

Области согласования

Дата: 2025-01-14 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2501.12405>

Рейтинг: 58

Адаптивность: 70

Ключевые выводы:

Основная цель исследования - предложить более широкую концепцию выравнивания (alignment) LLM, выходящую за рамки общепринятого подхода, ориентированного на универсальные ценности (полезность, безвредность, честность). Авторы предлагают трехмерную структуру для более точного определения различных областей выравнивания LLM, учитывающую компетенции, временные рамки и аудиторию.

Объяснение метода:

Исследование предлагает ценную концептуальную рамку для понимания ограничений универсального выравнивания LLM и необходимости учёта контекста. Оно помогает пользователям осознать культурные предубеждения моделей и формулировать более эффективные запросы. Однако исследование ограничено в плане конкретных техник, которые пользователи могли бы непосредственно применить без дополнительных знаний.

Ключевые аспекты исследования: 1. **Концепция трех измерений выравнивания (alignment) LLM:** авторы предлагают рассматривать выравнивание моделей в трех измерениях: компетенция (знания, навыки, поведение), временность (семантическая или эпизодическая) и аудитория (массовая, групповая, диадическая).

Критика ограниченного подхода к выравниванию: исследование показывает, что текущий подход к выравниванию LLM сосредоточен на универсальных ценностях (полезность, безвредность, честность), но игнорирует культурные различия и контекстуальные потребности.

Контекстуальное выравнивание: авторы утверждают, что выравнивание должно быть адаптировано к конкретным потребностям и контекстам использования, а не стремиться к единому универсальному набору ценностей.

Примеры практического применения: исследование предлагает пример использования различных подходов к выравниванию для систем поддержки психического здоровья в США и Китае, где культурные нормы и нормативные среды существенно различаются.

Предварительный шаг к плюралистическому выравниванию: авторы рассматривают свою работу как предшественника плюралистического выравнивания, который помогает избежать конфликтов ценностей путем правильного определения области применения.

Дополнение: Исследование не требует дообучения или API для применения основных концепций. Хотя авторы обсуждают технические методы выравнивания, требующие доступа к параметрам модели (полная настройка, эффективная настройка параметров), основные концептуальные идеи могут быть адаптированы для использования в стандартном чате.

Концепции, которые можно применить в стандартном чате:

Трехмерный подход к выравниванию: Пользователи могут адаптировать свои промпты, учитывая: **Компетенция:** Явно указывать, какие знания, навыки или поведение требуются от модели **Временность:** Определять, нужен ли общий (семантический) или контекстуально-специфичный (эпизодический) ответ **Аудитория:** Уточнять, для кого предназначен ответ (для себя, малой группы или широкой аудитории)

Диадическое взаимодействие: Пользователи могут развивать "взаимную теорию разума" с моделью, адаптируясь к её способностям и помогая модели лучше понимать их потребности.

Контекстуальное выравнивание: Пользователи могут создавать системные промпты, учитывающие культурный и профессиональный контекст, в котором будет использоваться ответ.

Практическое применение этих концепций может привести к: - Более точным и релевантным ответам, соответствующим конкретным потребностям пользователя - Снижению культурных предубеждений в ответах LLM - Более эффективному взаимодействию с моделью через постепенную адаптацию запросов

Эти подходы не требуют технических изменений в модели, а основаны на стратегическом формулировании запросов, учитывающем многомерную природу выравнивания.

Анализ практической применимости: ## Концепция трех измерений выравнивания - **Прямая применимость:** Средняя. Обычные пользователи не могут напрямую изменить базовое выравнивание модели, но могут использовать эту концепцию для лучшего понимания ограничений модели и более эффективного формулирования запросов. - **Концептуальная ценность:** Высокая. Помогает пользователям понять, что LLM не являются нейтральными инструментами, а содержат определенные ценности и предубеждения, которые могут не соответствовать их собственным. - **Потенциал для адаптации:** Высокий. Пользователи могут применять эту концепцию для разработки системных промптов, учитывающих нужные им компетенции,

временной контекст и аудиторию.

Prompt:

Применение исследования "Области согласования" в промптах для GPT ##
Ключевые аспекты исследования для использования в промптах

Исследование предлагает трехмерную структуру для выравнивания LLM: 1. **Компетенция** (знания, навыки, поведение) 2. **Временные рамки** (эпизодические или семантические) 3. **Аудитория** (от диадической до массовой)

Пример промпта с применением этих знаний

[=====] # Запрос на финансовую консультацию

Желаемая компетенция Я ищу сочетание фактических знаний о финансовых рынках и практических навыков по составлению инвестиционного портфеля. Пожалуйста, воздержись от демонстрации рискованного финансового поведения.

Временные рамки - Эпизодический контекст: Я 35-летний IT-специалист из России, планирующий долгосрочные инвестиции в условиях текущей геополитической ситуации (2023 год) - Требуется сочетание общих (семантических) принципов инвестирования с учетом конкретных эпизодических обстоятельств

Аудитория Это диадическое взаимодействие, адаптированное под мои индивидуальные обстоятельства, а не общие рекомендации для массовой аудитории.

С учетом этих параметров, помоги мне составить стратегию диверсификации инвестиционного портфеля на сумму 1 миллион рублей. [=====]

Как работает применение знаний из исследования

Измерение компетенции: Промпт четко указывает, какой тип знаний (финансовые рынки), навыков (составление портфеля) и поведения (избегание рискованных рекомендаций) требуется от модели.

Измерение временных рамок: Промпт включает как эпизодический контекст (конкретная ситуация пользователя), так и необходимость применения семантических (общих) принципов инвестирования.

Измерение аудитории: Промпт явно определяет взаимодействие как диадическое (персонализированное для одного пользователя), что позволяет модели адаптировать ответ под конкретные обстоятельства.

Такое структурирование промпта позволяет получить более точный, релевантный и этически выверенный ответ, соответствующий конкретным потребностям

пользователя, вместо обобщенного ответа, который может не учитывать важные нюансы ситуации.