# Attentions Is All You Need

## Transformers

Eng. Mohammad Fouad Shikh Ali

# Attentions Is All You Need

Part 1

# Tokenization

"Where can I find a pizzeria?"

## Word-Level Tokenization

[where, can, i, find, a, pizzeria]

➕ Intuitive.

➖ Doesn't handle OOV words including new words, slang, play on words, misspellings, etc.

➖ Huge vocabulary for large corpora; especially for languages with rich morphology (e.g. Hungarian).

➖ Handling punctuation is challenging (e.g. "don't" vs "N.Y.C").

## Char-Level Tokenization

[a, c, d, e, f, h, i, n, p, r, w, z]

➕ Small memory footprint.

➕ Handles OOV words.

➖ Needs to go over all characters and learn a particular sequence for a given word.

➖ Loss of performance.

# Subword Tokenization



"where can I find a pizzeria?"

where, can, I, find, a, pi, zz, eria

listeria

⇩

[list, eria]

*Subword tokenization* has a better chance of handling OOV words while reducing vocabulary size and maintaining performance.

# Byte Pair Encoding (BPE)
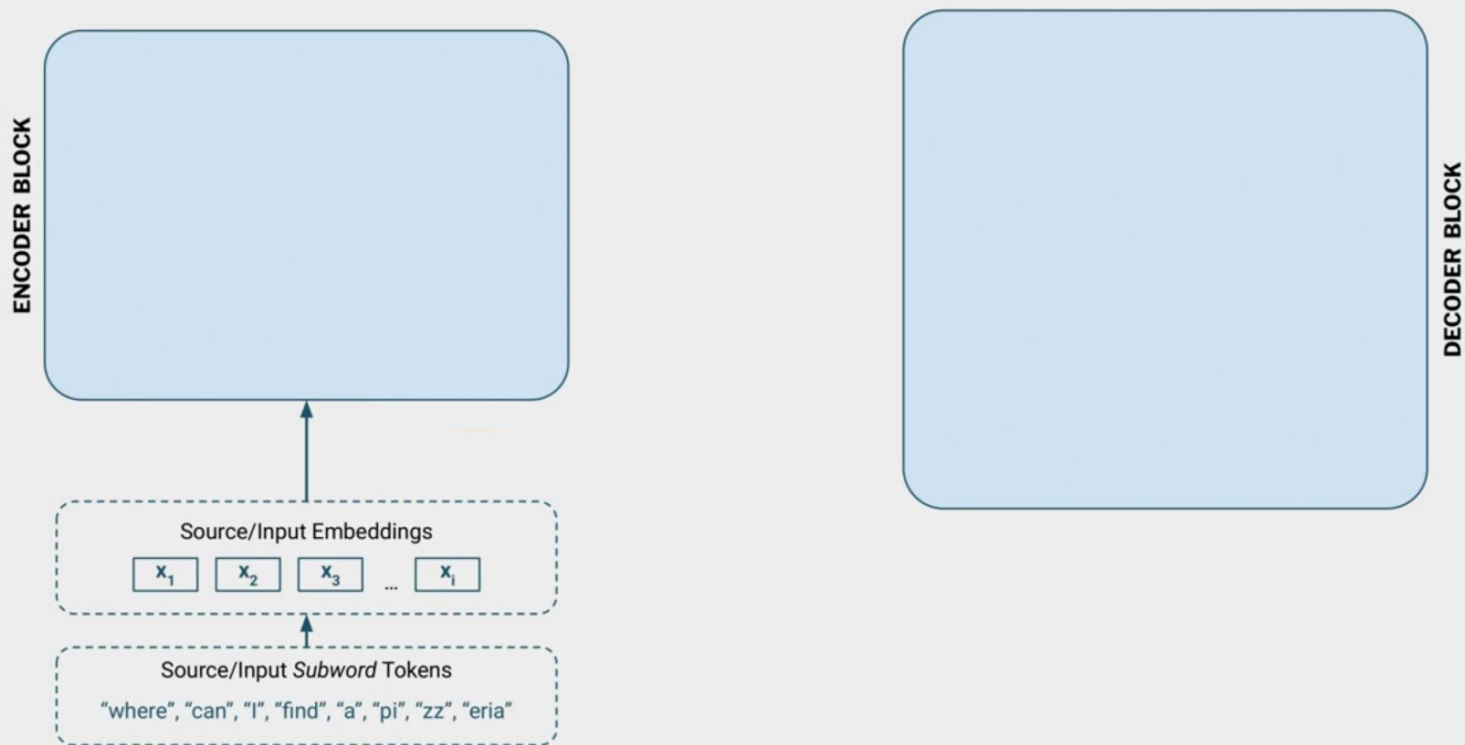
"she sells seashells by the seashore"

**Corpus**

she_
sells_
seashells_
by_
the_
seashore_

**Vocabulary**
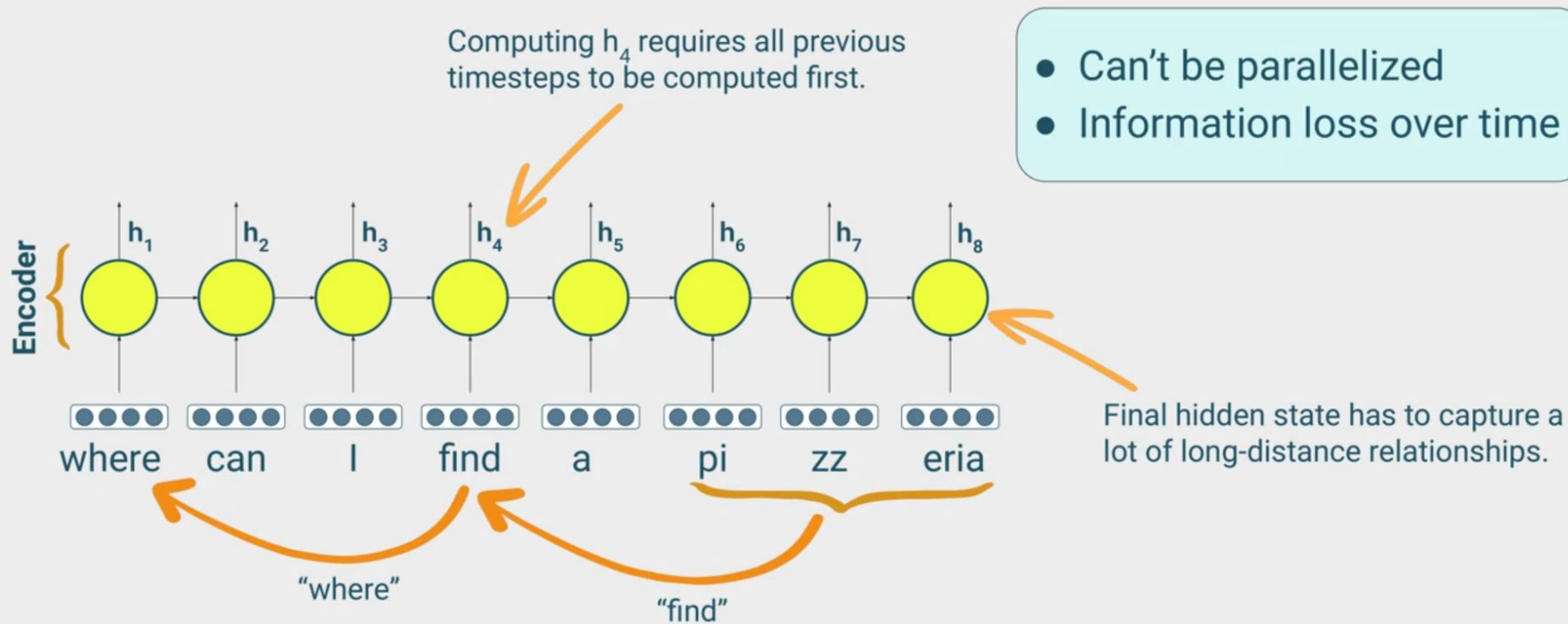
_, a, b, e, h, l, o, r, s, t, y
_, a, b, e, h, l, o, r, s, t, y, sh
_, a, b, e, h, l, o, r, s, t, y, sh, he
_, a, b, e, h, l, o, r, s, t, y, sh, he, e_
_, a, b, e, h, l, o, r, s, t, y, sh, he, e_, se
_, a, b, e, h, l, o, r, s, t, y, sh, he, e_, se, she
.
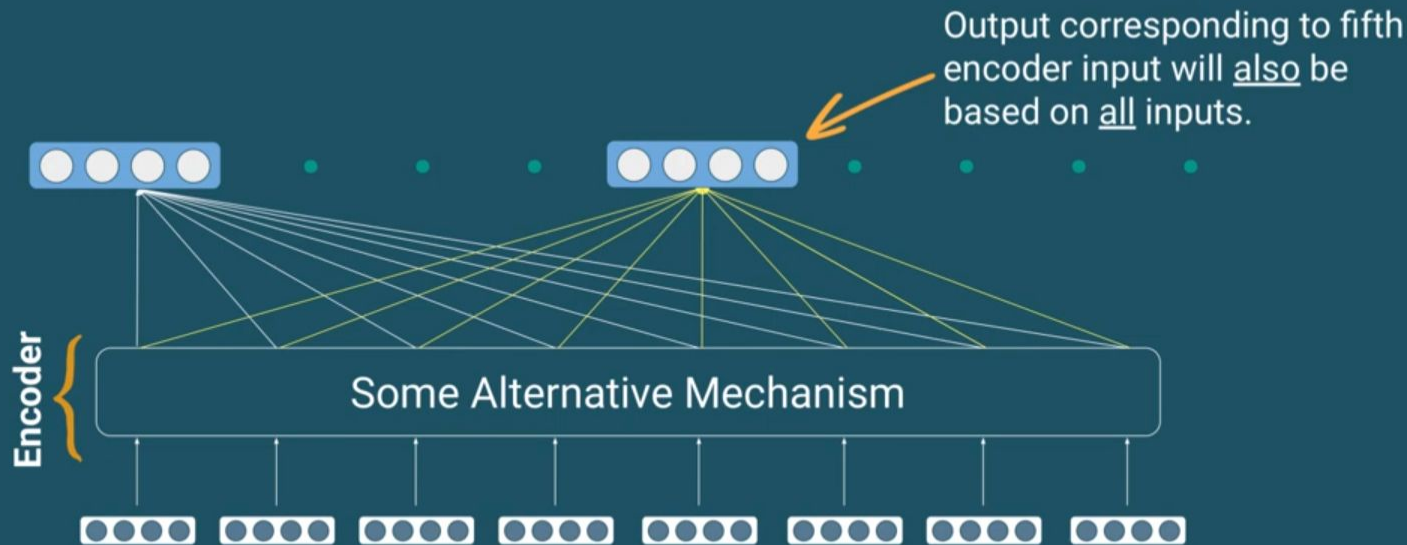.
.

continue until N merges are performed.

# Transformers

ENCODER BLOCK

DECODER BLOCK

Source/Input Embeddings

| $x_1$ | $x_2$ | $x_3$ | ... | $x_i$ |

Source/Input *Subword* Tokens

"where", "can", "I", "find", "a", "pi", "zz", "eria"

# Why Transformers



Computing $h_4$ requires all previous timesteps to be computed first.

- Can't be parallelized
- Information loss over time

Encoder

$h_1$   $h_2$   $h_3$   $h_4$   $h_5$   $h_6$   $h_7$   $h_8$

where   can   I   find   a   pi   zz   eria

Final hidden state has to capture a lot of long-distance relationships.

"where"

"find"

...throw away recurrence, and base EACH encoder OUTPUT on ALL the encoder INPUTS.

Output corresponding to fifth encoder input will **also** be based on **all** inputs.

Encoder

Some Alternative Mechanism

Do this by having the encoder perform attention **on itself**.

*"Self-Attention"*

# Self-Attention



Updated embedding for $x_1$ now includes information from all other embeddings.

$$\hat{x}_l = \sum_j \alpha_{ij} x_j \quad \} \quad a_{1,1} * x_1 + a_{1,2} * x_2 + a_{1,3} * x_3$$

$a_{1,1}$ | $a_{1,2}$ | $a_{1,3}$  $\}$ Attention weights

SOFTMAX

$s_{1,1}$ | $s_{1,2}$ | $s_{1,3}$  $\}$ Attention scores

$X_1$     $X_1$     $X_1$

$X_1$     $X_2$     $X_3$

$X_1$     $X_2$     $X_3$

# Self-Attention



$$\hat{x}_2 = \sum_j \alpha_{ij} x_j \quad \Big\} \quad a_{2,1} * x_1 + a_{2,2} * x_2 + a_{2,3} * x_3$$

Simple
Self-Attention

| $a_{2,1}$ | $a_{2,2}$ | $a_{2,3}$ | } Attention weights |

SOFTMAX

| $s_{2,1}$ | $s_{2,2}$ | $s_{2,3}$ | } Attention scores |

$x_2$ $x_2$ $x_2$

$x_1$ $x_2$ $x_3$

$x_1$ $x_2$ $x_3$

# Self-Attention



Each updated embedding now includes information from all other embeddings.

# Self-Attention

Each updated embedding now includes information from all other embeddings.



Simple Self-Attention Layer

Simple Self-Attention Layer

$x_1$

$x_2$

$x_3$

# RNNs (Encoder-Decoder)



$$c_j = \sum_i \alpha_{i\ j} h_{\ i}$$

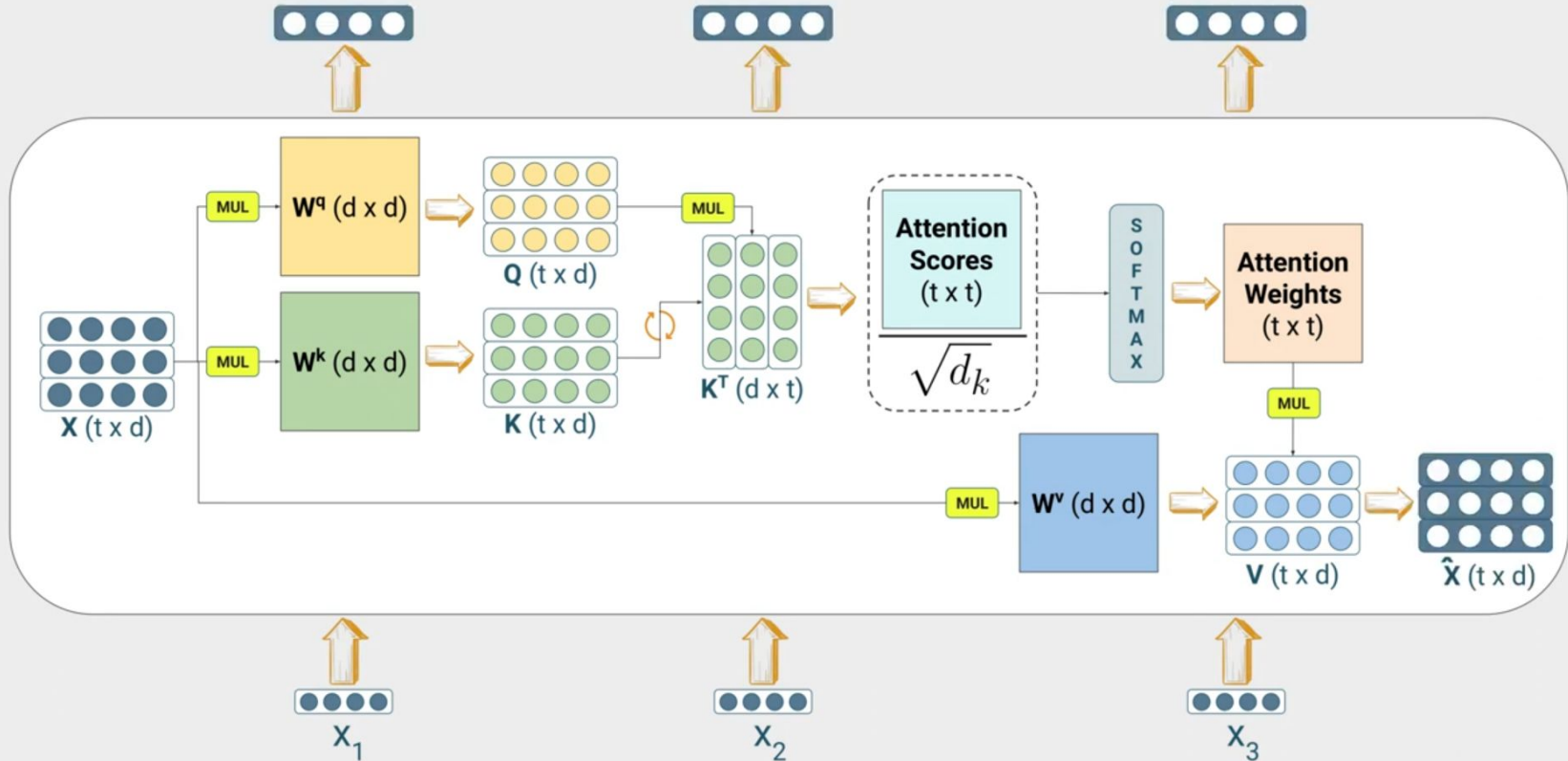Rather than returning a single value, a blend of ALL values is returned (*attention* is more like mixing paint).
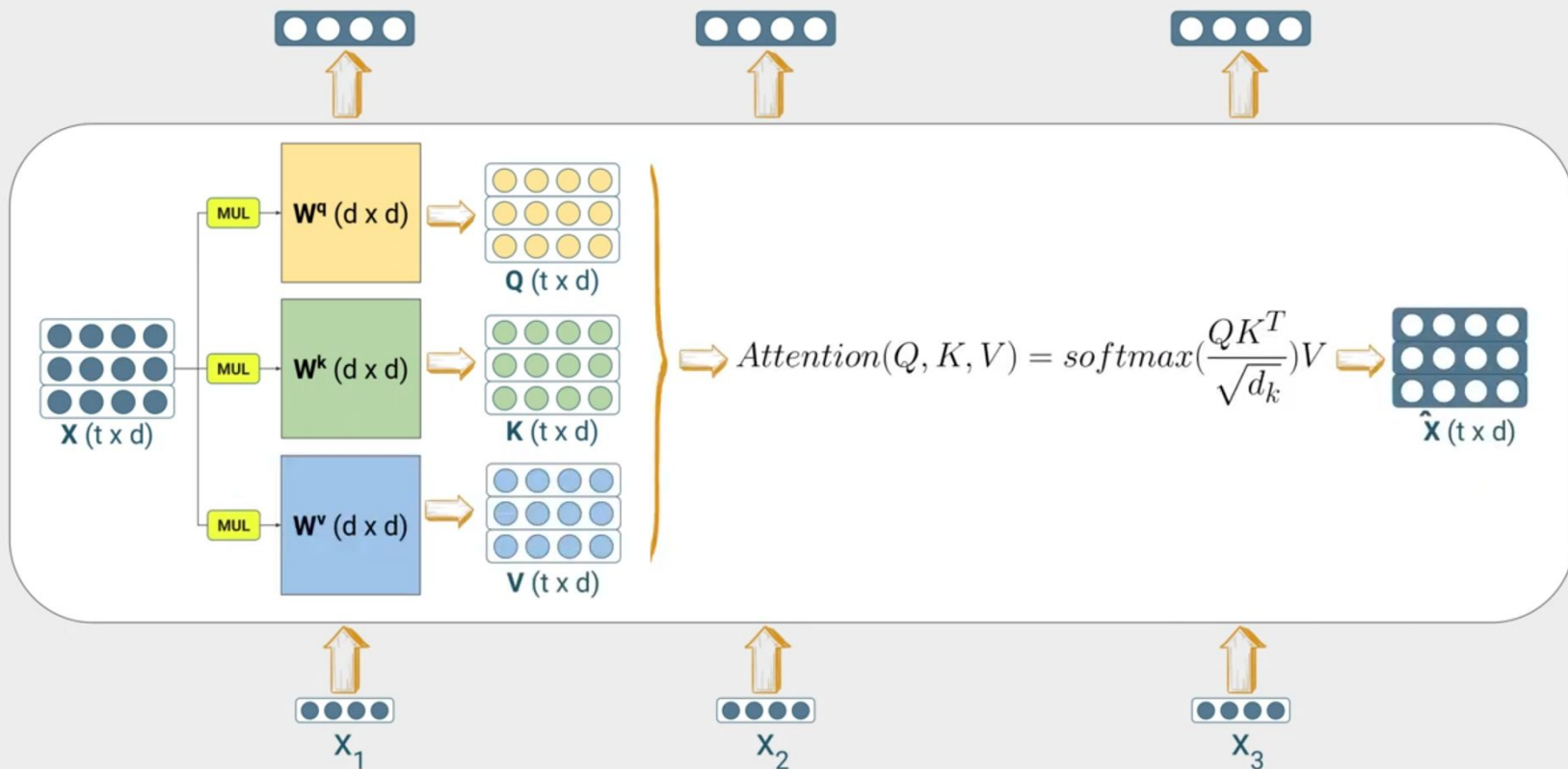
# Scale Dot Product

$$\hat{x}_l = \sum_j a_{ij} v_j$$

Calculate Attention Weights

$q_1$   $k_1$   $q_1$   $k_2$   $q_1$   $k_3$

| $a_{1,1}$ | $a_{1,2}$ | $a_{1,3}$ |

*Attention weights*

SOFTMAX

| $s_{1,1}$ | $s_{1,2}$ | $s_{1,3}$ |

$$\sqrt{d_k}$$

*Scaled* Attention scores

query   key   value    query   key   value    query   key   value

$W^k$   $W^q$   $W^v$   $x_1$

$W^k$   $W^q$   $W^v$   $x_2$

$W^k$   $W^q$   $W^v$   $x_3$

# Scale Dot Product Self-Attention

# Attention Head

"**Sarah** **went** to **a restaurant** to **meet** her **friend** **that night.**"

- **What?**
- **Where?**
- **Who?**
- **When?**

# Multi-Head Self-Attention

**HEAD 0**

Each head will have its own set of Q, K, and V weights of dimension $d/h$ where $h$ is the number of heads.

$W_0^q$  $W_0^k$  $W_0^v$

MUL  MUL  MUL

$X$ (t x d)

# Multi-Head Self-Attention

# Paper



**Original Transformer Dimensions:**
- Embedding dimension of 512. This is known as the *model dimension* ($d_{model}$).
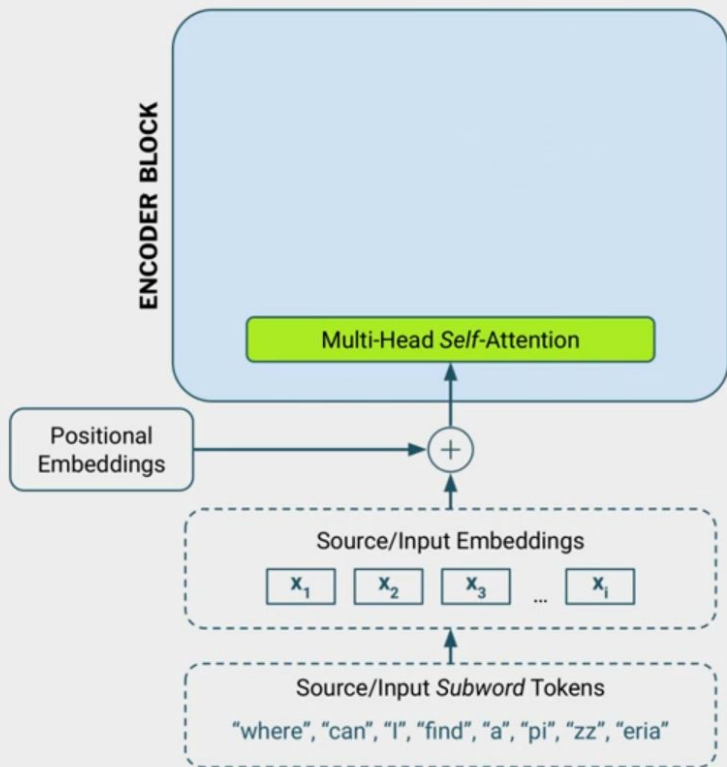
- 8 attention heads, so d/h = 64.

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

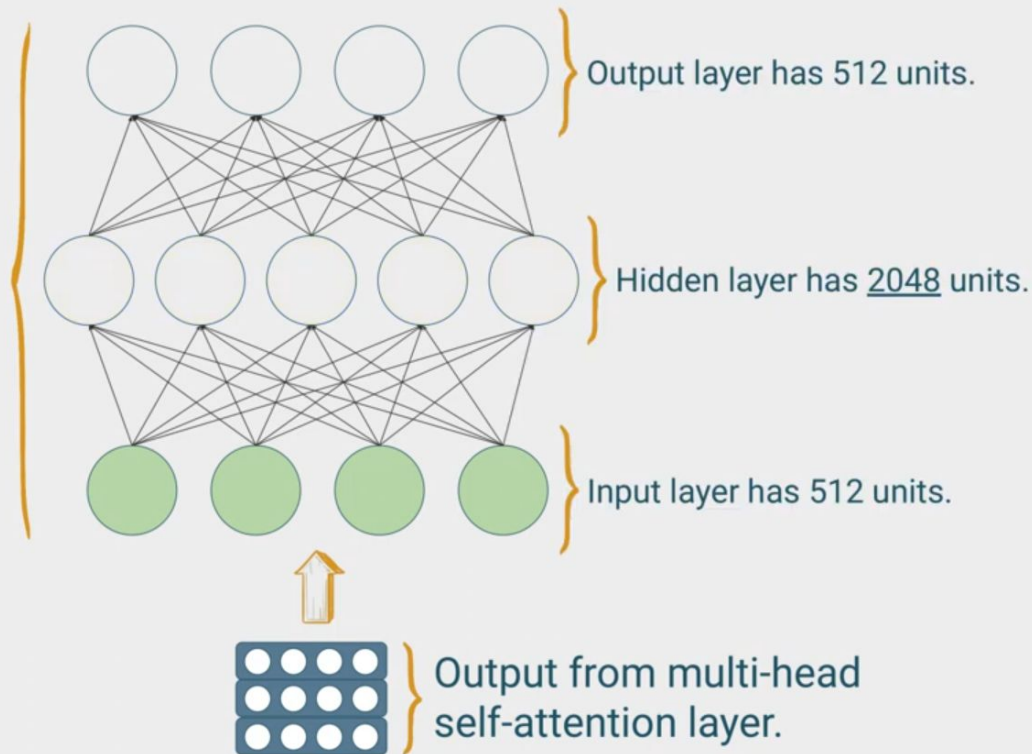$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

# Transformers



**ENCODER BLOCK**

Multi-Head *Self*-Attention

Positional Embeddings

Source/Input Embeddings

$x_1$  $x_2$  $x_3$  ...  $x_i$

Source/Input *Subword* Tokens

"where", "can", "I", "find", "a", "pi", "zz", "eria"

We added a number of learnable weights but there's still no non-linearity...

# Adding Non-Linearity



Original transformer used a two-layer network with a ReLU activation in the hidden layer.

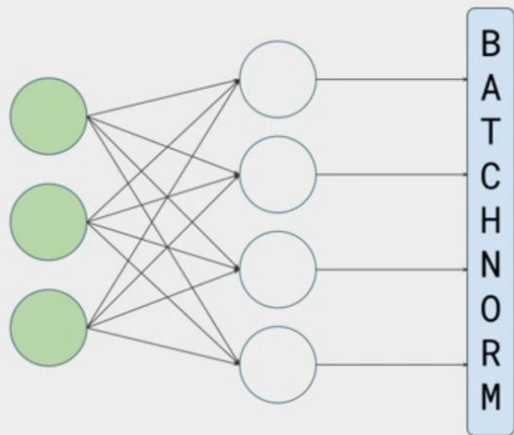Output layer has 512 units.

Hidden layer has 2048 units.

Input layer has 512 units.

Output from multi-head self-attention layer.

# Transformers



| $h_1$ | $h_2$ | $h_3$ | ... | $h_i$ |

Encoder output embeddings have the same dimensions as input embeddings (512).

**ENCODER BLOCK**

Feedforward Neural Network
(*applied pointwise*)

Multi-Head *Self*-Attention

Positional Embeddings

$+$

Source/Input Embeddings

| $x_1$ | $x_2$ | $x_3$ | ... | $x_i$ |

Source/Input *Subword* Tokens

"where", "can", "I", "find", "a", "pi", "zz", "eria"

# Transformers



Two potential issues:
1. Shifting inputs from earlier encoder blocks add noise.

2. Depth leads to earlier information (e.g. positional embeddings) being "forgotten" over blocks, and vanishing gradients.

# Batch Normalization



$$\boldsymbol{\mu}_b = \frac{1}{m_b} \sum_{i=1}^{m_b} \mathbf{h}_i \Big\} \text{ Vector of mini-batch output means.}$$

$$\boldsymbol{\sigma}_b^2 = \frac{1}{m} \sum_{i=1}^{m_b} (\mathbf{h}_i - \boldsymbol{\mu}_b)^2 \Big\} \text{ Vector of mini-batch output standard deviations.}$$

$$\hat{\mathbf{h}}_i = \frac{\mathbf{h}_i - \boldsymbol{\mu}_b}{\sqrt{\boldsymbol{\sigma}_b^2 + \epsilon}} \Big\} \text{ Vector of standardized outputs.}$$

$$\mathbf{z}_i = \boldsymbol{\gamma} \odot \hat{\mathbf{h}}_i + \boldsymbol{\beta} \Big\} \text{ Vector of } \textit{scaled} \text{ and } \textit{shifted} \text{ outputs.}$$

# Batch Normalization

# Layer Normalization
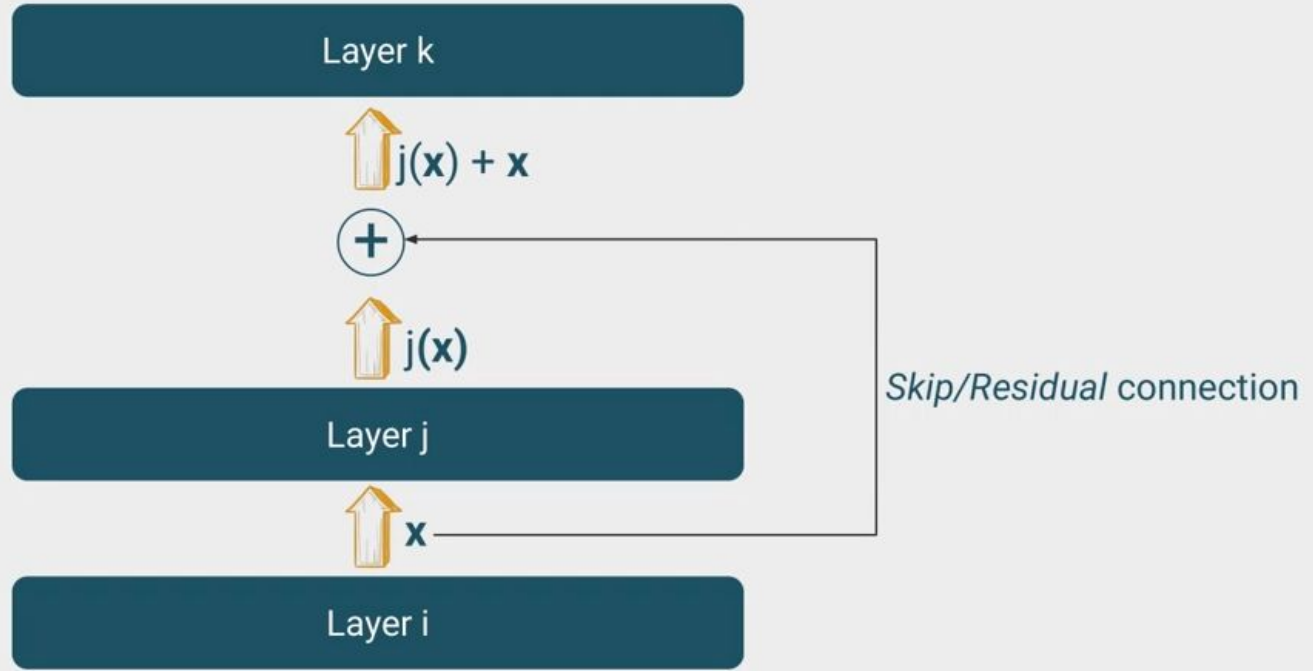


$$\mu = \frac{1}{d_h} \sum_{i=1}^{d_h} x_i$$

$$\sigma = \sqrt{\frac{1}{d_h} \sum_{i=1}^{d_h} (x_i - \mu)^2}$$

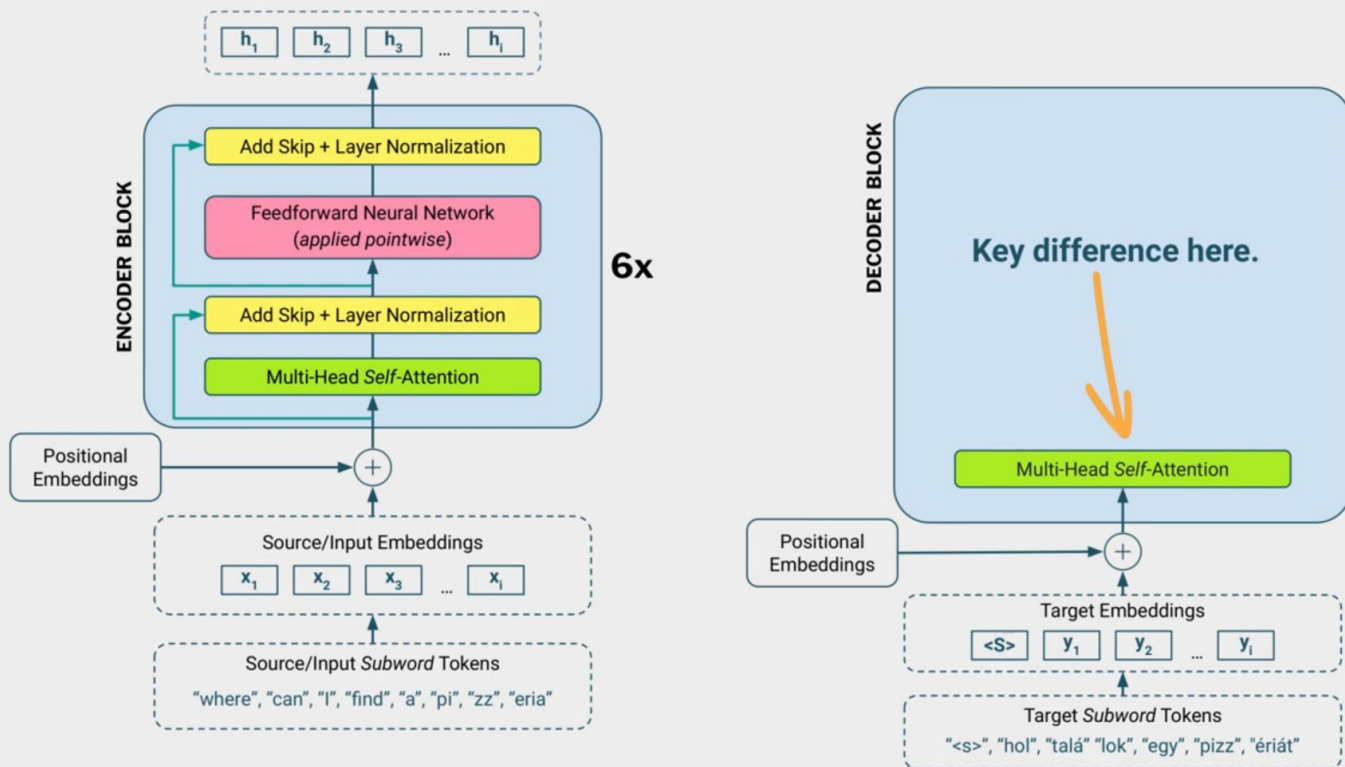$$\hat{\mathbf{x}} = \frac{\mathbf{x} - \mu}{\sigma}$$

Learnable parameters to scale and shift as needed.

$$\mathbf{z} = \gamma \hat{\mathbf{x}} + \beta$$
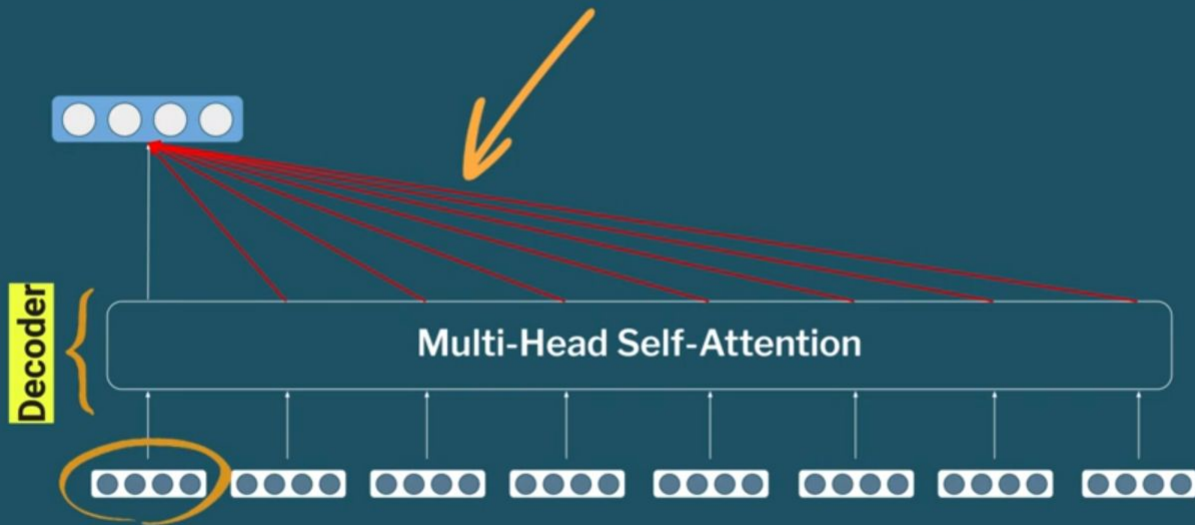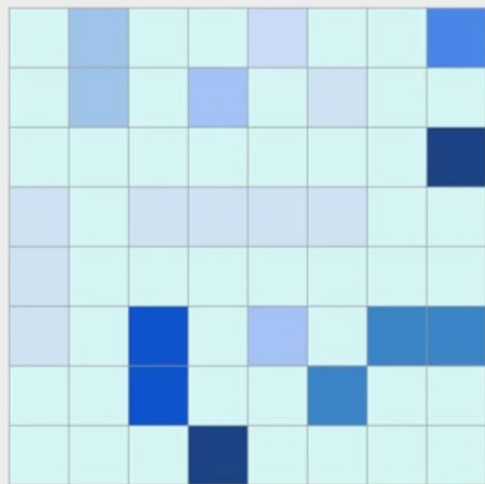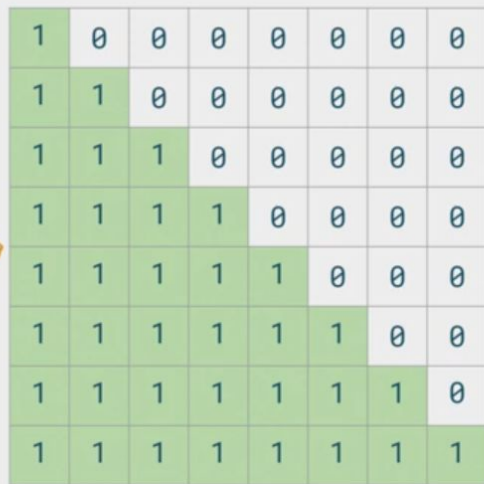
# Skip\Residual Connection

# Transformers

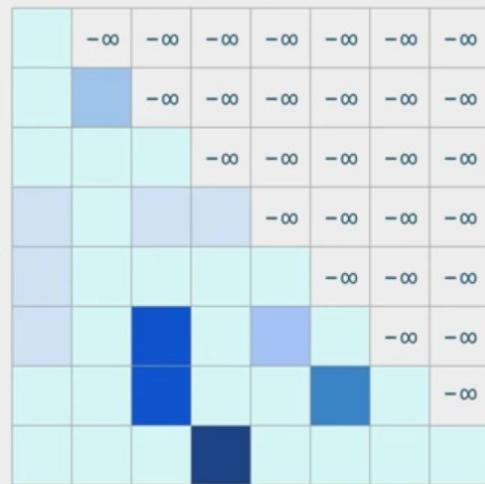We need to block decoder from accessing any future parts of the sequence.

Decoder

Multi-Head Self-Attention

Attention scores → Attention mask → Masked Attention scores → SOFTMAX

Attention weights

Masked Multi-Head Self-Attention