

Аналитика Больших данных: технологии

С.Д. Белов, ЛИТ ОИЯИ
belov@jinr.ru

Осенняя Школа по информационным технологиям ОИЯИ 2022
14-19 ноября 2022г., ЛИТ ОИЯИ, Дубна

Современные тенденции в анализе и управлении данными

Gartner Top Data and Analytics Trends for 2022



Activate Dynamism and Diversity

- Adaptive AI Systems
- Data-Centric AI
- Metadata-Driven Data Fabric
- Always Share Data



Augment People and Decisions

- Context-Enriched Analysis
- Business-Composed D&A
- Decision-Centric D&A
- Skills and Literacy Shortfall



Institutionalize Trust

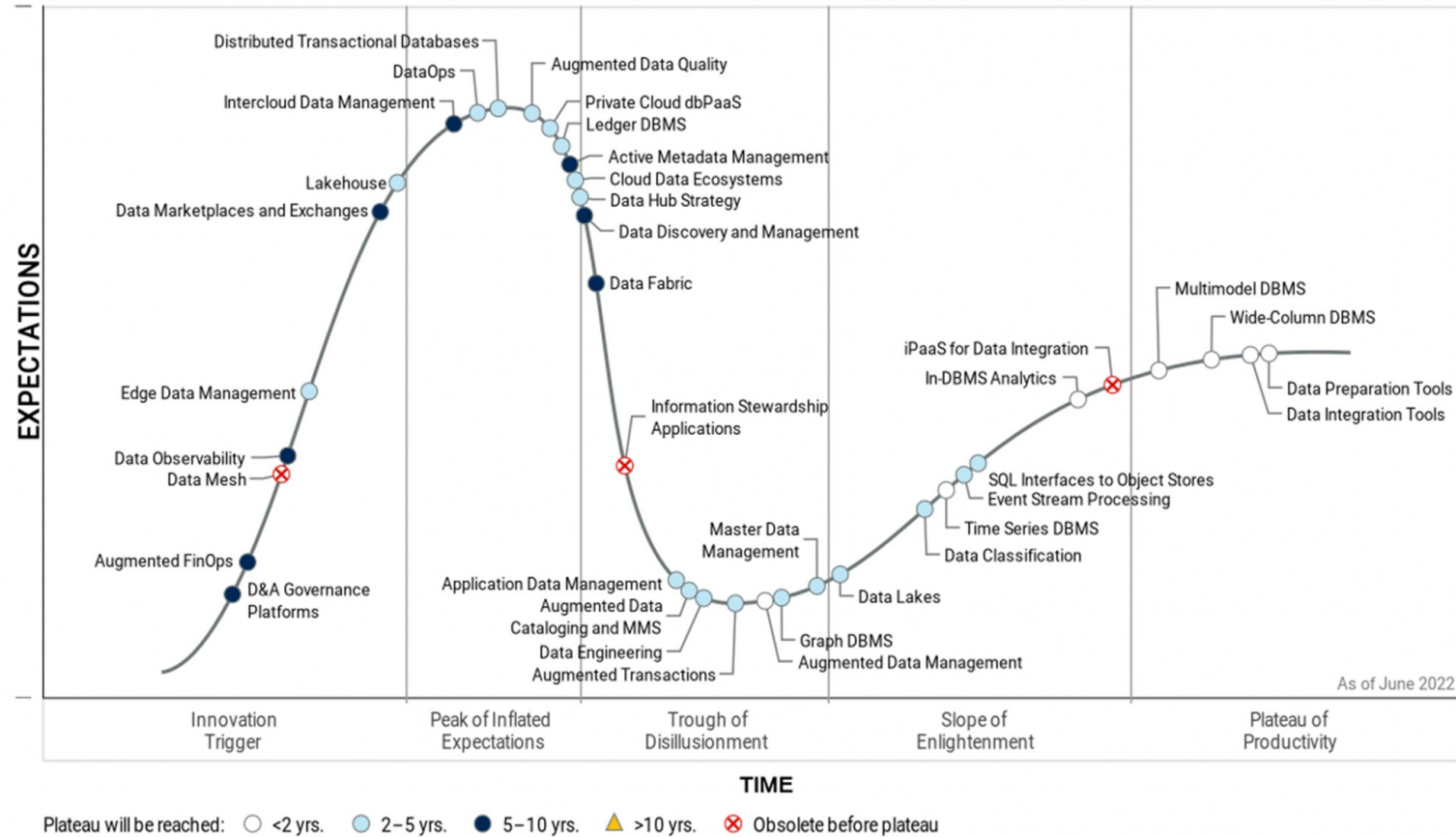
- Connected Governance
- AI Risk Management
- Vendor and Regional Ecosystems
- Expansion to the Edge

gartner.com

Source: Gartner
© 2022 Gartner, Inc. All rights reserved. CM_GTS_1740785



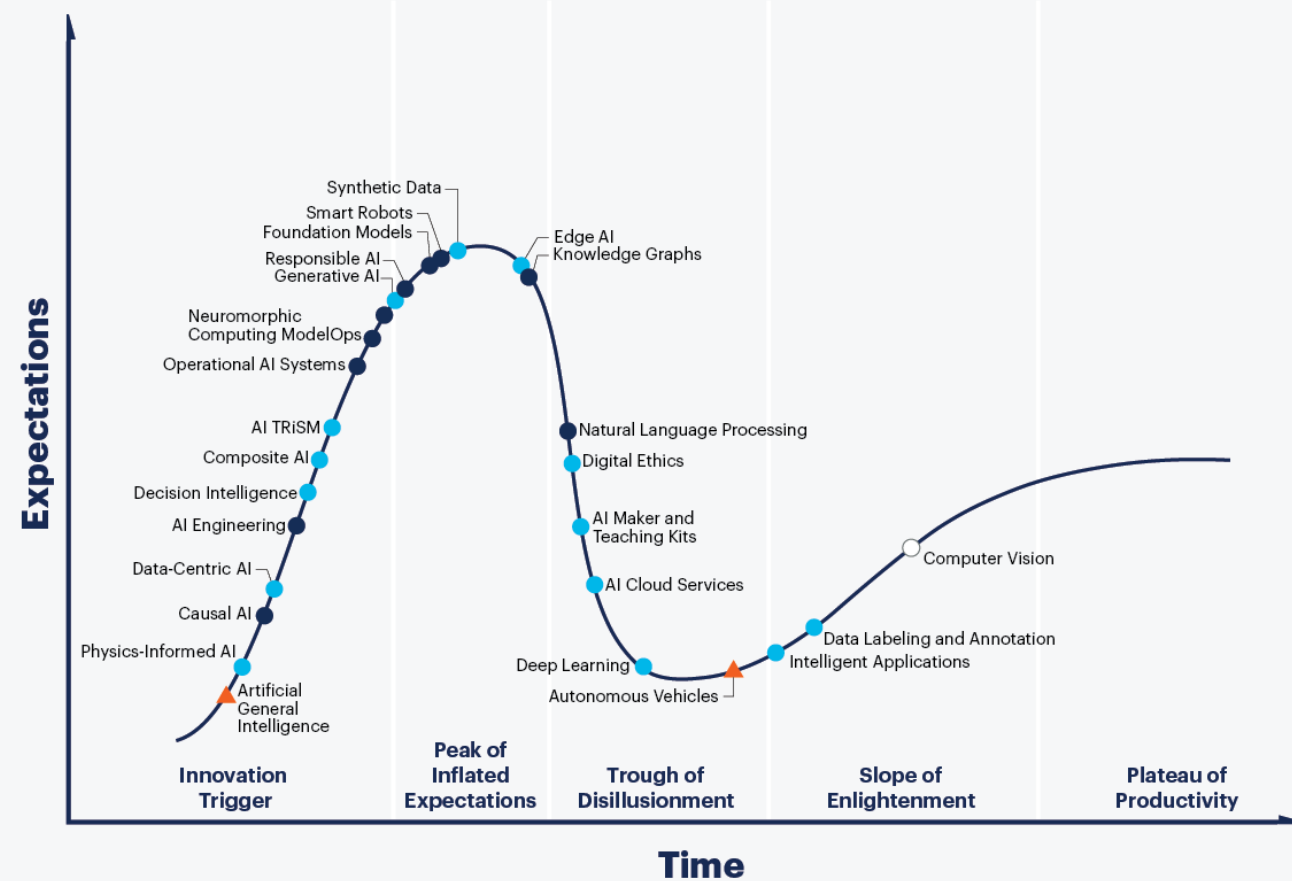
Hype Cycle for Data Management, 2022



Source: Gartner (June 2022)
<https://www.denodo.com/en/document/analyst-report/2022-gartner-hype-cycle-data-management>

Современные
тенденции
развития
технологий
искусственного
интеллекта

Hype Cycle for Artificial Intelligence, 2022



Plateau will be reached:

○ less than 2 years

● 2 to 5 years

● 5 to 10 years

▲ more than 10 years

⊗ obsolete before plateau

As of July 2022

[gartner.com](https://www.gartner.com)

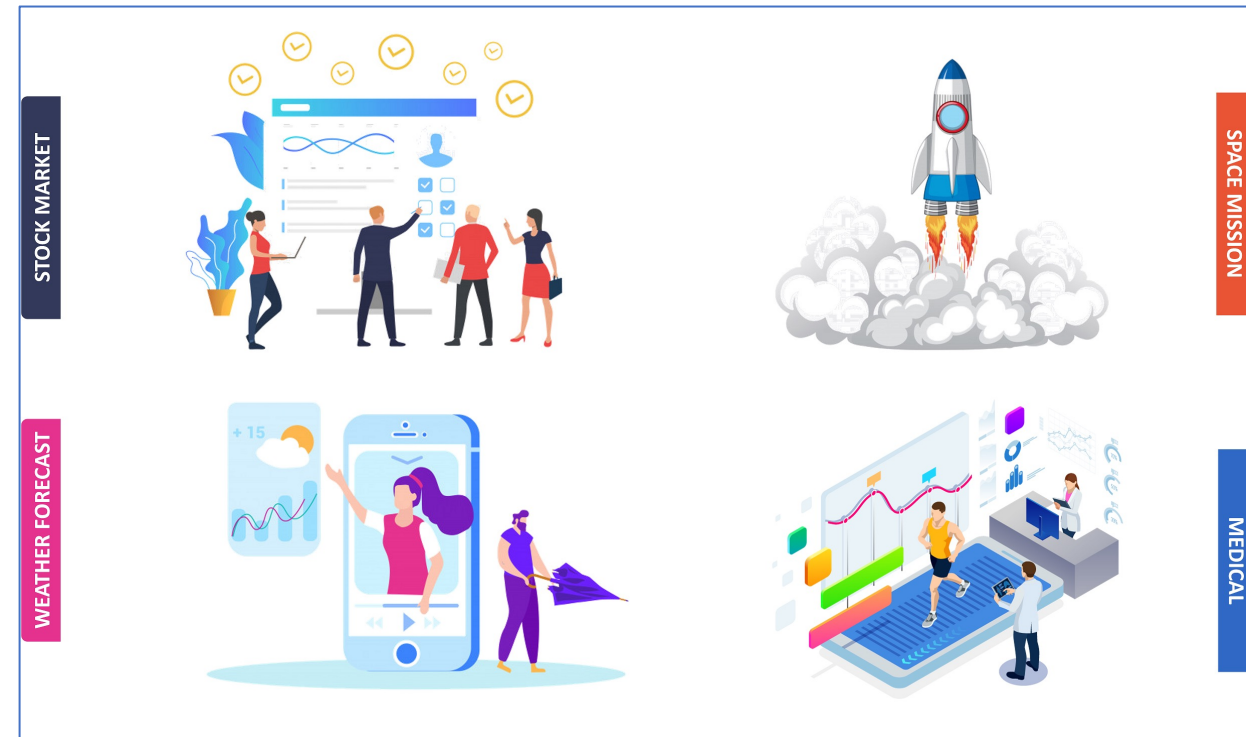
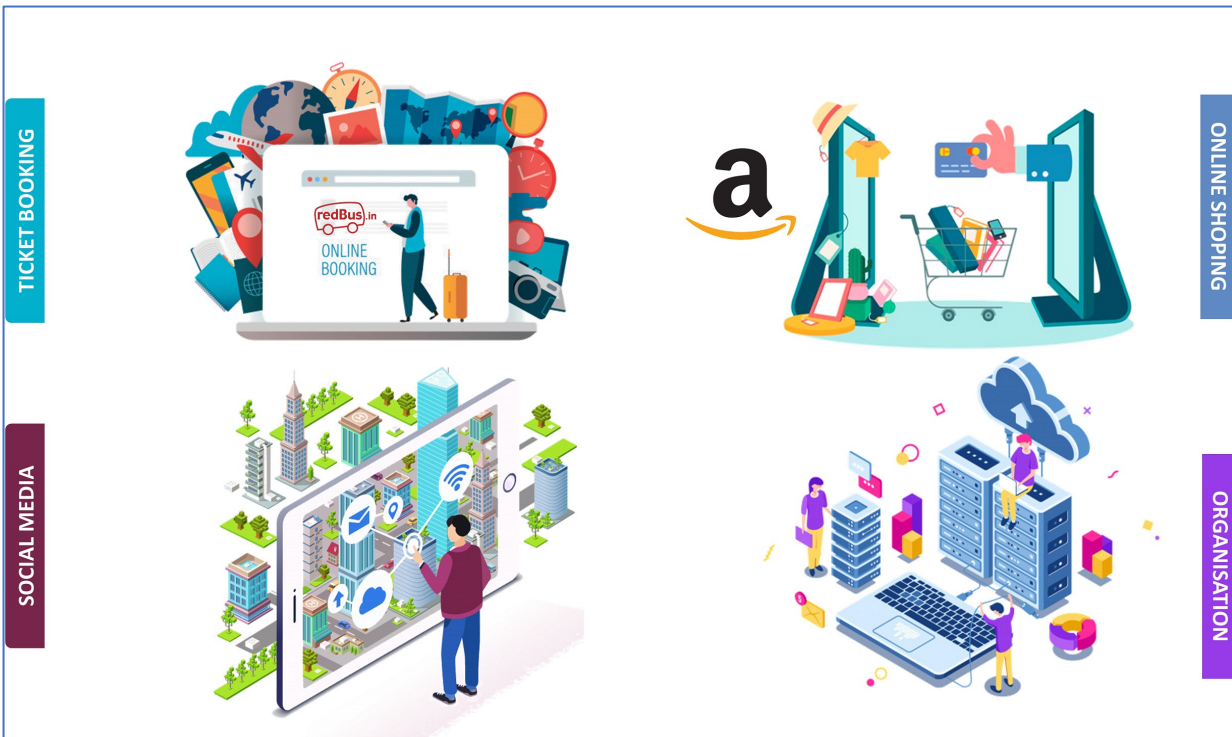
Source: Gartner
© 2022 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner and Hype Cycle are registered trademarks of Gartner, Inc. and its affiliates in the U.S. 1957302

Gartner

Классификация анализа Больших данных по типам использования

Операционная деятельность

Аналитика

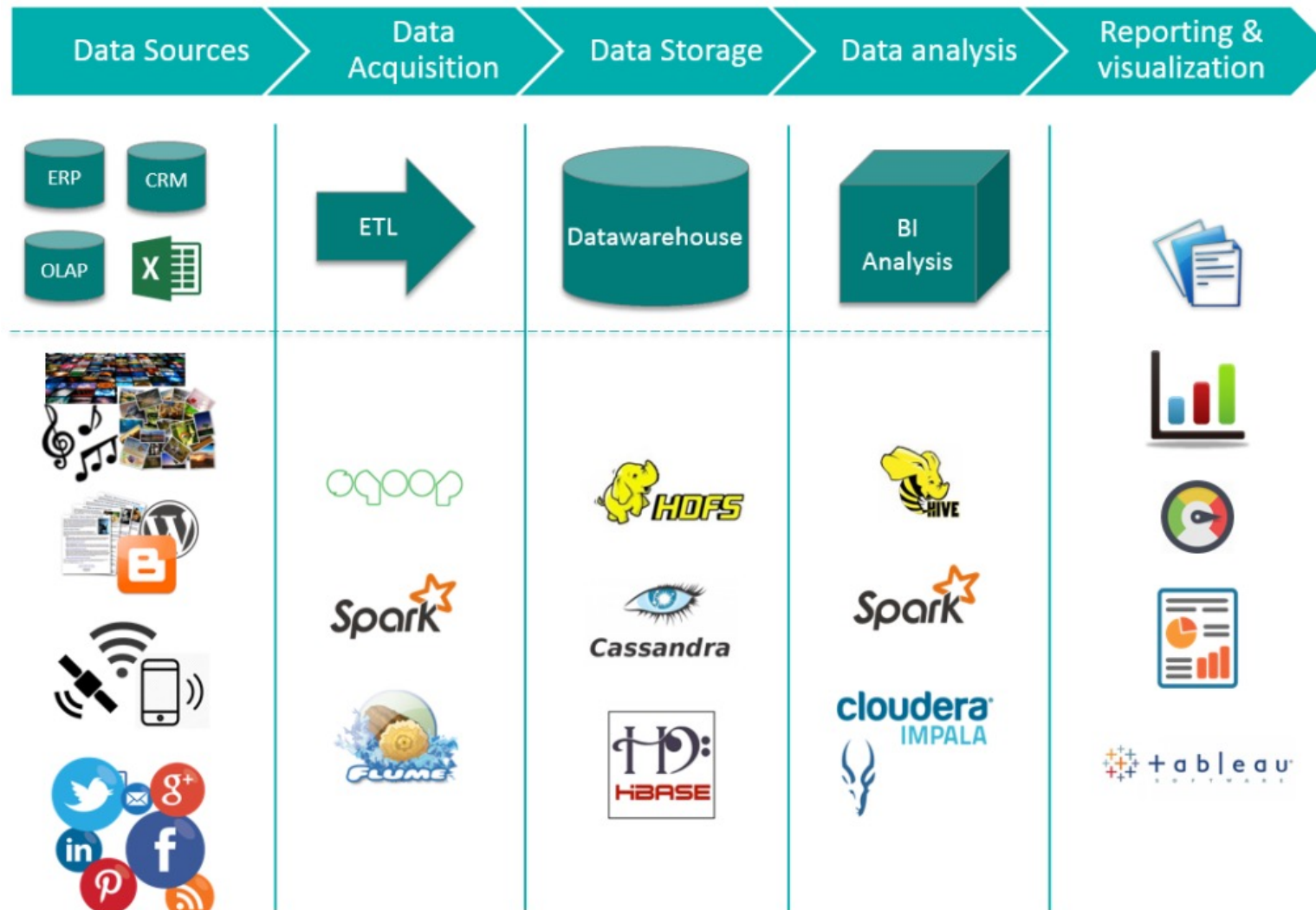


Классификация технологий по функционалу

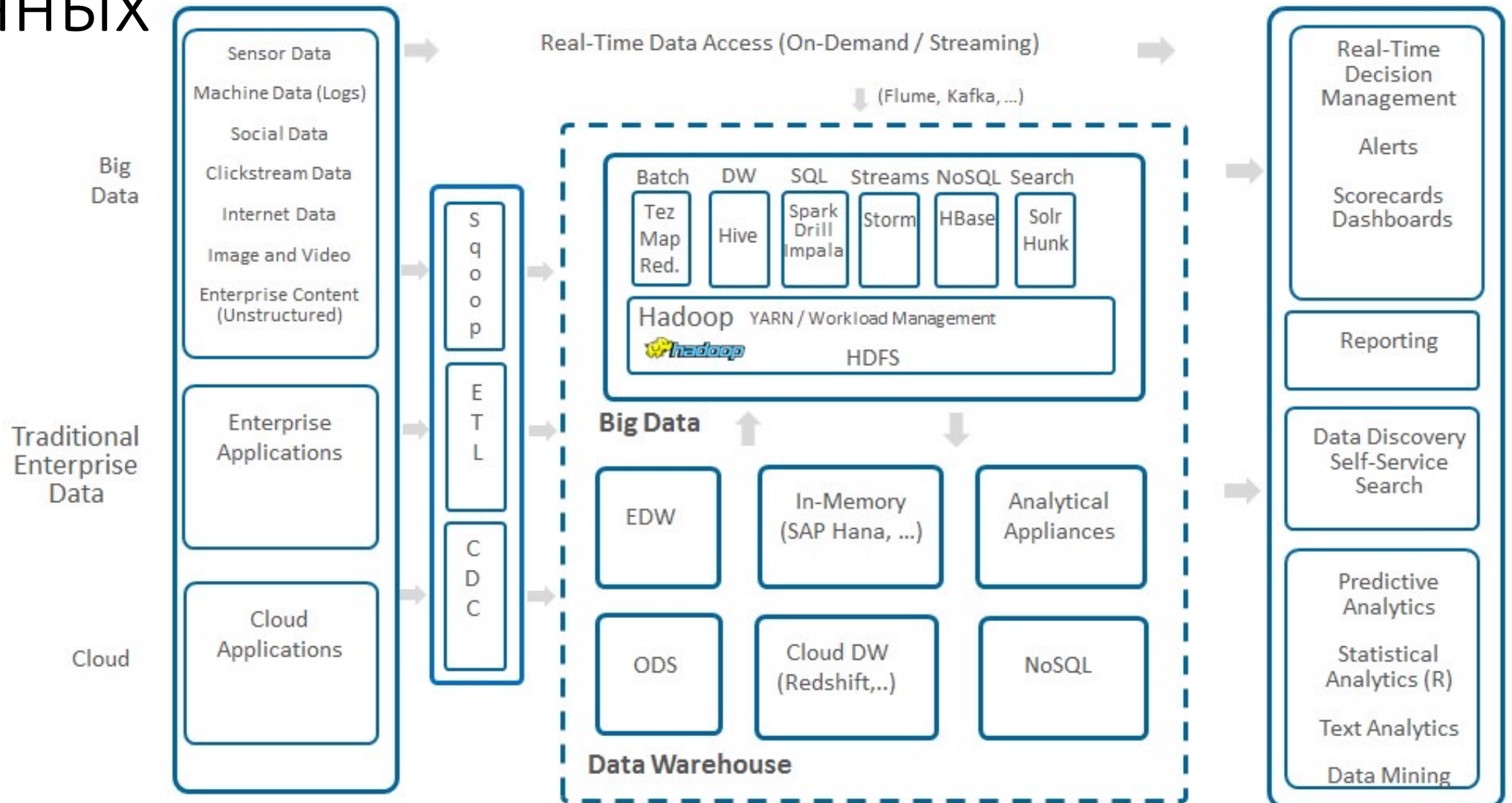
- Основные группы технологий Больших данных
- Хранилища данных (Data Storage)
- Извлечение/добыча данных (Data Mining)
- Аналитика (Data Analytics)
- Визуализация (Data Visualization)



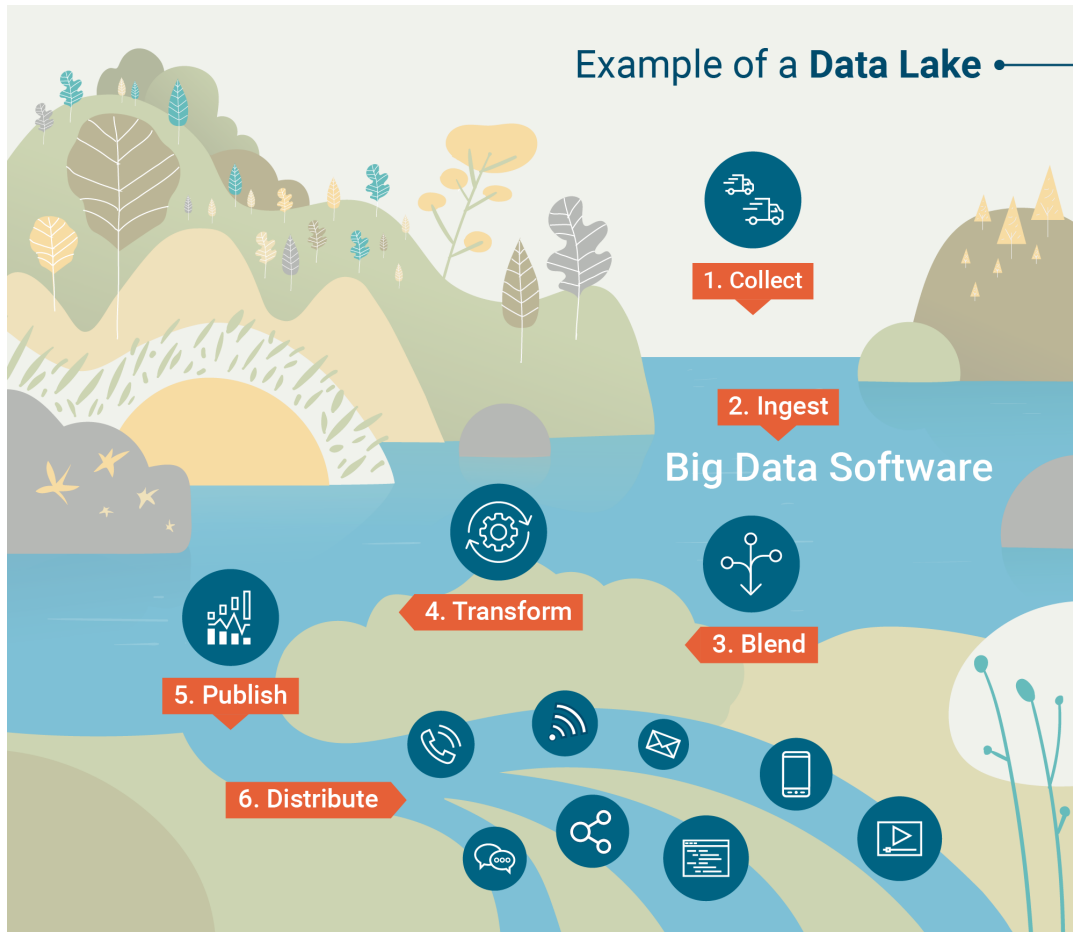
Пример цепочки анализа Больших данных



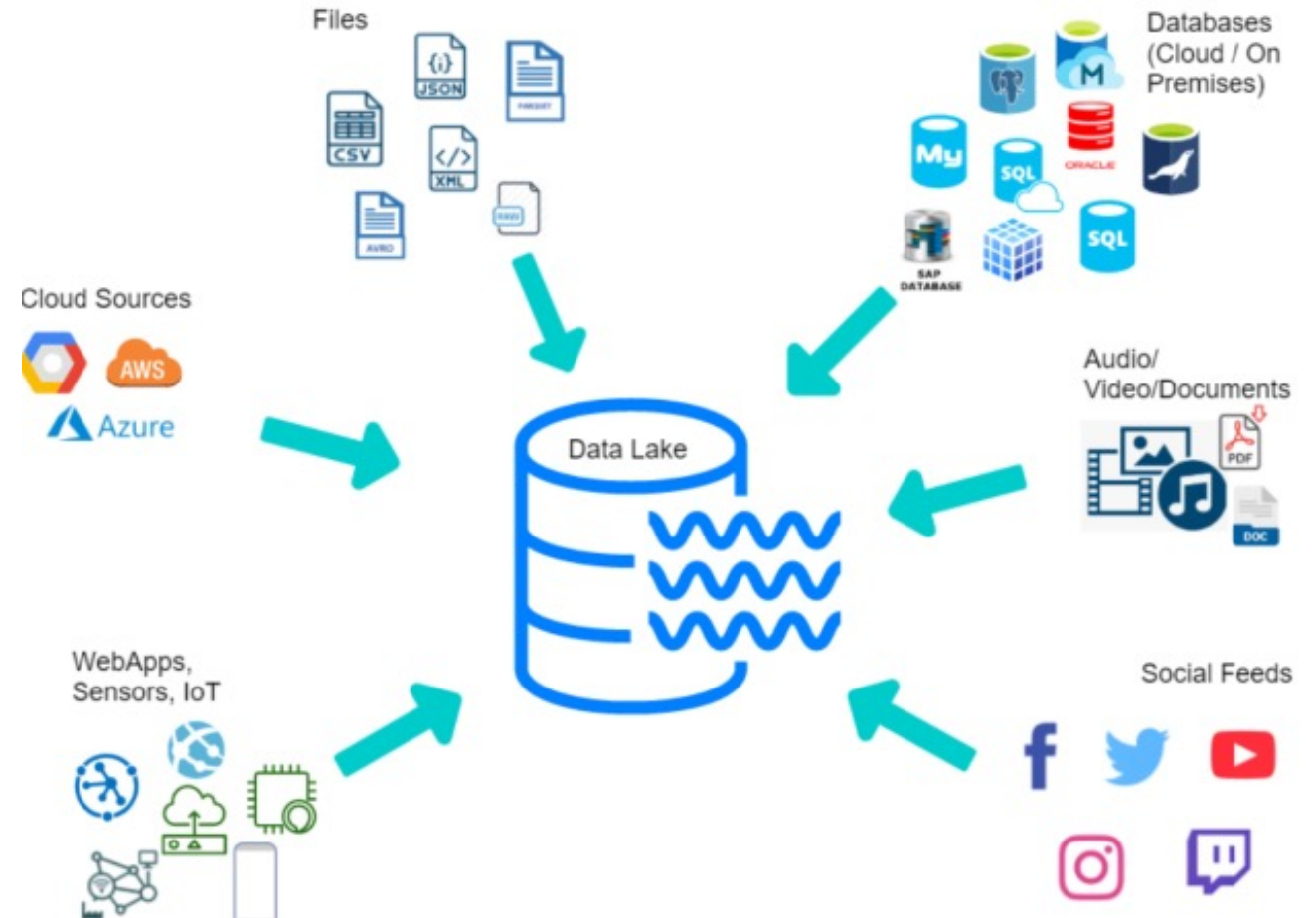
Пример архитектуры обработки Больших данных



Озёра данных (Data Lakes)



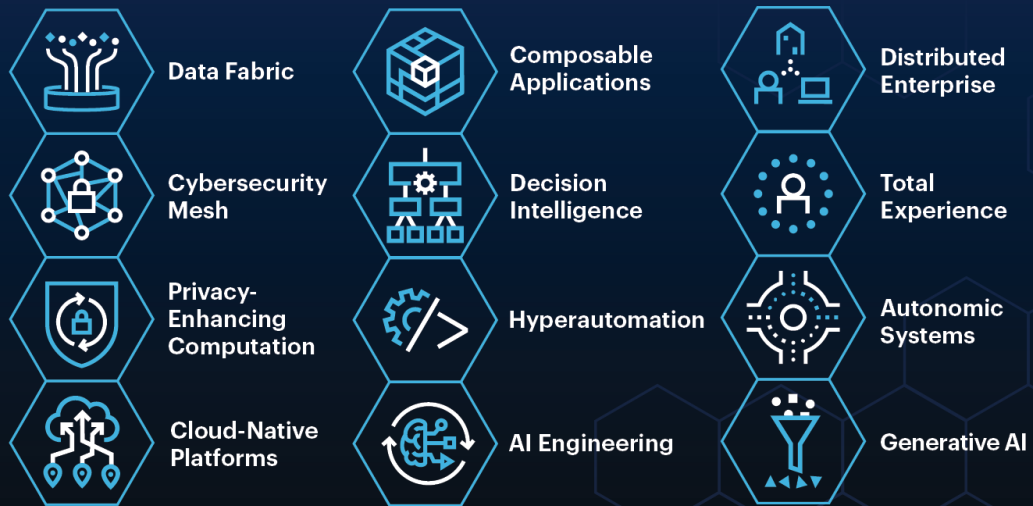
<https://www.g2.com/articles/what-is-a-data-lake>



<https://vitalflux.com/data-lake-design-principles-best-practices/>

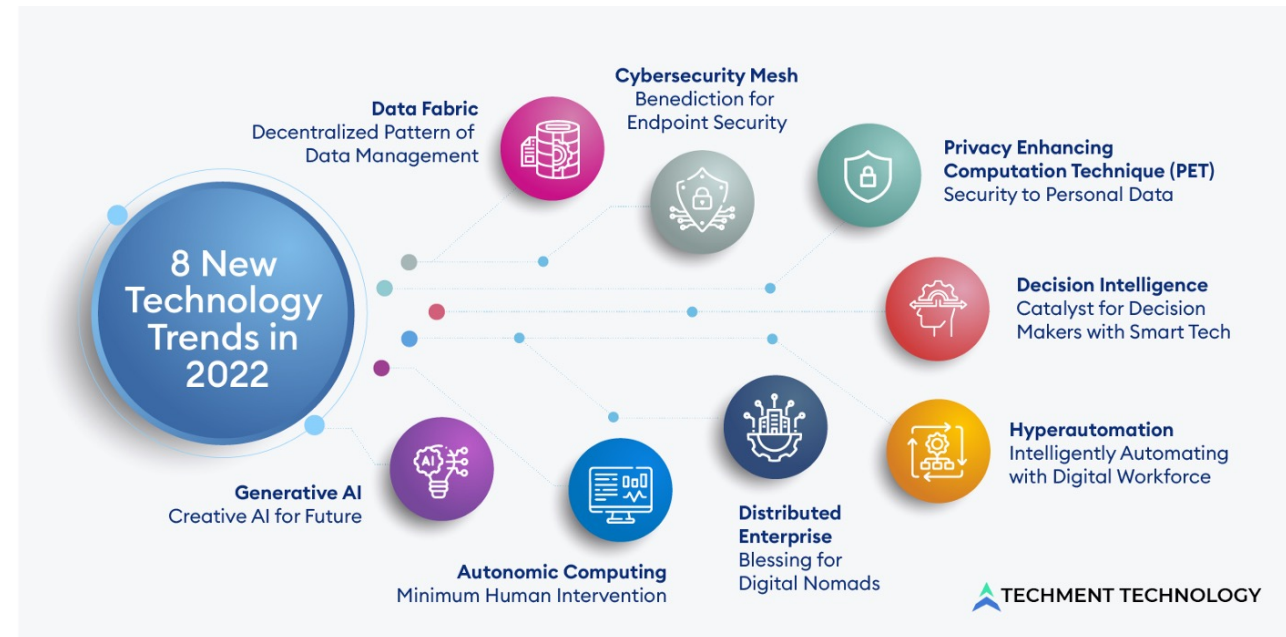
Некоторые технологические тренды

Top Strategic Technology Trends for 2022



© 2021 Gartner, Inc. and/or its affiliates. All rights reserved. 1397600

Gartner



Data Fabric



DATA CONSUMERS

DATA DELIVERY (ODBC/JDBC, APIS, SOAP, REST, ETC)

DATA INTEGRATION (ETL/ELT, DATA VIRTUALIZATION, MESSAGING, STREAMING, REPLICATION, SYNCHRONIZATION, IPAAS)

MASTER AND REFERENCE DATA MANAGEMENT

METADATA MANAGEMENT & DATA CATALOG

DATA SOURCES



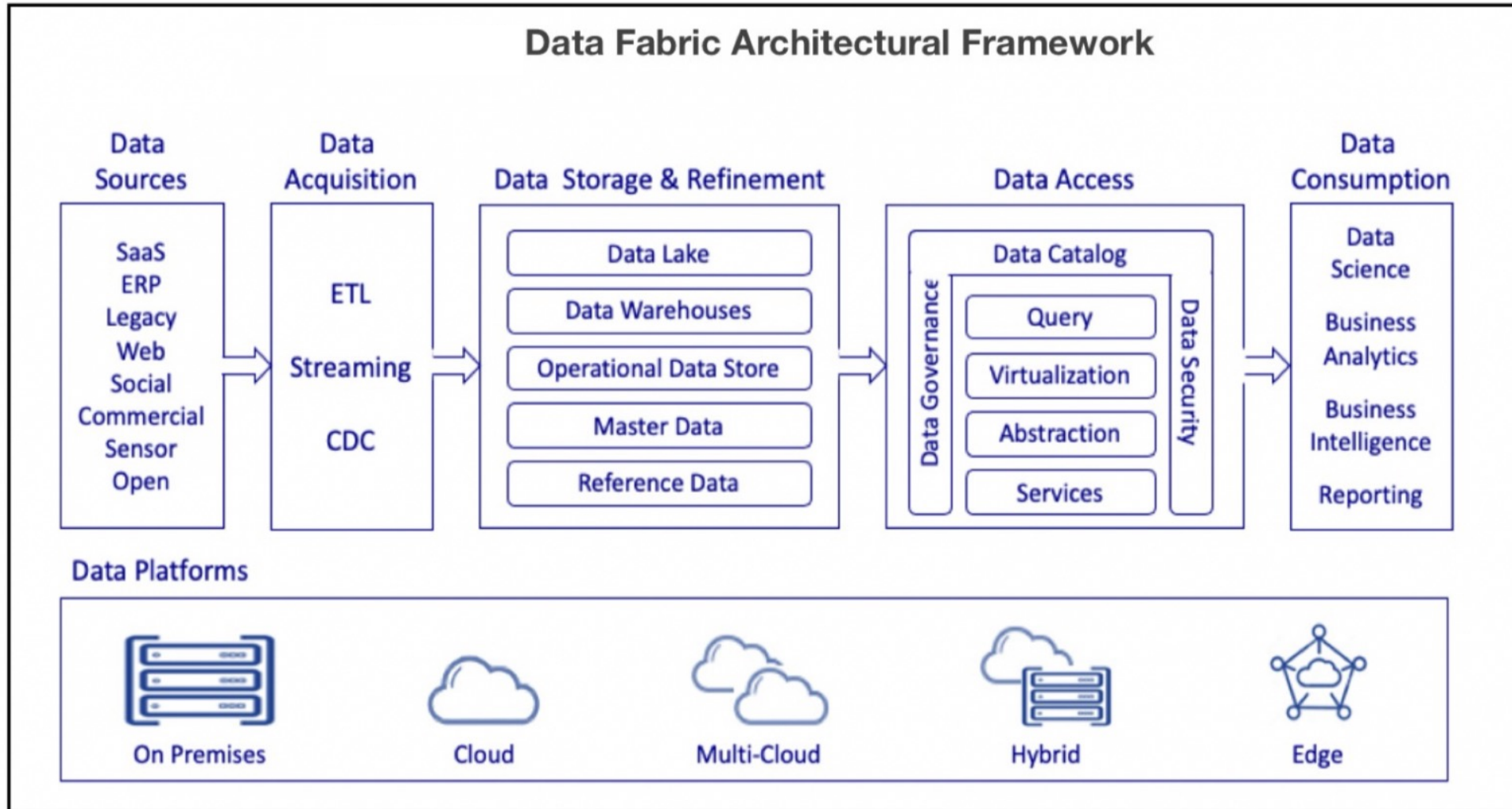
Data Fabric Delivers Integrated Data To All Data Consumers



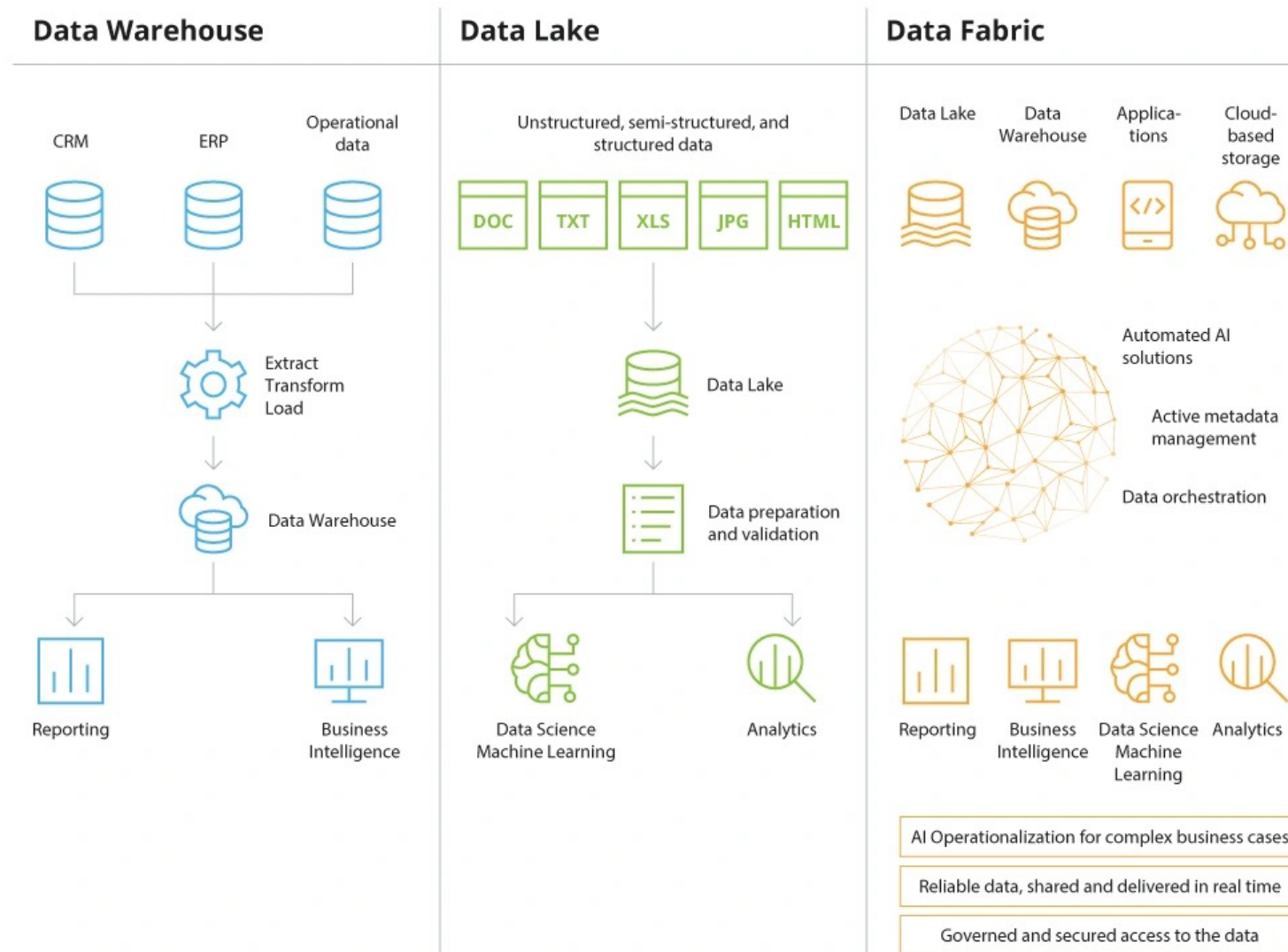
4 © 2021 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner is a registered trademark of Gartner, Inc. and its affiliates.



Архитектура Data Fabric



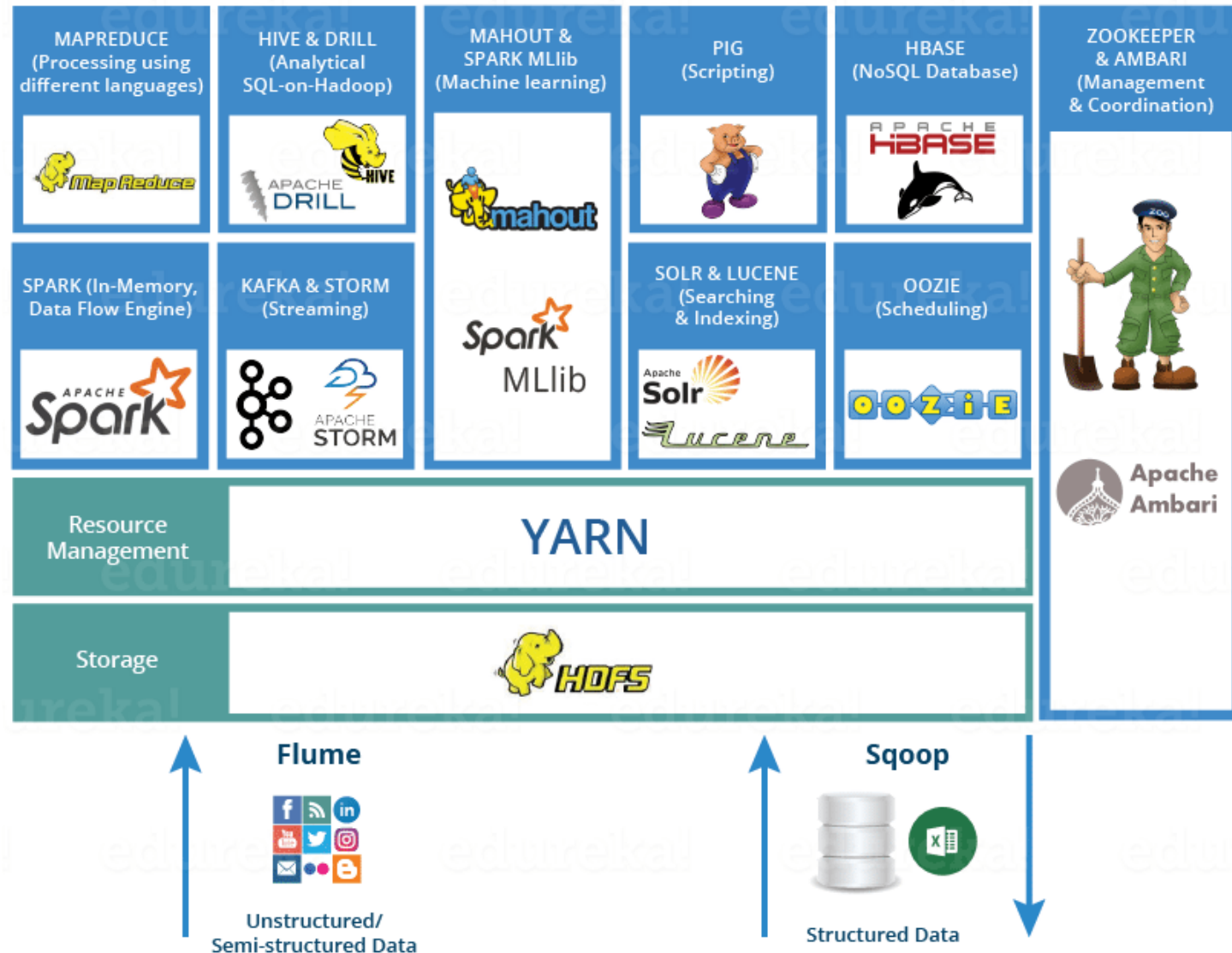
Сравнение Data Warehouse, Data Lake и Data Fabric



Пакеты с открытым исходным кодом



Экосистема Hadoop



Представление и хранение данных



<https://hadoop.apache.org/>



<https://www.mongodb.com/>



<https://cassandra.apache.org/>

- Фреймворк для разработки и выполнения распределенных программ
- Документоориентированная нереляционная (NoSQL) СУБД
- Распределённая нереляционная СУБД

Представление и хранение данных



<https://kylin.apache.org/>



<https://calcite.apache.org/>



<https://ignite.apache.org/>

- Интеграция различных источников данных, предоставление SQL-интерфейса и аналитика
- Динамическое управление разнородными данными, интеграция данных
- Распределённая СУБД с хранением данных в памяти (для высокопроизводительных вычислений)

Представление и хранение данных



<https://neo4j.com/>

- Графовая база данных



<https://hbase.apache.org/>

- Масштабируемая распределённая база данных для работы с большими данными

Извлечение/добыча данных – Data Mining



<https://prestodb.io/>



RAPIDMINER

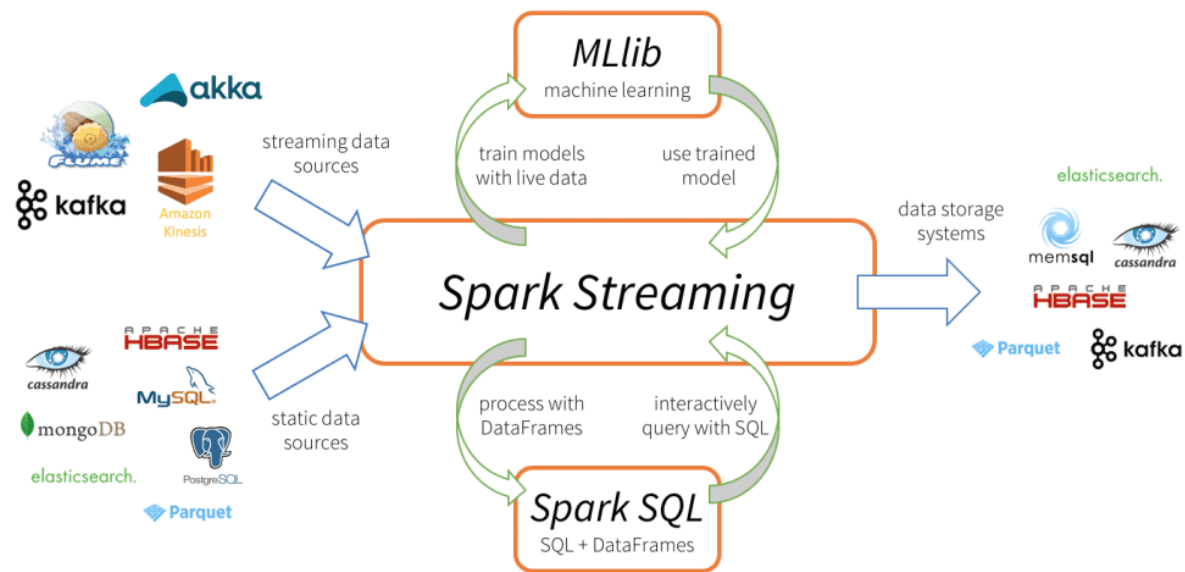
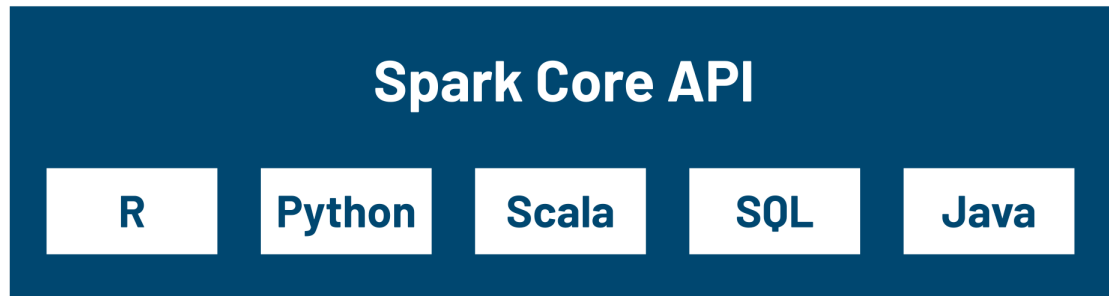
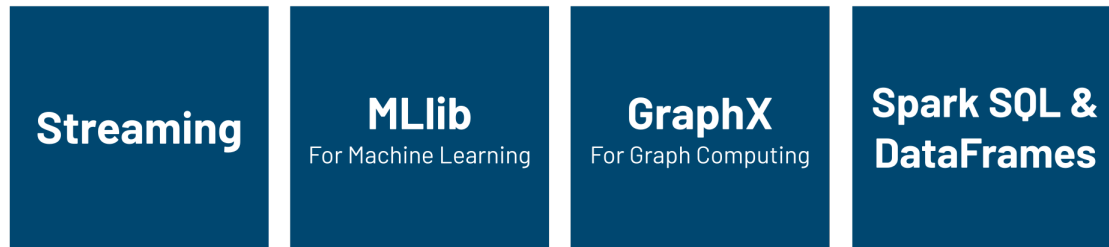
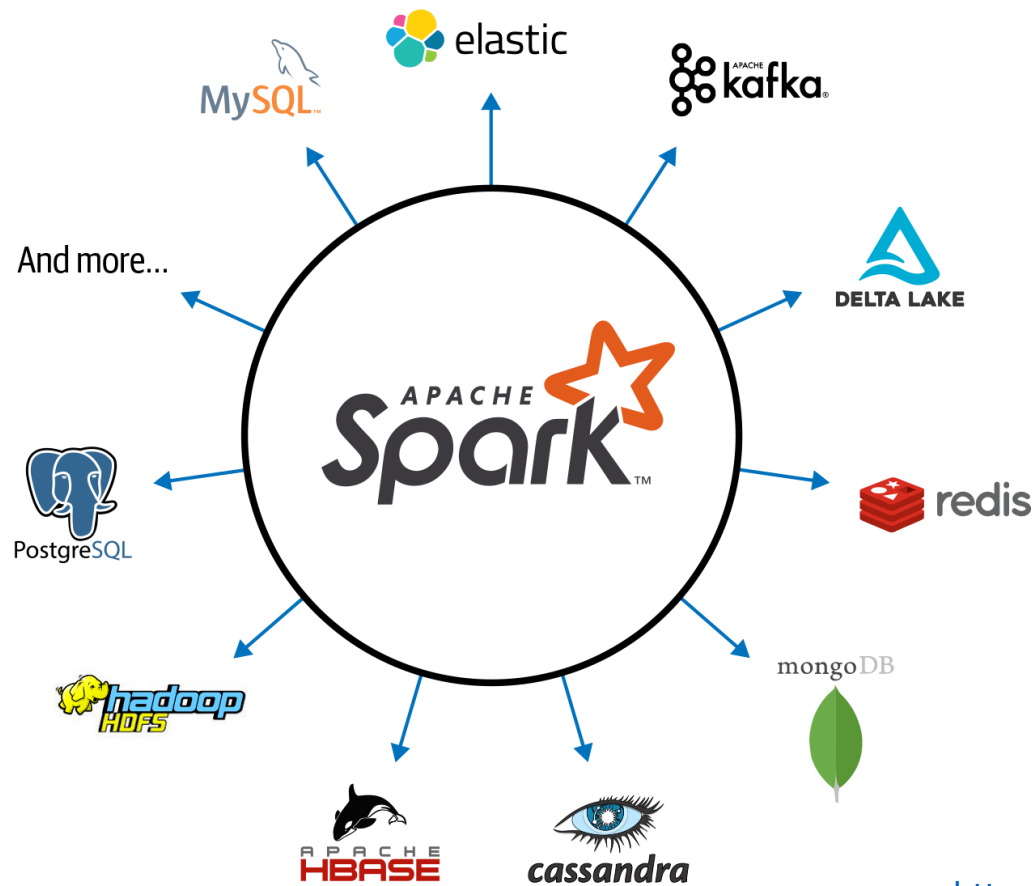
<https://github.com/rapidminer/>



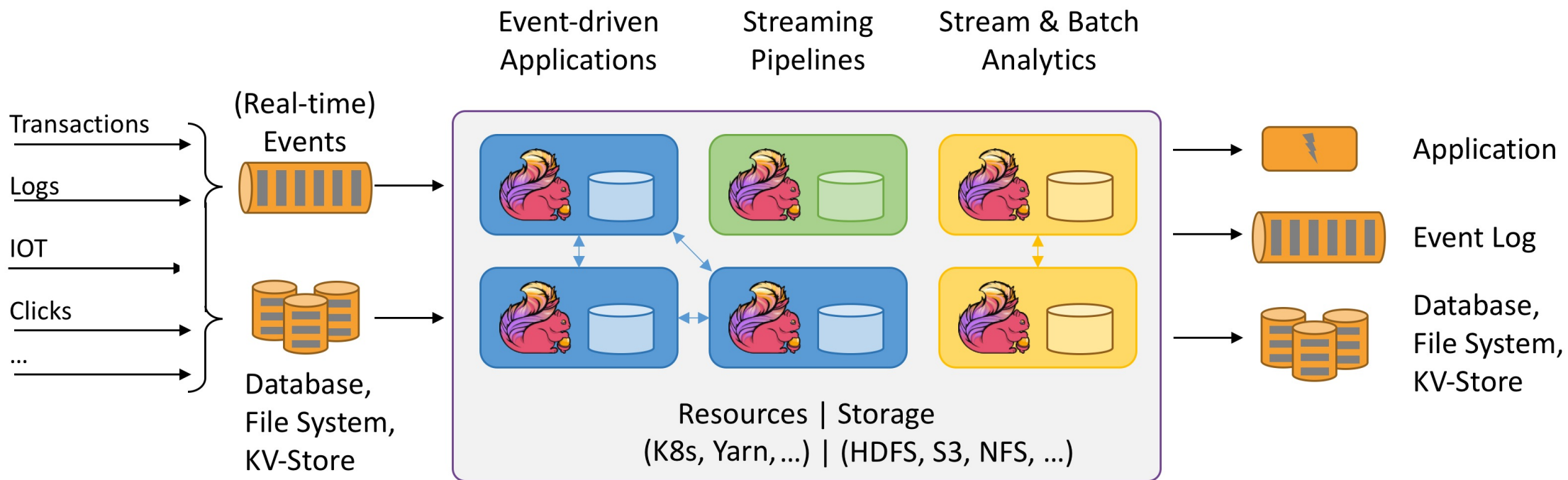
<https://www.elastic.co/>

- Система SQL-запросов к различным СУБД (в том числе нереляционным)
- Платформа для анализа данных
- Программная поисковая система

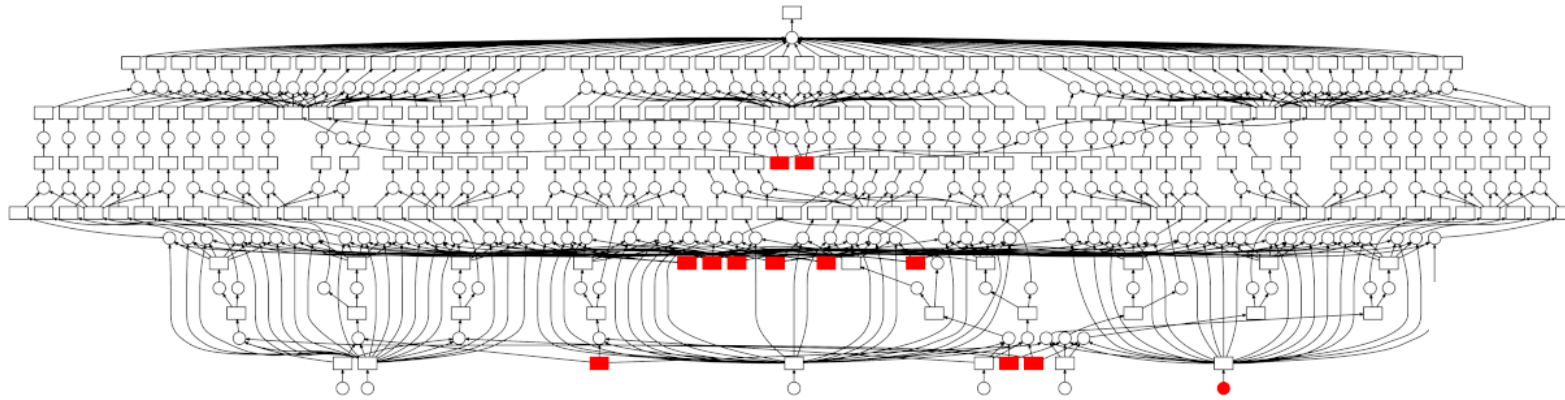
Анализ данных – Apache Spark



Анализ данных – Apache Flink

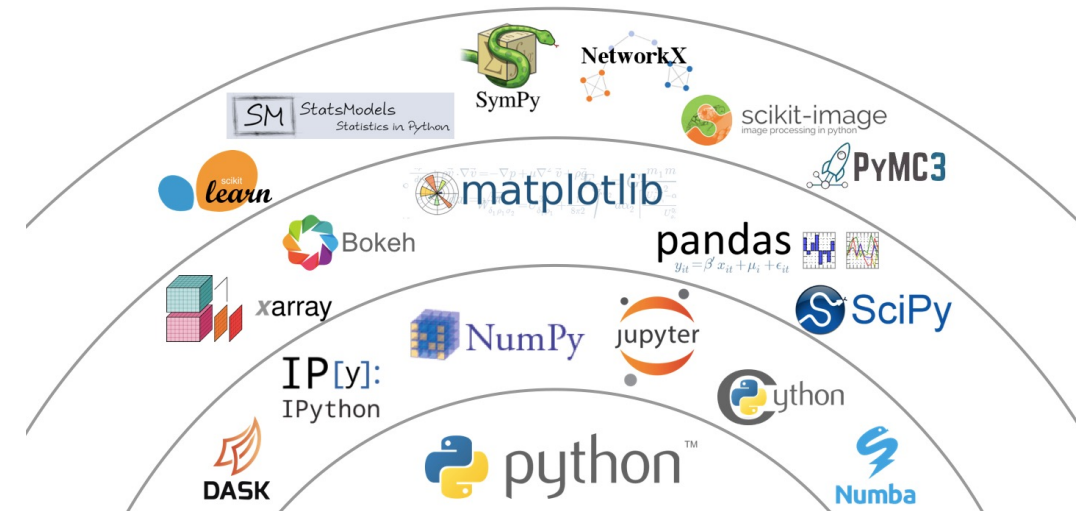


Анализ данных - dask



Python's Scientific Stack

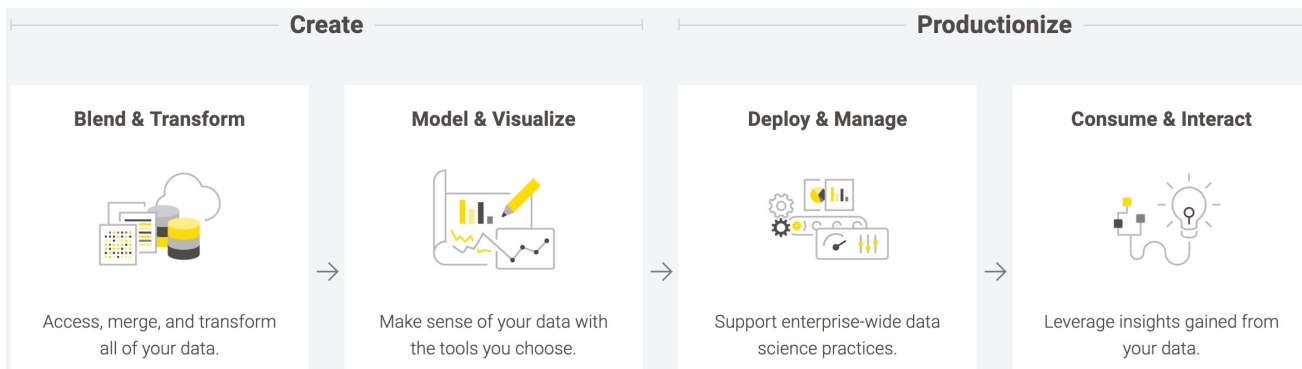
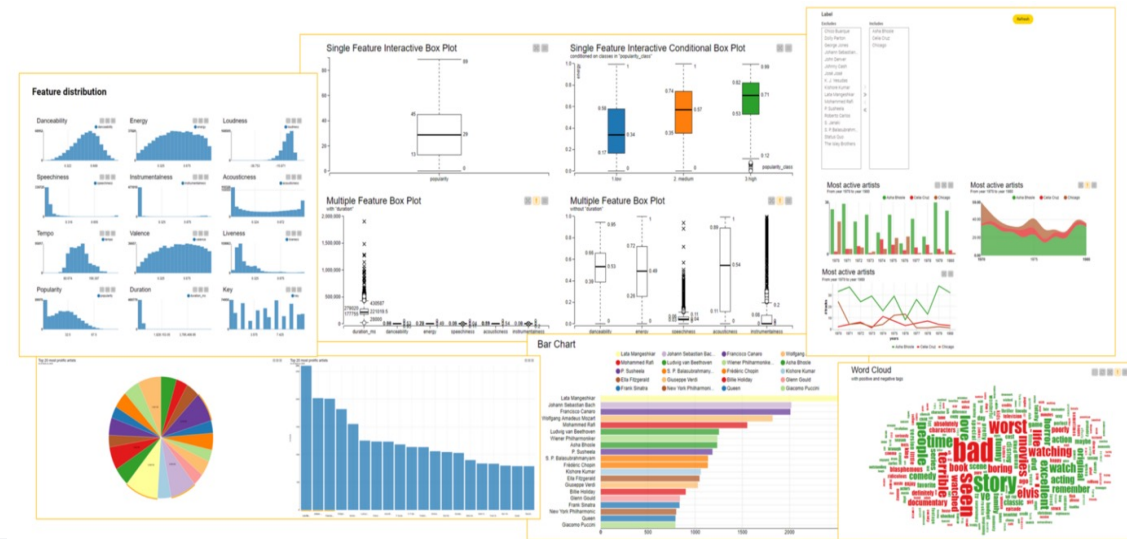
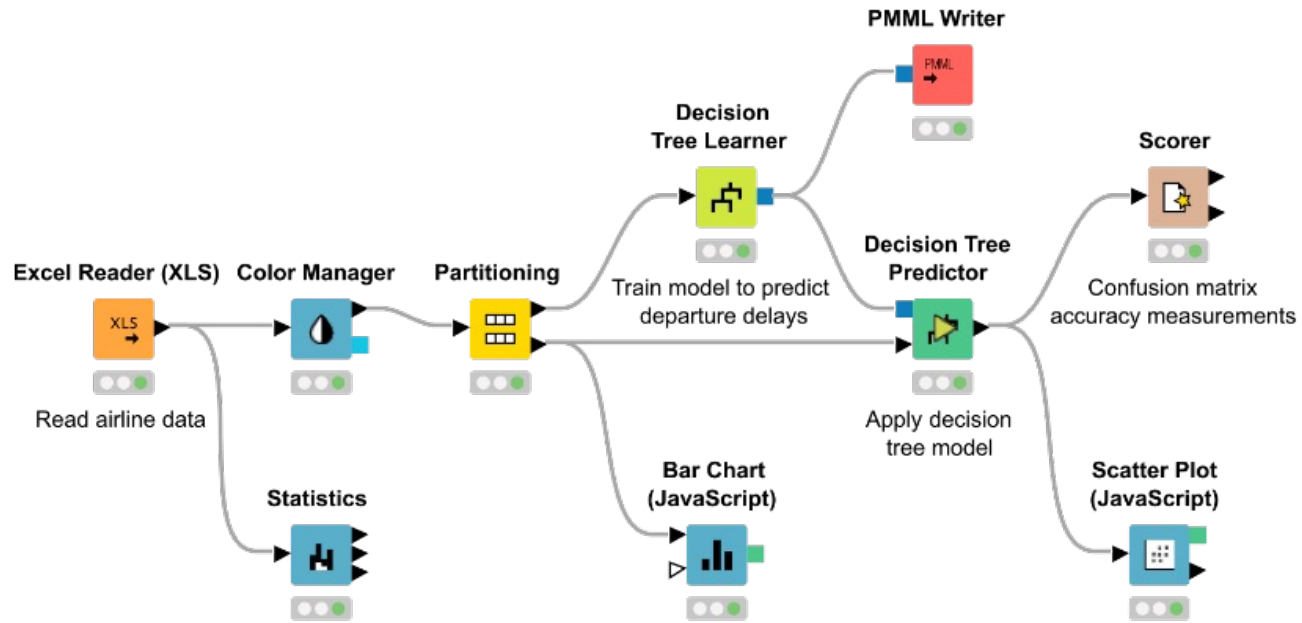
Jake Vanderplas PyCon 2017 Keynote



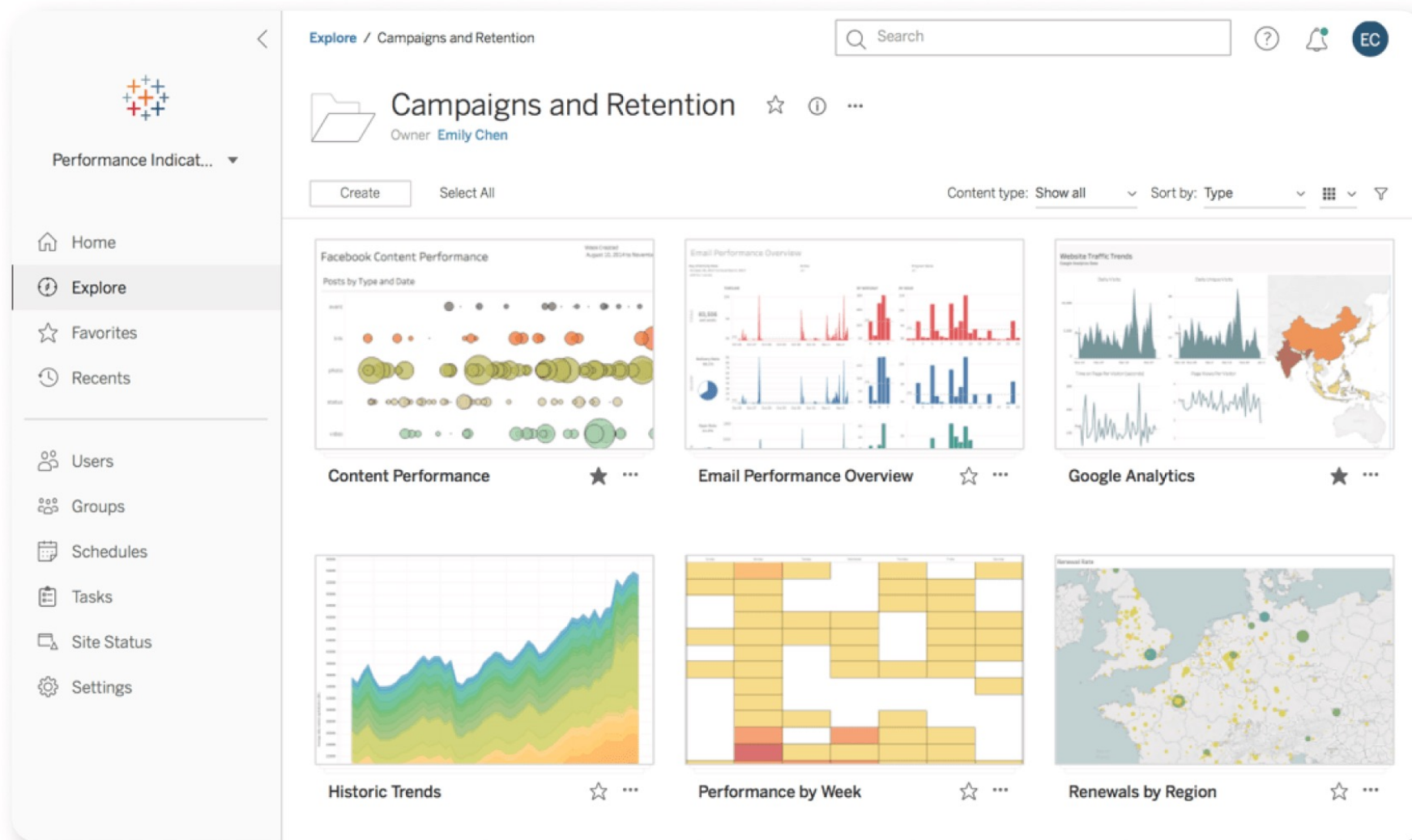
<https://www.dask.org/>

Масштабируемые параллельные вычисления для Python¹

Аналитика - KNIME



Анализ, визуализация - Tableau

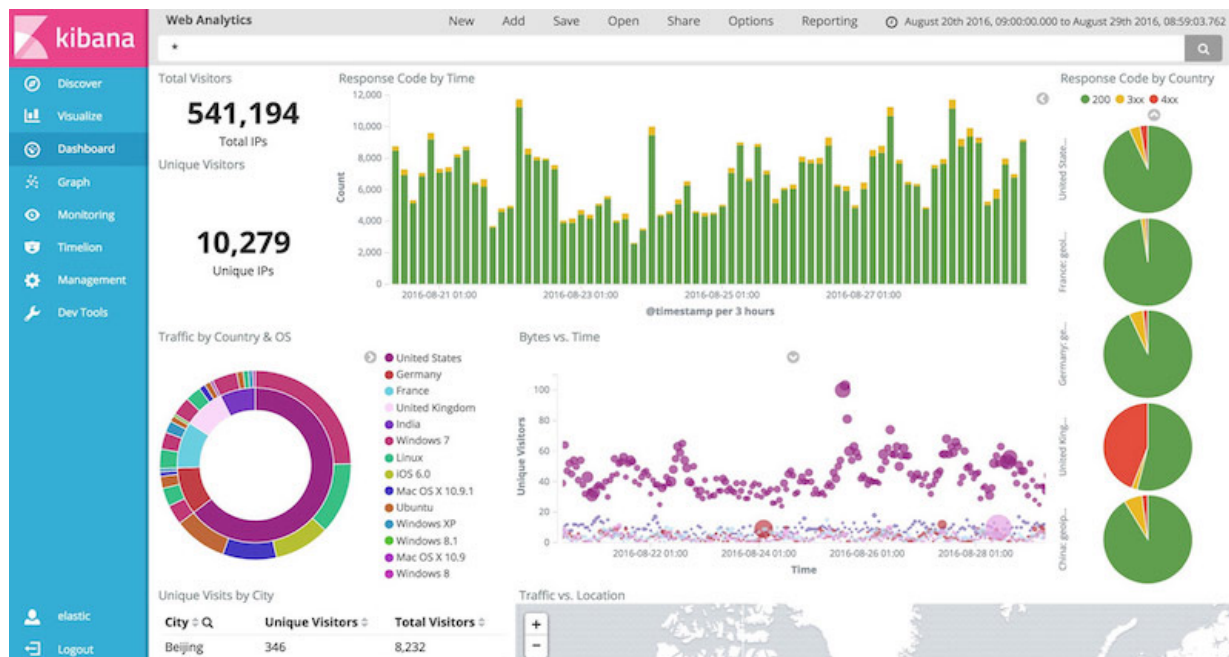


Визуализация, анализ – Apache Superset

The screenshot displays the Apache Superset dashboard interface. At the top, there is a dropdown menu for 'CHOOSE A DATASET' with 'video_game_sales' selected. Below this is a search bar for charts and a list of recommended tags such as 'Highly-used', 'ECharts', and 'Advanced-Analytics'. The main area is a grid of 21 different chart types, including Area Chart, Time-series Bar Chart, Big Number with Trendline (showing 215 +7.0% WoW), Big Number (showing 80.7M), Box Plot, Bubble Chart, Calendar Heatmap, Time-series Percent Change, Country Map, deck.gl Arc, deck.gl Geojson, deck.gl Grid, deck.gl 3D Hexagon, deck.gl Multiple Layers, deck.gl Polygon, deck.gl Scatterplot, deck.gl Screen Grid, Bar Chart, Time-series Area Chart, Time-series Chart, and Time-series Bar Chart v2.

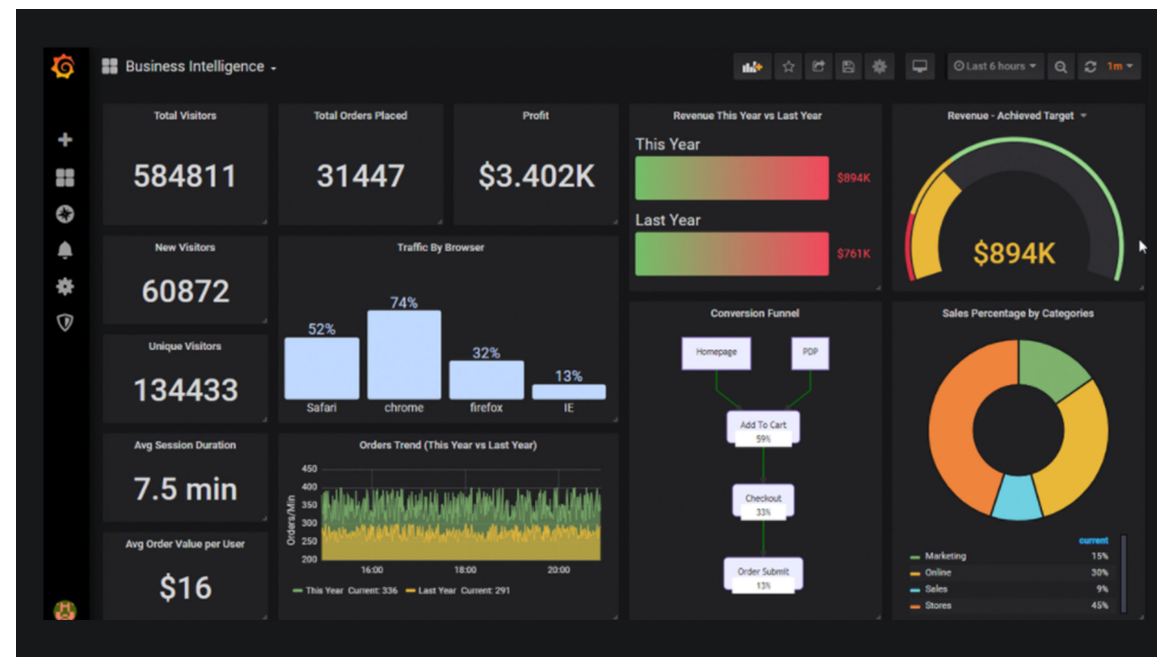


Визуализация – Grafana, Kibana



+ Исследование данных

<https://www.elastic.co/kibana/>



+ Временные ряды

<https://github.com/grafana/grafana>

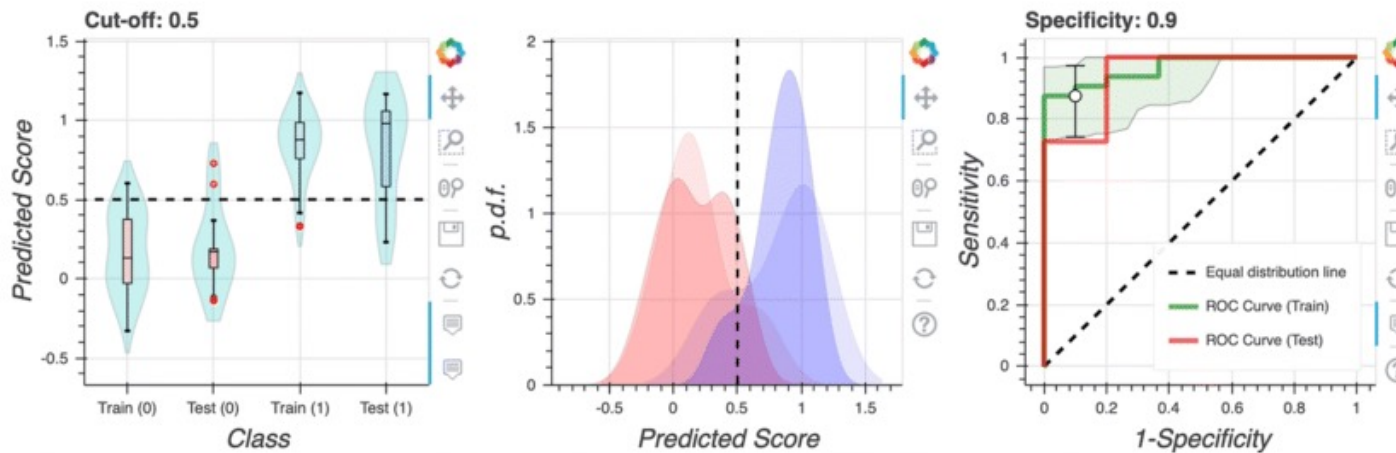
Исследование данных, визуализация (I)

```
# Calculate Ypredicted score using modelPLS.test
YVpred = modelPLS.test(XVknn)

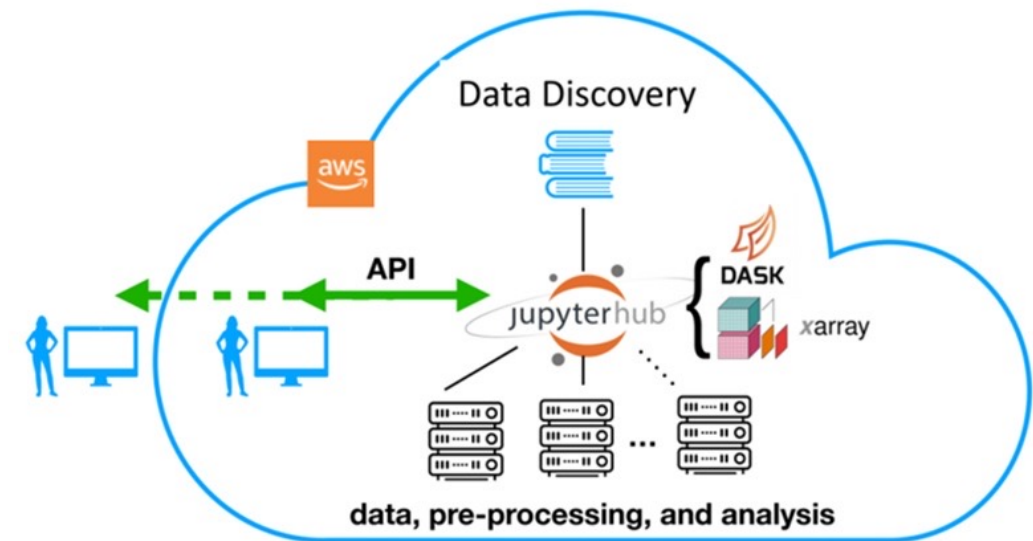
# Evaluate Ypred against Ytest
evals = [Ytest, YVpred] # alternative formats: (Ytest, Ypred) or np.array([Ytest, Ypred])
#modelPLS.evaluate(evals, specificity=0.9)
modelPLS.evaluate(evals, cutoffscore=0.5)
```

BokehJS 1.1.0 successfully loaded.

Score cut-off fixed to: 0.5

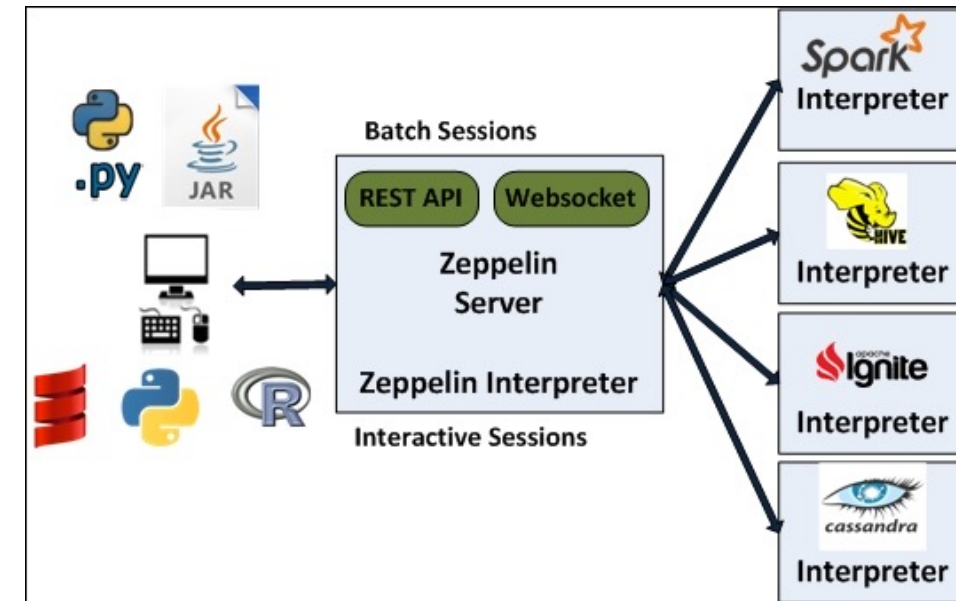
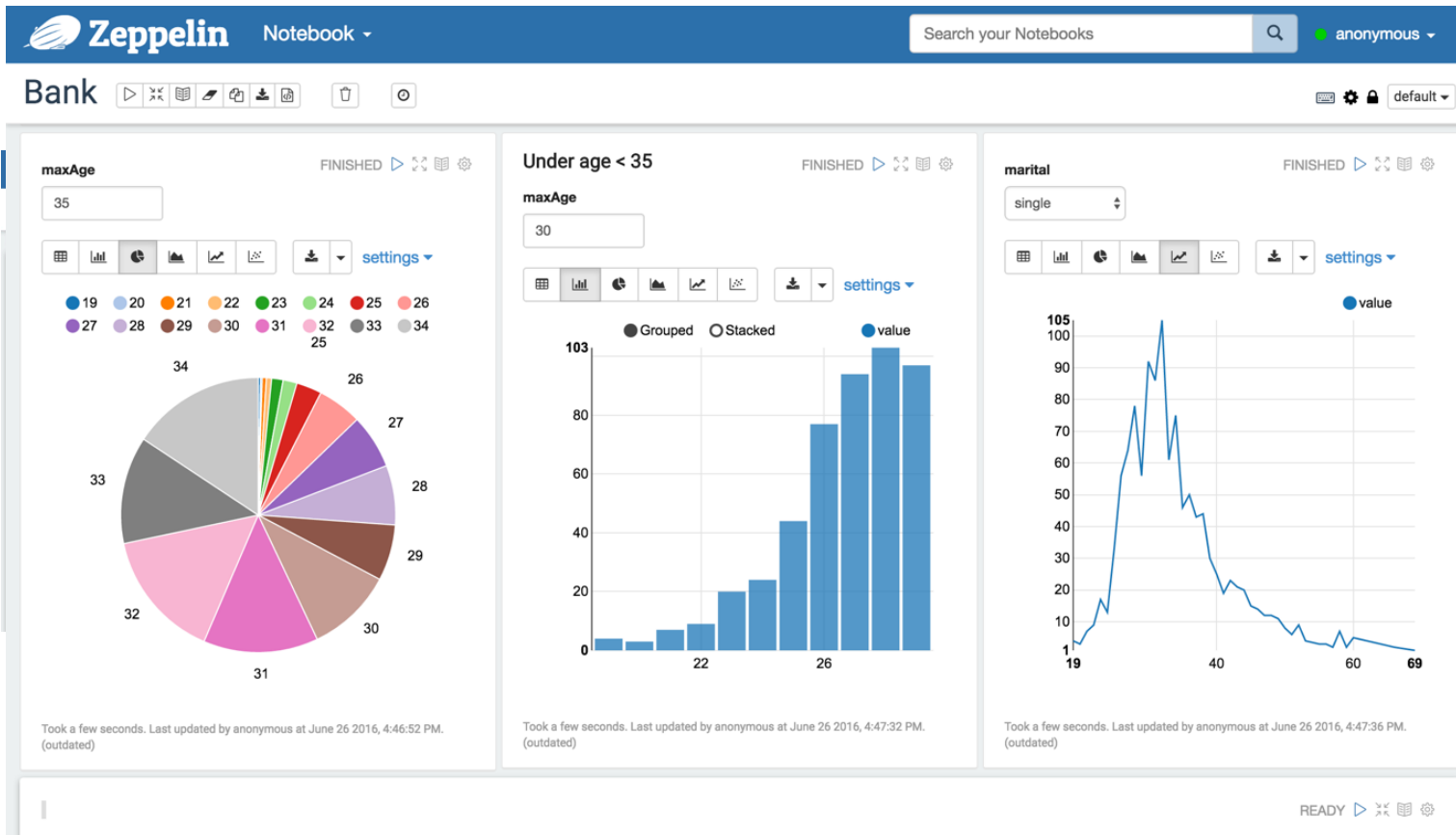


#	Evaluate	MW-U Pvalue	R2	AUC	Accuracy	Precision	Sensitivity	F1score
0	Train	2.66e-10	0.67 (0.5, 0.77)	0.97 (0.91, 1.0)	0.89 (0.81, 0.94)	0.9 (0.82, 0.95)	0.88 (0.74, 0.97)	0.89 (0.8, 0.95)
1	Test	6.37e-04	0.52	0.95	0.76	0.8	0.73	0.76



<https://jupyter.org/hub>

Исследование данных, визуализация (II)



Вычислительные эксперименты и машинное обучение – MLflow



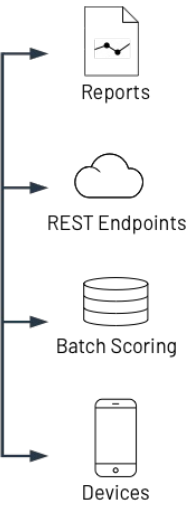
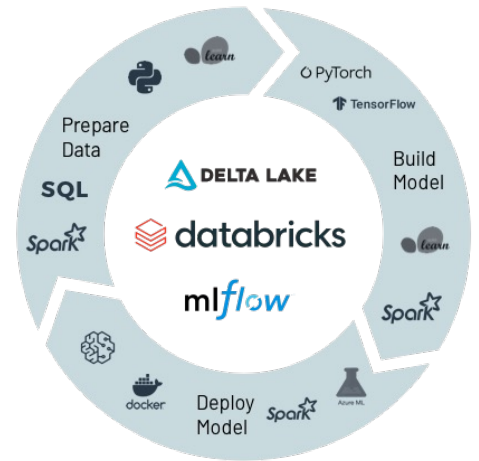
Files



Big Data



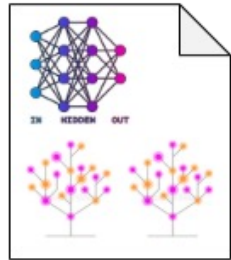
Streams



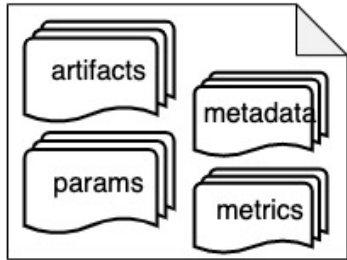
XGBoost



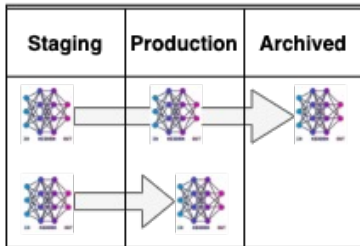
Models



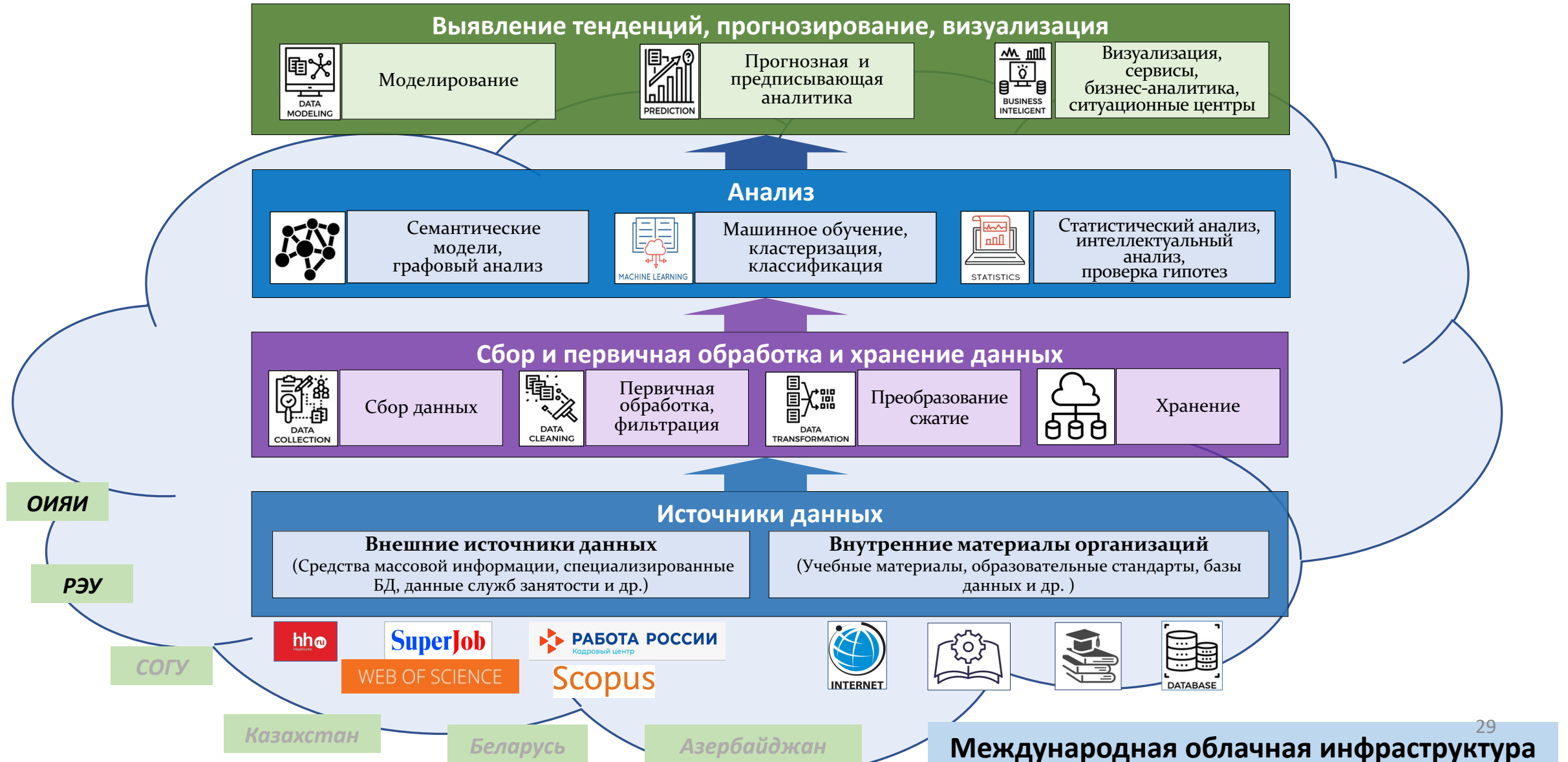
Tracking



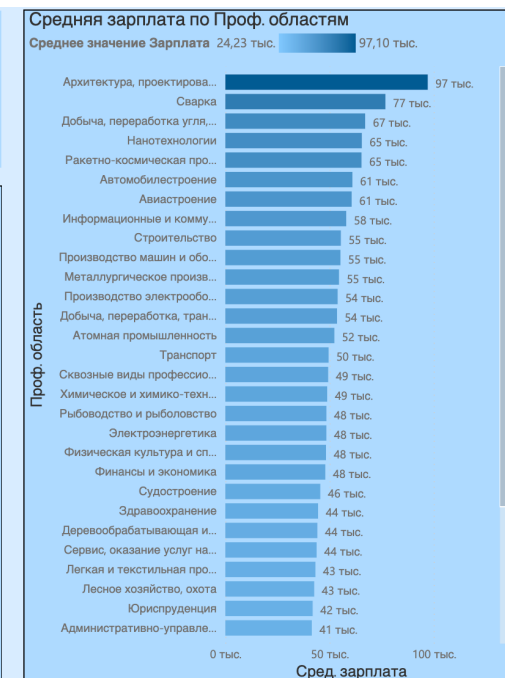
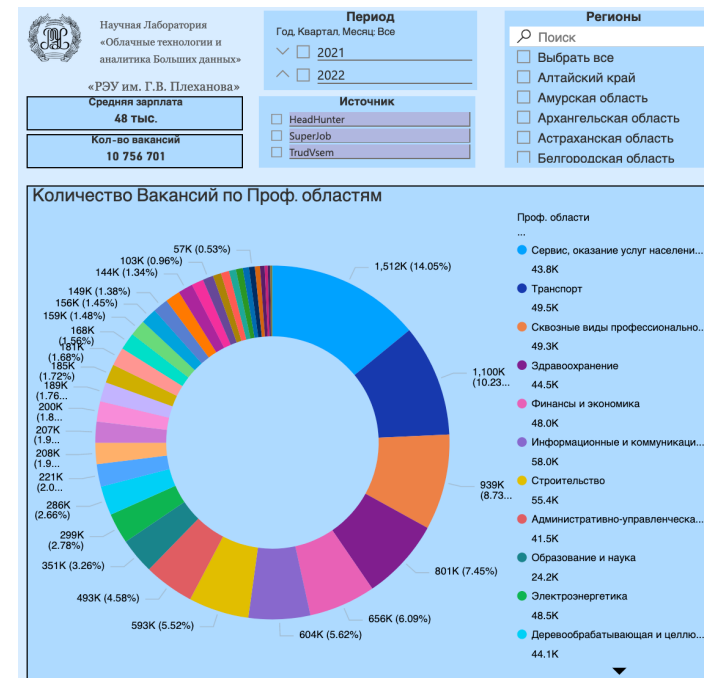
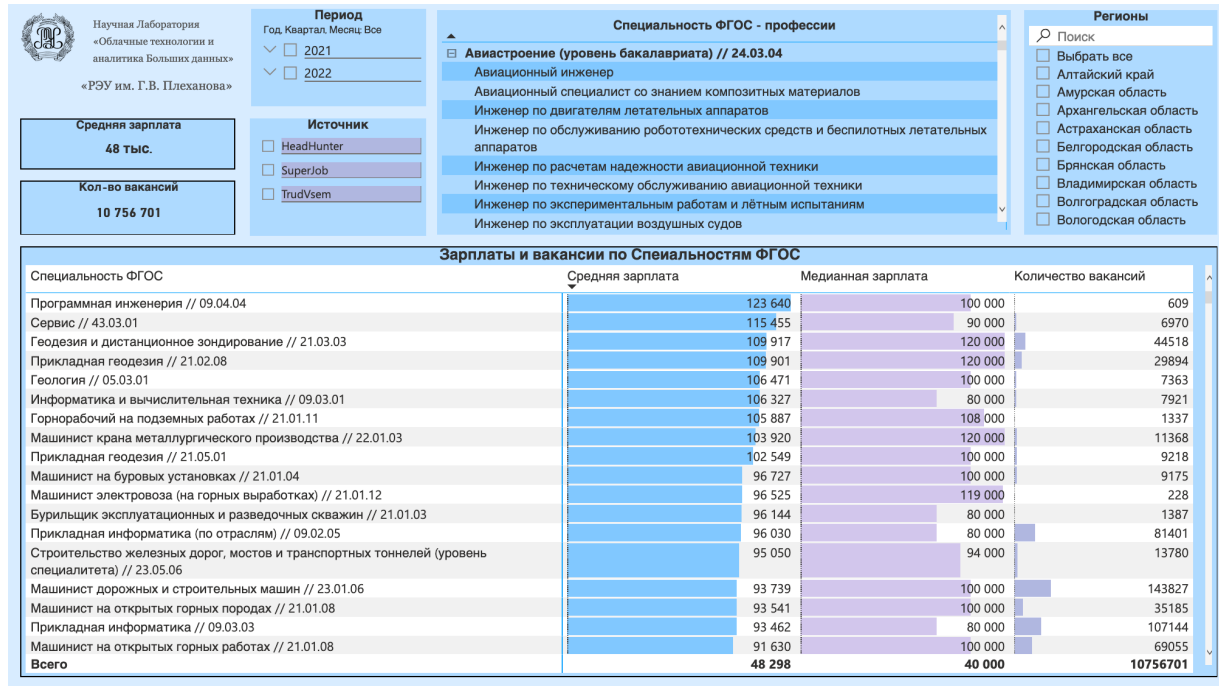
Registry



Пример реализации платформы Больших данных



Приложение: анализ рынка труда



Приложение: анализ рынка труда



Экосистема программных компонентов

Уровень	Программные пакеты
Визуализация и интерфейсы для доступа к системе	<ul style="list-style-type: none"> Zeppelin, Jupyter (пользовательский интерфейс) Graphana (создание отчетов и графическое представление результатов) KrakenD (организация программных шлюзов для различных компонентов) 
Распределенная аналитика Больших данных	<ul style="list-style-type: none"> Apache Kylin 
Вычислительные эксперименты в области машинного обучения	<ul style="list-style-type: none"> MLflow 
Вычисления в памяти	<ul style="list-style-type: none"> Apache Spark, Dask, Hadoop 
Организация процесса управлением потоками данных и сбором данных	<ul style="list-style-type: none"> Apache Kafka, Apache Flume, Apache Airflow, Celery, Scrapy 
Хранилища и специализированные базы данных	<ul style="list-style-type: none"> СЕРН, NFS (хранение и доступ к файлам) Elasticsearch (хранение и анализ структурированных данных) Apache Ignite (база с хранением данных в памяти для быстрого доступа и кеширования) Российское «озеро данных» Apache Calcite (интеграция хранилищ) 
Аутентификация и сквозная авторизация, безопасность	<ul style="list-style-type: none"> Free IPA, Vault 
Компьютерная инфраструктура, управление ресурсами	<ul style="list-style-type: none"> OpenNebula, Kubernetes, Puppet Docker, Git 

