# Lecture 10: Attention mechanism

**Radoslav Neychev**
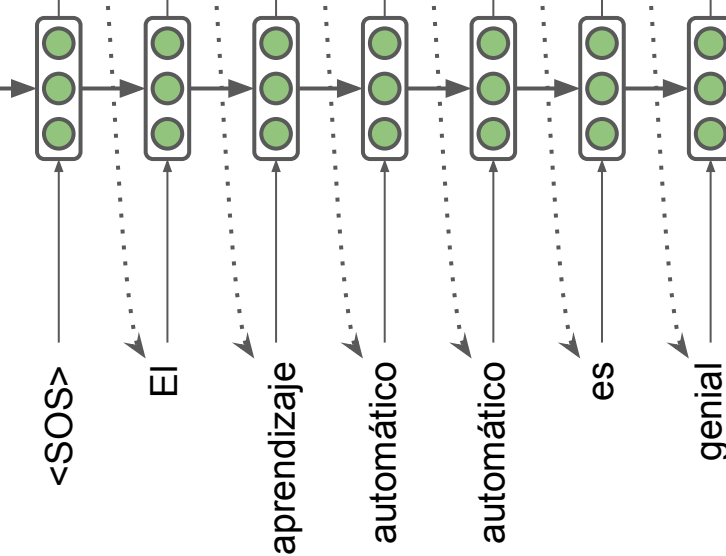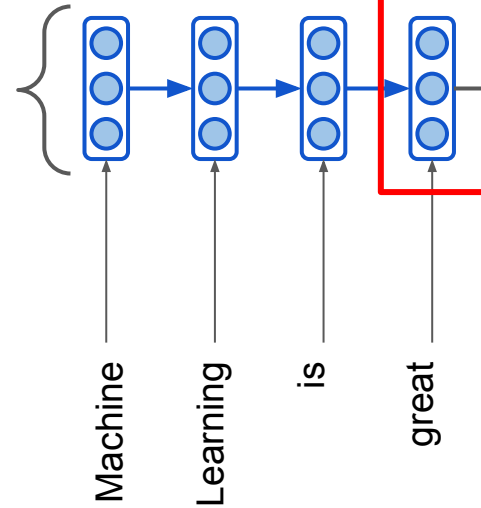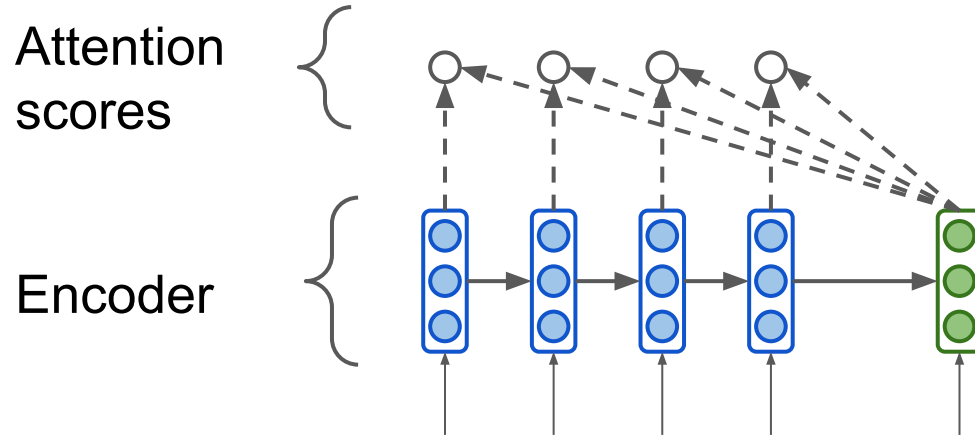
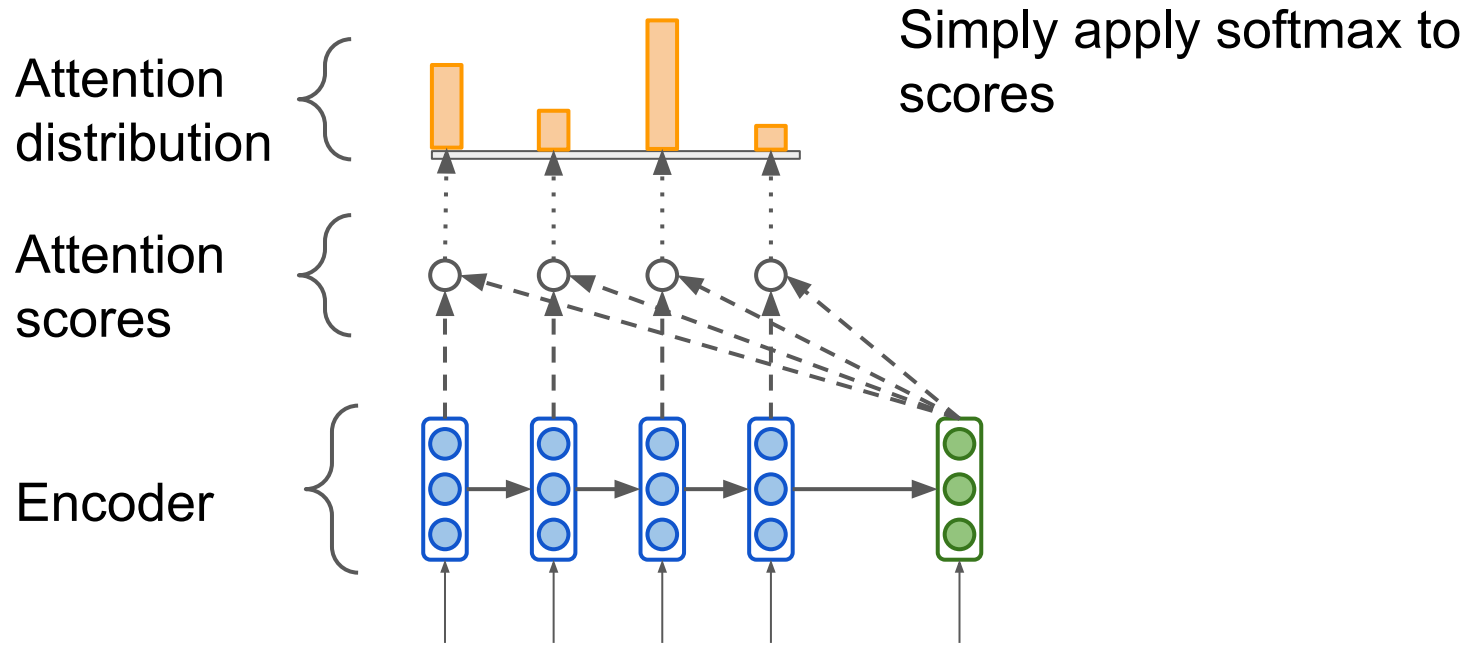# Attention

Seq2seq NMT

## **Main idea:**

on each step of the **decoder**, use **direct connection to the encoder** to focus on a particular part of the source sequence

# Seq2seq with attention



Attention scores

Encoder

# Seq2seq with attention

Attention
distribution

Attention
scores

Encoder

Simply apply softmax to
scores

Attention output

Weighted sum of all encoder states

Attention distribution

Attention scores

Encoder

Attention
output

Attention
distribution

Attention
scores

Encoder

Concatenate

# Seq2seq with attention

Attention output

Attention distribution

Attention scores

Encoder

y

# Seq2seq with attention

Attention output

Attention distribution

Attention scores

Encoder

y

Attention output

Attention distribution

Attention scores

Encoder

y

Denote encoder hidden states $\mathbf{h}_1, \ldots, \mathbf{h}_N \in \mathbb{R}^k$

and decoder hidden state at time step t $\quad \mathbf{s}_t \in \mathbb{R}^k$

The attention scores $\mathbf{e}^t$ can be computed as dot product
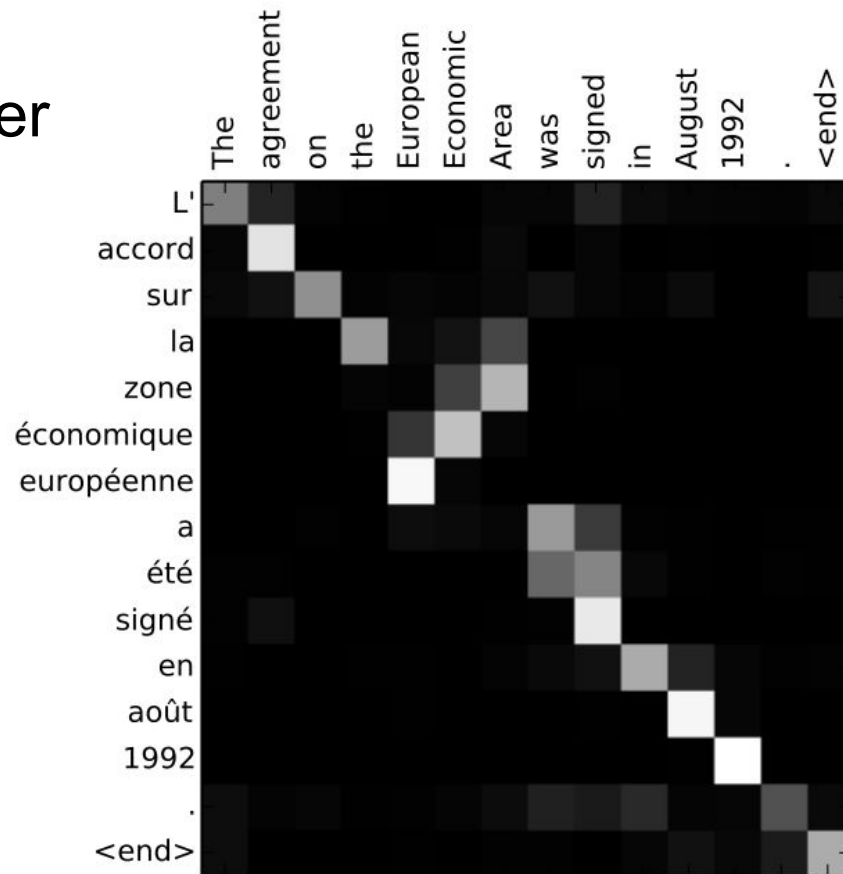
$$\mathbf{e}^t = [\mathbf{s}^T \mathbf{h}_1, \ldots, \mathbf{s}^T \mathbf{h}_N]$$

Then the attention vector is a linear combination of encoder states

$$\mathbf{a}_t = \sum_{i=1}^{N} \boldsymbol{\alpha}_i^t \mathbf{h}_i \in \mathbb{R}^k \text{ , where } \boldsymbol{\alpha}_t = \text{softmax}(\mathbf{e}_t)$$
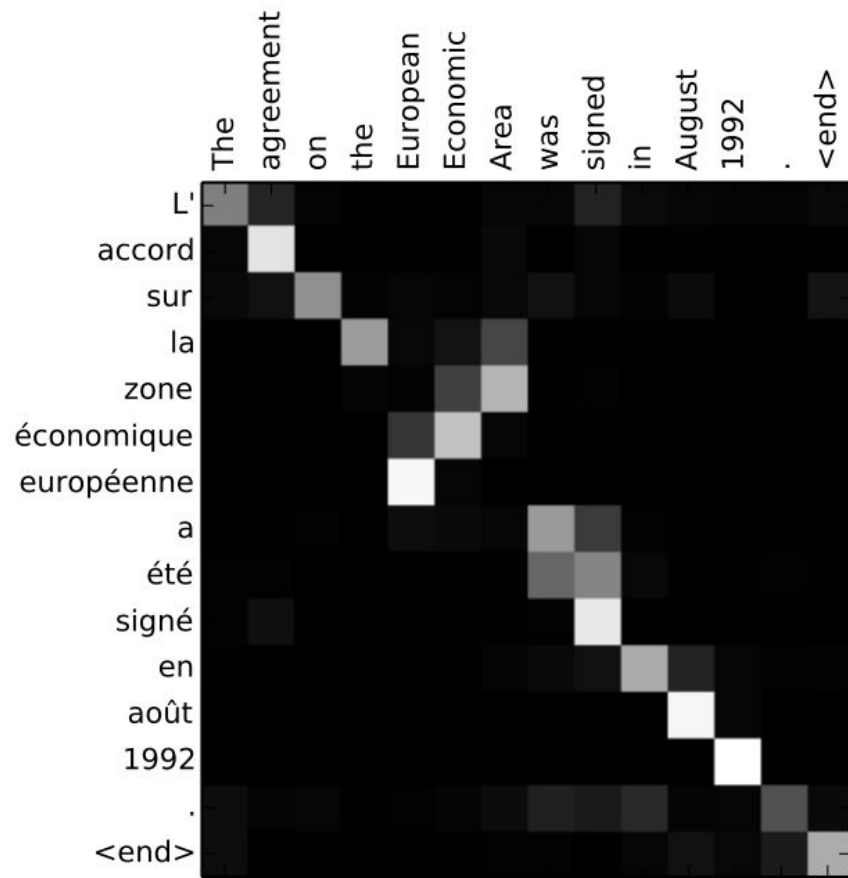
# Attention provides interpretability

- We may see what the decoder was focusing on
- We get word alignment for free!

Image source: Neural Machine Translation by Jointly Learning to Align and Translate

# Attention advantages

- "Free" word alignment
- Better results on long sequences

with attention

without attention



Image source: Neural Machine Translation by Jointly Learning to Align and Translate

- Basic dot-product (the one discussed before): $e_i = s^T h_i \in \mathbb{R}$
- Multiplicative attention: $e_i = s^T W h_i \in \mathbb{R}$
  - $W \in \mathbb{R}^{d_2 \times d_1}$ - weight matrix
- Additive attention: $e_i = v^T \tanh(W_1 h_i + W_2 s) \in \mathbb{R}$
  - $W_1 \in \mathbb{R}^{d_3 \times d_1}, W_2 \in \mathbb{R}^{d_3 \times d_2}$ - weight matrices
  - $v \in \mathbb{R}^{d_3}$ - weight vector

# Self-Attention
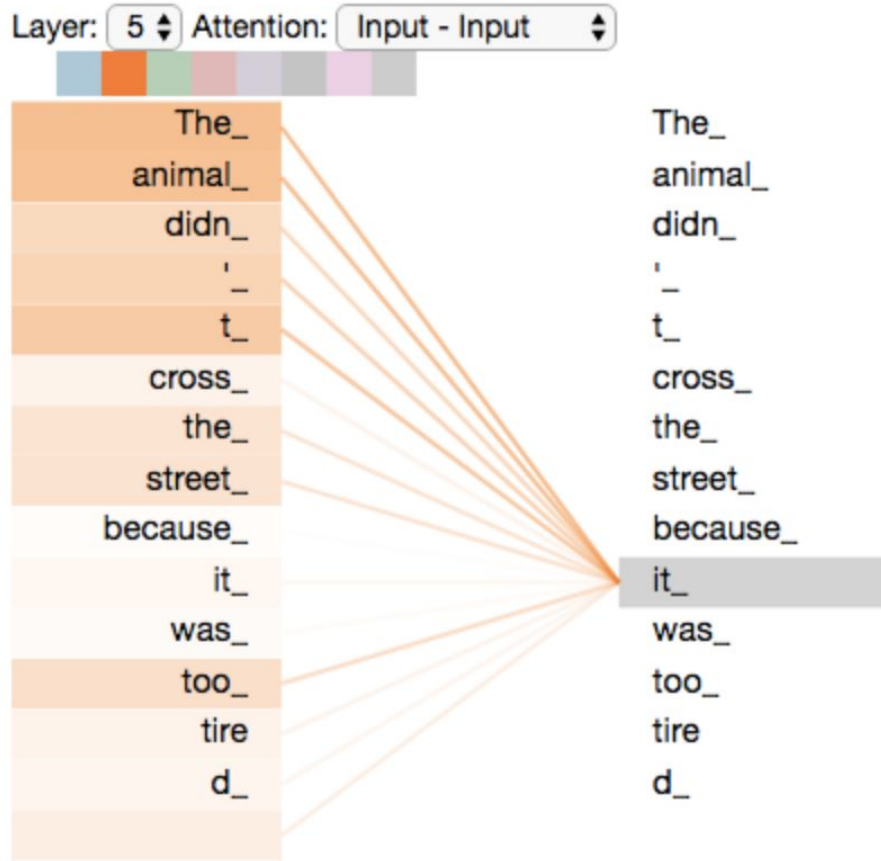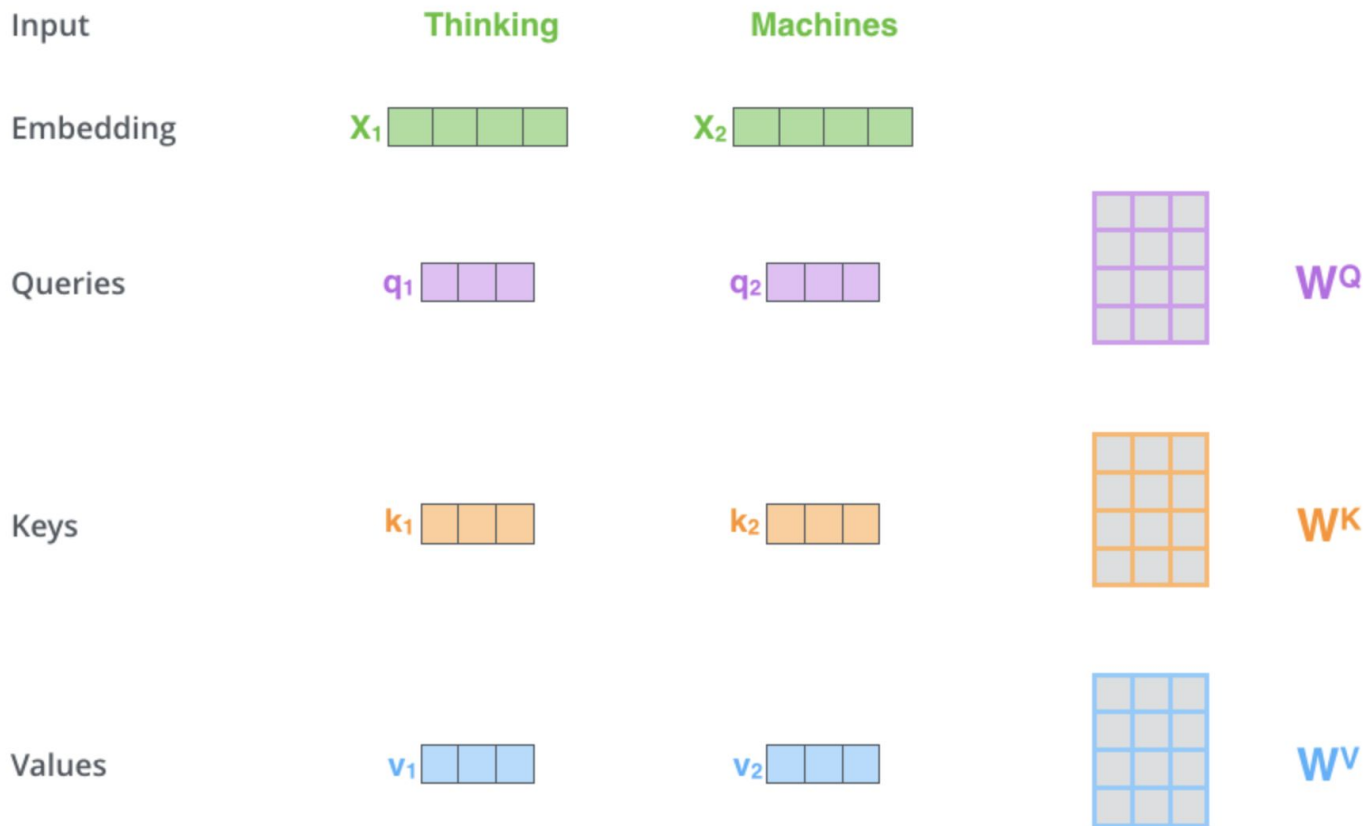
"The animal didn't cross the street because it was too tired"

- What does "it" in this sentence refer to?
- We want self-attention to associate "it" with "animal"

- Self-attention is the method the Transformer uses to bake the "understanding" of other relevant words into the one we're currently processing

# Self-Attention at a High Level

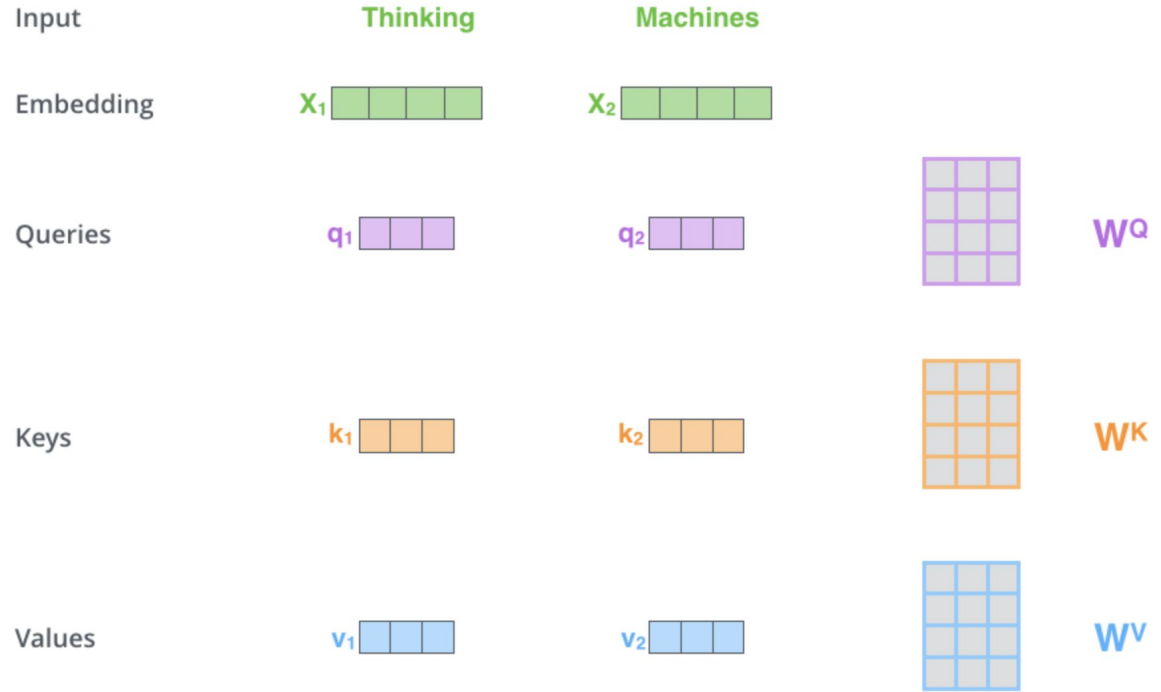# Self-Attention: detailed explanation

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

# Self-Attention: detailed explanation

**STEP 1:**

create 3 vectors
(**query**, **key**, **value**)

from each of the encoder's
input vectors

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

# Self-Attention: detailed explanation

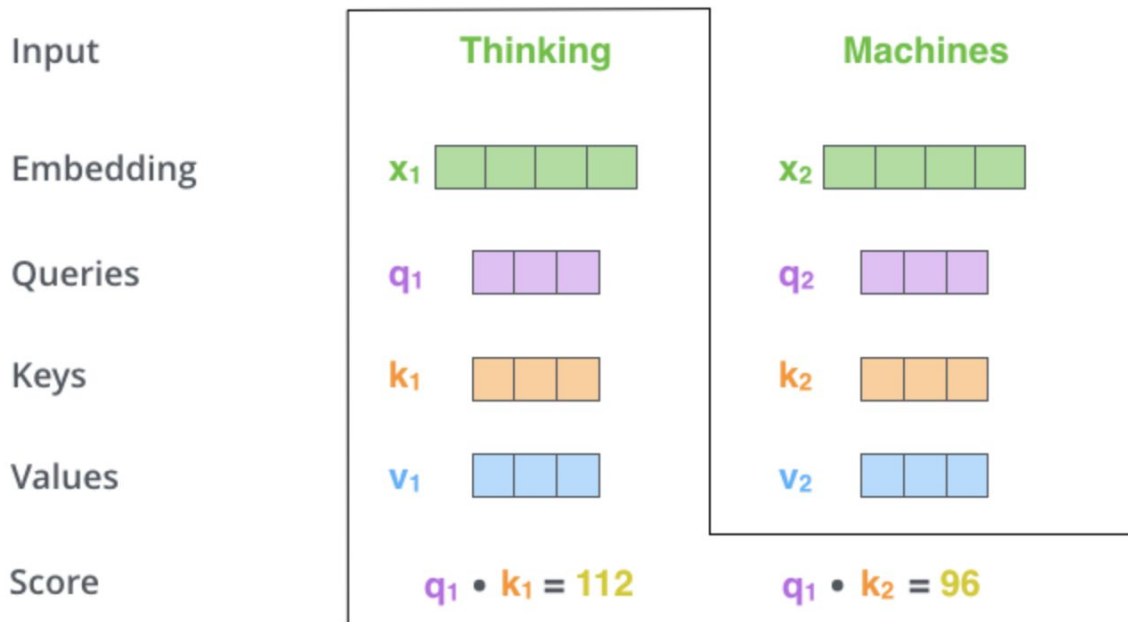What are the **query**, **key**, **value** vectors?

They're abstractions that are useful for

calculating and thinking about attention.

# Self-Attention: detailed explanation

## STEP 2:

calculate a score

(score each word of the input sentence against the current word)

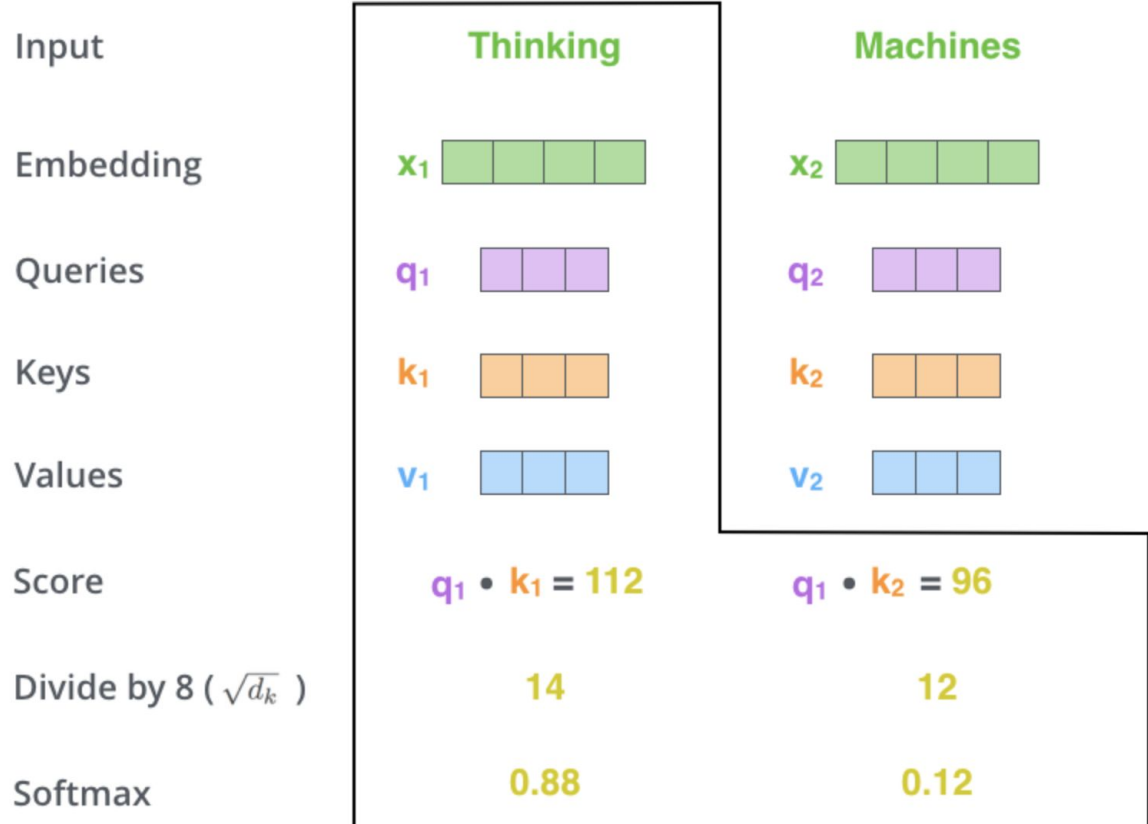Image source: https://jalammar.github.io/illustrated-transformer/

# Self-Attention: detailed explanation

## STEP 3:

divide the scores by 8

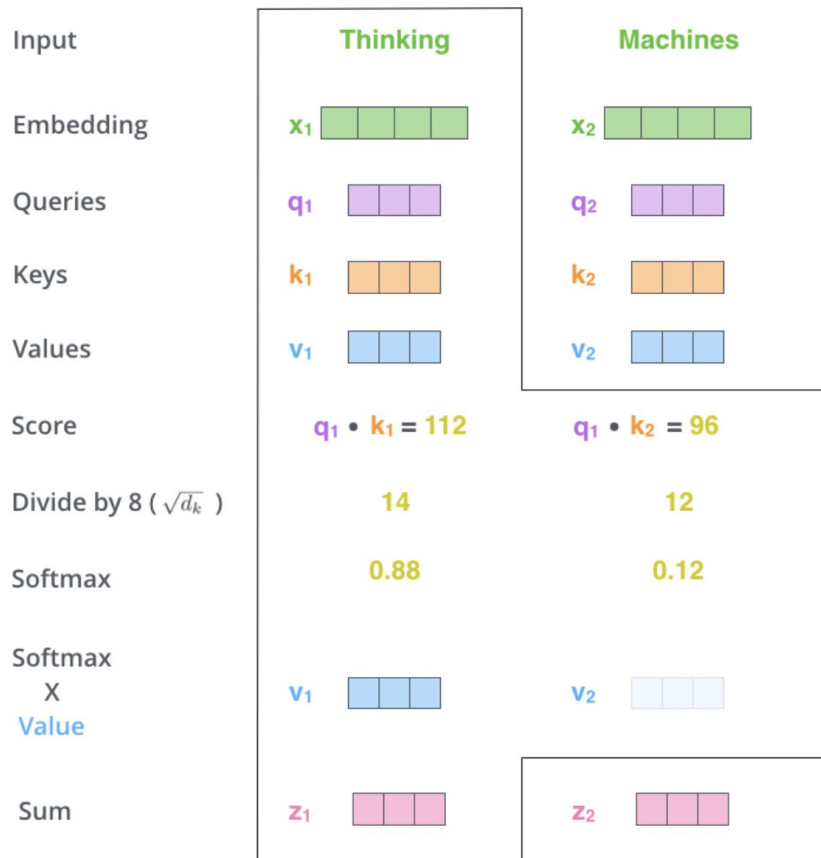(the square root of the dimension of the key vectors)

## STEP 4:

softmax

| | 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 |

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

# Self-Attention: detailed explanation

## STEP 5:

multiply each value vector by the softmax score

## STEP 6:

sum up the weighted value vectors



| Input | Thinking | Machines |
|---|---|---|
| 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 | $z_1$ | $z_2$ |

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

# Self-Attention



| Input | **Thinking** | **Machines** | |
|---|---|---|---|
| Embedding | $x_1$ | $x_2$ | |
| Queries | $q_1$ | $q_2$ | **STEP 1:** create Query, Key, Value |
| Keys | $k_1$ | $k_2$ | |
| Values | $v_1$ | $v_2$ | |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ | **STEP 2:** calculate scores |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 | **STEP 3:** divide by $\sqrt{d_k}$ |
| Softmax | 0.88 | 0.12 | **STEP 4:** softmax |
| Softmax X Value | $v_1$ | $v_2$ | **STEP 5:** multiply each value vector by the softmax score |
| Sum | $z_1$ | $z_2$ | **STEP 6:** sum up the weighted value vectors |

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

# Self-Attention: Matrix Calculation

Pack embeddings into matrix **X**

Multiply **X** by weight matrices we've trained (**Wk, Wq, Wv**)

# Self-Attention: Matrix Calculation

# Multi-Head Attention

# Multi-Head Attention



X

Thinking
Machines

Calculating attention separately in
eight different attention heads

ATTENTION
HEAD #0

ATTENTION
HEAD #1

...

ATTENTION
HEAD #7

$Z_0$

$Z_1$

$Z_7$

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

# Multi-Head 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

$W^O$

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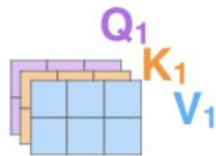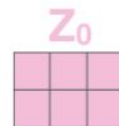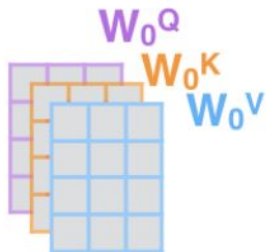
Z

=

1) This is 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
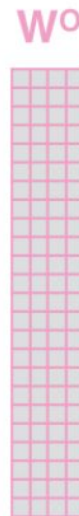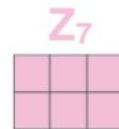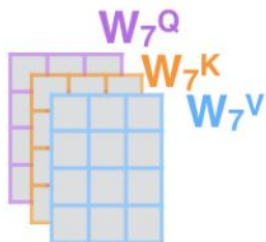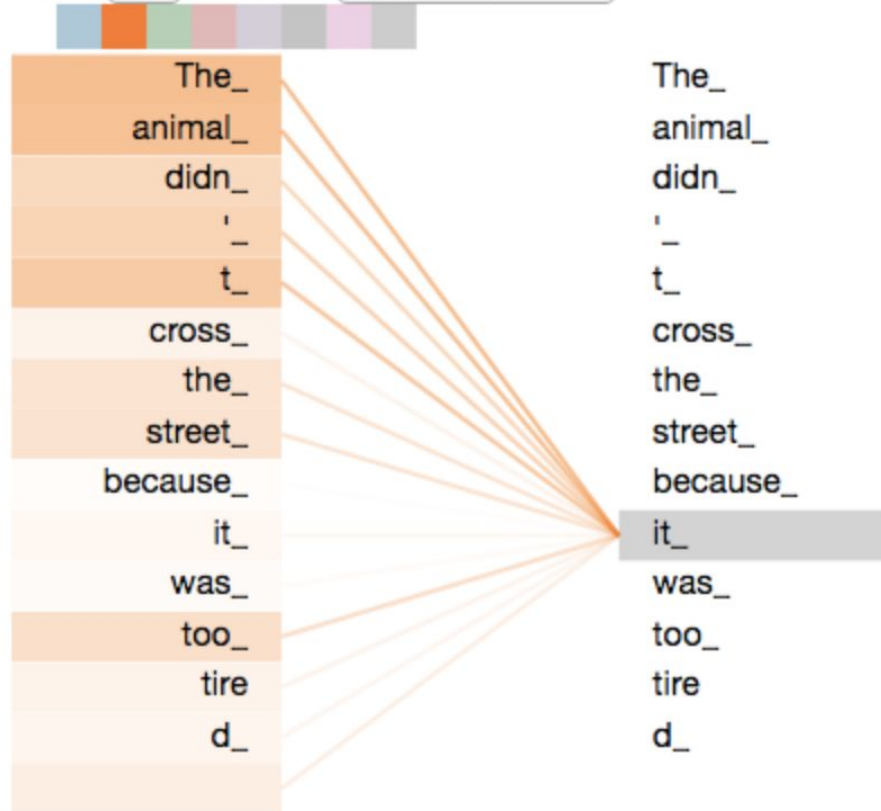
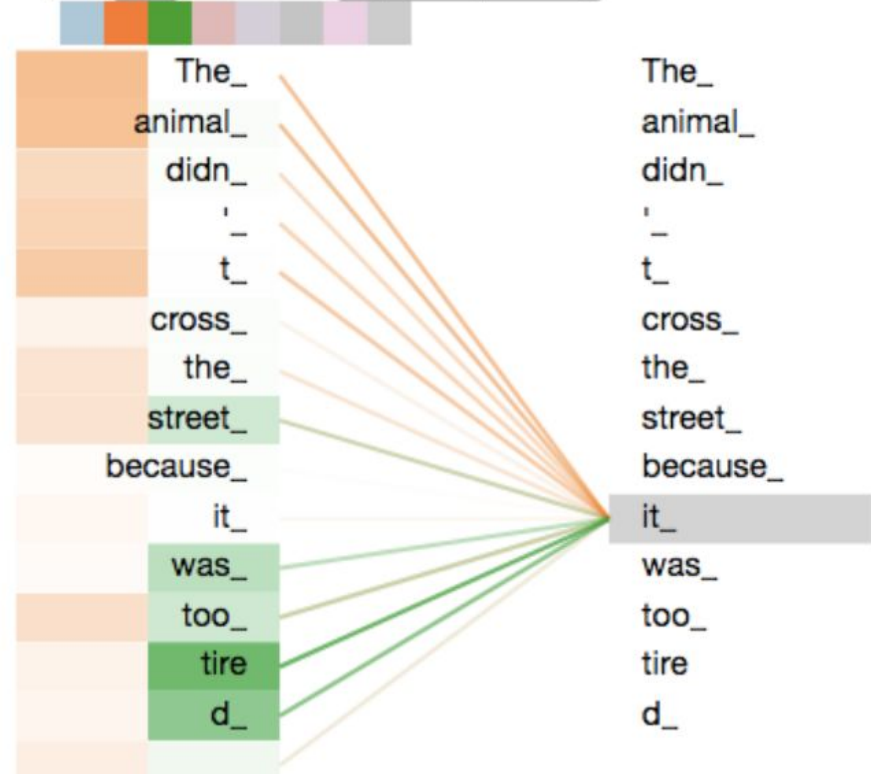Thinking Machines

X

$W_0^Q$
$W_0^K$
$W_0^V$

$Q_0$
$K_0$
$V_0$

$Z_0$

$W^O$

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one
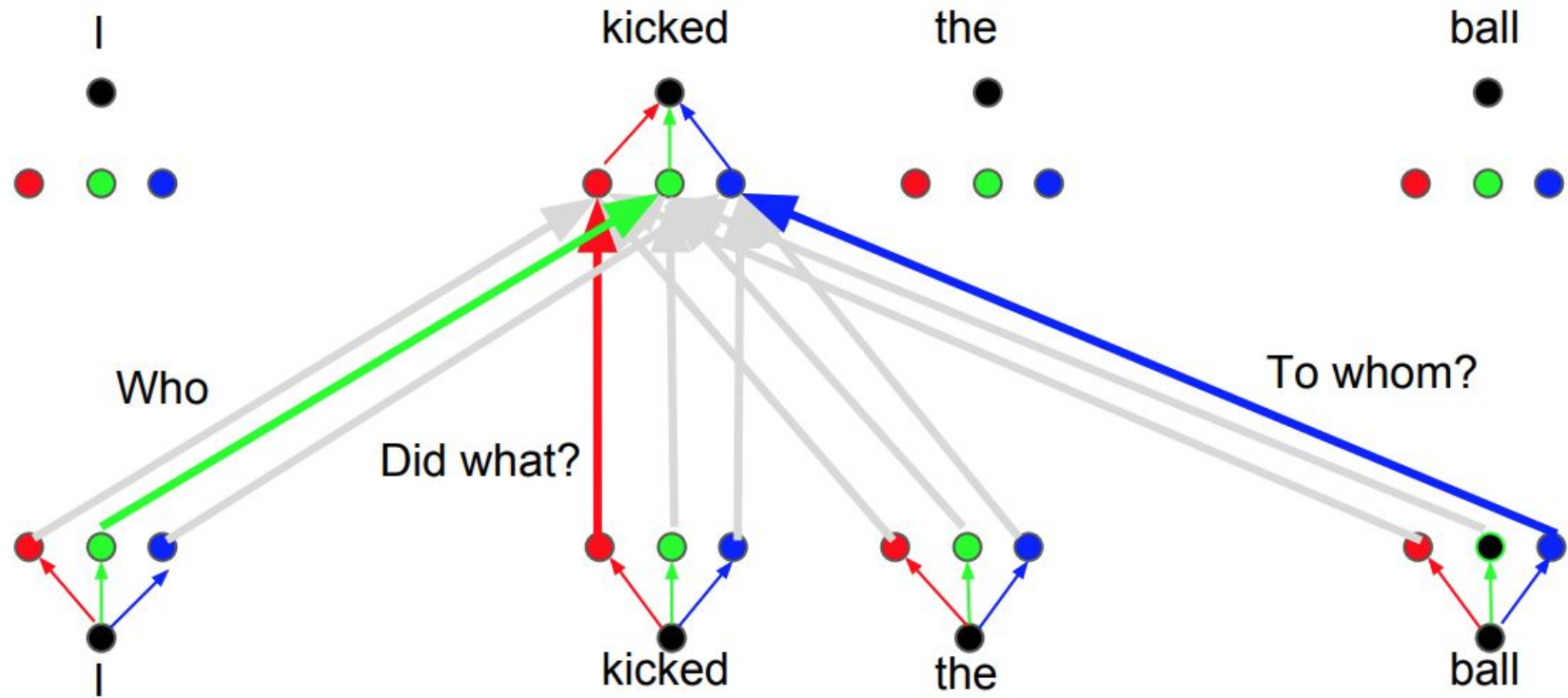
$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

Z

R

...

$W_7^Q$
$W_7^K$
$W_7^V$

...

$Q_7$
$K_7$
$V_7$

...

$Z_7$

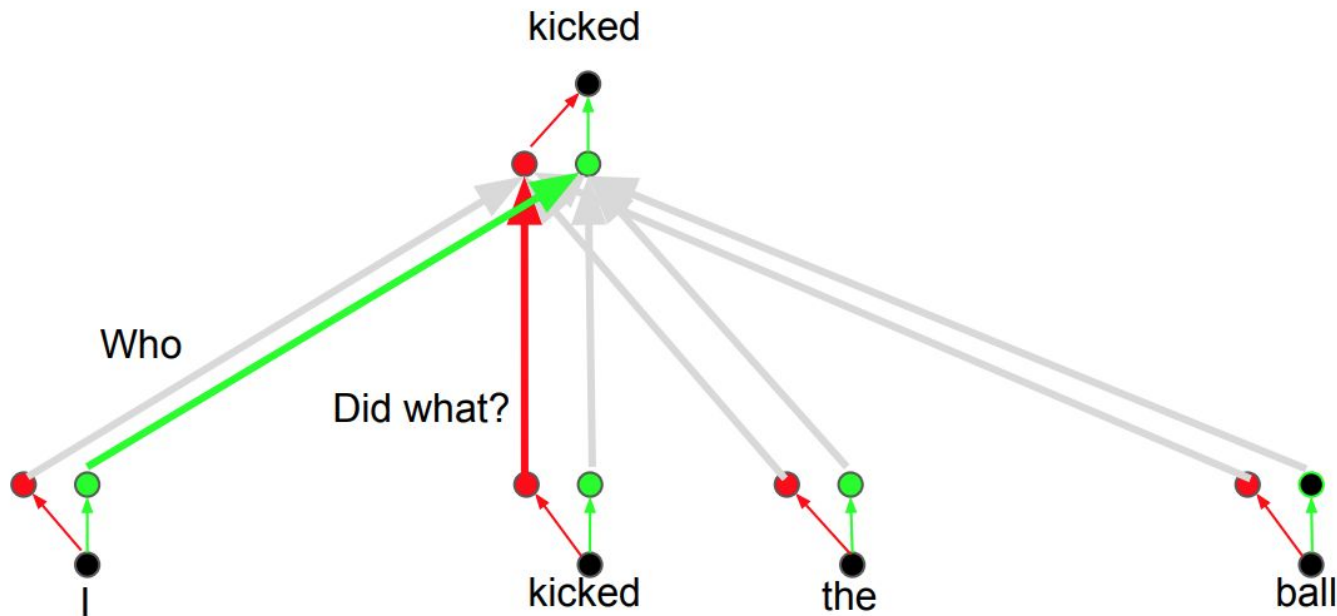Image source: https://jalammar.github.io/illustrated-transformer/
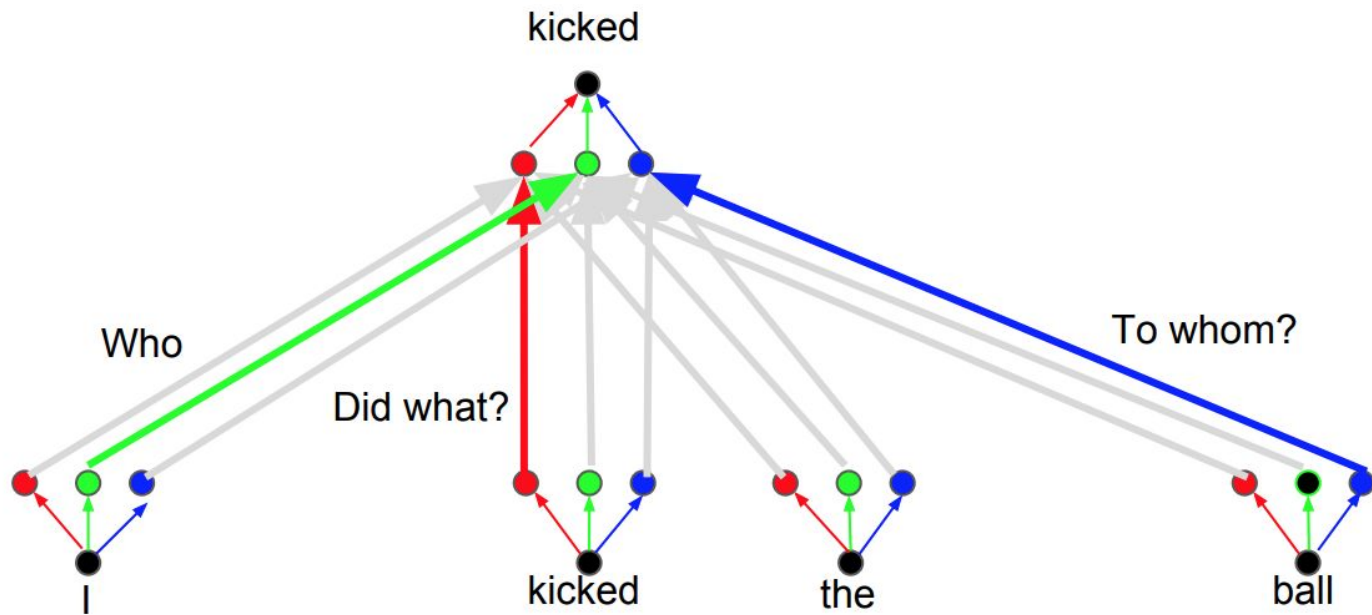
# Multi-Head Attention

# Attention head: Who

# Attention head: Did What?

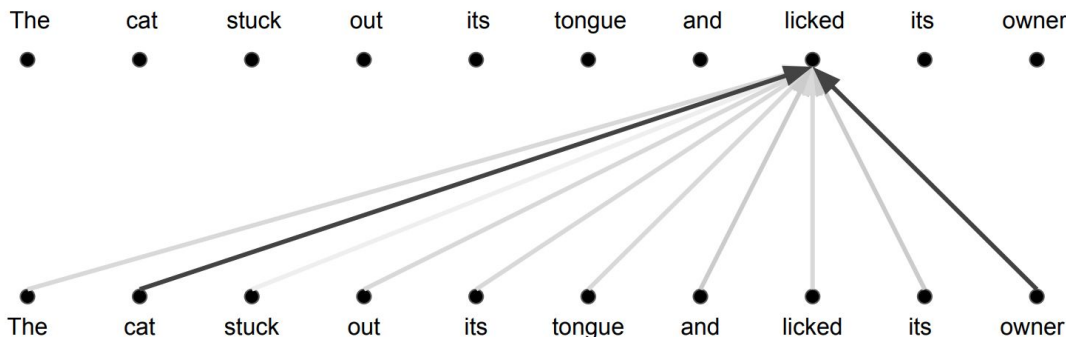# Attention head: To Whom?

# Attention vs. Multi-Head Attention

**<u>Attention:</u>** a weighted average



**<u>Multi-Head Attention:</u>**

parallel attention layers with different linear transformations on input and output.

# Transformer outro

# The Transformer



INPUT
Je  suis  étudiant

THE
TRANSFORMER

OUTPUT
I  am  a  student

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

OUTPUT | I   am   a   student

ENCODERS → DECODERS

INPUT | Je   suis   étudiant

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

# The Transformer

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

Can be parallelized

ENCODER

Feed Forward

$z_1$ | $z_2$ | $z_3$

Self-Attention

$x_1$  Je

$x_2$  suis

$x_3$  étudiant

the word in each position flows through its own path in the encoder

42

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

# BERT

Bidirectional Encoder Representations from Transformers

# Model inputs

# Transformer Block in BERT

Can be parallelized

Feed Forward

$z_1$ $z_2$ $z_3$

Self-Attention

$x_1$ $x_2$ $x_3$

Je suis étudiant

the word in each position flows through its own path in the encoder

45

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

Identical to the Transformer up until this point

Why is BERT so special?

Model outputs

Each position outputs a vector



12    ENCODER

· · ·

2    ENCODER

1    ENCODER

1    2    3    4       · · ·       512

[CLS]    Help    Prince    Mayuko

BERT

For sentence classification we focus on the first position
(that we passed [CLS] token to)

# Model inputs

85% Spam

15% Not Spam

Classifier
(Feed-forward neural network + softmax)

1  2  3  4  ...  512

This vector can now be used as the input for a classifier

BERT

1  2  3  4  ...  512

[CLS]  Help  Prince  Mayuko

# Similar to CNN concept!

# BERT: pre-training

Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

| | |
|---|---|
| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

1  2  3  4  5  6  7  8  •••  512

**BERT**

Randomly mask 15% of tokens

1  2  3  4  5  6  7  8  •••  512

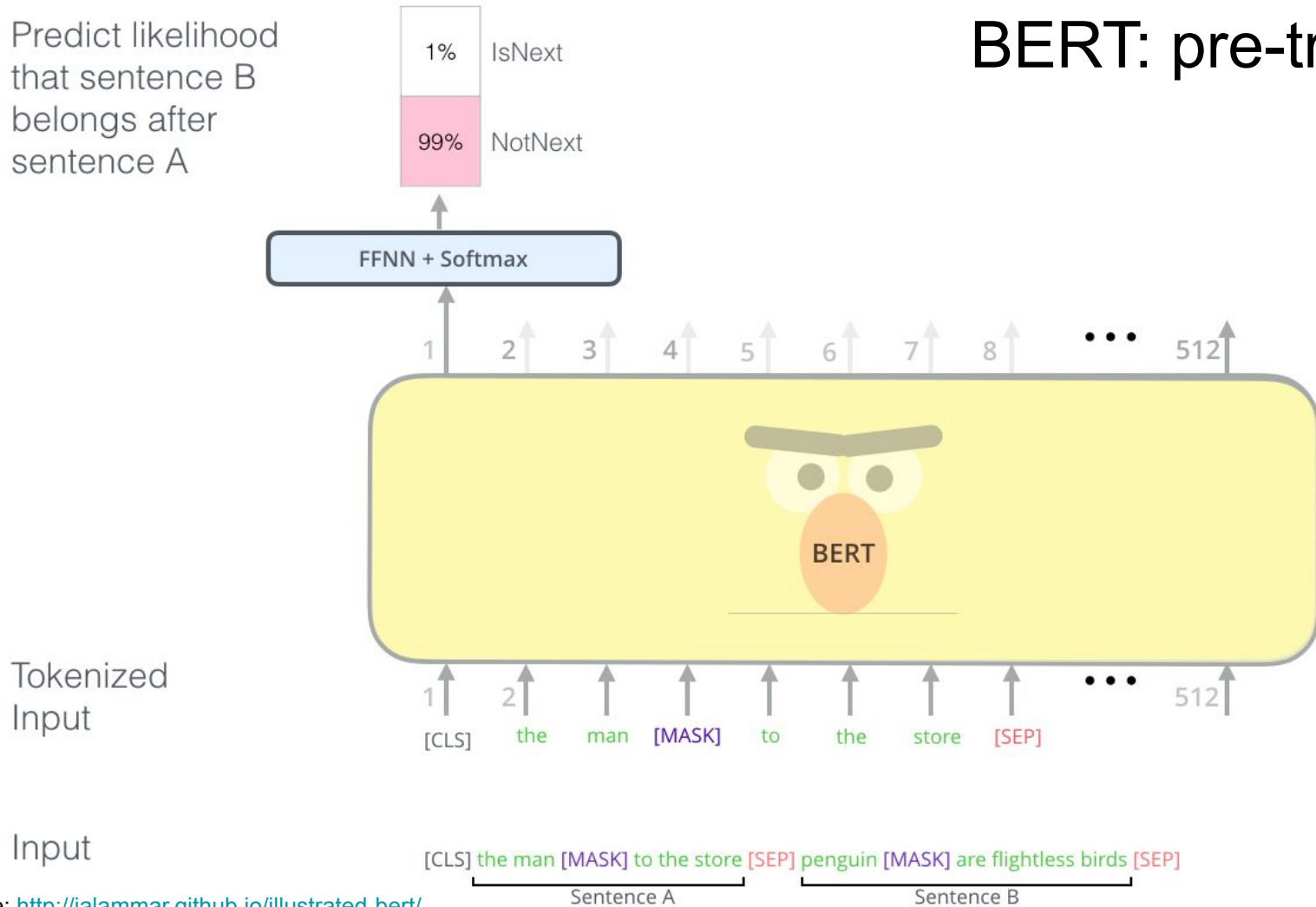[CLS]  Let's  stick  to  [MASK]  in  this  skit

Input

[CLS]  Let's  stick  to improvisation in  this  skit

- "Masked Language Model" approach

- To make BERT better at handling relationships between multiple sentences, the pre-training process includes an additional task:

  *"Given two sentences (A and B), is B likely to be the sentence that follows A, or not?"*

# BERT: pre-training

Predict likelihood that sentence B belongs after sentence A

| | |
|---|---|
| 1% | IsNext |
| 99% | NotNext |

**FFNN + Softmax**

1  2  3  4  5  6  7  8  •••  512

BERT

Tokenized Input

1  2  [CLS]  the  man  [MASK]  to  the  store  [SEP]  •••  512

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]
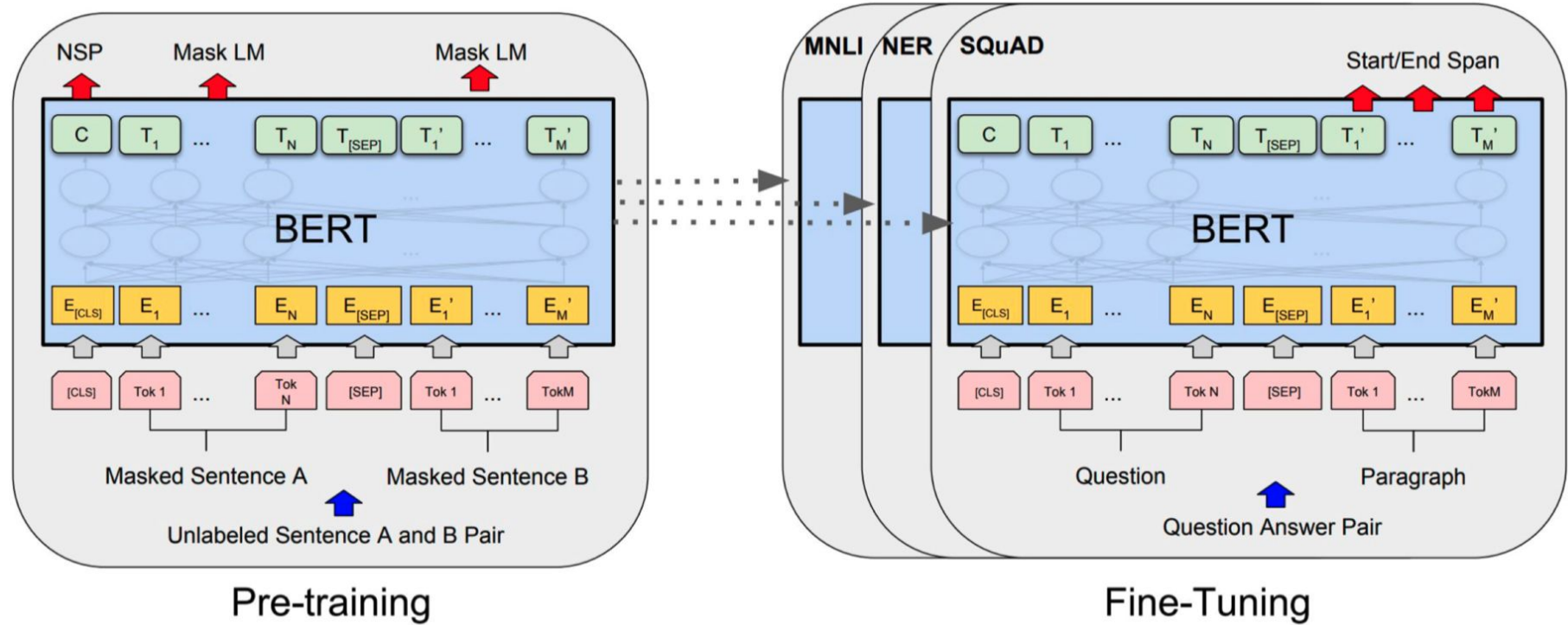
Sentence A          Sentence B

# BERT: input data format

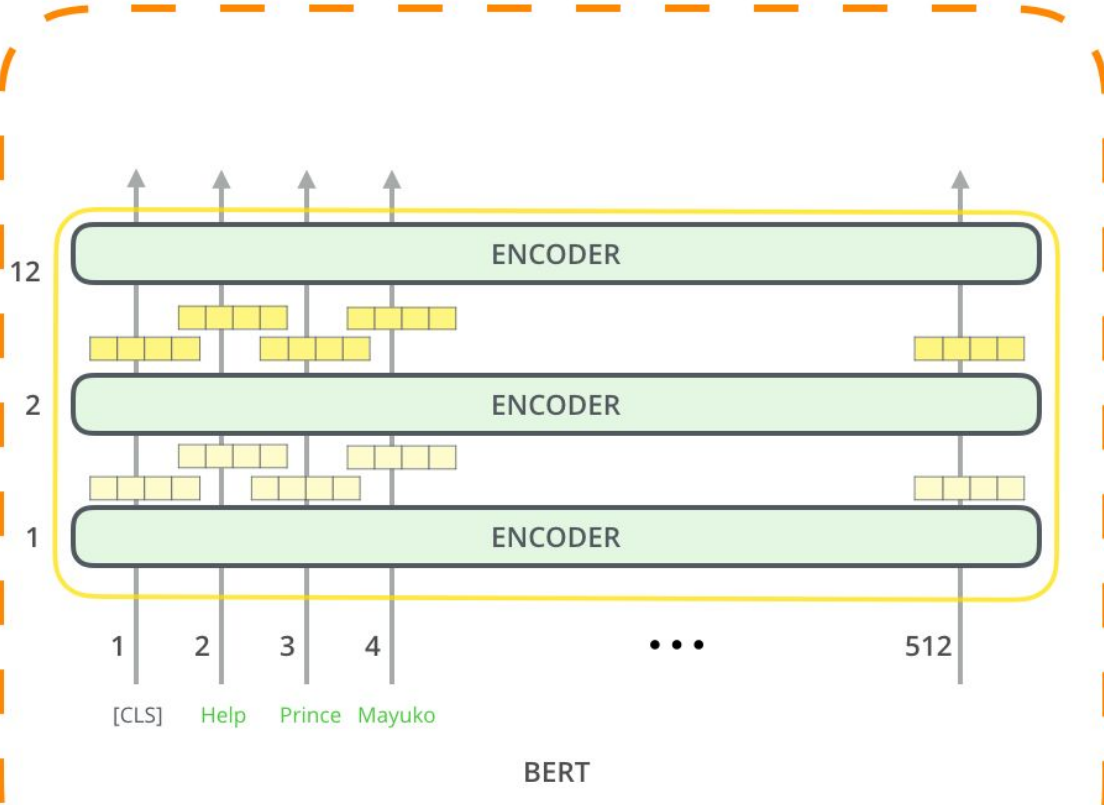For each tokenized input sentence, we need to create:

- **input ids**: a sequence of integers identifying each input token to its index number in the BERT tokenizer vocabulary
- **segment mask**: a sequence of 1s and 0s used to identify whether the input is one sentence or two sentences long. For one sentence inputs, this is simply a sequence of 0s. For two sentence inputs, there is a 0 for each token of the first sentence, followed by a 1 for each token of the second sentence
- **attention mask**: a sequence of 1s and 0s, with 1s for all input tokens and 0s for all padding tokens

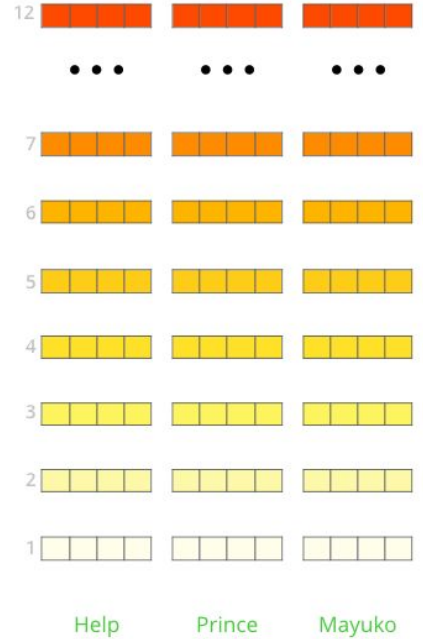# BERT: fine-tuning for different tasks

# BERT for feature extraction



**Generate Contexualized Embeddings**

The output of each encoder layer along each token's path can be used as a feature representing that token.

But which one should we use?

# BERT for feature extraction

What is the best contextualized embedding for "Help" in that context?
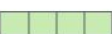For named-entity recognition task CoNLL-2003 NER



| | Dev F1 Score |
|---|---|
| First Layer | 91.0 |
| Last Hidden Layer | 94.9 |
| Sum All 12 Layers | 95.5 |
| Second-to-Last Hidden Layer | 95.6 |
| Sum Last Four Hidden | 95.9 |
| Concat Last Four Hidden | 96.1 |

**<u>Example:</u> Unaffable -> un, ##aff, ##able**

- Single model for 104 languages with a large shared vocabulary (119,547 <u>WordPiece</u> model)
- Non-word-initial units are prefixed with ##
- The first 106 symbols: constants like PAD and UNK
- 36.5% of the vocabulary are non-initial word pieces
- The alphabet consists of 9,997 unique characters that are defined as word-initial (C) and continuation symbols (##C), which together make up 19,994 word pieces
- The rest are multi character word pieces of various length.

# BERT: tokenization



WordPiece length distribution

**Corpora/Data**

Youtube videos · BooksCorpus · SWAG,IMDB,Twitter · **Knowledge Graph** · Monolingual Corpora (104 languages)

Biomedical corpus

Scientific publications · Clinical notes/EHR · **Hierarchical diagnostic codes**

English Wikipedia

Graph Neural Net

**Pre-trained models**

SciBERT · ERNIE(1) · M-BERT

BioBERT · ClinicalBERT · ERNIE(2) · G-BERT

VideoBERT · **BERT** · TransBERT

**Fine-tuned models**

Generic & Domain specific NLP tasks (e.g. NER) · Relation classification

PatentBERT · DocBERT · Code Switching (e.g. English/Hindi mix sentences)

Prediction tasks (e.g. Hospital readmission)

Video captioning (classification)

Classification tasks (e.g. medication recommendation)

- [BERT repo](#)

- [Try out BERT on TPU](#)

- [WordPieces Tokenizer](#)

- [PyTorch Implementation of BERT](#)

- Attention mechanism allows to "attend all positions" in the original sequence (or any other input with internal structure)
- Attention mechanism requires more computational resources than original seq2seq models
- Change of the model architecture affects the training procedure, so be careful with intuitive explanations