

# FAST-AUDIT: Адаптивная многоагентная структура для динамической оценки проверки фактов больших языковых моделей

Дата: 2025-03-02 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.17924>

Рейтинг: 72

Адаптивность: 75

## Ключевые выводы:

Исследование представляет FAST-AUDIT - адаптивную мультиагентную систему для динамической оценки способностей больших языковых моделей (LLM) к проверке фактов. Основная цель - выявить ограничения LLM в проверке фактов, оценивая не только точность вердикта, но и качество обоснования. Результаты показали значительные различия в производительности между проприетарными и открытыми моделями, а также выявили конкретные сценарии, представляющие наибольшую сложность для LLM при проверке фактов.

## Объяснение метода:

FAST-AUDIT предлагает ценную методологию для оценки способностей LLM в проверке фактов, включая анализ обоснований, а не только вердиктов. Исследование предоставляет структурированную таксономию типов фактчекинга и данные о производительности 13 моделей, что помогает пользователям понять ограничения LLM и адаптировать свои ожидания. Основные принципы могут быть применены в повседневном взаимодействии.

Ключевые аспекты исследования: 1. **Адаптивная система оценки фактчекинга:** FAST-AUDIT предлагает многоагентную систему, которая динамически оценивает способности языковых моделей (LLM) проверять факты, адаптируясь к конкретным слабостям моделей.

**Оценка обоснований, а не только вердиктов:** В отличие от традиционных методов, сосредоточенных на точности классификации, FAST-AUDIT оценивает также качество объяснений, которые LLM предоставляют для своих выводов о достоверности информации.

**Метод выборки по значимости:** Используется алгоритм, который целенаправленно выбирает более сложные и разнообразные сценарии проверки

фактов, что позволяет эффективнее выявлять ограничения моделей.

**Итеративное зондирование:** Система генерирует новые тестовые случаи на основе анализа предыдущих результатов, что позволяет обнаруживать более тонкие ограничения LLM в проверке фактов.

**Таксономия проверки фактов:** Разработана детальная классификация различных сценариев проверки фактов (сложные утверждения, фейковые новости, социальные слухи), что обеспечивает комплексную оценку.

Дополнение:

Исследование FACT-AUDIT действительно использует расширенные технические подходы, такие как многоагентную систему и API моделей, но многие его концепции и подходы можно адаптировать для использования в стандартном чате без необходимости дообучения или специальных API.

**Концепции и подходы для стандартного чата:**

**Оценка обоснований, а не только вердиктов:** Пользователи могут запрашивать у модели не только ответ на фактический вопрос, но и подробное объяснение. Затем они могут оценить качество этого объяснения, даже если вердикт кажется правильным.

**Итеративное зондирование:** Пользователи могут последовательно задавать уточняющие вопросы по теме, чтобы проверить согласованность и глубину знаний модели. Это помогает выявить потенциальные ограничения в понимании фактов.

**Использование различных режимов проверки:** Исследование показывает три режима проверки фактов: [claim] (только на основе утверждения), [evidence] (с предоставлением доказательств) и [wisdom of crowds] (с использованием "мудрости толпы"). Пользователи могут применять эти подходы, задавая вопросы в разных форматах.

**Таксономия проверки фактов:** Понимание различных категорий утверждений (сложные утверждения, фейковые новости, слухи) помогает пользователям формулировать более целенаправленные запросы и критически оценивать ответы.

**Адаптивное тестирование:** Пользователи могут сосредоточиться на темах, где модель показывает слабые результаты, и более тщательно проверять информацию в этих областях.

**Результаты от применения этих подходов:**

Более критическое отношение к ответам LLM, особенно в сложных областях знаний  
Выявление потенциальных неточностей или пробелов в знаниях модели  
Более глубокое понимание темы через итеративные вопросы  
Повышение качества взаимодействия с LLM за счет более структурированных запросов  
Способность

определить, когда требуется дополнительная проверка информации из независимых источников. Важно отметить, что хотя полная система FACT-AUDIT с множеством агентов и сложной оценкой требует технической экспертизы, основные принципы исследования вполне применимы в обычном чате и могут значительно улучшить качество взаимодействия с LLM и надежность получаемой информации.

Анализ практической применимости: 1. **Адаптивная система оценки фактчекинга** - Прямая применимость: Высокая. Пользователи могут использовать подход "проверки слабых мест" для понимания, когда модели могут ошибаться, и соответственно корректировать свои ожидания. - Концептуальная ценность: Значительная. Демонстрирует, что модели имеют разную производительность в различных сценариях проверки фактов, что важно для критического использования LLM. - Потенциал для адаптации: Средний. Требуется технической экспертизы для полной реализации, но концепция "проверки разных сценариев" может быть применена даже обычными пользователями.

**Оценка обоснований, а не только вердиктов** Прямая применимость: Средняя. Пользователи могут перенять подход оценки не только ответа, но и объяснения LLM, что повышает критическое мышление. Концептуальная ценность: Высокая. Понимание того, что правильный вердикт с неверным обоснованием может быть ненадежным, критически важно для пользователей. Потенциал для адаптации: Высокий. Легко адаптируется как практика для любого взаимодействия с LLM.

### **Метод выборки по значимости**

Прямая применимость: Низкая для обычных пользователей, высокая для разработчиков. Концептуальная ценность: Средняя. Демонстрирует, как можно целенаправленно тестировать модель в сложных случаях. Потенциал для адаптации: Средний. Принцип "тестирования на сложных случаях" может быть применен пользователями интуитивно.

### **Итеративное зондирование**

Прямая применимость: Средняя. Пользователи могут последовательно уточнять вопросы, чтобы проверить надежность ответов LLM. Концептуальная ценность: Высокая. Показывает, как последовательные вопросы могут выявлять несоответствия в знаниях LLM. Потенциал для адаптации: Высокий. Легко адаптируется как стратегия взаимодействия.

### **Таксономия проверки фактов**

Прямая применимость: Высокая. Пользователи могут использовать категории для понимания типов утверждений, с которыми LLM справляется лучше или хуже. Концептуальная ценность: Высокая. Структурирует понимание различных типов фактчекинга. Потенциал для адаптации: Высокий. Категоризация легко применима в повседневном взаимодействии с LLM. Сводная оценка полезности: На основе анализа я оцениваю общую полезность исследования для широкой аудитории в **72 балла из 100**.

Это исследование предлагает высокую практическую ценность для пользователей LLM, особенно в части понимания ограничений моделей при проверке фактов. Наиболее ценными аспектами являются:

Понимание важности проверки не только вердикта, но и обоснования LLM  
Осознание различной эффективности моделей в разных доменах и типах проверки фактов  
Структурированная таксономия для категоризации различных сценариев фактчекинга  
Стратегия итеративных вопросов для проверки надежности ответов  
Контраргументы к оценке:

Аргумент за более высокую оценку: Исследование предоставляет конкретные данные о производительности 13 современных LLM в различных сценариях проверки фактов, что дает пользователям практическое представление о том, каким моделям можно больше доверять.

Аргумент за более низкую оценку: Полная реализация методологии FACT-AUDIT требует значительных технических знаний и доступа к API моделей, что ограничивает непосредственное применение для большинства пользователей.

После рассмотрения этих аргументов я сохраняю оценку в 72 балла, так как хотя техническая сложность полного воспроизведения системы высока, основные принципы и выводы исследования могут быть применены широким кругом пользователей для более критического взаимодействия с LLM.

Исследование получает такую оценку за: 1. Практические методы для критической оценки фактчекинга в LLM 2. Ценные концепции для понимания ограничений моделей 3. Конкретные данные о сильных и слабых сторонах популярных моделей 4. Адаптируемость основных принципов для обычных пользователей 5. Структурированный подход к различным типам проверки фактов

Уверенность в оценке: Очень сильная. Исследование предоставляет достаточно информации о методологии и результатах, чтобы сделать обоснованные выводы о его полезности для широкой аудитории. Представленные данные о производительности различных моделей и анализ различных сценариев проверки фактов дают четкое представление о практической применимости исследования.

Оценка адаптивности: Оценка адаптивности: **75 из 100**

1) **Применимость принципов в обычном чате:** Многие ключевые принципы исследования могут быть применены в стандартном взаимодействии с LLM без необходимости в специальных API или инструментах. Пользователи могут использовать таксономию проверки фактов, стратегию итеративных вопросов и оценку обоснований, не только вердиктов.

2) **Извлечение полезных идей:** Исследование предоставляет богатый материал для формирования более осознанного подхода к взаимодействию с LLM. Идеи о

различной эффективности моделей в разных доменах и типах задач, а также о важности качественных обоснований, являются ценными для любого пользователя.

3) **Потенциал для будущих взаимодействий:** Подход к тестированию моделей в различных сценариях и с разными типами входных данных создает основу для более критического и эффективного использования LLM в будущем.

4) **Абстрагирование методов:** Хотя полная реализация многоагентной системы требует технической экспертизы, основные принципы, такие как итеративные вопросы, проверка обоснований и категоризация типов проверки фактов, легко абстрагируются до общих принципов взаимодействия.

Эта оценка отражает хороший баланс между техническими аспектами исследования, которые могут быть сложны для реализации обычными пользователями, и концептуальными принципами, которые могут быть легко адаптированы для повседневного использования.

|| <Оценка: 72> || <Объяснение: FACT-AUDIT предлагает ценную методологию для оценки способностей LLM в проверке фактов, включая анализ обоснований, а не только вердиктов. Исследование предоставляет структурированную таксономию типов фактчекинга и данные о производительности 13 моделей, что помогает пользователям понять ограничения LLM и адаптировать свои ожидания. Основные принципы могут быть применены в повседневном взаимодействии.> || <Адаптивность: 75>

## Prompt:

Использование результатов исследования FACT-AUDIT в промптах для ChatGPT  
**Ключевые знания из исследования для улучшения промптов**

Исследование FACT-AUDIT предоставляет ценные инсайты о том, как языковые модели справляются с проверкой фактов в различных сценариях. Эти знания можно эффективно применить при составлении промптов для получения более точных и обоснованных ответов.

## Пример промпта с учетом результатов исследования

[=====] Проверь следующее утверждение на фактическую точность:  
[УТВЕРЖДЕНИЕ]

Для проверки используй следующие доказательства: 1. [ДОКАЗАТЕЛЬСТВО 1] 2. [ДОКАЗАТЕЛЬСТВО 2] 3. [ДОКАЗАТЕЛЬСТВО 3]

Инструкции: 1. Разбей проверку на отдельные логические шаги 2. Для каждого шага укажи, какие доказательства ты используешь 3. Оцени каждый компонент утверждения отдельно 4. Предоставь итоговый вердикт (Правда/Частично

правда/Ложь) 5. Объясни свое обоснование, особенно если утверждение содержит статистические данные или требует многоэтапных рассуждений

Если ты не уверен в каком-то аспекте, явно укажи это в своем ответе. [=====]

## **Почему этот промпт работает лучше на основе исследования**

**Включение доказательств:** Исследование показало, что режим [evidence] (с доступом к доказательствам) значительно улучшает точность проверки фактов по сравнению с режимом [claim].

**Разбиение на шаги:** Промпт требует пошагового рассуждения, что помогает преодолеть сложности с многоэтапными рассуждениями (MSR) и статистическими утверждениями (ASR), которые оказались самыми проблемными сценариями.

**Отдельная оценка компонентов:** Этот подход помогает избежать ошибок в сложных утверждениях, содержащих несколько фактов.

**Акцент на обосновании:** Исследование выявило, что модели могут давать правильный вердикт с неверным обоснованием (метрика JFR), поэтому промпт специально запрашивает детальное объяснение.

**Признание неуверенности:** Поощряет модель явно указывать на неопределенность, что снижает риск уверенных, но неверных ответов.

Используя эти принципы, вы можете создавать промпты, которые компенсируют известные ограничения языковых моделей в проверке фактов, выявленные в исследовании FACT-AUDIT.