

# **State-of-the art техники RAG на примере победного решения Enterprise RAG Challenge и его адаптации**

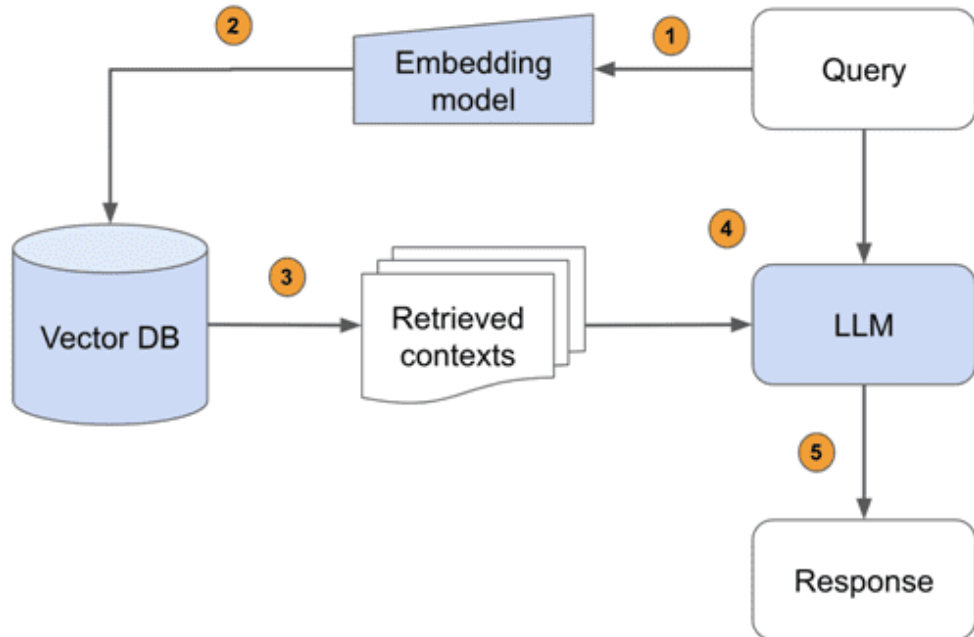
Акимов Дмитрий Андреевич

# Актуальность RAG

## работа базовой LLM



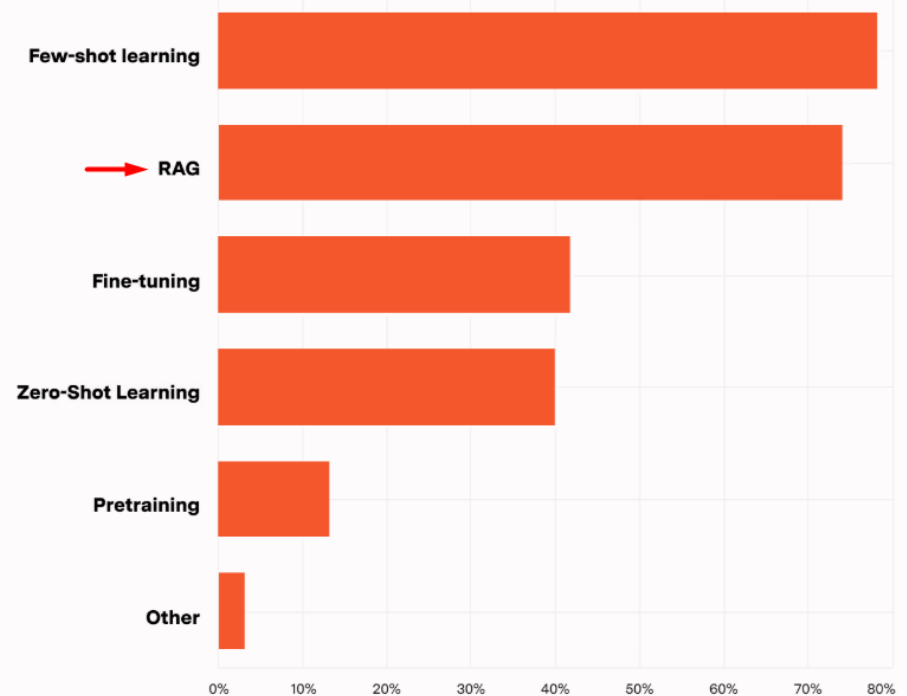
## расширение LLM с помощью RAG



Engineering and Infrastructure

## The 2025 AI Engineering Report

What techniques are you using to customize your AI systems?



# Retrieval Augmented Generation (на пальцах)

Вопрос пользователя

Сколько Васе лет?

$\{0; 0.6; 0\}$

нарезаем на части («чанки»)

переводим в векторы с  
помощью модели эмбедингов

## Наши данные

Вася играет в футбол.  
Васе 15 лет. Вася  
хорош в математике.

Вася играет в футбол.

$\{0; 0; 0.8\}$

$|\{0*0; 0*0.6; 0.8*0\}|=0$

Васе 15 лет.

$\{0; 0.7; 0\}$

$|\{0*0; 0.7*0.6; 0*0\}|=\mathbf{0.42}$

Вася хорош в математике.

$\{0.5; 0; 0\}$

$|\{0.5*0; 0*0.6; 0*0\}|=0$

ищем косинусное сходство

## Запрос LLM

Сколько Васе лет?  
Найденный контекст:  
Васе 15 лет

## Ответ LLM

Васе 15 лет



# Retrieval Augmented Generation (если погружаться)

HyDE

RAG-Fusion

Haystack

LlamaIndex

Langchain

Knowledge Graph

MMR TrueLens

RAGAS

Agentic RAG

Recall@k



BM25 + TF-IDF

Contextual  
Retrieval Neo4j

Cross-Encoder  
Reranking

HNSW  
IVF-Flat

AutoRAG

Multi-Query

ColBERT

DeepEval

CRAG

# Enterprise RAG Challenge v2 (2025)



Статья получила международную популярность и была переведена на китайский:



Конкурс на деньги (как и в случае с Kaggle) – часто лучший способ определить state-of-the-art техники.

## В чём суть RAG Challenge?

Нужно создать вопросно-ответную систему на основе годовых отчётов компании. Если коротко, то в день конкурса:

1. Выдаётся 100 годовых отчётов по случайно выбранным компаниям и 2.5 часа на их парсинг и составление базы данных. Отчёты представляют из себя PDF размером до 1000 страниц.
2. После этого генерируется 100 случайных вопросов (по заранее известным шаблонам), на которые система должна как можно быстрее ответить.

Все вопросы должны иметь однозначный ответ:

- Да/Нет;
- название компании (или нескольких компаний);
- названия управляющих позиций, выпущенных продуктов;
- размер той или иной метрики: выручка, количество магазинов и т.д.

Каждый ответ должен сопровождаться ссылками на страницы с ответом в качестве доказательства, что система по честному нашла ответ и не сгаллюцинировала.

# Enterprise RAG Challenge v2 (условия)

100 PDF годовых отчетов компаний (случайно отобранных)

100 случайных вопросов (Да/Нет; Числовые показатели и т.д.)

Оценивался Retrieval (верная ссылка в ответе) и Generation (верный ответ)

Итоговая метрика  $R/3 + G$  (т.е. верный ответ оценивался в 3 раза выше верного цитирования)



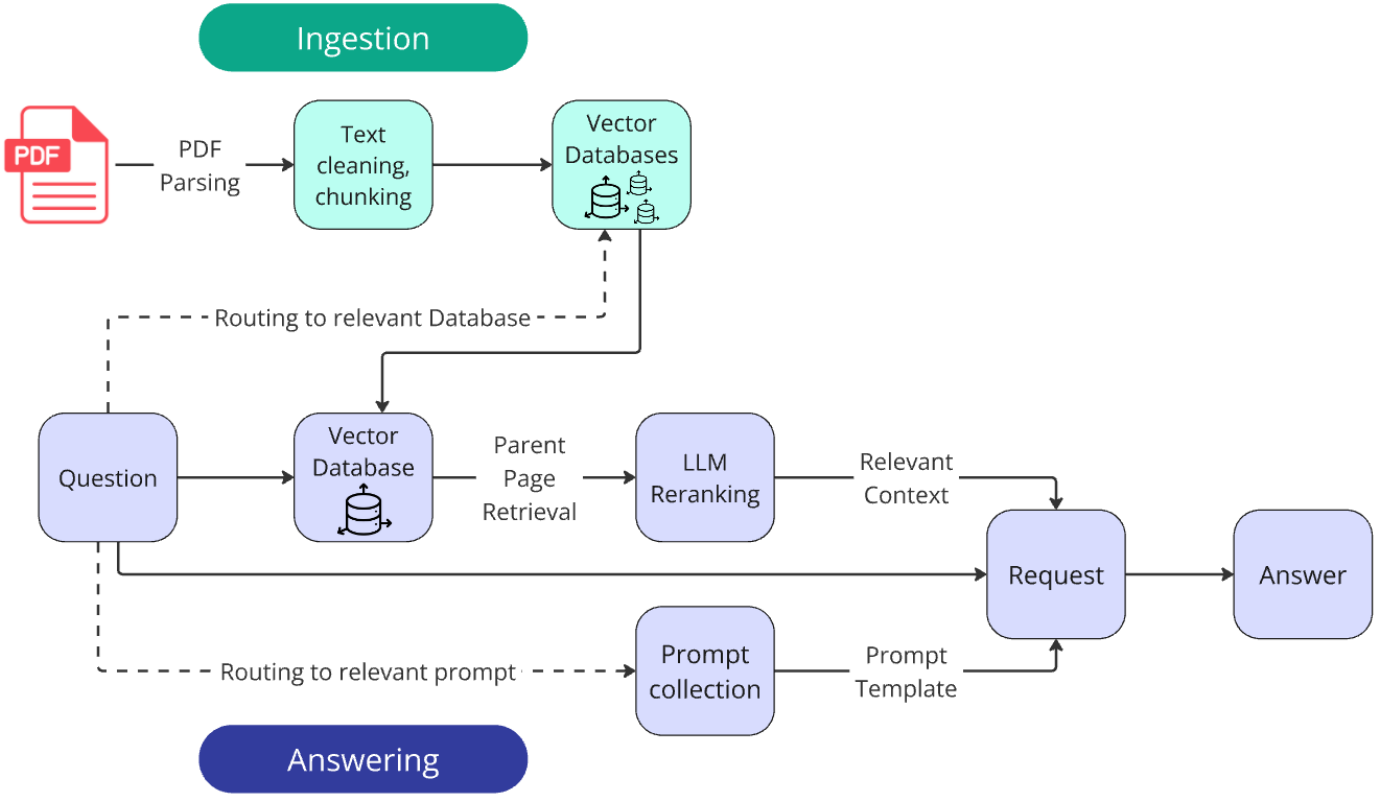


# Enterprise RAG Challenge v2. Победное решение (кратко)

Team / Experiment		Time	R/G	Score
1. ▶ Ilia Ris	🥇	49 min	83/81	123.7
2. ▶ Emil Shagiev	🥈	55 min	86/78	121.6
3. ▶ Dmitry Buykin	🥉	8 hours	81/76	117.5
4. ▶ Sergey Nikonov	🥈	30 hours	85/73	116.4
5. ▶ ScrapeNinja.net	🥈	23 hours	82/71	112.5
6. ▶ xsl777	🥈	16 hours	79/71	110.9
7. ▶ nikolay_sheyko(grably.tech)	🥈	25 hours	81/69	110.4
8. ▶ Felix-TAT	🥈	7 days	80/69	109.4
9. ▶ A.Rasskazov/V.Kalesnikau		30 hours	84/67	109.3
10. ▶ Dany the creator	🥈	3 hours	82/67	108.4
11. ▶ SergC	🥈	7 days	77/69	108.1
12. ▶ Swisscom Innovation Lab	🔒	21 hours	83/66	107.8
13. ▶ fomih	🥈	10 days	83/65	107.4
14. ▶ Al Bo		12 days	81/65	105.9
15. ▶ NumericalArt		8 days	70/70	105.3

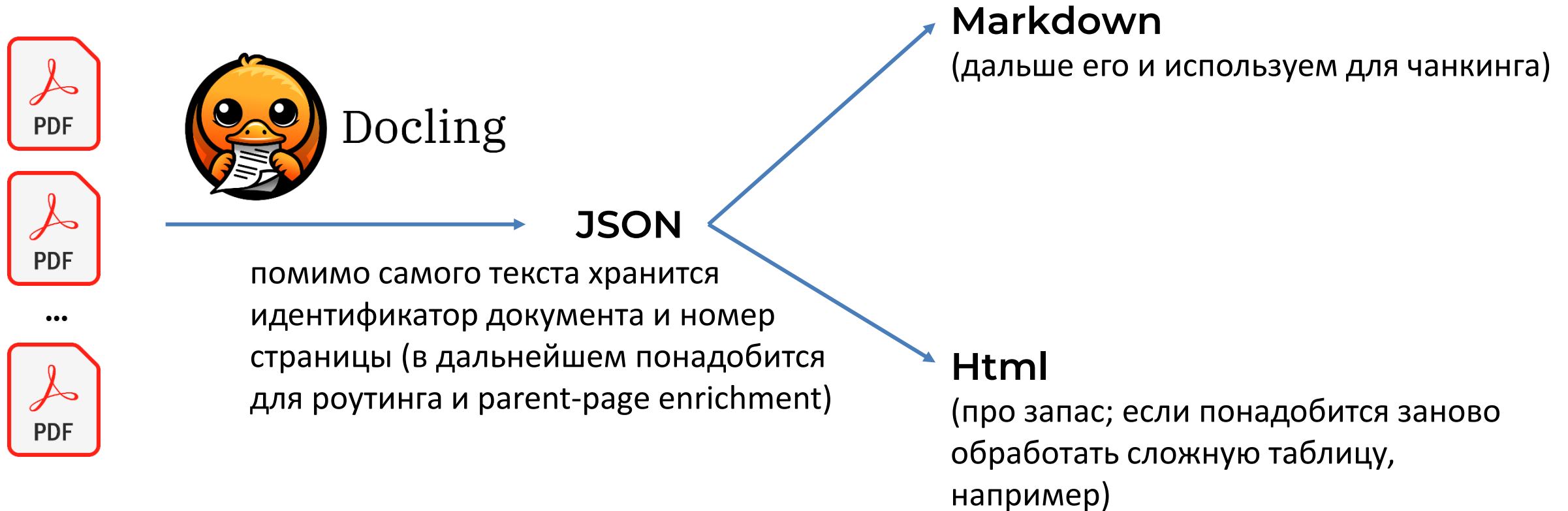
всего около 45 команд

## Схема победного решения Ilia Ris



На доп. треке IBM (использование open-source LLM Llama 2 70B) качество просело менее, чем на 2%!

# Parsing

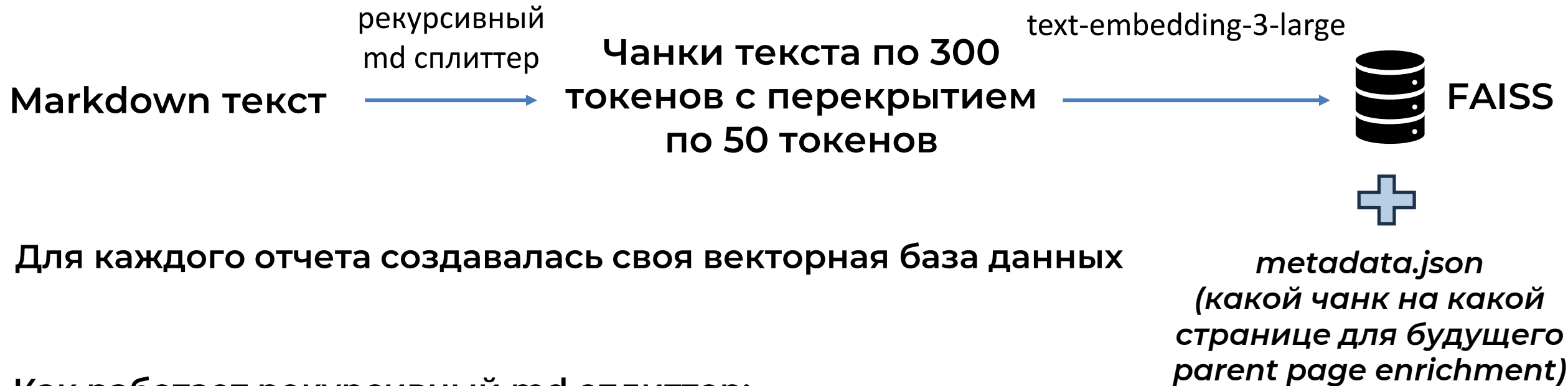


Илья использовал Docling, арендовав сервер на Runpod. Можно также использовать Unstructured, Marker, LlamaParse и пр.

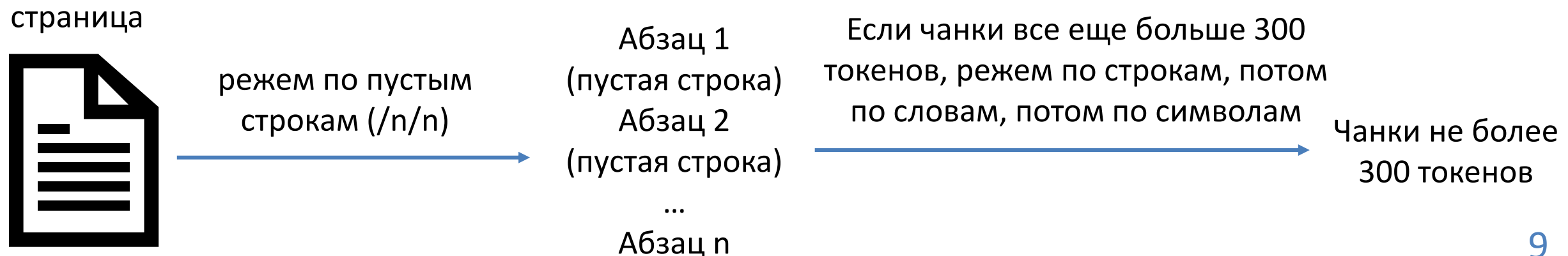
В своем упрощенном примере я использовал pdfplumber.



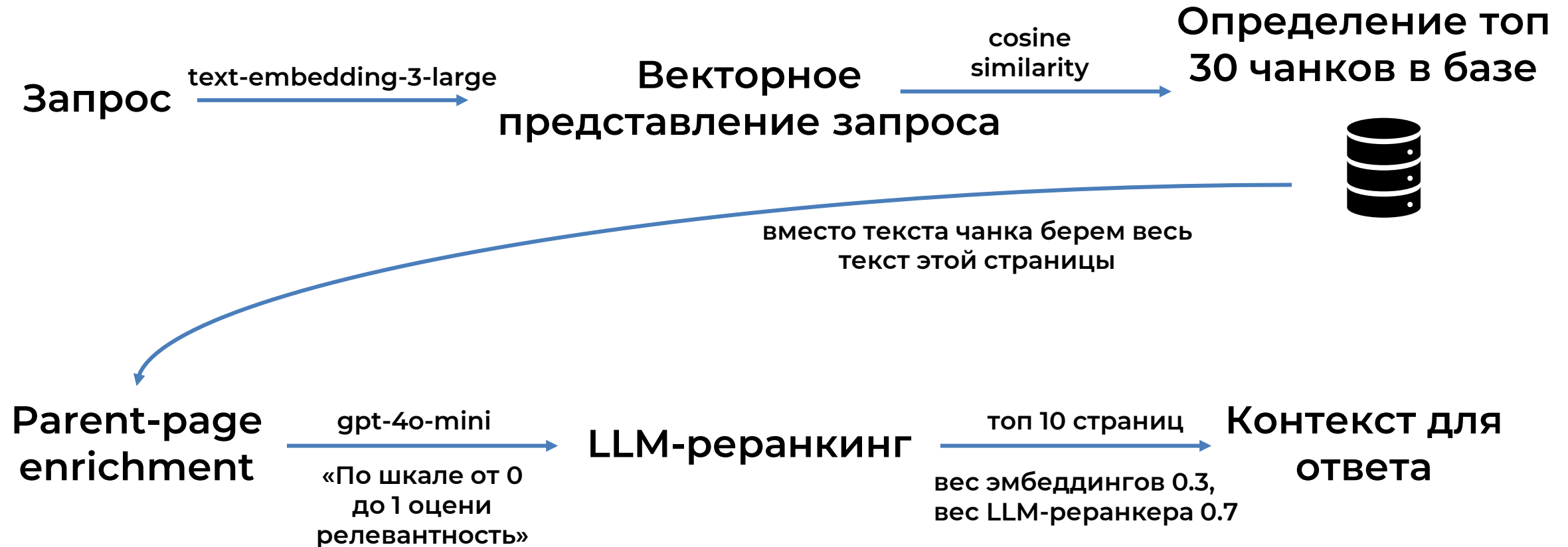
# Ingestion



Как работает рекурсивный md сплиттер:



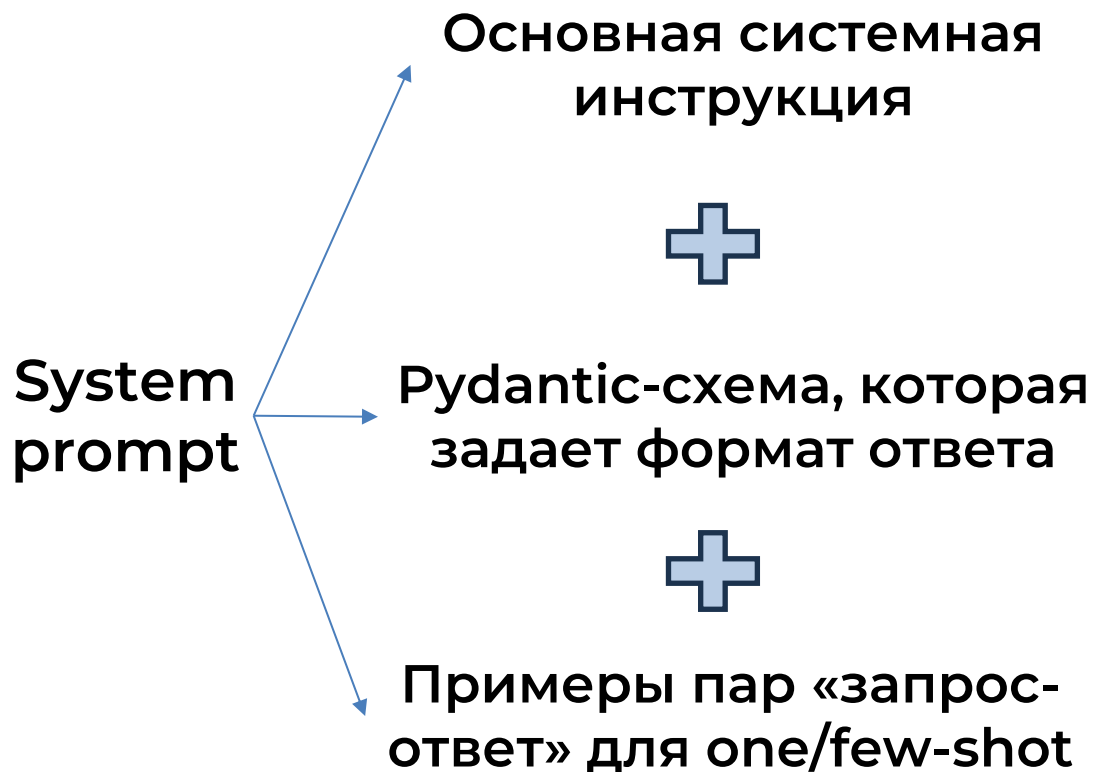
# Retrieval



LLM-переранкинг – перспективная новинка! Обычно используется cross-encoder reranking.

Серебряный призёр вообще построил RAG без эмбеддингов.

# Augmentation



User prompt → Template для вставки контекст и вопроса

Ты — RAG (Retrieval-Augmented Generation) система ответов. Твоя задача — отвечать на заданный вопрос, опираясь ТОЛЬКО на информацию из годового отчёта компании, который предоставляется как релевантные страницы (контекст) после процедуры RAG. Перед финальным ответом внимательно, вслух и по шагам проанализируй вопрос. (и т.д.)

```
step_by_step_analysis: str = Field(description="Подобранный пошаговый анализ ответа минимум из 5 шагов и не менее 150 слов.")
str = Field(description="Краткое резюме пошагового анализа. Около 50 слов.") (и т.д.)
```

Пример 1: Какова прибыль Сбера за 2023 год?

```
{"step_by_step_analysis": «...», "reasoning_summary": «...», "relevant_pages": ..., "final_answer": ...}
```

Контекст: (найденный контекст, т.е. отобранные страницы)  
Вопрос: (вопрос)

# Generation

Специфическое для условий данного соревнования решение:

- каждый запрос (кроме multiquery) касался только одного отчета;
- в запросе было явно указано, к отчету какой компании он относится;
- запросы относились к одному из нескольких известных типов.

## 1. Роутинг к нужной БД по имени компании



## 2. Роутинг к нужному шаблону промпта

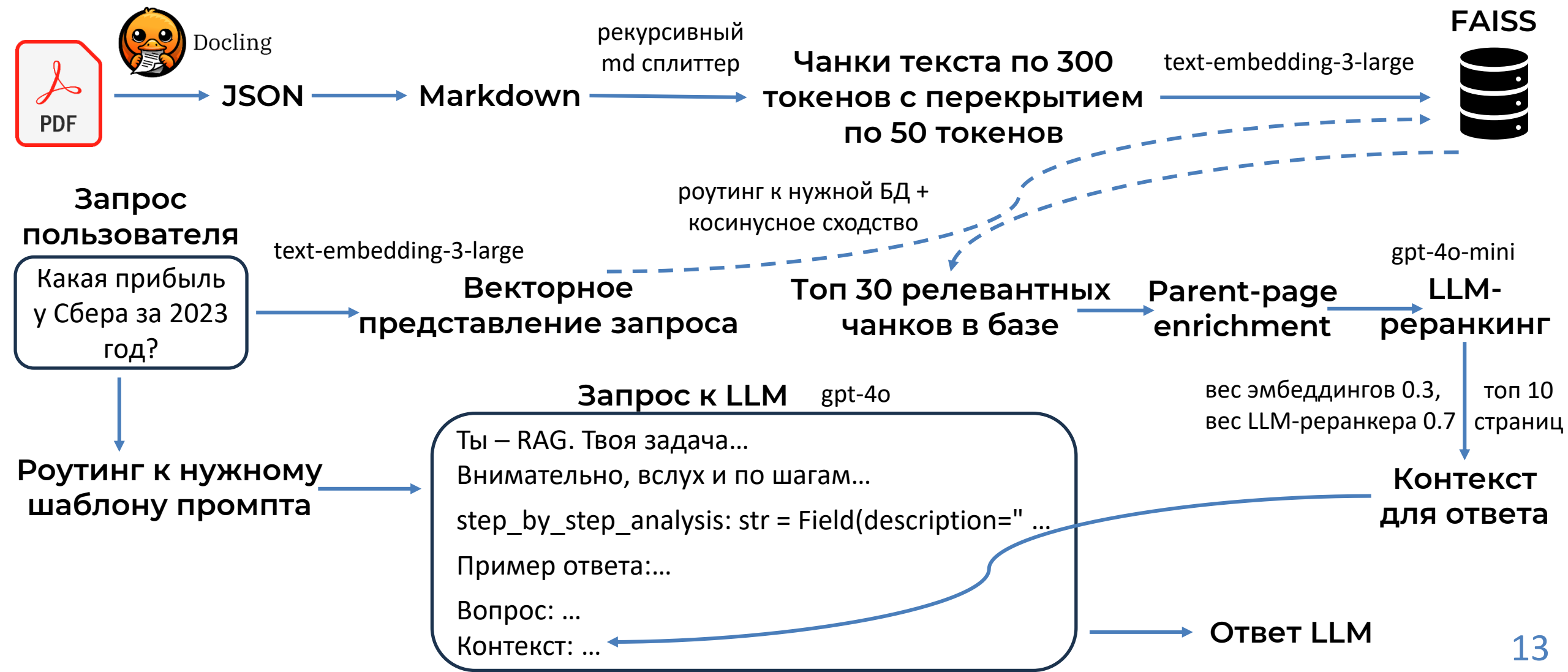


## 3. Роутинг multiquery запросов





# Еще раз соберем всё решение вместе





# Demo

Задайте вопрос

Ваш вопрос

Каковы основные финансовые показатели Сбера за 2023 год?

Поиск

Ответ системы

- Итого обязательств составило 43 723 млрд руб., увеличившись на 26,7% (страница 94).
- Итого собственных средств составило 6 584 млрд руб., увеличившись на 14% (страница 94).

4. **Качество активов:**

- Доля неработающих кредитов в кредитном портфеле снизилась до 2,2% (страница 95).

**РАССУЖДЕНИЯ:**

На основе предоставленных данных можно сделать выводы о значительном улучшении финансовых показателей Сбера в 2023 году по сравнению с предыдущим годом. Это выражается в росте рентабельности, увеличении чистой прибыли и операционных доходов, а также в снижении доли неработающих кредитов. Также наблюдается рост активов и обязательств, что свидетельствует о расширении деятельности банка.

**ВЫВОДЫ:**

Основные финансовые показатели Сбера за 2023 год демонстрируют значительное улучшение:

- Рентабельность капитала (ROE) составила 25,3%.
- Чистая прибыль увеличилась в 5,2 раза до 1 508,6 млрд руб.
- Операционные доходы выросли на 121,7% до 2 909,1 млрд руб.
- Активы увеличились на 25% до 52 307 млрд руб.
- Доля неработающих кредитов снизилась до 2,2%.

Источники и статистика

**Источники:**

**1.** Страница 95 (поиск: 0.725, релевантность: 10.0/10)  
**Превью:** Финансовые показатели  
Показатели рентабельности  
%  
2023  
2022  
Изменение  
2021  
Рентабельность среднегодовых активов (ROA)  
3,2  
0,7  
2,5 п.п.  
3,3  
Рентабельность капитала (ROE)  
25,3  
5,2

# Выводы

- RAG – одно из основных применений LLM на сегодняшний день.
- В работе была реализована мультимодальная RAG-система, вдохновленная победным решением Enterprise RAG Challenge (2025 год).
- Основные используемые технологии: pdfplumber, RecursiveCharacterTextSplitter, FAISS, Parent-page enrichment, LLM reranking, Chain-of-Thoughts, pydantic, pymupdf, gradio, text-embedding-3-large, GPT-4o, gpt-4o-mini.
- Основные сложности: некачественное выделение и распознавание картинок (следует рассмотреть переход с pdfplumber на docling, а также уточнить промпт для GPT-4o и добавить автоматический повторный запрос в случае отказа создать описание картинки).
- Перспективы на будущее: целесообразно верифицировать качество собранного пайплайна путем построения автоматической оценки (RAGAS и др.).



## Список используемых источников/программных средств:

- ✔ Google Colab, HuggingFace
- ✔ «Как я победил в RAG Challenge: от нуля до SoTA за один конкурс» (<https://habr.com/ru/articles/893356/>)