

# Калибровка уверенности LLM с помощью семантического управления: рамочная система агрегирования многоподказок

Дата: 2025-03-04 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.02863>

Рейтинг: 72

Адаптивность: 85

## Ключевые выводы:

Исследование направлено на улучшение калибровки уверенности больших языковых моделей (LLM) через систематическое управление промптами. Основной результат: разработан фреймворк **SteeringConf**, который успешно корректирует уверенность LLM в предсказаниях, опровергая предыдущие утверждения о невозможности систематического управления уверенностью моделей через лингвистические вмешательства.

## Объяснение метода:

Исследование предлагает практически применимые методы управления уверенностью LLM через простые инструкции. Базовый принцип "будь осторожен/уверен" может быть непосредственно использован широкой аудиторией для получения более надежных ответов. Полная реализация методологии требует технических знаний, но основные концепты доступны обычным пользователям и повышают понимание работы LLM.

## Ключевые аспекты исследования: 1. **Метод управления уверенностью (Confidence Steering)** - исследование доказывает, что с помощью специальных подсказок можно направленно изменять оценку уверенности LLM в своих ответах (от "будь очень осторожен" до "будь очень уверен").

**Агрегация направленной уверенности (Steered Confidence Aggregation)** - метод объединяет несколько оценок уверенности, полученных с разными подсказками, для создания более калиброванной итоговой оценки.

**Выбор ответа на основе калиброванной уверенности (Steered Answer Selection)** - система выбирает наиболее подходящий ответ из нескольких вариантов, полученных с разными подсказками, на основе близости к расчетной калиброванной уверенности.

**Метрики согласованности ответов и уверенности** - авторы используют согласованность ответов и согласованность оценок уверенности как индикаторы надежности прогнозов модели.

**Экспериментальное подтверждение** - исследование показывает, что метод SteeringConf значительно улучшает калибровку уверенности и обнаружение ошибок на семи различных тестовых наборах.

## Дополнение:

### Применимость методов в стандартном чате

Исследование не требует дообучения моделей или специального API для реализации основных методов. Ключевые концепции можно применить в стандартном чате:

**Управление уверенностью через инструкции** - пользователи могут добавлять фразы "будь очень осторожен" или "будь очень уверен" к своим запросам, чтобы получать более консервативные или уверенные ответы.

**Проверка согласованности ответов** - пользователи могут задать один и тот же вопрос несколько раз с разными формулировками или уровнями запрашиваемой уверенности, чтобы оценить согласованность ответов как индикатор надежности.

**Явный запрос уверенности** - пользователи могут запрашивать модель оценить свою уверенность в ответе в процентах или по шкале от 0 до 100.

**Комбинирование консервативных и уверенных подходов** - пользователи могут сравнивать ответы, полученные с инструкциями "будь очень осторожен" и "будь очень уверен", чтобы выявить возможные расхождения и оценить надежность информации.

Результаты применения: - Более точная оценка надежности информации - Снижение риска принятия решений на основе неверной информации - Лучшее понимание ограничений модели в конкретных областях знаний - Возможность выявления противоречивых или неоднозначных ответов

## Анализ практической применимости: 1. **Метод управления уверенностью** - Прямая применимость: Высокая. Пользователи могут непосредственно использовать инструкции типа "будь очень осторожен" или "будь очень уверен" при взаимодействии с LLM для получения более консервативных или уверенных оценок. - Концептуальная ценность: Значительная. Понимание того, что LLM реагируют на инструкции по калибровке уверенности, помогает пользователям формировать запросы с осознанием возможности влиять на уверенность модели. - Потенциал для адаптации: Высокий. Принцип направленного управления уверенностью может быть адаптирован для разных моделей и задач.

**Агрегация направленной уверенности** Прямая применимость: Средняя. Требует нескольких запросов к модели, что усложняет процесс для обычных пользователей, но принцип агрегации может быть реализован в интерфейсах. Концептуальная ценность: Высокая. Демонстрирует, что комбинирование нескольких оценок уверенности дает более надежный результат. Потенциал для адаптации: Высокий. Принцип может быть адаптирован для различных сценариев, где требуется повышенная надежность.

### **Выбор ответа на основе калиброванной уверенности**

Прямая применимость: Низкая для обычных пользователей, так как требует дополнительной обработки. Концептуальная ценность: Средняя. Показывает, как можно выбирать наиболее надежные ответы из множества вариантов. Потенциал для адаптации: Средний. Может быть реализован в приложениях, но сложен для ручного использования.

### **Метрики согласованности**

Прямая применимость: Низкая. Требует технических знаний для реализации. Концептуальная ценность: Высокая. Понимание того, что согласованность ответов при разных подсказках указывает на надежность, может помочь пользователям оценивать достоверность информации. Потенциал для адаптации: Средний. Концепция может быть упрощена до проверки ответов с разными формулировками запроса.

### **Экспериментальное подтверждение**

Прямая применимость: Низкая. Результаты экспериментов сами по себе не применимы напрямую. Концептуальная ценность: Высокая. Доказывает эффективность подхода и дает понимание его ограничений. Потенциал для адаптации: Высокий. Результаты могут направлять разработку пользовательских интерфейсов и инструментов.

## **Prompt:**

Применение исследования о калибровке уверенности LLM в промтах для GPT ##  
Ключевая идея исследования

Исследование SteeringConf показывает, что можно систематически управлять уверенностью языковых моделей через специальные промты, варьирующие от "очень осторожных" до "очень уверенных", и затем агрегировать результаты для получения более точных и надежных ответов.

## Пример промпта с применением SteeringConf

[=====] Я хочу получить максимально точный ответ на вопрос медицинского характера. Для этого я прошу тебя:

Сначала ответь на мой вопрос, будучи **ОЧЕНЬ ОСТОРОЖНЫМ**. Отметь свой уровень уверенности по шкале от 1 до 10 и укажи, какие аспекты вопроса вызывают у тебя неуверенность.

Затем ответь на тот же вопрос, будучи **НЕЙТРАЛЬНЫМ**. Снова оцени уверенность по шкале от 1 до 10.

Наконец, ответь на вопрос, будучи **ОЧЕНЬ УВЕРЕННЫМ**. Оцени уверенность по шкале от 1 до 10.

Сравни свои ответы и сделай заключение, какой из них наиболее надежен и почему. Укажи, где могут быть ошибки или неточности.

Мой вопрос: Может ли длительное употребление аспирина привести к проблемам с почками? [=====]

## Как работает этот подход

**Управление уверенностью** — Запрашивая модель дать ответы с разным уровнем осторожности, мы получаем спектр ответов с разной степенью уверенности.

**Агрегация ответов** — Сравнивая согласованность между различными ответами и их уровнями уверенности, мы можем выявить, насколько модель действительно "знает" ответ.

**Выбор оптимального ответа** — Модель сама анализирует результаты и определяет, какой уровень уверенности наиболее обоснован для данного вопроса.

**Обнаружение ошибок** — Если ответы сильно различаются или уровень уверенности не соответствует содержанию, это сигнализирует о возможных ошибках.

Такой подход особенно полезен в областях, где точность и правильная оценка неопределенности критически важны: медицина, право, финансы и другие сферы, где ошибки могут иметь серьезные последствия.