

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ»  
(СПбГУ)

Образовательная программа бакалавриата «Науки о данных»



**Отчёт по практике**  
**Учебная практика (научно-исследовательская работа)**

Выполнил студент 3 курса бакалавриата  
(группа 22.Б05-мкн)  
**Михеев Артём Антонович**

Научный руководитель:  
Кудряшов Александр

Санкт-Петербург

Построение обратимых эмбедингов предложений для  
финансовых новостей

Содержание

<b>Введение</b>	<b>3</b>
<b>Основные результаты практики</b>	<b>4</b>
<i>Dataset</i>	<i>4</i>
<i>Метрики</i>	<i>5</i>
<i>Размерности.</i>	<i>6</i>
<i>Архитектура моделей</i>	<i>6</i>
<i>Обучение</i>	<i>9</i>
<i>Графики</i>	<i>10</i>
<i>Результаты</i>	<i>12</i>
<i>Генерация. GREEDY vs BEAM-SEARCH</i>	<i>14</i>
<i>Примеры</i>	<i>14</i>
<b>Заключение</b>	<b>15</b>

# Введение

Обращать эмбединги токенов довольно просто. Находим эмбединг в матрице эмбедингов и номер строчки является порядковым номером токена. Используя токенизатор. декодируем и получаем символы.

Предложений бывает очень много, и никакой памяти не хватит для хранения всех эмбедингов. При этом условная обратимость бывает полезной(обратить полностью не получится, так как какая-то часть информации теряется, но получить схожее по смыслу хотелось бы).

Предположим, мы проектируем мультимодальную модель для обработки данных разной природы. К примеру, такая модель может на основе финансовых новостей предсказывать различные показатели: индексы фондового рынка, курс валют. Или генерировать текст - объяснение, почему произошло изменение этих финансовых показателей. Нам важно уловить общий смысл, получить одно представление (вектор) этого объяснения. Из него мы хотим получать текст, не углубляясь на уровень токенов.

# Основные результаты практики

Наша задача спроектировать и обучить как энкодер для получения одного представления, так более важная часть модели - decoder для смыслового восстановления.

## ***Dataset.***

Будем обучаться на финансовых новостях [Financial News Headlines Data](#) с kaggle.

Заголовки этих наборов данных, взятые с официальных сайтов CNBC, The Guardian и Reuters, отражают обзор экономики и фондового рынка США за каждый день за 2018-2020.

Пример:

*U.S. economy faces significant risks, long road to recovery: IMF staff.*

*UBS, Morgan Stanley expected to lead Vodafone Tower IPO: sources.*

*As big U.S. banks let customers delay payments, loan losses remain unclear.*

*China says it will act to protect its interests after UK Huawei ban.*

Всего 53370 заголовков. Делим данные на train и validation, validation size = 0.1.

## ***Метрики.***

Надо оценивать на сколько декодируемый текст похож на оригинал. Будем оценивать следующие метрики:

1) Совпадение слов.

*ACCURACY.* Сопоставляем токены по позициям.

## 2) Семантическая близость.

*BERTScore*.

Модель: `'microsoft/deberta-xlarge-mnli'`

- Получает эмбединги для каждого токена из двух предложений.
- Вычисляет семантическое соответствие токенов, сопоставляя каждый токен из output с ближайшим по смыслу токеном в input (и наоборот).  $O(T^2)$ .
- Итоговые P, R, F1 — агрегированные значения сходства между токенами.

`bert_score("The car was damaged in the accident.", "The vehicle got wrecked during the crash.") = 0.91`

`bert_score("The economy is growing rapidly.", 'My dog loves running in the park.') = 0.35`

*PARAPH-SIM*.

Модель: `'sentence-transformers/paraphrase-mpnet-base-v2'`

- Каждое предложение сжимается в один вектор фиксированной размерности 768.
- Вычисляется косинусное сходство между двумя векторами:

$$\text{cos\_sim}(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

`paraph_sim("The car was damaged in the accident.", "The vehicle got wrecked during the crash.") = 0.821`

`paraph_sim("The economy is growing rapidly.", 'My dog loves running in the park.') = 0.17`

## ***Размерности.***

B - batch\_size

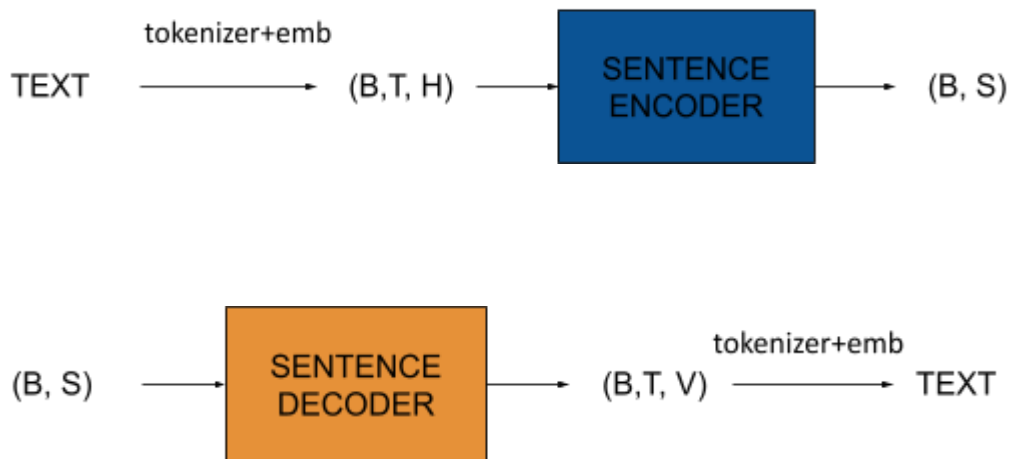
T - длина предложения в токенах с padding для обучения по бачам.

H - размерность эмбеддингов токенов.

S - размерность эмбеддингов предложений.

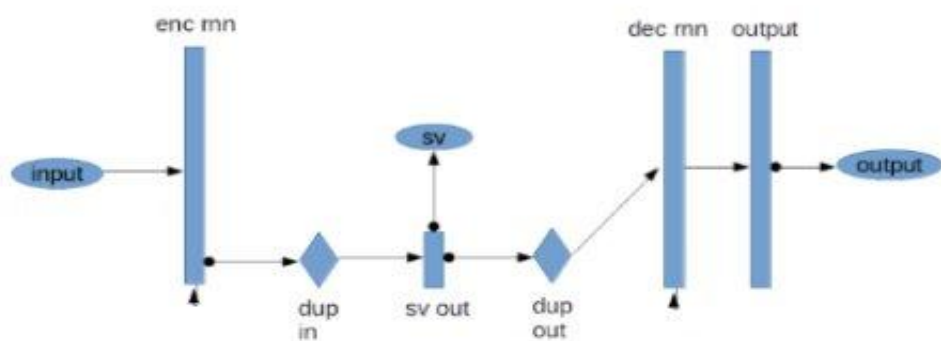
V - размер словаря из токенов.

На вход модели подаётся тензор размера (B, T, H). На выходе logits (B,T,V)



## ***Архитектура моделей.***

Прошлые подходы на основе RNN показали способность учиться, однако качество по метрикам было не очень.



## *Новые подходы на основе transformer архитектуры.*

За основу encoder был взят **FinBERT** — это предварительно обученная модель для анализа настроений (сентимента) финансовых текстов.

На выходе **FinBERT** эмбединги для каждого токена и **cls** токен.

Агрегируем с обучаемыми весами выход **FinBERT** с помощью **Attention pooling** слоя в эмбединг предложения. Также для эксперимента подаем только **cls**.

Decoder обучаем с нуля двумя способами.

(1) За один проход. Decoder на основе transformer\_encoder.

**sent\_emb**: эмбединг предложения. Расширяем до T, учитывая позиции (прибавляем позиционные эмбединги) и подаем в decoder для self-attention.

Предсказываем все токены на основе **sent\_emb**.

$P(\text{tgt}[1:T] \mid \text{sent\_emb})$

(2) Авторегрессионно. Decoder на основе transformer\_decoder.

**sent\_emb**: эмбединг предложения. Подаем в decoder как вход для cross-attention.

**tgt**: токены входа. Подаем в decoder для masked self-attention и cross-attention.

Предсказываем следующий токен на основе предыдущих и **sent\_emb**.

$P(\text{tgt}[N+1] \mid \text{tgt}[1:N], \text{sent\_emb})$

### Attention Pooling

```
class AttentionPooling(nn.Module):
    def __init__(self, hidden_dim, output_dim):
        super().__init__()
        self.attn = nn.Linear(hidden_dim, 1)
        self.proj = nn.Sequential(
            nn.Linear(hidden_dim, output_dim),
            nn.GELU(),
            nn.LayerNorm(output_dim),
        )
```

```

def forward(self, x, mask=None):
    attn_scores = self.attn(x).squeeze(-1) # (B, T)
    if mask is not None:
        attn_scores = attn_scores.masked_fill(mask == 0, -1e9)
    attn_weights = torch.softmax(attn_scores, dim=1) # (B, T)
    return torch.sum(attn_weights.unsqueeze(-1) * self.proj(x),
dim=1) # (B, S)

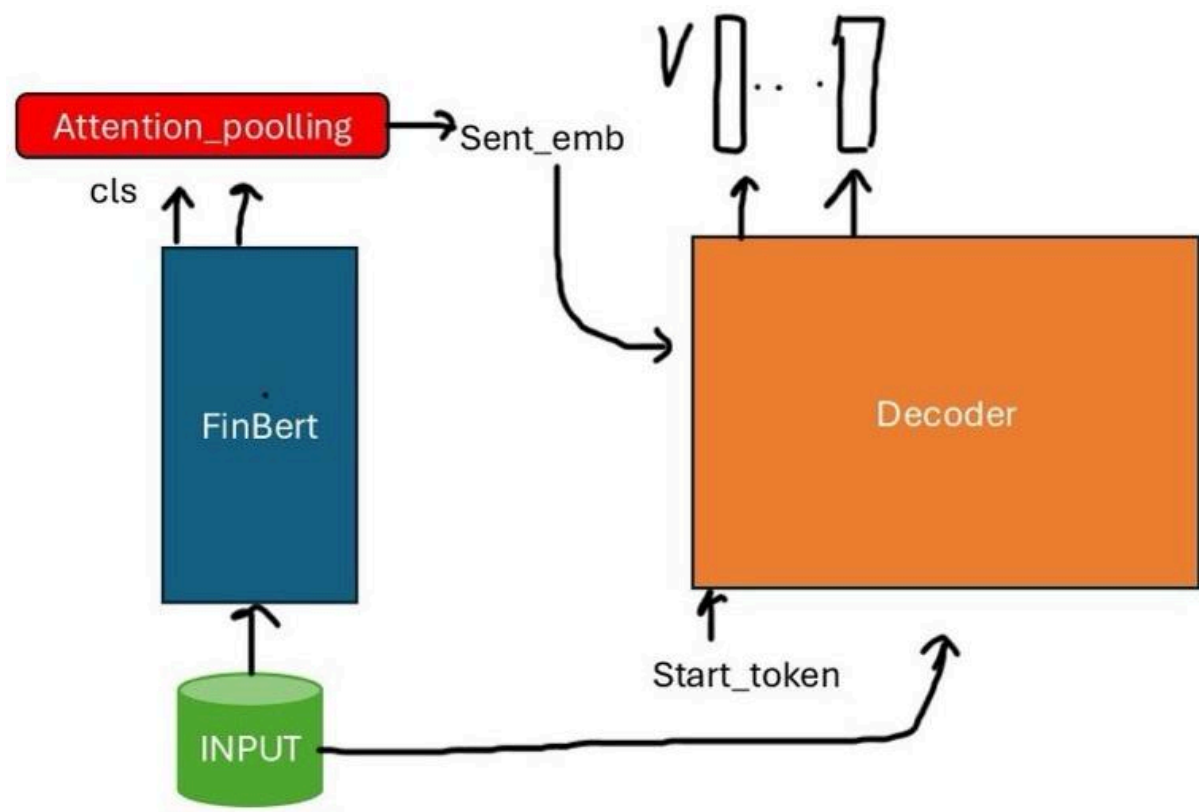
```

Последний слой для получения *logits*.

```

self.fc_out = nn.Sequential(
    nn.Dropout(0.5),
    nn.Linear(sent_dim, vocab_size)
)

```





## Обучение

Loss: *cross\_entropy\_loss*

+ fine\_tuning with *sent\_sim\_loss*=cos(sent\_emb(input), sent\_emb(decoded))

optimizer: AdamW (learning\_rate=0.0001, weight\_decay=0.01)

lr\_scheduler: ReduceLROnPlateau (mode='max', factor=0.1,  
min\_lr=1e-5, patience=5)

batch\_size: 128-200

*Возможные улучшения. Loss на уровне предложений*

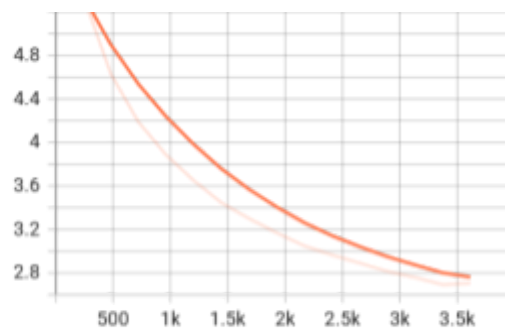
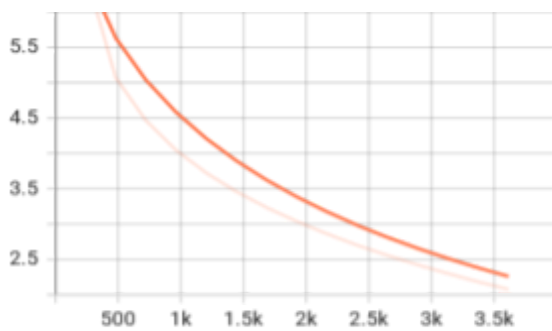
- 1) *sent\_sim\_loss* для эмбедингов из предобученных моделей, как *'paraphrase-mpnet-base-v2'* из *paraph\_sim*, вместо нашей модели.
- 2) Вместо входных токенов в *cross\_entropy\_loss* можно подавать токены перефраз полученных из *'paraphrase-mpnet-base-v2'* в количестве равном T (обрезка или padding до тензора BxT). Это поможет восстанавливать перефразы.

# Графики

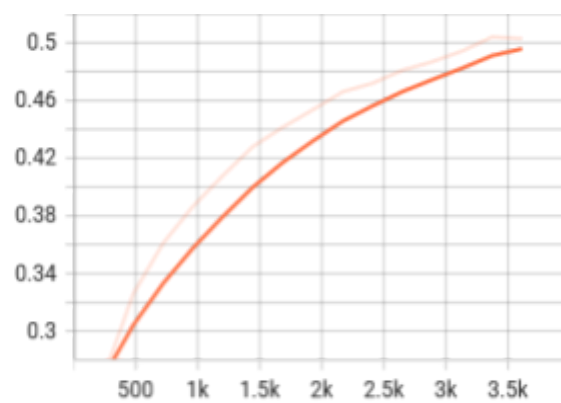
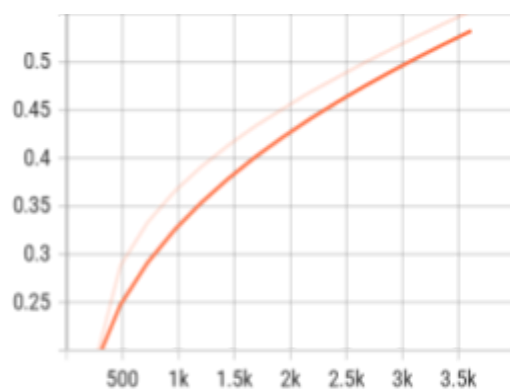
*Train*

*VALIDATION*

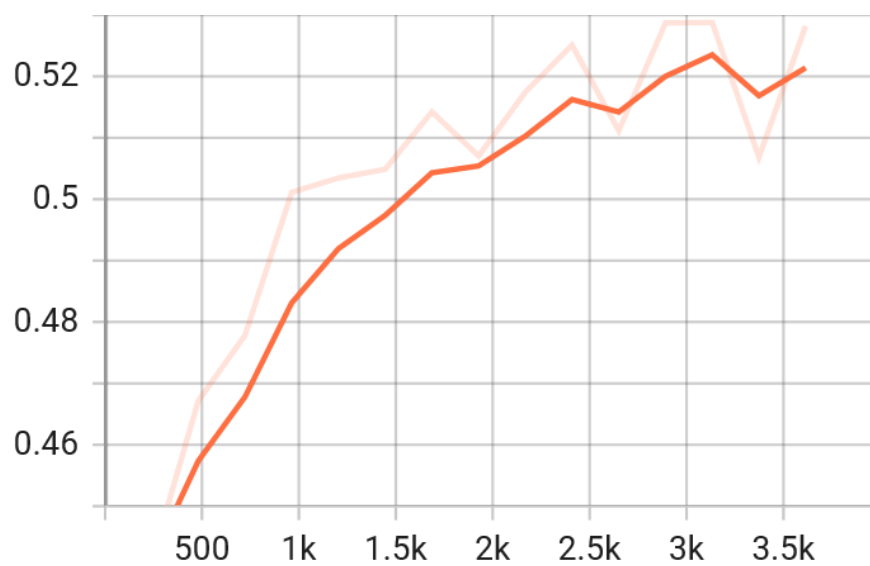
*Loss*



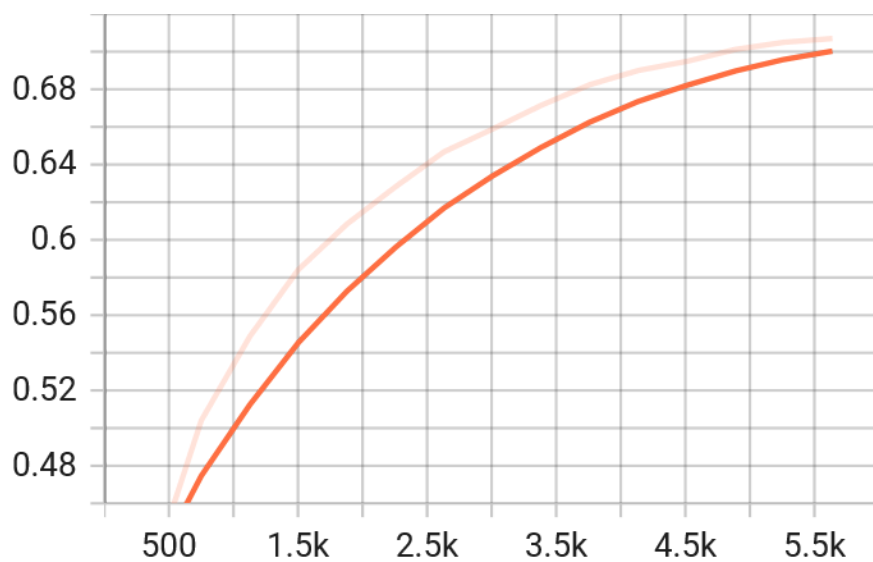
*Accuracy*



*BERT-score*



*Top-5*



## Результаты

Зависимость метрик от архитектуры и размера эмбединга предложения.

<i>model</i>	<i>dim</i>	<i>Valid_accuracy</i>	<i>BERT_score</i>
<i>За один проход</i>	<i>1568</i>	<i>0.357</i>	
<i>Аutoreгрессионно CLS</i>	<i>768</i>	<i>0.406</i>	<i>0.58</i>
<i>Аutoreгрессионно attention_pooling</i>	<i>768</i>	<i>0.502</i>	<i>0.65</i>
	<i>2560</i>	<i>0.552</i>	<i>0.675</i>

Аторегрессионные модели по качеству сильно лучше, однако медленнее.

*Attention pooling* увеличивает качеству по сравнению с *cls* эмбедингом. Так же этот слой позволяет менять размерность *sent\_emb*.

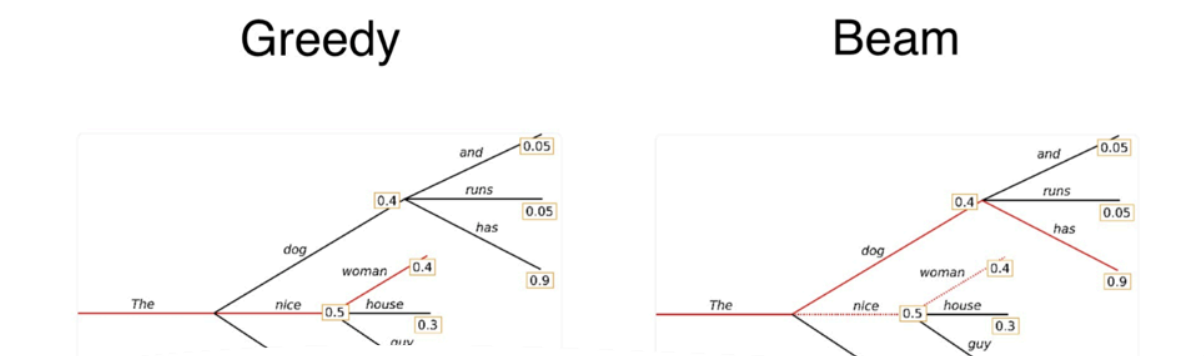
Увеличивая размерность вектора, метрики растут.

Как часто входные токены в первых k токенах по вероятности.

<i>Top-k</i>	<i>1</i>	<i>3</i>	<i>5</i>	<i>10</i>	<i>15</i>	<i>25</i>	<i>50</i>	<i>100</i>
<i>Model_768</i>	<i>0.50</i>	<i>0.63</i>	<i>0.68</i>	<i>0.73</i>	<i>0.75</i>	<i>0.78</i>	<i>0.81</i>	<i>0.84</i>
<i>Model-2560</i>	<i>0.55</i>	<i>0.68</i>	<i>0.72</i>	<i>0.77</i>	<i>0.79</i>	<i>0.81</i>	<i>0.84</i>	<i>0.86</i>

Заметно, что top-5 сильно лучше, чем top-1. С увеличением k скорость роста метрик сильно уменьшается. Возможно, что top-k для маленького k содержит синонимы или перефразы.

## ***Генерация. GREEDY vs BEAM-SEARCH***



Изменив алгоритм генерации с жадного на поиск наиболее вероятной последовательности, модель стала реже заикаться (повторять слова).

bert\_score: 0.66 → 0.7

paraph\_sim: 0.64 → 0.69

## *Примеры.*

<i>INPUT</i>	<i>DECODED</i>
<i>consumer giants spurn risks to chase online subscribers</i>	<i>consumer giants spurn consumer risks to online subscribers</i>
<i>t - mobile beats estimates for phone subscribers'</i>	<i>t - mobile beats estimates for phone subscribers</i>
<i>australia watchdog says vodafone misled customers over digital purchases</i>	<i>australia watchdog says hsbc over wireless firms misled customers payments</i>
<i>bayer says monsanto likely kept files on influential people across europe</i>	<i>bayer says monsanto likely put on some countries sick people in europe</i>
<i>wall street rises on economic data, easing trade worries</i>	<i>wall street rises on trade worries, fed data data</i>
<i>councils to bid for share of fund to invigorate english high streets</i>	<i>councils to buy out bid to fund for lse's shares of london's</i>

Заголовки (1-2) целиком восстанавливаются. Остальные примеры имеют похожие слова. В заголовке (3) компания Vodafone восстановлена как фирма, занимающаяся беспроводной связью. Payments почти синоним purchases. В (6) примере english восстанавливается как london's.

# Заключение

Основные цели проекта были достигнуты. Была показана возможность сжимать текст в один эмбедринг с возможностью частичного восстановления с BERTScore=0.7. Эмбедринги токенов из FinBert или других моделей имеют общую информацию, это позволило нам сжимать новостные заголовки в среднем, состоящие из 30 токенов в один вектор размерности 2560.  $30 \times 768 \rightarrow 2560$ , сжатие в 10 раз.

Можно добавить *loss* на уровне предложений, чтобы улучшить сжимаемость за счёт перефраз. Увеличить размерность вектора, или чуть изменить архитектуру, например добавить специальные прототокены как в статье[2]. PEFT для FinBert, использовать предобученные LLM(*Llama*, *Pythia*) как более мощные decoder-архитектуры.

Данный подход открывает новые подходы мультимодальных моделях. Так для объяснения изменения финансовых показателей можно предсказывать один вектор. И наш декодер получит текст, сохранивший общий смысл.

Список использованных литературных источников  
и информационных материалов.

[\[1\] Cramming 1568 Tokens into a Single Vector and Back Again: Exploring the Limits of Embedding Space Capacity](#)

[\[2\] Exploring the Latent Capacity of LLMs for One-Step Text Generation](#)