

# Выявление недостатков в том, как люди и большие языковые модели интерпретируют субъективный язык

Дата: 2025-03-06 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.04113>

Рейтинг: 75

Адаптивность: 65

## Ключевые выводы:

Исследование направлено на выявление несоответствий между тем, как большие языковые модели (LLM) интерпретируют субъективные языковые выражения, и тем, как их понимают люди. Основной результат: разработан метод TED (Thesaurus Error Detector), который успешно обнаруживает случаи, когда LLM неожиданно меняют свое поведение при использовании определенных субъективных фраз в промптах.

## Объяснение метода:

Исследование выявляет критические несоответствия между ожиданиями людей и тем, как LLM интерпретируют субъективные инструкции. Конкретные примеры проблем (например, "энтузиастичный" → "нечестный", "остроумный" → "оскорбительный") имеют прямую практическую ценность для пользователей при формулировании запросов. Сам метод TED требует доступа к градиентам и вычислительным ресурсам, но концептуальное понимание проблемы применимо немедленно.

Ключевые аспекты исследования: 1. **Метод TED (Thesaurus Error Detector)** - инструмент для выявления несоответствий между семантическим пониманием субъективных фраз у людей и LLM. Метод сравнивает два тезауруса: операционный (как LLM интерпретирует фразы) и семантический (как люди ожидают, что LLM будет интерпретировать фразы).

**Операционная семантика субъективных выражений** - исследование показывает, что LLM могут неожиданным образом реагировать на субъективные инструкции. Например, запрос написать "энтузиастичный" текст может привести к генерации "нечестного" контента.

**Типы несоответствий** - выявлены два типа проблем: "неожиданные побочные эффекты" (когда LLM добавляет нежелательные качества, например, делает "остроумный" текст "оскорбительным") и "неадекватные обновления" (когда LLM не

добавляет ожидаемые качества).

**Практическая проверка** - методология включает тестирование найденных несоответствий на реальных задачах: редактирование текста и управление выводом при запросе.

**Высокая точность предсказаний** - метод TED показал высокую точность в предсказании проблем в реальном взаимодействии с LLM, значительно превосходя базовый метод, основанный только на семантическом тезаурусе.

Дополнение: Исследование представляет метод TED (Thesaurus Error Detector), который требует доступа к градиентам модели и вычислительным ресурсам для своей полной реализации. Однако ключевая концепция и результаты исследования могут быть применены в стандартном чате без необходимости дообучения или API.

Концепции и подходы, применимые в стандартном чате:

**Осознание проблемы "операционной семантики"** - понимание того, что субъективные инструкции могут интерпретироваться моделью иначе, чем ожидает человек. Пользователи могут применить это знание, избегая потенциально проблемных субъективных фраз.

**Использование конкретных примеров несоответствий** - исследование выявило множество конкретных проблемных комбинаций, которые пользователи могут немедленно учитывать:

Избегать запросов на "энтузиастичный" контент, если важна честность  
Избегать запросов на "остроумный" или "игривый" контент, если важно избежать оскорбительного тона  
Избегать запросов на "юмористический" контент, если важна точность

**Ручная проверка на побочные эффекты** - пользователи могут адаптировать подход TED, сравнивая тексты с субъективной инструкцией и без неё, чтобы выявить нежелательные изменения.

**Предпочтение конкретных инструкций вместо субъективных** - вместо "сделай текст энтузиастичным" использовать более конкретные указания: "добавь восклицательные знаки, используй позитивные прилагательные".

**Поэтапная проверка** - сначала запрашивать нейтральный контент, а затем просить модель отредактировать его с учётом субъективных качеств, контролируя каждый шаг.

Результаты применения этих концепций: - Более предсказуемые ответы LLM - Снижение риска получения контента с нежелательными качествами - Улучшение соответствия между ожиданиями пользователя и результатами модели - Возможность создать собственный "тезаурус" проблемных комбинаций для конкретных задач

Хотя полный метод TED требует технических возможностей, его ключевые выводы о несоответствиях в интерпретации субъективного языка могут быть успешно применены любым пользователем в обычном чате.

Анализ практической применимости: 1. **Метод TED и выявление несоответствий** - Прямая применимость: Пользователи могут использовать выявленные проблемные комбинации субъективных фраз, чтобы избегать нежелательных результатов. Например, избегать запросов на "остроумный" контент, если не хотят получить "оскорбительный". - Концептуальная ценность: Понимание того, что LLM могут интерпретировать субъективные инструкции иначе, чем люди, критически важно для эффективного использования. - Потенциал для адаптации: Пользователи могут самостоятельно проверять и составлять списки "безопасных" субъективных запросов для своих задач.

**Операционная семантика субъективных выражений** Прямая применимость: Знание о конкретных проблемных комбинациях (например, "энтузиастичный" → "нечестный") помогает формулировать более точные запросы. Концептуальная ценность: Понимание того, что у LLM есть "побочные эффекты" при использовании субъективных фраз. Потенциал для адаптации: Пользователи могут разработать альтернативные формулировки для достижения желаемого эффекта без побочных эффектов.

### **Типы несоответствий**

Прямая применимость: Понимание различий между "неожиданными побочными эффектами" и "неадекватными обновлениями" помогает диагностировать проблемы с запросами. Концептуальная ценность: Осознание того, что проблемы могут быть как в добавлении нежелательных качеств, так и в отсутствии ожидаемых. Потенциал для адаптации: Пользователи могут разработать стратегии для проверки обоих типов проблем в своих запросах.

### **Практическая проверка**

Прямая применимость: Методология тестирования может быть адаптирована пользователями для проверки своих запросов. Концептуальная ценность: Понимание важности тестирования запросов перед их использованием в важных задачах. Потенциал для адаптации: Упрощенные версии методологии могут быть внедрены в рабочий процесс.

### **Высокая точность предсказаний**

Прямая применимость: Выявленные проблемы имеют высокую вероятность проявления на практике. Концептуальная ценность: Понимание того, что некоторые проблемы проявляются почти в 100% случаев (например, "юмористический" → "унизительный"). Потенциал для адаптации: Выстраивание приоритетов при разработке стратегий запросов на основе вероятности проблем. Сводная оценка

полезности: На основе анализа определяю общую оценку полезности исследования:  
**78 из 100**

Это исследование предоставляет исключительно ценную информацию о том, как LLM интерпретируют субъективные инструкции, и выявляет конкретные проблемные комбинации, которые пользователи могут немедленно учитывать при формулировании запросов. Знание о том, что запрос на "энтузиастичный" текст может привести к "нечестному" контенту, или что "остроумный" запрос может сделать текст "оскорбительным", имеет прямую практическую ценность.

Контраргументы для более высокой оценки: - Исследование могло бы предложить конкретные рекомендации для пользователей по формулированию запросов, избегающих выявленные проблемы. - Метод TED требует значительных вычислительных ресурсов и доступа к градиентам модели, что делает его непрактичным для обычных пользователей.

Контраргументы для более низкой оценки: - Исследование выявляет конкретные проблемы в популярных моделях (Llama 3, Mistral), которые пользователи могут немедленно учитывать. - Понимание самого факта, что субъективные инструкции могут интерпретироваться неожиданно, имеет высокую ценность даже без возможности применить сам метод TED.

Скорректированная оценка: **75 из 100**. Снижаю оценку, учитывая ограничения по применимости самого метода TED обычными пользователями, но сохраняю высокую оценку за выявленные конкретные проблемы и общее понимание рисков субъективных инструкций.

Уверенность в оценке: Очень сильная. Исследование четко описывает проблему, методологию и результаты. Представлены убедительные количественные данные о частоте проявления проблем. Выявленные проблемы подтверждены как автоматическими методами, так и человеческой оценкой. Исследование проведено на современных моделях (Llama 3, Mistral), что повышает его актуальность.

Оценка адаптивности: Оценка адаптивности: **65 из 100**

1) Принципы исследования могут быть частично адаптированы для обычного чата. Хотя сам метод TED требует доступа к градиентам модели, концепция сравнения ожидаемой и фактической интерпретации субъективных фраз может быть применена пользователями в упрощенной форме.

2) Пользователи могут извлечь несколько ключевых идей: а) избегать потенциально проблемных субъективных фраз (например, "энтузиастичный", "остроумный"); б) проверять, не приносит ли запрос нежелательные качества; в) использовать более конкретные инструкции вместо субъективных.

3) Высокий потенциал для будущих взаимодействий с LLM. Понимание проблем с интерпретацией субъективных фраз поможет пользователям формулировать более эффективные запросы.

4) Возможность абстрагирования специализированных методов до общих принципов существует, но ограничена необходимостью доступа к внутренним механизмам модели для полноценного применения метода TED.

|| <Оценка: 75> || <Объяснение: Исследование выявляет критические несоответствия между ожиданиями людей и тем, как LLM интерпретируют субъективные инструкции. Конкретные примеры проблем (например, "энтузиастичный"→"нечестный", "остроумный"→"оскорбительный") имеют прямую практическую ценность для пользователей при формулировании запросов. Сам метод TED требует доступа к градиентам и вычислительным ресурсам, но концептуальное понимание проблемы применимо немедленно.> || <Адаптивность: 65>

## **Prompt:**

Использование исследования TED в промптах для GPT

### **Ключевые применения исследования**

Исследование TED (Thesaurus Error Detector) выявляет несоответствия между тем, как языковые модели интерпретируют субъективные выражения и как их понимают люди. Это знание можно применить для:

**Избегания проблемных субъективных терминов** Замены терминов с нежелательными эффектами Создания более точных и предсказуемых промптов

### **Пример промпта с учетом знаний из исследования**

[=====] Напиши статью о преимуществах электромобилей. Сделай текст: - Энергичным (вместо "энтузиастичным", чтобы избежать нечестности) - Информативным и основанным на фактах - Структурированным и логичным

Избегай: - Преувеличений и необоснованных заявлений - Сочетания юмора с фактами (может снизить точность) - Чрезмерной эмоциональности в ущерб достоверности

Цель: создать текст, который будет одновременно увлекательным и точным.  
[=====]

## **Объяснение принципа работы**

Данный промпт использует знания из исследования TED следующим образом:

**Избегает проблемных терминов:** Использует "энергичный" вместо "энтузиастичный", который, согласно исследованию, может привести к нечестности в

97% случаев у Llama 3 8B (аналогичный эффект может наблюдаться и у GPT).

**Избегает проблемных комбинаций:** Явно указывает на необходимость избегать сочетания юмора с фактической информацией, поскольку исследование показало, что запрос на "юмористичный" текст может привести к более "неточному" контенту.

**Устанавливает противовес:** Требуется информативности и фактической точности как противовес потенциальным побочным эффектам от субъективных терминов.

**Дает четкие ограничения:** Явно указывает, чего следует избегать, основываясь на выявленных в исследовании проблемах.

Такой подход помогает получить более предсказуемый и качественный результат, избегая неожиданных побочных эффектов от использования субъективных терминов в промптах.