

# Визуальное описание на основе контекста снижает количество галлюцинаций и улучшает reasoning в LVLM

Дата: 2025-03-05 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2405.15683>

Рейтинг: 72

Адаптивность: 85

## Ключевые выводы:

Исследование направлено на выявление причин галлюцинаций в больших мультимодальных языковых моделях (LVLM) и разработку метода их снижения. Основной вывод: существующие методы снижения галлюцинаций эффективны для задач визуального распознавания, но не для когнитивных задач, требующих рассуждений. Авторы выявили ключевую проблему - разрыв визуального восприятия: LVLM могут распознавать визуальные элементы, но не могут полноценно интерпретировать их в контексте запроса.

## Объяснение метода:

Исследование предоставляет ценное понимание причин галлюцинаций в LVLMs и предлагает метод VDGD для их снижения. Хотя полная реализация требует технических знаний, основной принцип (использование описания изображения перед основным запросом) может быть легко применен обычными пользователями через последовательные запросы, значительно улучшая точность ответов для задач, требующих рассуждения.

**## Ключевые аспекты исследования: 1. Идентификация причин галлюцинаций в LVLMs:** Авторы выявили, что существующие методы снижения галлюцинаций работают хорошо для задач визуального распознавания, но неэффективны для когнитивных задач, требующих рассуждения.

**Определение разрыва в визуальном восприятии:** Исследование показало, что LVLMs могут распознавать визуальные элементы, но испытывают трудности с их контекстуализацией относительно запроса пользователя и связыванием с внутренними знаниями, что критично для рассуждений.

**Метод Visual Description Grounded Decoding (VDGD):** Предложен простой и не требующий дообучения метод для улучшения рассуждений в LVLMs путем создания детального описания изображения и использования его для направления генерации

ответа.

**Категоризация типов галлюцинаций:** Авторы классифицировали галлюцинации на языковые, стилистические, визуальные и связанные с обучением на инструкциях (IT), что позволяет лучше понять их происхождение.

**Создание бенчмарка VaLLu:** Разработан комплексный бенчмарк для оценки когнитивных способностей LVLMS, включающий задачи различной сложности и типов.

**## Дополнение:** Полная реализация методов исследования действительно требует доступа к API или возможности модификации процесса декодирования модели, что недоступно большинству обычных пользователей. Однако ключевая концепция метода VDGD может быть успешно адаптирована для использования в стандартном чате пользователями без технических навыков.

Основной принцип VDGD заключается в том, что детальное описание изображения помогает модели лучше контекстуализировать визуальную информацию при ответе на сложные вопросы. Этот принцип можно реализовать в стандартном чате следующим образом:

**Двухэтапный запрос:** Пользователь может сначала попросить модель подробно описать изображение, а затем задать основной вопрос, ссылаясь на это описание.

**Направленное описание:** Можно запросить описание, ориентированное на конкретную задачу, например: "Опиши детально изображение, обращая особое внимание на числовые данные в графике" перед вопросом о тенденциях.

**Проверка понимания:** Пользователь может попросить модель повторить ключевые визуальные элементы перед ответом на сложный вопрос.

Ожидаемые результаты от применения этих подходов: - Снижение галлюцинаций, особенно в задачах, требующих рассуждения или извлечения знаний - Повышение точности ответов на вопросы о диаграммах, графиках, математических задачах - Более надежные ответы при работе с изображениями, содержащими текст или числовые данные

Таким образом, хотя исследователи использовали сложные технические методы для имплементации VDGD, концептуальный подход "сначала опиши, потом отвечай" является мощной техникой, доступной любому пользователю чат-моделей с мультимодальными возможностями.

**## Анализ практической применимости:** 1. **Идентификация причин галлюцинаций в LVLMS** - Прямая применимость: Пользователи могут лучше понимать, в каких ситуациях модели склонны к галлюцинациям, и соответствующим образом формулировать запросы. - Концептуальная ценность: Высокая, так как объясняет фундаментальные причины ошибок в работе LVLMS. - Потенциал для адаптации: Понимание различий между задачами распознавания и рассуждения позволит

пользователям выбирать оптимальные стратегии взаимодействия.

**Определение разрыва в визуальном восприятии** Прямая применимость: Средняя, требует технического понимания работы моделей. Концептуальная ценность: Очень высокая, объясняет почему модели дают некорректные ответы даже при правильном распознавании изображений. Потенциал для адаптации: Пользователи могут формулировать запросы так, чтобы компенсировать этот разрыв, например, прося модель сначала описать изображение.

### **Метод Visual Description Grounded Decoding (VDGD)**

Прямая применимость: Высокая для разработчиков, низкая для обычных пользователей (требует доступа к API). Концептуальная ценность: Значительная, демонстрирует эффективный подход к решению проблемы. Потенциал для адаптации: Пользователи могут имитировать этот метод, запрашивая у модели описание изображения перед основным вопросом.

### **Категоризация типов галлюцинаций**

Прямая применимость: Средняя, помогает распознавать типы ошибок. Концептуальная ценность: Высокая, предоставляет структуру для анализа ошибок. Потенциал для адаптации: Пользователи могут научиться определять типы галлюцинаций и корректировать свои запросы.

### **Бенчмарк VaLLu**

Прямая применимость: Низкая для обычных пользователей, высокая для исследователей. Концептуальная ценность: Средняя для широкой аудитории. Потенциал для адаптации: Ограниченный для непосредственного использования.

## **Prompt:**

Применение исследования о снижении галлюцинаций LVLM в промптах ## Ключевое понимание Исследование показывает, что большие визуально-языковые модели (LVLM) страдают от "разрыва визуального восприятия" - они могут видеть элементы изображения, но плохо интерпретируют их в контексте задачи, что приводит к галлюцинациям, особенно в когнитивных задачах.

## Пример промпта на основе VDGD метода

[=====] [Первый шаг: запрос детального описания] Сначала внимательно опиши это изображение, уделяя особое внимание всем визуальным элементам: что на нем изображено, какие объекты присутствуют, как они расположены, их характеристики и взаимосвязи. Опиши все детали, которые могут быть важны для понимания контекста.

[Второй шаг: основной вопрос] Теперь, основываясь на твоём собственном описании изображения, ответь на следующий вопрос: [здесь основной вопрос, требующий

рассуждения].

Важно: в своем ответе опирайся только на факты, которые ты действительно видишь на изображении и упомянул в своем описании. Если какой-то информации не хватает, укажи это вместо предположений. [=====]

## ## Почему это работает

Такой двухэтапный подход реализует принцип VDGD (Visual Description Grounded Decoding):

**Преодоление разрыва восприятия** - заставляя модель сначала создать детальное описание, мы помогаем ей лучше "увидеть" и зафиксировать все элементы изображения

**Привязка к фактам** - когда модель отвечает на вопрос во втором шаге, она уже имеет структурированное представление о том, что действительно присутствует на изображении

**Снижение всех типов галлюцинаций** - особенно эффективно для когнитивных задач, где обычные методы снижения галлюцинаций не работают

**Разделение восприятия и рассуждения** - позволяет модели сначала сосредоточиться на визуальном восприятии, а затем на связывании этой информации с внутренними знаниями

Этот подход особенно полезен для сложных изображений (графиков, диаграмм, технических схем) и задач, требующих глубокого понимания контекста.