

Московский государственный университет имени М.В. Ломоносова Факультет
Вычислительной Математики и Кибернетики
Кафедра Математических Методов Прогнозирования

Левыкин Александр Михайлович

Детекция галлюцинаций больших языковых моделей

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА

Научный руководитель:
д.ф.-м.н., профессор
Воронцов К.В.

Москва, 2025

Содержание

1	Введение	2
1.1	Цель работы	4
2	Постановка задачи	5
3	Обзор существующих методов	5
4	Предложенный метод	6
4.1	Instruction-based подход	6
4.2	Обучение рекуррентной нейросети LSTM	7
4.3	Дообучение NER моделей	7
4.4	Методика оценивания модели	8
5	Вычислительный эксперимент	11
5.1	Данные	11
5.1.1	SemEval-2025	11
5.2	Instruction-based подход	11
6	Заключение	13
7	Приложение	16

Аннотация

В данной работе рассматривается задача детекции галлюцинаций больших языковых моделей. Основное внимание уделяется решению задачи в token classification постановке, в которой требуется классифицировать на наличие галлюцинаций каждый token ответа модели. Анализируются подходы instruction-based, дообучение NER моделей, RAG подход, а также анализ временных рядов. Дополнительно рассматривается задача, в которой для фрагментов размечена степень уверенности, с которой он был отнесён к галлюцинации, и соответствующие такому подходу критерии оценки качества моделей детекции.

1 Введение

Задача детекции галлюцинаций в больших языковых моделях (Large Language Models, LLM) является одной из подзадач обработки естественного языка (NLP), которая направлена на выявление ложных или недостоверных ответов, генерируемых моделями. Галлюцинации могут проявляться в виде неправдивых утверждений, вымышленных данных и фактов, что ограничивает точность и применимость LLM в ряде критически важных задач. Детекция таких ошибок критически важна в областях, где на основе необходима для обеспечения надежности и повышения доверия к моделям в реальных приложениях.

Применение детекции галлюцинаций может быть полезно в следующих областях:

- В создании автоматических систем поддержки решений (например, в медицине или праве) идентификация ложных данных в ответах модели помогает минимизировать риски и избегать принятия решений на основе недостоверной информации [1].
- В задачах генерации текстов и чат-ботах применение детекции галлюцинаций позволяет улучшить качество и достоверность предоставляемых ответов, помогая пользователям получать более точную информацию [2].
- В аналитике социальных медиа и маркетинговых исследованиях детекция галлюцинаций позволяет более эффективно анализировать настроения и тренды, исключая из анализа ложные и искаженные данные [3].

Галлюцинация LLM — это ответ (или его часть), сгенерированный моделью, который не соответствует входным данным (промпту) или ранее сгенерированному контексту, либо противоречит общеизвестным фактам о мире.

Задачи детекции галлюцинаций в больших языковых моделях можно классифицировать по нескольким основным признакам, включая подходы к детекции, использование эталонных данных и уровень детализации анализа.

1. Детекция галлюцинаций с использованием эталонных данных (reference-based hallucination detection) Одним из распространённых методов является эталонная (reference-based) детекция галлюцинаций, при которой сгенерированный текст сравнивается с имеющимся эталоном. Этот подход успешно применяется в таких задачах, как абстрактное суммирование текста (Maynez et al., 2020), машинный перевод (Wang and Sennrich, 2020), генерация текста на основе данных (Rebuffel

et al., 2021) и создание подписей к изображениям (Rohrbach et al., 2018). Однако для многих задач генерации текста общего формата эталонные данные могут быть недоступны, что ограничивает применимость таких методов. Например, в диалоговых системах реального времени, таких как чат-боты, модель часто генерирует ответы без возможности сравнения с эталоном.

2. Безэталонная детекция галлюцинаций (reference-free hallucination detection) В условиях, где эталонные данные отсутствуют, используются методы безэталонной детекции галлюцинаций, которые основываются только на контексте и правилах, встроенных в модель. Такие подходы становятся актуальными в системах, работающих в реальном времени (например, чат-боты и системы автозаполнения текста), где сравнение с эталоном затруднено из-за невозможности получить эталонные данные. Эти методы направлены на определение неконсистентных и потенциально ложных утверждений исключительно по информации, доступной модели в контексте текущего сеанса взаимодействия.

Галлюцинации можно классифицировать по следующей структуре:

- Intrinsic галлюцинации (внутренние галлюцинации) внутри себя делятся на:
 - Input-conflict галлюцинации — это случай, когда генерируемый текст не соответствует запросу пользователя.
 - Context-conflict галлюцинации - случай, когда модель генерирует высказывание, противоречащее ранее сгенерированной информации.
- Extrinsic (fact-conflict) галлюцинации (внешние галлюцинации) — эти галлюцинации определяются различно в зависимости от задачи:
 - генерируемый текст не соответствует фактам реального мира (для referenced-free постановки задачи),
 - генерируемый текст не может быть подтвержден информацией, содержащейся во входных данных - эталонном документе (для reference-based постановки).

Кроме того, подходы к решению задачи детекции галлюцинации можно поделить на два вида по степени детальности обнаружения:

- ****Уровень предложений или документов****: Многие подходы, такие как детекция фейковых новостей (Zellers et al., 2019) или факт-чекинг (Thorne и Vlachos, 2018), анализируют галлюцинации на уровне предложения или всего документа. Однако такой подход может быть недостаточно точным для выявления конкретных ложных утверждений внутри текста, так как он оценивает весь текст в целом.
- ****Детекция на уровне токенов****: В альтернативных методах детекция галлюцинаций осуществляется на уровне отдельных токенов. Такой подход позволяет более точно идентифицировать ложные утверждения в тексте и, в некоторых случаях, корректировать сгенерированный текст в реальном времени (например, управляя вероятностью появления тех или иных токенов). Этот метод более эффективен для задач, требующих высокую точность детекции и исправления, таких как системы генерации текста реального времени.

1.1 Цель работы

Reference-free постановка задачи являются более сложной с той точки зрения, что в ней не задан документ-образец, относительно которого идёт поиск галлюцинаций. Из-за этого качество таких моделей ниже. В данной работе исследуются различные решения для reference-free задач (Question Answering и детекция фактологических ошибок) и предлагается подход со сведением её к reference-based постановке. Кроме того, задача решается на уровне токенов, в token-classification постановке, что дополнительно усложняет решение задачи и делает исследование особенно ценным и комплексным.

2 Постановка задачи

Пусть заданы:

- Запрос пользователя (query, промпт, задача) - последовательность токенов $Q = \{q_1, \dots, q_m\}$, где m - длина последовательности.
- Ответ модели - последовательность токенов $X = \{x_1, x_2, \dots, x_n\}$, где n - длина последовательности.

Задача детекции галлюцинаций заключается в поиске всех пар вида (i_k, j_k) , где i_k и j_k , такие что $i_k \leq j_k$ - индексы начала и конца текстового фрагмента $\overline{x_{i_k} x_{i_k+1} \dots x_{j_k}}$, являющегося галлюцинацией.

Предполагается отсутствие вложенных сущностей (плоская задача NER), поэтому постановка выше эквивалентна следующей:

Задача детекции галлюцинаций заключается в классификации каждого токена ответа $X = \{x_1, x_2, \dots, x_n\}$ на один из 3-х классов $y_i \in \{B, I, O\}$, где B обозначает начало галлюцинации, I — продолжение галлюцинации, а O — токен, не являющийся частью галлюцинации.

3 Обзор существующих методов

dfgdfgdfgdfg[1]fkgdfgdfgdfg [2]

4 Предложенный метод

4.1 Instruction-based подход

В предложенном подходе для детекции галлюцинаций используется instruction-based метод, при котором большая языковая модель (LLM) генерирует текст, явно указывая на галлюцинации в ответе. В данном эксперименте вместо вывода вероятностей или меток на уровне токенов модель на естественном языке выделяет фрагменты текста, которые, по её мнению, являются галлюцинациями.

Одним из ключевых свойств крупных языковых моделей является способность решать задачи, на которые они не были специально обучены, используя инструкции (промпты), формирующие поведение модели в ходе генерации текста.

Исследования по подбору промпта широко развились [5, 3] и получили термин “prompt engineering”, а способность генеративных моделей проявлять умения, которые от них не требовались при обучении, получила название эмерджентным поведением (emergent behavior) и повлекло множество исследований [6, 4].

Формирование промпта для явной индикации галлюцинаций: Данные в задаче представляют из себя наборы пар: запрос пользователя $Q = \{q_1, q_2, \dots, q_m\}$ и ответ модели $X = \{x_1, x_2, \dots, x_n\}$. Задача состоит в создании промпта P , который направит модель на явное указание фрагментов, являющихся галлюцинациями.

Преимущества и недостатки instruction-based метода

- **Преимущества:**

- **Адаптивность:** позволяет модели решать задачи, на которые она не была специально обучена, благодаря использованию промптов.
- **Экономичность:** отсутствует необходимость в дообучении модели, что снижает затраты на дополнительные данные и вычислительные ресурсы.
- **Универсальность:** подходит для широкого спектра задач, где можно управлять поведением модели через промпты.
- **Интерпретируемость:** модель может выдавать галлюцинации в явном виде, делая выводы более прозрачными для пользователей.

- **Недостатки:**

- **Чувствительность к формулировке промпта:** результат может сильно зависеть от точной формулировки запроса, что требует времени на подбор оптимальных инструкций.
- **Отсутствие гарантий точности:** модель может ошибаться в детекции галлюцинаций, так как метод не обеспечивает вероятностных оценок надёжности.
- **Ограниченная гибкость:** сложно управлять степенью уверенности модели в предсказаниях без дополнительных вероятностных данных.
- **Зависимость от масштабов модели:** для эффективного выполнения задачи могут требоваться крупные модели, что увеличивает вычислительные затраты.

4.2 Обучение рекуррентной нейросети LSTM

4.3 Дообучение NER моделей

Задача выделения фрагментов галлюцинаций в тексте представляет собой задачу классификации токенов, в которой каждому токenu присваивается метка, указывающая, является ли он частью галлюцинации или нет. Для решения этой задачи используется формат BIO-разметки (отмечающий начало, середину и конец фрагментов, содержащих галлюцинации). Формат BIO позволяет более точно определять начало и границы фрагментов и показывает лучшие результаты в сравнении с форматом бинарной классификации каждого токена. В данном разделе мы подробно опишем задачу обучения модели и введем основные обозначения.

Формализация задачи

Рассмотрим текст, состоящий из N токенов, представленных последовательностью $X = (x_1, x_2, \dots, x_N)$. Модель должна классифицировать каждый токен x_i как принадлежащий галлюцинации или нет, используя BIO-формат, где:

- **B** обозначает начало фрагмента галлюцинации,
- **I** обозначает токен, продолжающий галлюцинацию,
- **O** обозначает токен, не связанный с галлюцинацией.

Для этого модель обучается предсказывать вероятности \hat{S} для каждого токена текста.

Обозначения матриц меток и предсказаний:

$S \in \{0, 1\}^{N \times 3}$: бинарная матрица, содержащая истинные метки для каждого токена, где:

- $S_{i,0} = 1$, если токен x_i отмечен как **B**,
- $S_{i,1} = 1$, если токен x_i отмечен как **I**,
- $S_{i,2} = 1$, если токен x_i отмечен как **O**.

$\hat{S} \in [0, 1]^{N \times 3}$: матрица предсказаний модели, где $\hat{S}_{i,t}$ обозначает вероятность того, что токен x_i относится к классу $t \in \{B, I, O\}$.

Ограничения нормировки:

Для каждого токена модель должна предсказывать распределение вероятностей по трем меткам, что выражается следующим условием:

$$\sum_t \hat{S}_{i,t} = 1, \quad \forall i \in \{1, \dots, N\}$$

Функция потерь — кросс-энтропия:

Для оптимизации параметров span-detection моделей используется функция потерь на основе кросс-энтропии, которая минимизируется по параметрам модели θ' :

$$CE(S, \hat{S}) = -\frac{1}{N} \sum_{i=1}^N \sum_{t \in \{B, I, O\}} S_{i,t} \log(\hat{S}_{i,t})$$

Здесь: - $S_{i,t}$ — бинарный индикатор для истинной метки токена x_i по классу t , - $\hat{S}_{i,t}$ — вероятность, предсказанная моделью для метки t токена x_i .

И задача обучения модели заключается в поиске параметров, доставляющих минимум функции потерь:

$$CE(S, p(S | \theta)) = CE(S, \hat{S}) = -\frac{1}{N} \sum_{i=1}^N \sum_{t \in \{B, I, O\}} S_{i,t} \log(\hat{S}_{i,t}) \longrightarrow \min_{\theta}$$

4.4 Методика оценивания модели

Обозначения для i -го класса:

- True Positive (TP_i): количество токенов или фрагментов, которые модель правильно классифицировала как принадлежащие классу i . Определяется как количество случаев, в которых предсказанный класс и истинный класс равны i :

$$TP_i = \sum_{j=1}^N \mathbb{I}(\hat{y}_j = i \wedge y_j = i)$$

где N — общее количество токенов, \hat{y}_j — предсказанный класс для токена j , y_j — истинный класс для токена j , а $\mathbb{I}(\cdot)$ — индикаторная функция, принимающая значение 1, если условие выполняется, и 0 в противном случае.

- False Positive (FP_i): количество токенов или фрагментов, которые модель ошибочно классифицировала как принадлежащие классу i . Определяется как количество случаев, в которых предсказанный класс равен i , но истинный класс не равен i :

$$FP_i = \sum_{j=1}^N \mathbb{I}(\hat{y}_j = i \wedge y_j \neq i)$$

- False Negative (FN_i): количество токенов или фрагментов, которые модель ошибочно не отнесла к классу i . Определяется как количество случаев, в которых истинный класс равен i , но предсказанный класс не равен i :

$$FN_i = \sum_{j=1}^N \mathbb{I}(\hat{y}_j \neq i \wedge y_j = i)$$

Метрики для i -го класса:

1. Точность (Precision) для i -го класса: Точность показывает, какая доля токенов, предсказанных моделью как относящиеся к классу i , действительно относятся к этому классу.

$$\text{Precision}_i = \frac{TP_i}{TP_i + FP_i}$$

2. Полнота (Recall) для i -го класса: Полнота показывает, какую долю токенов класса i модель правильно предсказала, по отношению к общему числу токенов данного класса в данных.

$$\text{Recall}_i = \frac{TP_i}{TP_i + FN_i}$$

3. F1-score для i -го класса: F1-score является гармоническим средним точности и полноты и даёт общую оценку качества предсказания для класса i .

$$F1_i = 2 \cdot \frac{\text{Precision}_i \cdot \text{Recall}_i}{\text{Precision}_i + \text{Recall}_i}$$

Общие метрики для всех классов:

4. F1-macro (макроусреднённый F1-score): F1-macro вычисляется как среднее значение F1-score по всем классам, где каждый класс имеет одинаковый вес. Это полезно, когда важно учитывать качество предсказаний для каждого класса независимо от его размера.

$$F1\text{-macro} = \frac{1}{C} \sum_{i=1}^C F1_i$$

где C — количество классов.

5. F1-micro (микроусреднённый F1-score): F1-micro учитывает общие суммы TP , FP и FN по всем классам и вычисляет F1 на основе общей точности и полноты. Он полезен, когда важен вклад каждого примера, вне зависимости от размера класса.

$$F1\text{-micro} = \frac{2 \cdot \sum_{i=1}^C TP_i}{2 \cdot \sum_{i=1}^C TP_i + \sum_{i=1}^C FP_i + \sum_{i=1}^C FN_i}$$

Дополнительно рассматривается задача, в которой для каждого фрагмента текста указана степень уверенности, с которой он был отнесён к галлюцинации. Чтобы оценить качество модели в таких случаях, используется корреляция Спирмана, которая позволяет измерить степень согласованности между ранжированными предсказаниями модели и истинными значениями.

Пусть имеются следующие значения для N фрагментов:

- r_i — ранг степени уверенности для i -го фрагмента, предсказанный моделью,
- s_i — истинный ранг степени уверенности для i -го фрагмента.

Корреляция Спирмана (ρ) рассчитывается по формуле:

$$\rho = 1 - \frac{6 \sum_{i=1}^N (r_i - s_i)^2}{N(N^2 - 1)}$$

где N — количество фрагментов. Данная формула основывается на квадрате разницы между рангами предсказанных и истинных значений и позволяет оценить, насколько предсказанные уверенности соответствуют истинным.

Критерий **IoU** (Intersection over Union, пересечение над объединением) используется для оценки схожести предсказанных и истинных фрагментов. IoU измеряет долю пересечения между предсказанным фрагментом и истинным фрагментом от их объединения. Формула IoU определяется как:

$$\text{IoU} = \frac{|A \cap B|}{|A \cup B|}$$

где:

- A — множество токенов, принадлежащих истинному фрагменту галлюцинации.
- B — множество токенов, принадлежащих предсказанному фрагменту галлюцинации.
- $|A \cap B|$ — количество токенов, которые находятся одновременно в истинном и предсказанном фрагментах.
- $|A \cup B|$ — количество токенов, которые находятся хотя бы в одном из фрагментов (их объединение).

Для задачи с несколькими фрагментами галлюцинаций метрика IoU может быть усреднена по всем парам предсказанных и истинных фрагментов, чтобы получить средний IoU.

5 Вычислительный эксперимент

5.1 Данные

5.1.1 SemEval-2025

В рамках активного соревнования Semeval-2025(<https://helsinki-nlp.github.io/shroom/>) на момент написания статьи доступна размеченная валидационная выборка. Она состоит из 10 файлов на разных языках. В каждом файле 50 объектов-троек: вопрос-ответ-размеченные галлюцинации.

5.2 Instruction-based подход

Для проведения эксперимента выбрана модель Meta-Llama-3.1-8B-Instruct. Модель была представлена в 2024 году, обучалась на данных до декабря 2023 года и содержит одни из самых актуальных знаний среди больших языковых моделей, а также показывает высокие результаты на бенчмарках. Кроме того, платформа kaggle с видеокартой Nvidia P100 с 16 Gb видеопамати добавляет ограничения на размер модели, и Llama-3.1-8B занимает всю доступную видеопамать. Так, для экспериментов была выбрана эта модель как наиболее актуальная и мощная из возможных для запуска на kaggle.

Для эксперимента взяты 50 объектов из валидационной выборки соревнования SemEval-2025 и проведено множество экспериментов с промптами. Наилучший результат достигнут для следующего промпта:

```
messages = [ {'role': 'system', 'content': 'You are a fact-checking assistant. Your task is to identify fragments of the response that are hallucinations - parts of the text that are factually incorrect or made up by model. Pay attention to facts, dates, numbers, places. Detect only hallucination words, without neighbour words. Give me only a list of fragments-hallucinations you found in model output.'}, {'role': 'user', 'content': formatted_str}, ]
```

В formatted_str подавалась строка следующего вида:

```
formatted_str = "query":"{data['model_input']}" "model_output":  
"{data['model_output_text']}"
```

Для оценки качества по регламенту конкурса используется метрика IoU.

Результаты: распределение IoU на объектах в данном эксперименте представлено на рис. 1.

Средний IoU составил **39.7%**.

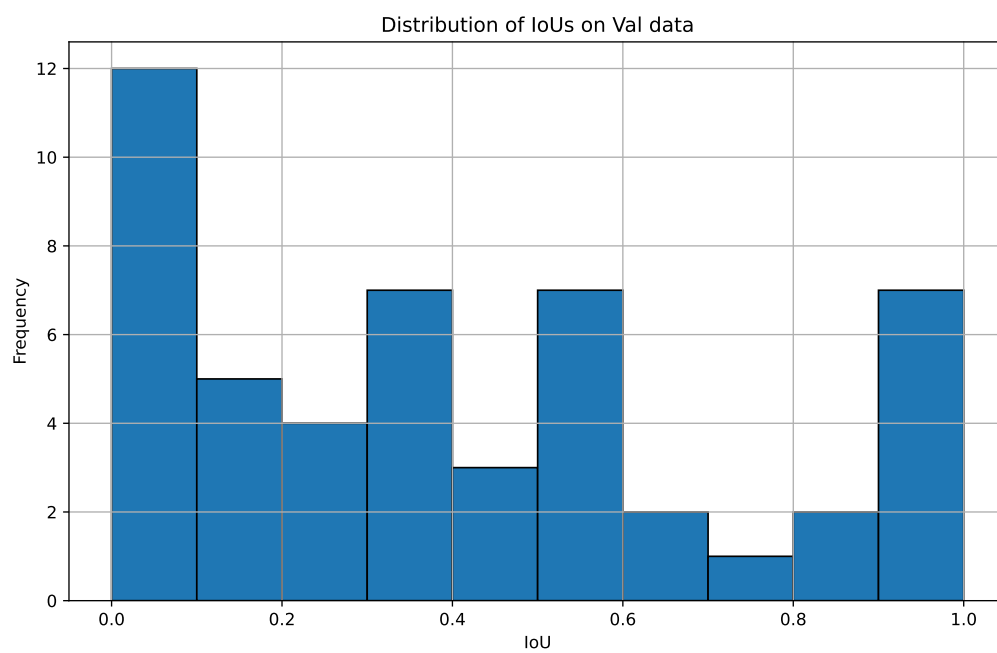


Рис. 1: распределение IoU на объектах SemEval-2025.

6 Заключение

В данной работе были рассмотрены разные векторные представления фрагментов для задачи классификации связей, а также рассмотрен подход MRC. В результате экспериментов самым лучшим по F_1 с микро усреднением оказался метод MRC с векторным представлением фрагментов, которое строится усреднением через механизм внимания.

Был предложен критерий оценки качества моделей для совместного выделения сущностей и фрагментов, который сглаживает F-меру и позволяет лучше оценить качество разметки. Лучший результат показал подход MRC. Основная проблема методов - низкое качество поиска релевантных пар фрагментов.

Дальнейшая работа связана с улучшением качества совместного выделения фрагментов и связей между ними. [1]

Список литературы

- [1] Abaho, M., Bollegala, D., Williamson, P., Dodd, S.: Detect and classify—joint span detection and classification for health outcomes. arXiv preprint arXiv:2104.07789 (2021)
- [2] Baly, R., Martino, G.D.S., Glass, J., Nakov, P.: We can detect your bias: Predicting the political ideology of news articles. arXiv preprint arXiv:2010.05338 (2020)
- [3] Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., Neubig, G.: Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing (2021), <https://arxiv.org/abs/2107.13586>
- [4] OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H.W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S.P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S.S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N.S., Khan, T., Kilpatrick, L., Kim, J.W., Kim, C., Kim, Y., Kirchner, J.H., Kiros, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C.M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S.M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H.P., Michael, Pokornyy, Pokrass, M., Pong, V.H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F.P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M.B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J.F.C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J.J., Wang, A., Wang, B., Ward, J., Wei, J.,

- Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., Zoph, B.: Gpt-4 technical report (2024), <https://arxiv.org/abs/2303.08774>
- [5] Reynolds, L., McDonell, K.: Prompt programming for large language models: Beyond the few-shot paradigm (2021), <https://arxiv.org/abs/2102.07350>
- [6] Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E.H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J., Fedus, W.: Emergent abilities of large language models (2022), <https://arxiv.org/abs/2206.07682>

7 Приложение