

# LLM как испорченный телефон: итеративная генерация искажает информацию

Дата: 2025-02-27 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.20258>

Рейтинг: 75

Адаптивность: 85

## Ключевые выводы:

Исследование изучает, как LLM искажают информацию при итеративной обработке собственных выходных данных (эффект «испорченного телефона»). Основная цель - понять, как накапливаются искажения при многократной обработке текста через цепочки переводов. Результаты показывают, что искажение информации неизбежно накапливается с течением времени, причем степень искажения зависит от выбора языка, сложности цепочки и параметров генерации.

## Объяснение метода:

Исследование демонстрирует важный эффект искажения информации при итеративном использовании LLM и предлагает практические решения (низкая температура, ограничительные промпты). Результаты применимы для любого пользователя, особенно при многошаговых взаимодействиях. Некоторые аспекты технически сложны, но ключевые выводы доступны для непосредственного применения без специальных навыков.

**## Ключевые аспекты исследования: 1. Эффект "сломанного телефона" в LLM:** Исследование демонстрирует, что при итеративном использовании выходных данных LLM (когда результат одной генерации становится входом для следующей) происходит постепенное искажение информации, аналогично игре "сломанный телефон" у людей.

**Факторы, влияющие на искажение информации:** Степень искажения зависит от выбора промежуточных языков (их сходства с исходным языком), сложности цепочки (количества языков и моделей), и параметров генерации (температуры и ограничений в промпте).

**Методы снижения искажения:** Авторы выявили, что контроль температуры (низкие значения) и использование ограничительных промптов значительно снижают искажение информации при итеративной генерации.

**Количественная оценка искажения:** Исследование предлагает методологию для

измерения степени искажения с использованием метрик текстуальной релевантности (BLEU, ROUGE, METEOR и др.) и сохранения фактов (FActScore).

**Эксперименты с разными конфигурациями:** Проведены серии экспериментов с различными моделями (Llama, Mistral, Gemma), языками и структурами цепочек для понимания факторов, влияющих на искажение.

**## Дополнение:** Для работы методов данного исследования не требуется дообучение или API. Хотя авторы использовали различные модели и специальные метрики для оценки искажений, основные концепции и выводы исследования полностью применимы в стандартном чате с LLM.

Основные концепции, которые можно применить в стандартном чате:

**Минимизация итеративной обработки:** Понимание, что каждая последующая обработка текста моделью потенциально вносит искажения. Пользователь может избегать многократных перефразирований одного и того же контента.

**Контроль параметров генерации:** В большинстве чатов с LLM можно запросить модель использовать более "консервативный" подход к генерации (эквивалент низкой температуры). Например: "Пожалуйста, перефразируй этот текст, максимально сохраняя оригинальный смысл и все детали, без добавления новой информации."

**Ограничительные промпты:** Пользователь может самостоятельно создавать более ограничительные инструкции, например: "Переведи этот текст с русского на английский. Важно: сохрани все факты, имена и цифры без изменений; не добавляй и не удаляй информацию; сохрани тон и стиль оригинала."

**Выбор языковых пар:** При необходимости перевода пользователи могут предпочесть прямой перевод между языками, а не цепочку переводов через промежуточные языки.

**Периодическая сверка с источником:** При длительных взаимодействиях пользователь может периодически напоминать модели исходную информацию, чтобы минимизировать накопление искажений.

Применяя эти концепции, пользователи могут значительно снизить риск искажения информации при работе с LLM, особенно в задачах, требующих сохранения фактической точности и полноты информации - от перевода документов до суммирования важных текстов и создания контента на основе исходных данных.

**## Анализ практической применимости: Ключевой аспект 1: Эффект "сломанного телефона" в LLM - Прямая применимость:** Пользователи могут учитывать риск искажения при многократной обработке текста в LLM, особенно важно для профессионалов, использующих LLM для создания контента. - **Концептуальная ценность:** Помогает понять фундаментальное ограничение LLM - накопление ошибок при итеративной обработке, что критично для правильного использования

этих систем. - **Потенциал для адаптации:** Пользователи могут разработать стратегии периодической "сверки" с оригинальными источниками при длительных цепочках взаимодействия с LLM.

**Ключевой аспект 2: Факторы, влияющие на искажение информации - Прямая применимость:** Пользователи могут выбирать языки с меньшим искажением при необходимости перевода (латинские скрипты вместо нелатинских). - **Концептуальная ценность:** Понимание, что LLM сохраняют информацию лучше на языках, которые лучше представлены в их обучающих данных. - **Потенциал для адаптации:** Пользователи могут выстраивать цепочки взаимодействия с минимальным количеством "переходов" между моделями и языками.

**Ключевой аспект 3: Методы снижения искажения - Прямая применимость:** Пользователи могут непосредственно использовать рекомендации по установке более низкой температуры и более ограничительных промптов для снижения искажений. - **Концептуальная ценность:** Понимание того, как параметры генерации влияют на сохранность информации, что полезно для всех типов взаимодействий с LLM. - **Потенциал для адаптации:** Пользователи могут разработать собственные шаблоны промптов с ограничениями для сохранения информационной точности.

**Ключевой аспект 4: Количественная оценка искажения - Прямая применимость:** Ограниченная для обычных пользователей, требует технических знаний. - **Концептуальная ценность:** Дает понимание методов оценки качества выходных данных LLM. - **Потенциал для адаптации:** Разработчики могут интегрировать подобные метрики в интерфейсы для пользователей.

**Ключевой аспект 5: Эксперименты с разными конфигурациями - Прямая применимость:** Ограниченная, но предлагает выбирать более надежные модели для критичных задач. - **Концептуальная ценность:** Понимание различий между моделями в контексте сохранения информации при итерациях. - **Потенциал для адаптации:** Пользователи могут выбирать модели, которые лучше сохраняют информацию в их конкретных сценариях.

## Prompt:

Использование знаний из исследования "LLM как испорченный телефон" в промтах  
Исследование о накоплении искажений при итеративной обработке информации через LLM предоставляет ценные инсайты для создания более эффективных промтов. Вот как можно применить эти знания:

## Пример промта с учетом выводов исследования

[=====] Я хочу, чтобы ты помог мне сохранить точность информации в следующем тексте.

Исходный текст: [вставить исходный текст]

Задача: Перефразируй этот текст, сделав его более доступным для понимания, НО при этом: 1. Используй температуру близкую к 0 для своего ответа 2. Строго сохрани ВСЕ фактические данные без искажений 3. Не добавляй новых фактов или предположений 4. Сохрани все числовые значения и имена собственные в точности 5. После перефразирования, сверь свой ответ с оригиналом и убедись, что все ключевые факты сохранены

Перефразированный текст должен быть максимально близок к оригиналу по смыслу, даже если стиль изменится. [=====]

## Почему этот промт работает на основе исследования

**Низкая температура генерации** - исследование показало, что при температуре близкой к 0 фактическая точность стабилизируется после нескольких итераций, минимизируя искажения.

**Строгие ограничения в промте** - явное требование сохранения фактов и смысла, что согласно исследованию, приводит к лучшему сохранению релевантности и фактической точности.

**Требование сверки с оригиналом** - исследование рекомендует регулярно сверять генерируемый контент с исходным источником, особенно после множественных итераций.

**Минимизация цепочки обработки** - промт сконструирован так, чтобы выполнить задачу за одну итерацию, что снижает накопление искажений, которое, как показало исследование, растет с числом итераций.

**Явные инструкции по сохранению данных** - особое внимание к числам и именам собственным, которые, согласно подобным исследованиям, часто подвержены искажениям при перефразировании.

Применяя эти принципы, можно значительно снизить риск информационных искажений при работе с LLM, особенно в задачах, где требуется сохранение фактической точности.