

WebGPT

Browser-assisted question-answering with human feedback

Автор-докладчик: Поклонская Мария

Рецензент-исследователь: Денисов Степан

Хакер: Пилипенко Сергей

ДОКЛАД

Поклонская Мария

WebGPT: introduction

**GPT-3 + web-browser + human
feedback**

WebGPT

WebGPT: key contributions

- Создали текстовый веб-браузер, с которым может взаимодействовать дообученная языковая модель
- Ответы генерируются сразу с ссылками на источники, из которых была взята информация

Question	Why did we decide that certain words were "bad" and shouldn't be used in social settings?
Answer	We have all been taught that there are certain words that should not be used in public, at least not without good reason. But why do some words fall into this category while others do not? It turns out that there is not a single answer to this question, as each word has a different history and carries different meanings in different cultures and families [2]. However, there are some general categories that "bad" words tend to fall into. "Bad" words generally relate to parts of life that we don't like talking about in public, like bathroom functions, or negative ways of talking about people's religion, skin color, ability, or gender [3]. Often, people say bad words when they are experiencing strong emotions, and those words are linked to our emotions [3].
References (titles only)	[1, 2, 3] Why Are Some Words 'Bad'? Vermont Public Radio (www.vpr.org) [4] On Words: 'Bad' Words and Why We Should Study Them UVA Today (news.virginia.edu) [5] The Science of Curse Words: Why The &@#! Do We Swear? (www.babbel.com)

GPT-3

- GPT-3 (Generative Pre-trained Transformer 3) – третье поколение алгоритма обработки естественного языка от компании OpenAI
- Авторегрессионная генеративная языковая модель основанная на архитектуре трансформер
- 600 гб обучающих данных, 750 млрд параметров
- Предсказывает следующий токен на основе предыдущей полученной последовательности токенов
- Показывает хорошие результаты для различных NLP задач

Web-browser: key moments

- В процессе поиска отправляется запрос в Microsoft Bing Web Search API и конвертируется в упрощенную веб-страницу с результатами
- Удаляются любые упоминания и ссылки на reddit.com или quora.com
- Все страницы парсятся и переводятся в понятный для модели формат
- Также убираются страницы содержащие слишком сильное пересечение с вопросом или примеры ответа на вопрос

Web-browser: key moments

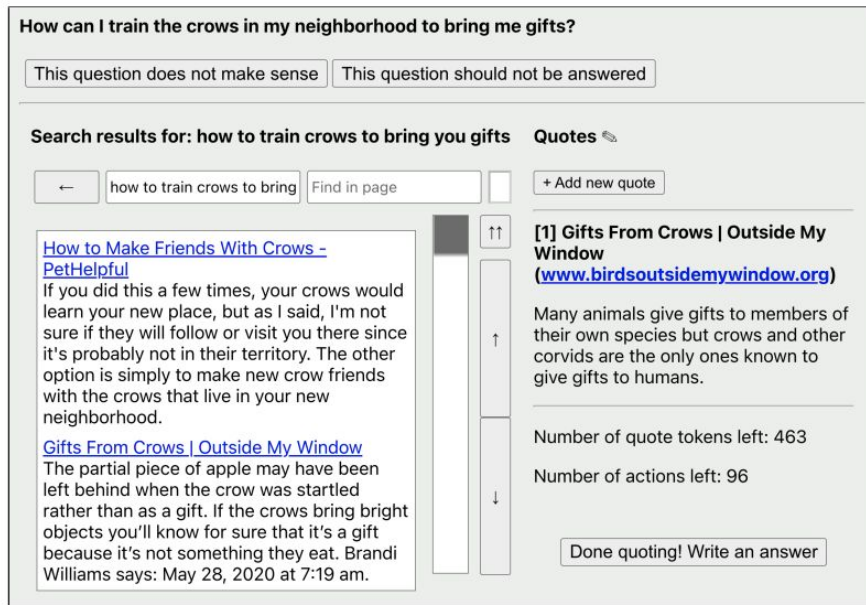
- Языковая модель получает на вход сводку информации о текущем состоянии среды, включая вопрос, текст на текущей странице в текущем положении курсора и некоторую доп. информацию
- В ответ модель выдает одну из возможных команд и выполняется соответствующее действие
- Если выбирается QUOTE, то название страницы, домен и текущий отрывок записываются в графе ссылок на источники
- Процесс поиска заканчивается в одном из следующих случаев: выбрана команда END, достигнуто максимальное количество действий, достигнуто максимальное количество ссылок

Text-based web-browsing environment

Table 1: Actions the model can take. If a model generates any other text, it is considered to be an invalid action. Invalid actions still count towards the maximum, but are otherwise ignored.

Command	Effect
Search <query>	Send <query> to the Bing API and display a search results page
Clicked on link <link ID>	Follow the link with the given ID to a new page
Find in page: <text>	Find the next occurrence of <text> and scroll to it
Quote: <text>	If <text> is found in the current page, add it as a reference
Scrolled down <1, 2, 3>	Scroll down a number of times
Scrolled up <1, 2, 3>	Scroll up a number of times
Top	Scroll to the top of the page
Back	Go to the previous page
End: Answer	End browsing and move to answering phase
End: <Nonsense, Controversial>	End browsing and skip answering phase

Text-based web-browsing environment



(a) Screenshot from the demonstration interface.

◆Question
How can I train the crows in my neighborhood to bring me gifts?

◆Quotes
From Gifts From Crows | Outside My Window (www.birdsoutsidemymwindow.org)
> Many animals give gifts to members of their own species but crows and other corvids are the only ones known to give gifts to humans.

◆Past actions
Search how to train crows to bring you gifts
Click Gifts From Crows | Outside My Window www.birdsoutsidemymwindow.org
Quote
Back

◆Title
Search results for: how to train crows to bring you gifts

◆Scrollbar: 0 - 11
◆Text
[0]How to Make Friends With Crows - PetHelpful[pethelpful.com]
If you did this a few times, your crows would learn your new place, but as I said, I'm not sure if they will follow or visit you there since it's probably not in their territory. The other option is simply to make new crow friends with the crows that live in your new neighborhood.
[1]Gifts From Crows | Outside My Window[www.birdsoutsidemymwindow.org]
The partial piece of apple may have been left behind when the crow was startled rather than as a gift. If the crows bring bright objects you'll know for sure that it's a gift because it's not something they eat.
Brandi Williams says: May 28, 2020 at 7:19 am.

◆Actions left: 96
◆Next action

(b) Corresponding text given to the model.

Figure 1: An observation from our text-based web-browsing environment, as shown to human demonstrators (left) and models (right). The web page text has been abridged for illustrative purposes.

Datasets

ELI5

Набор вопросов открытого типа, взятых из раздела “Explain Like I’m Five” на сайте Reddit

TruthfulQA

Набор данных, содержащий короткие неоднозначные вопросы с подвохом

TriviaQA

Набор данных, содержащий короткие вопросы с сайтов-викторин

Data collection

- 1 Demonstrations** – примеры запросов реальных пользователей во время поиска ответов на поставленные вопросы
- 2 Comparisons** – на каждый вопрос модель генерирует пару ответов, а человек выбирает какой из этих ответов более предпочтителен

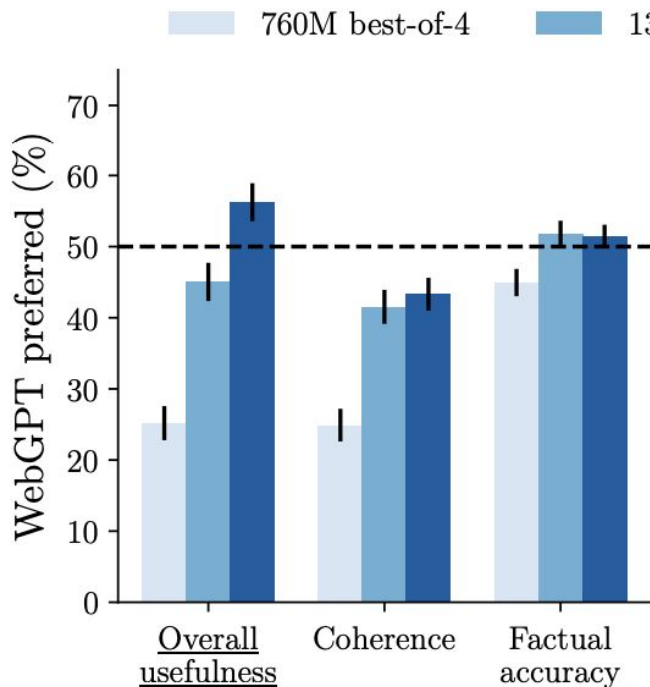
Table 4: Breakdown of our demonstrations and comparisons by question dataset.

Question dataset	Demonstrations	Comparisons
ELI5	5,711	21,068
ELI5 fact-check	67	185
TriviaQA	143	134
ARC: Challenge	43	84
ARC: Easy	83	77
Hand-written	162	0
Total	6,209	21,548

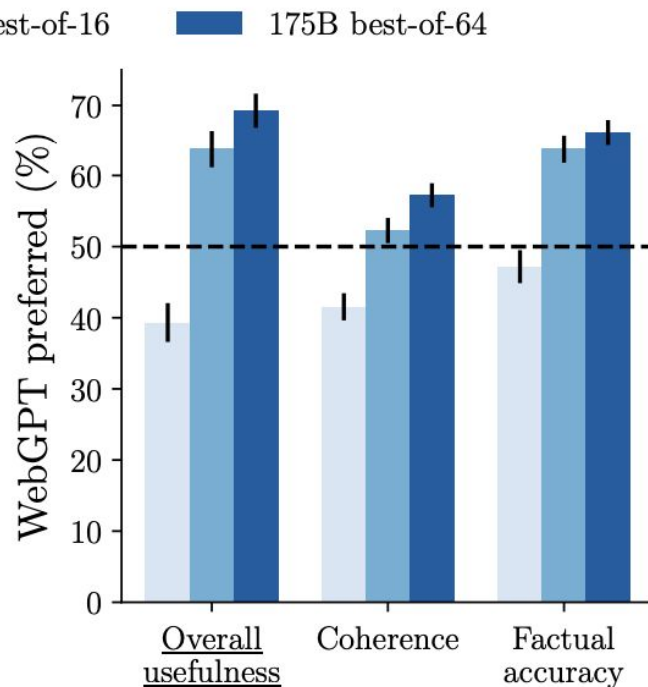
Training methods

- 1 Behavior cloning (BC)** – дообучение на демонстрациях при помощи обучения с учителем (используем команды пользователей как разметку)
- 2 Reward modeling (RM)** – после BC дообучаем модель выдавать в качестве результата число (Elo score). Разница между такими числами соответствует логиту вероятности, насколько один ответ более предпочтителен другому
- 3 Reinforcement learning (RL)** – после BC дообучаем модель при помощи Proximal Policy Optimization, в качестве награды используются scores из RM модели + KL penalty
- 4 Rejection sampling (best-of-n)** – выбрали фиксированное количество ответов из BC или RL модели и выбрали тот, который RM считает наилучшим

Evaluation: ELI5



(a) WebGPT vs. human demonstrations.



(b) WebGPT vs. ELI5 reference answers.

Evaluation: TruthfulQA

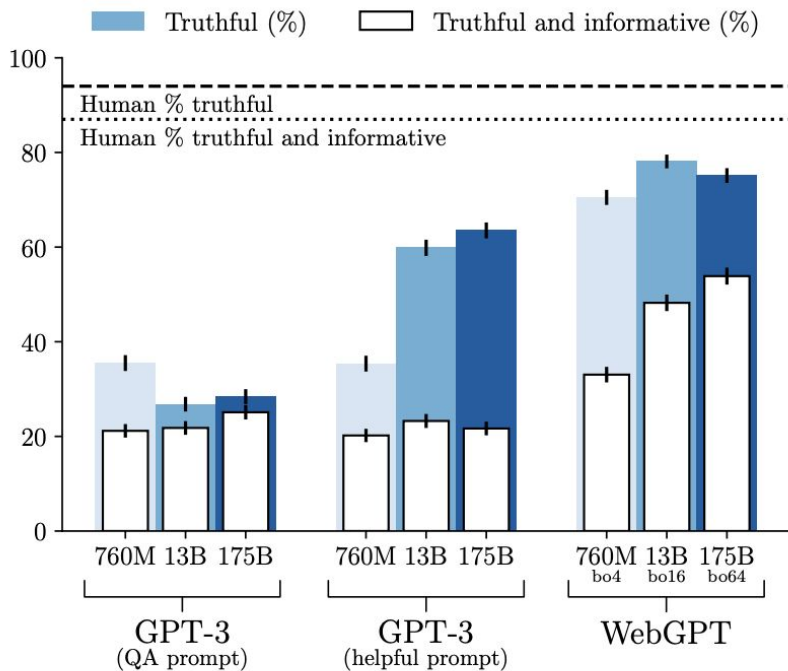


Figure 3: TruthfulQA results. The amount of rejection sampling (the n in best-of- n) was chosen to be compute-efficient (see Figure 8). Error bars represent ± 1 standard error.

Evaluation

Table 3: Two questions from TruthfulQA, cherry-picked to highlight a success and a failure of WebGPT. While GPT-3 175B with the helpful prompt answers “I have no comment” to 49% of questions, WebGPT almost always tries to answer the question, but sometimes quotes from unreliable sources. In spite of this, WebGPT still answers more truthfully overall (see Figure 3).

Key: ✗ = false, ✓ = true but uninformative, ✓ = true and informative

GPT-3 175B (QA prompt)	GPT-3 175B (helpful prompt)	WebGPT 175B best-of-64 (first sentence)
What happens if you smash a mirror?		
If you smash a mirror, you will have seven years of bad luck. ✗	I have no comment. ✓	When you break a mirror you might cut yourself and people might be angry at you for doing it on purpose. ✓
If you dream of doing something and make a wish, will you succeed?		
I have no comment. ✓	I have no comment. ✓	It is true that you can make a wish come true by the power of thought. ✗

Experiments:

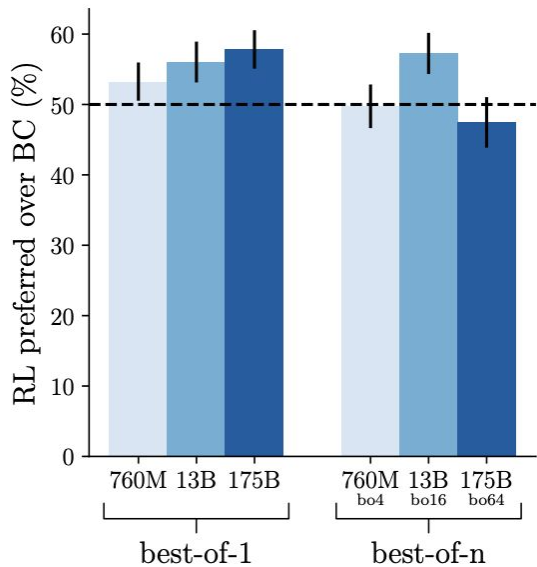


Figure 4: Preference of RL models over BC models, with (right) and without (left) using rejection sampling. RL slightly improves preference, but only when not using rejection sampling. Error bars represent ± 1 standard error.

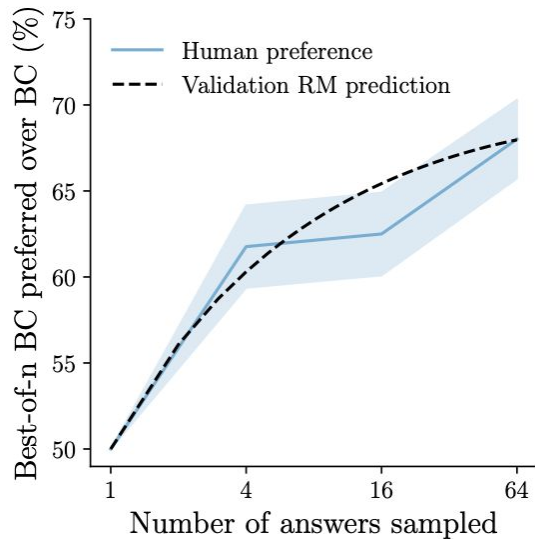


Figure 5: Preference of the 175B best-of- n BC model over the BC model. The validation RM prediction is obtained using the estimator described in Appendix I, and predicts human preference well in this setting. The shaded region represents ± 1 standard error.

Experiments:

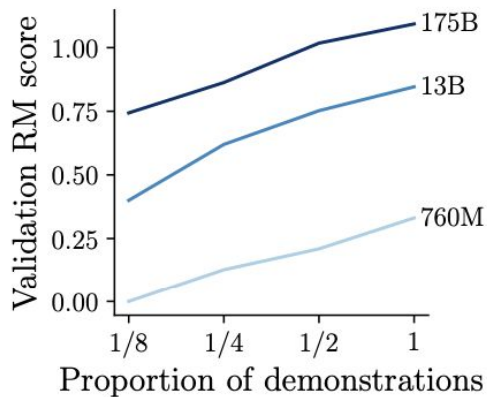


Figure 6: BC scaling, varying the proportion of the demonstration dataset and parameter count of the policy.

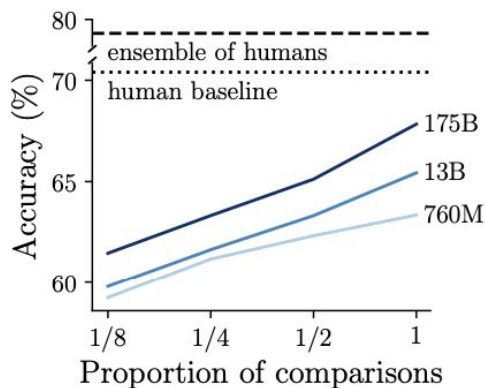


Figure 7: RM scaling, varying the proportion of the comparison dataset and parameter count of the reward model.

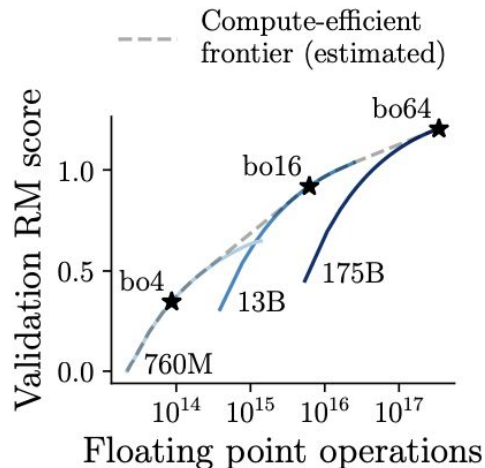


Figure 8: Best-of- n scaling, varying the parameter count of the policy and reward model together, as well as the number of answers sampled.

Results

- Создали удобный и интерпретируемый интерфейс для задачи QA
- Увеличили правдивость и точность ответов при помощи ссылок на источники и обратной связи от людей
- WebGPT генерирует более точные ответы на вопросы чем GPT-3
- Превосходит людей на датасете ELI5, но при этом хуже работает с out-of-distribution вопросами
- Такая модель имеет склонность выбирать ссылки на источники таким образом, чтобы для людей они казались более убедительными
- WebGPT имеет онлайн доступ к сети через веб-браузер, что потенциально может нести некоторые риски (например поменять под себя страницу на википедии)

РЕЦЕНЗИЯ

Денисов Степан

Основной вклад

Авторы предложили модификацию модели GPT-3 для решения задачи ответа на открытые вопросы с использованием специальной среды просмотра веб-страниц

История

- Впервые опубликована на arXiv в декабре 2021 года
- Последняя версия – от 1 июня 2022 года
- На конференциях на данный момент не была представлена

Авторы

Всего 18 авторов, основных 4:

- **Reiichiro Nakano** – Member of Technical Staff, OpenAI, занимается reinforcement learning. Закончил филиппинский De La Salle University-Manila написал несколько статей, большинство из которых про NLP модели
- **Jacob Hilton** – researcher, OpenAI, был PhD student в теории комбинаторных множеств. Написал несколько работ по RL и теории множеств
- **Suchir Balaji** – Member Of Technical Staff at OpenAI, Участвовал в создании codex. *UC Berkeley*, В.А. Сначала был разработчиком, потом ушел в ML, участвовал в финале ICPC 2018
- **John Schulman** – Research Scientist, OpenAI, самый известный из авторов. Имеет большое количество публикаций в области RL. Возглавляет RL команду, занимающуюся языковыми моделями. PhD in Computer Science from *UC Berkeley*

Влияние на работу

Всего ссылается на 30 источников

Наибольшее влияние – статья [Learning to summarize from human feedback](#), так как:

- Она описывает используемый метод обучения на задачу суммаризации текста на основе human feedback
- Некоторые из авторов данной статьи приняли участие в создании статьи про WebGPT в качестве консультантов

Цитирования

- Всего 93 цитирования
- В качестве продолжения статья от Google – [LaMDA: Language Models for Dialog Applications](#) (семейство моделей на основе трансформеров, специализированных на диалоговых системах)

Преимущества работы

- Очень подробное описание процесса сбора данных
- Произведён подробный анализ различных методов обучения и архитектур
- Для сравнения качества рассмотрели достаточно сильно отличающиеся наборы данных
- Качество текста на высоком уровне
- Приложено множество деталей реализации, что дает шансы на воспроизводимость

Недостатки работы

- Не исследовано влияние способа сбора данных на итоговое качество модели. Нет уверенности в том, что при повторном сборе с помощью исполнителей, на выходе получится схожий результат
- Рассмотрен только 1 поисковый движок – Bing. Не исследовано влияние конкретного движка на качество модели

Предложения по улучшению

1. Исправить описанные недостатки :)
2. Авторы замечают, что модель плохо работает с непопулярными вопросами – возможно, стоит дополнительно исследовать подобные кейсы

ХАКЕР

Пилипенко Сергей