

# RuozhiBench: Оценка LLM с помощью логических ошибок и вводящих в заблуждение предпосылок

Дата: 2025-02-18 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.13125>

Рейтинг: 68

Адаптивность: 75

## Ключевые выводы:

Исследование представляет RuozhiBench - двуязычный набор данных из 677 вопросов, содержащих логические ошибки и обманчивые предпосылки, для оценки способности языковых моделей (LLM) распознавать и правильно рассуждать о логических ошибках. Даже лучшая модель (Claude 3 Haiku) достигла только 62% точности по сравнению с человеческим результатом более 90%.

## Объяснение метода:

Исследование предоставляет ценную таксономию логических ошибок и методологию их обнаружения, что помогает пользователям критически оценивать ответы LLM. Категоризация типов обманчивых вопросов и метод парных запросов могут быть адаптированы для повседневного использования. Однако требуется определенная адаптация технической методологии для обычных пользователей.

## Ключевые аспекты исследования: 1. RuozhiBench - это двуязычный набор данных из 677 вопросов, содержащих логические ошибки и вводящие в заблуждение предпосылки, созданный для оценки способности LLM распознавать обманчивый контент.

Исследование включает комплексную оценку 17 LLM с использованием как открытого формата, так и формата с выбором из двух вариантов, показывая ограниченные способности моделей в обнаружении логических ошибок.

Методология исследования включает создание "нормальных" парных вопросов для сравнения производительности моделей на обманчивых и необманчивых входных данных.

Исследователи разработали и сравнили два формата оценки: генеративный (RuozhiBench-Gen) и с множественным выбором (RuozhiBench-MC), выявив их преимущества и ограничения.

Результаты показывают, что даже лучшая модель достигла только 62% точности, что значительно ниже человеческого уровня (более 90%), демонстрируя существенный разрыв в способности моделей обрабатывать обманчивый контент.

## Дополнение:

Исследование RuozhiBench не требует дообучения или специального API для применения его ключевых концепций в стандартном чате. Хотя авторы использовали API для формальной оценки моделей, основные подходы и методы могут быть адаптированы обычными пользователями.

Концепции и подходы, применимые в стандартном чате:

**Таксономия логических ошибок** - пользователи могут научиться распознавать шесть типов ошибок (логические ошибки, ошибки здравого смысла, ошибочные предположения, научные заблуждения, абсурдные фантазии и другие) и использовать это знание для оценки ответов LLM.

**Метод парных вопросов** - пользователи могут формулировать один и тот же запрос разными способами (с потенциальной логической ошибкой и без неё) для проверки надёжности ответов.

**Подход с множественным выбором** - пользователи могут предлагать LLM выбрать из нескольких вариантов ответа, что часто приводит к более надёжным результатам, чем открытая генерация.

**Проверка позиционных предубеждений** - исследование показало, что модели часто предпочитают первый вариант ответа, что можно использовать для проверки надёжности их выбора.

Ожидаемые результаты от применения этих подходов: - Повышенная способность распознавать ненадёжные ответы LLM - Улучшенное качество ответов через более структурированные запросы - Более критический подход к оценке информации, предоставляемой моделями - Возможность проверки логической согласованности ответов без технических знаний

Эти методы не требуют специальных инструментов и могут быть применены в любом стандартном интерфейсе чата с LLM.

## Анализ практической применимости: 1. **Распознавание обманчивого контента:** - Прямая применимость: Пользователи могут использовать понимание типов логических ошибок для более критической оценки ответов LLM в повседневных взаимодействиях. - Концептуальная ценность: Высокая - исследование демонстрирует фундаментальные ограничения современных LLM в логическом мышлении. - Потенциал для адаптации: Пользователи могут научиться формулировать запросы, которые помогут LLM избегать логических ошибок.

**Методология оценки (RuozhiBench-Gen и RuozhiBench-MC):** Прямая применимость: Ограниченная для обычных пользователей, больше подходит для разработчиков и исследователей. Концептуальная ценность: Средняя - представляет различные способы оценки ответов LLM, но требует технических знаний. Потенциал для адаптации: Пользователи могут адаптировать подход с множественным выбором для проверки ответов LLM, предлагая моделям выбрать из альтернатив.

#### **Категоризация типов обманчивых вопросов:**

Прямая применимость: Пользователи могут использовать эту таксономию для идентификации потенциально проблематичных запросов. Концептуальная ценность: Высокая - структурированное понимание различных типов логических ошибок помогает пользователям распознавать их. Потенциал для адаптации: Классификация может быть использована как руководство для формулирования более точных запросов.

#### **Сравнение нормальных и обманчивых вопросов:**

Прямая применимость: Пользователи могут переформулировать запросы для проверки надежности ответов LLM. Концептуальная ценность: Высокая - демонстрирует, как небольшие изменения в формулировке могут значительно влиять на качество ответов. Потенциал для адаптации: Метод парных вопросов может быть адаптирован для проверки надежности ответов на критически важные запросы.

#### **Результаты сравнения различных моделей:**

Прямая применимость: Пользователи могут использовать эту информацию для выбора более надежных моделей для задач, требующих логического мышления. Концептуальная ценность: Средняя - показывает различия в способностях моделей, но быстро устаревает с выпуском новых моделей. Потенциал для адаптации: Методология сравнения может быть адаптирована для оценки новых моделей.

### **Prompt:**

Применение знаний из RuozhiBench в промтах для GPT ## Ключевые выводы из исследования

Исследование RuozhiBench показывает, что даже лучшие языковые модели (включая GPT) имеют ограниченную способность распознавать логические ошибки и обманчивые предпосылки, достигая максимум 62% точности по сравнению с человеческим результатом более 90%.

## Пример промта с учетом результатов исследования

[=====] Проанализируй следующий аргумент на наличие логических ошибок:

[ТЕКСТ АРГУМЕНТА]

Инструкции: 1. Внимательно рассмотри предпосылки и заключение аргумента 2. Определи, есть ли в аргументе логические ошибки (например, ложные дихотомии, круговые рассуждения, ложные предпосылки) 3. Если обнаружишь ошибку, объясни её точную природу 4. Предложи исправленную версию аргумента 5. Оцени аргумент по шкале от 1 до 5, где: - 1: содержит критические логические ошибки - 5: логически безупречен

Формат ответа: - Анализ предпосылок: - Выявленные логические ошибки: - Исправленная версия: - Оценка (1-5):

Важно: Перед ответом тщательно проверь свои рассуждения на наличие логических противоречий. [=====]

## Как работают знания из исследования в данном промпте

**Структурированный подход:** Промпт разбивает задачу на четкие шаги, что помогает модели последовательно анализировать логические конструкции, компенсируя обнаруженную в исследовании слабость моделей в распознавании логических ошибок.

**Явные инструкции:** Исследование показало, что модели нуждаются в явных указаниях для анализа логической структуры, поэтому промпт содержит конкретные инструкции по поиску противоречий.

**Формат множественного выбора:** Шкала оценки от 1 до 5 использует принцип множественного выбора, который, согласно исследованию, повышает точность ответов модели.

**Проверка самоанализа:** Финальное напоминание проверить собственные рассуждения учитывает тенденцию моделей не замечать логические ошибки в своих собственных выводах.

**Структурированный вывод:** Формат ответа с четкими разделами помогает модели систематизировать анализ, что особенно важно для задач с логическим рассуждением, где модели показывают ограниченную эффективность.

Такой подход к составлению промптов позволяет компенсировать выявленные в исследовании RuozhiBench ограничения языковых моделей в области логического рассуждения.