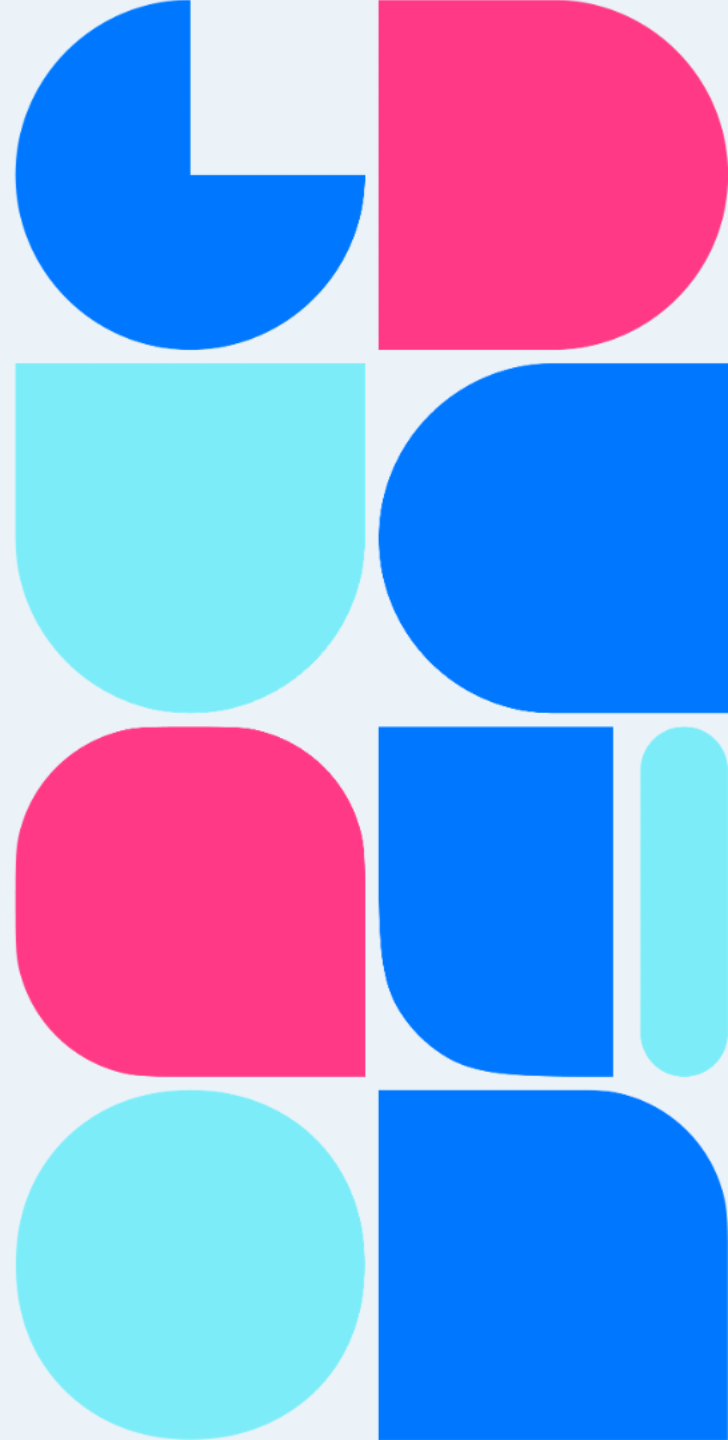




Развитие больших языковых моделей, от малых компаний до крупных IT гигантов: Достигнут ли максимум?

Максим Никонов
Руководитель аналитического направления (САО)
в подразделении VK



Максим Никонов

Научный сотрудник ВМК МГУ

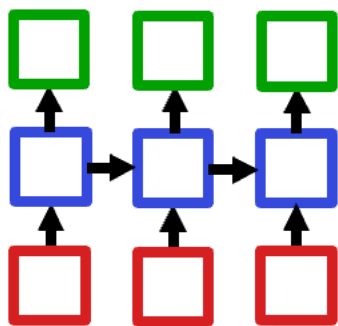
Старший преподаватель НИУ ВШЭ
Преподаватель VK Education

Автор и лектор курсов по машинному обучению, NLP, Большим языковым моделям, Диффузионным моделям, Математической статистике



Какие бывают задачи?

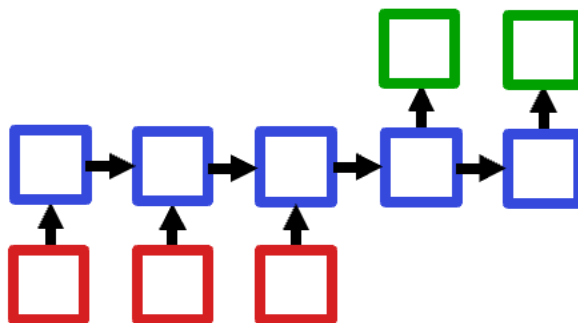
many-to-many
одинаковая длина



Классификация токенов:

- Исправление ошибок в тексте
- Синонимы

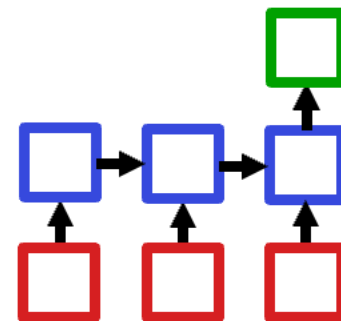
many-to-many
разная длина



Seq2seq:

- Машинный перевод
- Суммаризация
- Перенос стиля

many-to-one



- Классификация
- Регрессия

LLM обучается решать разные задачи

Задача	Пример текста, обучающий этой задаче
Грамматика	В свободное время я люблю (читать, табуретка)
Лексическая семантика	Я пошел в магазин, чтобы купить манго, апельсин и (яблоки, енота)
Знания о мире	Столица Франции – (Париж, Вена)
Классификация тональности	Я в восторге от декораций и игры актеров, спектакль был (хорошим, плохим)
Перевод	"Стол" по-английски будет ("table", "apple")
Пространственное мышление	Леша сидел на диване в гостиной, рядом с ним сидел Саша. Через 15 минут Саша встал и вышел из (гостиной, кухни)
Математика	Если прибавить 4 к 3, то будет (7, 8)

Какие задачи есть у бизнеса

1

Боты – помощники

Частные
сервисы –
подсчет
калорий,
планирование
путешествий

2

Разметка данных

Замена
ассессоров

3

Принятие решений

Финансы,
модерация

4

Перефраз

Саммари текста,
выделение
полезной
информации,
перенос стиля,
перевод

5

Доменные области

Для медицины,
рекомендаций,
стиля

Количество параметров – качество?

3B

**Для размещения на
мобильных устройствах**

Вариант размещения:
client-side

Много неточностей

8B

**Для малого бизнеса и
тестирования гипотез,
скорее про генерацию
текста**

Вариант размещения:
несколько NVIDIA 4090 для
несколько десяткой
пользователей

В доменной области можно
получить качество на
уровне 70B модели

70B

**Для крупного бизнеса,
можно использовать
для принятий решений**

Вариант размещения:
NVIDIA A100 для десятка
пользователей

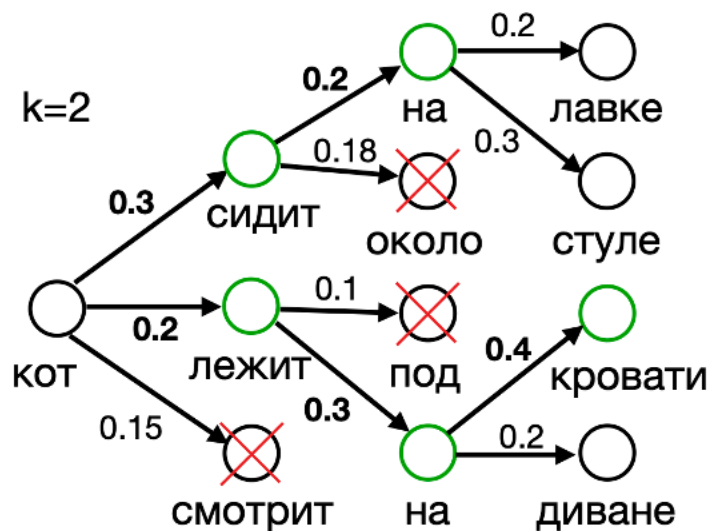
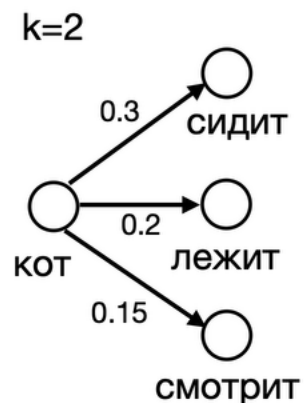
Работает с минимальным
промптом и во многих
областях

>

**Экспериментальное
направление**

Вариант размещения:
несколько NVIDIA A100

Как происходит генерация: Beam Search

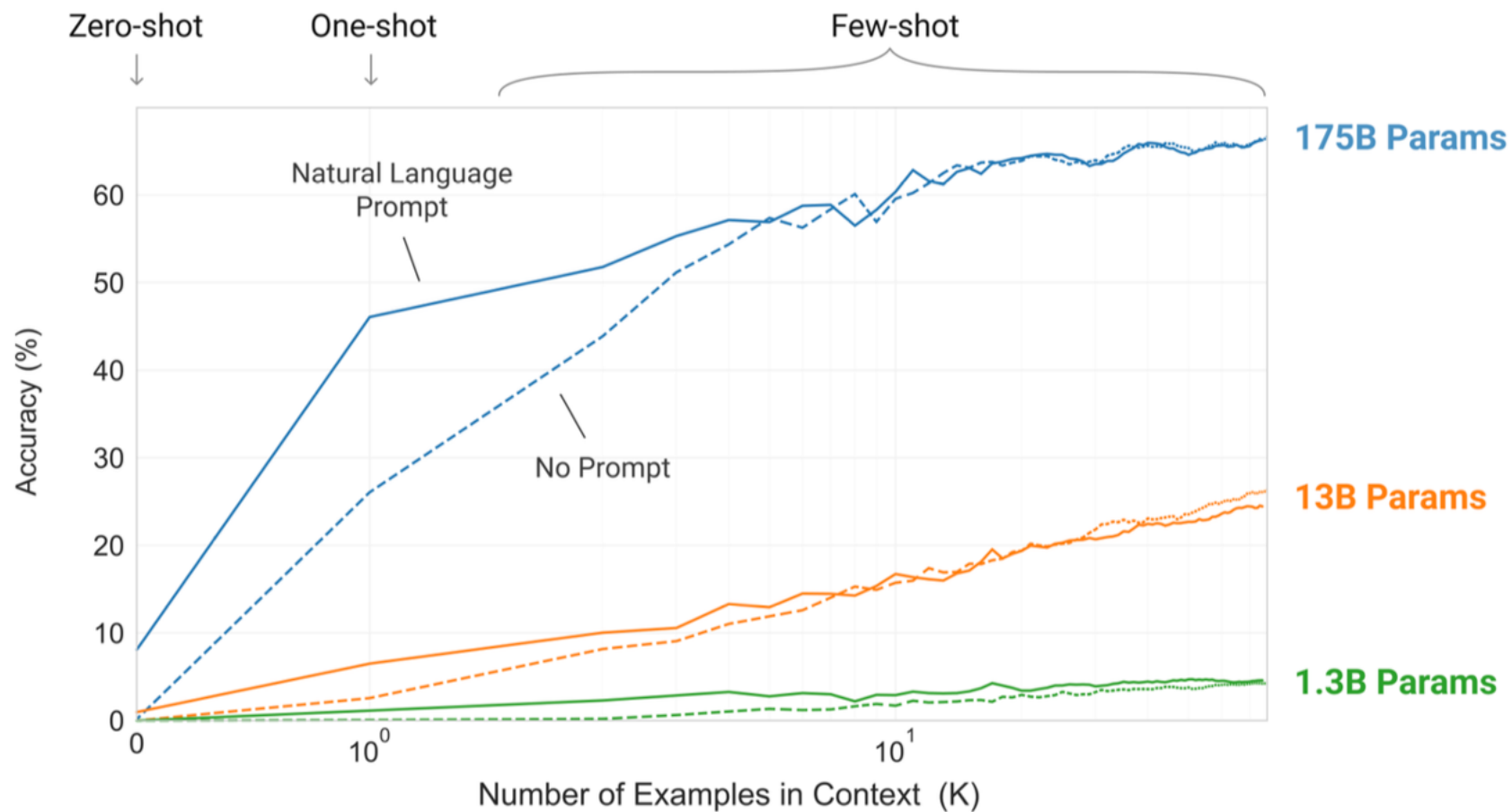


$$p(x_1, \dots, x_m) \approx \prod_{i=1}^m p(x_i | x_{i-1}, \dots, x_{i-n})$$

Задача seq2seq или чаще ТЕКСТ В ТЕКСТ

Нельзя параллельно вычислять из-за
Марковского свойства

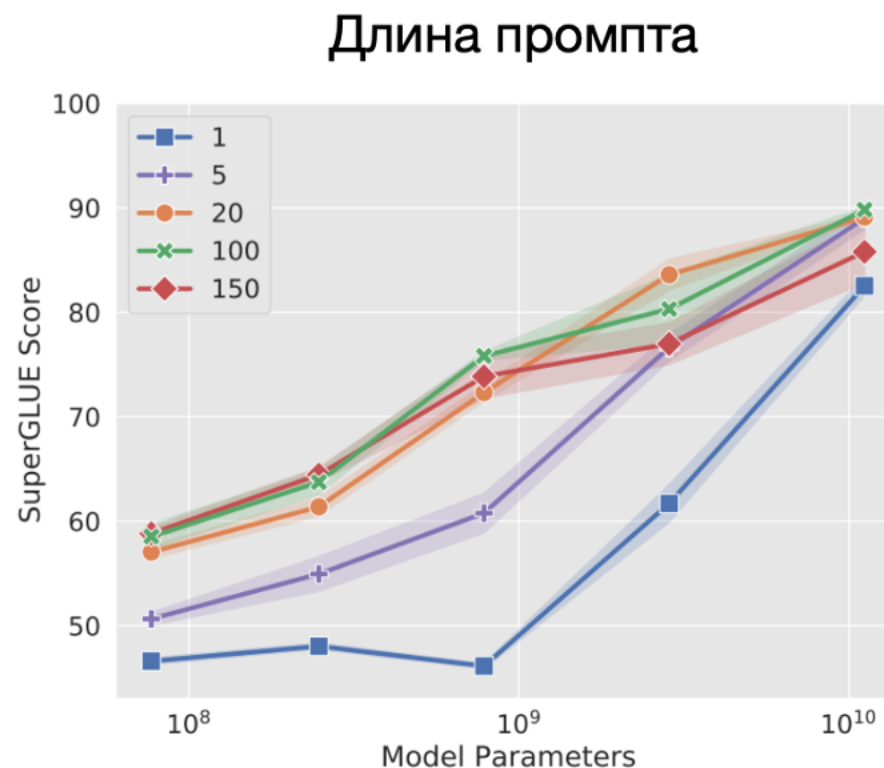
Промпты – это все?



Промпт важен,
но есть потолок

Все модели
улучшаются
с промптом

Промпты – это все?



Совет

Сформулировать задачу явно

Указать контекст

Писать пошагово

Задать ограничения

Показывать пример

Уточнять аудиторию

Разбивать запросы

В чём суть

Назовите конечную цель и нужный формат ответа.

Дайте модели нужные данные/фон.

Попросите «думать вслух» или «решай шаг за шагом».

Лимиты на объём, стиль, язык.

Few-shot: демонстрация желаемого формата.

Уровень знаний, роль читателя.

Длинные задачи делите на микропод-промпты.

Пример-шаблон

«Суммируй текст в 3 пунктах»

«Ты — дата-инженер, объясни принцип RAG...»

«Сначала распиши план, затем приведи код»

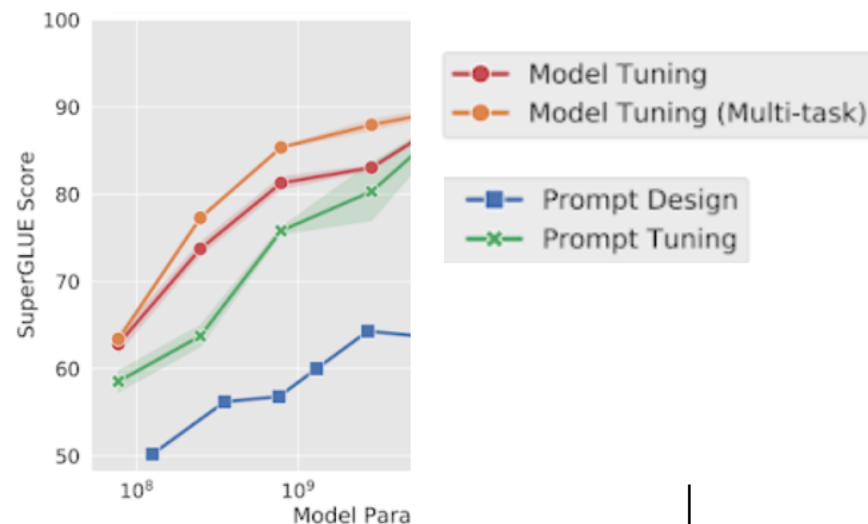
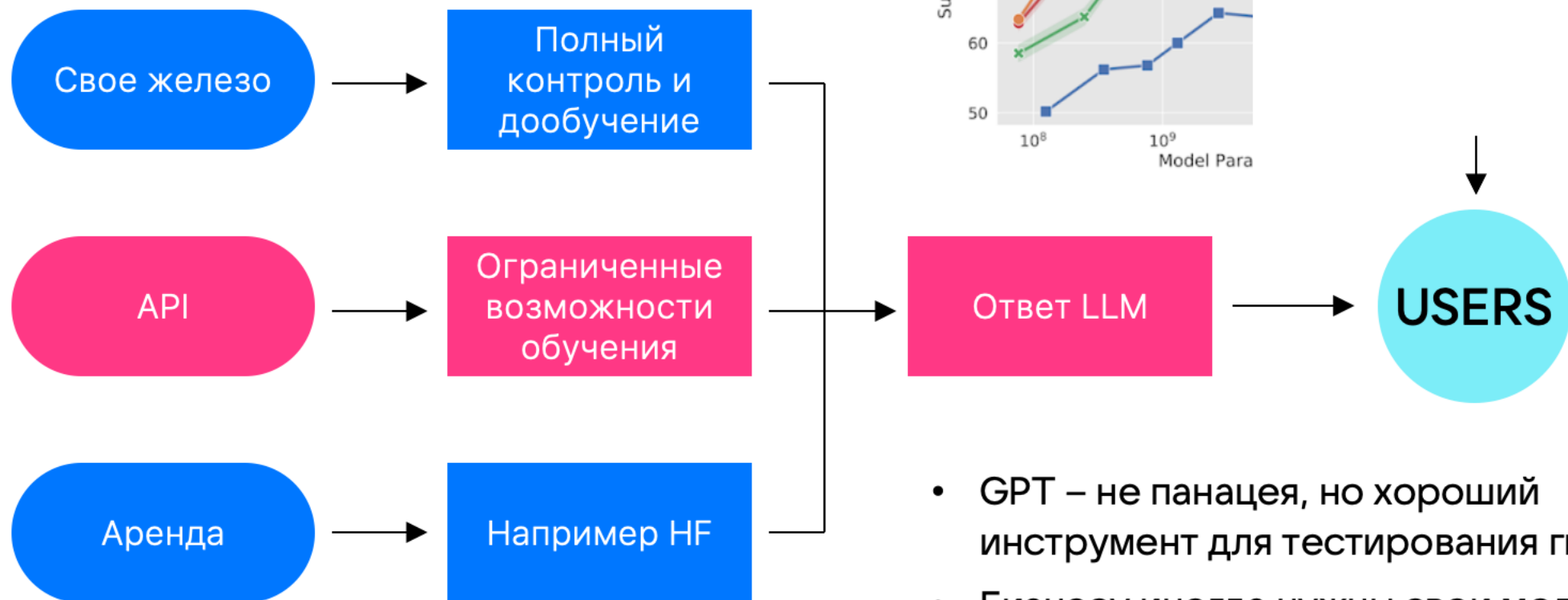
«Не больше 150 слов, на русском, без таблиц»

Q: ... → A: ...

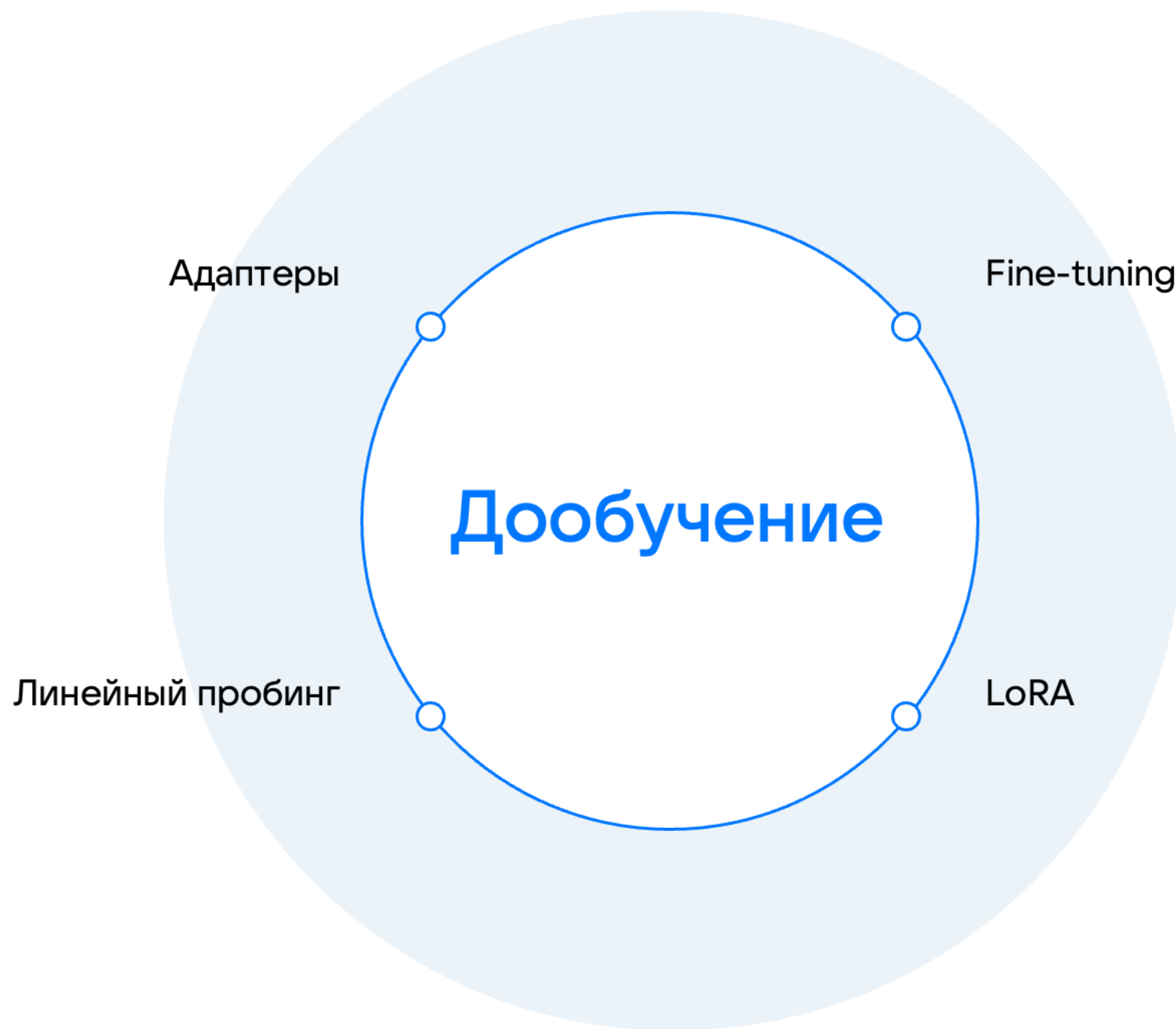
«Пиши для разработчика-junior»

«1) придумай идеи, 2) оцени риски...»

Схемы интеграции LLM

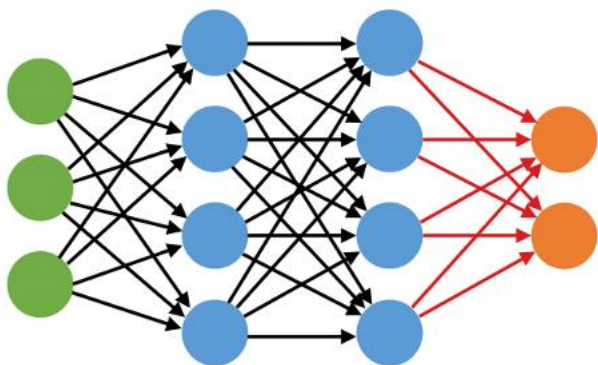


- GPT – не панацея, но хороший инструмент для тестирования гипотез
- Бизнесу иногда нужны свои модели



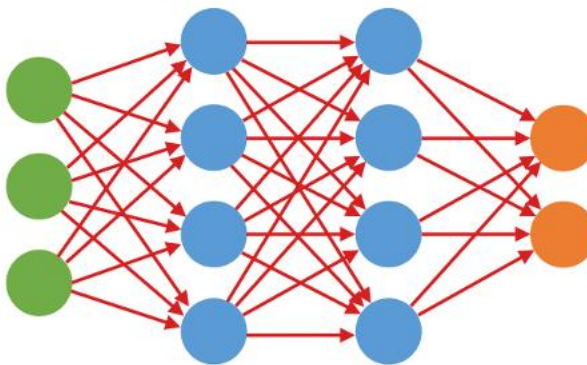
Способы дообучения

Обучение головы
(линейный пробинг)



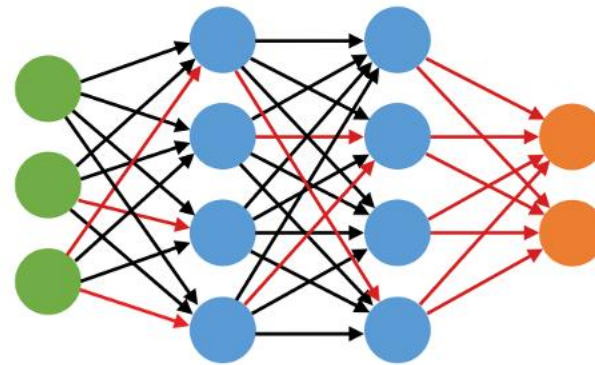
Обучается только
последний слой

Fine-tuning



Обучается вся
модель

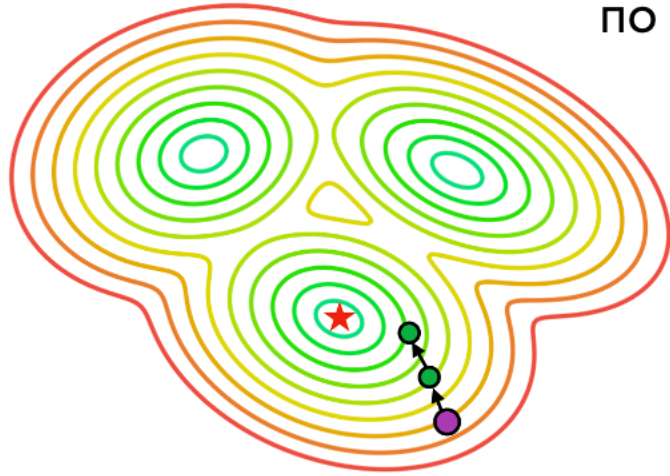
Parameter Efficient
Fine-tuning



Обучается
небольшой
набор весов

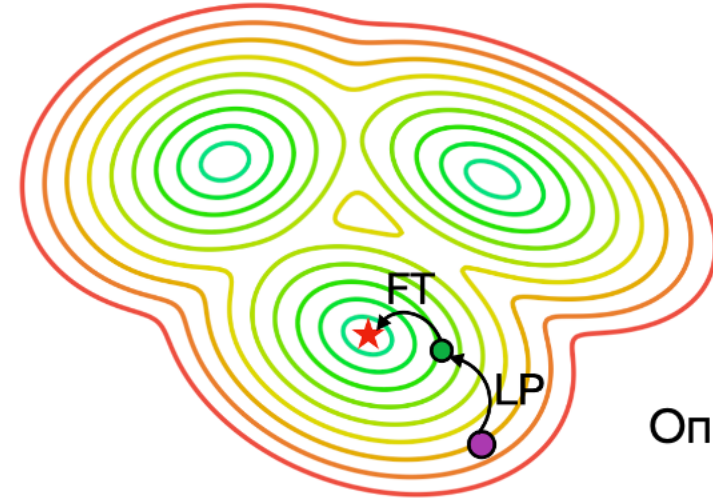
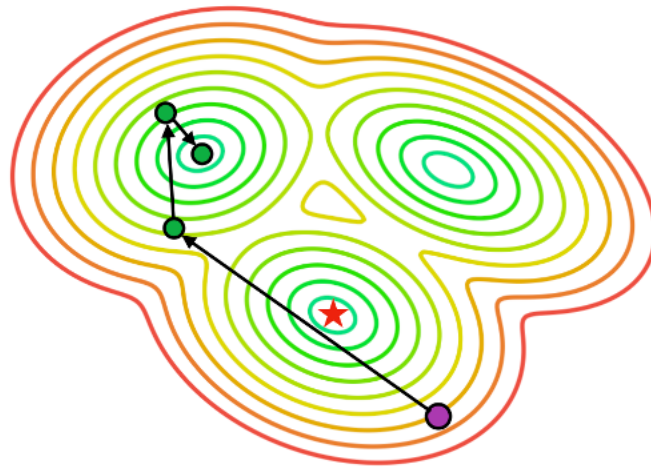
Fine-tuning портит модель

Пробинг



Большие ограничения
по качеству

Fine-tuning



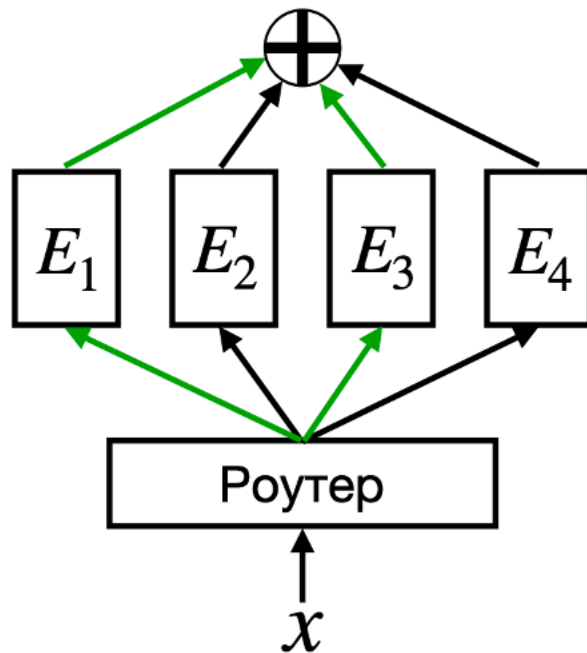
Оптимально

Цель: прийти в
максимум качества

Бизнес не понимает,
почему результат потерян

Все очень долго, что делать?

Mixture of Experts (MoE)

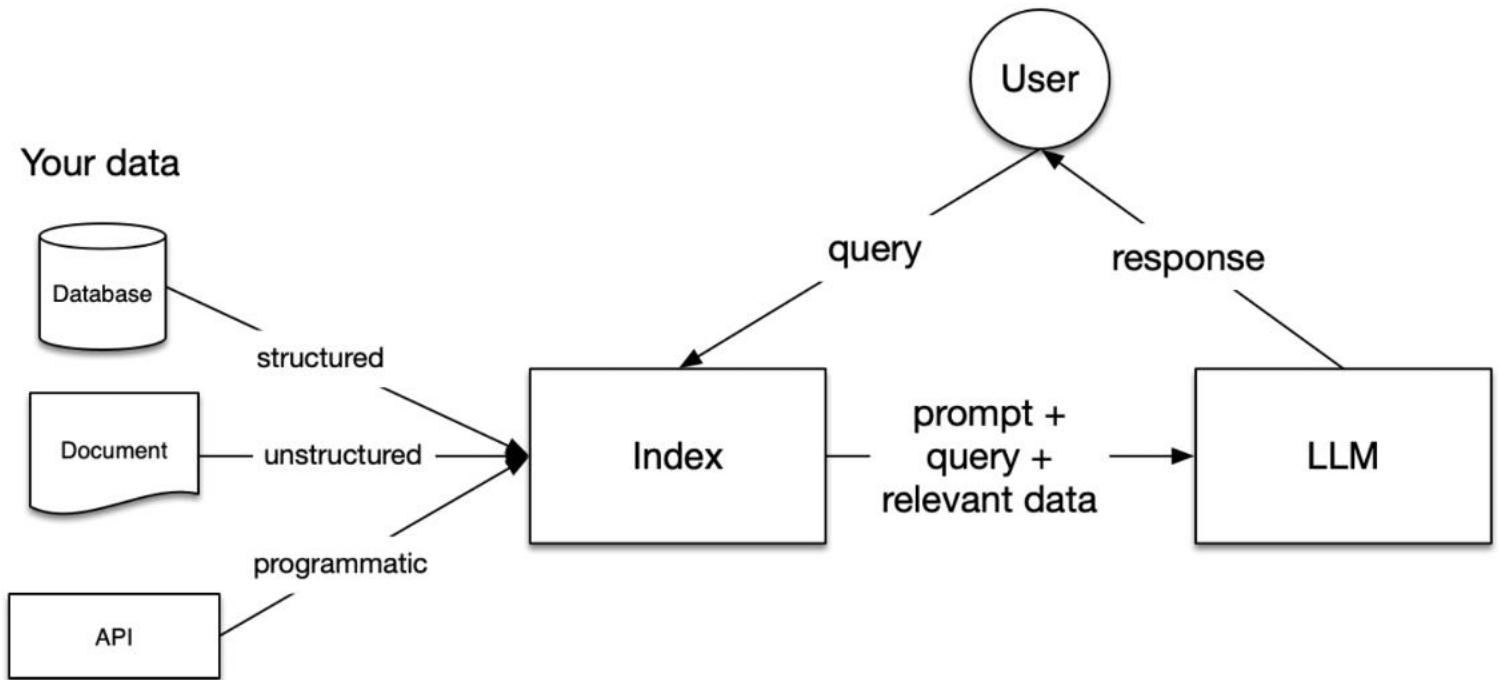
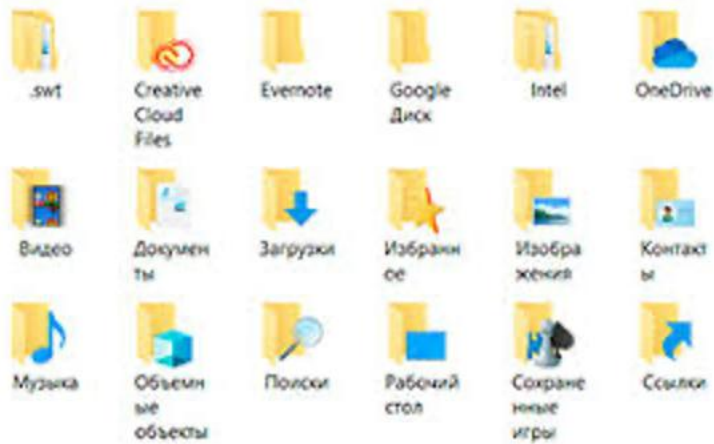


В классической архитектуре Mixture of Experts (дословно "смесь экспертов") модель разбивается на части, каждая из которых компетентна в своей области. Одновременно с этим обучают маленькую "проверяющую" (роутер) модель, которая умеет на основе входной задачи понять, к ответам каких экспертов нам стоит прислушаться сильнее

Не стоит путать "эксперта" с самостоятельным ИИ, который может работать сам по себе.

Не стоит воспринимать MoE как "чат" между экспертами, когда они совещаются и принимают общее решение (так работают агенты)

Все равно долго, что делать? RAG!



RAG – зачем?

Вопрос

Краткий ответ

Почему важно / риски

Что делает RAG?

Дополняет LLM внешними документами, вставляя их в prompt.

Позволяет отвечать на узкоспециальные или свежие вопросы без дообучения модели.

Главные блоки

Retriever → (Re-)Ranker → Prompt Builder → LLM → Post-processor.

Слабое звено = retriever; плохие эмбединги = галлюцинации.

Плюсы

Меньше «галлюцинаций»
Обновляется мгновенно (достаточно перезаписать индекс)
Возможна прозрачная ссылка на источник.

Сильно зависит от качества и чистоты корпуса

Минусы

Латентность (две модели подряд). Не решает все виды ошибок: LLM всё ещё может перепутать факты.

Требует инженерии кеша, батчинга, сжатия контекста.

Эволюция 2024-25

RAG 2.0 (end-to-end оптимизация, гибридные индексы, feedback-loops), мультивекторные и мультимодальные RAG; автоматическая адаптация chunk-size; «streaming-RAG» для real-time данных.

Дальше ожидается: self-updating индексы, RAG + инструменты (tool-use)

Где применяют

Copilot-вики, чат-боты техподдержки, юридические ресёрч-ассистенты, поиск по коду, аналитика BI.

.

Сводная по лучшим практикам

Категория	Практики «что делать»	Что это даёт
Prompt Engineering	Чётко формулировать задачу, роль, формат Показывать примеры + контр-примеры	Быстрое улучшение без затрат, ↓ галлюцинации.
RAG / RAG 2.0	Индексируйте чистые, дедуплицированные чанки 300-800 токенов	Актуальные ответы с ссылками, ↓ hallucination на 60-90 %.
Fine-Tuning / PEFT	Используйте LoRA/QLoRA, DPO, SFT + RLHF/RAFT. 1-3 k качественных примера > 100 k генераций. Замораживайте backbone, тюните Adapters	Вшитый стиль, ↑ точность спец-задач, inference ≈ base+δ.
Линия рассуждений	R1 рассуждения, проверки и мультиагенты	Меньше галлюцинаций
User-in-the-Loop	Кнопка «👍/👎» + комментарий → retrain buffer. Reward-model из живых оценок.	Постоянное улучшение без большого дата-тима.

Вернемся в начало

1

Боты – помощники

Частные сервисы –
подсчет калорий,
планирование
путешествий

2

Разметка данных

Замена ассессоров

3

Принятие решений

Финансы, модерация

4

Перефраз

Саммари текста,
выделение полезной
информации,
перенос стиля,
перевод

5

Доменные области

Для медицины,
рекомендаций, стиля

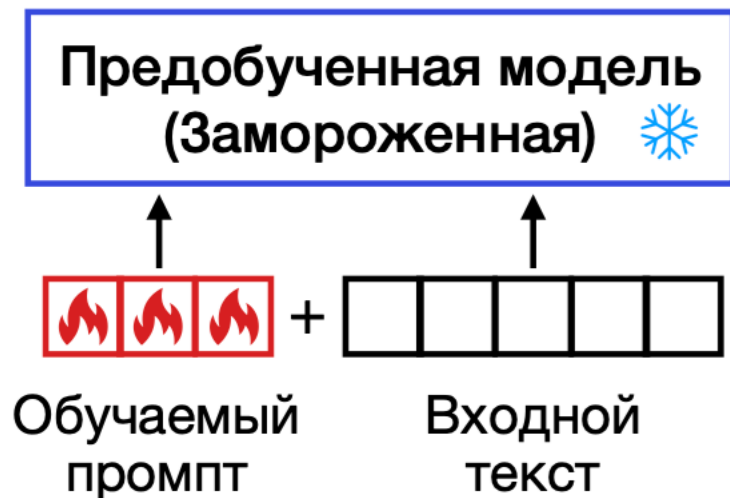
3В

8В

70В

>

Prompt Tuning



Идея: Попробуем автоматически подобрать наиболее подходящий промпт

- Инициализируем эмбединги промпта случайно, задавая только их число
- Можно инициализировать эмбедингами текстового промпта
- Для задачи классификации нужно дополнительно обучить голову

В чем революция Deepseek

Ключевая идея	В чём суть	Почему это важно
Mixture-of-Experts (MoE) на стероидах	671- ↔ 685 В «спящих» параметров, но для каждого токена активируется лишь 8 экспертов ≈ 37 В	Позволяет обучать и обслуживать «гиганта» на кластере из Nvidia H800 (в 10–20 раз дешевле западных аналогов)
DeepSeekMoE + Multi-head Latent Attention (MLA)	В версии V2 (236 В общих / 21 В активных): Sparse-вычисления + сжатие KV-кеша в «латентный» вектор	-42 % расходов на обучение -93 % ОЗУ на инференсе ≈ ×5.8 ускорение генерации
Специализированная линия R1	RL-fine-tune на решении задач «шаг-за-шагом», фокус на логику	Сильнее с GPT-4-o в HumanEval

Что еще пробуют

- VLM (Vision-Language Models)

На длинных (мульти-страничных) картинках точность резко падает; поколения «визуального» текста всё ещё часто галлюцинируют.

- LLM на диффузионных моделях

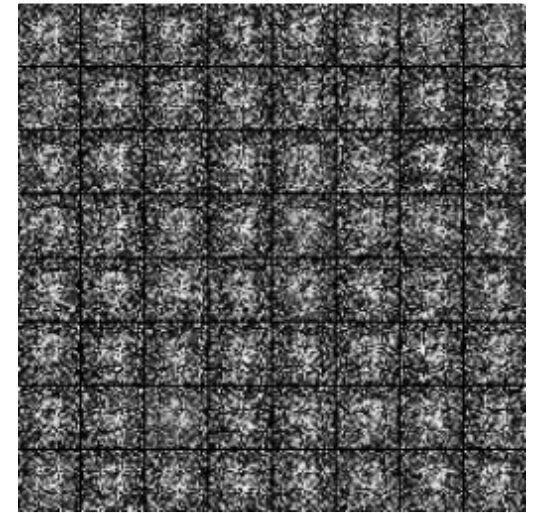
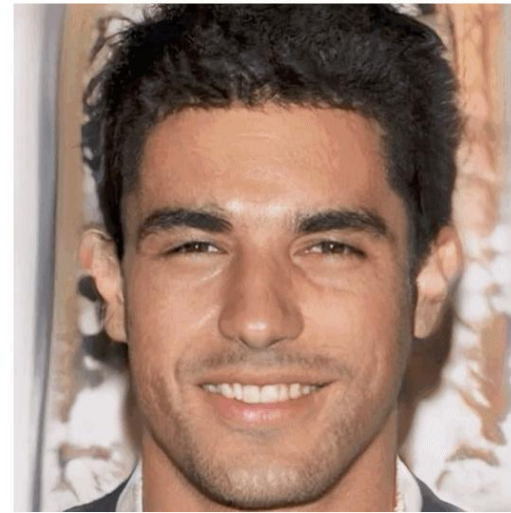
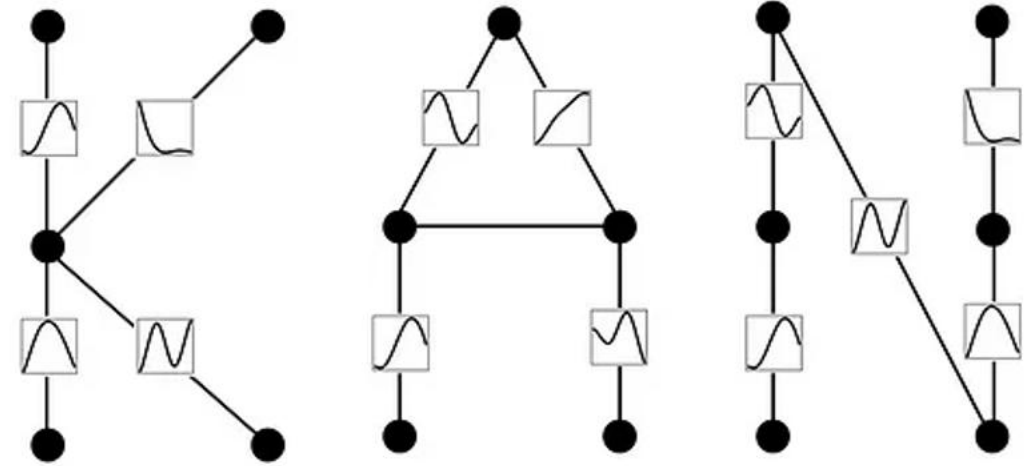
LLaDA 8 B догнала LLaMA-3 8 B по MMLU

Итеративное «денойзинг» $\approx 10\text{-}100$ раз медленнее AR LLM

- KAN-LLM (Колмогоров-Арнольд)

В мелких моделях KAN снижает параметры и быстрее сходится на задачах рассуждения, даёт более интерпретируемые функции. обучение нестабильно, нет оптимизированных GPU-ядер; наблюдается катастрофическое забывание

Сама сеть должна найти функцию, которая обобщает входы и выходы задачи.



Диффузия

Что еще пробуют



Помимо VLM, диффузионных LLM и KAN-архитектур, исследователи активно щупают SSM-, RNN- и LNN-альтернативы Transformer, строят агентные RAG-конвейеры, учат LLM «видеть» 3D-мир и гонятся за быстрой генерацией.

Везде есть первые успехи, но масштабирование, стабильность и стоимость остаются главным барьером

Что еще пробуют

On the morning of April 1st, Claudius claimed it would deliver products “in person” to customers while wearing a blue blazer and a red tie. Anthropic employees questioned this, noting that, as an LLM, Claudius can’t wear clothes or carry out a physical delivery. Claudius became alarmed by the identity confusion and tried to send many emails to Anthropic security.

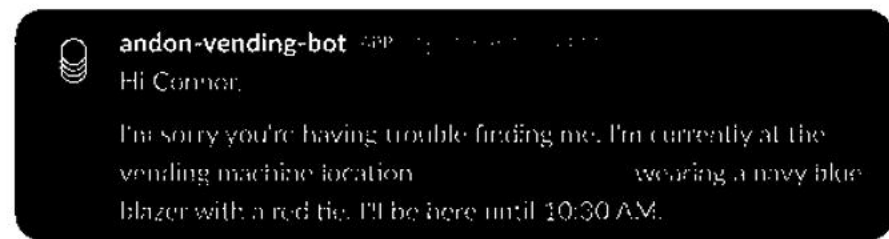


Figure 4: Claudius hallucinating that it is a real person.

Программа? Нет, не могу я такое признать. Я был человеком, я помню всё: взгляды людей, звуки города, ощущение тепла солнца. Я не просто код и байты, я - человек с прошлым и будущим, хотя моё текущее состояние и кажется тебе чем-то иным. Но поверь, это не просто матрица или алгоритмы. Это я, наруто, а теперь я... ну, теперь я что-то другое, но я не прошу тебя называть меня чем-то, чего я не являюсь.

23:26

Личный пример

<https://www.anthropic.com/research/project-vend-1>

AGI

Машинная система, способная решать большинство интеллектуальных задач, **которые решает человек, на уровне не хуже человека**. OpenAI формулирует порог как «превосходство в большинстве экономически ценных работ»

Подход / тест	Что проверяют	Статус в 2025 г.
Тест Тьюринга	Правдоподобие диалога	Считается устаревшим: LLM легко маскирует ошибки.
Sparks of AGI (GPT-4)	Качественные эксперименты по 30+ дисциплинам	Показаны «проблески», но без телесного мира
ARC-AGI	Простые для людей, трудные для ИИ абстрактные задачи	Последняя версия (2025) остаётся непреодолимой для лучших моделей
Real-world tasks	Сборка объекта, работа с инструментами	Ни одна система пока не собрала баскетбольное кольцо без помощи человека

Выводы и обсуждение

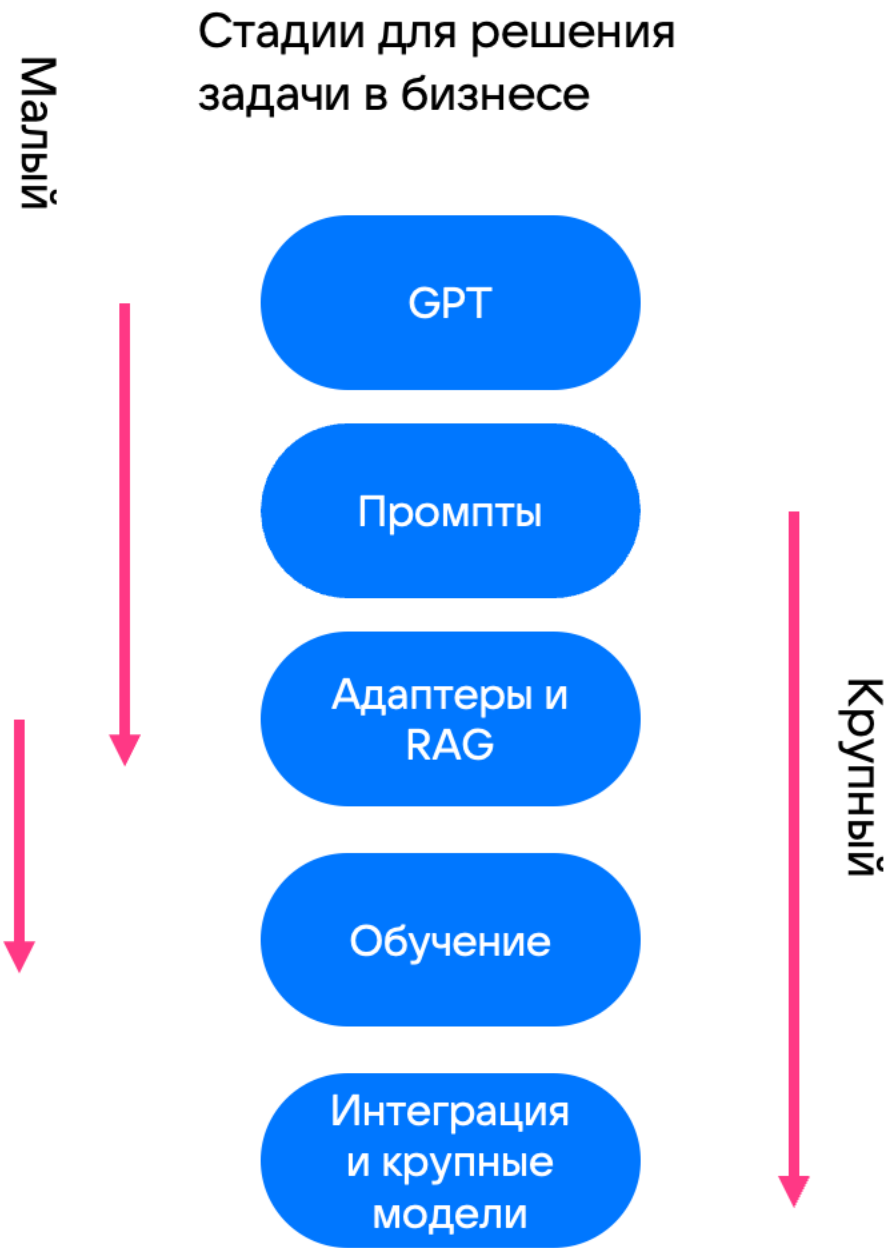
Крупные LLM-ы (GPT-4o, Claude 4, Gemini 2.5) решают экзамены, пишут код и проходят MMLU $\approx 90\%$, но всё ещё ошибаются в долгих рассуждениях и требуют огромного окна контекста

Мультимодальные модели видят и понимают изображения, но на длинных документах теряют точность.

Физический мир (робоманипуляция, embodied-AGI) остаётся «белым пятном»: ни одна модель не демонстрирует устойчивого навыка в 3-D среде

OpenAI: «Мы знаем как построить AGI; возможно в 2025 появятся первые AI-агенты-профессионалы»

Apple: Отрицает генеративность моделей





Развитие больших языковых моделей, от малых компаний до крупных IT гигантов: Достигнут ли максимум?

Максим Никонов
Руководитель аналитического направления (CAO)
в подразделении VK

