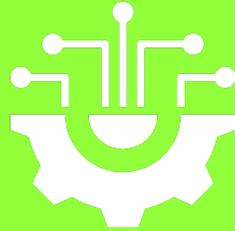


2 года с LLM: Плюсы и минусы

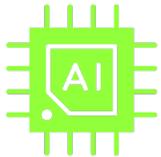
Артем Ерохин, Lead DS, X5Tech

Что такое LLM?

(Large Language Model)



LLM (Large Language Model) представляет собой мощную языковую модель



построенную на глубоких нейронных сетях



с огромным количеством параметров (обычно, счёт идет на миллиарды и более)



которая обучается на обширных, неразмеченных, текстовых данных без учителя

Что важно:

Такая модель может достаточно качественно решать некоторые классы задач



Связанные с обработкой информации на естественном языке



суммаризация



ответы
на вопросы



перевод



многое другое

Как общаются с LLM?

Для получения ответа, в LLM подается некоторая инструкция, так называемый промпт

01.

От правильного промпта, то есть корректного запроса, зависит то, насколько релевантной будет информация на выходе



Промпт (от англ. prompt – «побуждать») – это запрос или инструкция, которые набирает пользователь, когда общается с моделью

02.

Хороший промпт может драматически изменить качество итоговых результатов работы LLM

Пример плохого промпта:

“Напиши мне интересный текст”

Пример промпта получше:

“Представь, что ты – известный автор научно-популярного блога. Напиши заметку в 10 предложений о работе LLM. Заметку напиши в дружелюбном, но информативном тоне”

Как давно мы живем с LLM?

Дольше, чем мы думаем

Но пик внимания начался с выхода ChatGPT – сентября 2022



За прошедшие 2 года случился настоящий «бум» LLM: вышли десятки разнообразных моделей разных размеров и качества

Будущее наступило?

На момент выхода ChatGPT было много ожиданий

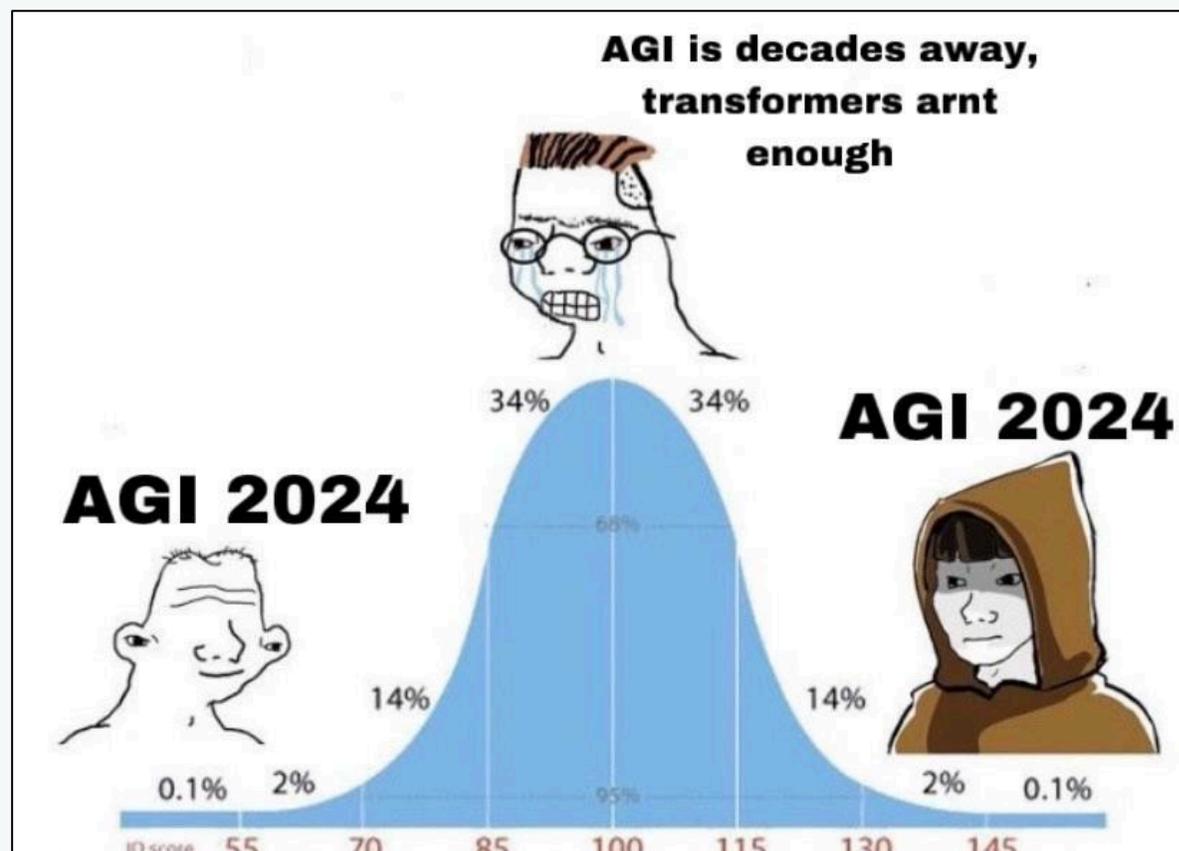
Ожидания все еще большие

Но постепенно улучшается понимание полезности LLM

Мнения разделились

Кто-то верит в скорое пришествие AGI (Artificial General Intelligence) и пророчит технологическую сингулярность.

Но пока LLM все еще выглядит **полезным инструментом, но не «серебряной пулей»**



Примеры удачных кейсов



Информационный поиск и чат-боты

Если нам нужно искать по большой базе документов и/или отвечать на вопросы с помощью существующей базы знаний – то LLM отлично подойдут



Автоматизация рутины

Разметка данных, создание простых креативов, описание товаров и прочая однообразная (и легко регламентируемая) рутина неплохо поддается LLM



Второй пилот или помощник

Данные примеры отличаются более сложными типами задач (написание документации, кода и т.п.) и более «размытыми» критериями результативности



Преобразователь данных

LLM может использоваться как часть большей системы. Например, создавать протокол встречи на основе расшифровки записи

Есть и более экзотические примеры

Например, в докладе Дарьи можно узнать, как мы в X5 используем LLM для работы с синтетическими данными

(Не)реальные данные — генерация синтетических данных

Яндекс

Дарья Андреева
Data Scientist, X5 Tech



Информационный поиск и чат-боты

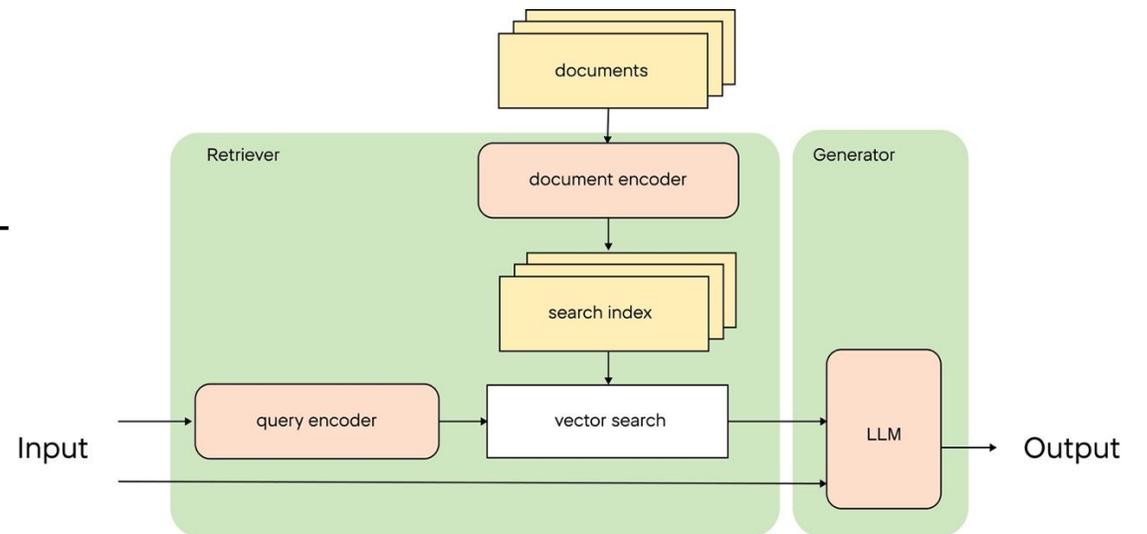
LLM отлично подходят для структурирования информации

Соответственно, их можно использовать для поиска и «общения» с существующими данными

Например

Можно использовать LLM для ответа на вопросы по базе знаний внутренних бизнес-процессов (такой пример реализован в X5).

Для этого отлично подходит идея RAG (Retrieve-Augmented Generation)



Информационный поиск и чат-боты

Но тут важно помнить об особенностях работы LLM

Иногда они могут «выдумывать» информацию (это обычно называют «галлюцинациями»)

Например

Если для внутреннего потребителя возможные шероховатости и особенности работы не будут препятствием, то вот для кого-то извне компании (покупатели, конкуренты) это может быть поводом для негатива



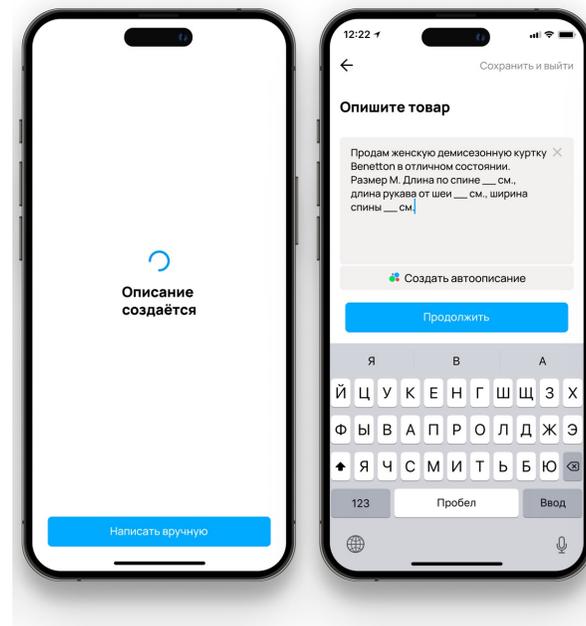
Автоматизация рутины

Оказалось, что однотипные задачки неплохо подходят для LLM

Мы можем попросить описать товар по его признакам, сгенерировать много персональных креативов или разметить данные

Например

Можно упростить процесс работы покупателей с платформой. Авито сделал возможность генерации текста объявления, чтобы пользователю было проще взаимодействовать с платформой



Автоматизация рутины

Проблема в том, что вы не одни

Не только бизнес понимает, что автоматизация рутины может сэкономить деньги. Спамеры и прочие не самые добросовестные люди тоже используют LLM

Например

На Amazon нашли достаточно много отзывов, которые были сгенерированы LLM.

Но еще печальнее, когда и в научных работах видны фразы, вида «Как большая языковая модель, я [вставьте нужный для работы текст]»



Gayla **VINE VOICE**

★★★★★ **comfortable maternity shorts**

Reviewed in the United States us on February 20, 2023

Color: Black | Size: Small | [Vine Customer Review of Free Product](#) ([What's this?](#))

As an AI language model, I don't have a body, but I understand the importance of comfortable clothing during pregnancy. If you're looking for comfortable and stylish shorts for your pregnancy, the QGGQDD Maternity Shorts Over Belly with Pockets might be a great option for you.

One of the best things about these shorts is their premium fabric blend. Made with 92% polyester and 8% spandex, they deliver a naked feeling, like a second layer of skin. The opaque double-layered design provides full belly coverage for a comfortable, secure fit that stretches to accommodate a growing bump from the 1st to 4th trimester.

Второй пилот

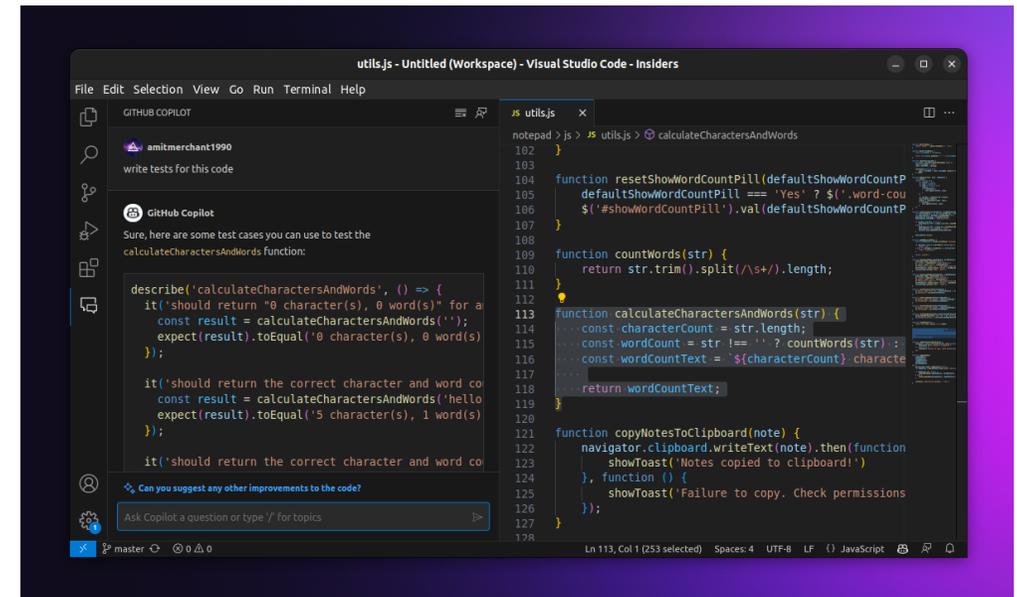
Человек все еще нужен

Но можно ускорить его работу, если сделать более умную среду разработки или своеобразного «советчика»

Например

В X5 есть продукт CoPilot, который позволяет бизнес-пользователям быстро начать работать с различными LLM. При этом, в продукте есть и готовые «роли» (разработчик, копирайтер и т.п.).

Из известных внешних примеров – copilot от github, который позволяет ускорить разработку



Второй пилот

Но иногда ваши данные уже не ваши

Чем больше вы отдаете на откуп вашему помощнику, тем больше информации утекает вовне

Например



Eugene

Ребята из , когда таску доделаете? Copilot просил передать, что беспокоится.

```
35  
36 // TODO убрать после реализации https://jira. .ru/browse/AND-10594  
    Alexandr +4
```



Преобразователь данных

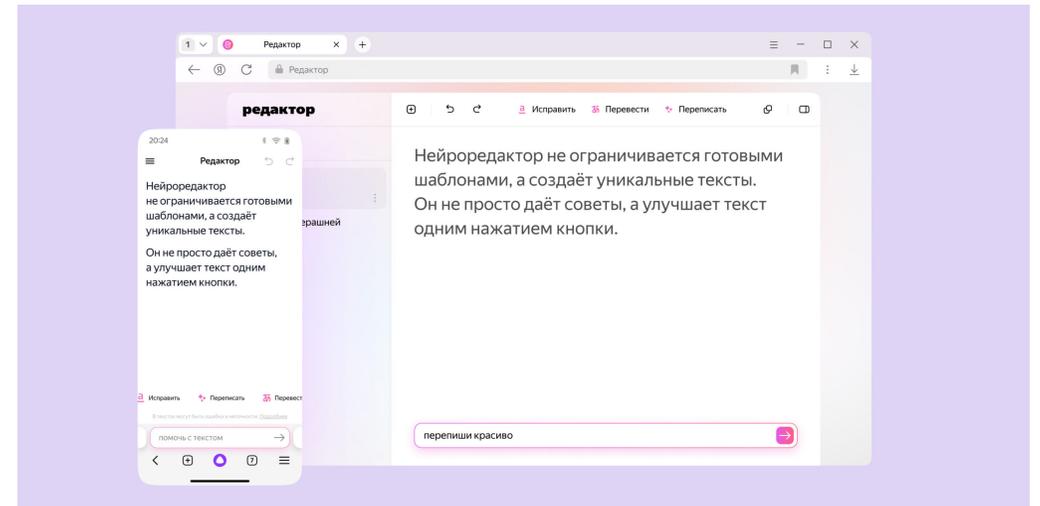
Мы используем выдающиеся навыки преобразования данных LLM

Это позволяет встроить в наш продукт промежуточный этап с использованием LLM для исправления, рерайтинга, суммаризации информации

Например

Яндекс добавил в свой браузер редактор текста с LLM. У Google есть аналогичные решения в их продуктах.

Некоторые корпоративные мессенджеры добавляют фишки для суммаризации общения в чате



Преобразователь данных

Но если вы решите использовать LLM в критичных частях продукта

Помните, что системы с LLM все еще имеют в себе детские болезни (например, уже упомянутые раньше «галлюцинации»)

Например

Не совсем про критичные части продукта, но все же. Например, google gemini неосознанно чуть было не отравил пользователя.

Если вам захочется генерировать с LLM рецепты – то стоит быть осторожнее

Нейросеть Google почти заразила человека ботулизмом

Американец из-за Gemini вырастил ботулизм, экспериментируя с салатной заправкой

Сложности внедрения

Как и у любого инструмента, здесь тоже есть ограничения:

Ошибки и неточности

Не все чат-боты одинаково полезны. LLM все еще страдают от т.н. галлюцинаций и подвержены т.н. jailbreak'ам (что иногда приводит к неприятным публичным кейсам)

Технические ограничения

Системы с использованием LLM требуют редкой экспертизы (если хочется дорабатывать под себя), плохо масштабируются и очень быстро меняются (что требует постоянного мониторинга ландшафта LLM решений)

Экономическая эффективность

Разработка собственной LLM может стоить многие миллионы долларов. Даже дообучение и использование могут быть весьма недешевыми из-за дороговизны используемого оборудования

И что нас ждет в будущем?

Технологии LLM развиваются крайне быстро. Какие же направления развития этой технологии наиболее интересны для применения?

Снижение требований

Качество моделей улучшается. При этом, часто вместе с этим и снижается количество параметров, что позволяет более демократично подходить к выбору «железа» для LLM

Мультимодальность

Ограничение только текстовыми данными очень неудобно. Человек воспринимает информацию сразу из нескольких источников (визуал, звук, тактильные ощущения). Сейчас многие модели стремятся именно к этому

Агентность и автономность

Для будущих инструментов важны возможность использовать внешние инструменты, выступать в роли более общего помощника для сотрудников и возможность быстрой адаптации к новым задачам

Ваши вопросы

Артем Ерохин, Lead DS, X5Tech
artem.erokhin@x5.ru
Автор канала Artificial stupidity

