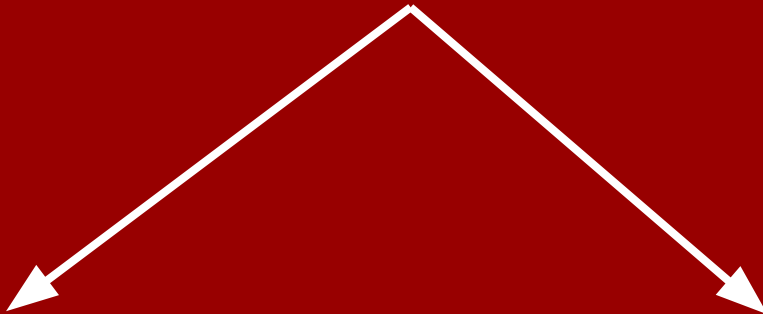


LANGUAGE MODELS + THE TRANSFORMER

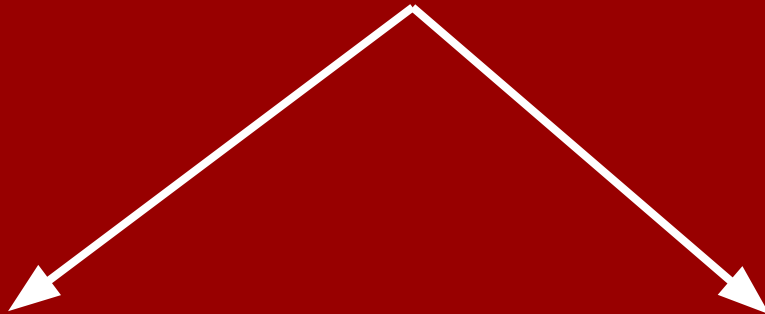
NLP

NLU

NLG



NLP



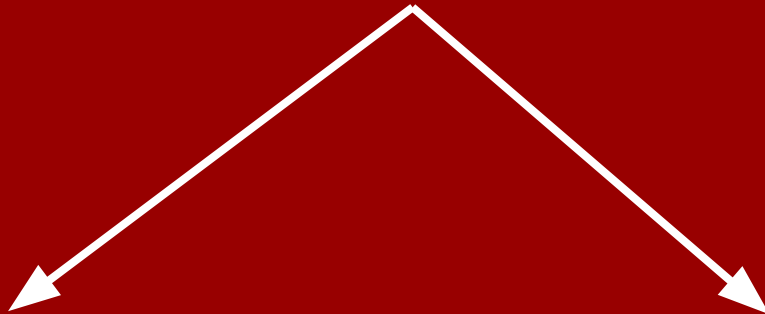
NLU

understanding

NLG

generation

NLP



NLU

understanding

NLG

generation



SYNTAX

SEMANTICS

NLP

```
graph TD; NLP --> NLU; NLP --> NLG; NLU --> SYNTAX; NLU --> SEMANTICS
```

NLU

understanding

NLG

generation

SYNTAX

SEMANTICS

**STATISTICAL
MODELING**

**WHAT EXACTLY IS A
LANGUAGE MODEL??**

LANGUAGE REPRESENTATION

LANGUAGE REPRESENTATION

```
graph TD; A[LANGUAGE REPRESENTATION] --> B[FEATURE EXTRACTION]; A --> C[EMBEDDING];
```

**FEATURE
EXTRACTION**

Classical ML

EMBEDDING

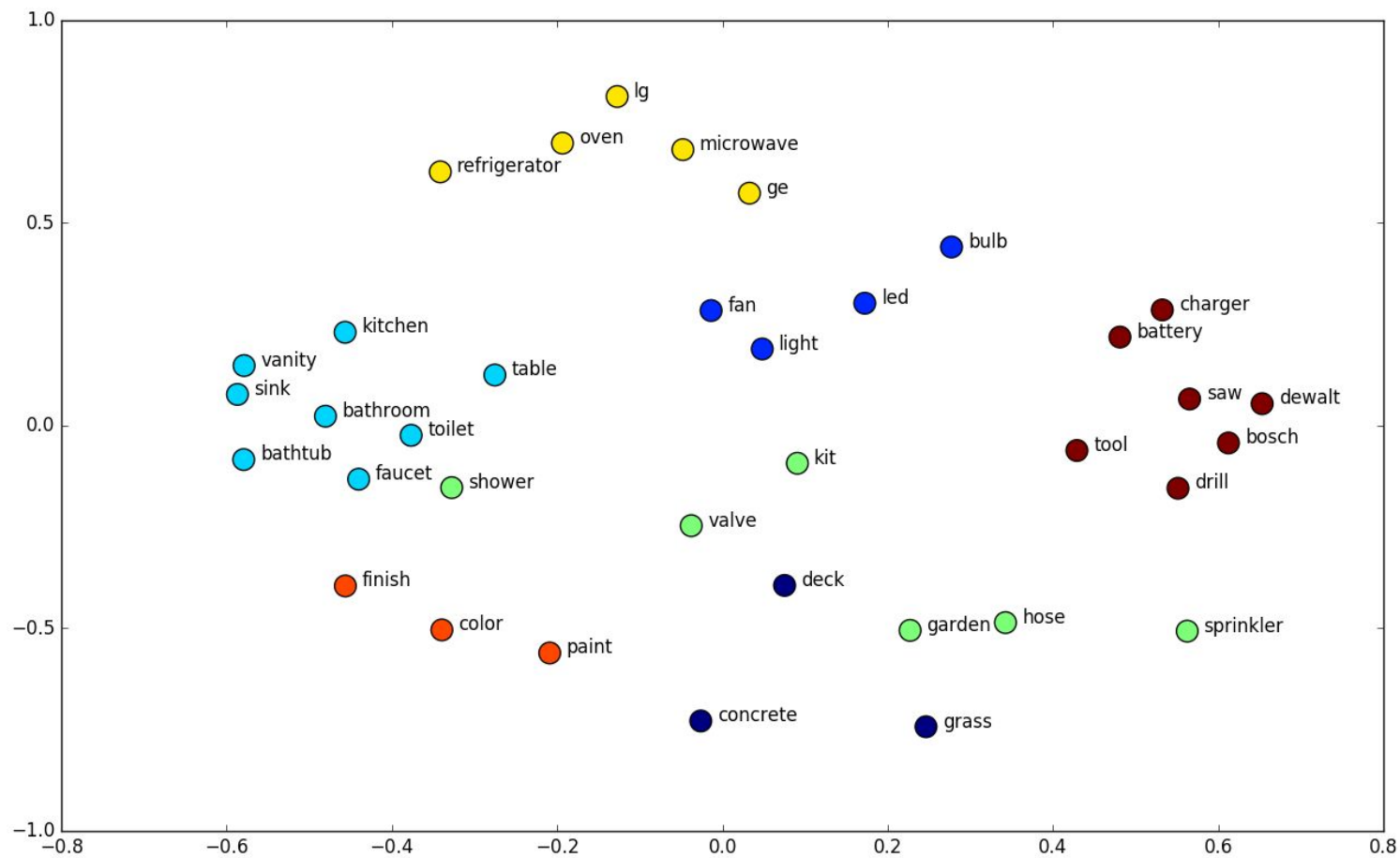
Deep & LLMs

<i>cat</i> →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<i>kitten</i> →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<i>dog</i> →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<i>houses</i> →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8
<i>man</i> →	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i> →	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
<i>king</i> →	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i> →	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

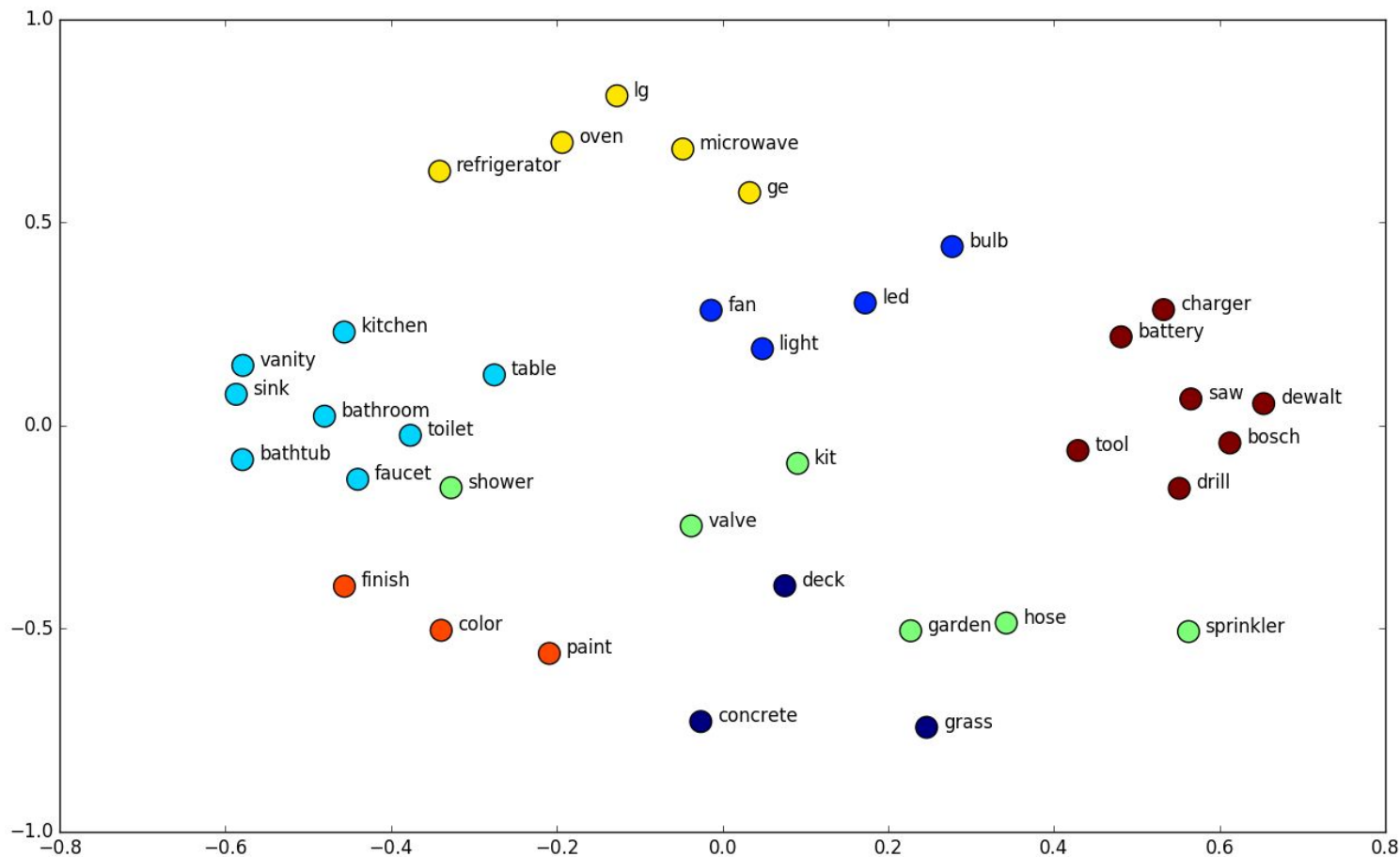
WORD EMBEDDINGS

Representing words by a
vector of numbers

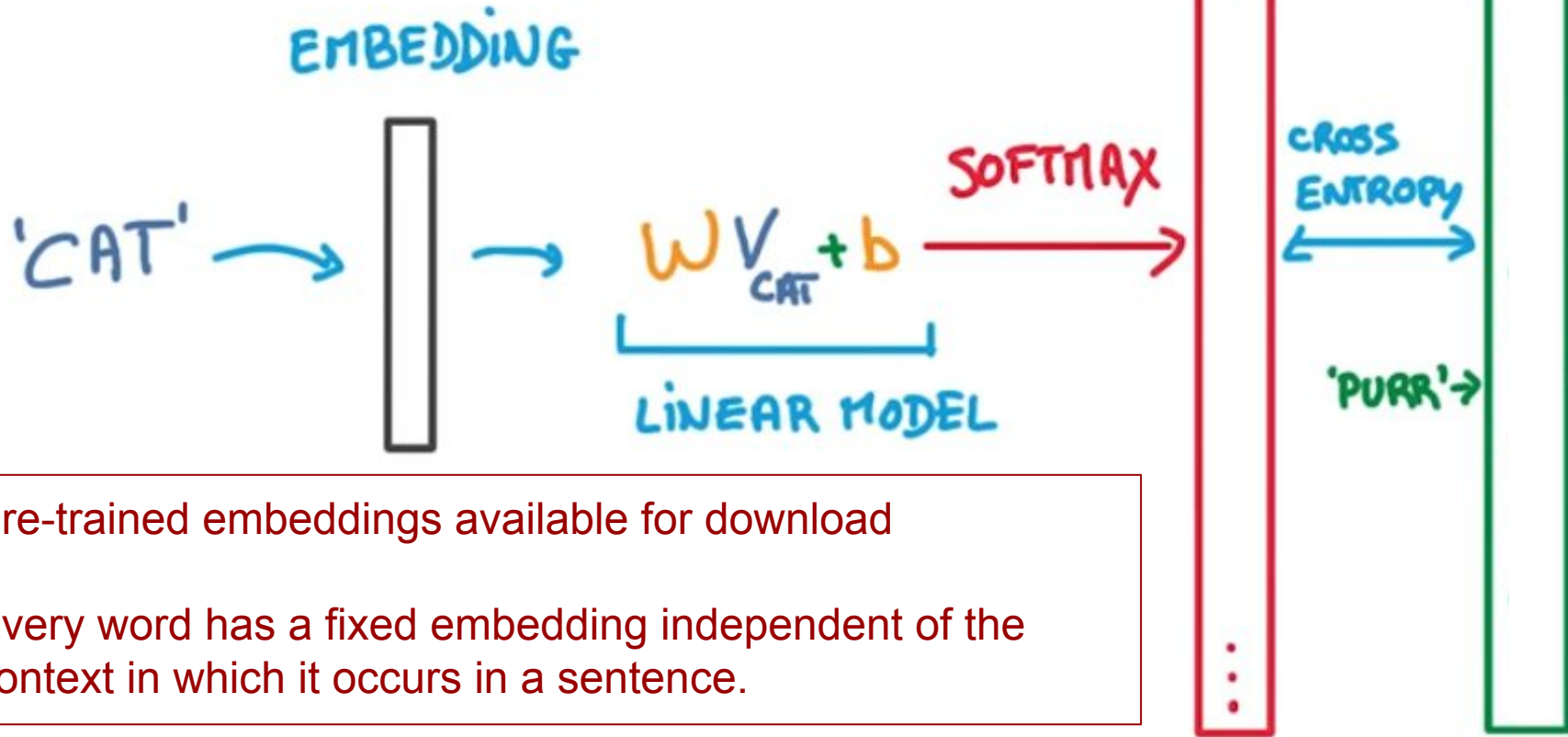
How long should
these vectors be?



How do we know that the vector representations are correct?

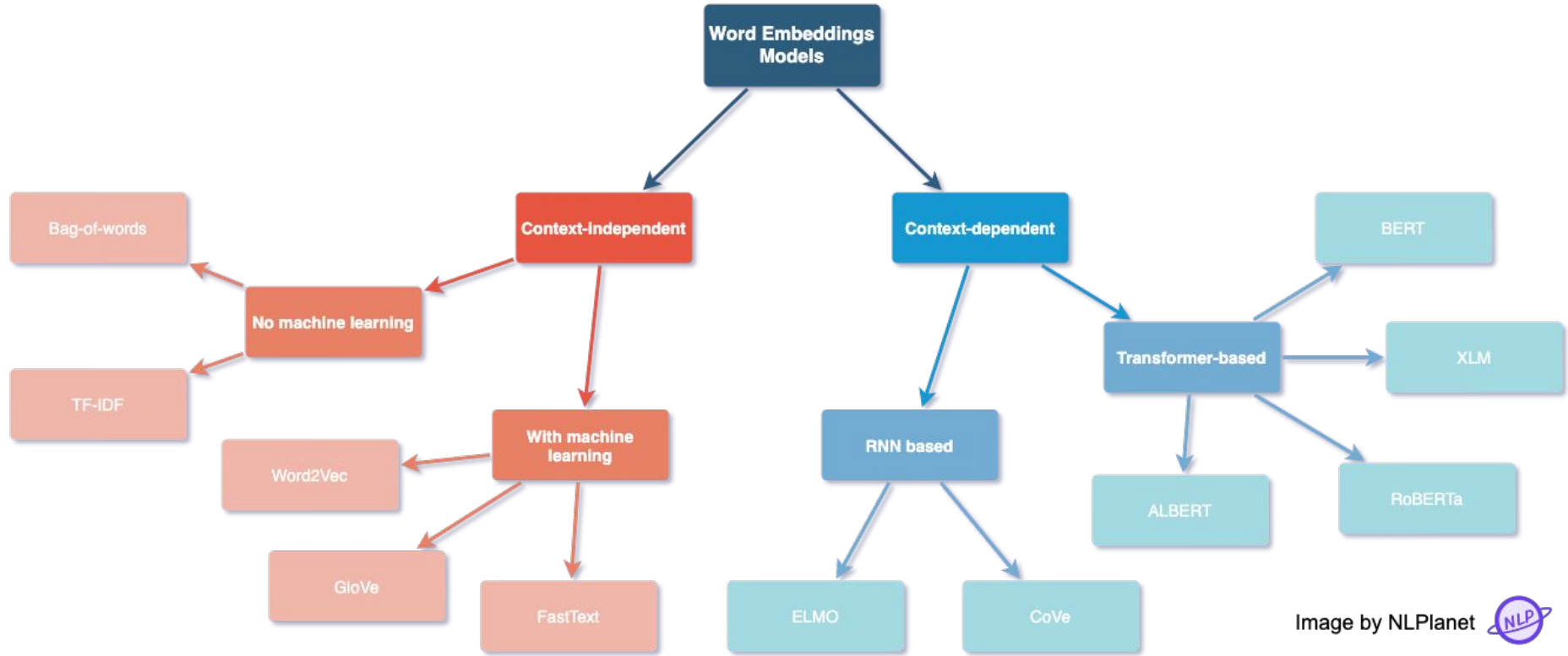


WORD 2 VEC.



Pre-trained embeddings available for download

Every word has a fixed embedding independent of the context in which it occurs in a sentence.





THE TRANSFORMER MODEL



What's the difference between these two devices in terms of how they treat the incoming information and data?



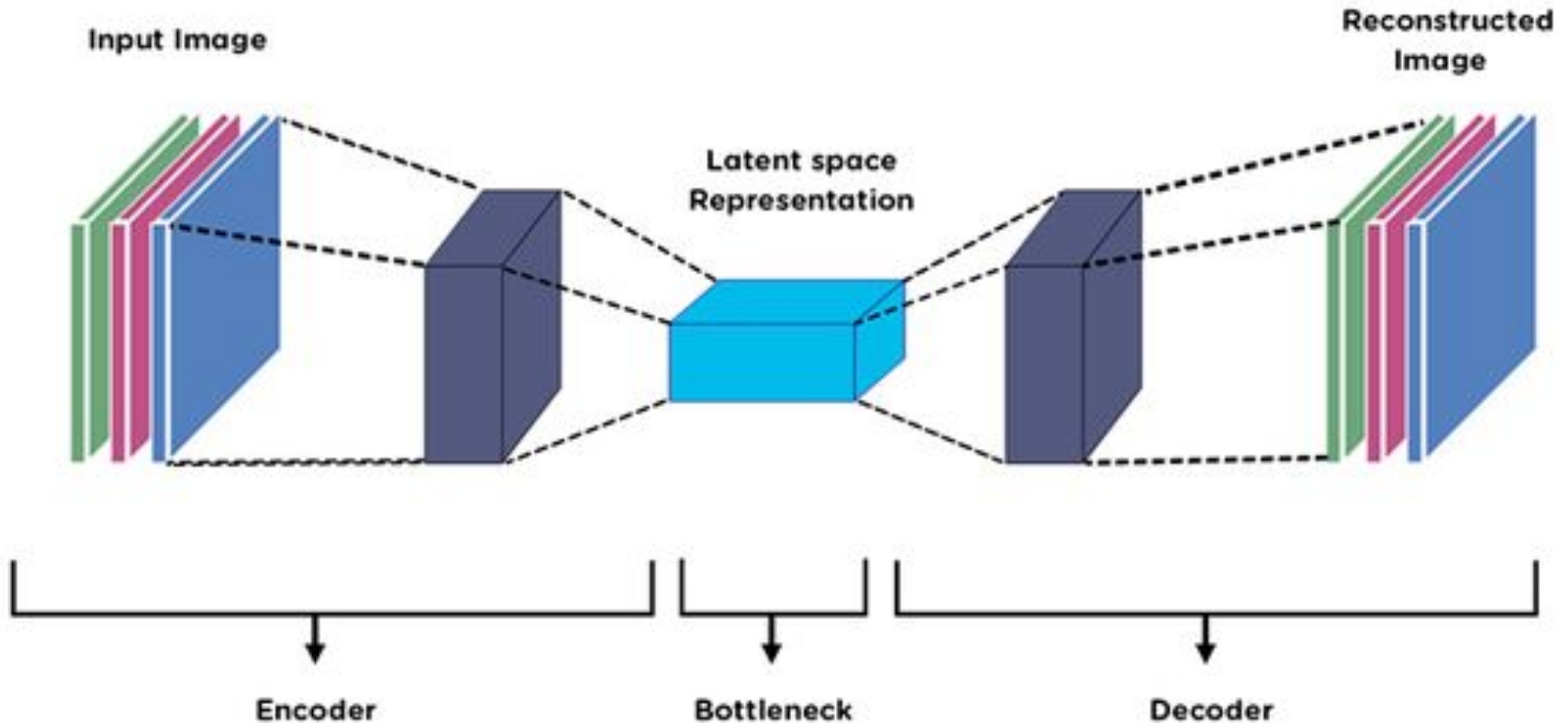
Why is the one on the left considered to be intelligent,
and the one on the right considered to be dumb?



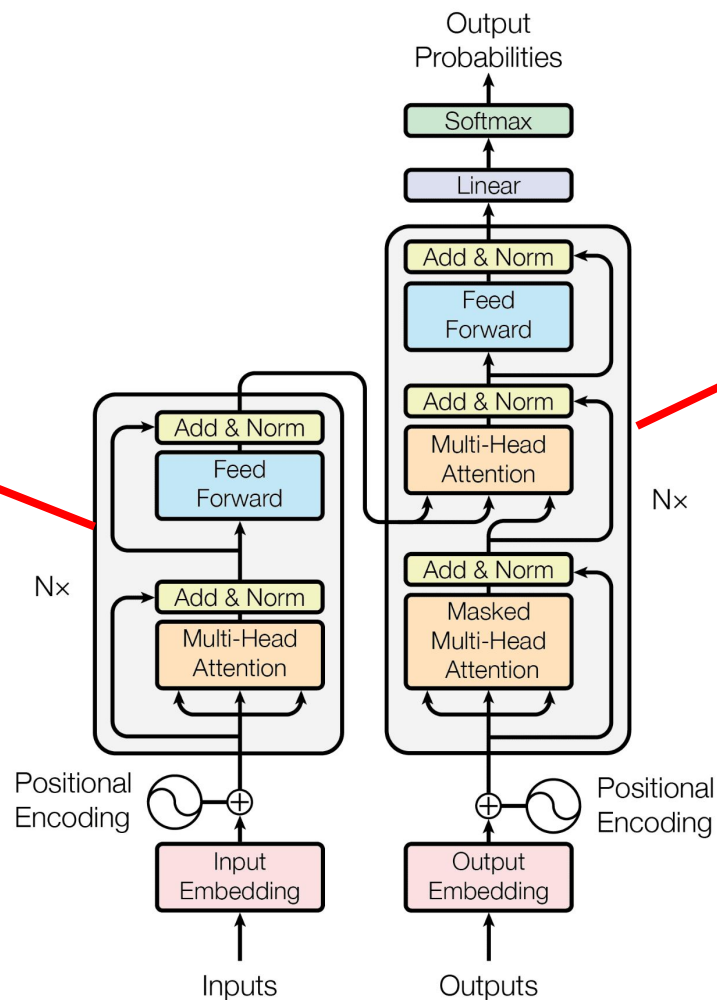
Intelligence is about being able to figure out the **essence** of a topic,
and not just memorizing facts.

AUTO-ENCODERS

[lossy compression]



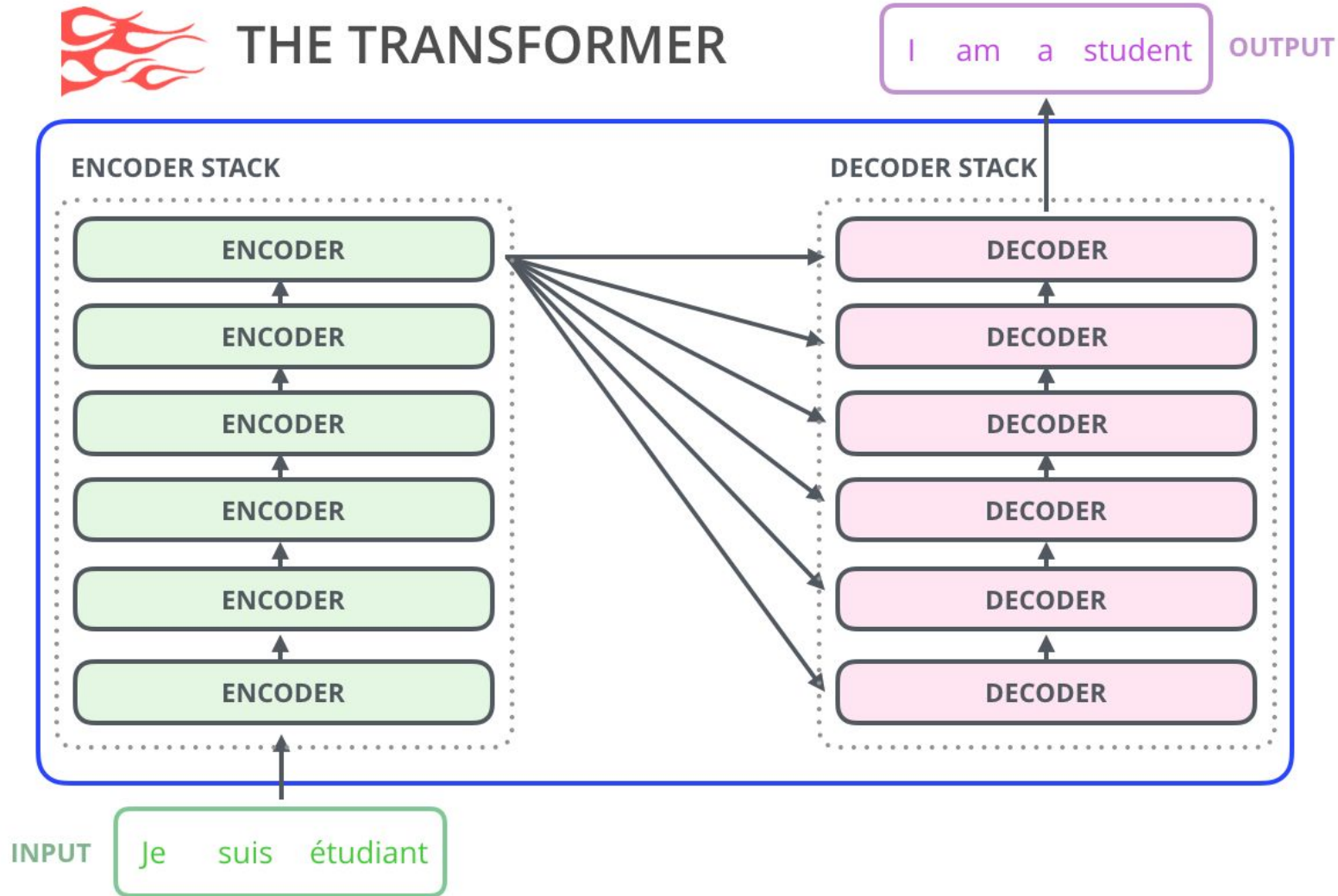
Encoder



Decoder



THE TRANSFORMER



" I like strawberries ", 3 words

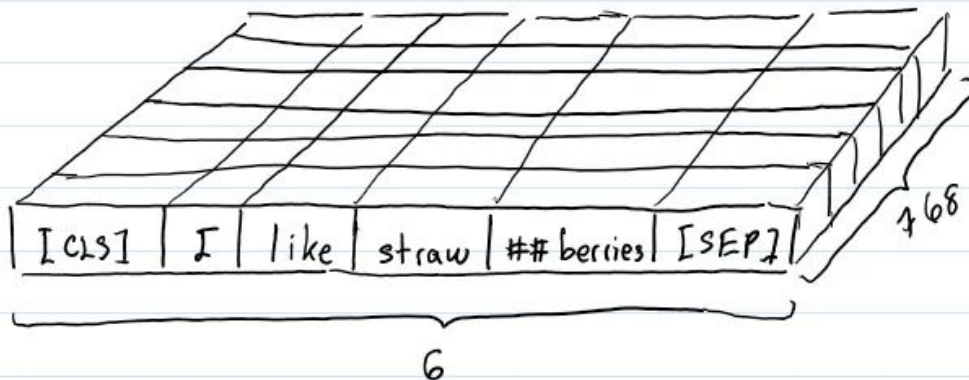
↓ ①

"[CLS]", "I", "like", "straw", "##berries", "[SEP]", 6 tokens

↓ ②



↓ result



**Non-contextual
Token Embeddings**

" I like strawberries ", 3 words

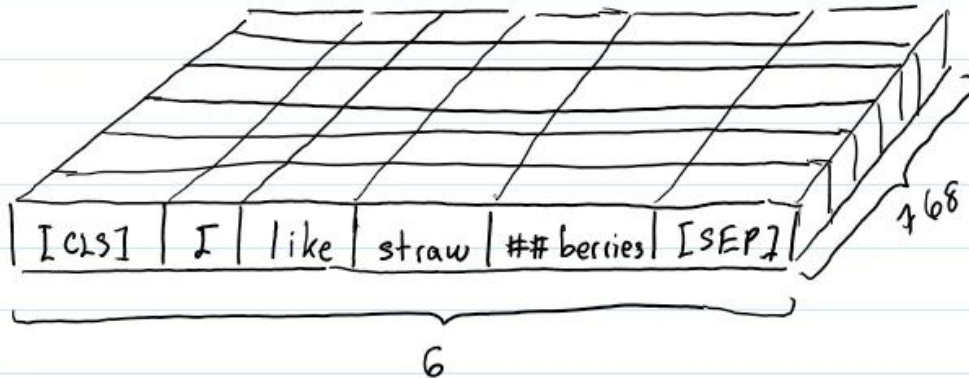
↓ ①

"[CLS]", "I", "like", "straw", "##berries", "[SEP]", 6 tokens

↓ ②



↓ result



**Non-contextual
Token Embeddings**



THE TRANSFORMER

ENCODER BLOCK

Feed Forward Neural Network

Self-Attention

robot

must

obey

orders

<eos>

<pad>

...

<pad>

1

2

3

4

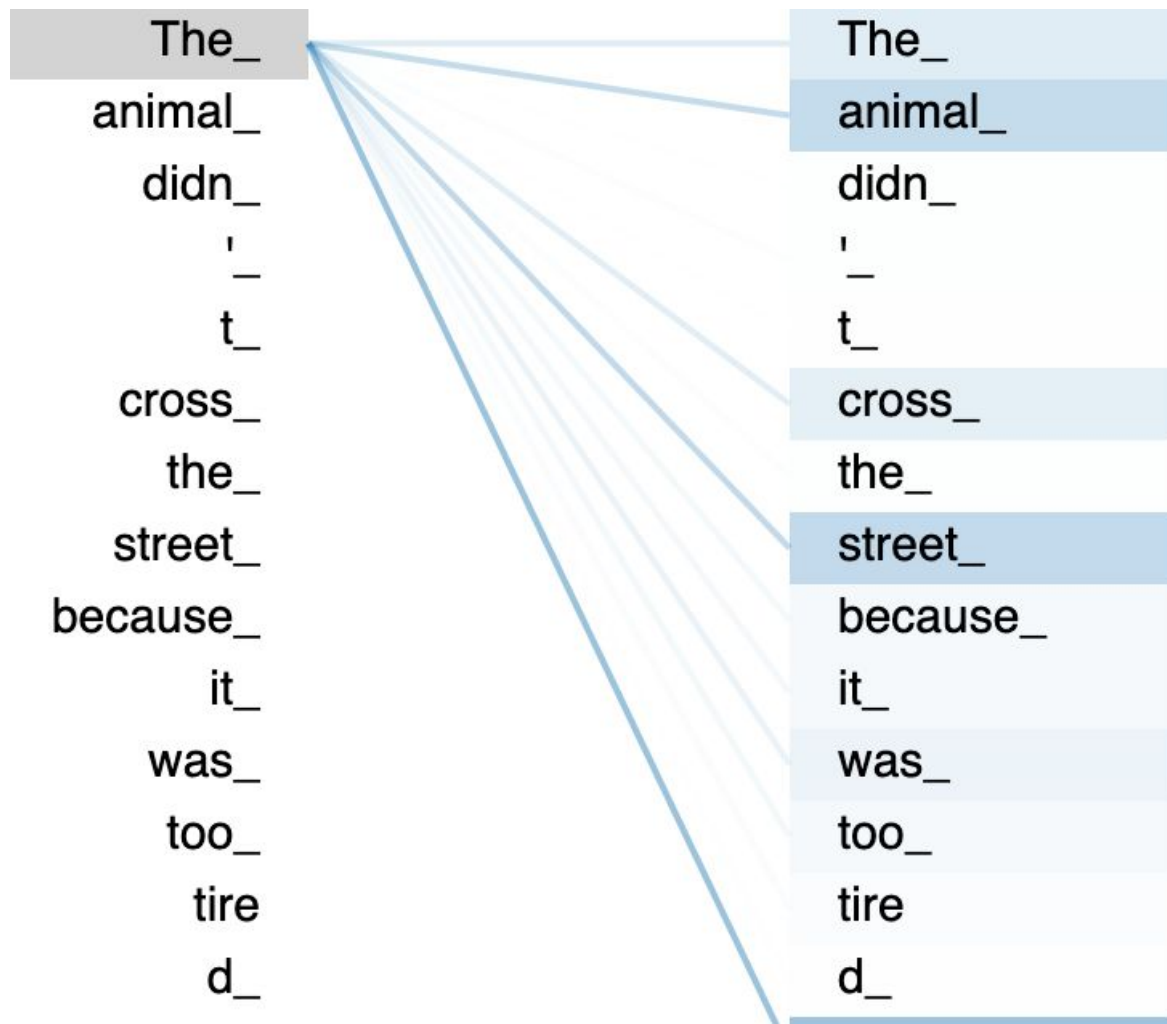
5

6

512

The Attention Mechanism

Taking care of the word sequence is important, but there are also long range dependencies between words.

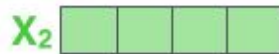
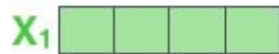


Input

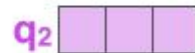
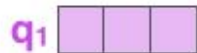
Thinking

Machines

Embedding

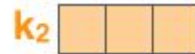
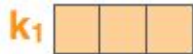


Queries



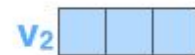
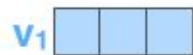
W^Q

Keys



W^K

Values



W^V

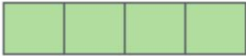
Input

Thinking

Machines

Embedding

x_1 

x_2 

Queries

q_1 

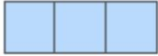
q_2 

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

Keys



Values



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



Sum





THE TRANSFORMER

DECODER BLOCK

Feed Forward Neural Network

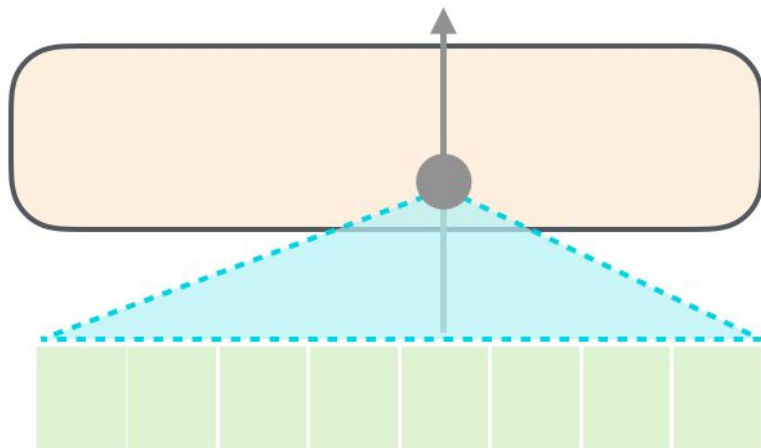
Encoder-Decoder Self-Attention

Masked Self-Attention

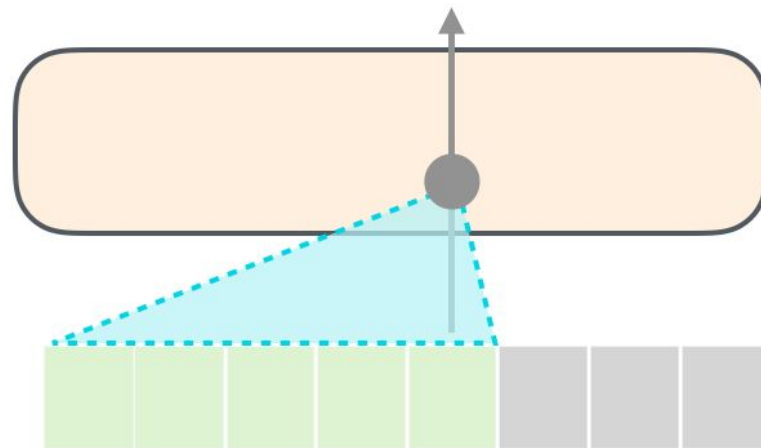
Input



Self-Attention



Masked Self-Attention





THE TRANSFORMER

DECODER BLOCK

Feed Forward Neural Network

Encoder-Decoder Self-Attention

Masked Self-Attention

Input





THE TRANSFORMER

Important for tasks
like translation

DECODER BLOCK

Feed Forward Neural Network

Encoder-Decoder Self-Attention

Masked Self-Attention

Input

<s>

robot

must

obey

1

2

3

4

5

6

512



GPT

DECODER

...

DECODER

DECODER



BERT

ENCODER

...

ENCODER

ENCODER

**WHAT EXACTLY IS A
LANGUAGE MODEL?**

HOW ARE TRANSFORMER EMBEDDINGS
DIFFERENT FROM WORD2VEC?

HOW ARE TRANSFORMER EMBEDDINGS
DIFFERENT FROM ELMo EMBEDDINGS?

**WHAT EXACTLY IS AN
AUTO-ENCODER?**

**IN THE TRANSFORMER MODEL,
WHAT DOES AN ENCODER DO?**

**IN THE TRANSFORMER MODEL,
WHAT DOES A DECODER DO?**

**WHAT DOES THE
ENCODER BLOCK CONTAIN?**

**WHAT DOES THE
DECODER BLOCK CONTAIN?**

HOW IS THE SELF-ATTENTION OF
DECODER BLOCK DIFFERENT FROM
THAT OF ENCODER BLOCK?

**WHY DOES NOT THE GPT MODEL
HAVE ANY ENCODER BLOCK?**

**WHY DOES NOT THE BERT MODEL
HAVE ANY DECODER BLOCK?**

" I like strawberries ", 3 words

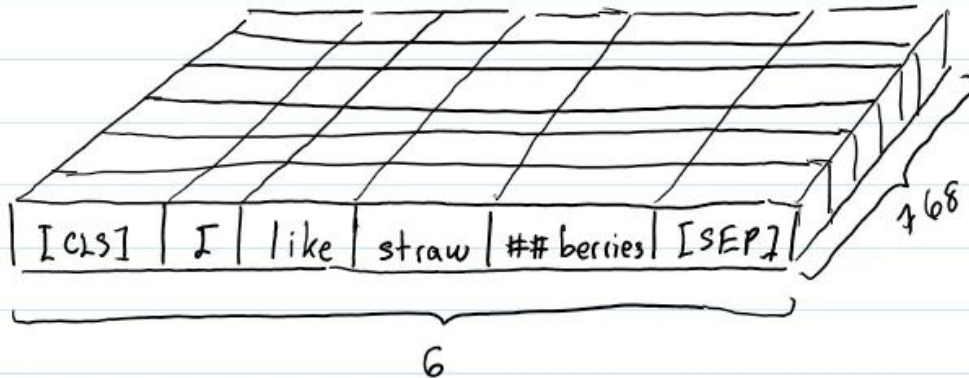
↓ ①

"[CLS]", "I", "like", "straw", "##berries", "[SEP]", 6 tokens

↓ ②



↓ result



How does
tokenization work?

SUB-WORD TOKENIZERS for TRANSFORMERS

Tokenizer	By	Used In	Merge Criteria	Advantage
WordPiece	Google	BERT	Normalized Score	More context
Byte Pair Encoding (BPE)	Philip Gage	GPT	Sub-word frequency	Faster training
SentencePiece	Google	Llama, XLNet, T5, PaLM	Same as BPE	Language independent

Corpus

1 huggingface
1 hugging
1 face
1 hug
1 hugger
2 learning
2 learner
2 learners
1 learn

Splits

h ##u ##g ##g ##i ##n ##g ##f ##a ##c ##e
h ##u ##g ##g ##i ##n ##g
f ##a ##c ##e
h ##u ##g
h ##u ##g ##g ##e ##r
l ##e ##a ##r ##n ##i ##n ##g
l ##e ##a ##r ##n ##e ##r
l ##e ##a ##r ##n ##e ##r ##s
l ##e ##a ##r ##n

Vocab

##a	##s	
##c	##u	fa
##e	f	fac
##f	h	hug
##g	l	##gfac
##i	hu	hugg
##n	##fa	huggi
##r	##fac	

Vocab

##a	##s	
##c	##u	fa
##e	f	fac
##f	h	hug
##g	l	##gfac
##i	hu	hugg
##n	##fa	huggi
##r	##fac	

h u g g i n g f a c e

Vocab

##a

##c

##e

##f

##g

##i

##n

##r

##s

##u

f

h

l

hu

##fa

##fac

fa

fac

hug

##gfac

hugg

huggi

h u g g i

##n

##g f a c e