

Decompose-ToM: Улучшение рассуждений о Теории Разума в больших языковых моделях через симуляцию и декомпозицию задач

Дата: 2025-01-15 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2501.09056>

Рейтинг: 75

Адаптивность: 85

Ключевые выводы:

Исследование направлено на улучшение способностей больших языковых моделей (LLM) к рассуждению с использованием теории разума (Theory of Mind, ToM). Авторы предлагают алгоритм Decompose-ToM, который значительно улучшает производительность LLM на сложных задачах ToM, особенно на задачах высшего порядка, без необходимости дополнительного обучения моделей.

Объяснение метода:

Исследование предлагает ценные принципы декомпозиции сложных задач и симуляции перспектив, которые могут быть адаптированы обычными пользователями для улучшения взаимодействия с LLM. Методы не требуют дополнительного обучения моделей, но полная реализация алгоритма технически сложна.

Ключевые аспекты исследования: 1. **Метод Decompose-ToM** – алгоритм, улучшающий способности языковых моделей (LLM) в задачах "теории разума" (Theory of Mind, ToM) путем декомпозиции сложных задач на более простые подзадачи.

Рекурсивная симуляция перспектив – техника, при которой модель последовательно "представляет себя" на месте каждого персонажа, чтобы понять их осведомленность и убеждения.

Декомпозиция задачи на подзадачи – разбиение сложных задач ToM на более простые компоненты: идентификация субъекта, переформулировка вопроса, обновление модели мира и оценка доступности информации.

Управление знаниями агентов – метод определения, какой информацией владеет каждый агент в зависимости от контекста (где они находились, что видели и слышали).

Значительное улучшение результатов – особенно в сложных задачах высокого порядка ToM (когда нужно определить, что думает А о том, что думает В о том, что думает С).

Дополнение: Для работы методов этого исследования **не требуется** дообучение или специальное API. Авторы используют стандартные LLM через обычное API без дополнительного обучения, что явно указано в статье: "...не требуя минимальной настройки промптов для разных задач и без дополнительного обучения модели".

Концепции, которые можно применить в стандартном чате:

Последовательная симуляция перспектив – можно просить чат-модель "представить себя" на месте определенного человека перед ответом на вопрос. Например: "Представь, что ты Алиса, которая не знает, что Боб переложил конфету. Где бы ты искала конфету?"

Декомпозиция сложного вопроса – разбивка сложных вопросов на серию простых шагов. Например, вместо "Что думает А о том, что думает В о намерениях С?", можно сначала спросить "Что знает В о намерениях С?", а затем "Зная это, что может думать А о мнении В?"

Явное отслеживание доступности информации – можно просить модель перечислить, какой информацией владеет каждый персонаж в истории, прежде чем делать выводы об их убеждениях.

Обновление модели мира – можно явно просить модель отслеживать, где находятся персонажи и какой информацией они владеют на каждом этапе истории.

Результаты применения этих концепций: - Улучшение ответов в задачах, требующих понимания разных точек зрения - Более точное моделирование убеждений людей с неполной информацией - Лучшее понимание мотивов персонажей в историях - Более эмпатичные и нюансированные анализы межличностных ситуаций - Повышение точности в задачах, требующих отслеживания, кто какой информацией владеет

Анализ практической применимости: **1. Метод Decompose-ToM - Прямая применимость:** Средняя. Пользователи могут применить концепцию декомпозиции сложных запросов на более простые, но полная реализация алгоритма требует технических навыков. - **Концептуальная ценность:** Высокая. Демонстрирует, что сложные рассуждения можно разбить на простые шаги, что улучшает понимание пользователями принципов эффективного взаимодействия с LLM. - **Потенциал для адаптации:** Высокий. Принцип "разбивай сложное на простое" можно применять к широкому спектру задач.

2. Рекурсивная симуляция перспектив - Прямая применимость: Средняя. Пользователи могут напрямую просить LLM "поставить себя на место"

определенного человека или персонажа для лучшего понимания его точки зрения. - **Концептуальная ценность:** Очень высокая. Помогает понять, что LLM может эффективнее решать задачи с множеством перспектив, если явно просить его рассматривать ситуацию с разных точек зрения. - **Потенциал для адаптации:** Высокий. Подход применим для анализа мотивов персонажей, прогнозирования реакций людей, улучшения эмпатии в текстах.

3. Декомпозиция задачи на подзадачи - Прямая применимость: Высокая. Пользователи могут использовать пошаговый подход при формулировке сложных запросов к LLM. - **Концептуальная ценность:** Высокая. Демонстрирует, что сложные задачи лучше решаются при последовательной декомпозиции. - **Потенциал для адаптации:** Очень высокий. Принцип применим практически ко всем сложным запросам к LLM.

4. Управление знаниями агентов - Прямая применимость: Низкая. Требует технических навыков для полной реализации. - **Концептуальная ценность:** Средняя. Полезна для понимания, как LLM может отслеживать доступность информации для разных персонажей. - **Потенциал для адаптации:** Средний. Может быть адаптирована для задач, где важно понимать, кто какой информацией владеет.

5. Улучшение результатов в сложных задачах ТоМ - Прямая применимость: Средняя. Обычные пользователи могут применять упрощенные версии метода для улучшения ответов LLM в задачах, требующих понимания разных точек зрения. - **Концептуальная ценность:** Высокая. Показывает, что даже сложные задачи "теории разума" решаемы при правильном подходе. - **Потенциал для адаптации:** Высокий. Методы могут быть адаптированы для различных задач социального анализа и моделирования.

Prompt:

Применение исследования Decompose-ToM в промптах для GPT ## Ключевые концепции для эффективных промптов

Исследование Decompose-ToM предлагает несколько мощных техник для улучшения способности языковых моделей рассуждать с использованием теории разума (ToM). Вот как можно применить эти знания в промптах:

Рекурсивная симуляция перспектив - моделирование точки зрения каждого участника **Декомпозиция сложных задач** - разбиение на подзадачи **Отслеживание доступности знаний** - учет того, что знает каждый участник **Символическое представление состояния мира** - структурированное отслеживание изменений **Пошаговое рассуждение** - применение Chain-of-Thought ## Пример промпта с применением Decompose-ToM

[=====] # Задача анализа сложной социальной ситуации

Контекст Алиса рассказала Борису, что планирует сделать сюрприз Виктории на

день рождения. Борис случайно упомянул при Денисе о подготовке сюрприза, но не сказал для кого. Позже Денис встретил Викторию и сказал: "Я слышал, тебя ждет какой-то сюрприз".

Инструкции для анализа (используя Decompose-ToM)

Идентификация субъектов и их перспектив: Перечисли всех участников ситуации
Для каждого участника определи исходную информацию

Символическое представление знаний:

Для каждого участника создай структуру: [что знает X о том, что знает Y о Z]

Симуляция перспектив по порядку:

Симулируй мысли Алисы: "Я знаю о сюрпризе, Борис знает, Виктория не знает"
Симулируй мысли Бориса: "Я знаю о сюрпризе, Алиса знает, Виктория не знает, Денис знает о сюрпризе, но не знает для кого"
Симулируй мысли Дениса: "Я знаю о каком-то сюрпризе, но не знаю для кого"
Симулируй мысли Виктории после разговора с Денисом

Анализ последствий:

Оцени, как изменилась ситуация после действий каждого участника
Определи, испорчен ли сюрприз и почему
Проведи анализ этой ситуации, используя вышеописанный подход. [=====]

Как работает этот подход

Данный промпт применяет ключевые принципы Decompose-ToM:

Декомпозиция - разбивает сложную социальную ситуацию на четкие этапы анализа
Симуляция перспектив - предлагает модели "встать на место" каждого участника
Отслеживание знаний - явно моделирует, кто что знает на каждом этапе
Символическое представление - структурирует знания в виде вложенных представлений
Такой подход позволяет модели более точно анализировать сложные социальные взаимодействия, особенно когда речь идет о нескольких уровнях понимания (например, "что X думает о том, что Y знает о Z").