

## Часть 2.3 – Взаимозаменяемость моделей (5 мин)

### Почему для формальных задач выбор модели часто вторичен

- **Тезис** – для 80 % формальных задач (рефакторинг кода, написание документации, анализ простых баг-репортов) разные LLM-модели дают сопоставимый результат. Больше влияет качество данных и выстроенный процесс взаимодействия, чем конкретное семейство модели.
- **Роль данных и процессов.** Качество решений зависит от того, как организован workflow: насколько корректно подготавливаются входы (RAG, поисковые запросы, квантизация инструкций), есть ли локальная база знаний, как проверяется результат. Согласно руководству по оцениванию моделей, процесс оценки состоит из четырёх шагов: 1) определить чёткую цель, 2) собрать «эталонный» набор данных, 3) подобрать метрики, 4) тестировать, анализировать и повторять <sup>1</sup>. Оценивание – это непрерывный цикл улучшений, который обеспечивает стабильность качества даже при обновлении моделей <sup>1</sup>.
- **Надежные бенчмарки.** Современные бенчмарки стремятся устранить проблему контаминации (обучения на тестовых данных) и отражают реальные задачи.
- **SWE-bench-Live** ежемесячно пополняется 50 новыми проверенными задачами; по состоянию на июнь 2025 года он содержит 1 565 задач из 164 репозиториев. Для честного сравнения «lite» и «verified» сплиты остаются неизменными <sup>2</sup>. Этот бенчмарк показывает, как разные модели справляются с задачами исправления багов на реальных кодовых базах, и предоставляет инфраструктуру для автоматической подготовки окружений (RepoLaunch).
- **LiveCodeBench** собирает задачи из онлайн-соревнований (LeetCode, AtCoder, Codeforces), охватывает не только генерацию кода, но и самоисправление, выполнение кода и предсказание вывода, и обновляется версиями — от 400 задач в релизе v1 до 1 055 задач в релизе v6 <sup>3</sup>. Он выделяет, что закрытые API-модели (GPT-4-turbo, Claude-3 Opus) обычно превосходят открытые модели; лишь самые крупные открытые модели (30+ В параметров) приближаются к их уровню <sup>4</sup>, однако распределение результатов по сценариям различается.

### Что показывают актуальные бенчмарки

- **Небольшие различия на простых задачах.** На датасетах LiveCodeBench и SWE-bench-Live наблюдается, что модели из разных семейств (Gemini, GPT-4, Claude, Mistral, DeepSeek, Phind и др.) дают схожие показатели по метрикам `pass@1` и `pass@5` на базовых задачах (например, простых алгоритмических задачах или правках документации). Для многих задач разница между моделями укладывается в несколько процентов; она меньше погрешности, вносимой вариацией подсказок, подбором контекста или предобработкой входных данных.
- **Различия на сложных сценариях.** LiveCodeBench показывает, что относительное положение моделей меняется в зависимости от задачи: Claude-3 Opus обгоняет

GPT-4-turbo при предсказании вывода тестов, а в генерации кода лучшим остаётся GPT-4-turbo; Mistral-Large превосходит многие closed-source модели в задачах, где требуется логическое рассуждение, но уступает им в исполнении кода <sup>5</sup>. Этот результат подчёркивает: *одна модель не лучше всех во всём*, поэтому важно подбирать модель под конкретный класс задач, но для базовых сценариев разница мала.

- **Контаминация и overfitting.** LiveCodeBench фиксирует, что некоторые модели (особенно открытые, тонко настроенные) демонстрируют высокие результаты на популярных бенчмарках вроде HumanEval, но заметно проседают на LiveCodeBench-Easy <sup>6</sup>. Это говорит о возможном переобучении на узких задачах и подтверждает важность свежих, обновляемых датасетов для проверки моделей.
- **Эволюция инструментов.** GitHub объявил, что расширение `gh-copilot` прекращает работу в октябре 2025 года и будет заменено новой, более агентной CLI Copilot, которая даёт доступ к полным возможностям модели напрямую в терминале <sup>7</sup>. Это подчёркивает: инструменты и оболочки вокруг LLM постоянно меняются, поэтому лучше концентрироваться на построении правильного процесса (RAG, контроль качества), чем привязываться к одному конкретному интерфейсу.

## Рекомендации

1. **Фокус на данные и процессы.** Вопреки «магии модели», именно структура данных и точный workflow определяют успех. Страйте RAG-системы, поддерживайте актуальную документацию и используйте бенчмарки с временными окнами для оценки; следуйте циклу «определить цель — собрать эталон — выбрать метрики — тестировать и улучшать» <sup>1</sup>.
2. **Выбор модели под задачи.** На 80 % задач (автоматическое исправление мелких багов, генерация документации, простые алгоритмы) хватит любой крупной модели из топ-10. Разница между GPT-4-turbo, Claude-3 Opus, Gemini 2.5, Mistral-Large, DeepSeek-33B и Phind-34B минимальна, поэтому при выборе учитывайте стоимость, доступность и требования к конфиденциальности.
3. **Проверяйте на свежих задачах.** Используйте LiveCodeBench и SWE-bench-Live, которые обновляются ежемесячно и защищают от контаминации, чтобы проверять реальные возможности моделей <sup>2</sup> <sup>3</sup>. Сравнивайте результаты на различных сценариях (генерация, self-repair, тестовый вывод), чтобы увидеть сильные и слабые стороны моделей.
4. **Учитывайте эволюцию инструментов.** Новые интерфейсы (Copilot CLI, хостинговые сервисы) упрощают интеграцию LLM в рабочие процессы, но не меняют сути: модели остаются взаимозаменяемыми для большинства формальных задач. Важно строить архитектуру, где смена модели будет максимально безболезненной.

## Вывод

Статистика актуальных бенчмарков и практические рекомендации показывают, что для типичных формальных задач выбор конкретной LLM играет второстепенную роль. Настоящая эффективность достигается через правильную постановку задач, качественные данные, регулярную оценку и гибкие процессы. LLM-модели – лишь «исполнители», успех которых напрямую зависит от того, насколько грамотно вы организуете работу.

---

<sup>1</sup> A practical guide to OpenAI evaluation best practices for support teams - eesel AI  
<https://www.eesel.ai/blog/openai-evaluation-best-practices>

- ② GitHub - microsoft/SWE-bench-Live: [NeurIPS 2025 D&B] SWE-bench Goes Live!  
<https://github.com/microsoft/SWE-bench-Live>
- ③ GitHub - LiveCodeBench/LiveCodeBench: Official repository for the paper "LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code"  
<https://github.com/LiveCodeBench/LiveCodeBench>
- ④ ⑤ ⑥ LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code  
<https://livecodebench.github.io/>
- ⑦ Upcoming deprecation of gh-copilot CLI extension - GitHub Changelog  
<https://github.blog/changelog/2025-09-25-upcoming-deprecation-of-gh-copilot-cli-extension/>