

# Charformer: Fast Character Transformers via Gradient-based Subword Tokenization

Орлов Александр  
Солодуха Мария  
Княжевский Владимир

# Charformer

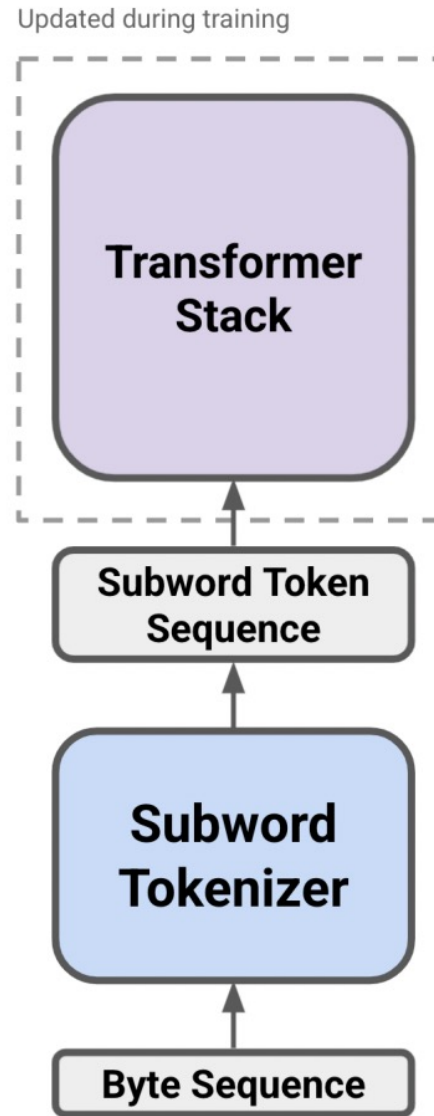
- Трансформер с токенайзером на основе градиентов
- Пайплайн обычного токенайзера
  1. Нормализация
  2. Пре-токенизация
  3. Токенизация (предобученная модель)
  4. Постпроцессинг (токены)

# Проблемы обычных токенайзеров

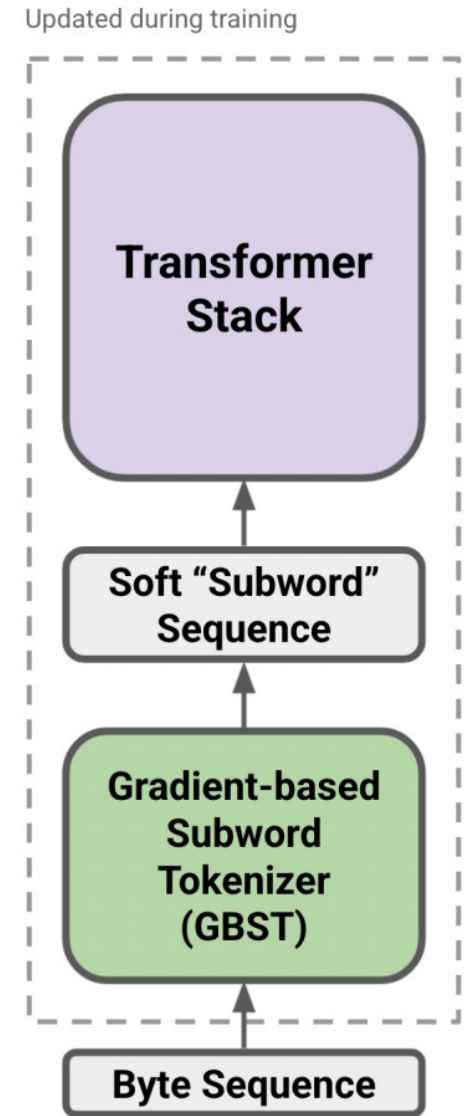
- Разделение на токены по частоте, без учета лексики
- Некоторые языки не имеют разделения на слова
- Проблема мультязычных моделей – разные языки имеют разное количество токенов
- Проблема с файнтюнингом – распределение слов на обучении может отличаться

# Charformer

- Charformer пробует решить эти проблемы
- Работает с байтами (символами)
- Разделяет слова с помощью блоков подслов-кандидатов
- Выучивает интерпретируемые подслова
- Вычислительно эффективен



**Subword Model**



**Charformer**

# Gradient-based subword tokenizer (GBST)

1. Текст
2. Байты
3. Эмбеддинги (фиксированный словарь размера 256)
4. Матрица  $X \in \mathbb{R}^{L \times d}$ , где L - кол-во байтов, d - размерность эмбеддинга

Основная идея: научить модель выполнять сегментацию данных по подсловам, выбирая по ней наиболее подходящий для всех символов блок подслов

# Gradient-based subword tokenizer (GBST)

- Блок подслов
- Функция проекция  $F$   
(Avg Pooling)
- Объединяем все блоки подслов
- Авторы предлагают брать  $s$  равным  $b$
- Считаем  $X_b$  для различных  $b$

$$X_{i:i+b} \in \mathbb{R}^{b \times d}$$

$$F : \mathbb{R}^{b \times d} \rightarrow \mathbb{R}^d$$

$$X_b = [F(X_{i:i+b}); F(X_{(i+s):(i+s)+b}); \dots]$$

# Проблемы

1. В таком пайплайне обучения мы рассматриваем не все возможные подслова
  - Можно уменьшить  $s$ , но это заметно увеличит объем вычислений
  - Вместо этого авторы предлагают применить 1D свертку к матрице  $X$
2. Мы не учитываем порядок символов
  - При применении свертки мы бонусом получаем учет порядка в функции  $F$

# Получение представлений

- Пусть  $X_{b,i}$  - блок длины  $b$ , в который попадает  $i$ -ый символ

Тогда введем модель, которая будет оценивать этот блок.	$F_R : \mathbb{R}^d \rightarrow \mathbb{R}$
С помощью неё можем оценивать все наши блоки.	$p_{b,i} = F_R(X_{b,i})$
Отсюда можем получить матрицу вероятностей	$P_i = \text{softmax}(p_{1,i}, p_{2,i}, \dots, p_{M,i}),$ $P \in \mathbb{R}^{L \times M}$
Также добавим self-attention блок, чтобы учесть вероятности соседних подслов	$\hat{P} = \text{softmax}(PP^T)P$
Наконец получаем новое представление входных данных	$\hat{X} = \sum_b^M \hat{P}_{b,i} X_{b,i}$



# Усечение

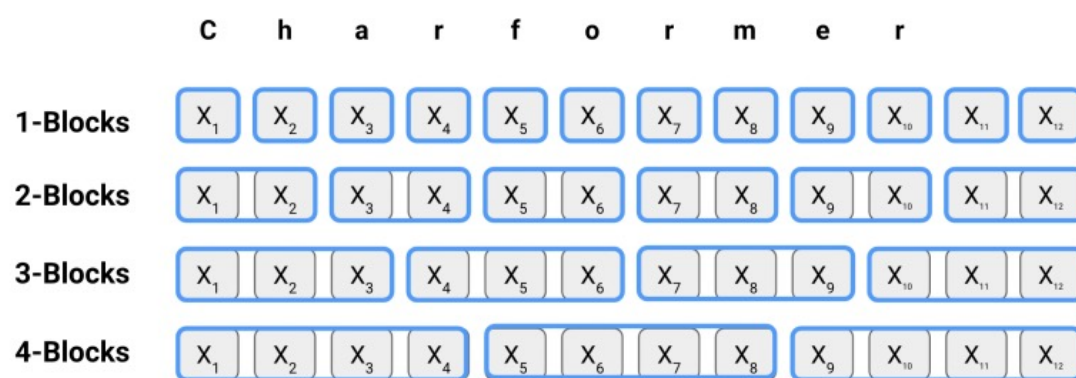
- С помощью Avg Pooling получаем векторные представления блоков символов

$$F_D : \mathbb{R}^{L \times d} \rightarrow \mathbb{R}^{\frac{L}{d_s} \times d}$$

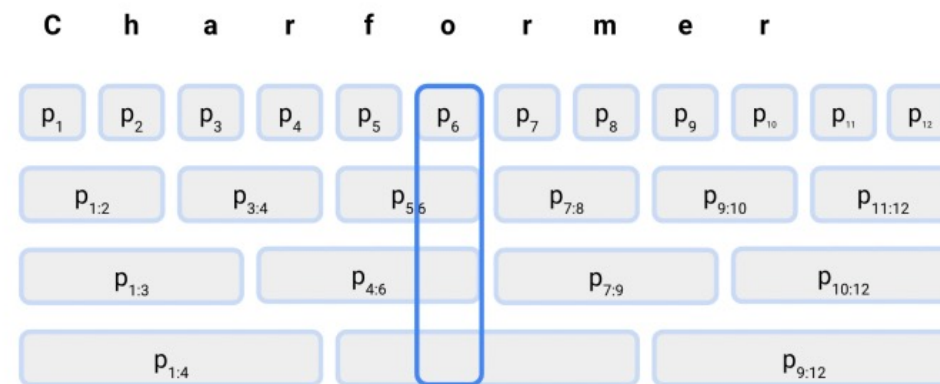
- В итоге на вход трансформеру мы будем подавать:

$$F_D(\hat{X}), \quad \hat{X} = [\hat{X}_1, \dots, \hat{X}_L]$$

# Пример



(a) Formation of subword blocks to be scored by  $F_R$ . Offsets and/or pre-GBST convolutions not shown.



(b) Block scores that have been expanded back to length  $L$ . Softmax is taken over block scores at each position  $i$  to form block weights for constructing latent subword representations.

Figure 2: Illustration of subword block formation and scoring.

# Пример

- Далее получившиеся на предыдущей картинке вероятности передаем по построенному пайплайну

$$P_i = \text{softmax}(p_{1,i}, p_{2,i}, \dots, p_{M,i})$$

$$\hat{P} = \text{softmax}(PP^T)P$$

$$\hat{X} = \sum_b^M \hat{P}_{b,i} X_{b,i}$$

Выход GBST

$$F_D(\hat{X})$$

# Эксперименты

## Стандартные английские датасеты

Model	$ \theta $	SST-2	MNLI	QNLI	MRPC	QQP	STSB	COLA	AVG
BERT <sub>Base, Subword</sub>	110M	<u>92.7</u>	<u>84.4/-</u>	88.4	86.7/-	-	-	-	-
T5 <sub>Base, Subword</sub>	220M	<u>92.7</u>	84.2/84.6	<u>90.5</u>	<u>88.9/92.1</u>	91.6/88.7	88.0	53.8	84.3
Byte-level T5 <sub>Base</sub>	200M	<u>91.6</u>	82.5/ <u>82.7</u>	88.7	<u>87.3/91.0</u>	90.9/87.7	84.3	45.1	<u>81.5</u>
Byte-level T5+Conv <sub>Base</sub>	205M	89.8	81.1/82.5	<u>89.2</u>	83.6/89.2	90.7/87.7	85.0	<u>47.1</u>	81.2
Byte-level T5+LASC <sub>Base</sub>	205M	90.0	80.0/80.8	87.1	82.8/88.1	89.0/85.4	83.7	25.3	77.0
CHARFORMER <sub>Base</sub>	203M	<u>91.6</u>	<u>82.6/82.7</u>	89.0	<u>87.3/91.1</u>	<u>91.2/88.1</u>	<u>85.3</u>	42.6	81.4
Byte-level T5 <sub>SBase</sub>	133M	91.2	<u>83.9/83.7</u>	90.9	85.5/89.2	91.1/88.1	85.7	49.3	82.6
CHARFORMER <sub>SBase</sub>	134M	<u>91.5</u>	83.7/ <u>84.4</u>	<u>91.0</u>	<u>87.5/91.4</u>	<u>91.4/88.5</u>	<u>87.3</u>	<u>51.8</u>	<u>83.6</u>

# Эксперименты

## Классификация комментариев с Wiki

Model	Civil Comments	Wiki Comments
$T5_{Base, Subword}$	81.2 / -	91.5 / -
Byte-level $T5_{Base}$	82.8 / 78.7	<u>93.2</u> / 75.4
Byte-level $T5+LASC_{Base}$	82.9 / 78.2	93.0 / 75.0
$CHARFORMER_{Base}$	<u>83.0</u> / <u>78.8</u>	92.7 / <u>79.7</u>
$CHARFORMER_{SBase}$	83.0 / 78.9	93.5 / 75.5

# Эксперименты

## Классификация больших документов

Model	IMDb	News
$T5_{Base, Subword}$	94.2	93.5
Byte-level $T5_{Base}$	<u>91.5</u>	93.6
Byte-level $T5+LASC_{Base}$	91.1	93.5
$CHARFORMER_{Base}$	<u>91.5</u>	<u>94.0</u>
$CHARFORMER_{SBase}$	94.4	94.1

# Эксперименты

## Многоязычный перевод текстов

Model	$ \theta $	In-Language	Translate-Train-All				Zero-Shot	
		TyDiQA-GoldP	XQuAD	MLQA	XNLI	PAWS-X	XNLI	PAWS-X
mBERT <sub>Base</sub> (Subword)	179M	77.6/68.0	-/-	-/-	-	-	65.4	81.9
mT5 <sub>Base</sub> (Subword)	582M	<u>80.8/70.0</u>	75.3/59.7	67.6/48.5	75.9	89.3	<u>75.4</u>	<u>86.4</u>
Byte-level T5 <sub>Base</sub>	200M	75.6/65.4	68.6/54.3	61.8/44.4	69.4	87.1	57.4	80.9
Byte-level T5+LASC <sub>Base</sub>	205M	70.6/59.7	66.8/52.1	58.8/41.1	67.9	84.8	55.2	79.0
CHARFORMER <sub>Base</sub>	203M	<u>75.9/65.6</u>	<u>70.2/55.9</u>	<u>62.6/44.9</u>	<u>71.1</u>	<u>87.2</u>	<u>57.6</u>	<u>81.6</u>
CHARFORMER <sub>SBase</sub>	134M	79.1/68.8	73.6/59.0	66.3/48.5	72.2	88.2	66.6	<u>85.2</u>
CHARFORMER <sub>SBase,LongPT</sub>	134M	<u>81.2/71.3</u>	<u>74.2/59.8</u>	<u>67.2/49.4</u>	<u>72.8</u>	<u>88.6</u>	<u>67.8</u>	83.7

# Производительность

## Модели на основе подслов

Model	Batch Size	$L$	$d_s$	$ \theta $	Speed (steps/s)	FLOPS
mT5 <sub>Base</sub> (Subword)	1024	1024	-	582M	1.54	$1.3 \times 10^{15}$
CHARFORMER <sub>SBase</sub>	1024	2048	2	134M	1.98	$4.3 \times 10^{14}$
CHARFORMER <sub>SBase, LongPT</sub>	2048	2048	2	134M	1.01	$4.3 \times 10^{14}$



# Производительность

## Модели на основе байтов

Model	$L$	$d_s$	$ \theta $	Speed (steps/s)	FLOPS	Peak Mem.
T5 <sub>Base</sub> (Subword)	512	-	220M	9.3	$1.1 \times 10^{13}$	-
Byte-level T5 <sub>Base</sub>	1024	1	200M	8.2	$2.9 \times 10^{13}$	3.09GB
Byte-level T5+LASC <sub>Base</sub>	1024	4	205M	15	$9.9 \times 10^{12}$	1.62GB
CHARFORMER <sub>Base</sub>	1024	2	206M	11	$1.6 \times 10^{13}$	1.95GB
CHARFORMER <sub>Base</sub>	1024	3	203M	15	$1.1 \times 10^{13}$	1.63GB
CHARFORMER <sub>SBase</sub>	1024	2	134M	14	$1.3 \times 10^{13}$	1.73GB
CHARFORMER <sub>SBase</sub>	1024	3	134M	20	$8.7 \times 10^{12}$	1.34GB

Спасибо за внимание!