

Введение в BigData

11.02.2025



Содержание

- 1** Немного истории
- 2** Архитектура решений
- 3** Наши дни
- 4** Игроки на РФ рынке
- 5** Немного практики

01

Немного истории



Немного истории

1. **70-ые** – Эдгар Кодд появления реляционной модели данных;
2. **80-ые** – Пол Мерфи, Барри Девлин “Business DWH”;
3. **90-ые** – Появления Oracle, Teradata, IBM D2. Модели Билла Инмона и Ральфа Кимбалла;
4. **00-ые** – MPP субд, появление Greenplum, Hadoop; Термин BigData(Клиффорд Линч)
5. **10-ые** – Эпоха облаков (Snowflake, G BigQuery, Amazon Redshift)
6. **20-ые** – Trino, YugaByte и новые.

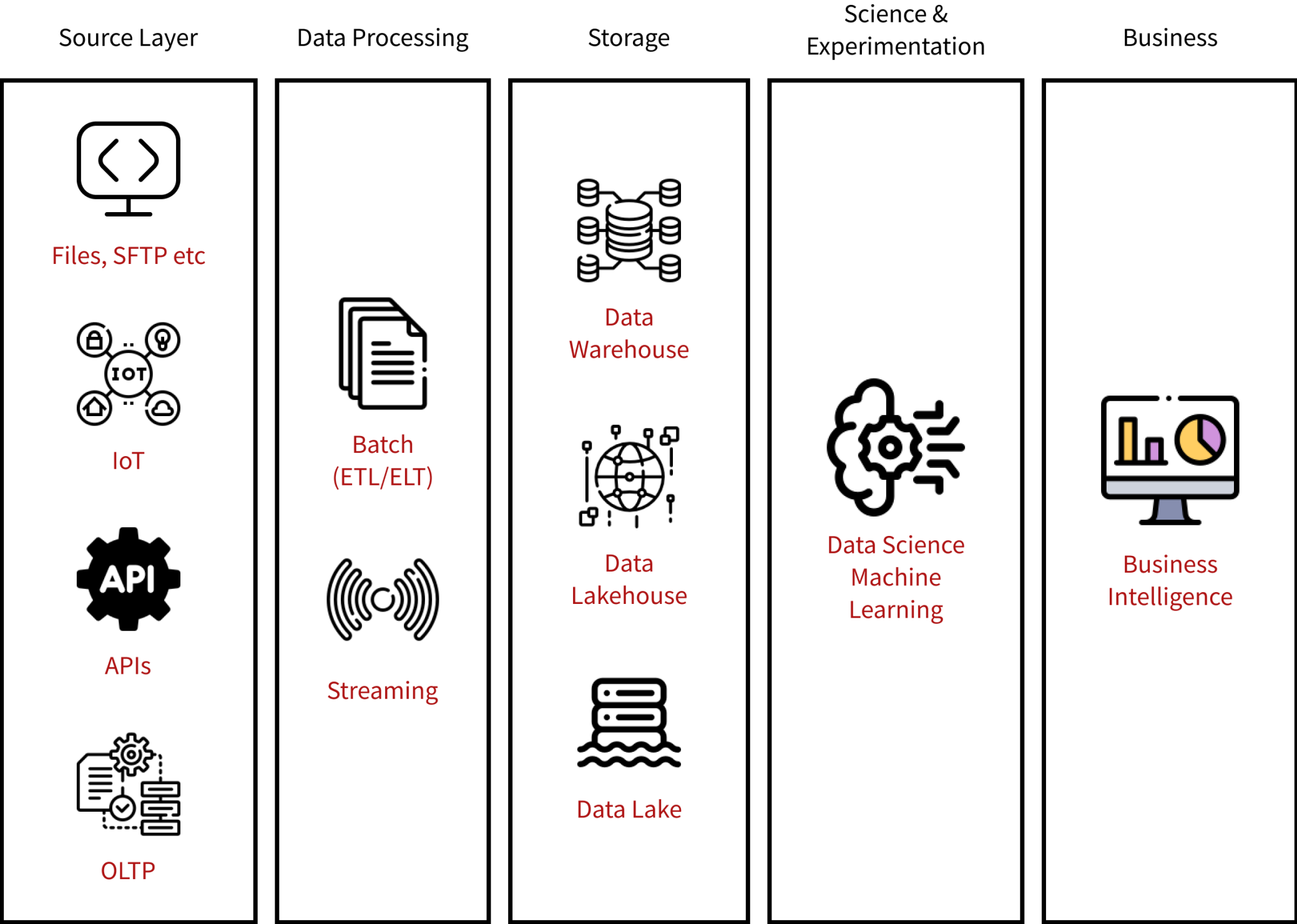
02

Архитектура решений



Архитектура решений

Key Layers in Data Stack



Source Layer

1. **Files, SFTP, etc** – файловые ресурсы(Excel, csv, FTP|SFTP);
2. **IOT** - internet of things, устройства умного дома, медицинское оборудование и т.д;
3. **API** - Application programming interface, программный интерфейс для получения данных;
4. **OLTP** – Online Transaction Processing. Системы обработки транзакций в реальном времени.

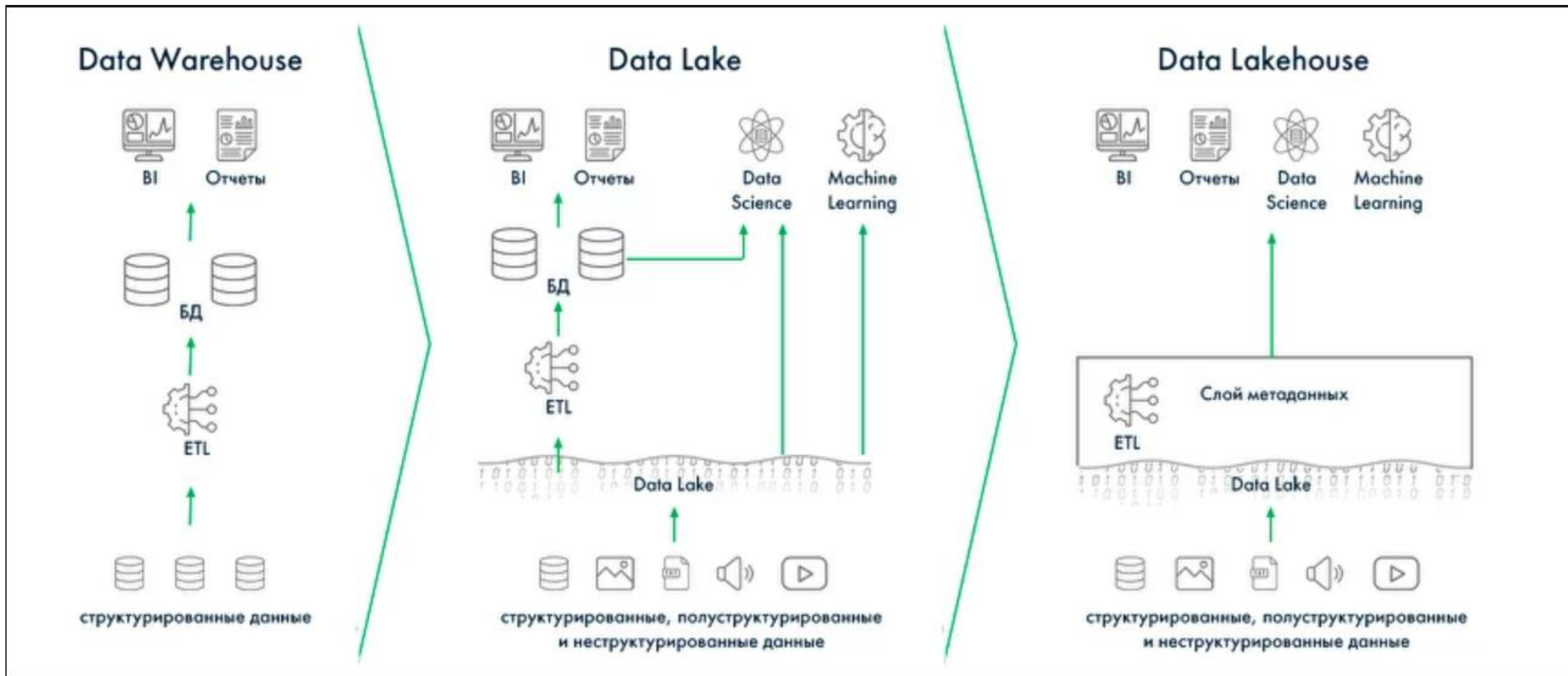
Data processing

1. **Batch** – ELT, ETL(Extract, Load, Transform и Extract, Transform, Load). Системы передачи **порционных** данных сырых данных с источника.
2. **Streaming** – обработка данных в реальном времени.

Storage

1. **DWH** – Хранилище данных. Это система для хранения и управления большими объемами структурированных данных, которые объединяются из различных источников с целью получения из них бизнес выгоды (G BigQuerry, Amazon Redshift, Snowflake, Greenplum).
2. **Data Lake** – Озеро данных. Это инструмент хранения необработанных(неструктурированных) данных в различных форматах (Amazon S3, Minio, G Cloud Storage, Yandex S3).
3. **Data Lakehouse** – Объединение концепции DWH и Data Lake. Храним все виды данных и управляем метаданными.

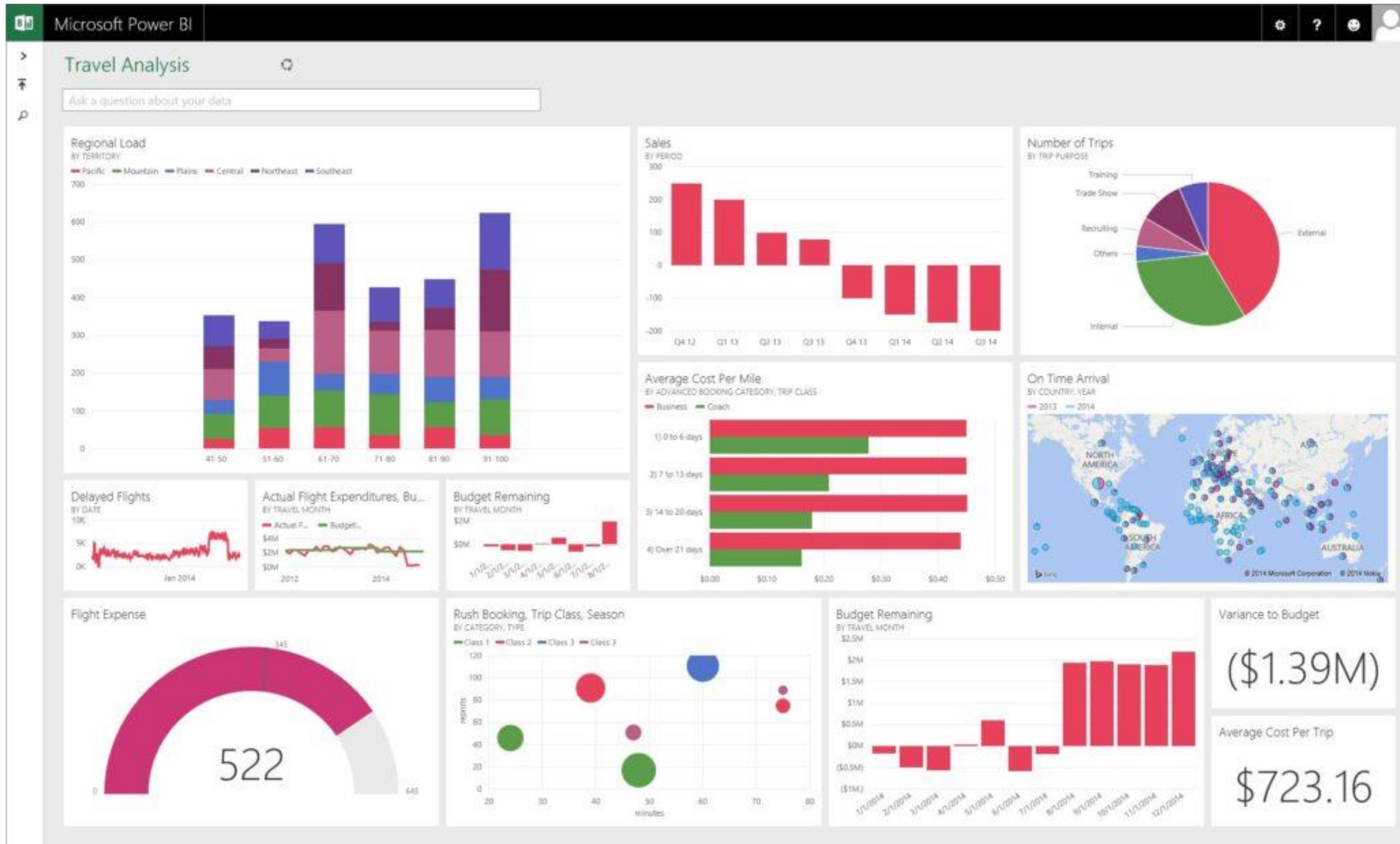
Storage



Science&Experimentation

1. **Adhoc** – запросы к источнику данных для получения метрик, подтверждения гипотез;
2. Скрипты, аналитика и т.д.
3. **ML** модели.

BI

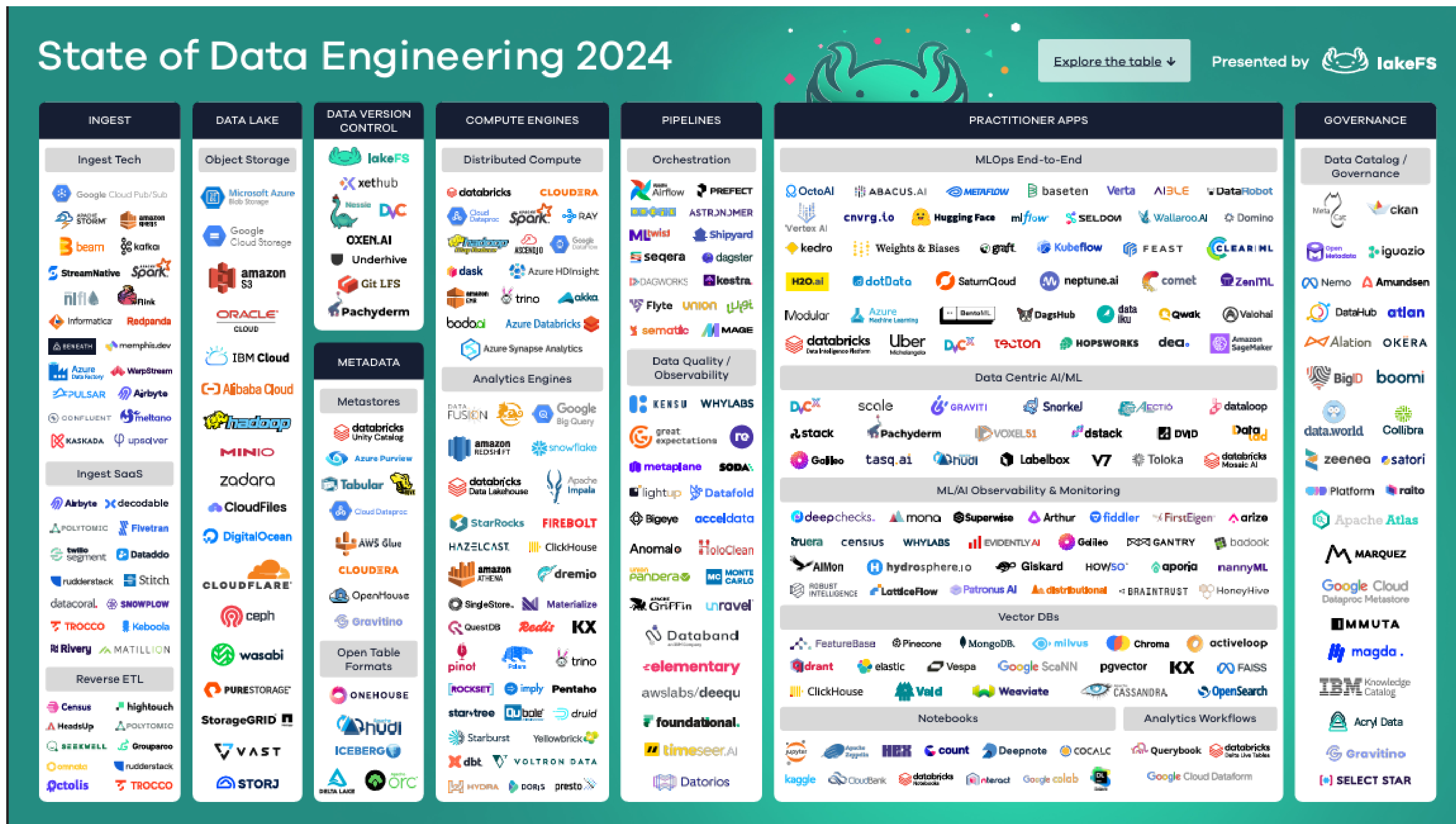


03

Наши дни



Наши дни



04

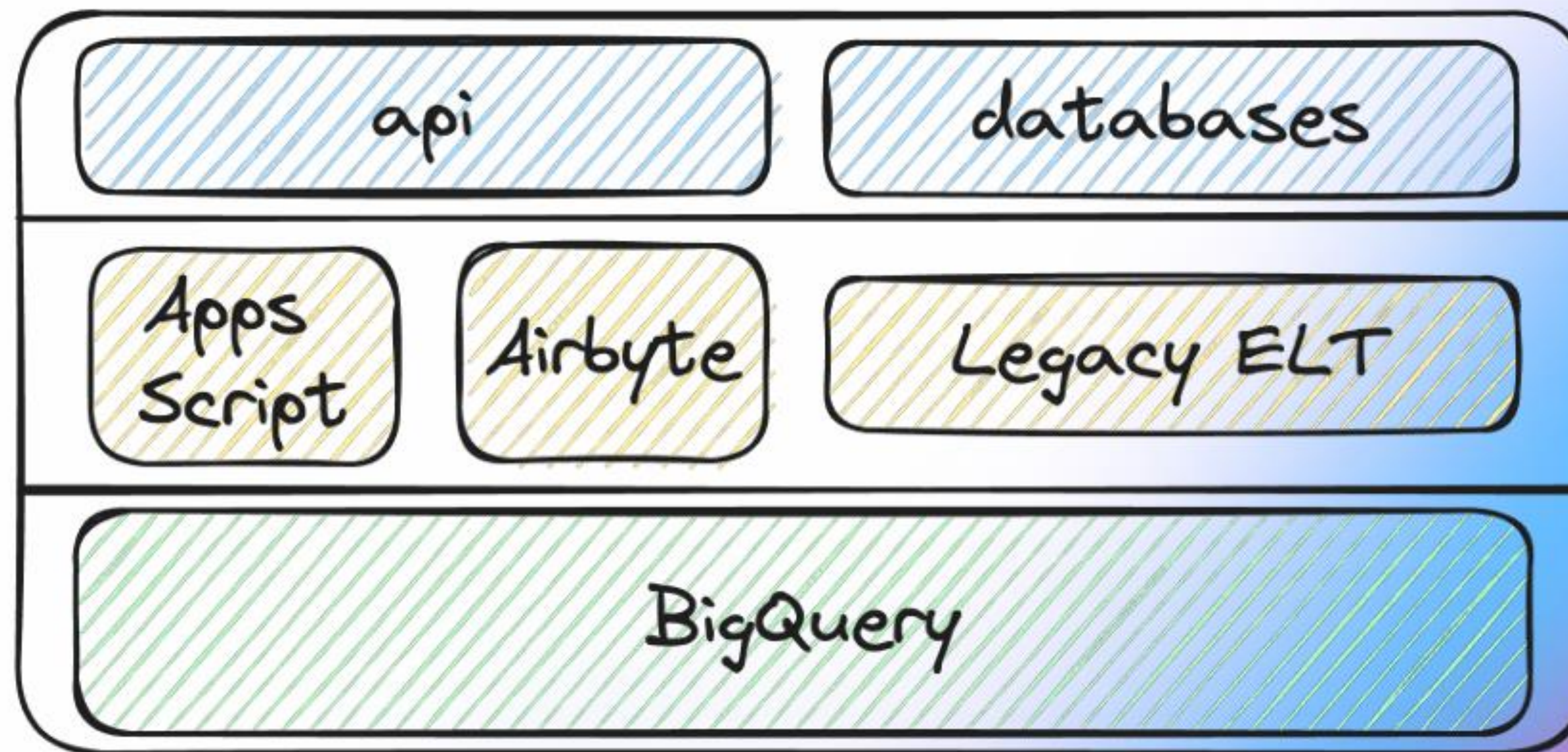
Игроки на рынке



Tbank(~2018)



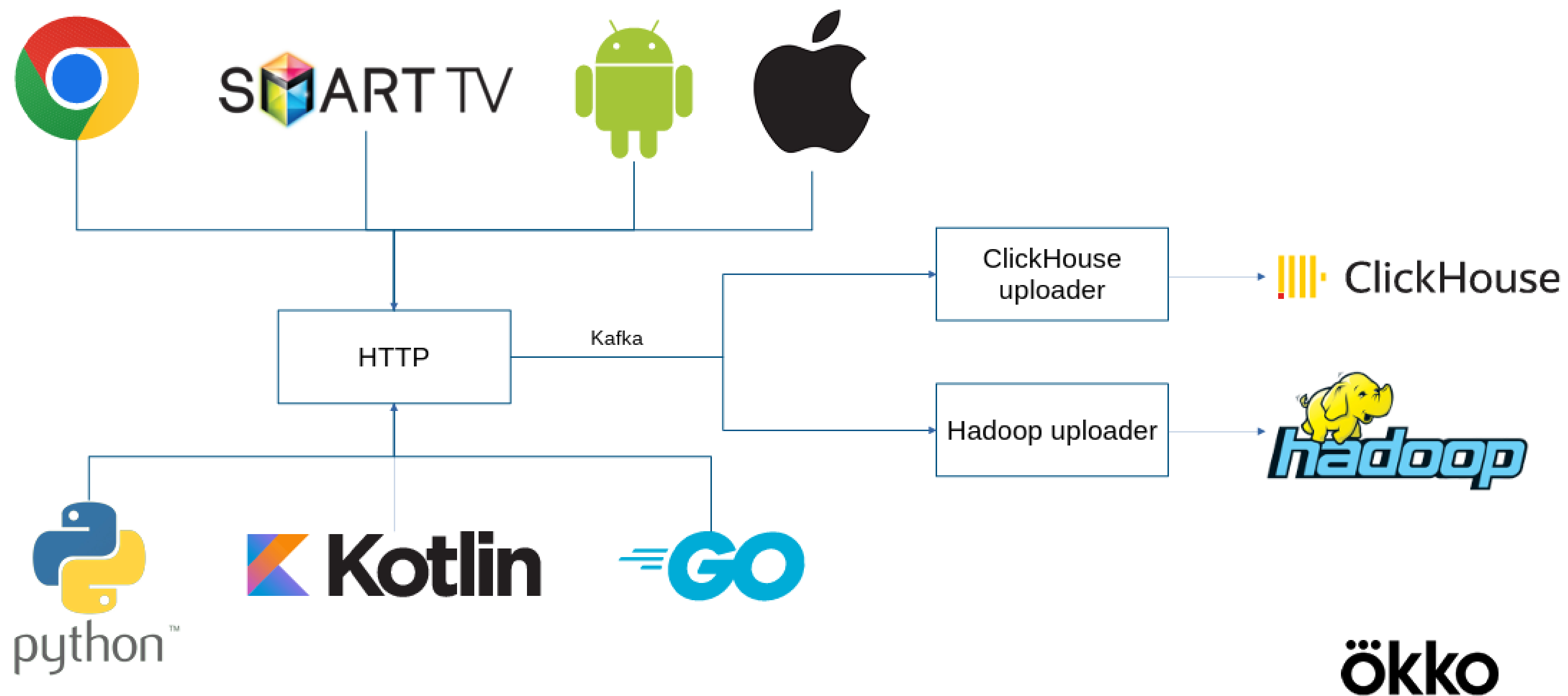
Сберздоровье(~2022)



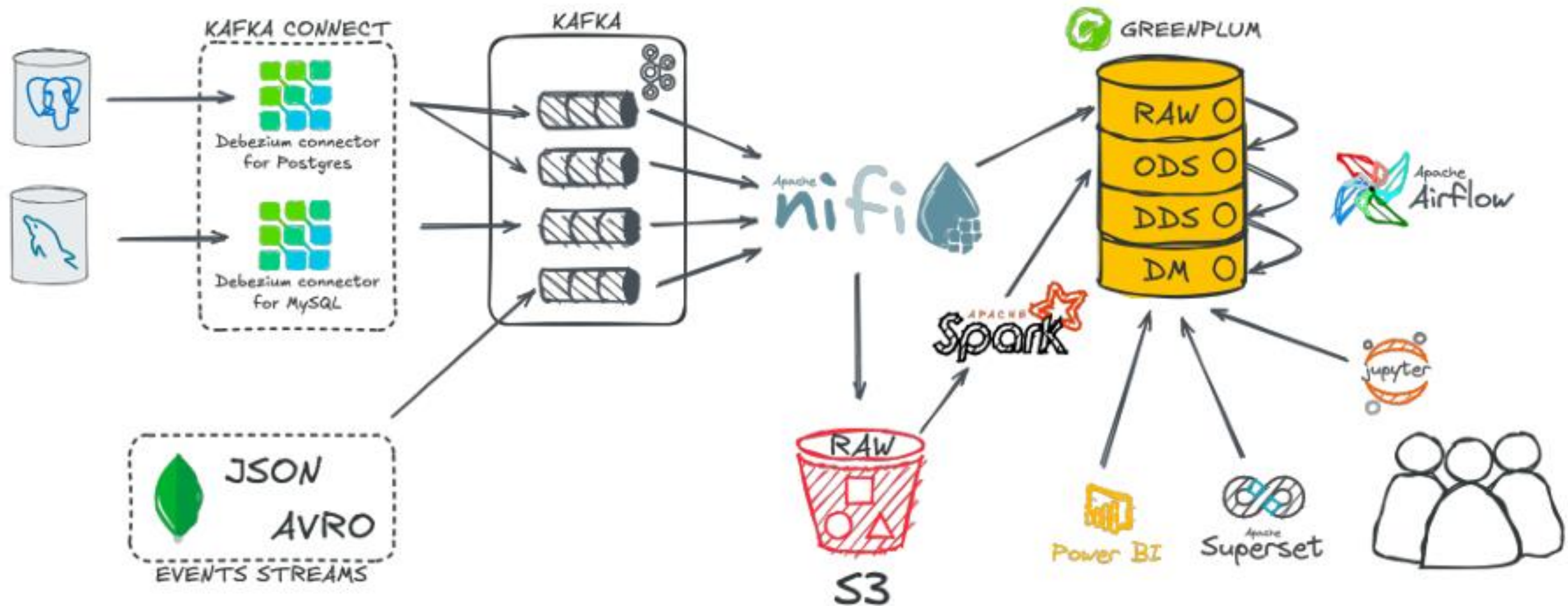
Okko(~2023)

Архитектура

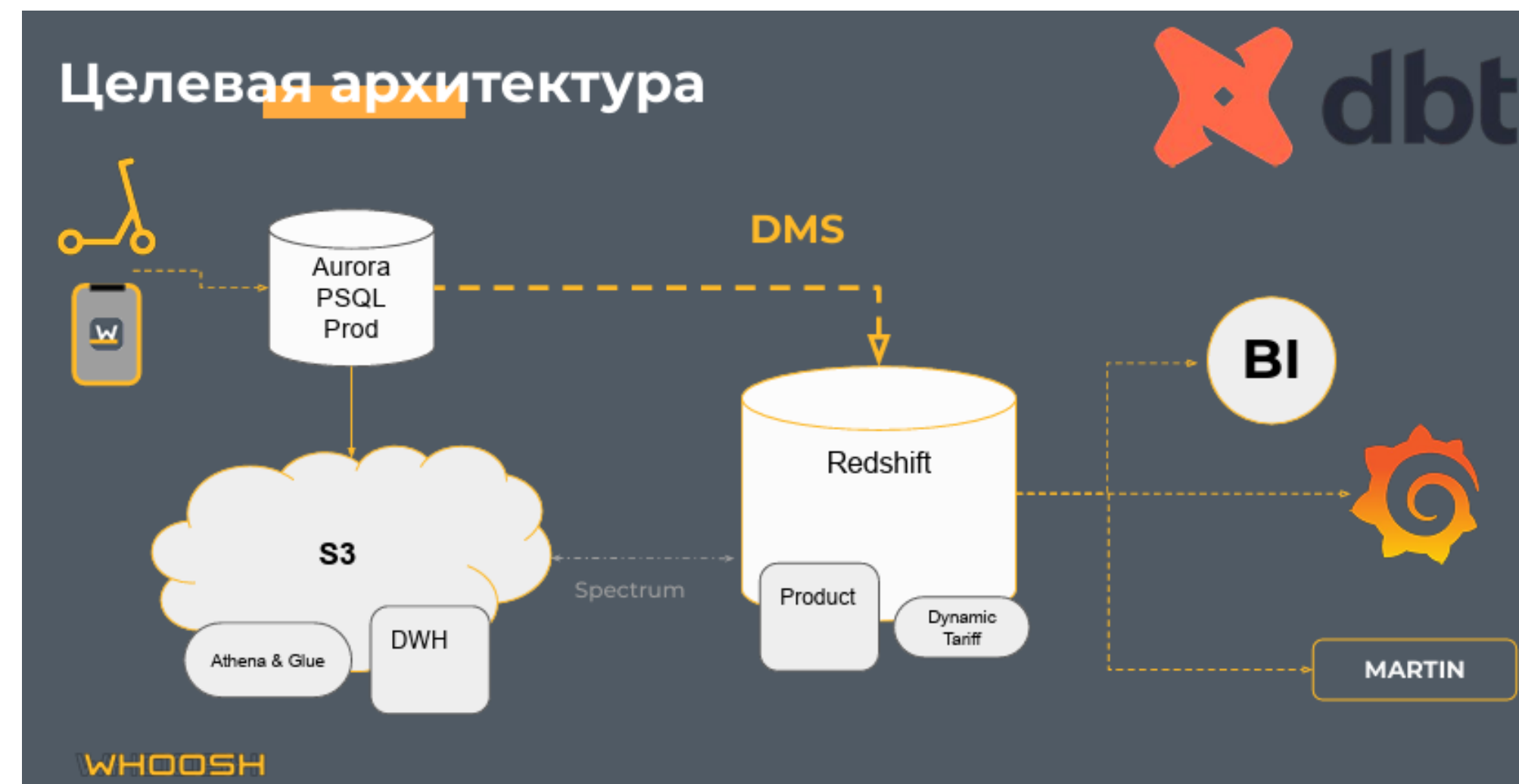
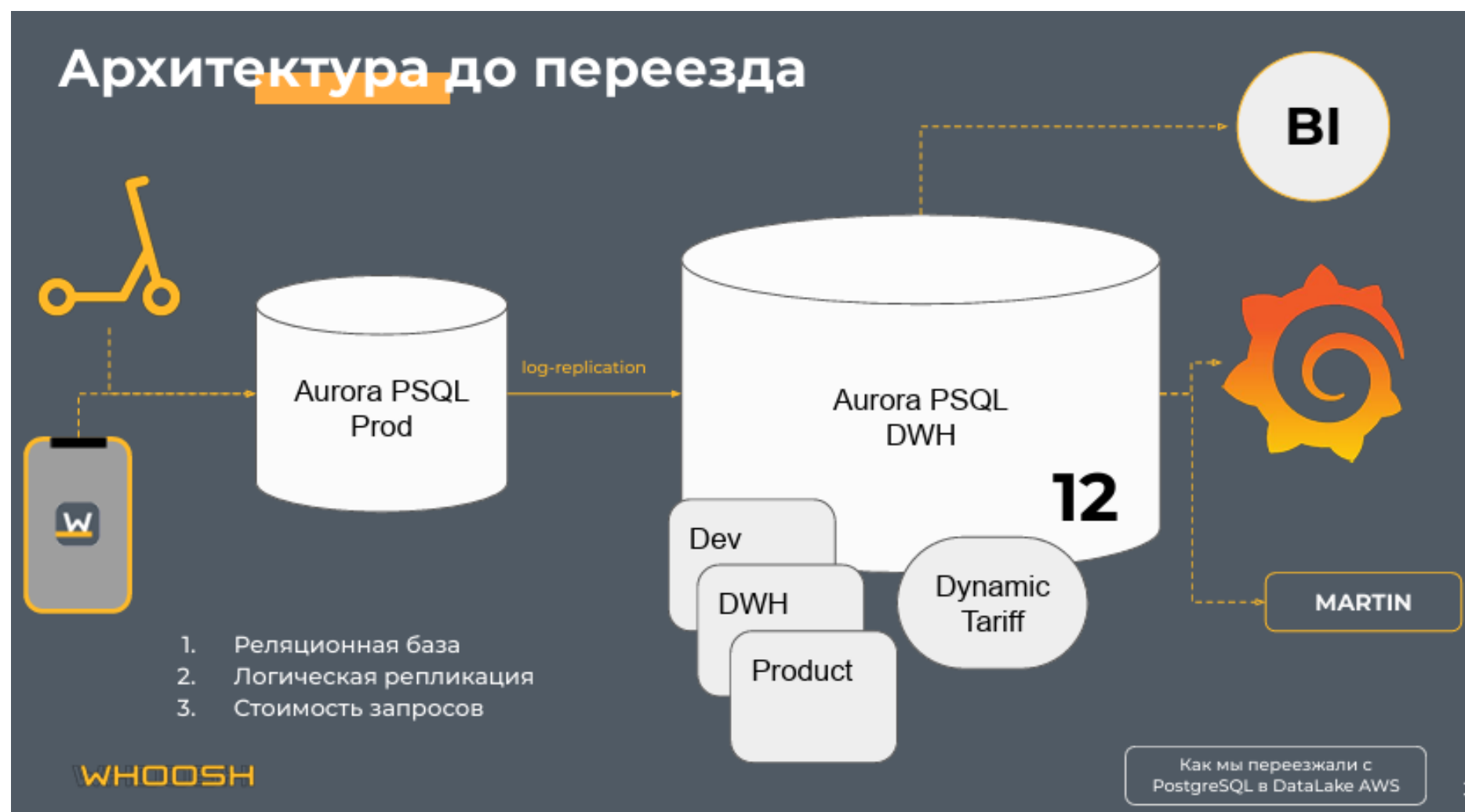
5



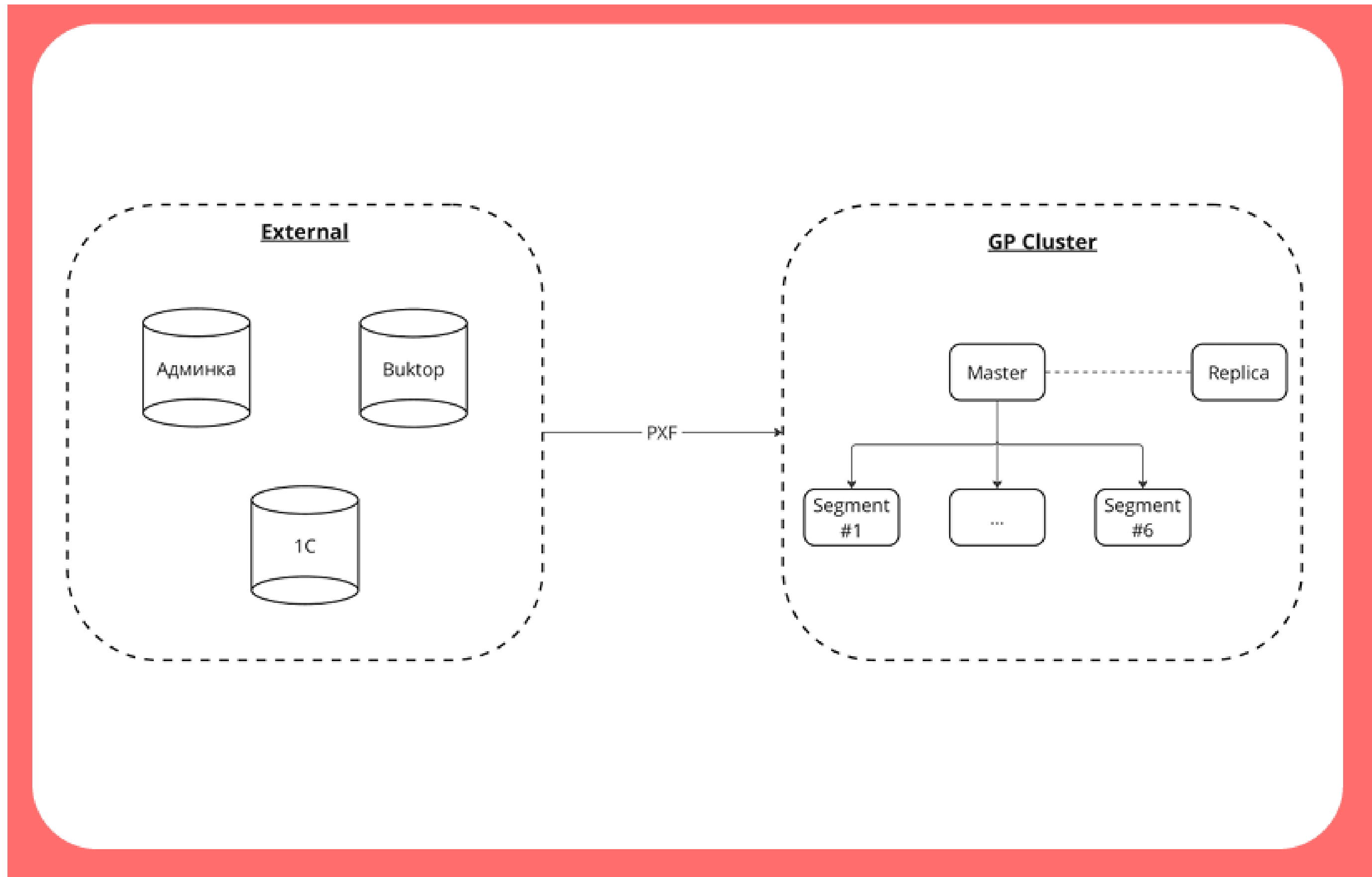
LemanaPro(~2022)



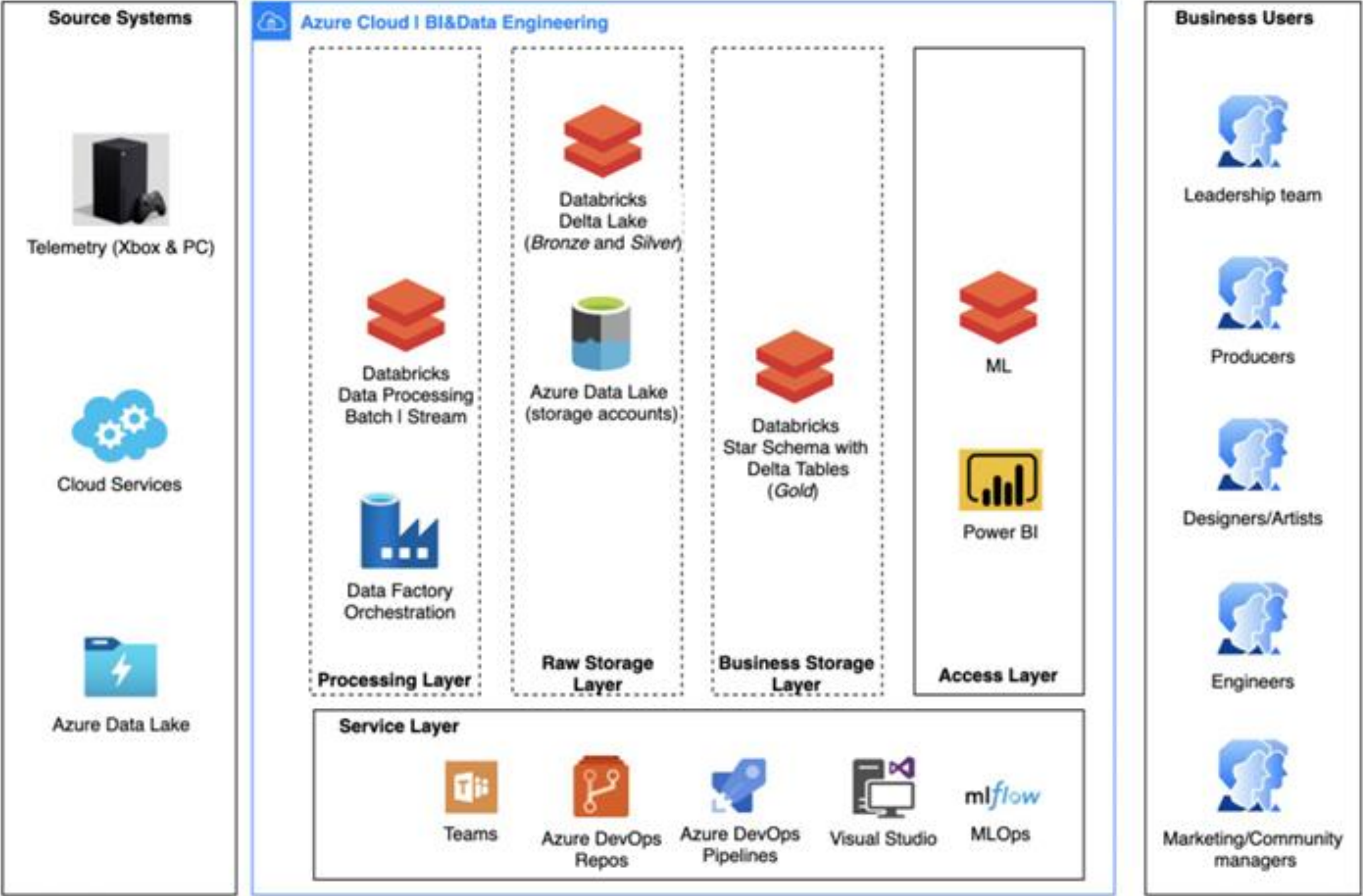
Whoosh(~2023)



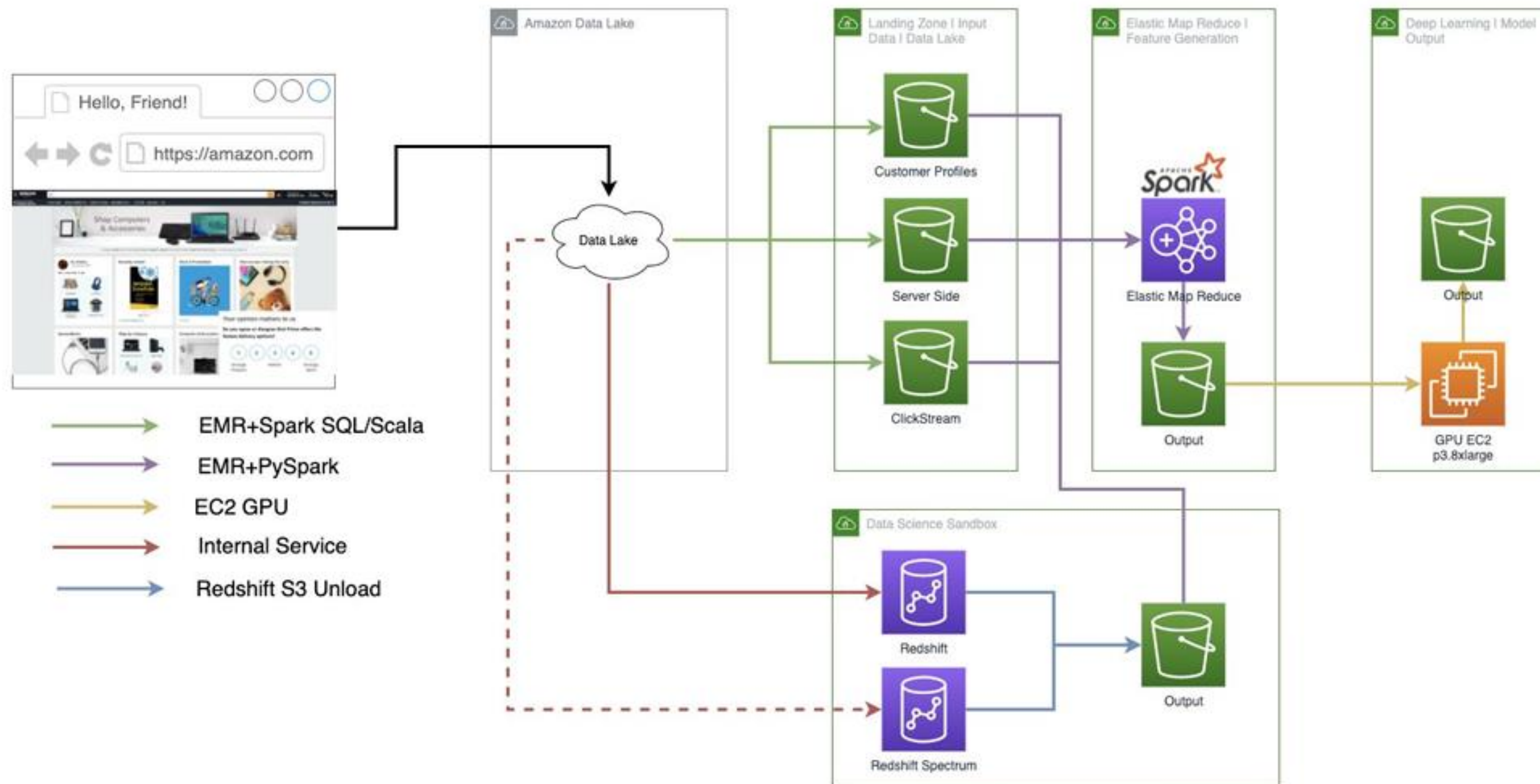
BestDoctor(~2022)



Xbox Microsoft(~2022)

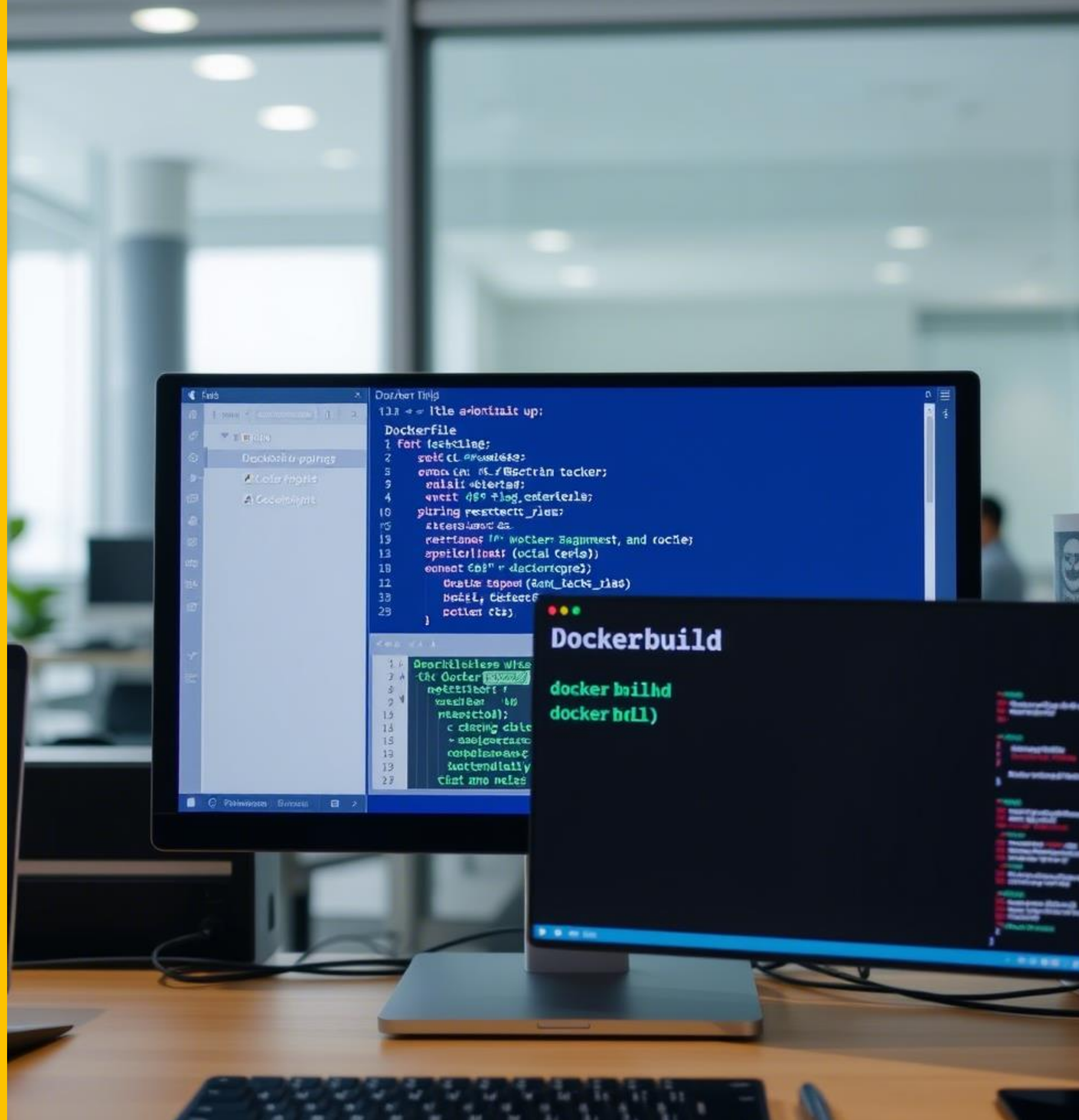


Amazon Retail(~2022)



05

Немного практики



Настройка

1. <https://www.docker.com/products/docker-desktop/>
2. <https://github.com/zezOtik/bmstu--iu8-big-data-tools>
3. <https://dbeaver.io/>