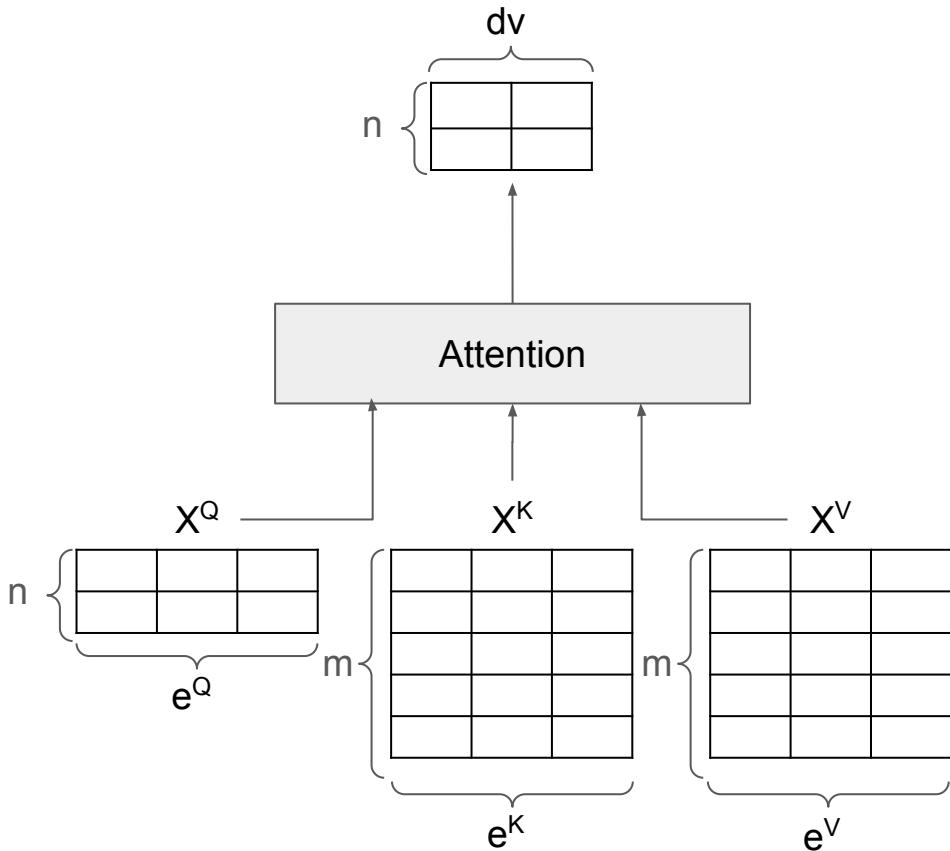


Attention

Attention przyjmuje trzy wejścia 2D $\{\mathbb{R}^{n \times e^q}, \mathbb{R}^{m \times e^k}, \mathbb{R}^{m \times e^v}\}$ i zwraca pojedyncze wyjście 2D $\{\mathbb{R}^{n \times d^v}\}$.

Jak za chwilę zobaczymy, attention przypomina sumę ważoną. Każde *query* $x^Q \in X^Q$ jest dopasowywane do wszystkich *keys* X^K . Te x^K , które lepiej pasują dokładają więcej *values* x^V do wyniku.

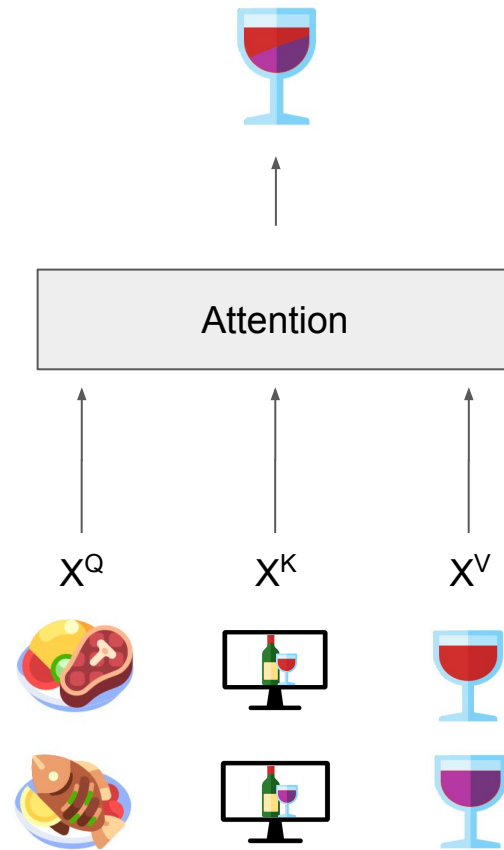
Zazwyczaj $e^Q = e^K = e^V$.



Attention: przykład

Powiedzmy, że prowadzimy restaurację. Chcemy podawać wyjątkowe wina pasujące do pozycji w naszym menu, ale dysponujemy jedynie małym zbiorem win z Lidl.

Kiedy klient prosi o wino pasujące do dania, porównujemy przepis na to danie X^Q do informacji o naszych winach znalezionych w internecie X^K . Następnie zlewamy razem wszystkie wina X^V , biorąc więcej tych, które zdawały się pasować lepiej.



Attention: implementacja (dot-product attention)

Attention(X^Q, X^K, X^V) =

$$\text{softmax}\left(\frac{X^Q W^Q (X^K W^K)^T}{\sqrt{d_k}}\right) X^V W^V$$

Dzielenie przez $\sqrt{d_k}$ czyni wariancję argumentu softmax równą jeden, co pomaga przy trenowaniu.

n : number of query tokens

m : number of key/value tokens

e : embedding length

d_k : dimension of key representation

d_v : dimension of value representation

$X^Q \in \mathbb{R}^{n \times e}$: queries

$X^K \in \mathbb{R}^{m \times e}$: keys

$X^V \in \mathbb{R}^{m \times e}$: values

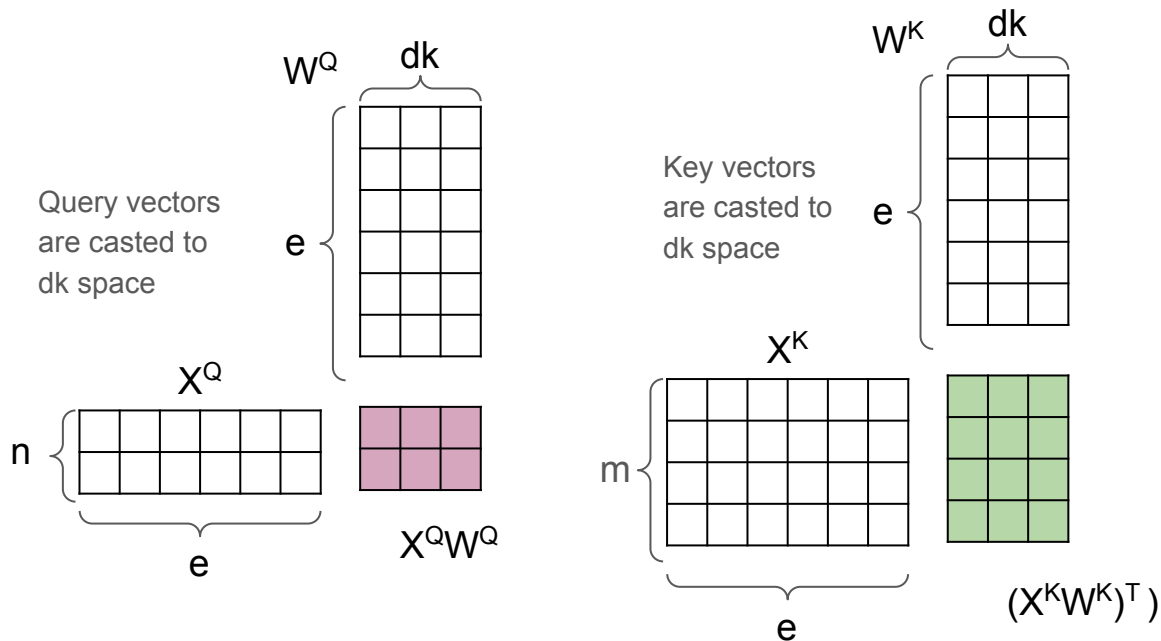
$W^Q \in \mathbb{R}^{e \times d_k}$: query weights (trainable)

$W^K \in \mathbb{R}^{e \times d_k}$: key weights (trainable)

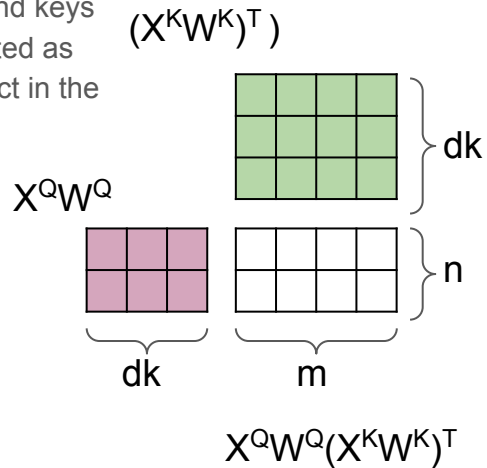
$W^V \in \mathbb{R}^{e \times d_v}$: value weights (trainable)

Attention: query-key matching

$$\text{softmax}\left(\frac{X^Q W^Q (X^K W^K)^T}{\sqrt{dk}}\right) X^V W^V$$



Match between queries and keys is calculated as dot product in the dk space



Attention: czy oba W^Q i W^K są potrzebne?

$$X^Q W^Q (X^K W^K)^T = X^Q W^Q (W^K)^T (X^K)^T$$

$W^Q (W^K)^T$ to przekształcenie liniowe, więc zamiast trenować osobno W^Q i W^K , moglibyśmy użyć jednej macierzy W^{QK} .

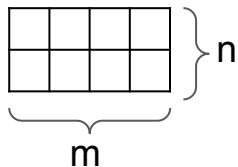
Jednakże, $W^Q \in \mathbb{R}^{e \times dk}$, $W^K \in \mathbb{R}^{e \times dk}$, więc $W^{QK} \in \mathbb{R}^{e \times e}$

W oryginalnym *Transformer* z 2017, $e=512$, $dk=64$, więc użycie W^{QK} zwiększyłoby liczbę parametrów czterokrotnie. Prawdziwy powód, dla którego używamy obu W^Q i W^K to redukcja złożoności.

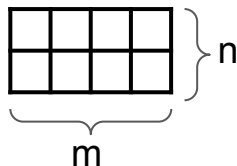
Attention: suma ważona

$$\text{softmax}\left(\frac{X^Q W^Q (X^K W^K)^T}{\sqrt{dk}}\right) X^V W^V$$

$$\frac{X^Q W^Q (X^K W^K)^T}{\sqrt{dk}}$$

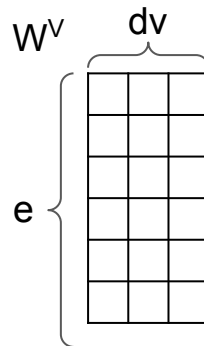
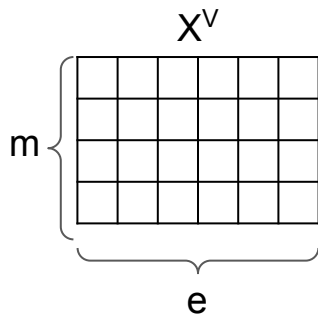


Every row sums up to one after softmax



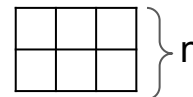
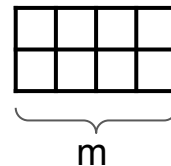
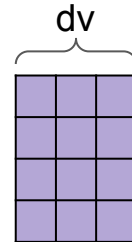
$$\text{softmax}\left(\frac{X^Q W^Q (X^K W^K)^T}{\sqrt{dk}}\right)$$

Value vectors are casted to dv space



$$X^V W^V$$

A weighted sum of $X^V W^V$ becomes the output



Attention: kolejność tokenów na wejściu nie ma znaczenia

Attention(X^Q, X^K, X^V) =

$$\text{softmax}\left(\frac{X^Q W^Q (X^K W^K)^T}{\sqrt{dk}}\right) X^V W^V$$

Permutacje X^Q zmieniają jedynie kolejność wyjść, a permutacje (X^K, X^V) nie zmieniają nic.

Sprzeczne z popularnym przekonaniem: skoro attention jest stosowane na ciągach tokenów, to na pewno zna ich kolejność?

Okazuje się, że nie: jeśli pozycja tokenu jest ważna, musi być w nim zapisana.

X^Q to zbiór, a (X^K, X^V) to zbiór par.

Multi-head attention

Wiele modułów attention przyjmujących te same dane; ich wyniki są konkatelowane.

MultiheadAttention(X^Q, X^K, X^V) =

$$\text{concat}_{i \in \{1, h\}} \left(\text{softmax} \left(\frac{X_i^Q W_i^Q (X_i^K W_i^K)^T}{\sqrt{d_k}} \right) X_i^V W_i^V \right) W^O$$

W typowych implementacjach Transformer, $d_k = d_v = e/h$, $d_o = e$.

h : number of heads

n : number of query tokens

m : number of key/value tokens

e : embedding length

d_k : dimension of key representation

d_v : dimension of value repr.

d_o : dimension of output repr.

$X^Q \in \mathbb{R}^{n \times e}$: query tokens

$X^K \in \mathbb{R}^{m \times e}$: key tokens

$X^V \in \mathbb{R}^{m \times e}$: value tokens

$W_i^Q \in \mathbb{R}^{e \times d_k}$: query weights

$W_i^K \in \mathbb{R}^{e \times d_k}$: key weights

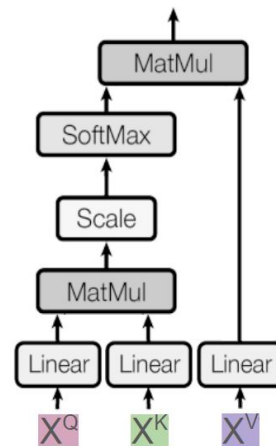
$W_i^V \in \mathbb{R}^{e \times d_v}$: value weights

$W^O \in \mathbb{R}^{d_v \times h \times d_o}$: output weights

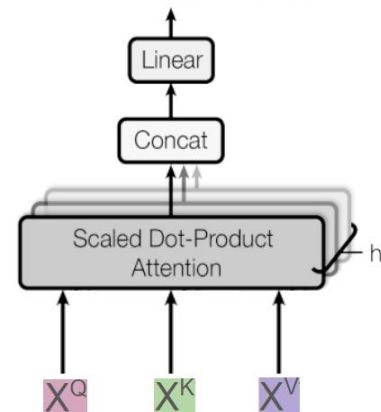
Multi-head attention: diagram

$$\text{concat}_{i \in \{1, h\}} \left(\text{softmax} \left(\frac{X^Q W_i^Q (X^K W_i^K)^T}{\sqrt{dk}} \right) X^V W_i^V \right) W^O$$

Scaled Dot-Product Attention



Multi-Head Attention



Multi-head attention: czy W_i^V są niezbędne?

Podobnie jak W^Q and W^K , zmniejszają złożoność

$$\underbrace{\text{concat}_{i \in \{1, h\}} \left(\text{softmax} \left(\frac{X^Q W_i^Q (X^K W_i^K)^T}{\sqrt{dk}} \right) X^V W_i^V \right)}_{\in \mathbb{R}^{n \times dv \times h}} W^O =$$

$$\underbrace{\text{concat}_{i \in \{1, h\}} \left(\text{softmax} \left(\frac{X^Q W_i^Q (X^K W_i^K)^T}{\sqrt{dk}} \right) X^V \right)}_{\in \mathbb{R}^{n \times e \times h}} YW^O$$

$$Y = \left\{ \begin{array}{c|ccc} & \overbrace{\hspace{1.5cm}}^{dv \times h} & & & \\ \hline W_1^V & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & \underbrace{\hspace{1.5cm}}_{W_2^V} \\ \hline 0 & 0 & 0 & \end{array} \right\} e \times h$$

$$\begin{aligned} W_i^V &\in \mathbb{R}^{e \times dv} & W^O &\in \mathbb{R}^{dv \times h \times do} \\ Y &\in \mathbb{R}^{e \times h \times dv \times h} & YW^O &\in \mathbb{R}^{e \times h \times do} \end{aligned}$$

YW^O można byłoby zastąpić macierzą $Z \in \mathbb{R}^{e \times h \times do}$, ale to zmieniłoby liczbę parameterów z $h \times e \times dv + h \times dv \times do$ na $e \times h \times do$. W oryginalnym Transformer to $h/2$ -krotny wzrost.

Warianty attention

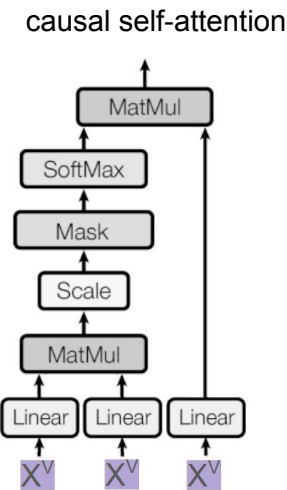
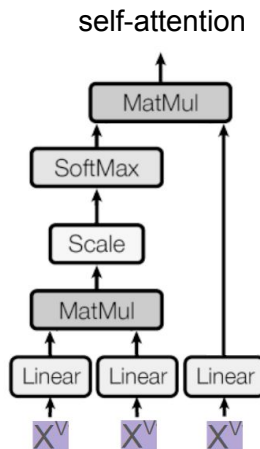
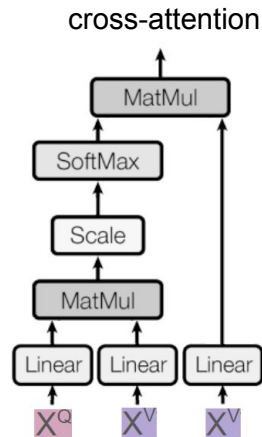
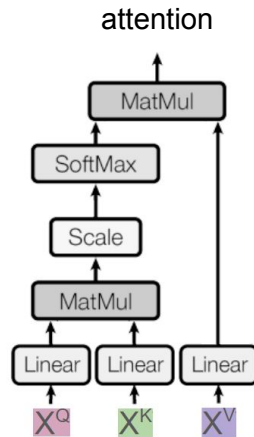
W praktyce najczęściej używane są trzy ograniczone warianty *attention*.

Cross-attention: $X^Q \neq X^K = X^V$

Self-attention: $X^Q = X^K = X^V$

Causal self-attention: $X^Q = X^K = X^V$ & “wcześniejsze” tokeny nie dopasowują się do “późniejszych”.

Trzy architektury: *encoder-only*, *decoder-only* i *encoder-decoder*.



Encoder-only Transformer

Sieć używająca self-attention,
fully-connected layers, residual connections
i layer normalization.

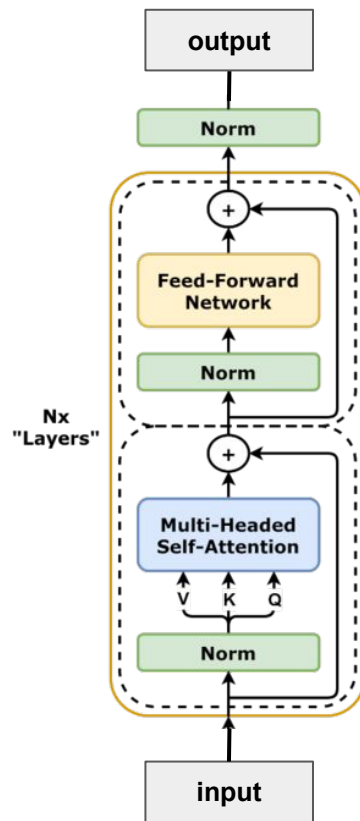
FC operują na pojedynczych tokenach.

BERT(liczba warstw L /embedding size e)

BERT-Large(24/1024): 340M parametrów

BERT-Base(12/768): 110M parametrów

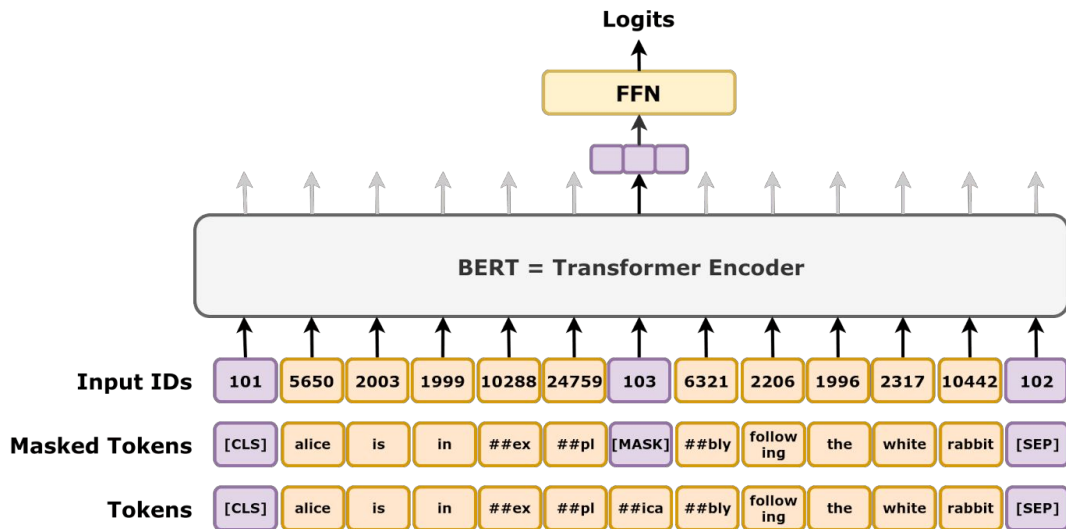
Sieci FC mają dwie warstwy ukryte: jedną
wielkości $4 \cdot e$ i jedną wielkości e .



Pre-training encoder-only Transformers

BERT jest pre-trained na dużym korpusie (3,3B słów). Jego zadanie to zamiana wszystkich tokenów [MASK] na przewidywane tokeny. Pozostałe wyjścia nie wpływają na funkcję straty.

Ucząc się wykonywać to proste zadanie, BERT poznaje zależności występujące w tekście.



use Adam with learning rate of $1e-4$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, L2 weight decay of 0.01, learning rate warmup over the first 10,000 steps, and linear decay of the learning rate. We use a dropout probability of 0.1 on all layers. We use a `gelu` acti-

We train with batch size of 256 sequences (256 sequences * 512 tokens = 128,000 tokens/batch) for 1,000,000 steps, which is approximately 40 epochs over the 3.3 billion word corpus.

<https://dvgodoy.github.io/dl-visuals/BERT/>

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (2018)

Fine-tuning encoder-only Transformers

Pre-trained BERT może być fine-tuned do innego zadania, zwykle na znacznie mniejszym zbiorze treningowym. W zależności od typu zadania, niewielkie warstwy FC mogą zostać dodane na koniec.

parameter	pre-train	fine-tune
lr	1e-4	3e-5
batch size	256	16, 32
epochs	40	2, 3, 4

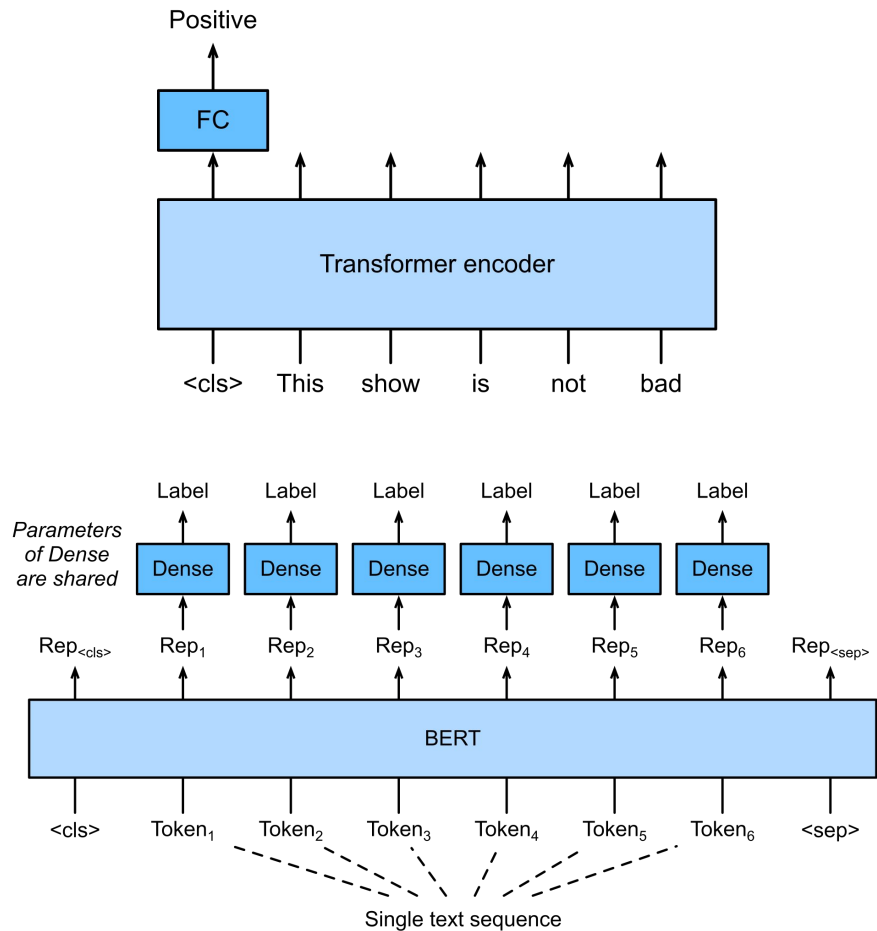
SST-2 The Stanford Sentiment Treebank is a binary single-sentence classification task consisting of sentences extracted from movie reviews with human annotations of their sentiment (Socher et al., 2013).

CoLA The Corpus of Linguistic Acceptability is a binary single-sentence classification task, where the goal is to predict whether an English sentence is linguistically “acceptable” or not (Warstadt et al., 2018).

Fine-tuning: przykłady

Klasyfikacja: zamiast opierać się na całym wyjściu, klasyfikacja bazuje jedynie na wyjściu tokena specjalnego <cls>.

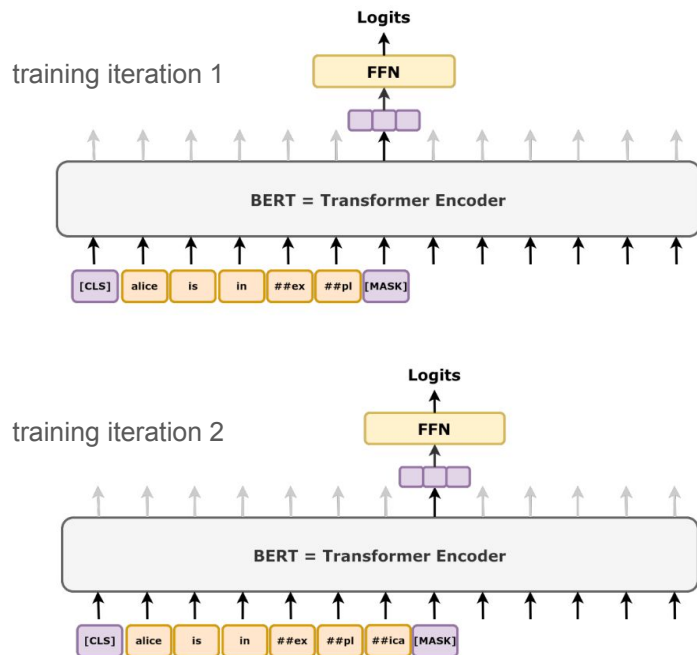
Part-of-speech tagging: wszystkie wyjścia są używane do przewidzenia części mowy odpowiadających poszczególnym tokenom.



Encoder-only Transformers kiepsko uczą się generacji

W teorii, encoder Transformer może być użyty do generacji tekstu; wystarczy dodać token [MASK] na końcu. W praktyce, przy standardowym *pre-training* prowadzi to do kiepskich wyników z powodu distribution shift.

Z drugiej strony, *pre-training* jedynie na tekstach kończących się na [MASK] jest zbyt kosztowny, bo wtedy tylko jeden token może być ewaluowany podczas jednej iteracji.



Causal self-attention

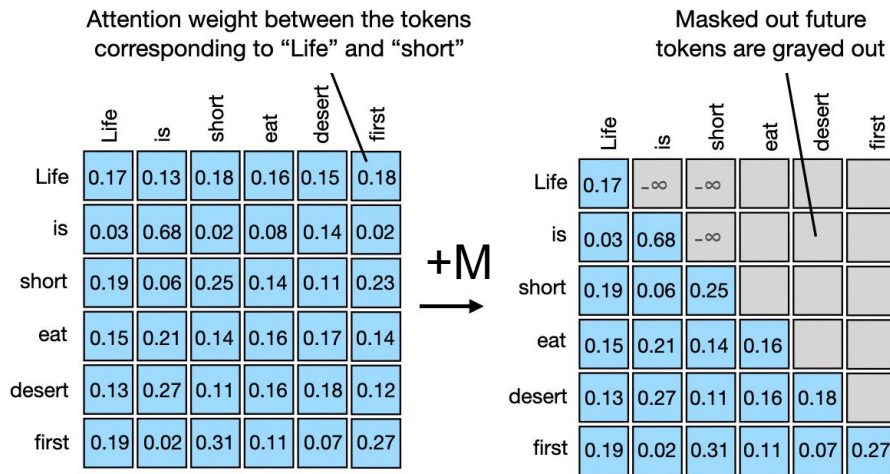
Modyfikacja attention, która ułatwia uczenie generowania.

Zamiast maskować wejścia, argumenty softmax są maskowane: każdy token może dopasować się tylko do siebie i poprzednich tokenów.

W trakcie trenowania wiele predykcji może być ewaluowanych jednocześnie.

Wprowadzona jest kolejność tokenów i jednokierunkowy przepływ informacji.

$$\text{CausalAttention}(X^Q, X^K, X^V) = \text{softmax}\left(\frac{X^Q W^Q (X^K W^K)^T}{\sqrt{dk}} + M\right) X^V W^V$$



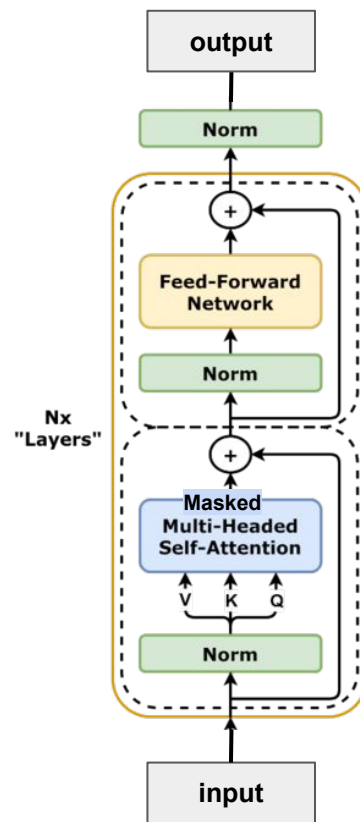
Decoder-only Transformer

To samo, co encoder-only Transformer, tylko z *causal* (a.k.a. *masked*) *self-attention*.

Pozwala to na znacznie szybsze uczenie, ale może pogarszać wyniki.

Llama3: 8B, 70B, 405B parametrów

	8B	70B	405B
Layers	32	80	126
Model Dimension	4,096	8192	16,384
FFN Dimension	14,336	28,672	53,248
Attention Heads	32	64	128
Key/Value Heads	8	8	8
Peak Learning Rate	3×10^{-4}	1.5×10^{-4}	8×10^{-5}
Activation Function	SwiGLU		
Vocabulary Size	128,000		
Positional Embeddings	RoPE ($\theta = 500,000$)		

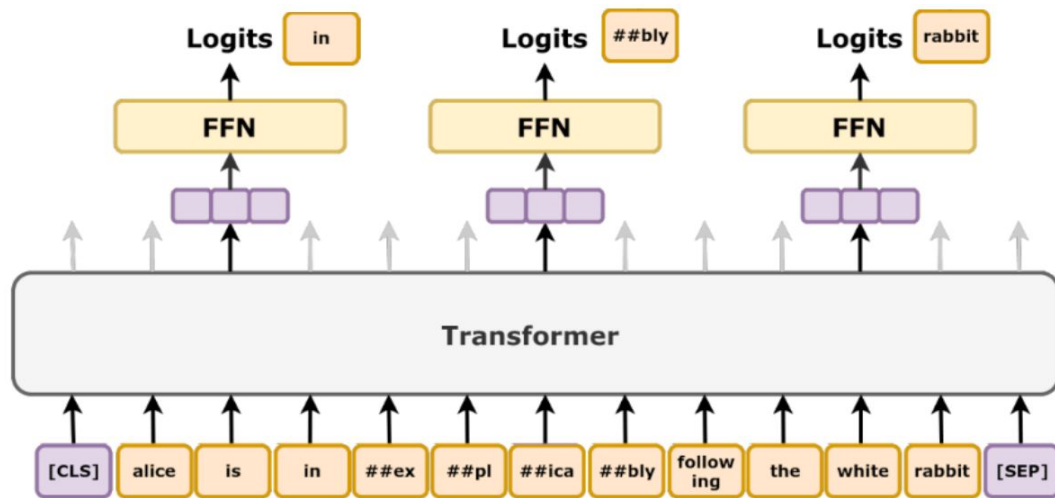


warstwy FC operują na pojedynczych tokenach, więc nie wymagają maskowania

Pre-training decoder-only Transformers

Przewidywanie następnego tokenu. Masked attention gwarantuje, że przewidywanie tokenu korzysta tylko z tokenów przed nim. Ręczne maskowanie nie jest wymagane.

Llama3: 15.6T tokens. Uczenie 405B zajęło 54 dni na 16K H100 GPU. AdamW, lr 3e-4, lr schedule. Tylko jedna epoka (wyjątek: Wikipedia i BookCrawl, gdzie dwie).



Encoder-decoder Transformer

Jedno stałe wejście jest analizowane przez encoder. Drugie wejście jest analizowane i generowane przez decoder. Architektura dobra do seq2seq: tłumaczenia, odpowiadanie na pytania, streszczenia. Wolniejsza niż decoder-only i przez to mniej popularna. Jej przewaga nad decoder-only zdaje się maleć ze skalą.

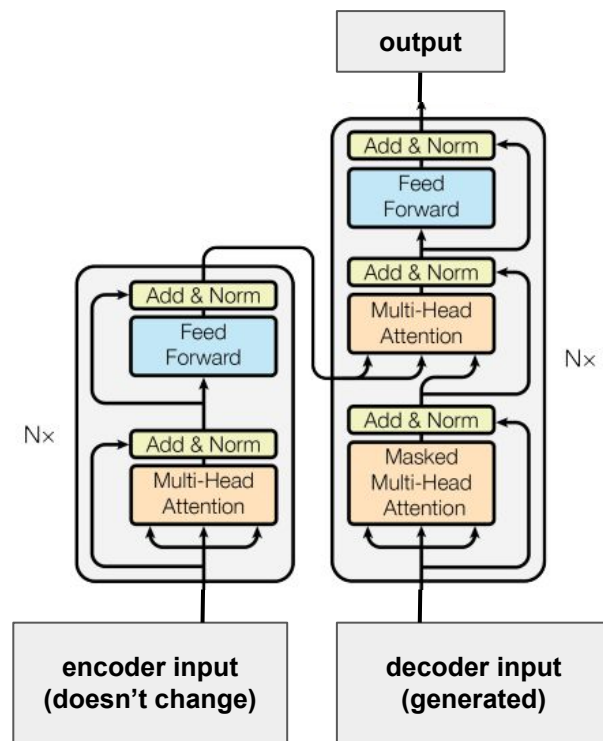
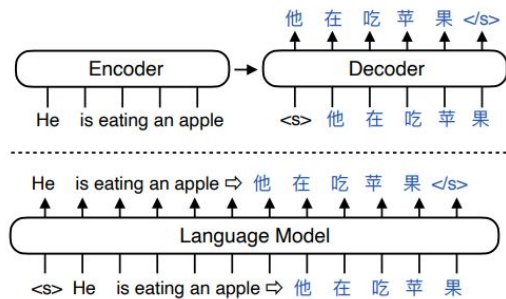


Figure 1: Encoder-Decoder (ED) framework and decoder-only Language Model (LM).

Pre-training encoder-decoder Transformers

T5 było uczone na przewidywaniu usuniętych fragmentów tekstu.

Dane: CommonCrawl, 34B tokenów.

T5: ~2.8B parametrów, L=24, e=1024, FC layer: 16384, h=32.

Może być fine-tuned do każdego zadania sentence-to-sentence, np. tłumaczenie, generacja, pisanie podsumowań itd.

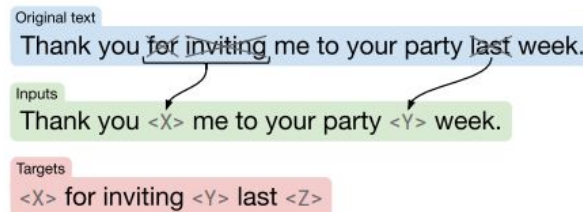


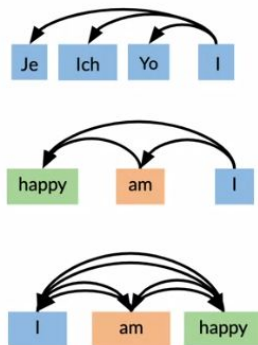
Figure 2: Schematic of the objective we use in our baseline model. In this example, we process the sentence “Thank you for inviting me to your party last week.” The words “for”, “inviting” and “last” (marked with an \times) are randomly chosen for corruption. Each consecutive span of corrupted tokens is replaced by a sentinel token (shown as $\langle X \rangle$ and $\langle Y \rangle$) that is unique over the example. Since “for” and “inviting” occur consecutively, they are replaced by a single sentinel $\langle X \rangle$. The output sequence then consists of the dropped-out spans, delimited by the sentinel tokens used to replace them in the input plus a final sentinel token $\langle Z \rangle$.

Transformer: podsumowanie architektur

Architektura	encoder-only	decoder-only	encoder-decoder
Attention	self-	masked self-	self-, masked self-, cross-
Zastosowanie*	analiza	generacja	tłumaczenie
Przykłady (parametry)	BERT (1.3B)	GPT-4 (1.8T?), Llama (405B)	T5 (2.8B), BART (406M)

Three ways of attention

- **Encoder/decoder attention:** One sentence (decoder) looks at another one (encoder)
- **Causal (self) attention:** In one sentence, words look at previous words (used for generation)
- **Bi-directional self attention:** In one sentence, words look at both previous and future words



<https://blog.sailor.plus/deep-learning/attention/>

* nie wyłącznie: można używać encoder-decoders do generacji, decoders do analizy itd.

Przekazywanie tekstu

Transformery przyjmują nieuporządkowane (z wyjątkiem causal self-attention) zbiory wektorów. Każde dane muszą być zamienione na ten format.

Text: tokenization + embedding + positional encoding. Embedding to przekształcenie liniowe $\mathbb{R}^{t \times e}$, gdzie t to rozmiar alfabetu.

universally

embeddings, choć generowane za pomocą pojedynczego przekształcenie liniowego, zawierają informacje na temat relacji tokenów

Word	universally	uniformly	explicitly	widely	invariably	duly
Cosine Similarity	1.00	0.73	0.70	0.69	0.69	0.65

Tiktokenizer

It is a truth universally acknowledged, that a single man in possession of a good fortune, must be in want of a wife.

meta-llama/Meta-Llama-3-70B

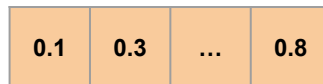
Token count

26

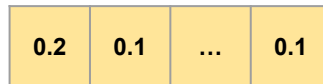
It is a truth universally acknowledged, that a single man in possession of a good fortune, must be in want of a wife.

universally

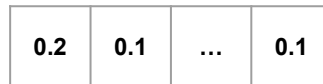
embedding



+ positional encoding (5th token)



=



Transformer input

e

Przekazywanie obrazu

Podobne do przekazywania tekstu:
obraz dzielimy na kawałki + embedding
+ positional encoding.

Positional encoding może zawierać info
o położeniu w 2D, ale 1D działa okej.

Vision Transformer to encoder-only
Transformer trenowany do klasyfikacji
(86M parametrów, potem powiększony
do 22B).

