

# Исследование и контроль разнообразия в беседе с LLM-агентом

Дата: 2025-02-21 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2412.21102>

Рейтинг: 72

Адаптивность: 80

## Ключевые выводы:

Исследование направлено на изучение и контроль разнообразия в диалогах между агентами на основе LLM. Основная цель - разработать метод, позволяющий балансировать между стабильностью в структурированных задачах и вариативностью в творческих сценариях. Главный результат - создание метода Adaptive Prompt Pruning (APP), который позволяет контролировать разнообразие диалогов через единый параметр  $\lambda$ , динамически удаляя компоненты промпта на основе их весов внимания.

## Объяснение метода:

Исследование предлагает практичный метод контроля разнообразия в диалогах с LLM через управление содержимым промпта. Хотя полная реализация APP требует доступа к весам внимания, основные принципы (удаление избыточной информации, порядок блоков) легко адаптируются к обычному использованию. Исследование дает глубокое понимание факторов, влияющих на разнообразие ответов, что ценно для любого пользователя LLM.

## Ключевые аспекты исследования: 1. **Адаптивное прореживание промпта (APP)** - метод для контроля разнообразия диалогов в симуляциях LLM-агентов путем динамического удаления компонентов промпта на основе их весов внимания.

**Модуляризация промпта** - исследователи разделили промпт на блоки (базовая информация, память, предыдущие диалоги, окружение и текущий диалог), что позволило изучить влияние каждого компонента на разнообразие.

**Параметр  $\lambda$  для контроля разнообразия** - единый параметр, позволяющий плавно регулировать степень разнообразия диалогов: чем выше  $\lambda$ , тем больше компонентов удаляется из промпта.

**Процесс проверки и исправления** - метод для устранения несоответствий, возникающих при удалении информации из промпта, что позволяет сохранять связность диалога.

**Анализ влияния порядка блоков и предварительных знаний модели** - исследование показало, что порядок блоков и частота имен существенно влияют на разнообразие диалогов.

**## Дополнение:**

**### Применимость методов в стандартном чате без дообучения или API**

Исследование не требует дообучения модели или специального API для применения его ключевых концепций. Хотя полная реализация APP с использованием весов внимания недоступна в стандартных интерфейсах, основные принципы и выводы могут быть адаптированы для использования в обычном чате.

**Применимые концепции и подходы:**

**Модуляризация промпта и выборочное включение информации:** Пользователи могут структурировать свои запросы по блокам (контекст, предыстория, инструкции) Целенаправленно исключать определенные блоки информации для повышения разнообразия

**Управление порядком информации:**

Размещение наиболее важной информации в начале промпта Избегание размещения текущего контекста в самом начале промпта

**Использование известных имен и концепций:**

При необходимости увеличить разнообразие - использовать общеизвестные имена/концепции При необходимости более предсказуемых ответов - использовать малоизвестные имена

**Двухэтапный подход с проверкой:**

Генерация ответа с ограниченной информацией для разнообразия Проверка ответа на соответствие важным исключенным деталям При необходимости - запрос на корректировку **Ожидаемые результаты:**

- Повышение разнообразия ответов при разных запусках с аналогичными запросами
- Лучший контроль над степенью креативности модели
- Более глубокое понимание причин однотипности ответов
- Возможность сознательно балансировать между разнообразием и согласованностью информации

Важно отметить, что исследователи использовали специальные техники (доступ к весам внимания) не потому, что это необходимо для работы метода, а для более точной количественной оценки и автоматизации процесса, который в упрощенном виде доступен любому пользователю.

### ## Анализ практической применимости: **Адаптивное прореживание промпта (APP)**

- Прямая применимость: Высокая. Пользователи могут непосредственно использовать принцип удаления частей контекста для получения более разнообразных ответов. Параметр  $\lambda$  позволяет точно настраивать уровень разнообразия. - Концептуальная ценность: Высокая. Исследование демонстрирует, как избыточная информация в промпте может ограничивать разнообразие ответов, что помогает пользователям лучше понимать причины однотипности ответов LLM. - Потенциал для адаптации: Высокий. Хотя полная реализация APP требует доступа к весам внимания, пользователи могут применять упрощенный подход, удаляя определенные блоки информации из своих запросов.

**Модуляризация промпта** - Прямая применимость: Средняя. Разделение промпта на блоки - полезная техника для структурирования запросов, которую пользователи могут применять в повседневном взаимодействии с LLM. - Концептуальная ценность: Высокая. Понимание того, какие блоки информации больше всего влияют на разнообразие (например, блок "Память"), дает пользователям инструмент для контроля над ответами LLM. - Потенциал для адаптации: Высокий. Пользователи могут экспериментировать с добавлением или удалением определенных блоков информации для достижения желаемого уровня разнообразия.

**Процесс проверки и исправления** - Прямая применимость: Средняя. Пользователи могут внедрить дополнительный шаг проверки ответов на соответствие удаленной информации. - Концептуальная ценность: Высокая. Понимание компромисса между разнообразием и согласованностью информации. - Потенциал для адаптации: Средний. Требует дополнительных запросов к модели, но может быть реализован через простые инструкции.

**Влияние порядка блоков и предварительных знаний** - Прямая применимость: Высокая. Пользователи могут экспериментировать с порядком информации в промптах. - Концептуальная ценность: Очень высокая. Понимание того, как порядок информации и известность имен влияют на ответы, дает глубокое понимание работы LLM. - Потенциал для адаптации: Высокий. Легко применимо в любом контексте использования LLM.

## **Prompt:**

Использование знаний из исследования разнообразия диалогов в промптах для GPT

### ## Ключевые применимые знания из исследования

Исследование APP (Adaptive Prompt Pruning) показывает, что:

Разнообразие диалогов можно контролировать через удаление определенных компонентов промпта. Блок памяти больше всего ограничивает разнообразие ответов. Порядок блоков в промпте значительно влияет на разнообразие (хронологический порядок лучше). Комбинирование методов (APP + настройка температуры) дает синергетический эффект. Использование популярных имен активирует параметрические знания модели. ## Пример промпта с применением знаний из исследования

[=====] # Творческая дискуссия о будущем технологий

## Инструкции для GPT ( $\lambda=0.7$ , модификация по методу APP): - Ты эксперт по футурологии по имени Гарри Поттер - Веди диалог в творческом формате, предлагая неожиданные, но обоснованные идеи - [УДАЛЕНО: блок памяти о предыдущих обсуждениях] - Используй последние 2-3 реплики для контекста, но не ограничивай себя только ими - Информация в хронологическом порядке: сначала базовые знания, потом текущий контекст - Температура генерации: 0.8

## Вопрос: Как ты думаешь, как изменится роль социальных сетей в обществе через 15 лет? [=====]

## Объяснение применения знаний из исследования

**Удаление блока памяти** (согласно методу APP с  $\lambda=0.7$ ) - намеренно убираем элемент, который больше всего ограничивает разнообразие

**Хронологический порядок информации** - структурируем промпт так, чтобы информация шла в хронологическом порядке, что способствует разнообразию

**Использование популярного имени** ("Гарри Поттер") - активирует параметрические знания модели для более разнообразных ответов

**Комбинирование методов** - используем и структурные модификации промпта (APP), и настройку температуры (0.8) для синергетического эффекта

**Ограничение контекста** - используем только последние 2-3 реплики вместо всей истории диалога, что уменьшает "якорение" и способствует разнообразию

Такой промпт позволяет получить более творческие и разнообразные ответы без потери связности и релевантности, что особенно ценно для креативных задач, мозговых штурмов и исследовательских дискуссий.