

# Картирование надежности в больших языковых моделях: библиометрический анализ, связывающий теорию с практикой

Дата: 2025-02-27 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.04785>

Рейтинг: 62

Адаптивность: 75

## Ключевые выводы:

Исследование направлено на преодоление разрыва между теоретическими дискуссиями о доверии к LLM и практической реализацией. Основная цель - картировать и систематизировать понимание доверия к LLM через библиометрический анализ 2006 публикаций (2019-2025) и систематический обзор 68 ключевых работ. Главные результаты показывают, что доверие к LLM часто формулируется через существующие организационные структуры доверия, но существует значительный разрыв между теоретическими принципами и конкретными стратегиями разработки.

## Объяснение метода:

Исследование предоставляет ценную концептуальную основу для понимания доверия к LLM (модель ABI) и 20 стратегий повышения доверия. Несмотря на то, что большинство стратегий ориентированы на разработчиков, некоторые (инженерия промптов, человек-в-цикле) доступны обычным пользователям. Концепции калибровки доверия помогают избегать как чрезмерного, так и недостаточного доверия к LLM. Требуется адаптация для практического применения.

Ключевые аспекты исследования: 1. **Библиометрический анализ исследований по доверию к LLM:** Исследование анализирует 2006 публикаций (2019-2025) о доверии и этике в крупных языковых моделях через сетевой анализ соавторства, совместное появление ключевых слов и отслеживание тематической эволюции.

**Определения доверия к LLM:** Авторы систематизировали различные определения доверия к LLM, выявив, что треть определений основана на организационной модели доверия Мейера (способность, доброжелательность, целостность) с адаптацией к контексту AI.

**Практические стратегии повышения доверия:** Исследование предлагает 20 конкретных методов для повышения доверия к LLM на разных этапах жизненного

цикла, включая RAG, методы объяснимости и аудит после обучения.

**Разрыв между теорией и практикой:** Авторы выявили значительный разрыв между теоретическими принципами и практическими рекомендациями по обеспечению доверия к LLM, что затрудняет их реальное внедрение.

**Эволюция исследований доверия к LLM:** Исследование показывает, как дискуссии об этике ИИ трансформировались в обсуждения доверия к LLM, при этом часто не предлагая конкретных реализаций.

Дополнение:

## **# Применимость методов в стандартном чате**

Большинство методов, описанных в исследовании, не требуют дообучения или API для базового применения в стандартном чате. Хотя разработчики используют более сложные реализации для максимальной эффективности, пользователи могут адаптировать ключевые концепции:

### **## Адаптируемые концепции в стандартном чате:**

**Retrieval-Augmented Generation (RAG):** Пользователь может имитировать RAG, предоставляя модели дополнительный контекст в промпте. Пример: "Используя следующую информацию [вставка текста/данных], ответь на вопрос..."

#### **Инженерия промптов:**

Непосредственно применима в стандартном чате. Техники: цепочка рассуждений (Chain of Thought), разбиение задачи на подзадачи.

#### **Объяснимость и прозрачность:**

Запрос модели объяснить свои рассуждения. Пример: "Объясни свой процесс мышления при формировании этого ответа".

#### **Выражение неопределенности:**

Запрос модели указать уровень уверенности. Пример: "Укажи степень уверенности в каждом пункте твоего ответа".

#### **Человек в цикле:**

Пользователи могут итеративно улучшать ответы, давая обратную связь. Пример: "Твой ответ содержит неточность в пункте X. Пожалуйста, исправь и улучши его".

### **## Результаты применения:**

Эти адаптированные подходы помогают: - Повысить надежность ответов через предоставление дополнительного контекста - Улучшить понимание процесса рассуждения модели - Выявить потенциальные области неопределенности - Калибровать соответствующий уровень доверия к ответам LLM

Хотя эти адаптации менее мощны, чем полноценные технические реализации, они значительно улучшают взаимодействие с LLM в стандартном чате.

Анализ практической применимости: **Библиометрический анализ исследований по доверию к LLM**: - Прямая применимость: Низкая. Сам анализ представляет академический интерес, но не дает готовых инструментов для пользователей. - Концептуальная ценность: Высокая. Дает представление о тенденциях и ключевых темах в исследованиях доверия к LLM, помогая понять, на что обращать внимание. - Потенциал для адаптации: Средний. Понимание исследовательских тенденций может помочь пользователям отслеживать наиболее перспективные направления.

**Определения доверия к LLM**: - Прямая применимость: Средняя. Понимание компонентов доверия (способность, доброжелательность, целостность) помогает критически оценивать ответы LLM. - Концептуальная ценность: Высокая. Предоставляет основу для оценки надежности взаимодействия с LLM. - Потенциал для адаптации: Высокий. Пользователи могут применять эти критерии для оценки ответов LLM и соответствующей корректировки своих запросов.

**Практические стратегии повышения доверия**: - Прямая применимость: Средняя. Большинство стратегий (16 из 20) ориентированы на разработчиков, но некоторые доступны пользователям (например, инженерия промптов). - Концептуальная ценность: Высокая. Понимание того, какие техники повышают доверие, позволяет пользователям задавать более эффективные запросы. - Потенциал для адаптации: Высокий. Знание о техниках, таких как RAG, объяснимость и инженерия промптов, может быть адаптировано даже обычными пользователями.

**Разрыв между теорией и практикой**: - Прямая применимость: Низкая. Понимание этого разрыва само по себе не дает практических инструментов. - Концептуальная ценность: Высокая. Осознание ограничений текущих подходов к доверию помогает пользователям быть более критичными. - Потенциал для адаптации: Средний. Понимание ограничений может мотивировать пользователей разрабатывать собственные стратегии взаимодействия.

**Эволюция исследований доверия к LLM**: - Прямая применимость: Низкая. Историческое развитие исследований имеет ограниченную практическую ценность. - Концептуальная ценность: Средняя. Контекст помогает понять текущее состояние исследований. - Потенциал для адаптации: Низкий. Историческая эволюция мало что дает для практического применения.

Сводная оценка полезности: Предварительная оценка: 65 из 100

Исследование имеет высокую полезность для широкой аудитории благодаря: 1)

Систематизации 20 практических стратегий повышения доверия к LLM, некоторые из которых доступны обычным пользователям (инженерия промптов, человек-в-цикле)  
2) Предоставлению четкой концептуальной основы для оценки доверия к LLM через модель ABI (способность, доброжелательность, целостность) 3) Выявлению проблемы калибровки доверия, помогающей пользователям избегать как чрезмерного, так и недостаточного доверия к LLM

Контраргументы к оценке:

Почему оценка могла бы быть выше: - Исследование предоставляет комплексный обзор проблемы доверия к LLM, что может существенно улучшить понимание пользователями ограничений и возможностей этих моделей - Некоторые стратегии (RAG, объяснимость) могут быть адаптированы пользователями через правильно сформулированные запросы

Почему оценка могла бы быть ниже: - Большинство стратегий (16 из 20) ориентированы на разработчиков, а не на конечных пользователей - Исследование имеет преимущественно академический характер и не предлагает готовых инструментов для немедленного применения - Многие концепции требуют значительной технической адаптации для использования обычными пользователями

После рассмотрения этих аргументов, скорректированная оценка: 62 из 100.

Эта оценка отражает высокую концептуальную ценность исследования при ограниченной прямой применимости для широкой аудитории. Исследование дает ценную основу для понимания доверия к LLM, но требует дополнительной адаптации для практического применения.

Уверенность в оценке: Очень сильная. Исследование предоставляет четкую структуру определений доверия к LLM и конкретные стратегии повышения доверия, что позволяет точно оценить его практическую применимость. Библиометрический анализ дает надежную основу для понимания текущего состояния исследований в этой области. Четкое разделение стратегий на те, что доступны разработчикам и пользователям, позволяет точно оценить их полезность для широкой аудитории.

Оценка адаптивности: Оценка адаптивности: 75 из 100

Исследование демонстрирует высокую адаптивность благодаря:

1) Концептуальной модели доверия (ABI: способность, доброжелательность, целостность), которую пользователи могут применять для оценки ответов LLM в любом чате

2) Описанию стратегий, четыре из которых могут быть непосредственно использованы пользователями: инженерия промптов, человек-в-цикле, обратная связь от заинтересованных сторон и обучение с подкреплением от человеческой обратной связи

3) Понятию калибровки доверия, которое может быть применено пользователями для оценки их собственного уровня доверия к LLM (избегание как чрезмерного, так и недостаточного доверия)

4) Возможности адаптировать стратегии, ориентированные на разработчиков (например, RAG, объяснимость), через специфические запросы, требующие от модели поиска дополнительной информации или объяснения своих ответов

Хотя многие технические аспекты исследования требуют специальных знаний, основные концепции и некоторые стратегии могут быть адаптированы для использования в стандартном чате, что делает исследование достаточно перспективным для широкой аудитории.

|| <Оценка: 62> || <Объяснение: Исследование предоставляет ценную концептуальную основу для понимания доверия к LLM (модель ABI) и 20 стратегий повышения доверия. Несмотря на то, что большинство стратегий ориентированы на разработчиков, некоторые (инженерия промптов, человек-в-цикле) доступны обычным пользователям. Концепции калибровки доверия помогают избегать как чрезмерного, так и недостаточного доверия к LLM. Требуется адаптация для практического применения.> || <Адаптивность: 75>

## **Prompt:**

Использование исследования о доверии к LLM в промптах для GPT

### **Ключевые аспекты исследования для промптов**

Исследование "Картирование надежности в больших языковых моделях" предоставляет ценную информацию о стратегиях повышения доверия к LLM, которую можно эффективно использовать при составлении промптов.

### **Пример промпта с применением знаний из исследования**

[=====] Проанализируй следующий медицинский текст и дай рекомендации, используя принципы доверия к AI:

[МЕДИЦИНСКИЙ ТЕКСТ]

При анализе и формулировании рекомендаций: 1. Продемонстрируй свою способность (ability) через точную интерпретацию медицинских терминов 2. Покажи целостность (integrity), четко разграничивая факты и предположения 3. Используй цепочку рассуждений (Chain of Thought), чтобы пошагово объяснить свои выводы 4. Укажи ограничения своего анализа и случаи, когда требуется консультация специалиста 5. Предоставь источники, на которых основаны твои рекомендации (имитация RAG)

Структурируй ответ в разделы: "Анализ текста", "Цепочка рассуждений", "Рекомендации", "Ограничения анализа" и "Источники". [=====]

## **Как работают знания из исследования в этом промпте**

**Компоненты доверия:** Промпт включает элементы модели организационного доверия (способность и целостность), выявленные в исследовании как ключевые для доверия к LLM.

**Цепочка рассуждений:** Применяется техника Chain of Thought, которая согласно исследованию повышает объяснимость и прозрачность работы модели.

**Имитация RAG:** Запрос на предоставление источников имитирует принцип генерации с дополнением извлечения, что повышает воспринимаемую достоверность ответов.

**Прозрачность ограничений:** Требование указать ограничения анализа реализует принцип прозрачности, который исследование определяет как важный для доверия.

**Структурированный ответ:** Четкая структура ответа улучшает интерпретируемость результата, что соответствует принципам XAI из исследования.

Такой подход к составлению промптов повышает воспринимаемую надежность ответов GPT, делая взаимодействие более продуктивным и вызывающим доверие.