

# Информационный поиск



LLM and RAG

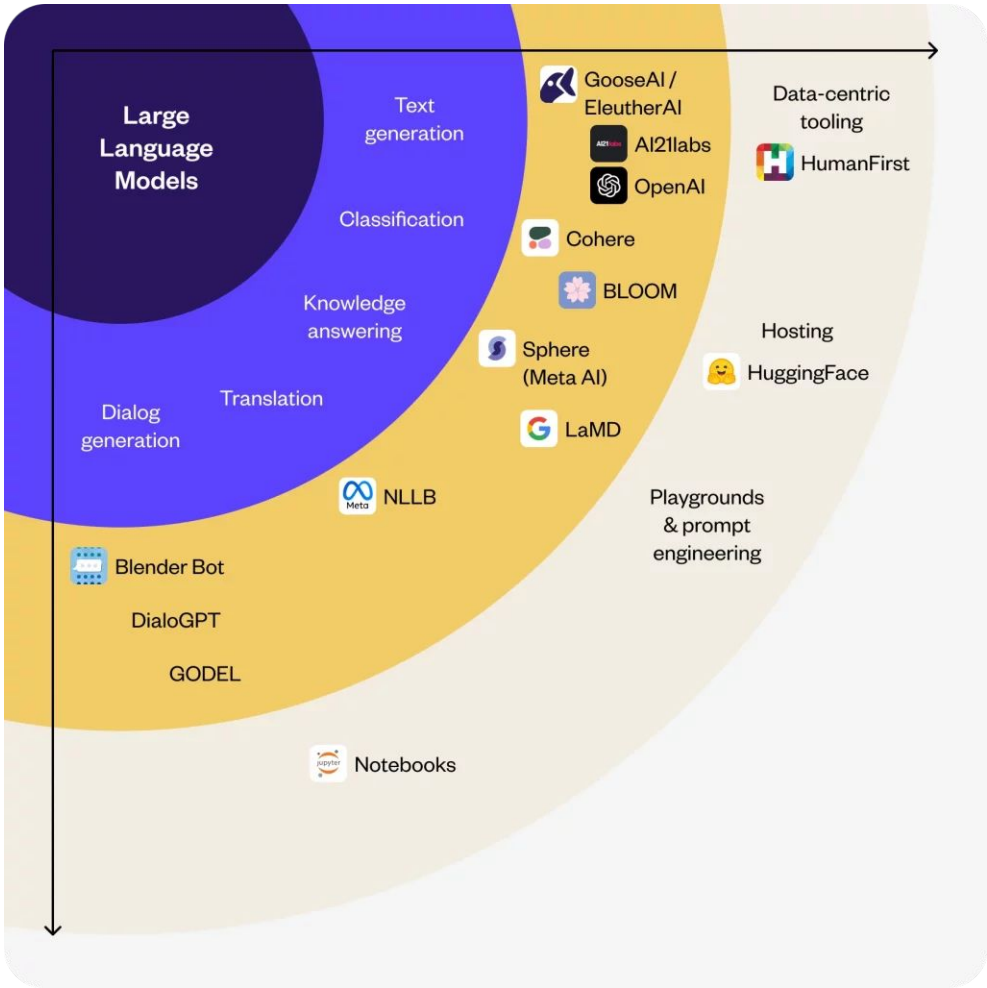


# Large Language Models

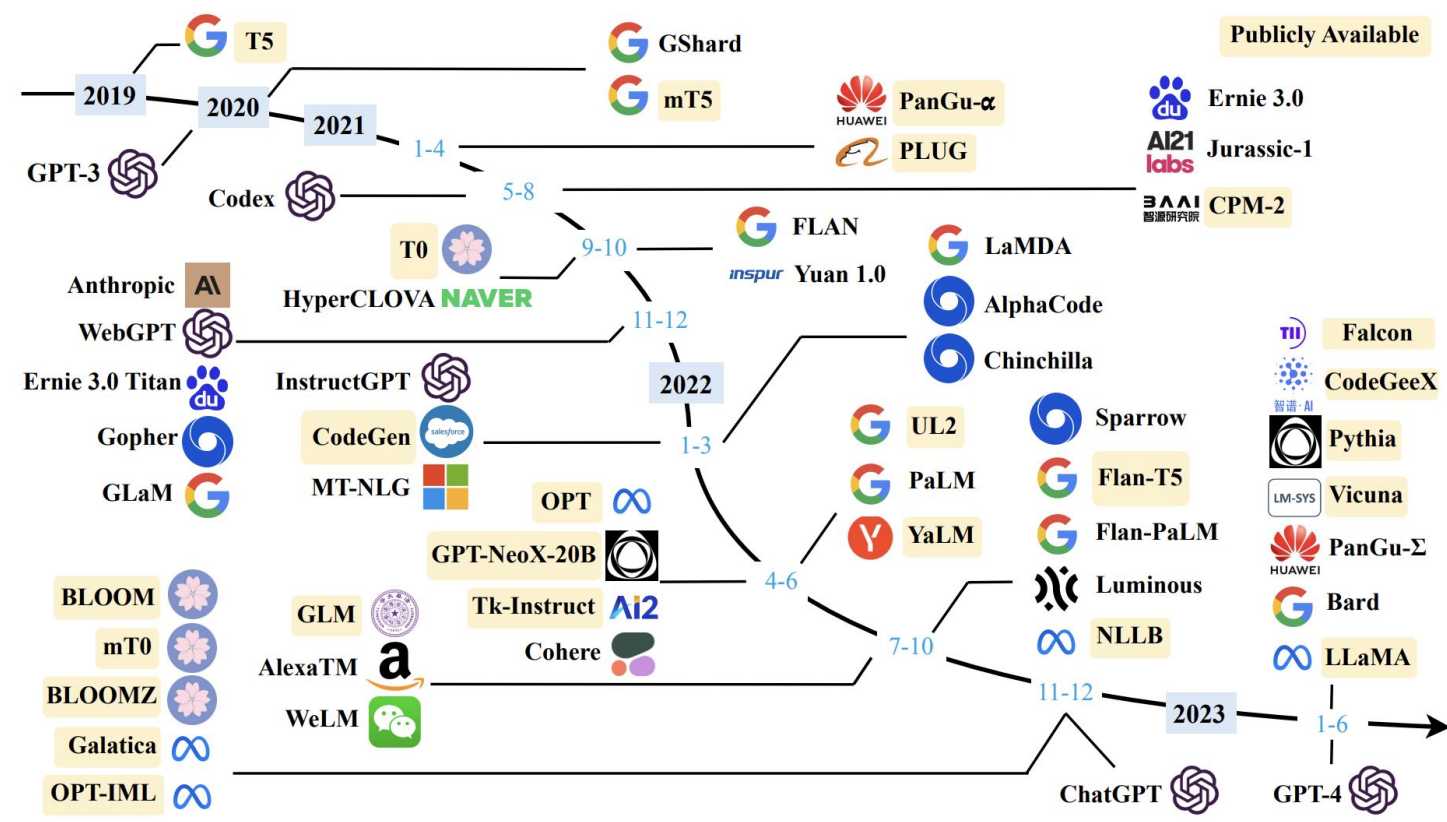
На данный момент большие языковые модели являются State-of-the-Art решениями во многих областях NLP (и не только).

При этом вокруг них активно выстраивается целая экосистема: API, специальные библиотеки, приемы использования и дообучения на малых мощностях.

Мы сегодня поговорим про LLM в целом, а также про приемы промптинга и RAG.



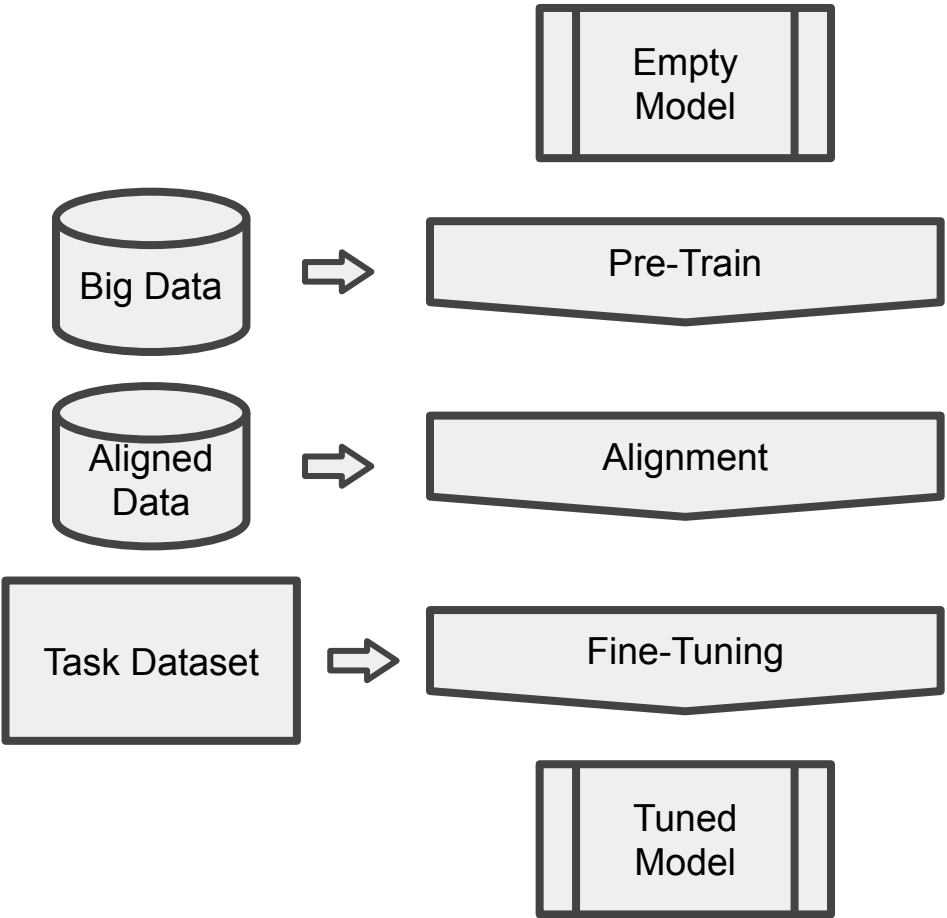
# Large Language Models



# Процесс обучения

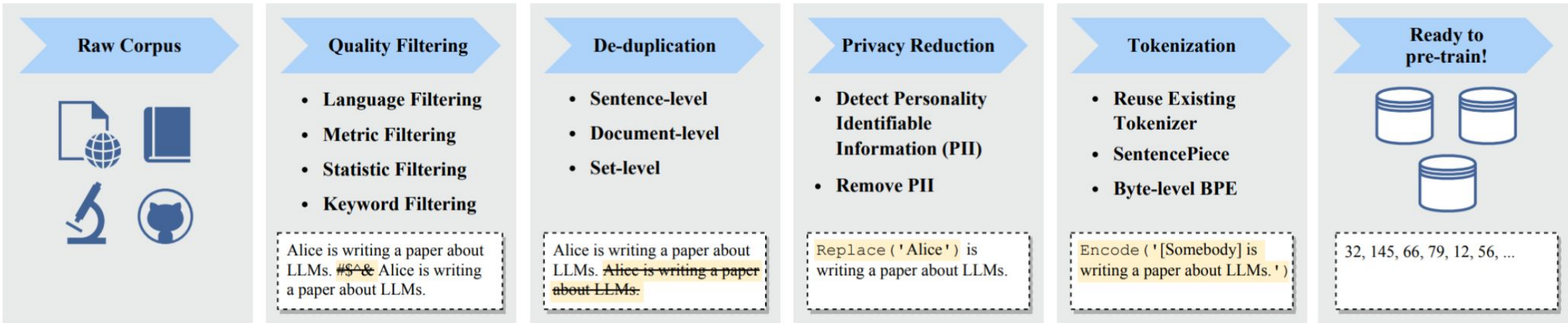
Обычно процесс разделяют на три больших этапа:

1. **Pre-train:** на данном этапе происходит основное обучение, о котором мы говорили выше. Длинный и вычислительно дорогой процесс
2. **Alignment** (обычно выделяют для llm): здесь модель учится соблюдать определенный набор правил, “выравнивается” под некоторое распределение
3. **Fine-tuning:** дообучение модели на конкретную задачу



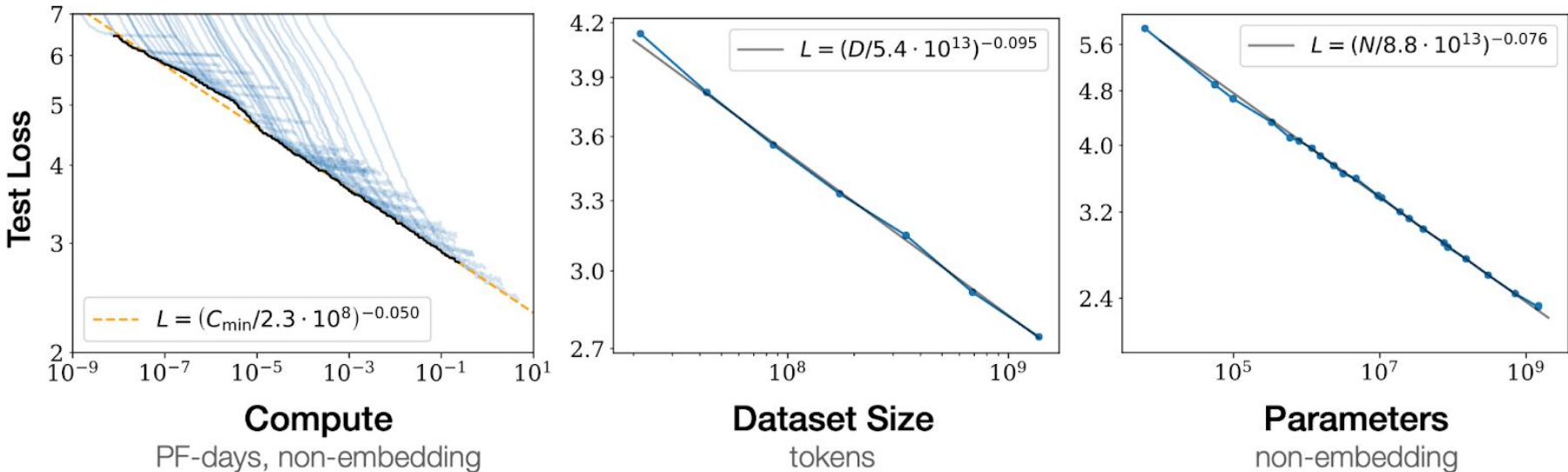
# Процесс обучения: pre-train

1. Обычно говорят, что такие модели обучаются на сырых данных, однако препроцессинг все-таки нужен:
2. Бывают дубли в данных, которые увеличивают время обучения, но не дают информации
3. Персональные данные лучше не включать в такие модели
4. Предложения на других языках тоже



# Параметры обучения

1. Обучение модели обычно существенно зависит от трех параметров: количества итераций обучения, размера датасета и размера самой модели.
2. И, конечно, архитектуры, но тут чаще всего выбор между Трансформером и Трап





# Вот как-то так



# Процесс обучения: alignment

На данный момент есть целая активно развивающаяся область на стыке AI, информационной безопасности и этики, которая занимается вопросами элаймента.

## Безопасность и этика

Учим модель правильно реагировать на запросы:

- Никакой симметричной реакции на негатив
- Нельзя разглашать персональные данные

## Реакция на инструкции

Современная llm должна уметь следовать произвольным набором инструкций:

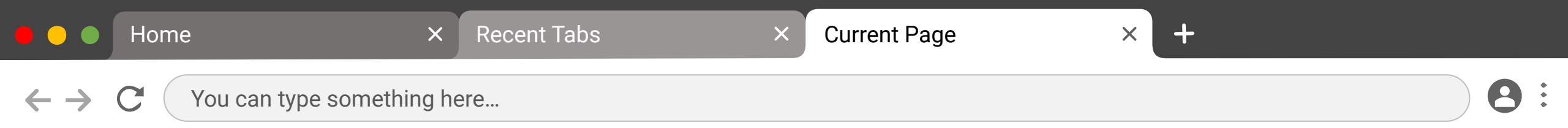
- "Напиши код для задачи"
- "Расскажи рецепт тирамису"

## Выбросы и смещения

Модель не должна галлюцинировать и выходить из строя на странных примерах:

- Мир меняется
- Пользователи не всегда пишут что-то осмысленное.





# Процесс обучения: fine-tuning

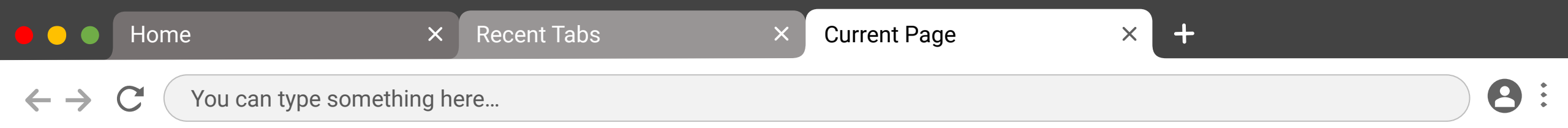
Существует множество задач, на которые можно дообучать модель.

## Encoder + Classification Head

- Natural Language Inference
- Sentiment Analysis
- Spam Detection
- Hate Speech Recognition
- Topic classification

## Decoder | Encoder-Decoder

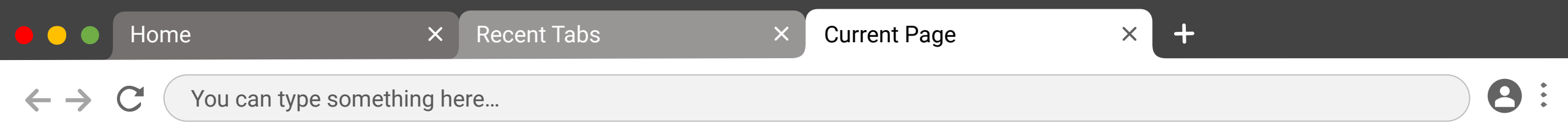
- Summarization
- Question Answering
- Machine Translation
- Style Transfer
- Text Generation in general



# Промптинг

Если модель была дообучена на следование инструкциям, то ее называют instruction-tuned (чаще всего обозначается -it в названии модели). Такие модели лучше всего реагируют на промнты и письменные инструкции, так что для всего, что будет обсуждаться дальше, лучше использовать их.

Промптинг - это процесс подбора затравки (промпта) для модели так, чтобы результат генерации лучше всего соответствовал ожиданиям. Он бывает ручной, когда мы просто придумываем промпт, и автоматический, куда относятся техники RAG, обучение моделей-промптеров и прочие похожие подходы.



# Промптинг: виды

## Самые частые виды промптинга:

- Zero-shot
- Few-shot
- Chain-of-Thought (CoT)
- Self-Consistency
- Retrieval Augmented Generation (RAG)

*Prompt:*

```
Classify the text into neutral, negative or positive.  
Text: I think the vacation is okay.  
Sentiment:
```

*Output:*

```
Neutral
```

*Prompt:*

```
This is awesome! // Negative  
This is bad! // Positive  
Wow that movie was rad! // Positive  
What a horrible show! //
```

*Output:*

```
Negative
```

# Промптинг: Chain-of-Thought

Вместо того, чтобы сразу требовать ответ на поставленный вопрос, мы даем модели пример рассуждений и просим от нее думать похожим образом.

На практике такой метод оказывается дороже (почему?), но эффективнее, особенно в математических задачах.

### Standard Prompting

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

### Chain-of-Thought Prompting

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✅

# Промптинг: Self-Consistency

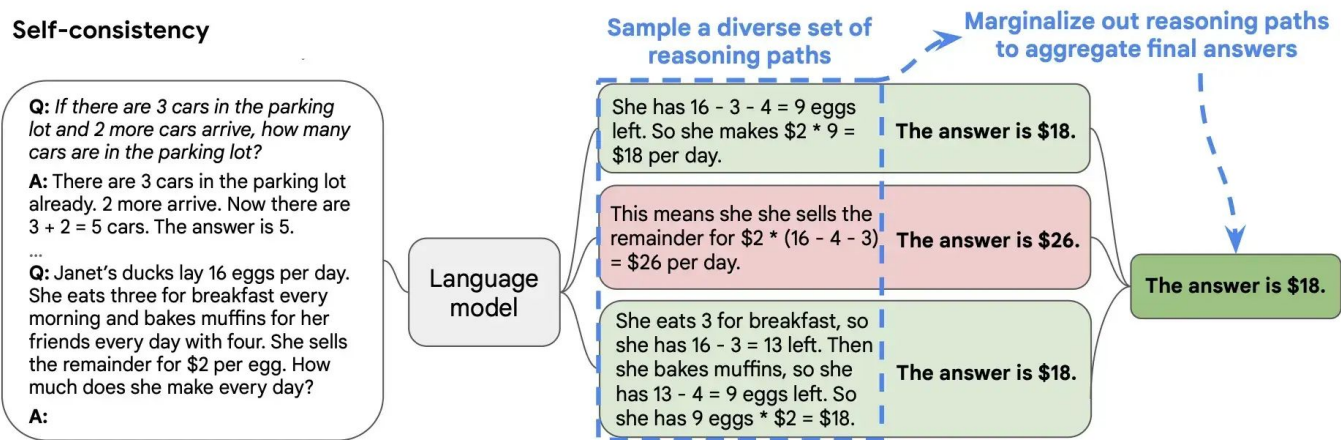
Похоже на цепочки рассуждения, но здесь мы генерируем несколько вариантов таких цепочек и потом берем средний (или самый частотный) результат их полученных.

Плюсы:

- Можем избежать принятия галлюцинации как правильного ответа, так как она менее вероятно повторится

Минусы:

- Еще более дорогой Self-consistency



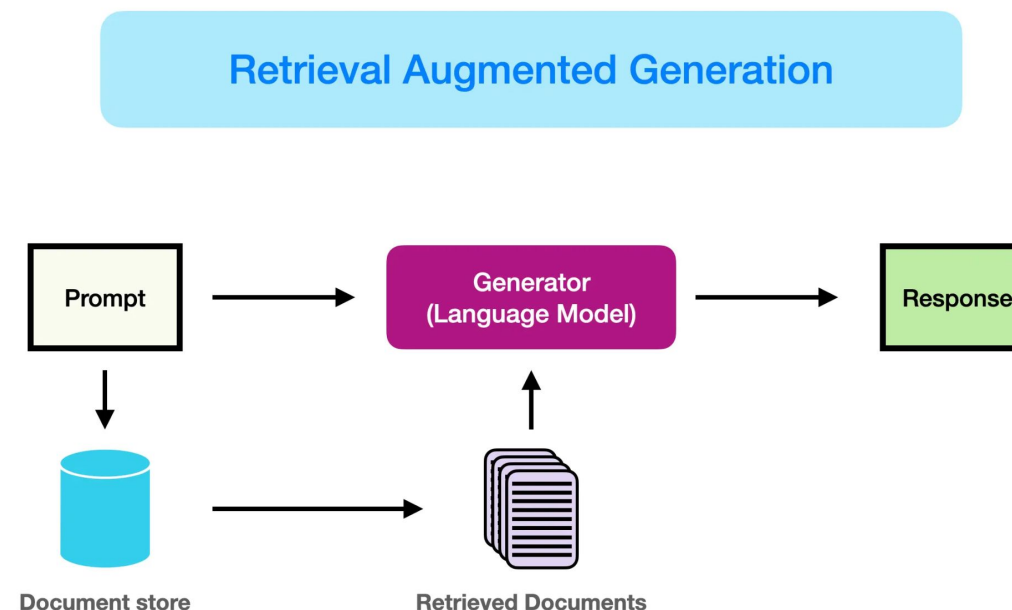
# Промптинг: Retrieval Augmented Generation

Здесь мы улучшаем генерацию с помощью стороннего источника знаний.

Как это работает:

- Векторизуем все документы в нашей базе
- Когда нам поступает промпт, ищем топ-n ближайших документов
- На вход LLM подаем промпт плюс найденные документы

Какие здесь есть трудности?

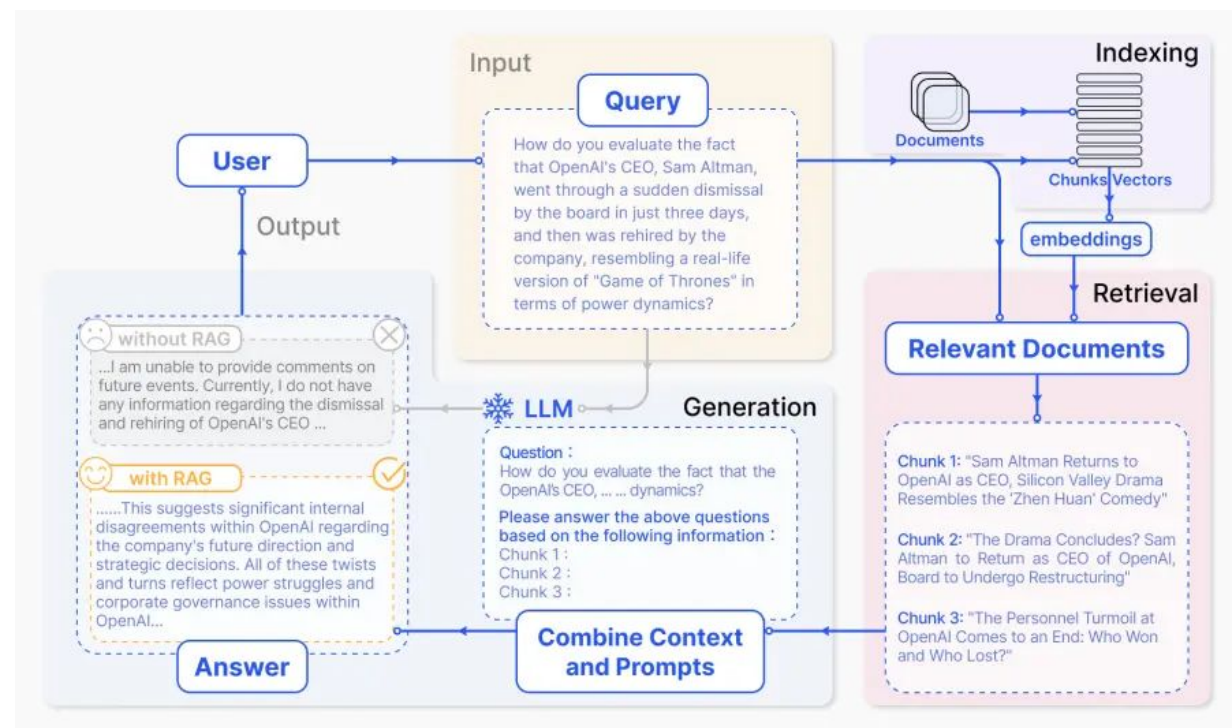




# Промптинг: Retrieval Augmented Generation

Хитрости:

- Можно делить тексты на чанки так, чтобы не разрывать информацию, но при этом сократить размер промпта
- Если использовать для векторизации саму LLM, то можно попробовать пропустить какое-то кол-во вычислений (почему это на практике сложно реализовать?)
- Иногда можно давать не все параграфы, а заголовки/результаты суммаризации



# Если нет LLM

В чем-то все описанное похоже на задачу QA, просто используется большая и тяжелая модель, которую не надо дообучать. Но если у вас есть время и данные, но нет LLM, то можно обучить генеративную модель вроде T5, чтобы она генерировала вам ответы по нескольким абзацам, которые вы будете подбирать по логике RAG.

