# Transformers Part 1

Lecture 17
Subir Varma

# So Far …

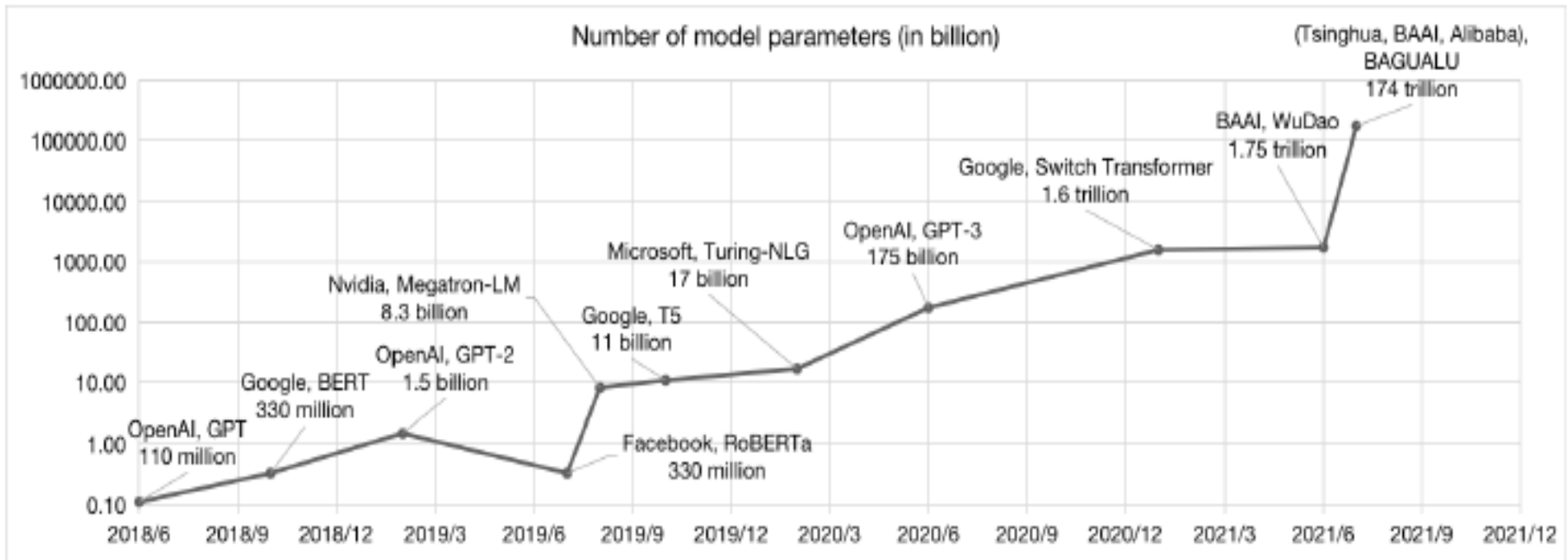                     

Linear → Dense Feed Forward → ConvNets → RNNs → LSTMs → Attention

(vectors)             (images)            (sequences)

▼

## Next: Transformers

Also used for Sequences

# Transformers

Number of model parameters (in billion)

1000000.00 ·
100000.00 ·
10000.00 ·
1000.00 ·
100.00 ·
10.00 ·
1.00 ·
0.10 ·

2018/6  2018/9  2018/12  2019/3  2019/6  2019/9  2019/12  2020/3  2020/6  2020/9  2020/12  2021/3  2021/6  2021/9  2021/12

OpenAI, GPT
110 million

Google, BERT
330 million

OpenAI, GPT-2
1.5 billion

Nvidia, Megatron-LM
8.3 billion

Facebook, RoBERTa
330 million

Google, T5
11 billion

Microsoft, Turing-NLG
17 billion

OpenAI, GPT-3
175 billion

Google, Switch Transformer
1.6 trillion

BAAI, WuDao
1.75 trillion

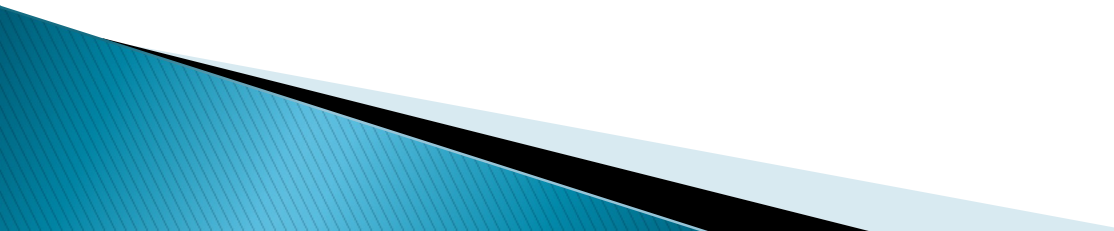(Tsinghua, BAAI, Alibaba),
BAGUALU
174 trillion

- Has made possible much larger models (higher capacity)
- Can be trained using Self Supervised Learning, so very large datasets are available
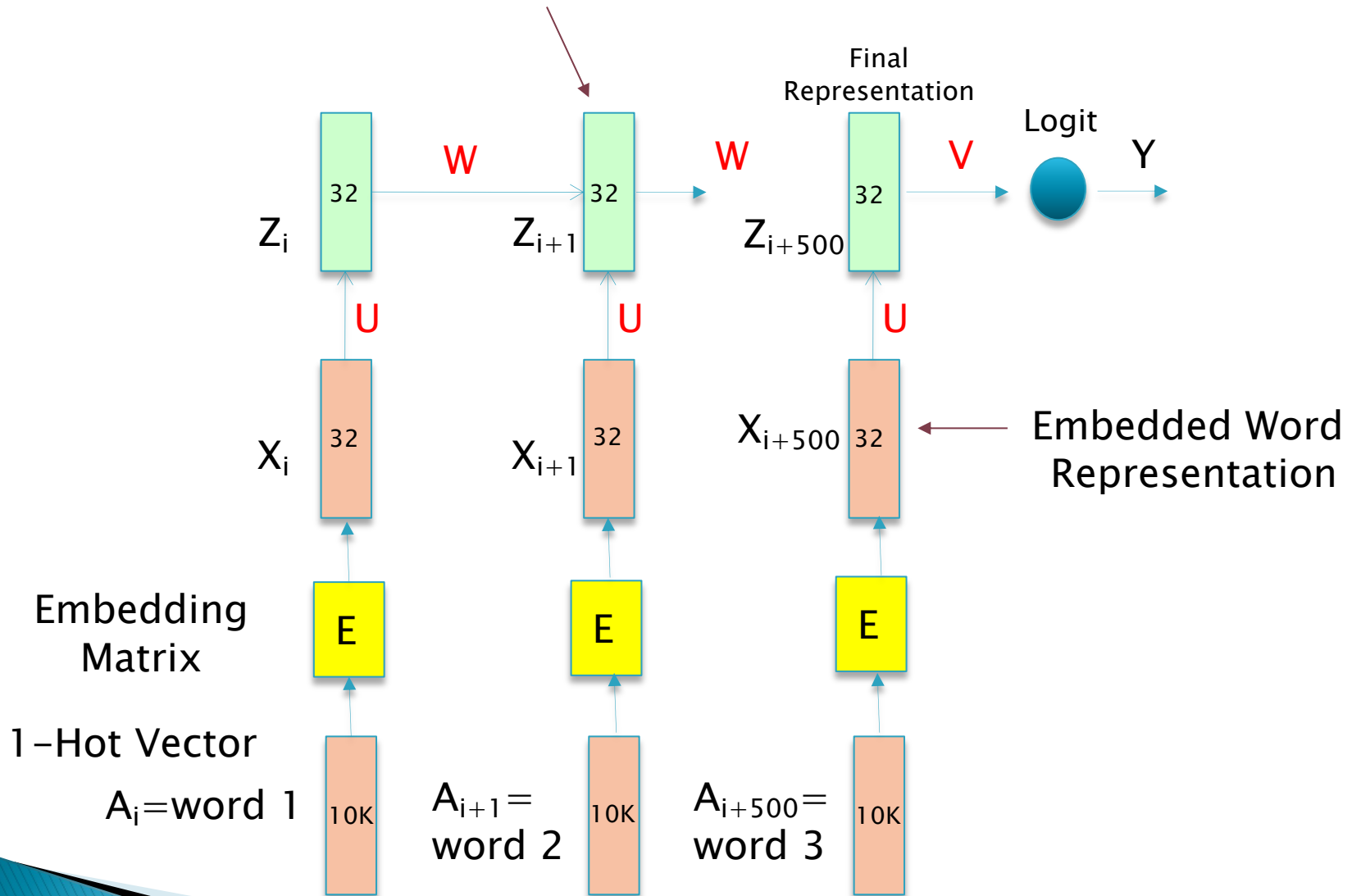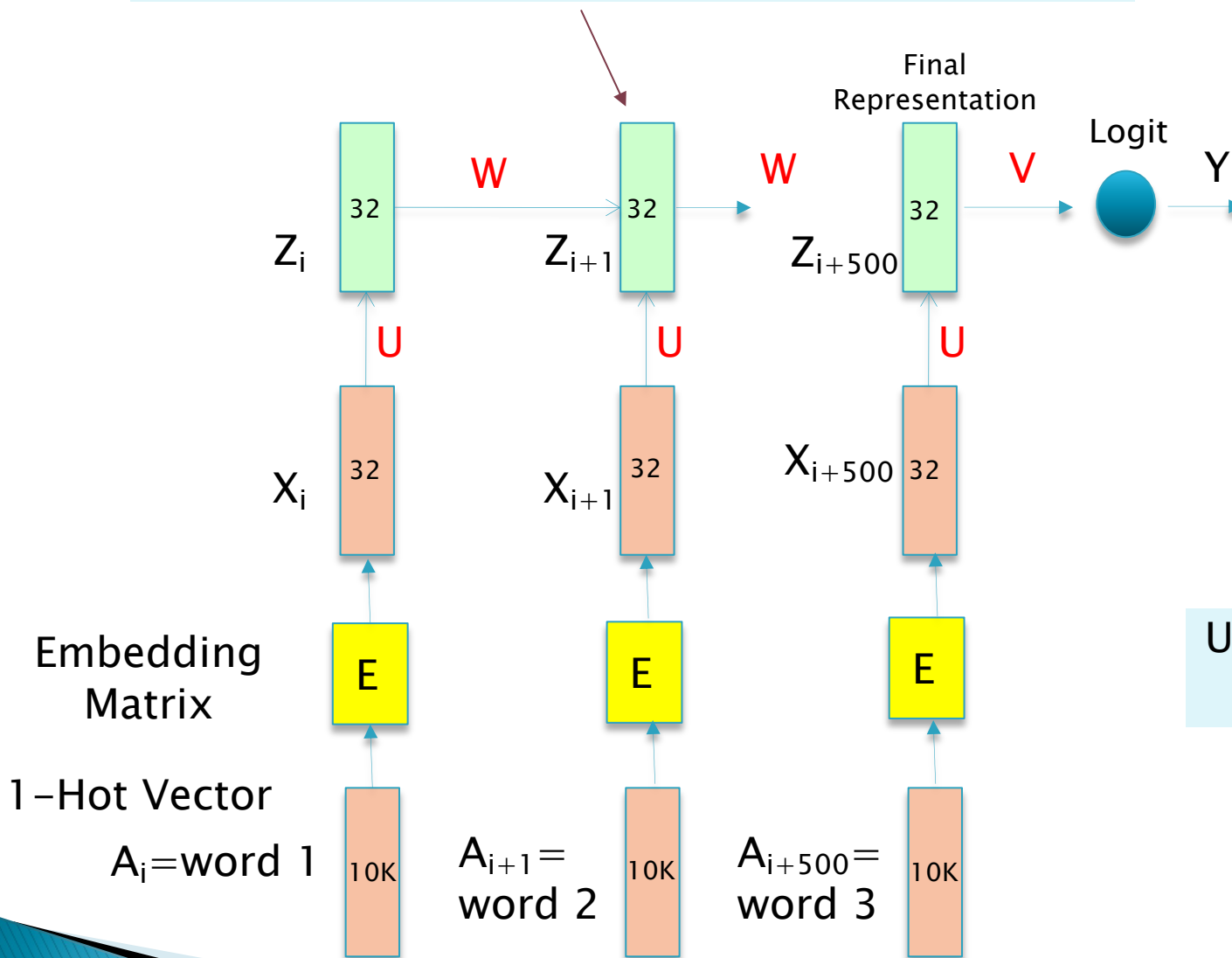
Huge Models ⟷ Huge Datasets

# Problems with RNNs/LSTMs

- Sequential computation prevents parallelization
- Despite GRUs and LSTMs, RNNs still need the attention mechanism to deal with long range dependencies – path length for codependent computation between states grows with sequence
- But if attention gives us access to any state... maybe we don't need the RNN?

But this representation is only a function of the words that came before it

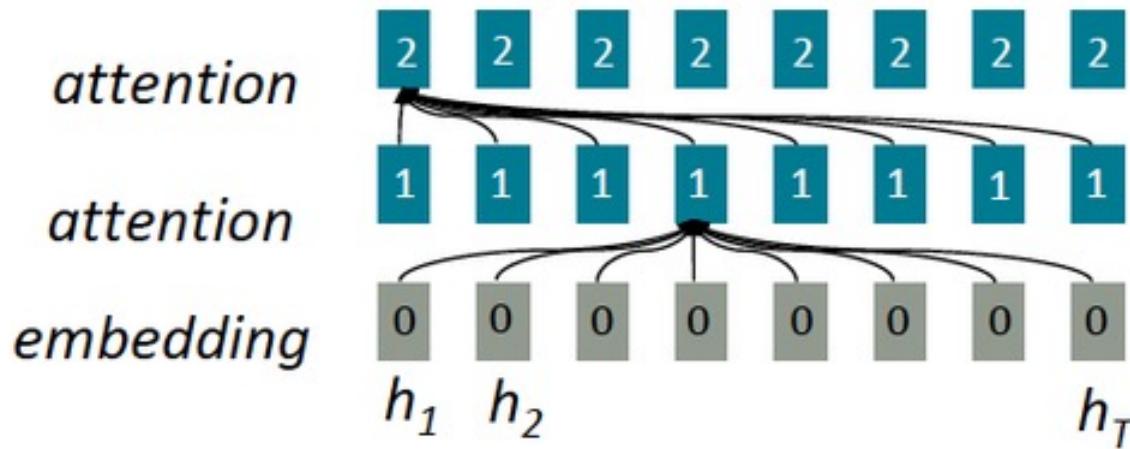Why not make it a function of all the words in the sentence!

Final Representation

Logit

$Z_i$    32    W    $Z_{i+1}$    32    W    $Z_{i+500}$    32    V    Y

U    U    U

$X_i$    32    $X_{i+1}$    32    $X_{i+500}$    32

How?

Embedding Matrix

E    E    E

Use the Attention Mechanism

1-Hot Vector

$A_i$=word 1    10K    $A_{i+1}$= word 2    10K    $A_{i+500}$= word 3    10K

# Self Attention

# Self Attention



attention

attention

embedding

$h_1$  $h_2$  $h_T$

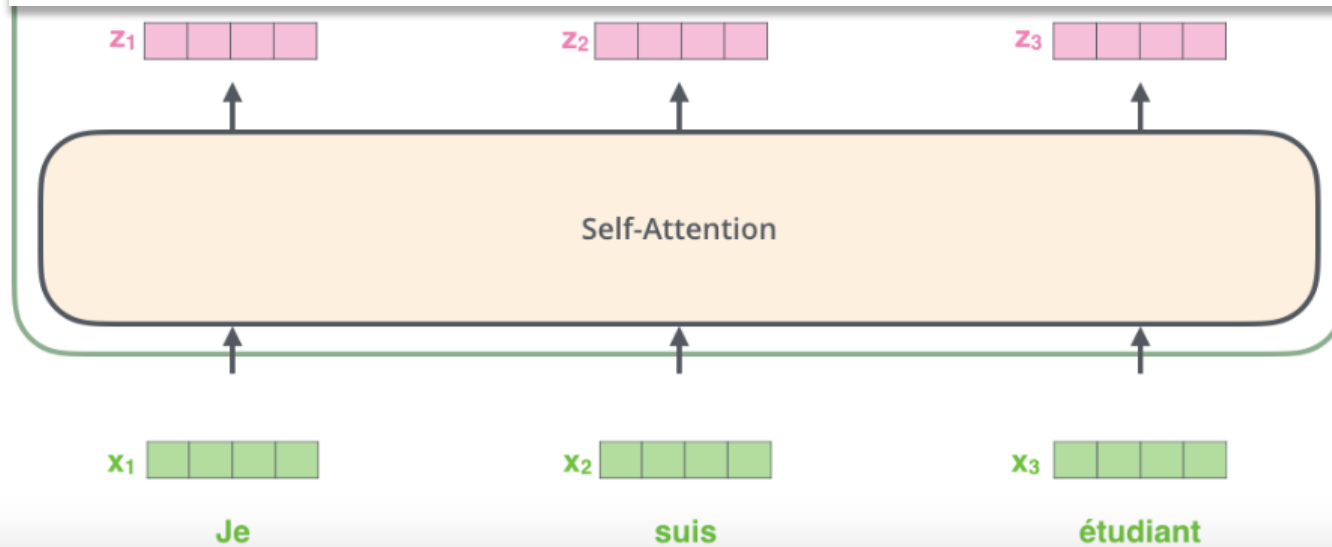All words attend to all words in previous layer; most arrows here are omitted

- The representation of each of the words is modified by every other word in the sequence by using the Self Attention mechanism
- The idea behind this architecture is that after several layers of Self Attention, each word develops a representation that takes into account all the other words that exist in the sentence.

# Self Attention

- Meaning of a word is context specific:
  - Example: Mark a Date vs Going on a date vs buying date at the market
- A Smart Embedding technique would provide a different vector representation for a word depending upon the other words surrounding it

- Self Attention: A way to make Word Representations Context Aware
- It does so by modulating the the word representation by using the representations of other words in the sentence.

# Self Attention (cont)

How to Compute the Self Attention Scores?
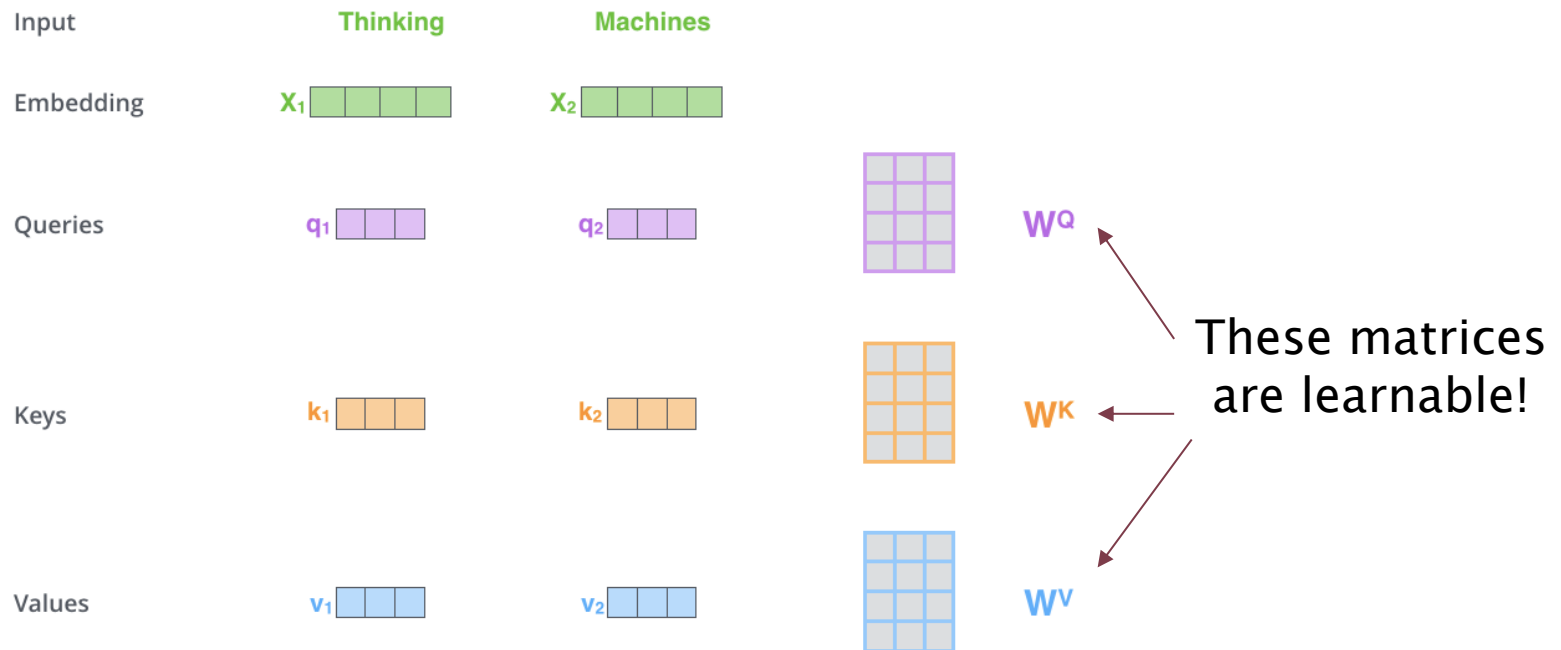
# Computing Self Attention Scores

▸ Simplest way:

$$a_{ik} = x_i \circ x_{k,} \, 1 \leq k \leq n$$
$$w_{ik} = softmax_k(a_{ik})$$

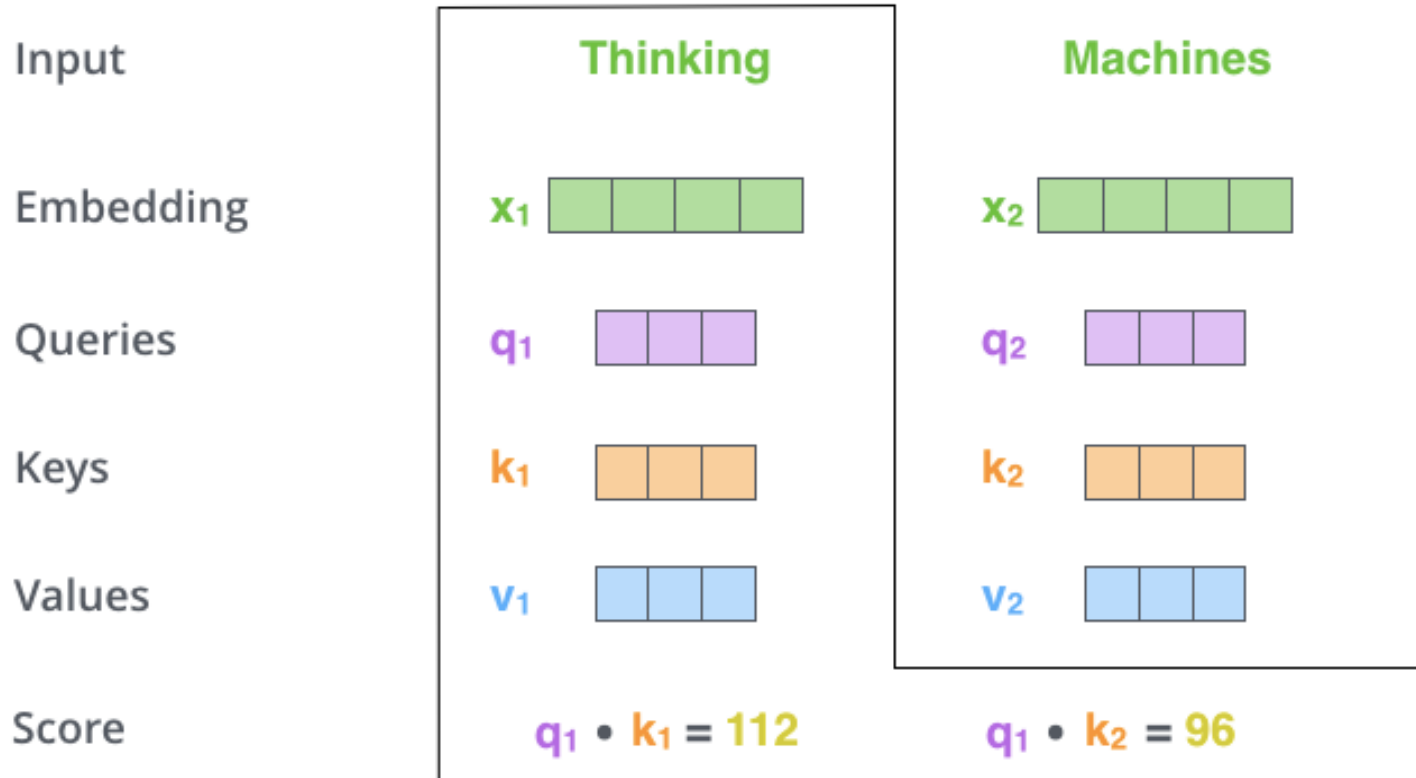This technique does not involve any learning since there are no parameters

# Generalized Self Attention

▸ The initial word embeddings are sent through 3 independent sets of dense projections, resulting in 3 separate vectors (per word)
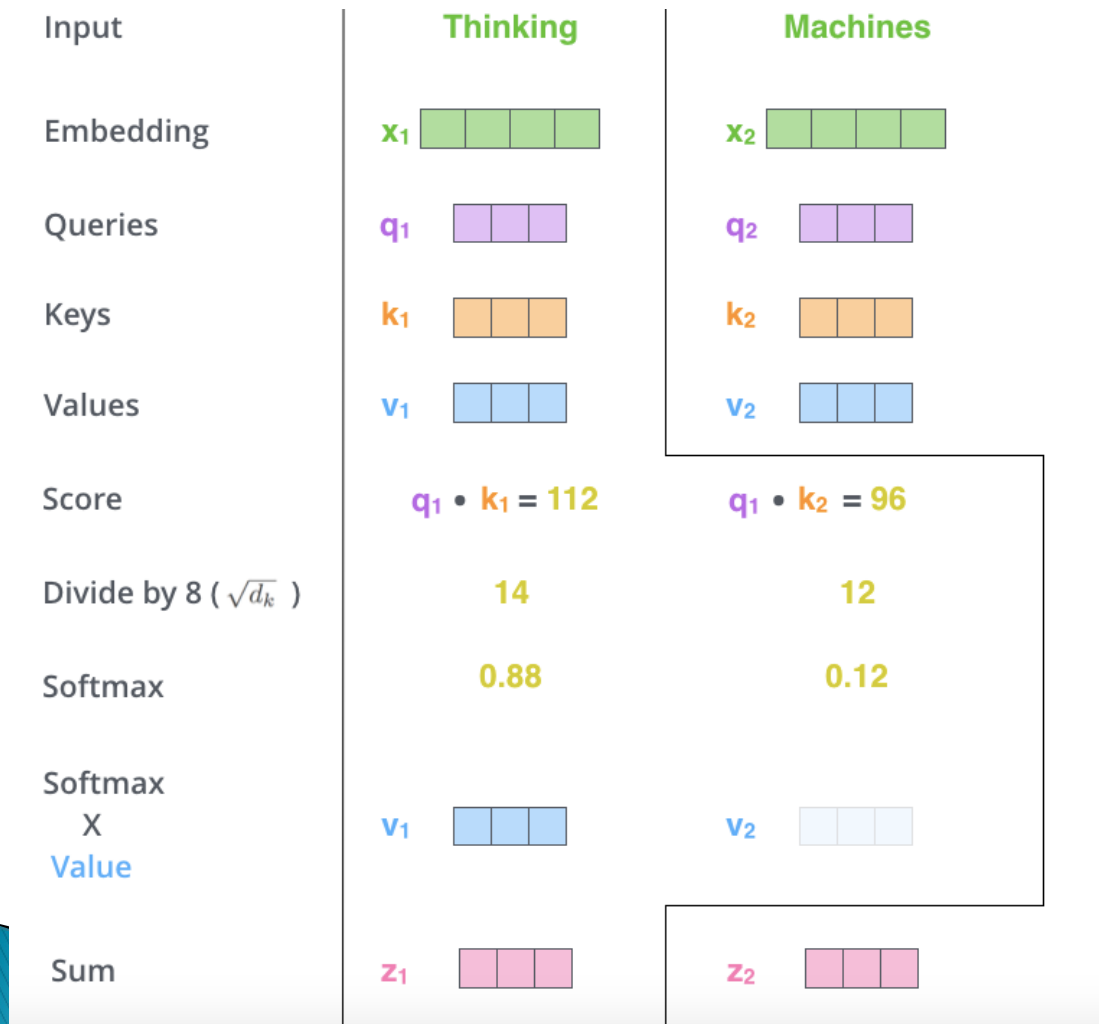


Multiplying x1 by the WQ weight matrix produces q1, the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

# Generalized Self Attention (cont)
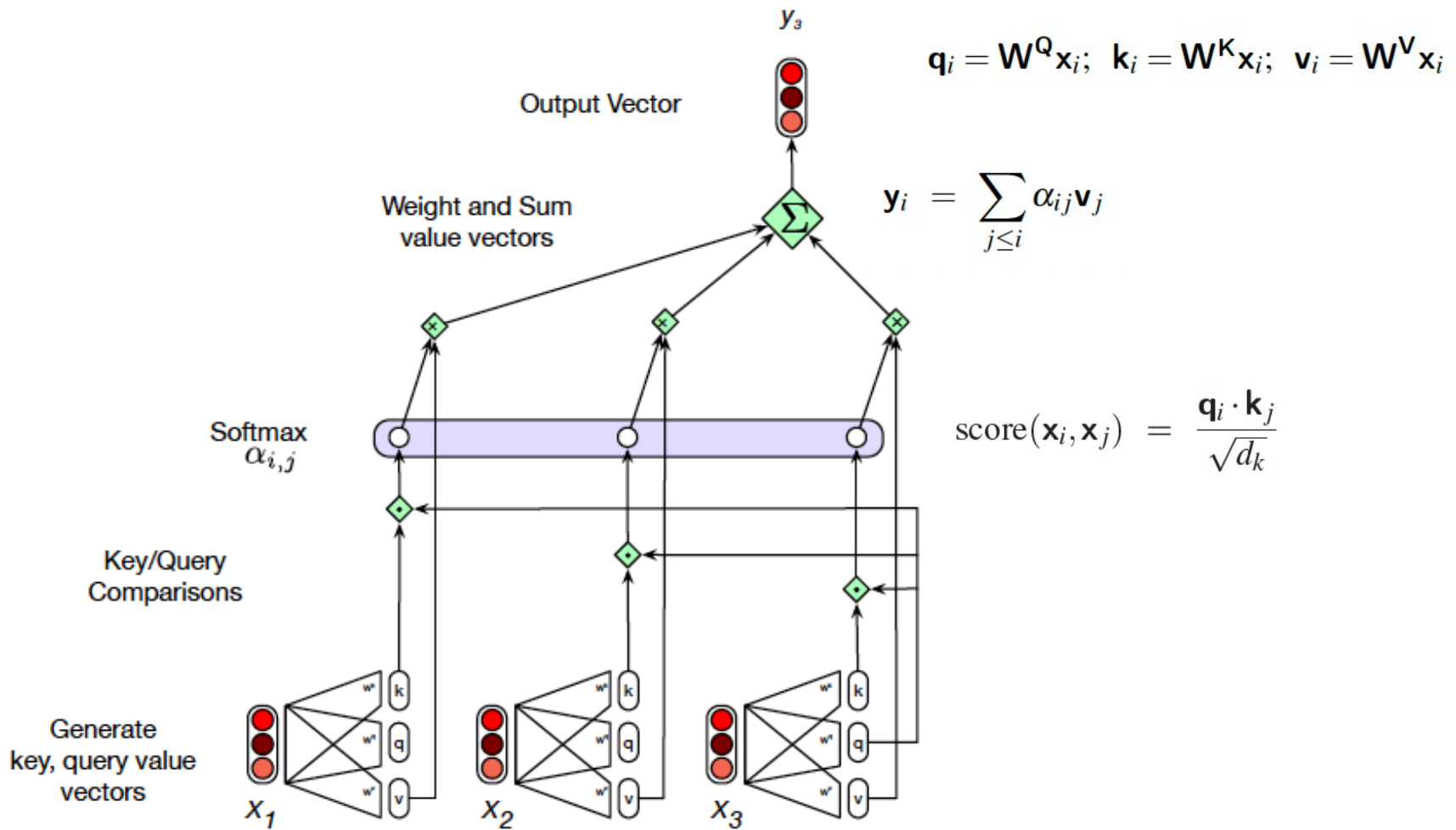
# Generalized Self Attention (cont)

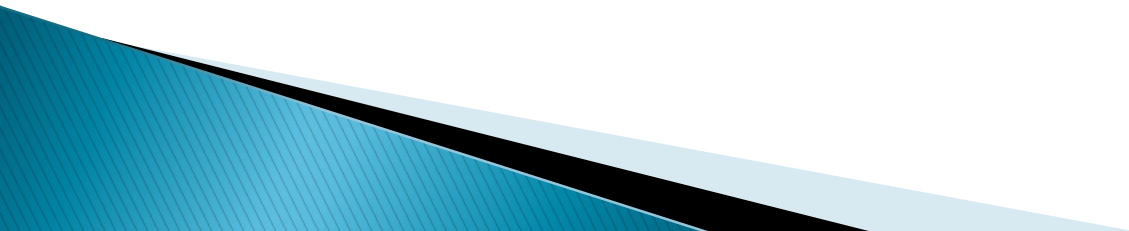| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ ▢▢▢▢ | $x_2$ ▢▢▢▢ |
| Queries | $q_1$ ▢▢▢ | $q_2$ ▢▢▢ |
| Keys | $k_1$ ▢▢▢ | $k_2$ ▢▢▢ |
| Values | $v_1$ ▢▢▢ | $v_2$ ▢▢▢ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ ▢▢▢ | $v_2$ ▢▢▢ |
| Sum | $z_1$ ▢▢▢ | $z_2$ ▢▢▢ |

$Q=W^QX; \ K=W^KX; \ V=W^VX$

Summarized as:

$$Z = softmax(\frac{QK^T}{\sqrt{d}})V$$

# Generalized Self Attention (cont)



$$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i; \quad \mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i; \quad \mathbf{v}_i = \mathbf{W}^V \mathbf{x}_i$$

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

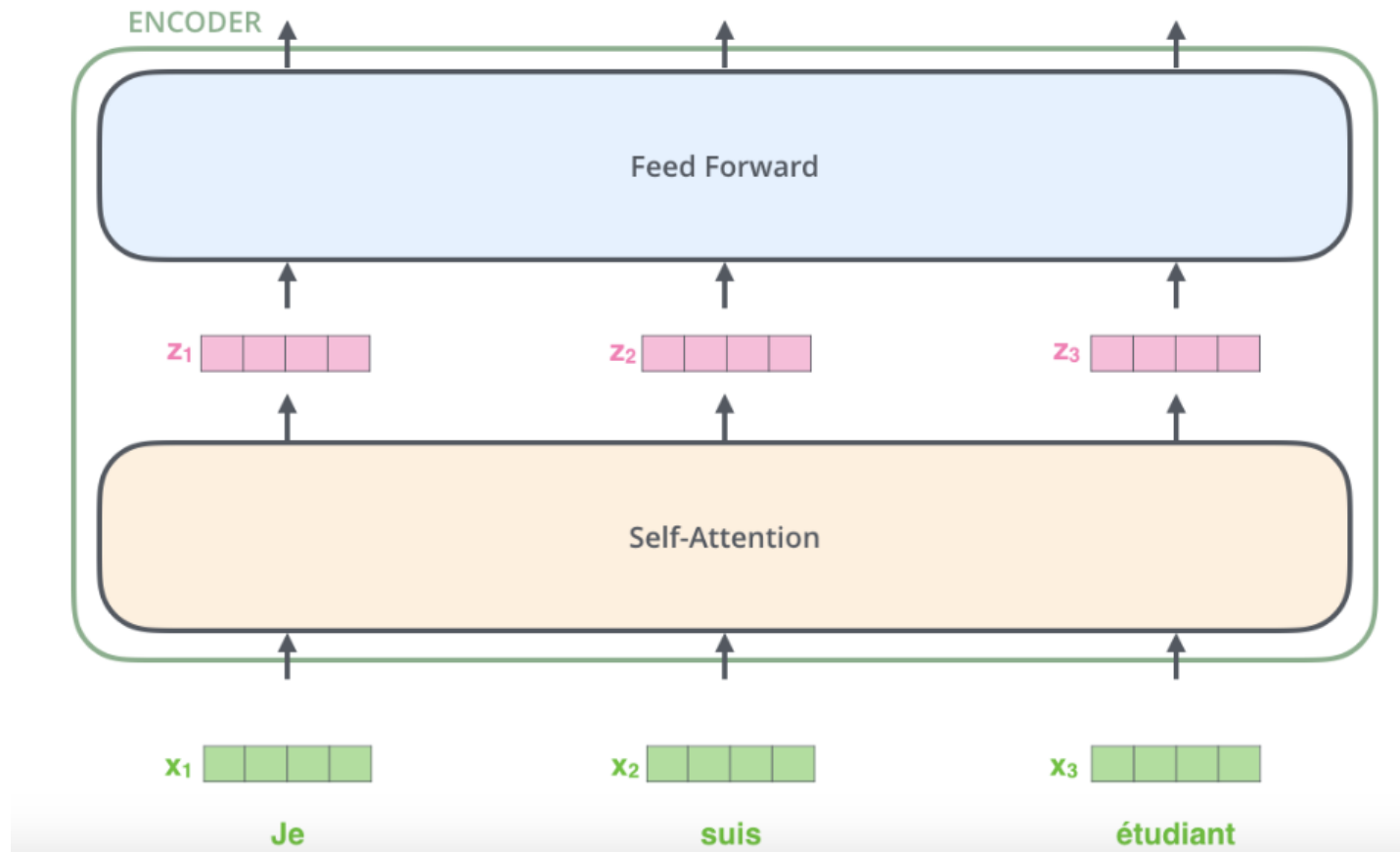$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}}$$
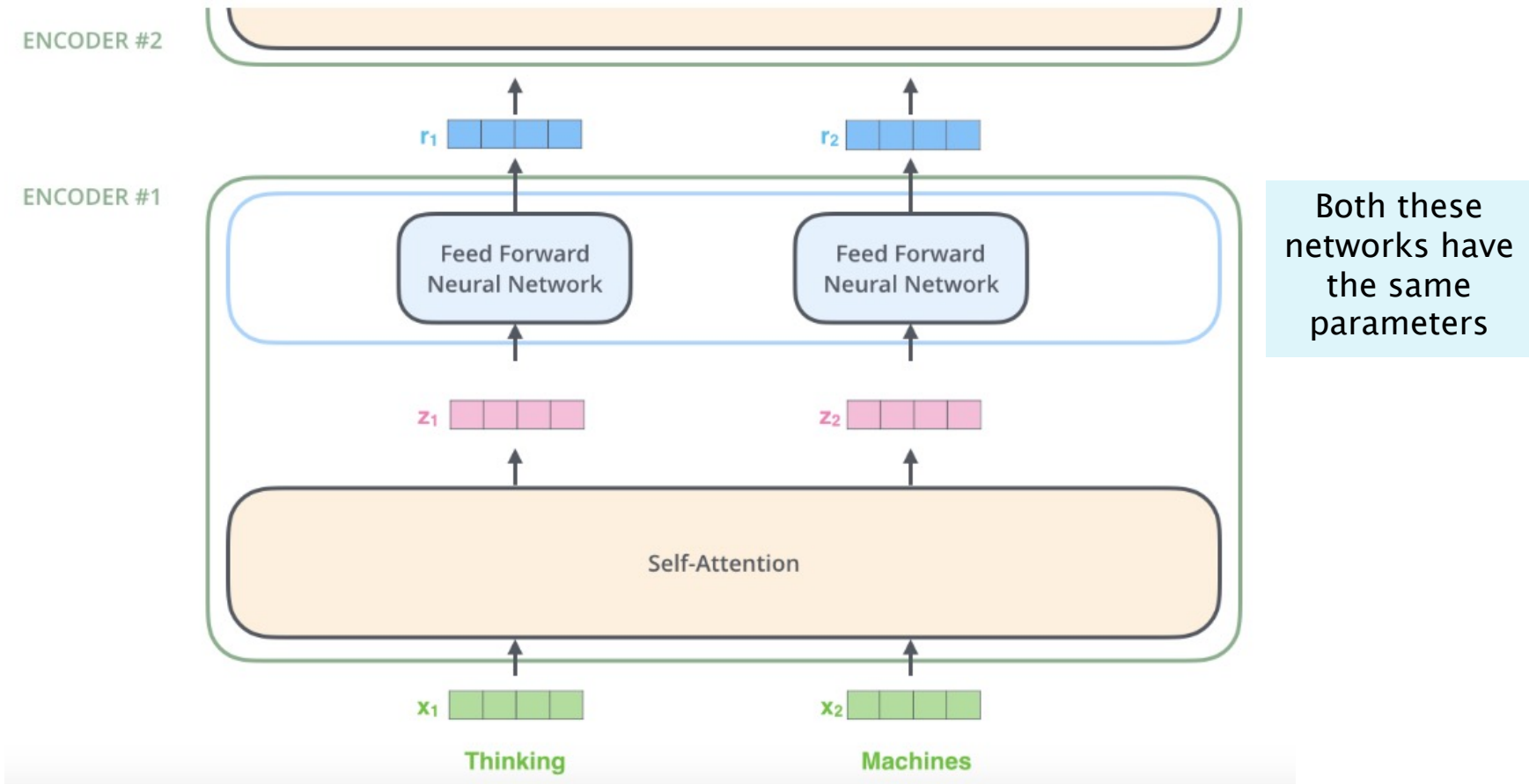
# Transformer Encoders

# Transformers – Encoders



Here we begin to see one key property of the Transformer, which is that the word in each position flows through its own path in the encoder. There are dependencies between these paths in the self-attention layer. The feed-forward layer does not have those dependencies, however, and thus the various paths can be executed in parallel while flowing through the feed-forward layer.
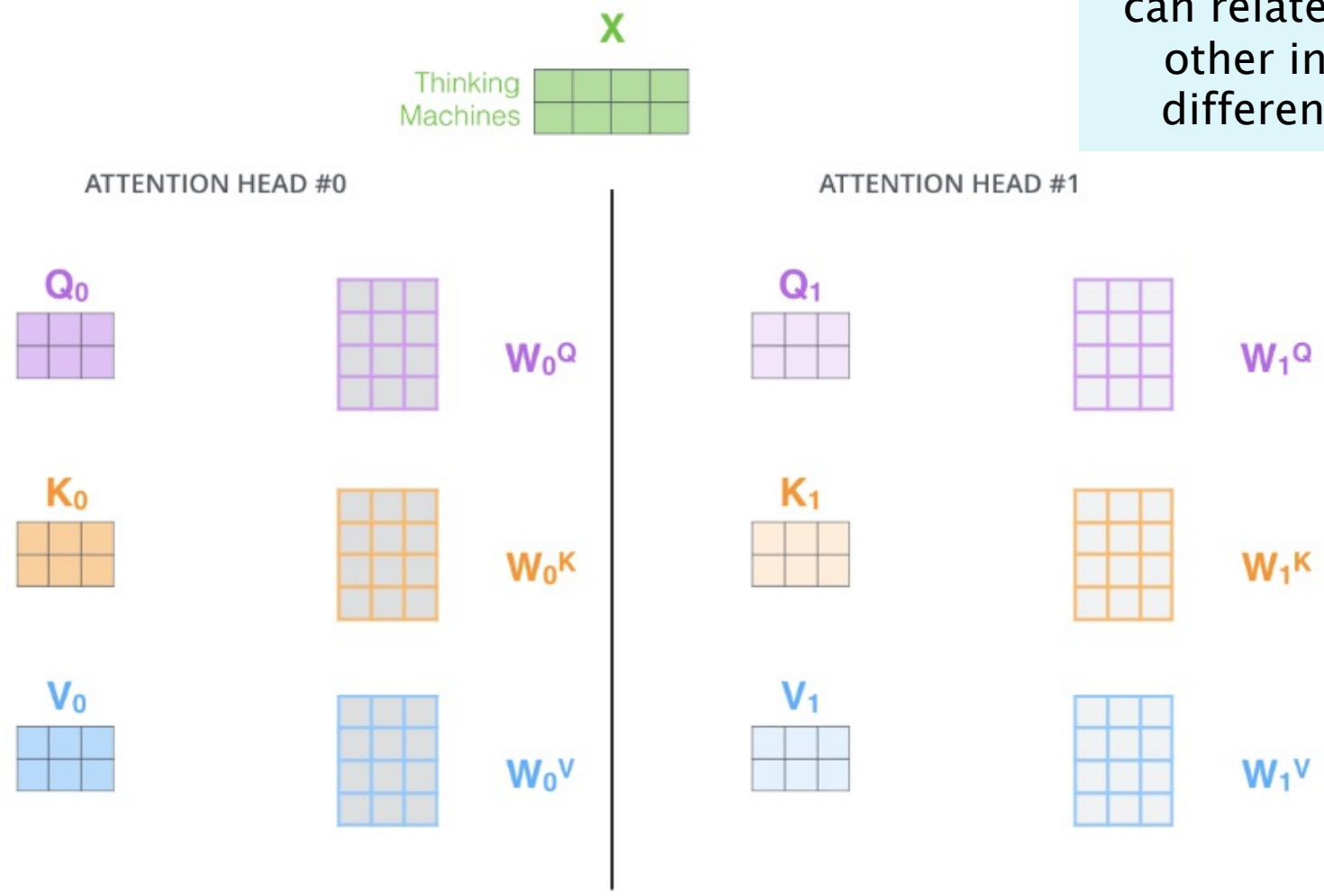
# Transformers – Encoders



The word at each position passes through a self-attention process. Then, they each pass through a feed-forward neural network -- the exact same network with each vector flowing through it separately.

# Multi-Headed Attention

Words in a sentence can relate to each other in many different ways

X

Thinking Machines

ATTENTION HEAD #0

$Q_0$

$W_0^Q$

$K_0$

$W_0^K$

$V_0$

$W_0^V$

ATTENTION HEAD #1

$Q_1$

$W_1^Q$

$K_1$

$W_1^K$

$V_1$

$W_1^V$

With multi-headed attention, we maintain separate Q/K/V weight matrices for each head resulting in different Q/K/V matrices. As we did before, we multiply X by the WQ/WK/WV matrices to produce Q/K/V matrices.
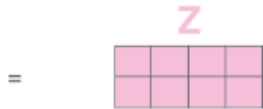
# Multi-Headed Attention

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

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

=

$W^O$

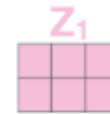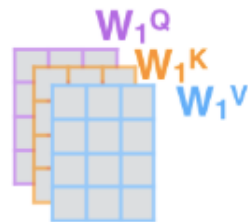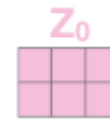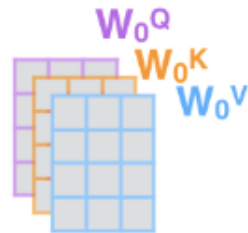# Attention – Summary



1) ...s our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer
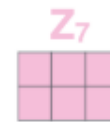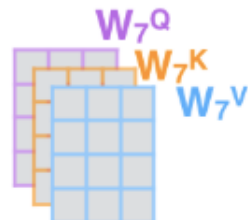
# Dense Feed Forward Layer

$$R_i = ReLU(Z_i W_1 + b_1) W_2 + b_2, \quad i = 1, \ldots, N$$



W₁       W₂

Z           R

Size d      Size 4d      Size d

# Add Residual Connections and Layer Normalization

A Single
Layer

# Layer Normalization

▸ Batch Normalization vs Layer Normalization

Batch Number

|       | B1 | B2 | B3 |
|-------|----|----|----|

Node Activations

A1      Batch Normalization

A2

A3

Layer Normalization

$$\mu_L = \frac{1}{d} \sum_{m=1}^{d} a(m)$$

$$\sigma_L^2 = \frac{1}{d} \sum_{m=1}^{d} (a(m) - \mu_L)^2$$

$$\hat{a}(m) = \frac{a(m) - \mu_L}{\sqrt{\sigma_L^2 + \epsilon}}$$

$$c(m) = \gamma \hat{a}(m) + \beta$$

# Multiple Layers

$$(W_3^Q, W_3^K, W_3^V), W_3^{FF}$$

$$(W_2^Q, W_2^K, W_2^V), W_2^{FF}$$

Parameters
not shared
Between layers

$$(W_1^Q, W_1^K, W_1^V), W_1^{FF}$$

Modularity: All elements of input sequence share the same parameters:
RNN like property

# Counting Number of Parameters

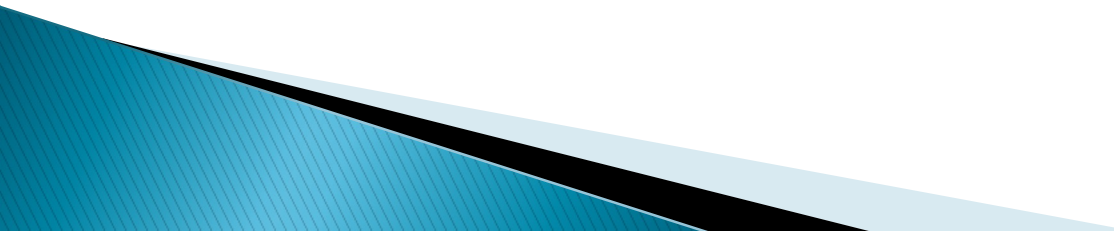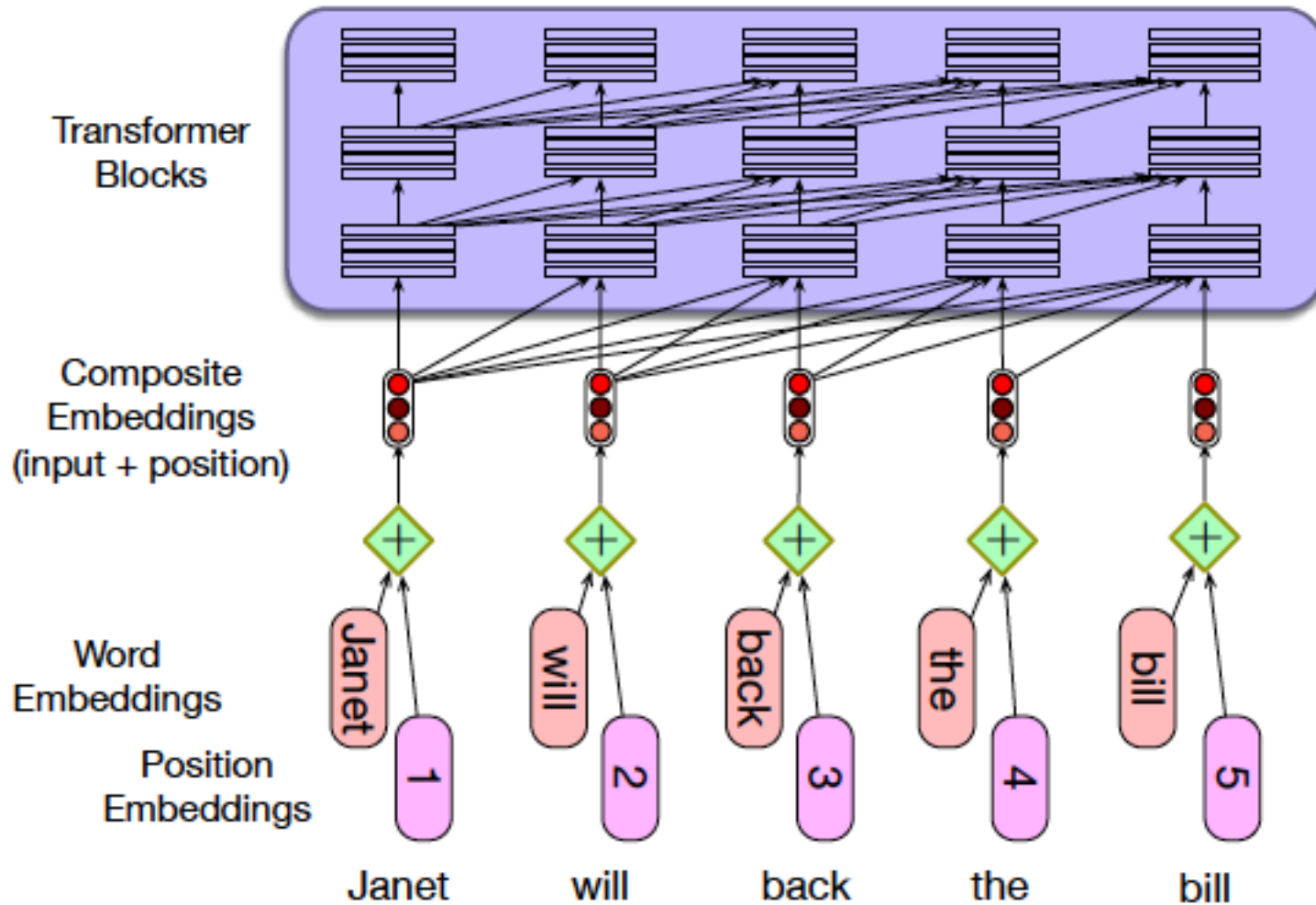| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | **Parameters in the Self Attention Block** | | | | # Parms | | Embedding Dimension | | | 32 |
| 2 | | | | | | | Number of Attention Heads | | | 1 |
| 3 | W^Q | | | 32 x 32 + 32 | 1,056 | | Dense Dimension | | | 32 |
| 4 | W^K | | | 32 x 32 + 32 | 1,056 | | Number of Blocks | | | 1 |
| 5 | W^V | | | 32 x 32 + 32 | 1,056 | | Sequence Length | | | 600 |
| 6 | | | | | | | | | | |
| 7 | Sub Total | | | | 3,168 | | | | | |
| 8 | | | | | | | | | | |
| 9 | **Number of Attention Heads** | | | | 2 | | | | | |
| 10 | | | | | | | | | | |
| 11 | Sub Total | | | | 6,336 | | | | | |
| 12 | | | | | | | | | | |
| 13 | **Parameters in the Projection Block** | | | 2 X 32 X 32 + 32 | 2,080 | | | | | |
| 14 | | | | | | | | | | |
| 15 | Sub Total | | | | 8,416 | | | | | |
| 16 | | | | | | | | | | |
| 17 | **Parameters in Dense Feed Forward Block** | | | | | | | | | |
| 18 | | | | | | | | | | |
| 19 | DFN1 | | | 32 X 32 + 32 | 1,056 | | | | | |
| 20 | DFN2 | | | 32 X 32 + 32 | 1,056 | | | | | |
| 21 | | | | | | | | | | |
| 22 | Sub Total | | | | 10,528 | | | | | |
| 23 | | | | | | | | | | |
| 24 | **Layer Normalization 1** | | | 32 X 2 | 64 | | | | | |
| 25 | | | | | | | | | | |
| 26 | **Layer Normalization 2** | | | 32 X 2 | 64 | | | | | |
| 27 | | | | | | | | | | |
| 28 | **Grant Total** | | | | **10,656** | | | | | |
| 29 | | | | | | | | | | |
| 30 | | | | | | | | | | |
| 31 | **Self Attention Computations** | | | | | | | | | |
| 32 | | | | | | | | | | |
| 33 | For a single element | | | 32 X 32 x 600 | 614,400 | | | | | |
| 34 | For Entire Sequence | | | ( ) X 600 | 368,640,000 | | | | | |
| 35 | | | | | | | | | | |

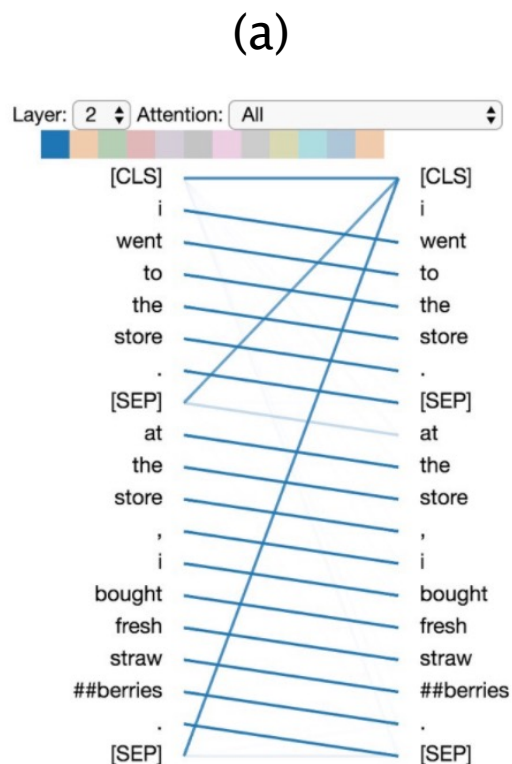# Positional Encoding

- Positional information arises naturally in RNN/LSTMs
- Transformer Architecture is invariant to permutations of the input sequence
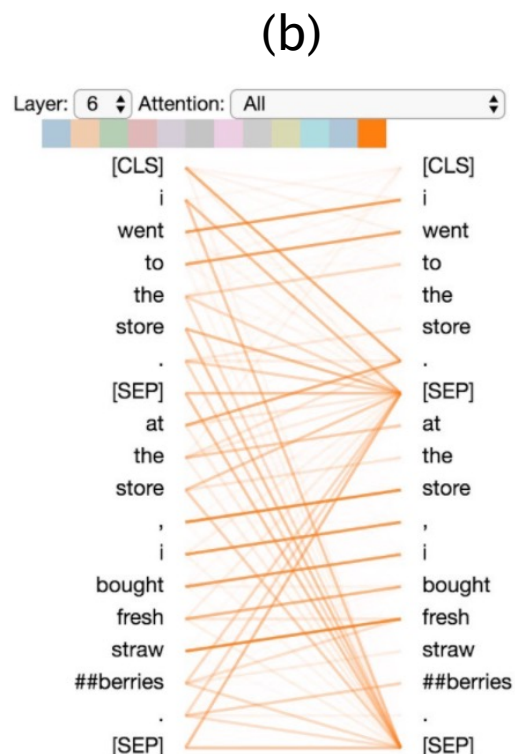- This is a problem if the position is important Example: NLP

# Positional Encoding

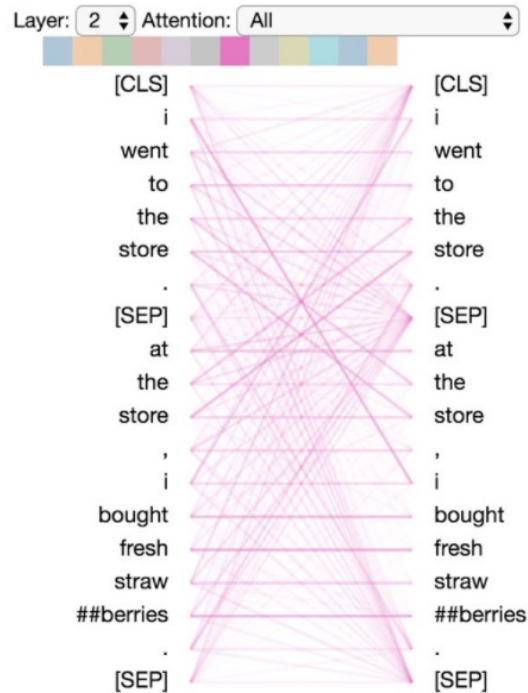# Visualizing Attention in Transformers

(a)

(b)



Attention to Next Word
Similar to a Backwards RNN

Attention to Prior Word
Similar to a Forward RNN
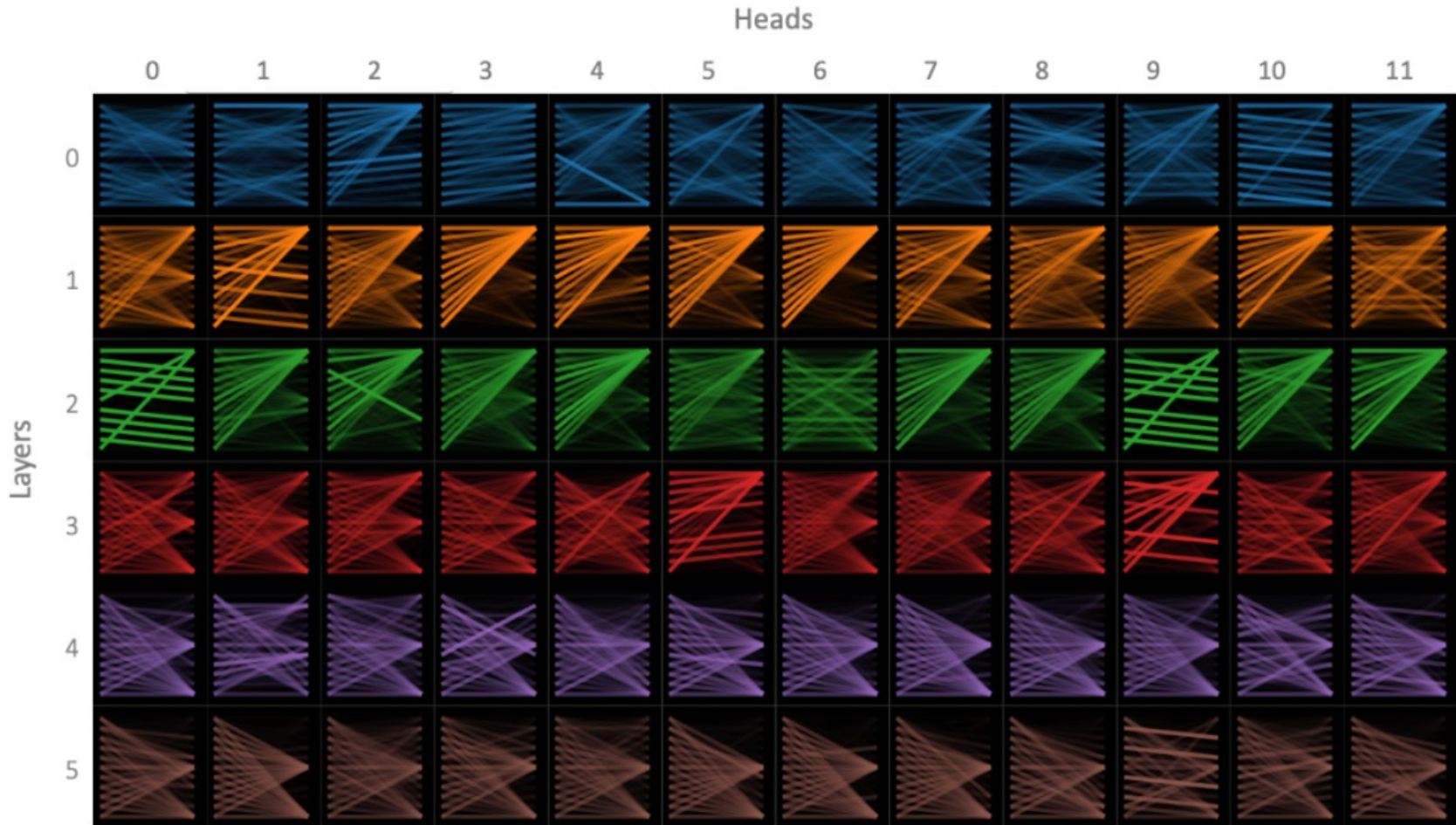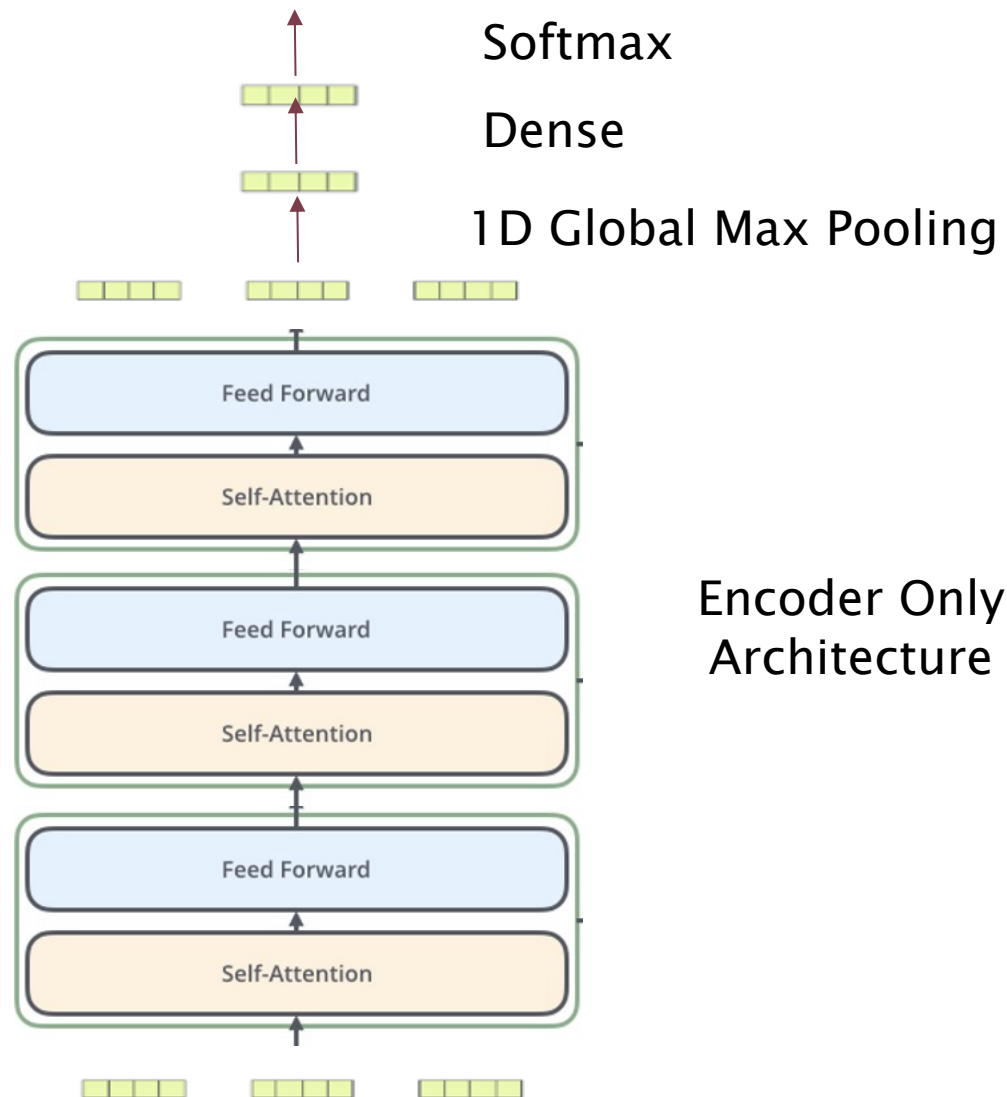
# Visualizing Attention



(c)

(d)

Attention to Similar Words

# Visualizing Attention

# Classification using Transformers

Softmax

Dense

1D Global Max Pooling

Feed Forward

Self-Attention

Encoder Only
Architecture

Feed Forward

Self-Attention

Feed Forward

Self-Attention

# Further Reading

- Das and Varma: Chapter Transformers
- Chollet (2nd Edition): Chapter 11, Section11.4
- http://jalammar.github.io/illustrated-transformer/