

МакГайвер: Являются ли большие языковые модели креативными решателями проблем?

Дата: 2025-02-22 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2311.09682>

Рейтинг: 85

Адаптивность: 90

Ключевые выводы:

Исследование направлено на оценку способностей современных языковых моделей (LLM) к творческому решению проблем в условиях ограничений. Авторы создали датасет MACGYVER, содержащий более 1600 реальных проблем, требующих нестандартного использования предметов. Основной вывод: хотя современные LLM (особенно GPT-4) демонстрируют определенные способности к творческому решению проблем, они все еще значительно отстают от коллективного человеческого интеллекта, особенно в понимании физических свойств предметов и их возможного применения.

Объяснение метода:

Исследование предоставляет готовые стратегии промптинга (итеративная рефлексия и дивергентно-конвергентное мышление), которые могут быть немедленно применены в стандартных чатах. Детальный анализ типичных ошибок LLM в физическом рассуждении дает пользователям концептуальную основу для критической оценки ответов. Особенно ценны выводы о дополняющих возможностях человека и LLM, способствующие более эффективному взаимодействию.

Ключевые аспекты исследования: 1. **MACGYVER Dataset** - новый набор данных из более 1600 практических проблем, требующих нестандартного использования предметов в ограниченных условиях. Задачи разработаны для оценки творческого решения проблем как у людей, так и у языковых моделей.

Сравнение людей и LLM - исследование показывает, что люди и LLM демонстрируют разные сильные стороны в решении задач: люди лучше справляются с задачами из знакомых областей, а LLM имеют более широкие, но менее глубокие знания в специализированных областях.

Анализ ошибок LLM - выявлены типичные ошибки моделей: предложение физически невыполнимых действий, неправильное использование инструментов, галлюцинации и игнорирование ограничений.

Техники улучшения производительности - предложены две стратегии промптинга: итеративная пошаговая рефлексия (проверка выполнимости каждого шага) и дивергентно-конвергентное мышление (анализ возможностей каждого предмета перед решением).

Оценка эффективности промптинг-стратегий - эксперименты показывают, что предложенные стратегии значительно улучшают способность LLM решать творческие задачи с ограничениями.

Дополнение: Методы, представленные в исследовании, не требуют дообучения моделей или специального API для их использования. Хотя авторы использовали API для своих экспериментов, предложенные подходы могут быть полностью реализованы в стандартном чате с LLM.

Основные концепции и подходы, которые можно применить в стандартном чате:

Дивергентно-конвергентное мышление. Эта стратегия включает два этапа: Сначала попросить LLM проанализировать каждый доступный предмет и его возможные применения (дивергентное мышление) Затем попросить модель использовать этот анализ для формирования решения (конвергентное мышление) Пример промпта: "Проанализируй возможные применения каждого из этих предметов: [список предметов]. Затем предложи решение проблемы [описание проблемы], используя эти предметы."

Итеративная пошаговая рефлексия. Этот метод можно реализовать через серию сообщений: Получить первоначальное решение от LLM Попросить модель проверить каждый шаг на физическую выполнимость Попросить модель улучшить решение с учетом выявленных проблем

Распознавание типичных ошибок. Пользователи могут проактивно запрашивать модель проверить свой ответ на наличие:

Физически невыполнимых действий Неправильного использования инструментов Использования недоступных инструментов Нарушения указанных ограничений В исследовании показано, что метод дивергентно-конвергентного мышления помогает улучшить результаты для всех протестированных моделей, включая менее мощные, чем GPT-4. Это делает его особенно ценным для широкого круга пользователей.

Применяя эти подходы, пользователи могут: - Получать более практически выполнимые решения - Снизить количество галлюцинаций в ответах LLM - Улучшить эффективность предлагаемых решений - Развивать более систематический подход к формулировке запросов

Анализ практической применимости: 1. **MACGYVER Dataset** - Прямая применимость: средняя. Хотя сам датасет имеет ограниченную прямую пользу для рядовых пользователей, примеры из него могут быть полезны для обучения пользователей формулировать нестандартные запросы. - Концептуальная ценность:

высокая. Понимание типов задач, которые бросают вызов LLM, помогает пользователям формировать более эффективные запросы. - Потенциал для адаптации: высокий. Принципы составления проблемных сценариев могут быть использованы для создания собственных задач.

Сравнение людей и LLM Прямая применимость: высокая. Понимание сильных и слабых сторон LLM помогает пользователям эффективнее использовать модели, зная, в каких областях LLM могут ошибаться. Концептуальная ценность: высокая. Осознание дополняющих возможностей человека и LLM способствует формированию более эффективных стратегий взаимодействия. Потенциал для адаптации: средний. Выводы о различиях между людьми и LLM могут быть применены к различным задачам.

Анализ ошибок LLM

Прямая применимость: очень высокая. Знание типичных ошибок моделей позволяет пользователям предвидеть и корректировать потенциальные проблемы. Концептуальная ценность: высокая. Понимание ограничений LLM в физическом рассуждении помогает реалистично оценивать возможности моделей. Потенциал для адаптации: высокий. Пользователи могут разрабатывать собственные стратегии проверки предложений LLM на основе выявленных типов ошибок.

Техники улучшения производительности

Прямая применимость: очень высокая. Предложенные стратегии промптинга могут быть непосредственно использованы пользователями в повседневных взаимодействиях с LLM. Концептуальная ценность: высокая. Методы демонстрируют важность итеративного подхода и систематического анализа при работе с LLM. Потенциал для адаптации: очень высокий. Стратегии могут быть модифицированы для различных типов задач.

Оценка эффективности промптинг-стратегий

Прямая применимость: высокая. Количественные результаты помогают пользователям выбрать наиболее эффективную стратегию для конкретной LLM. Концептуальная ценность: средняя. Данные подтверждают эффективность структурированного подхода к промптингу. Потенциал для адаптации: высокий. Методология оценки может быть применена к оценке других стратегий промптинга.

Prompt:

Применение исследования "МакГайвер" в промптах для GPT ## Ключевые выводы из исследования

Исследование показывает, что хотя современные языковые модели обладают некоторыми способностями к творческому решению проблем, они всё ещё значительно отстают от людей, особенно в понимании физических свойств предметов и их применения. Однако существуют стратегии промптинга, которые

могут значительно улучшить результаты.

Пример эффективного промпта на основе исследования

[=====] # Задача: Творческое решение проблемы с ограниченными ресурсами

Контекст Я оказался в следующей ситуации: [описание проблемы]. Доступные предметы: [список предметов]. Ограничения: [описание ограничений].

Инструкции (на основе исследования "МакГайвер")

Дивергентный этап: Для каждого доступного предмета перечисли 3-5 возможных нестандартных способов его использования, учитывая физические свойства.

Конвергентный этап: На основе перечисленных возможностей, предложи 3 различных решения проблемы.

Итеративная рефлексия: Для каждого решения:

Разбей его на конкретные шаги Проверь физическую выполнимость каждого шага
Укажи потенциальные проблемы Модифицируй решение для устранения этих проблем

Финальная оценка: Оцени каждое решение по:

Выполнимости (учитывая физические законы) Эффективности Безопасности
Надежности

Рекомендация: Выбери наиболее оптимальное решение и объясни свой выбор.
[=====]

Как это работает

Данный промпт использует ключевые стратегии, выявленные в исследовании:

Дивергентно-конвергентное мышление - сначала исследуются все возможные применения предметов, затем формулируются конкретные решения, что повышает эффективность на 6.5%.

Итеративная пошаговая рефлексия - проверка физической выполнимости каждого шага и модификация решения, что снижает количество невыполнимых решений на 9.7%.

Коллективное решение проблем - запрос нескольких вариантов решения, что имитирует коллективный интеллект.

Проверка физической выполнимости - акцент на физические свойства предметов, что помогает избежать основного источника ошибок (42.4% ошибок

связаны с неправильным использованием инструментов).

Такая структура промпта компенсирует слабые стороны LLM и максимизирует их творческий потенциал при решении практических задач.