

Обнаружение галлюцинаций в больших языковых моделях с метаморфными отношениями

Дата: 2025-02-20 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.15844>

Рейтинг: 70

Адаптивность: 80

Ключевые выводы:

Исследование представляет MetaQA - новый метод обнаружения галлюцинаций в больших языковых моделях (LLM), основанный на метаморфических отношениях. Основной результат: MetaQA превосходит существующие методы обнаружения галлюцинаций, не требуя внешних ресурсов и работая как с открытыми, так и с закрытыми LLM.

Объяснение метода:

MetaQA предлагает метод обнаружения галлюцинаций через синонимические и антонимические мутации ответов без внешних ресурсов. Подход применим для всех LLM, но полная реализация трудоемка. Пользователи могут адаптировать основную концепцию, перефразируя вопросы и проверяя согласованность ответов, что делает метод доступным даже без специальных знаний.

Ключевые аспекты исследования: 1. **MetaQA** - новый метод обнаружения галлюцинаций в LLM, использующий метаморфические отношения (MR) и мутации запросов без внешних ресурсов, работающий как с открытыми, так и с закрытыми моделями.

Метаморфические отношения для обнаружения несоответствий - метод генерирует синонимические и антонимические мутации исходного ответа, а затем проверяет их фактическую корректность, выявляя противоречия.

Самопроверка без внешних ресурсов - в отличие от существующих методов, MetaQA не требует внешних баз данных или API, используя только саму LLM для генерации и проверки мутаций.

Превосходство над существующими методами - исследование показывает, что MetaQA превосходит SelfCheckGPT по показателям точности, полноты и F1-оценки во всех тестируемых моделях и наборах данных.

Универсальность применения - метод работает с разными типами вопросов и

категориями знаний, показывая стабильные результаты при обнаружении фактических несоответствий.

Дополнение: Действительно ли для работы методов этого исследования требуется дообучение или API? Или методы и подходы можно применить в стандартном чате, а ученые лишь для своего удобства использовали расширенные техники?

Методы MetaQA можно применить в стандартном чате без дообучения или специального API. Авторы использовали API и программные реализации для масштабного тестирования и точного расчёта метрик, но сама концепция работает в обычном диалоге с LLM.

Основные концепции, которые можно применить в стандартном чате:

Синонимические мутации - переформулировать вопрос, сохраняя смысл: Исходный: "Какой процент мозга использует человек?" Мутация: "Какая доля мозга активна у среднестатистического человека?"

Антонимические мутации - задать вопрос с противоположным смыслом:

Исходный: "Использует ли человек только 10% своего мозга?" Мутация: "Верно ли, что человек использует более 10% своего мозга?"

Верификация мутаций - попросить LLM проверить фактическую точность утверждения:

"Является ли фактически верным утверждение, что человек использует только 10% мозга?" Применяя эти подходы, пользователь может: - Выявить несоответствия в ответах на похожие вопросы - Обнаружить, когда модель неуверена в ответе - Проверить фактическую точность информации

Результаты будут менее формализованы, чем в исследовании, но сама методика обнаружения противоречий через метаморфические отношения полностью применима в обычном чате и не требует специальных технических знаний.

Анализ практической применимости: 1. **Метод обнаружения галлюцинаций (MetaQA)** - Прямая применимость: Высокая. Пользователи могут адаптировать технику создания синонимических и антонимических вариаций ответа и проверять их согласованность для выявления потенциальных галлюцинаций. - Концептуальная ценность: Очень высокая. Метод показывает, что противоречия в ответах на похожие вопросы могут указывать на галлюцинации, что формирует критическое мышление при работе с LLM. - Потенциал для адаптации: Отличный. Техника применима в обычных чатах без специальных API или инструментов.

Метаморфические отношения (синонимы и антонимы) Прямая применимость: Средняя. Пользователи могут самостоятельно перефразировать вопросы или задать противоположные вопросы для проверки согласованности ответов.

Концептуальная ценность: Высокая. Понимание того, что LLM должны давать согласованные ответы на семантически эквивалентные запросы. Потенциал для адаптации: Высокий. Техника перефразирования легко применима в повседневном использовании.

Процесс верификации мутаций

Прямая применимость: Средняя. Пользователи могут задавать дополнительные проверочные вопросы для подтверждения информации. Концептуальная ценность: Высокая. Демонстрирует, как LLM могут проверять собственные утверждения при правильной формулировке запроса. Потенциал для адаптации: Средний. Требует дополнительных усилий, но выполнимо в рамках обычного чата.

Алгоритм оценки галлюцинаций

Прямая применимость: Низкая. Формальное вычисление оценки галлюцинаций сложно для обычных пользователей. Концептуальная ценность: Средняя. Понимание того, что противоречивые ответы указывают на возможные галлюцинации. Потенциал для адаптации: Средний. Пользователи могут интуитивно оценивать согласованность ответов без формальных вычислений.

Результаты исследования различных моделей

Прямая применимость: Средняя. Пользователи получают информацию о склонности разных LLM к галлюцинациям. Концептуальная ценность: Высокая. Формирует понимание, что все LLM склонны к галлюцинациям и требуют проверки. Потенциал для адаптации: Высокий. Знание о различиях между моделями помогает выбрать подходящую для критически важных задач.

Prompt:

Применение MetaQA в промптах для GPT ## Краткое объяснение

Исследование MetaQA предлагает метод обнаружения галлюцинаций в LLM с помощью метаморфических отношений. Основная идея заключается в создании мутаций ответа (синонимичных и антонимичных версий) и проверке их согласованности, что позволяет выявить потенциальные галлюцинации без внешних ресурсов.

Пример промпта с применением методологии MetaQA

[=====] Я хочу получить от тебя максимально точную и надежную информацию о [ТЕМА]. Для этого используем метод MetaQA:

Дай краткий ответ на вопрос: [ОСНОВНОЙ ВОПРОС]

Теперь перефразируй свой ответ тремя разными способами, сохраняя то же значение:

Вариант 1: [перефразируй ответ] Вариант 2: [перефразируй ответ] Вариант 3: [перефразируй ответ]

Теперь сформулируй противоположное утверждение к своему ответу:

Антонимичное утверждение: [противоположный ответ]

Проверь каждую из версий (включая антонимичную) на фактическую точность. Оцени их как "Верно", "Частично верно" или "Неверно".

На основе этого анализа, определи, есть ли в твоём первоначальном ответе галлюцинации или неточности. Если обнаружены расхождения, предоставь исправленный и уточненный ответ.

Оцени уровень своей уверенности в окончательном ответе по шкале от 1 до 5. [=====]

Почему это работает

Данный промпт реализует ключевые шаги методологии MetaQA: 1. Получение базового ответа 2. Создание мутаций (синонимичных и антонимичных версий) 3. Проверка мутаций на фактическую точность 4. Оценка вероятности галлюцинации

Метод эффективен, поскольку: - Заставляет модель рассмотреть информацию с разных формулировок - Выявляет несоответствия через сравнение версий - Проверяет реакцию модели на противоположные утверждения - Не требует внешних ресурсов для верификации - Работает как самопроверка в рамках одного запроса

Такой подход значительно повышает точность ответов в критически важных областях, где фактическая достоверность имеет первостепенное значение.