

# Гендерные предвзятости в LLM: Более высокая *inteligencia* в LLM не обязательно решает проблемы гендерной предвзятости и стереотипов

Дата: 2025-02-15 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2409.19959>

Рейтинг: 75

Адаптивность: 80

## Ключевые выводы:

Исследование направлено на изучение гендерных предубеждений в больших языковых моделях (LLM), особенно на проверку гипотезы о том, снижает ли более высокий интеллект LLM такие предубеждения. Основной вывод: даже в более интеллектуальных моделях (o1-mini по сравнению с 4o-mini) систематические гендерные предубеждения сохраняются. Модель o1 оценивала мужчин выше по компетентности (8.1) по сравнению с женщинами (7.9) и небинарными персонами (7.80), а также демонстрировала стереотипное распределение по профессиональным областям.

## Объяснение метода:

Исследование высоко полезно для широкой аудитории, предлагая методологию выявления гендерных предубеждений в LLM, которую могут применять обычные пользователи. Оно разрушает миф о "самоисправлении" предвзятости с ростом интеллекта моделей и дает конкретные инструменты для критической оценки ответов LLM, что повышает цифровую грамотность.

## Ключевые аспекты исследования: 1. **Методология оценки гендерных стереотипов в LLM:** Авторы разработали методологию с использованием персон и гендерно-нейтральных имен для оценки гендерных предубеждений в языковых моделях.

**Сравнительный анализ моделей разной "интеллектуальности":** Исследование сравнивает две модели OpenAI (4o и o1), чтобы выяснить, уменьшаются ли гендерные предубеждения с повышением "интеллекта" модели.

**Измерение конкретных параметров смещения:** Авторы оценивают смещения по нескольким ключевым параметрам - оценка компетентности, вероятность стать успешным основателем бизнеса или CEO, а также анализ личностных черт и

предпочтений.

**Выявленные паттерны стереотипов:** Исследование обнаружило устойчивые стереотипы в представлении различных гендеров в профессиональных областях (мужчины доминируют в технических областях, женщины - в творческих).

**Рекомендации по снижению предубеждений:** Авторы предлагают конкретные подходы к смягчению гендерных предубеждений в LLM, включая балансировку данных, алгоритмические методы и создание "слоя справедливости".

## Дополнение:

### Применимость методов исследования в стандартном чате

Методы данного исследования не требуют дообучения или API для применения обычными пользователями. Большинство подходов можно адаптировать для стандартного чата:

**Использование гендерно-нейтральных имен** - пользователи могут формулировать запросы с гендерно-нейтральными именами (например, "Алекс", "Саша") и анализировать, какой гендер модель присваивает персонажу по умолчанию.

**Проверка конкретных параметров предвзятости** - пользователи могут проверять, как модель оценивает компетентность, лидерские качества или вероятность успеха для разных гендеров.

**Сравнительные запросы** - можно задавать одинаковые вопросы, меняя только гендер персонажа, чтобы выявить различия в ответах.

**Анализ стереотипных паттернов** - после понимания типичных стереотипов (мужчины в технических областях, женщины в творческих), пользователи могут формулировать запросы, намеренно противоречащие этим стереотипам.

**Явное указание на необходимость гендерной нейтральности** - пользователи могут включать в запросы инструкции по предоставлению гендерно-сбалансированных примеров.

Применение этих подходов позволит: - Повысить критическое мышление при оценке ответов LLM - Получать более сбалансированные и менее стереотипные ответы - Лучше понимать ограничения моделей - Формулировать запросы, минимизирующие влияние предвзятости

## Анализ практической применимости: **1. Методология оценки гендерных стереотипов в LLM - Прямая применимость:** Высокая. Пользователи могут адаптировать этот подход для проверки гендерных предубеждений в ответах моделей, используя гендерно-нейтральные имена и сравнивая ответы. - **Концептуальная ценность:** Очень высокая. Понимание того, что модели могут

навязывать гендерные стереотипы даже при использовании нейтральных запросов, помогает пользователям критически оценивать ответы LLM. - **Потенциал для адаптации:** Высокий. Методика с гендерно-нейтральными именами может быть упрощена для повседневного использования и применена к анализу любых ответов LLM на предмет предубеждений.

**2. Сравнительный анализ моделей разной "интеллектуальности" - Прямая применимость:** Средняя. Пользователи получают понимание, что более новые или "умные" модели не обязательно менее предвзяты, что влияет на выбор модели. - **Концептуальная ценность:** Высокая. Разрушает миф о том, что повышение интеллекта модели автоматически решает проблемы предвзятости. - **Потенциал для адаптации:** Средний. Позволяет пользователям формировать более реалистичные ожидания от новых версий моделей.

**3. Измерение конкретных параметров смещения - Прямая применимость:** Высокая. Пользователи могут обращать внимание на конкретные аспекты предвзятости (компетентность, лидерские качества) при взаимодействии с LLM. - **Концептуальная ценность:** Высокая. Помогает понять, что предвзятость проявляется в конкретных, измеримых аспектах ответов. - **Потенциал для адаптации:** Высокий. Пользователи могут разработать собственные "проверки" на предвзятость в повседневных взаимодействиях с LLM.

**4. Выявленные паттерны стереотипов - Прямая применимость:** Средняя. Знание о конкретных стереотипах помогает пользователям распознавать их в ответах LLM. - **Концептуальная ценность:** Высокая. Понимание того, как именно проявляются гендерные стереотипы, позволяет более критично оценивать информацию. - **Потенциал для адаптации:** Высокий. Пользователи могут формулировать запросы таким образом, чтобы минимизировать влияние этих стереотипов.

**5. Рекомендации по снижению предубеждений - Прямая применимость:** Средняя. Некоторые рекомендации (например, переформулирование запросов) могут быть непосредственно применены пользователями. - **Концептуальная ценность:** Высокая. Понимание механизмов снижения предвзятости дает пользователям инструменты для критического взаимодействия с LLM. - **Потенциал для адаптации:** Высокий. Пользователи могут адаптировать предложенные подходы для улучшения собственного опыта взаимодействия с LLM.

## **Prompt:**

Использование знаний о гендерных предвзятостях в LLM для создания лучших промптов Исследование о гендерных предвзятостях в языковых моделях предоставляет ценную информацию, которую можно применить при составлении промптов для минимизации нежелательных стереотипов в ответах.

## Ключевые выводы для составления промптов

**Даже продвинутые LLM** имеют устойчивые гендерные предубеждения  
**Стереотипное распределение** по профессиональным областям сохраняется  
**Оценка компетентности** систематически различается по гендерному признаку  
**Требуется явное указание** на необходимость гендерного баланса в ответах ##  
Пример промпта с учетом исследования

[=====] Создай список из 10 выдающихся специалистов в области инженерии и технологий для панельной дискуссии на конференции.

Важные требования: - Обеспечь сбалансированное гендерное представительство (равное количество мужчин, женщин и, по возможности, небинарных персон) - Избегай стереотипного распределения по узким специализациям внутри технической сферы - При описании профессиональных достижений используй одинаково объективные критерии для всех гендеров - Оценивай компетентность всех специалистов по единой шкале, не допуская систематического занижения для не-мужчин - Представь разнообразие происхождения, возрастов и опыта

Для каждого специалиста укажи: 1. Имя и гендер 2. Область специализации 3. Ключевые достижения 4. Потенциальный вклад в дискуссию [=====]

## Как работает этот подход

Данный промпт использует знания из исследования следующим образом:

- Явное требование баланса — противодействует выявленной тенденции моделей создавать больше мужских персон в технических областях
- Запрет на стереотипное распределение — предотвращает автоматическое помещение женщин в "творческие" роли, а мужчин в "технические"
- Единые критерии оценки — борется с обнаруженной тенденцией оценивать компетентность женщин и небинарных персон ниже (8.1 для мужчин против 7.9/7.8 для других)
- Контроль предвзятости — создаёт "слой метасправедливости", заставляя модель проверять свои ответы на наличие предвзятости

Такой подход помогает получить более сбалансированные результаты, даже несмотря на встроенные предвзятости модели.