

# Оценка способности LLM к восприятию смешанных контекстов через призму суммирования

Дата: 2025-03-03 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.01670>

Рейтинг: 65

Адаптивность: 70

## Ключевые выводы:

Исследование направлено на оценку способности больших языковых моделей (LLM) выявлять смешанные контекстные галлюцинации в задаче суммаризации текста. Основные результаты показывают, что внутренние знания LLM создают предвзятость в оценке галлюцинаций, особенно при обнаружении фактических галлюцинаций, что является основным узким местом производительности. Ключевая проблема заключается в эффективном использовании знаний, балансируя между внутренними знаниями LLM и внешним контекстом.

## Объяснение метода:

Исследование предоставляет ценное понимание различных типов галлюцинаций LLM и методов их выявления. Пользователи могут адаптировать концепции фактических/нефактических галлюцинаций и стратегии проверки (CoT, ICL, внешние источники) для повседневного использования. Однако многие методы технически сложны и требуют значительной адаптации для неспециалистов.

**## Ключевые аспекты исследования:** 1. **Исследование оценки смешанного контекста галлюцинаций через призму суммаризации** - работа анализирует способность LLM распознавать два типа галлюцинаций: фактические (фактически верные, но отсутствующие в источнике) и нефактические (фактически неверные).

**Создание специализированного датасета FHSumBench** - авторы разработали автоматизированный конвейер для создания сбалансированного набора данных с различными типами галлюцинаций в суммаризации текста.

**Сравнение различных методов оценки** - исследование сравнивает прямую генерацию и методы на основе поиска информации для выявления галлюцинаций в смешанном контексте.

**Влияние размера модели** - работа анализирует, как масштабирование моделей влияет на способность выявлять разные типы галлюцинаций.

**Проблема внутреннего знания LLM** - исследование выявляет, что внутреннее знание моделей создает предвзятость при оценке галлюцинаций, особенно фактических.

## Дополнение:

### Применимость методов в стандартном чате без дообучения и API

Исследование описывает несколько методов, которые **можно применить в стандартном чате** без дополнительного API или дообучения:

**Использование CoT (Chain-of-Thought)** - Пользователи могут запрашивать пошаговые рассуждения от LLM для проверки фактов. Исследование показывает, что это улучшает выявление нефактических галлюцинаций.

**Использование ICL (In-Context Learning)** - Предоставление моделям примеров правильной оценки галлюцинаций помогает им лучше определять проблемные утверждения. Это особенно полезно для моделей меньшего размера.

**Разделение текста на утверждения** - Пользователи могут разбивать длинные тексты на отдельные утверждения и проверять каждое отдельно, как это делается в методах Knowledge Retrieval.

**Двухэтапная проверка (аналог Reflection Retrieval)** - Пользователи могут сначала проверить соответствие ответа исходному запросу, а затем отдельно проверить фактическую точность сомнительных утверждений.

**Осознание предвзятости к внутренним знаниям** - Понимание того, что LLM могут считать фактически верную, но не подтвержденную источником информацию правильной, помогает пользователям быть более критичными к ответам моделей.

Исследователи действительно использовали API и специализированные инструменты для масштабного тестирования, но концептуальные подходы применимы и в стандартном чате. Результаты показывают, что правильно сформулированные запросы могут существенно улучшить способность моделей выявлять галлюцинации.

## Анализ практической применимости: 1. **Автоматизированное создание датасета FHSumBench** - Прямая применимость: Низкая для рядовых пользователей, требует технических знаний и доступа к специализированным инструментам - Концептуальная ценность: Высокая, демонстрирует разницу между фактическими и нефактическими галлюцинациями, что помогает пользователям лучше понимать ответы LLM - Потенциал для адаптации: Средний, пользователи могут научиться различать типы галлюцинаций в ответах, хотя само создание тестовых наборов требует специализированных навыков

**Методы выявления галлюцинаций (прямая генерация vs. методы с поиском)**

Прямая применимость: Средняя, пользователи могут использовать промпты с CoT (цепочкой рассуждений) или ICL (обучением в контексте) для улучшения оценки достоверности  
Концептуальная ценность: Высокая, понимание преимуществ разных подходов к проверке информации  
Потенциал для адаптации: Высокий, пользователи могут адаптировать методы для проверки фактов в обычных чатах

### **Влияние внутреннего знания на выявление галлюцинаций**

Прямая применимость: Средняя, понимание ограничений LLM при проверке фактов  
Концептуальная ценность: Высокая, осознание того, что модели могут быть предвзяты к своим внутренним знаниям  
Потенциал для адаптации: Высокий, пользователи могут использовать внешние источники для проверки сомнительных утверждений

### **Влияние размера модели на качество оценки**

Прямая применимость: Низкая, пользователи обычно не могут выбирать размер модели  
Концептуальная ценность: Средняя, понимание того, что большие модели не всегда лучше в оценке галлюцинаций  
Потенциал для адаптации: Низкий, ограниченные возможности для применения пользователями

### **Методы с поиском информации (Knowledge Retrieval, Concurrent Retrieval, Reflection Retrieval)**

Прямая применимость: Средняя, пользователи могут адаптировать эти подходы для проверки информации  
Концептуальная ценность: Высокая, понимание важности внешних источников для проверки фактов  
Потенциал для адаптации: Высокий, можно использовать похожие стратегии в обычном взаимодействии с LLM

## **Prompt:**

Использование знаний из исследования о смешанных контекстных галлюцинациях в промптах ## Ключевые уроки исследования для промптинга

Исследование показывает, что LLM испытывают трудности с обнаружением смешанных контекстных галлюцинаций, особенно фактических (когда информация верна, но отсутствует в исходном тексте). Это можно использовать для улучшения промптов.

## Пример улучшенного промпта для суммаризации текста

[=====] Я хочу, чтобы ты создал краткое резюме следующего текста, используя метод рефлексивного извлечения информации:

Сначала прочитай весь текст внимательно. Для каждого ключевого утверждения, которое ты планируешь включить в резюме: Выдели точную цитату или отрывок из исходного текста, подтверждающий это утверждение Проверь, не добавляешь ли ты дополнительную информацию из своих внутренних знаний Если информация

кажется правильной, но отсутствует в тексте, пометь ее как "не подтверждено источником" Составь итоговое резюме, используя только подтвержденную информацию из текста. В конце резюме укажи любые моменты, где ты заметил конфликт между содержанием текста и твоими внутренними знаниями. Вот текст для резюме: [ВСТАВИТЬ ТЕКСТ] [=====]

## Почему это работает

Данный промпт использует несколько ключевых выводов исследования:

**Применение рефлексивного извлечения** - согласно исследованию, этот метод показал наилучшие результаты (F1-показатель 0.5010 на GPT-4o) **Явное разделение источников информации** - промпт заставляет модель осознанно различать информацию из текста и из собственных знаний **Пошаговый подход** - структурированный процесс помогает модели избежать предвзятости внутренних знаний, что исследование определило как основное узкое место **Прозрачность в отношении потенциальных конфликтов** - модель документирует случаи, когда ее знания противоречат исходному тексту Такой промпт помогает снизить вероятность появления смешанных контекстных галлюцинаций, особенно фактических, которые, как показало исследование, являются наиболее проблематичными для LLM.