

# Plan eksperymentów

Michał Iwaniuk, Bartłomiej Borycki

November 2025

## 1 Rozpatrywane konfiguracje hiperparametrów

### 1.1 Architektura

Przyjmijmy następujące oznaczenia

- $l^t$  – liczba warstw enkodera (liczba bloków Transformer).
- $d^h$  – rozmiar wektora ukrytego (*hidden size*)
- $d^f$  – rozmiar warstwy pośredniej w feed-forward network (FFN), często nazywany *intermediate size*. W klasycznym BERT zwykle wynosi  $4 \cdot d^h$ .
- $h$  – liczba głowic w mechanizmie wielogłowicowej uwagi.
- $d^{q|k|v}$  – wymiar przestrzeni zapytań (*query*), kluczy (*key*) oraz wartości (*value*) w każdej głowicy uwagi. (W klasycznym BERT przyjmuje się zazwyczaj  $d^q = d^k = d^v = \frac{d^h}{h}$ )

### 1.2 Wspólne parametry treningowe

**Architektura:** MLP\_dropout: 0.1, Embedding\_dropout: 0.1, Pozy-cje: `rope` (base=10000, scale=1.0). Atencja: projection\_bias=true, attn\_out\_dropout=0.1, attn\_dropout=0.0.

**MLM head:** tie\_mlm\_weights=true, mask\_p=0.15, mask\_token\_p=0.8, ran-dom\_token\_p=0.1.

**Classification head:** num\_labels=2, classifier\_dropout=0.1, pooling=`cls`, pooler\_type=`bert`.

**Trening:** batch=32 (`seq_len`=512) / 2 (`seq_len`=8192), lr=2e-5, warmup=0.1, wd=0.01, max\_grad\_norm=1.0, grad\_accum=1, AMP=true, loss=cross\_entropy.

### 1.3 Architektura Bazowa

Jako standardowa architektura małego modelu BERT przyjmiemy  $BERT_{SMALL}$  opisana w artykule *Well-Read Students Learn Better: On the Importance of Pre-training Compact Models*

$BERT_{SMALL}$ :

$l^t$	$d^h$	$d^f$	$h$	$d^{q k v}$
4	512	2048	8	512

### 1.4 Szukanie optymalnych parametrow

W oparciu o architekturę bazową przeprowadzimy eksperymenty mające na celu wyznaczenie optymalnych hiperparametrów dla mechanizmów *FAVOR+* oraz LSH, które zostaną następnie wykorzystane w fazie właściwego treningu. Optymalizacja obejmie następujące parametry:

- **Mechanizm *FAVOR+*:** Optymalizacji poddany zostanie hiperparametr `nb_features`, określający liczbę ortogonalnych wektorów projekcyjnych. Rozpatrzono wartości proporcjonalne do wymiaru  $d_k$  głowy atencji:  $0.5d_k$ ,  $1.0d_k$ ,  $2.0d_k$  oraz  $4.0d_k$ .
- **Mechanizm LSH (*Locality-Sensitive Hashing*):** Analizie poddane zostaną trzy kluczowe parametry:
  - **Rozmiar fragmentu ( $m$ ):** Testowane wartości to 32, 64 oraz 128.
  - **Liczba funkcji haszujących (`num_hashes`):** Sprawdzone zostaną konfiguracje z 4 oraz 8 haszami.
  - **Maskowanie wewnątrz fragmentu (`mask_within_chunks`):** Parametr logiczny, gdzie:
    - \* **True:** zapytania (*queries*) wewnątrz fragmentu mogą zwracać uwagę (*attend*) wyłącznie na klucze z tego samego kubelka LSH w obrębie okna.
    - \* **False:** zapytania mogą zwracać uwagę na dowolny klucz w obrębie okna.

## 2 Architektura właściwa (do klasyfikacji)

Wykorzystamy architektury modeli wyłonione w artykule *AutoTinyBERT: Automatic Hyper-parameter Optimization for Efficient Pre-trained Language Models* jako podstawa do treningu właściwych modeli klasyfikacyjnych.

Na powyższych architekturach przeprowadzone zostaną procesy uczenia dla

- **SDPA,**

$l^t$	$d^h$	$d^f$	$h$	$d^{q k v}$
5	564	1054	8	512
4	396	624	6	384
4	432	384	4	256
3	320	608	4	256

- **FAVOR** (z wykorzystaniem optymalnych parametrów wyznaczonych w poprzednim etapie),
- **LSH** (z wykorzystaniem optymalnych parametrów wyznaczonych w poprzednim etapie).

Dodatkowo, w celu stworzenia punktu odniesienia (*baseline*), wszystkie trzy wymienione typy atencji (SDPA, FAVOR, LSH) zostaną wytrenowane również na architekturze bazowej  $BERT_{SMALL}$ .

Taki sposób przeprowadzenia eksperymentów ma na celu weryfikacje czy wskazane architektury faktycznie zapewniają najkorzystniejszy stosunek szybkości działania do jakości predykcji również w przypadku metod FAVOR i LSH, analogicznie jak ma to miejsce w standardowych modelach BERT.

### 3 Trening

#### 3.1 Korpus Wikipedii do pretreningu MLM

Wybrano podzbior ok. 200 000 artykułów.

Wybierane dokumenty zawierające przynajmniej jedno ze słów kluczowych:

```
["film", "sport", "business", "science", "technology", "news"]
```

#### 3.2 TAPT + Fine-tuning(klasyfikacja): IMDB, AG News i korpus ArXiv

- **IMDB:** seq\_len=512
- **AG News:** seq\_len=512
- **ArXiv Classification:** seq\_len=8192

### 4 Podsumowanie

- Każdy z trzech mechanizmów atencji (SDPA, FAVOR, LSH) zostanie pretrenowany w pięciu wariantach architektury (bazowy  $BERT_{SMALL}$  oraz cztery konfiguracje z AutoTinyBERT). Co daje łącznie **15 modeli**.

- Po procesie TAPT + Fine-tuning trzymamy łącznie **45 finalnych modeli**, co przeloży się na 45 wyników klasyfikacji podlegających późniejszej analizie.