

# Эффективное управление SteerLLM для соблюдения предпочтений путем создания уверенных направлений

Дата: 2025-03-04 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.02989>

Рейтинг: 60

Адаптивность: 75

## Ключевые выводы:

Исследование направлено на разработку метода эффективного управления выводом больших языковых моделей (LLM) в соответствии с предпочтениями пользователей. Авторы предложили теоретическую основу для понимания методов управления моделями и разработали алгоритм CONFST (Confident Direction Steering), который позволяет направлять вывод LLM путем модификации активаций во время вывода без необходимости дополнительной настройки модели.

## Объяснение метода:

Исследование предлагает метод CONFST, позволяющий управлять выводом LLM через модификацию внутренних активаций на основе истории пользователя. Высокая концептуальная ценность: показывает, как модель может адаптироваться к стилю и тематическим предпочтениям без явных инструкций. Ограниченная прямая применимость: требует доступа к внутренним параметрам модели, недоступным в обычных API.

Ключевые аспекты исследования: 1. **Теоретическая модель механизма управления LLM:** Исследование представляет математическую основу для понимания того, как работает "model steering" - метод управления выходными данными LLM через модификацию внутренних активаций модели.

**Метод CONFST (Confident Direction Steering):** Предложен новый алгоритм, который находит "уверенные направления" в пространстве активаций модели, представляющие предпочтения пользователя, и использует их для управления выводом LLM.

**Многонаправленное управление:** В отличие от существующих методов, которые обычно работают только с двумя направлениями (например, правдивое vs неправдивое), CONFST способен работать с множественными направлениями предпочтений одновременно.

**Простота внедрения:** Метод не требует перебора всех слоёв и головок внимания для выбора наиболее подходящих, а работает с фиксированным слоем, что значительно упрощает реализацию.

**Неявное управление без инструкций:** В отличие от многих других методов, CONFST не требует явных инструкций от пользователя, а выводит предпочтения из истории взаимодействия.

Дополнение: Методы исследования действительно требуют доступа к внутренним активациям модели и, в представленной форме, не могут быть напрямую применены в стандартном чате без API-доступа или дообучения. Однако ключевые концепции и подходы могут быть адаптированы для использования в стандартных чатах.

Адаптируемые концепции и подходы:

**Использование истории взаимодействий для неявного управления:** Пользователи могут создать "профиль предпочтений" через серию взаимодействий в одной сессии, формируя у модели понимание их стиля. Например, начав разговор с нескольких примеров предпочитаемого стиля (краткого, технического, разговорного), пользователь может "настроить" модель.

**Многонаправленное управление:** Вместо технического управления активациями, пользователи могут явно указывать несколько параметров в своих запросах: "Ответь кратко, технически точно и с фокусом на практическом применении".

**Отбор "уверенных" примеров:** Пользователи могут предоставлять только самые яркие и однозначные примеры предпочитаемого стиля/тематики в начале разговора, следуя идее отбора "уверенных направлений".

**Постепенное обучение предпочтениям:** Пользователи могут давать обратную связь после ответов модели ("Это хорошо, но сделай следующий ответ более кратким"), постепенно направляя модель в нужную сторону.

Ожидаемые результаты от адаптации этих подходов:

- Более персонализированные ответы, соответствующие стилистическим предпочтениям пользователя
- Возможность неявного управления фокусом ответов на определенные тематики без явного указания в каждом запросе
- Улучшенный пользовательский опыт за счет адаптации модели к индивидуальным потребностям

Хотя эти адаптированные подходы не будут столь же эффективны, как прямое

управление активациями, они могут значительно улучшить взаимодействие пользователей с LLM в стандартных чатах, применяя основные принципы исследования CONFST.

Анализ практической применимости: 1. **Теоретическая модель механизма управления LLM:** - Прямая применимость: Низкая. Математическая модель полезна для исследователей, но не для обычных пользователей. - Концептуальная ценность: Высокая. Объясняет, почему и как работает управление моделью, что может помочь пользователям понять, как LLM реагируют на разные типы запросов. - Потенциал для адаптации: Средний. Теоретическое понимание может помочь пользователям более эффективно формулировать запросы, зная концептуально, как модель обрабатывает информацию.

**Метод CONFST:** Прямая применимость: Средняя. Требуется доступ к внутренним активациям модели, что обычно недоступно в общедоступных API, но принципы могут быть адаптированы. Концептуальная ценность: Высокая. Показывает, что модель может "запоминать" стиль пользователя из истории взаимодействий. Потенциал для адаптации: Высокий. Пользователи могут применить принцип "обучения" модели своему стилю через последовательные взаимодействия.

#### **Многонаправленное управление:**

Прямая применимость: Низкая для обычных пользователей, высокая для разработчиков. Концептуальная ценность: Высокая. Понимание того, что модель может одновременно учитывать несколько аспектов предпочтений. Потенциал для адаптации: Высокий. Пользователи могут задавать запросы, учитывающие несколько аспектов одновременно (например, "краткий, но информативный ответ").

#### **Простота внедрения:**

Прямая применимость: Высокая для разработчиков, низкая для обычных пользователей. Концептуальная ценность: Средняя. Упрощение технической реализации не особенно важно для конечных пользователей. Потенциал для адаптации: Средний. Упрощенные методы могут быть быстрее внедрены в пользовательские интерфейсы.

#### **Неявное управление без инструкций:**

Прямая применимость: Высокая. Позволяет модели адаптироваться к пользователю без явных инструкций. Концептуальная ценность: Очень высокая. Демонстрирует, что модели могут "учиться" на истории взаимодействий. Потенциал для адаптации: Высокий. Пользователи могут формировать последовательность запросов так, чтобы модель "уловила" их предпочтения. Сводная оценка полезности: Предварительная оценка: 65 из 100.

Исследование предлагает метод CONFST, который может значительно улучшить персонализацию взаимодействия с LLM без необходимости дорогостоящей дополнительной тренировки. Метод позволяет модели адаптироваться к

предпочтениям пользователя по стилю (краткость, полезность) и тематике контента на основе истории взаимодействий.

Однако реализация требует доступа к внутренним активациям модели, что обычно недоступно через стандартные API. Несмотря на это, концептуальные идеи (использование истории взаимодействий для адаптации, одновременное управление несколькими аспектами) могут быть применены пользователями в модифицированном виде.

Контраргументы к оценке: 1. Оценка могла бы быть выше (75-80), так как исследование предлагает неявное управление моделью без необходимости явных инструкций, что очень ценно для обычных пользователей, которые могут не знать, как эффективно формулировать запросы. 2. Оценка могла бы быть ниже (45-50), поскольку реализация метода требует специальных технических навыков и доступа к внутренним параметрам модели, что делает прямое применение невозможным для большинства пользователей.

После рассмотрения этих аргументов, я корректирую оценку до 60 из 100, поскольку, несмотря на высокую концептуальную ценность, практическая применимость для широкой аудитории ограничена техническими требованиями.

Оценка в 60 баллов отражает высокую полезность исследования для понимания принципов персонализации LLM, но при этом учитывает значительные ограничения для непосредственного применения этих принципов обычными пользователями.

Уверенность в оценке: Очень сильная. Исследование четко описывает метод CONFST, его преимущества и ограничения. Эксперименты подтверждают эффективность метода для управления различными аспектами вывода LLM (тематика, стиль). Ограничения для широкого применения также ясны: необходимость доступа к внутренним активациям модели.

Оценка адаптивности: Оценка адаптивности: 75 из 100.

Принципы и концепции исследования имеют высокий потенциал для адаптации:

1) Идея использования истории взаимодействий для неявного управления моделью может быть адаптирована в обычных чатах путем последовательного формирования определенного "профиля" пользовательских предпочтений.

2) Концепция одновременного управления несколькими аспектами вывода (например, тематика + стиль) может быть применена через явное указание множественных параметров в промпте.

3) Понимание того, что модель может идентифицировать "уверенные направления" в пространстве активаций, может помочь пользователям формировать более последовательные и однозначные запросы.

4) Принцип выбора "уверенных" примеров для обучения может быть адаптирован в

виде предоставления модели только наиболее репрезентативных образцов предпочитаемого стиля/тематики.

Ограничением является то, что полная реализация метода требует технического доступа к внутренним активациям модели, но общие принципы могут быть адаптированы через техники формирования промптов и контекста.

|| <Оценка: 60> || <Объяснение: Исследование предлагает метод CONFST, позволяющий управлять выводом LLM через модификацию внутренних активаций на основе истории пользователя. Высокая концептуальная ценность: показывает, как модель может адаптироваться к стилю и тематическим предпочтениям без явных инструкций. Ограниченная прямая применимость: требует доступа к внутренним параметрам модели, недоступным в обычных API.> || <Адаптивность: 75>

## Prompt:

Использование SteerLLM и CONFST в промптах для GPT

### Ключевые знания из исследования

Исследование предлагает метод CONFST (Confident Direction Steering), который позволяет управлять выводом языковых моделей без дополнительной настройки, используя "уверенные направления" из истории взаимодействия с пользователем. Хотя сам метод требует доступа к внутренним активациям модели (что недоступно обычным пользователям GPT), принципы исследования можно адаптировать для создания эффективных промптов.

### Пример промпта с использованием принципов SteerLLM

[=====] Я хочу, чтобы ты действовал как эксперт по кибербезопасности. Вот несколько примеров моего предпочтительного стиля коммуникации:

"Проблема уязвимости SQL-инъекций решается применением параметризованных запросов." "При анализе инцидента важно сохранять все логи и доказательства в неизменном виде." "Многофакторная аутентификация снижает риск несанкционированного доступа на 99,9%." Обрати внимание на технический, лаконичный стиль с конкретными цифрами и рекомендациями.

Используя этот стиль, объясни, пожалуйста, основные принципы защиты от атак типа "человек посередине" для обычного пользователя. [=====]

### Объяснение применения знаний из исследования

В этом промпте используются ключевые принципы из исследования CONFST:

**Создание "уверенных направлений"** — предоставление нескольких примеров предпочтительного стиля, что создает четкие векторы для модели.

**Неявное управление на основе истории** — вместо прямых инструкций ("будь кратким"), демонстрируются примеры желаемого поведения, что соответствует подходу CONFST извлекать предпочтения из истории.

**Комбинирование нескольких направлений** — примеры одновременно задают несколько параметров: техническую точность, лаконичность и конкретность.

**Явное указание на ключевые аспекты** — обращение внимания модели на конкретные характеристики примеров усиливает "уверенные направления".

Такой подход помогает получить более персонализированные и соответствующие предпочтениям ответы без прямого доступа к внутренним механизмам GPT.