

# FB-Bench: Тонкий многозадачный бенчмарк для оценки отклика LLM на человеческую обратную связь

Дата: 2025-02-16 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2410.09412>

Рейтинг: 65

Адаптивность: 75

## Ключевые выводы:

Исследование представляет FB-Bench - новый многозадачный бенчмарк для оценки отзывчивости больших языковых моделей (LLM) на обратную связь от пользователей в реальных сценариях использования на китайском языке. Основные результаты показывают, что разрыв в производительности между открытыми и закрытыми LLM сокращается, большинство моделей демонстрируют сбалансированную способность исправлять ошибки и поддерживать ответы, а открытые модели показывают превосходные возможности поддержания ответов.

## Объяснение метода:

Исследование предлагает ценную таксономию типов обратной связи и анализ их влияния на ответы LLM. Пользователи могут применять эти знания для более эффективного взаимодействия, особенно используя подсказки и руководство. Однако техническая направленность и китайский язык исследования требуют некоторой адаптации для широкого применения.

## Ключевые аспекты исследования: 1. **Создание FB-Bench** - разработка детального многозадачного бенчмарка для оценки отзывчивости языковых моделей (LLM) на обратную связь от пользователей в реальных сценариях использования на китайском языке.

**Трехуровневая иерархическая таксономия** - классификация взаимодействий человека и LLM по трем компонентам: запросы пользователей (8 типов задач), ответы модели (5 типов недостатков) и обратная связь от пользователей (9 типов).

**Два ключевых сценария взаимодействия** - исследование фокусируется на исправлении ошибок и сохранении ответа как основных сценариях взаимодействия человека с LLM.

**Метод оценки на основе чек-листа** - разработка взвешенного чек-листа для детальной оценки каждого образца, с использованием GPT-4o в качестве судьи.

**Выявление факторов, влияющих на отзывчивость моделей** - анализ того, как типы задач, типы обратной связи и недостатки предыдущих ответов влияют на способность моделей реагировать на обратную связь.

## Дополнение:

Исследование FB-Bench не требует дообучения или API для практического применения его основных концепций обычными пользователями. Хотя сами исследователи использовали технически сложные методы для создания бенчмарка и оценки моделей, ключевые выводы и подходы могут быть адаптированы для стандартного чата без дополнительных инструментов.

**Концепции и подходы, применимые в стандартном чате:**

**Типы эффективной обратной связи:** Использование "подсказок и руководства" (hinting guidance) значительно улучшает качество ответов LLM "Указание на ошибки" (pointing out errors) помогает моделям исправлять неточности "Разъяснение намерений" (clarifying intent) улучшает релевантность ответов

**Понимание двух сценариев взаимодействия:**

В сценарии "исправления ошибок" важно четко указывать на ошибки и направлять модель В сценарии "сохранения ответа" следует избегать предоставления дезинформации или необоснованных претензий

**Учет типов задач:**

Для сложных математических и логических задач может потребоваться более детальная обратная связь Для задач создания и перевода текста эффективны разные типы обратной связи **Ожидаемые результаты от применения:** - Более точные и релевантные ответы от LLM - Сокращение количества итераций для получения желаемого результата - Лучшее понимание, как формулировать эффективную обратную связь в различных контекстах - Повышение способности направлять модель к желаемому результату без необходимости в технических знаниях

## Анализ практической применимости: 1. **Создание FB-Bench** - Прямая применимость: Ограниченная для обычных пользователей, поскольку это инструмент для исследователей и разработчиков. - Концептуальная ценность: Высокая, демонстрирует важность учета различных типов обратной связи при взаимодействии с LLM. - Потенциал для адаптации: Пользователи могут адаптировать понимание различных типов обратной связи для более эффективного взаимодействия с LLM.

**Трехуровневая иерархическая таксономия** Прямая применимость: Средняя, пользователи могут использовать понимание различных типов задач и обратной связи для формирования более эффективных запросов. Концептуальная ценность:

Высокая, помогает понять, как структурировать взаимодействие с LLM. Потенциал для адаптации: Пользователи могут адаптировать таксономию для своих нужд, понимая, какие типы обратной связи более эффективны.

### **Два ключевых сценария взаимодействия**

Прямая применимость: Высокая, пользователи могут осознанно применять различные стратегии при исправлении ошибок и подтверждении правильности ответов. Концептуальная ценность: Значительная для понимания того, как модели обрабатывают различные ситуации. Потенциал для адаптации: Пользователи могут разработать персональные стратегии взаимодействия для каждого сценария.

### **Метод оценки на основе чек-листа**

Прямая применимость: Низкая для обычных пользователей, в основном полезна для исследователей. Концептуальная ценность: Средняя, демонстрирует методы оценки ответов LLM. Потенциал для адаптации: Пользователи могут создавать мысленные чек-листы при оценке ответов LLM.

### **Выявление факторов, влияющих на отзывчивость моделей**

Прямая применимость: Высокая, пользователи могут применять наиболее эффективные типы обратной связи. Концептуальная ценность: Значительная, помогает понять сильные и слабые стороны LLM. Потенциал для адаптации: Пользователи могут адаптировать свои стратегии обратной связи в зависимости от типа задачи и модели.

## **Prompt:**

Использование знаний из FB-Bench в промтах для GPT ## Ключевые применения исследования

Исследование FB-Bench предоставляет ценные инсайты о том, как большие языковые модели реагируют на обратную связь пользователей. Эти знания можно эффективно применить при составлении промтов для улучшения взаимодействия с GPT.

## Пример промпта с использованием находок исследования

[=====] Я хочу, чтобы ты помог мне решить математическую задачу по оптимизации.

Следуя выводам исследования FB-Bench, я буду структурировать свой запрос так:

Вот полная формулировка задачи: [описание задачи] Мне нужно пошаговое решение с пояснениями каждого этапа. Если я замечу ошибку, я укажу на конкретный шаг и предоставлю дополнительную информацию. Пожалуйста, сохраняй правильные части решения при внесении исправлений. При необходимости я могу предоставить

дополнительные подсказки для направления решения Начни, пожалуйста, с анализа условий задачи и определения метода решения. [=====]

## Объяснение применения знаний из исследования

В этом промпте использованы следующие находки из FB-Bench:

**Структурированные подсказки и руководства** - исследование показало, что все модели достигают оценок выше 80% при получении конкретных подсказок

**Фокус на математической задаче** - учитывая, что модели показывают более низкую производительность в математических областях, промпт предусматривает пошаговое решение

**Подготовка к предоставлению обратной связи** - заранее указано, что будут даваться конкретные указания на ошибки, что по данным исследования помогает моделям эффективнее корректировать ответы

**Сохранение правильных частей** - исследование показало, что открытые модели лучше справляются с поддержанием верных частей ответов, поэтому промпт явно запрашивает это

**Готовность предоставить дополнительные подсказки** - учитывая, что модели лучше адаптируются при получении разъяснений намерений пользователя

Такой подход к составлению промтов, основанный на эмпирических данных FB-Bench, повышает вероятность получения качественных и точных ответов от GPT.