

Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - 0 x +

colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=BvWYUYhHYEE7

Update

+ Code + Text

Connect |



- which are then merged together so as to produce a compact representation

$$z_i = w_i * v_i,$$

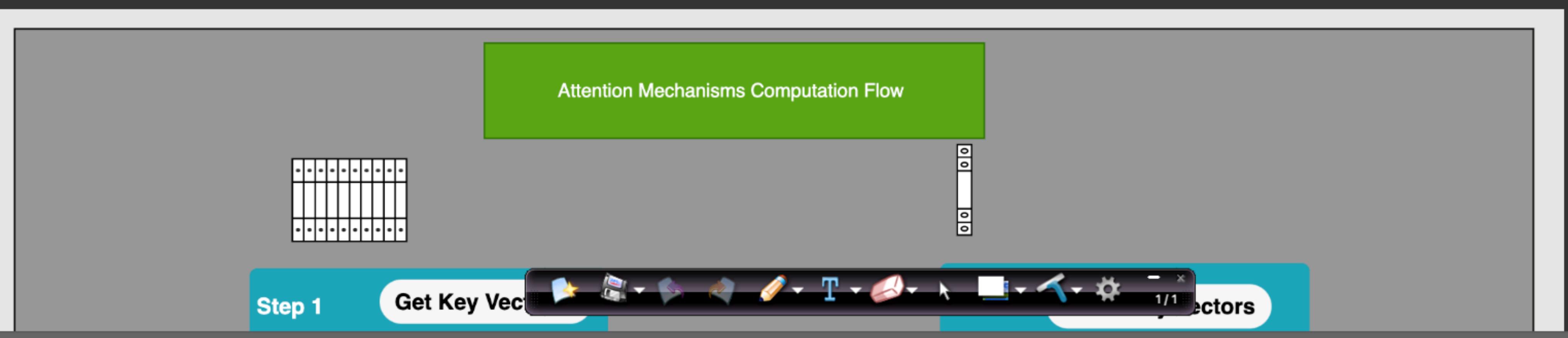
$$C = \sum_{i=1}^{d_k} (z_i),$$

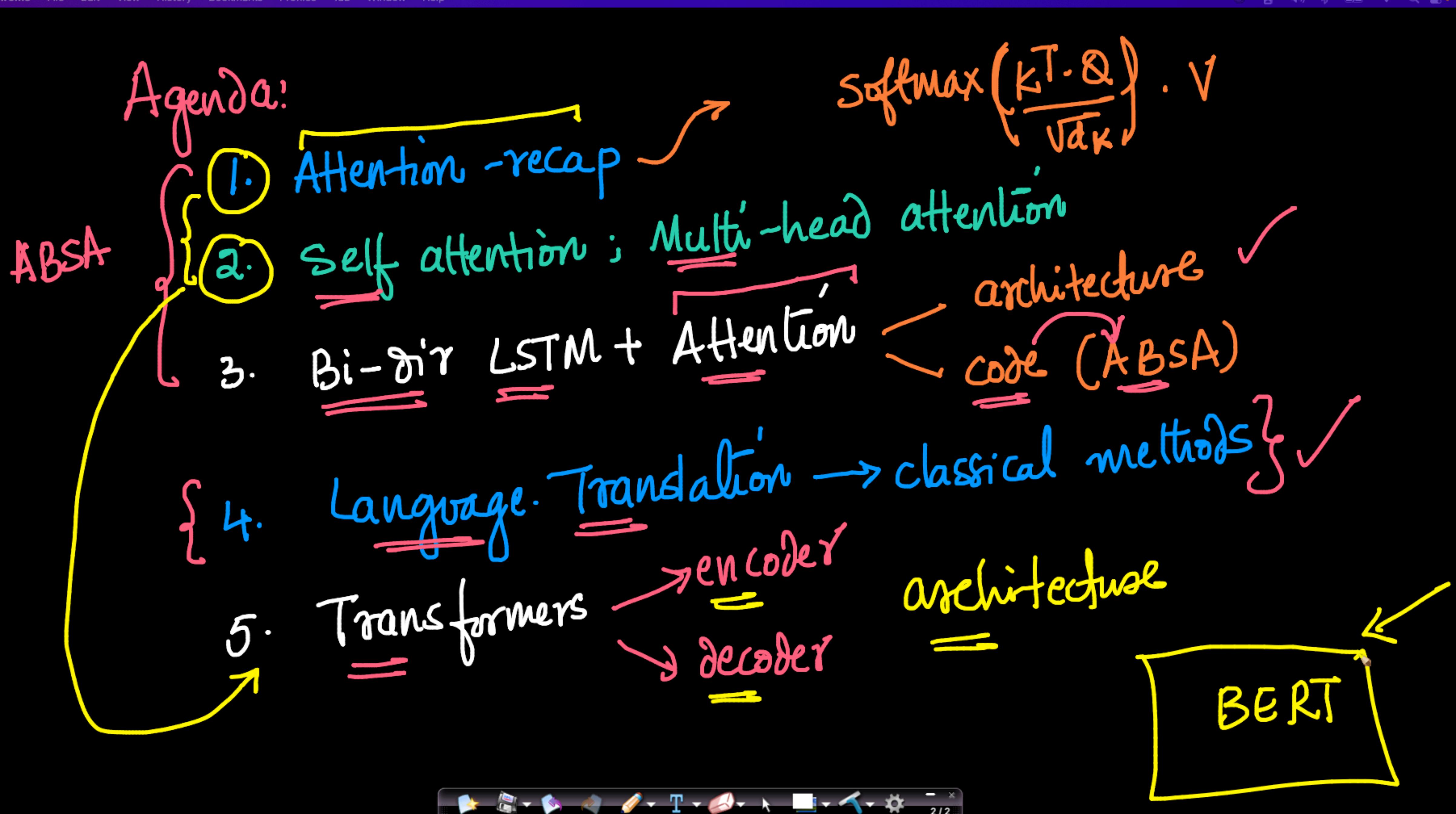
Attention: $\text{Softmax}\left(\frac{K^T Q}{\sqrt{d}}\right) \cdot V$

How to compute attention?

Ex - we have source sentence of 6 words and a target word , and we have to calculate attention score and context vector for target word.

- Assume source sentence food is delicious but price is high
- target words tasty
- word 3 is similar in context to target words .





colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=jPWKmLPRYECT

+ Code + Text

ATTENTION_FUNCTION

Equation $z(Q, K)$

Location_Basedt $f(Q)$

$f(Q, K) = [e_0; e_1; ; ; ; e_{d_k}]$

Convolution_Based $e_j = \frac{1}{l} \sum_{i=j-l}^k e_{i,j}$

$activation(W \cdot [k_i; k_{i+1}; ; ; k_{i+l}] + b)$

$z = \text{Softmax} \left(\frac{K^T \cdot Q}{\sqrt{d_K}} \right) \cdot \sqrt{d_K}$

sent \rightarrow *aspects* \rightarrow $= K$

$d_K = 5$ dim of each word vector

What are different type of Attention Mechanisms?

Attention Mechanisms can be broadly categorized into four category:

1. Number of Sequences.
2. Number of Abstraction.
3. Number of Position.
4. Number of Representation.

ATTENTION MECHANISMS

+ Code + Text

Connect |  

What are attention mechanisms based on Number of Sequences?

{x}

1. Distinctive Attention Mechanisms

- In Distinctive Attention Mechanisms, Both key and Query text Sequence is different.
- KEY VECTOR == VALUE VECTOR!= QUERY VECTOR
- How to Compute Distinctive AM?
 - Compute KEY, QUERY and VALUE vectors
 - Compute Energy function
 - Compute Attention Score

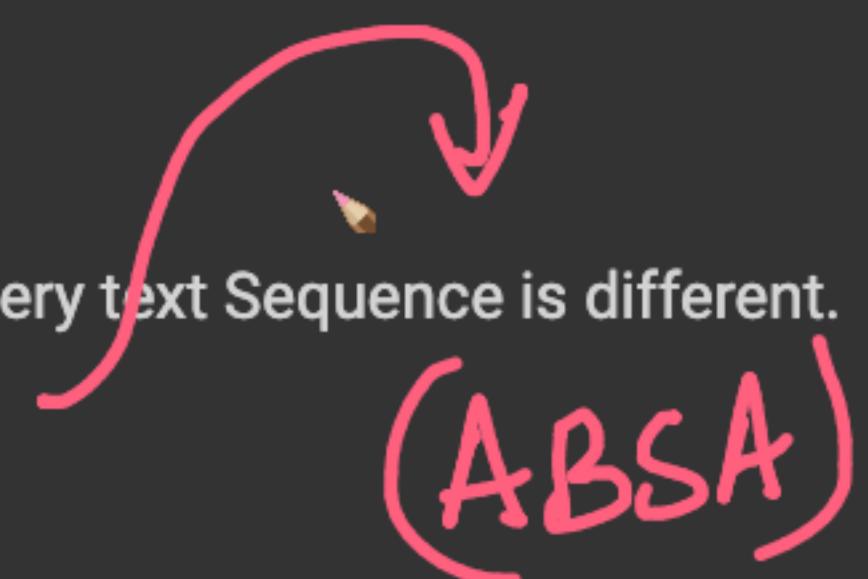


Fig Below - Distinctive Attention Mechanisms.

Distinctive Attention

Context



colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=jPWKmLPRYECT

+ Code + Text Connect |

• HOW TO COMPUTE SELF ATT:

- Compute KEY, QUERY and VALUE Vectors
- Compute Energy function
- Compute Attention Score

W_1 W_2 W_3

Self Attention

$Z = \text{softmax} \left(\frac{K^T \cdot Q}{\sqrt{d_K}} \right) \cdot V$

$K = Q = \sqrt{d_K}$ \Leftrightarrow Self attention

5 / 5

+ Code + Text

KEY = [h1, h2, h3...h]
QUERY=VALUE = KEY

What are attention mechanisms based on Number of Position?

1. Soft Attention Mechanisms

- In soft attention, we compute attention weight $w_1, w_2, w_3 \dots w_n$, where $\sum W = 1$.
- And, It uses a weighted average of all hidden states of the input sequence to build the context vector.

$$z_i = w_i * v_i,$$

$$C = \sum_{i=1}^{d_k} (z_i),$$

- The usage of the soft weighing method makes the neural network amenable to efficient learning through backpropagation,
- Pro:** the model is smooth and differentiable.
- Con:** expensive when the source input is large.

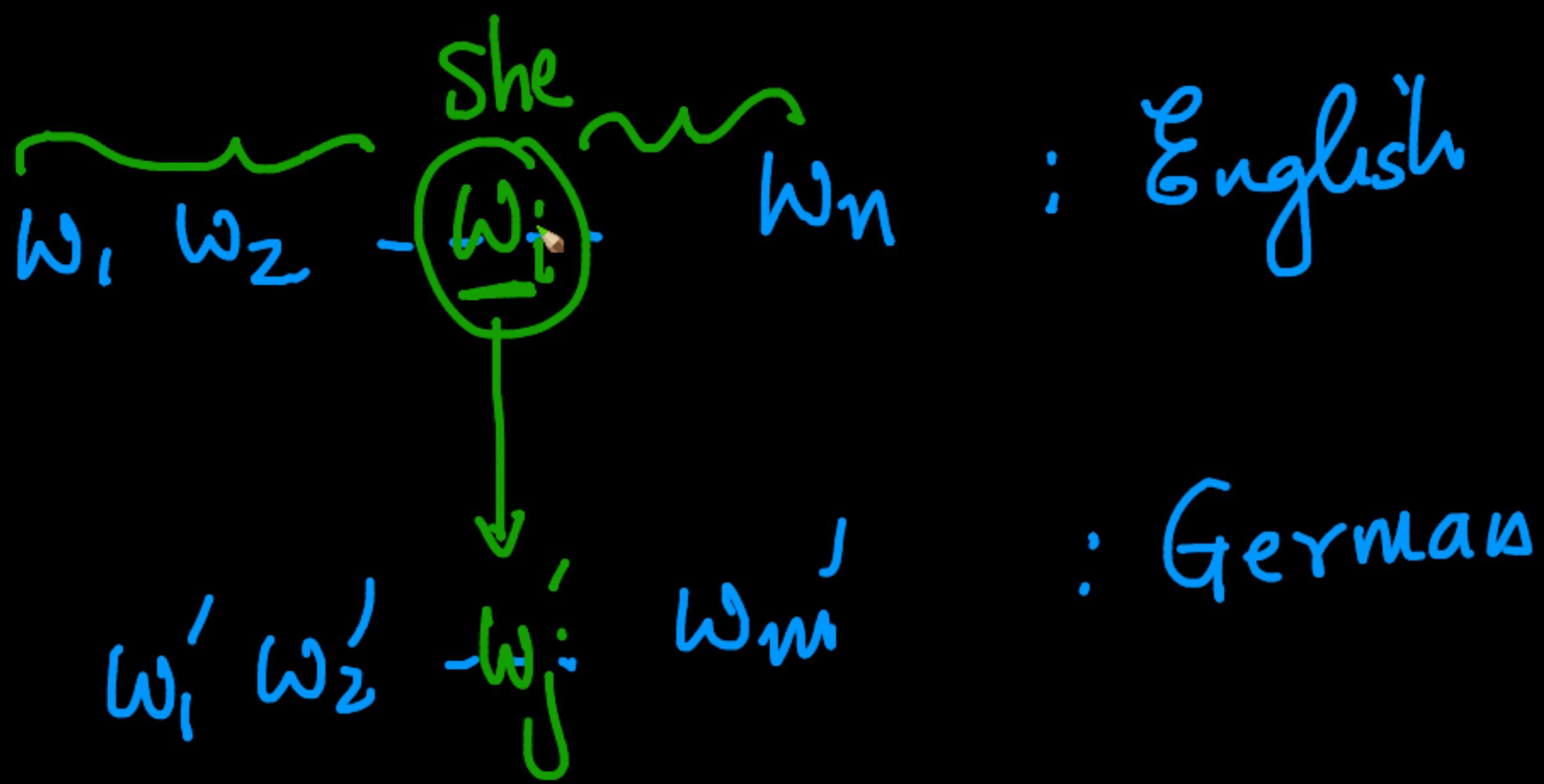
$$Z = \text{Softmax} \left(\frac{K^T \cdot Q}{\sqrt{d_k}} \right) \cdot V$$

\downarrow

$$w_1, w_2, w_3 \dots w_n$$

$\leftarrow \rightarrow$

Self-attention:



+ Code + Text

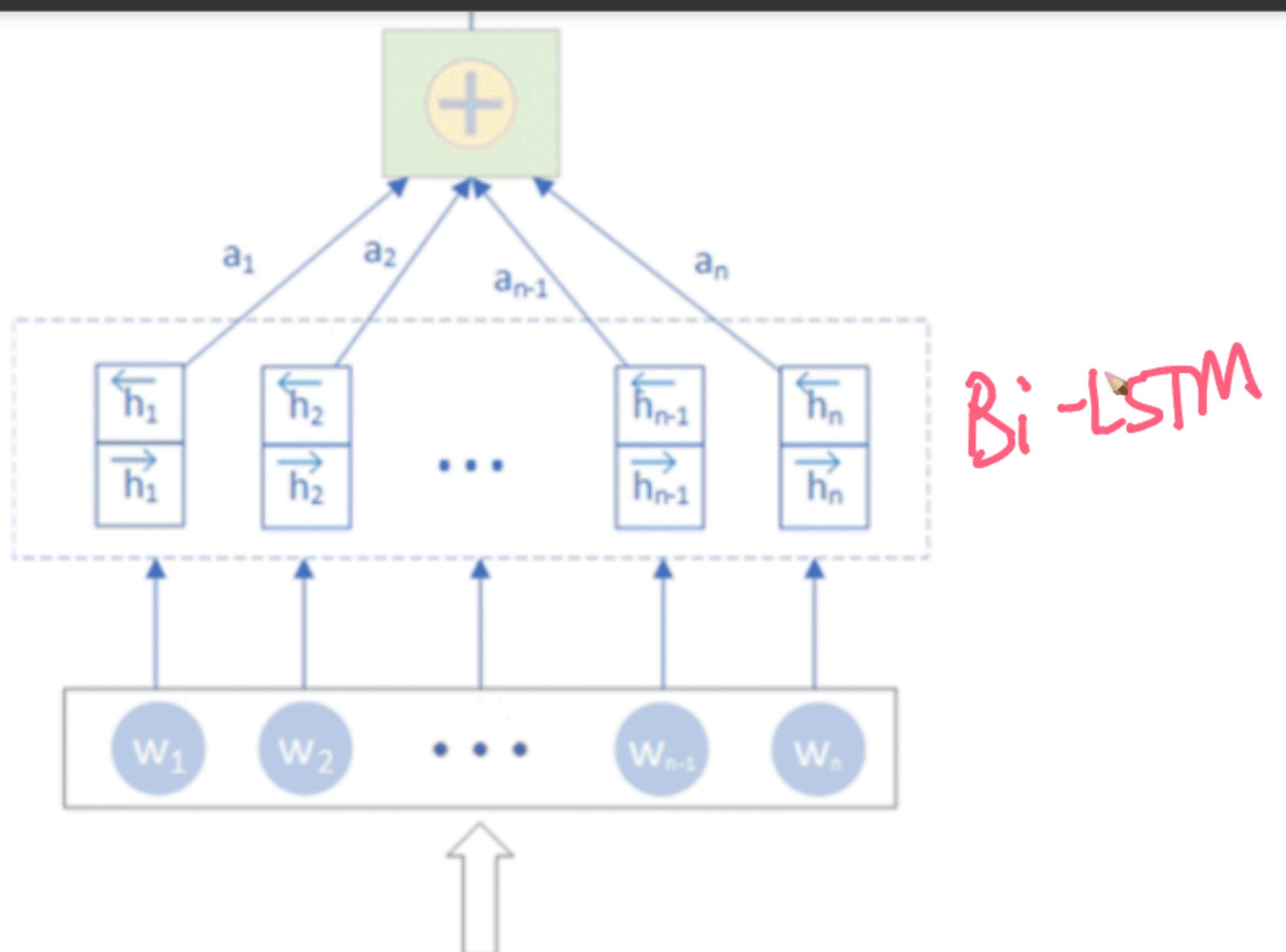
Word-Level Attention

Attention layer

Word encoder

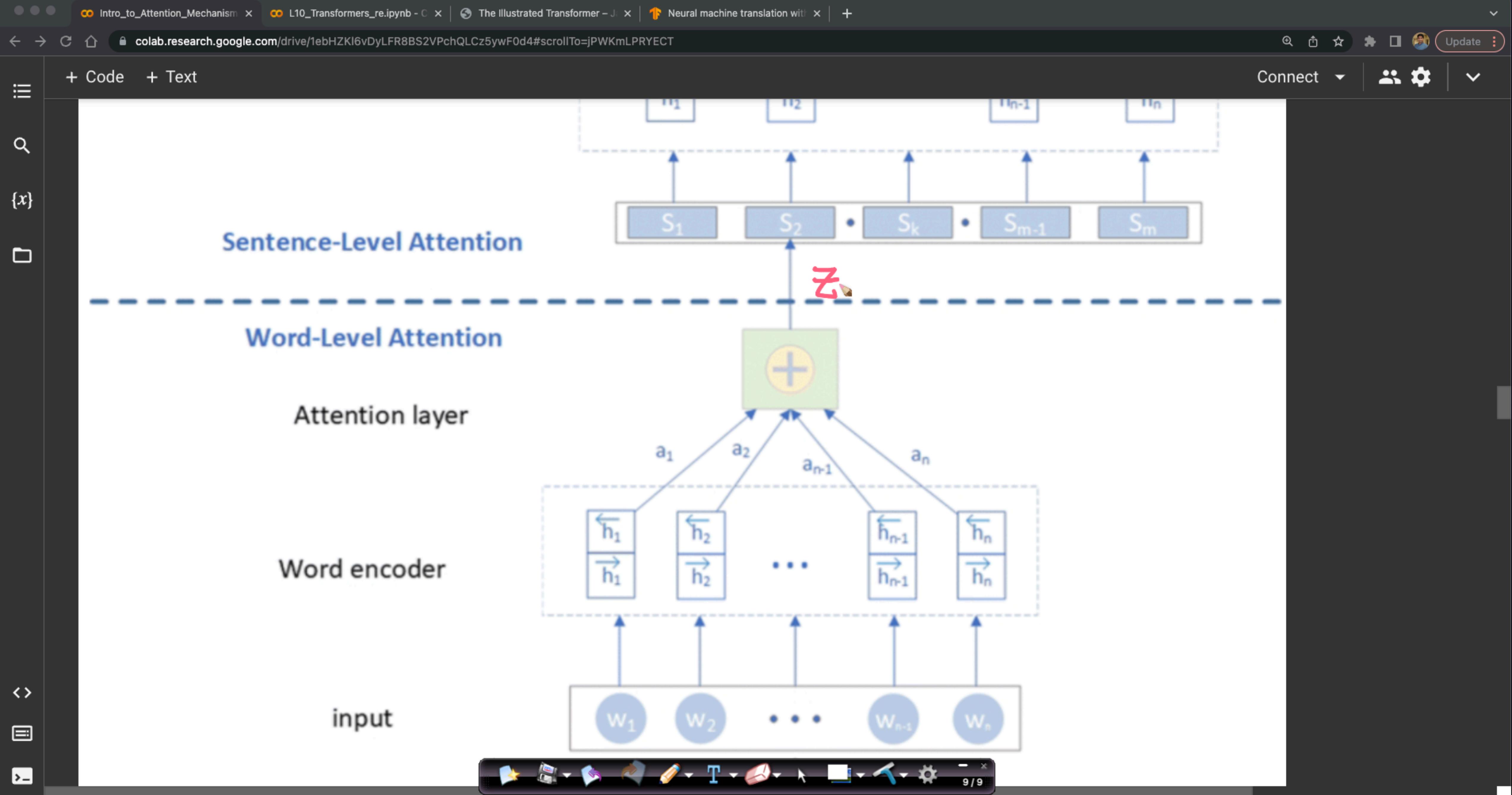
input

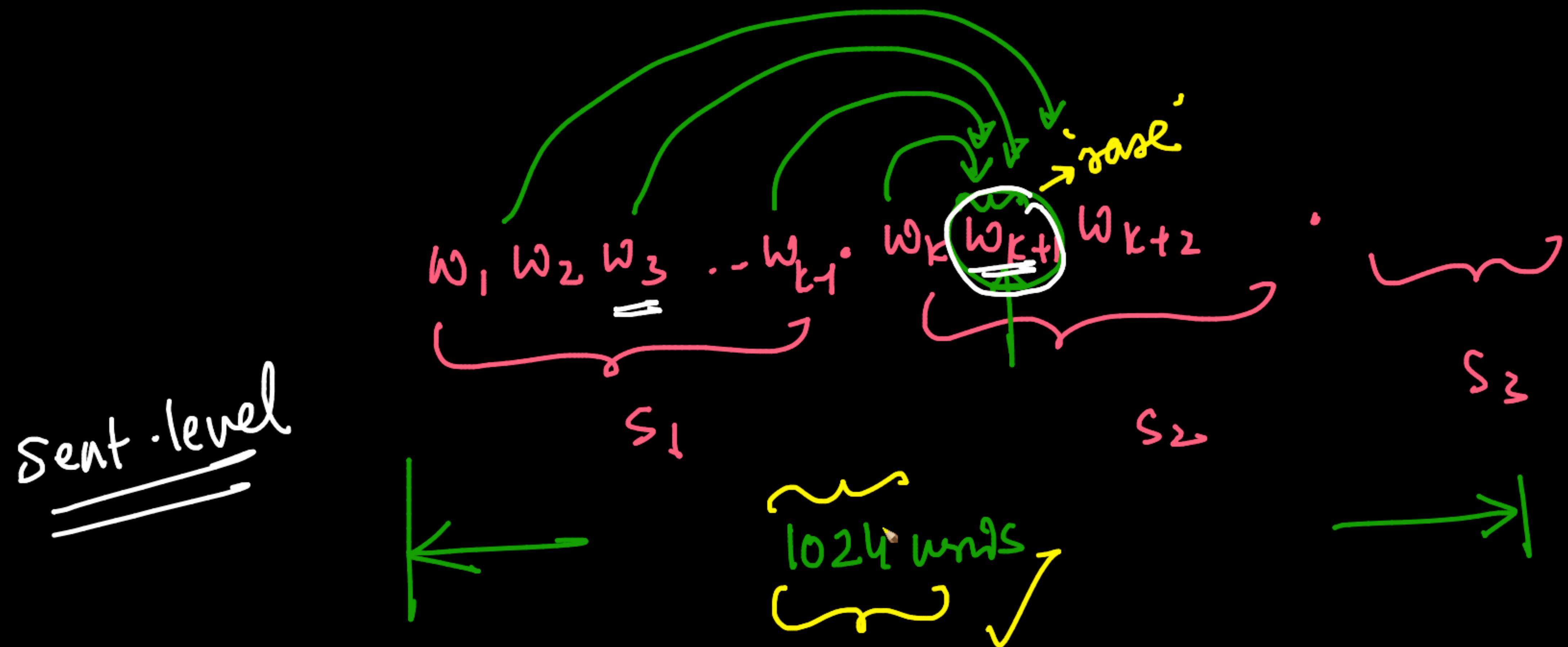
text

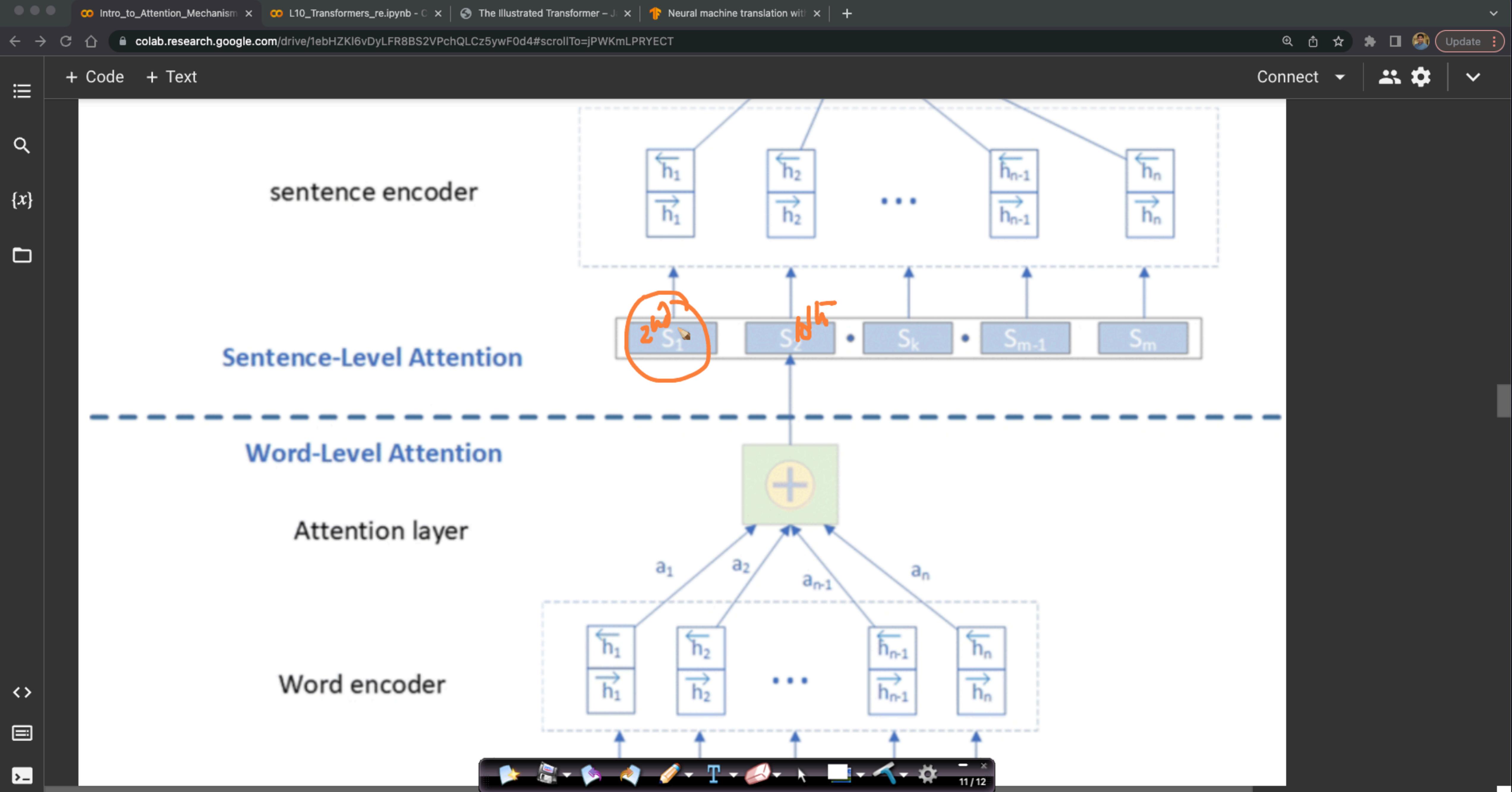


Bi-LSTM

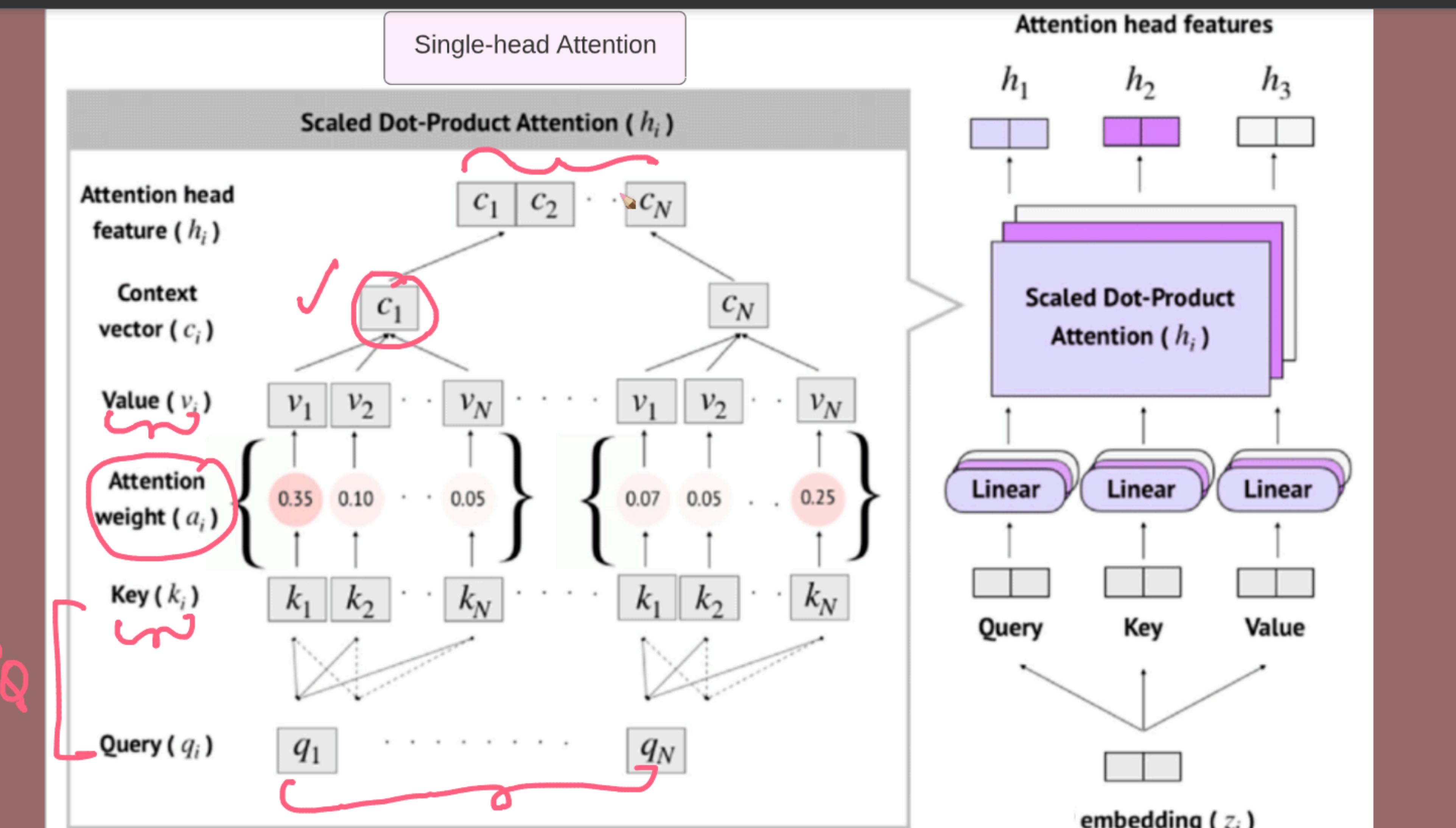
The hierarchical attention mechanism architecture. We connect the word-level and sentence-level features by attention mechanism to get a feature vector with more information and the Softmax function is used to normalize the result



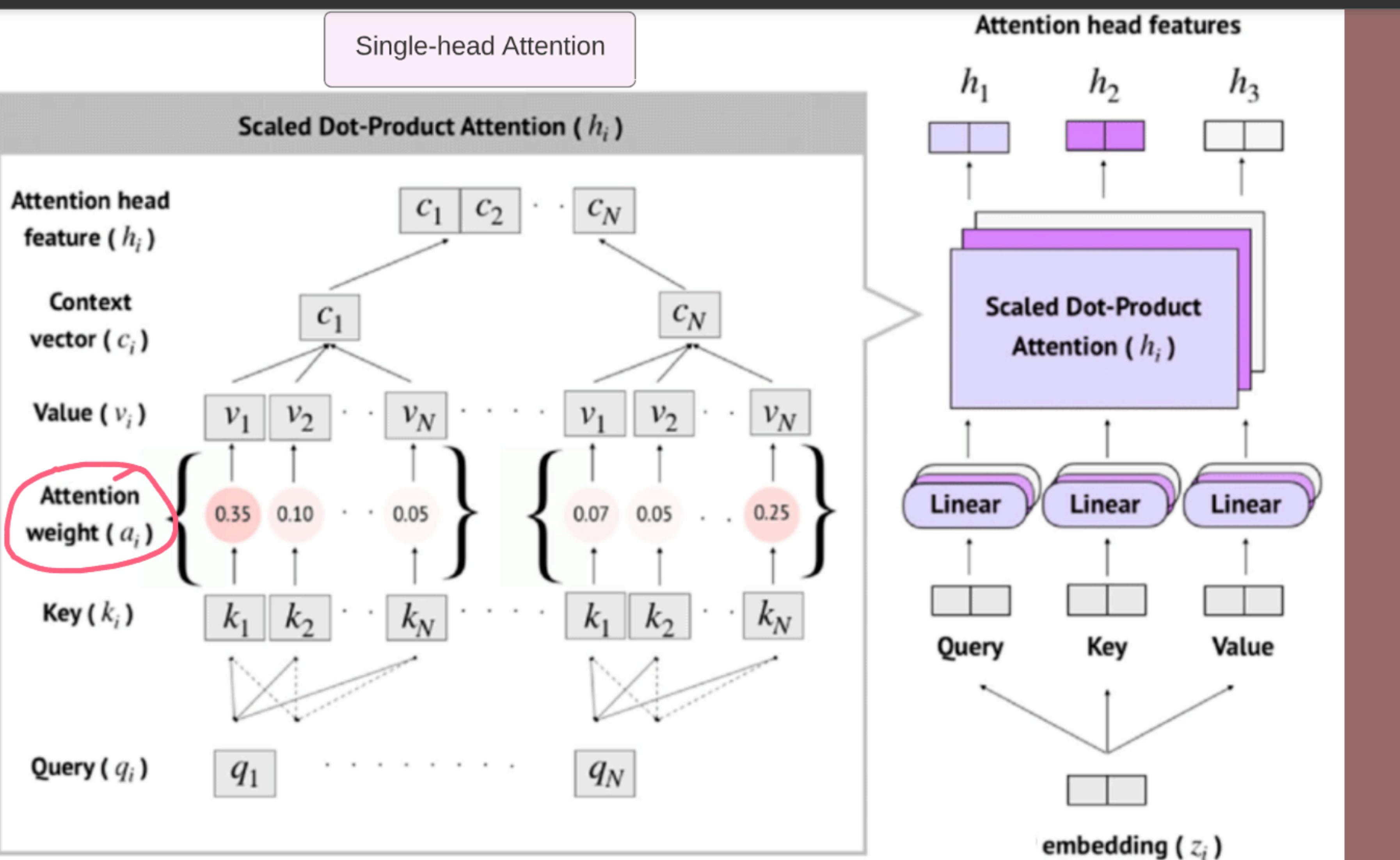


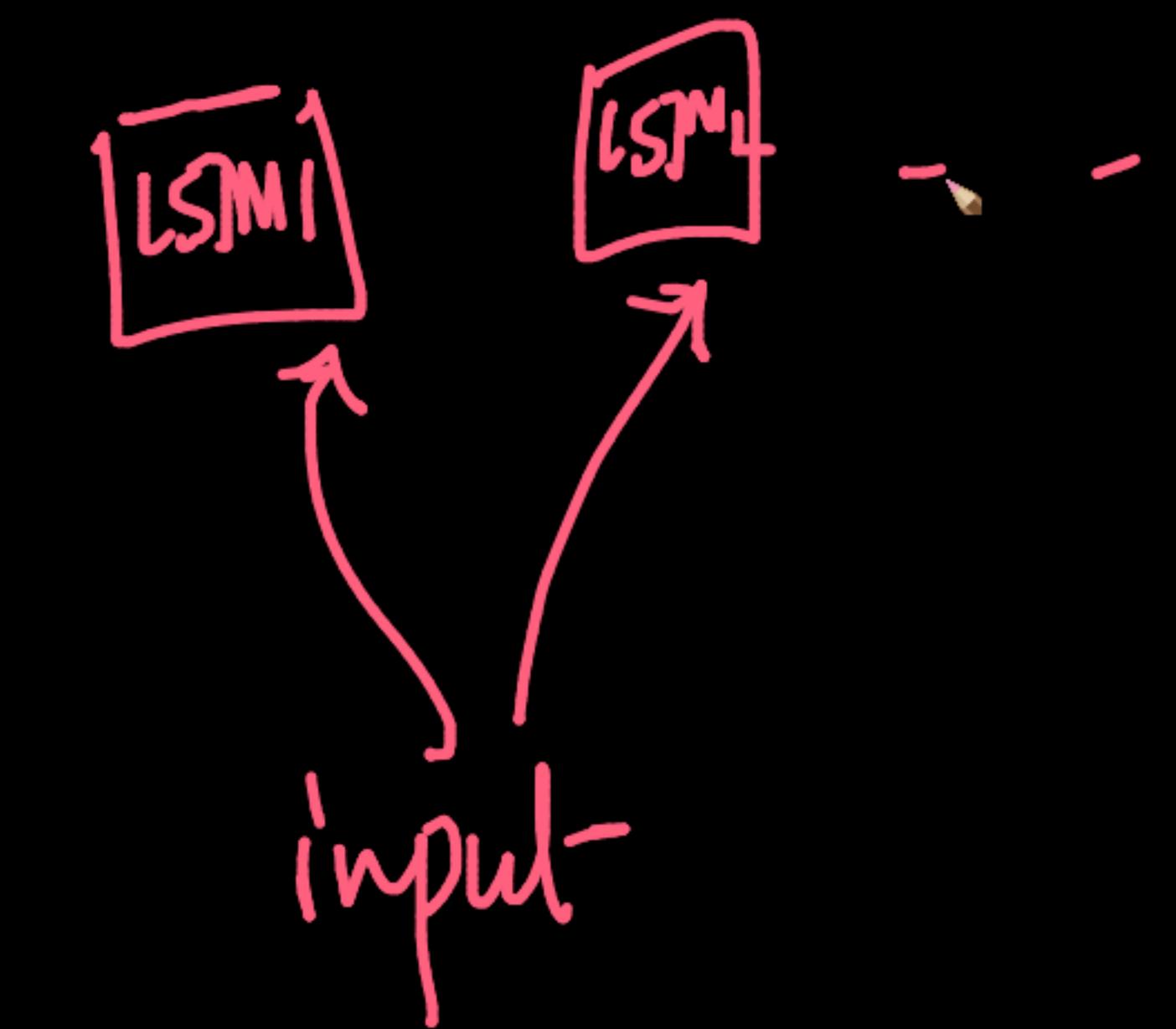


+ Code + Text

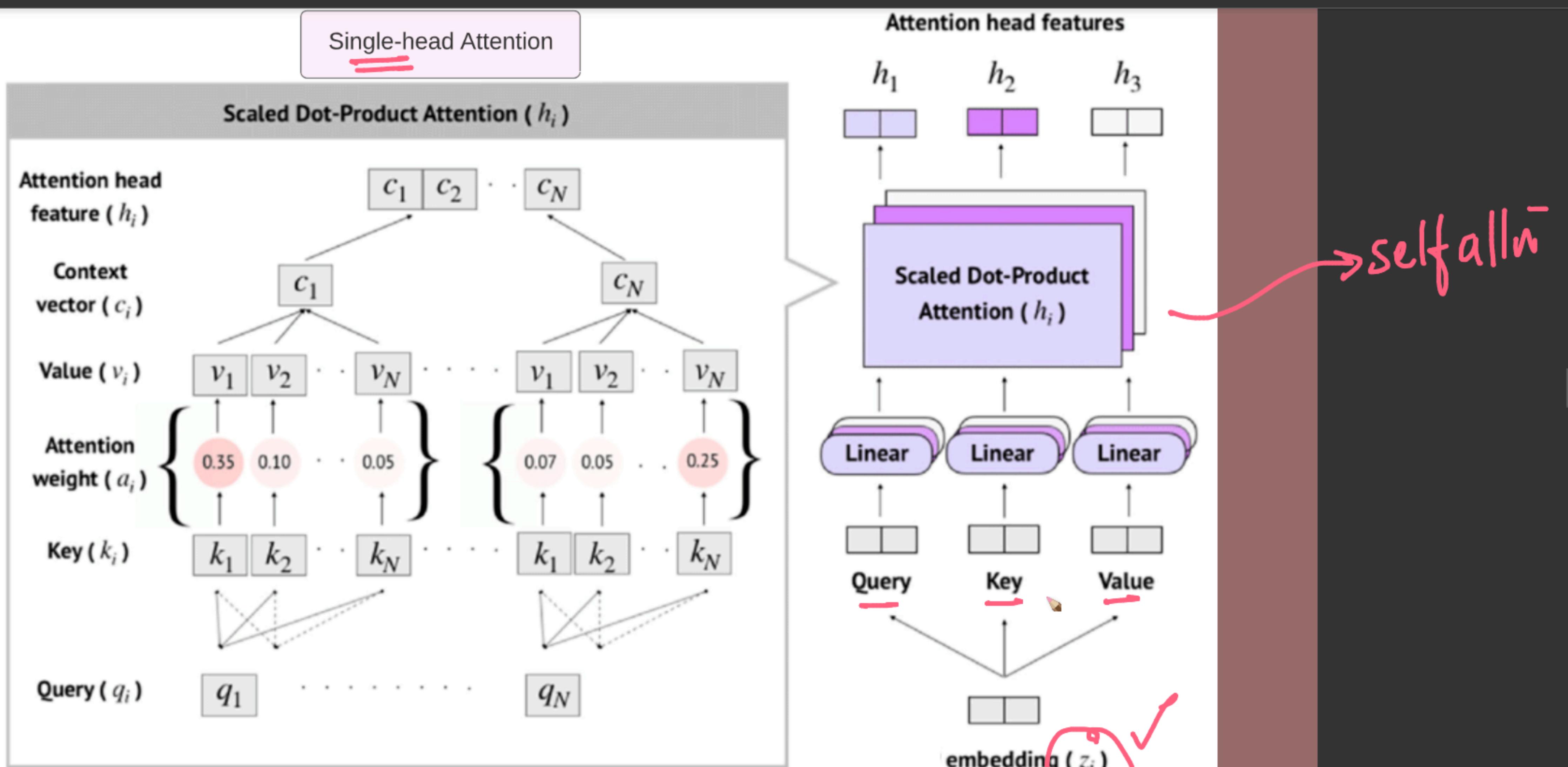


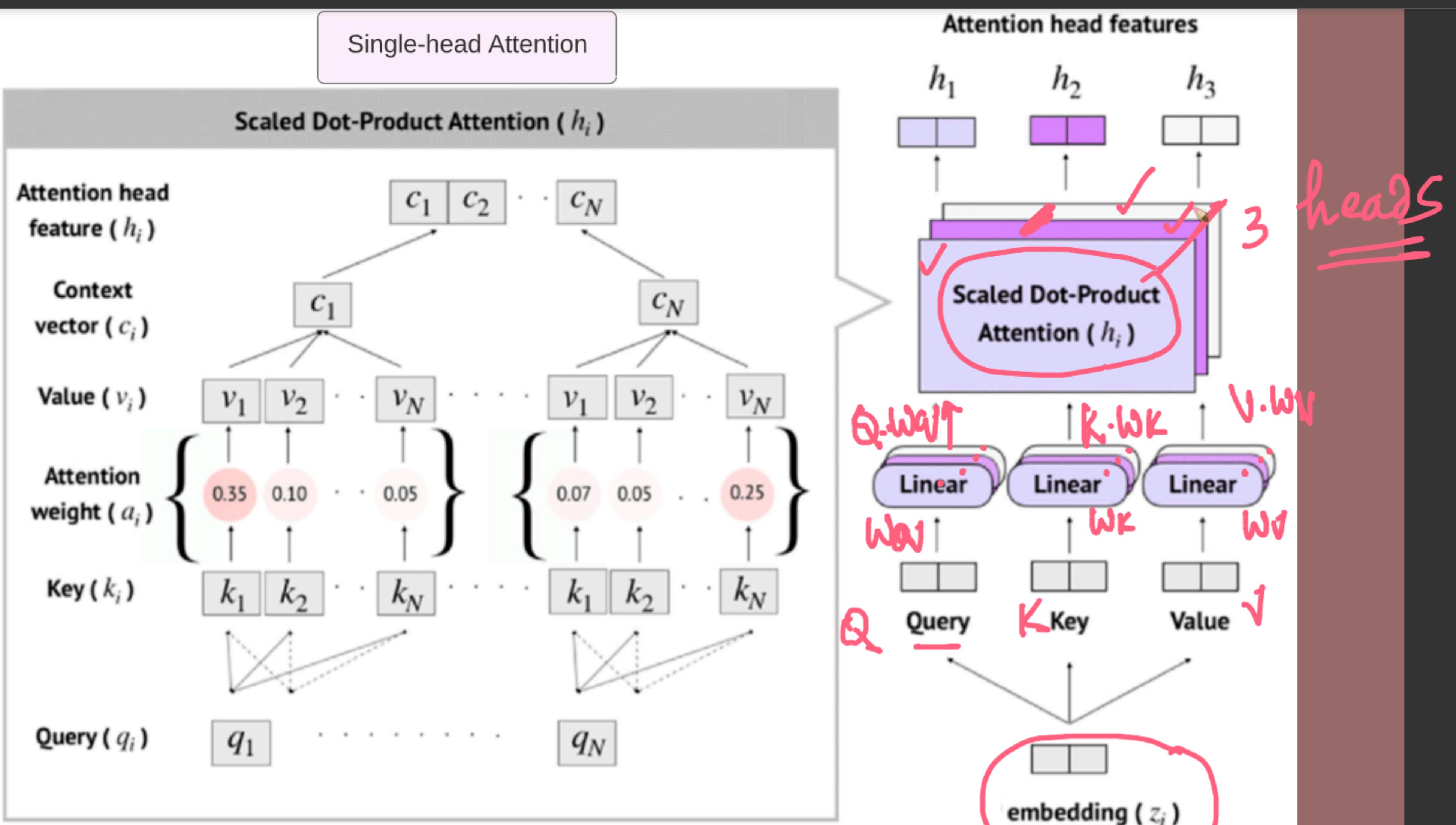
+ Code + Text



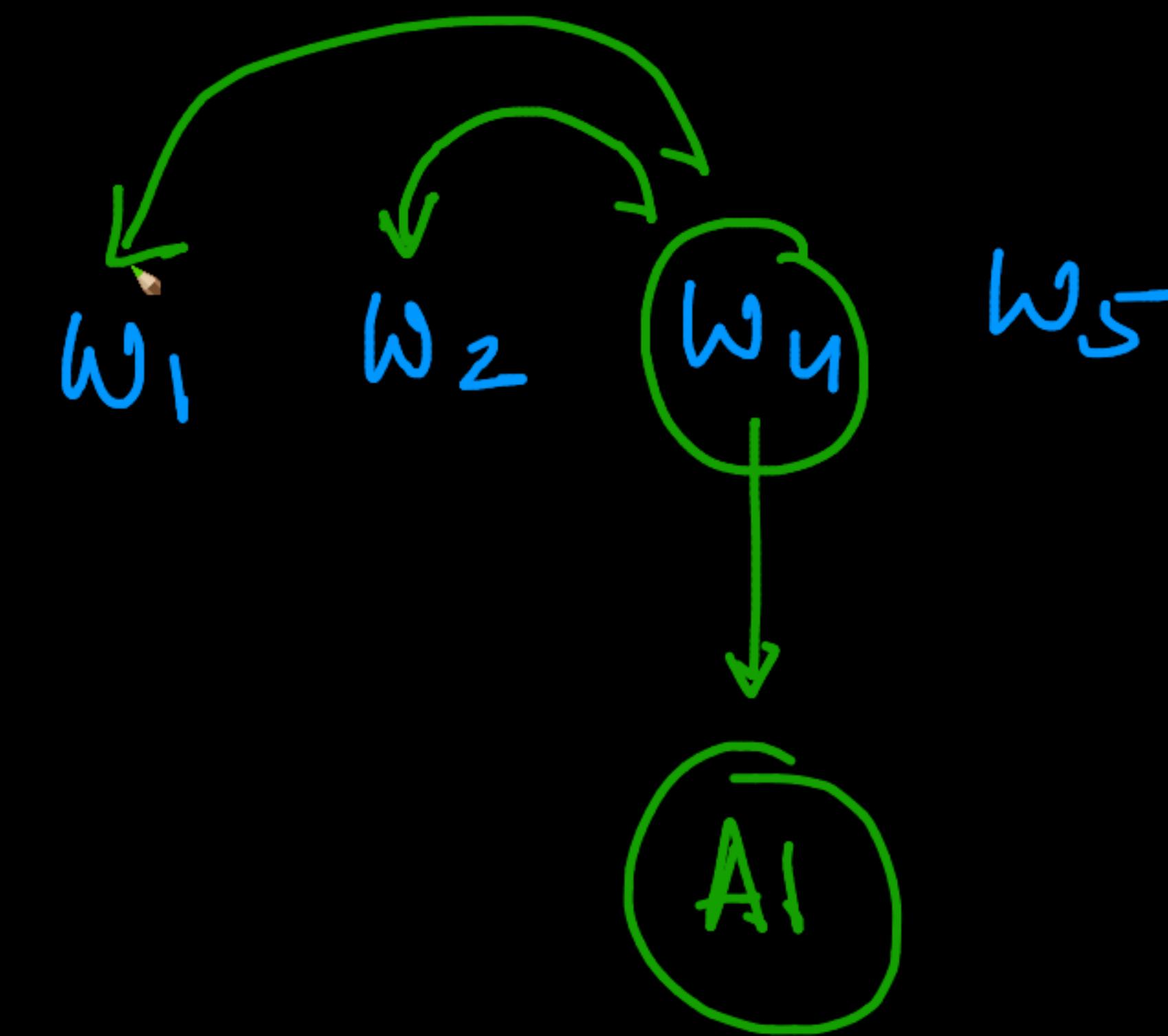


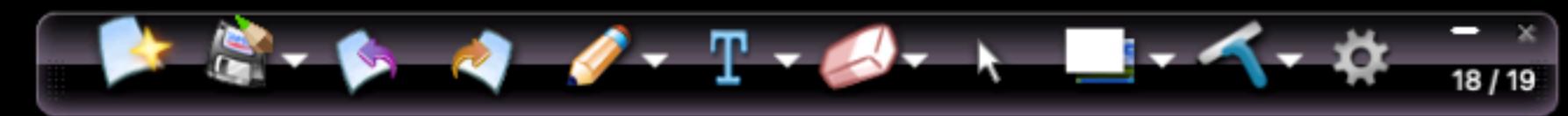
+ Code + Text

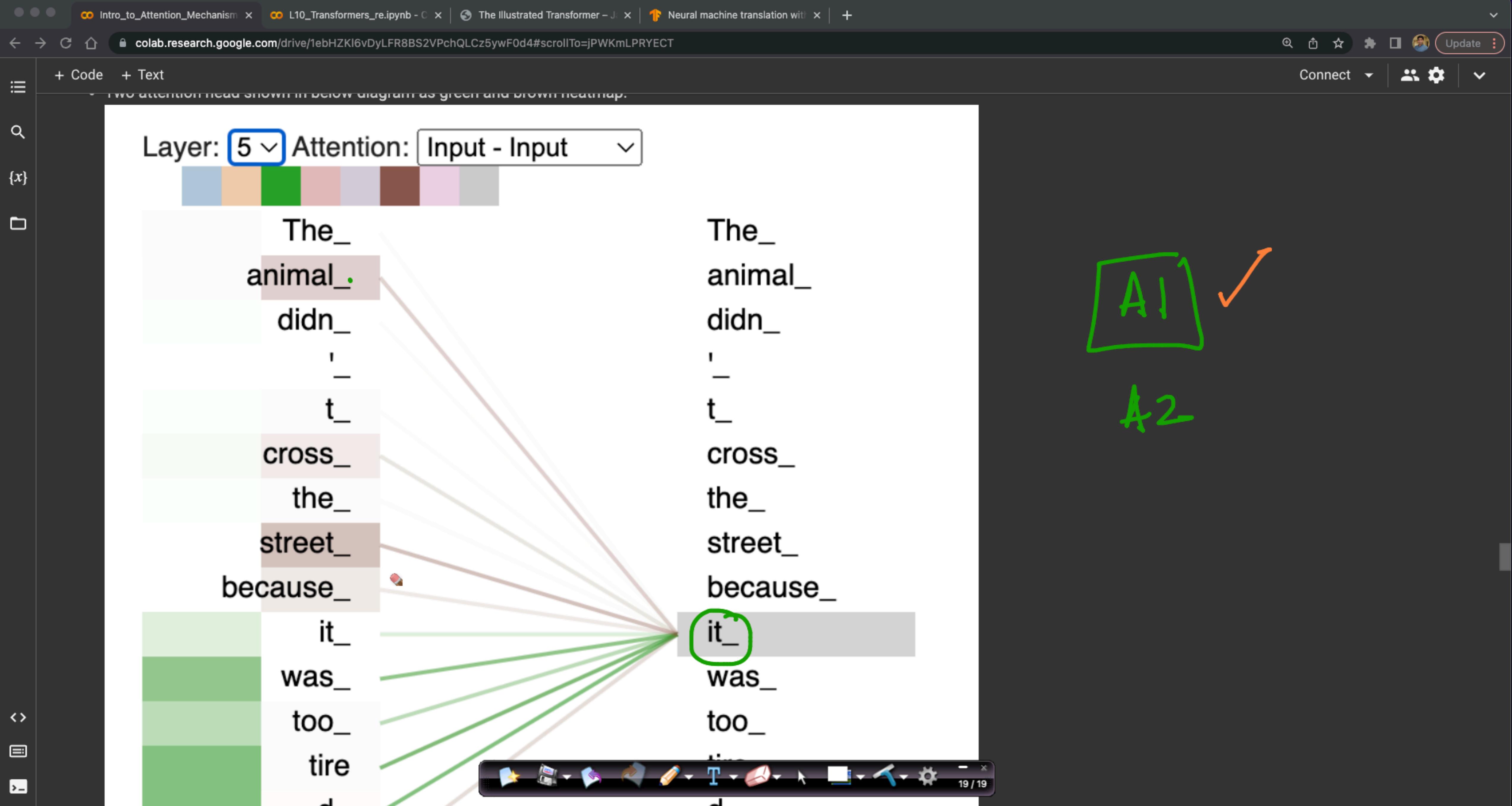


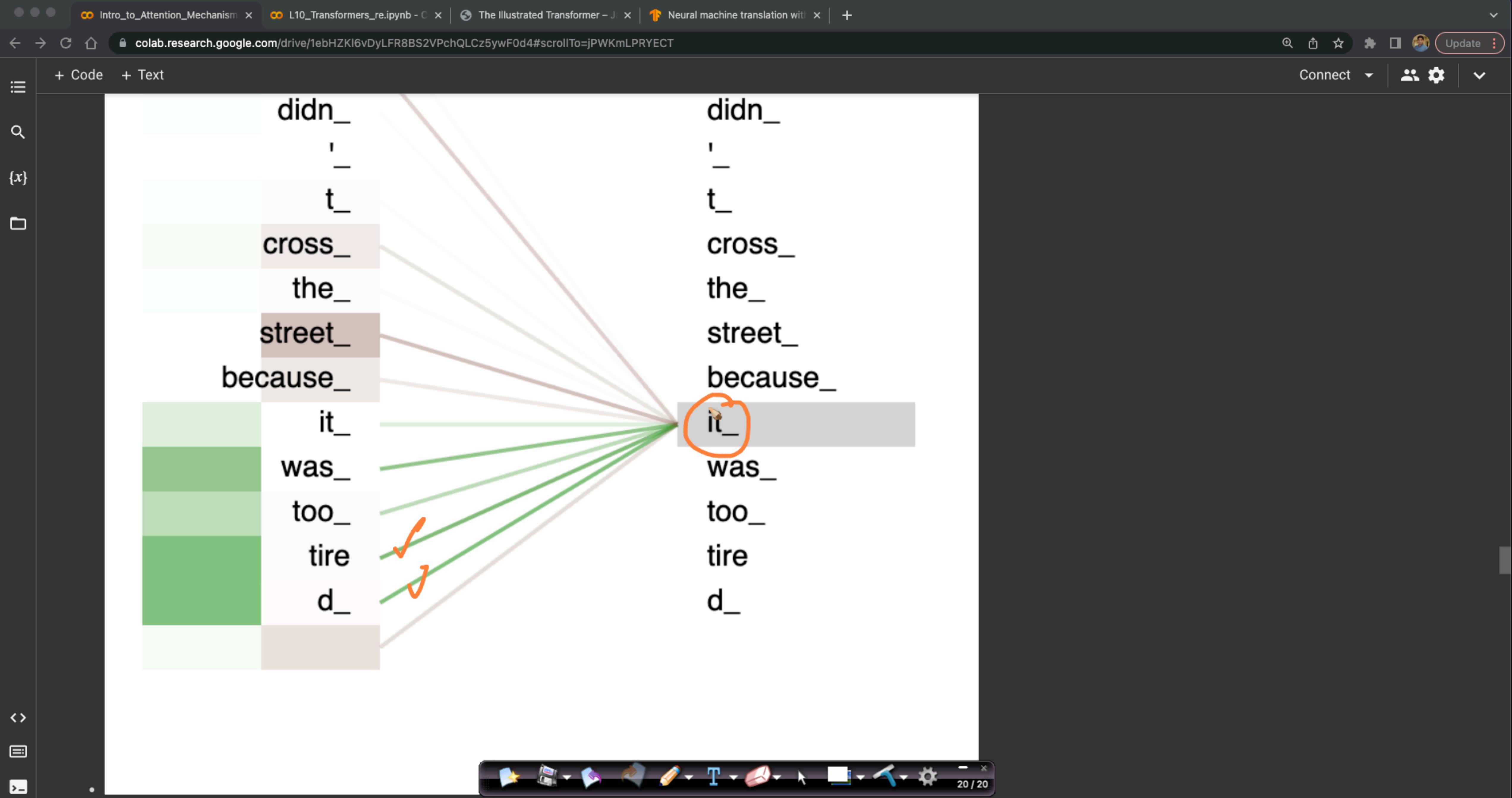


e.g:
~~A1~~









+ Code + Text

```
Y = [glove_vectors[i].tolist() for i in target]
Y = np.stack([pca.transform(i) for i in Y]).reshape(11, 5)
```

Lets Compute one of the popular Multi-headed attention in NLP.

- Multi-headed Self Attention

1. Step-1 Randomly initialized wq, wk and wv MATRIX.

- Each word in X is represented using a vector of size 5.
- Randomly Initialize 3 matrices, $w_q, w_v, w_k = [5 \times 9]$ matrix

```
▶ wk = np.ones((5, 9))
    vv = np.ones((5, 9))
    wq = np.ones((5, 9))
    display(HTML('<h3>WQ</h3>'))
    np_display(wq)
```



WQ

	0	1	2	3	4	5	6	7	8
0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
2	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
3	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
4	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0

Bi-directional LSTM ✓
Bi-directional LSTM + attention
Attention - Enc-Dec
↳ Transformers

Attention - Enc
↳ BERT

+ Code + Text

```
y = [glove_vectors[i].tolist() for i in target]
Y = np.stack([pca.transform(i) for i in Y]).reshape(11, 5)
```

{x} Lets Compute one of the popular Multi-headed attention in NLP.

- Multi-headed Self Attention

▼ 1. Step-1 Randomly initialized wq, wk and wv MATRIX.

- Each word in X is represented using a vector of size 5.
- Randomly Initialize 3 matrices, $w_q, w_v, w_k = [5 \times 9]$ matrix

```
▶ wk = np.ones((5, 9))
wv = np.ones((5, 9))
wq = np.ones((5, 9))
display(HTML('<h3>WQ</h3>'))
np_display(wq)
```



WQ

	0	1	2	3	4	5	6	7	8
0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
2	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
3	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
4	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0



+ Code + Text

0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
2	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
3	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
4	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0

W_q, W_v, W_k

3 heads

2. Define number of Attention Head

- For our task assumes N=3

```
[ ] N_head = 3
```

3. Compute Query, Value and Key Matrix

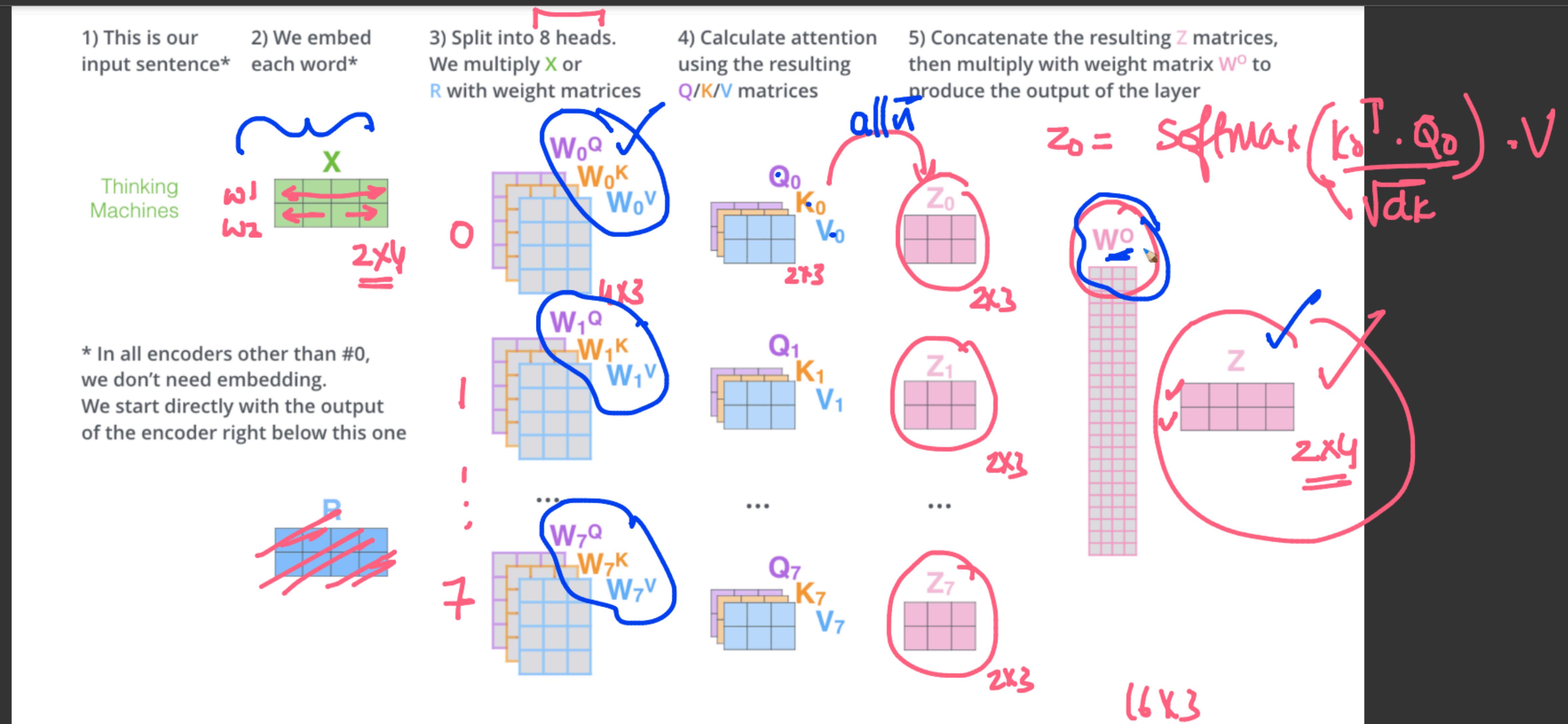
- Three new matrix can be generated by dot product on each weight matrices with seq_embedding,

Query Matrix = $Q = X \bullet W_q$

Key Matrix = $K = X \bullet W_k$

Value Matrix = $V = X \bullet W_v$

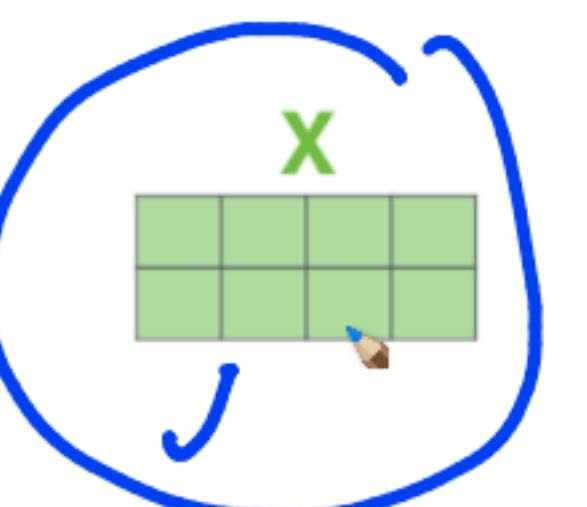
+ Code + Text



+ Code + Text

Connect ▾ |   | ✓

- 1) This is our input sentence*
- 2) We embed each word*



- In all encoders other than #0, we don't need embedding.
We start directly with the output of the encoder right below this one

3) Split into 8 heads.
We multiply \mathbf{X} or
 \mathbf{R} with weight matrices



The diagram illustrates three parallel convolutional layers, W_1^Q , W_1^K , and W_1^V , being applied to an input tensor. The input tensor is represented by a 4x4 grid of light blue squares. Each layer consists of a 2x2 kernel (indicated by colored boxes) applied across the input grid. The output of each layer is a 2x2 grid of darker squares. The layers are color-coded: W_1^Q is purple, W_1^K is orange, and W_1^V is blue.

The diagram illustrates three parallel convolutional layers, W_7^Q , W_7^K , and W_7^V , being applied to an input image. The input image is represented by a grid of 16 light-gray squares. Each layer produces a feature map consisting of 4 larger squares. The first layer, W_7^Q , is shown in purple and produces a feature map where all four squares are light gray. The second layer, W_7^K , is shown in orange and produces a feature map where the top-left square is dark orange and the others are light gray. The third layer, W_7^V , is shown in blue and produces a feature map where the bottom-right square is dark blue and the others are light gray. Ellipses above the layers indicate they are part of a larger network.

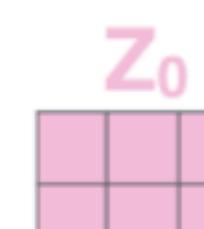
4) Calculate attention using the resulting $O/K/V$ matrices



$$\begin{matrix} Q_1 \\ K_1 \\ V_1 \end{matrix}$$

Q₇ **K₇** **V₇**

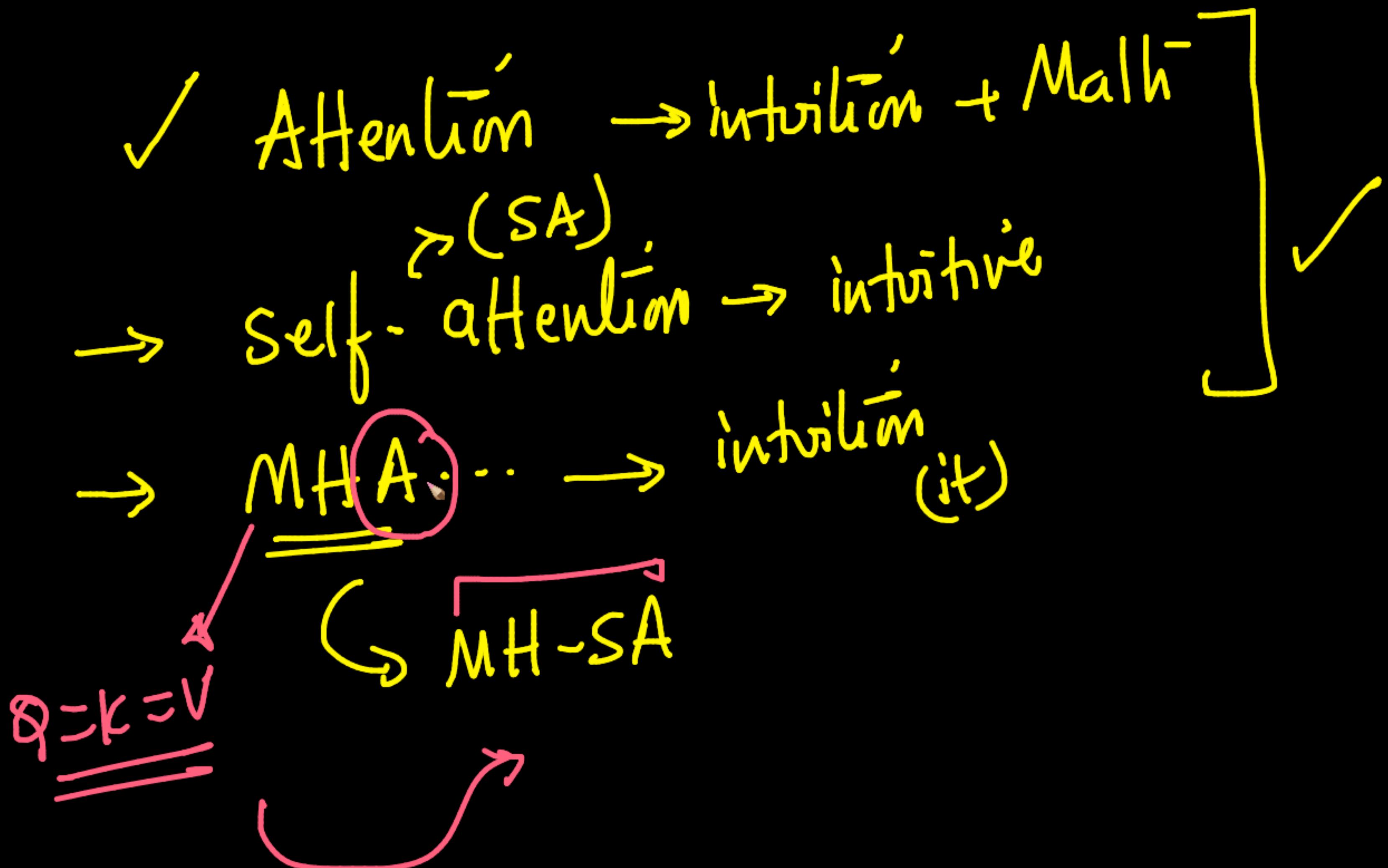
5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer



W

A vertical column of 40 small, light blue squares arranged in a grid pattern. The grid is bounded by a thick pink border.

Z



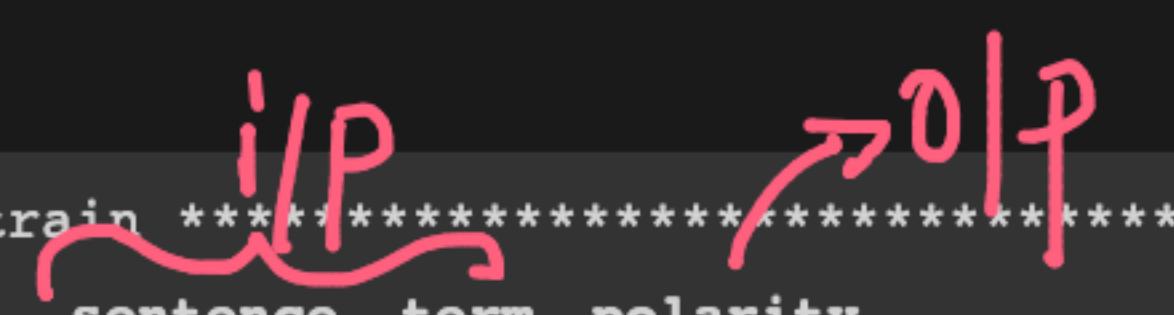
colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=jPWKmLPRYECT

+ Code + Text Connect |  

```
print(f"len(val_pairs) validation pairs")
```

***** Train and Test Split *****
2202 total pairs
1872 training pairs
330 validation pairs

[] print('*'*30,'train', '*'*30)
train_pairs.head()

***** train *****

sentence term polarity

	sentence	term	polarity
432	your a sushi fan, you love expertly cut fish, ...	eel	2
291	the coffe is very good, too.	coffe	2
1197	during the course of the past 3 months, the ch...	staff	0
2390	the bread was stale, the salad was overpriced ...	salad	0
212	best restaurant in the world, great decor, gre...	decor	2

[] print('*'*30,'Train Sentiment label Distribution', '*'*30)
train_pairs['polarity'].value_counts()

***** Train Sentiment label Distribution *****

polarity	count
2	1345
0	453
1	74

27 / 27

8 / 100 | https://www.udemy.com/course/the-art-of-transformers/ | L10_The_Illustrated_Transformer.ipynb | The Illustrated Transformer - Jupyter Notebook

atural machine translation with x | +

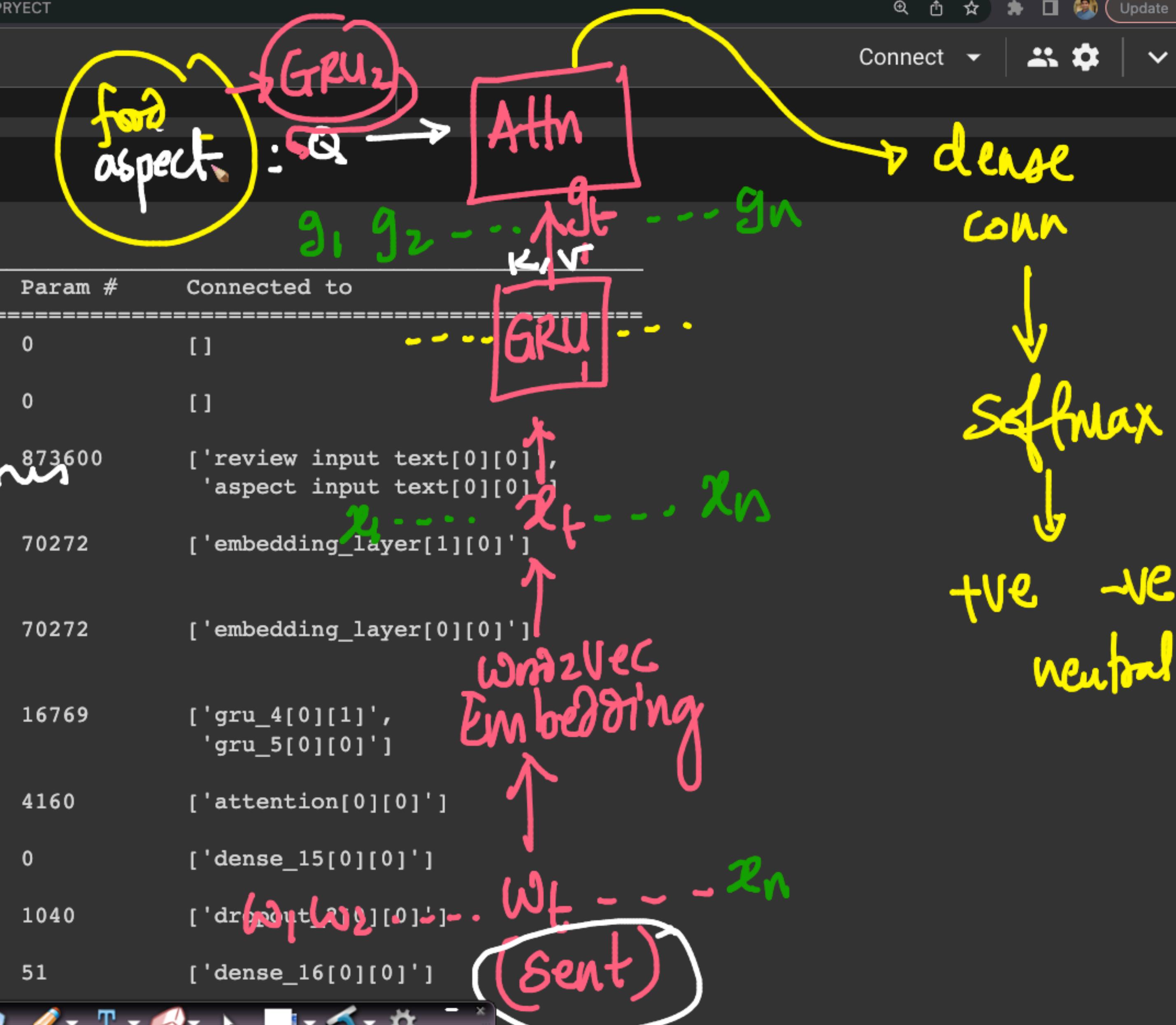
Update

+ Code + Text

CARTRIDGES - Save money, use!

```
[ ] model.summary()
```

Layer (type)	Output Shape	Param #	Connected to
aspect input text (InputLayer)	[(None, 2)]	0	[]
review input text (InputLayer)	[(None, 20)]	0	[]
embedding_layer (Embedding)	multiple	873600	['review input text[0][0]', 'aspect input text[0][0]']
gru_4 (GRU)	[(None, 2, 64), (None, 64)]	70272	['embedding_layer[1][0]']
gru_5 (GRU)	[(None, 20, 64), (None, 64)]	70272	['embedding_layer[0][0]']
attention (Attention)	((None, 64), (None, 20, 1))	16769	['gru_4[0][1]', 'gru_5[0][0]']
dense_15 (Dense)	(None, 64)	4160	['attention[0][0]']
dropout_2 (Dropout)	(None, 64)	0	['dense_15[0][0]']
dense_16 (Dense)	(None, 16)	1040	['dropout_2[0][0]']
dense_17 (Dense)	(None, 3)	51	['dense_16[0][0]']



Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - 0 x The Illustrated Transformer - Jupyter Notebook x Neural machine translation with attention x +

colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=jPWKmLPRYECT

Update

+ Code + Text

Connect |  

```
aspect_input = tf.keras.Input(shape=(2, ), dtype='float64', name = 'aspect input text')

review_embeddings = embedding_layer(review_input)
aspect_embeddings = embedding_layer(aspect_input)

{x}

    _, aspect = tf.keras.layers.GRU(64,
        return_sequences=True, # return the full output sequence
        return_state=True, # return the last hidden state
    )(aspect_embeddings)
    lstm = tf.keras.layers.GRU(64,
        return_sequences=True,
        return_state=True,
    )
whole_seq_output, _ = lstm(review_embeddings)

context_vector, attention_weights = Attention(128)(aspect, whole_seq_output)
dense_con = tf.keras.layers.Dense(64, activation="selu")(context_vector)
dropout = tf.keras.layers.Dropout(0.5)(dense_con) # Adding dropout layer to reduce overfitting
dense_output = tf.keras.layers.Dense(16, activation="selu")(dropout)
output = tf.keras.layers.Dense(3, activation="softmax")(dense_output)

model = tf.keras.Model(inputs=[review_input, aspect_input], outputs=output)
```

Train and Val set Preparation



Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - o The Illustrated Transformer - J x Neural machine translation with x +

colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=jPWKmLPRYECT

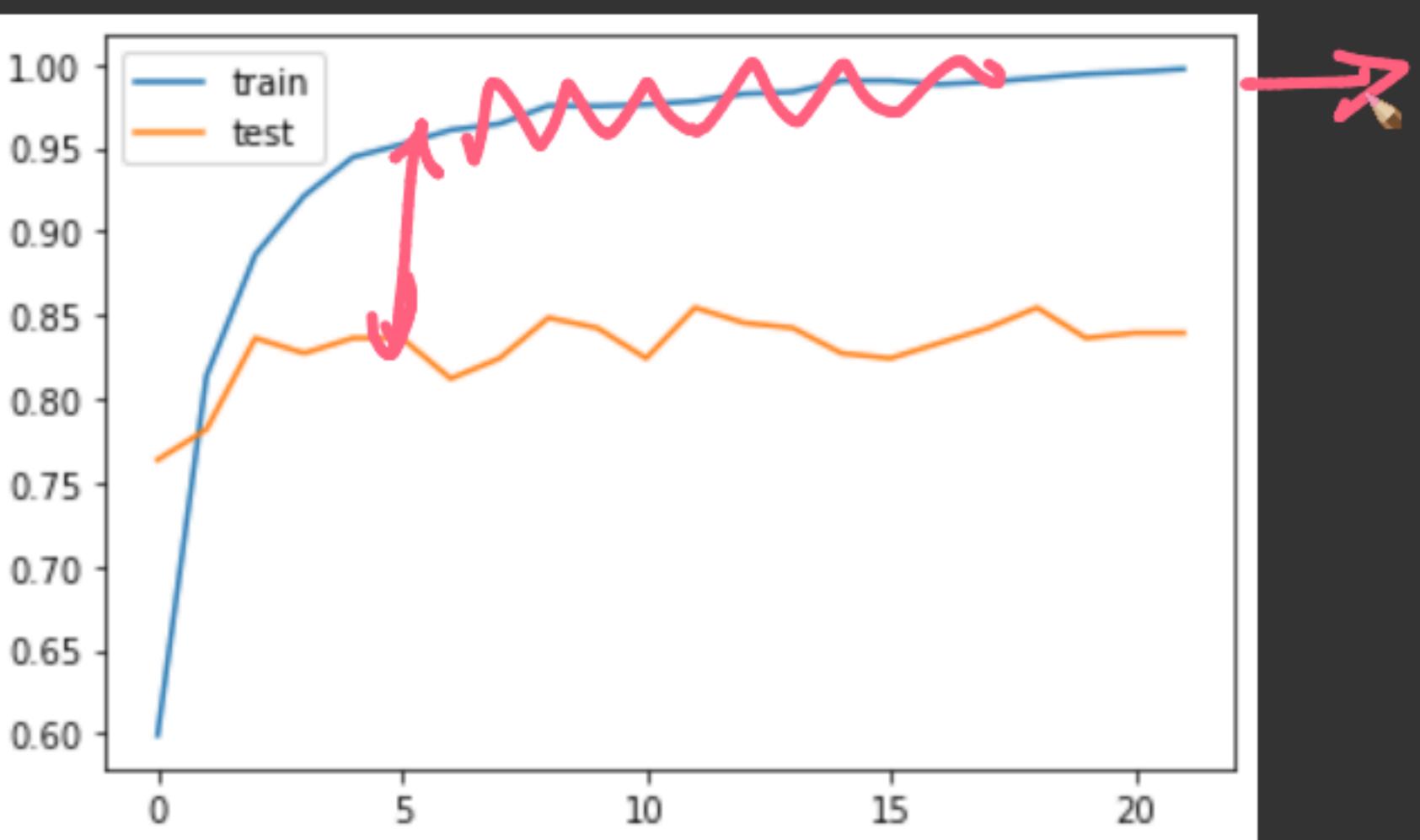
Update

+ Code + Text

Connect |

```
[ ] from matplotlib import pyplot

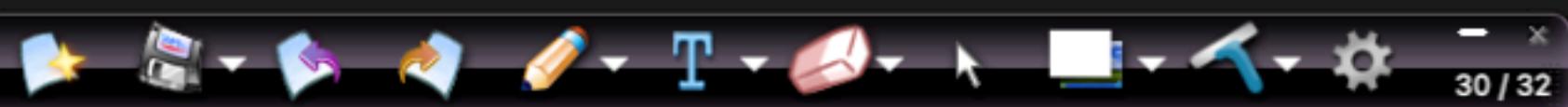
pyplot.plot(history.history['acc'], label='train')
pyplot.plot(history.history['val_acc'], label='test')
pyplot.legend()
pyplot.show()
```



Inferencing

```
[ ] # load the most accurate model, following val_acc curve
model.load_weights('aspect_based_sa.h5')

label_map = {
```



Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - c x The Illustrated Transformer - J x Neural machine translation with x + colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=ZnoWCOOyu0bZ

+ Code + Text Connect |

Model: "model_4"

Layer (type)	Output Shape	Param #	Connected to
aspect input text (InputLayer)	[(None, 2)]	0	[]
review input text (InputLayer)	[(None, 20)]	0	[]
embedding_layer (Embedding)	multiple	873600	['review input text[0][0]', 'aspect input text[0][0]']
gru_4 (GRU)	[(None, 2, 64), (None, 64)]	70272	['embedding_layer[1][0]']
gru_5 (GRU)	[(None, 20, 64), (None, 64)]	70272	['embedding_layer[0][0]']
attention (Attention)	((None, 64), (None, 20, 1))	16769	['gru_4[0][1]', 'gru_5[0][0]']
dense_15 (Dense)	(None, 64)	4160	['attention[0][0]']
dropout_2 (Dropout)	(None, 64)	0	['dense_15[0][0]']
dense_16 (Dense)	(None, 16)	1040	['dropout_2[0][0]']
dense_17 (Dense)	(None, 3)	51	['dense_16[0][0]']

Total params: 1,036,164
Trainable params: 162,564
Non-trainable params: 873,600

gru_4 *gru_5* { *gru_4* *gru_5* } *embedding_layer* *attention* *dense_15* *dropout_2* *dense_16* *dense_17*

+ Code + Text

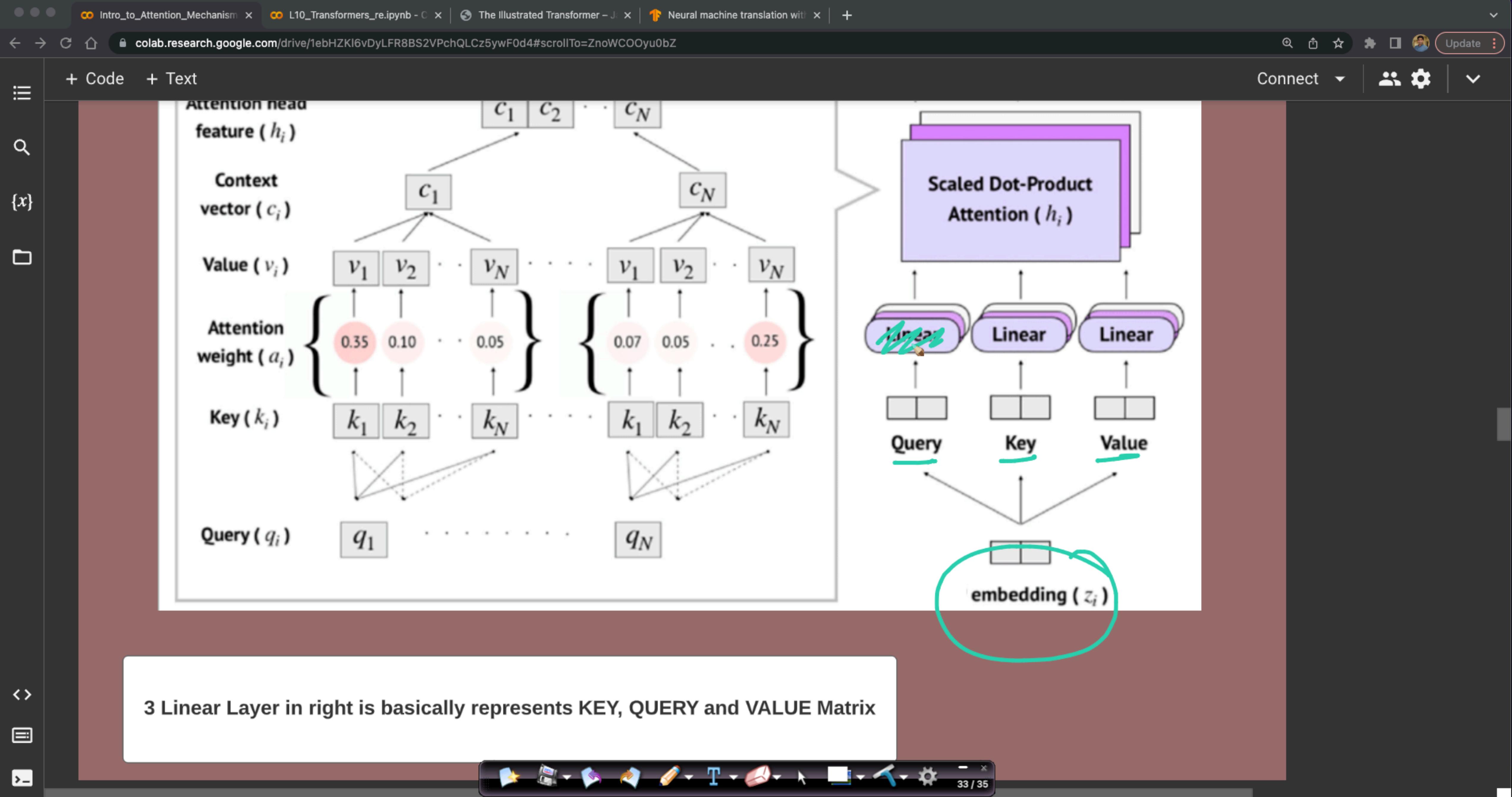
Connect

Model: "model_4"

Layer (type)	Output Shape	Param #	Connected to
aspect input text (InputLayer)	[(None, 2)]	0	[]
review input text (InputLayer)	[(None, 20)]	0	[]
embedding_layer (Embedding)	multiple	873600	['review input text[0][0]', 'aspect input text[0][0]']
gru_4 (GRU)	[(None, 2, 64), (None, 64)]	70272	['embedding_layer[1][0]']
gru_5 (GRU)	[(None, 20, 64), (None, 64)]	70272	['embedding_layer[0][0]']
attention (Attention)	((None, 64), (None, 20, 1))	16769	['gru_4[0][1]', 'gru_5[0][0]']
dense_15 (Dense)	(None, 64)	4160	['attention[0][0]']
dropout_2 (Dropout)	(None, 64)	0	['dense_15[0][0]']
dense_16 (Dense)	(None, 16)	1040	['dropout_2[0][0]']
dense_17 (Dense)	(None, 3)	51	['dense_16[0][0]']

Total params: 1,036,164
Trainable params: 162,564
Non-trainable params: 873,600

dense { w^Q w^K w^V



Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - o The Illustrated Transformer - J x Neural machine translation with x +

colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=ZnoWCOOyu0bZ

Update

+ Code + Text

Connect



```
x1 = tokenizer.texts_to_sequences([['service']])
x1 = tf.keras.preprocessing.sequence.pad_sequences(x1, maxlen=2)
preds = model.predict([x,x1])
print('*'*30,'aspect-term : service', '*'*30)
print(label_map[np.argmax(preds)])
```

✓ Input String : the service was worst , but ambience was amazing

***** aspect-term : ambience *****

WARNING:tensorflow:5 out of the last 13 calls to <function Model.make_predict_function.<locals>.predict_function at 0x7fcff424f440> triggered

1/1 [=====] - 1s 1s/step

✓ Positive

1/1 [=====] - 0s 23ms/step

***** aspect-term : service *****

Negative

▼ Attention Visualization

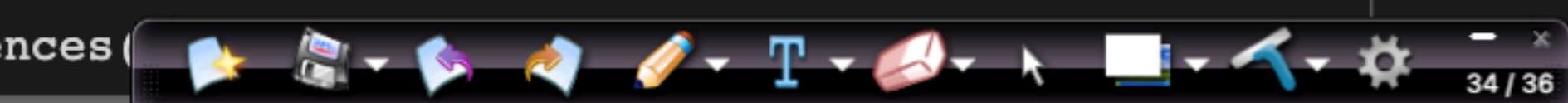
```
[ ] attention_model = tf.keras.Model(inputs=model.input, outputs=[model.output, model.get_layer('attention').output, model.input])
```

```
[ ] def plot_attention(target):
```

```
    # plot attention curve
```

```
    #define input output
```

```
    x1 = tokenizer.texts_to_sequences([target])
```



Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - c x The Illustrated Transformer - J x Neural machine translation with x +

← → ⌂ ⌂ colab.research.google.com/drive/1ebHZKI6vDyLFR8BS2VPchQLCz5ywF0d4#scrollTo=ZnoWCOOyu0bZ

Update

+ Code + Text

Connect |  

```
preds = model.predict([x,x1])
print('*'*30,'aspect-term : service', '*'*30)
print(label_map[np.argmax(preds)])
```

{x} Input String : the service was worst , but ambience was amazing

***** aspect-term : ambience *****

WARNING:tensorflow:5 out of the last 13 calls to <function Model.make_predict_function.<locals>.predict_function at 0x7fcff424f440> triggered warnings.

1/1 [=====] - 1s 1s/step

Positive

1/1 [=====] - 0s 23ms/step

***** aspect-term : service *****

Negative

▼ Attention Visualization

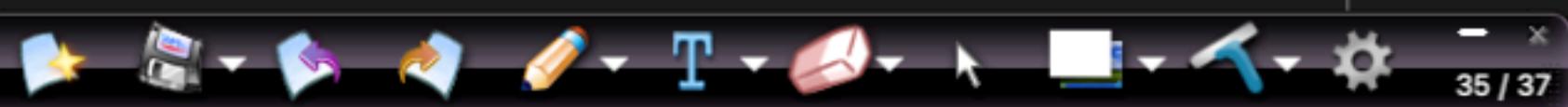
```
[ ] attention_model = tf.keras.Model(inputs=model.input, outputs=[model.output, model.get_layer('attention').output, model.input])
```

```
[ ] def plot_attention(target):

    # plot attention curve

    #define input output
    x1 = tokenizer.texts_to_sequences([[target]])
    x1 = tf.keras.preprocessing.sequence.pad_sequences(x1, maxlen=2)

    preds = model.predict([x,x1])
```



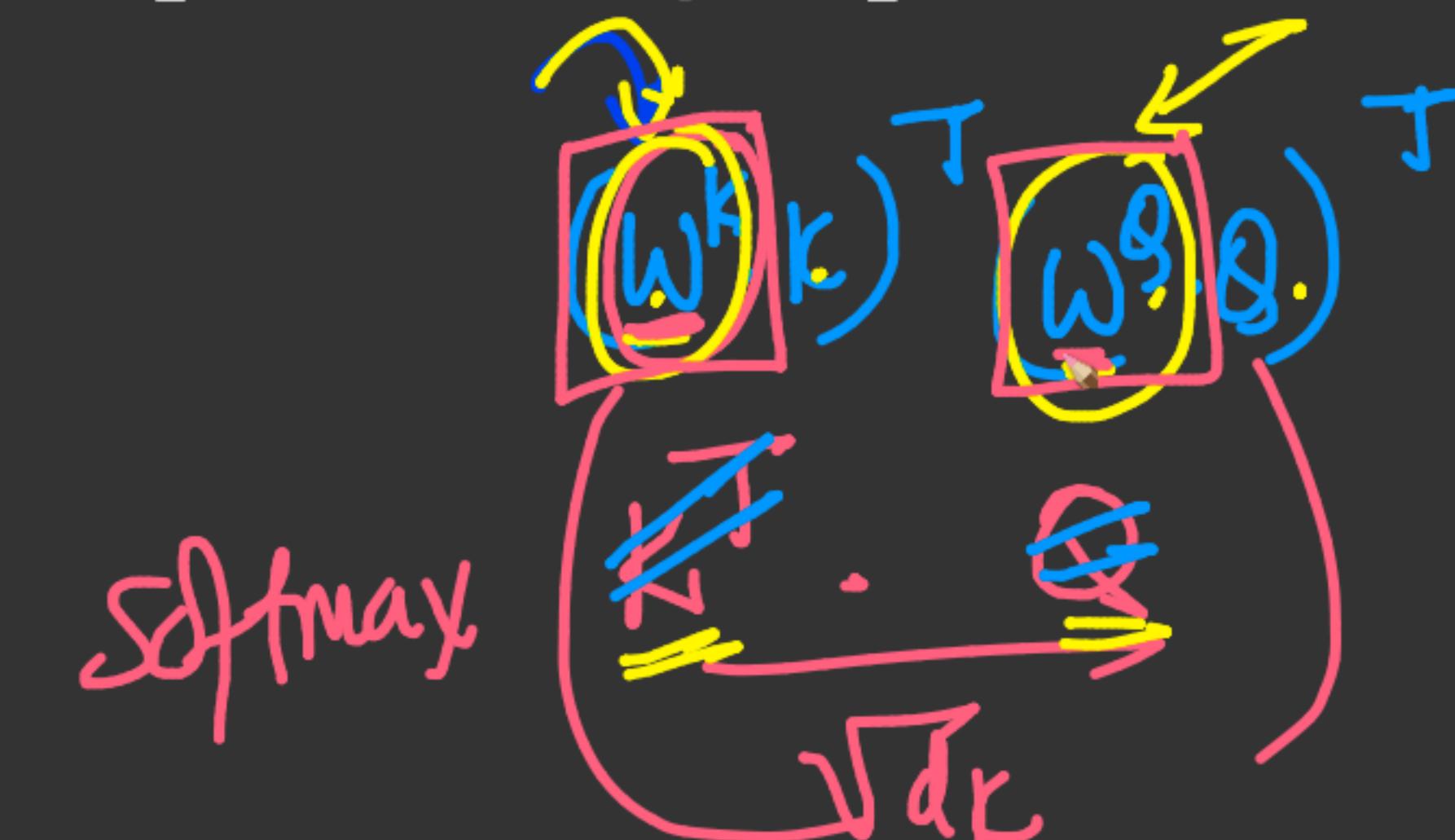
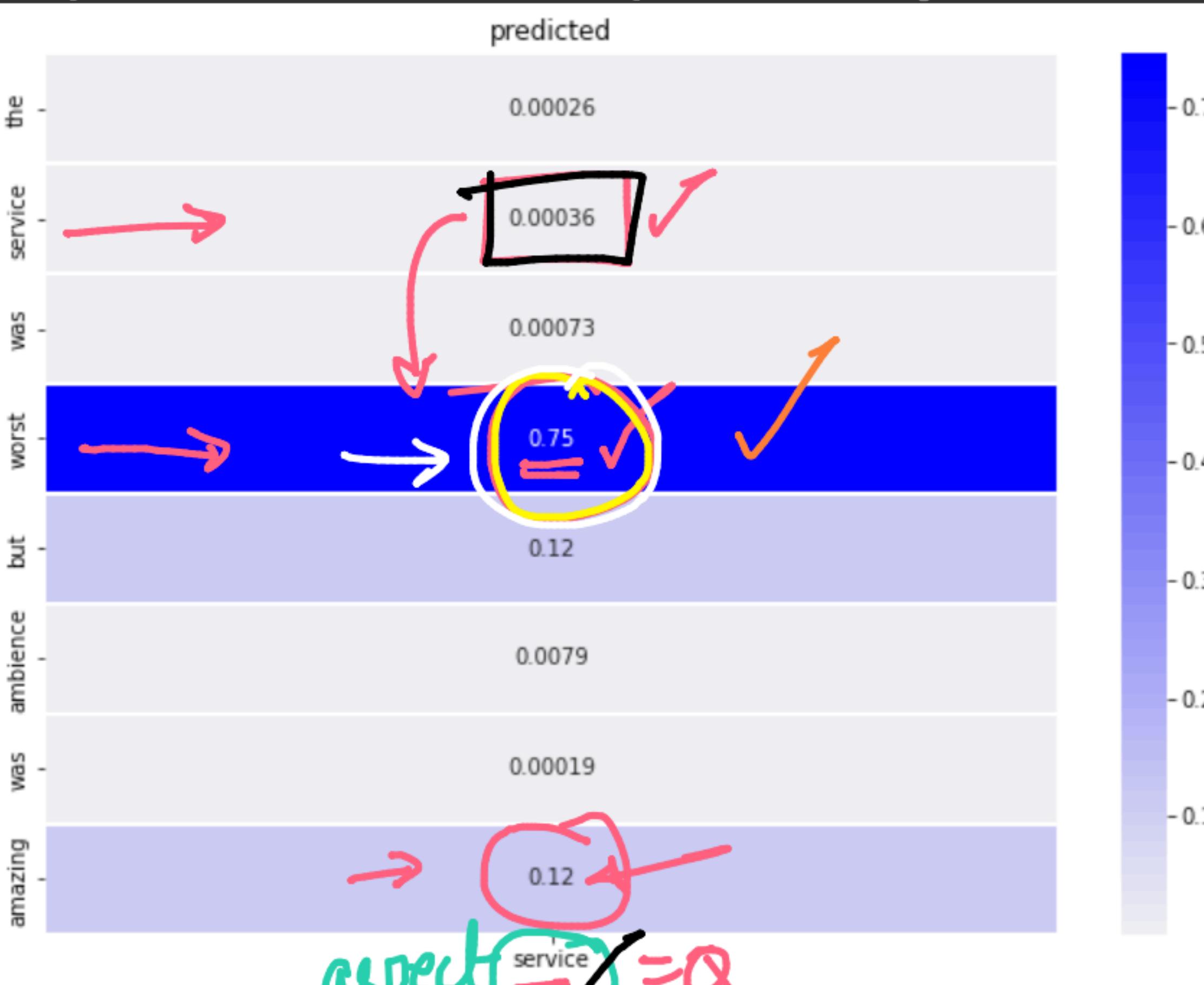
+ Code + Text

***** aspect-term : service *****

Negative ✓

WARNING:tensorflow:5 out of the last 13 calls to <function Model.make_predict_function.<locals>.predict_function at 0x7fcff40560e0> triggered warnings.

1/1 [=====] - 1s 959ms/step



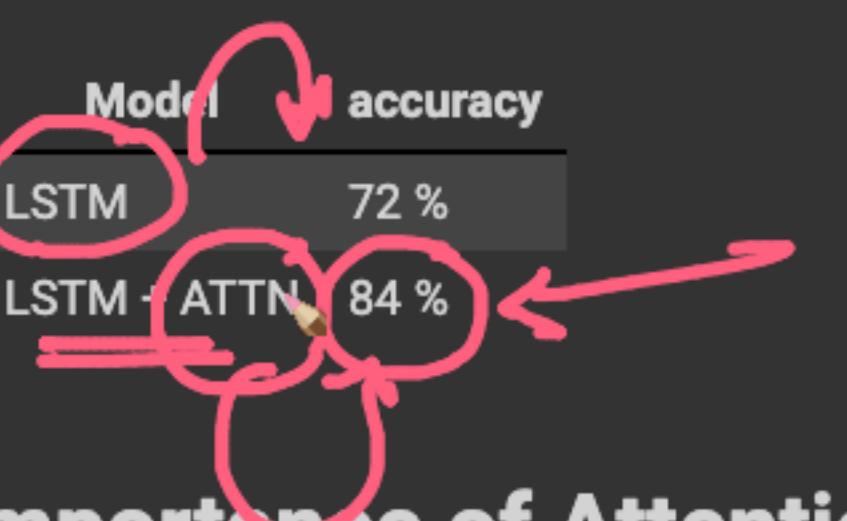
+ Code + Text

Connect |

Explainability

Model	Explainability
LSTM	Not possible
LSTM + ATTN	possible to visualize

Performance



Importance of Attention Mechanisms in training Popular NLP Model

- helps in learning good contextual representation of a text
- Adding Explainability
- helps in performance improvement
- Facilitates transfer learning

TRANSFORMERS

- Language translation Model
- based on MULTI-HEAD SELF-ATTENTION MECHANISMS
- Get rid of recurrent architecture to understand



+ Code + Text

Connect |

Explainability

Model Explainability

Model	Explainability
LSTM	Not possible
LSTM + ATTN	possible to visualize



Performance

Model accuracy

Model	accuracy
LSTM	72 %
LSTM + ATTN	84 %

Importance of Attention Mechanisms in training Popular NLP Model

- helps in learning good contextual representation of a text
- Adding Explainability
- helps in performance improvement
- Facilitates transfer learning

$w_1 \dots w_t \dots w_{l_w}$

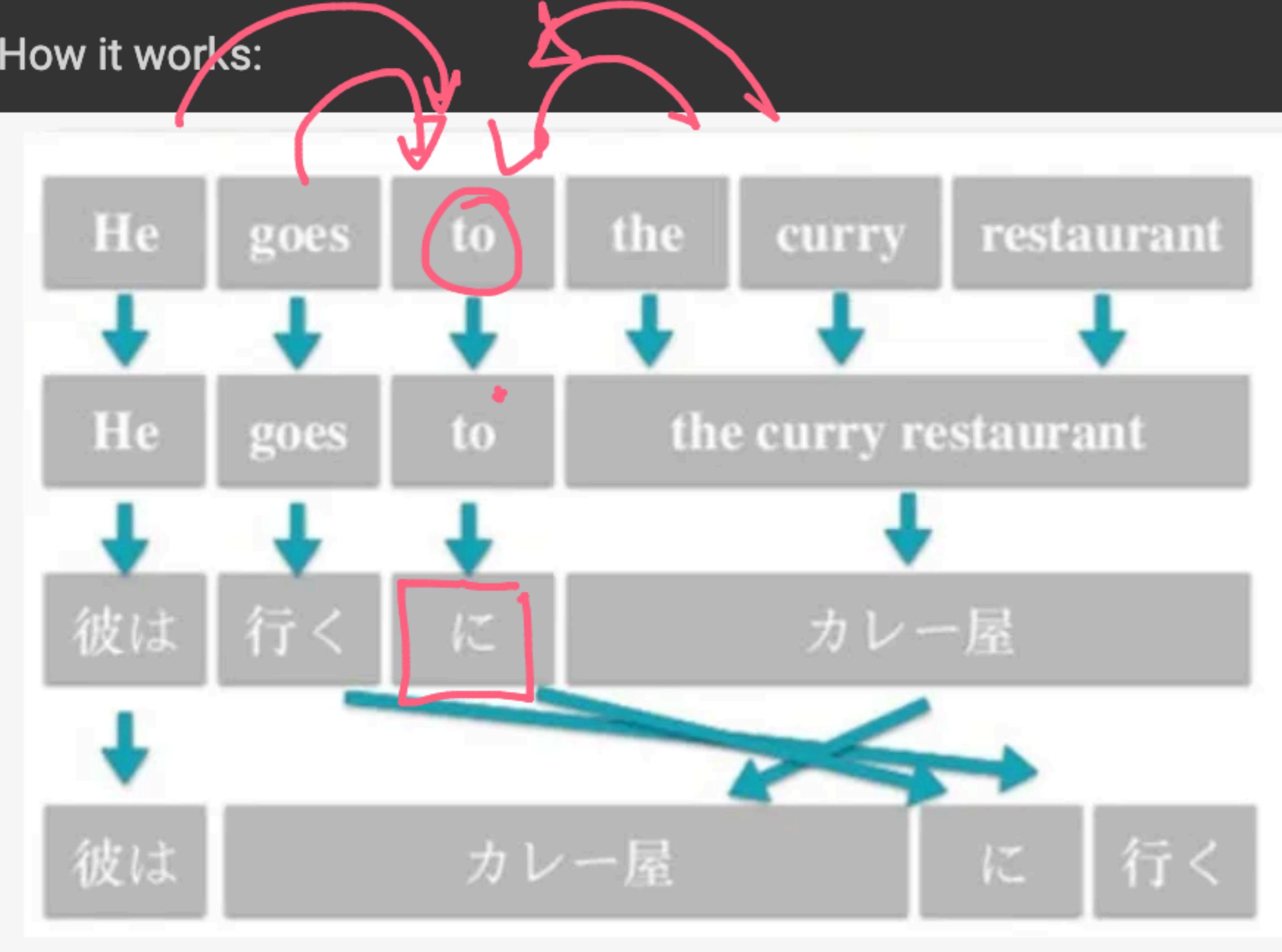
TRANSFORMERS

- Language translation Model
- based on MULTI-HEAD SELF-ATTENTION MECHANISMS
- Get rid of recurrent architecture to understand text

+ Code + Text Last edited on 9 November

Connect | User Settings | More

- Sentence is divided into small units of words and phrases.
- Statistical Models learn how a specific source word is translated into different word in target language.
- Models give more weightage to the source-target pair which comes maximum times in corpus.

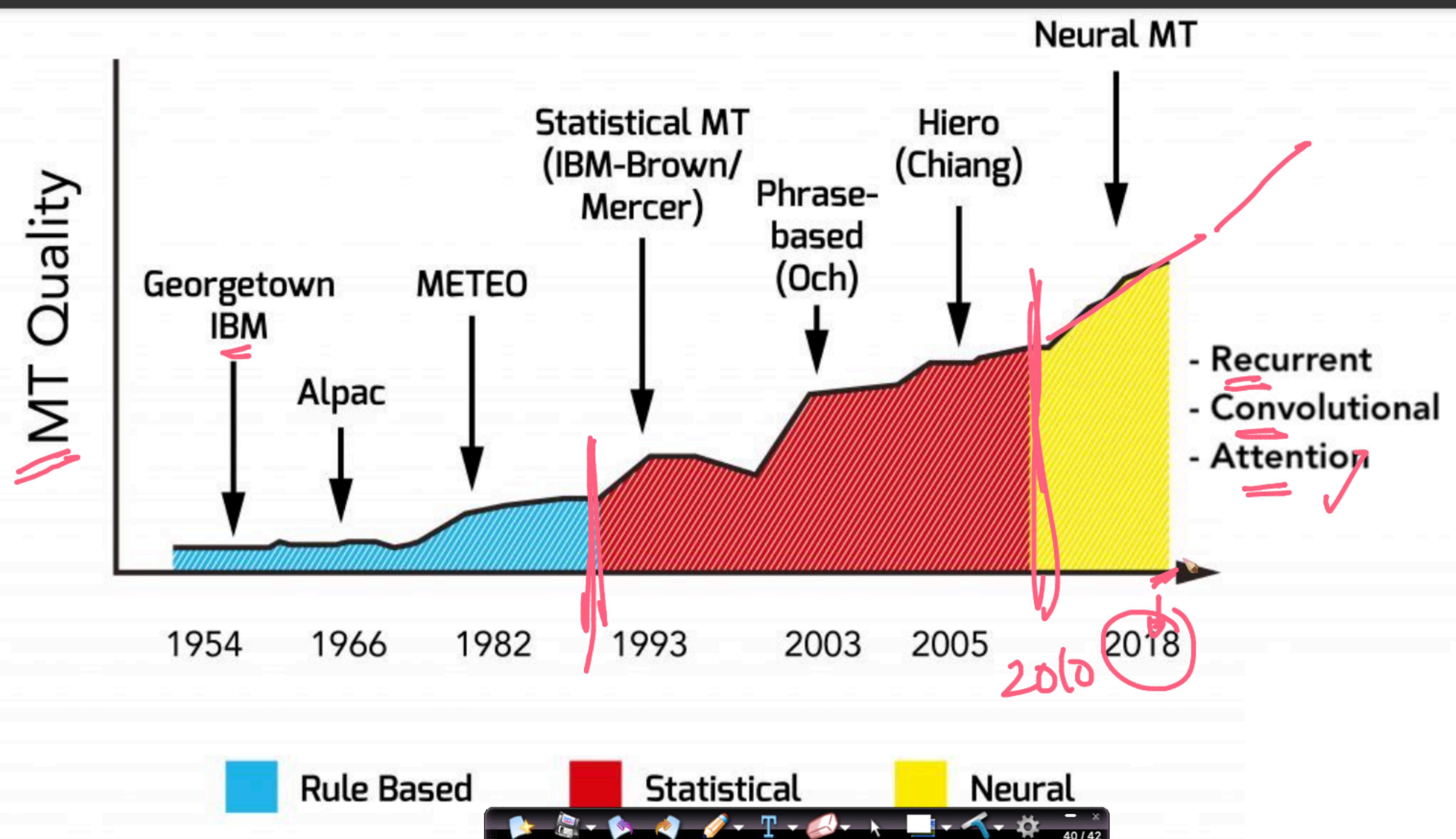


$$P(y_t \mid x_1, x_2, \dots)$$

- Corpus is created which contains very large data set of approved previous translations.
- Translation Model(TM) is build using the f

+ Code + Text Last edited on 9 November

Connect | People | Settings | More



Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - x The Illustrated Transformer - Jupyter Notebook x Neural machine translation with x + colab.research.google.com/drive/1n4mT8SduOrSfIPUzIL5Q9OPf8aFhhtNN#scrollTo=U6t5GYld5lLw

+ Code + Text Last edited on 9 November Connect |

What is Encoder-Decoder Architecture?

{x}

Encoder-Decoder Architecture

Input Sequence
 $x = \{x_1, x_2, \dots, x_n\}$

Encoder

Intermediate representation z

Decoder

Predicted Sequence
 $y = \{y_1, y_2, \dots, y_m\}$

Consist of:

- Encoder Component
- Encoder Vector/Hidden State Vector/Context Vector
- Decoder Component

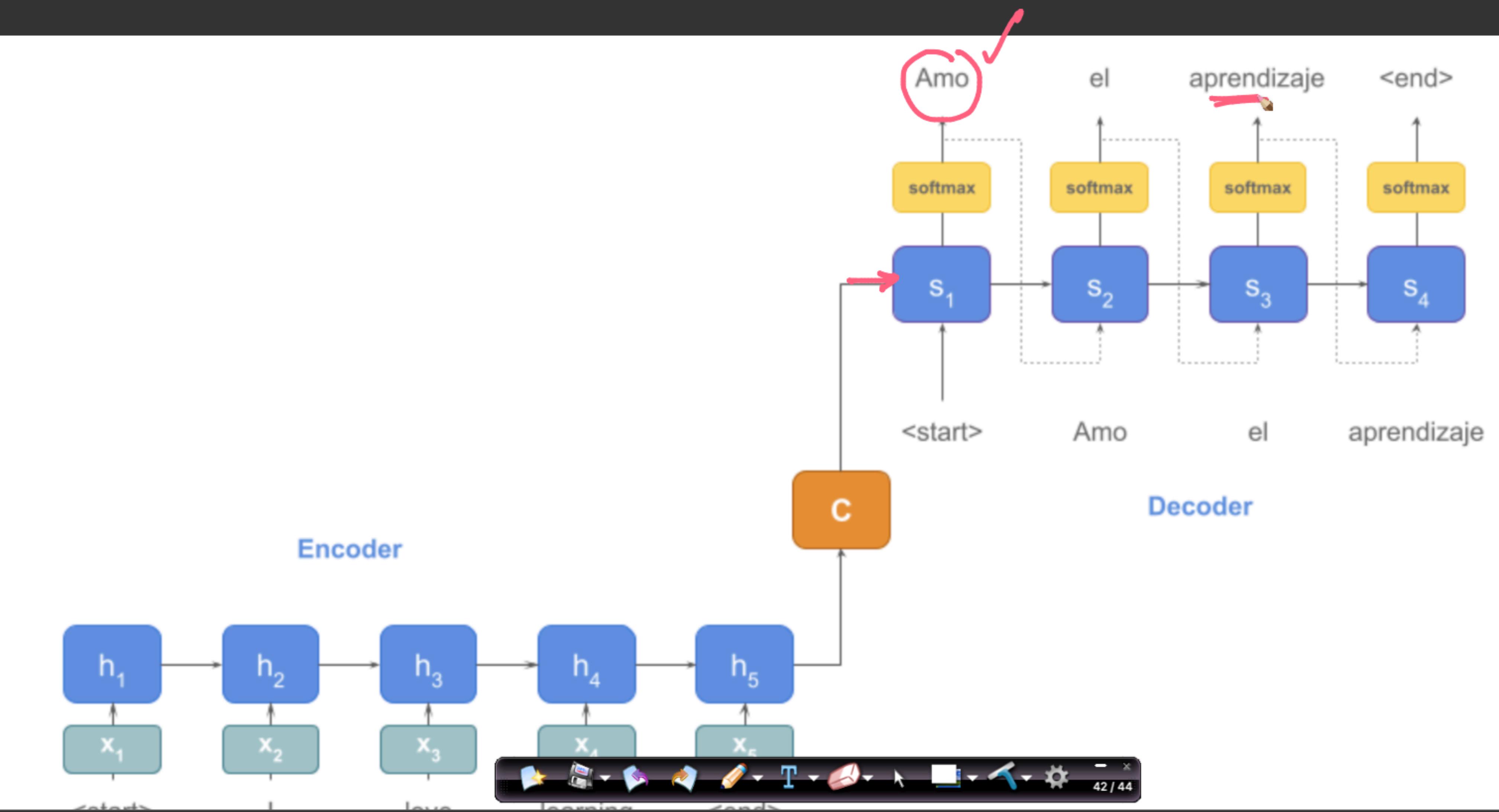
Encoder

```
graph LR; x["Input Sequence  
x = {x1, x2, ..xn}"] --> Encoder[Encoder]; Encoder -- "z" --> Decoder[Decoder]; Decoder --> y["Predicted Sequence  
y = {y1, y2, ..ym}"]
```

+ Code + Text Last edited on 9 November

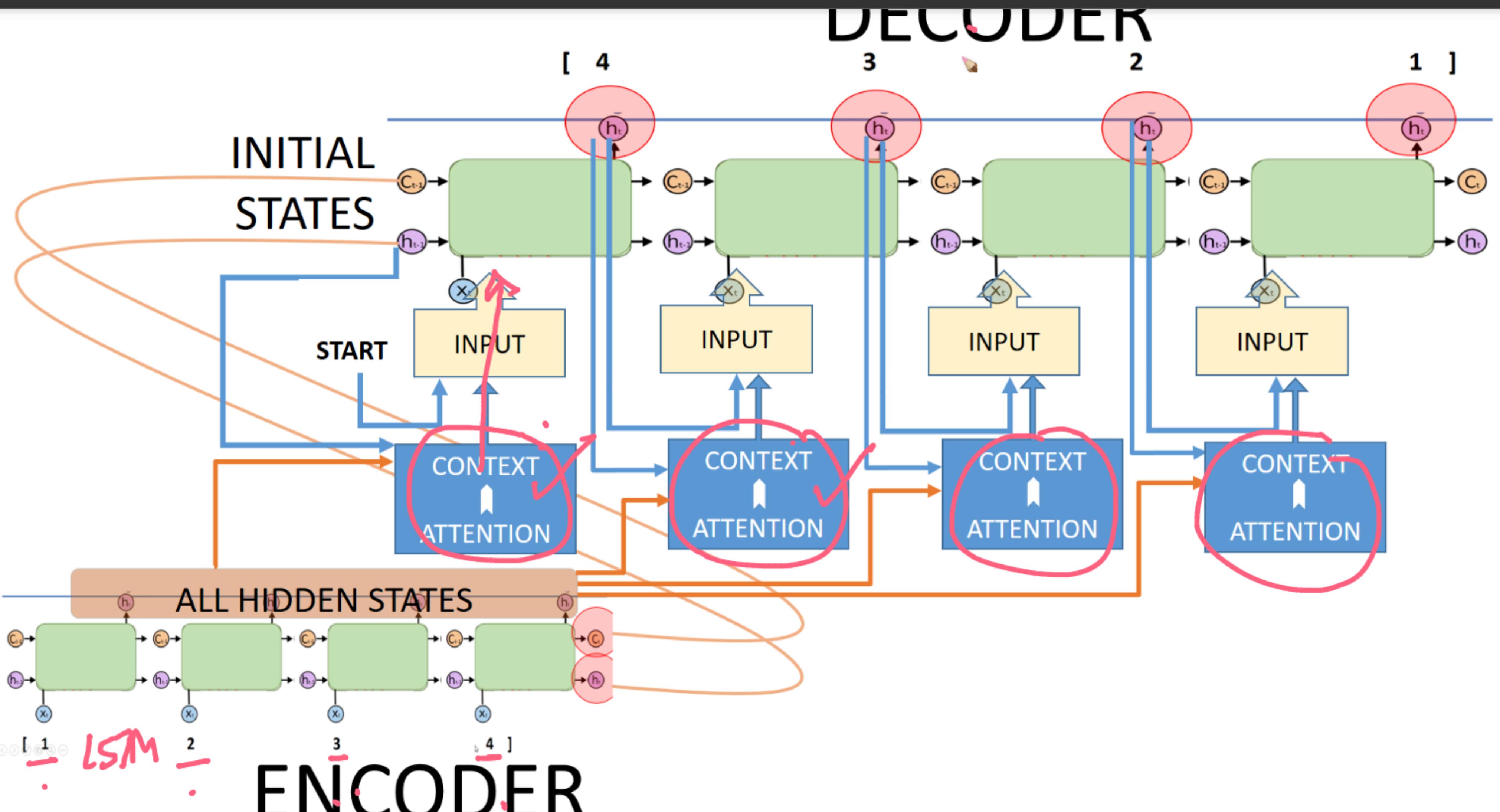
Connect | People | Settings

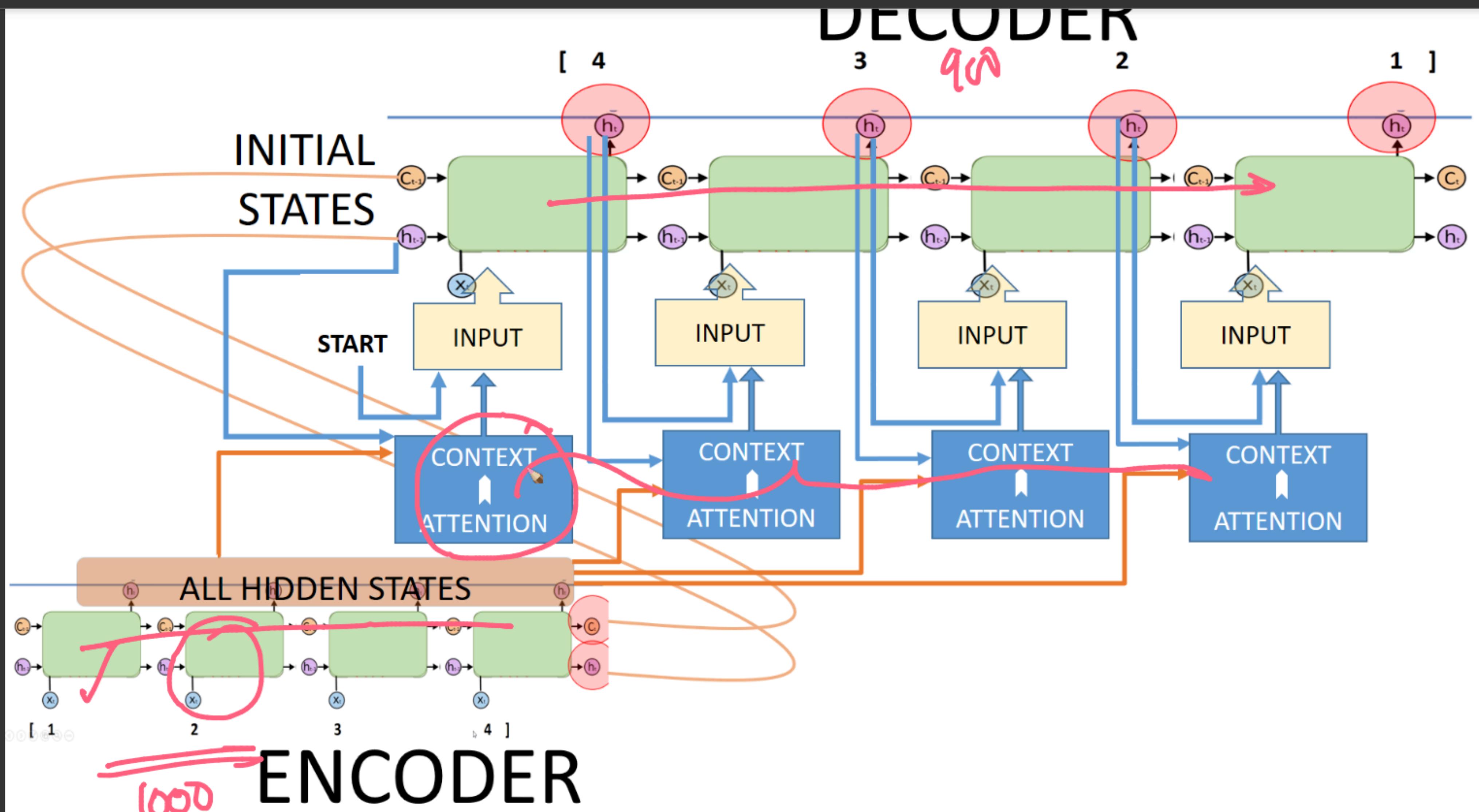
- The Encoder summarizes the input sequence into encoded vectors
- Decoder uses Encoded vectors and Generates output sequences using softmax



+ Code + Text Last edited on 9 November

Connect | People | Settings | More

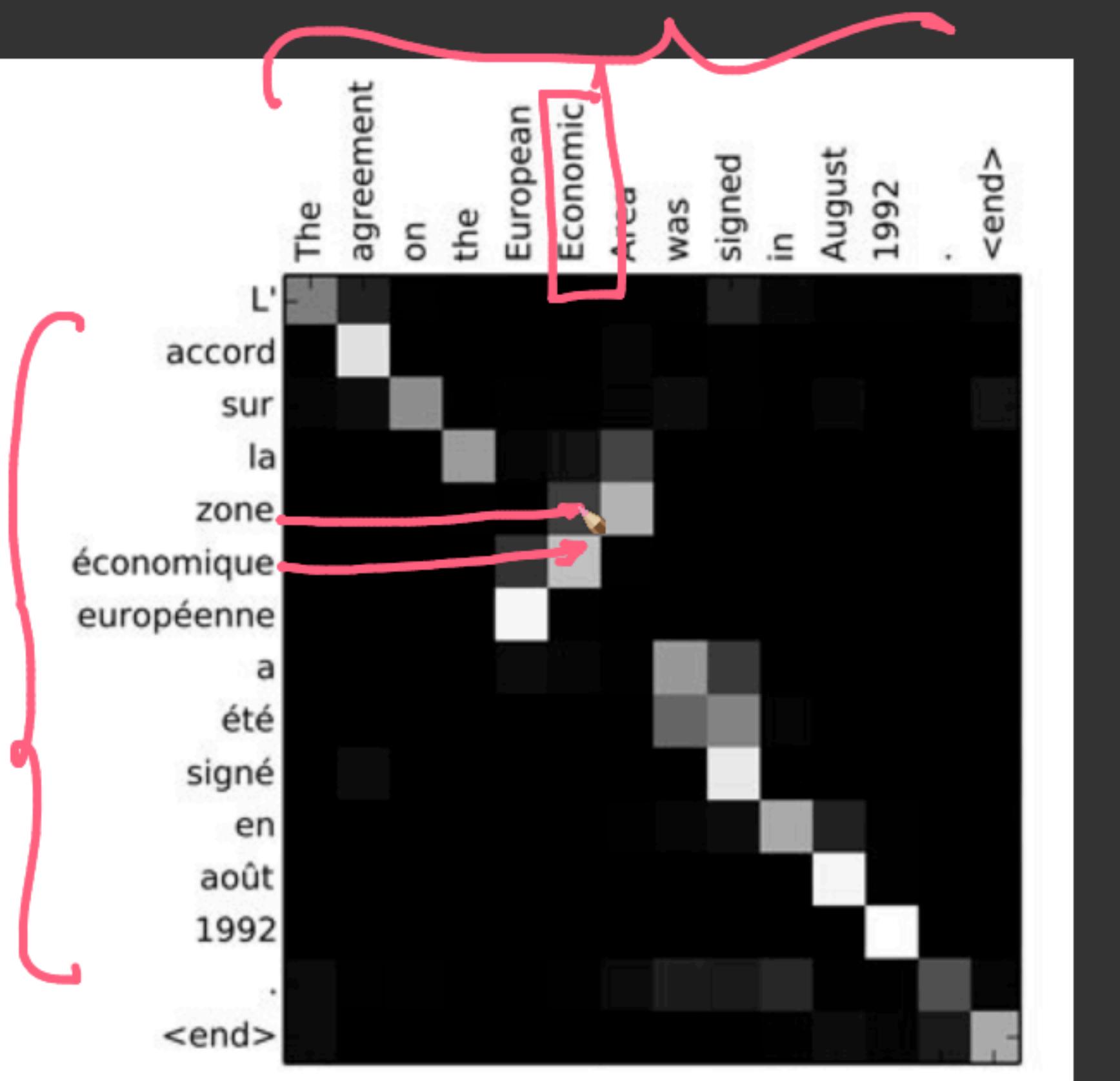




Connect |  

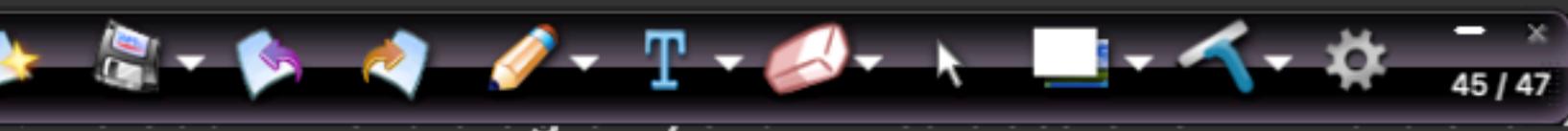
+ Code + Text Last edited on 9 November

The Attention Mechanism can be viewed as a method for making the RNN work better by letting the network know where to look as it is performing its task.



Let's See our translation problems:

- Translating a sentence from English to French



network is “paying attention” as it

colab.research.google.com/drive/1n4mT8SduOrSfIPUzIL5Q9OPf8aFhhtNN#scrollTo=x7MWDTU1Vgno

+ Code + Text Connect |

Property RNN/LSTM RNN/LSTM+Attention Transformers

{x}	Recurrence blocks	Yes	Yes	No
Attention	No	Yes	Yes	
context vector	last timestamp state vector	calculated using attention mechanisms	using attention mechanisms	
attention alignment function	NA	Dot product	Scaled Dot Product	
attention type	NA	single head	Multi head	
operation	sequential	sequential	parallel	

What are Transformers Components?

Following are the components of Transformers:

1. Position Encodings
2. Transformers Encoders
3. Multi-headed Attention
4. Transformers Decoders

46 / 48

Intro_to_Attention_Mechanism x L10_Transformers_re.ipynb - Colab Notebooks The Illustrated Transformer - Jupyter Notebook Neural machine translation with x + colab.research.google.com/drive/1n4mT8SduOrSfIPUzIL5Q9OPf8aFhhtNN#scrollTo=x7MWDTU1Vgno

+ Code + Text Connect |

Property RNN/LSTM RNN/LSTM+Attention Transformers

Recurrence blocks	Yes	Yes	No
Attention	No	Yes	Yes
context vector	last timestamp state vector	calculated using attention mechanisms	using attention mechanisms
attention alignment function	NA	Dot product	Scaled Dot Product
attention type	NA	single head	Multi head
operation	sequential	sequential	parallel

{x}

What are Transformers Components?

Following are the components of Transformers:

1. Position Encodings
2. Transformers Encoders
3. Multi-headed Attention
4. Transformers Decoders

47 / 49

Intro_to_Attention_Mechanism x

L10_Transformers_re.ipynb - C x

The Illustrated Transformer - J x

Neural machine translation with x

+

colab.research.google.com/drive/1n4mT8SduOrSfIPUzIL5Q9OPf8aFhhtNN#scrollTo=x7MWDTU1Vgno

Update

:

+ Code + Text

Connect



▼

attention alignment function NA

Dot product

Scaled Dot Product

↑ ↓ ↪



⋮

attention type NA

single head

Multi head

operation sequential

sequential

parallel

{x}

What are Transformers Components?

Following are the components of Transformers:

- ✓ 1. Position Encodings
- ✓ 2. Transformers Encoders
- ✓ 3. Multi-headed Attention
- ✓ 4. Transformers Decoders
- ✓ 5. FFN

<>

≡

▶

As Transformer does not use recurrence layer like RNN/LSTM, how it understand the relative position of sequence in the inputs?



Intro_to_Attention_Mechanism x

L10_Transformers_re.ipynb - C x

The Illustrated Transformer - J x

Neural machine translation with x

colab.research.google.com/drive/1n4mT8SduOrSfIPUzIL5Q9OPf8aFhhtNN#scrollTo=x7MWDTU1Vgno

Update

+ Code + Text

Connect



attention alignment function	NA	Dot product	Scaled Dot Product
attention type	NA	single head	Multi head
operation	sequential	sequential	parallel

{x}

What are Transformers Components?

Following are the components of Transformers:

1. Position Encodings
2. Transformers Encoders
3. Multi-headed Attention
4. Transformers Decoders
5. FFN

As Transformer does not use recurrence layer like RNN/LSTM, how it understand the relative position of sequence in the inputs?



+ Code + Text

Connect



attention alignment function NA

Dot product

Scaled Dot Product



attention type NA

single head

Multi head

operation sequential

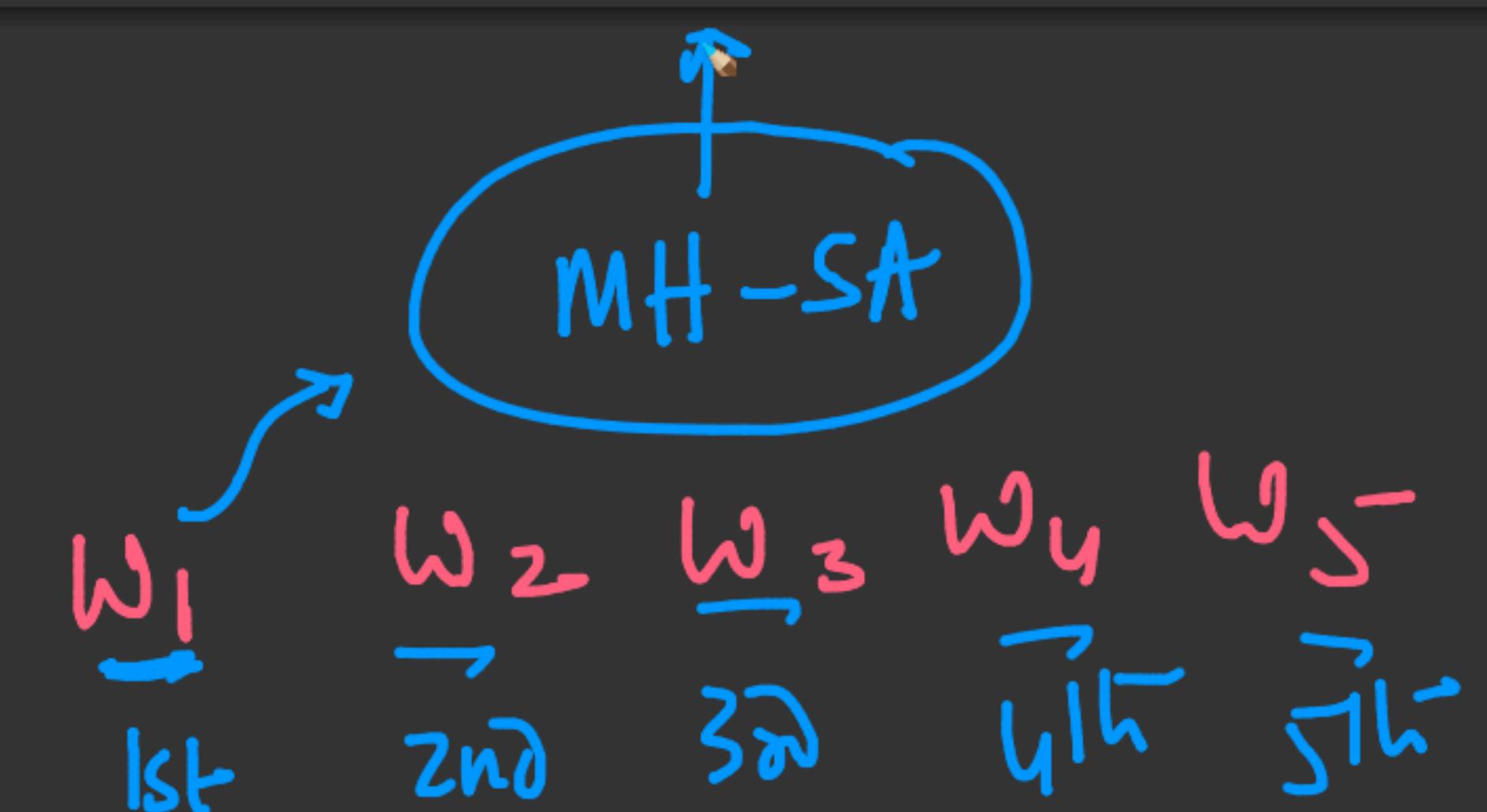
sequential

parallel

What are Transformers Components?

Following are the components of Transformers:

1. Position Encodings
2. Transformers Encoders
3. Multi-headed Attention
4. Transformers Decoders
5. FFN



As Transformer does not use recurrence layer like RNN/LSTM, how it understand the relative position of sequence in the inputs?

The Illustrated Transformer - Jalammar

Intro_to_Attention_Mechanism x | L10_Transformers_re.ipynb - C x | The Illustrated Transformer - Jalammar x | Neural machine translation with x | +

jalamar.github.io/illustrated-transformer/ Update

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

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

R

W₁
W₂

X

W₀^Q W₀^K W₀^V

Q₀ K₀ V₀

Z₀

W₁^Q W₁^K W₁^V

Q₁ K₁ V₁

Z₁

...

W₇^Q W₇^K W₇^V

Q₇ K₇ V₇

Z₇

...

W^O

Z

Now that we have touched upon attention heads, let's revisit our example from before to see where the different attention heads are focusing as we encode the word "it" in our example sentence:

Layer: 5 Attention: Input - Input

didn_

51 / 53

+ Code + Text

Connect



attention alignment function NA

Dot product

Scaled Dot Product

attention type

NA

single head

Multi head

operation

sequential

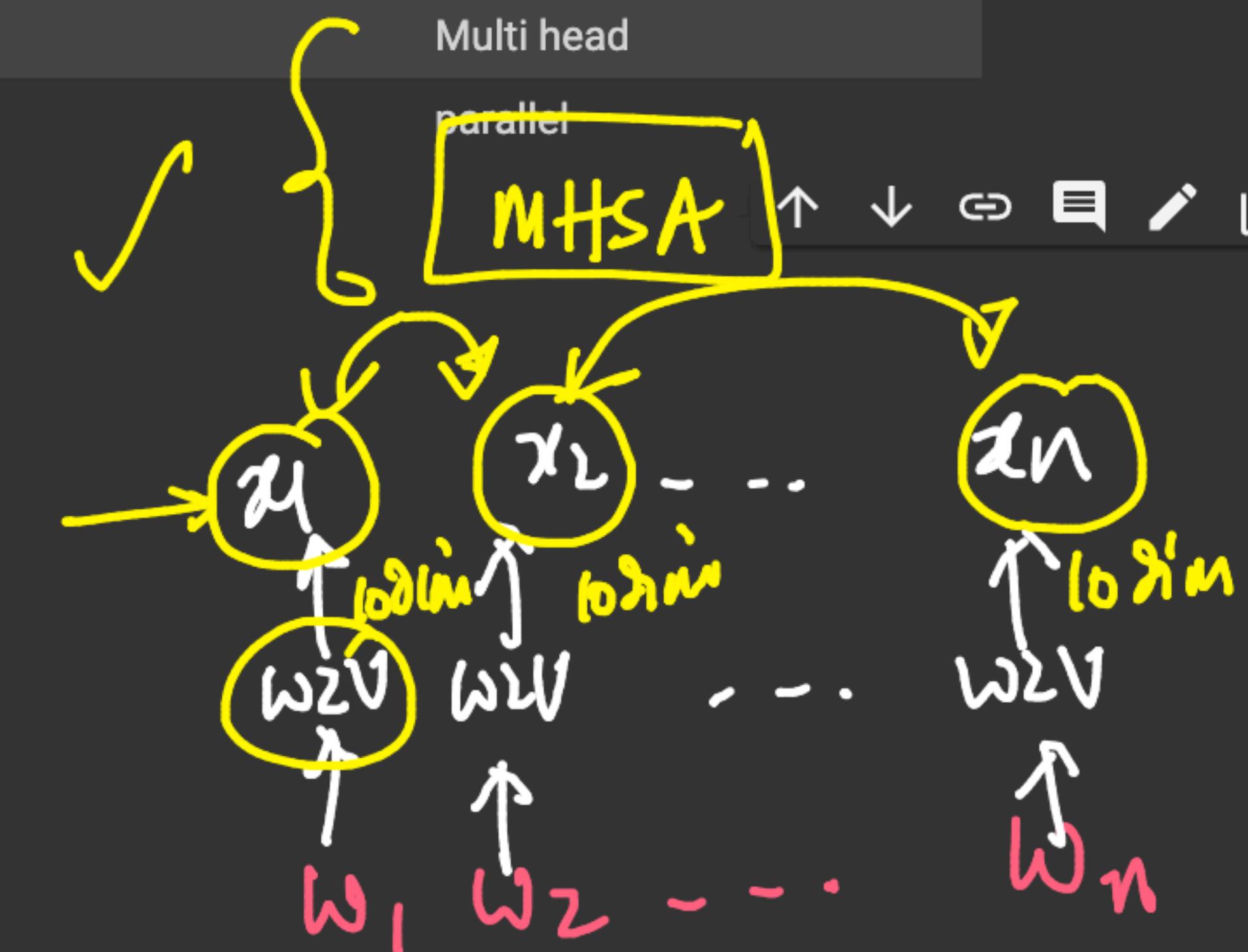
sequential

parallel

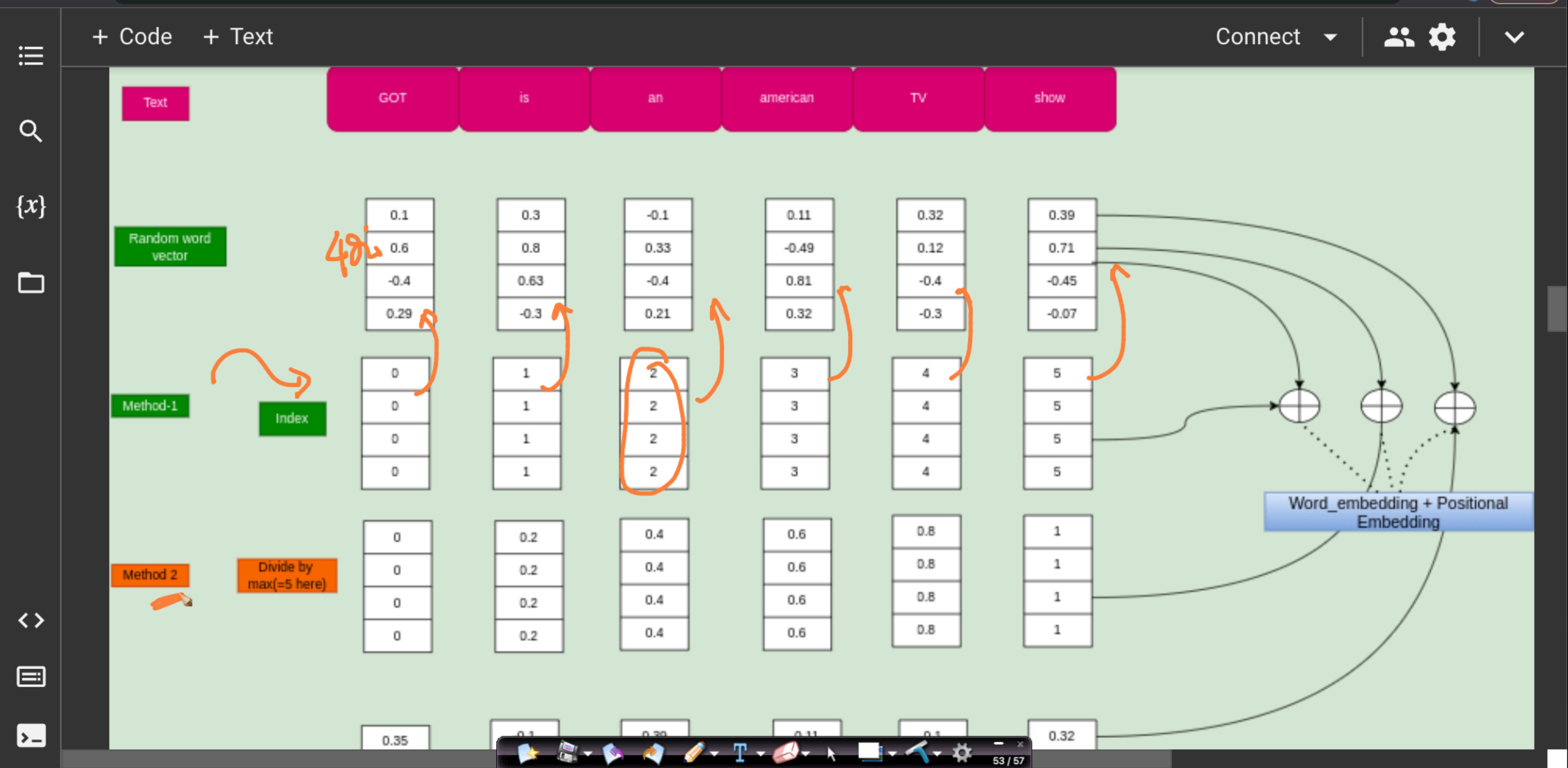
What are Transformers Components?

Following are the components of Transformers:

- ✓ 1. Position Encodings
- 2. Transformers Encoders
- 3. Multi-headed Attention
- 4. Transformers Decoders
- 5. FFN

{ positional
enc }

As Transformer does not use recurrence layer like RNN/LSTM, how it understand the relative position of sequence in the inputs?



+ Code + Text

Connect

Text GOT is an american TV show

{ x }

Random word vector

0.1	0.3	-0.1	0.11	0.32	0.39
0.6	0.8	0.33	-0.49	0.12	0.71
-0.4	0.63	-0.4	0.81	-0.4	-0.45
0.29	-0.3	0.21	0.32	-0.3	-0.07

Method-1

Index

0	1	2	3	4	5
0	1	2	3	4	5
0	1	2	3	4	5
0	1	2	3	4	5

Word_embedding + Positional Embedding

Method 2

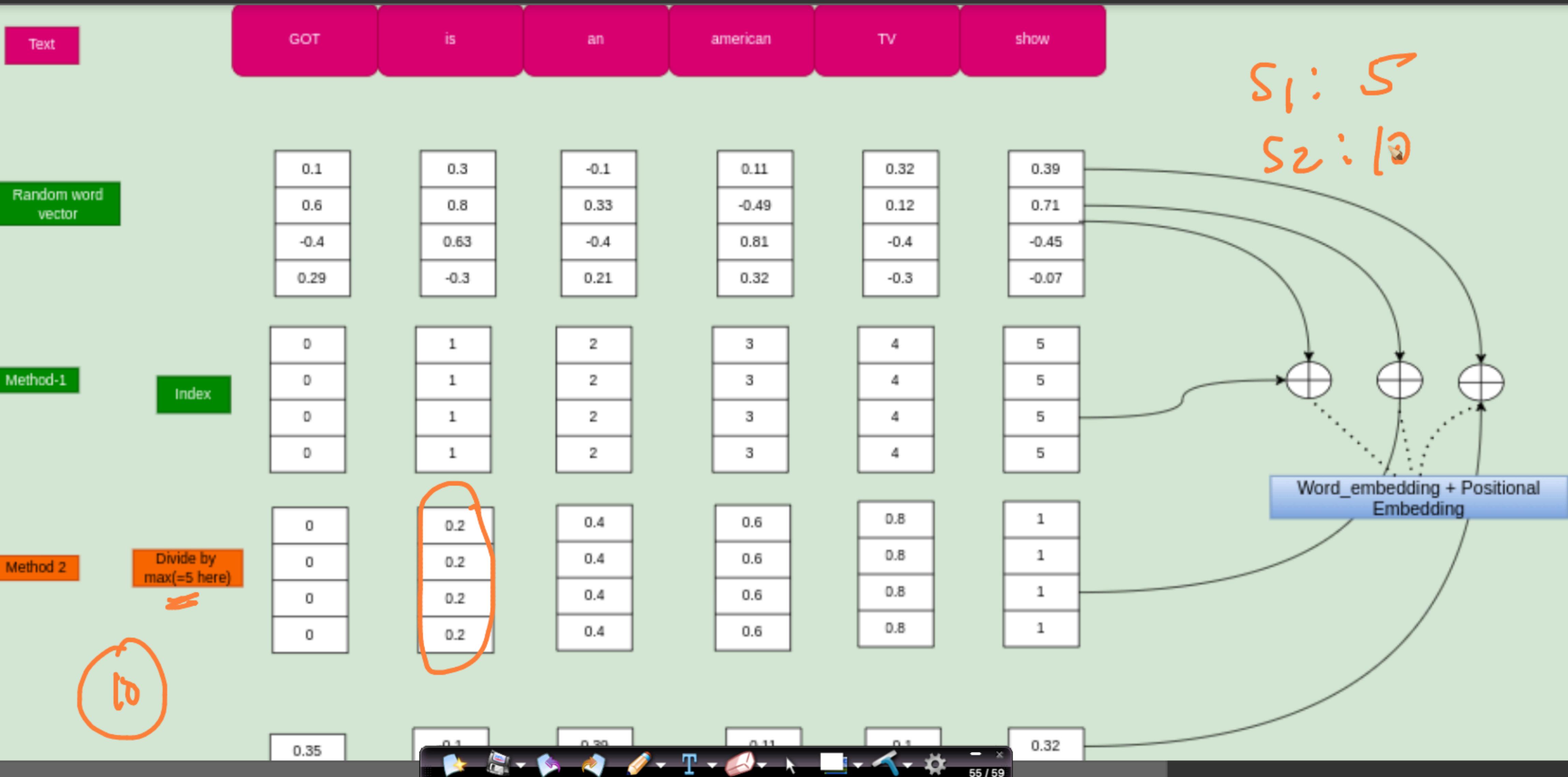
Divide by max(=5 here)

0	0.2	0.4	0.6	0.8	1
0	0.2	0.4	0.6	0.8	1
0	0.2	0.4	0.6	0.8	1
0	0.2	0.4	0.6	0.8	1

0.35 0.1 0.29 0.11 0.32 0.32

54 / 58

+ Code + Text

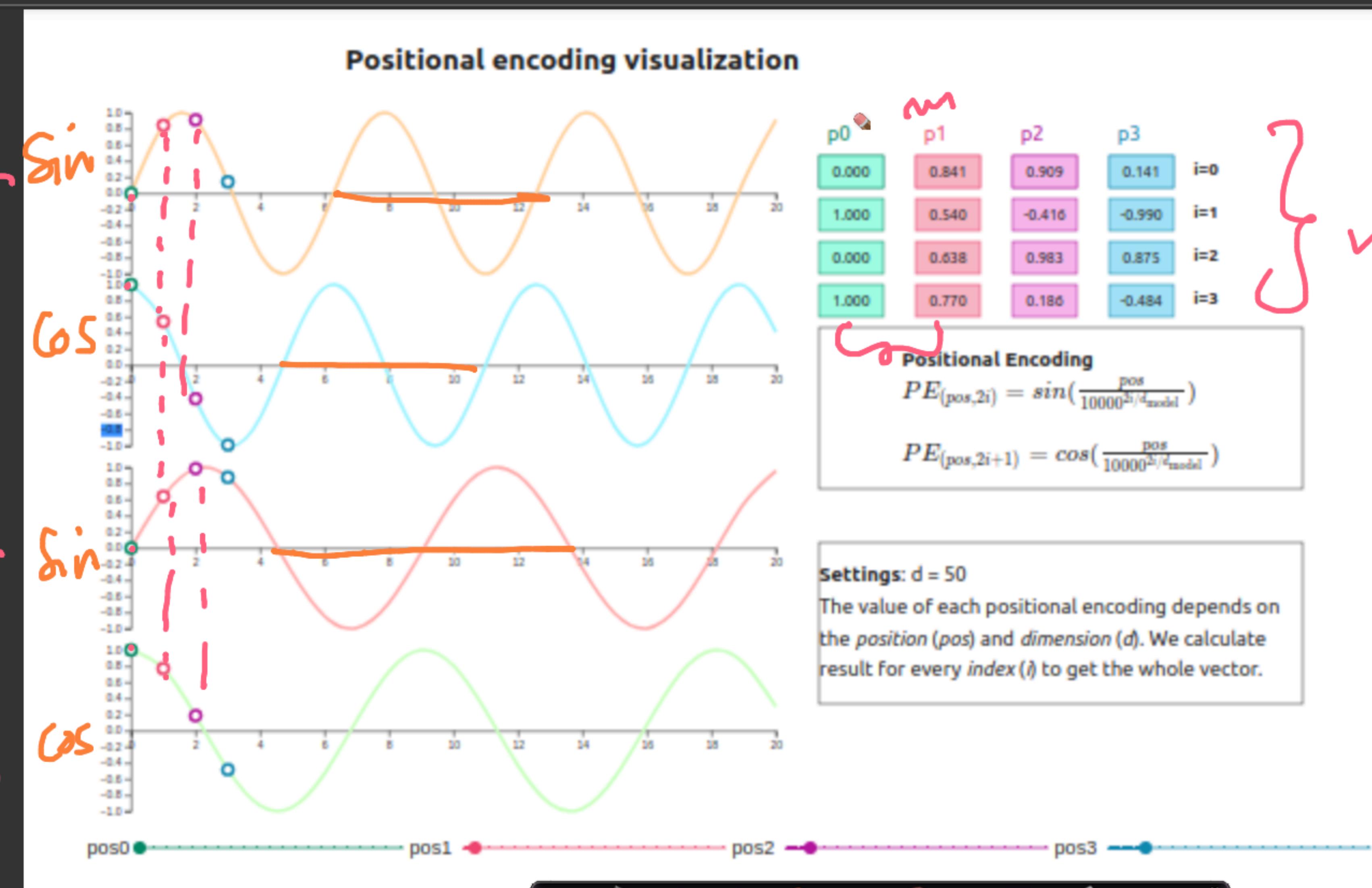
 $s_1 : S$
 $s_2 : [0]$ 

+ Code + Text

Connect



▼



+ Code + Text

Connect

Positional encoding visualization

Settings: $d = 50$
The value of each positional encoding depends on the *position (pos)* and *dimension (d)*. We calculate result for every *index (i)* to get the whole vector.

Positional Encoding

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

pos	p0	p1	p2	p3
0	0.000	0.841	0.909	0.141
1	1.000	0.540	-0.416	-0.990
2	0.000	0.638	0.983	0.875
3	1.000	0.770	0.186	-0.484

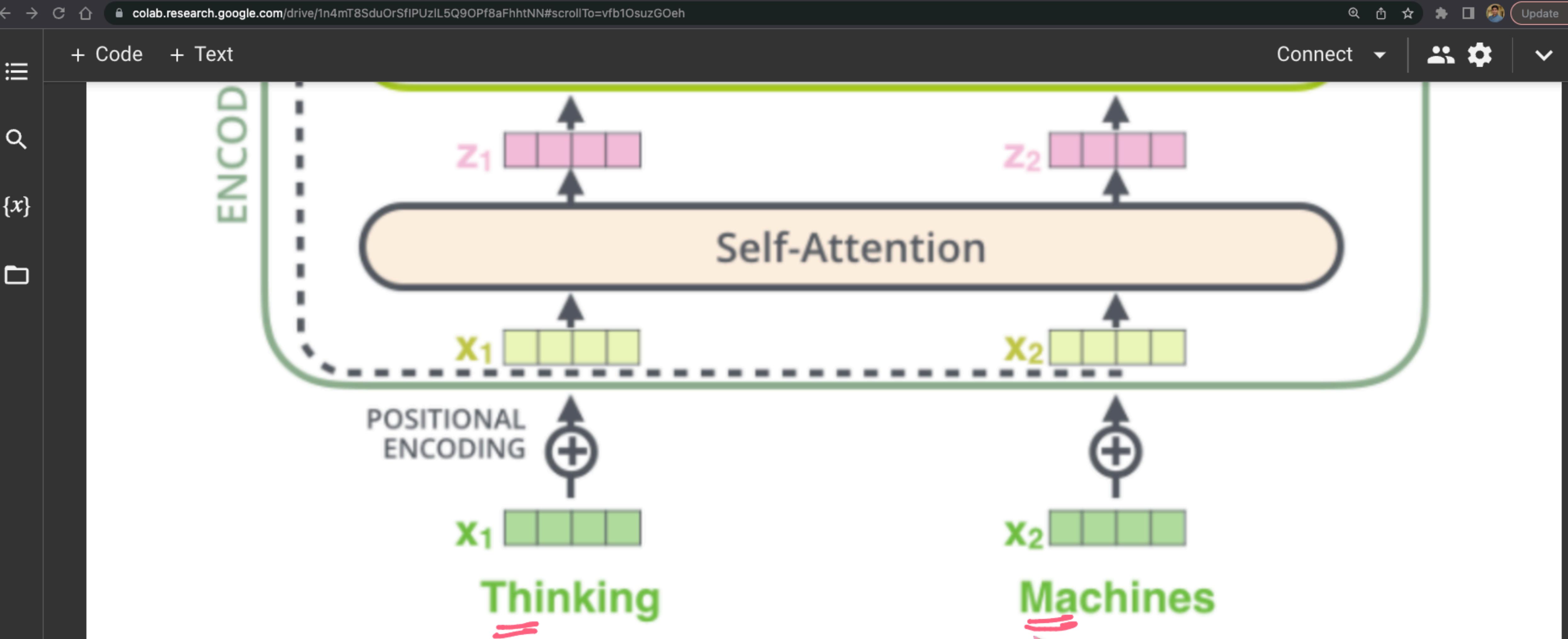
i=0
i=1
i=2
i=3

① \rightarrow pos
② \rightarrow dim
③ any i/p length
④ rel pos

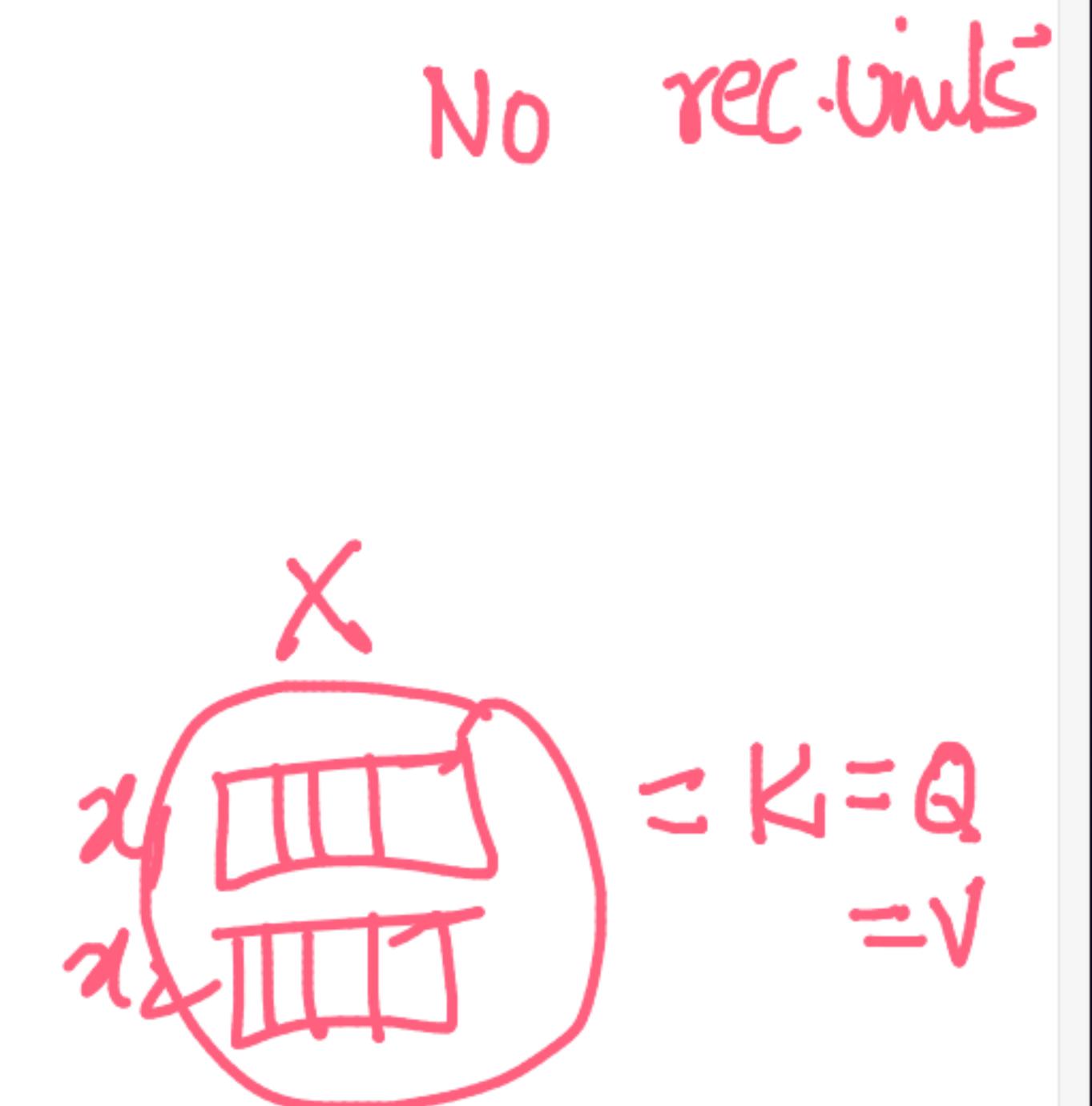
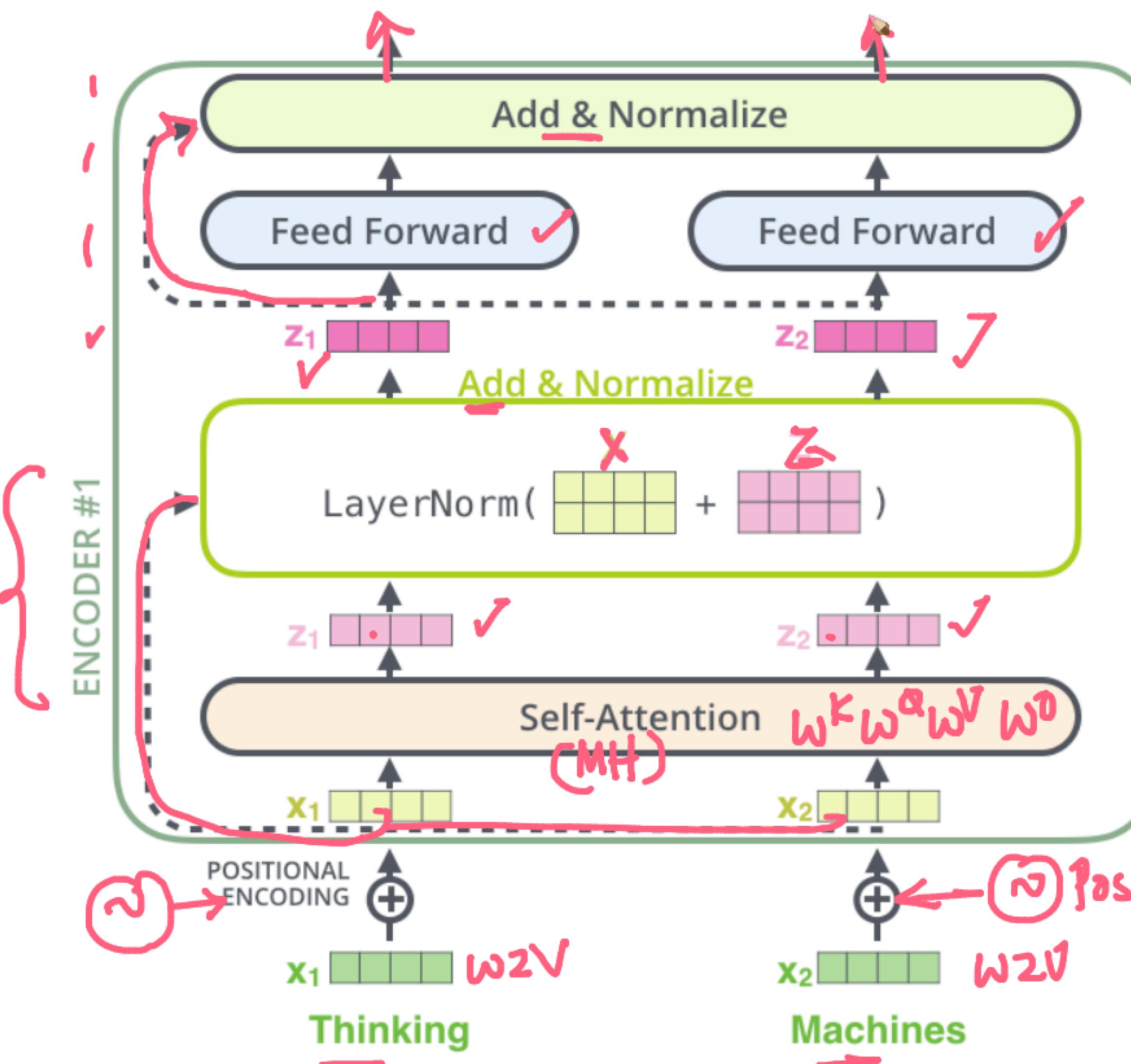
pos $\rightarrow x_1, x_2, \dots, x_i, x_{i+1}, \dots, x_n$

$w_1, w_2, w_3, w_4, w_5, w_6$

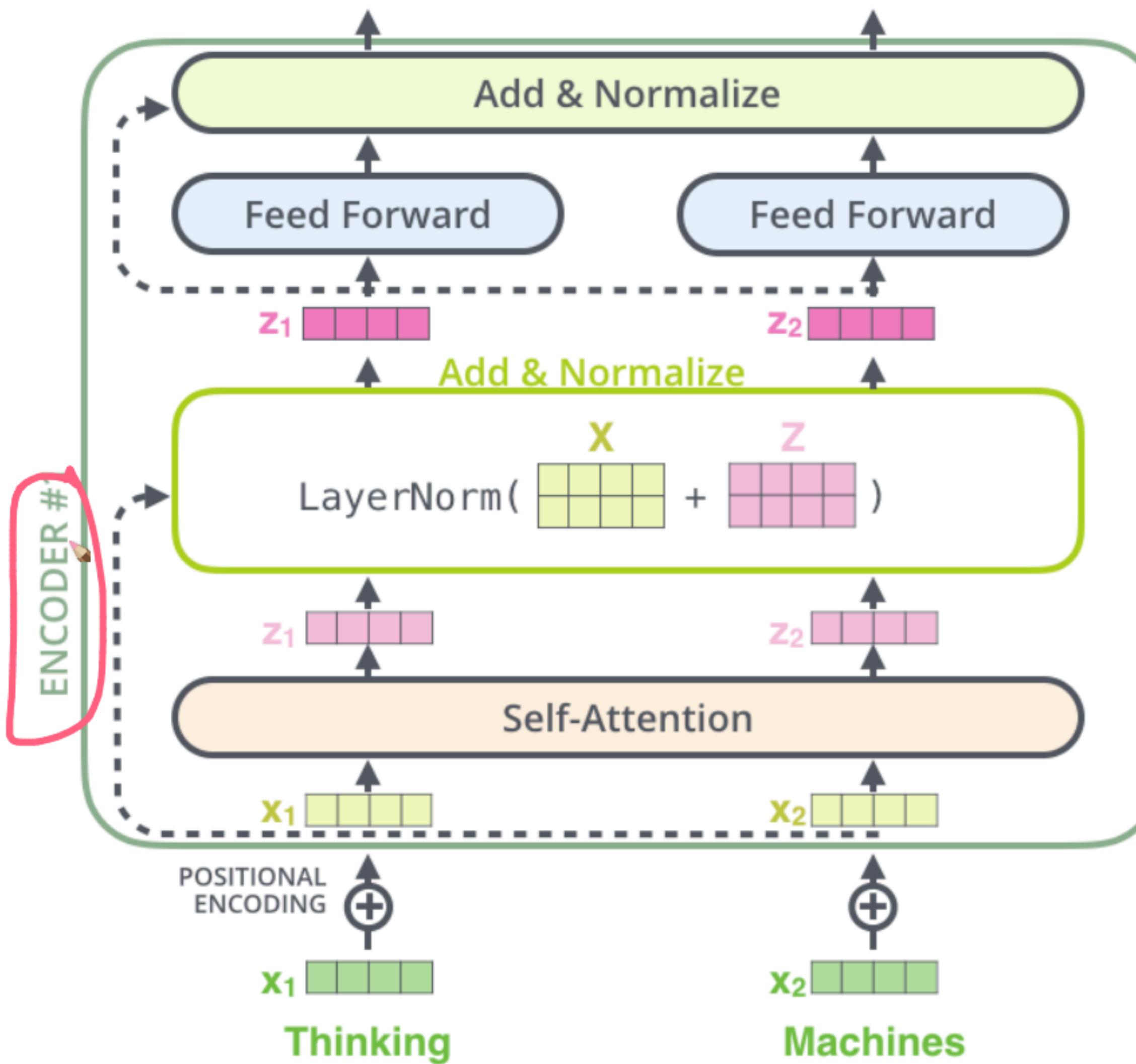
57 / 61



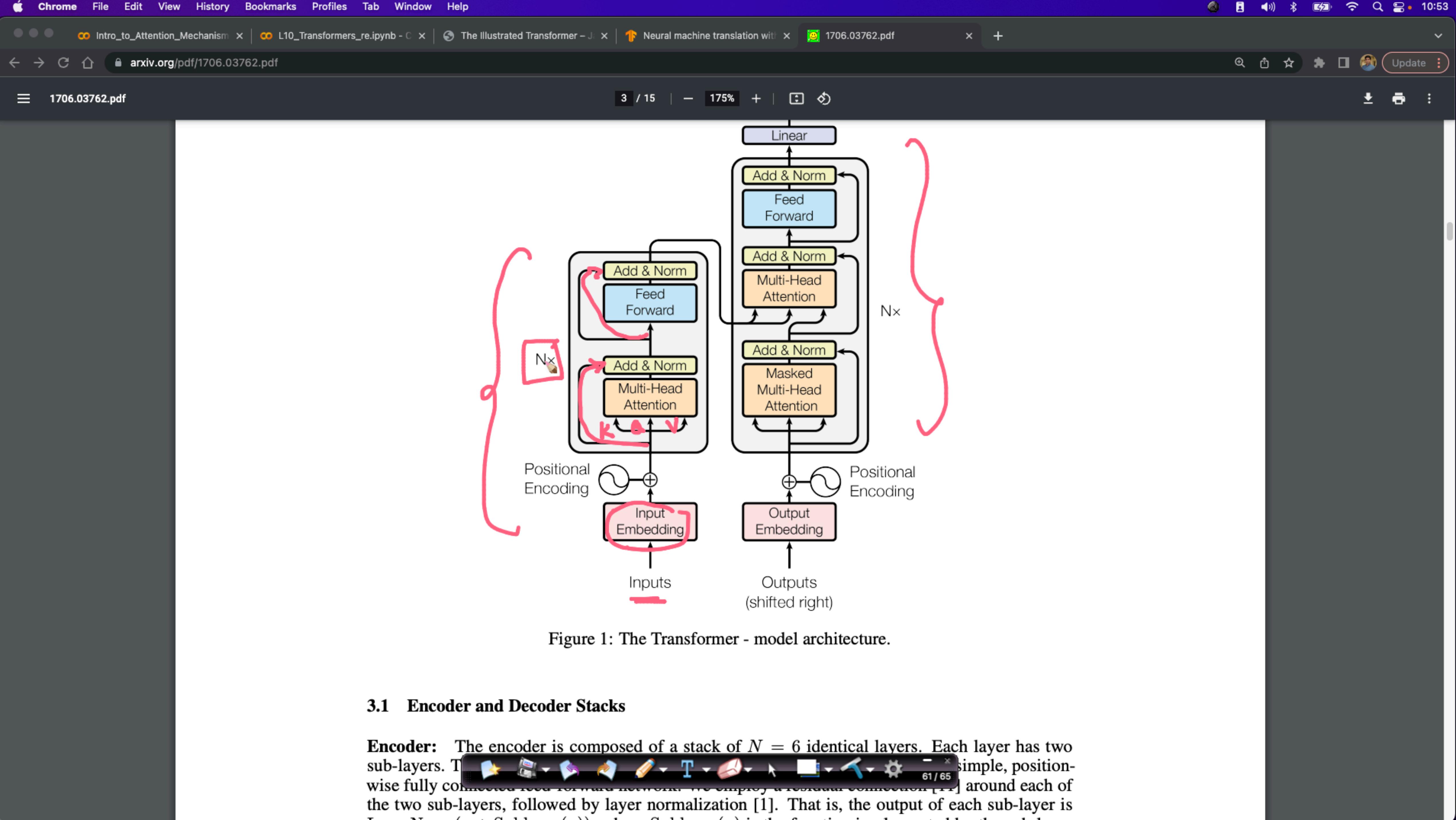
Transformer Decoder

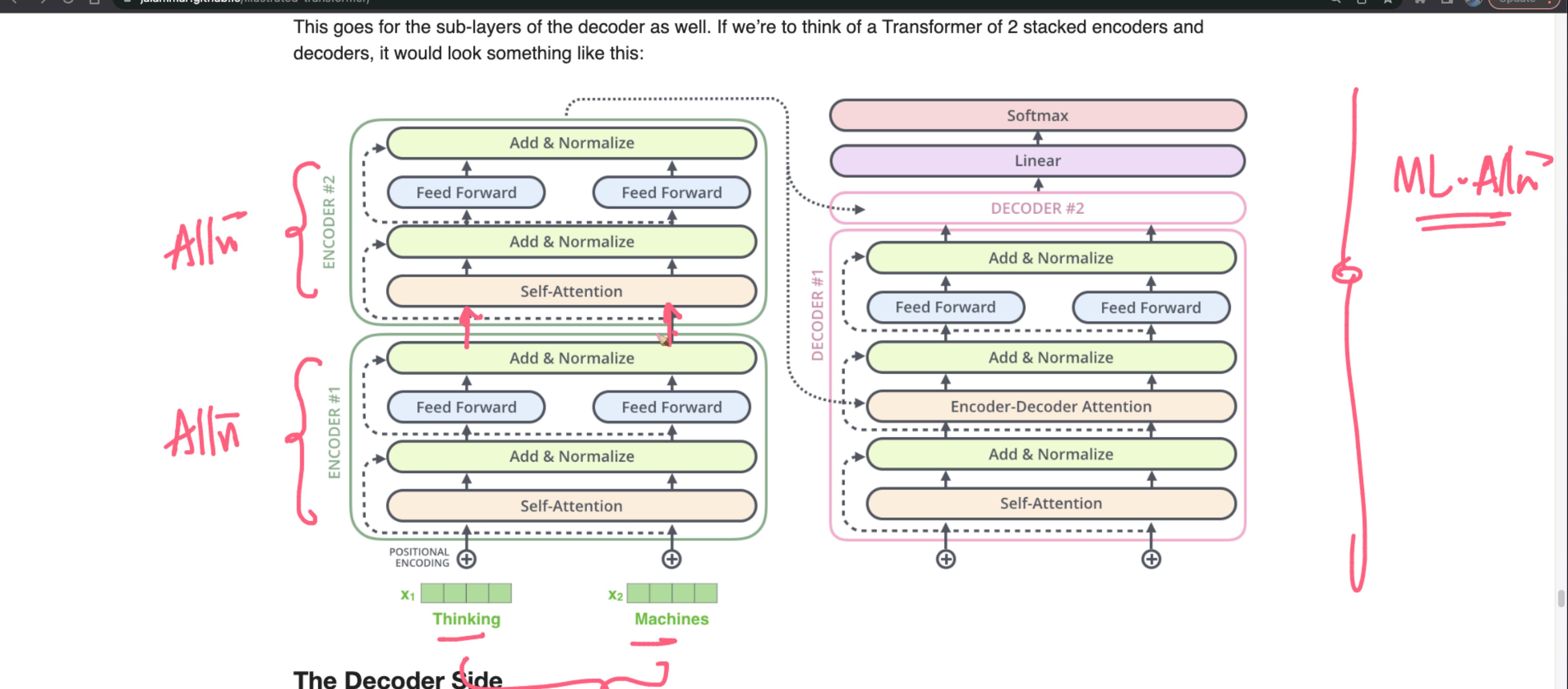


This goes for the sub-layers of the decoder as well. If we do the same for the Transformer, if stacked encoders and decoders, it would look something like this:

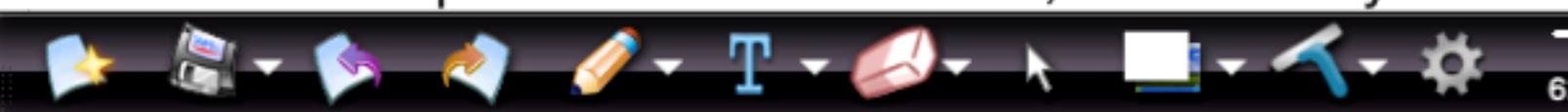


This goes for the sub-layers of the encoder as well. If we do the same for the decoder, it would look something like this:

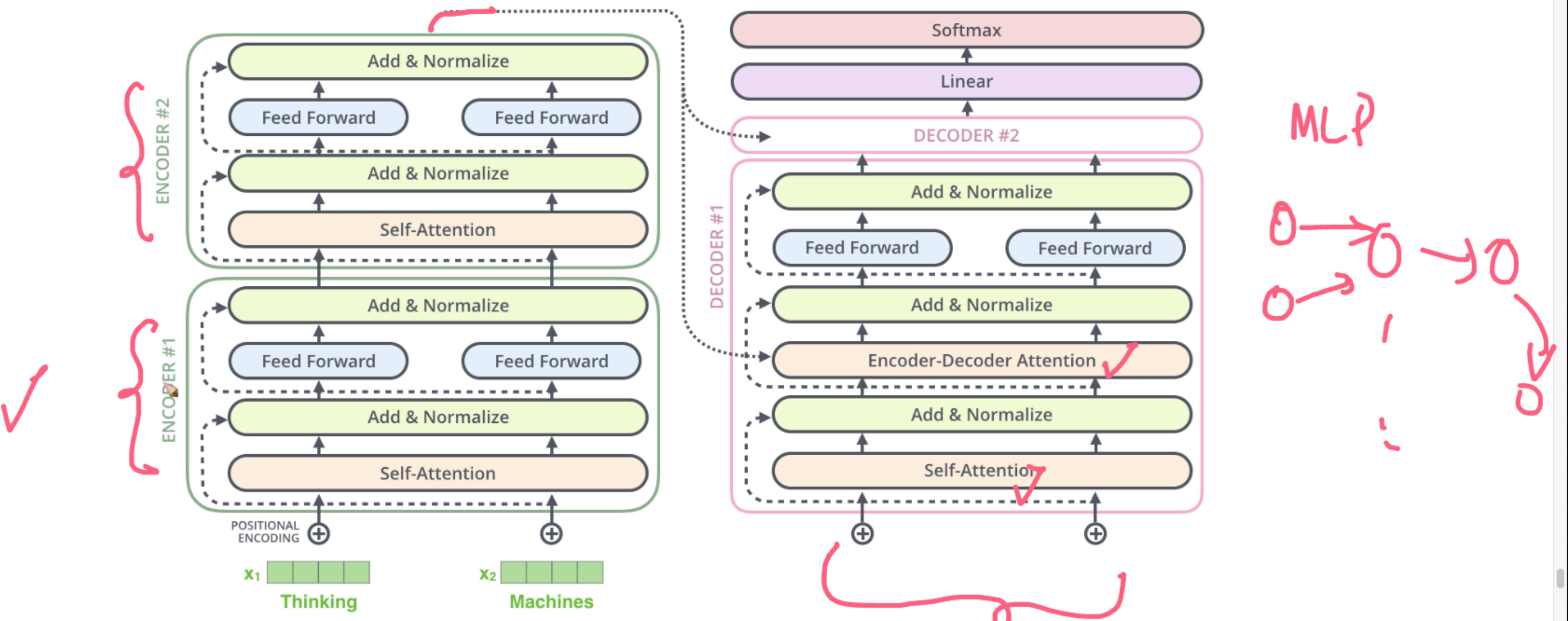




Now that we've covered most of the concepts on the encoder side, we basically know how the components of decoders work as well. B



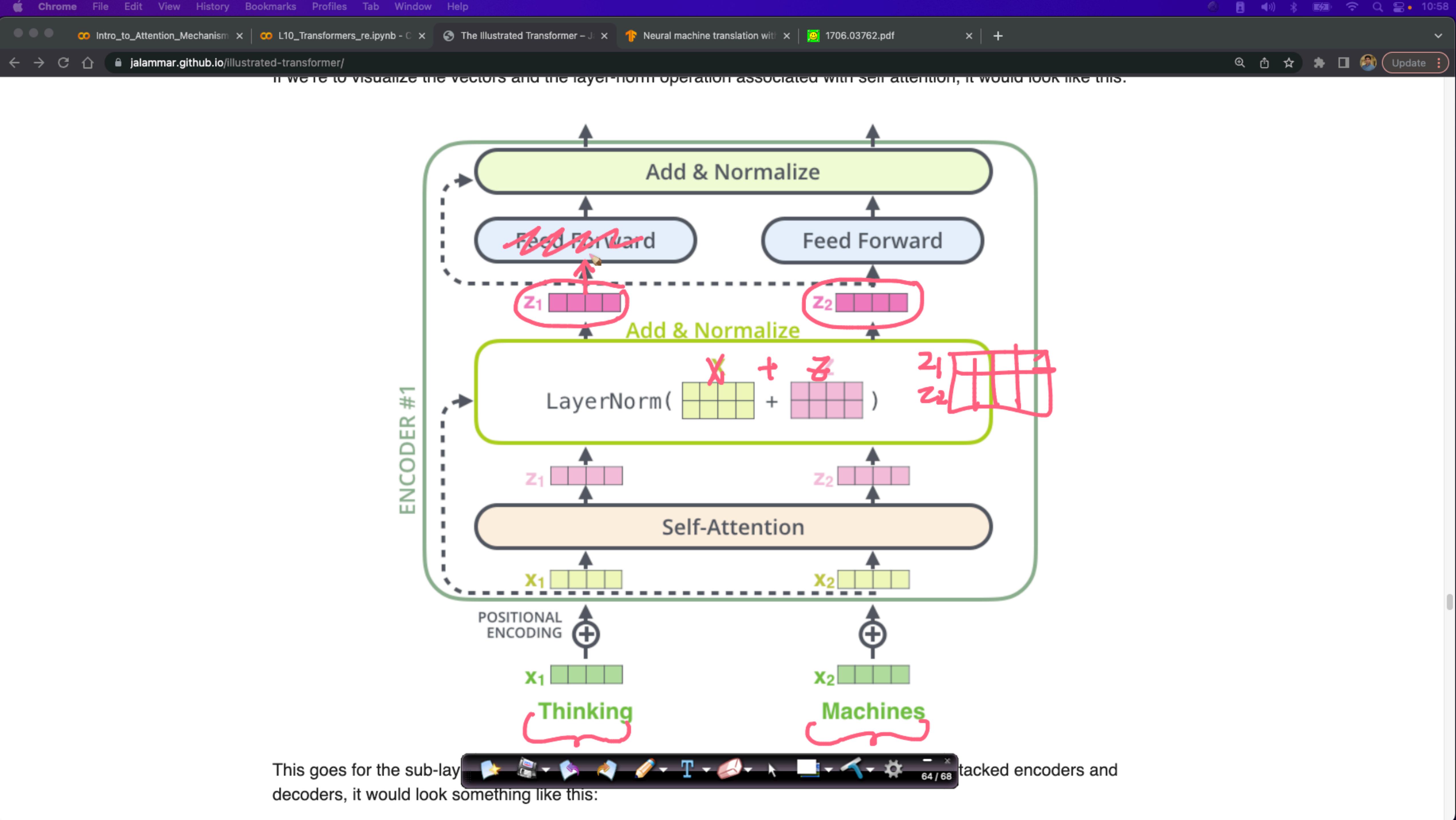
This goes for the sub-layers of the decoder as well. If we're to think of a Transformer of 2 stacked encoders and decoders, it would look something like this:

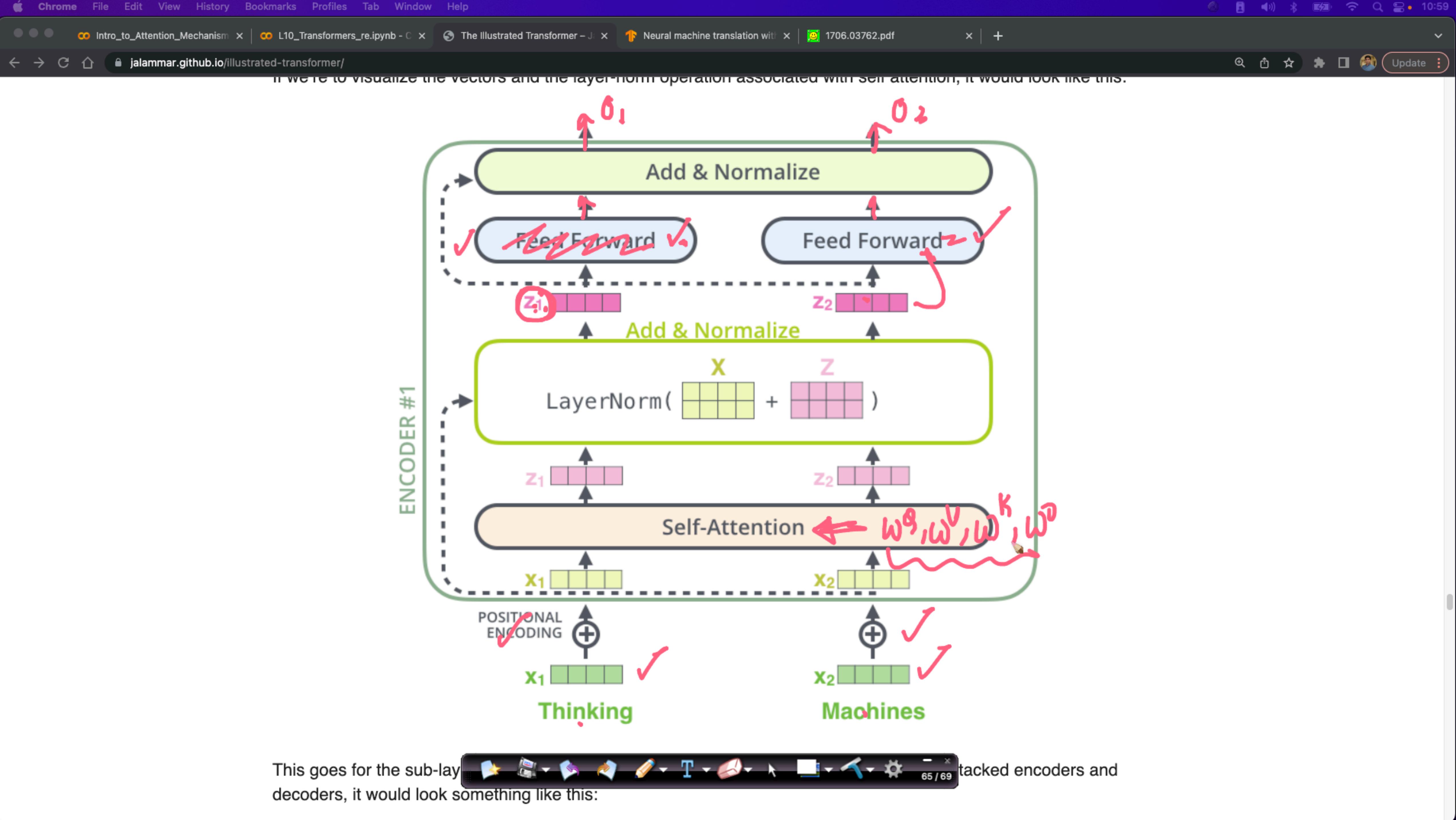


The Decoder Side

Now that we've covered most of the concepts on the encoder side, we basically know how the components of decoders work as well. Below is a toolbar with various icons for file operations.

The encoder starts by processing the input sequence. The output of the top encoder is then transformed into a set of





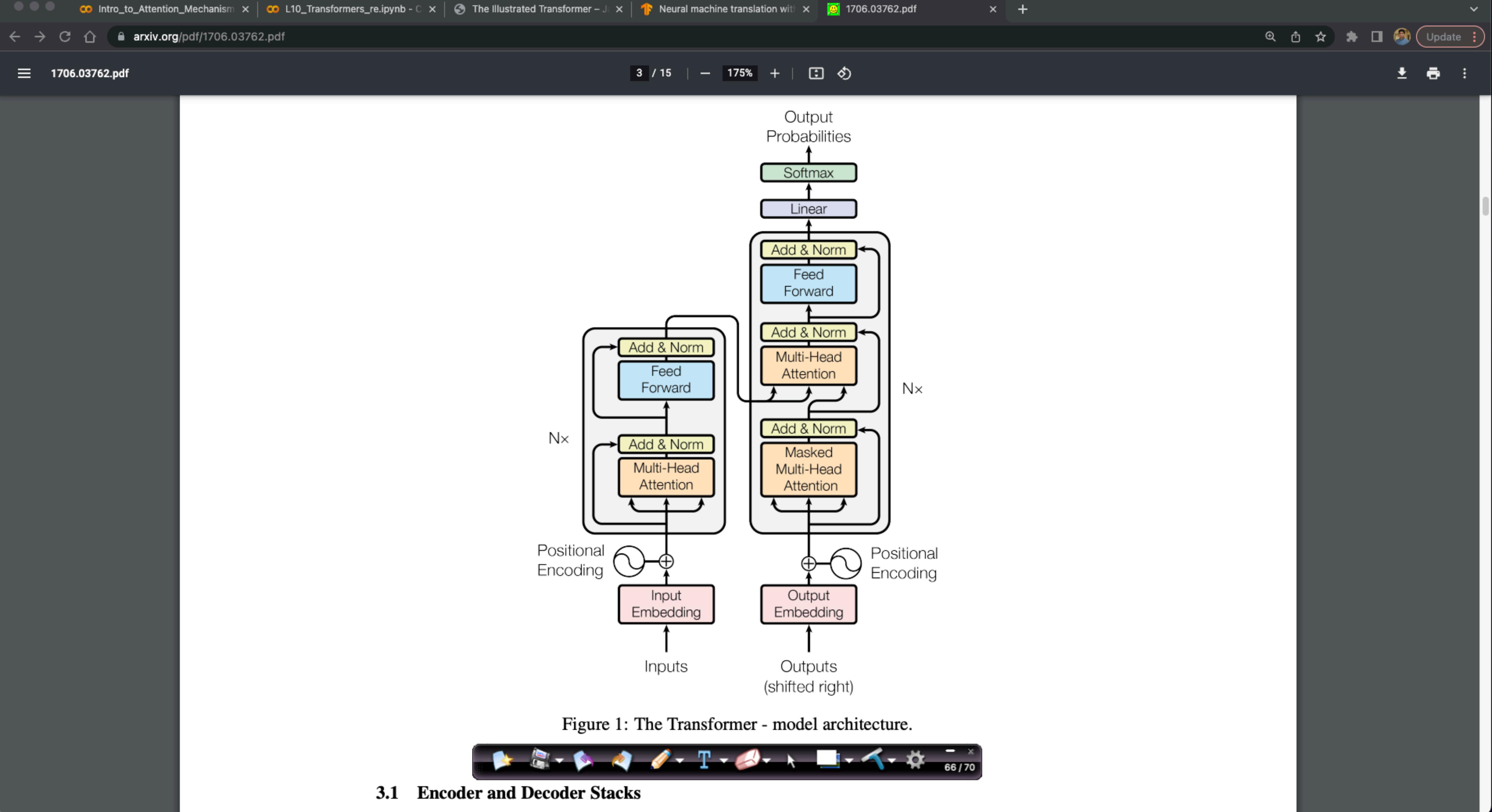


Figure 1: The Transformer - model architecture.

3.1 Encoder and Decoder Stacks