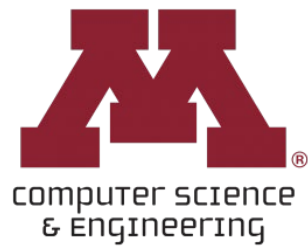# CSCI 5541: Natural Language Processing

**Lecture 10: Deep Dive on Transformers**

Dongyeop Kang (DK), University of Minnesota

dongyeop@umn.edu | twitter.com/dongyeopkang | dykang.github.io

Courtesy of Paramount Pictures

Using some slides borrowed from Anna Goldie (Google Brain) and John Hweitt (Stanford)
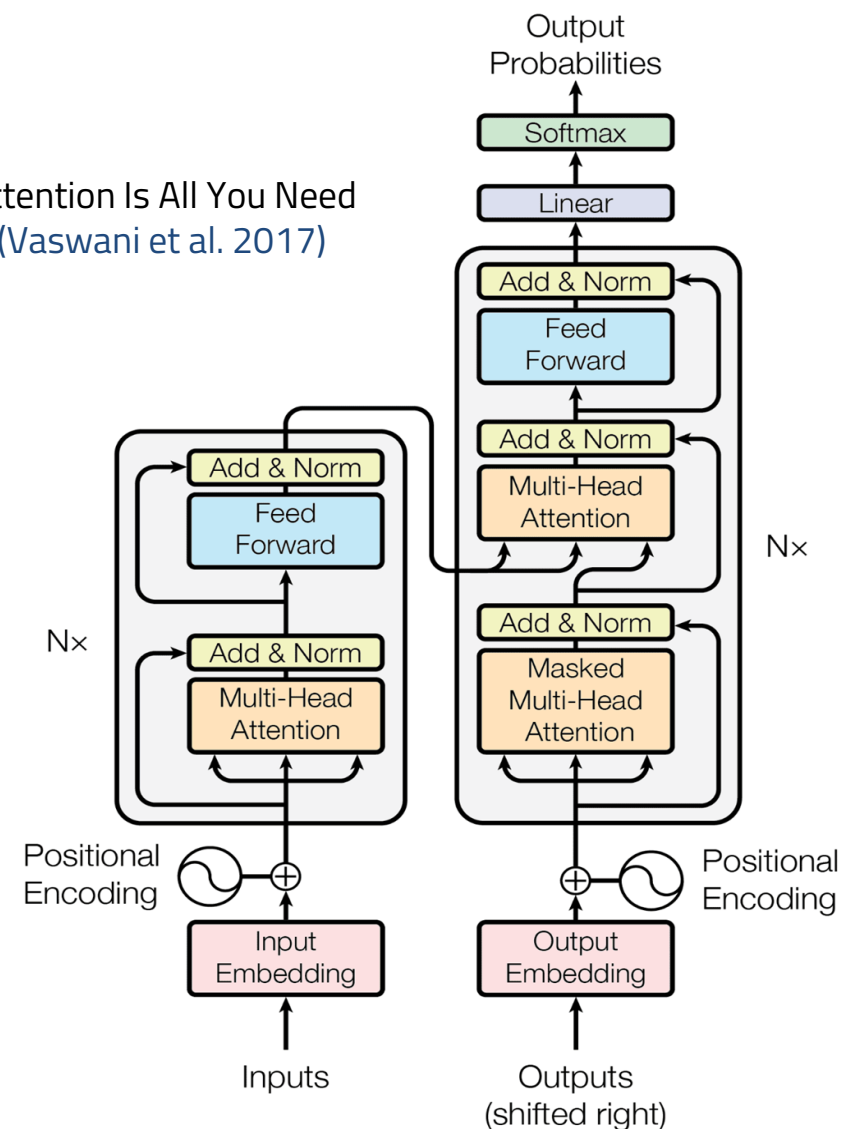
# Summary of Transformers

❑ A sequence-to-sequence model based entirely on attention

❑ Strong results on translation and a wide variety of other tasks

❑ Fast: only matrix multiplications

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

[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

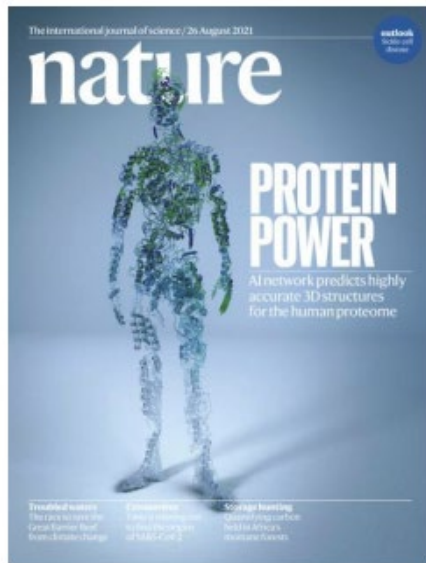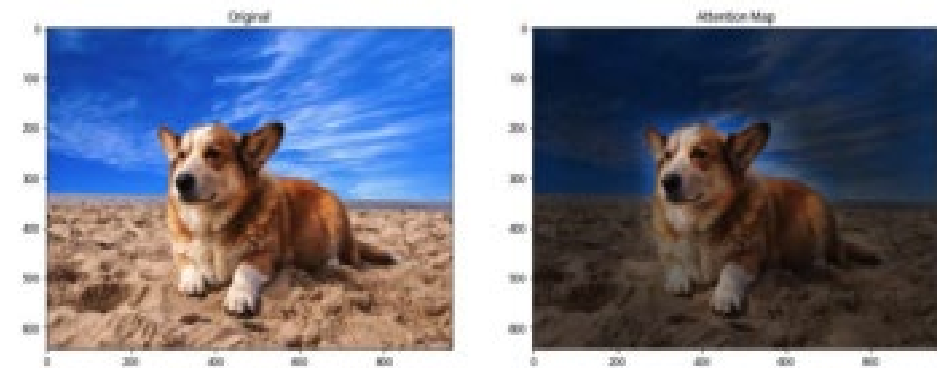| Rank | Model | Accuracy↑ | Paper | Code | Result | Year | Tags ✎ |
|------|-------|-----------|-------|------|--------|------|--------|
| 1 | SMART-RoBERTa Large | 97.5 | SMART: Robust and Efficient Fine-Tuning for Pre-trained Natural Language Models through Principled Regularized Optimization | ○ | → | 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
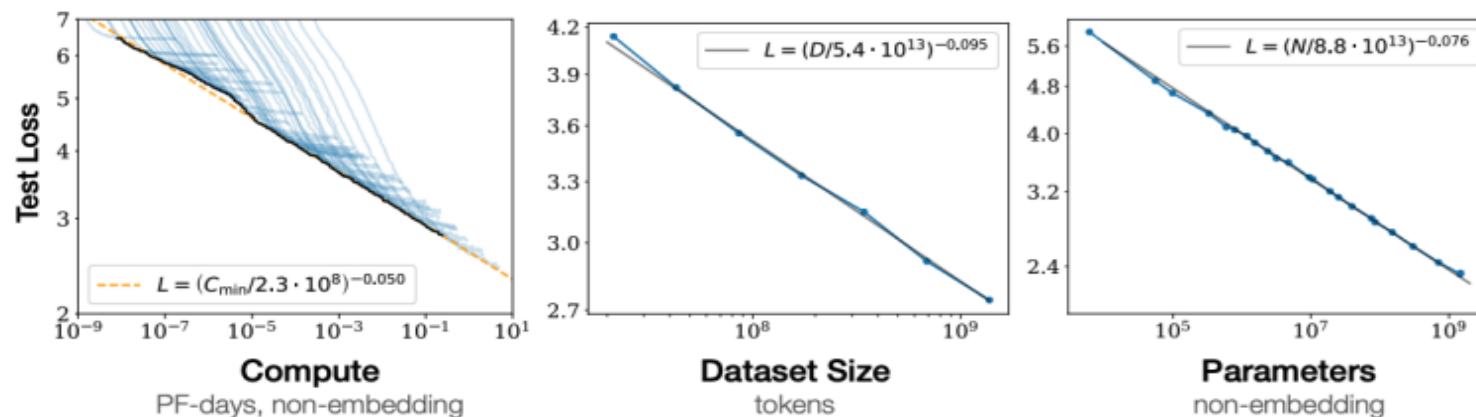


Alpha?Fold2 (Jumper et al., 2021)

## 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!

❑ If we keep scaling up these models (with no change to the architecture), could they eventually match or exceed human-level performance?
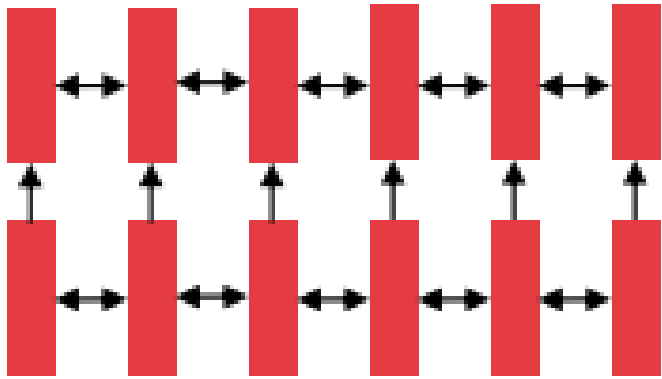


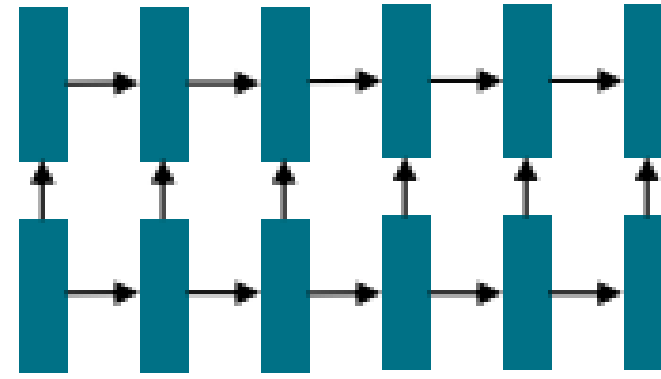Kaplan et al., 2020, *Scaling Laws for Neural Language Models*

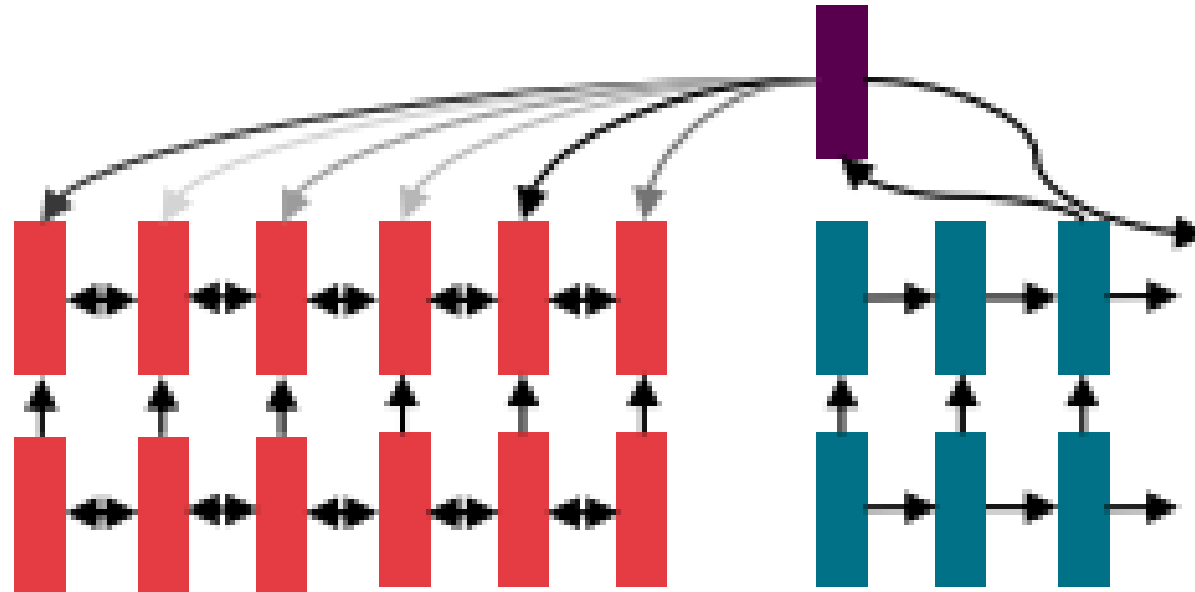# Why self-attention?

# 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.
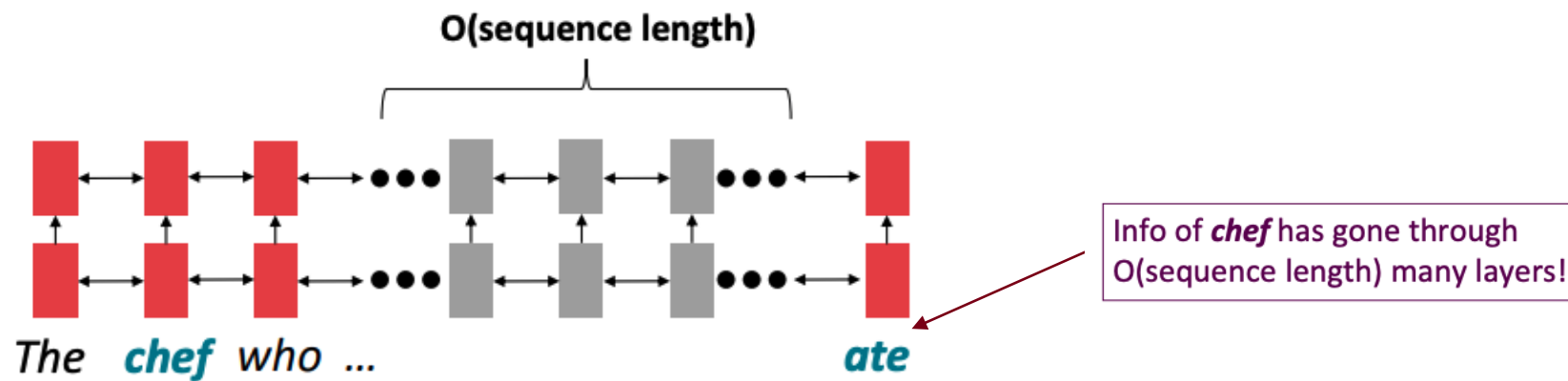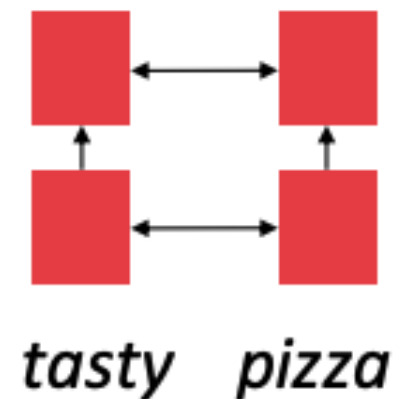
# Sequence-to-sequence with attention



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

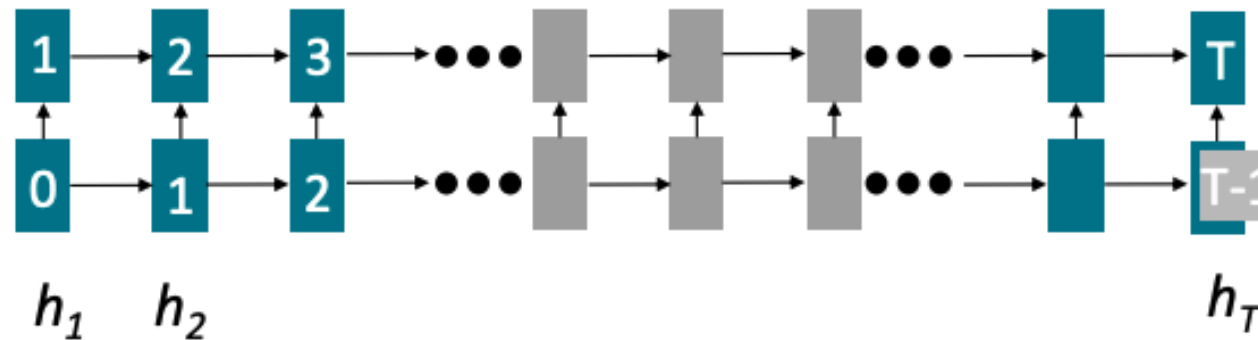# Issues with recurrent models: **Linear interaction distance**

❑ Forward RNNs are unrolled "left-to-right".

❑ It encodes linear locality:

  o Nearby words often affect each other's meanings

❑ **Problem**: RNNs take **O(sequence length) steps** for distant word pairs to interact



tasty   pizza



O(sequence length)

The   *chef*   who  …                                      ate

Info of *chef* has gone through
O(sequence length) many layers!

# 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!
- o But future RNN hidden states can't be computed fully before past RNN hidden states have been computed
- o Particularly problematic as sequence length increases, as we can no longer batch many examples together due to memory limitations
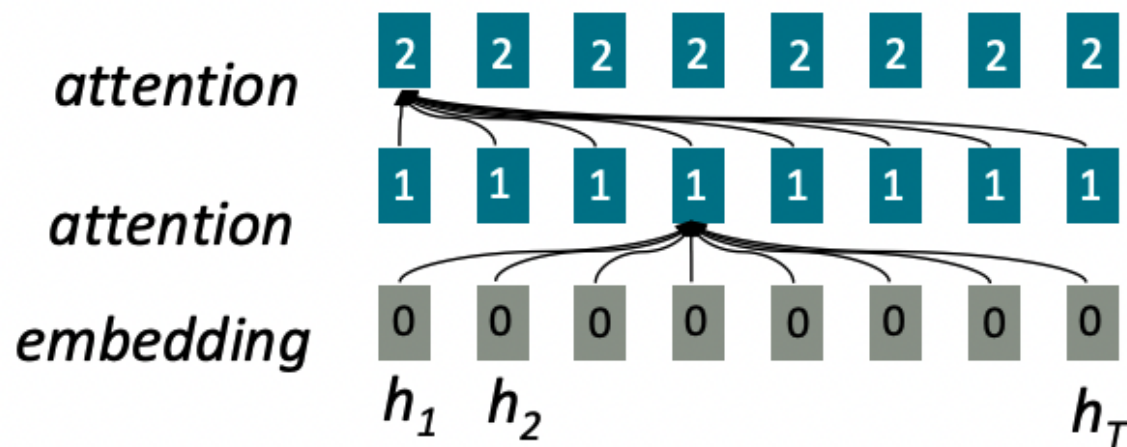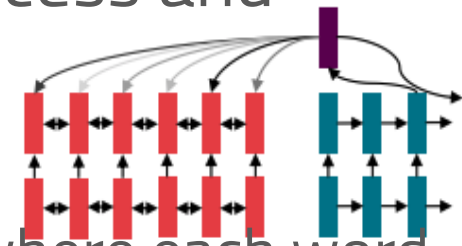
Numbers indicate min # of steps before a state can be computed

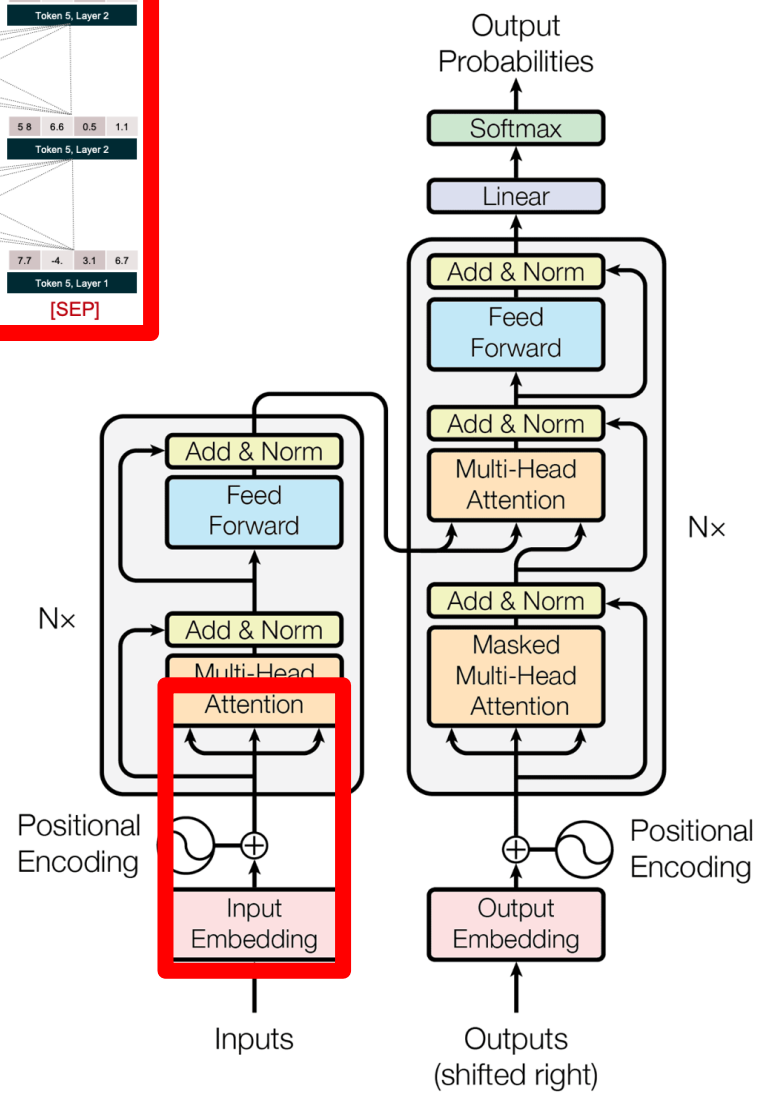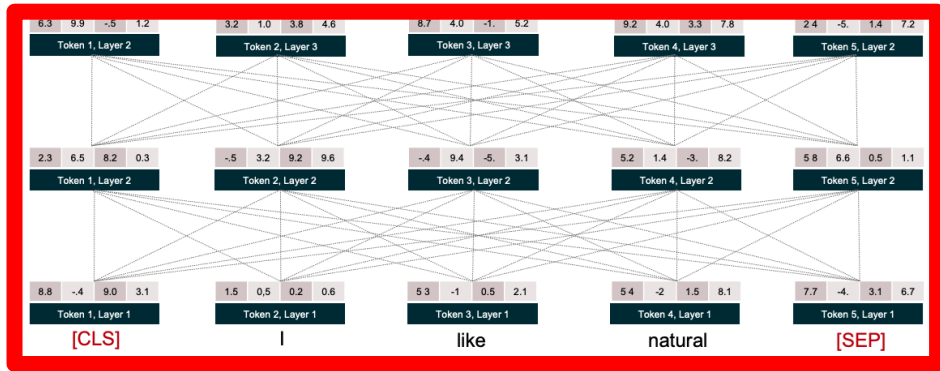# If not recurrence, then what? How about (self) attention?

❑ **Attention** treats each word's representation as a **query** to access and incorporate information from a **set of values**.
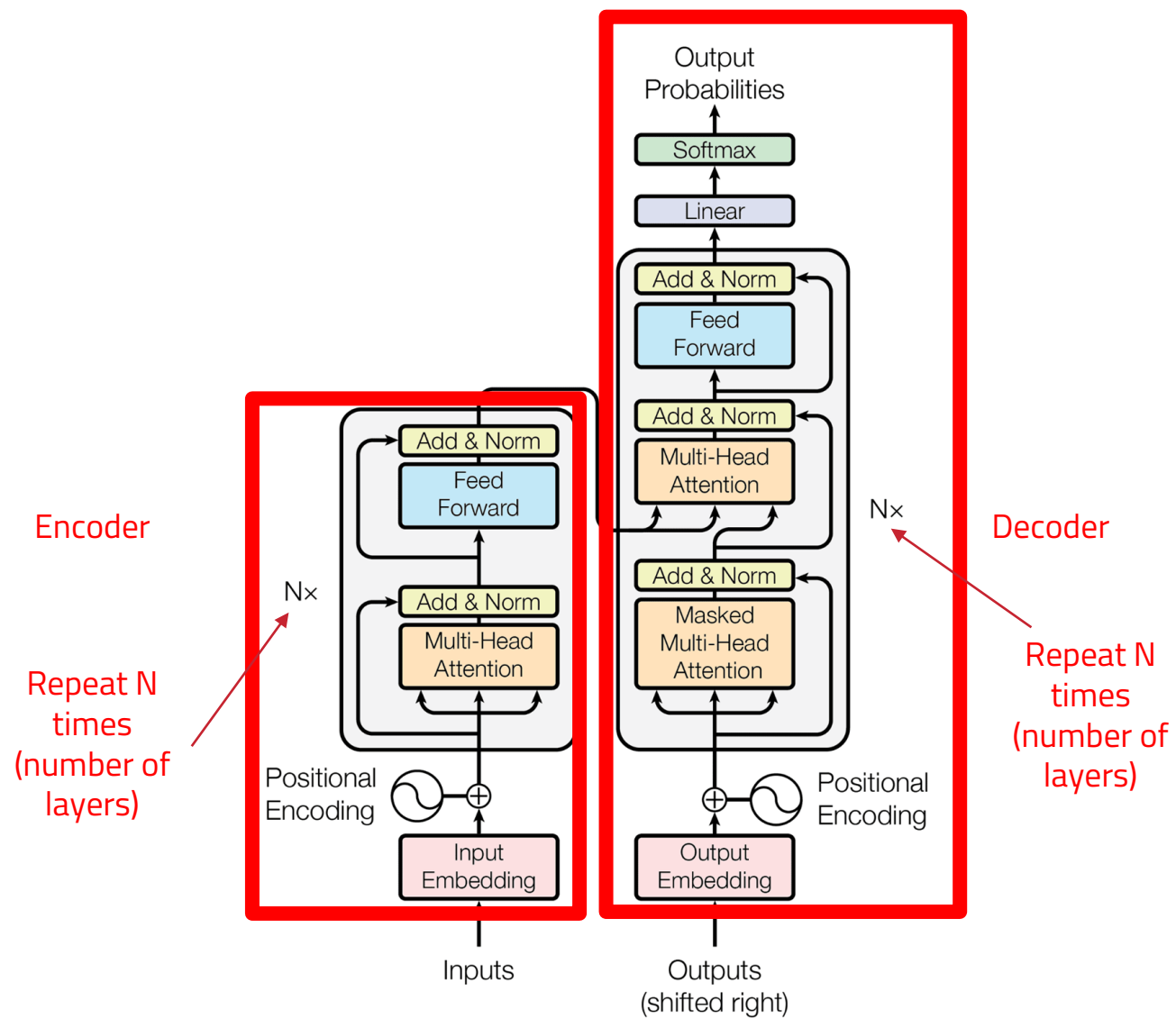
- ○ We saw attention from the decoder to the encoder;
- ○ Self-attention is encoder-encoder (or decoder-decoder) attention where each word attends to each other word within the input (or output).



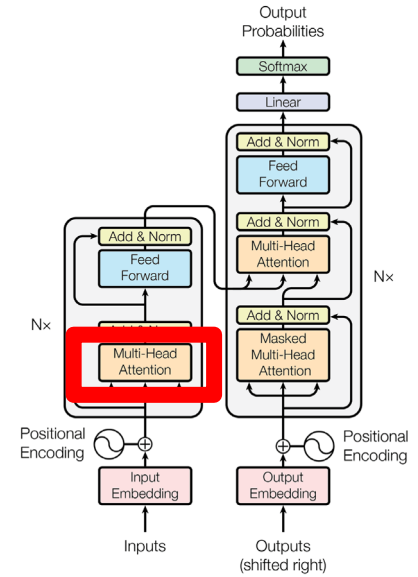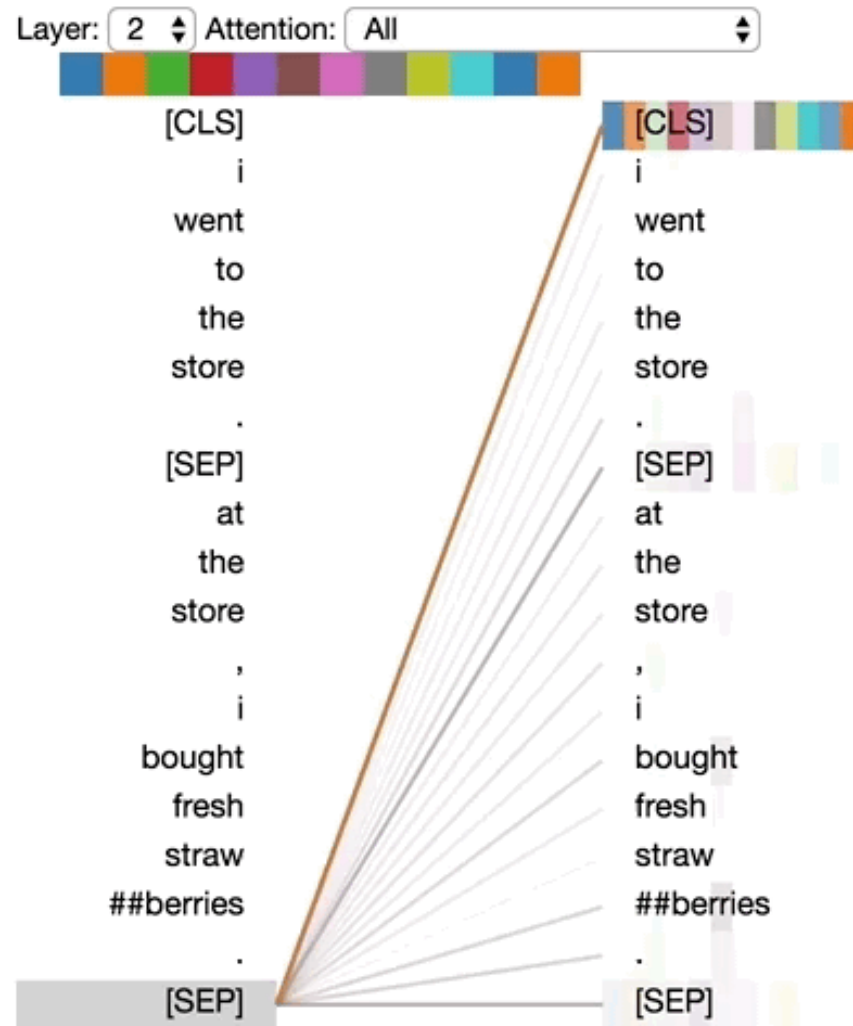All words attend to all words in previous layer; most arrows are omitted
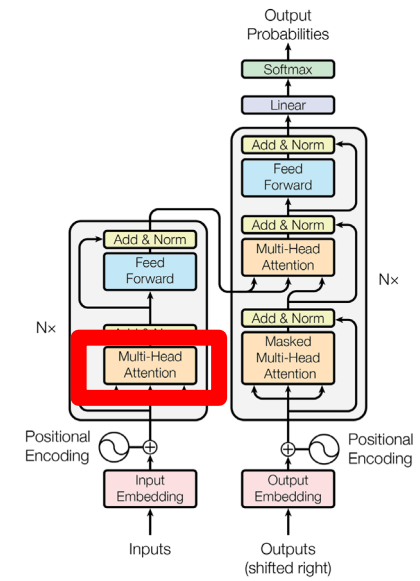
"I went to the store. At the store, I bought fresh strawberries."

https://github.com/jessevig/bertviz
https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb

**Heads**



https://github.com/jessevig/bertviz
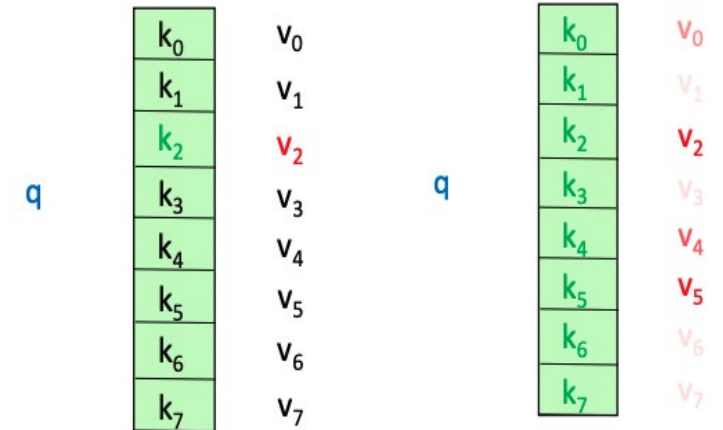
https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb

# Encoder: Self-Attention

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

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

# Recipe for Self-Attention in the Transformer Encoder

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

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$
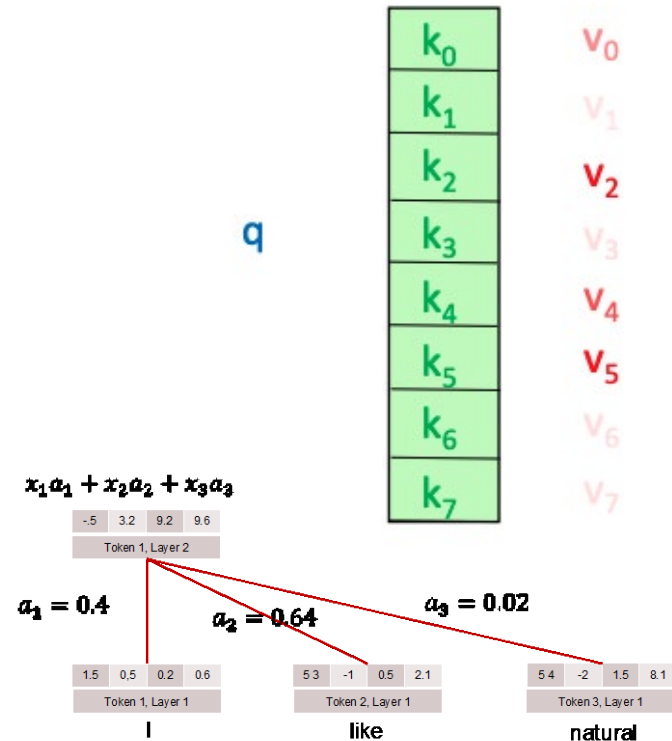
❑ Step 2: Calculate attention score between query and keys.

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

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

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

❑ Step 4: Take a weighted sum of values.

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



$x_1 a_1 + x_2 a_2 + x_3 a_3$

| -5 | 3.2 | 9.2 | 9.6 |

Token 1, Layer 2

$a_1 = 0.4$   $a_2 = 0.64$   $a_3 = 0.02$

| 1.5 | 0.5 | 0.2 | 0.6 |

Token 1, Layer 1

| 5 3 | -1 | 0.5 | 2.1 |

Token 2, Layer 1

| 5 4 | -2 | 1.5 | 8.1 |

Token 3, Layer 1

I     like     natural

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

❏ Step 1: For each word , 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.

$$E = QK^T$$

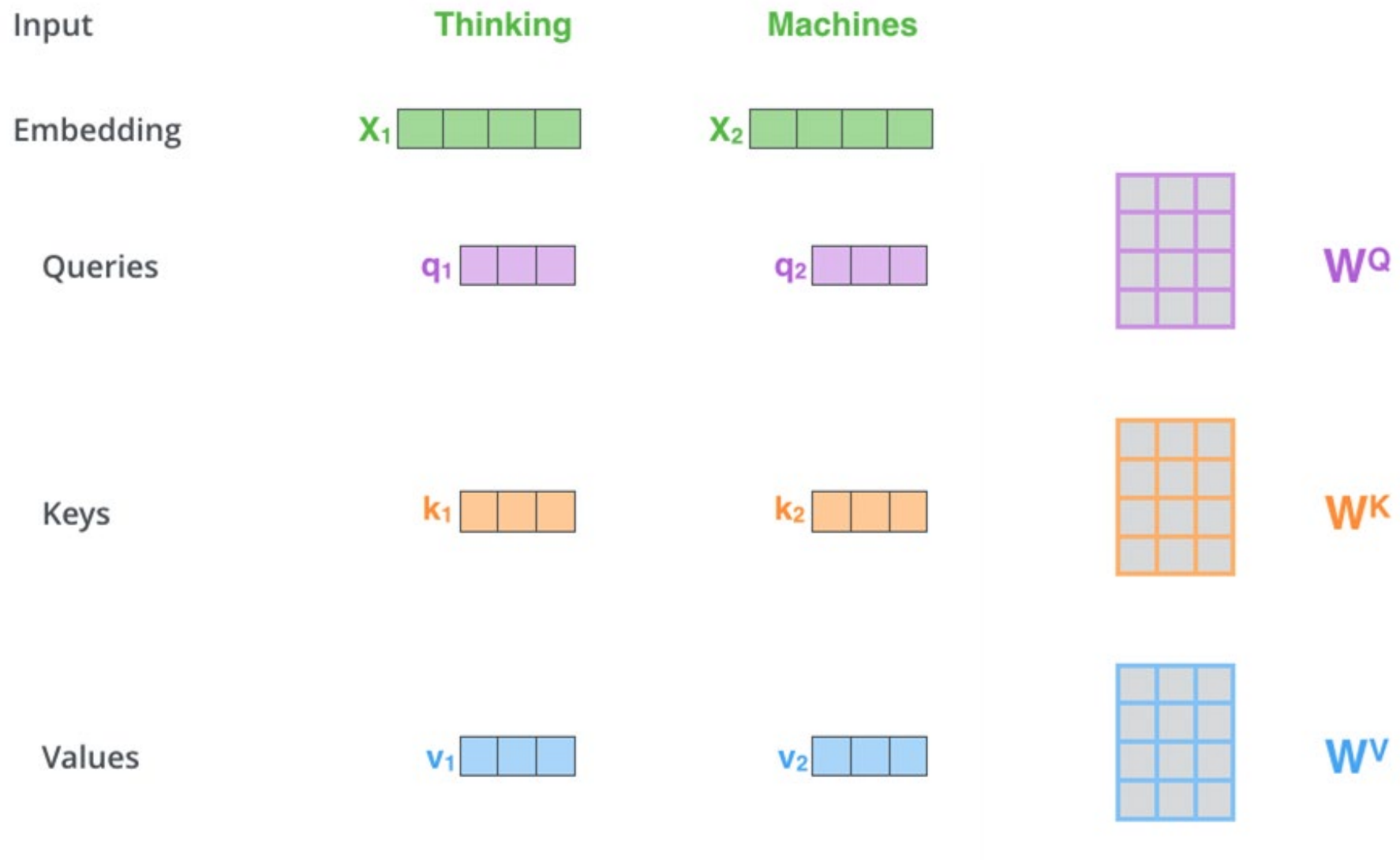❏ Step 3: Take the softmax to normalize attention scores.

$$A = softmax(E)$$

❏ Step 4: Take a weighted sum of values.

$$Output = AV$$

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

Input     **Thinking**     **Machines**

Embedding   $X_1$    $X_2$

Queries   $q_1$    $q_2$    $W^Q$

Keys   $k_1$    $k_2$    $W^K$

Values   $v_1$    $v_2$    $W^V$

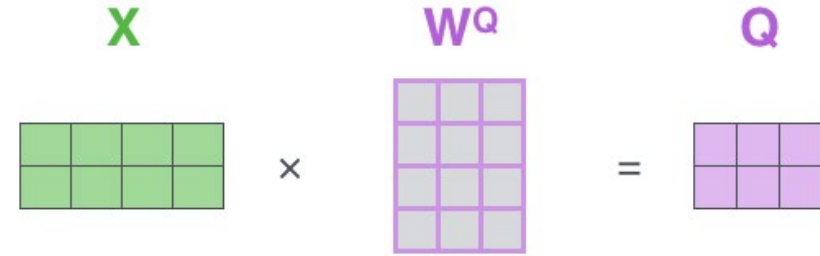https://jalammar.github.io/illustrated-transformer/

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

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

❑ Step 2: Calculate attention score between query and keys.
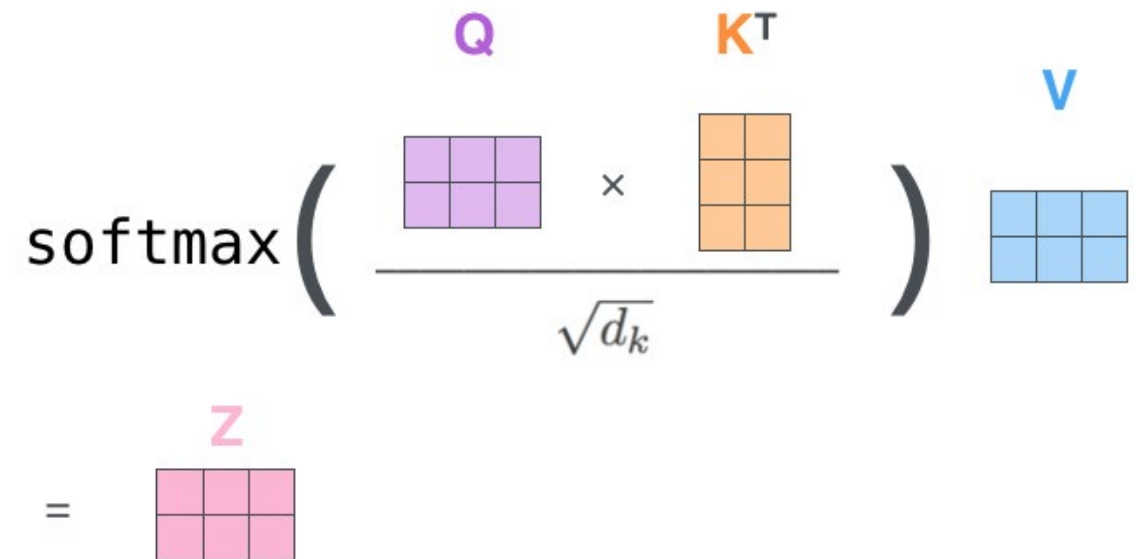
$$E = QK^T$$

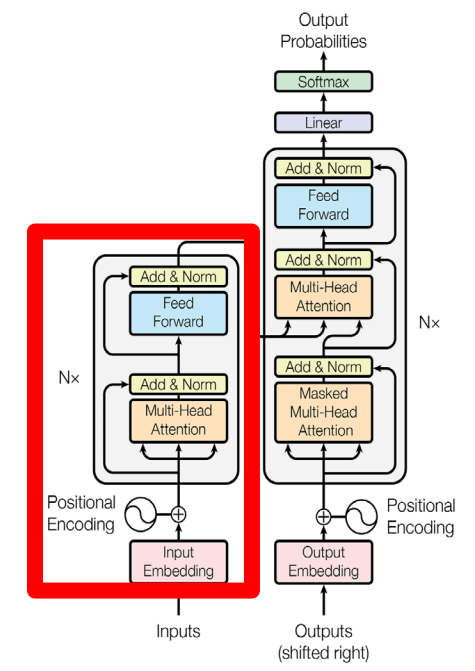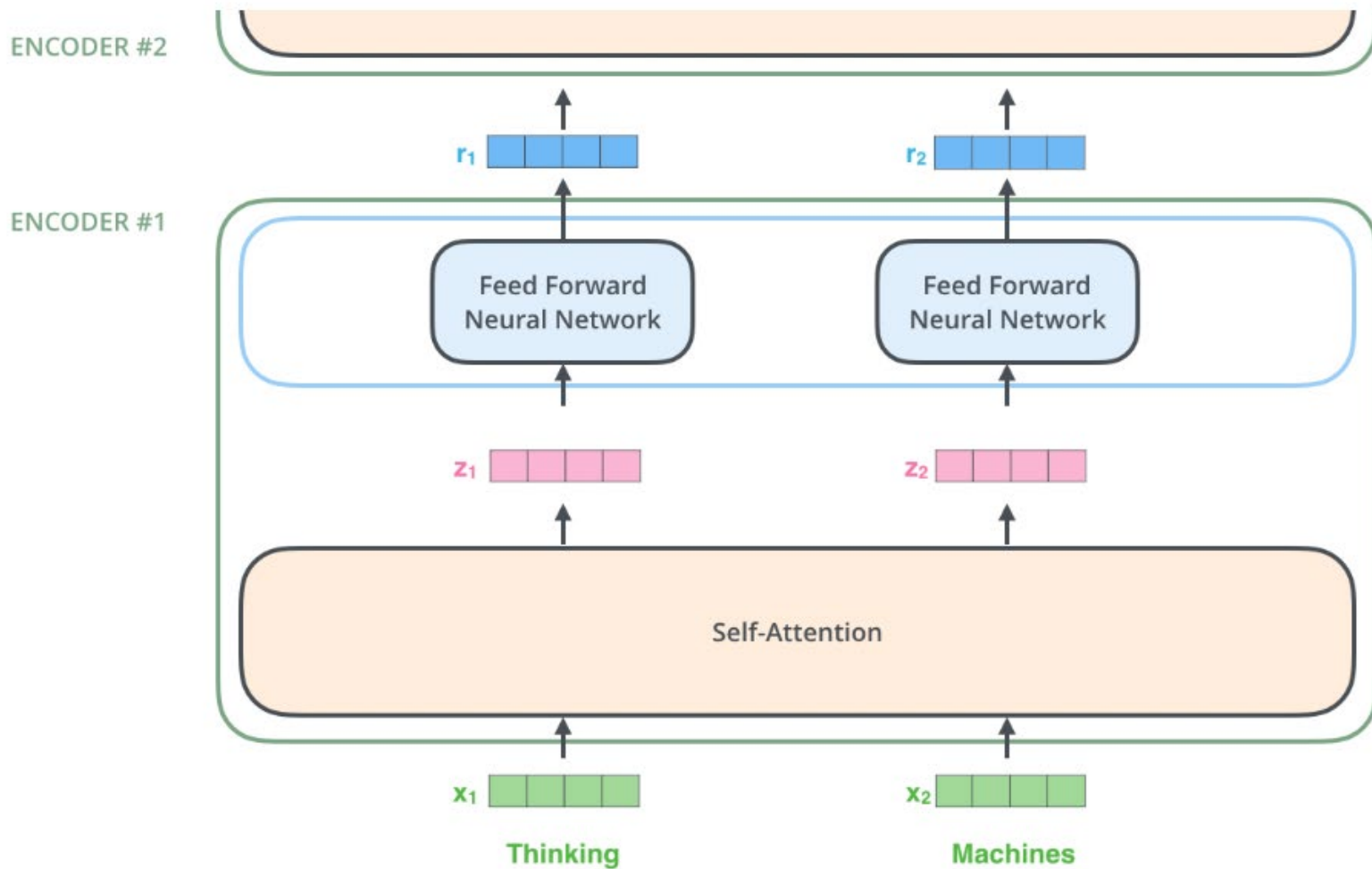❑ Step 3: Take the softmax to normalize attention scores.

$$A = softmax(E)$$

❑ Step 4: Take a weighted sum of values.

$$Output = AV$$

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

ENCODER #2

ENCODER #1

$r_1$

$r_2$

Feed Forward Neural Network

Feed Forward Neural Network

$z_1$

$z_2$

Self-Attention

$x_1$ Thinking

$x_2$ Machines

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Add & Norm

Feed Forward

Multi-Head Attention

Add & Norm

Add & Norm

Multi-Head Attention

Masked Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

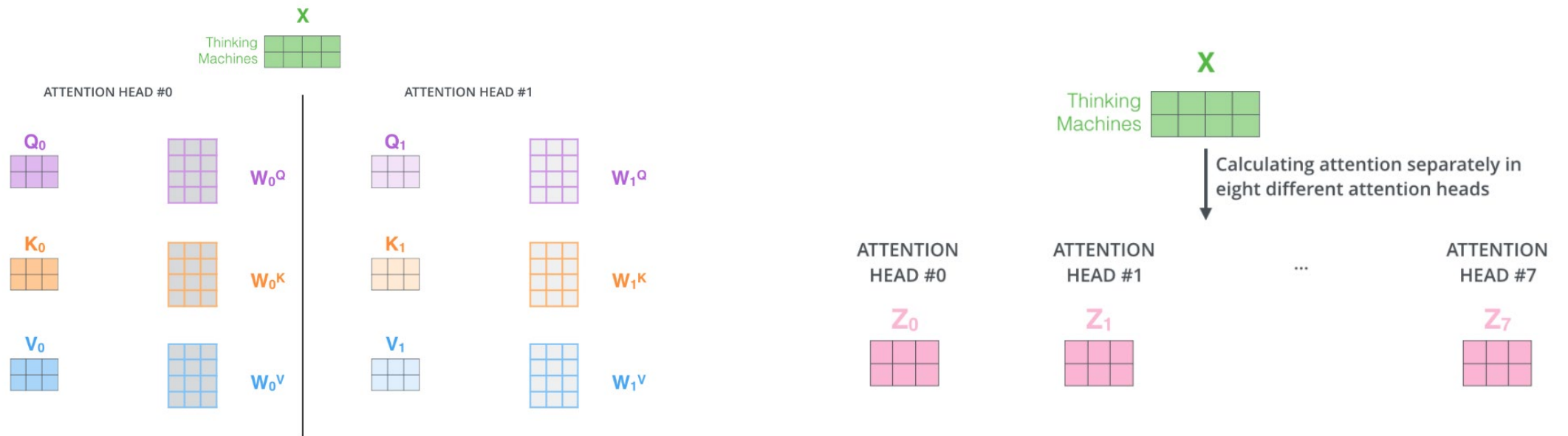https://jalammar.github.io/illustrated-transformer/

# 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.

# Condensing multi-head attentions into a single matrix
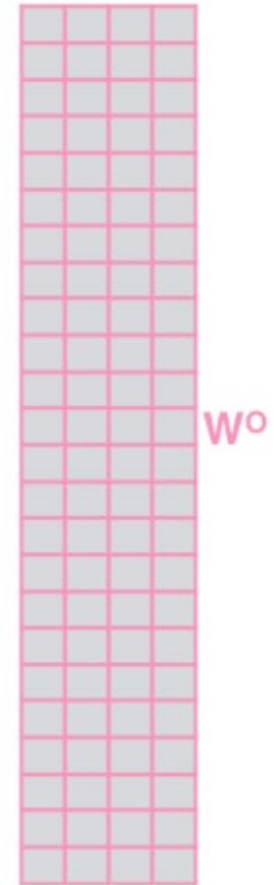
1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

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

X

$W^o$

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

$Z$

=

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

https://github.com/jessevig/bertviz
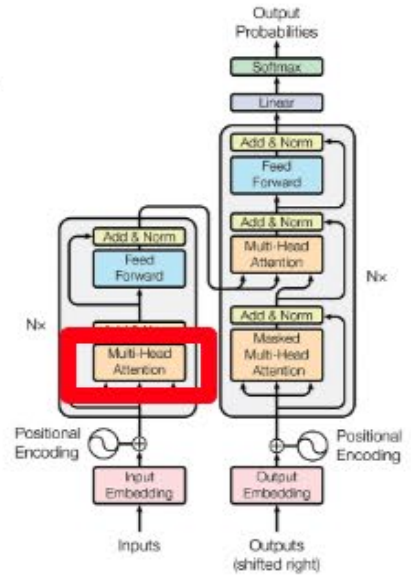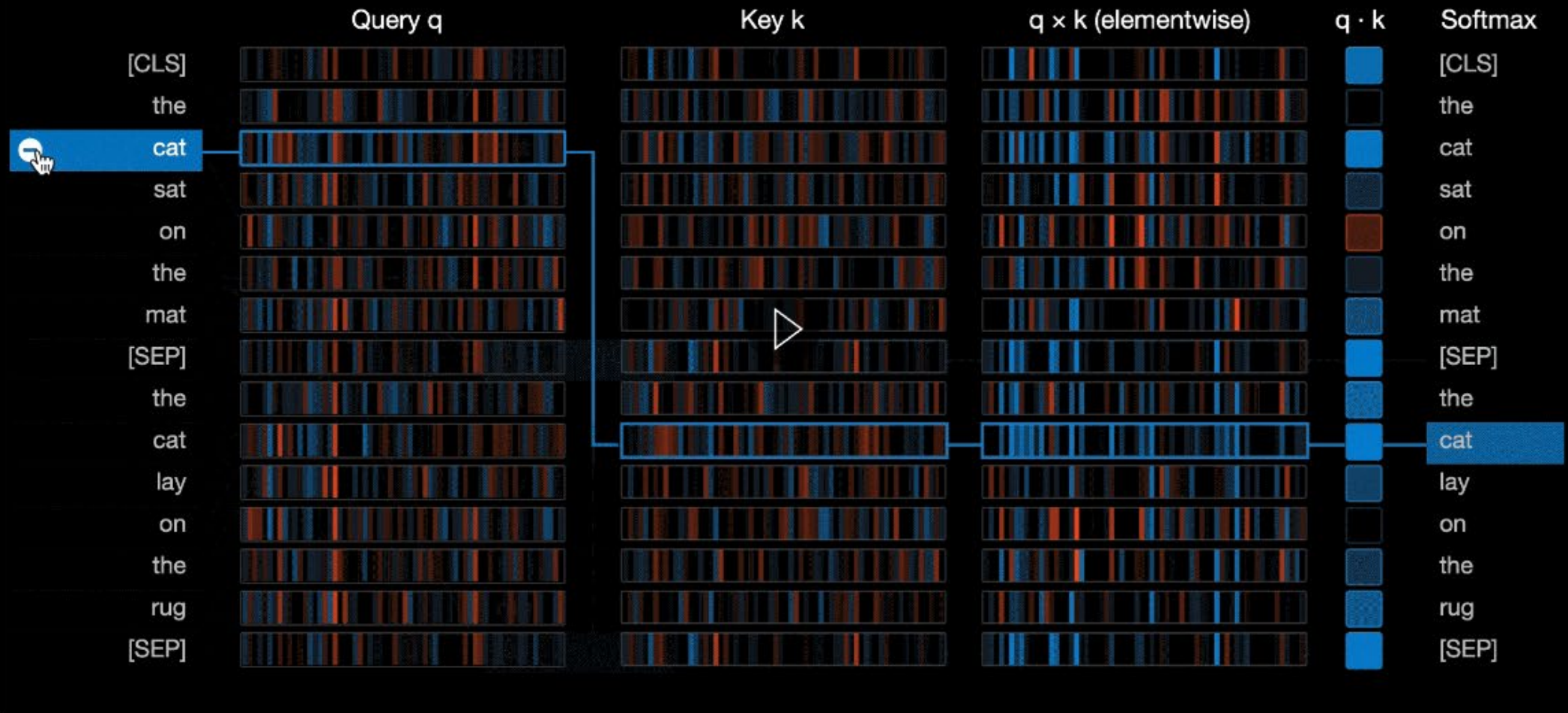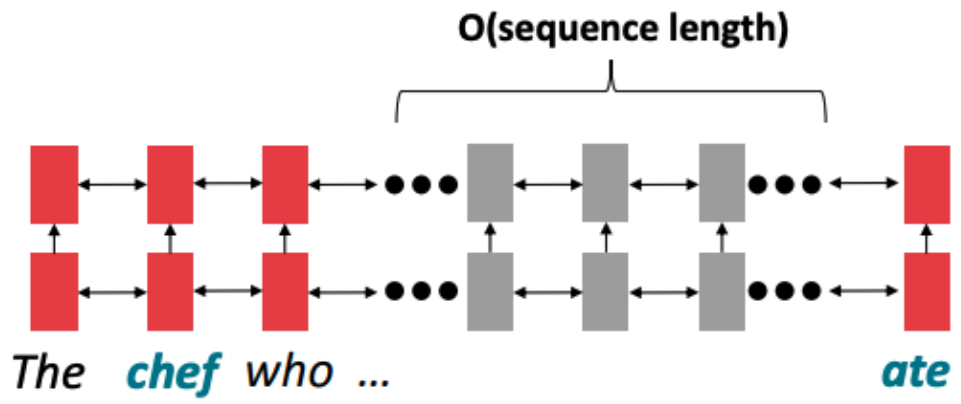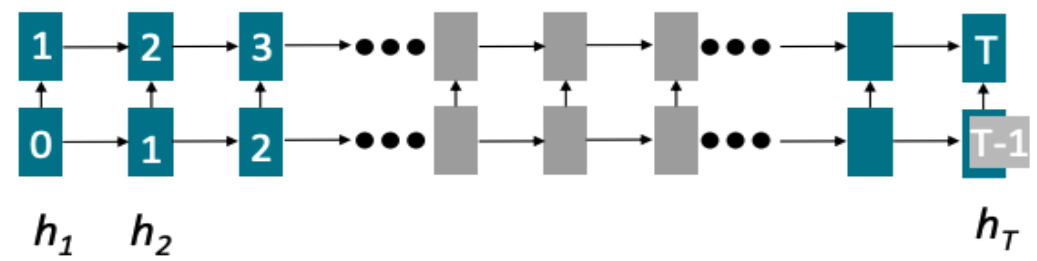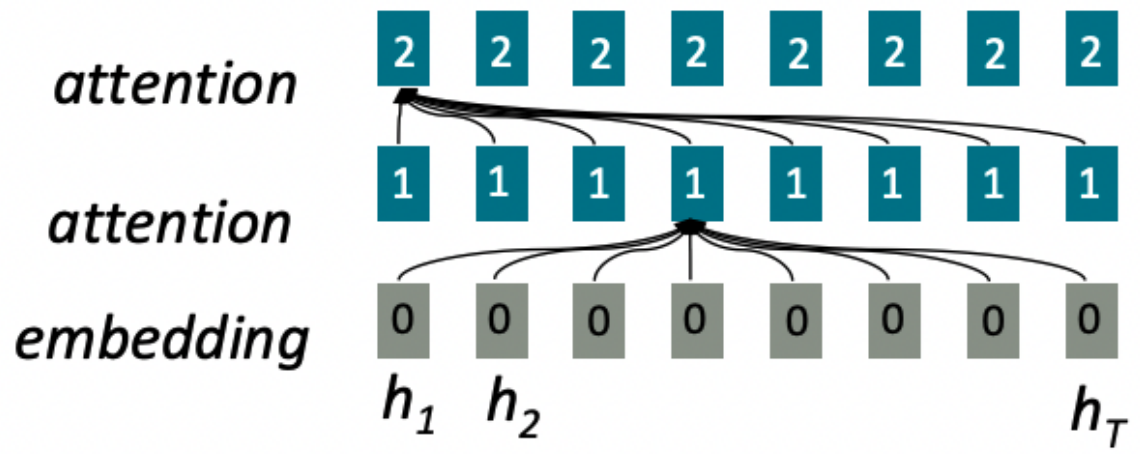https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb
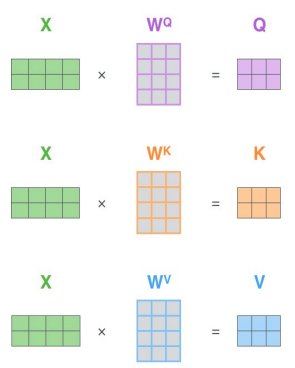
Linear interaction distance: O(T)
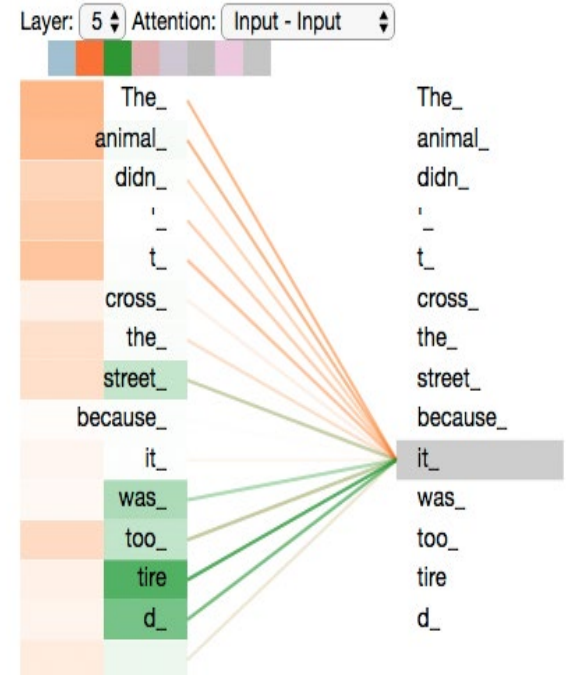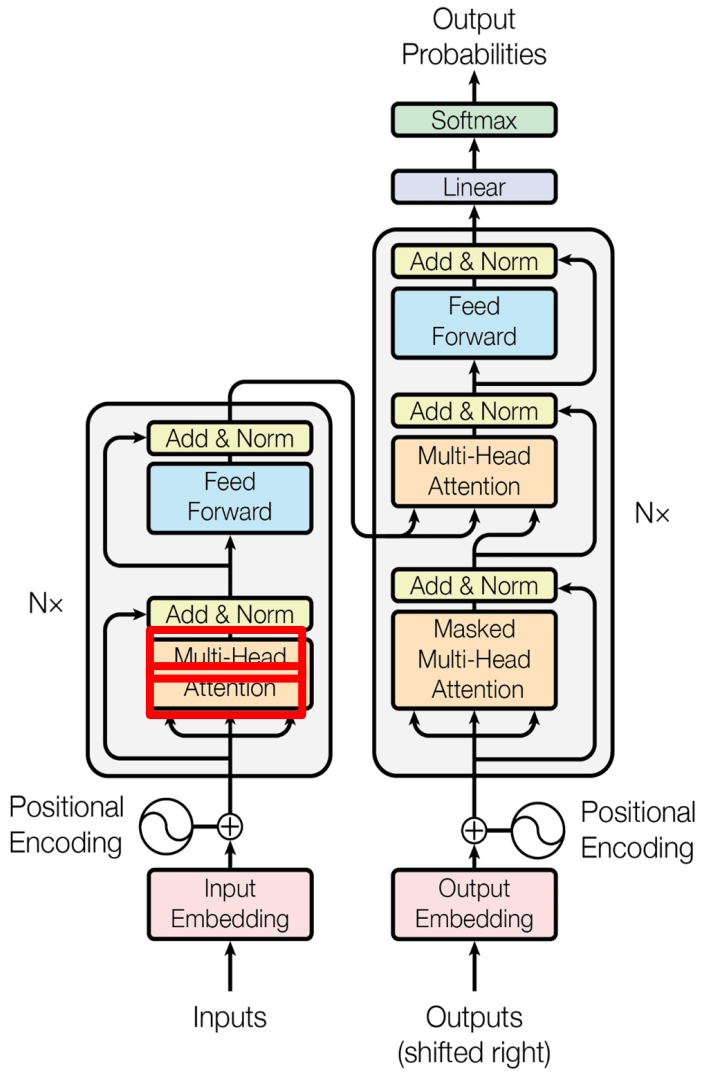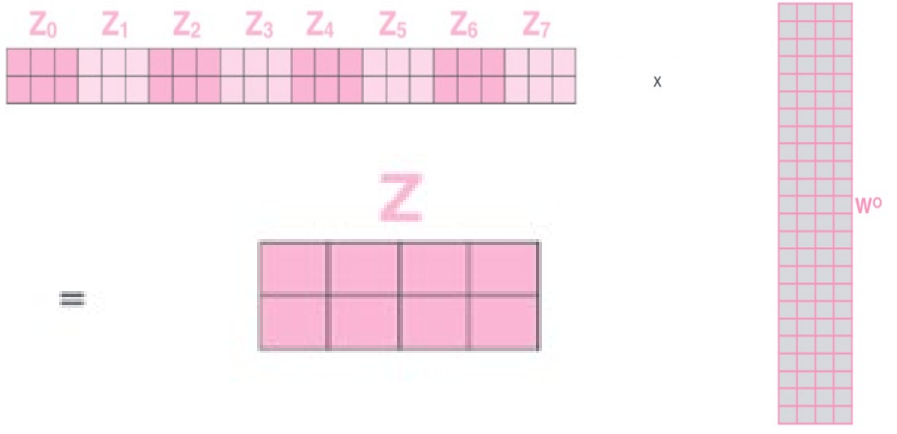
Lack of parallelization: O(T)

Self-attention: O(L)

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



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

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



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

N×

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

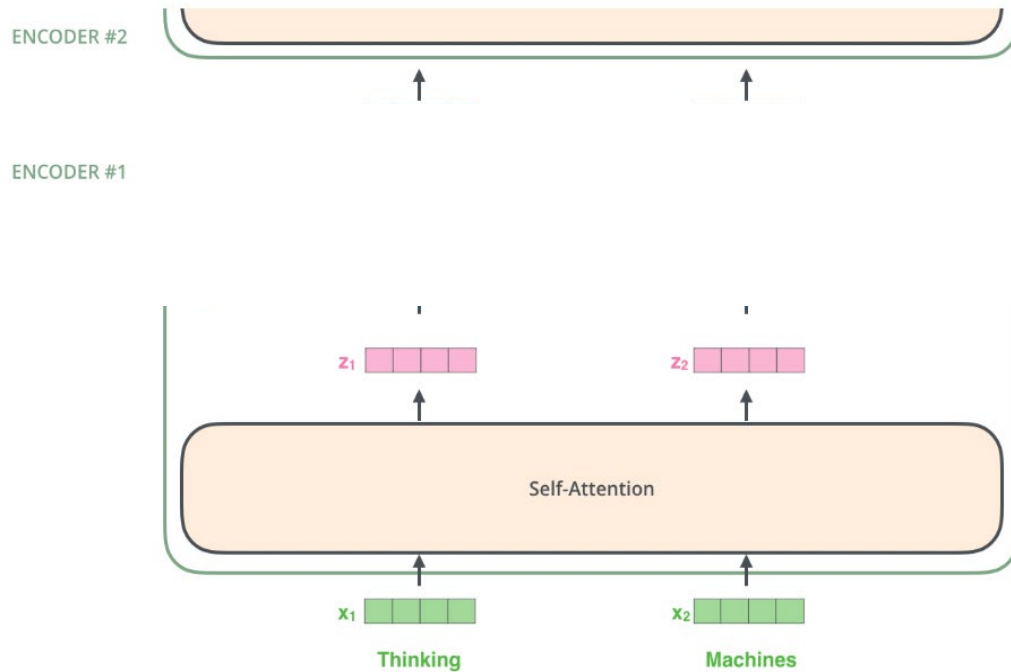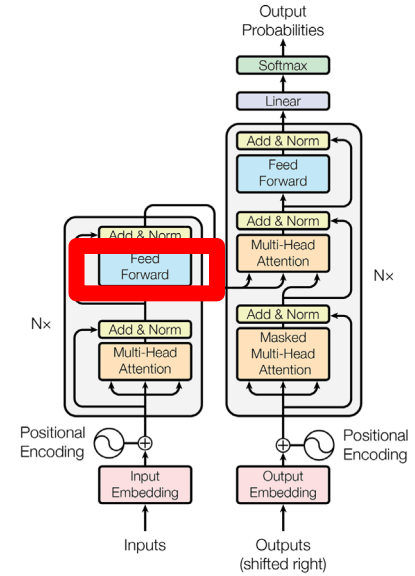Layer: 5 ⬦ Attention: Input - Input ⬦

# 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).
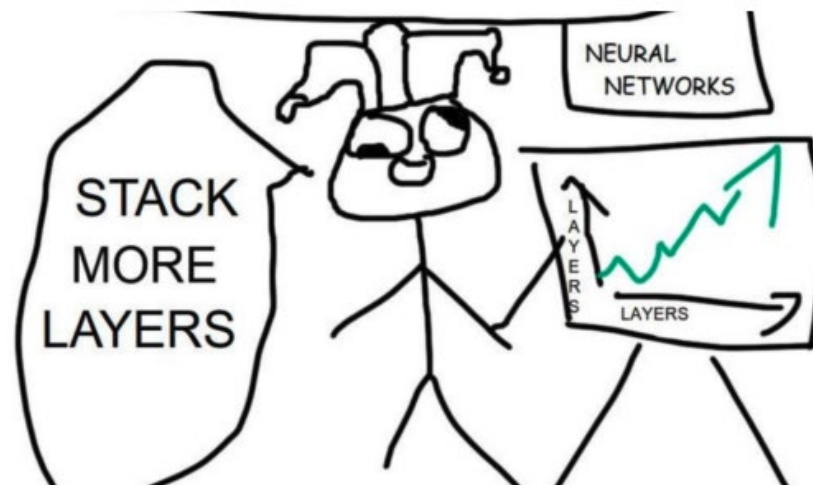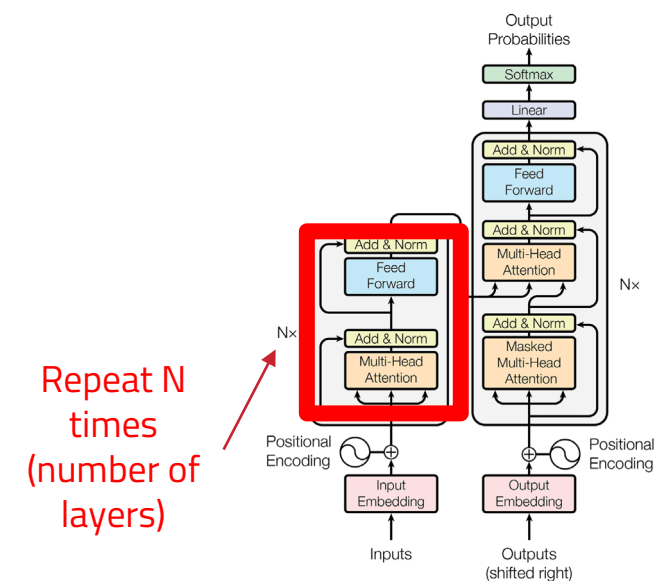


Equation for Feed-Forward layer

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

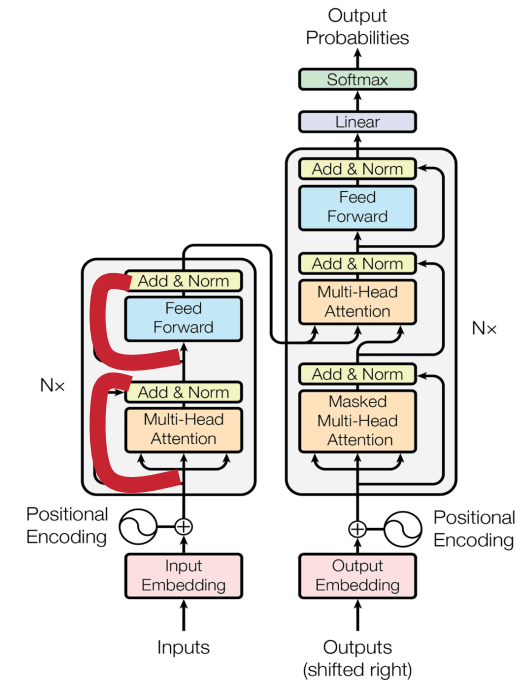# Stacking deep neural nets

❑ Training trick #1: Residual Connections

❑ Training trick #2: LayerNorm
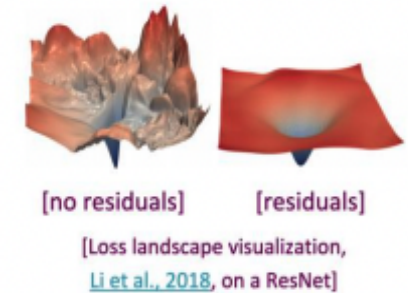
❑ Training trick #3: Scaled Dot Product Attention



Repeat N times (number of layers)

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

❏ Residual connections are a simple but powerful technique from computer vision.

❏ Directly passing "raw" embeddings to the next layer prevents the network from "forgetting" or distorting important information as it is processed by many layers.
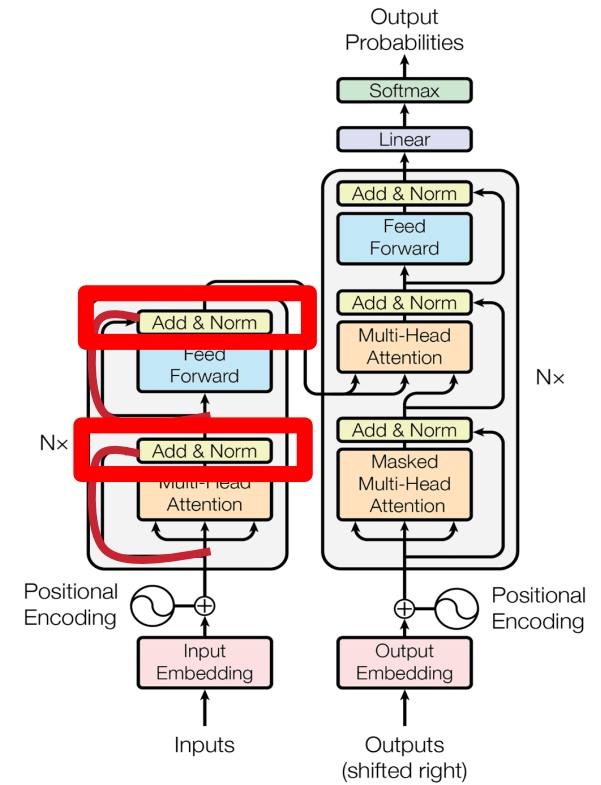


$$x_\ell = F(x_{\ell-1}) + x_{\ell-1}$$

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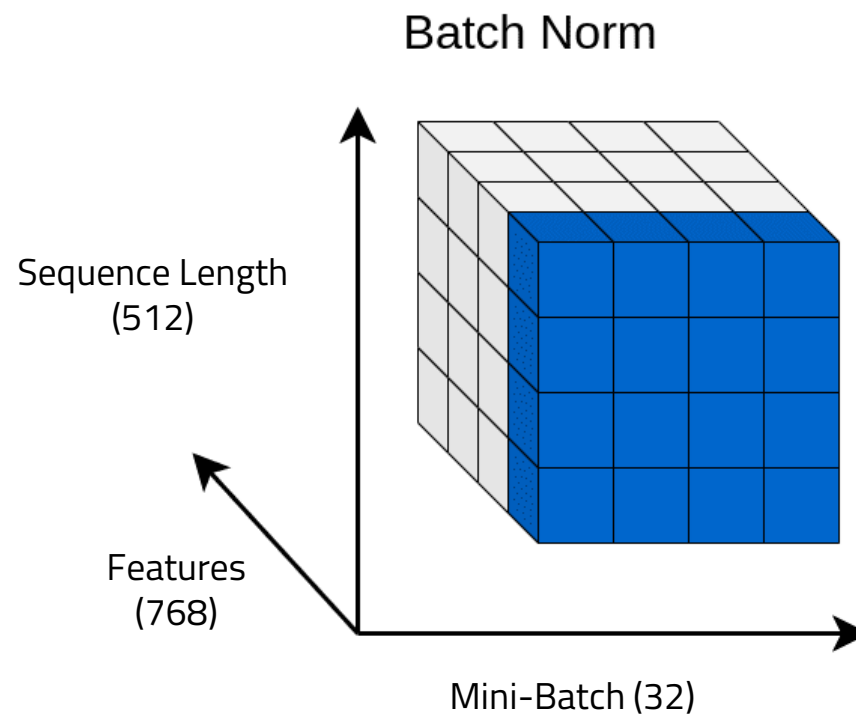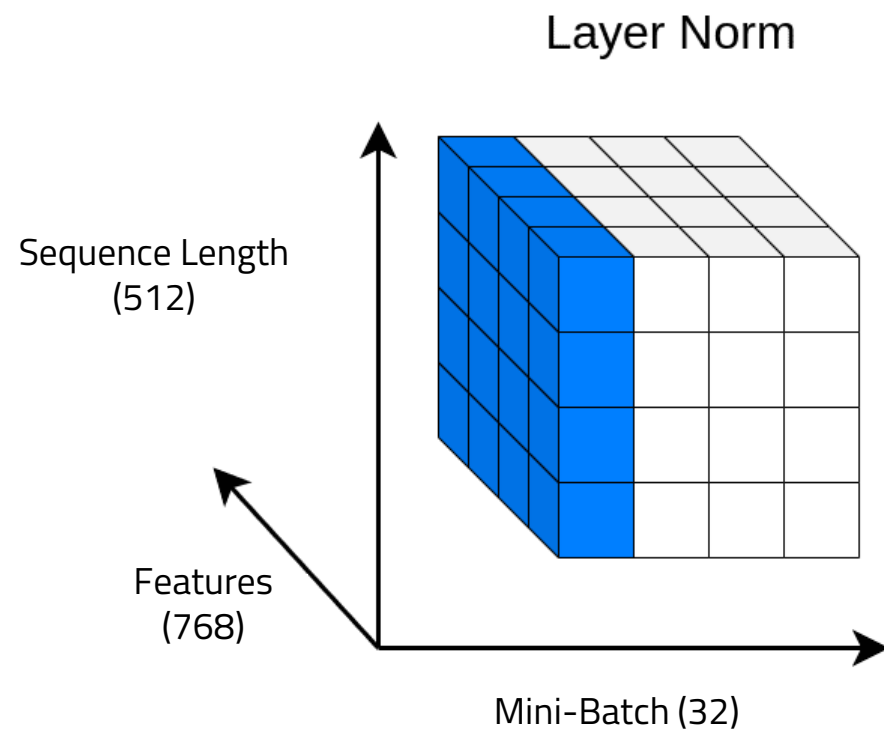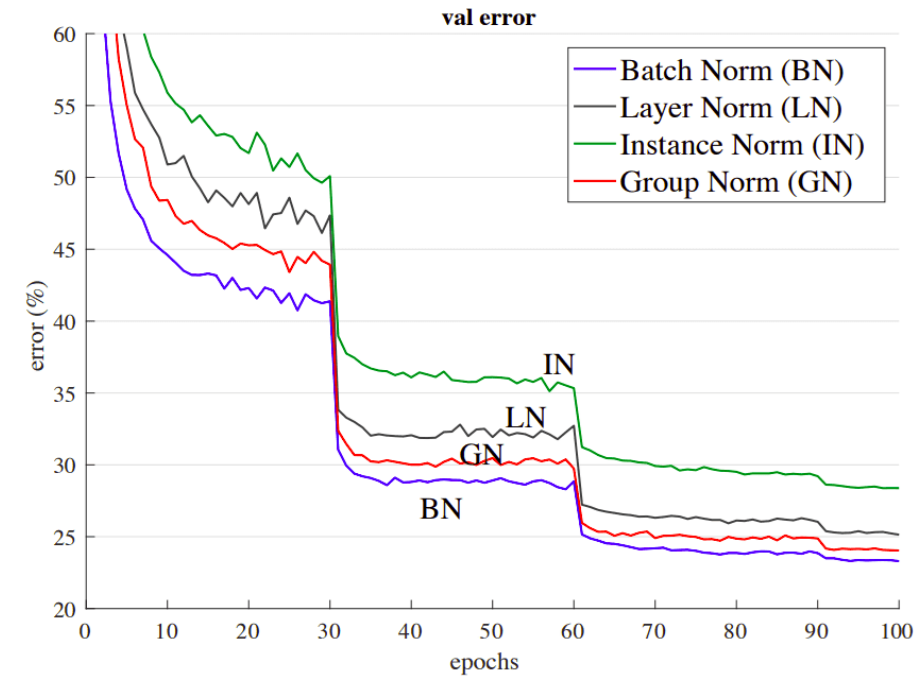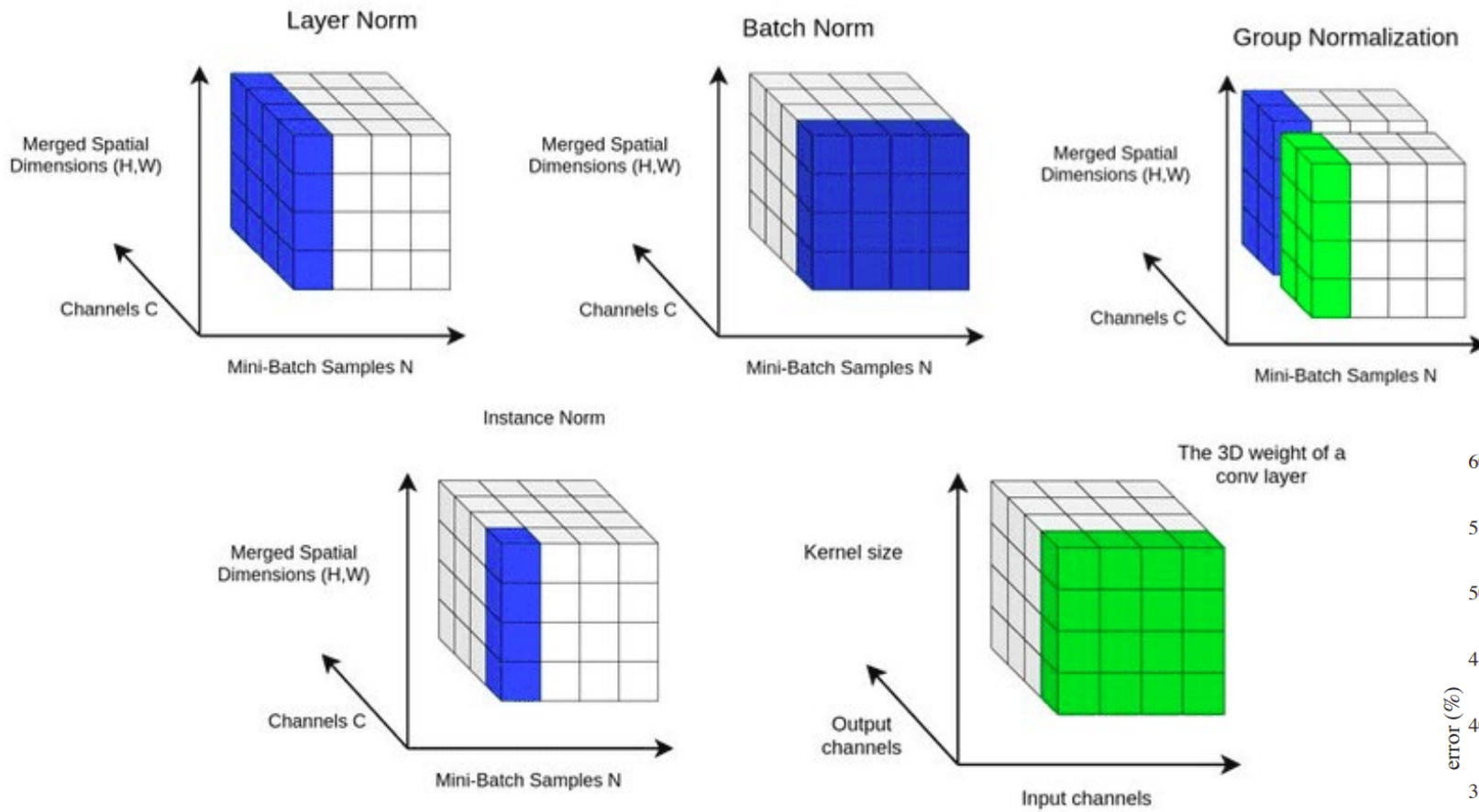
**Mean:** $\mu^l = \frac{1}{H} \sum_{i=1}^{H} a_i^l$    **Standard Deviation:** $\sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^{H} (a_i^l - \mu^l)^2}$

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

# Layer norm vs Batch norm



Layer Norm

Sequence Length
(512)

Features
(768)

Mini-Batch (32)

Batch Norm

Sequence Length
(512)

Features
(768)

Mini-Batch (32)

Layer Norm

Batch Norm

Group Normalization

Instance Norm

The 3D weight of a conv layer
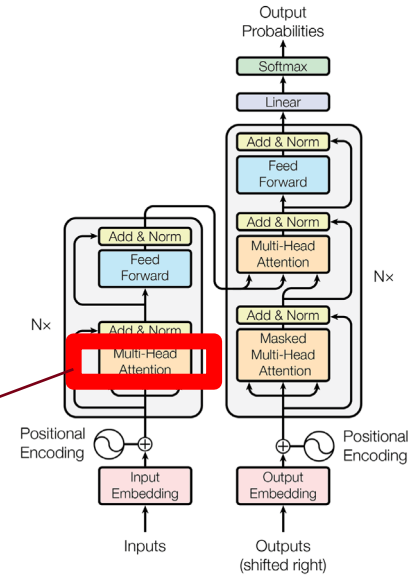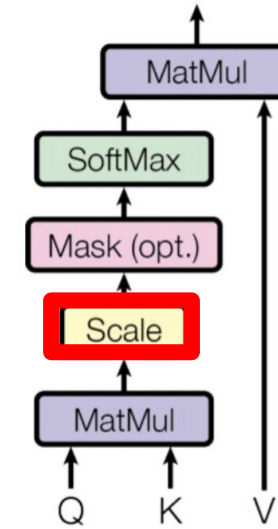
val error

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 $i \in \{1, 2, \ldots, T\}$ are position vectors
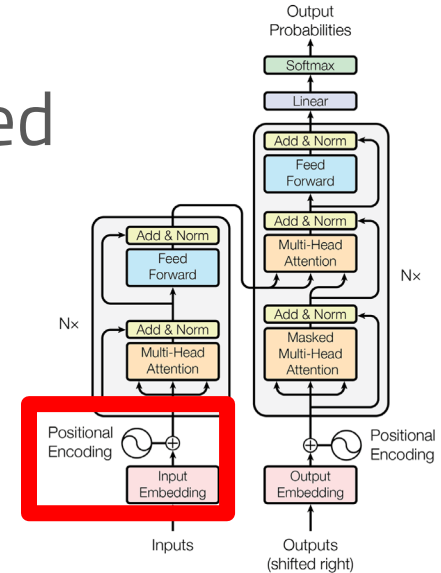
- ❏ Easy to incorporate this info into our self-attention block: just add the $pi$ to our inputs!

$$v_i = \tilde{v}_i + p_i$$
$$q_i = \tilde{q}_i + p_i$$
$$k_i = \tilde{k}_i + p_i$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add…
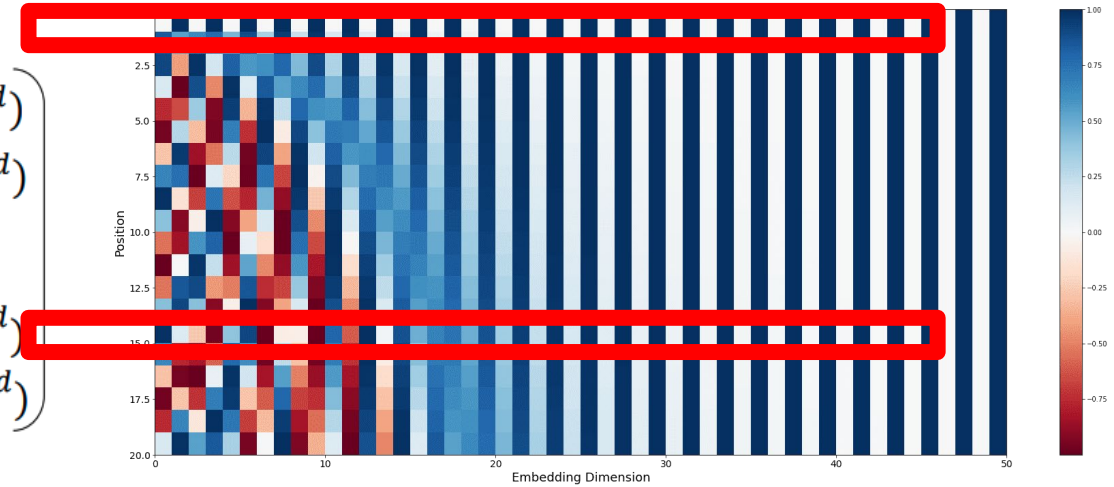
# Position representation vectors through sinusoids

❑ Sinusoidal position representations: concatenate sinusoidal functions of varying periods:
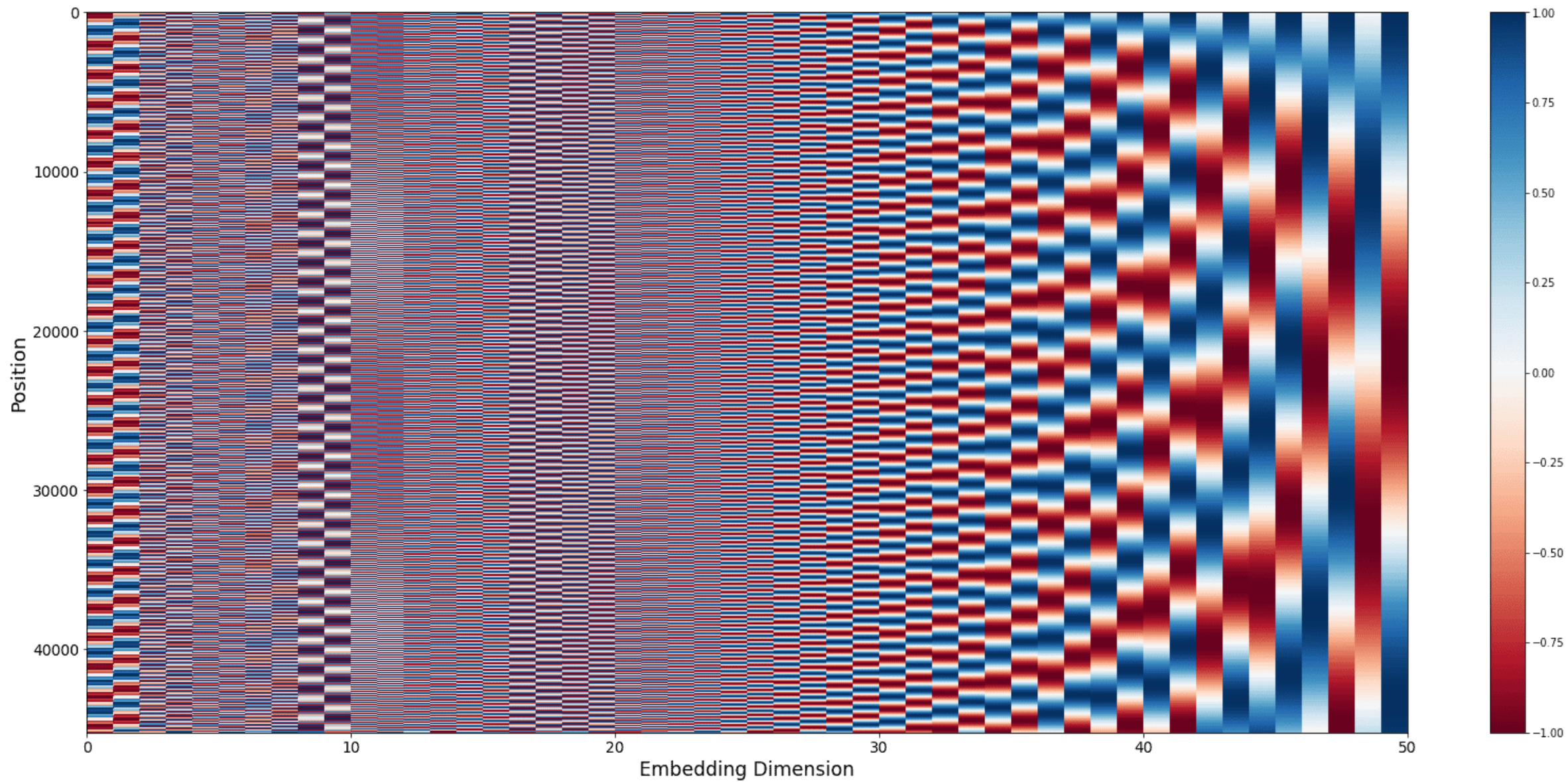
$$\overrightarrow{p_t}^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k.t), & \text{if } i = 2k \\ \cos(\omega_k.t), & \text{if } i = 2k+1 \end{cases} \qquad p_i = \b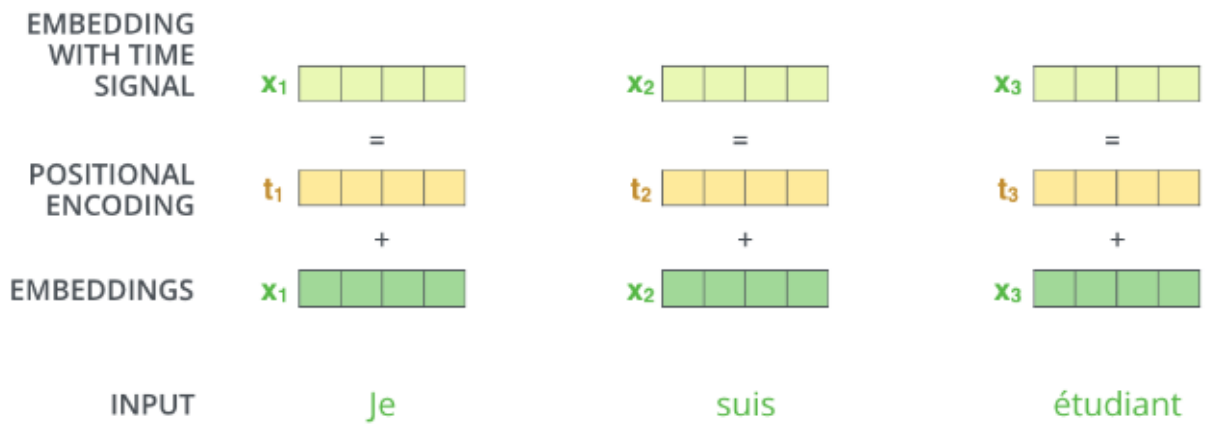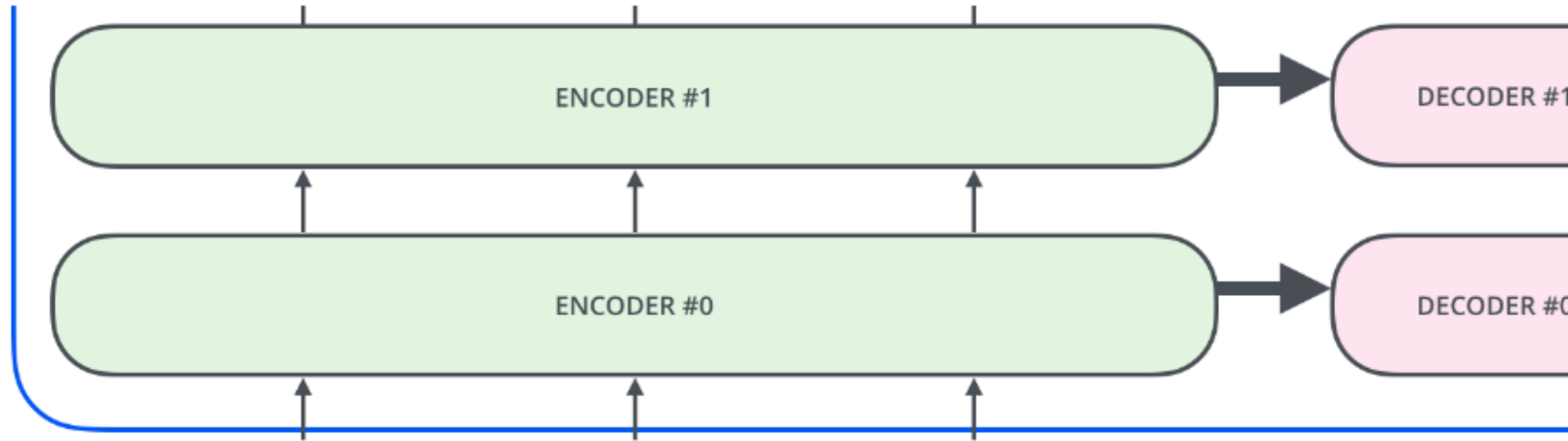egin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$
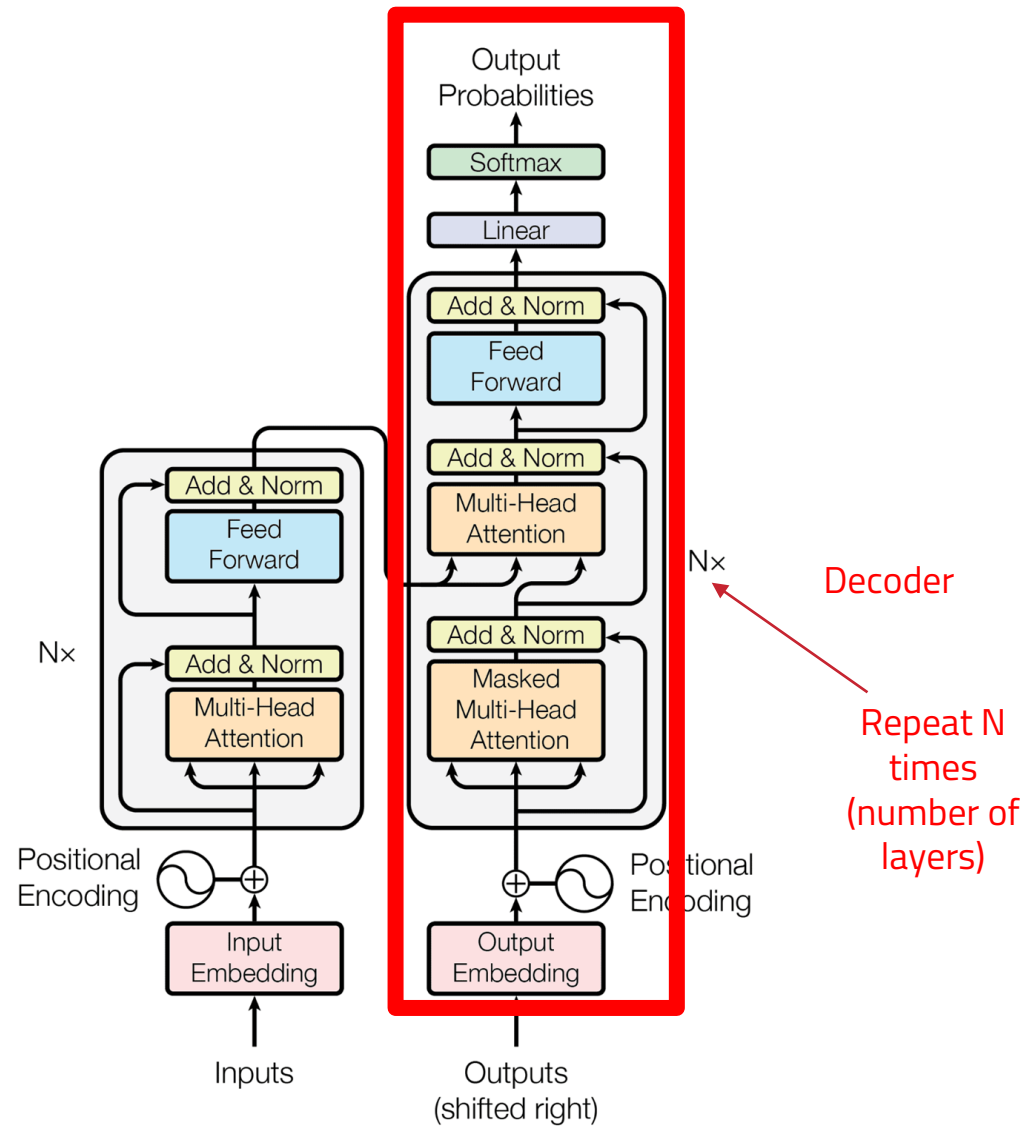


❑ Pros: Periodicity indicates that maybe "absolute position" isn't as important
❑ Cons: Not learnable; also the extrapolation doesn't really work

ENCODER #1 → DECODER #1

ENCODER #0 → DECODER #0

EMBEDDING WITH TIME SIGNAL: $x_1$, $x_2$, $x_3$

=

POSITIONAL ENCODING: $t_1$, $t_2$, $t_3$

+

EMBEDDINGS: $x_1$, $x_2$, $x_3$

INPUT: Je  suis  étudiant
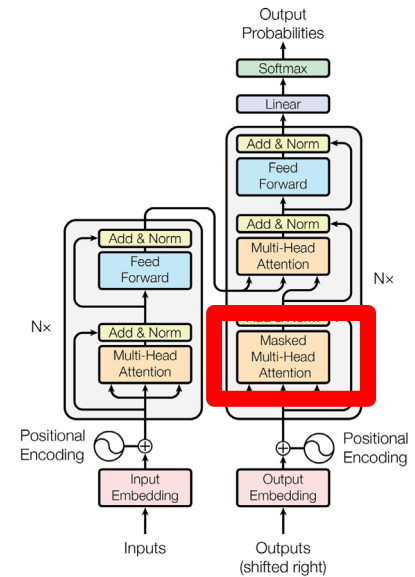
Decoder

Repeat N times (number of layers)

# 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.
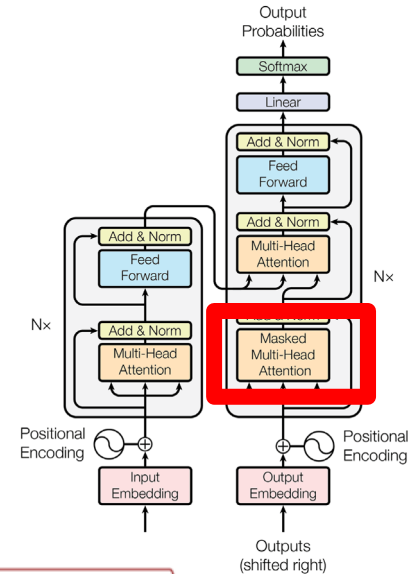
# Masking the future in self-attention

❑ 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 $-\infty$

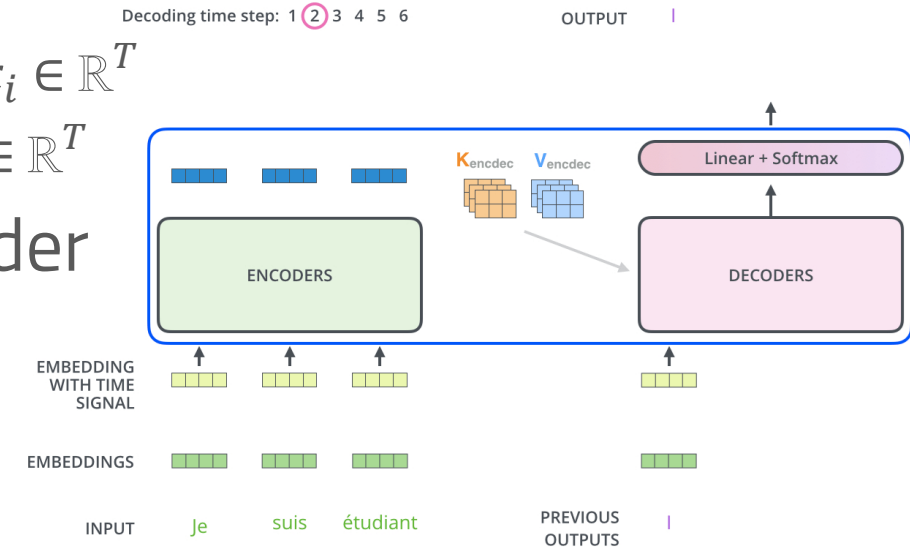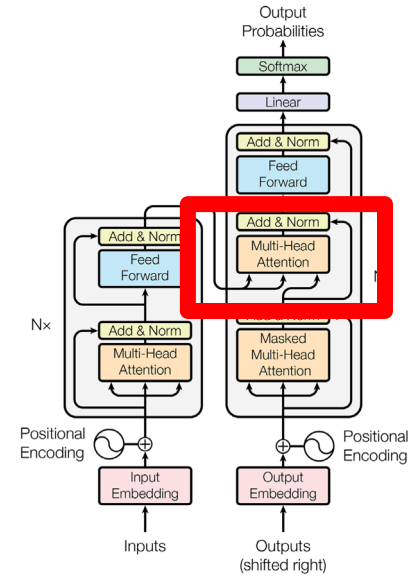$$e_{ij} = \begin{cases} q_i^\mathsf{T} k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

We can look at these (not greyed out) words

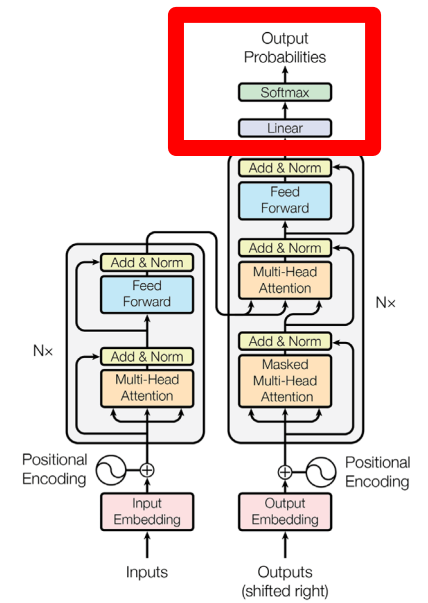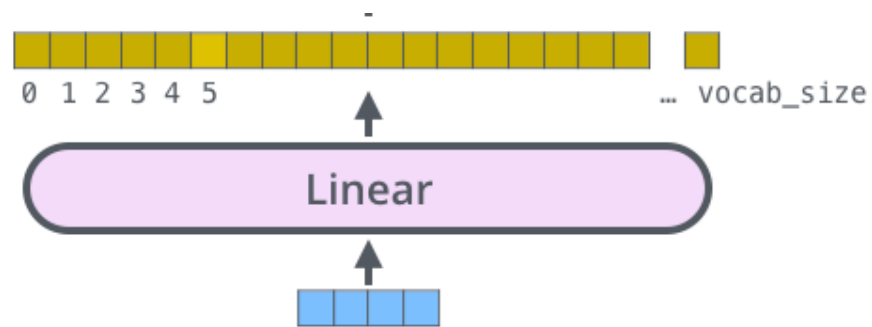|        | [START] | The | chef | who |
|--------|---------|-----|------|-----|
| [START] | $-\infty$ | $-\infty$ | $-\infty$ | $-\infty$ |
| The    |         | $-\infty$ | $-\infty$ | $-\infty$ |
| chef   |         |     | $-\infty$ | $-\infty$ |
| who    |         |     |      | $-\infty$ |

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,

  ○ $q_i = Q z_i$

logits

0 1 2 3 4 5        … vocab_size

Linear

Decoder stack output

Which word in our vocabulary is associated with this index?

am

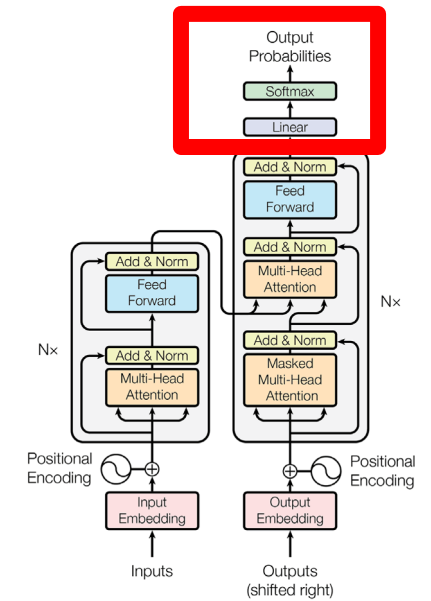Get the index of the cell with the highest value (`argmax`)

5

log_probs

0 1 2 3 4 5 … vocab_size

**Softmax**

logits

0 1 2 3 4 5 … vocab_size

**Linear**

Decoder stack output

Output Probabilities

Softmax

Linear

Add & Norm
Feed Forward

Add & Norm
Feed Forward

Add & Norm
Multi-Head Attention

Add & Norm
Multi-Head Attention

N×

N×

Add & Norm
Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Drawback of Transformer

# Drawback of Transformer

❏ **Quadratic compute in self-attention:**

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

❏ **Position representations:**

- o  Are simple absolute indices the best we can do to represent position?
- o  Relative linear position attention [Shaw et al., 2018]
- o  Dependency syntax-based position [Wang et al., 2019]

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

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

# Reduce $O(T\text{^}2)$ 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.
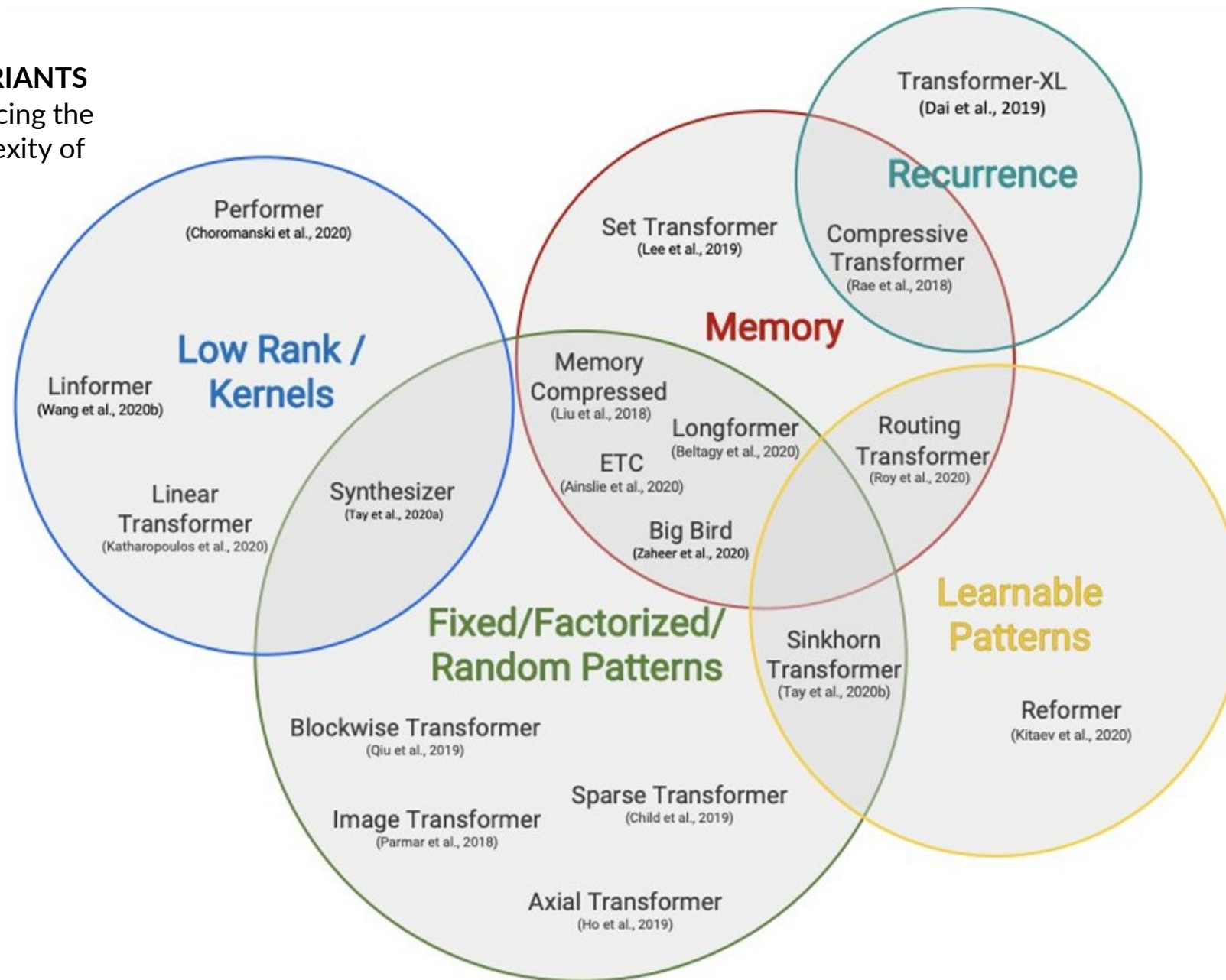


(a) Random attention    (b) Window attention    (c) Global Attention    (d) BIGBIRD

**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.'*



## 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) |

# 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 | |

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 | |
| 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.

  o Why? On our way back to fully-connected models, throwing out the inductive bias of CNNs and RNNs.