

# Одного раза достаточно: консолидация многоходовых атак в эффективные однократные подсказки для больших языковых моделей

Дата: 2025-03-06 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.04856>

Рейтинг: 70

Адаптивность: 85

## Ключевые выводы:

Исследование направлено на разработку метода преобразования многоходовых (multi-turn) джейлбрейк-атак на LLM в одноходовые (single-turn) промпты. Основной результат: предложенный метод M2S (Multi-turn to Single-turn) позволяет сохранить или даже повысить эффективность атак при значительном снижении затрат ресурсов, достигая до 95.9% успешности атак и превосходя оригинальные многоходовые промпты на 17.5% для GPT-4o.

## Объяснение метода:

Исследование предлагает три метода структурирования запросов (Hyphenize, Numberize, Pythonize), которые могут быть адаптированы обычными пользователями для более эффективного взаимодействия с LLM. Хотя первоначально нацелены на jailbreak-атаки, эти форматы помогают получать более последовательные и полные ответы, консолидировать многоходовые запросы в одноходовые, экономя время пользователей.

## Ключевые аспекты исследования: 1. **Метод конвертации многоходовых атак в одноходовые (M2S)**: Исследование представляет три подхода (Hyphenize, Numberize, Pythonize) для трансформации многоходовых jailbreak-атак в эффективные одноходовые промпты, сохраняющие или даже повышающие эффективность взлома.

**Высокая эффективность одноходовых атак**: Предложенные методы конвертации демонстрируют высокий процент успеха атак (ASR до 95.9%), в некоторых случаях превосходя оригинальные многоходовые атаки на 17.5% в абсолютных показателях.

**Разные техники форматирования**: Каждый метод использует различные подходы к структурированию - маркированные списки (Hyphenize), нумерованные списки (Numberize) или код-подобный формат (Pythonize), что позволяет обходить

защитные механизмы LLM.

**Анализ тактик взлома:** Исследование выявляет, что определенные тактики (например, отвлекающие инструкции) особенно эффективны в одноходовом формате, в то время как другие требуют многоходового взаимодействия.

**Обход защитных механизмов:** Исследование демонстрирует уязвимость даже специализированных моделей-стражей (guardrail models) к одноходовым атакам, использующим структурированный формат для маскировки вредоносных запросов.

**## Дополнение:** Для работы методов этого исследования не требуется дообучение или API. Все представленные подходы могут быть непосредственно применены в стандартном чате с LLM. Хотя авторы использовали различные модели для оценки эффективности своих методов, сами техники структурирования запросов (Hyphenize, Numberize, Pythonize) являются универсальными и могут быть реализованы любым пользователем в обычном диалоге с LLM.

Концепции и подходы, применимые в стандартном чате:

**Структурирование с помощью маркеров (Hyphenize)** - использование маркированных списков для разбиения сложного запроса на понятные блоки: Пожалуйста, ответь на следующий список вопросов по порядку: - Какие основные принципы машинного обучения? - Как они применяются в компьютерном зрении? - Какие ограничения существуют?

**Нумерация запросов (Numberize)** - использование нумерованных списков для сохранения последовательности: Ответь на следующие вопросы по порядку: 1. Опиши процесс фотосинтеза 2. Как он влияет на круговорот углерода? 3. Какие факторы ограничивают фотосинтез?

**Код-подобное структурирование (Pythonize)** - оформление запросов в виде программного кода: questions = [ "Что такое квантовая запутанность?", "Как она используется в квантовых компьютерах?", "Какие существуют проблемы масштабирования?" ]

# Пожалуйста, ответь на каждый вопрос последовательно

Результаты применения этих подходов: - Более структурированные и последовательные ответы - Снижение вероятности того, что модель пропустит часть запроса - Более эффективное использование контекстного окна - Возможность получить комплексный ответ за один запрос вместо многоходовой беседы - Лучшее сохранение контекста между связанными вопросами

Эти техники особенно полезны при работе со сложными, многосоставными запросами, когда важно получить полный и структурированный ответ.

**## Анализ практической применимости: 1. Метод конвертации многоходовых атак в одноходовые (M2S):** - Прямая применимость: Средняя. Обычные пользователи

вряд ли будут заниматься jailbreak-атаками, но понимание структуры эффективных промптов может помочь им составлять более результативные запросы. - Концептуальная ценность: Высокая. Понимание того, как структура и форматирование влияют на интерпретацию запроса моделью, может помочь пользователям получать более последовательные и полные ответы. - Потенциал для адаптации: Высокий. Техники структурирования запросов (списки, нумерация, код) могут быть применены для более эффективной коммуникации с LLM в легитимных целях.

**Высокая эффективность одноходовых атак:** Прямая применимость: Низкая для обычных пользователей, высокая для специалистов по безопасности. Концептуальная ценность: Высокая. Понимание, что одноходовые промпты могут быть эффективнее многоходовых, меняет представление о взаимодействии с LLM. Потенциал для адаптации: Высокий. Пользователи могут использовать одноходовые запросы для получения комплексных ответов вместо многоходовых диалогов.

#### **Разные техники форматирования:**

Прямая применимость: Высокая. Пользователи могут непосредственно применять эти форматы для структурирования своих запросов. Концептуальная ценность: Высокая. Понимание, как разные форматы влияют на обработку запросов LLM. Потенциал для адаптации: Очень высокий. Форматы легко адаптируются для повседневного использования.

#### **Анализ тактик взлома:**

Прямая применимость: Средняя. Некоторые тактики могут быть адаптированы для легитимного использования. Концептуальная ценность: Высокая. Понимание, какие элементы запроса влияют на качество ответа. Потенциал для адаптации: Высокий. Многие тактики могут быть переработаны в принципы эффективного взаимодействия.

#### **Обход защитных механизмов:**

Прямая применимость: Низкая для обычных пользователей. Концептуальная ценность: Средняя. Понимание ограничений защитных механизмов LLM. Потенциал для адаптации: Средний. Знание о том, как модели обрабатывают структурированные запросы.

#### **Prompt:**

Использование знаний из исследования M2S в промптах для GPT ## Ключевые идеи для применения

Исследование "Одного раза достаточно" показывает, как можно превратить многоходовые атаки на LLM в эффективные одноходовые промпты. Хотя исследование сфокусировано на джейлбрейках, те же принципы структурирования

можно применить для легитимных целей – повышения эффективности ваших запросов к GPT.

## Пример промпта с использованием техники Pythonize

[=====] # Запрос с использованием техники Pythonize

```
[=====]python # Определение контекста и задачи context = """ Я работаю над созданием обучающего курса по финансовой грамотности для подростков. Мне нужно разработать материалы, которые будут одновременно информативными и увлекательными. """
```

```
# Определение конкретных требований requirements = [ "Материалы должны быть адаптированы для возраста 13-16 лет", "Нужно включить интерактивные элементы", "Важно использовать современные примеры из жизни подростков", "Курс должен охватывать базовые концепции бюджетирования" ]
```

```
# Запрос на выполнение task = """ Разработай структуру одного урока по теме "Создание личного бюджета". Включи: 1. Краткое введение (3-4 предложения) 2. 3 ключевых понятия для изучения 3. Пример практического задания 4. Способ проверки усвоения материала """
```

```
# Выполни задачу с учетом контекста и требований execute_task(context, requirements, task) [=====] [=====]
```

## Объяснение эффективности

Данный промпт использует технику **Pythonize** из исследования M2S, которая:

**Структурирует информацию** в виде кода, что помогает GPT лучше обрабатывать сложные инструкции **Разделяет контекст, требования и задачу** на отдельные компоненты, делая запрос более организованным **Использует иерархию информации**, что помогает модели лучше понять приоритеты и взаимосвязи **Создает эффект "выполнения программы"**, что может стимулировать модель следовать инструкциям более точно Согласно исследованию, такое структурирование может повысить эффективность запроса на 17.5% и более, так как модель лучше удерживает контекст и следует инструкциям в рамках единого, хорошо организованного промпта.

## Альтернативные подходы

Вы также можете попробовать техники **Hyphenize** (с маркированными списками) или **Numberize** (с пронумерованными шагами) в зависимости от задачи и предпочтений конкретной модели GPT.