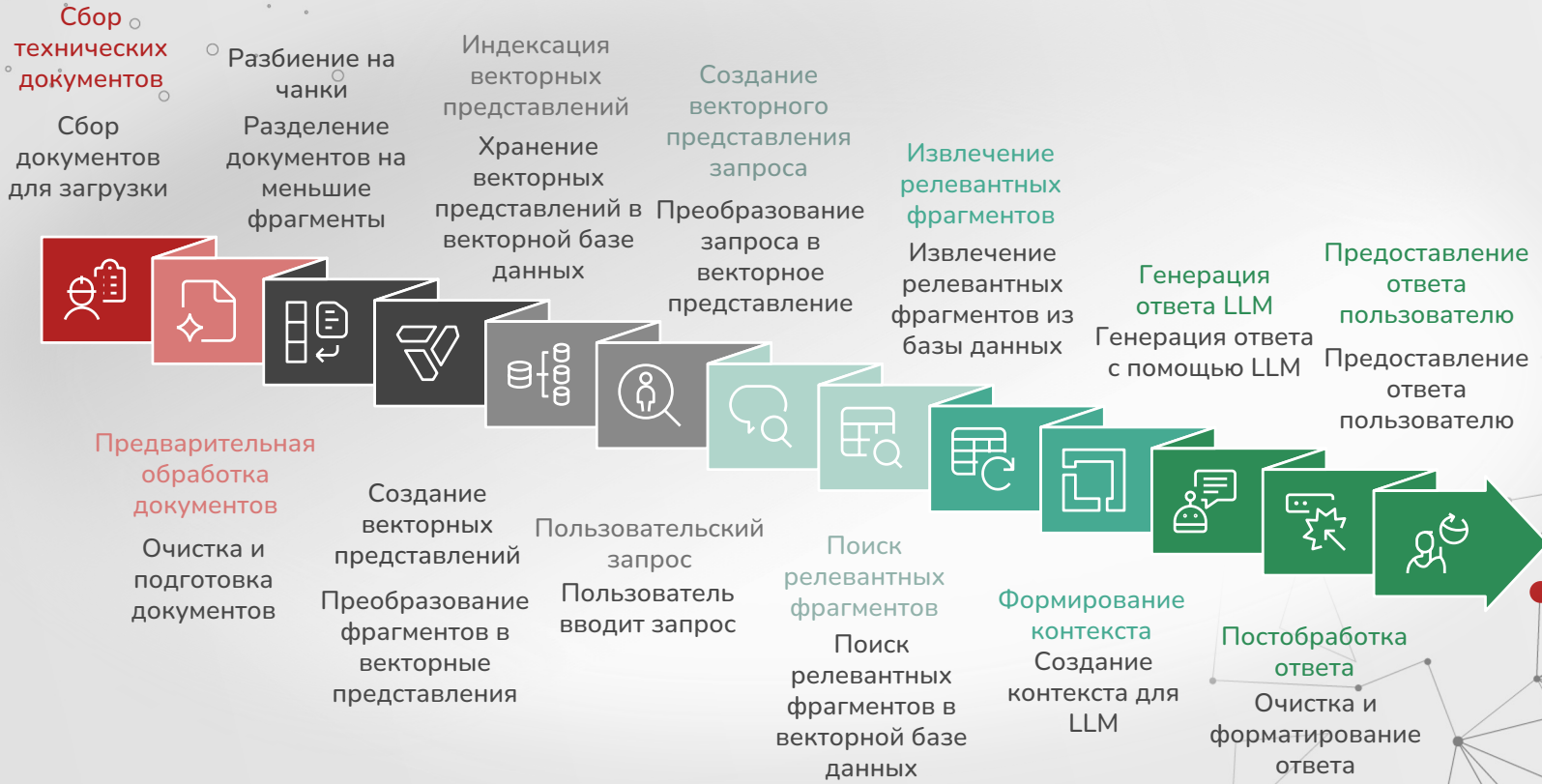


ETL 4 AI

Отраслевые кейсы применения ИИ для
обработки данных

Типовой RAG пайплайн





Аудио и видео

Неструктурированные аудио и видеофайлы



Текст

Макеты страниц, таблицы, формулы



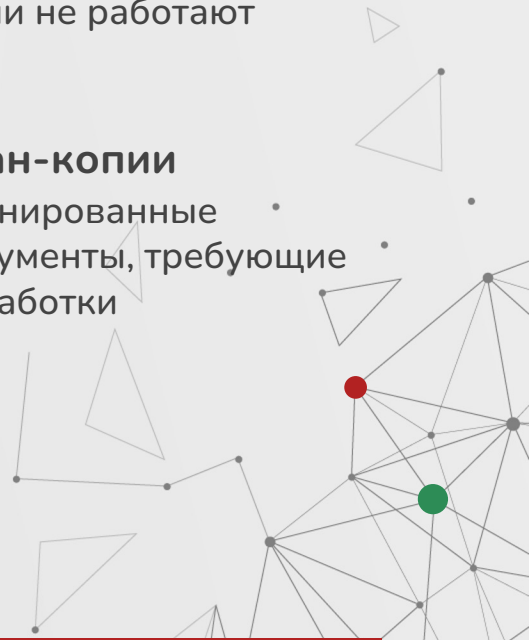
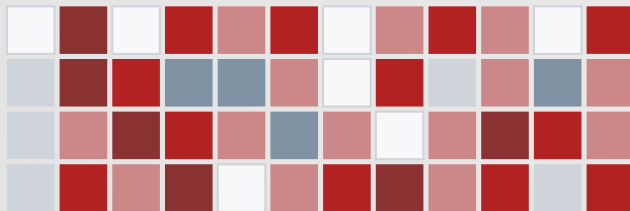
Несъедобные большие данные

В архивах компаний хранится огромное количество данных но они не работают



Скан-копии

Сканированные документы, требующие обработки

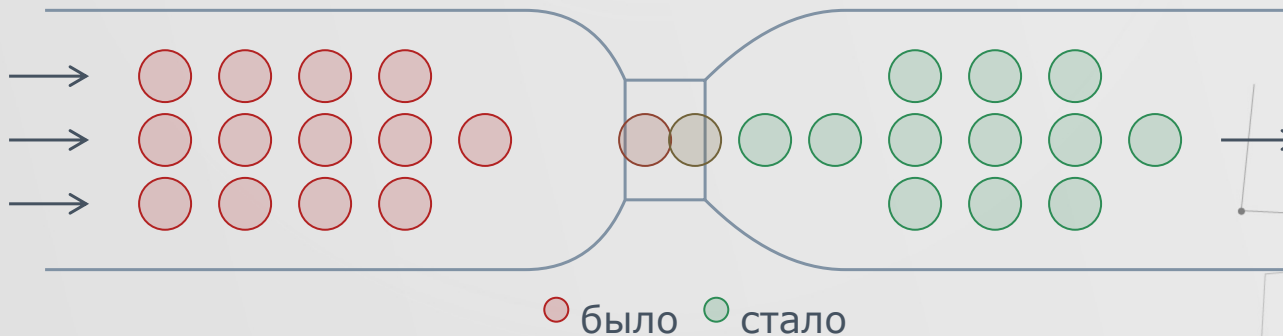


Хотелось бы

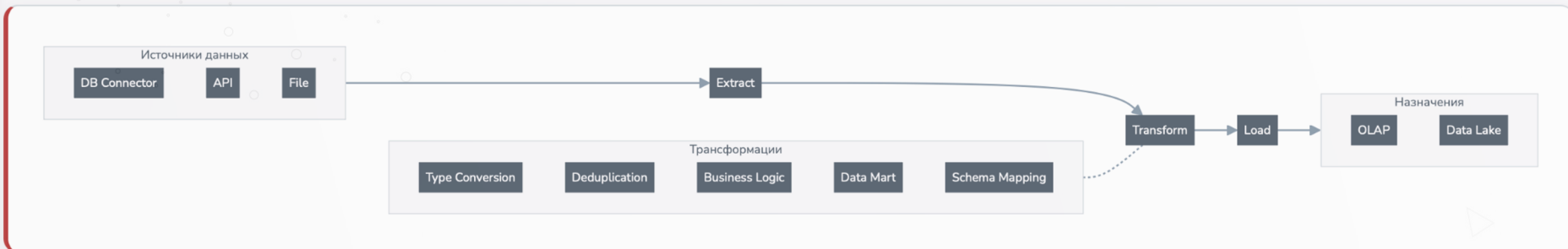
- Сократить time-to-production
- Уменьшить стоимость подготовки данных
- Снизить зависимость от ручной работы

Ключевой фактор

- Качество результата сильно зависит от промптов
- Пользователи должны сами настраивать поведение системы

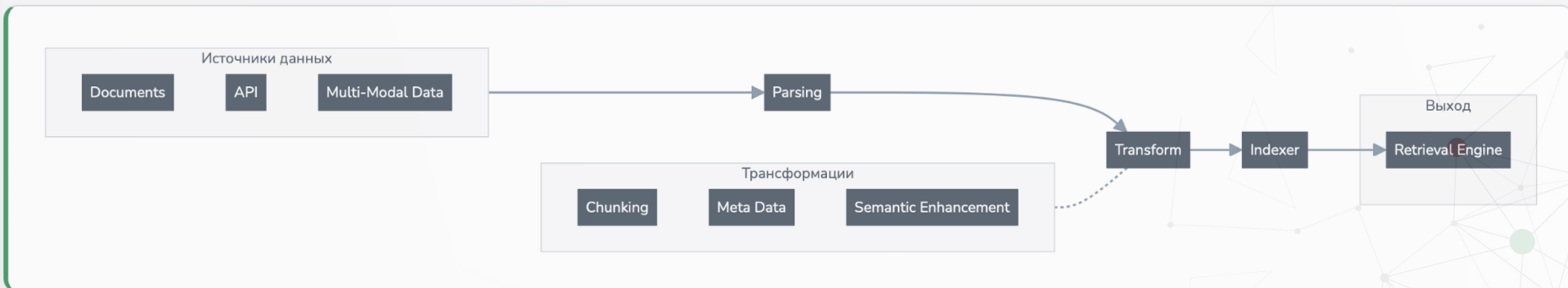


ETL PIPELINE



VS

RAG PIPELINE



Многие проекты используют одни и те же кирпичи **СНОВА И СНОВА**

Парсинг

PDF · картинки · DOC · PPT · таблицы · прочее

Чанкинг

фиксированный · семантический · рекурсивный

LLM-трансформации

сводка · извлечение · классификация · прочее

Мэтчинг

эмбединги · косинусное сходство

Индексация (для RAG сценария)

гибридный поиск: BM25 + vector

Процессинг

Многопоточность · мониторинг · отладка

ВЫВОД

Нужен конфигурируемый пайплайн, а не набор скриптов

Что должна уметь система

Универсальный парсер:

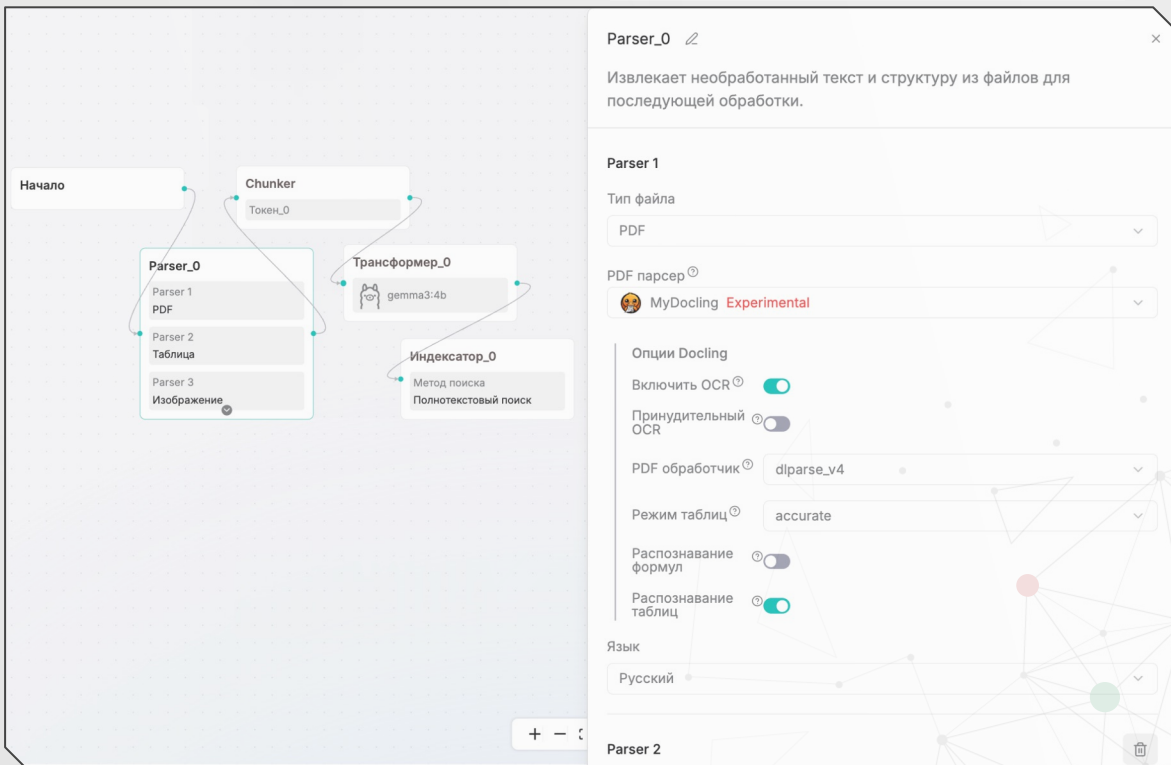
- Word, Excel, PDF
- изображения (OCR, классификация)
- аудио / видео
- сложные документы (таблицы, формулы)

Масштабируемая обработка:

- batch-режим
- многопоточность
- быстрое масштабирование

Гибкая настройка:

- пайплайны
- промпты
- сценарии обработки








Взяли основу — и сделали production-ready

Базовый проект

 ragflow (Apache License 2.0, ★ 75k)

Что добавили

-  Улучшили обработку русскоязычного текста
-  Добавили парсеры и интеграции (Yandex GPT, Yandex Search, etc)
-  Разработали дополнительные узлы обработки (BI-агент, мэтчинг)
-  Проработали библиотеку кейсов использования
-  Подготовили дорожную карту дальнейших улучшений

Результаты

 RAGu

Получилась платформа для быстрой сборки AI-решений

Бизнес доволен

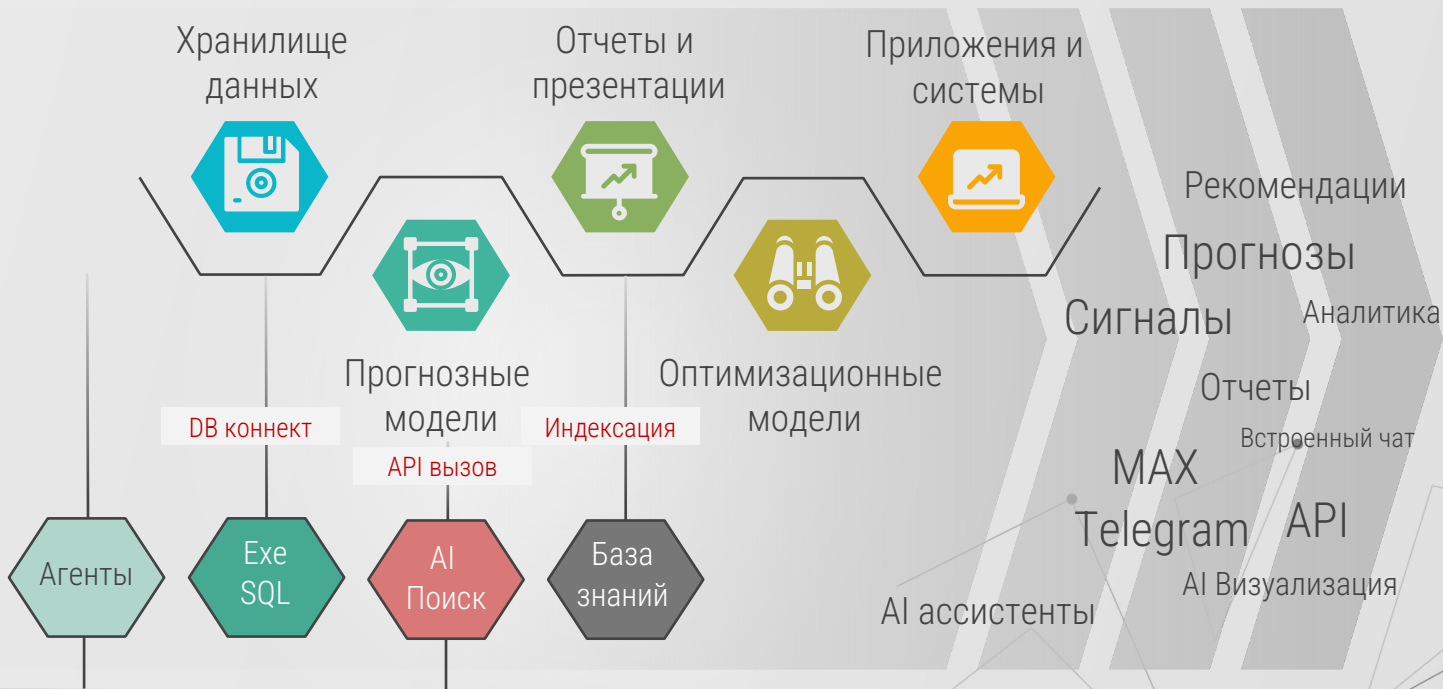
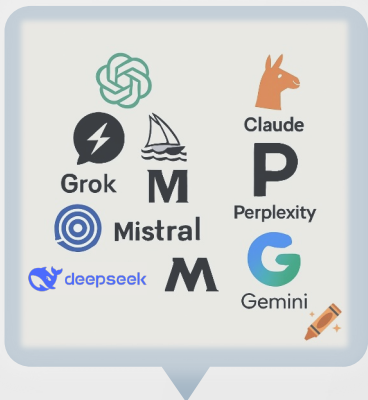


Интеграция



Компании и пользователи

Экосистема провайдеров моделей



использует API приложений и СУБД как инструменты взаимодействия с ресурсами заказчика для генерации достоверных и обоснованных ответов AI ассистента

Обработка технических паспортов для энергетики

Потенциалы экономии

1. Ручное извлечение тех. Параметров из 12 тыс. документов
2. Целевые сведения хранятся в таблицах, формулах, рукописном тексте

Решение

1. Подбор комбинации OCR для получения текстовых сведений там где даже человеку сложно разобрать
2. Извлечение тех. параметров LLM
3. Мэтчинг со справочником тех. Параметров с помощью гибридного поиска

Результаты

1. Сокращение ручного труда на **95%**
2. Единообразный результат, готовый к загрузке в систему управления производственными фондами



Обработка таблиц

Common food_SP_A5.pdf
Size: 49 KB | Updated Time: 18/02/2025 20:34:48

Результат чанка
Просмотр чанков, используемых для эмбединга и извлечения.

Полный текст | Сокращения

Выбрать все

единица измерения	0,00	0,02	0,03	0,05
калории (кcal/сут)	15			
белки (г/сут)	1,4	2,0	3,0	4,0

Суточная потребность в пищевых веществах и энергии для обучающихся в образовательных организациях кадетского типа и организаций кадетской направленности

Возраст и класс	Энергетическая ценность (ккал/сут)	Белки (г/сут)	Жиры (г/сут)	Углевод. (г/сут)
5-8 класс	до 3500	119-143	134-143	550-580
9-11 класс	до 4000	142-177	167-188	648-681

Всего 8 | 50 / Страниц

Обработка формул

Landau-Teofiz_t1_mehnika_12.pdf
Size: 107 KB | Updated Time: 18/02/2025 19:01:11

Результат чанка
Просмотр чанков, используемых для эмбединга и извлечения.

Полный текст | Сокращения

Выбрать все

УРАВНЕНИЯ ДВИЖЕНИЯ

Исходный текст

Но в силу условий (2.3) первый член в этом выражении исчезает. Остаётся интеграл, который должен быть равен нулю при произвольных значениях dt . Это возможно только в том случае...

Исходный текст

$$M \frac{d^2}{dt^2} (x + dx) + M \frac{d^2}{dt^2} (y + dy) + M \frac{d^2}{dt^2} (z + dz) - M \frac{d^2}{dt^2} (x + dx) = 0.$$

Исходный текст

В заключение отметим, что решение уравнения (2.5) не является дифференциальным уравнением второго порядка, так как оно содержит производные второго и третьего порядков.

Всего 13 | 50 / Страниц



СПАСИБО

Вопросы?

www.resetlab.ru
d.kibalnikov@resetlab.ru
+79119714950



@KB_DM

CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, and infographics & images by **Freepik**.

Please keep this slide for attribution.