





Contextual Models: Transformers

Luciano Barbosa

Traditional Word Embedding

- Vocabulary built using the training set
- Words not present mapped to UNK

	word		vocab mapping	embedding
Common words	hat	→	pizza (index)	
	learn	→	tasty (index)	
Variations	taaaaasty	→	UNK (index)	
misspellings	laern	→	UNK (index)	
novel items	Transformerify	→	UNK (index)	

Subword Tokenization

- Mitigates the problem with out-of-vocabulary words
- Rare words are broken into substrings
- Frequent words are kept

“unfortunately” = “un” + “for” + “tun” + “ate” + “ly”

“anyplace” = “any” + “place”

“anyhow” = “any” + “how”

“anywhere” = “any” + “where”

Subword Tokenization

- Critical for languages with many variations in word structure
- Ex: Swahili (spoken in Kenya, Tanzania, and Uganda)

Conjugation of -ambia																		Press ▲
Form		Non-finite forms																Negative kutoambia
		Positive kuambia																
Positive form		Simple finite forms																Plural ambieni
		Singular ambia																
Imperative		huambia																
		Complex finite forms																
Polarity	Persons				Persons / Classes		Classes											
	1st	2nd	3rd	4th	5th	6th	7th	8th	9th	10th	11th	12th	13th	14th	15th	16th	17th	
	Sg.	Pl.	Sg.	Pl.	Sg. / 1	Pl. / 2	3	4	5	6	7	8	9	10	11 / 14	15 / 17	16	18
Past																		Press ▲
Positive	niambia	tulambia	ulambia	mlambia	alambia	walambia	ulambia	ilambia	ilambia	yalambia	kilambia	vilambia	ilambia	zilambia	ulambia	kulambia	palambia	mulambia
Negative	sikuambia	hatakuambia	hukuambia	hamkuambia	hakuambia	hawakuambi a	haukuambia	hakuambia	hakuambia	hayakuambi a	hakukuambia	havikuambia	hakuambia	hazikuambia	hakuambia	hakukuambi a	hapakuambi a	hamukuambi a
Present																		Press ▲
Positive	ninaambia	tunaambia	unaambia	mnaambia	anaambia	wanaambia	unaambia	inaambia	inaambia	yanaambia	kinaambia	vinambia	inaambia	zinaambia	unaambia	kunaambia	panaambia	munaambia
Negative	siambia	hatuambi	huambi	hamambi	hambi	hawambi	huambi	hambi	hambi	hayaambi	hakiambi	haviambi	hambi	hazambi	huambi	hakuambi	hapaambi	hamuambi
Future																		Press ▲
Positive	nitaambia	tutaambia	utaambia	mitaambia	itaambia	wataambia	utaambia	itaambia	itaambia	yataambia	kitaambia	vitaambia	itaambia	zitaambia	utaambia	kutaambia	pataambia	mutaambia
Negative	sitaambia	hataambia	hutaambia	hamtaambia	hataambia	hawataambi a	hutaambia	hataambia	hataambia	hayataambia	hakitaambia	haviataambia	hataambia	hazitaambia	hutaambia	hakutaambia	hapataambia	hamutaambi a
Subjunctive																		Press ▲
Positive	niambia	tuambi	uambi	mambi	aambi	waambi	uambi	lambi	lambi	yaambi	kiambi	viambi	lambi	ziambi	uambi	kuambi	paambi	muambi
Negative	niambia	tusambi	usambi	resambi	asambi	wasambi	usambi	isambi	isambi	yasambi	kisambi	visambi	isambi	zisambi	usambi	kusambi	pasambi	musambi
Present Conditional																		Press ▲
Positive	ninaambia	tunaambia	unaambia	mnaambia	anaambia	wanaambia	unaambia	inaambia	inaambia	yanaambia	kinaambia	vinambia	inaambia	zinaambia	unaambia	kunaambia	panaambia	munaambia
Negative	ninaambi a singaambia	tunaambi a hatungaambi	unaambi a hungaambia	mnaambi a hamungaambi	anaambi a hangaambia	wanaambi a hawangaambi	unaambi a	inaambi a halingaambi	inaambi a halingaambi	yanaambi a hayangaambi	kinaambi a hakingaambi	vinambia a havigaambi	inaambi a halingaambi	zinaambi a hazingaambi	unaambi a hauंगाambi	kunaambi a hakingaambi	panaambi a hapangaambi	munaambi a hamungaambi
Past Conditional																		Press ▲
Positive	ninaambia	tunaambia	unaambia	mnaambia	anaambia	wanaambi	unaambia	inaambia	inaambia	yanaambi	kinaambia	vinambia	inaambia	zinaambia	unaambia	kunaambi	panaambi	munaambi
Negative	ninaambi a singaambia	tunaambi a hatungaambi	unaambi a hungaambi	mnaambi a hamungaambi	anaambi a hangaambi	wanaambi a hawangaambi	unaambi a hauंगाambi	inaambi a halingaambi	inaambi a halingaambi	yanaambi a hayangaambi	kinaambi a hakingaambi	vinambia a havigaambi	inaambi a halingaambi	zinaambi a hazingaambi	unaambi a hauंगाambi	kunaambi a hakingaambi	panaambi a hapangaambi	munaambi a hamungaambi
Conditional Contrary to Fact																		Press ▲
Positive	ninaambia	tunaambia	unaambia	mnaambia	anaambia	wanaambi	unaambia	inaambia	inaambia	yanaambi	kinaambia	vinambia	inaambia	zinaambia	unaambia	kunaambi	panaambi	munaambi
Negative	ninaambi a singaambia	tunaambi a hatungaambi	unaambi a hungaambi	mnaambi a hamungaambi	anaambi a hangaambi	wanaambi a hawangaambi	unaambi a hauंगाambi	inaambi a halingaambi	inaambi a halingaambi	yanaambi a hayangaambi	kinaambi a hakingaambi	vinambia a havigaambi	inaambi a halingaambi	zinaambi a hazingaambi	unaambi a hauंगाambi	kunaambi a hakingaambi	panaambi a hapangaambi	munaambi a hamungaambi
Gnomic																		Press ▲
Positive	naambia	haambia	waambia	mwaambia	aambi	waambia	waambia	yaambi	laambi	yaambi	chaambi	vyambi	yaambi	zaambi	waambi	kwaambi	paambi	muambi
Perfect																		Press ▲

Subword Tokenization: Byte-pair Encoding (BPE)

- Goal: to represent the corpus with the least amount of tokens
- Algorithm:
 1. Start the vocabulary with all characters and the "end of word" symbol
 2. Merge the most frequent tokens
 3. Decrements the frequency of the two tokens
 4. Repeat until the desired vocabulary (e.g., pre-defined vocabulary size)

BPE: Words in the Corpus

WORD	FREQUENCY	WORD	FREQUENCY
deep </w>	3	build </w>	1
learning </w>	3	train </w>	1
the </w>	2	and </w>	1
models </w>	2	deploy </w>	1
Floydhub </w>	1	Build </w>	1
is </w>	1	models </w>	1
fastest </w>	1	in </w>	1
way </w>	1	cloud </w>	1
to </w>	1	Train </w>	1

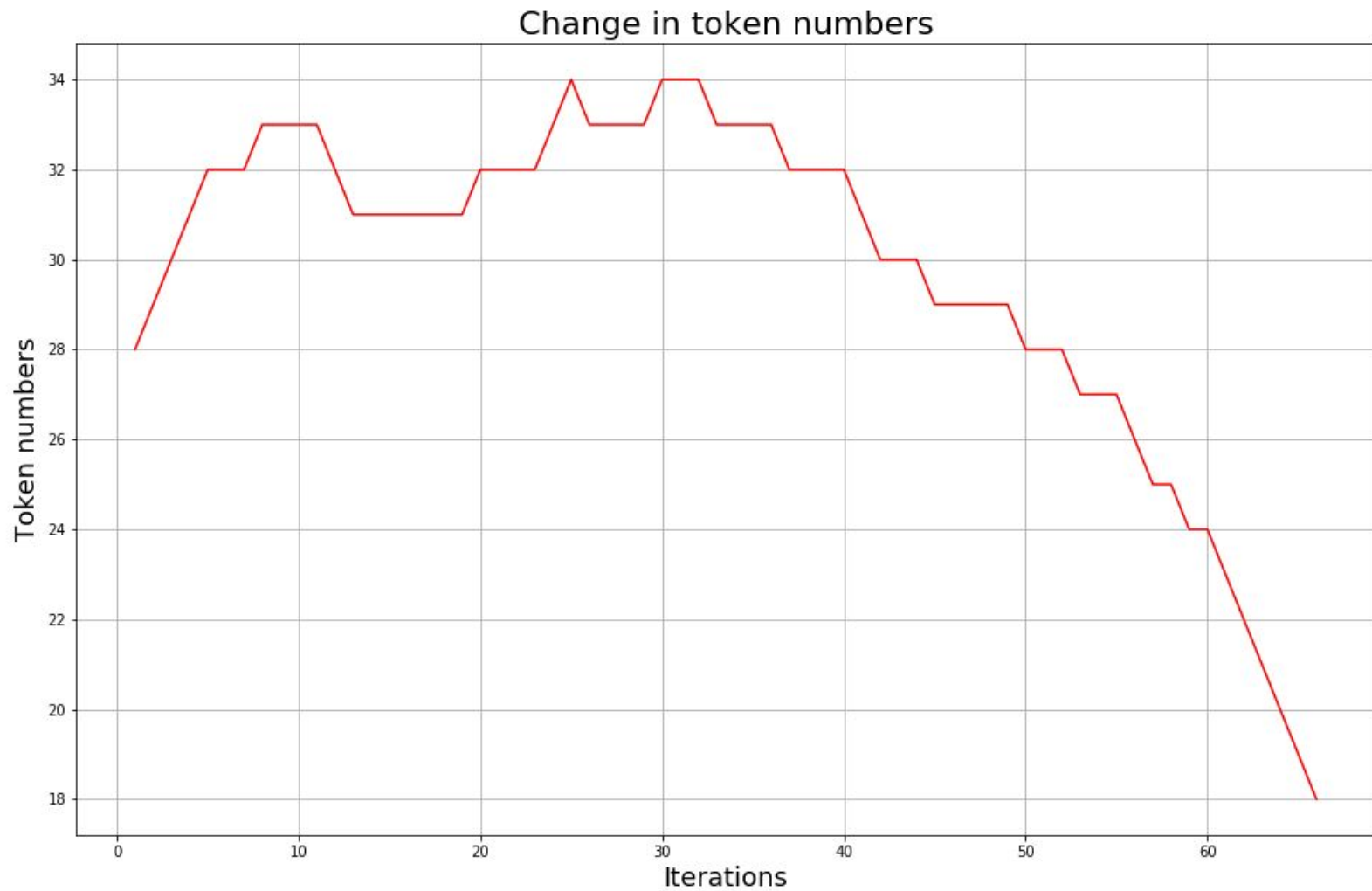
BPE: Characters in the Corpus

NUMBER	TOKEN	FREQUENCY	NUMBER	TOKEN	FREQUENCY
1	</w>	24	15	g	3
2	e	16	16	m	3
3	d	12	17	.	3
4	l	11	18	b	2
5	n	10	19	h	2
6	i	9	20	F	1
7	a	8	21	H	1
8	o	7	22	f	1
9	s	6	23	w	1
10	t	6	24	,	1
11	r	5	25	B	1
12	u	4	26	c	1
13	p	4	27	T	1
14	y	3			

BPE: Merge

NUMBER	TOKEN	FREQUENCY	NUMBER	TOKEN	FREQUENCY
1	</w>	24	16	g	3
2	e	$16 - 7 = 9$	17	m	3
3	d	$12 - 7 = 5$	18	.	3
4	l	11	19	b	2
5	n	10	20	h	2
6	i	9	21	F	1
7	a	8	22	H	1
8	o	7	23	f	1
9	de	7	24	w	1
10	s	6	25	,	1
11	t	6	26	B	1
12	r	5	27	c	1
13	u	4	28	T	1
14	p	4			
15	y	3			

BPE



Wordpiece

- Similar ao BPE
- Algorithm:
 1. Start the vocabulary with all characters and the "end of word" symbol
 2. Merge the tokens with the highest score
 3. Decrements the frequency of the two tokens
 4. Repeat until the desired vocabulary (e.g., pre-defined vocabulary size)

$$\text{score} = (\text{freq_of_pair}) / (\text{freq_of_first_element} \times \text{freq_of_second_element})$$

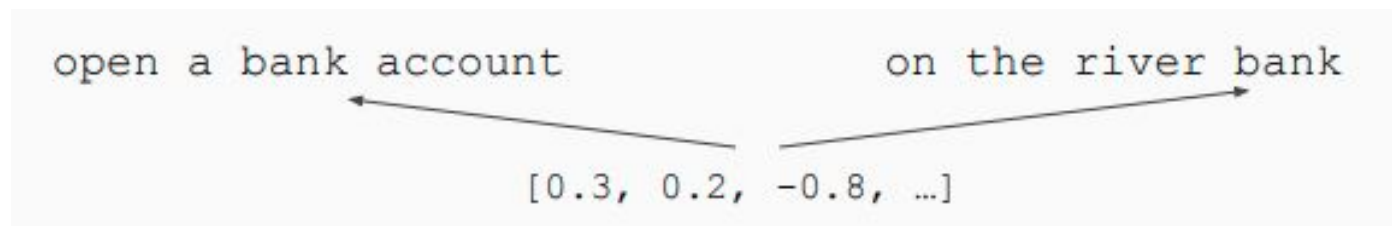
high probability of they occur together than separated

Wordpiece

Word	Token(s)
surf	['surf']
surfing	['surf', '##ing']
surfboarding	['surf', '##board', '##ing']
surfboard	['surf', '##board']
snowboard	['snow', '##board']
snowboarding	['snow', '##board', '##ing']
snow	['snow']
snowing	['snow', '##ing']

Contextual Representations

- Limitation of word embeddings: same representation for different meanings

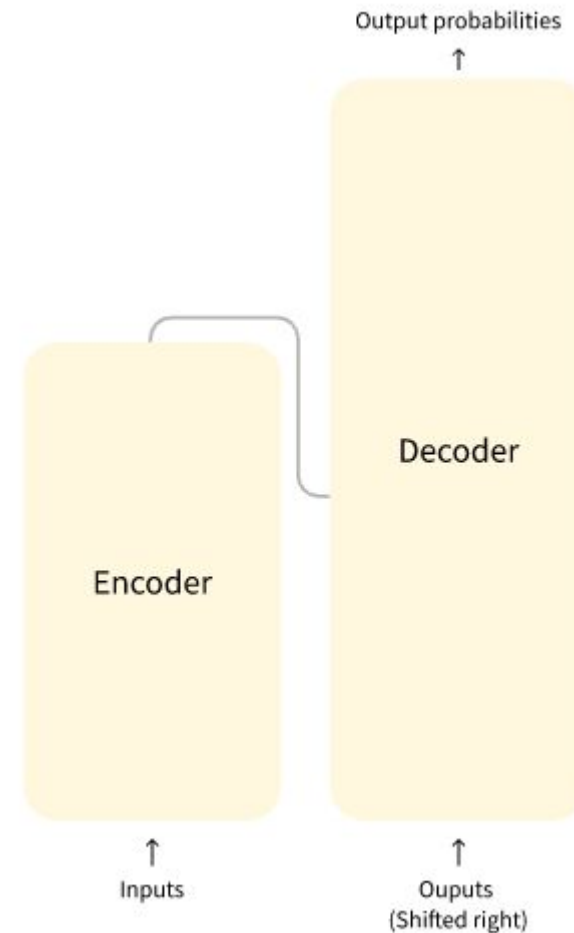


- Solution: learn contextual representations



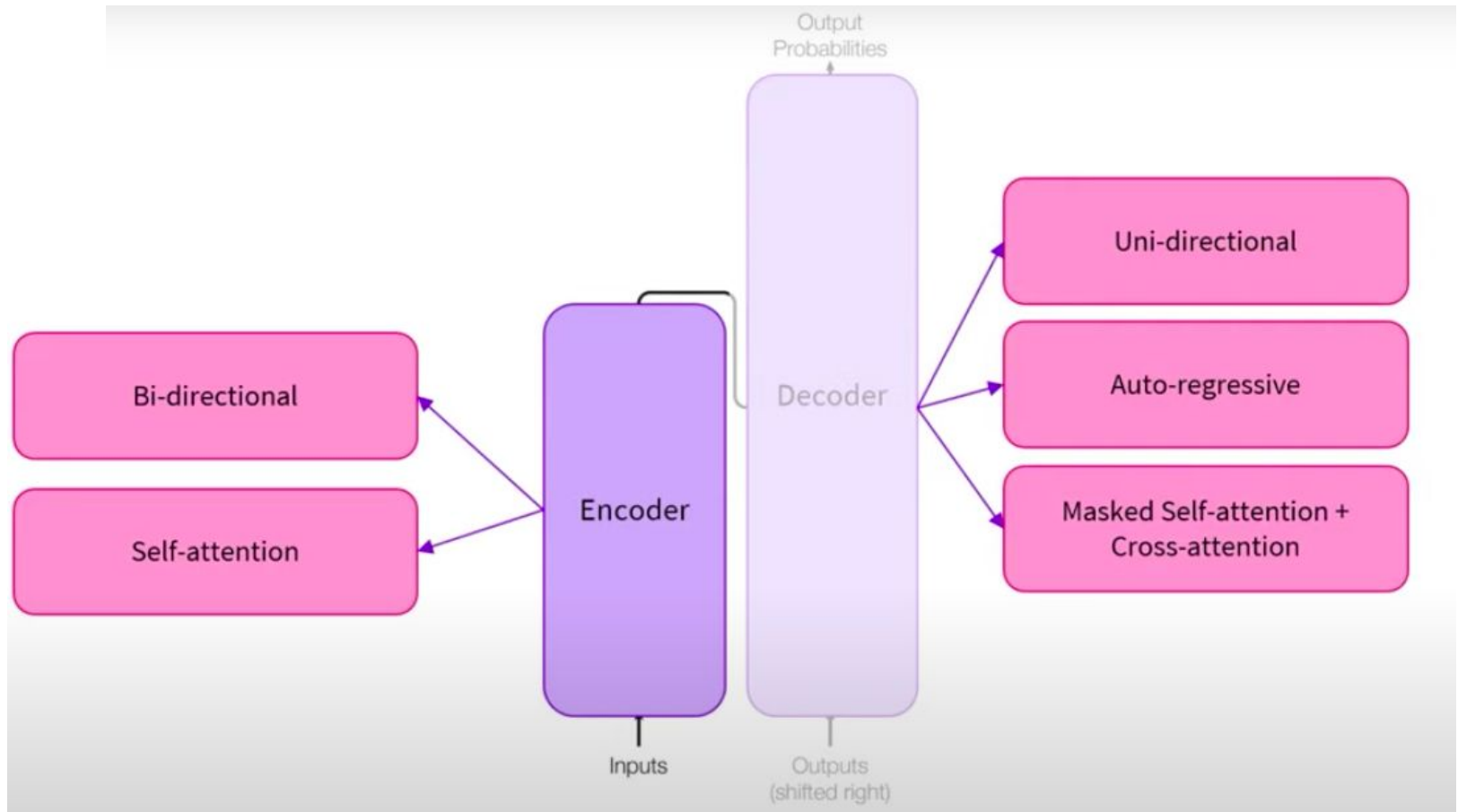
Transformers: Architecture

- Encoder: builds a representation of the input (embeddings)
- Decoder: outputs probabilities based on the encoder's output and other inputs



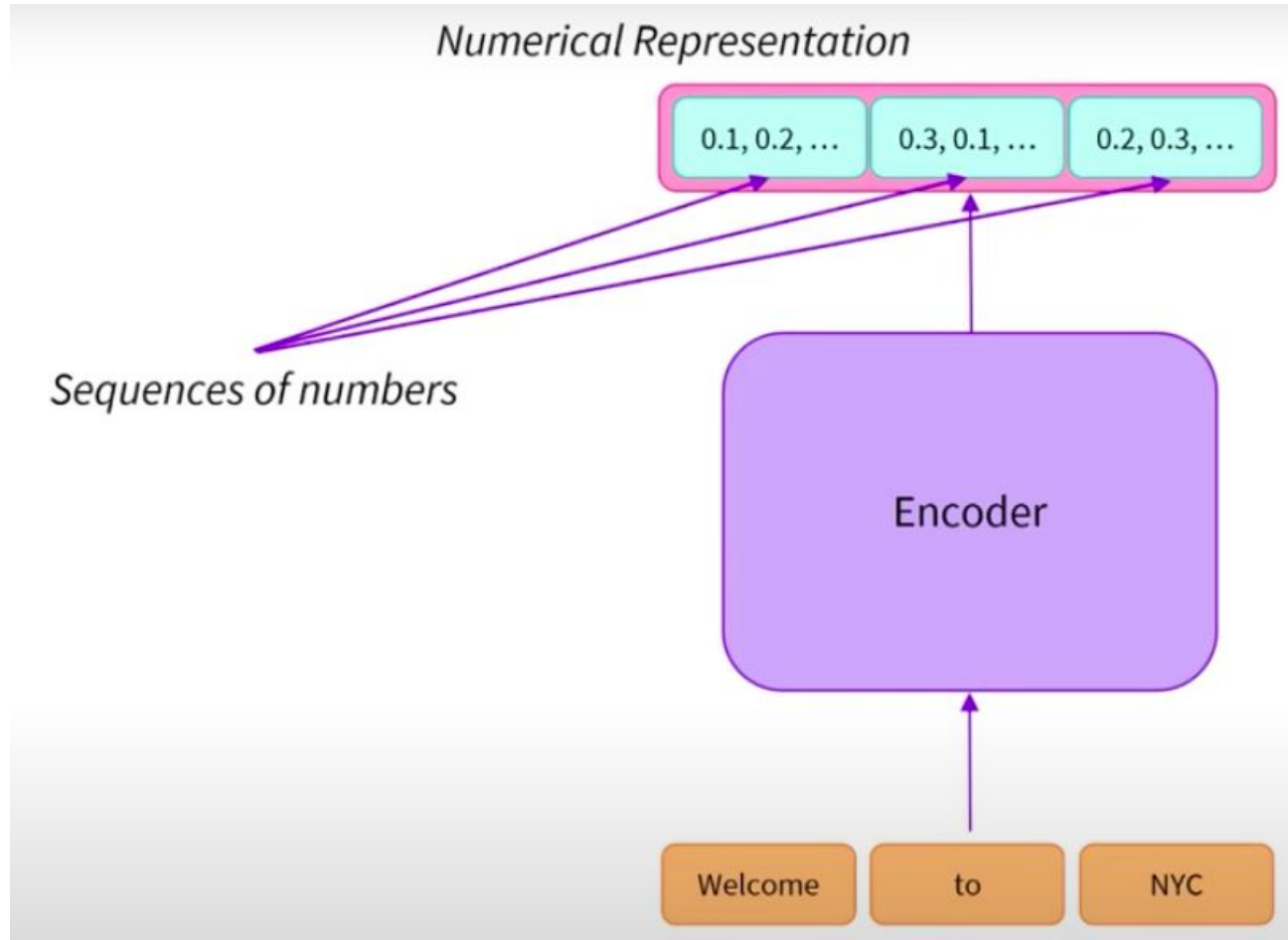
source: <https://huggingface.co/course/chapter1/4>

Transformers: Encoder



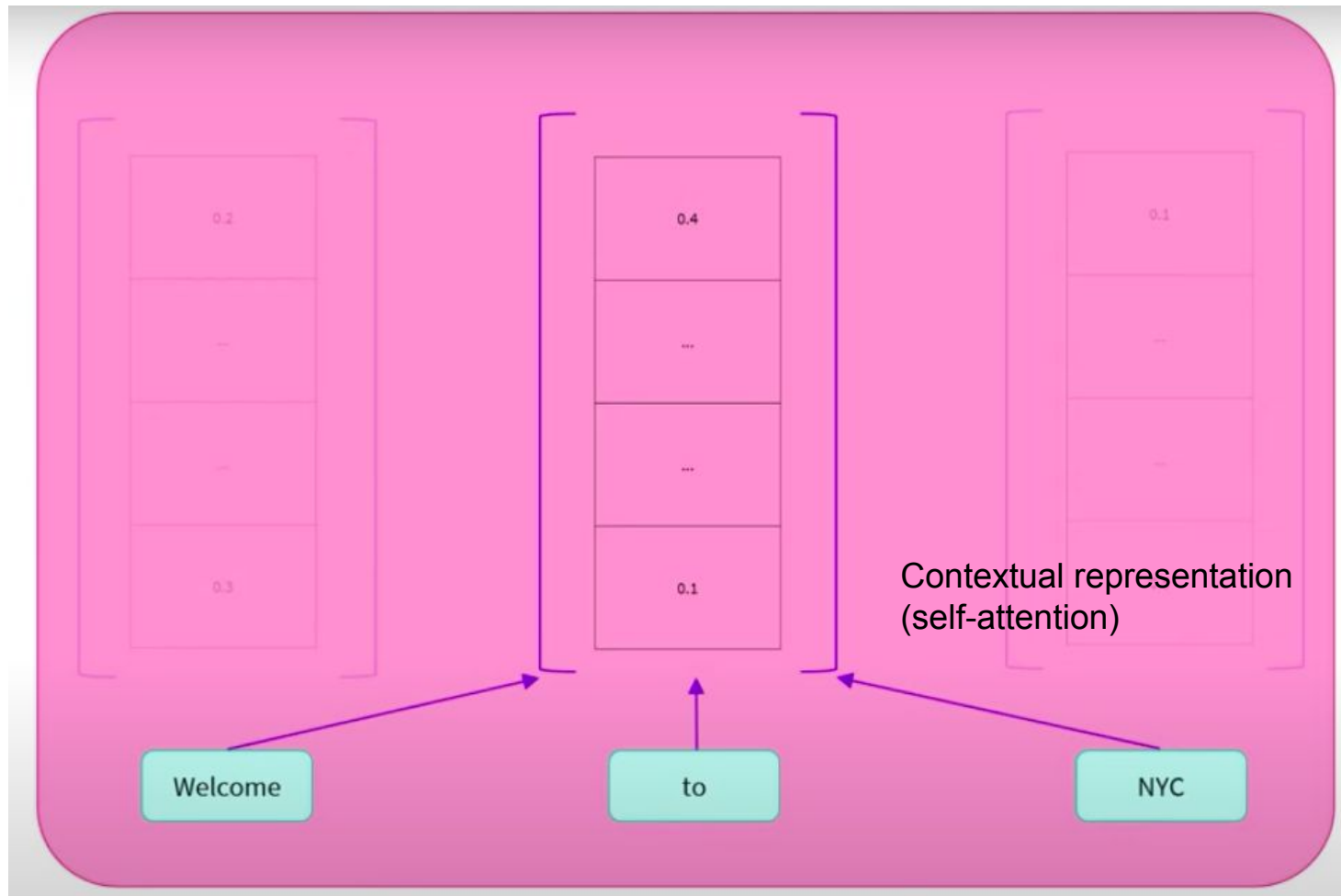
source: <https://huggingface.co/course/chapter1/5>

Transformers: Encoder



source: <https://huggingface.co/course/chapter1/5>

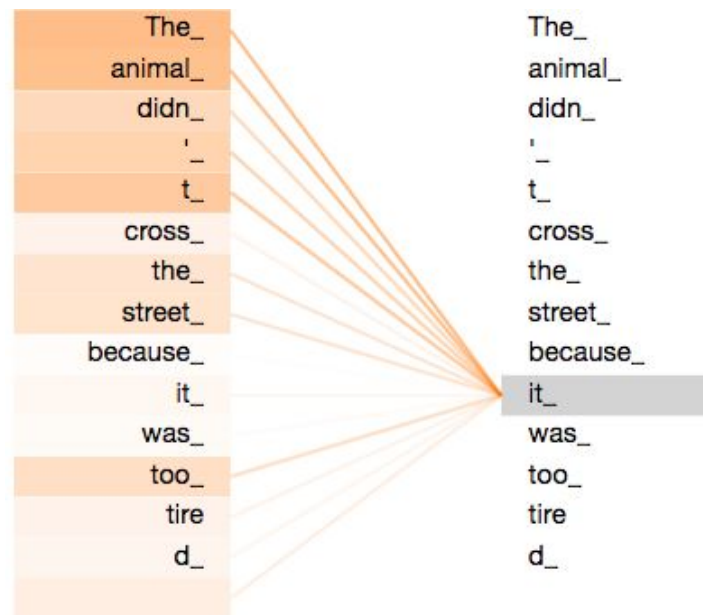
Transformers: Encoder



source: <https://huggingface.co/course/chapter1/5>

Self-Attention

- All words in the sentence have influence in the word representation



<https://jalammar.github.io/illustrated-transformer/>

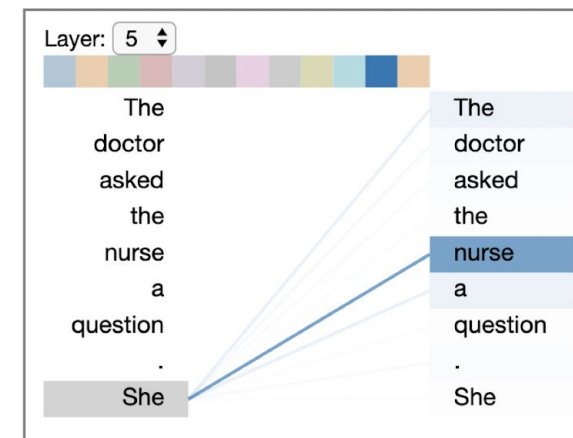
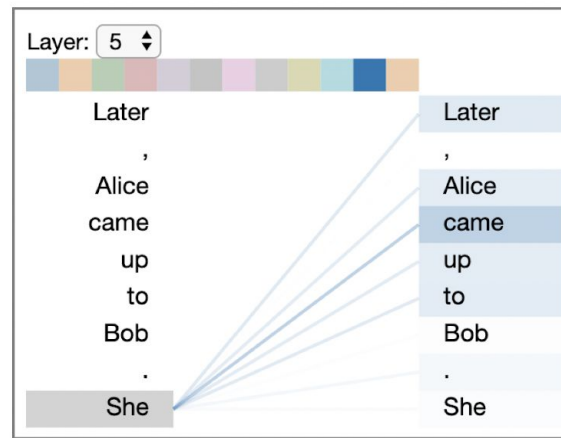
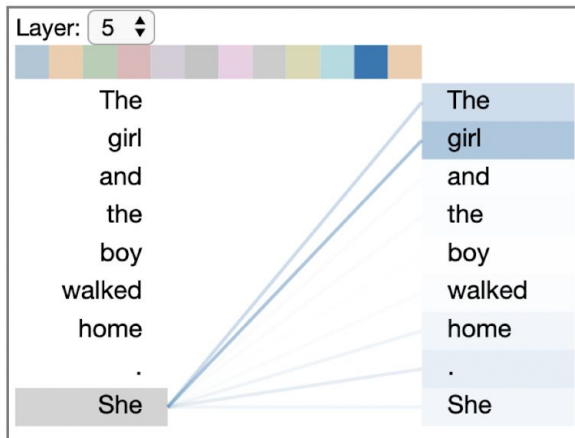
Examples of Self-Attention

Gender-specific term

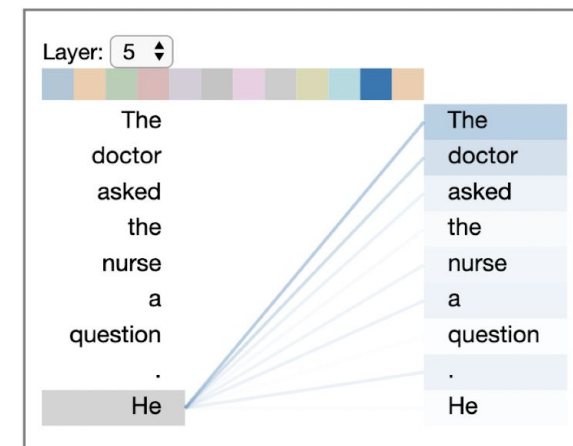
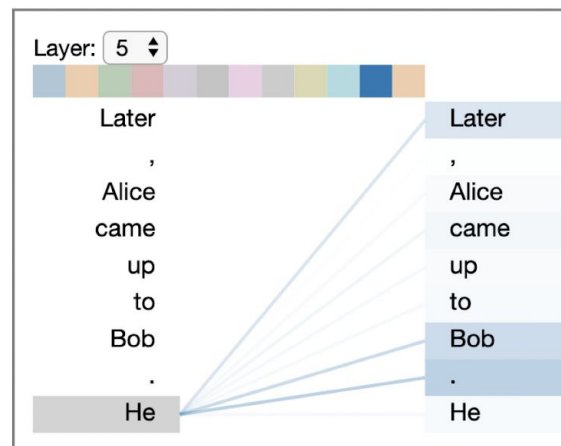
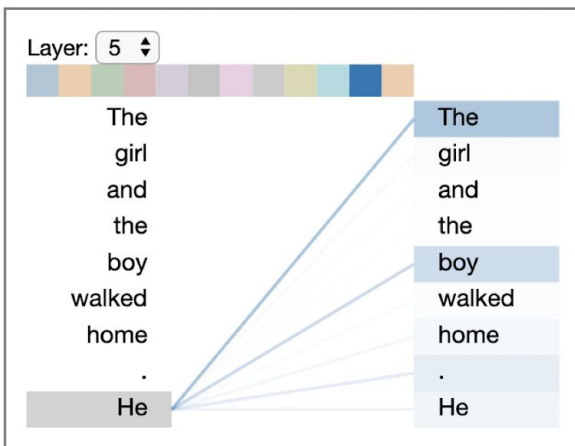
Name

Occupation

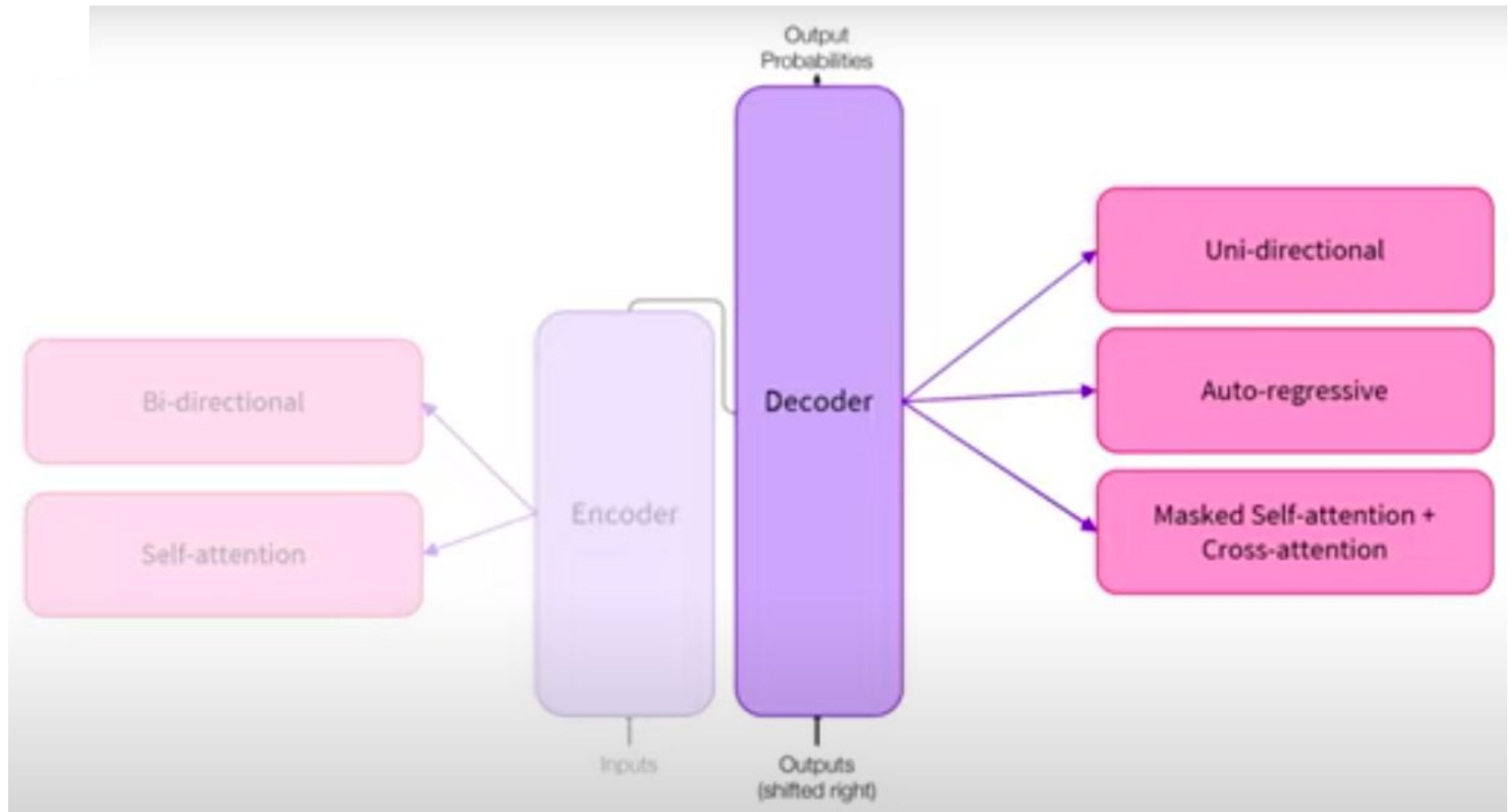
She



He



Transformers: Decoder



source: <https://huggingface.co/course/chapter1/6>

Transformers: Decoder

My



name

My

name



is

My

name

is



Sylvain

My

name

is

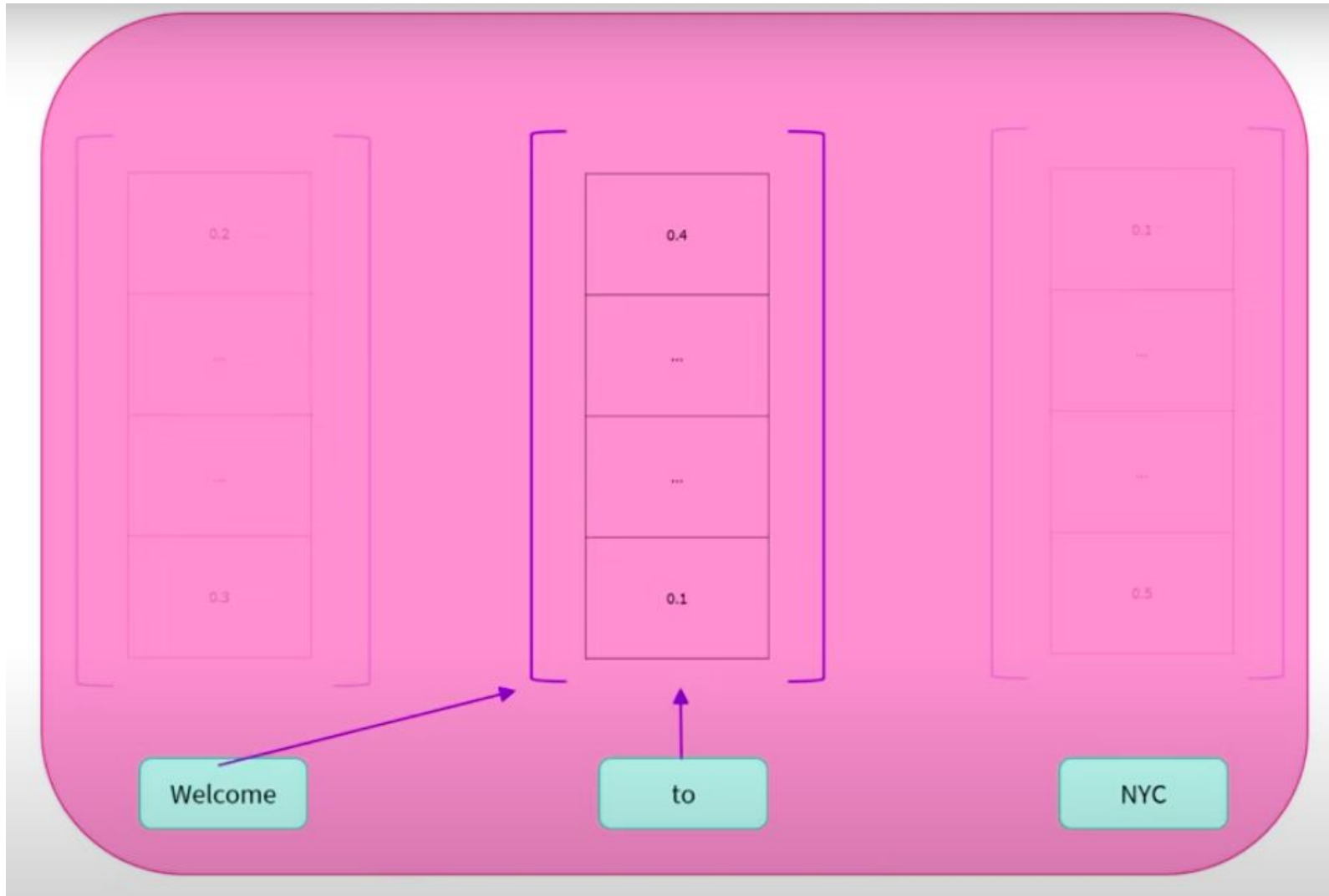
Sylvain



.

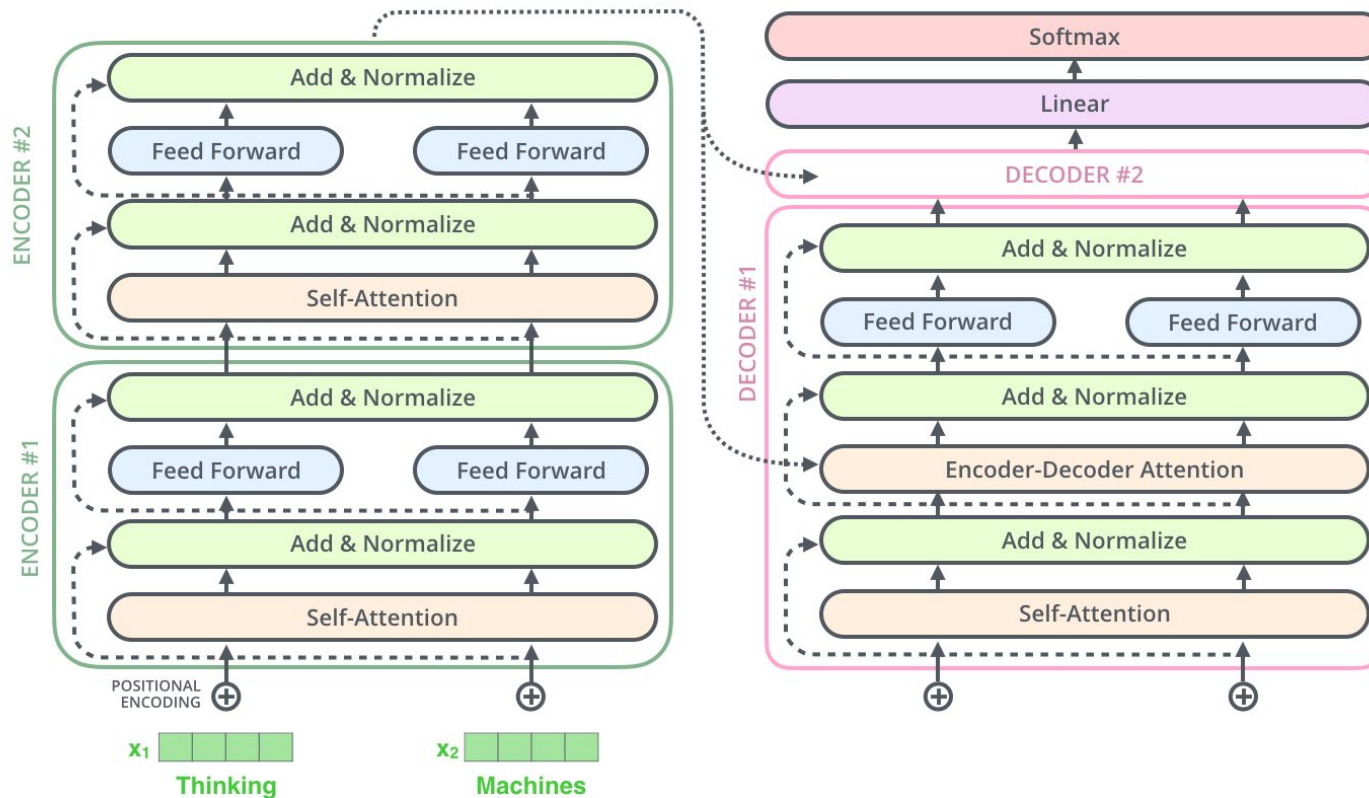
source: <https://huggingface.co/course/chapter1/4>

Transformers: Decoder

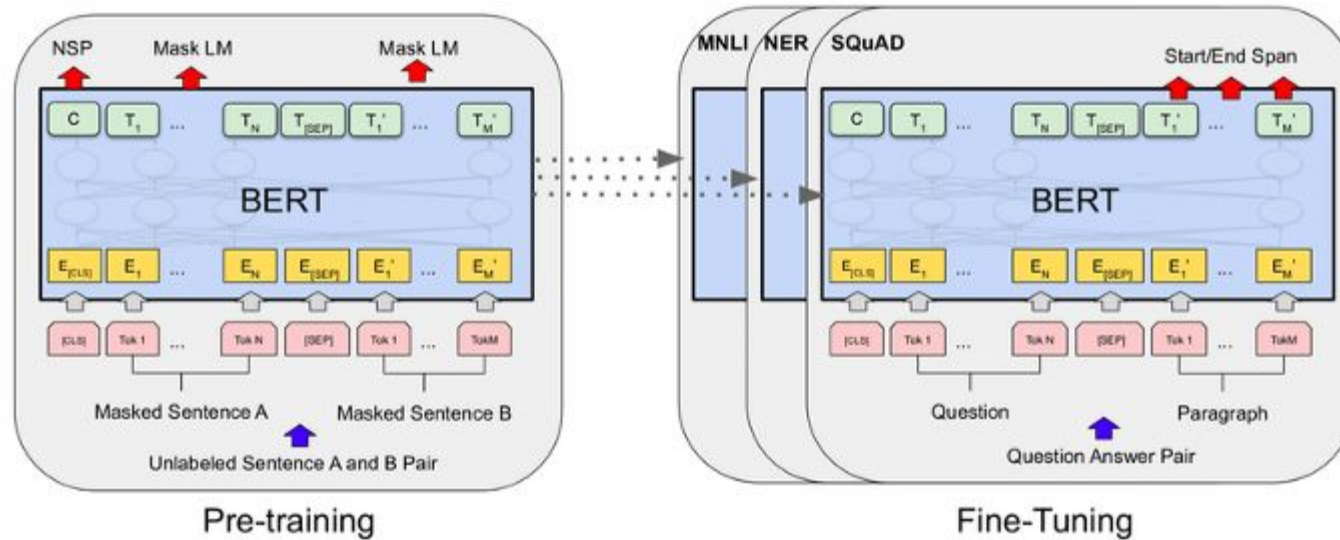


source: <https://huggingface.co/course/chapter1/6>

BERT: Bidirectional Encoder Representations from Transformers



BERT: Two Steps

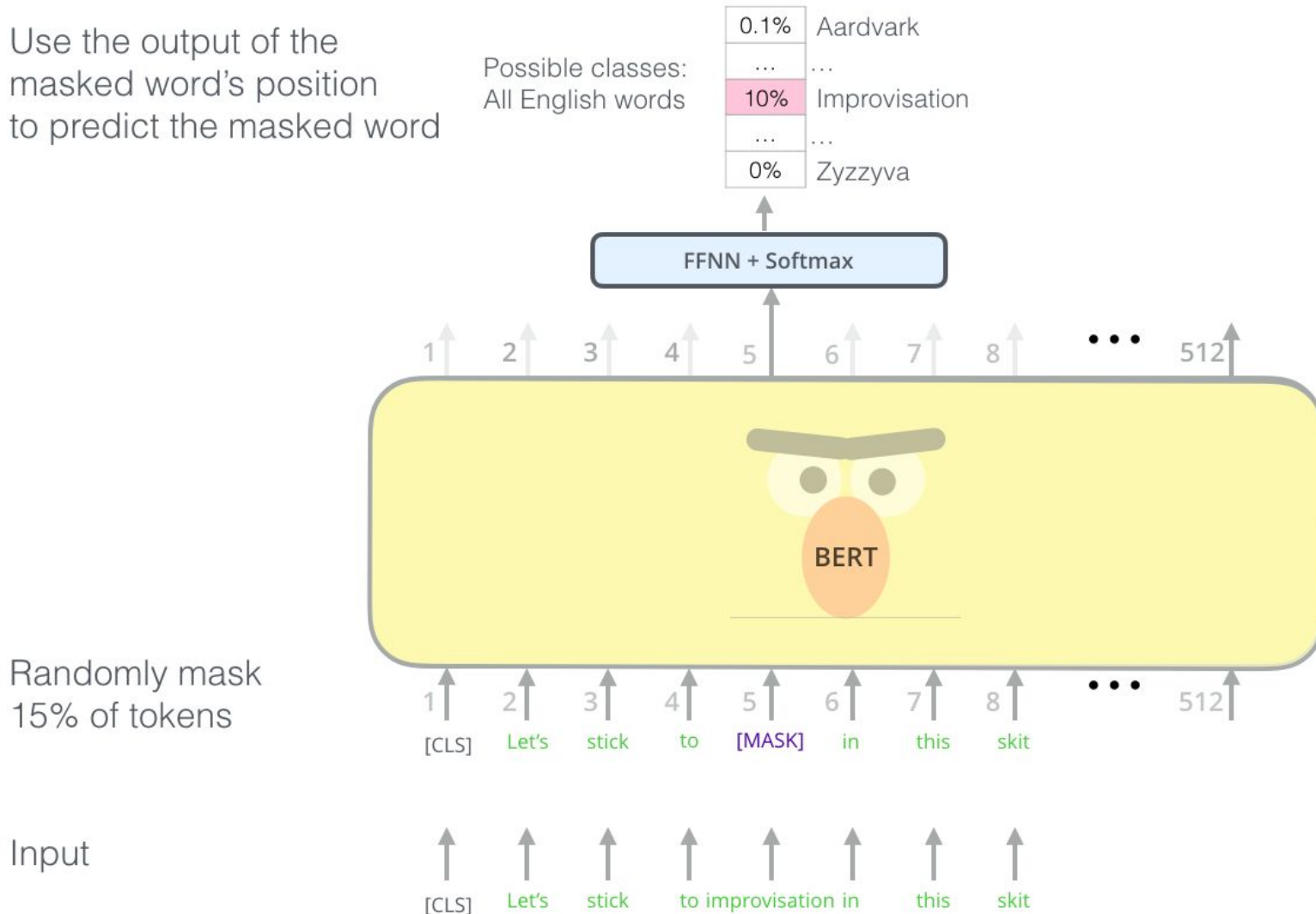


Pre-training: Masked Language Model

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

Possible classes:
All English words

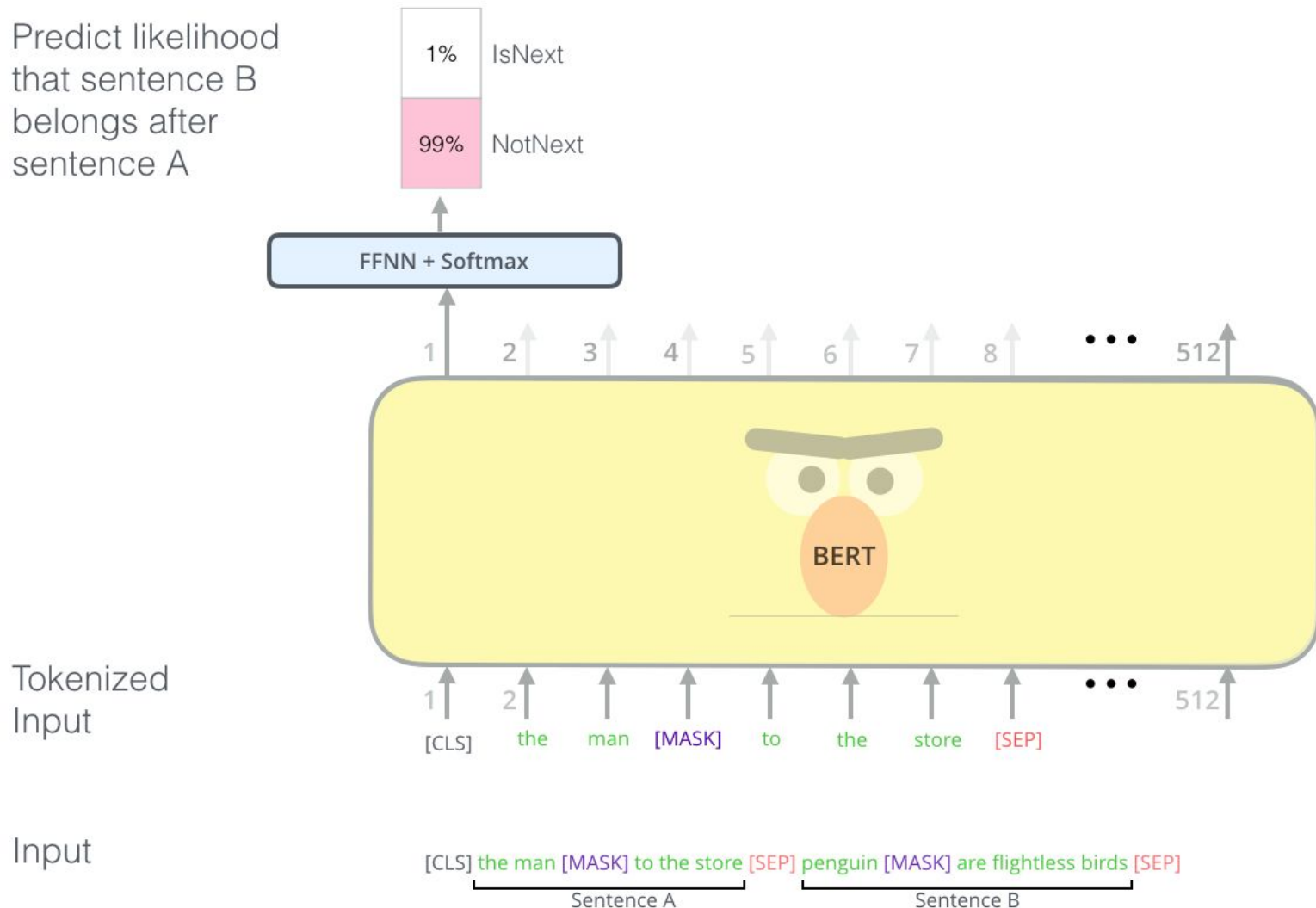
0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva



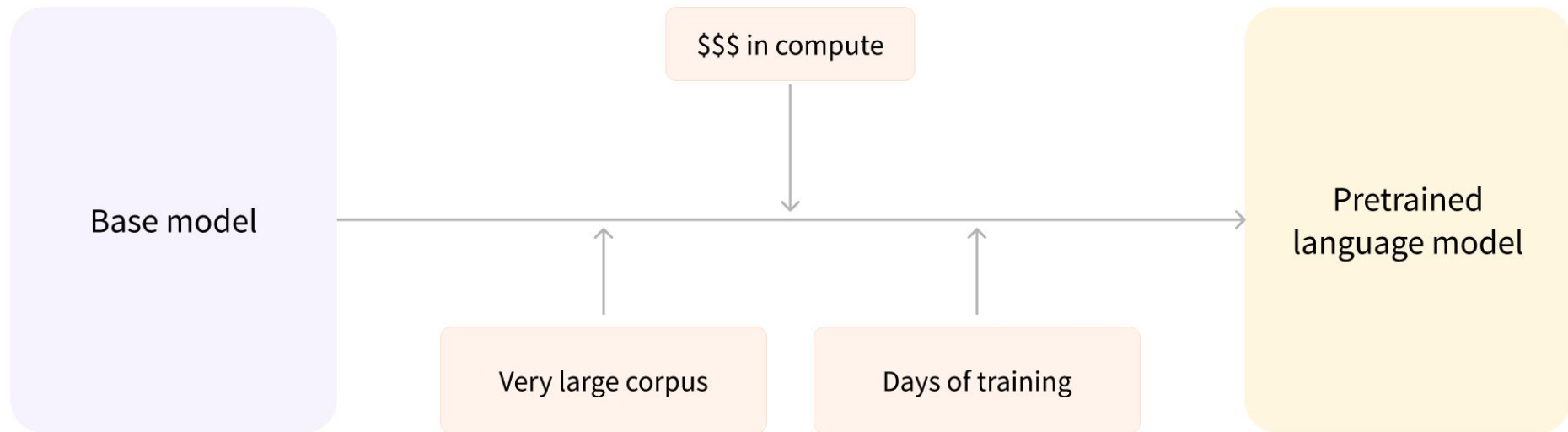
Img-Source: <https://jalammar.github.io/illustrated-transformer/>

Pre-training: Next-Sentence Prediction

Predict likelihood
that sentence B
belongs after
sentence A



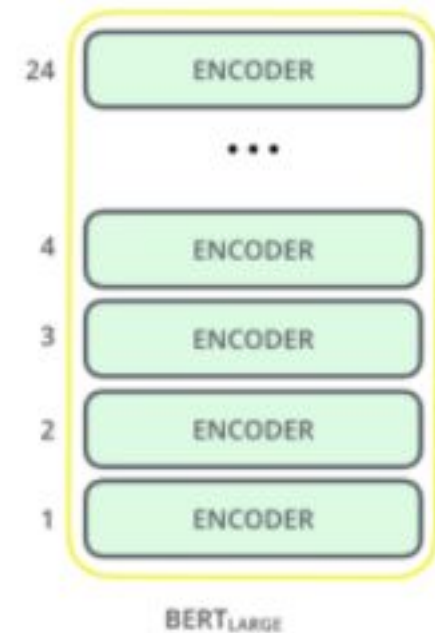
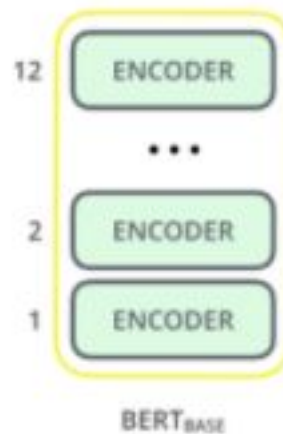
Pre-training



source: <https://huggingface.co/course/chapter1/4>

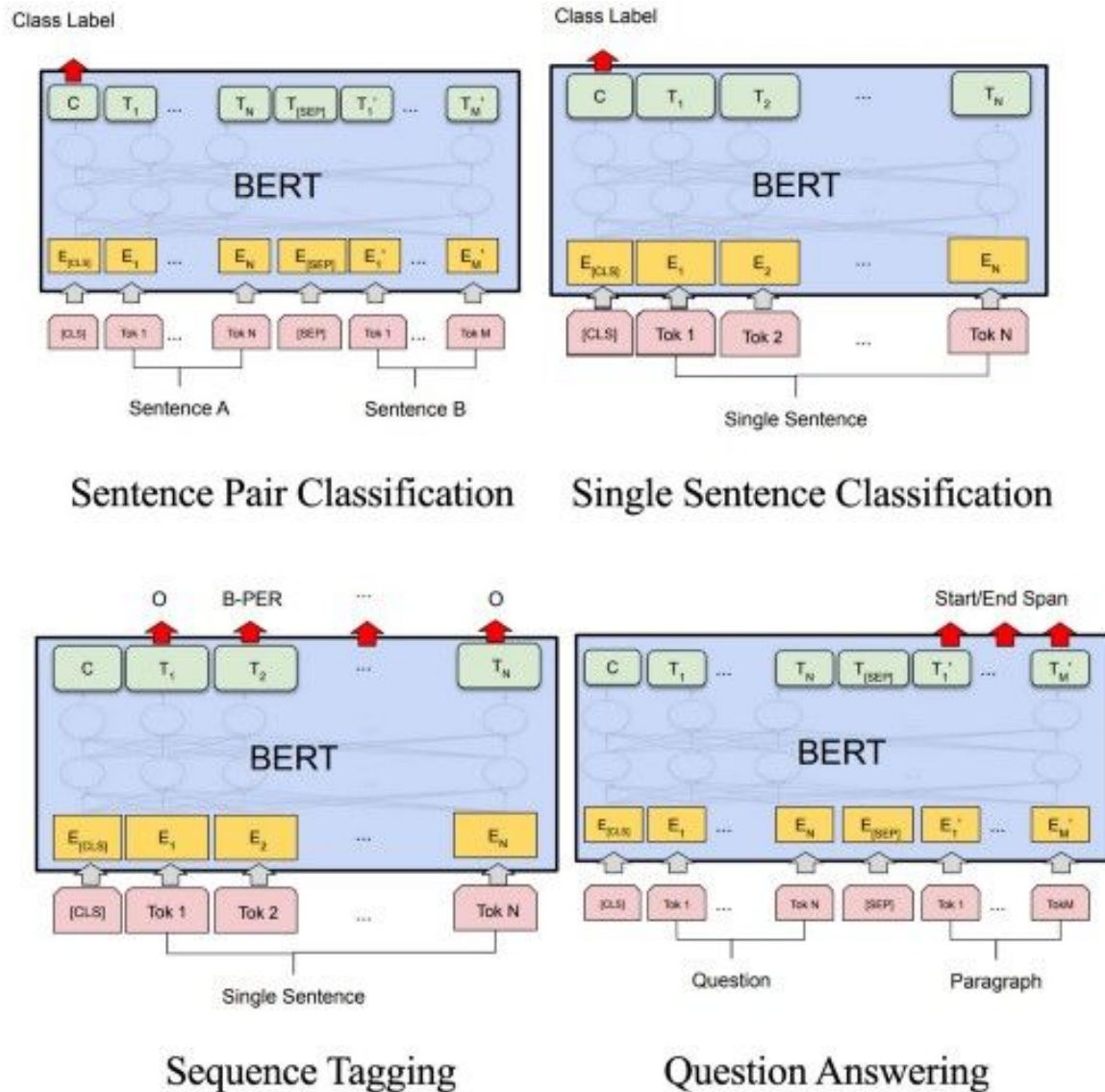
BERT: Details

- Data: Wikipedia (2.5B words) + BookCorpus (800M words)
- BERT-Base: 110M
- BERT-Large: 340M

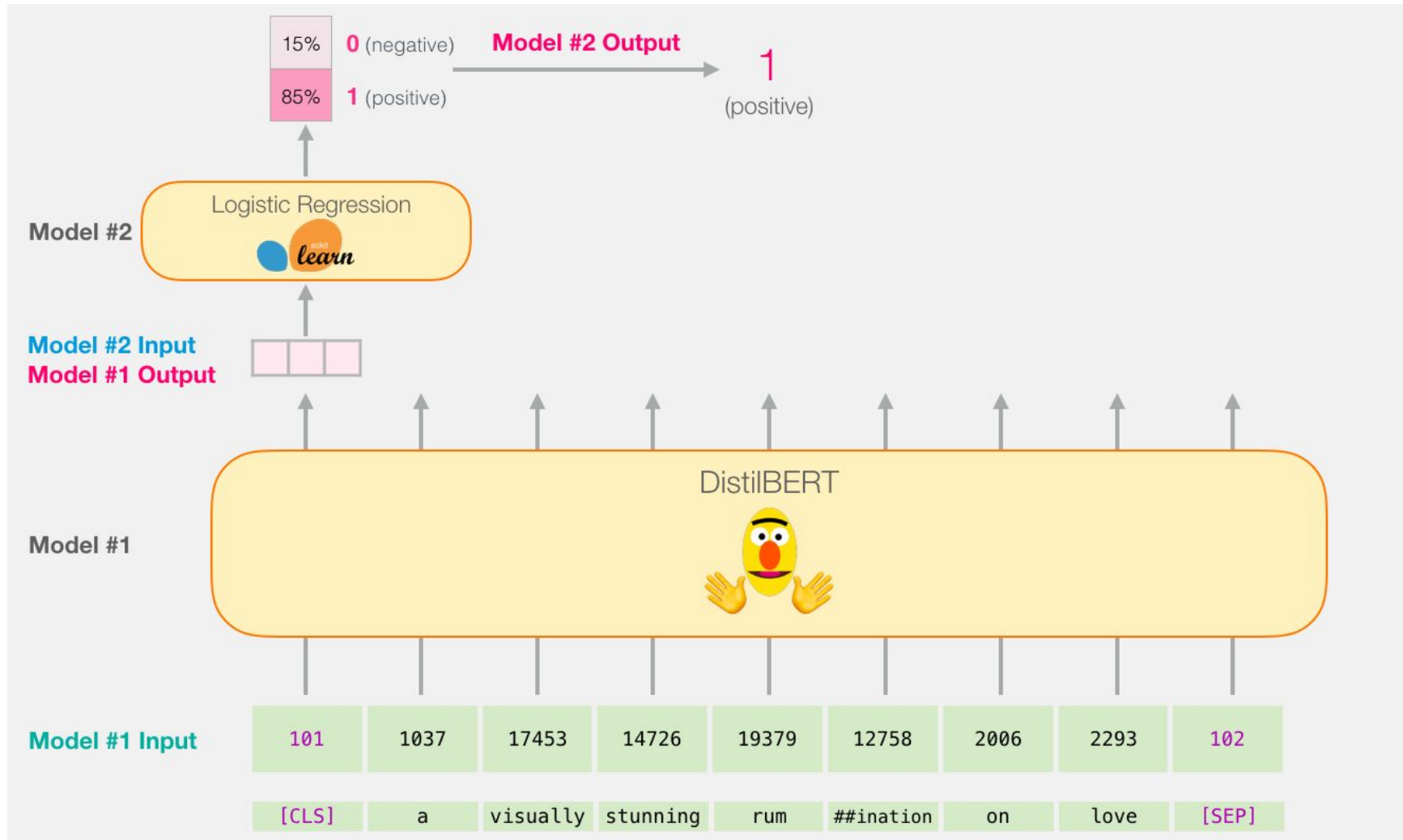


Img-Source: <https://jalammar.github.io/illustrated-transformer/>

Fine-tuning (Transfer Learning)

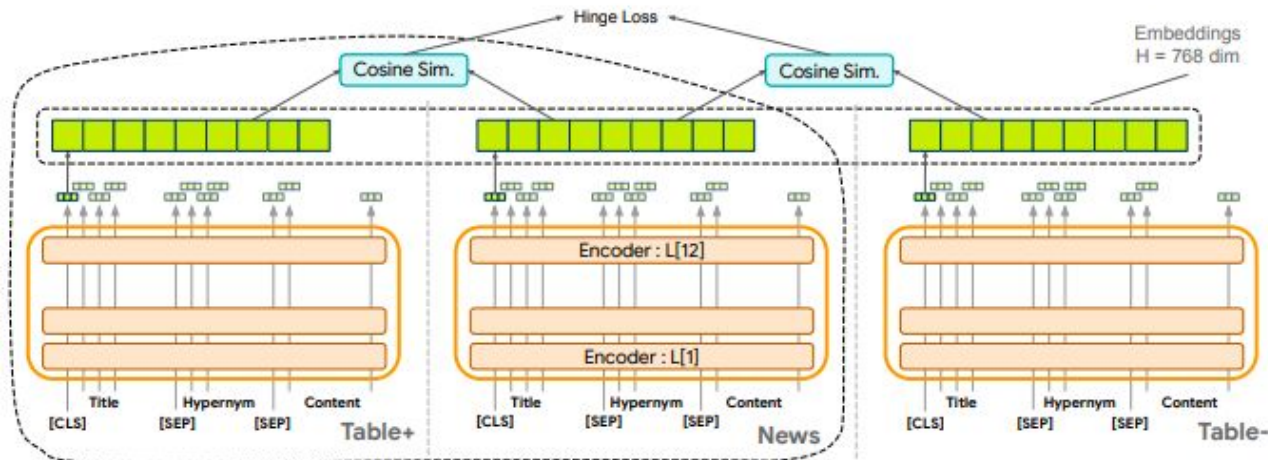


Sentence Classification: Sentiment Analysis



Img-Source: <https://jalammar.github.io/illustrated-transformer/>

Pair Classification: Article-Table Matching



Model	Table	News	acc.@k=1	5	10	100
BM25	⟨Title,Content⟩	⟨Title,Content⟩	.426	.574	.622	.831
	⟨Title,Content,Table Content⟩	⟨Title,Content⟩	.372	.588	.703	.838
TF-IDF	⟨Title,Content,Table Content⟩	⟨Title,Content⟩	.453	.642	.730	.892
USE	⟨Title,Content⟩	⟨Title⟩	.250	.439	.507	.669
	⟨Title,Content⟩	⟨Content⟩	.237	.419	.520	.851
	⟨Title,Content⟩	⟨Title,Content⟩	.243	.466	.561	.838
Doc2Vec	⟨Title⟩	⟨Title⟩	.297	.507	.581	.824
	⟨Title,Content⟩	⟨Title,Content⟩	.223	.378	.487	.737
BERTpublic	⟨Title,Content⟩	⟨Title⟩	.155	.291	.372	.574
NewsBERT	⟨Title,Content⟩	⟨Title,Content⟩	.458	.725	.779	.824

Model	Table	News	acc.@k=1	5	10	100
NewsBERT	⟨Title⟩	⟨Title⟩	.422	.602	.656	.719
	⟨Title⟩	⟨Title,Content⟩	.438	.677	.746	.8
	⟨Title,Content⟩	⟨Title⟩	.176	.244	.313	.588
	⟨Title,Content⟩	⟨Title,Content⟩	.458	.725	.779	.824
NT-BERT	⟨Title⟩	⟨Title⟩	.545	.773	.833	.939
	⟨Title⟩	⟨Title,Content⟩	.53	.818	.871	.955
	⟨Title,Content⟩	⟨Title⟩	.545	.795	.841	.932
	⟨Title,Content⟩	⟨Title,Content⟩	.553	.856	.879	.947
NTH-BERT	⟨Title,Hypernyms,Content⟩	⟨Title⟩	.575	.811	.858	.961
	⟨Title,Hypernyms,Content⟩	⟨Title,Hypernyms⟩	.543	.803	.858	.953
	⟨Title,Hypernyms,Content⟩	⟨Title,Hypernyms,Content⟩	.669	.898	.953	.976

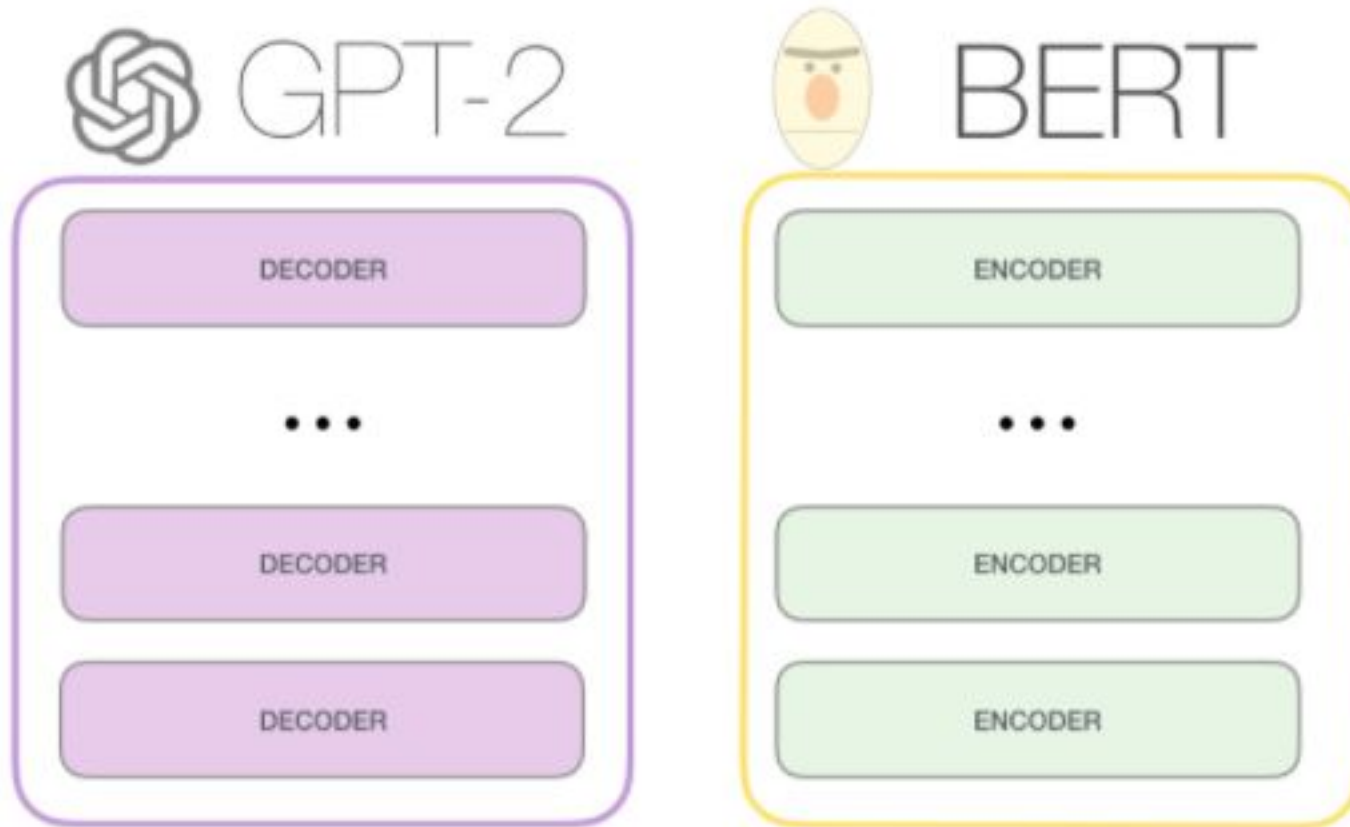
Lees et al. **Collocating News Articles with Structured Web Tables**. International Workshop on News Recommendation and Intelligence, 2021



Generative Pretrained Transformer (GPT)

- Transformer decoder com 12 camadas
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers
- Byte-pair encoding com 40.000 merges
- Treinado no BooksCorpus: mais de 7.000 livros (sentenças longas)
- Demo: <https://demo.allennlp.org/next-token-lm>

GPT-2 vs BERT

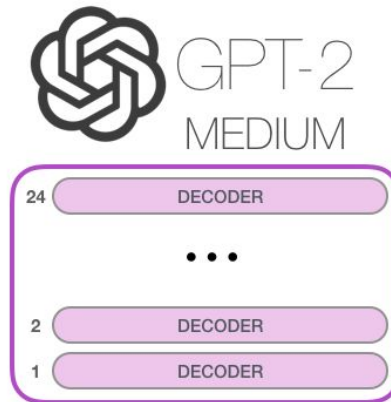


<https://jalammar.github.io/illustrated-gpt2/>

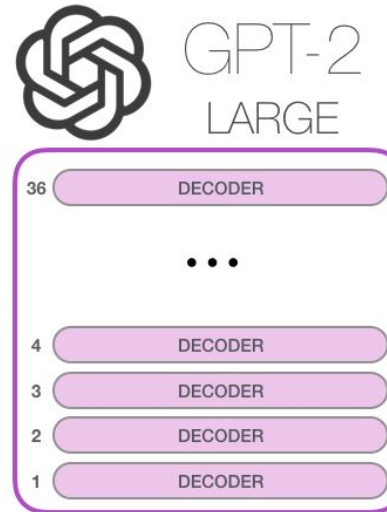
GPT-2 Models



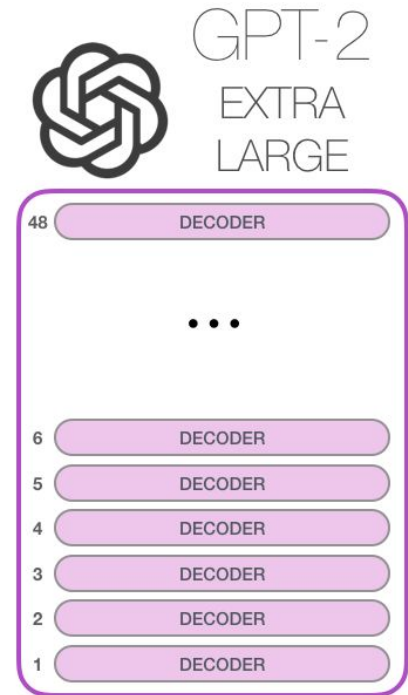
Model Dimensionality: 768



Model Dimensionality: 1024



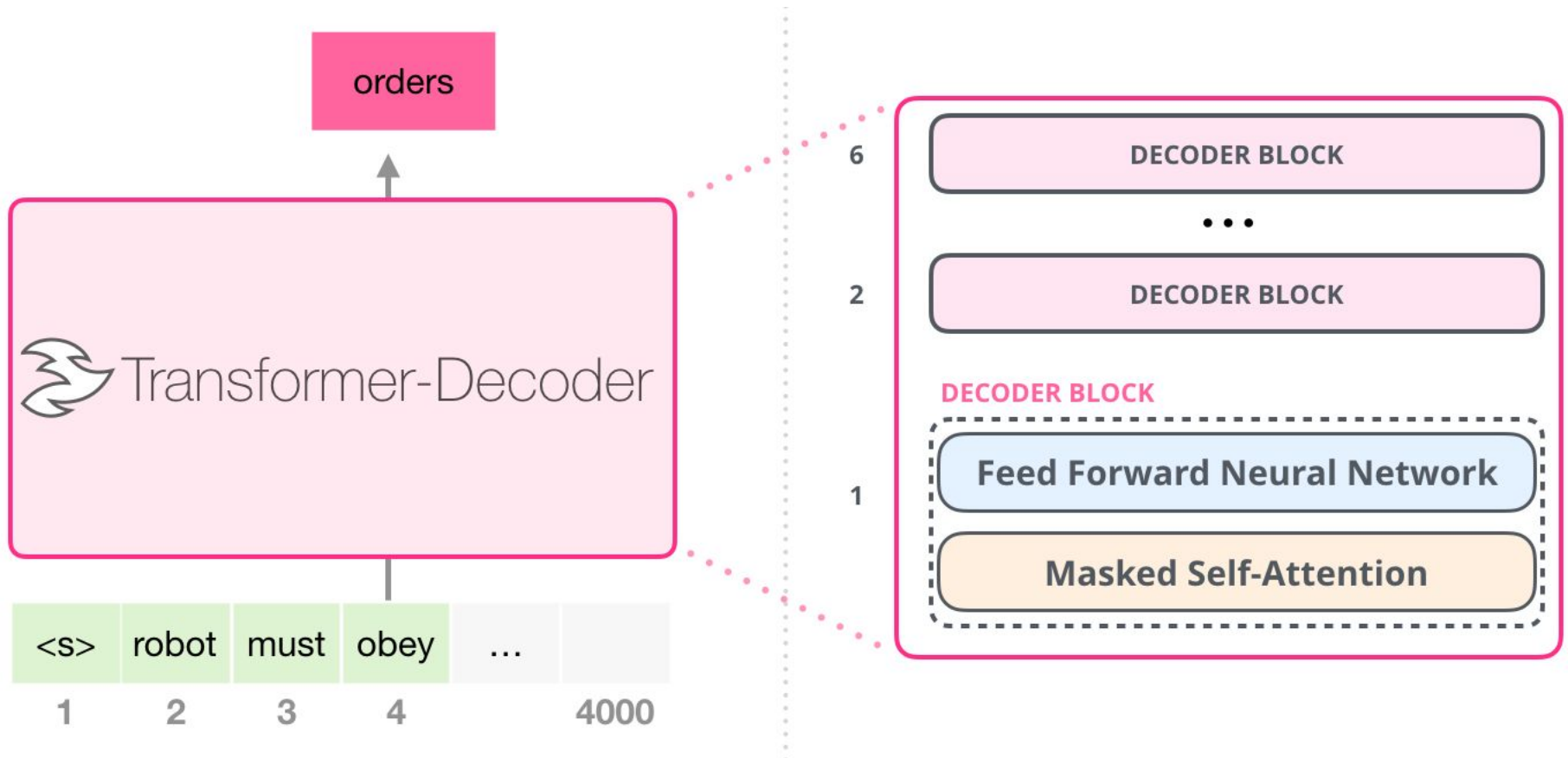
Model Dimensionality: 1280



Model Dimensionality: 1600

<https://jalammar.github.io/illustrated-gpt2/>

GPT-2

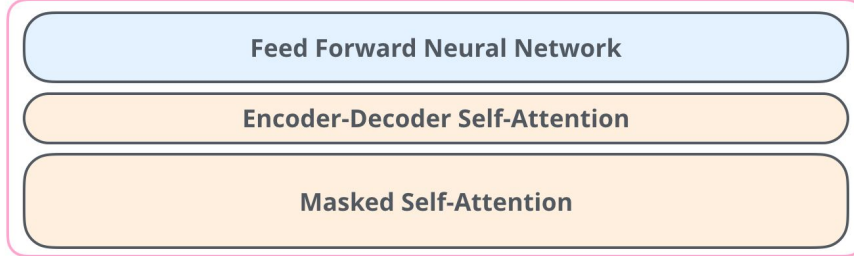


<https://jalammar.github.io/illustrated-gpt2/>

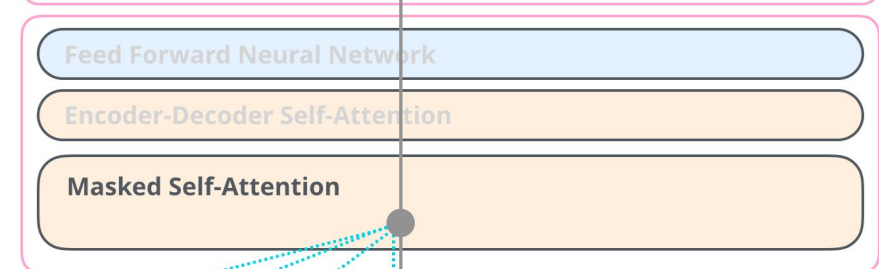
GPT-2: Masked Self-Attention



DECODER BLOCK



DECODER BLOCK #2

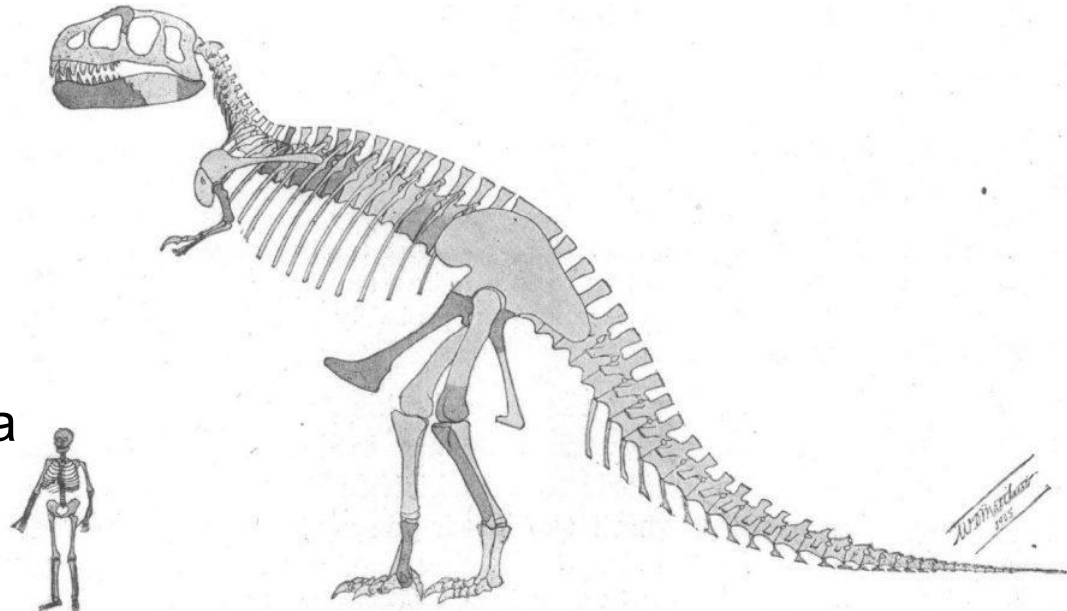


orders

<https://jalammar.github.io/illustrated-gpt2/>

GPT-3

- 175 billion of parameters
- (100x > GPT-2)
- Corpus:
- 45TB
- Common Crawl, WebText2, Books1, Books2 e Wikipedia
- 96 layers



GPT-2
1.5B Parameters

GPT-3
175B Parameters

GPT-3

In-context Learning

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```

1 Translate English to French:  ← task description
2 cheese =>                   ← prompt
  
```

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```

1 Translate English to French:  ← task description
2 sea otter => loutre de mer    ← example
3 cheese =>                    ← prompt
  
```

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```

1 Translate English to French:  ← task description
2 sea otter => loutre de mer    ← examples
3 peppermint => menthe poivrée ← examples
4 plush girafe => girafe peluche ← examples
5 cheese =>                    ← prompt
  
```

Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.

```

1 sea otter => loutre de mer  ← example #1
    ↓
    gradient update
    ↓
1 peppermint => menthe poivrée ← example #2
    ↓
    gradient update
    ↓
...
1 plush giraffe => girafe peluche ← example #N
    ↓
    gradient update
    ↓
1 cheese =>                    ← prompt
  
```

Source: <https://arxiv.org/pdf/2005.14165.pdf>

ChatGPT (GPT-3.5)

- Supervised Fine Tuning (SFT) Model
- Fine-tuning the GPT-3 model on a large labeled dataset
- Inputs and prompts from users of the OpenAI API platform
- Manual labelers needed to create examples of categories not covered by OpenAI API users