

# От инструментов к товарищам по команде: оценка LLM в многосессионных взаимодействиях при кодировании

Дата: 2025-02-19 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.13791>

Рейтинг: 78

Адаптивность: 85

## Ключевые выводы:

Исследование оценивает способность языковых моделей (LLM) эффективно сотрудничать в долгосрочных взаимодействиях. Авторы создали датасет MEMORYCODE для тестирования способности моделей отслеживать и выполнять простые инструкции по кодированию в многосессионных диалогах. Результаты показывают, что даже современные модели (включая GPT-4o) значительно теряют эффективность при необходимости извлекать и интегрировать информацию из длинных цепочек инструкций, распределенных по нескольким сессиям.

## Объяснение метода:

Исследование выявляет критическое ограничение LLM — неспособность эффективно использовать информацию в длительных взаимодействиях, даже если задачи просты. Предоставляет конкретные данные о падении эффективности с ростом контекста (GPT-4o теряет 67% точности) и анализирует причины. Пользователи могут применить стратегии "освежения памяти", структурирования взаимодействий и явного указания на обновления инструкций.

## Ключевые аспекты исследования: 1. **Многосессионное взаимодействие с LLM:** Исследование представляет MEMORYCODE - набор данных, моделирующий многосессионное взаимодействие между человеком (ментором) и LLM (учеником), имитирующий реальные рабочие ситуации, где информация накапливается и обновляется со временем.

**Проблема долгосрочной памяти:** Авторы выявили фундаментальное ограничение современных LLM - их неспособность эффективно отслеживать и применять инструкции, полученные в течение длительного взаимодействия, даже когда сами инструкции просты.

**Оценка моделей:** Исследование тестирует различные модели (GPT-4o, Llama 3.1 и др.) на способность запоминать и выполнять инструкции по кодированию,

полученные в разных сессиях, среди нерелевантной информации.

**Падение производительности с ростом контекста:** Эксперименты показывают, что даже лучшие модели значительно теряют в эффективности при увеличении количества сессий - GPT-4o показывает падение точности на 67% при полной истории диалога.

**Анализ причин неэффективности:** Исследование выявляет, что проблема связана не только с извлечением информации из контекста, но и с неспособностью моделей рассуждать о цепочке инструкций и их обновлениях.

## Дополнение: Исследование не требует дообучения или специального API для применения его методов. Ученые использовали стандартные модели через API только для удобства тестирования, но выявленные ограничения и предложенные подходы полностью применимы в стандартном чате с LLM.

Основные концепции, которые можно применить в стандартном чате:

**Стратегия периодического резюмирования** - пользователи могут периодически просить модель суммировать ключевые инструкции и договоренности, достигнутые в предыдущих сессиях.

**Явное обновление инструкций** - при изменении требований, пользователи должны явно указывать, что это обновление предыдущей инструкции, а не новое требование.

**Структурирование взаимодействия** - разбивка длинных сессий на более короткие, с четкими задачами и минимумом нерелевантной информации.

**Проактивное напоминание** - периодическое напоминание модели о ключевых инструкциях, особенно при выполнении задач, где эти инструкции должны применяться.

**Фреймворк проверки** - пользователи могут создать систему проверки, регулярно тестируя, помнит ли модель важные детали из предыдущих сессий.

Применение этих концепций позволит значительно повысить эффективность длительных взаимодействий с LLM, преодолевая выявленное исследованием фундаментальное ограничение моделей. Это особенно ценно для рабочих сценариев, где взаимодействие с моделью происходит на протяжении длительного времени и требует сохранения контекста.

## Анализ практической применимости: 1. **Многосессионное взаимодействие с LLM** - Прямая применимость: Высокая. Исследование выявляет серьезные ограничения LLM в долгосрочных взаимодействиях, что критически важно для пользователей, работающих с моделями в течение продолжительного времени. - Концептуальная ценность: Очень высокая. Понимание, что LLM теряют эффективность при длительных взаимодействиях, помогает пользователям

корректировать свои ожидания и стратегии работы. - Потенциал для адаптации: Высокий. Пользователи могут разработать стратегии "освежения памяти" моделей, периодически напоминая ключевую информацию.

**Проблема долгосрочной памяти** Прямая применимость: Средняя. Пользователи могут минимизировать потерю информации, разбивая взаимодействие на короткие сессии с четкими инструкциями. Концептуальная ценность: Высокая. Понимание, что простые инструкции могут быть забыты в длинном контексте, помогает пользователям структурировать свои запросы. Потенциал для адаптации: Высокий. Пользователи могут создавать "сводки" предыдущих инструкций и периодически их предоставлять модели.

### **Оценка моделей**

Прямая применимость: Высокая. Результаты показывают, что даже самые продвинутые модели (GPT-4o) имеют ограничения, что важно учитывать при выборе модели для длительных взаимодействий. Концептуальная ценность: Средняя. Понимание разницы в производительности моделей разного размера помогает пользователям делать более информированный выбор. Потенциал для адаптации: Средний. Пользователи могут выбирать модели с учетом длительности планируемого взаимодействия.

### **Падение производительности с ростом контекста**

Прямая применимость: Очень высокая. Пользователи должны учитывать, что эффективность моделей значительно снижается при длинных диалогах, даже если контекст помещается в окно модели. Концептуальная ценность: Высокая. Осознание того, что проблема не в размере контекстного окна, а в способности обрабатывать информацию, меняет подход к использованию моделей. Потенциал для адаптации: Высокий. Пользователи могут разработать стратегии повышения эффективности, например, резюмирование предыдущих инструкций.

### **Анализ причин неэффективности**

Прямая применимость: Средняя. Понимание, что проблема связана с рассуждением, а не только с извлечением информации, помогает пользователям формулировать более эффективные запросы. Концептуальная ценность: Высокая. Это объясняет, почему простое предоставление информации не гарантирует ее использование моделью. Потенциал для адаптации: Высокий. Пользователи могут разработать методы упрощения рассуждений для модели, например, явно указывая на противоречия или обновления инструкций.

### **Prompt:**

Применение результатов исследования о многосессионных взаимодействиях с LLM  
## Ключевые выводы для промптинга

Исследование "От инструментов к товарищам по команде" ясно показывает, что

даже продвинутые LLM (включая GPT-4o) значительно теряют эффективность при работе с информацией, распределенной по длинным диалогам. Производительность моделей падает с >90% до примерно 10% точности при увеличении количества сессий.

## Пример эффективного промпта с учетом исследования

[=====] # Запрос на разработку функции для обработки данных

## Контекст и предыдущие инструкции - Мы разрабатываем систему обработки финансовых данных - Ранее мы договорились использовать pandas для обработки таблиц - Все функции должны включать подробные комментарии - Обработка ошибок должна быть реализована с помощью try-except - Производительность критична для больших наборов данных

## Текущая задача Необходимо создать функцию для очистки финансовых данных со следующими требованиями: 1. Функция должна принимать DataFrame с финансовыми транзакциями 2. Удалять дубликаты транзакций 3. Заменять отсутствующие значения в поле 'amount' на медианное значение 4. Конвертировать даты в стандартный формат ISO

## Формат ответа - Предоставьте полный код функции - Добавьте краткое описание работы функции - Укажите на возможные ограничения вашего решения [=====]

## Почему этот промпт эффективен с учетом исследования

**Консолидация информации:** Вместо того, чтобы полагаться на память модели о предыдущих инструкциях, промпт включает все важные ранее оговоренные требования.

**Структурированность:** Четкое разделение на секции помогает модели организовать свой ответ и не пропустить важные детали.

**Самодостаточность:** Промпт содержит всю необходимую информацию для выполнения задачи, не требуя от модели обращения к предыдущим сессиям.

**Явные ожидания:** Раздел "Формат ответа" устанавливает четкие критерии для генерируемого контента.

## Практическое применение

Данный подход особенно полезен при длительной работе над проектом, когда важно, чтобы модель следовала установленным ранее соглашениям и требованиям. Вместо ожидания, что LLM будет помнить все предыдущие инструкции (что, как показало исследование, ненадежно), лучше создавать промпты, содержащие сводку всей важной контекстной информации.