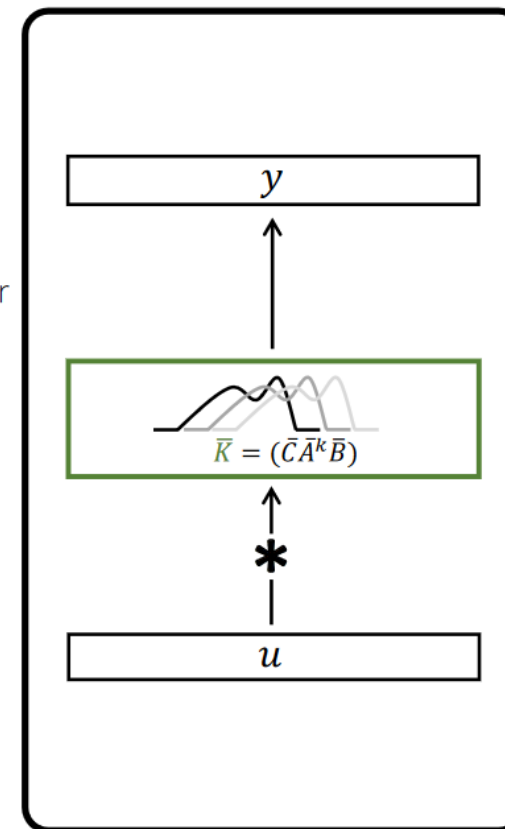
$\Delta t$



Recurrent

- ✓ unbounded context
- ✓ efficient inference

or



Convolutional

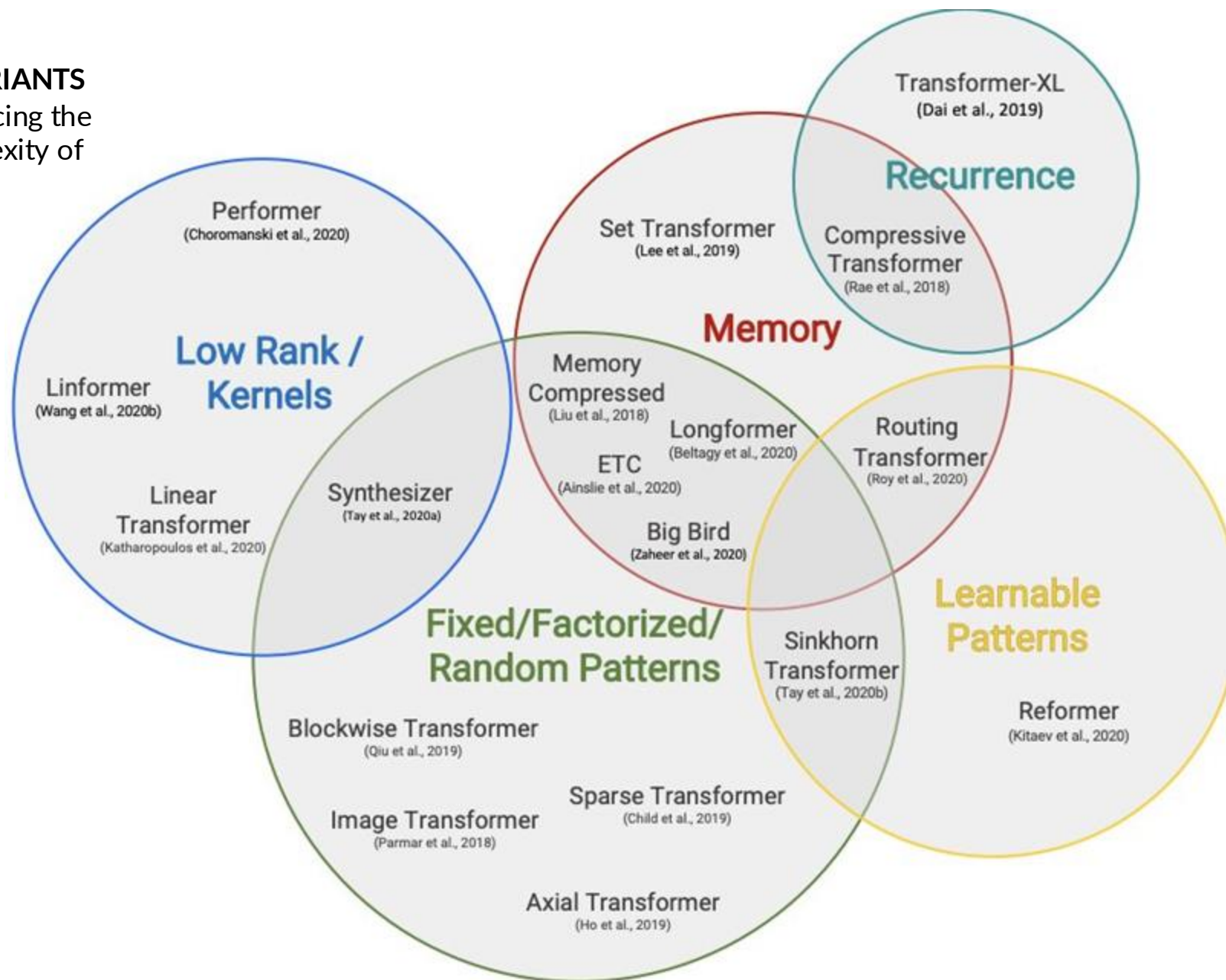
- ✓ local information
- ✓ parallelizable training

<https://huggingface.co/blog/lbourdois/get-on-the-ssm-train>



## TRANSFORMER VARIANTS

Lots of focus on reducing the computational complexity of transformer models.



# Do Transformer Modifications Transfer?

□ 'Surprisingly, we find that most modifications do not meaningfully improve performance.'

Model	Params	Ops	Steps/s	Early loss	Final loss	SGEUE	XSens	WikiQ	WSMT EnDe
Vanilla Transformer	223M	11.3T	3.50	2.182 ± 0.005	1.838	71.66	17.78	23.02	26.62
GeLU	223M	11.3T	3.58	2.179 ± 0.003	1.838	<b>75.79</b>	<b>17.86</b>	<b>25.13</b>	26.47
Swish	223M	11.3T	3.62	2.186 ± 0.003	1.847	<b>73.77</b>	17.74	<b>24.34</b>	<b>26.75</b>
ReLU	223M	11.3T	3.56	2.270 ± 0.007	1.932	67.83	16.73	23.02	26.38
GLU	223M	11.3T	3.59	2.174 ± 0.003	<b>1.814</b>	<b>74.39</b>	<b>17.42</b>	<b>24.34</b>	<b>27.22</b>
GeGLU	223M	11.3T	3.55	2.130 ± 0.006	<b>1.793</b>	<b>75.86</b>	<b>18.27</b>	<b>24.87</b>	<b>26.87</b>
NeGLU	223M	11.3T	3.57	2.145 ± 0.004	<b>1.803</b>	<b>76.17</b>	<b>18.96</b>	<b>24.87</b>	<b>27.62</b>
SeLU	223M	11.3T	3.55	2.315 ± 0.004	1.948	68.76	16.78	22.75	25.89
SwGeLU	223M	11.3T	3.53	2.127 ± 0.003	<b>1.789</b>	<b>76.00</b>	<b>18.20</b>	<b>24.34</b>	<b>27.62</b>
LiGLU	223M	11.3T	3.59	2.149 ± 0.005	<b>1.798</b>	<b>75.34</b>	<b>17.97</b>	<b>24.34</b>	26.53
Sigmoid	223M	11.3T	3.63	2.291 ± 0.019	1.967	74.81	17.51	23.02	26.30
Sofplus	223M	11.3T	3.47	2.297 ± 0.011	1.850	<b>73.45</b>	17.65	<b>24.34</b>	<b>26.89</b>
RMS Norm	223M	11.3T	3.68	2.167 ± 0.008	<b>1.821</b>	<b>75.45</b>	<b>17.94</b>	<b>24.97</b>	<b>27.14</b>
Resers	223M	11.3T	3.51	2.282 ± 0.003	1.939	61.69	15.64	29.50	26.37
Resers + LayerNorm	223M	11.3T	3.26	2.223 ± 0.006	1.858	70.42	17.58	23.02	26.29
Resers + RMS Norm	223M	11.3T	3.34	2.221 ± 0.009	1.875	70.33	17.33	23.02	26.19
Flang	223M	11.3T	3.95	2.392 ± 0.012	2.087	68.56	14.42	23.02	26.31
24 layers, $d_v = 1536, H = 6$	224M	11.3T	3.33	2.290 ± 0.007	1.843	<b>74.89</b>	17.75	<b>25.13</b>	<b>26.89</b>
18 layers, $d_v = 2048, H = 8$	223M	11.3T	3.38	2.185 ± 0.005	<b>1.831</b>	<b>76.45</b>	16.53	<b>24.34</b>	<b>27.19</b>
8 layers, $d_v = 4096, H = 18$	223M	11.3T	3.69	2.180 ± 0.005	1.847	<b>74.58</b>	17.69	<b>23.28</b>	<b>26.85</b>
6 layers, $d_v = 6144, H = 24$	223M	11.3T	3.79	2.291 ± 0.019	1.837	<b>73.55</b>	17.59	<b>23.60</b>	<b>26.68</b>
Block sharing	45M	11.3T	3.91	2.497 ± 0.027	2.164	64.50	14.53	29.96	25.48
+ Factorized embeddings	45M	9.4T	4.31	2.631 ± 0.205	2.383	60.84	14.00	19.84	25.27
+ Factorized & shared embeddings	26M	9.1T	4.27	2.987 ± 0.213	2.385	53.85	11.37	19.84	25.19
Decoder only block sharing	176M	11.3T	3.68	2.298 ± 0.023	1.929	69.60	16.23	23.02	26.23
Encoder only block sharing	144M	11.3T	3.79	2.352 ± 0.029	2.092	67.30	16.13	<b>23.81</b>	26.08
Factorized Embedding	227M	9.4T	3.80	2.298 ± 0.008	1.855	70.41	15.92	22.75	26.50
Factorized & shared embeddings	262M	9.1T	3.92	2.320 ± 0.019	1.952	69.69	16.33	22.22	26.41
Trid. encoder/decoder input embeddings	248M	11.3T	3.55	2.192 ± 0.002	1.840	<b>73.79</b>	17.72	<b>24.34</b>	26.49
Trid. decoder input and output embeddings	248M	11.3T	3.57	2.197 ± 0.007	<b>1.827</b>	<b>74.88</b>	17.74	<b>24.87</b>	<b>26.67</b>
Unid. embeddings	223M	11.3T	3.53	2.195 ± 0.005	<b>1.834</b>	<b>72.99</b>	17.58	<b>23.28</b>	26.48
Adaptive input embeddings	261M	9.2T	3.55	2.250 ± 0.002	1.899	66.57	16.21	<b>24.97</b>	<b>26.68</b>
Adaptive softmax	261M	9.2T	3.60	2.384 ± 0.005	1.992	<b>72.85</b>	16.67	25.16	25.36
Adaptive softmax without projection	223M	10.8T	3.43	2.229 ± 0.009	1.914	<b>71.82</b>	17.18	23.02	25.72
Mixture of softmaxes	223M	16.3T	2.26	2.227 ± 0.017	<b>1.821</b>	<b>76.77</b>	17.62	22.75	<b>26.82</b>
Transparent attention	223M	11.3T	3.23	2.181 ± 0.014	1.876	64.21	16.40	25.16	<b>26.80</b>
Dynamic convolution	257M	11.8T	2.65	2.483 ± 0.009	2.047	68.30	12.67	25.16	17.83
Lightweight convolution	224M	10.4T	4.07	2.370 ± 0.019	1.969	63.87	14.96	23.02	24.73
Evolved Transformer	223M	9.8T	3.69	2.220 ± 0.003	1.863	<b>73.87</b>	18.78	<b>24.97</b>	26.58
Synthesizer (dense)	224M	11.4T	3.47	2.334 ± 0.021	1.962	61.83	14.27	16.14	<b>26.63</b>
Synthesizer (dense plus)	243M	12.8T	3.22	2.191 ± 0.019	1.840	<b>73.88</b>	16.96	<b>23.81</b>	<b>26.71</b>
Synthesizer (dense plus alpha)	243M	12.8T	3.01	2.180 ± 0.007	<b>1.828</b>	<b>74.23</b>	17.02	<b>23.28</b>	26.61
Synthesizer (factorized)	207M	10.2T	3.96	2.341 ± 0.017	1.968	62.78	15.39	<b>23.55</b>	26.42
Synthesizer (random)	251M	10.2T	4.98	2.326 ± 0.012	2.009	64.27	18.45	19.56	26.41
Synthesizer (random plus)	282M	12.8T	3.63	2.189 ± 0.004	1.842	<b>73.52</b>	17.04	<b>24.87</b>	26.43
Synthesizer (random plus alpha)	292M	12.8T	3.42	2.196 ± 0.007	<b>1.828</b>	<b>75.34</b>	17.08	<b>24.06</b>	26.39
Universal Transformer	64M	40.8T	0.88	2.496 ± 0.036	2.053	70.13	14.09	19.05	23.91
Mixture of experts	649M	11.7T	3.29	2.148 ± 0.008	<b>1.785</b>	<b>74.55</b>	<b>18.13</b>	<b>24.06</b>	<b>26.84</b>
Swish Transformer	1166M	11.7T	3.18	2.133 ± 0.005	<b>1.758</b>	<b>73.28</b>	<b>18.82</b>	<b>26.19</b>	<b>26.85</b>
Fixed Transformer	223M	1.9T	4.30	2.359 ± 0.008	1.918	67.34	16.28	22.75	23.20
Weighted Transformer	286M	71.8T	0.59	2.578 ± 0.022	1.989	69.94	16.98	23.02	26.30
Product key memory	421M	306.6T	0.25	2.155 ± 0.003	<b>1.798</b>	<b>75.16</b>	17.04	<b>23.55</b>	<b>26.73</b>

## Do Transformer Modifications Transfer Across Implementations and Applications?

Sharan Narang\* Hyung Won Chung Yi Tay William Fedus  
 Thibault Fevry† Michael Matena† Karishma Malkan† Noah Fiedel  
 Noam Shazeer Zhenzhong Lan† Yanqi Zhou Wei Li  
 Nan Ding Jake Marcus Adam Roberts Colin Raffel†



# Scaling up Transformer

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hrs)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)

<http://hal.cse.msu.edu/teaching/2020-fall-deep-learning/14-nlp-and-transformers/#/22/0/9>





# Scaling up Transformer

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hrs)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13GB	
BERT-Large	24	1024	16	340M	13GB	



# Scaling up Transformer

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hrs)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13GB	
BERT-Large	24	1024	16	340M	13GB	
XLNet-Large	24	1024	16	340M	126GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160GB	1024x V100 (1 day)
GPT-2	48	1600	?	1.5B	40GB	
Megatron-LM	72	3072	32	8.3B	174GB	512x V100 (9 days)
Turing-NLG	78	4256	28	17B	?	256x V100
GPT-3	96	12288	96	175B	694GB	?

Brown et al, "Language Models are Few-Shot Learners", arXiv 2020

<http://hal.cse.msu.edu/teaching/2020-fall-deep-learning/14-nlp-and-transformers/#/22/0/9>



# Summary

- ❑ Transformers are a new neural network model that only uses attention (and many other training tricks!!)
- ❑ However, the models are extremely expensive
- ❑ Improvements (unfortunately) seem to mostly come from even more expensive models and more data
- ❑ If you can afford large data and large compute, transformers are the go to architecture, instead of CNNs, RNNs, etc.
  - Why? On our way back to fully-connected models, throwing out the inductive bias of CNNs and RNNs.

