

# Закон затмения знаний: к пониманию, прогнозированию и предотвращению галлюцинаций LLM

Дата: 2025-02-22 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.16143>

Рейтинг: 68

Адаптивность: 75

## Ключевые выводы:

Исследование направлено на понимание, предсказание и предотвращение галлюцинаций в больших языковых моделях (LLM). Авторы выявили феномен 'затенения знаний' (knowledge overshadowing), при котором доминирующие знания в модели подавляют менее распространенные знания, что приводит к генерации неточных фактов. Исследование установило логарифмически-линейный закон, позволяющий количественно предсказывать вероятность галлюцинаций в зависимости от популярности знаний, их длины и размера модели.

## Объяснение метода:

Исследование предлагает ценную концепцию "знаниевого затенения", объясняющую причины галлюцинаций в LLM, и лог-линейный закон для их предсказания. Высокая концептуальная ценность для понимания ограничений LLM, но техническая сложность метода CoDA и математическое обоснование ограничивают прямое применение обычными пользователями. Требуется адаптация для широкой аудитории.

Ключевые аспекты исследования: 1. **Концепция "знаниевого затенения" (knowledge overshadowing)** - ключевой механизм галлюцинаций в LLM, при котором доминирующие знания подавляют менее распространенные, что приводит к искажению фактов при генерации текста.

**Лог-линейный закон галлюцинаций** - исследователи установили, что вероятность фактических галлюцинаций линейно растет с логарифмической шкалой: (а) относительной популярности знаний, (б) относительной длины знаний и (в) размера модели.

**Теоретическое обоснование эффекта затенения** - авторы связывают эффект с механизмами обобщения в LLM, показывая, что доминирующие знания имеют лучшие границы обобщения, что приводит к подавлению редких знаний.

**Метод CoDA (Contrastive Decoding to Amplify overshadowed knowledge)** - предложенный метод декодирования, который выявляет затененные знания и усиливает их влияние, что значительно снижает галлюцинации без необходимости дополнительного обучения.

**Экспериментальное подтверждение** - авторы провели масштабные эксперименты на предобученных и дообученных моделях разных размеров, подтвердив закономерности эффекта затенения знаний.

Дополнение:

Исследование не требует дообучения или API для применения основных концепций, хотя метод CoDA в полной форме использует доступ к вероятностям токенов, что недоступно в стандартном чате.

Концепции и подходы, применимые в стандартном чате:

**Понимание знаниевого затенения** - пользователи могут избегать смешивания в запросе популярных и редких концепций, разбивая сложные вопросы на части.

**Применение лог-линейного закона** - можно упреждать галлюцинации, учитывая:

Избегать длинных контекстов при запросе о редких фактах Для редких тем предпочитать более крупные модели Разделять запросы с несколькими условиями на отдельные

**Упрощённая версия CoDA** - можно реализовать через:

Запрос модели проверить свой ответ, выделяя потенциально спорные части  
Использование контрастных запросов, где одни элементы маскируются для выявления влияния других

**Стратегия разделения условий** - когда запрос содержит несколько условий (например, "известные женщины-учёные в области ИИ"), разбивать его на последовательность уточняющих вопросов.

Результаты от применения: снижение количества фактических ошибок, особенно в сложных запросах с несколькими условиями или при запросах о редких фактах в контексте более популярных тем.

Анализ практической применимости: 1. **Концепция знаниевого затенения - Прямая применимость:** Средняя. Пользователи могут осознанно избегать ситуаций, когда в запросе смешиваются разные по популярности концепции, особенно при формулировке вопросов. - **Концептуальная ценность:** Высокая. Понимание механизма затенения помогает пользователям осознать, почему LLM могут искажать факты даже при высококачественном обучении. - **Потенциал для**

**адаптации:** Высокий. Пользователи могут разработать стратегии формулировки запросов, подчеркивая менее популярные аспекты или разбивая запросы на части.

**Лог-линейный закон галлюцинаций** **Прямая применимость:** Низкая. Рядовым пользователям сложно применить математическую формулу для оценки вероятности галлюцинаций. **Концептуальная ценность:** Высокая. Понимание факторов, влияющих на вероятность галлюцинаций, помогает пользователям предвидеть проблемные ситуации. **Потенциал для адаптации:** Средний. Знание о влиянии длины контекста и размера модели может помочь выбрать подходящую модель и формат запроса.

## Метод CoDA

**Прямая применимость:** Низкая для обычных пользователей, так как требует доступа к вероятностям токенов и техническую реализацию. **Концептуальная ценность:** Средняя. Идея контрастного декодирования может быть адаптирована в упрощенной форме. **Потенциал для адаптации:** Высокий. Принципы метода могут быть реализованы в интерфейсах чат-ботов или инструментах проверки фактов.

## Экспериментальные результаты

**Прямая применимость:** Средняя. Пользователи могут использовать выявленные закономерности для критической оценки ответов LLM. **Концептуальная ценность:** Высокая. Результаты наглядно демонстрируют, когда и почему возникают галлюцинации. **Потенциал для адаптации:** Высокий. Знание о влиянии популярности, длины и размера модели может быть преобразовано в практические рекомендации.

## Теоретическое обоснование

**Прямая применимость:** Низкая. Математические выкладки малополезны для обычных пользователей. **Концептуальная ценность:** Средняя. Понимание связи между обобщением и галлюцинациями углубляет понимание работы LLM. **Потенциал для адаптации:** Средний. Теоретические insights могут быть переформулированы в более доступные принципы работы с LLM. Сводная оценка полезности: На основе проведенного анализа определяю общую оценку полезности исследования: **68/100**.

Исследование предлагает глубокое понимание механизма галлюцинаций в LLM через концепцию знаниевого затенения. Эта концепция имеет высокую практическую ценность для пользователей, позволяя им лучше понимать, почему и когда LLM могут выдавать неверную информацию. Лог-линейный закон помогает предсказать вероятность галлюцинаций в зависимости от трех ключевых факторов, что может быть адаптировано в практические рекомендации.

Однако техническая сложность метода CoDA и математическое обоснование ограничивают прямое применение результатов обычными пользователями. Метод требует доступа к внутренним вероятностям модели, что недоступно в стандартных

чат-интерфейсах.

### **Контраргументы к оценке:**

*Почему оценка могла бы быть выше:* Концепция знаниевого затенения фундаментально меняет понимание того, как работают LLM, и может значительно повысить осведомленность пользователей о возможных проблемах. Также метод CoDA показывает значительное улучшение фактологической точности на нескольких датасетах.

*Почему оценка могла бы быть ниже:* Исследование в значительной степени академическое, с акцентом на теоретическое обоснование и технические детали, которые малодоступны для обычных пользователей. Применение метода CoDA требует технических знаний и программных навыков.

После рассмотрения этих аргументов я считаю, что оценка 68/100 справедливо отражает баланс между высокой концептуальной ценностью и ограниченной прямой применимостью для широкой аудитории.

Уверенность в оценке: Очень сильная. Исследование было проанализировано комплексно, с учетом всех ключевых аспектов и их применимости для широкой аудитории. Оценка основана на тщательном анализе как прямой практической применимости, так и концептуальной ценности работы.

### **Оценка адаптивности: Оценка адаптивности: 75/100**

Исследование представляет высокоадаптивные концепции и принципы, которые могут быть трансформированы для использования широкой аудиторией:

Концепция знаниевого затенения легко адаптируется в практические рекомендации по формулировке запросов к LLM, помогая избегать смешивания концепций разной популярности.

Лог-линейный закон может быть преобразован в простые эвристики для оценки надежности ответов LLM (например, "длинные запросы с редкими концепциями имеют высокий риск галлюцинаций").

Принципы метода CoDA, хотя и технически сложные, могут быть реализованы в пользовательских интерфейсах как опции проверки фактов или альтернативных формулировок запросов.

Экспериментальные результаты могут быть адаптированы в образовательные материалы, помогающие пользователям критически оценивать ответы LLM.

Понимание влияния размера модели может помочь пользователям выбирать подходящие модели для своих задач, особенно когда важна фактологическая точность.

Однако метод CoDA в его текущей форме требует доступа к внутренним вероятностям модели, что ограничивает его прямую адаптацию в стандартных чат-интерфейсах.

|| <Оценка: 68> || <Объяснение: Исследование предлагает ценную концепцию "знаниевого затенения", объясняющую причины галлюцинаций в LLM, и лог-линейный закон для их предсказания. Высокая концептуальная ценность для понимания ограничений LLM, но техническая сложность метода CoDA и математическое обоснование ограничивают прямое применение обычными пользователями. Требуется адаптация для широкой аудитории.> || <Адаптивность: 75>

## Prompt:

Использование закона затмения знаний в промптах для GPT

### Ключевые принципы из исследования

Исследование выявило феномен **затмения знаний (knowledge overshadowing)**, при котором: - Доминирующие знания подавляют редкие знания - Частота галлюцинаций зависит от: - Относительной популярности знаний (P) - Относительной длины знаний (L) - Размера модели (S)

### Пример промпта с учетом закона затмения знаний

[=====]

Запрос о малоизвестном историческом факте Мне нужна информация о малоизвестном историческом событии - восстании в городе [название редкого события].

ВАЖНО: - Я знаю, что это событие менее известно, чем [популярное событие того же периода] - Пожалуйста, сосредоточься именно на запрашиваемом событии, а не на более известных событиях того периода - Если ты не уверен в фактах, укажи это явно и не пытайся заполнить пробелы предположениями - Приведи все доступные тебе детали именно об этом конкретном событии, его датах, участниках и последствиях

Дополнительно: если возможно, сравни это событие с более известным [популярное событие], указав ключевые различия. [=====]

### Почему это работает

**Борьба с низкой популярностью (P):** Промпт явно указывает на редкость запрашиваемой информации и предупреждает модель не подменять ее более

популярными знаниями

**Компенсация относительной длины (L):** Промпт делает акцент на важных элементах, выделяя их структурно и повторяя ключевые моменты

**Снижение риска галлюцинаций:** Промпт содержит прямую инструкцию указывать на неуверенность вместо генерации потенциально неверных фактов

**Контрастное усиление:** Запрос на сравнение с более известным событием действует подобно методу CoDA из исследования, помогая модели лучше дифференцировать знания

Такой подход помогает "вытащить" затененные знания на передний план и снизить вероятность галлюцинаций, особенно при работе с редкими фактами.