

# Обучение в контексте против настройки инструкций: случай малых и многоязычных языковых моделей

Дата: 2025-03-03 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2503.01611>

Рейтинг: 75

Адаптивность: 85

## Ключевые выводы:

Исследование сравнивает эффективность обучения в контексте (ICL) и инструктивной настройки (instruction tuning) для многоязычных и малых языковых моделей. Основной вывод: ICL значительно уступает инструктивной настройке в многоязычных сценариях и на малых моделях, даже при применении оптимизации прямых предпочтений (DPO).

## Объяснение метода:

Исследование предоставляет практически применимый метод URIAL для улучшения следования инструкциям базовыми моделями без специальной настройки. Результаты о влиянии языка и размера модели дают пользователям ценное понимание ограничений LLM. Анализ критических ошибок помогает избегать проблемных ситуаций. Большинство выводов могут быть применены с минимальной адаптацией.

**## Ключевые аспекты исследования:** 1. **Сравнение методов инструктирования LLM:** Исследование сравнивает три подхода к выполнению инструкций моделями: нулевой шот (zero-shot), обучение в контексте (ICL/URIAL) с использованием нескольких примеров и полноценная инструктивная настройка модели (instruction tuning).

**Мультиязычное сравнение:** Авторы оценивают эффективность данных методов не только на английском, но и на французском и испанском языках, выявляя разницу в качестве выполнения инструкций в зависимости от языка.

**Влияние размера модели:** Исследование анализирует, как размер модели (от 1.7B до 18B параметров) влияет на способность следовать инструкциям без специальной настройки.

**Применение DPO:** Авторы исследуют, насколько метод Direct Preference

Optimisation (DPO) может улучшить способность базовых моделей следовать инструкциям без полноценной инструктивной настройки.

**Анализ критических ошибок:** В работе проводится детальный анализ критических ошибок моделей (бесконечные циклы, генерация нерелевантного кода), которые могут сделать ответы полностью непригодными для использования.

## Дополнение:

### Применимость методов без дообучения или API

Исследование демонстрирует, что методы ICL (In-Context Learning), особенно URIAL с тремя стилистическими примерами, могут быть применены в стандартном чате **без какого-либо дообучения или API**. Это ключевое преимущество данного исследования для обычных пользователей.

### Концепции и подходы для стандартного чата

**Метод URIAL:** Добавление трех примеров взаимодействия в промпт: Два примера стандартных ответов на обычные запросы Один пример ответа на чувствительный запрос, демонстрирующий отказ от вредного контента Системный промпт, описывающий желаемое поведение

**Структурирование примеров:** Исследование показывает, что формат примеров (Query/Answer) важен для эффективности метода.

**Языковая адаптация:** Примеры должны быть на том же языке, что и запрос пользователя. Исследование предоставляет шаблоны для английского, испанского и французского.

**Обнаружение критических ошибок:** Пользователи могут идентифицировать признаки проблемных ответов (повторения, нерелевантный код).

### Ожидаемые результаты применения

**Улучшение следования инструкциям:** Применение URIAL может повысить качество ответов базовой модели на 0.5-1 балл по 5-балльной шкале.

**Повышение безопасности:** URIAL значительно улучшает способность модели отклонять вредные запросы.

**Улучшение языковой согласованности:** URIAL повышает вероятность получения ответа на том же языке, что и запрос.

**Снижение критических ошибок:** Применение URIAL существенно снижает вероятность бесконечных циклов и генерации нерелевантного кода.

Важно отметить, что эффективность метода снижается для маленьких моделей и на

неанглийских языках, но все равно дает заметное улучшение по сравнению с прямым запросом к базовой модели.

**## Анализ практической применимости: 1. Сравнение методов инструктирования LLM:** - Прямая применимость: Высокая. Пользователи могут непосредственно применять метод URIAL (добавление 3 примеров в промпт) для улучшения качества ответов базовых моделей. - Концептуальная ценность: Высокая. Исследование демонстрирует, что даже без специальной настройки можно значительно улучшить следование инструкциям. - Потенциал для адаптации: Высокий. Пользователи могут адаптировать примеры под свои конкретные задачи.

**Мультиязычное сравнение:** Прямая применимость: Средняя. Пользователи неанглоязычных стран получают понимание, что модели могут работать хуже на их родном языке. Концептуальная ценность: Высокая. Исследование показывает, что качество следования инструкциям снижается на неанглийских языках. Потенциал для адаптации: Средний. Пользователи могут адаптировать примеры на своем языке, но результаты будут хуже, чем на английском.

#### **Влияние размера модели:**

Прямая применимость: Высокая. Пользователи получают понимание, что маленькие модели существенно хуже следуют инструкциям без специальной настройки. Концептуальная ценность: Высокая. Исследование ясно демонстрирует, что для маленьких моделей инструктивная настройка критически важна. Потенциал для адаптации: Средний. Пользователи могут выбирать модели подходящего размера в зависимости от своих задач.

#### **Применение DPO:**

Прямая применимость: Низкая. Большинство пользователей не смогут самостоятельно применить DPO к моделям. Концептуальная ценность: Средняя. Исследование показывает, что даже без инструктивной настройки можно улучшить качество моделей. Потенциал для адаптации: Низкий. Технология требует специальных навыков и ресурсов.

#### **Анализ критических ошибок:**

Прямая применимость: Высокая. Пользователи могут идентифицировать типичные критические ошибки моделей. Концептуальная ценность: Высокая. Понимание типичных ошибок помогает избегать проблемных случаев. Потенциал для адаптации: Высокий. Знание о типичных ошибках позволяет корректировать промпты для их избегания.

#### **Prompt:**

Использование выводов исследования ICL vs Instruction Tuning в промптах ##  
Ключевые знания из исследования

Исследование показывает, что: - Инструктивная настройка (instruction tuning) превосходит обучение в контексте (ICL) для многоязычных и малых моделей - Для малых моделей (<2B параметров) разрыв между подходами особенно велик - На неанглийских языках ICL работает значительно хуже - Базовые модели чаще допускают критические ошибки

## Пример промпта с учетом этих знаний

[=====] Я работаю с языковой моделью Llama 3 размером 8B параметров на французском языке. На основе исследования об эффективности различных подходов к обучению:

Я знаю, что для неанглийских языков инструктивно настроенные модели работают лучше, чем ICL-подходы. Поэтому я предпочту использовать прямые инструкции вместо предоставления нескольких примеров. Задача: Создай краткое резюме следующего текста о климатических изменениях. [ТЕКСТ]

Пожалуйста, сделай резюме четким, структурированным и сохрани ключевые идеи оригинала. [=====]

## Объяснение эффективности

Этот промпт эффективен, потому что:

**Избегает ICL для неанглийского языка** — исследование показало, что на французском языке инструктивная настройка значительно превосходит ICL (оценки 4.45 vs 3.98) **Использует прямые четкие инструкции** вместо предоставления примеров, что оптимально для многоязычных моделей **Формулирует конкретные ожидания** от результата (четкость, структурированность), что снижает вероятность критических ошибок **Учитывает размер модели** — для 8B модели разрыв между подходами существенен, но не критичен, как для моделей <2B Если бы мы работали с моделью меньшего размера (например, 1.7B) на неанглийском языке, разница была бы еще более значительной, и использование инструктивного подхода стало бы критически важным.