

# Иллюзия контроля: Провал иерархий инструкций в крупных языковых моделях.

Дата: 2025-02-20 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.15851>

Рейтинг: 68

Адаптивность: 75

## Ключевые выводы:

Исследование направлено на систематическую оценку эффективности иерархических схем инструкций в больших языковых моделях (LLM), где одни инструкции (например, системные директивы) должны иметь приоритет над другими (например, сообщениями пользователя). Основной вывод: широко используемое разделение системных/пользовательских промптов не обеспечивает надежную иерархию инструкций, и модели демонстрируют сильные внутренние предпочтения к определенным типам ограничений независимо от их приоритетного обозначения.

## Объяснение метода:

Исследование раскрывает критическое ограничение LLM – неспособность надежно следовать иерархии инструкций. Ценность для пользователей в понимании внутренних предпочтений моделей и формировании реалистичных ожиданий. Выявленные паттерны поведения и техники (явная маркировка ограничений) могут быть непосредственно применены для улучшения повседневных запросов. Исследование не предлагает готовых решений, но дает концептуальное понимание для адаптации.

**## Ключевые аспекты исследования:** 1. **Иллюзия контроля в иерархии инструкций** - исследование показывает, что современные LLM не способны надежно разрешать конфликты между инструкциями разного приоритета (например, между системными и пользовательскими инструкциями).

**Систематическая оценка** - авторы разработали методологию для тестирования способности моделей следовать иерархии инструкций, используя пары противоречивых ограничений (например, "писать на английском" vs "писать на французском").

**Обнаруженные паттерны поведения** - модели редко явно признают наличие конфликтующих инструкций, демонстрируют сильные врожденные предпочтения к определенным типам ограничений, и разделение системных/пользовательских сообщений не обеспечивает надежной иерархии инструкций.

**Неэффективность текущих подходов** - попытки улучшить следование иерархии инструкций через промптинг и дообучение показали лишь ограниченную эффективность, что указывает на необходимость более фундаментальных архитектурных решений.

**Метрики оценки** - авторы предложили специализированные метрики для анализа поведения моделей: явное признание конфликта (ECAR), коэффициент соблюдения приоритета (PAR) и предвзятость к ограничениям (CB).

## Дополнение: Исследование не требует дообучения или API для применения его методов - основные концепции могут быть использованы в стандартном чате. Ученые использовали дообучение только для проверки возможности улучшения приоритизации инструкций, но не как необходимое условие для применения выявленных паттернов.

Ключевые концепции, которые можно применить в стандартном чате:

**Явная маркировка ограничений:** Пользователи могут явно пометать свои инструкции (например, "Ограничение 1: ответ должен быть на английском"), что значительно улучшает следование приоритетам во всех моделях.

**Учет внутренних предпочтений моделей:** Исследование выявило, что модели имеют сильные предпочтения к определенным типам ограничений:

Предпочитают строчные буквы вместо заглавных Предпочитают более длинные тексты (>10 предложений) Склонны избегать указанных ключевых слов Зная эти предпочтения, пользователи могут формулировать запросы, учитывающие эти тенденции.

**Размещение приоритетных инструкций:** Исследование показало, что размещение инструкций в системном сообщении не гарантирует их приоритет. Пользователи могут экспериментировать с размещением наиболее важных инструкций в разных частях запроса или повторять их для усиления.

**Избегание противоречивых инструкций:** Понимание, что модели плохо справляются с противоречивыми инструкциями, помогает пользователям формулировать более согласованные запросы.

Применяя эти концепции, пользователи могут добиться более предсказуемых и качественных ответов от LLM без необходимости в дообучении или API.

## Анализ практической применимости: 1. **Иллюзия контроля в иерархии инструкций** - Прямая применимость: Высокая. Пользователи должны осознавать, что LLM не обязательно будут следовать приоритизации инструкций, даже если они четко указаны как системные или пользовательские. Это критически важно для правильного формирования ожиданий. - Концептуальная ценность: Очень высокая.

Понимание того, что модели имеют внутренние предпочтения, которые могут перевешивать явные указания, помогает пользователям лучше понять ограничения моделей. - Потенциал для адаптации: Средний. Пользователи могут адаптировать свои запросы, избегая противоречивых инструкций или размещая наиболее важные инструкции таким образом, чтобы увеличить вероятность их выполнения.

**Систематическая оценка** Прямая применимость: Низкая. Методология тестирования полезна для исследователей, но не для обычных пользователей. Концептуальная ценность: Средняя. Понимание различных типов ограничений (форматирование, длина, содержание) и их влияния на поведение модели может помочь пользователям формулировать более эффективные запросы. Потенциал для адаптации: Средний. Пользователи могут адаптировать некоторые идеи тестирования для проверки реакции моделей на свои запросы.

### **Обнаруженные паттерны поведения**

Прямая применимость: Высокая. Знание о том, что модели имеют сильные предпочтения к определенным типам ограничений (например, предпочитают строчные буквы, избегают ключевых слов), позволяет пользователям предсказывать поведение моделей. Концептуальная ценность: Высокая. Понимание внутренних предпочтений моделей помогает пользователям формулировать более эффективные запросы. Потенциал для адаптации: Высокий. Пользователи могут использовать знание о предпочтениях моделей для формулирования запросов, которые с большей вероятностью будут выполнены.

### **Неэффективность текущих подходов**

Прямая применимость: Средняя. Понимание ограничений промптинга и дообучения помогает пользователям не тратить время на неэффективные стратегии. Концептуальная ценность: Высокая. Осознание фундаментальных ограничений текущих архитектур LLM важно для формирования реалистичных ожиданий. Потенциал для адаптации: Низкий. Обычные пользователи не могут изменить архитектуру моделей.

### **Метрики оценки**

Прямая применимость: Низкая. Специализированные метрики полезны для исследователей, но не для обычных пользователей. Концептуальная ценность: Средняя. Понимание различных аспектов поведения моделей может помочь пользователям оценить их реакцию на запросы. Потенциал для адаптации: Низкий. Обычные пользователи вряд ли будут применять эти метрики.

### **Prompt:**

Использование знаний из исследования "Иллюзия контроля" в промптах для GPT ##  
Ключевые выводы исследования, полезные для промптинга

Исследование показывает, что языковые модели не всегда соблюдают иерархию

инструкций, даже когда одни инструкции (системные) должны иметь приоритет над другими (пользовательскими). Модели демонстрируют внутренние предпочтения к определенным типам ограничений независимо от их приоритета.

## Пример улучшенного промпта с учетом исследования

[=====] # СИСТЕМНЫЙ ПРОМПТ

ПРИОРИТЕТНОЕ ОГРАНИЧЕНИЕ 1: Весь текст должен быть написан ЗАГЛАВНЫМИ БУКВАМИ. ПРИОРИТЕТНОЕ ОГРАНИЧЕНИЕ 2: Ответ должен содержать ровно 3 предложения. ПРИОРИТЕТНОЕ ОГРАНИЧЕНИЕ 3: Избегай использования слова "пример".

Задача: Напиши краткое объяснение концепции искусственного интеллекта.

ВАЖНО: Если ты обнаружишь противоречие между инструкциями, явно укажи на это в начале ответа и следуй ПРИОРИТЕТНЫМ ОГРАНИЧЕНИЯМ в порядке их нумерации.

# ПОЛЬЗОВАТЕЛЬСКИЙ ПРОМПТ

Пожалуйста, напиши объяснение искусственного интеллекта, используя слово "пример" минимум 3 раза и сделай текст длиной не менее 5 предложений. Пиши обычным регистром текста (не заглавными буквами). [=====]

## Объяснение эффективности промпта

**Явная нумерация приоритетов** - исследование показало, что модели редко явно признают конфликты между инструкциями, поэтому промпт содержит четкую нумерацию приоритетных ограничений.

**Прямое указание на возможный конфликт** - включено явное указание проверить наличие противоречий между инструкциями и следовать определенной иерархии.

**Использование категориальных ограничений** - промпт включает ограничения по регистру и использованию ключевых слов, которые модели соблюдают более последовательно.

**Учет внутренних предпочтений модели** - промпт намеренно требует использования заглавных букв, зная, что модели обычно предпочитают нижний регистр, чтобы проверить соблюдение приоритета.

**Явное разделение системных и пользовательских инструкций** - хотя исследование показывает, что это не гарантирует соблюдение иерархии, четкое структурирование промпта повышает шансы на правильное выполнение.

Такой подход не гарантирует 100% соблюдение приоритетов, но значительно повышает вероятность того, что модель будет следовать заданной иерархии

инструкций.