

# Уверенность улучшает самосогласованность в больших языковых моделях

Дата: 2025-02-10 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.06233>

Рейтинг: 82

Адаптивность: 90

## Ключевые выводы:

Исследование направлено на улучшение метода Self-consistency для LLM путем внедрения оценки уверенности модели в собственных ответах. Основным результатом - предложенный метод Confidence Informed Self-Consistency (CISC) достигает сравнимой точности при снижении вычислительных затрат на более чем 40% в среднем.

## Объяснение метода:

CISC - практичный метод улучшения взаимодействия с LLM, требующий минимальных изменений в промптах. Позволяет сократить количество запросов на 40% при сохранении точности. Легко адаптируется к разным задачам и моделям. Предлагает несколько способов извлечения уверенности, включая простые вербальные оценки. Может быть реализован обычными пользователями без технических навыков через стандартные интерфейсы чатов.

## Ключевые аспекты исследования: 1. **Confidence-Informed Self-Consistency (CISC)** - исследование представляет новый метод повышения эффективности стандартного подхода Self-Consistency. CISC использует самооценку LLM о правильности своих рассуждений для взвешенного голосования при выборе итогового ответа.

**Снижение вычислительных затрат** - CISC позволяет достичь той же точности, что и стандартный Self-Consistency, но с использованием на 40% меньше вычислительных ресурсов (меньше сгенерированных цепочек рассуждений).

**Метод P(True)** - наиболее эффективный способ извлечения уверенности модели, когда LLM оценивает правильность своего ответа в бинарном формате (0 или 1).

**Within-Question Discrimination (WQD)** - новая метрика для оценки способности модели различать правильные и неправильные ответы на один и тот же вопрос.

**Температурное масштабирование** - нормализация уверенности с помощью

функции softmax и настраиваемого параметра температуры существенно улучшает эффективность CISC.

## Дополнение: Исследование CISC не требует дообучения или специального API для основной реализации. Ключевые элементы метода могут быть применены в стандартном чате с любой LLM без модификации самой модели.

Для реализации в стандартном чате можно использовать следующие концепции:

**Генерация нескольких решений** - пользователь может попросить модель решить одну задачу несколько раз (например, 5-10 раз) с инструкцией "решай эту задачу независимо каждый раз".

**Самооценка уверенности** - для каждого решения пользователь может запросить: "Оцени свою уверенность в этом ответе по шкале от 0 до 100" или "Считаешь ли ты этот ответ правильным? Ответь 'да' или 'нет'".

**Взвешенное голосование** - пользователь может представить модели все полученные решения с оценками уверенности и попросить: "На основе этих решений и оценок уверенности, какой ответ наиболее вероятно правильный?"

Хотя исследование использовало софтмакс-нормализацию с оптимальной температурой для взвешивания, даже простое взвешенное голосование без сложной нормализации дает значительное улучшение над стандартным подходом Self-Consistency.

Основной результат, который можно получить от применения этих концепций - более точные ответы при меньшем количестве запросов к модели, что экономит время и ресурсы. Например, вместо генерации 30 ответов для надежного результата может оказаться достаточным 10-15 ответов с оценкой уверенности.

Метод особенно полезен для сложных задач рассуждения, таких как математические задачи, логические головоломки и задачи, требующие многошаговых рассуждений.

## Анализ практической применимости: **1. Confidence-Informed Self-Consistency (CISC) - Прямая применимость:** Высокая. Пользователи могут сразу внедрить этот метод в свои взаимодействия с LLM, запрашивая модель оценить свою уверенность после решения задачи. - **Концептуальная ценность:** Значительная. Демонстрирует, что LLM способны оценивать качество своих собственных рассуждений. - **Потенциал для адаптации:** Высокий. Метод применим к широкому спектру задач рассуждения и может быть реализован через простые изменения в запросах.

**2. Снижение вычислительных затрат - Прямая применимость:** Очень высокая. Позволяет пользователям получать более качественные ответы при меньшем количестве запросов к API, что экономит время и деньги. - **Концептуальная ценность:** Средняя. Показывает компромисс между количеством запросов и

качеством результата. - **Потенциал для адаптации:** Высокий. Принцип может быть применен к различным задачам и моделям.

**3. Метод P(True) - Прямая применимость:** Высокая. Простой и эффективный способ получения оценки уверенности модели. - **Концептуальная ценность:** Значительная. Показывает, что бинарная оценка уверенности может быть эффективнее, чем более детализированные шкалы. - **Потенциал для адаптации:** Высокий. Легко реализуется через дополнительные запросы к модели.

**4. Within-Question Discrimination (WQD) - Прямая применимость:** Низкая. Это скорее исследовательская метрика, чем практический инструмент для обычных пользователей. - **Концептуальная ценность:** Высокая. Помогает понять, как оценивать способность модели к самопроверке. - **Потенциал для адаптации:** Средний. Может быть использован для выбора наилучшего метода извлечения уверенности.

**5. Температурное масштабирование - Прямая применимость:** Средняя. Требуется дополнительная настройка, но может быть реализовано. - **Концептуальная ценность:** Высокая. Показывает важность нормализации оценок уверенности. - **Потенциал для адаптации:** Высокий. Подход может быть использован в различных контекстах взвешенного голосования.

## Prompt:

Применение исследования CISC в промптах для GPT ## Ключевая идея исследования

Исследование показывает, что использование **оценки уверенности модели** при генерации нескольких ответов позволяет достичь лучших результатов с меньшими вычислительными затратами. Вместо простого подсчета частоты ответов (стандартный Self-consistency), метод CISC применяет взвешенное голосование, учитывая уверенность модели в каждом ответе.

## Пример промпта с применением CISC

[=====] Решите следующую математическую задачу: [ТЕКСТ ЗАДАЧИ]

Пожалуйста, выполните следующие шаги: 1. Предложите 5 различных подходов к решению этой задачи 2. Для каждого подхода: - Подробно опишите ход рассуждений - Получите ответ - Оцените вашу уверенность в правильности этого ответа по шкале от 0 до 100% - Объясните, почему вы присвоили именно такую оценку уверенности 3. В конце сделайте взвешенное голосование: - Перечислите все полученные ответы - Умножьте частоту каждого ответа на среднюю уверенность в этом ответе - Выберите ответ с наивысшим взвешенным показателем

Это очень важная задача, поэтому, пожалуйста, будьте максимально точны в своих рассуждениях и честны в оценке уверенности. [=====]

## ## Почему это работает

**Множественные пути решения:** Промпт запрашивает несколько независимых решений одной задачи (Self-consistency) **Оценка уверенности:** Для каждого решения модель оценивает свою уверенность (ключевой элемент CISC) **Взвешенное голосование:** Итоговый ответ выбирается не просто по частоте, а с учетом уверенности модели в каждом решении ## Преимущества подхода

- Экономия ресурсов: Согласно исследованию, можно получить такую же точность с меньшим количеством генераций (на 40-46% меньше)
- Улучшение качества: Взвешенное голосование помогает отфильтровать случайные ошибки, отдавая предпочтение решениям, в которых модель более уверена
- Прозрачность: Пользователь видит не только итоговый ответ, но и уровень уверенности модели в различных подходах

Такой подход особенно эффективен для задач, требующих сложных рассуждений, таких как математические задачи, логические головоломки или многошаговые процессы принятия решений.