

# Методы предобработки текстовых данных для ускорения обучения языковых моделей

Сурков Максим Константинович

Научный руководитель: Ямщиков Иван Павлович

Санкт-Петербургская школа физико-математических и компьютерных наук  
НИУ ВШЭ СПб

17 марта 2021 г.

# Обработка естественного языка в реальной жизни

- социальные сети
- электронная почта
- службы доставки
- голосовые помощники
- переводчики
- чат боты



- ❶ классификация последовательностей
  - спам
  - грубая речь<sup>1</sup>
- ❷ генерация выходной последовательности из исходной
  - машинный перевод
  - ответы на вопросы
- ❸ выделение информации из последовательностей
  - выделение именованных сущностей<sup>2</sup>

---

<sup>1</sup>G. H. Paetzold et al., SemEval'19 Task 5: Hate Speech Identification with RNN.

<sup>2</sup>Vikas Yadav et al., SemEval'19 Task 12: Deep-Affix Named Entity Recognition of Geolocation Entities. ACL'19

# Современные методы решения задач обработки естественного языка

- ❶ Механизм внимания<sup>1</sup>
- ❷ **BERT** (Google)<sup>2</sup>
- ❸ GPT-3 (OpenAI)<sup>3</sup>

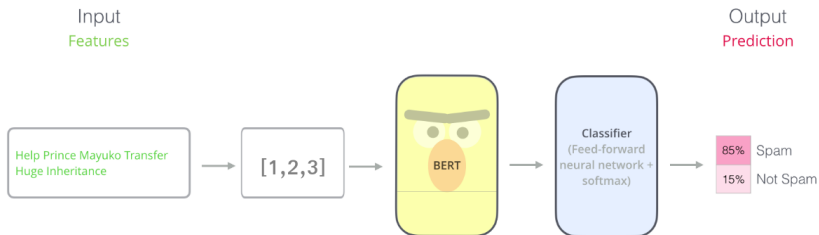
---

<sup>1</sup>Ashish Vaswani et al., Attention Is All You Need, 2017

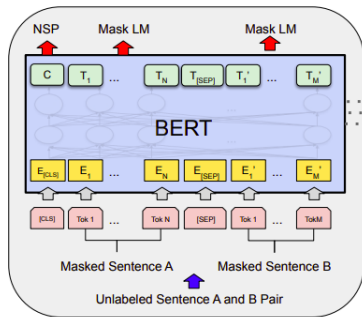
<sup>2</sup>Jacob Devlin et al., BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, 2019

<sup>3</sup>Tom B. Brown et al., Language Models are Few-Shot Learners, 2020

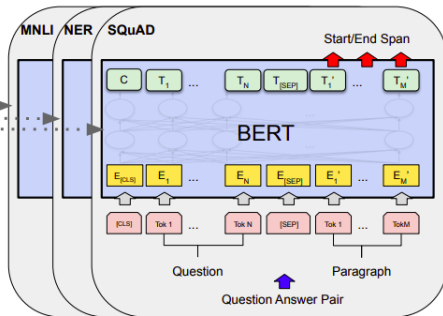
# BERT. Использование



# BERT. Обучение



Pre-training



Fine-Tuning

# BERT. Требуемые ресурсы

- количество параметров: 110M – 340M
- время на предобучение: от 2-4 дней до 1-2 недель<sup>1</sup>
  - мировой рекорд: 47 минут на **1472** V100 GPU<sup>2</sup>
- время на дообучение: 1-2 дня
- размеры данных:

Датасет	Размер
Wikipedia	3-600M
HND	600k-2M
s140	1.6M
IWSLT	200-230k
QQP	364k
MNLI	393k

<sup>1</sup>При использовании 1x-4x GPU Nvidia Tesla V100 32Gb

<sup>2</sup><https://developer.nvidia.com/blog/training-bert-with-gpus>

# BERT. Существующие методы оптимизации

- квантизация<sup>1</sup>
- дистилляция<sup>2</sup>
- прунинг<sup>3</sup>

---

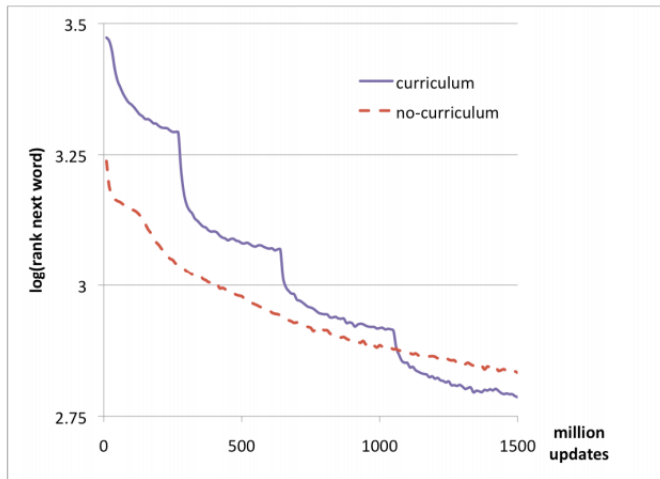
<sup>1</sup>Sheng Shen et al., Q-BERT: Hessian Based Ultra Low Precision Quantization of BERT, 2019

<sup>2</sup>Victor Sanh et al., DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter, 2020

<sup>3</sup>Hassan Sajjad et al., Poor Man's BERT: Smaller and Faster Transformer Models, 2020

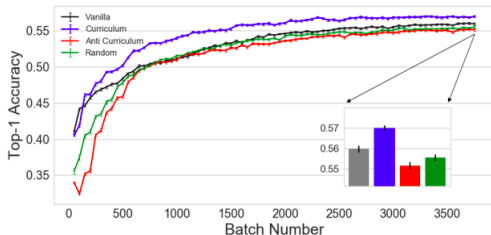


# Обучение с расписанием. Начало



Y. Bengio et al., Curriculum learning, 2009

- компьютерное зрение<sup>1</sup>



- обучение с подкреплением<sup>2</sup>

- глубокое обучение<sup>3</sup>

<sup>1</sup>Guy Hach Cohen, Daphna Weinshall, On The Power of Curriculum Learning in Training Deep Networks, 2019

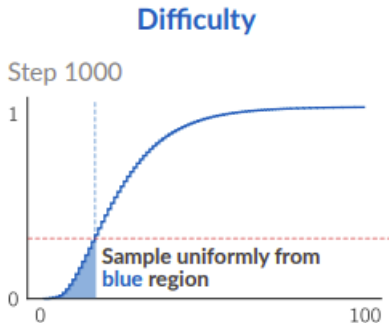
<sup>2</sup>Sanmit Narvekar et al., Curriculum Learning for Reinforcement Learning Domains: A Framework and Survey, 2020

<sup>3</sup>Mermer et al., Scalable Curriculum Learning for Artificial Neural Networks, 2017

# Обучение с расписанием в обработке языка

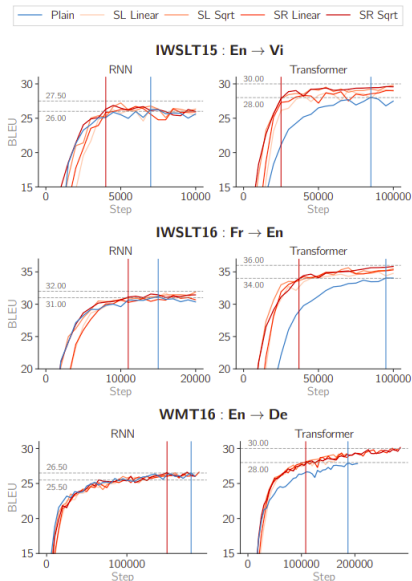
- Задача: машинный перевод
- Модель: BERT, LSTM
- Датасеты: IWSLT'15, IWSLT'16, WMT'16
- Алгоритм:

- 1 сортируем тексты по сложности (длина, логарифм вероятности правдоподобия)
- 2 в течение  $T$  шагов (рассмотрим шаг  $t$ )
  - считаем  $c(t) \in [0, 1]$
  - строим батч из  $c(t)$  **первых** текстов корпуса
  - шаг обучения



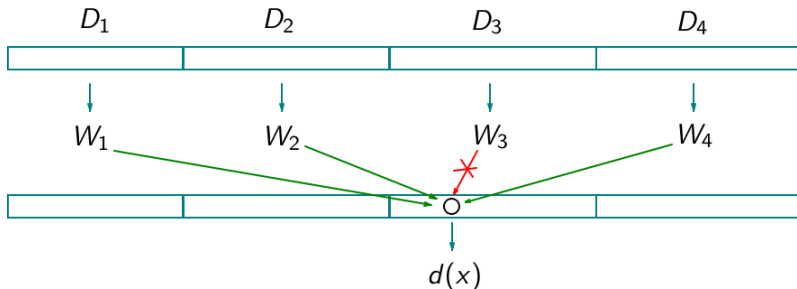
Е. А. Platanios et al., Competence-based Curriculum Learning for Neural Machine Translation, ACL'19

# Обучение с расписанием в обработке языка



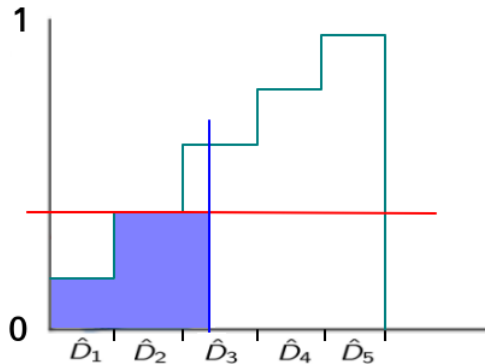
# Обучение с расписанием в обработке языка

- Задача: классификация
- BERT
- Датасеты: SQuAD 2.0, NewsQA, GLUE
- Алгоритм: в течение  $T$  шагов



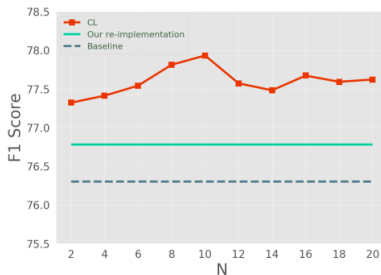
Benfeng Xu et al., Curriculum Learning for Natural Language Understanding, ACL'20

# Обучение с расписанием в обработке языка



# Обучение с расписанием в обработке языка

	MNLI-m	QNLI	QQP	RTE	SST-2	MRPC	CoLA	STS-B	Avg
<i>results on dev</i>									
BERT Large	<b>86.6</b>	92.3	91.3	70.4	93.2	88.0	60.6	90.0	84.1
BERT Large*	<b>86.6</b>	92.5	91.5	74.4	93.8	91.7	63.5	90.2	85.5
BERT Large+CL	<b>86.6</b>	<b>92.8</b>	<b>91.8</b>	<b>76.2</b>	<b>94.2</b>	<b>91.9</b>	<b>66.8</b>	<b>90.6</b>	<b>86.4</b>
<i>results on test</i>									
BERT Large	<b>86.7</b>	91.1	89.3	70.1	<b>94.9</b>	89.3	60.5	87.6	83.7
BERT Large*	86.3	92.2	<b>89.5</b>	70.2	94.4	89.3	60.5	87.3	83.7
BERT Large+CL	<b>86.7</b>	<b>92.5</b>	<b>89.5</b>	<b>70.7</b>	94.6	<b>89.6</b>	<b>61.5</b>	<b>87.8</b>	<b>84.1</b>



# Обучение с расписанием в обработке языка.

## Направления для исследований

- Много важных задач обработки естественного языка с **большими корпусами тренировочных данных**
- Решаются с помощью **тяжелых** моделей, которые **долго** учатся
- Не исследованы метрики оценки сложности текста (длина - текущий предел)
- Эксперименты проведены только на определенных задачах
  - ACL'19 - только задача машинного перевода
  - ACL'20 - только задача классификации<sup>1</sup>
- Не исследовано влияние обучения с расписанием на этапе предобучения

---

<sup>1</sup>Не совсем честное обучение с расписанием; Не ускоряет; Требуется еще больших ресурсов



**Цель:** ускорить обучение языковой модели BERT с помощью обучения с расписанием за счет метрики оценки сложности текстовых данных на задачах предобучения, классификации и машинного перевода

**Задачи:**

- 1 Найти эффективные<sup>1</sup> метрики оценки сложности текста
- 2 Реализовать механизм подсчета найденных метрик на больших датасетах
- 3 Исследовать влияние найденных метрик на скорость обучения языковой модели BERT
- 4 Сравнить найденные метрики с существующими метриками оценки сложности текста

---

<sup>1</sup>с точки зрения сокращения скорости обучения модели