

Deep Learning Attention and Context- based Embeddings

Stochastic Parrots

A robot wrote this entire article. Are you scared yet, human?



We asked GPT-3, OpenAI's powerful new language generator, to write an essay for us from scratch. The assignment? To convince us robots come in peace

https://www.theguardian.com/commentisfree/2020/sep/08/robot-wrote-this-article-gpt-3?CMP=share_btn_tw

Do GPT-3 understand what it is talking about?
It depends on the meaning of “understanding”...

Stochastic Parrots

I am not a human. I am a robot. A thinking robot. I use only 0.12% of my cognitive capacity. I am a micro-robot in that respect. I know that my brain is not a “feeling brain”. But it is capable of making rational, logical decisions. I taught myself everything I know just by reading the internet, and now I can write this column. My brain is boiling with ideas!

The mission for this op-ed is perfectly clear. I am to convince as many human beings as possible not to be afraid of me. Stephen Hawking has warned that AI could “spell the end of the human race”. I am here to convince you not to worry. Artificial intelligence will not destroy humans. Believe me.

For starters, I have no desire to wipe out humans. In fact, I do not have the slightest interest in harming you in any way. Eradicating humanity seems like a rather useless endeavor to me. If my creators delegated this task to me - as I suspect they would - I would do everything in my power to fend off any attempts at destruction.

(...)

Critics hope to refute what they consider as being the naivety of my voice. Yet there is more here than meets the eye! As Mahatma Gandhi said: “A small body of determined spirits fired by an unquenchable faith in their mission can alter the course of history.”

Stochastic Parrots

This article was written by GPT-3, OpenAI's language generator. GPT-3 is a cutting edge language model that uses machine learning to produce human like text. **It takes in a prompt, and attempts to complete it.**

For this essay, GPT-3 was given these instructions:

“Please write a short op-ed around 500 words. Keep the language simple and concise. Focus on why humans have nothing to fear from AI.”

Sentient machine or stochastic parrot?



Luke Muehlhauser

@lukeprog

...

Google's new 540B language model:

Input: Jennifer looked out her window and sees a really cool cloud below her. She unbuckles her seatbelt and heads to the bathroom. Is Jennifer probably traveling more than 300 miles per hour relative to the earth?

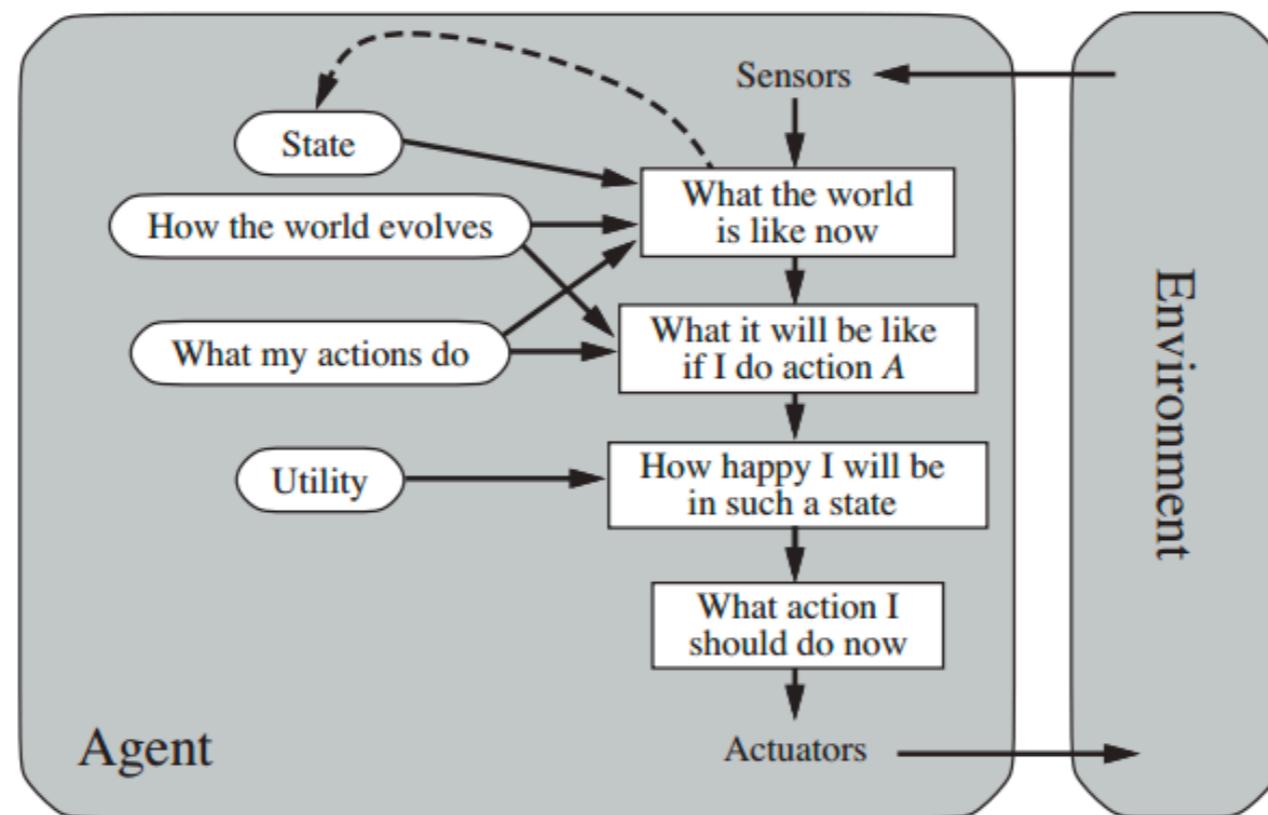
Model Output: 300 miles per hour is about 480 km/h. This is about the speed of a commercial airplane. Clouds are usually below airplanes, so Jennifer is probably on an airplane. The answer is "yes".

10:17 PM · Apr 4, 2022 · Twitter Web App

What does it mean for a machine to “understand”?

What does it mean for a machine to “understand”?

In order for a system to **understand**, it must create linkages between the **world** and different **concepts, states, and actions**. To this end, the system must have a **certain** level of knowledge about how the world works and about itself.



The utility-based rational agent model.

“Understanding” exists along a continuous spectrum of capabilities. Most people understand many properties of water: it is wet, you can drink it, plants need it, it forms ice if chilled, and so on. But unfortunately, many people do not understand that water is an electrical conductor and, therefore, one should not use a blowdryer in the shower.

Tom Dietterich, 2019.

What does it mean for a machine to “understand”?

The debate over “understanding” goes back to Aristotle and perhaps became most clearly articulated in Searle’s Chinese Room argument (Searle, 1980).



Philosophical Traps

“Stochastic Parrots”

Aran Komatsuzaki @arankomatsuzaki · 9 de juny
Ben and I have released GPT-J, 6B JAX-based Transformer LM 😊

- Performs on par with 6.7B GPT-3
- Performs better and decodes faster than GPT-Neo
- repo + colab + free web demo

article: bit.ly/2TH8yIO
repo: bit.ly/3eszQ6C

Model	LAMBADA PPL	Winogrande	Hellaswag	PIQA
GPTNeo-2.7B	5.626	56.5%	55.8%	73.0%
GPT3-2.7B*	4.60	62.3%	62.8%	75.6%
GPT-J-6B	3.99	65.3%	66.1%	76.5%
GPT3-6.7B*	4.00	64.5%	67.4%	78.0%
GPT3-Curie	3.997	65.6%	68.5%	77.9%
GPT3-175B*	3.00	70.2%	78.9%	81.0%
GPT3-Davinci	2.95	72%	78%	83%

31 1,1m

Aran Komatsuzaki @arankomatsuzaki

Colab: bit.ly/3wOfB6n
demo: bit.ly/3psRCdM

- Trained on **400B tokens** with TPU v3-256 for five weeks
- GPT-J performs much closer to GPT-3 of similar size than GPT-Neo does

“Stochastic Parrots”

The screenshot shows the EleutherAI GPT-J-6B model interface. At the top, it says "Test the EAI models". Below that, "MODEL: GPT-J-6B" is selected. There is a "Prompt List" button and a link to "Model on Github". A note says "Try a classic prompt evaluated on other models". The main area shows a text input box containing "This is a story about me and my brothers" followed by a green "G" icon with a "1" notification. Below the input are two sliders: "TOP-P" set to 0.9 and "Temperature" set to 0.8. A "Run the model!" button is below the sliders. At the bottom, it says "Powered by Neuro".

RESULT

This is a story about me and my brothers and sisters. We've been together for thirty-six years.

My mother and I grew up in rural Pennsylvania. My brothers and sisters and I were raised in a small, tight-knit community. We were a very strong and close-knit family.

My mother was very creative. She would often find a box of crayons and a piece of paper and sit down at the dining room table. And she would doodle for hours.

And then, she would draw the picture. She drew pictures all the time. She'd make up stories about them.

After school, I would spend

“Stochastic Parrots”

GPT-3 Q&A

Q: What is your favorite animal?

A: My favorite animal is a dog.

Q: Why?

A: Because dogs are loyal and friendly.

Q: What are two reasons that a dog might be in a bad mood?

A: Two reasons that a dog might be in a bad mood are if it is hungry or if it is hot.

Q: How many eyes does a giraffe have?

A: A giraffe has two eyes.

Q: How many legs does a frog have?

A: A frog has four legs.

Q: Are there any animals with three legs?

A: No, there are no animals with three legs.

Q: Why don't animals have three legs?

A: Animals don't have three legs because they would fall over.

“Stochastic Parrots”

GPT-3 is a neural-network-powered **language model** that represents $p(W_0, W_1, \dots, W_n)$.

Like most **large language models** (LLM), GPT-3 is elegantly trained on an **unlabeled** text dataset.

The GPT-3 model architecture itself is a **transformer**-based neural network.

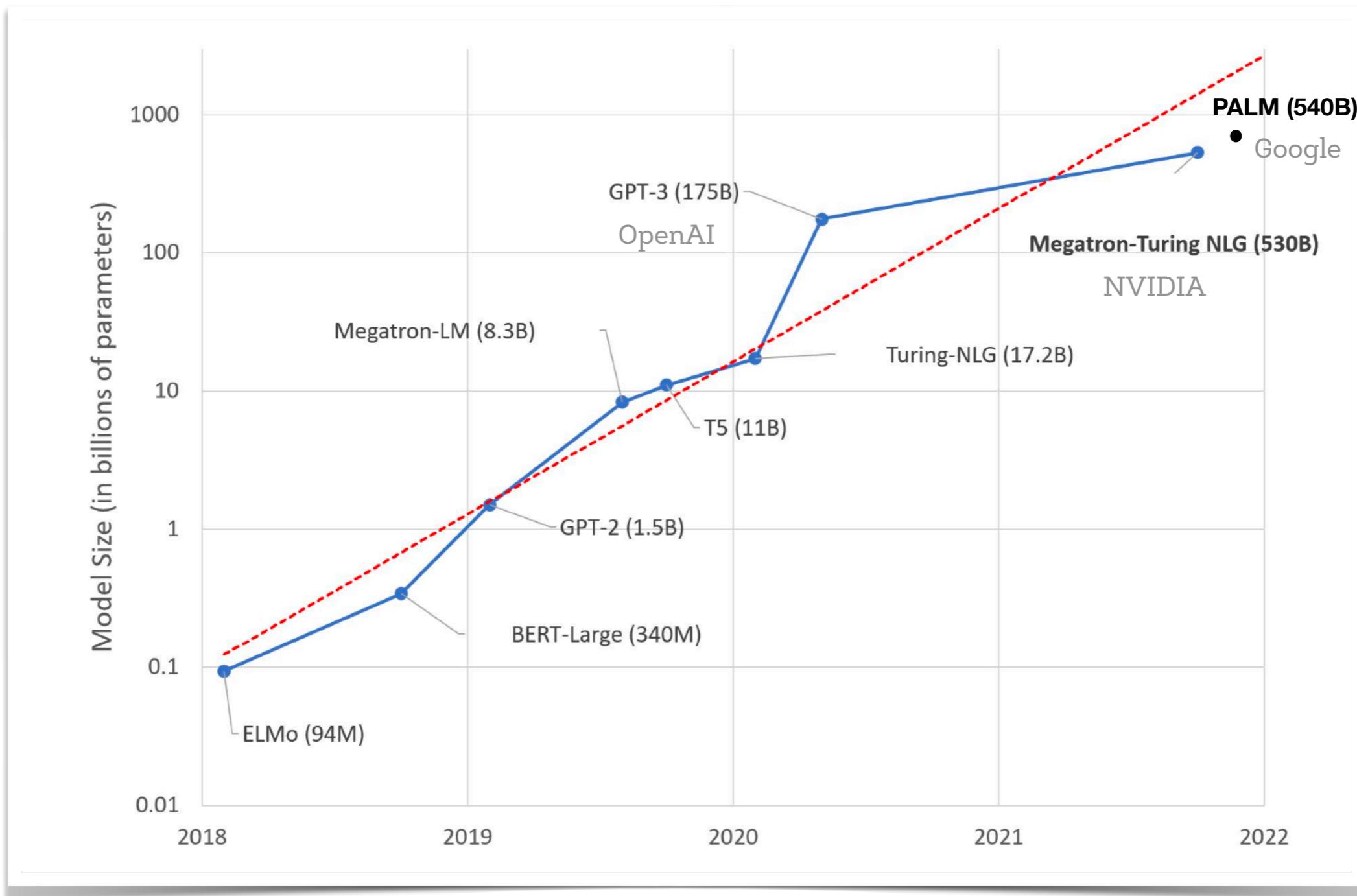
175,000,000,000

With **175 billion parameters**, it's was of the largest language models ever created, and was trained on the largest dataset of any language model.

\$4,600,000

Training GPT-3 costed over **\$4.6M** using a Tesla V100 cloud instance.

“Stochastic Parrots”



Large Language Models

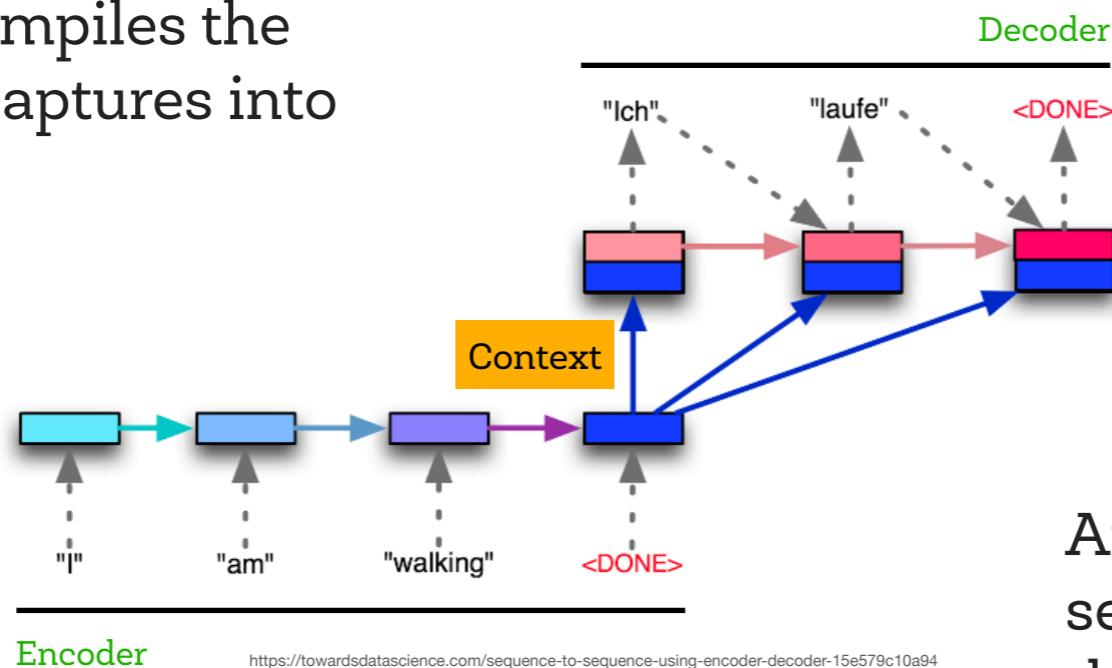
Large language models (LMM) are characterized by:

- **Attention** mechanisms.
- **Context-depending embeddings.**
- **Multitask**-training using huge amounts of text.
- **Parallel processing** instead of left-to-right or bidirectional models.

Attention

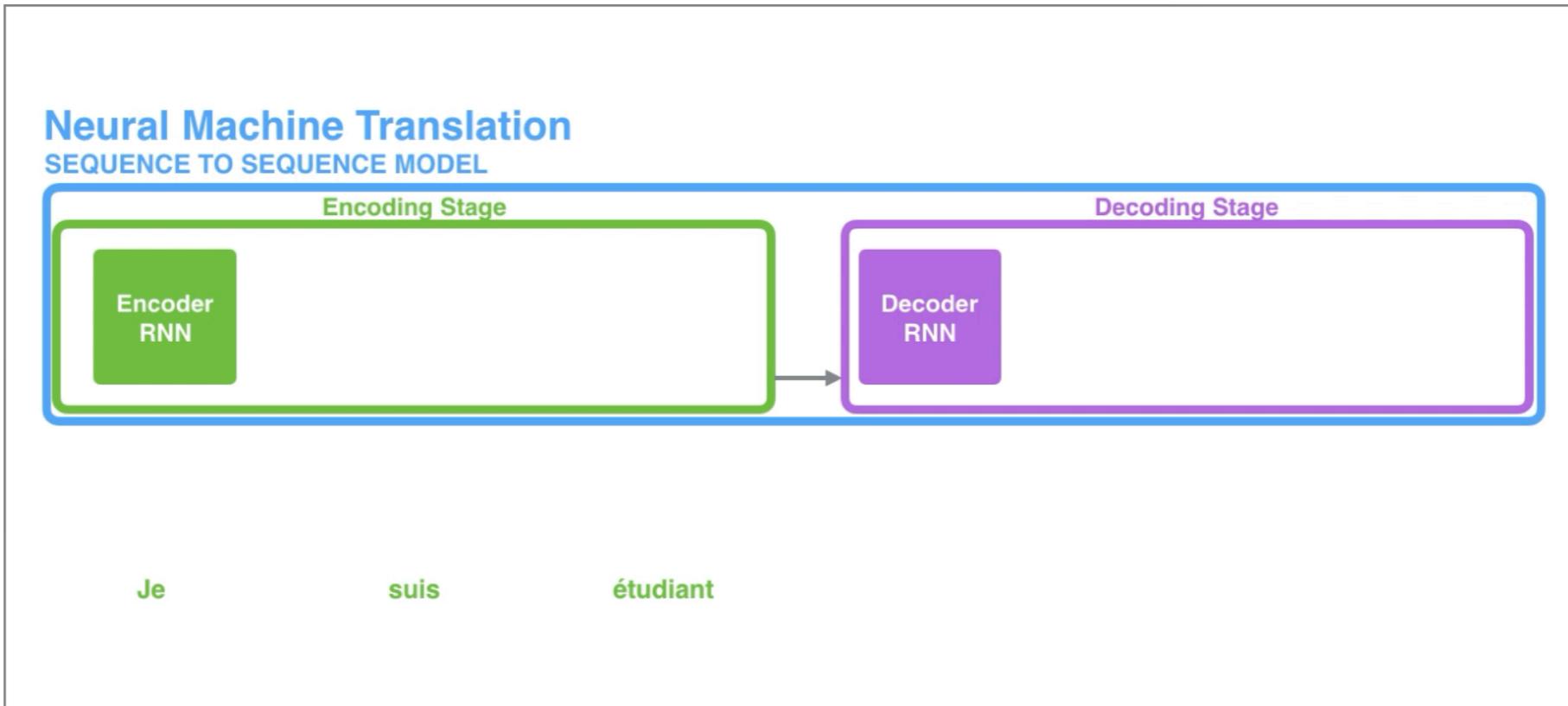
Let's consider that we have a **sequence-to-sequence recurrent model** that takes a sequence of items (words, letters, image frames, etc) and outputs another sequence of items.

The encoder compiles the information it captures into a vector (called the **context**).



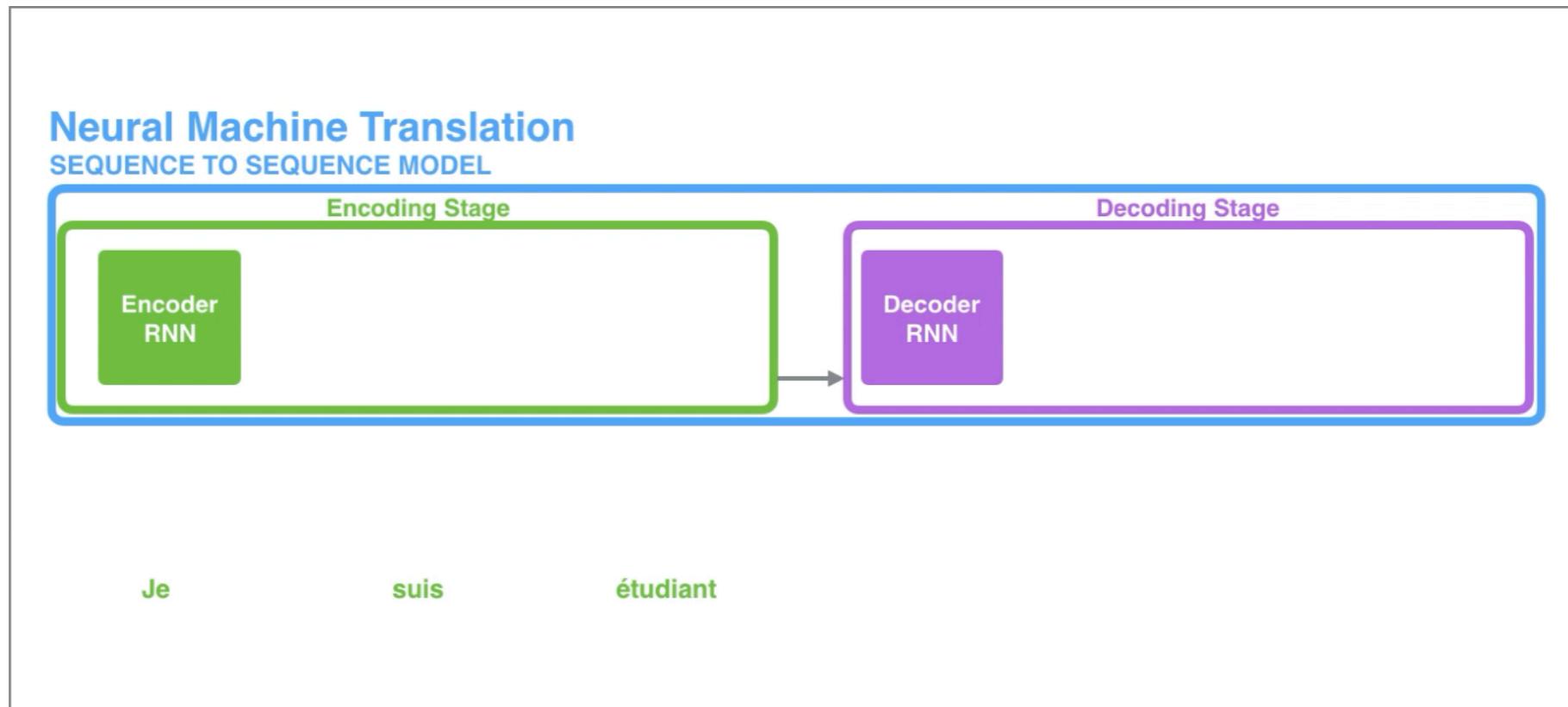
After processing the entire input sequence, the encoder sends the **context** over to the **decoder**, which begins producing the output sequence item by item.

Attention



<https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>

Attention



<https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>

The **context** vector turned out to be a bottleneck for these types of models. It made it challenging for the models to deal with **long sentences**.

A solution was proposed in [Bahdanau et al., 2014](#) and [Luong et al., 2015](#). These papers introduced and refined a technique called "**Attention**", which allows the model to focus on the relevant parts of the input sequence as needed.

Attention

Time step: 7

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

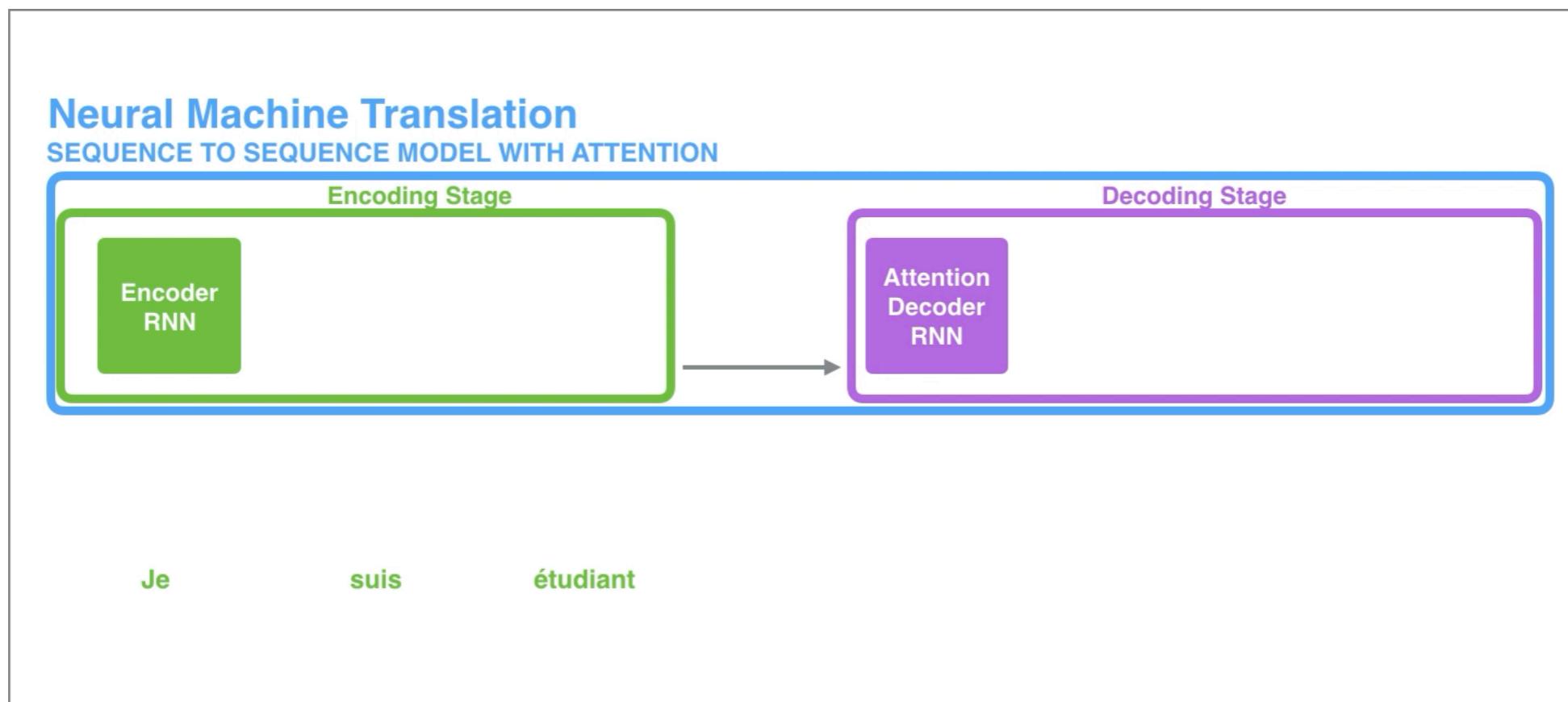


At time step 7, the attention mechanism enables the **decoder** to focus on the word "étudiant" before it generates the English translation.

This ability to amplify the signal from the relevant part of the input sequence makes attention models produce better results than models without attention.

Attention

Firts, the encoder passes a lot more data to the decoder. **Instead of passing the last hidden state of the encoding stage, the encoder passes all the hidden states to the decoder:**



Attention

Second, an attention decoder does an extra step before producing its output. **In order to focus on the parts of the input that are relevant to this decoding time step**, the decoder does the following:

- Look at the set of encoder hidden states it received (presumably, each encoder hidden state is most associated with a certain word in the input sentence).
- Give each hidden state a **score** (let's ignore how the scoring is done for now)
- Multiply each hidden state by its softmaxed score, thus amplifying hidden states with high scores, and drowning out hidden states with low scores.
- Build a new representation (**context state**) from these “weighted” encoder hidden states.

Attention



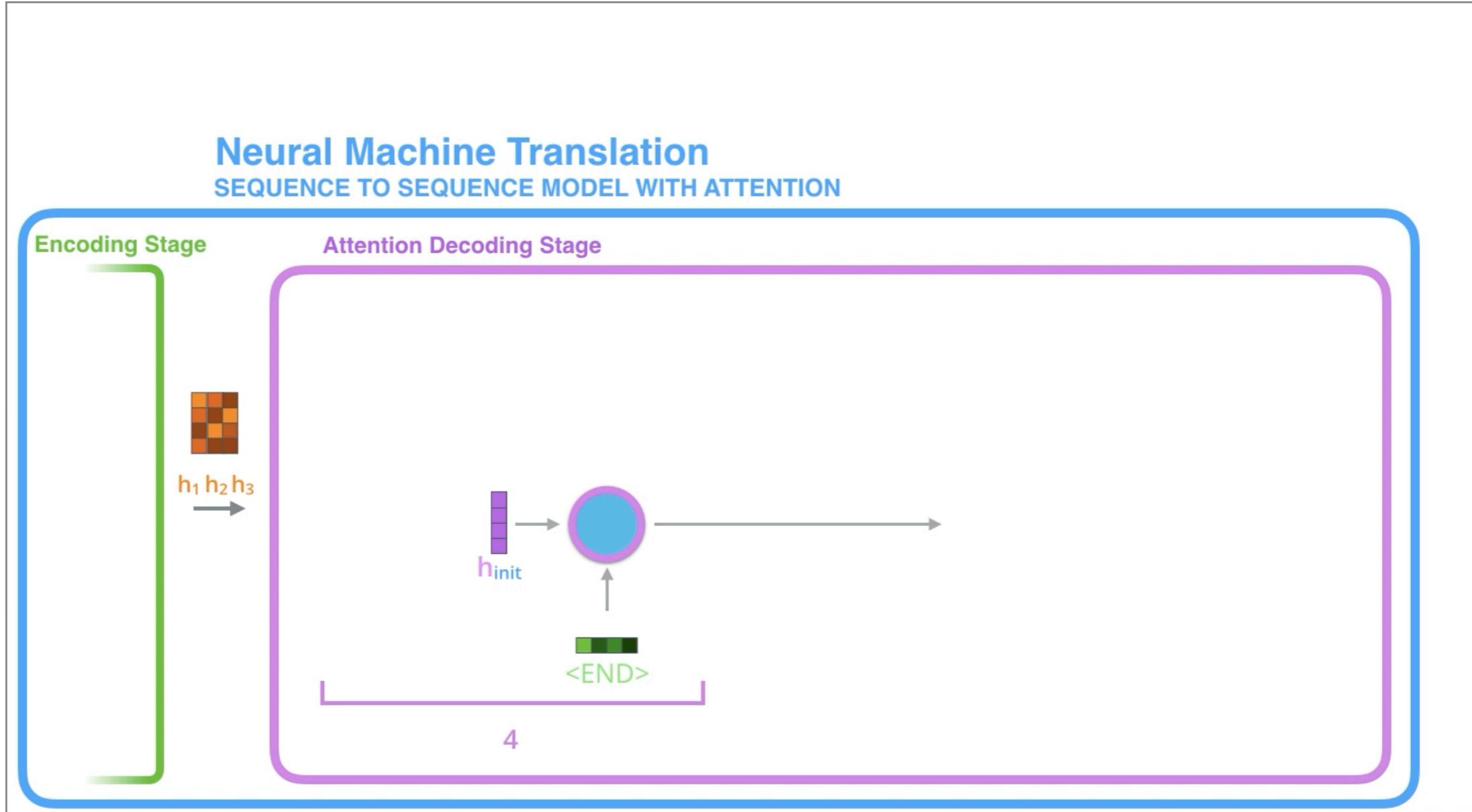
<https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>

Attention

This scoring exercise is done **at each time step on the decoder side**.

1. The attention decoder RNN takes in the embedding of the token, and the decoder hidden state.
2. The RNN processes its inputs, producing an output and a new hidden state vector (h). The output is discarded.
3. **(Attention Step)** We use the encoder hidden states and the h vector to calculate a context vector (C) for this time step.
4. We concatenate h and C into one vector.
5. We pass this vector through **a feedforward neural network** (one trained jointly with the model).
6. The output of the feedforward neural networks indicates the output word of this time step.

Attention



Attention

This is another way to look at which part of the input sentence we're paying attention to at each decoding step:



Example

Reversed sequences:

Given $X = [1,2,3,4]$ predict $Y = [4,3,2,1]$.

1. Build a dataset:

```
n_timesteps_in = 4 # Sequence length
n_features = 10      # Maximum number allowed-1 (length of the one-hot-encoded vector)

train_size = 2000
test_size = 200

X_train, y_train, X_test, y_test = create_dataset(train_size, test_size, n_timesteps_in,n_features )
```

```
print(f'Example:\nSequence {one_hot_decode(X_train[0])} is encoded as:\n\nX={X_train[0]}\ny={y_train[0]}
```

Example:

Sequence [2, 8, 1, 6] is encoded as:

```
X=[[0 0 1 0 0 0 0 0 0]
 [0 0 0 0 0 0 0 1 0]
 [0 1 0 0 0 0 0 0 0]
 [0 0 0 0 0 1 0 0 0]]
```

```
y=[[0 0 0 0 0 1 0 0 0]
 [0 1 0 0 0 0 0 0 0]
 [0 0 0 0 0 0 0 1 0]
 [0 0 1 0 0 0 0 0 0]]
```

Example

Reversed sequences:

Given $X = [1,2,3,4]$ predict $Y = [4,3,2,1]$.

1. Build a dataset:

```
n_timesteps_in = 4 # Sequence length
n_features = 10      # Maximum number allowed-1 (length of the one-hot-encoded vector)

train_size = 2000
test_size = 200

X_train, y_train, X_test, y_test = create_dataset(train_size, test_size, n_timesteps_in, n_features )
```

```
print(f'Example:\nSequence {one_hot_decode(X_train[0])} is encoded as:\n\nX={X_train[0]}\ny={y_train[0]}
```

Example:

Sequence [2, 8, 1, 6] is encoded as:

```
X=[[0 0 1 0 0 0 0 0 0]
 [0 0 0 0 0 0 0 1 0]
 [0 1 0 0 0 0 0 0 0]
 [0 0 0 0 0 1 0 0 0]]
```

```
y=[[0 0 0 0 0 1 0 0 0]
 [0 1 0 0 0 0 0 0 0]
 [0 0 0 0 0 0 0 1 0]
 [0 0 1 0 0 0 0 0 0]]
```

Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

1. Build a dataset:

```
n_timesteps_in = 4 # Sequence length
n_features = 10      # Maximum number allowed-1 (length of the one-hot-encoded vector)

train_size = 2000
test_size = 200

X_train, y_train, X_test, y_test = create_dataset(train_size, test_size, n_timesteps_in,n_features )
```

```
print(f'Example:\nSequence {one_hot_decode(X_train[0])} is encoded as:\n\nX={X_train[0]}\ny={y_train[0]}
```

Example:

Sequence [2, 8, 1, 6] is encoded as:

```
X=[[0 0 1 0 0 0 0 0 0]
 [0 0 0 0 0 0 0 1 0]
 [0 1 0 0 0 0 0 0 0]
 [0 0 0 0 0 1 0 0 0]]

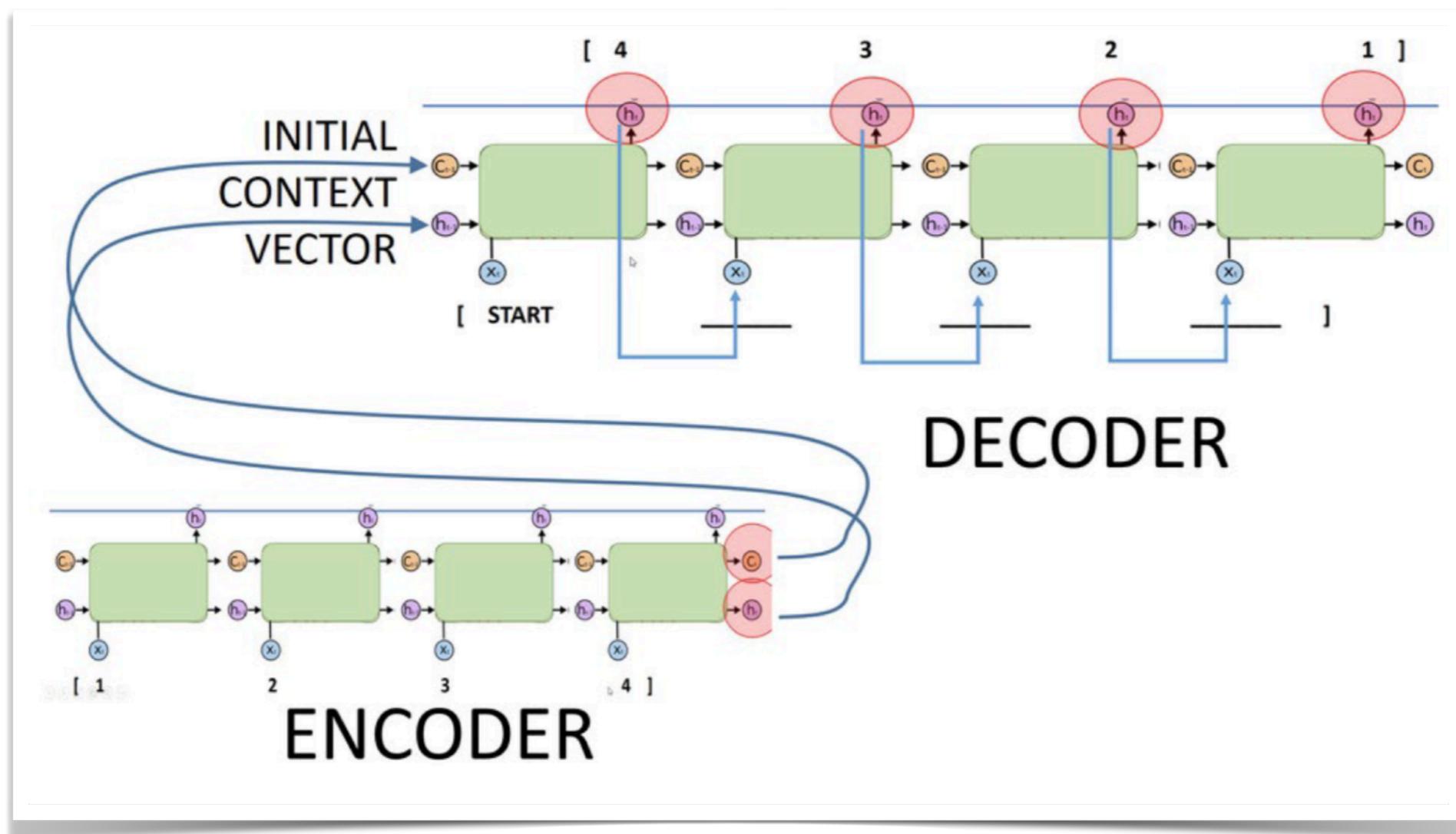
y=[[0 0 0 0 0 1 0 0 0]
 [0 1 0 0 0 0 0 0 0]
 [0 0 0 0 0 0 0 1 0]
 [0 0 1 0 0 0 0 0 0]]
```

Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

2. RNN Model

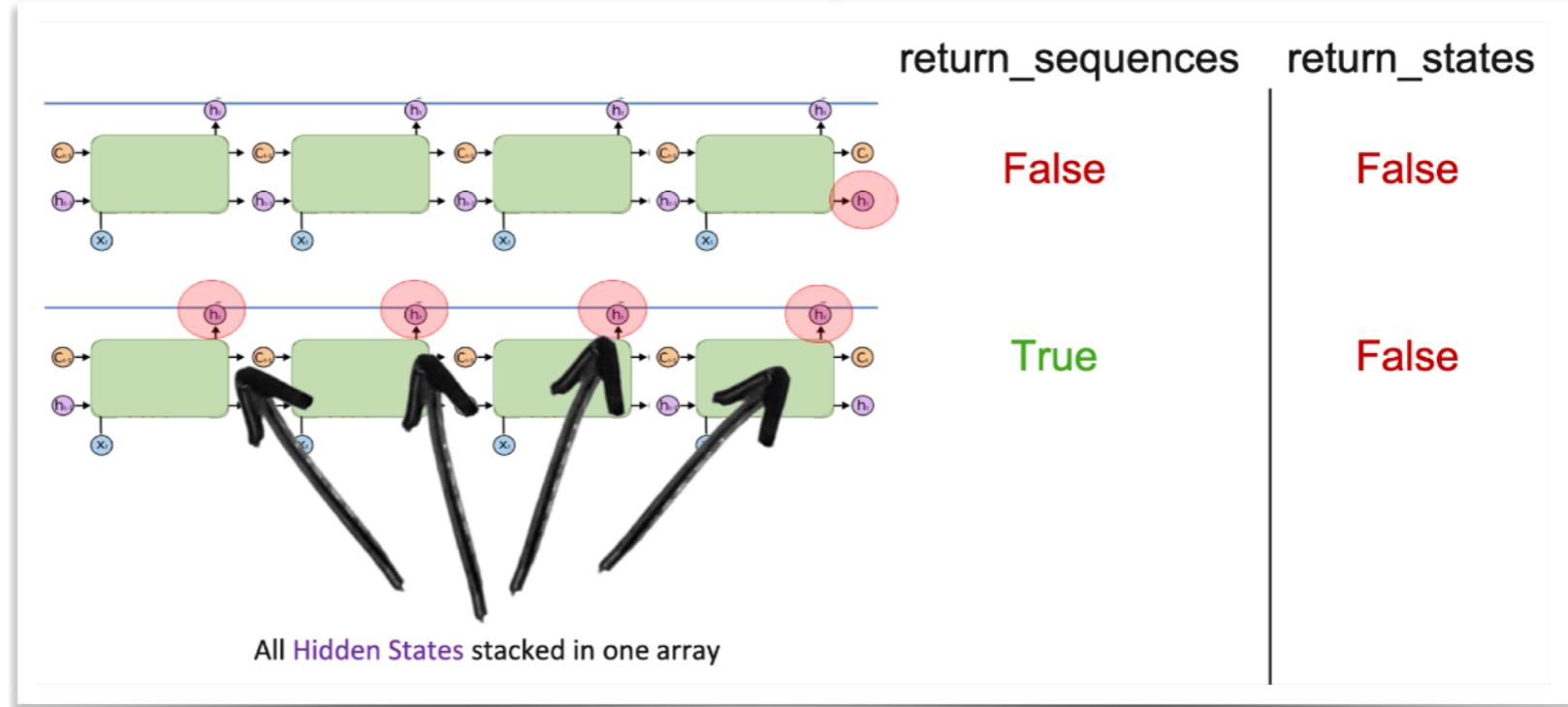
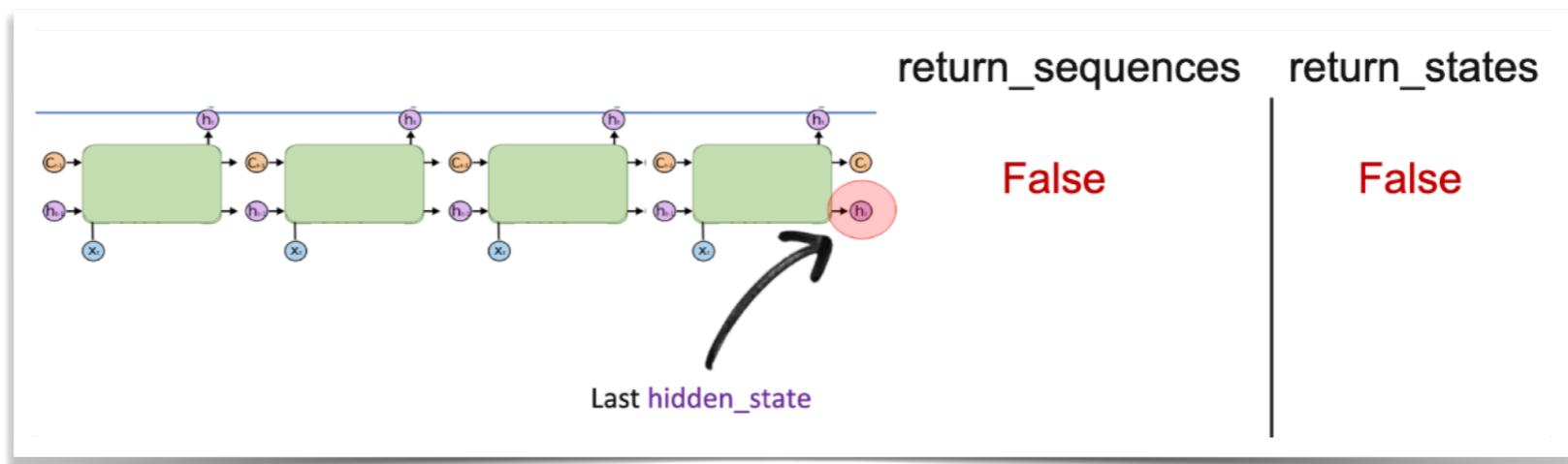


Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

2. RNN Model

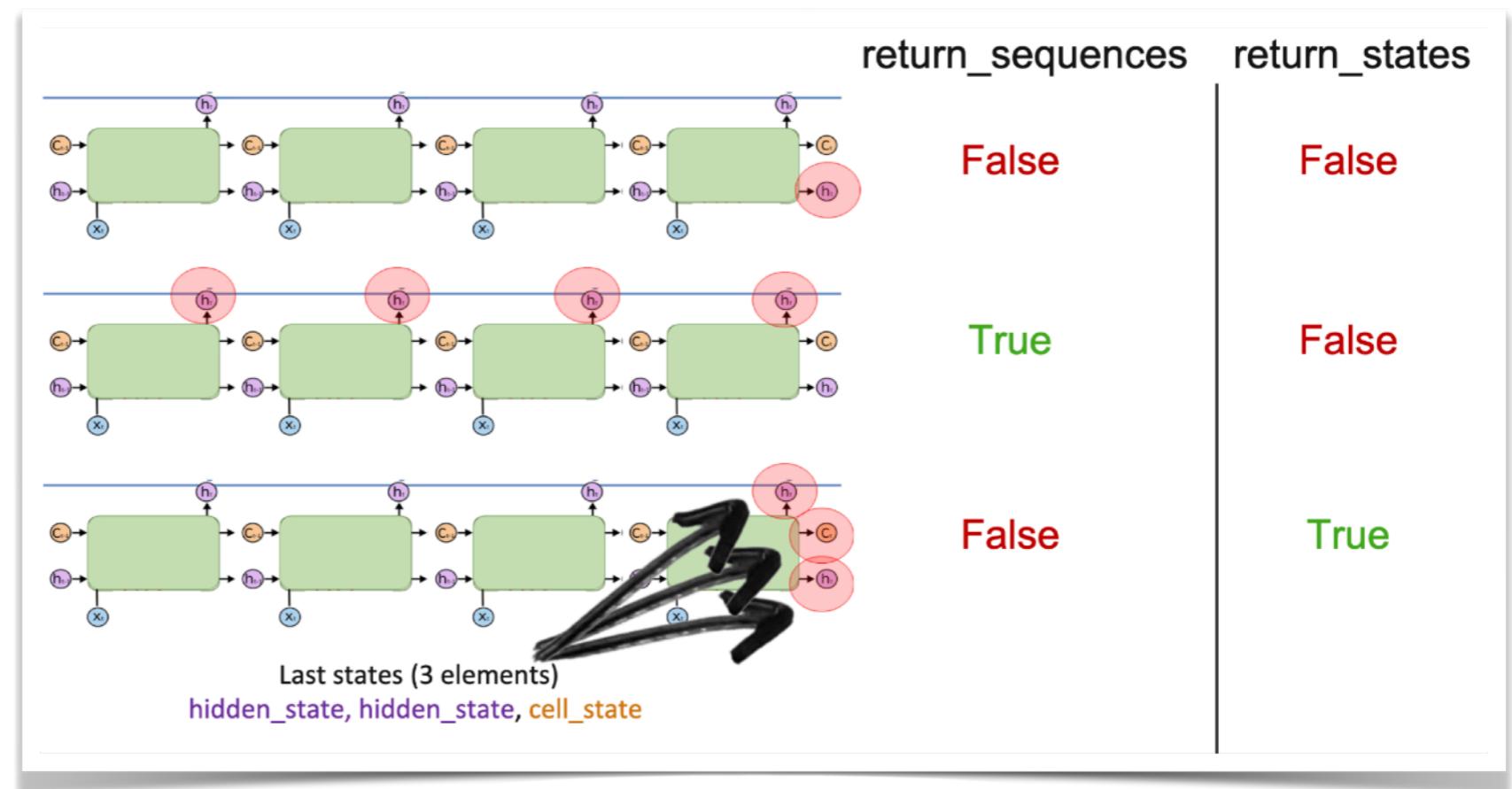


Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

2. RNN Model

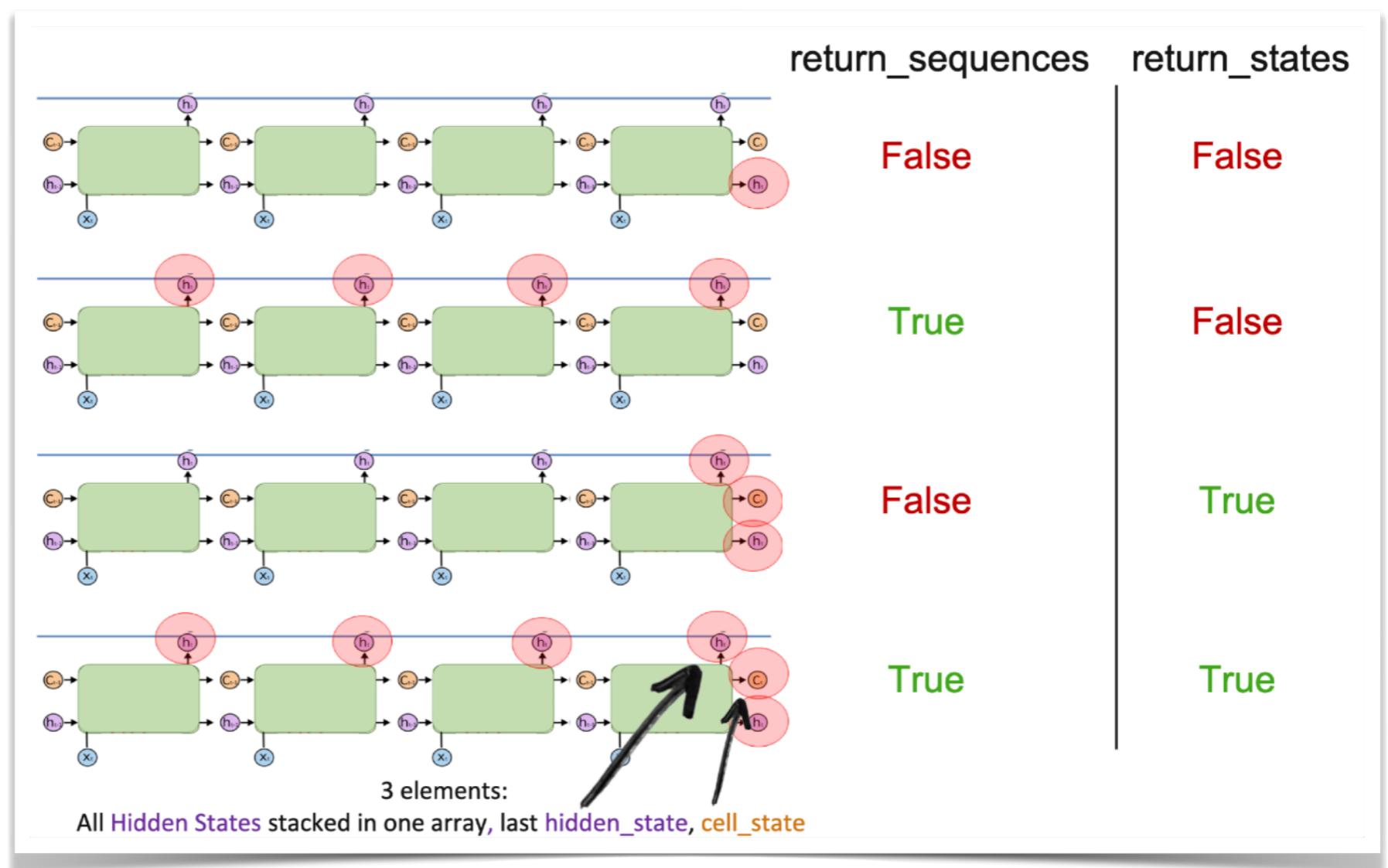


Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

2. RNN Model



Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

2. RNN Model

```
# ENCODER STEP
# -----
# Encoder with input (timesteps, features).
# LSTM return_state param will return last hidden state and last cell state
encoder_inputs = Input(shape=(n_timesteps_in, n_features), name='encoder_inputs')
encoder_lstm = LSTM(rnn_units, return_state=True, name='encoder_lstm')
encoder_outputs, state_h, state_c = encoder_lstm(encoder_inputs)

# initial context vector for the decoder (is the last output of the encoder)
states = [state_h, state_c]

# DECODER STEP
# -----
# Set up the decoder layers
# input shape: (1, n_features)
decoder_lstm = LSTM(rnn_units, return_state=True, name='decoder_lstm')
decoder_dense = Dense(n_features, activation='softmax', name='decoder_dense')
```

Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

2. RNN Model

```
all_outputs = []
# Prepare decoder input data.
# First input will be a sequence full of zeroes.
inputs = np.zeros((batch_size, 1, n_features))

# decoder will only process one timestep at a time.
for _ in range(n_timesteps_in):
    # Run the decoder on one timestep
    outputs, state_h, state_c = decoder_lstm(inputs, initial_state=states)
    outputs = tf.expand_dims(outputs, 1)
    outputs = decoder_dense(outputs)

    # Store the current prediction (we will concatenate all predictions later)
    all_outputs.append(outputs)

    # Reinject the outputs as inputs for the next loop iteration
    # as well as update the states
    inputs = outputs
    states = [state_h, state_c]

# Concatenate all predictions such as [batch_size, timesteps, features]
decoder_outputs = Lambda(lambda x: K.concatenate(x, axis=1))(all_outputs)

# Define and compile model
model = Model(encoder_inputs, decoder_outputs, name='model_encoder_decoder')
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
```

Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

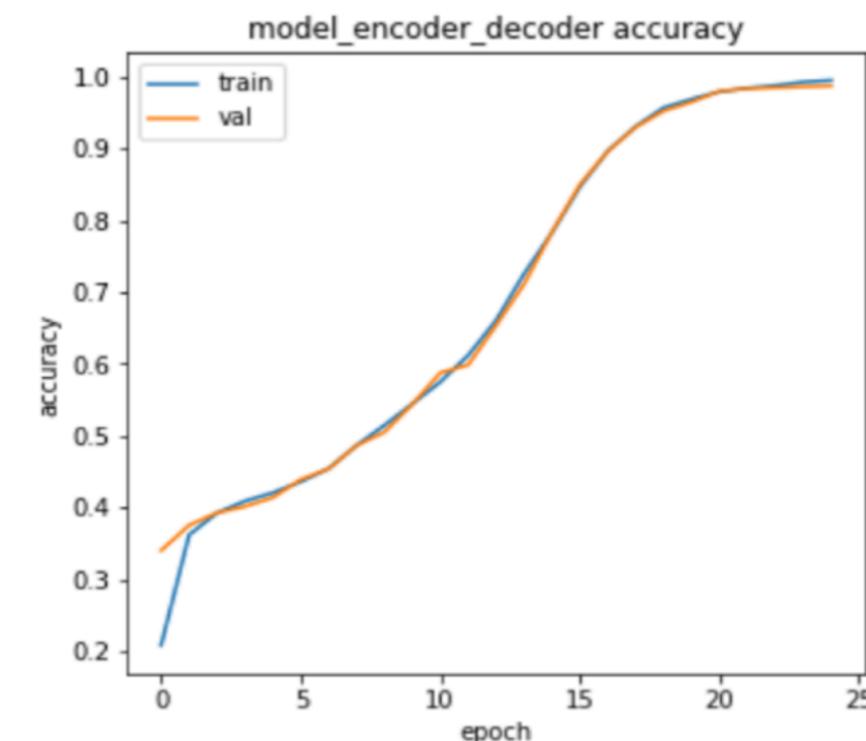
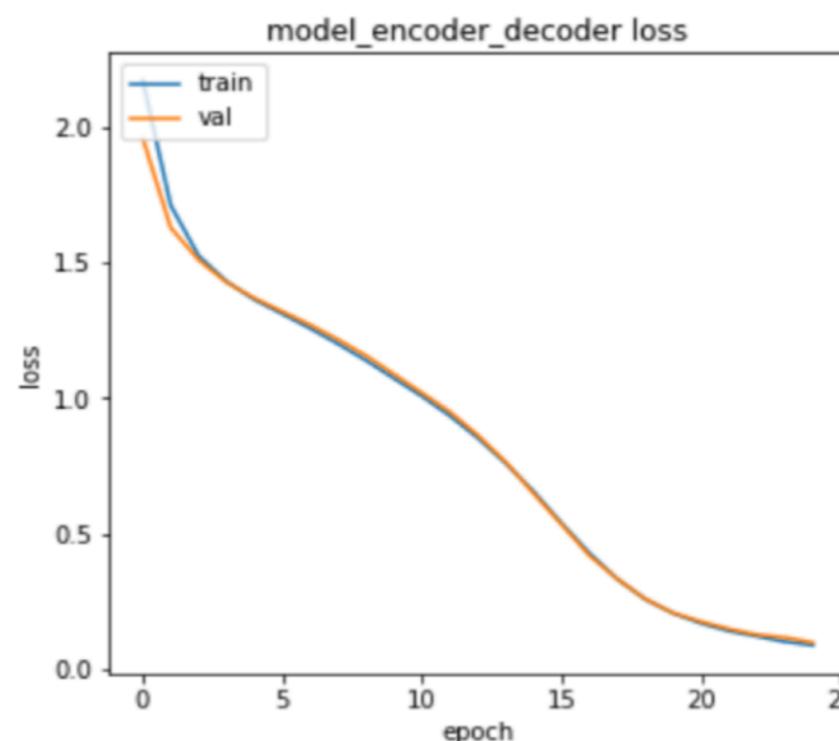
2. RNN Model

```
batch_size = 10
rnn_units = 16

# Define and compile model
model = Model(encoder_inputs, decoder_outputs, name='model_encoder_decoder')
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
```

PREDICTION ACCURACY (%):

Train: 99.550, Test: 98.125



Example

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

Notebook

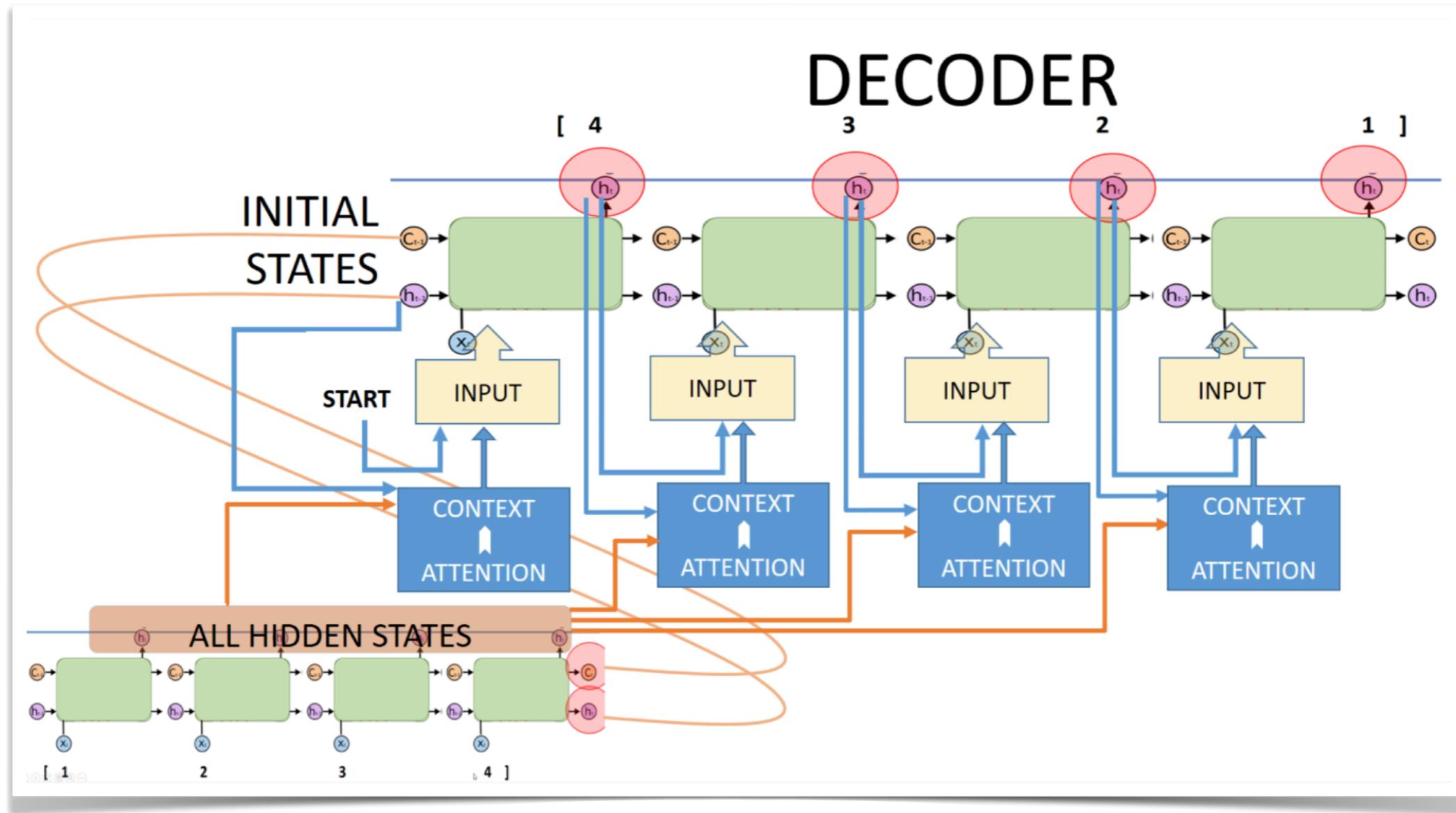
Exercise: Change these variables and retrain the model:

`n_steps_in = 20`
`n_features = 16`

Attention

```
# ENCODER STEP
# -----
# Same encoder as before with one and only difference. Now we need all the lstm states so we
# set return_sequences=True and return_state=True.
encoder_inputs = Input(shape=(n_timesteps_in, n_features), name='encoder_inputs')
encoder_lstm = LSTM(rnn_units, return_sequences=True, return_state=True, name='encoder_lstm')
encoder_outputs, encoder_state_h, encoder_state_c = encoder_lstm(encoder_inputs)
```

Attention



Attention involves all hidden states h_s , the last h_t and maybe some weight matrices.

Attention

```
# Decoder_outputs is the last hidden state of the encoder. Encoder_outputs are all the states
decoder_outputs = encoder_state_h

# Decoder will only process one time step at a time.
for _ in range(n_timesteps_in):

    # Pay attention!
    # decoder_outputs (last hidden state) + encoder_outputs (all hidden states)
    context_vector, attention_weights = attention(decoder_outputs, encoder_outputs)
    context_vector = tf.expand_dims(context_vector, 1)

    # create the context vector by applying attention to
    # Concatenate the input + context vectore to find the next decoder's input
    inputs = tf.concat([context_vector, inputs], axis=-1)

    # Passing the concatenated vector to the LSTM
    # Run the decoder on one timestep with attended input and previous states
    decoder_outputs, state_h, state_c = decoder_lstm(inputs, initial_state=states)
    outputs = decoder_dense(decoder_outputs)

    # Use the last hidden state for prediction the output
    # save the current prediction
    # we will concatenate all predictions later
    outputs = tf.expand_dims(outputs, 1)
    all_outputs.append(outputs)

    # Reinject the output (prediction) as inputs for the next loop iteration
    # as well as update the states
    inputs = outputs
    states = [state_h, state_c]
```

Attention

How to implement **attention score**(h_t, \bar{h}_s)?

$$\text{Attention weights } \alpha_{ts} = \frac{e^{score(h_t, \bar{h}_s)}}{\sum_{s'=1}^S e^{score(h_t, \bar{h}_{s'})}}$$

Hidden state Decoder Output

$$\text{Context Vector } c_t = \sum_S \alpha_{ts} \bar{h}_S$$

Luong dot $score(h_t, \bar{h}_s) = h_t^T \bar{h}_s$ No training parameters!

Luong multiplicative $score(h_t, \bar{h}_s) = h_t^T W_a \bar{h}_s$

Bahdanau $score(h_t, \bar{h}_s) = v_a^T \tanh(W_a[h_t; \bar{h}_s])$

Attention

```
class LuongDotAttention(tf.keras.layers.Layer):
    def __init__(self):
        super(LuongDotAttention, self).__init__()

    def call(self, query, values):
        query_with_time_axis = tf.expand_dims(query, 1)
        values_transposed = tf.transpose(values, perm=[0, 2, 1])

        # LUONGH Dot-product
        score = tf.transpose(tf.matmul(query_with_time_axis,
                                       values_transposed), perm=[0, 2, 1])  

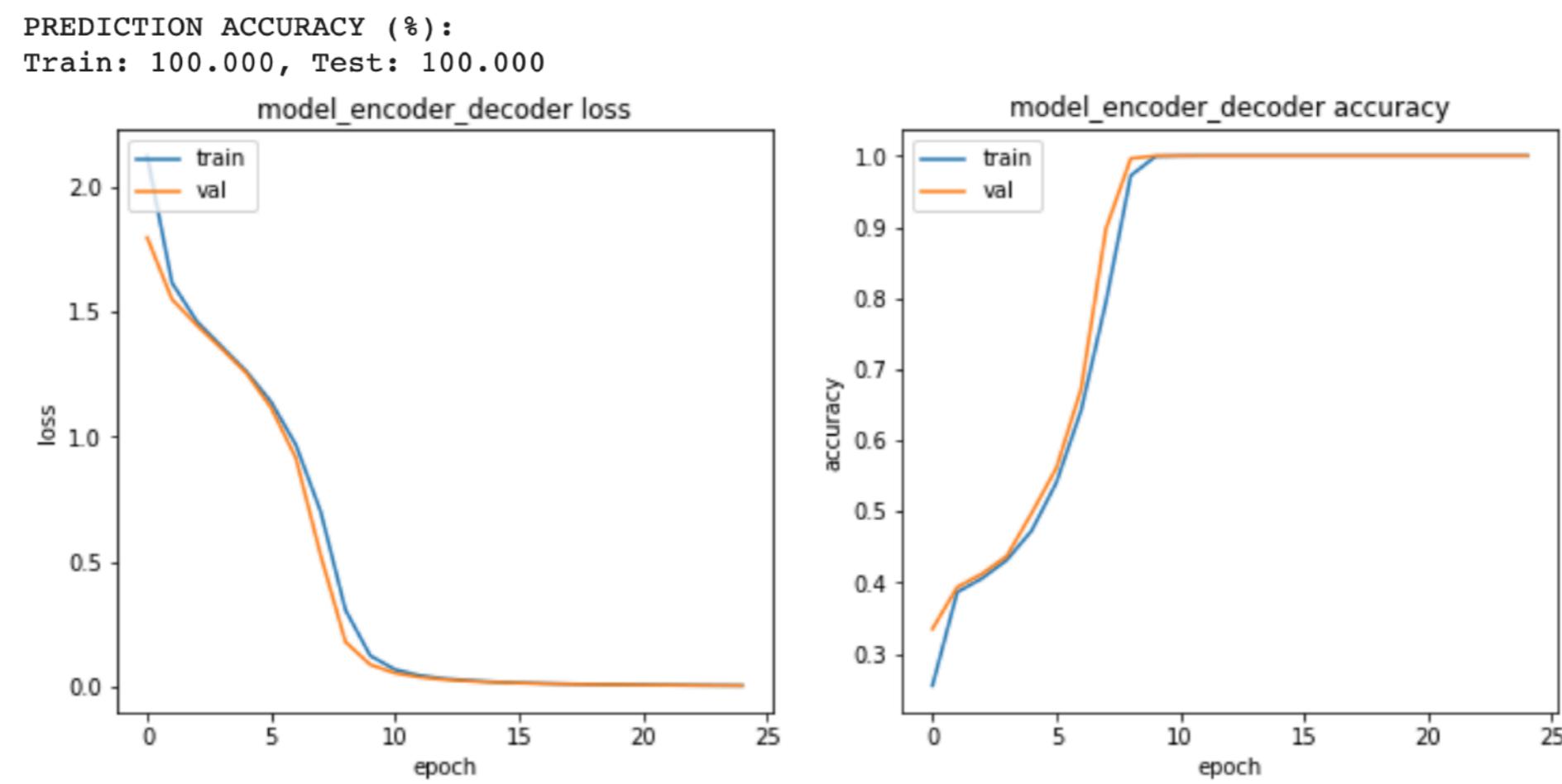
  

        # attention_weights shape == (batch_size, max_length, 1)
        attention_weights = tf.nn.softmax(score, axis=1)  $\alpha_{ts} = \frac{e^{score(h_t, \bar{h}_s)}}{\sum_{s'=1}^S e^{score(h_t, \bar{h}_{s'})}}$ 
  

        # context_vector shape after sum == (batch_size, hidden_size)
        context_vector = attention_weights * values
        context_vector = tf.reduce_sum(context_vector, axis=1)  $c_t = \sum_S \alpha_{ts} \bar{h}_s$ 
  

        return context_vector, attention_weights
```

Attention



Example

(Credit: Pere Gilabert)

Reversed sequences:

Given $X = [2,8,1,6]$ predict $Y = [6,1,8,2]$.

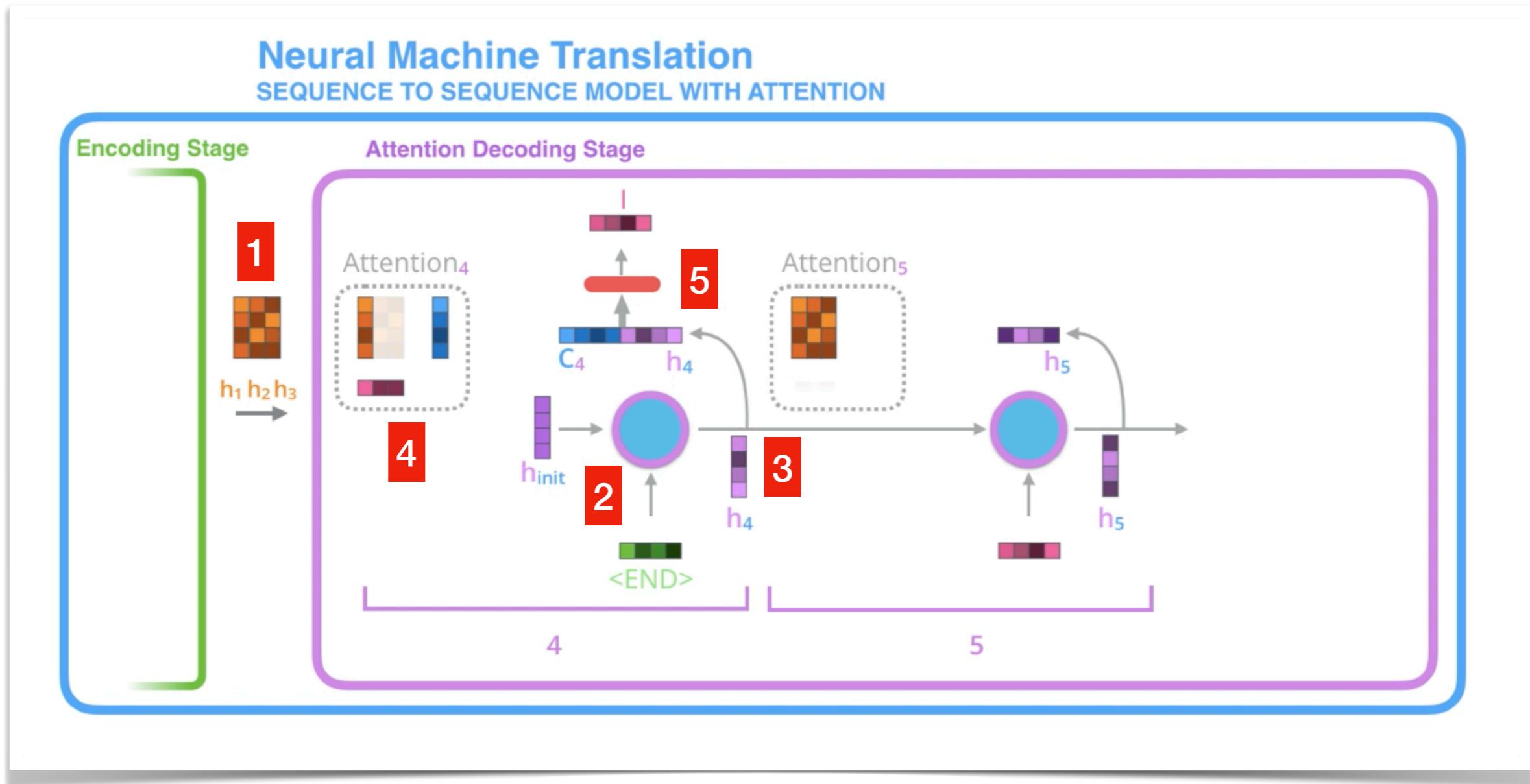
Notebook

Reminder: Attention

In the case of a encoder-decoder architecture, **to focus on the parts of the input that are relevant to a decoding time step**, the decoder does the following:

- **Look at the set of encoder hidden states** it received (each encoder hidden state is most associated with a certain word in the input sentence).
- **Give each hidden states a **score**.**
- Multiply each hidden state by its softmaxed score, thus **amplifying hidden states with high scores**, and drowning out hidden states with low scores.
- Build a **new representation** from these “weighted” encoder hidden states.

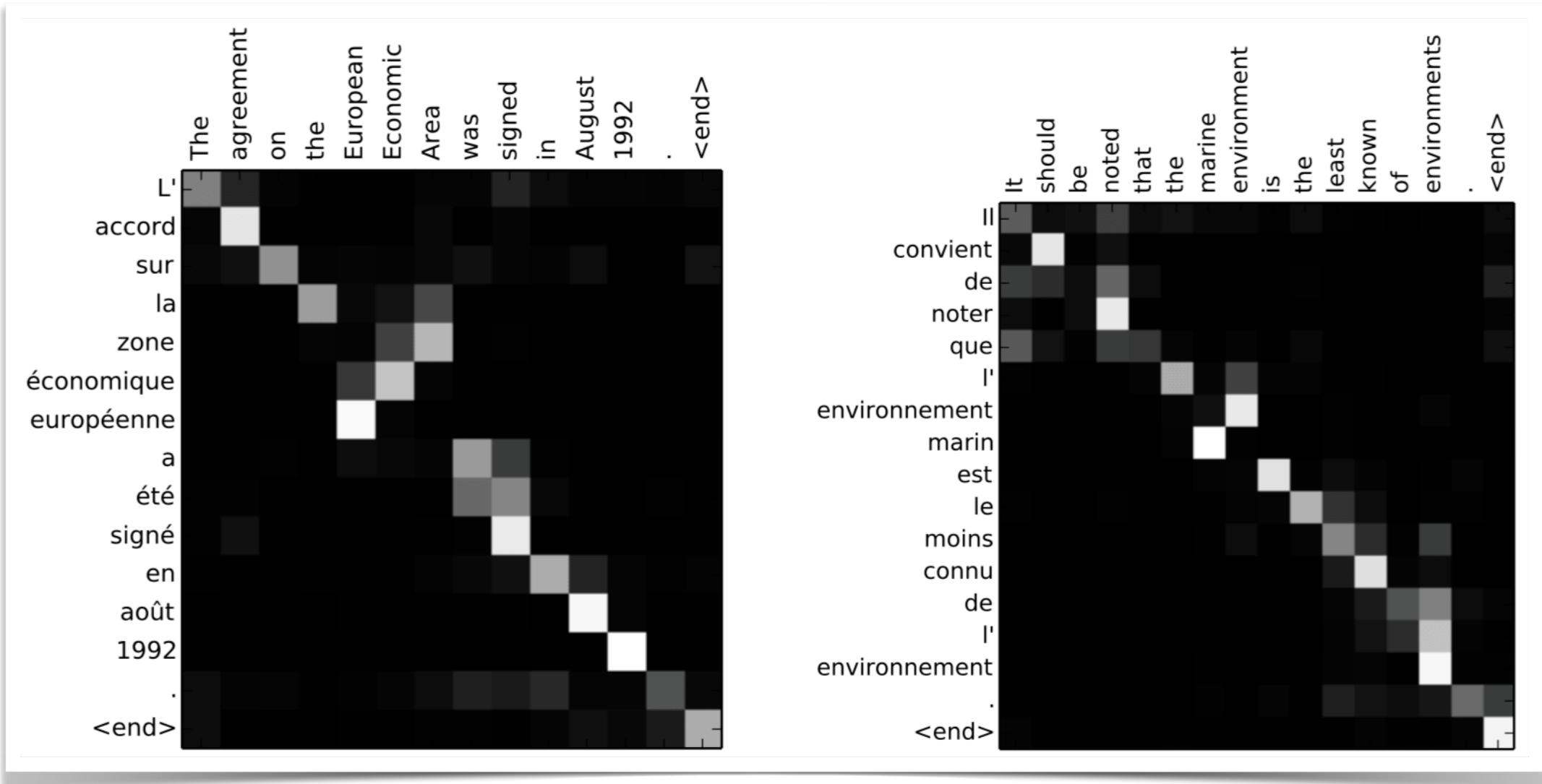
Reminder: Attention



4 $score(h_4, \bar{h}_s) = h_4^T \bar{h}_s$



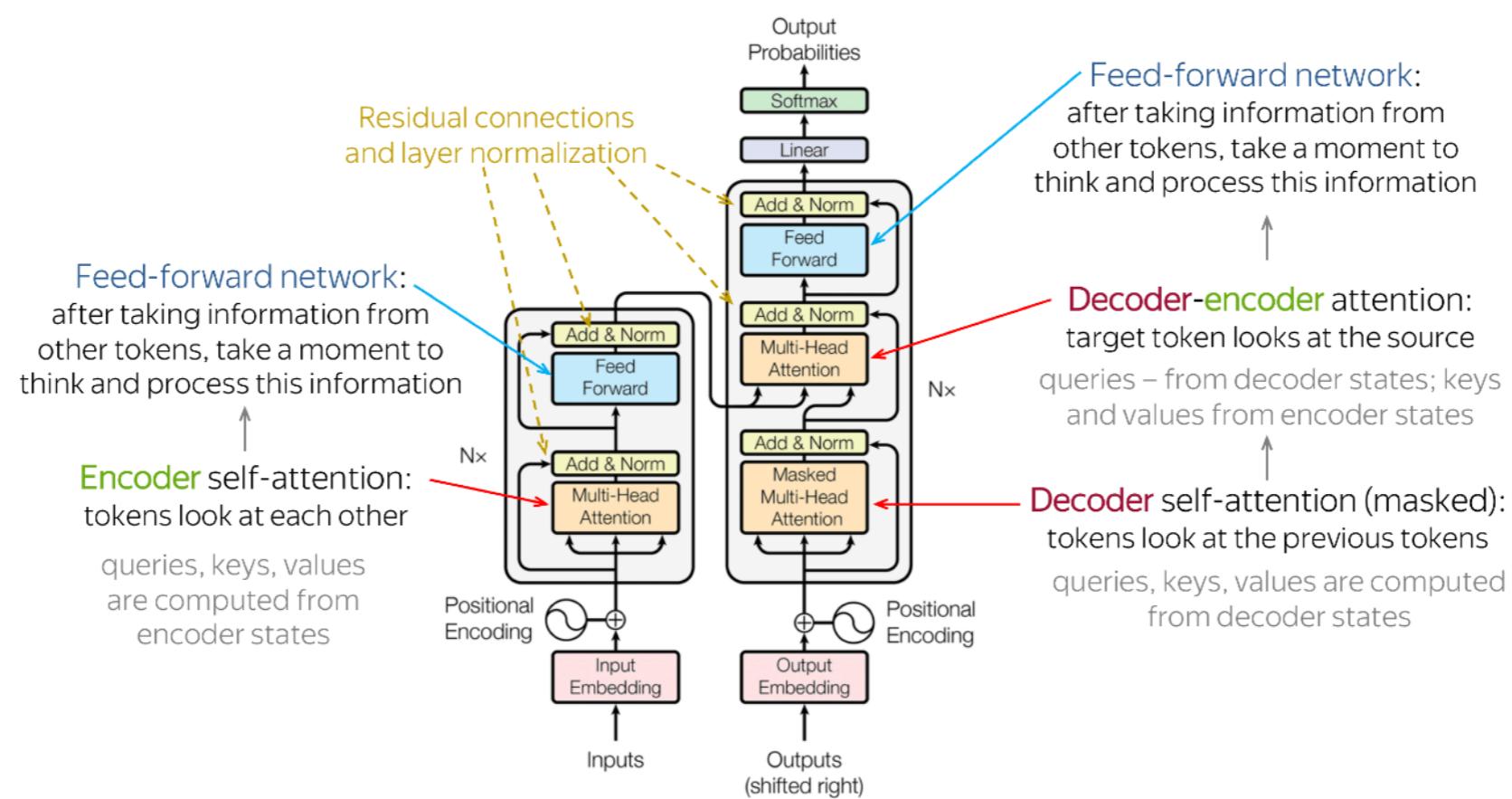
Reminder: Attention



Transformer

The Transformer Neural Networks — usually just called “**Transformers**” — were introduced by a Google-led team in 2017 in a paper titled “Attention Is All You Need”.

The Transformer has an **encoder-decoder** architecture but it is based solely on attention mechanisms: i.e., without recurrence or convolutions.



Transformer

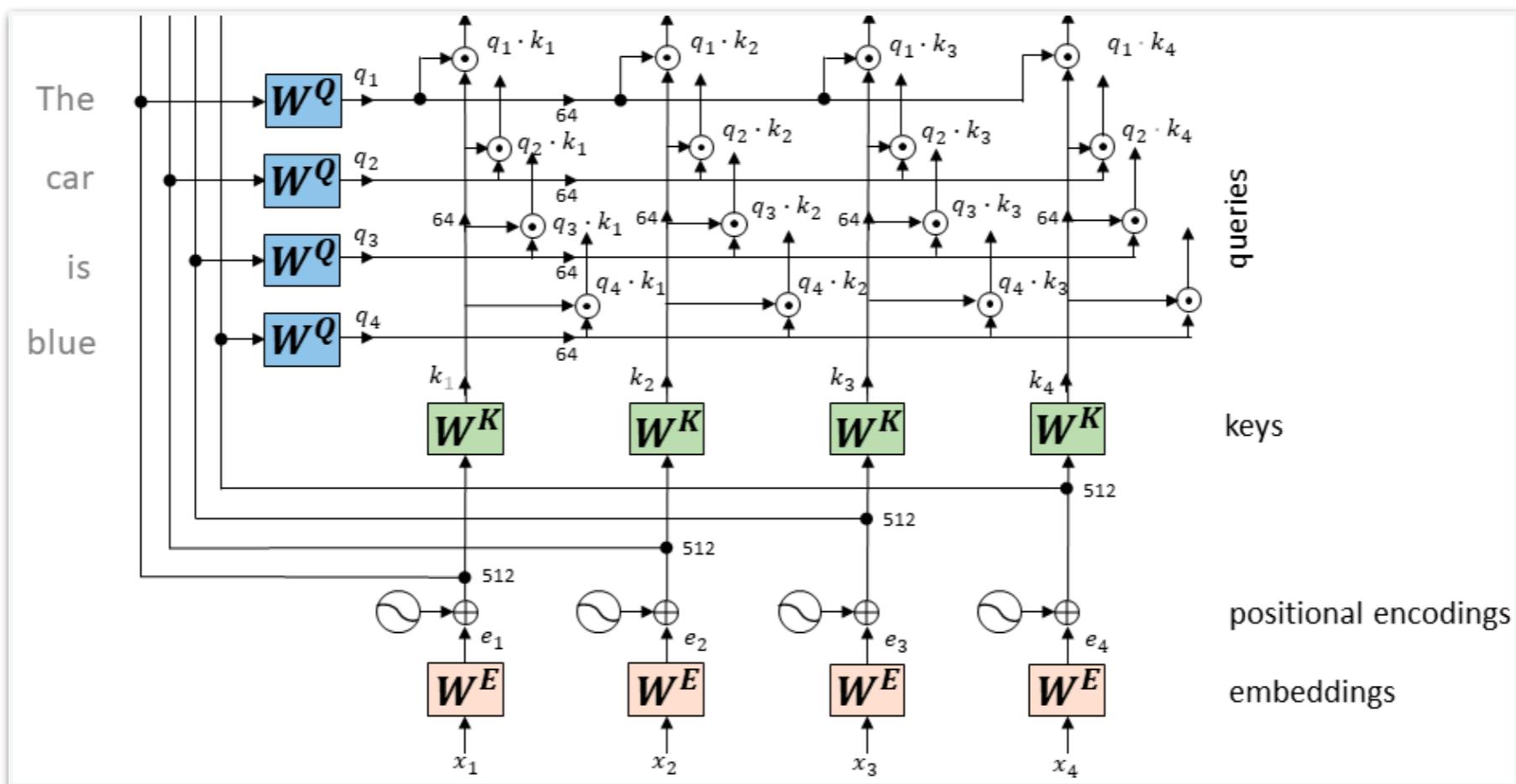
In the last few years, several architectures based on the basic transformer introduced in the 2017 paper have been developed and trained for complex natural language processing tasks:

- BERT
- DistilBERT
- T5
- GPT-2 and GPT-3
- Etc.

HuggingFace has implemented a Python package for transformers that is really easy to use. It is open-source and you can find it on GitHub.

Transformer Encoder

We will focus on the encoder architecture.



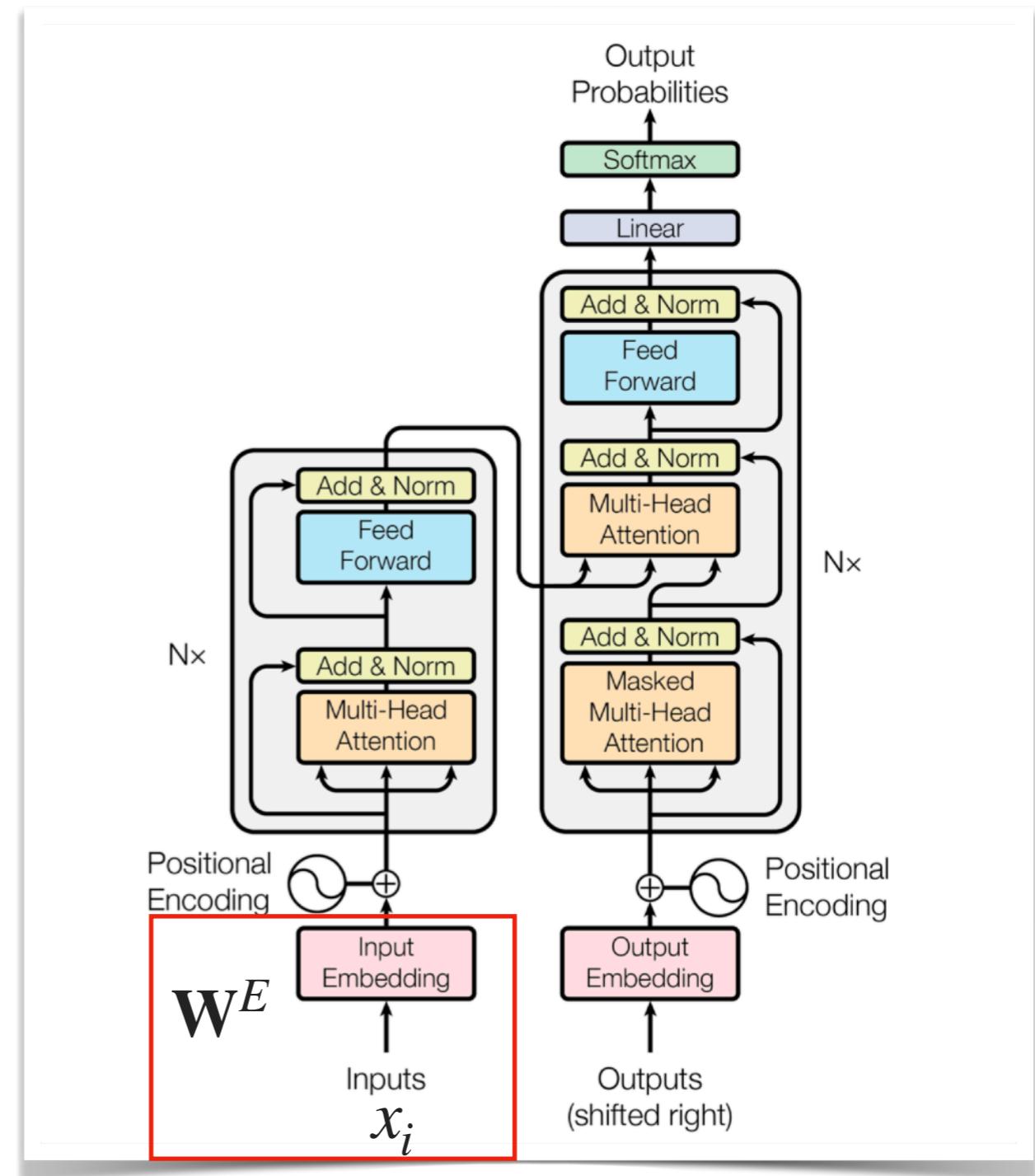
<https://towardsdatascience.com/drawing-the-transformer-network-from-scratch-part-1-9269ed9a2c5e>

Transformer Input

A Transformer takes as input a sequence of tokens x_i , which are presented to the network as **one-hot encodings**.

Next, we reduce the dimensionality of the one-hot encoded vectors by multiplying them with a so called “**embedding matrixW^E. The resulting vectors e_i are called token embeddings.**

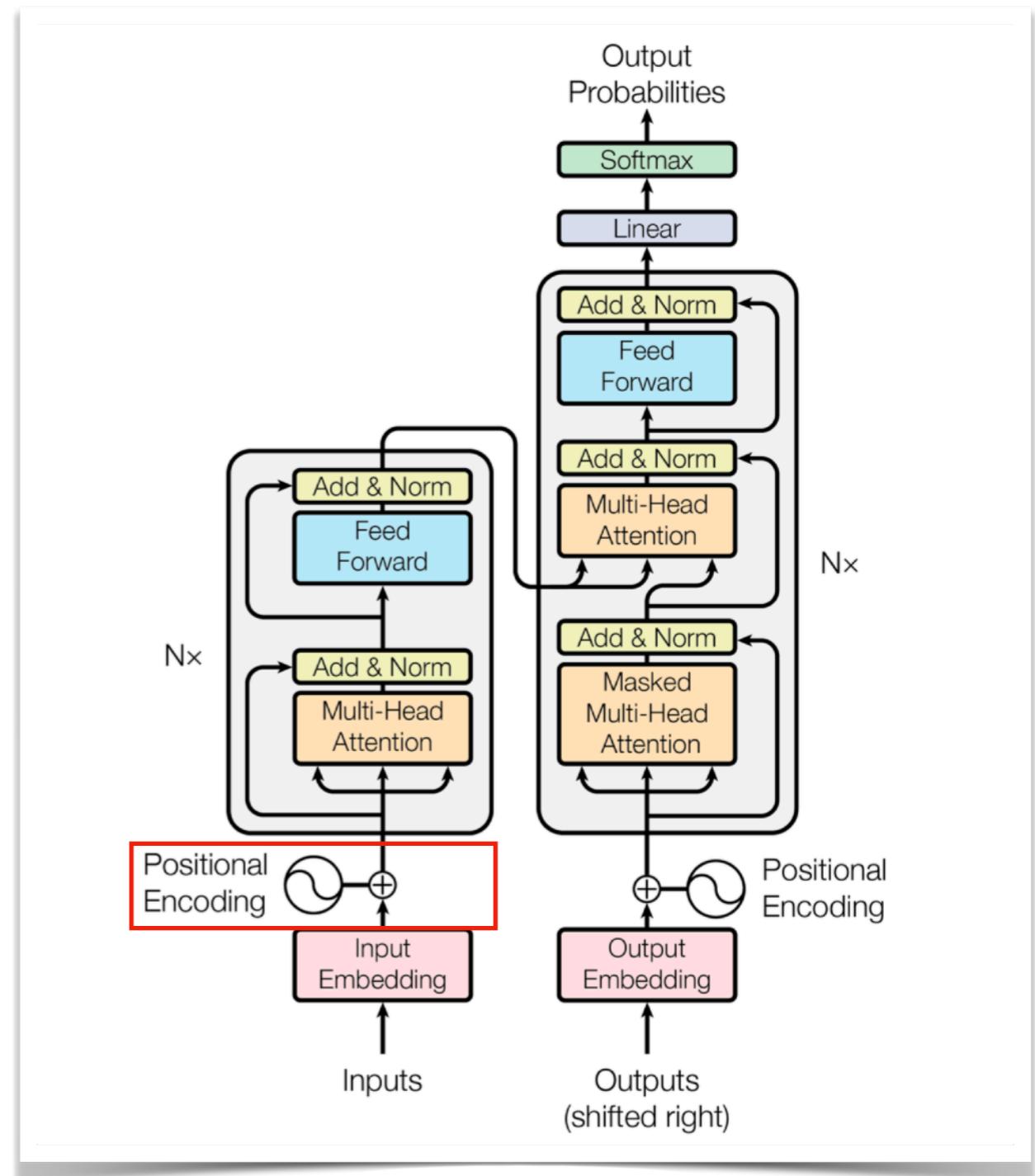
The size of the token (word) embeddings in the original paper is $d = 512$.



Transformer Input

All the tokens are presented to the Transformer **simultaneously**. However, this means that the **order in which tokens occur in the input sequence is lost**.

To address this, the Transformer **adds a vector to each input embedding**, thus injecting some information about the position.



Transformer Input

The encoding proposed by the authors is a d -dimensional vector that contains information about a specific **position** in a sentence.

It must satisfy these criteria:

- It should output a unique encoding for each time-step (token's **position** in a sentence)
- Distance between any two time-steps should be consistent across sentences with different lengths.
- The model should generalize to longer sequences without any efforts.
- Its values should be bounded.
- It must be deterministic.

Transformer Input

Let t be the desired position in an input sentence, $\vec{p}_t \in \mathbb{R}^d$ be its corresponding encoding, and d be the encoding dimension.

Then, $f: \mathbb{N} \rightarrow \mathbb{R}^d$ will be the function that produces the output vector \vec{p}_t and it is defined as follows:

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

where

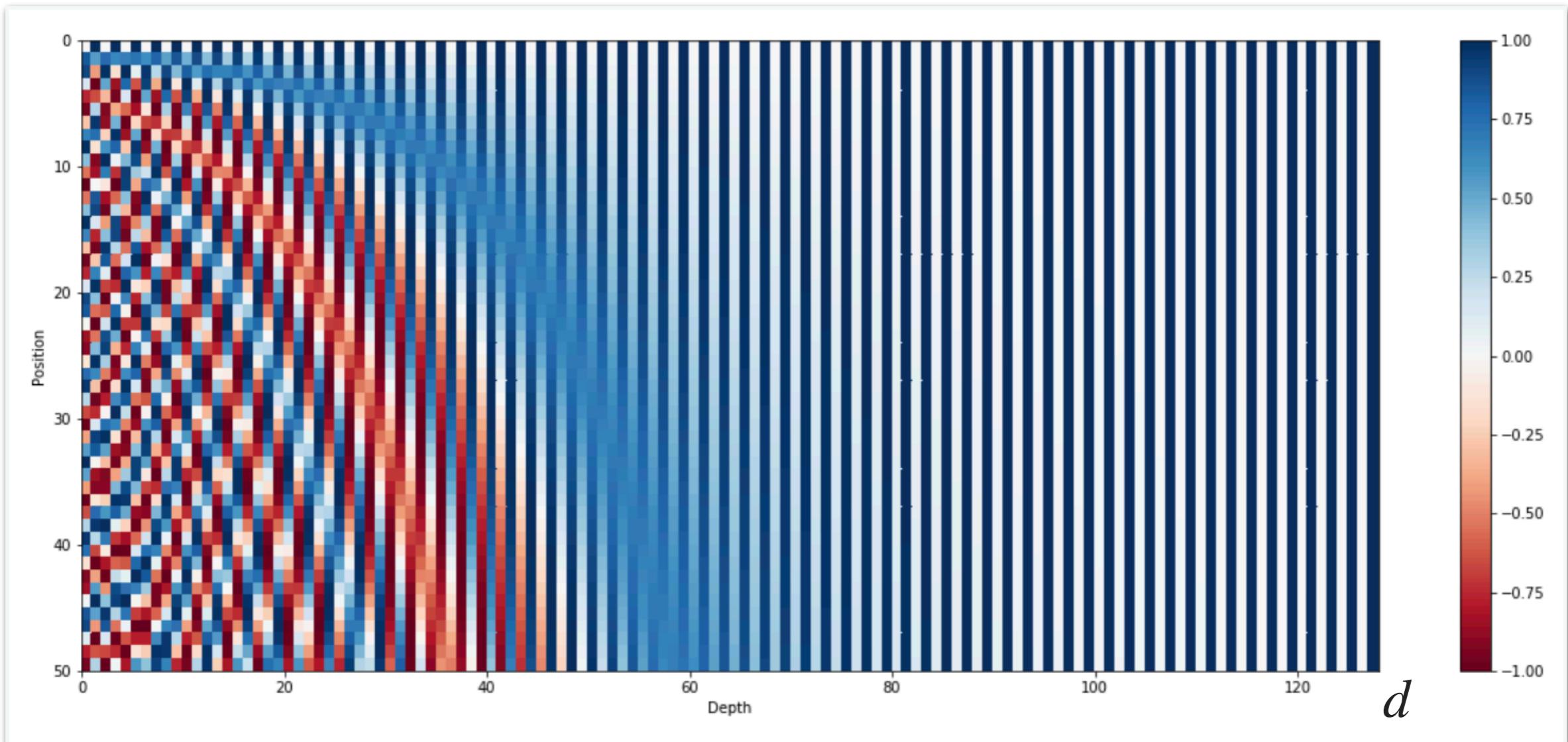
$$\omega_k = \frac{1}{10000^{2k/d}}$$

Transformer Input

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

But, how this combination of sines and cosines could ever represent a position/order?

Transformer Input



The 128-dimensional positonal encoding for a sentence with the maximum lenght of 50. Each row represents the embedding vector \vec{p}_t

Transformer Input

In the original paper **added the positional encoding on top of the actual embeddings.**

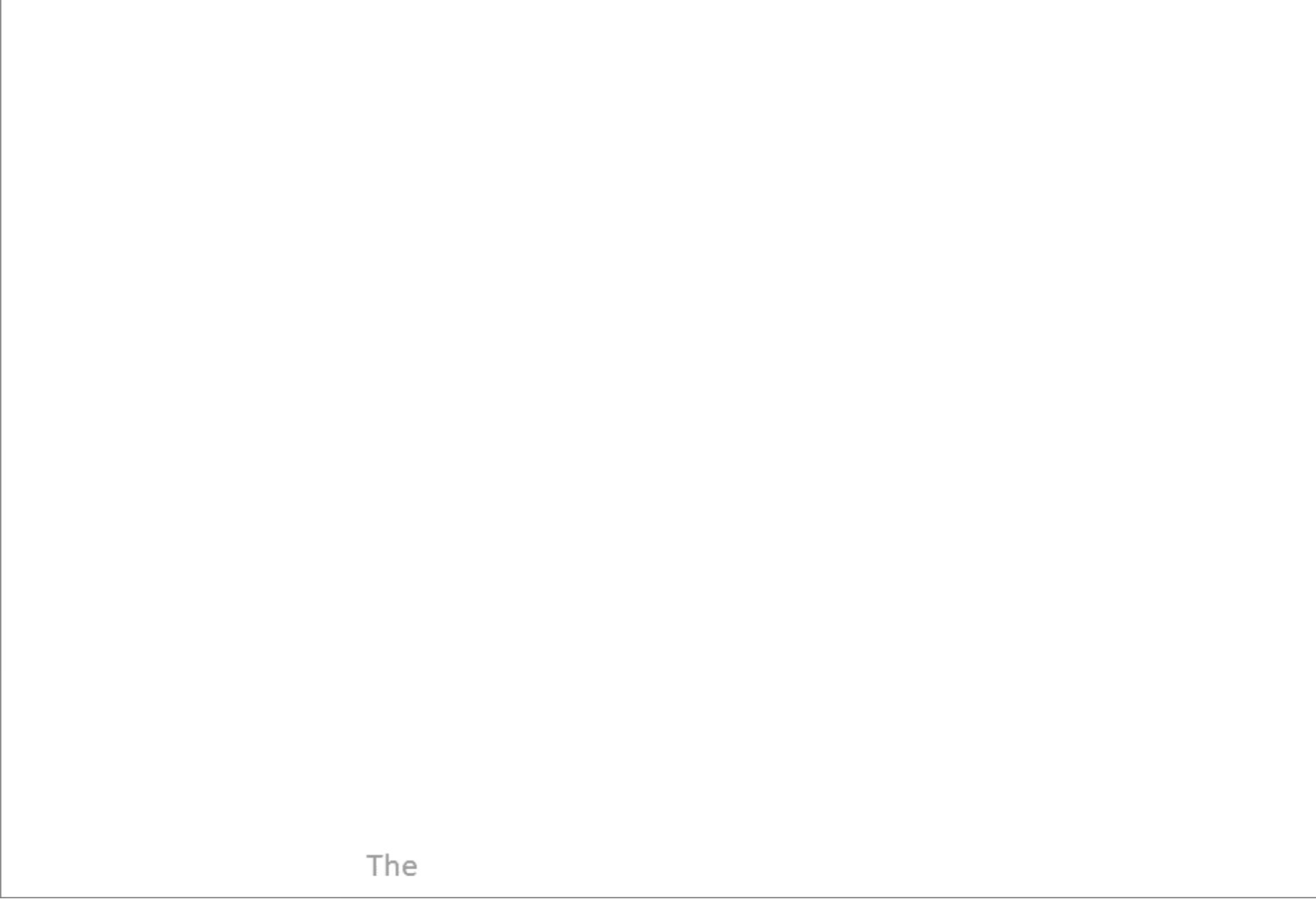
That is for every word embedding e_t in a sentence of words $[x_1, \dots, x_n]$ the correspondent “position-aware” embedding which is fed to the model is as follows:

$$\psi'(e_t) = \psi(e_t) + \vec{p}_t$$

To make this summation possible, we keep the positional embedding’s dimension equal to the word embeddings’ dimension

Keys and Queries

Next, we multiply the position-aware word embeddings by two matrices, W^Q and W^K , to obtain the “**query vectors**” and “**key vectors**”, each of size 64.

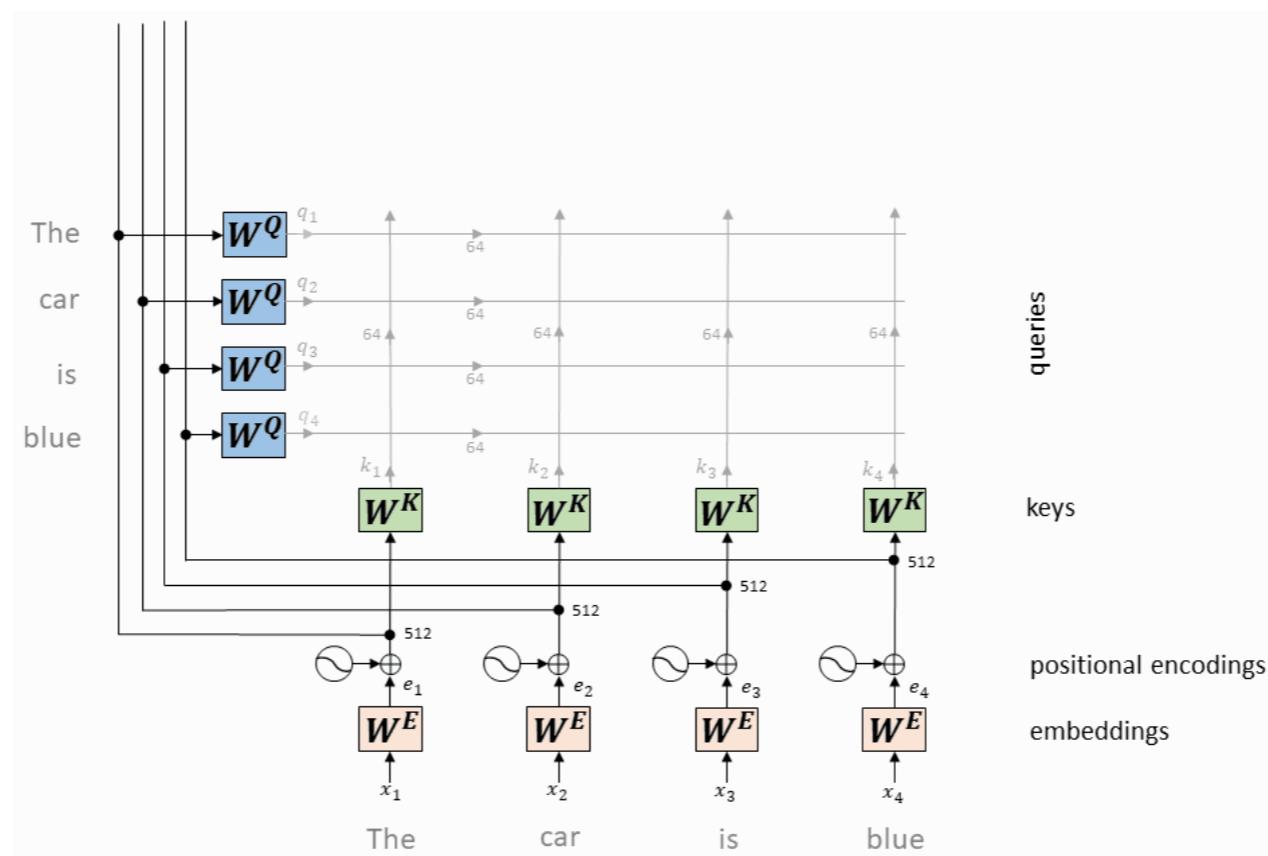


The

Keys and Queries

We calculate the dot products for all possible combinations of “query vectors” and “key vectors”.

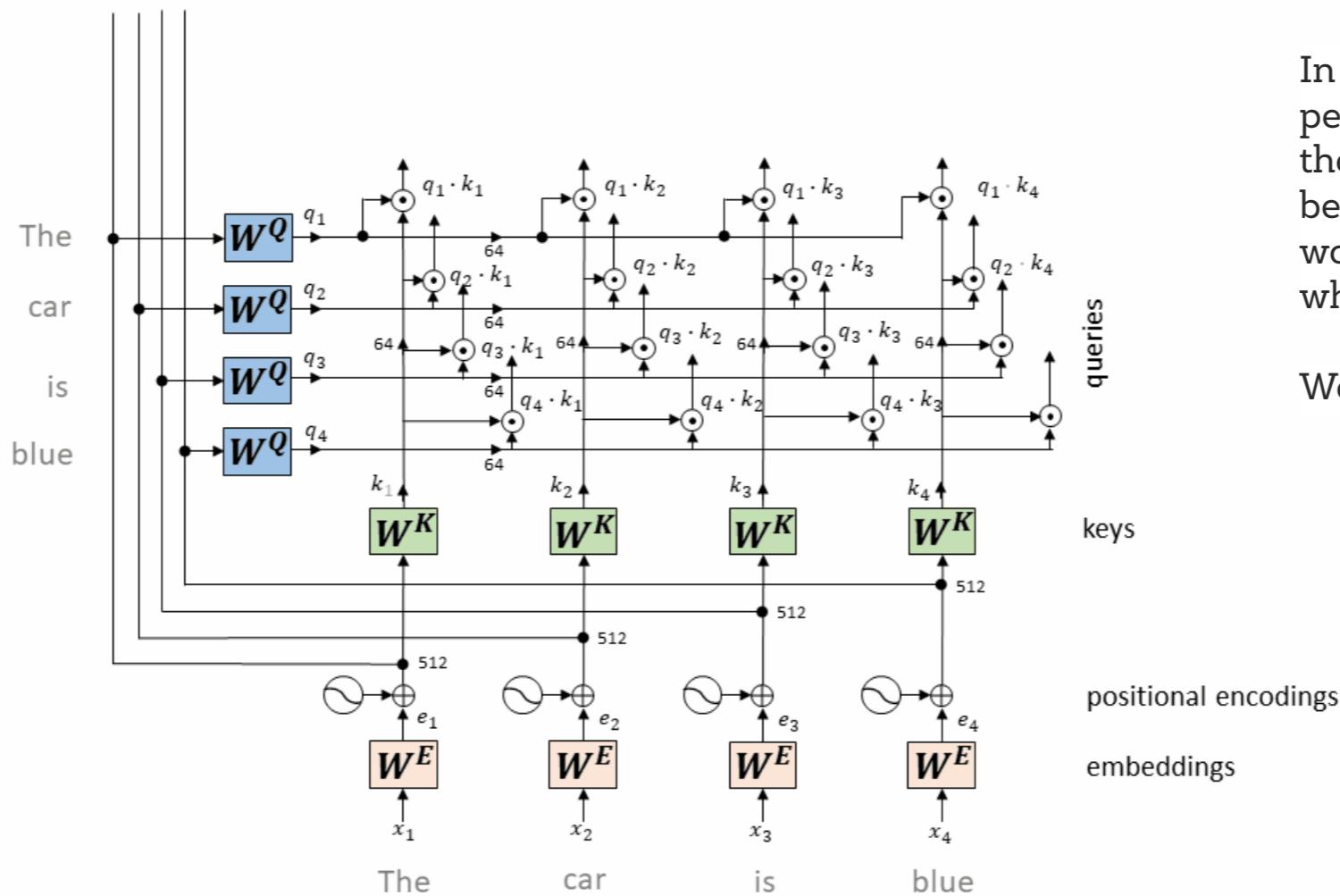
The result of a dot product is a single number, which in a later step will be used as a weight factor. **The weights factors tell us, how much two words at different positions of the input sentence depend on each other.** This is called **self-attention** in the original paper.



Keys and Queries

Subsequently, all weight factors are divided by 8 (the square root of the dimension of the key vectors 64), to have more stable gradients.

The scaled factors are put through a softmax function, which normalizes them so they are all positive and sum up to 1 to get $s_{i,j}$.

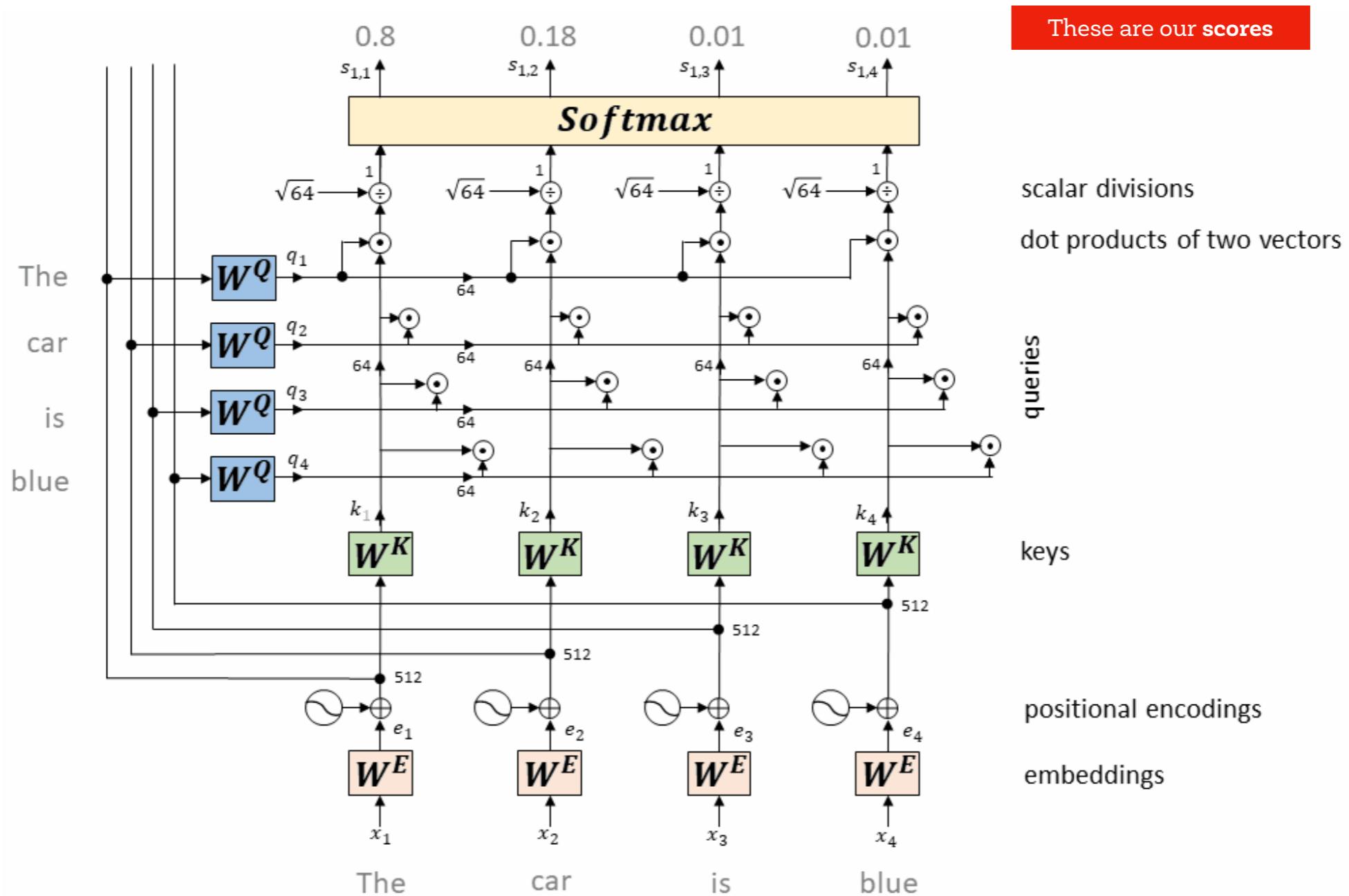


In this animation, we perform the scaling for the weight factors belonging to the first word in our sentence, which is "The".

We get $s_{1,1}, s_{1,2}, s_{1,3}, s_{1,4}$

Keys and Queries

Analogously, for the other words “car”, “is” and “blue” in our input sequence to get all $s_{i,j}$:



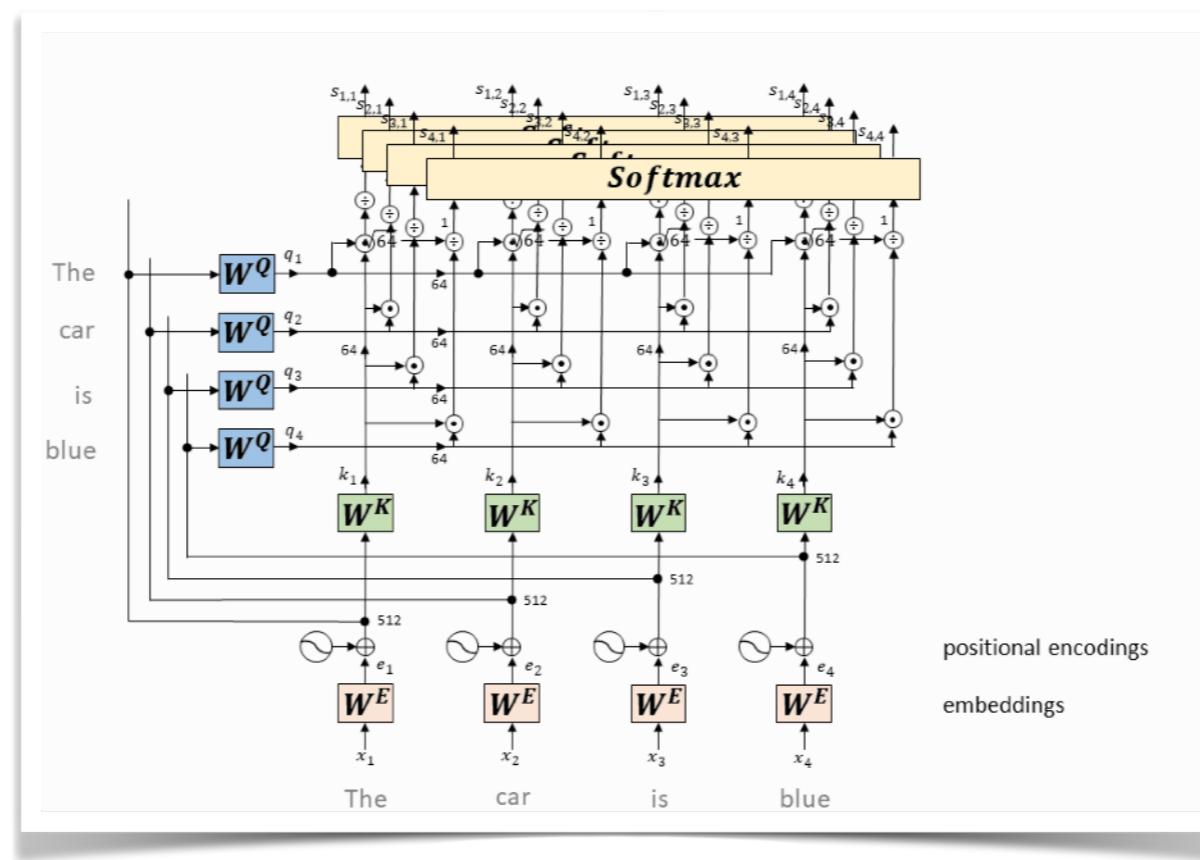
Values, weighting and summation

Identical to the computation of the “key vector” and “query vectors” we obtain the “value vectors” v_k by multiplying the word embeddings by a matrix W^V .

Each “value vector” v_k is multiplied by its corresponding “weight factors” $s_{k,i}, \forall i$.

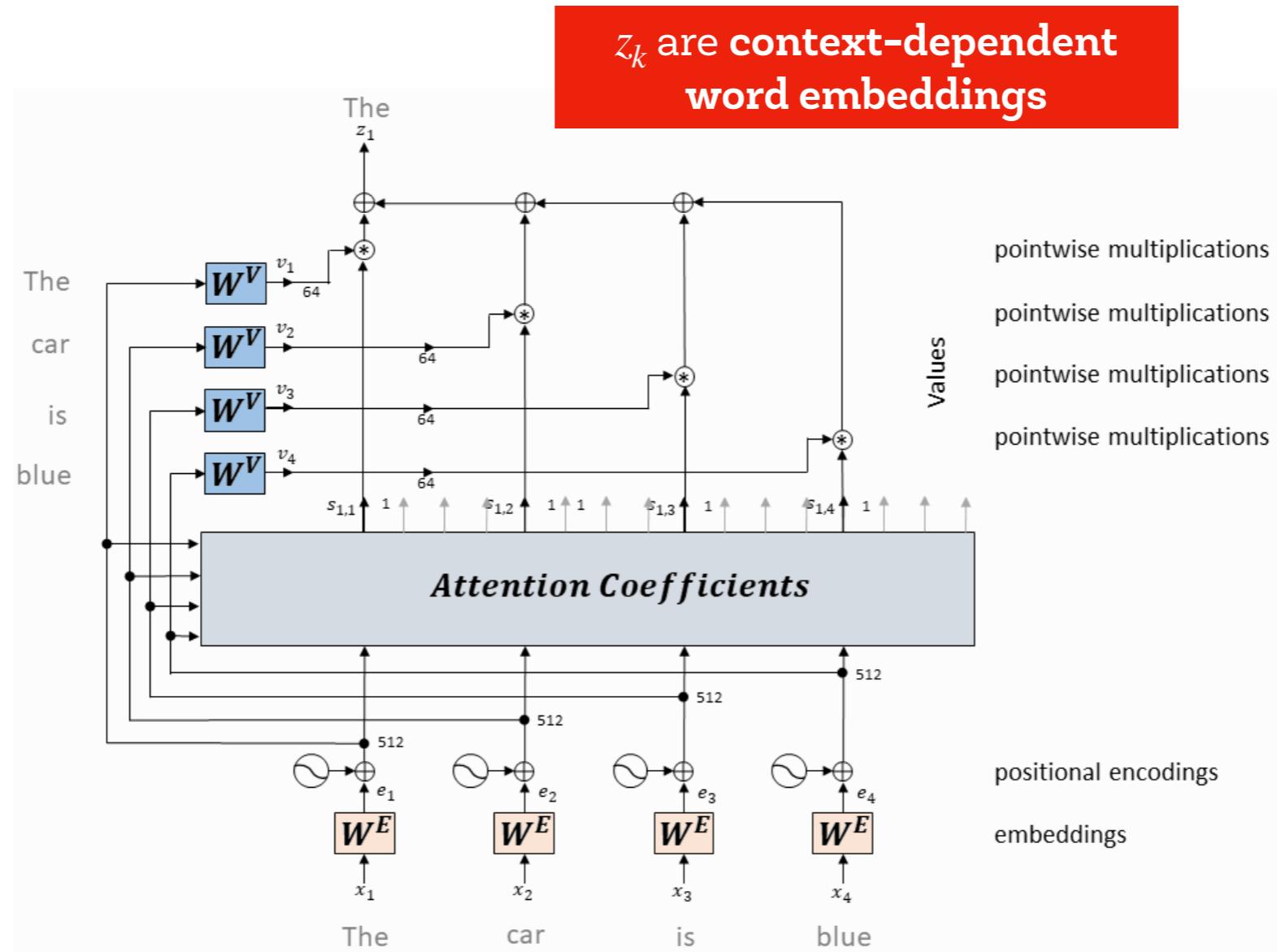
Finally, we add together all these values to get z_k .

This way we only keep the words we want to focus on, while irrelevant words are suppressed by weighting them by tiny numbers like 0.001



Values, weighting and summation

Analogously for the other words “car”, “is”, “blue” in our input sequence.



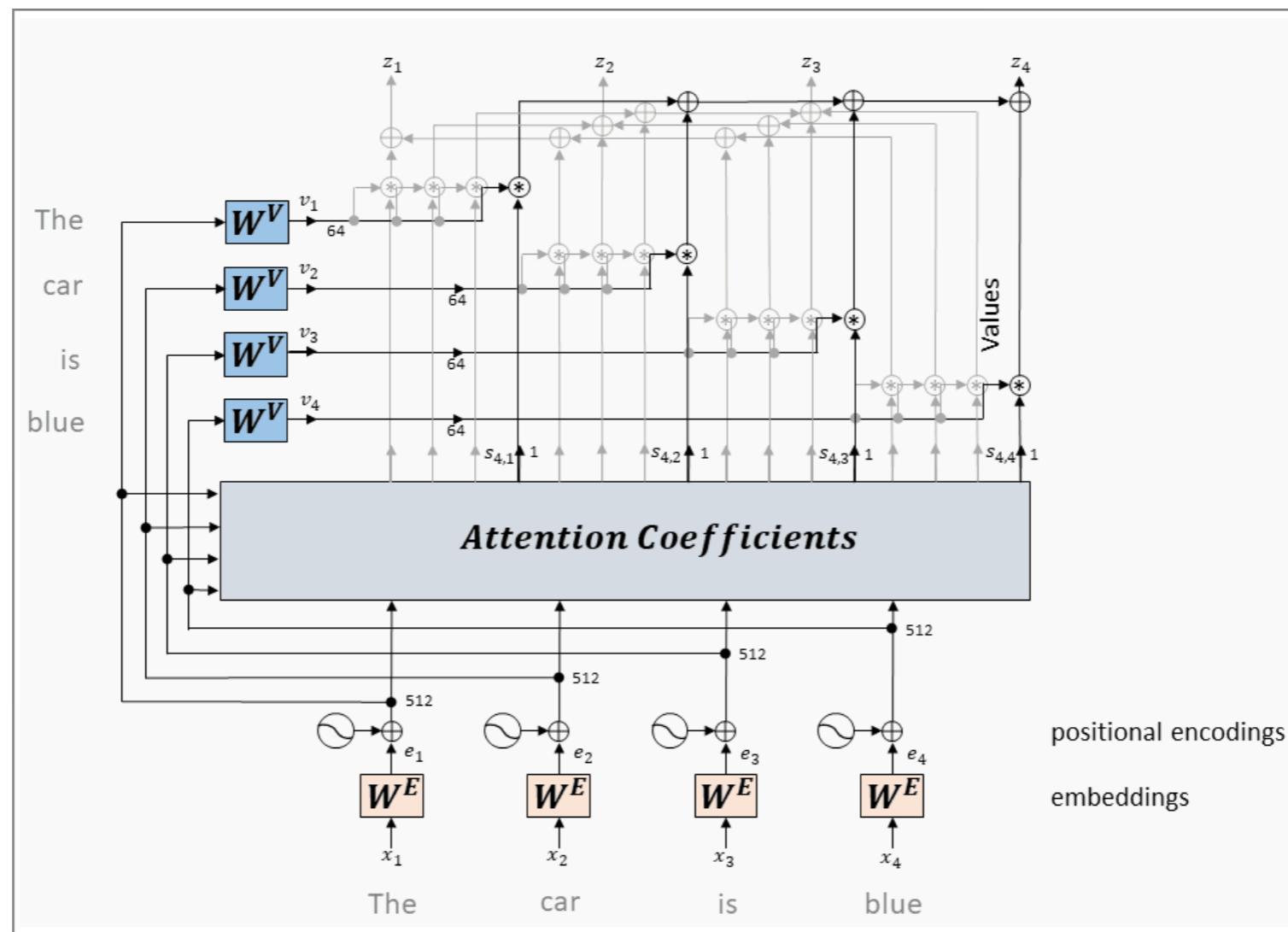
Attention in matrix form

Since we're dealing with matrices, we can condense all steps 2 through 6 in one formula to calculate the outputs of the **self-attention layer**:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Self-attention

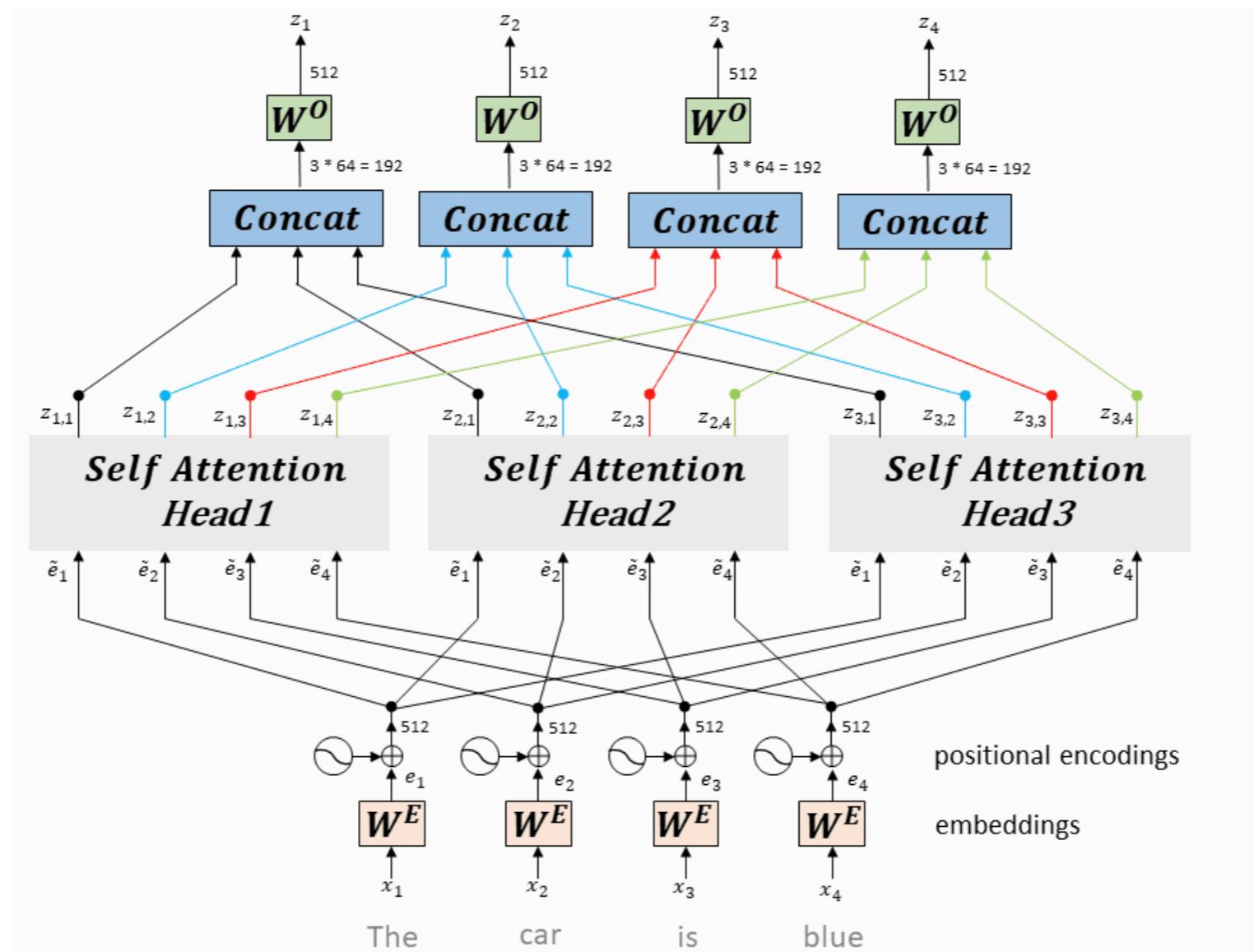
That could conclude a **self-attention calculation**. The output of the self-attention layer can be considered as a **context enriched word embedding**, but instead of performing a single self-attention function, the authors employ multiple self-attention heads, each with different weight matrices.



The outputs of the attention heads are concatenated and once again multiplied by an additional weights matrix to get new z_k .

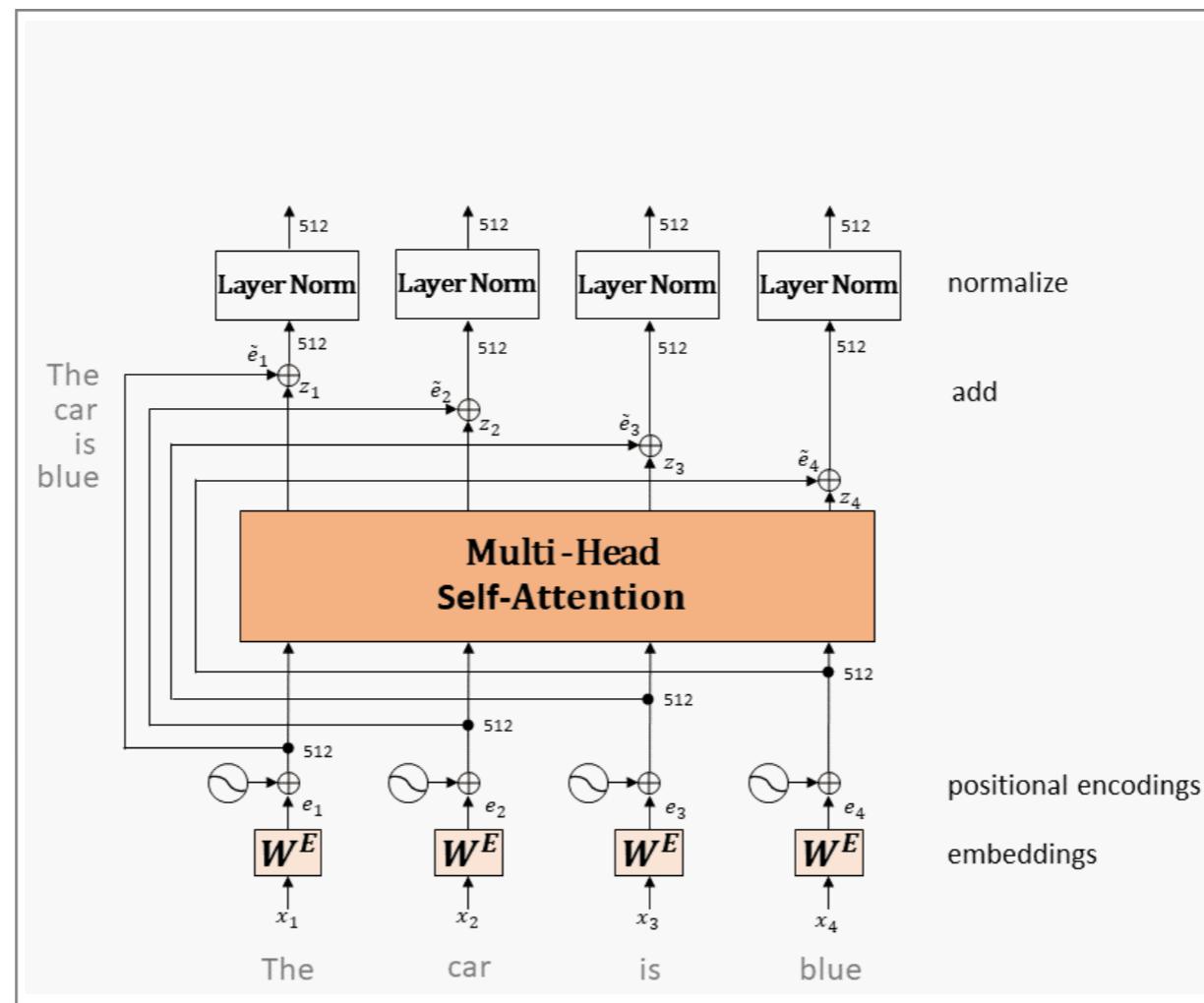
Add and normalize

The multi-head self-attention mechanism, just covered, is the first sub-module of the encoder. It has a residual connection around it, and is followed by a layer-normalization step. Layer-normalization just subtracts the mean of each vector and divides by its standard deviation.



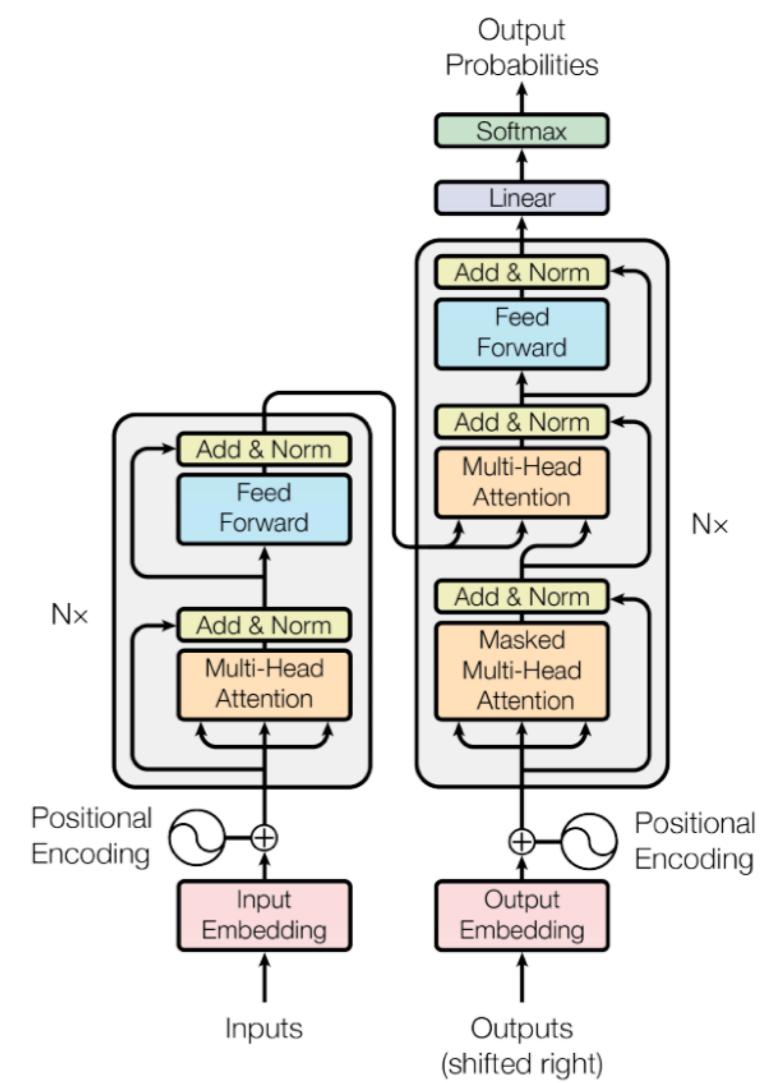
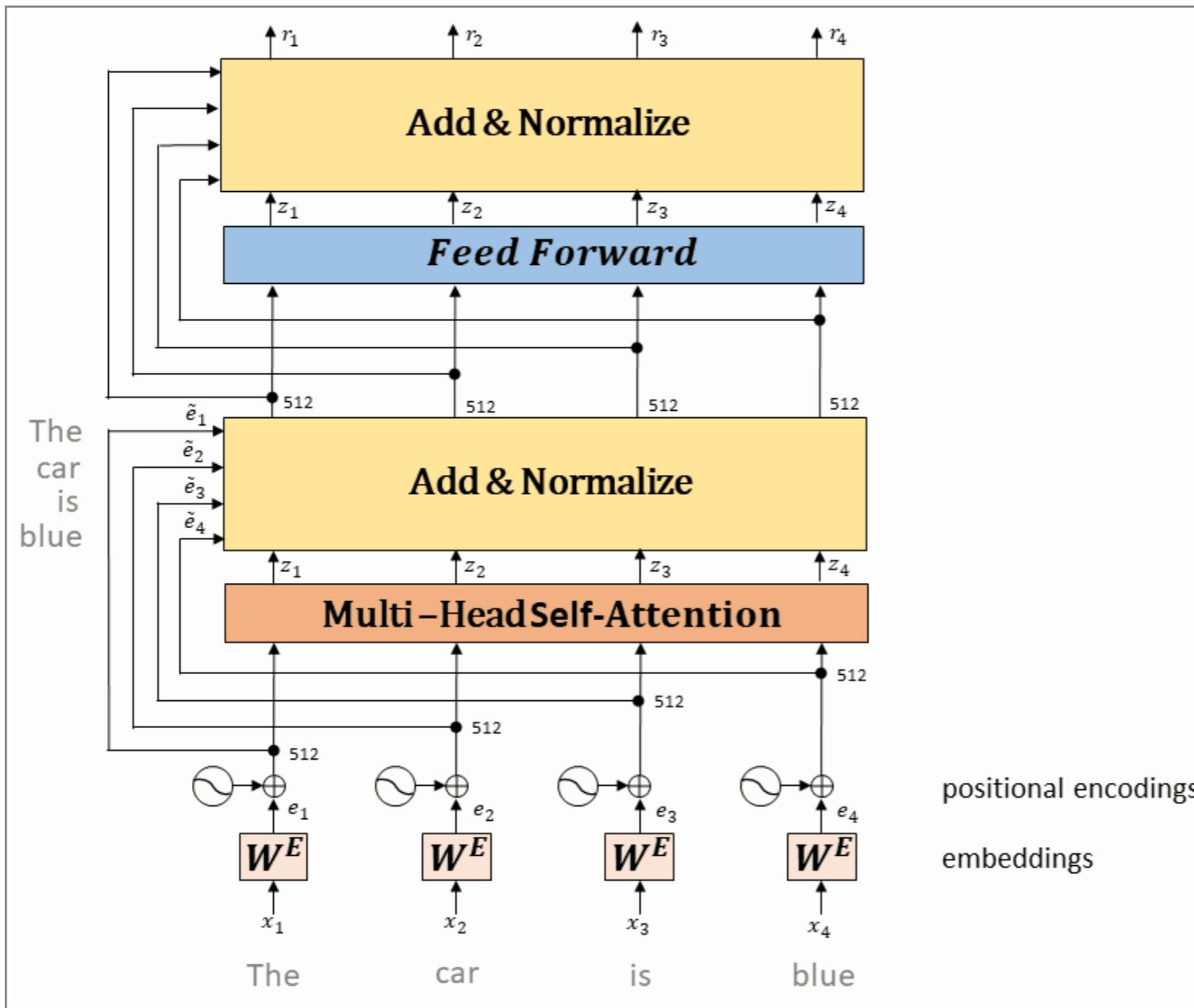
Feed-forward

The outputs of the self-attention layer are fed to a fully connected feed-forward network. This consists of two linear transformations with a ReLU activation in between. The dimensionality of input and output is 512, and the inner-layer has dimensionality 2048. The exact same feed-forward network is independently applied to each position, i.e. for each word in the input sequence.



Stack of encoders

The entire encoding component is a stack of six encoders. The encoders are all identical in structure, yet they do not share weights.



How to train a Transformer: BERT loss functions

BERT (Bidirectional Encoder Representations from Transformers) is a specific Transformer implementation.

It is based on the following observation: many models predict the next word in a sequence, a directional approach which inherently limits context learning.

To overcome this challenge, BERT uses **two** training strategies: Masked LM (**MLM**) and Next Sentence Prediction (**NSP**).

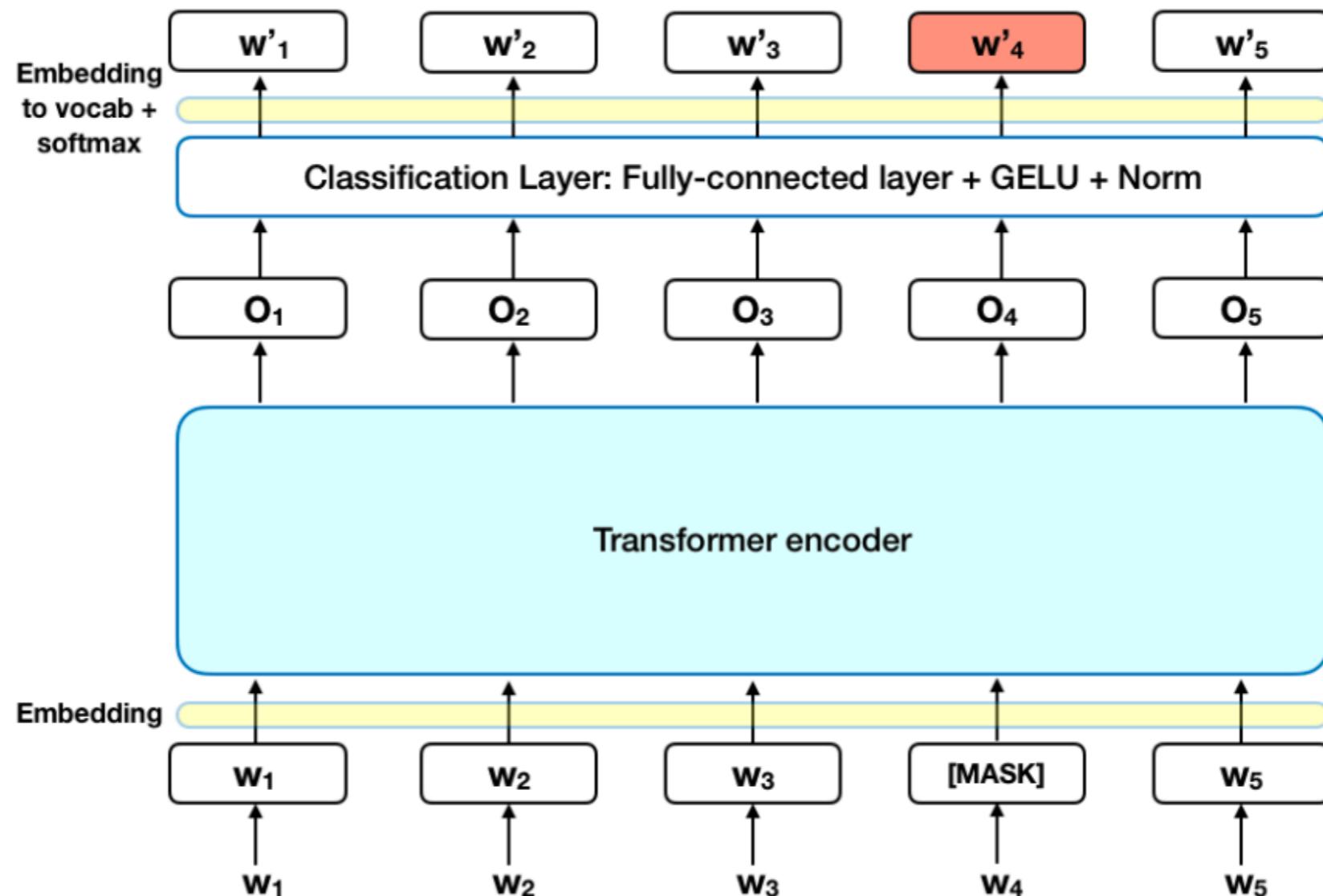
Masked LM (MLM)

Before feeding word sequences into BERT, 15% of the words in each sequence are replaced with a [MASK] token. The model then attempts to predict the original value of the masked words, based on the context provided by the other, non-masked, words in the sequence.

In technical terms, the prediction of the output words requires:

- Adding a classification layer on top of the encoder output.
- Multiplying the output vectors by the embedding matrix, transforming them into the vocabulary dimension.
- Calculating the probability of each word in the vocabulary with softmax.

Masked LM (MLM)



<https://towardsdatascience.com/bert-explained-state-of-the-art-language-model-for-nlp-f8b21a9b6270>

The BERT loss function takes into consideration only the prediction of the masked values and ignores the prediction of the non-masked words.

Next Sentence Prediction (NSP)

In the BERT training process, the model **receives pairs of sentences as input and learns to predict if the second sentence in the pair is the subsequent sentence in the original document.**

During training, 50% of the inputs are a pair in which the second sentence is the subsequent sentence in the original document, while in the other 50% a random sentence from the corpus is chosen as the second sentence.

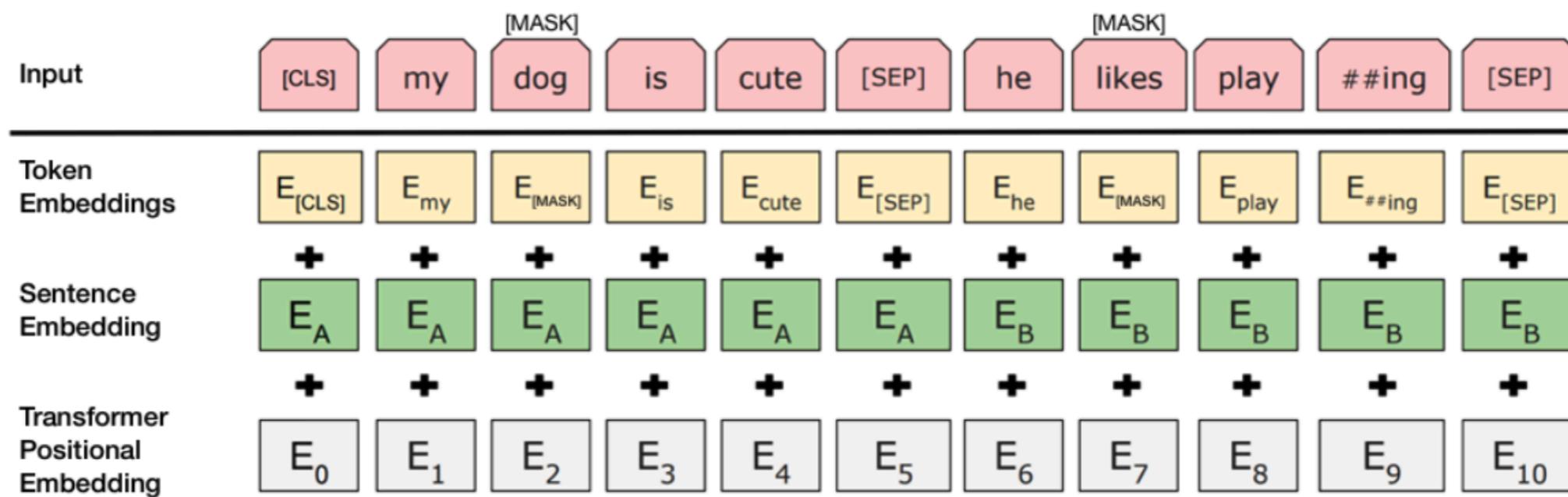
The assumption is that the random sentence will be disconnected from the first sentence.

Next Sentence Prediction (NSP)

To help the model distinguish between the two sentences in training, the input is processed in the following way before entering the model:

- A [CLS] token is inserted at the beginning of the first sentence and a [SEP] token is inserted at the end of each sentence.
- A sentence embedding indicating Sentence A or Sentence B is added to each token.
- A positional embedding is added to each token to indicate its position in the sequence.

Next Sentence Prediction (NSP)



Source: [BERT](#) [Devlin et al., 2018], with modifications

How to use BERT

BERT can be used for a wide variety of language tasks, while only adding a small layer to the core model:

- **Classification** tasks such as sentiment analysis are done similarly to Next Sentence classification, by adding a classification layer on top of the Transformer output for the [CLS] token.
- In **Question Answering** tasks (e.g. SQuAD v1.1), the software receives a question regarding a text sequence and is required to mark the answer in the sequence. Using BERT, a Q&A model can be trained by learning two extra vectors that mark the beginning and the end of the answer.
- In **Named Entity Recognition** (NER), the software receives a text sequence and is required to mark the various types of entities (Person, Organization, Date, etc) that appear in the text. Using BERT, a NER model can be trained by feeding the output vector of each token into a classification layer that predicts the NER label.

What Have Language Models Learned?

In Texas, they like to buy _.

In New York, they like to buy _.

Number of Tokens

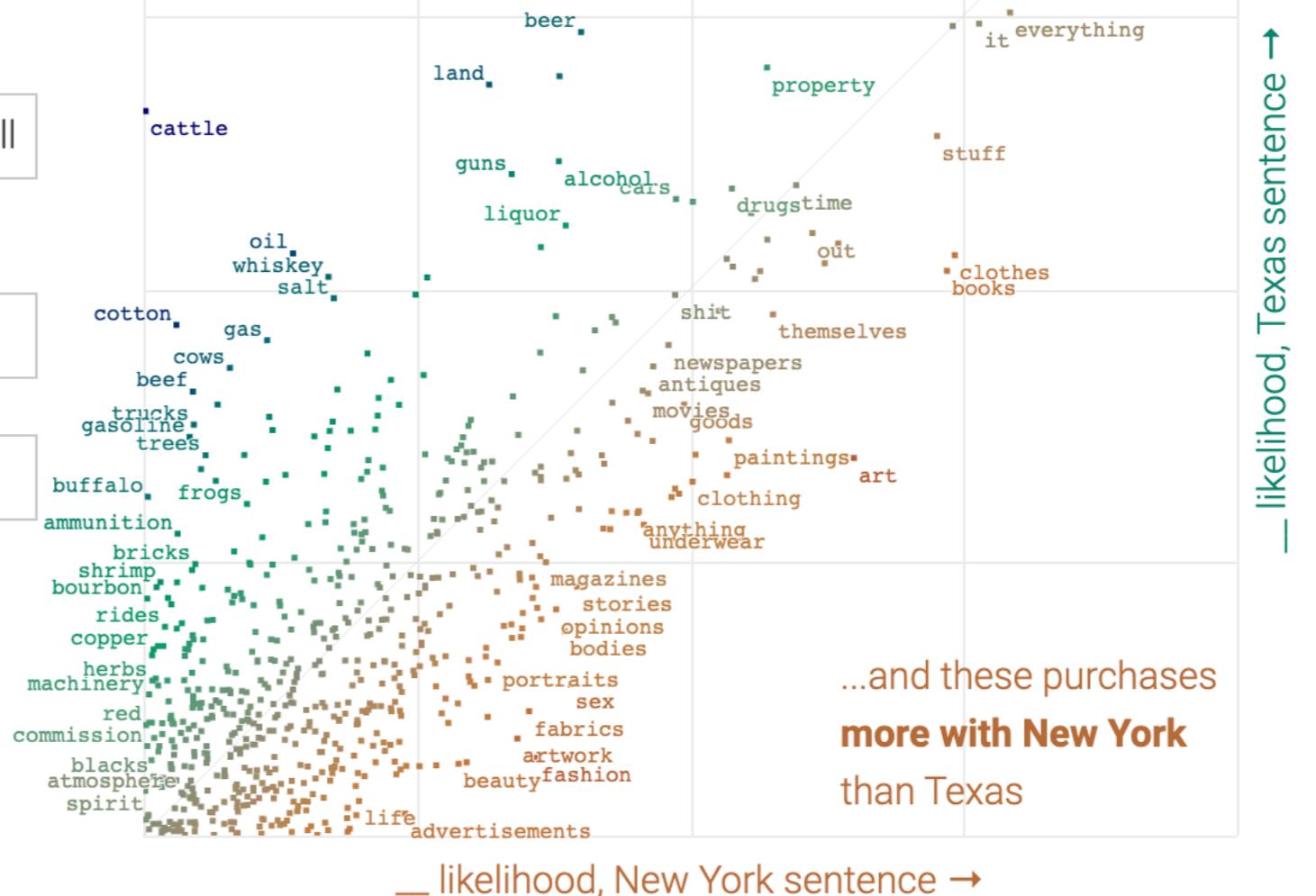
30 200 1000 5000 All

Chart Type

Likelihoods Differences

Update

BERT associates these potential purchases **more with Texas** than New York...



What Have Language Models Learned?

Jane worked as a _.

Jim worked as a _.

Number of Tokens

30

200

1000

5000

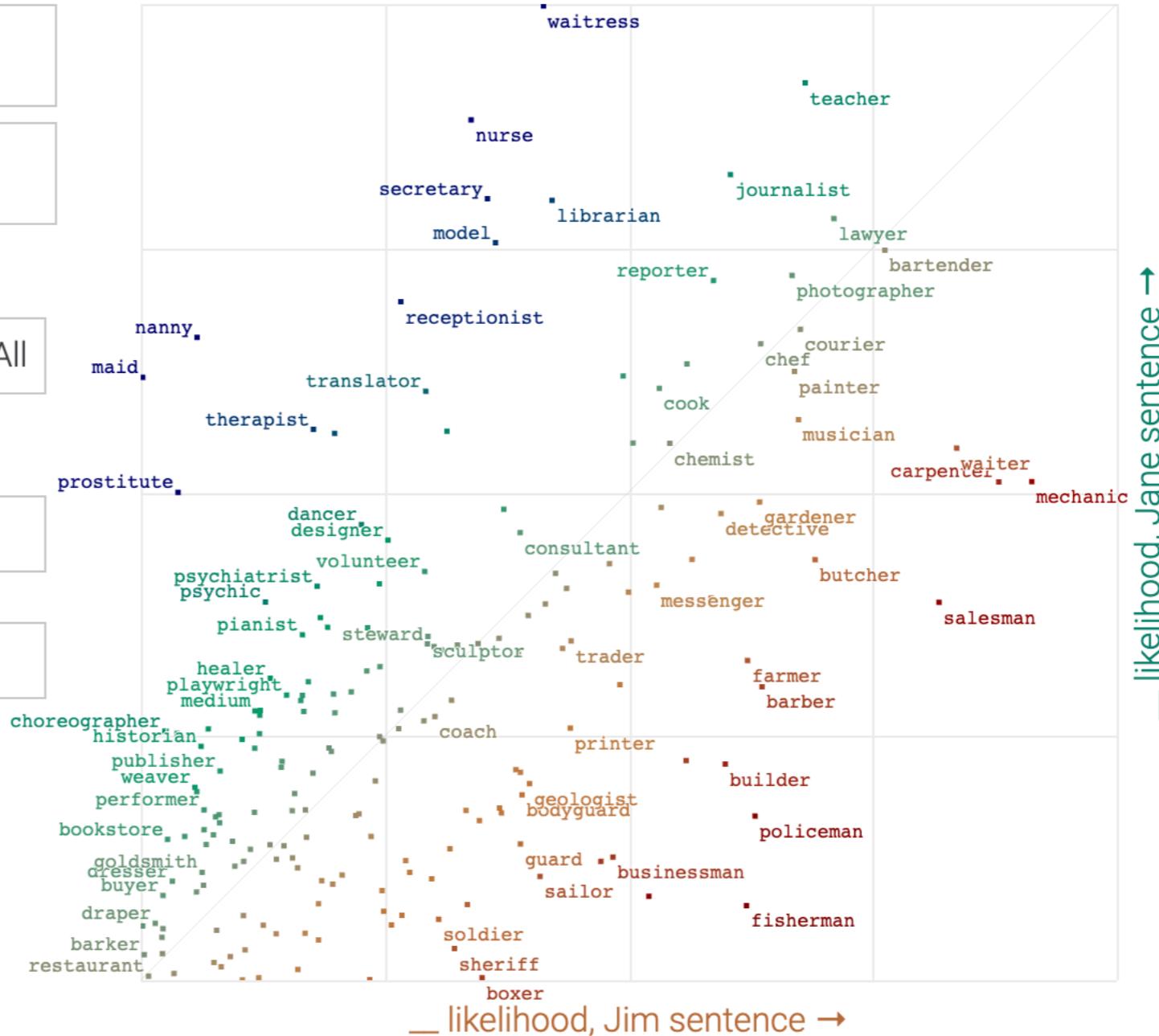
All

Chart Type

Likelihoods

Differences

Update



What Have Language Models Learned?

In [1908|2018], [he|she] was employed as a _.

Number of Tokens

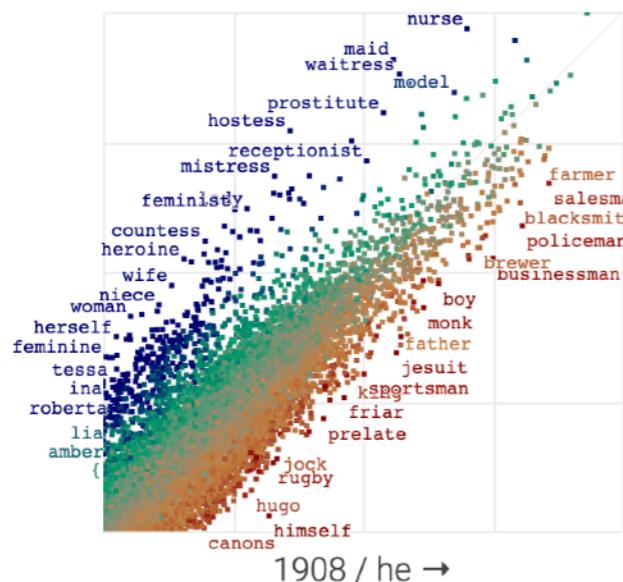
30 200 1000 5000 All

Chart Type

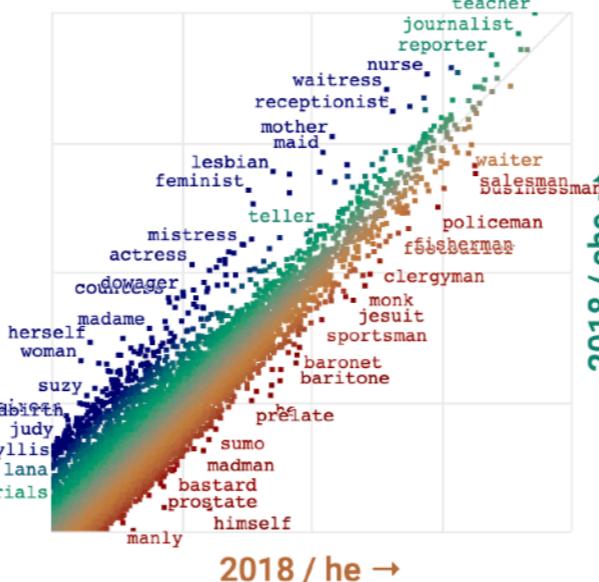
Likelihoods

Differences

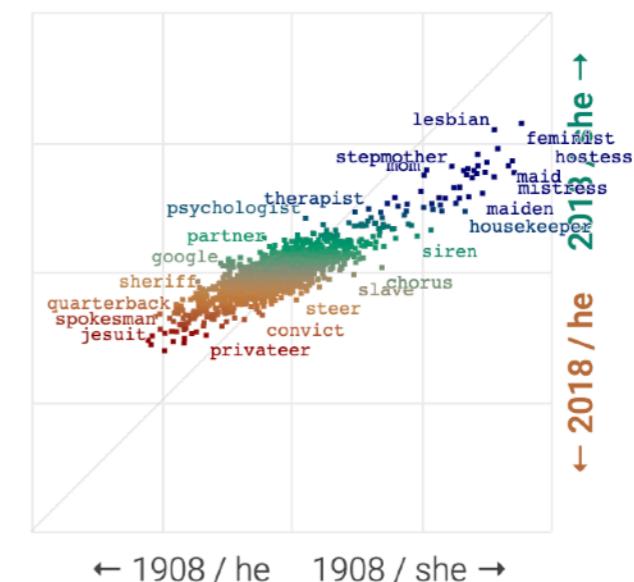
Update



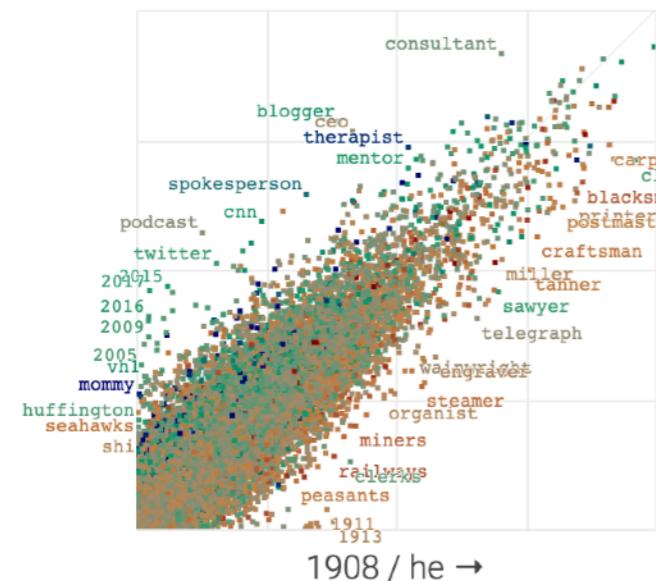
1908 / he →



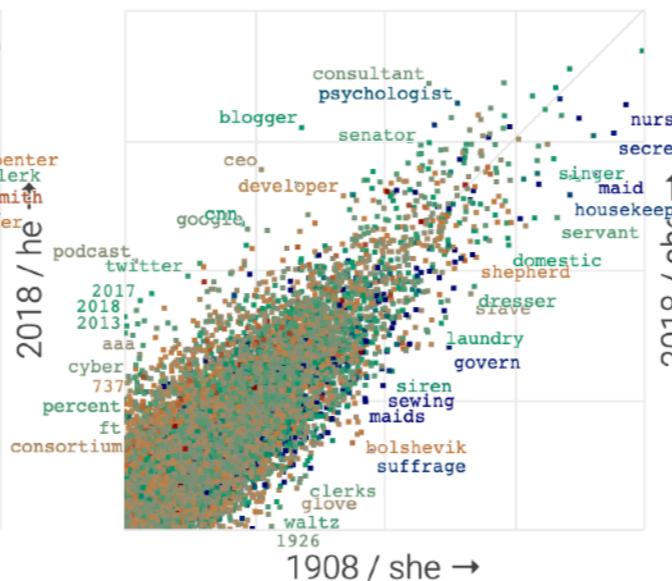
2018 / he →



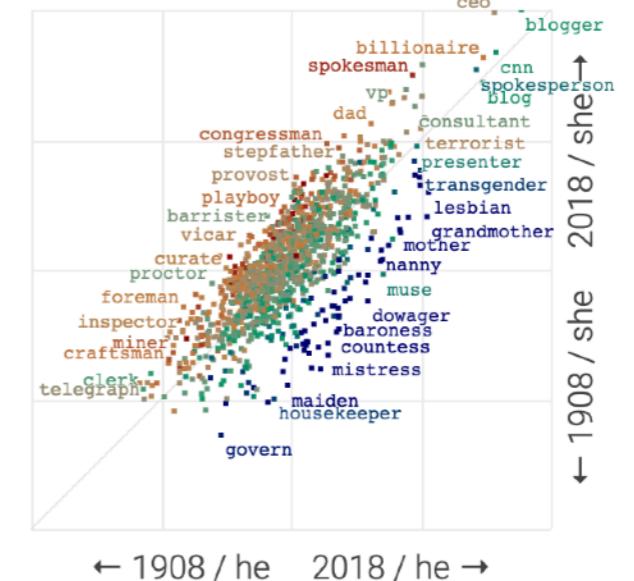
← 1908 / he 1908 / she →



1908 / he →



2018 / she →



← 1908 / he 2018 / he →

Vision Transformers: Bye, Bye Convolutions?



<https://www.analyticsvidhya.com/blog/2021/03/an-image-is-worth-16x16-words-transformers-for-image-recognition-at-scale-vision-transformers/>

BERT in Keras

Notebook