

Оценка персонализированных инструментов с поддержкой больших языковых моделей с точки зрения персонализации и проактивности

Дата: 2025-03-02 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.00771>

Рейтинг: 68

Адаптивность: 75

Ключевые выводы:

Исследование направлено на оценку персонализированных LLM-агентов, использующих инструменты, с точки зрения персонализации и проактивности. Авторы разработали новый бенчмарк ETAPP (Evaluation of Tool-augmented Agent from the Personalization and Proactivity Perspective) для оценки способности LLM использовать инструменты с учетом предпочтений пользователя.

Объяснение метода:

Исследование предлагает ценные концепции персонализации и проактивности, применимые при формулировке запросов к LLM. Метод E-ReAct и структура предпочтений пользователя могут быть адаптированы для повседневного использования. Однако многие технические аспекты (песочница, методы оценки) недоступны обычным пользователям без специальных навыков.

Ключевые аспекты исследования: 1. **Фреймворк ETAPP** - Авторы разработали новый бенчмарк для оценки персонализированного вызова инструментов (API) в языковых моделях с двух ключевых перспектив: персонализации и проактивности.

Архитектура памяти и предпочтений пользователя - Исследователи предложили структуру для хранения предпочтений пользователя, разделенную на долгосрочную (профиль пользователя и предпочтения инструментов) и краткосрочную память (текущее состояние пользователя).

Метод оценки на основе ключевых точек - Разработан подход, использующий заранее аннотированные ключевые точки для более точной оценки LLM в задачах персонализации, что значительно повышает согласованность с оценкой человека.

Анализ методов вызова инструментов - Исследование сравнивает различные методы вызова инструментов (Function Calling, ReAct, E-ReAct) и показывает, как

интеграция рассуждений перед вызовом инструментов улучшает персонализацию и проактивность.

Эксперименты по дообучению - Проведены эксперименты, демонстрирующие, как дообучение улучшает способность модели использовать инструменты, но имеет ограниченный эффект для новых типов инструкций.

Дополнение: Исследование не требует обязательного дообучения или специального API для применения его основных концепций. Хотя авторы использовали специальную среду и дообучение в своих экспериментах, ключевые идеи и подходы могут быть адаптированы для использования в стандартном чате с LLM.

Концепции, которые можно применить в стандартном чате:

Структура персонализации: Разделение информации о пользователе на долгосрочную (базовый профиль, общие предпочтения) и краткосрочную память (текущее состояние, контекст). Пользователи могут структурировать свои запросы, включая эту информацию в начале диалога или при необходимости.

Метод E-ReAct: Пользователи могут просить модель "подумать вслух" о персонализации и проактивности перед предоставлением ответа. Например: "Перед ответом проанализируй мои предпочтения и подумай, как ты можешь сделать ответ персонализированным и предвосхитить мои дополнительные потребности".

Критерии персонализации и проактивности: Пользователи могут явно указывать эти критерии в своих запросах, например: "Учти мои предпочтения в еде (перечисление) и предложи дополнительные варианты, которые могут мне понравиться, даже если я о них не спрашивал".

Структурированные профили предпочтений: Пользователи могут создавать и сохранять структурированные профили своих предпочтений по различным категориям (еда, музыка, путешествия и т.д.), которые можно включать в релевантные запросы.

Ожидаемые результаты от применения этих концепций: - Более персонализированные ответы, учитывающие предпочтения пользователя - Проактивные предложения, выходящие за рамки явного запроса - Лучший пользовательский опыт благодаря более глубокому пониманию моделью контекста и предпочтений - Более эффективное использование LLM для решения повседневных задач

Применение этих концепций не требует технических навыков и доступно любому пользователю стандартного чата с LLM.

Анализ практической применимости: 1. **Фреймворк ETAPP - Прямая применимость:** Средняя. Обычные пользователи не могут напрямую использовать бенчмарк, но могут применять критерии персонализации и проактивности при

формулировке запросов. - **Концептуальная ценность:** Высокая. Понимание двух ключевых аспектов персонализированных агентов (персонализация и проактивность) может помочь пользователям более эффективно формулировать запросы. - **Потенциал для адаптации:** Высокий. Критерии оценки могут быть превращены в рекомендации для пользователей о том, как получать более персонализированные ответы.

Архитектура памяти и предпочтений **Прямая применимость:** Средняя. Пользователи могут структурировать свои запросы, включая релевантную информацию о предпочтениях и текущем состоянии. **Концептуальная ценность:** Высокая. Понимание разделения на долгосрочные предпочтения и текущее состояние помогает пользователям давать более релевантный контекст. **Потенциал для адаптации:** Высокий. Можно создать шаблоны запросов, которые учитывают эту структуру.

Метод оценки на основе ключевых точек

Прямая применимость: Низкая. Это технический метод оценки, сложный для непосредственного использования обычными пользователями. **Концептуальная ценность:** Средняя. Понимание ключевых аспектов качественных ответов может помочь пользователям оценивать ответы LLM. **Потенциал для адаптации:** Средний. Можно преобразовать ключевые точки в чеклисты для оценки ответов.

Анализ методов вызова инструментов

Прямая применимость: Высокая. Пользователи могут запрашивать у модели "подумать вслух" перед выполнением действий, что улучшает персонализацию. **Концептуальная ценность:** Высокая. Понимание, что рассуждение перед действием улучшает результаты, может изменить подход к формулировке запросов. **Потенциал для адаптации:** Высокий. Легко адаптируется в простые промпты, стимулирующие размышление модели.

Эксперименты по дообучению

Прямая применимость: Низкая. Большинство пользователей не могут проводить дообучение моделей. **Концептуальная ценность:** Средняя. Понимание ограничений дообучения помогает иметь реалистичные ожидания. **Потенциал для адаптации:** Низкий. Требуются технические навыки и ресурсы для применения.

Prompt:

Использование знаний из исследования ETAPP в промтах для GPT ## Ключевые знания из исследования

Исследование ETAPP показывает, что: 1. Метод ReAct и Enhanced-ReAct превосходит простой Function Calling 2. Разделение предпочтений пользователя на высокоуровневые и низкоуровневые улучшает персонализацию 3. Включение этапа рассуждения перед вызовом инструментов повышает эффективность 4.

Проактивность требует предвосхищения неявных потребностей пользователя

Пример промта с применением этих знаний

[=====] # Запрос для персонализированного помощника по планированию питания

Контекст пользователя - Высокоуровневый профиль: женщина, 35 лет, работает офисным менеджером, вегетарианка, занимается йогой 3 раза в неделю - Низкоуровневые предпочтения: предпочитает органические продукты, избегает глютен, ограничивает потребление сахара, любит азиатскую кухню

Инструкции для модели 1. **Этап рассуждения (Enhanced-ReAct):** - Проанализируй профиль пользователя и определи ключевые диетические потребности - Учти недавнюю активность пользователя (последняя йога-сессия была вчера - высокая интенсивность) - Определи, какие инструменты потребуются для составления плана питания

Персонализация: Используй только релевантные для запроса предпочтения Адаптируй рецепты под вегетарианскую диету и ограничения по глютену Учитывай предпочтение азиатской кухни при выборе рецептов

Проактивность:

Предложи продукты, богатые белком, учитывая вчерашнюю интенсивную тренировку Проверь сезонность предлагаемых ингредиентов Предложи варианты для разных бюджетов без явного запроса об этом ## Задача Составь план питания на 3 дня, включая завтрак, обед и ужин, который подойдет пользователю. [=====]

Объяснение применения знаний из исследования

В этом промте я применил следующие принципы из исследования ETAPP:

Структура Enhanced-ReAct: Промт включает явный этап рассуждения перед действием, что согласно исследованию повышает качество персонализации и проактивности

Разделение предпочтений: Предпочтения разделены на высокоуровневые (общий профиль) и низкоуровневые (конкретные пищевые предпочтения), что снижает когнитивную нагрузку на модель

Явное указание на проактивность: Промт побуждает модель предвосхищать потребности пользователя (например, повышенная потребность в белке после тренировки), а не просто отвечать на явный запрос

Персонализация с учетом контекста: Промт направляет модель на использование только релевантных предпочтений для конкретного запроса

Такой подход, согласно исследованию, должен обеспечить более высокое качество

персонализированного ответа по сравнению с простым запросом без структуры Enhanced-ReAct.