

# Насколько эффективен код, сгенерированный LLM? Строгая и высокоуровневая оценка

Дата: 2025-02-18 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2406.06647>

Рейтинг: 60

Адаптивность: 55

## Ключевые выводы:

Исследование направлено на разработку строгого и высококачественного бенчмарка ENAMEL для оценки способности больших языковых моделей (LLM) генерировать эффективный код. Основные результаты показывают, что существующие LLM значительно отстают от экспертного уровня в генерации эффективного кода, испытывая трудности с разработкой продвинутых алгоритмов и оптимизацией реализации.

## Объяснение метода:

Исследование демонстрирует, что LLM генерируют функционально корректный, но неэффективный код. Предлагает методологию оценки и эталонные решения, но применение требует высокой технической подготовки. Ценно для понимания ограничений LLM в создании эффективного кода, но имеет ограниченную прямую применимость для широкой аудитории.

## Ключевые аспекты исследования: 1. **Новая метрика оценки эффективности кода** - исследователи разработали метрику `eff@k`, которая расширяет стандартную метрику `pass@k` и учитывает цензурированное время выполнения кода при превышении лимита времени.

**Эталонные эффективные решения** - эксперты разработали оптимальные по эффективности решения для 142 задач из наборов HumanEval и HumanEval+, многие из которых значительно эффективнее канонических решений.

**Сильные генераторы тестовых случаев** - исследователи создали строгие генераторы тестовых случаев, которые выявляют как неправильные, так и неоптимальные алгоритмы.

**Многоуровневая оценка** - разработана система с несколькими уровнями сложности входных данных, позволяющая дифференцировать код разной эффективности.

**Обширное тестирование 30 LLM** - результаты показывают, что даже самые

мощные модели (GPT-4) далеки от создания кода экспертного уровня эффективности ( $\text{eff}@1=0.454$ ).

**## Дополнение:** Действительно ли для работы методов этого исследования требуется дообучение или API? Или методы и подходы можно применить в стандартном чате, а ученые лишь для своего удобства использовали расширенные техники?

Для большинства методов этого исследования **не требуется дообучение или API**. Исследователи использовали API для некоторых моделей (например, GPT-4), но сама методология оценки эффективности и основные концепции могут быть применены в стандартном чате.

Концепции и подходы, которые можно применить в стандартном чате:

**Понимание разрыва между корректностью и эффективностью** - пользователи могут осознать, что LLM часто генерируют работающий, но неэффективный код, и более критично оценивать полученные решения.

**Многоуровневое тестирование** - пользователи могут создавать тесты разной сложности для проверки эффективности кода:

Начинать с малых входных данных для проверки корректности Постепенно увеличивать размер входных данных для оценки эффективности Искать граничные случаи, где неэффективные решения будут работать медленно

**Запросы на эффективные алгоритмы** - хотя исследование показало ограниченную эффективность простого запроса "самого эффективного алгоритма", пользователи могут:

Запрашивать оценку временной сложности решения Просить модель проанализировать и улучшить эффективность предложенного решения Использовать пошаговый промптинг, направляя модель к более эффективным решениям

**Изучение экспертных решений** - пользователи могут использовать примеры эффективных решений как образцы и просить модель объяснить их или создать подобные.

Результаты, которые можно получить от применения этих концепций: - Более эффективный код для решения задач - Лучшее понимание ограничений LLM в создании оптимальных алгоритмов - Развитие критического мышления при оценке генерируемого кода - Улучшение собственных навыков алгоритмического мышления

Хотя исследование показывает, что даже с прямыми указаниями LLM часто не могут генерировать алгоритмически оптимальный код, понимание этих ограничений и применение многоуровневого подхода к тестированию может значительно улучшить качество получаемых решений.

**## Анализ практической применимости: 1. Новая метрика оценки эффективности кода** - Прямая применимость: Пользователи могут адаптировать подход для оценки эффективности генерируемого кода в своих задачах - Концептуальная ценность: Высокая - помогает понять, что LLM часто генерируют функционально корректный, но неэффективный код - Потенциал для адаптации: Пользователи могут разработать собственные тесты с разными уровнями сложности для оценки эффективности решений

**Эталонные эффективные решения** Прямая применимость: Пользователи могут изучить эти решения как образцы оптимального кода Концептуальная ценность: Очень высокая - демонстрирует разрыв между типичными и экспертными решениями Потенциал для адаптации: Пользователи могут использовать эти образцы для обучения и улучшения своего кода

### **Сильные генераторы тестовых случаев**

Прямая применимость: Ограниченная для обычных пользователей, требует экспертных знаний Концептуальная ценность: Средняя - понимание важности тестирования худших случаев Потенциал для адаптации: Пользователи могут улучшить свои тестовые наборы, включив граничные случаи

### **Многоуровневая оценка**

Прямая применимость: Пользователи могут применить аналогичный подход для проверки своего кода Концептуальная ценность: Высокая - помогает понять, что эффективность кода нельзя оценивать бинарно Потенциал для адаптации: Хороший - применим при разработке систем тестирования

### **Обширное тестирование LLM**

Прямая применимость: Пользователи могут выбирать модели с лучшей эффективностью для своих задач Концептуальная ценность: Высокая - показывает ограничения LLM в создании эффективного кода Потенциал для адаптации: Пользователи могут использовать результаты для выбора моделей и стратегий промптинга

## **Prompt:**

Использование результатов исследования ENAMEL в промптах для GPT ##  
Ключевые инсайты из исследования

Исследование ENAMEL показывает, что даже современные LLM (включая GPT-4) хорошо справляются с генерацией корректного кода, но значительно хуже с генерацией *эффективного* кода. Этот разрыв можно использовать для создания более продуктивных промптов.

## Пример промпта для получения эффективного кода

[=====] # Задача: Оптимизация алгоритма поиска k-ой наименьшей суммы пар

## Контекст Мне нужно найти k-ую наименьшую сумму пар элементов из двух отсортированных массивов.

## Требования к решению 1. **Эффективность критична**: Решение должно иметь оптимальную временную сложность  $O(k \log k)$  или лучше. 2. **Несколько подходов**: Предложи 3 разных алгоритмических подхода к решению этой задачи. 3. **Анализ эффективности**: Для каждого подхода: - Укажи временную и пространственную сложность - Опиши потенциальные узкие места - Предложи оптимизации на уровне реализации

## Процесс решения 1. Сначала опиши общую идею каждого алгоритма 2. Затем реализуй наиболее эффективное решение с подробными комментариями 3. Проанализируй свое решение на сильных тестовых случаях, включая граничные условия и наихудшие сценарии 4. Предложи оптимизации реализации (например, использование кеширования, избегание ненужных вычислений)

## Формат вывода - Начни с краткого сравнения всех трех подходов - Предоставь полную реализацию лучшего подхода - Завершите анализом эффективности и обоснованием выбора [=====]

## Объяснение эффективности промпта

Данный промпт применяет основные выводы исследования ENAMEL:

**Генерация нескольких вариантов решения** (запрос трех разных подходов) - исследование показало, что метрика  $eff@k$  значительно улучшается с увеличением k.

**Явное указание требований к временной сложности** - исследование отмечает, что LLM часто не выбирают оптимальные алгоритмы без явных указаний.

**Разбиение задачи на подзадачи** (описание идеи → реализация → анализ → оптимизация) - помогает моделям справляться со сложными алгоритмическими задачами.

**Запрос на анализ эффективности** - заставляет модель критически оценить собственное решение.

**Акцент на сильных тестовых случаях** - исследование показало, что случайные тесты часто не выявляют субоптимальные алгоритмы.

Такой структурированный промпт значительно повышает шансы получить не просто корректное, но и эффективное решение, преодолевая ограничения LLM в генерации

оптимального кода, выявленные в исследовании ENAMEL.