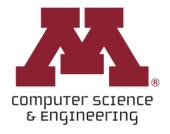
# CSCI 5541: Natural Language Processing

**Lecture 10: Deep Dive on Transformers** 





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



# Announcement (1014)

- Proposal Report (due: tonight)
- ☐ HW4 due (due: Oct 19, Sunday)
  - o Need extension?
- ☐ A series of exciting guest lectures are booked!
  - LLMs as agents (Oct 28) by Shuyu Gan
  - Artificial cognition (Nov 6) by Karin de Langis
  - Interpretability and Explainability (Nov 11) by Ryan Peters
  - Human-Al Collaboration (Nov 20) by Shirley A. Hayati





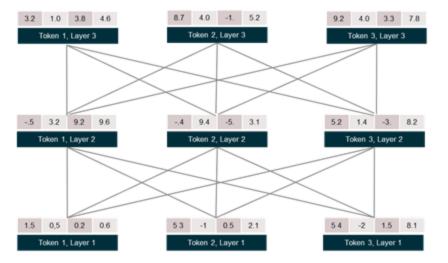




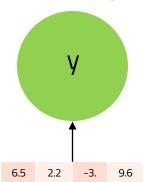
Stacked Bidirectional RNN trained to predict next word in language modeling task

53 85 -1 53 85 -1 21 87 ; 21 87 ; like

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

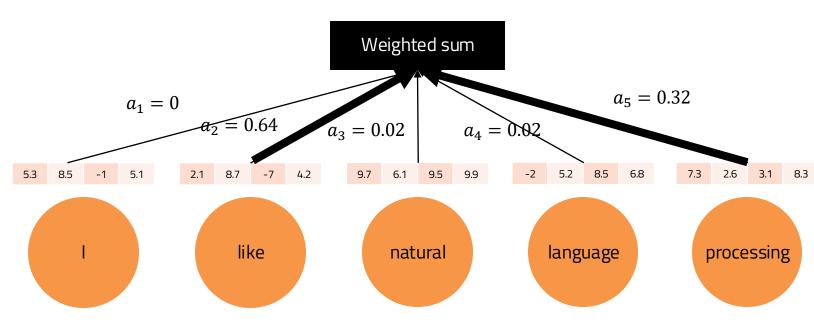


#### Positive / Negative



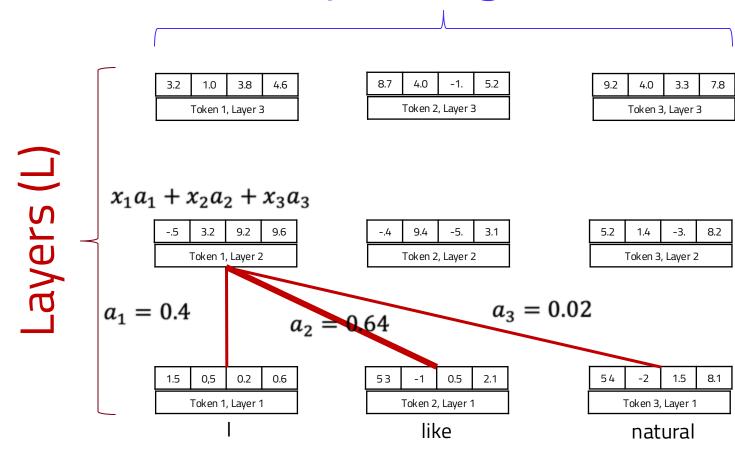
$$x_1a_1 + x_2a_2 + x_3a_3 + x_4a_4 + x_5a_5$$







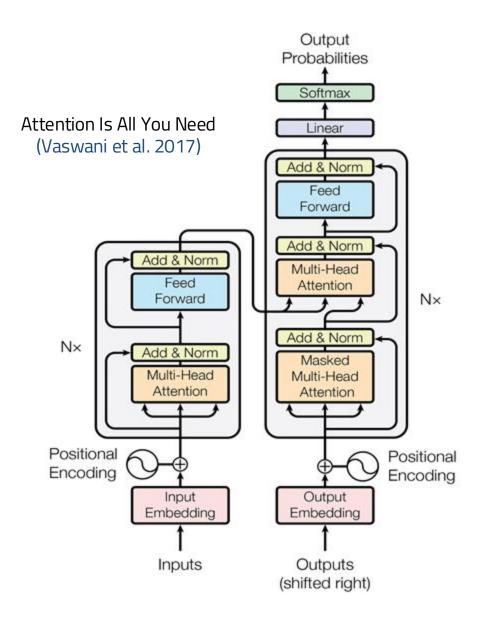
## Input Length (T)



# 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



# Strong results/findings and applications of Transformers

## Strong results with Transformers on machine translation

BLEU		Training Cost (FLOPs)		
EN-DE	EN-FR	EN-DE	EN-FR	
23.75				
	39.2		$1.0 \cdot 10^{20}$	
24.6	39.92	$2.3\cdot 10^{19}$	$1.4\cdot 10^{20}$	
25.16	40.46	$9.6\cdot 10^{18}$	$1.5\cdot 10^{20}$	
26.03	40.56	$2.0\cdot 10^{19}$	$1.2\cdot 10^{20}$	
	40.4		$8.0 \cdot 10^{20}$	
26.30	41.16	$1.8 \cdot 10^{20}$	$1.1\cdot 10^{21}$	
26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$	
27.3	38.1	3.3	$3.3 \cdot 10^{18}$ $2.3 \cdot 10^{19}$	
28.4	41.8	2.3		
	EN-DE 23.75 24.6 25.16 26.03 26.30 26.36 27.3	EN-DE EN-FR  23.75  39.2  24.6 39.92  25.16 40.46 26.03 40.56  40.4  26.30 41.16 26.36 41.29  27.3 38.1	EN-DE EN-FR EN-DE  23.75  39.2  24.6 39.92 25.16 40.46 9.6 · 10 <sup>18</sup> 26.03 40.56 20 · 10 <sup>19</sup> 40.4  26.30 41.16 26.36 41.29  27.3 38.1  3.3	

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

(Vaswani et al. 2017)

## Strong results with Transformers on document summarization

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

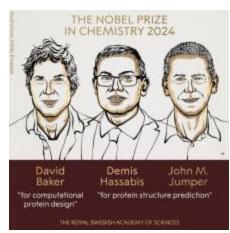
Rank	Model	Accuracy 🕈	Paper	Code	Result	Year	Tags 08'
1	SMART-RoBERTa Large	97.5	SMART: Robust and Efficient Fine-Tuning for Pre-trained Natural Language Models through Principled Regularized Optimization	C	Ð	2019	Transformer
2	T5-3B	97.4	Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer	C	Ð	2019	Transformer
3	MUPPET Roberta Large	97.4	Muppet: Massive Multi-task Representations with Pre- Finetuning	C	-0	2021	
4	ALBERT	97.1	ALBERT: A Lite BERT for Self-supervised Learning of Language Representations	C	Ð	2019	Transformer
5	T5-11B	97.1	Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer	0	Ð	2019	Transformer
6	StructBERTRoBERTa ensemble	97.1	StructBERT: Incorporating Language Structures into Pre- training for Deep Language Understanding		Ð	2019	Transformer
7	XLNet (single model)	97	XLNet: Generalized Autoregressive Pretraining for Language Understanding	0	-9	2019	Transformer
8	ELECTRA	96.9	ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators	C	-9	2020	
9	EFL	96.9	Entailment as Few-Shot Learner	C	-9	2021	Transformer
10	XLNet-Large (ensemble)	96.8	XLNet: Generalized Autoregressive Pretraining for Language Understanding	C	-10	2019	Transformer
11	RoBERTa	96.7	RoBERTa: A Robustly Optimized BERT Pretraining Approach	C	-20	2019	Transformer

https://paperswithcode.com/

## Transformers used outside of NLP

## Protein folding





AlphaFold2 (Jumper et al., 2021)

## Image Classification

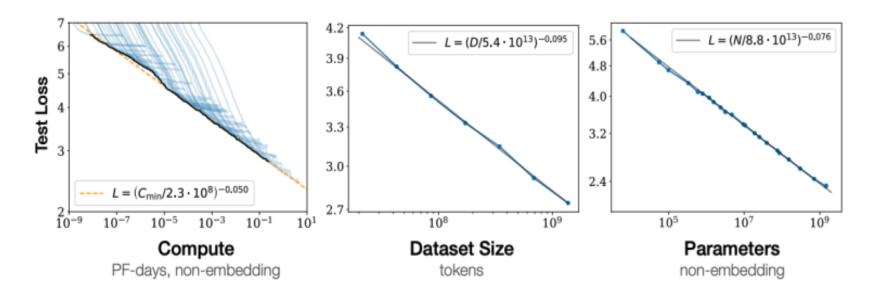


Vision Transformer (ViT) outperforms ResNet-based baselines with substantially less compute (Dosovitskiy et al. 2020)

11

# 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!

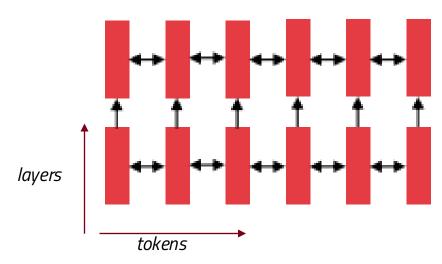


Kaplan et al., 2020, Scaling Laws for Neural Language Models

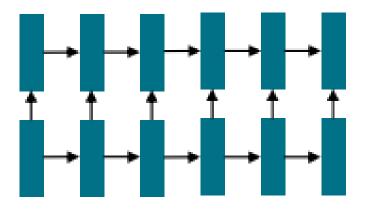
Why self-attention?

## Recurrence in RNNs





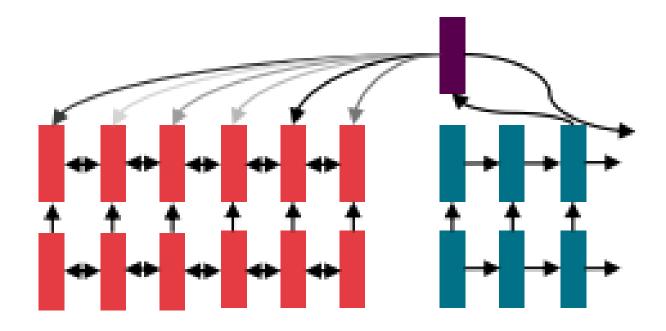
**Encoding**: Encode input sentences with bi-directional LSTM



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

# Sequence-to-sequence with attention

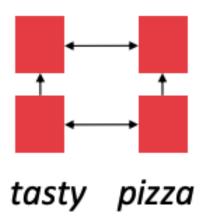


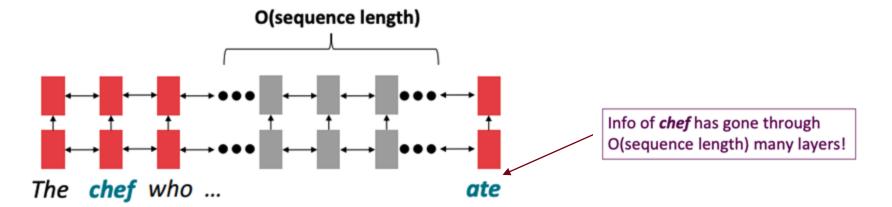


Use **attention** to allow flexible access to input memory

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

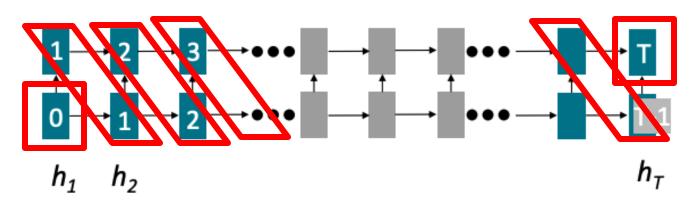
- Forward RNNs are unrolled "left-to-right".
- ☐ It encodes linear locality:
  - Nearby words often affect each other's meanings
- Problem: RNNs take **O(sequence length) steps** for distant word pairs to interact





## Issues with recurrent models: Lack of parallelizability

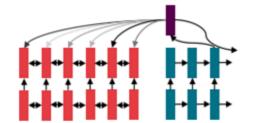
- ☐ Forward and backward passes have **O(seq length) un-parallelizable** operations
  - o GPUs (and TPUs) can perform many independent computations at once! But future RNN hidden states can't be computed fully before past RNN hidden states have been computed
  - Particularly problematic as sequence length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



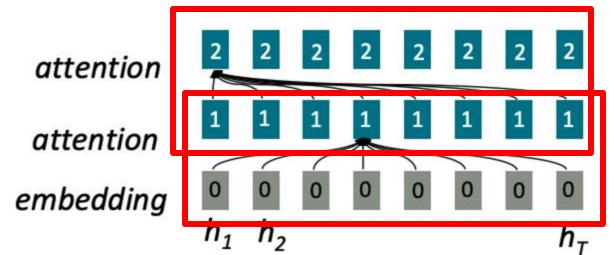
Numbers indicate min # of steps before a state can be computed

## 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.

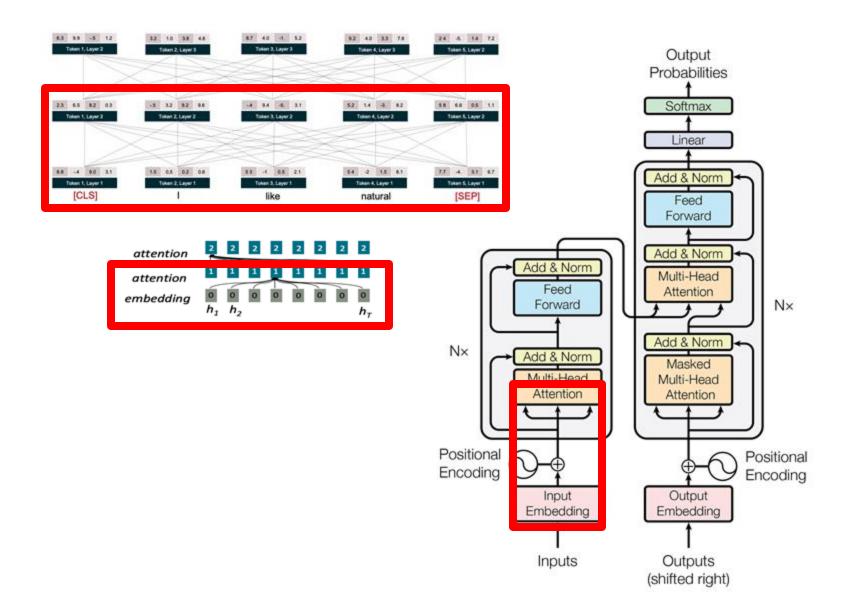


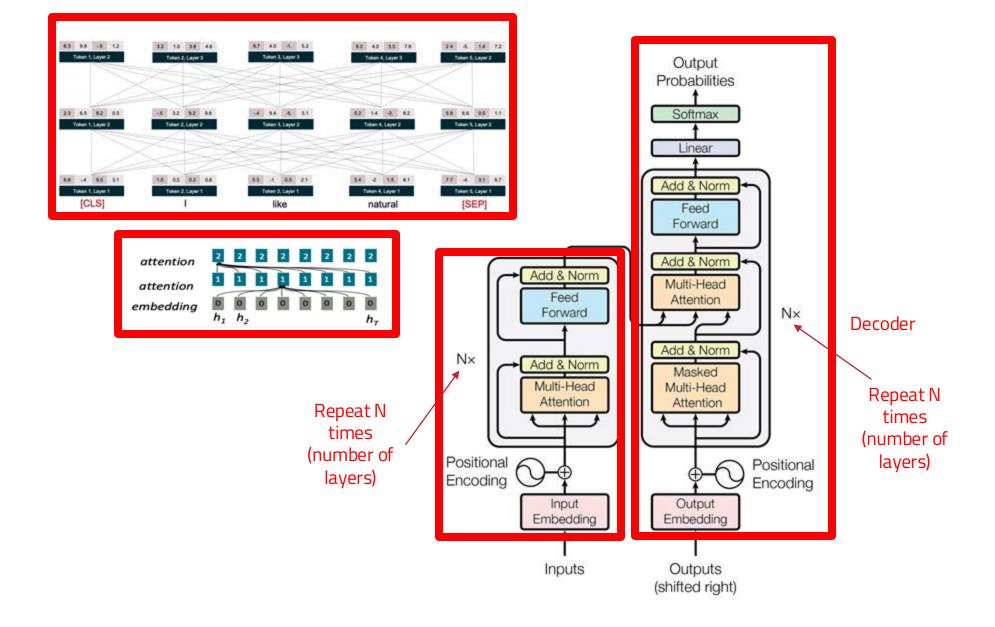
- We saw attention from the decoder to the encoder;
- Self-attention is encoder-encoder (or decoder-decoder) attention where each word attends to each other word within the input (or out



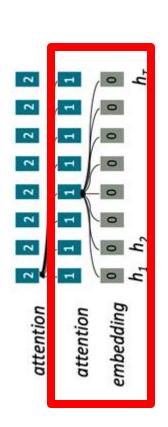
All words attend to all words in previous layer; most arrows are omitted

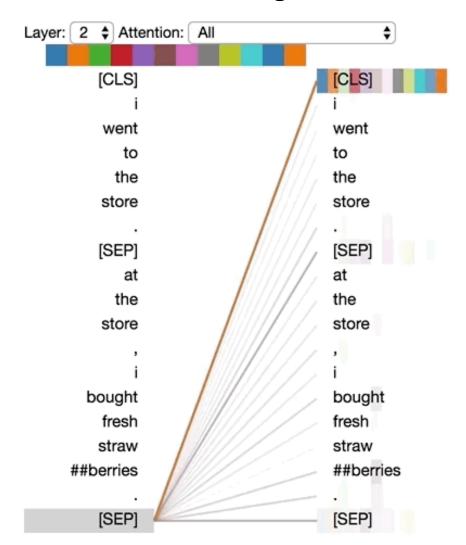
O(seq length) O(Layers)

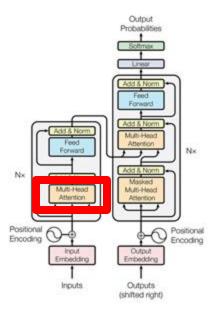




#### "I went to the store. At the store, I bought fresh strawberries."

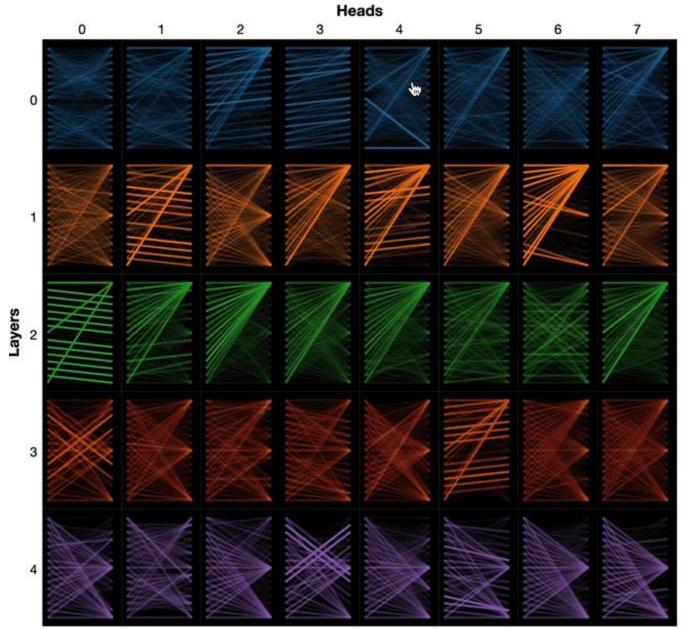


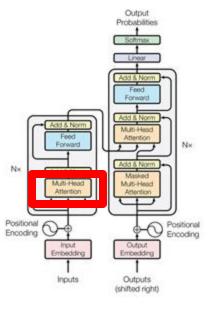




https://github.com/jessevig/bertviz

https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\_t2tipynb





HITTP://glulup.com/jessevig/peit/viz

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**.

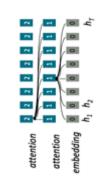
Let's think of attention as a "fuzzy" or approximate hashtable:

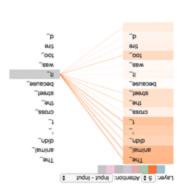


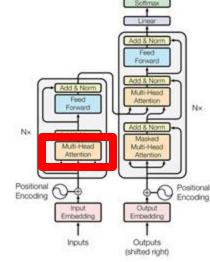


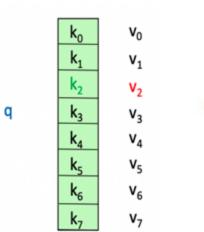
Each query (hash) maps to exactly one key-value pair.

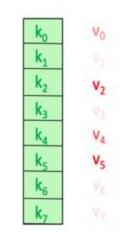
- In (self-)attention:
  - Each query (token in current layer) matches each key to varying degrees.
  - We return a sum of values (token in previous layer) weighted by the guery-key match (attention score).





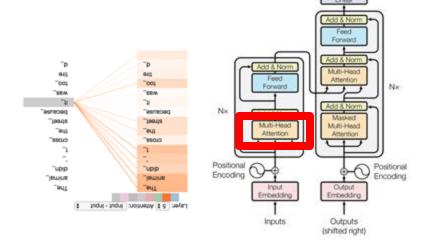




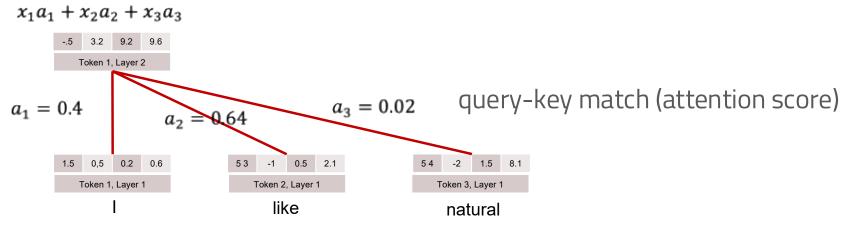


## **Encoder: Self-Attention**

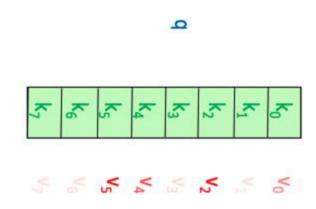
In (self-)attention: Each query (token in current layer) matches each key to varying degrees. We return a sum of values (token in previous layer) weighted by the query-key match (attention score).



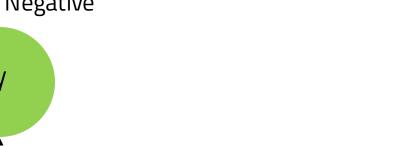
query (token in current layer)

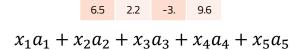


values (token in previous layer)



#### Positive / Negative



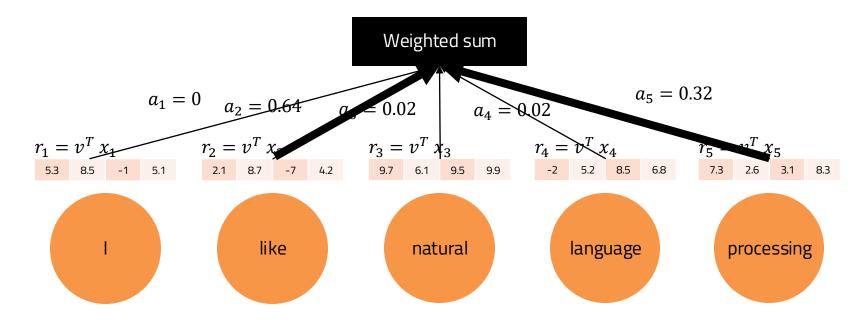




$$a = softmax(r)$$

 $v \in R^h$ 

6.5 2.2 -3. 9.6



## Recipe for Self-Attention in the Transformer Encoder

Model parameters to learn (randomly initialized)



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

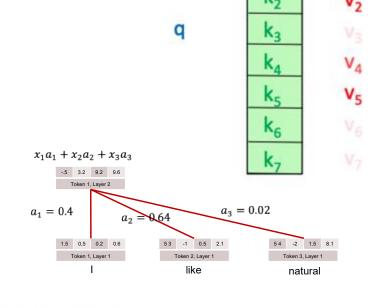
☐ Step 2: Calculate attention score between query and keys.

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

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

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

Step 4: Take a weighted sum of values.



$$Output_i = \sum_j \alpha_{ij} v_j$$

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

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

$$Q = XW^Q$$
  $K = XW^K$   $V = XW^V$ 

Step 2: Calculate attention score between query and keys.

$$E = QK^T$$

Step 3: Take the softmax to normalize attention scores.

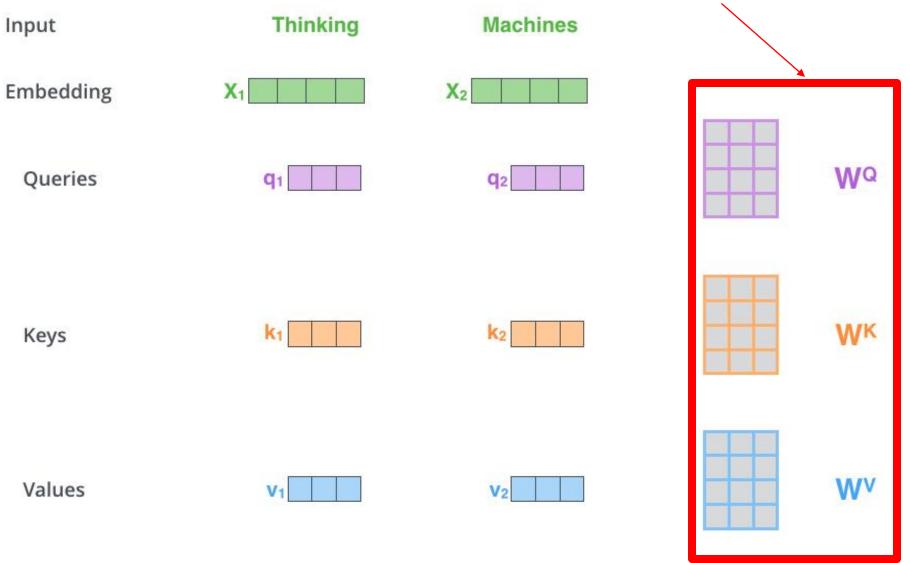
$$A = softmax(E)$$

Step 4: Take a weighted sum of values.

$$Output = AV$$

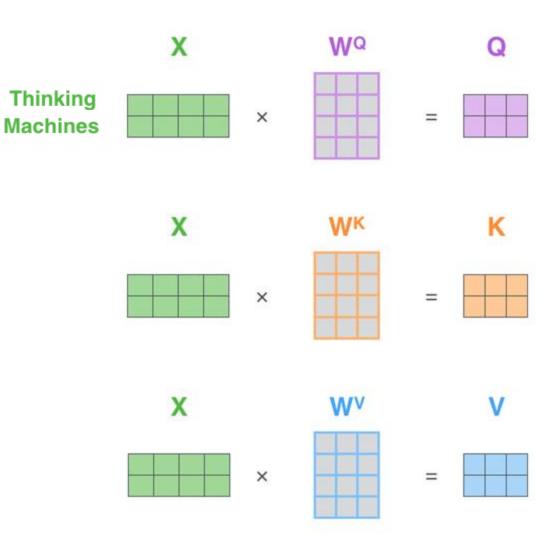
$$Output = softmax(QK^T)V$$

#### Model parameters to learn (randomly initialized)



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

$$Q = XW^Q$$
  $K = XW^K$   $V = XW^V$ 



Step 2: Calculate attention score between query and keys.

$$E = QK^T$$

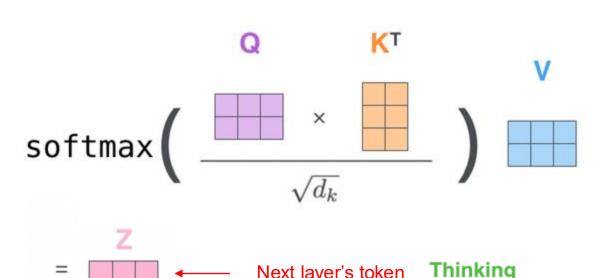
Step 3: Take the softmax to normalize attention scores.

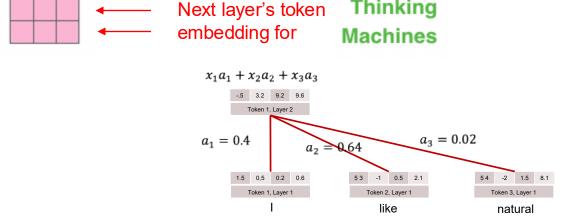
$$A = softmax(E)$$

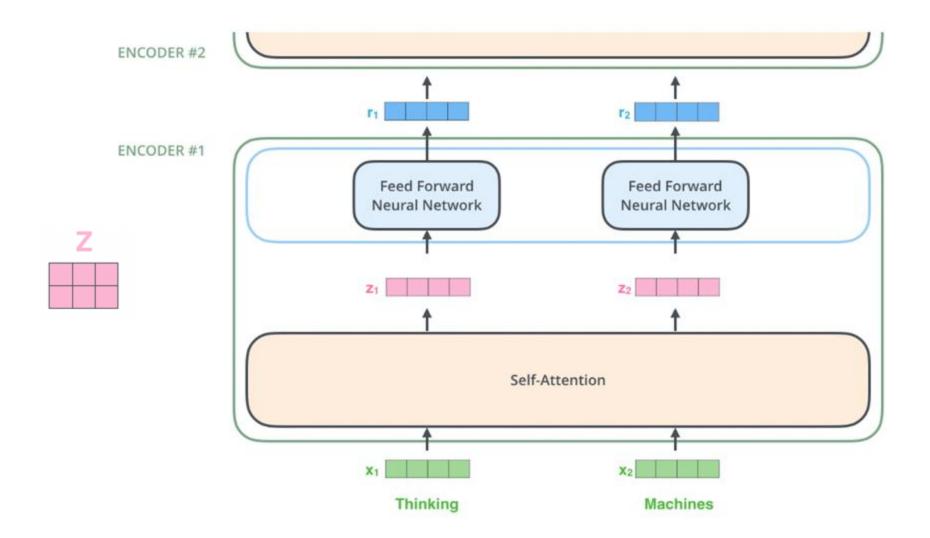
Step 4: Take a weighted sum of values.

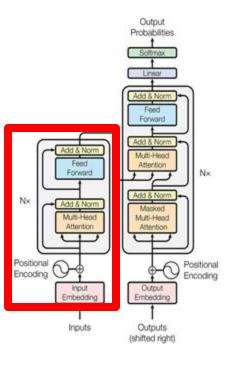
$$Output = AV$$

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



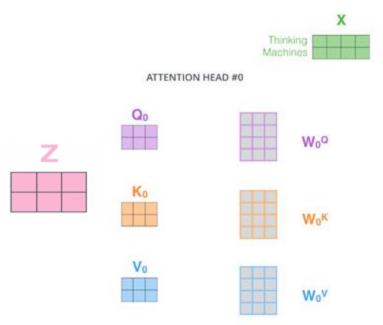






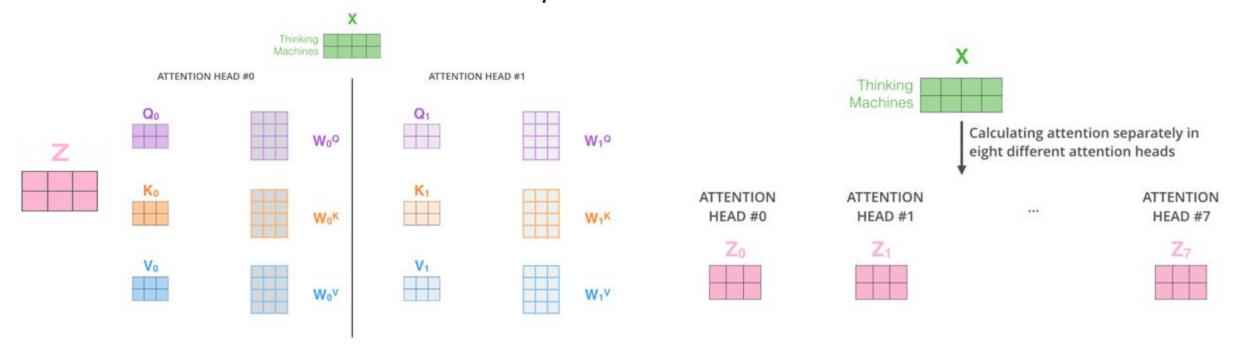
## Multi-headed self-attention

- ☐ 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.

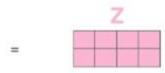


## Condensing multi-head attentions into a single matrix



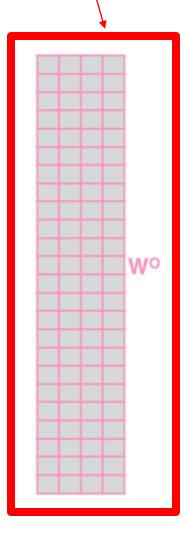


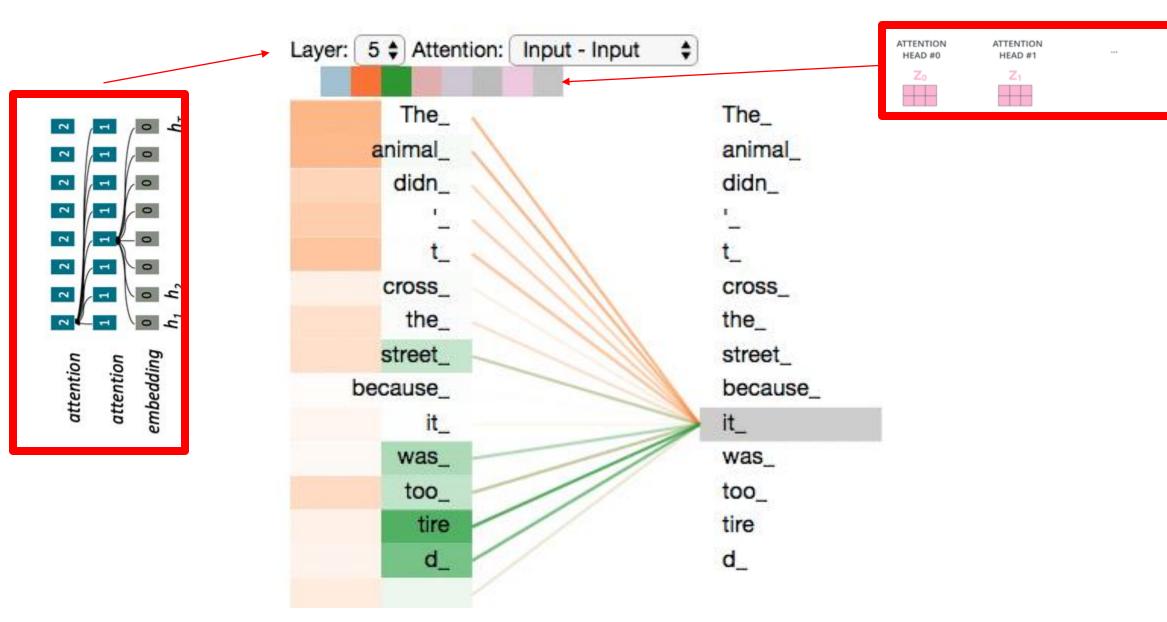
3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



 Multiply with a weight matrix W<sup>o</sup> that was trained jointly with the model

X



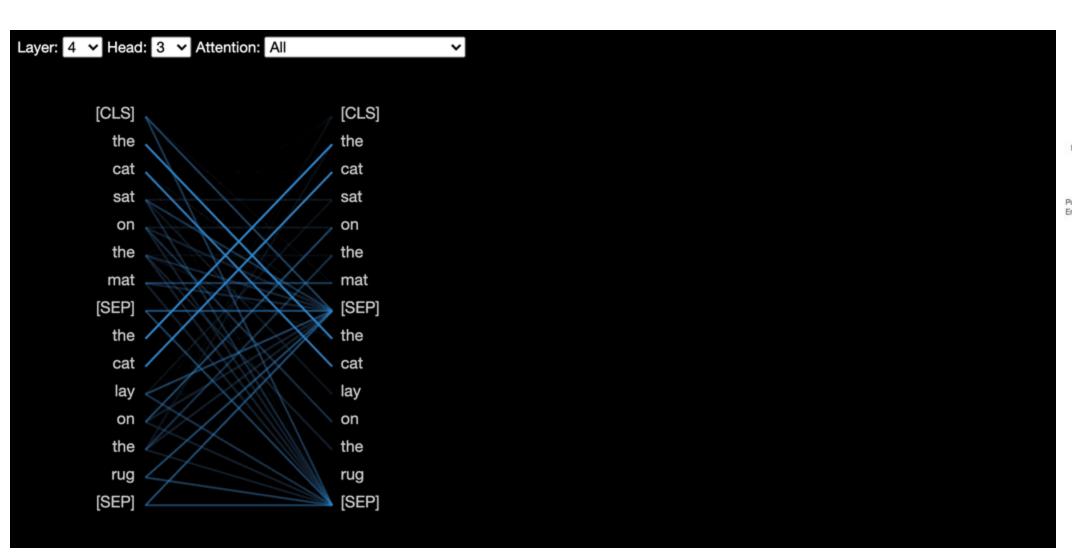


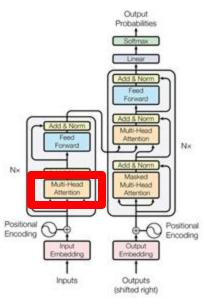
https://github.com/jessevig/bertviz

https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\_t2t.ipynb

ATTENTION

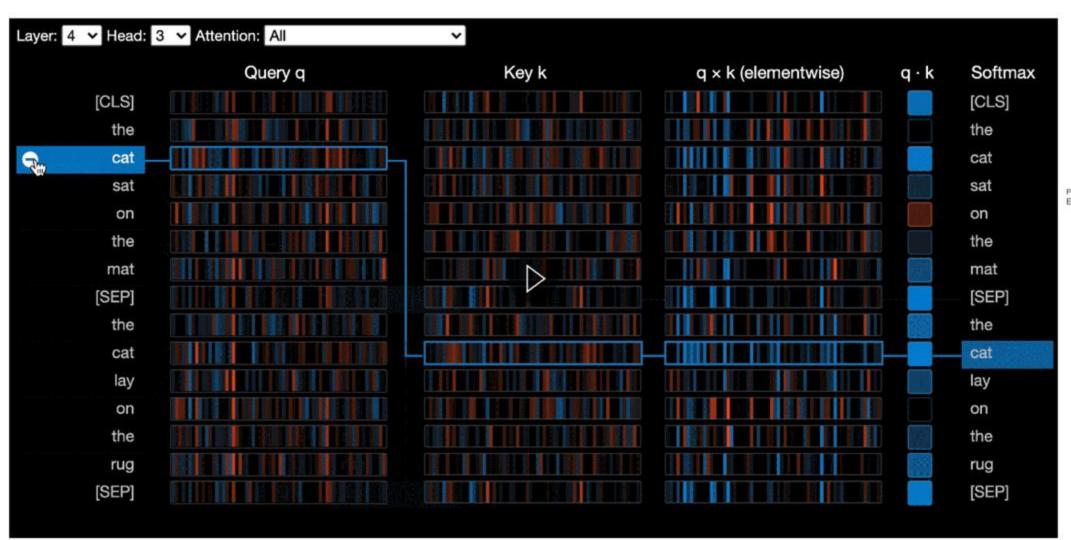
HEAD #7

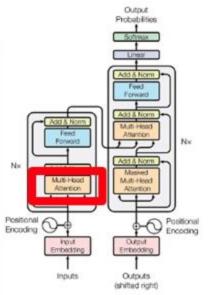




https://github.com/jessevig/bertviz

 $\underline{https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\_t2t.ipynb$ 





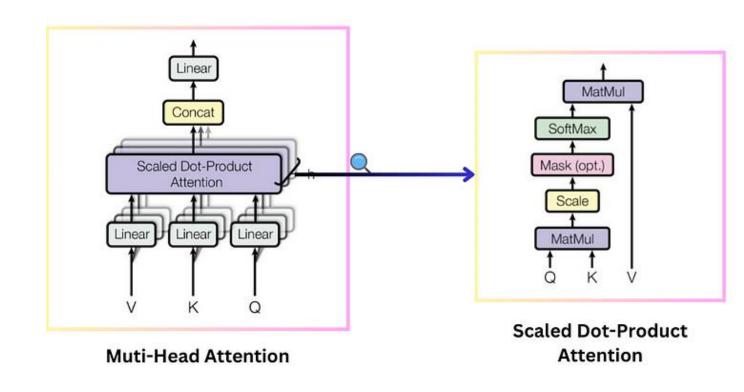
https://github.com/jessevig/bertviz

https://colab.research.google.com/github/tensorflow/tensor/blob/master/tensor/2tensor/notebooks/hello\_t2t.ipynb\_

## Recap

Multi-Head Attention

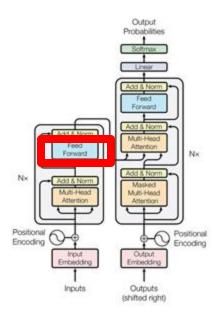
- □ Pass our input through the W<sub>v</sub>, W<sub>k</sub>, W<sub>q</sub> matrices for each head (corresponding to the 'Linear' boxes at right)
- □ Perform Scaled dot product attention for each head
- Concatenate the results for each head
- ☐ Use linear layer to project to original output dimension

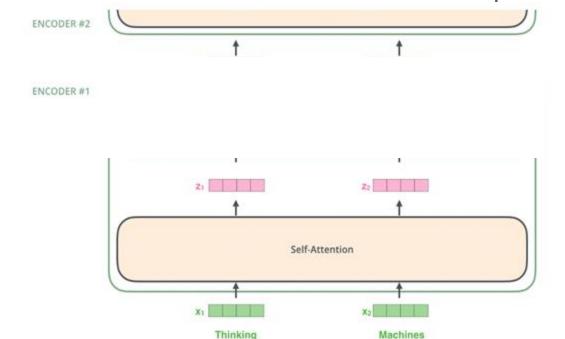


Other tricks than attention?

# 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).



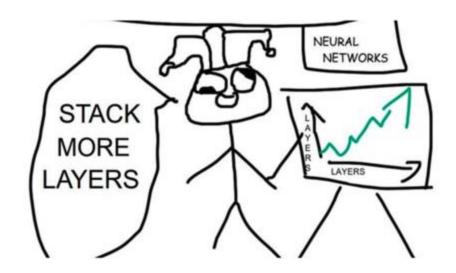


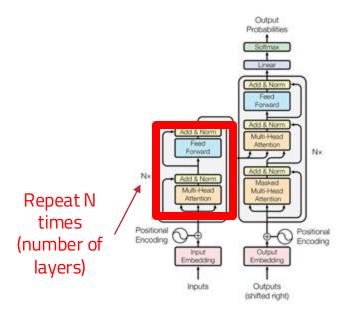
Equation for Feed-Forward layer

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

## Stacking deep neural nets

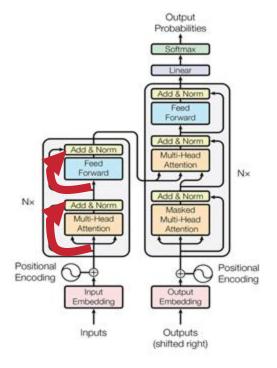
- ☐ Training trick #1: Residual Connections
- ☐ Training trick #2: LayerNorm
- ☐ Training trick #3: Scaled Dot Product Attention





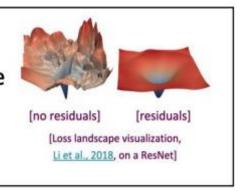
#### Trick #1: Residual Connections [He et al., 2016]

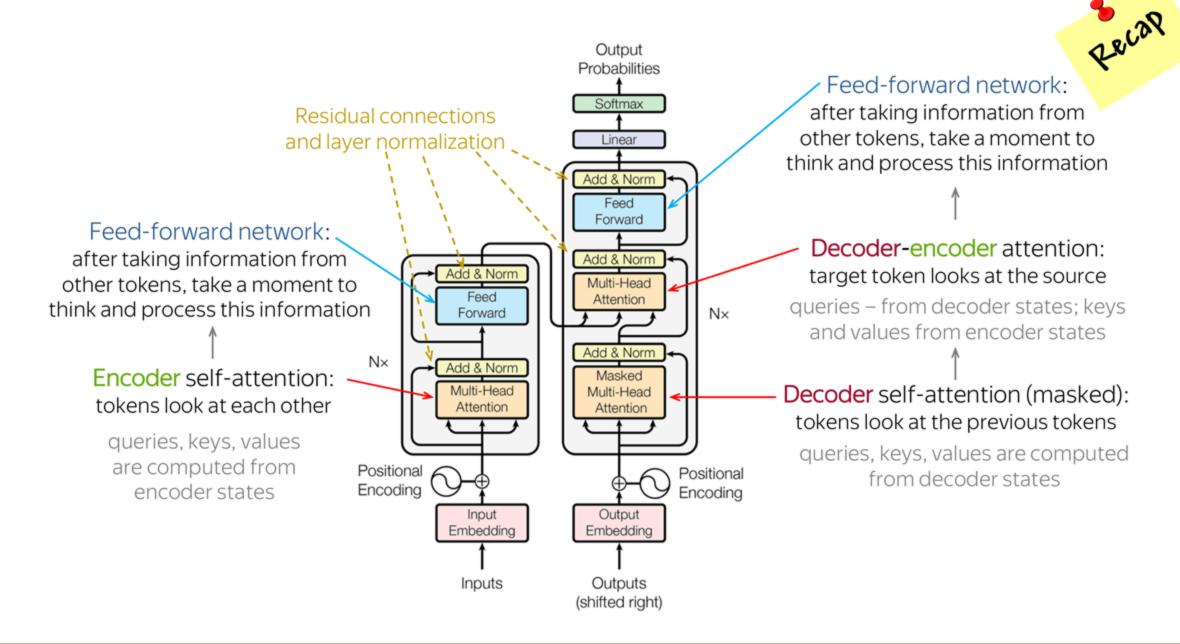
- Residual connections are a simple but powerful technique from computer vision.
- ☐ Similar to additive connection in LSTM
- ☐ Directly passing "raw" embeddings to the next layer prevents the network from "forgetting" or distorting important information as it is processed by many layers.



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

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

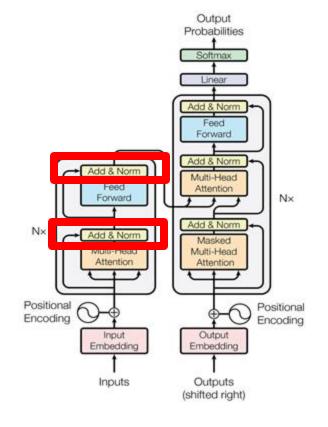




#### Trick #2: Layer Normalization [Ba et al., 2016]

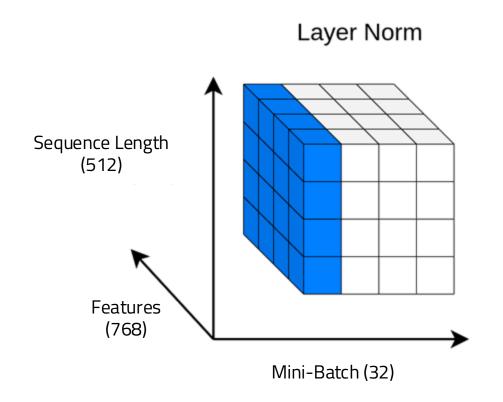
- ☐ **Problem**: Difficult to train the parameters of a given layer because its input from the layer beneath keeps shifting.
- ☐ **Solution**: Reduce uninformative variation by normalizing to *zero mean* and *standard deviation* of one within each layer.

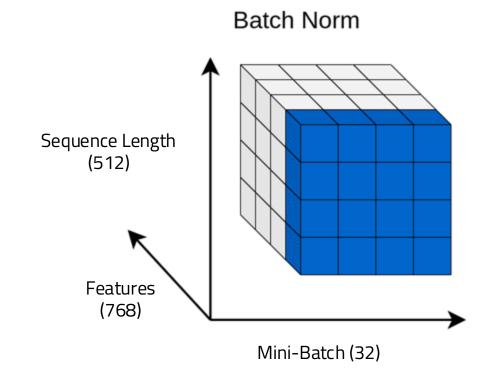
Mean: 
$$\mu^l = \frac{1}{H} \sum_{i=1}^{H} a_i^l$$
 Standard Deviation:  $\sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^{H} \left(a_i^l - \mu^l\right)^2}$ 

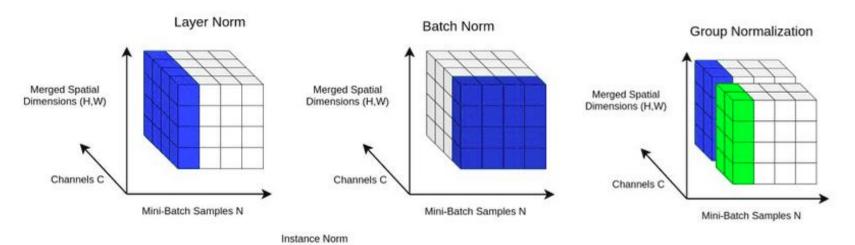


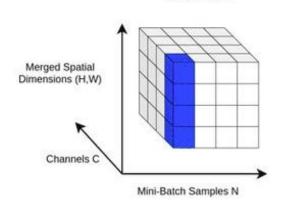
$$x^{\ell'} = \frac{x^{\ell} - \mu^{\ell}}{\sigma^{\ell} + \epsilon}$$

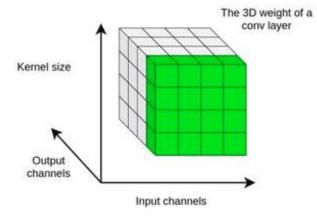
## Layer norm vs Batch norm

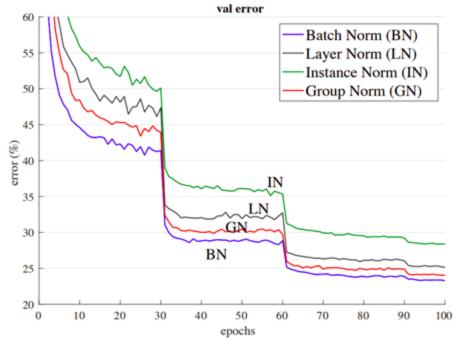








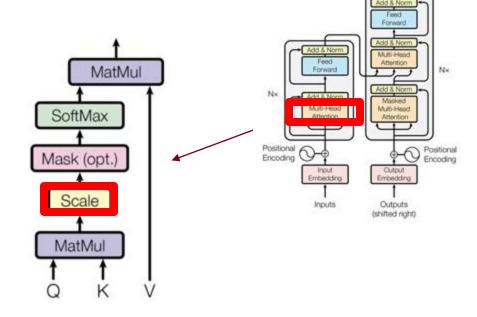




https://theaisummer.com/normalization

#### 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<sub>k</sub>



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

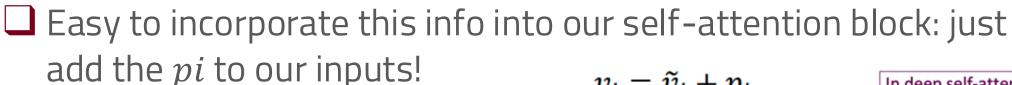


$$Output = softmax \left(QK^{T} / \sqrt{d_{k}}\right) V$$

#### 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  $i \in \{1, 2, ..., T\}$  are position vectors

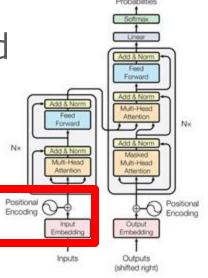


$$v_i = \tilde{v}_i + p_i$$

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

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

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

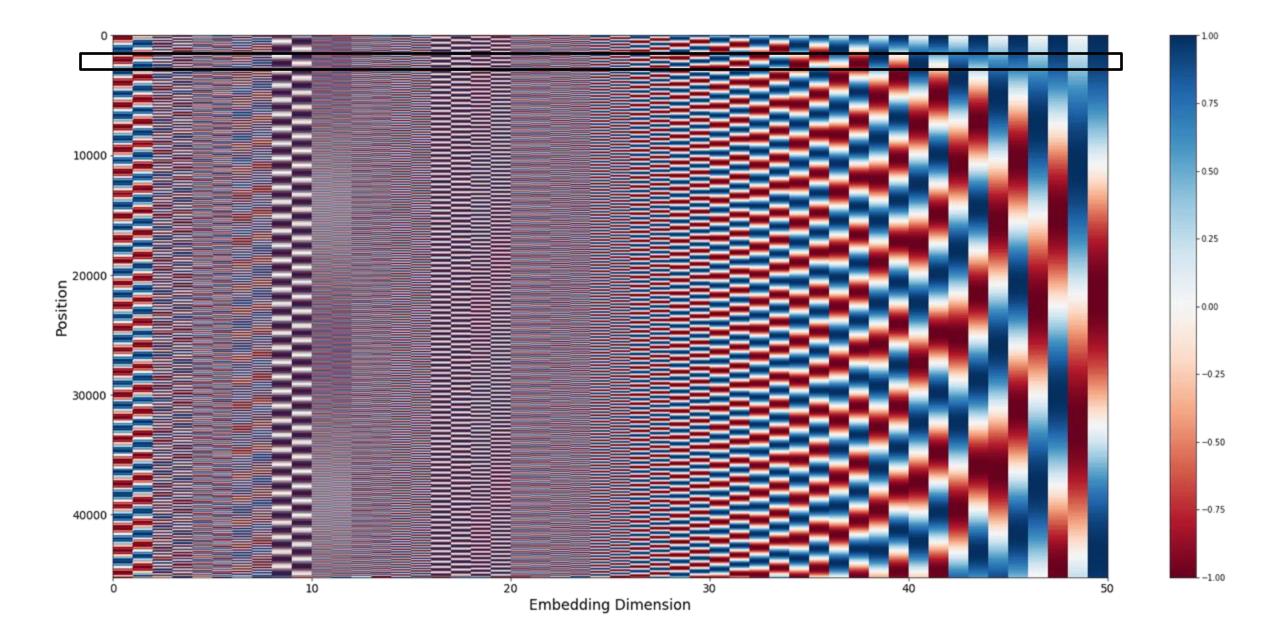


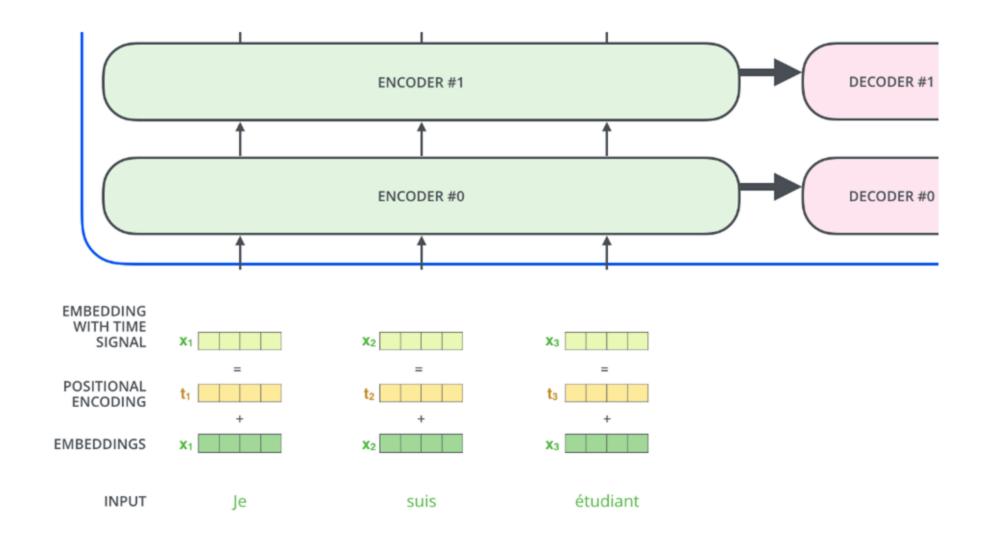
#### Position representation vectors through sinusoids

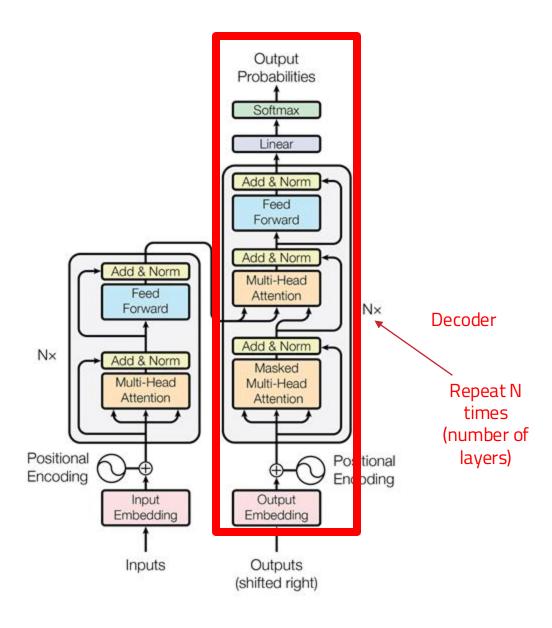
☐ Sinusoidal position representations: concatenate sinusoidal functions of varying periods:

$$\overrightarrow{p_t}^{(i)} = f(t)^{(i)} := egin{cases} \sin(\omega_k.t), & ext{if } i = 2k \ \cos(\omega_k.t), & ext{if } i = 2k + 1 \end{cases} \qquad p_i = egin{cases} \sin(i/10000^{2*1/d}) \ \cos(i/10000^{2*1/d}) \ \sin(i/10000^{2*rac{d}{2}/d}) \ \cos(i/10000^{2*rac{d}{2}/d}) \ \cos(i/10000^{2*rac{d}{2}/d}) \end{cases}$$

- ☐ Pros: Periodicity indicates that maybe "absolute position" isn't as important
- ☐ Cons: Not learnable; also the extrapolation doesn't really work

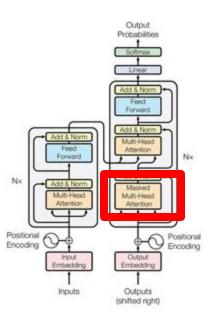






#### Decoder: Masked Multi-Head Self-Attention

☐ **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?



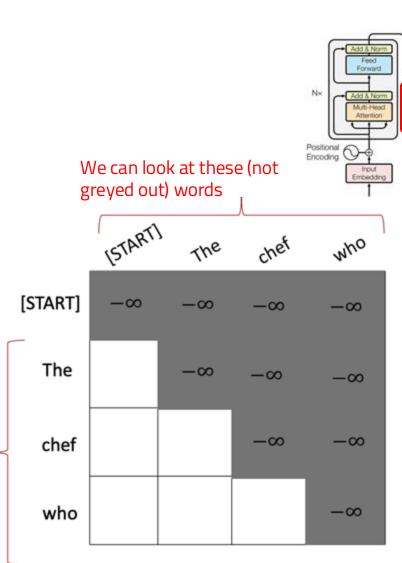
- □ Solution: Masked Multi-Head Attention.
  - At a high-level, we hide (mask) information about future tokens from the model.

# Masking the future in self-attention

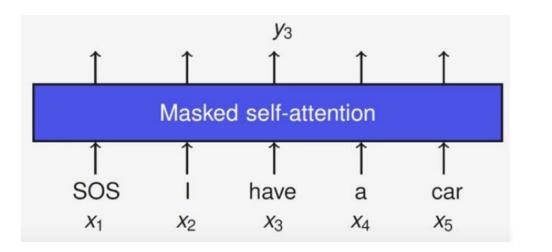
- To use self-attention in decoders, we need to ensure we can't peek at the future.
- At every timestep, we could change the set of keys and queries to include only past words. (Inefficient!)
- To enable parallelization, we mask out attention to future words by setting attention scores to -∞

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

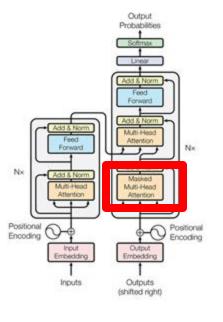
The chef For encoding these words who



## Masking the future in self-attention



$$Z_{13} = \frac{k_1^{\mathsf{T}} q_3}{\sqrt{d}}, \dots, Z_{53} = \frac{k_5^{\mathsf{T}} q_3}{\sqrt{d}}.$$

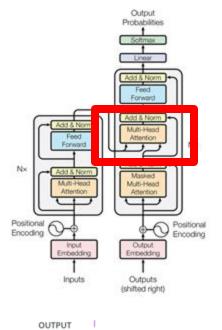


"Masked" weights:

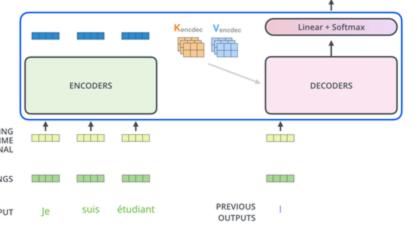
$$\begin{bmatrix} \tilde{W}_{13} \\ \tilde{W}_{23} \\ \tilde{W}_{33} \\ \tilde{W}_{43} \\ \tilde{W}_{53} \end{bmatrix} = \begin{bmatrix} \exp(Z_{13}) \\ \exp(Z_{23}) \\ \exp(Z_{33}) \\ 0 \\ 0 \end{bmatrix} \Rightarrow \begin{bmatrix} W_{13} \\ W_{23} \\ W_{33} \\ W_{43} \\ W_{53} \end{bmatrix} = \begin{bmatrix} \operatorname{softmax}(Z_{13}, Z_{23}, Z_{33}) \\ 0 \\ 0 \end{bmatrix} \Rightarrow y_3 = \sum_{i=1}^5 v_i W_{i3}$$

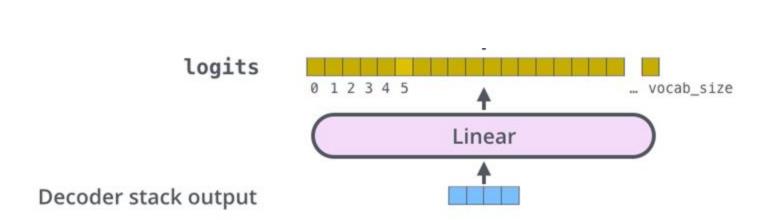
#### **Encoder-Decoder Attention**

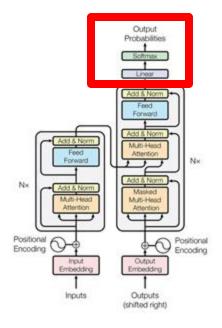
- ☐ We saw that self-attention is when keys, queries, and values come from the same source.
- ☐ In the decoder, we have attention that looks more like seq2seq with attention.
  - Let  $h_1$ ...  $h_T$  be output vectors from the Transformer encoder;  $x_i \in \mathbb{R}^T$
  - Let  $z_1$ ...  $z_T$ 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











Which word in our vocabulary is associated with this index?

am

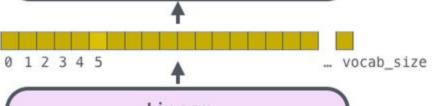
Get the index of the cell with the highest value (argmax)

5

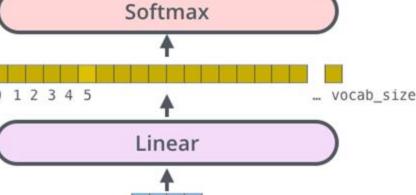
log\_probs

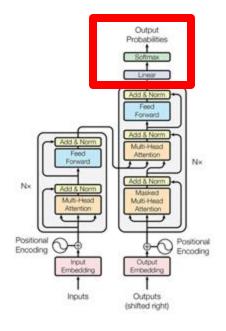


logits



Decoder stack output





#### Drawback of Transformer

#### Drawback of Transformer

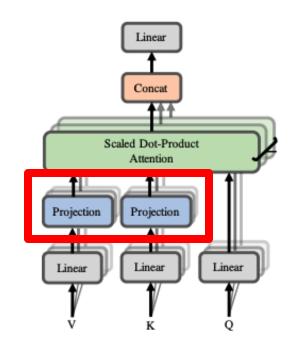
- ☐ Static positional embedding representations:
  - Are simple absolute indices the best we can do to represent position?
  - Relative linear position attention [Shaw et al., 2018]

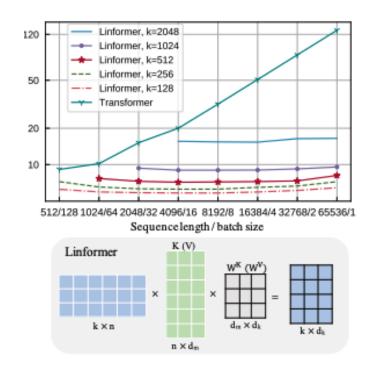
#### ☐ Quadratic compute in self-attention:

- $\circ$  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?

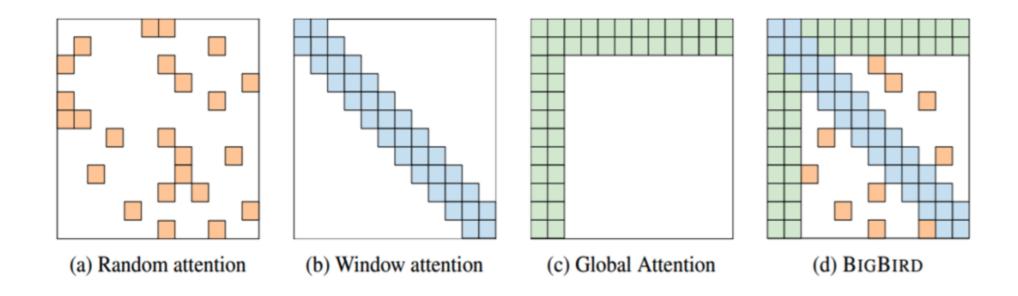
- ☐ LinFormer (Wang et al., 2020); O(T^2) -> O(T)
  - Map the sequence length dimension to a lower-dimensional space for values, keys

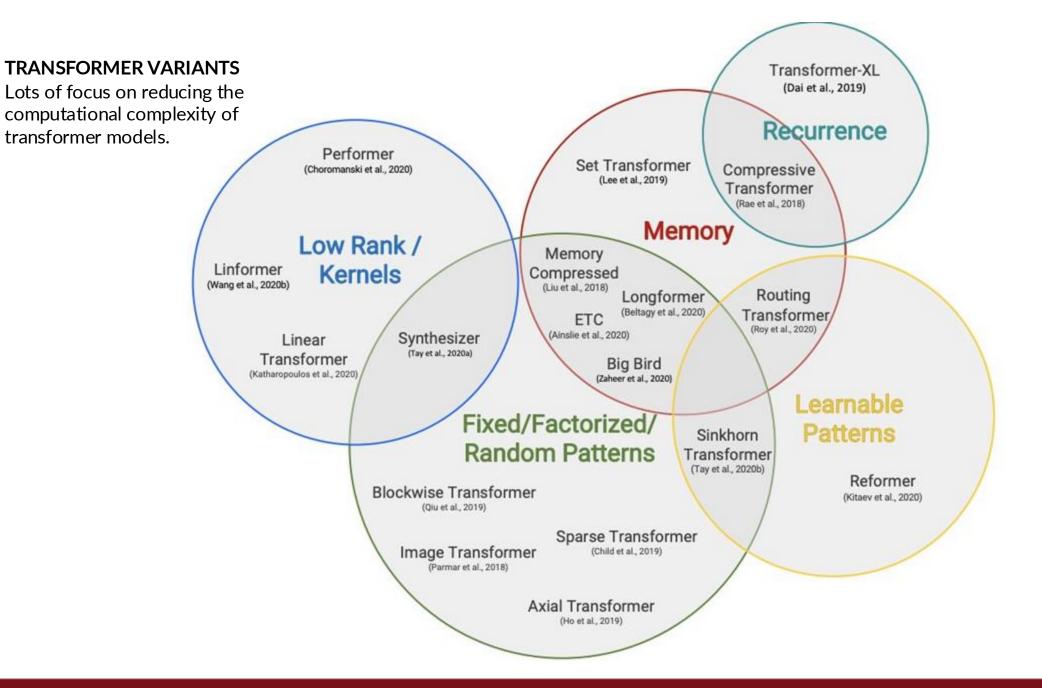




## Reduce $O(T^2)$ all-pairs self-attention cost?

- ☐ BigBird (Zaheer et al., 2021)
  - Replace all-pairs interactions with a family of other interactions, like local windows, looking at everything, and random interactions.





#### Do Transformer Modifications Transfer?

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

Model	Parama	Ops	Step/s	Early loss	First loss	SCLUE	X8um	Wild	WMT ExDe
Vanilla Transformer	223 M	11.17	3.50	$2.182 \pm 0.005$	1.636	71.66	17.79	29.02	26.62
GeLU	223 M	11.17	3.58	$2.179 \pm 0.803$	1.838	75.79	17.86	25.13	26.47
Swish	223M	11.07	3.62	$2.186 \pm 0.003$	1.947	79.77	17.74	24.34	26.75
BLU	223M	11.17	3.56	$2.270 \pm 0.807$	1.932	67.63	16.73	23.02	26.06
GLU	223M	11.07	3.59	$2.174 \pm 0.003$	1.614	74.00	17.42	24.34	27.12
GeGLU	22334	11.07	3.55	$2.130 \pm 0.006$	1.793	T1.96	18.27	24.87	36.87
BeGLU	223.M	11.17	3.07	$2.145 \pm 0.004$	1.803	76.17	18.36	34.87	37.63
SeLU	22334	11.17	3.55	$2.315 \pm 0.804$	1.948	68.76	16.76	22.75	25.99
SwiGLU	22334	11.47	3.53	$2.12T \pm 0.000$	1,789	76.00	18.20	24.34	27.62
LIGHT	223.0f	11.17	3.39	$2.149 \pm 0.005$	1,798	75.34	ST-ST	24.34	26.53
Signacid	223M	11.17	3.63	$2.291 \pm 0.019$	1.967	74.31	17.51	29.02	26.30
Softplus	223M	DAT	3.47	$2.297 \pm 0.011$	1.800	72.45	17.65	24.34	26.69
RMS Norm	223M	11.17	3.68	$2.167 \pm 0.006$	1.821	75.45	17.94	24.0T	27.14
		11.17	3.54			61.69	15.64	29.90	26.37
Bosero Bosero + LaverNorm	223M 223M	11.07	3.26	$2.262 \pm 0.008$ $2.223 \pm 0.006$	1.939	70.42	17.58	29.02	26.29
			3.34						
Benero + HMS Norm	22834	11.17		$2.221 \pm 0.000$	1.875	70.33	17.33	29.02	26.19
Fixep	2283d	11.17	2.95	$2.383 \pm 0.012$	3.007	58.56	14-42	20.02	26.31
26  layers, 4g = 1536, H = 6	224M	11.0T	3.33	$2.290 \pm 0.007$	1.943	74.00	17.75	25.18	26.69
18 layers, $4g = 3048$ , $H = 8$	2233M	DATE	3.38	$2.185 \pm 0.005$	1.651	76.45	16.83	24.34	37.30
8 lapers, $d_d = 4008, M = 18$	223M	11.17	3.69	$2.190 \pm 0.005$	1.847	74.58	17.69	23.28	36.85
6 layers, $d_{d} = 6164, H = 24$	2230I	11.17	3.79	$2.201 \pm 0.010$	1.807	73.55	17.59	24.60	26.66
Block sharing	65.M	DAT	3.91	$2.497 \pm 0.687$	2.164	64.50	14.53	21.96	20.46
+ Partorised embeddings	45.54	9.4T	4.71	$2.631 \pm 0.305$	2.183	60.84	14.00	19.84	20.27
+ Factorized & shared en-	20.0/	9.1T	4.37	$2.907 \pm 0.313$	2.385	53.95	11.37	19.84	25.19
beddings						-			
Enruder only block sharing	170M	11.17	3.68	$2.298 \pm 0.023$	1.929	69.60	16.23	23.02	26.23
Decader only block sharing	144M	11.17	3.79	$2.352 \pm 0.029$	2.092	6T.93	16.13	23.81	26.09
Factorized Embedding	22TM	9-4T	3.80	$2.298 \pm 0.006$	1.855	70.41	15.92	22.75	26.50
Factorized & shared embed-	292M	8.1T	3.92	$2.320 \pm 0.010$	1.902	69.69	16.33	22.22	26.44
dinge									
Tied encoder/decoder in-	246M	11.1T	3.55	$2.192 \pm 0.002$	1.940	71.79	17.72	24.34	26.49
pet enholdings									
Tied decoder input and out-	248M	11.17	3.07	$2.187 \pm 0.007$	1.827	74.86	17.74	34.87	36.67
pat embeddings									
Untied embeddings	273 M	11.17	3.53	$2.195 \pm 0.805$	1.834	72.99	17.58	23.28	26.48
Adaptive input embeddings	204M	8:2T	3.55	$2.250 \pm 0.002$	1.899	66.5T	16.20	24.0T	26.66
Adaptive soltmax	2043/	9:2T	3.60	$2.364 \pm 0.005$	1.982	72.91	16.67	25.16	25.56
Adaptive softmax without	223M	10.8T	3.43	$2.229 \pm 0.009$	1.914	71.82	17.10	23.02	25.72
projection									
Misture of softmanes	292M	16.37	2.26	$2.227 \pm 0.017$	1.821	76.77	17.62	22.75	26.62
	22334	11.17	3.33	2.181 ± 0.814	1.871	54.31	10.40	25.16	26.80
Transparent attention									
Dynamic convolution	25TM	11.8T	2.65	$2.403 \pm 0.009$	2.047	58.30	12.67	25.16	17.68
Lightweight convolution	224M	10.4T	4.07	$2.370 \pm 0.010$	1.989	63.07	14.56	23.02	24.73
Evolved Transformer	21/TM	9:9T	3.09	$2.220 \pm 0.803$	1.963	73.67	10.76	24.0T	26.56
Syntholiser (dense)	224M	HAT	3.47	$2.334 \pm 0.021$	1.962	61.68	14.27	16.14	26.63
Synthesiser (deuse plus)	243M	DAT	3.22	$2.191 \pm 0.010$	1.840	T3.96	16.96	23.81	36.71
Synthesizer (dense plus al-	24334	13.67	3.00	$2.180 \pm 0.007$	1.828	74.35	17.03	23.28	26.61
pha)									
Synthesizer (Tactorized)	28TM	10.1T	3.94	$2.341 \pm 0.017$	1.968	62.78	15.39	23.55	26.42
Synthesizer (random)	254M	10.1T	4.08	$2.326 \pm 0.012$	2:809	54.27	18.35	19.56	26.44
Synthesizer (random plue)	292M	12.87	3.63	$2.189 \pm 0.804$	1.842	73.32	17.04	24.8T	26.43
Synthosiser (random plus	29234	12.60	3.42	$2.186 \pm 0.807$	1.626	75.34	17.08	24.06	26.39
(lpha)									
Universal Transformer	8434	40.ET	0.88	$2.406 \pm 0.636$	2:053	70.13	14.09	19.05	23.91
Misture of experts	64834	11.77	3.30	$2.148 \pm 0.006$	1.780	74.55	18.13	24.08	36.94
Switch Transformer	1180M	11.7T	3.18	$2.135 \pm 0.807$	1.758	75.38	18.02	26.19	26.81
Funnel Transferore	223M	1.9T	4.30	$2.298 \pm 0.008$	1.919	67.34	16.26	22.75	23.20
Weighted Transfermer	2900M	71.8F	0.59	$2.376 \pm 0.821$	1.989	69.04	16.96	23.02	26.30
Product key memory	421.M	396.67	0.25	$2.155 \pm 0.000$	1,796	75.16	17.04	23.55	26.73

### Do Transformer Modifications Transfer Across Implementations and Applications?

Sharan Narang*	Hyung Won Chung	$\mathbf{Yi} \; \mathbf{Tay}$	William Fedus
Thibault Fevry <sup>†</sup>	${\bf Michael~Matena}^{\dagger}$	Karishma Malkan $^{\dagger}$	Noah Fiedel
Noam Shazeer	${\bf Zhenzhong}{\bf Lan}^{\dagger}$	Yanqi Zhou	Wei Li
Nan Ding	Jake Marcus	Adam Roberts	Colin Raffel $^{\dagger}$

# 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)

# 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

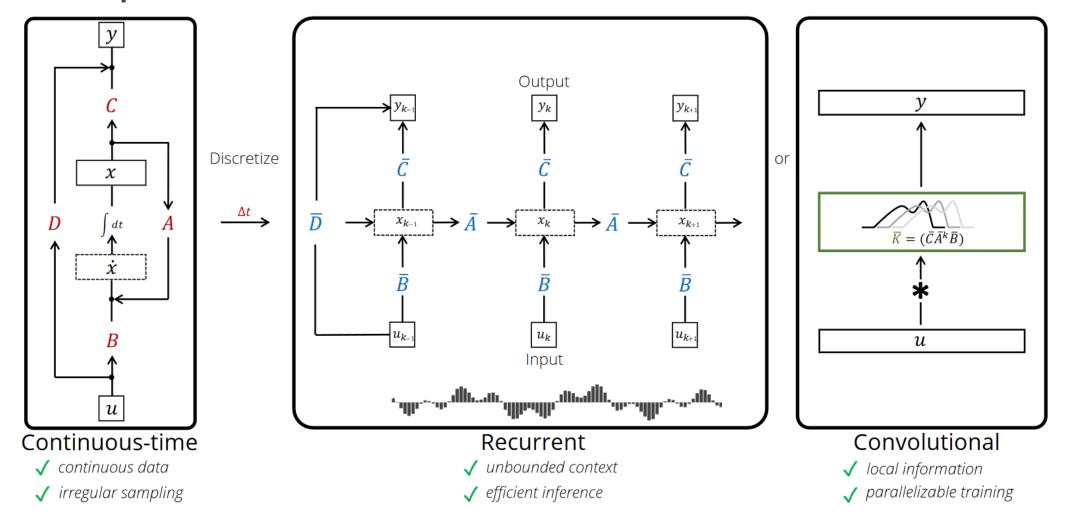
Model	Layers	Width	Heads	Params	Data	Training	
Transformer-Base	12	512	8	65M		8x P100 (12 hrs)	
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)	
BERT-Base	12	768	12	110M	13GB		
BERT-Large	24	1024	16	340M	13GB		
XLNet-Large	24	1024	16	340M	126GB	512x TPU-v3 (2.5 days)	
RoBERTa	24	1024	16	355M	160GB	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.

### State Space Models



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