

# Соединение исследований HCI и ИИ для оценки разговорных помощников в области программной инженерии

Дата: 2025-02-11 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.07956>

Рейтинг: 65

Адаптивность: 75

## Ключевые выводы:

Исследование направлено на разработку методов автоматической оценки LLM-ассистентов для разработки программного обеспечения (SE) с учетом человеческого фактора. Авторы предлагают объединить подходы из областей взаимодействия человека с компьютером (HCI) и искусственного интеллекта (AI) для создания комплексной системы оценки, которая сочетает симулированных пользователей и подход 'LLM as a judge'.

## Объяснение метода:

Исследование предлагает ценные концепции, которые могут быть адаптированы для повседневного использования LLM: "LLM как судья" для оценки ответов, учет разнообразия пользователей, многоходовые взаимодействия и критическое отношение к "эталонным" ответам. Хотя полная реализация методологии требует технических навыков, общие принципы доступны широкой аудитории.

**## Ключевые аспекты исследования:** 1. **Интеграция методов HCI и AI для оценки LLM-ассистентов:** Исследование предлагает объединить подходы из областей взаимодействия человека с компьютером (HCI) и искусственного интеллекта (AI) для автоматической оценки разговорных LLM-ассистентов в сфере разработки ПО.

**Симулированные пользователи:** Предлагается использовать LLM для создания симулированных пользователей, которые могут реалистично взаимодействовать с ассистентом, генерировать качественную обратную связь и выявлять проблемы с инклюзивностью.

**LLM как судья:** Подход, при котором LLM используется для оценки ответов других LLM по заданным критериям, что позволяет получать количественные метрики без необходимости проведения дорогостоящих исследований с участием людей.

**Персоны и разнообразие:** Важность создания репрезентативных персон для

симуляции разнообразных пользователей, что помогает выявлять "баги инклюзивности" — проблемы, которые возникают только у определенных групп пользователей.

**Ограничения существующих методов оценки:** Критика традиционных методов оценки LLM-ассистентов, основанных на сравнении с эталонными ответами, которые не отражают разнообразие возможных действительных ответов и не учитывают многоходовый характер реальных взаимодействий.

## Дополнение: Исследование не требует дообучения или специального API для применения основных концепций в стандартном чате. Хотя авторы используют более продвинутые технические подходы для систематической оценки, ключевые идеи могут быть адаптированы для использования в обычном диалоге с LLM.

Вот концепции, которые можно применить в стандартном чате:

**LLM как судья:** Пользователь может попросить LLM оценить свой предыдущий ответ по определенным критериям или сравнить несколько подходов к решению задачи. Например: "Я получил от тебя два решения моей задачи программирования. Оцени их по критериям эффективности, читаемости и следования лучшим практикам. Какое решение лучше и почему?"

**Персонализация запросов:** Пользователь может указать свой уровень опыта или предпочтительный стиль объяснения. Например: "Объясни мне, как работает рекурсия, как если бы я был начинающим программистом, который только изучает основы."

**Многоходовое взаимодействие:** Вместо попыток получить исчерпывающий ответ в одном запросе, пользователь может вести инкрементальный диалог, уточняя детали и постепенно двигаясь к решению. Это соответствует естественному процессу человеческого общения.

**Разнообразие перспектив:** Пользователь может запросить альтернативные точки зрения на проблему. Например: "Как бы эту задачу решил опытный разработчик Python? А как бы к ней подошел специалист по JavaScript?"

Применяя эти концепции, пользователь может получить: - Более точные и релевантные ответы, адаптированные к своему уровню подготовки - Многогранное понимание проблемы через различные перспективы - Более критическое отношение к ответам LLM - Улучшенное поэтапное решение сложных задач

Эти подходы не требуют технической реализации и могут быть использованы любым пользователем в рамках обычного диалогового интерфейса.

## Анализ практической применимости: **Аспект 1: Интеграция методов HCI и AI - Прямая применимость:** Средняя. Пользователи не могут напрямую реализовать эту интеграцию, но могут использовать идею о сочетании качественной и количественной оценки при взаимодействии с LLM. - **Концептуальная ценность:**

Высокая. Понимание необходимости оценивать LLM не только по точности ответов, но и по удобству взаимодействия помогает пользователям формировать более эффективные запросы. - **Потенциал для адаптации:** Высокий. Пользователи могут адаптировать идею многокритериальной оценки для собственных нужд, например, оценивая ответы LLM по нескольким параметрам одновременно.

**Аспект 2: Симулированные пользователи - Прямая применимость:** Низкая. Обычные пользователи вряд ли будут создавать симулированных пользователей для оценки взаимодействия с LLM. - **Концептуальная ценность:** Средняя. Понимание, что LLM может симулировать различные стили взаимодействия, может помочь пользователям лучше понять, как формулировать запросы. - **Потенциал для адаптации:** Средний. Пользователи могут попросить LLM симулировать разные подходы к решению задачи, что поможет получить более разнообразные ответы.

**Аспект 3: LLM как судья - Прямая применимость:** Высокая. Пользователи могут напрямую попросить LLM оценить другой ответ LLM или сравнить несколько ответов. - **Концептуальная ценность:** Высокая. Понимание, что LLM может критически анализировать собственные ответы, помогает пользователям формулировать запросы на оценку и улучшение ответов. - **Потенциал для адаптации:** Высокий. Пользователи могут разработать собственные критерии оценки и попросить LLM оценить ответы по этим критериям.

**Аспект 4: Персоны и разнообразие - Прямая применимость:** Средняя. Пользователи могут создавать запросы с учетом различных персон и сценариев использования. - **Концептуальная ценность:** Высокая. Понимание, что LLM может по-разному отвечать на запросы разных пользователей, помогает более эффективно формулировать запросы. - **Потенциал для адаптации:** Высокий. Пользователи могут адаптировать концепцию персон для получения более разнообразных и инклюзивных ответов.

**Аспект 5: Ограничения существующих методов оценки - Прямая применимость:** Средняя. Понимание ограничений помогает пользователям формировать более реалистичные ожидания. - **Концептуальная ценность:** Высокая. Осознание того, что один "правильный ответ" часто не существует, помогает пользователям лучше интерпретировать и использовать ответы LLM. - **Потенциал для адаптации:** Средний. Пользователи могут адаптировать критический подход к оценке ответов LLM в своих взаимодействиях.

## **Prompt:**

Использование знаний из исследования в промптах для GPT ## Ключевые инсайты исследования для промптов

Исследование предлагает комбинированный подход к оценке LLM-ассистентов, объединяющий симулированных пользователей и LLM-судей. Эти знания можно применить для создания более эффективных промптов.

## ## Пример промпта с использованием методологии исследования

[=====] Я разрабатываю помощника на базе GPT для junior-разработчиков, который помогает с задачами по JavaScript. Помоги мне улучшить мой промпт, используя следующий подход:

Создай 3 различные персоны пользователей (начинающий, средний уровень, с опытом в других языках) с разными потребностями и стилями взаимодействия.

Для каждой персоны смоделируй диалог с моим ассистентом, используя следующий базовый промпт: "Ты помощник по JavaScript. Твоя задача - объяснять концепции, помогать с отладкой и предлагать решения. Старайся давать понятные объяснения с примерами кода."

Оцени эффективность этого промпта по следующим критериям:

Точность предоставляемой информации (1-10) Понятность объяснений для конкретной персоны (1-10) Полезность примеров кода (1-10) Способность адаптироваться к уровню пользователя (1-10)

Предложи улучшения промпта, основываясь на результатах симуляции, добавив:

Конкретные инструкции по адаптации к разным уровням пользователей Структуру для предоставления ответов Стратегии для более эффективного объяснения сложных концепций [=====] ## Как работают знания из исследования в этом промпте

**Репрезентативные персоны** — промпт использует идею создания различных профилей пользователей для обеспечения инклюзивности и учета разнообразия аудитории.

**Симуляция диалогов** — применяется метод генерации взаимодействий между ассистентом и симулированным пользователем для тестирования эффективности промпта.

**Количественная оценка** — внедрена система оценки по конкретным метрикам (по шкале 1-10), что соответствует подходу "LLM as a judge" из исследования.

**Качественная обратная связь** — запрашиваются конкретные улучшения на основе выявленных проблем, что позволяет получить качественные выводы.

Такой подход позволяет итеративно улучшать промпты, основываясь на симулированном пользовательском опыте и структурированной оценке, что делает разработку более эффективной, не требуя постоянных реальных пользовательских тестов.