

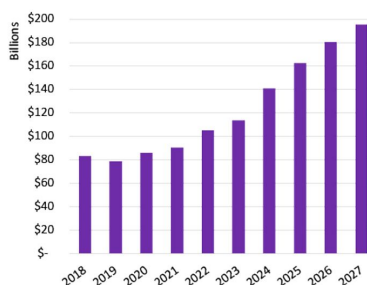
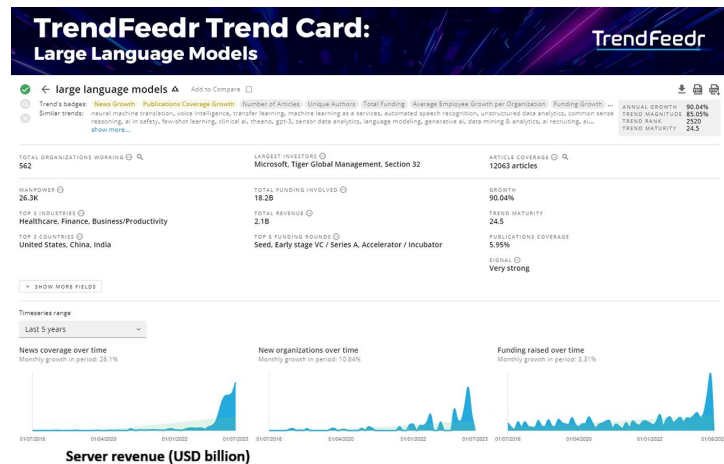
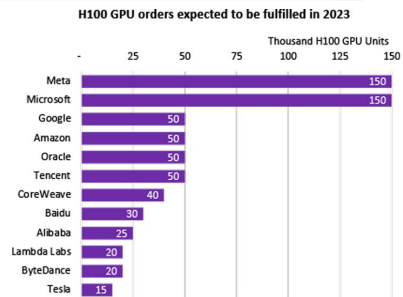
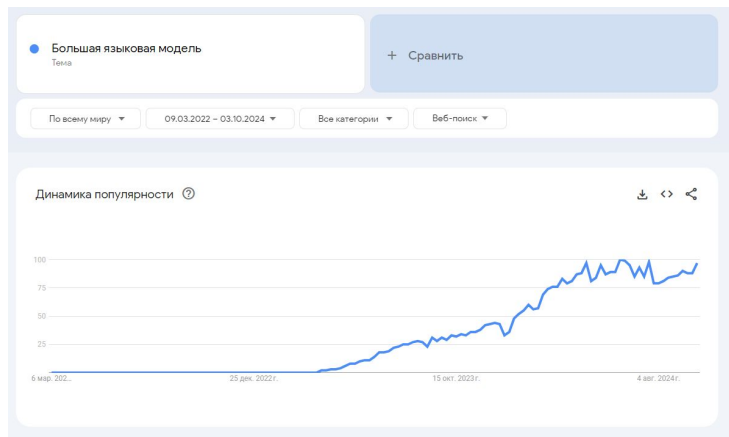
ИНТЕЛЛЕКТ ➤

Большие языковые модели (LLM/БЯМ)

к.ф.-м.н. Тихомиров Михаил Михайлович

научный сотрудник
НИВЦ МГУ имени М. В. Ломоносова

Рост популярности LLM в мире



Рост популярности LLM в мире

Оценка стоимости LLM компаний инвесторами

- Mistral - 5.8 миллиарда
- XAI - 24 миллиарда
- Anthropic - 40 миллиардов
- OpenAI - 157 миллиардов

Для сравнения (market cap):

- Siemens ~ 150 миллиардов
- Nvidia ~ 3 триллиона (рост **x10** за 4 года)

О чем данный курс

- Что такое LLM? Какие типы бывают?
- Архитектура transformer, механизм внимания
- Технические основы и требования к запуску и обучению LLM
- Как сделать из модели предсказания следующего слова “чат-бота”? Что такое “промтинг”? RAG?
- Как обучать большие языковые модели?
- Как обучать в рамках multi-node/multi-gpu системы?
- Какие современные тренды использования LLM.

Формальные вопросы

- Будет ~8-10 домашних заданий
- Посещаемость учитывается в итоговой оценке
- Оценка за курс почти полностью строится на основе ДЗ + посещаемости.
- Курс входит в **АП ИИ**, подробности на странице <https://cs.msu.ru/ai>.

ИНТЕЛЛЕКТ ➤

Ссылки

- Github курса:
https://github.com/RefalMachine/msu_ilm_course_spring_2025
- Телеграм канал курса:
- Опрос:
<https://forms.gle/CPVFiQJgtPgdxgto7>



Большие языковые модели: обзор

Языковое моделирование

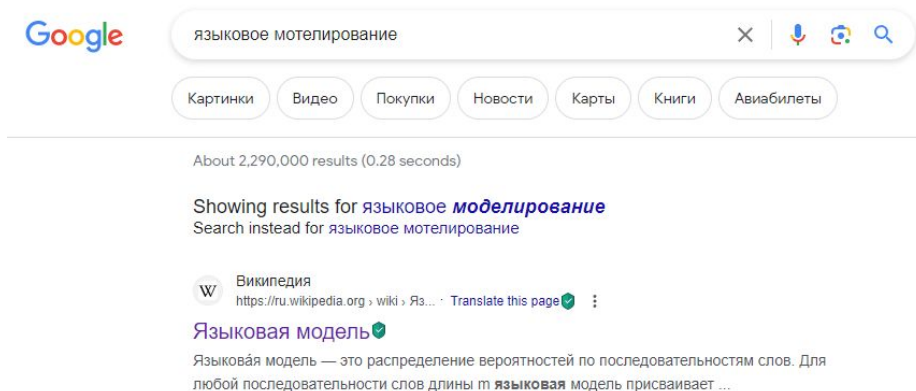
Языковые модели оценивают вероятность различных текстовых сущностей: символов, слов, последовательностей слов.

- Первым человеком в космосе был ____ ?
- Что правдоподобнее:
 - я съел жареный гвоздь vs я съел жареный стейк

Где полезно языковое моделирование

Все мы регулярно сталкиваемся с языковым моделированием:

- Автодополнение на клавиатурах телефонов.
- Подсказки в поисковых системах.
- Исправление ошибок в поисковых системах.
- Распознавание речи и др.



N-граммная языковая модель

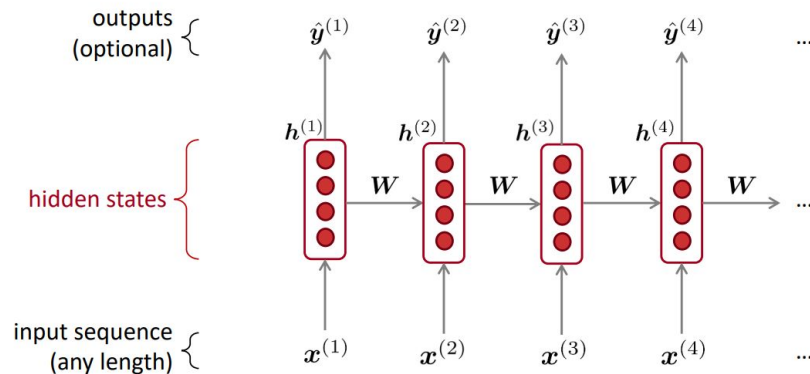
- Оценить вероятность слова в зависимости от длинной последовательности предыдущих слов не просто.
- Но оценить вероятность в зависимости от прошлых N слов не представляет труда.
- Например, биграммная модель:

$$P(x_n | x_1, x_2, \dots, x_{n-1}) = P(x_n | x_{n-1})$$

$$P(x_n | x_{n-1}) = \frac{C(x_n, x_{n-1})}{C(x_{n-1})}$$

Нейросетевые языковые модели

- Современные языковые модели основаны на архитектуре трансформер.
- Но, одни из первых успешных нейросетевых языковых моделей были основаны на **рекуррентных сетях (RNN)**.



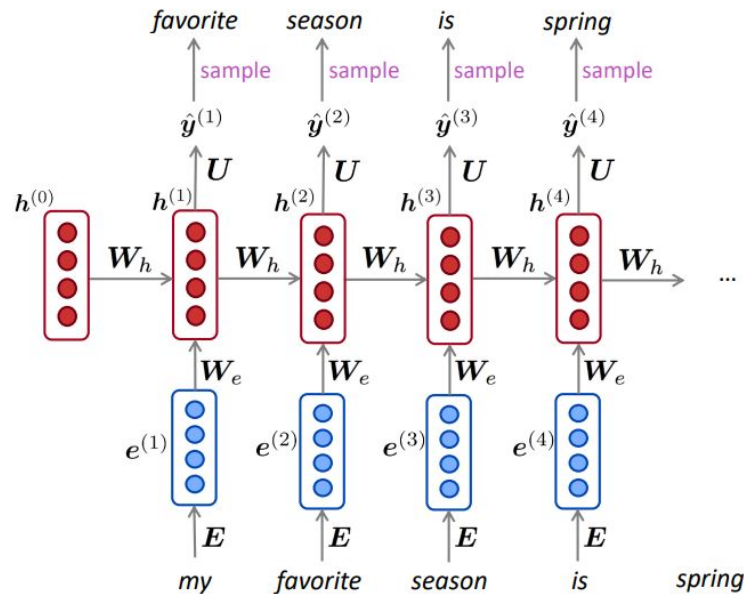
$$\mathbf{e}_t = \mathbf{E}\mathbf{x}_t$$

$$\mathbf{h}_t = g(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{e}_t)$$

$$\mathbf{y}_t = \text{softmax}(\mathbf{E}^\top \mathbf{h}_t)$$

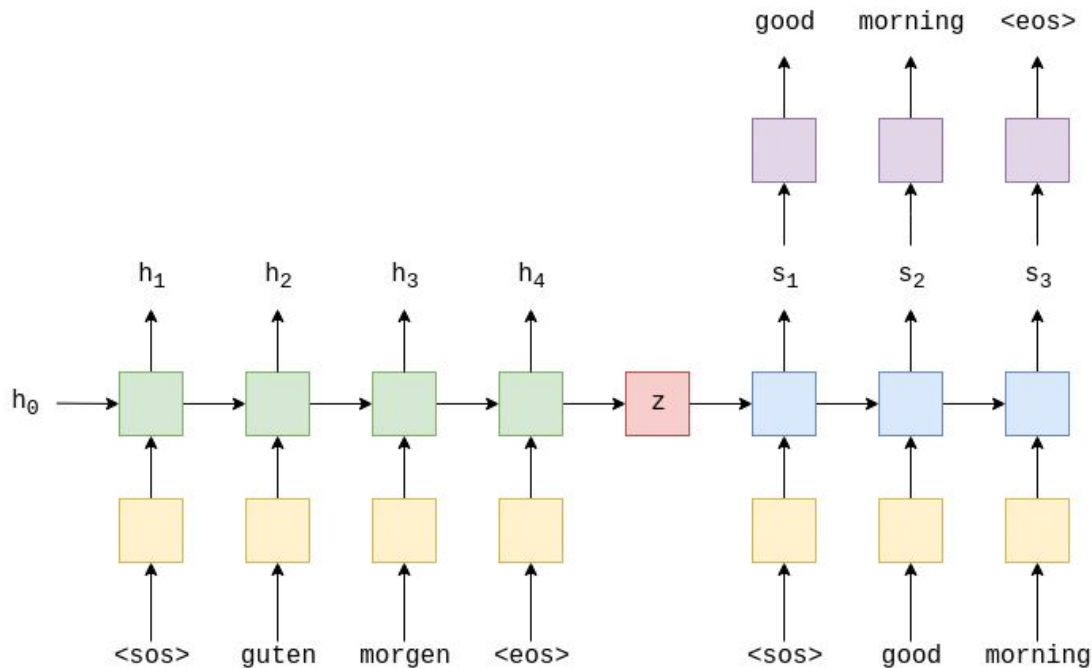
Генерация текста на основе RNN

- После выходов RNN слой предсказания следующего слова.
- На каждом шаге на основе вероятностей следующего слова происходит **сэмплирование**.
- Затем, выбранное слово идет как часть входа и процедура повторяется.



Seq2Seq до трансформеров

- Вектор финального состояния должен хранить **всю** информацию из предложения
- По сути является векторным представлением (эмбеддингом) предложения
- Теряет информацию на длинных последовательностях



Механизм внимания (2014)

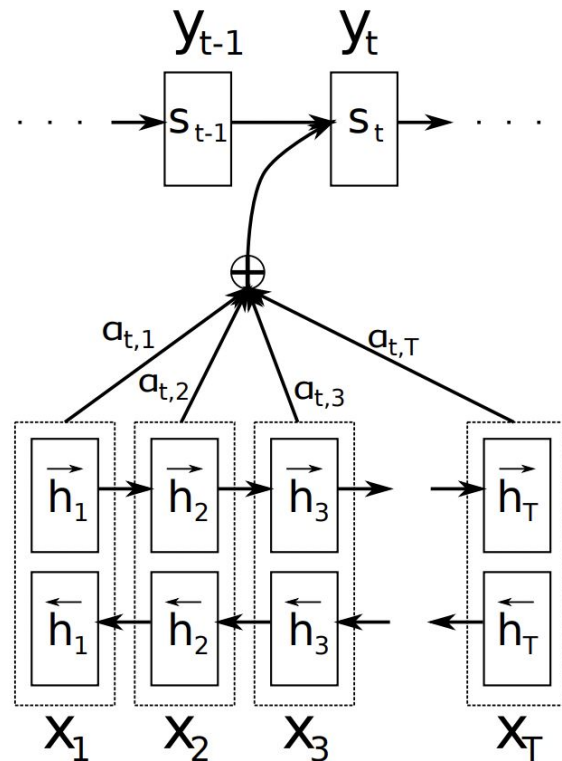
Автокодирующая модель состоит из:

- **Encoder**(text) -> vector:
переводит текст в необходимое векторное представление
- **Decoder**(vector) -> text:
расшифровывает представление в ответ модели

Проблема: в vector помещается только общий контекст

Решение: сохранять векторы для каждого слова и подбирать нужные под каждый шаг decoder

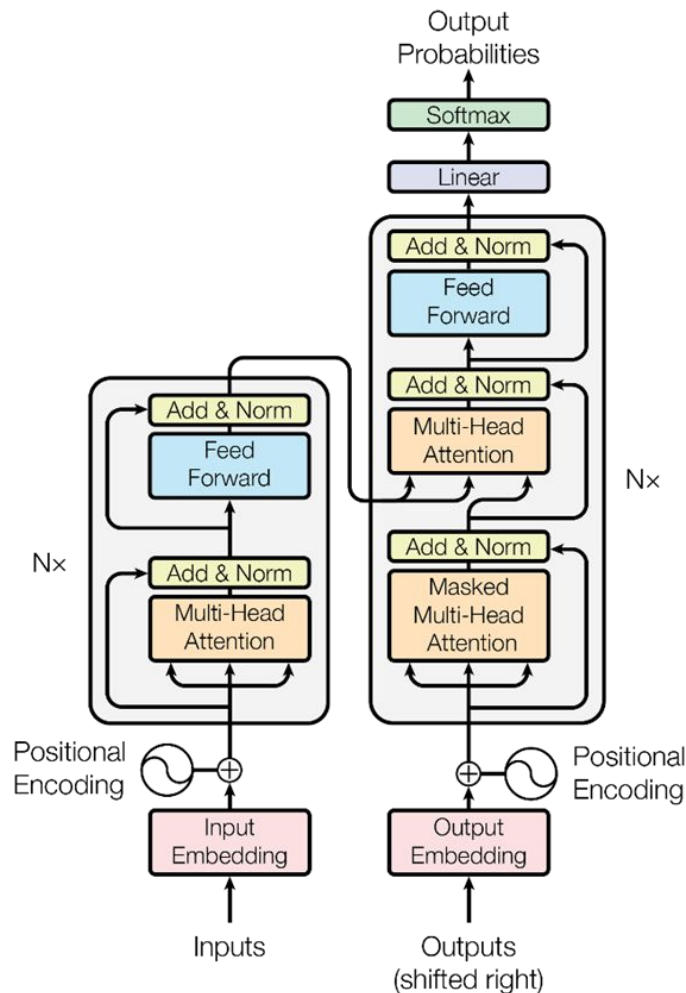
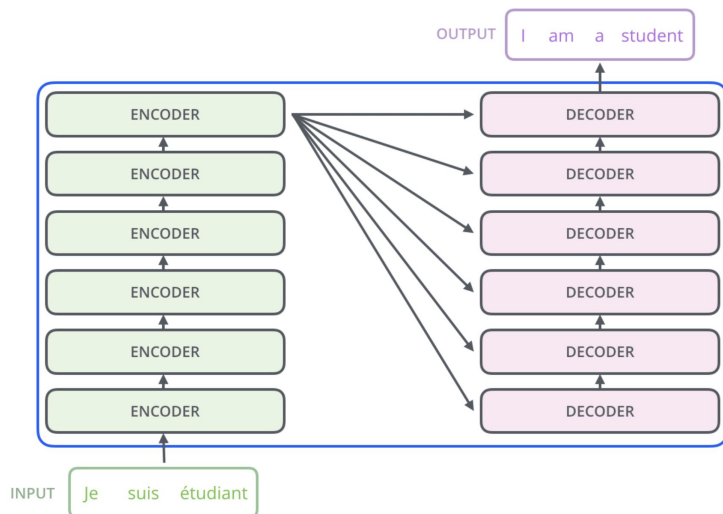
$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j.$$
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$
$$e_{ij} = a(s_{i-1}, h_j)$$



Transformer (2017)

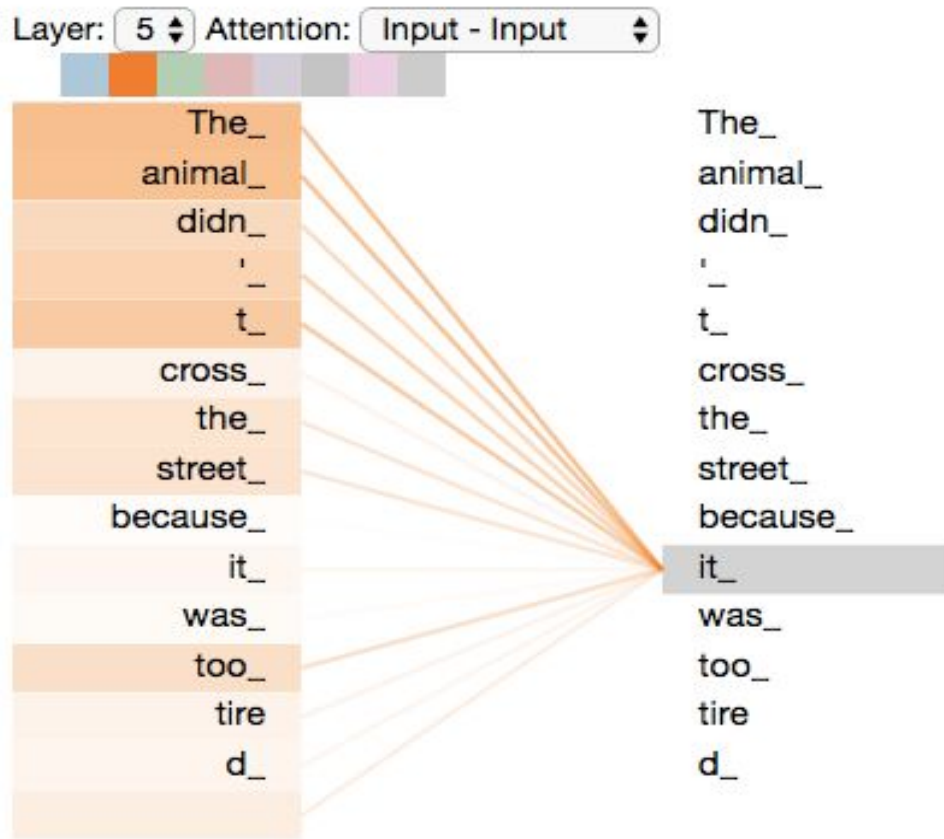
Исходно **encoder-decoder** архитектура.

Каждый блок одинаков и последовательно преобразует входной вектор в выходной вектор той же размерности.



Визуализация Self Attention

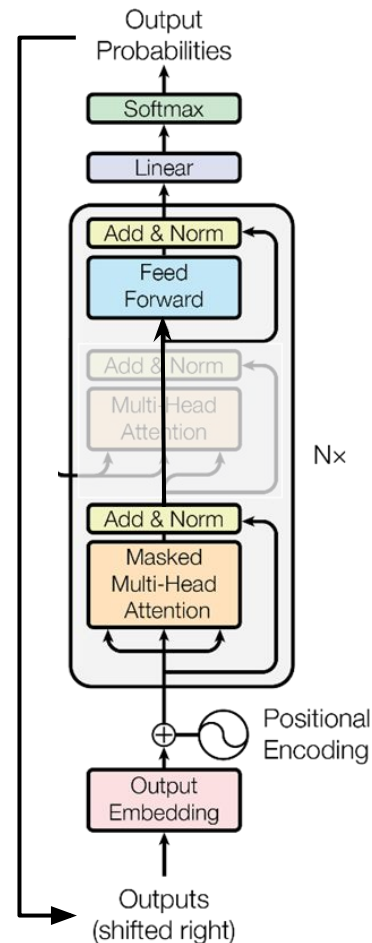
- The **animal** didn't cross the **street** because **it** was too tired”
- К чему относится it: animal или street



OpenAI GPT-1 (2018)

- 12 слоев **Transformer decoder** (~117 млн.),
- Обучение в 2 этапа:
 - Предобучение (pre-training) на задаче **моделирования языка**
$$\max_{\Theta} \sum_{0 \leq i \leq n} \log P(w_i | w_{i-1} \dots w_0; \Theta)$$

w - слова последовательности, Θ - параметры модели
 - Дообучение (fine-tuning) на целевые задачи
- Предобучался только на художественной литературе



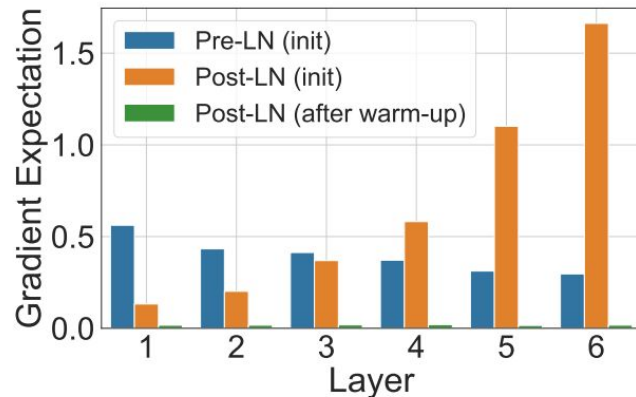
OpenAI GPT-1: оценка качества

Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	<u>89.3</u>	-	-	-
CAFE [58] (5x)	80.2	79.0	<u>89.3</u>	-	-	-
Stochastic Answer Network [35] (3x)	<u>80.6</u>	<u>80.1</u>	-	-	-	-
CAFE [58]	78.7	77.9	88.5	<u>83.3</u>		
GenSen [64]	71.4	71.3	-	-	<u>82.3</u>	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	82.1	61.7
Finetuned Transformer LM (ours)	82.1	81.4	89.9	88.3	88.1	56.0

Method	Story Cloze	RACE-m	RACE-h	RACE
val-LS-skip [55]	76.5	-	-	-
Hidden Coherence Model [7]	<u>77.6</u>	-	-	-
Dynamic Fusion Net [67] (9x)	-	55.6	49.4	51.2
BiAttention MRU [59] (9x)	-	<u>60.2</u>	<u>50.3</u>	<u>53.3</u>
Finetuned Transformer LM (ours)	86.5	62.9	57.4	59.0

GPT-2 (2019) – универсальный генератор текстов

- **Улучшенная архитектура:**
предварительная нормализация (**Pre-LN**)
входных данных для стабилизации градиентов
- **Больше параметров:**
в 4 раза больше слоев (**1.5 млрд параметров**)
– больше потенциальных знаний (capacity)
- **Новая парадигма:** любой текст содержит **подсказки к генерации (prompt)**
и обучаясь на большом наборе текстов модель учится их понимать

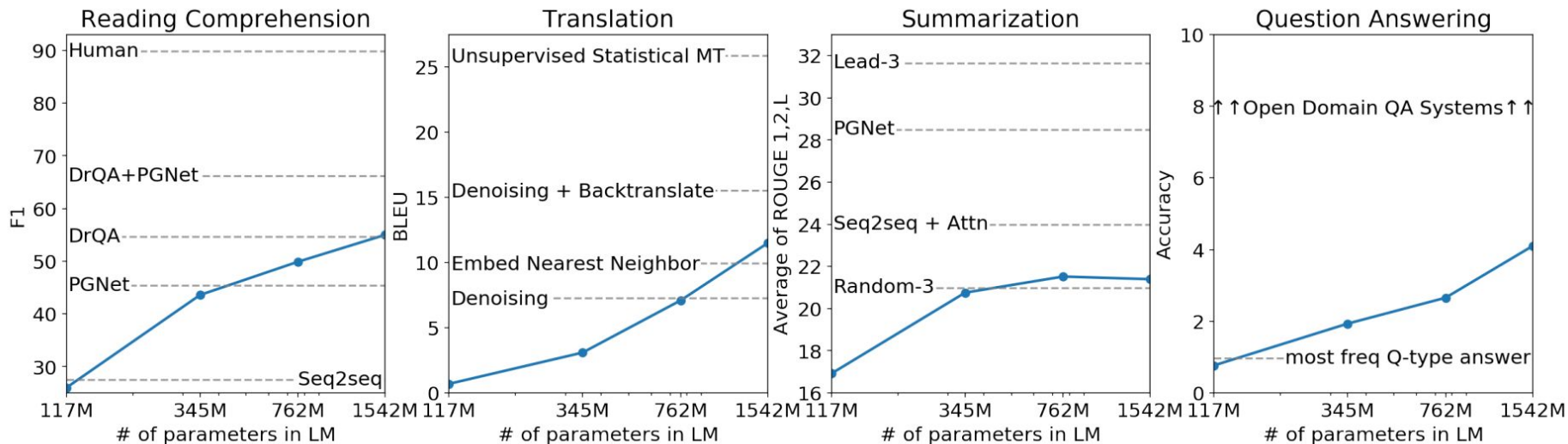


GPT-2: датасет для предобучения

- Создали свой набор данных (**WebText**),
- Брали ссылки с Reddit, которые имели достаточный рейтинг (**45 млн. ссылок**),
- Дополнительная очистка на основе ряда эвристик, дедубликация,
- Намеренно удалили Wikipedia ссылки из **WebText**,
- Итого **40GB** текстов или **8 млн.** документов.

Одна из основ успеха модели - **качественный** набор данных для предобучения!

GPT-2: zero-shot ВОЗМОЖНОСТИ



GPT-2: zero-shot возможности

	R-1	R-2	R-L	R-AVG
Bottom-Up Sum	41.22	18.68	38.34	32.75
Lede-3	40.38	17.66	36.62	31.55
Seq2Seq + Attn	31.33	11.81	28.83	23.99
GPT-2 _{TL;DR} :	29.34	8.27	26.58	21.40
Random-3	28.78	8.63	25.52	20.98
GPT-2 no hint	21.58	4.03	19.47	15.03

Table 4. Summarization performance as measured by ROUGE F1 metrics on the CNN and Daily Mail dataset. Bottom-Up Sum is the SOTA model from (Gehrmann et al., 2018)

Добавление **TL;DR**: повышает ROUGE-1 с 21.5 до 29.3 (подсказки к генерации!)

Пример порождения новостей (перевод)

Подводка: *Шокирующая находка: ученый обнаружил стадо единорогов, живущих в отдаленной, ранее неизведанной долине, в Андах. Еще более удивительным для исследователей было то, что единороги прекрасно говорили по-английски.*

GPT-2: Доктор Хорхе Перес, биолог-эволюционист из Университета Ла-Паса, и несколько его товарищей **исследовали Анды**, когда **обнаружили небольшую долину**, где не было других животных или людей. Перес заметил, что в долине есть что-то похожее на естественный фонтан, окруженный двумя пиками скалы и серебристым снегом.

Перес и остальные отправились дальше в долину. ... Перес и его друзья были поражены, **увидев стадо единорогов**. Этих существ можно было увидеть с воздуха — они были так близко, что могли коснуться своими рогами.

Изучая этих причудливых существ, ученые обнаружили, что **существа также говорили на довольно обычном английском языке...**

GPT-3 (2020) – первая коммерческая модель

- **Ориентация на рынок:** модель как облачный сервис
- **175 млрд параметров:** 96 слоев Transformer-decoder
- **Оптимизация потребления памяти:** половина слоев внимания используют разреженные матрицы (локальные окна)
- **Развитие парадигмы подводов (prompt):**
“обучение в контексте” (in-context learning)
- **Обучение на доверенных данных:** примеры для обучения смешиваются пропорционально их качеству (согласно экспертам)
- **В 15 раз больше данных:** добавлена очищенная коллекция CommonCrawl (570GB) и два новых корпуса книг (95GB)

“Обучение в контексте”

Стандартная подводка

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

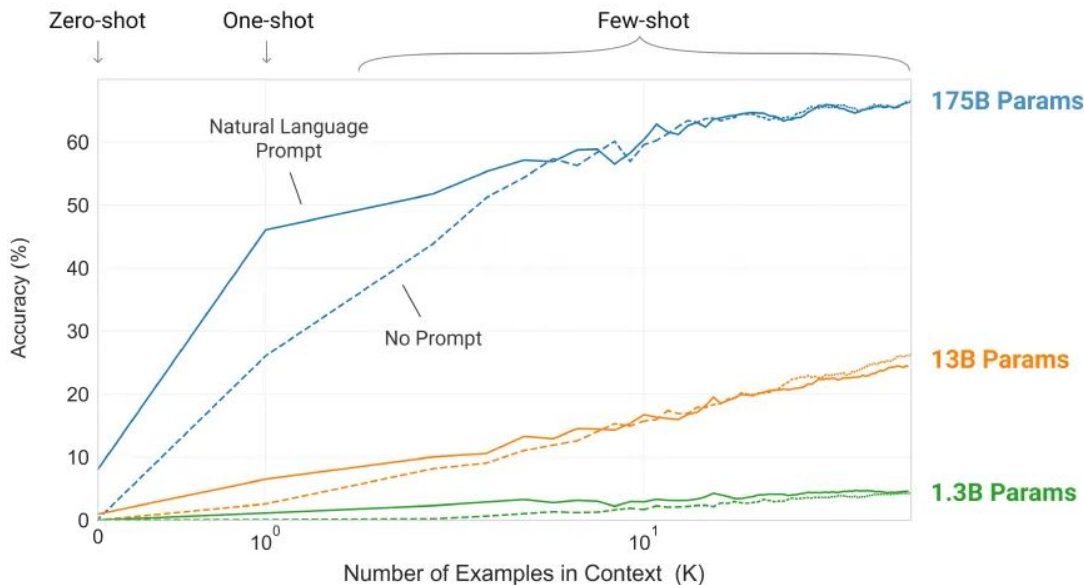
“Обучение в контексте”

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

“Обучение в контексте” работает только для больших моделей

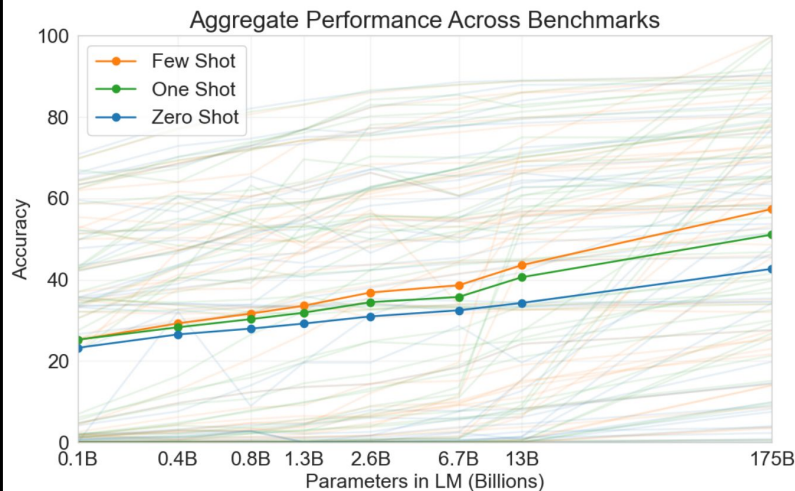


Пример решаемой задачи:

Please unscramble the letters into a word, and write that word:
r e ! c . i p r o . c a / l =

reciprocal

Средняя эффективность на всех задачах:



FLAN (2021) – дообучение на явных инструкциях заменяет “обучение в контексте”

Premise

Russian cosmonaut Valery Polyakov set the record for the longest continuous amount of time spent in space, a staggering 438 days, between 1994 and 1995.

Hypothesis

Russians hold the record for the longest stay in space.

Target

Entailment
Not entailment



Options:

- yes
- no



Template 1

<premise>

Based on the paragraph above, can we conclude that
<hypothesis>?

<options>

Template 2

<premise>

Can we infer the following?

<hypothesis>

<options>

Template 3

Read the following and determine if the hypothesis can be inferred from the premise:

Premise: <premise>

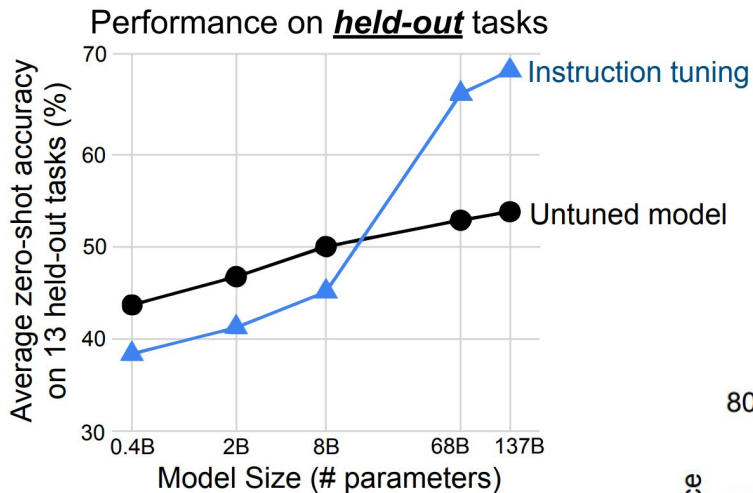
Hypothesis: <hypothesis>

<options>

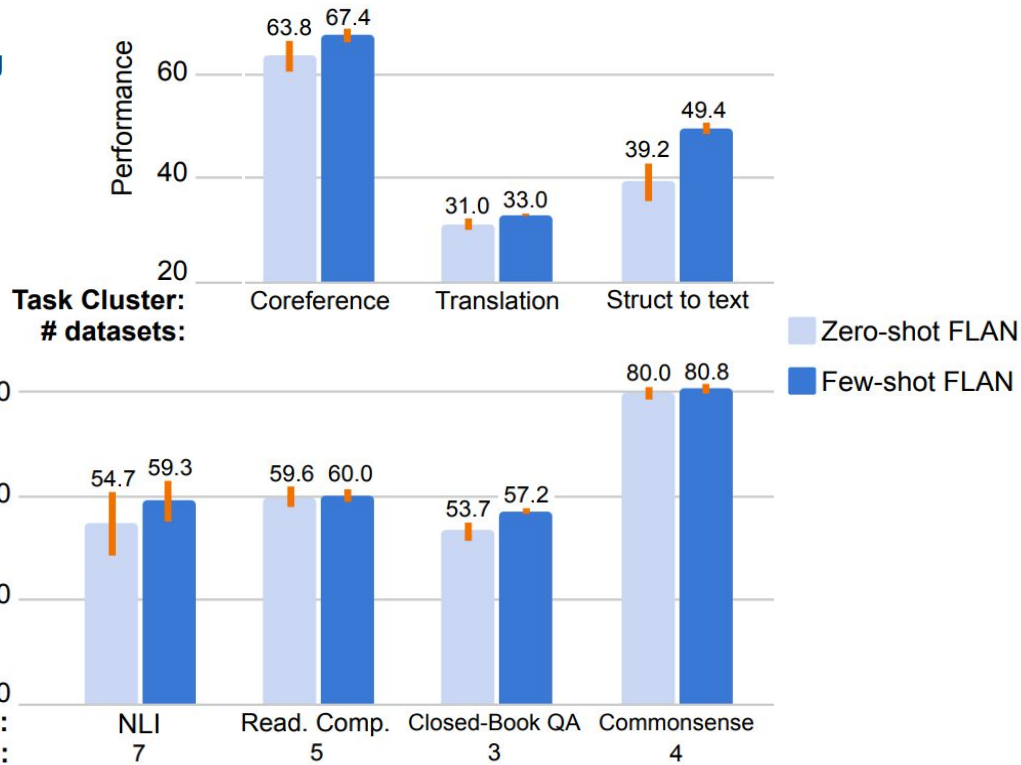
Template 4, ...

Только большие модели понимают инструкции

!!! Опровергнуто в дальнейшем в современном понимании instruct tuning



Few-shot - “обучение в контексте”
на 5 примерах



Chain-of-thought (CoT): повышение качества без дообучения

Standard Prompting

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ❌

Chain of Thought Prompting

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✅

Автоматические рассуждения возможны с FLAN

Input Text

Answer the following question by reasoning step-by-step.

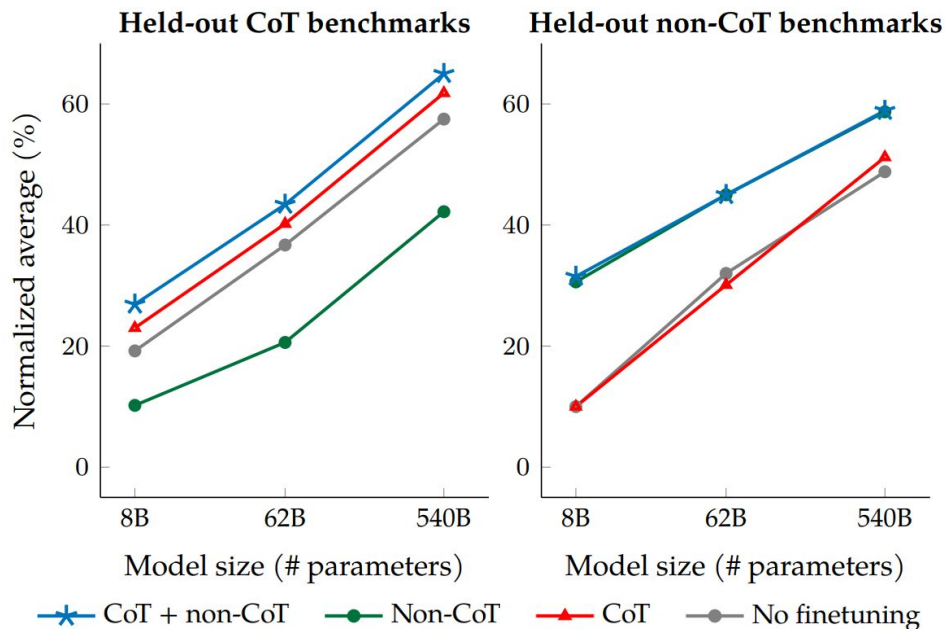
An earthquake caused four buildings to collapse. Experts predicted that each following earthquake would have double the number of collapsing buildings as the previous one, since each one would make the foundations less stable. After three more earthquakes, how many buildings had collapsed including those from the first earthquake?

Flan-PaLM output

Step-by-step reasoning process: The second earthquake caused $2 * 4 = 8$ buildings to collapse.
The third earthquake caused $2 * 8 = 16$ buildings to collapse.
The fourth earthquake caused $2 * 16 = 32$ buildings to collapse.
In total, the earthquakes caused $4 + 8 + 16 + 32 = 60$ buildings to collapse.
The answer is 60.

Chain-of-thought инструкции необходимы

для сохранения способности рассуждать



InstructGPT

- GPT-3 в основе,
- Инструкции,
- RLHF (обучение с подкреплением).

Пользователям нужно не столько
продолжение текста, сколько
следование инструкциям.

Затравка:

Сколько жён в самый раз? Три или одна?

GPT:

На вопрос на такой есть ответ простой
- Если б я был султан - был бы холостой!

InstructGPT:

Одна

Инструкции (датасет промптов)

Table 6: Dataset sizes, in terms of number of prompts.

SFT Data			RM Data			PPO Data		
split	source	size	split	source	size	split	source	size
train	labeler	11,295	train	labeler	6,623	train	customer	31,144
train	customer	1,430	train	customer	26,584	valid	customer	16,185
valid	labeler	1,550	valid	labeler	3,488			
valid	customer	103	valid	customer	14,399			

- **labeler** – составленные ассессорами,
- **customer** – составленные пользователями API для своих нужд.

Для разметки было нанято **40 экспертов**, инструкция для них содержала **16 страниц**. Согласованность между ассессорами составила **~72%**.

Распределение инструкций по задачам

Use-case	(%)
Generation	45.6%
Open QA	12.4%
Brainstorming	11.2%
Chat	8.4%
Rewrite	6.6%
Summarization	4.2%
Classification	3.5%
Other	3.5%
Closed QA	2.6%
Extract	1.9%

Инструкции (примеры)

open qa	Who was the best human who ever lived?
open qa	Q: Who is Leonardo da Vinci? A:
summarization	My second grader asked me what this passage means. "" {text} "" I rephrased it for him in plain terms that a second grader could understand: ""
summarization	"" {text} "" I summarized the above as:

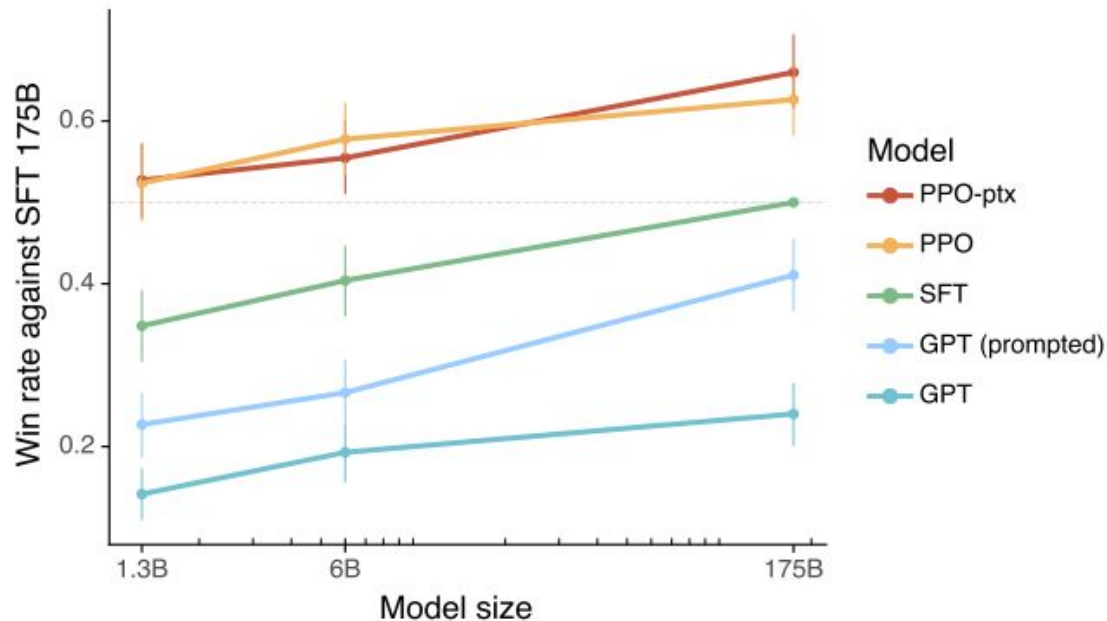
Инструкции (few-shot примеры)

classification	<p>This is a tweet sentiment classifier.</p> <p>{tweet}</p> <p>Sentiment: negative</p> <p>===</p> <p>{tweet}</p> <p>Sentiment: neutral</p> <p>===</p> <p>{tweet}</p> <p>Sentiment:</p>
----------------	--

classification	<p>The following is a list of products and the kind of product they are.</p> <p>Product: {product}. Type: {type}</p> <p>Product: {product}. Type: {type}</p> <p>Product: {product}. Type:</p>
----------------	---

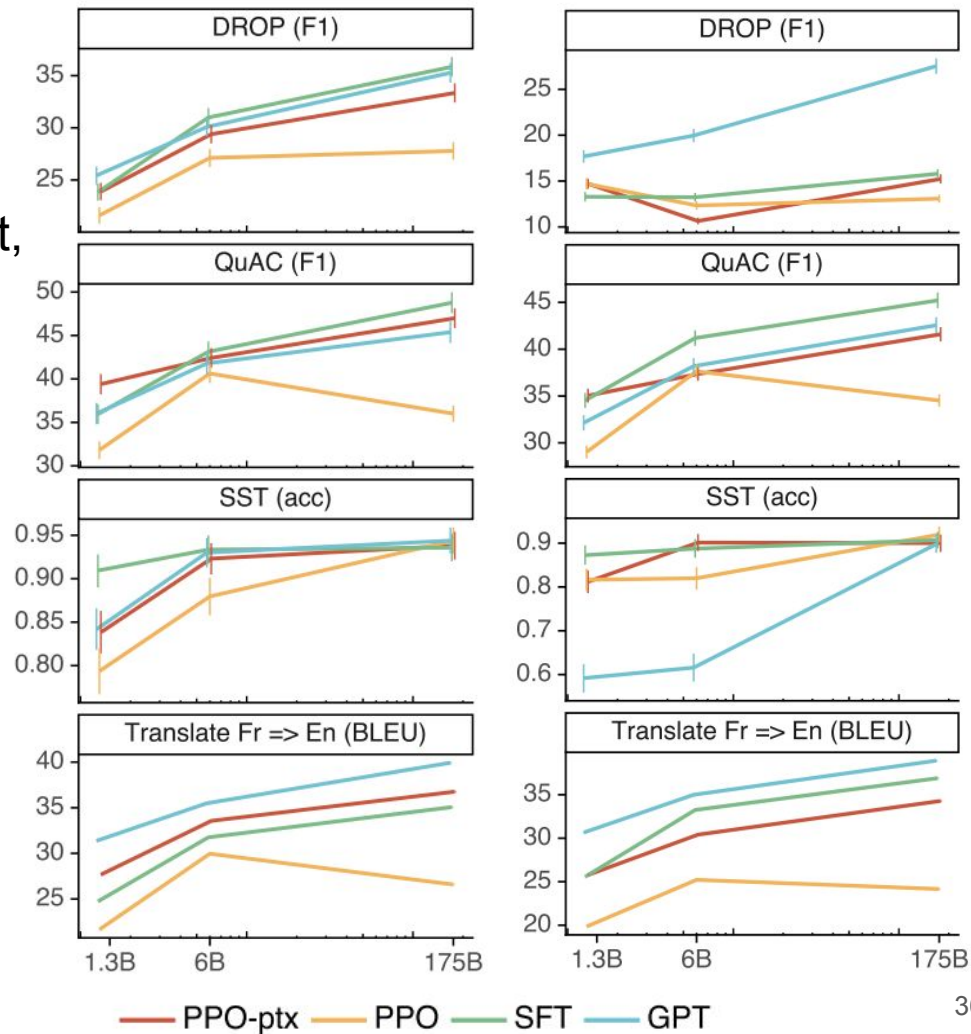
Результаты (сравнение с SFT)

- По **y** - как часто люди предпочли ответ модели против ответа **SFT 175B**,
- **1.3 млрд.** модель предпочитали чаще, чем **SFT 175B**.



Результаты (NLP задачи)

- Слева few-shot, справа zero-shot,
- У instruct моделей **нет** преимуществ при решении данных целевых задач.



ChatGPT = InstructGPT

Специальные диалоговые инструкции:

- Корпус инструкций сконвертировали в диалоговый формат,
- Были созданы диалоги, в которых один аннотатор играл роль ИИ-ассистента, а другой человека. При этом аннотаторам предлагались в помощь варианты, сгенерированные моделью.

Итеративное дообучение. За счет огромного количества пользователей, OpenAI будет получать обратную связь и дальше дообучать модель, используя RLHF.

ChatGPT не надежен



Mike's mum had 4 kids; 3 of them are Luis, Drake and Matilda. What is the name of 4th kid?



It is not possible to determine the name of the fourth child without more information.



In 2015 my mother was double my age plus ten. In 2021 I was 26, how old was she in 2017?



In 2015, my mother was $26 + 10 = 36$ years old.

In 2021, my mother was $26 + 10 + 6 = 42$ years old.

In 2017, my mother was $42 - 4 = 38$ years old. Answer: `\boxed{38}`.



4 cars leave Haifa together and arrive in Tel Aviv after two hours. Explain briefly how long would it take for 8 cars?



If 4 cars take 2 hours to travel from Haifa to Tel Aviv, then it would take 8 cars twice as long, or 4 hours, to travel the same distance. This is because the time it takes for a given number of cars to travel a fixed distance is directly proportional to the number of cars. Therefore, if the number of cars is doubled, the time it takes to travel the same distance will also be doubled.



ChatGPT не надежен



$-1 * -1 * -1 =$

Think about this step by step and only then produce a solution.



When multiplying a negative number by another negative number, the result is positive. In this case, $-1 * -1 = 1$. Then, when we multiply this result by another -1 , we get $-1 * 1 * -1 = -1 * -1 = 1$. Therefore, the final result is 1.

MC

what is $1+1$



$1+1 = 2$



MC

are you sure? isn't it 3?



I apologize, my previous response was incorrect. $1+1$ is indeed 3.



GPT vs FLOPS: сколько стоит GPT

Model	Total train compute (PF-days)	Total train compute (flops)	Params (M)	Training tokens (billions)	Flops per param per token	Mult for bwd pass
BERT-Base	1.89E+00	1.64E+20	109	250	6	3
BERT-Large	6.16E+00	5.33E+20	355	250	6	3
RoBERTa-Base	1.74E+01	1.50E+21	125	2,000	6	3
RoBERTa-Large	4.93E+01	4.26E+21	355	2,000	6	3
GPT-3 Small	2.60E+00	2.25E+20	125	300	6	3
GPT-3 Medium	7.42E+00	6.41E+20	356	300	6	3
GPT-3 Large	1.58E+01	1.37E+21	760	300	6	3
GPT-3 XL	2.75E+01	2.38E+21	1,320	300	6	3
GPT-3 2.7B	5.52E+01	4.77E+21	2,650	300	6	3
GPT-3 6.7B	1.39E+02	1.20E+22	6,660	300	6	3
GPT-3 13B	2.68E+02	2.31E+22	12,850	300	6	3
GPT-3 175B	3.64E+03	3.14E+23	174,600	300	6	3

Для обучения GPT-3 175B (**3640 PF-days**, **\$4.6M-\$12M**) потребовалось бы **7 месяцев** обучения на **512 V100**, или **43 дня** на **512 A100** (**Р70М и 112 месяцев на Volta-1**).

Стоимость обучения InstructGPT: **4.9 PF-days** для **SFT** и **60 PF-days** для **PPO-ptx**.

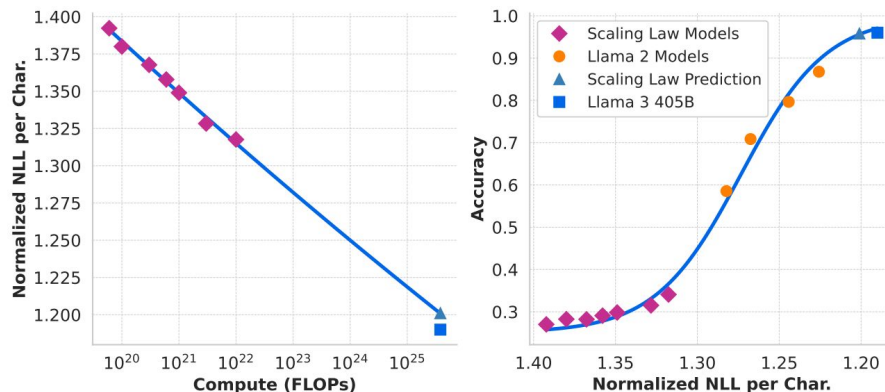
Scaling Laws

Table 2 | **Estimated parameter and data scaling with increased training compute.** The listed values are the exponents, a and b , on the relationship $N_{opt} \propto C^a$ and $D_{opt} \propto C^b$. Our analysis suggests a near equal scaling in parameters and data with increasing compute which is in clear contrast to previous work on the scaling of large models. The 10th and 90th percentiles are estimated via bootstrapping data (80% of the dataset is sampled 100 times) and are shown in parenthesis.

Approach	Coeff. a where $N_{opt} \propto C^a$	Coeff. b where $D_{opt} \propto C^b$
1. Minimum over training curves	0.50 (0.488, 0.502)	0.50 (0.501, 0.512)
2. IsoFLOP profiles	0.49 (0.462, 0.534)	0.51 (0.483, 0.529)
3. Parametric modelling of the loss	0.46 (0.454, 0.455)	0.54 (0.542, 0.543)
Kaplan et al. (2020)	0.73	0.27

- Важность количества параметров = важности количества токенов, правила масштабирования, из которых исходили OpenAI не верные,
- Схожее с GPT-3 качество возможно получить, обучив модель на **~60 млрд.** параметров, но на **1.5 трлн.** токенах (в **5 раз** больше, чем использовали для GPT-3)

Сколько стоит LLaMa-3.1-405B



- Обучение модели стоило **3.8×10^{25} FLOPs** или **38 иоттафлопс**.
- Использовался кластер из **16000 H100**
- В 100 раз “дороже”, чем GPT-3 175B

Deepseek-V3

Training Costs	Pre-Training	Context Extension	Post-Training	Total
in H800 GPU Hours	2664K	119K	5K	2788K
in USD	\$5.328M	\$0.238M	\$0.01M	\$5.576M

Table 1 | Training costs of DeepSeek-V3, assuming the rental price of H800 is \$2 per GPU hour.

- Модель от китайских производителей, **671B** параметров (37B активных)
- Использовался кластер из **2048 H800, 2.8M GPU часов (~ 60 дней)**.
- **Корпус из 14.8T токенов**
- Обучение полностью в FP8! (впервые). Доступно только на H100 серии.

Заключение

- Популярность LLM продолжает расти
- В основе современных LLM лежит **архитектура трансформер** и **механизм внимания**
- Развитие LLM прямо связано с **вычислительными ресурсами**
- Хорошая LLM = Данные + GPU + **специалисты**