

Метод Self Refine: Самопроверка и итеративное улучшение ответов

Self Refine (Самоуточнение) - это метод промпт-инжиниринга, использующий итеративное улучшение ответов через самокритику. Основная идея метода заключается в том, что модель сначала дает ответ, затем оценивает его недостатки и предлагает улучшенную версию.

Как работает Self Refine:

1. **Первоначальная генерация ответа:** модель создает начальный ответ на вопрос
2. **Критическая оценка:** модель анализирует свой ответ на наличие недостатков, неточностей или ошибок
3. **Улучшенная версия:** на основе выявленных недостатков модель генерирует улучшенный ответ
4. **Потенциально повторяющийся процесс:** при необходимости этапы 2 и 3 могут повторяться несколько раз

Пример промпта для Self Refine:

Сначала дай ответ на следующий вопрос: [ваш вопрос]

После того как ты дал первый ответ, пожалуйста:

1. Критически оцени свой ответ, выявив недостатки, неточности или пропущенную информацию
2. Предложи улучшенную версию ответа, учитывающую все выявленные недостатки
3. Укажи, какие конкретно коррекции были внесены и почему они важны

Почему это работает:

Self Refine работает за счет активации двух когнитивных систем мышления в языковых моделях:

- **System 1** (быстрое, интуитивное мышление) используется для начального ответа
- **System 2** (медленное, аналитическое мышление) активируется при самопроверке и улучшении

Исследования показывают, что такой подход значительно повышает качество ответов по сравнению с одноэтапной генерацией. Вы получаете:

- Повышение точности при решении сложных задач
- Снижение "галлюцинаций" и фактических ошибок
- Более структурированные и полные ответы
- Возможность видеть процесс самокоррекции модели, что повышает доверие к ответу

Self Refine особенно эффективен для задач, требующих высокой точности, детализации или логических рассуждений, таких как решение математических задач, написание кода или анализ сложных текстов.