

СИСТЕМЫ ОБРАБОТКИ БОЛЬШИХ ДАННЫХ

Разработка приложений



К.Т.Н.
Папулин Сергей Юрьевич
papulin_bmstu@mail.ru

Программа курса



- Hadoop: Hadoop Distributed File System (HDFS)
- Управление ресурсами и приложениями (YARN), платформа MapReduce
- Apache Spark. Распределенная координация с Zookeeper
- Системы потоковой обработки. Apache Storm. Spark Streaming.
- Системы потоковой обработки. Apache Flink. Kafka
- Системы обработки графов: Giraph, Spark GraphX, Spark GraphFrames
- NoSQL СУБД. Обработка больших данных с SQL/SQL-подобными запросами

- Средства разработки и развертывания виртуальных машин в облаке
- Развертывание кластера Cloudera под большие данные в облаке
- Основные команды HDFS, работа с API
- Использование парадигмы MapReduce для обработки данных. Особенности Hadoop Streaming
- Типы данных. Объединение данных с MapReduce
- Разработка приложений с использованием Tez и Oozie
- Развертывание Spark кластера. Настройка интерактивной среды разработки

РК1

- Spark. Основные операции над RDD
- Spark. Основные операции над Dataframe
- Spark. Взаимодействие с HDFS, Parquet, Avro
- Поточковая обработка. Разработка Storm приложений
- Поточковая обработка. Разработка приложений под Spark Streaming
- Обработка графов. Разработка приложений под Spark GraphX и GraphFrame
- NoSQL СУБД. Базовые операции манипулирования данными в HBase
- Обработка данных с SQL/SQL-подобными запросами. Базовые операции ETL для Hive, Pig

- MapReduce
- Spark + Kafka
- Spark + MLlib
- Spark GraphFrame

cloudera



Java, Scala, Python

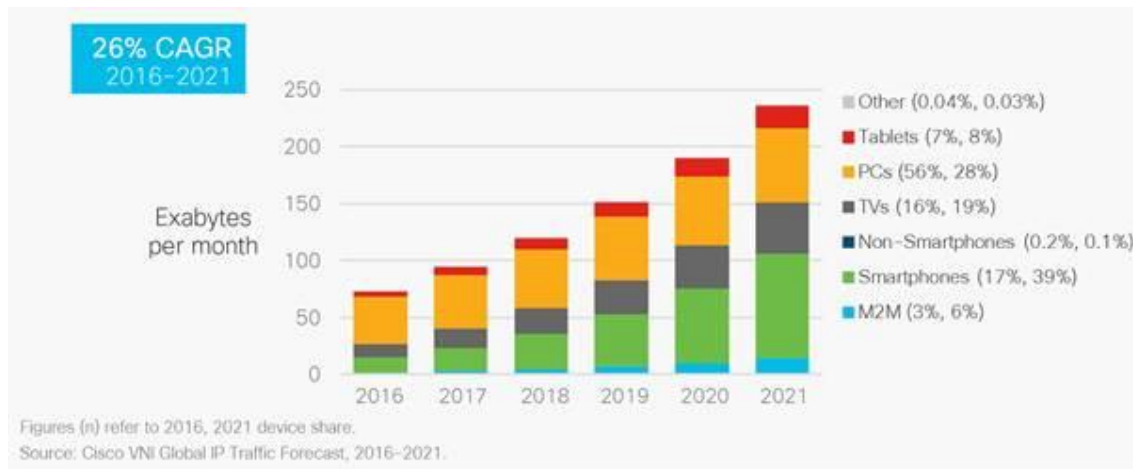
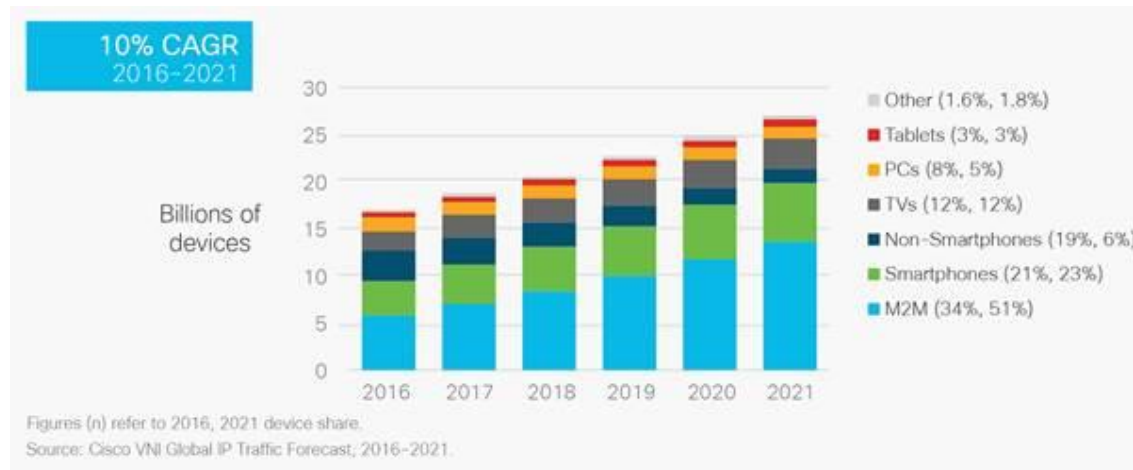
<https://github.com/bigdataprocsystems>

Лекция 1. Концепция Больших Данных



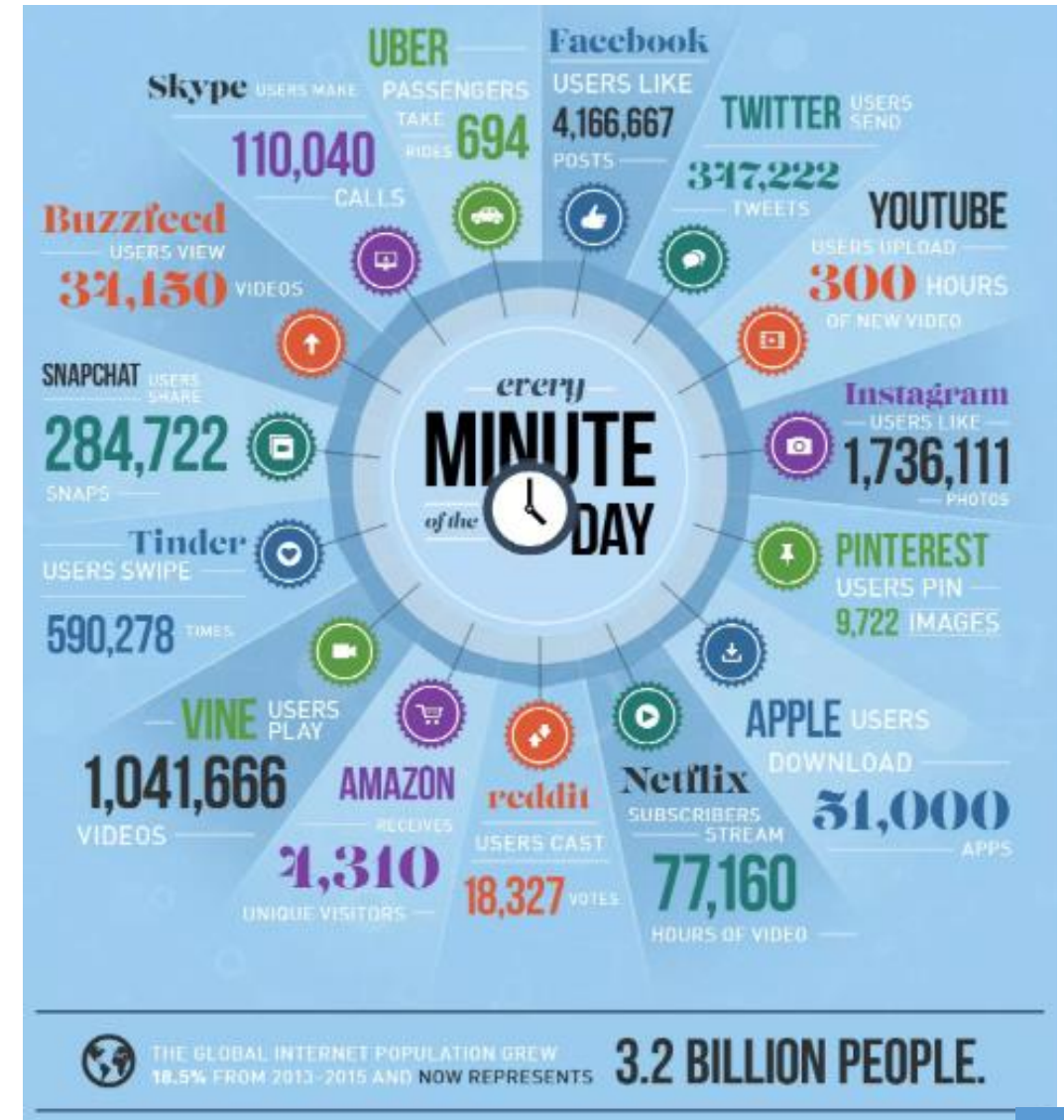
- Большие данные – 4V
- Параллельные и распределенные вычисления
- Системы обработки и хранения больших данных
- стек технологий
- Облачные ресурсы
- Примеры использования

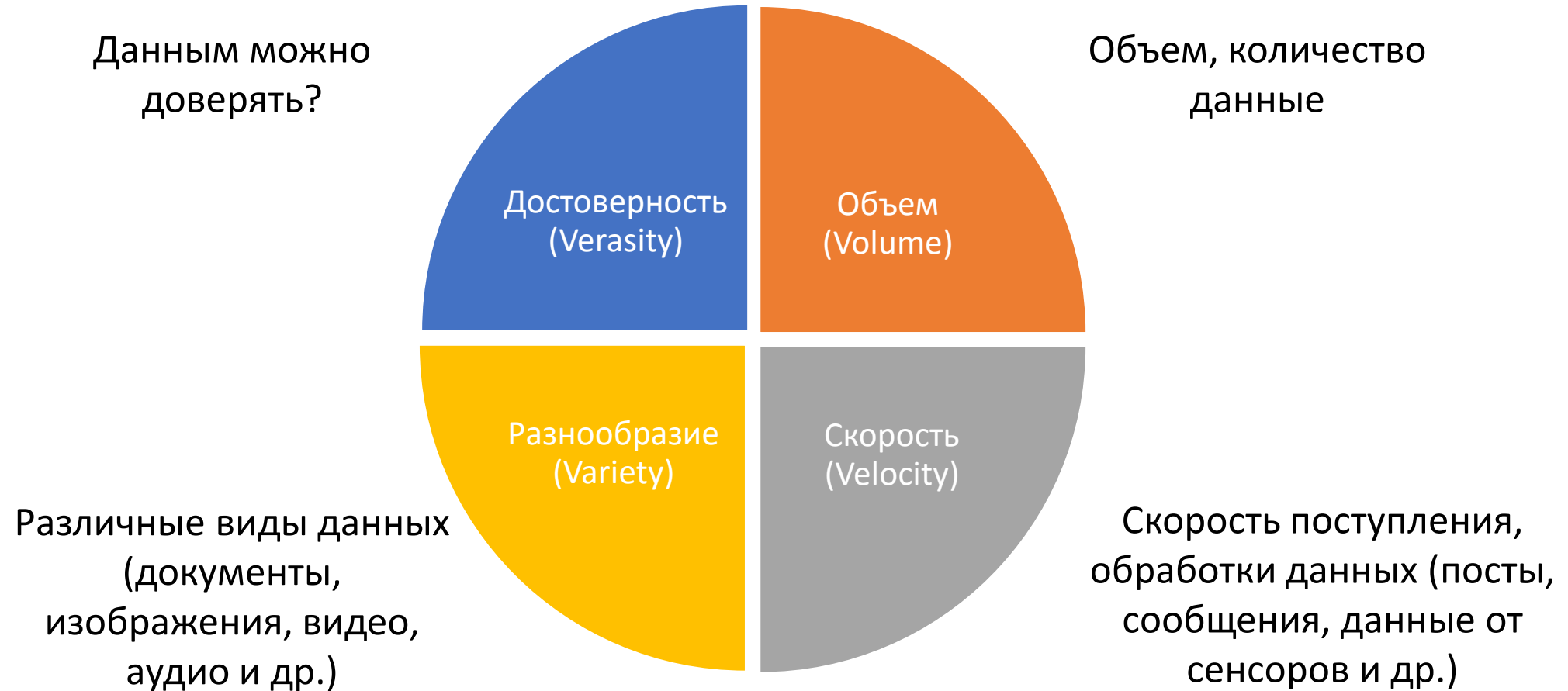
Большие Данные



Пример, Facebook

2.23 млрд. активных пользователей в месяц (2018)
90,032 постов в день (2018)





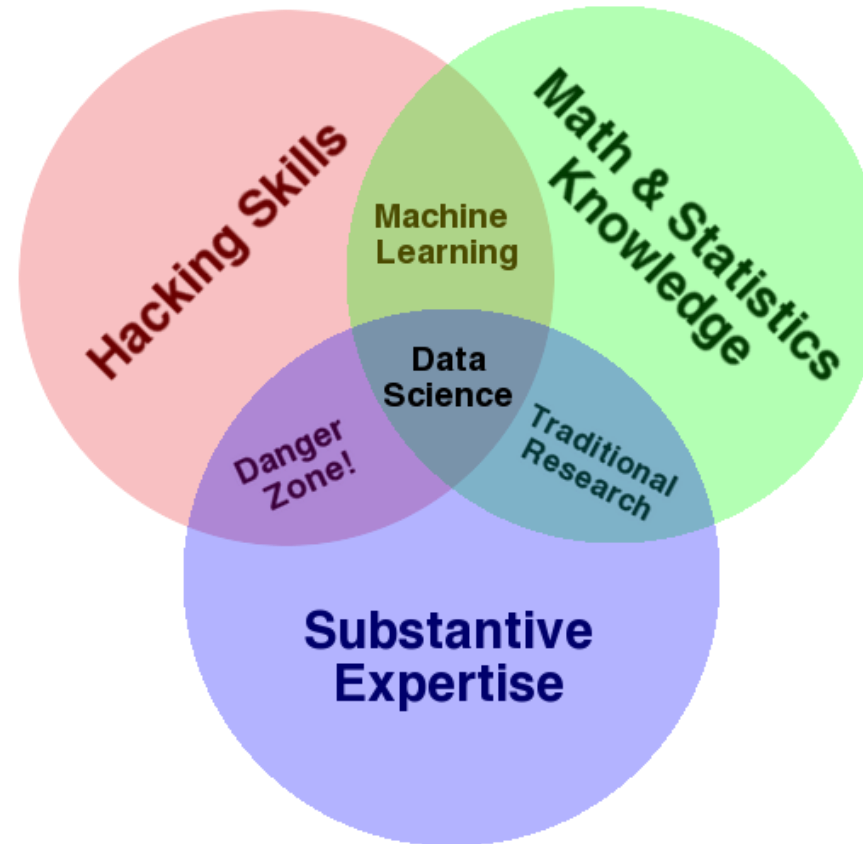
Чем больше данных у нас есть, тем больше знаний мы можем извлечь, лучшее решение можем принять

Чем быстрее обрабатываются поступающие данные, тем быстрее можно начать анализ

Чем более разнообразные источники данных (социальные сети, история просмотров, покупок и пр.), тем лучше можно составить портрет клиента

Чем более достоверные данные, тем точнее можно составить портрет клиента

Наука о данных (Data Science)



Анализ данных

Источники данных



Источники данных

Публичные данные

Экономические
Перепись
Гео-информация
Погода
Открытые данные

Коммерческие данные

Бизнес-информация
Исследования рынка
Кредитное бюро

Социальные сети

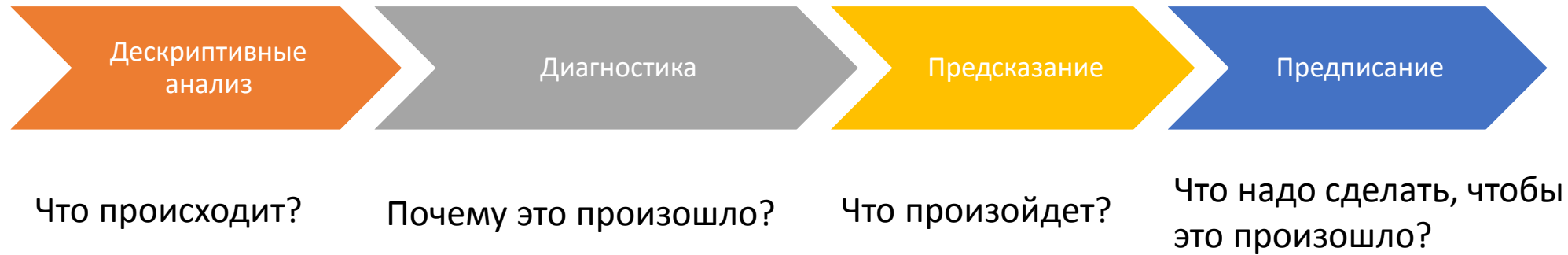
Сообщества
Блоги
Twitter, Facebook, LinkedIn, Tumblr

Операционные данные

Сенсоры
GPS
Транзакции

Корпоративные данные

Взаимодействия с клиентами
Отчеты
Логи
Контакты



Регрессия

Информационный
поиск

Классификация

Дескриптивный
анализ

Кластеризация

Выбор модели

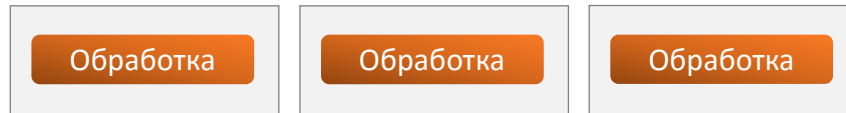
Уменьшение
размерности

Выбор признаков

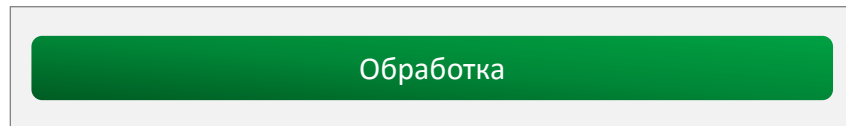
Архитектура систем обработки больших данных

Общие данным

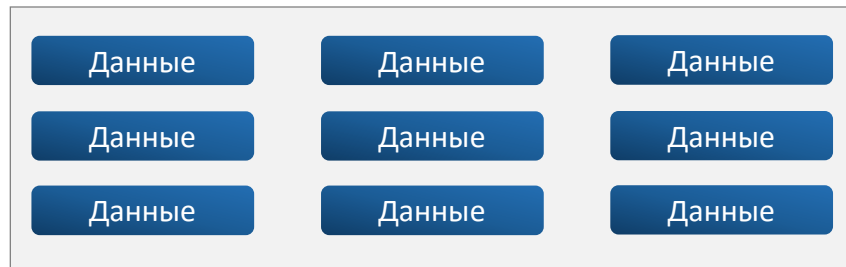
Приложения



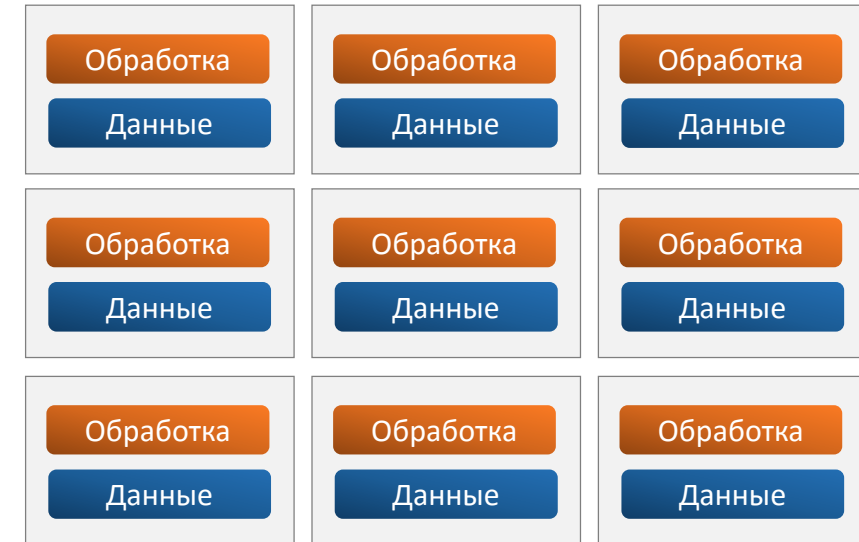
СУБД



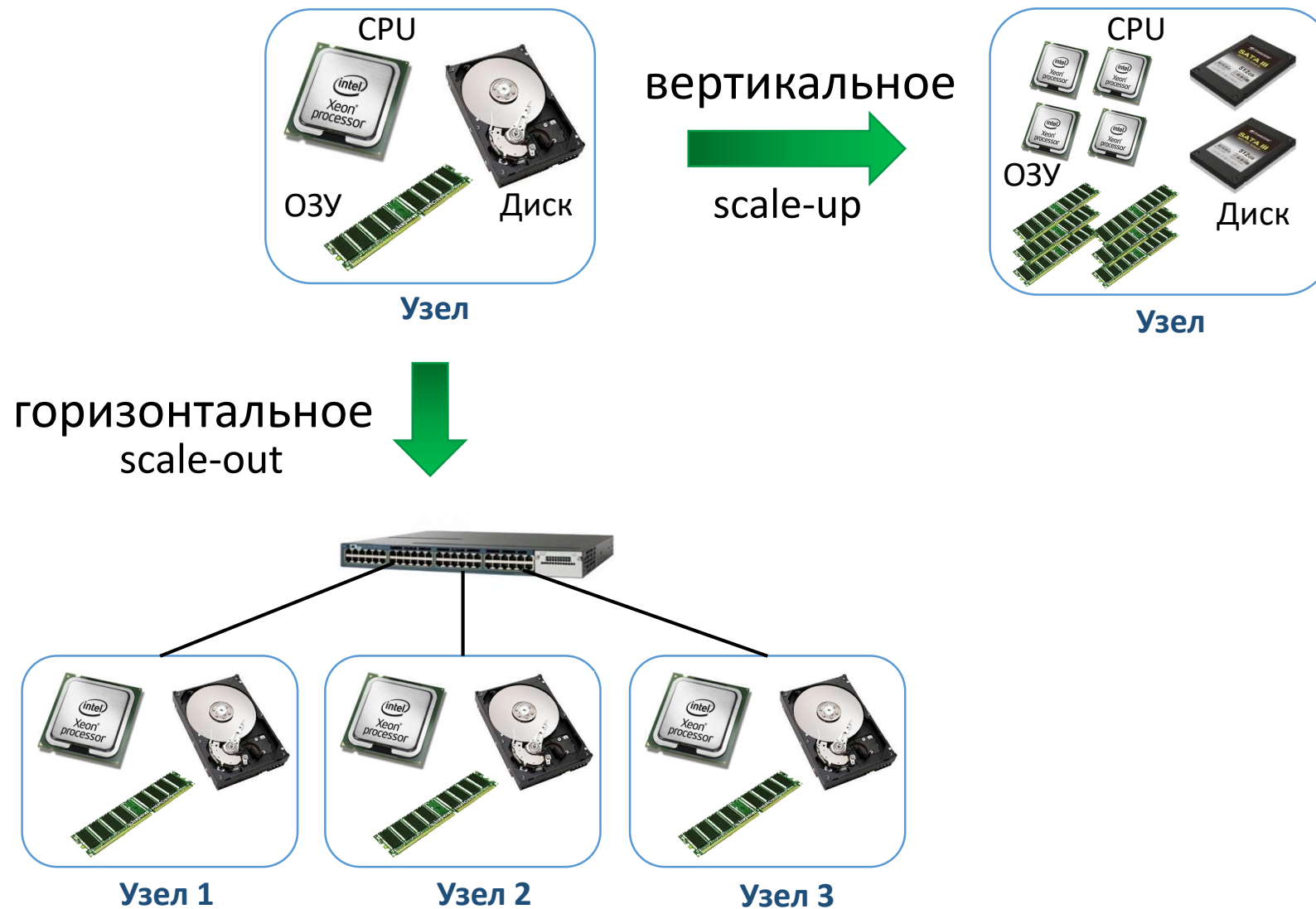
SAN/NAS



Данные обрабатываются там же, где они хранятся



Наращивание производительности



Кластер

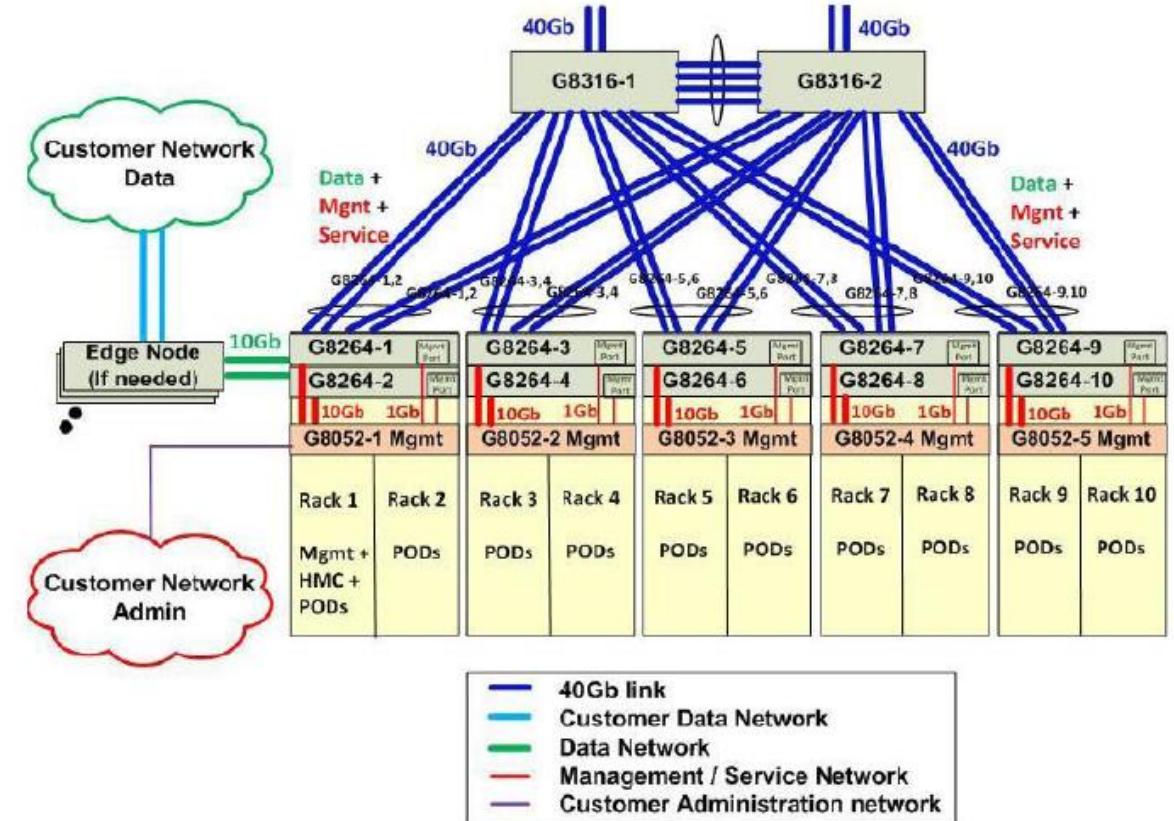
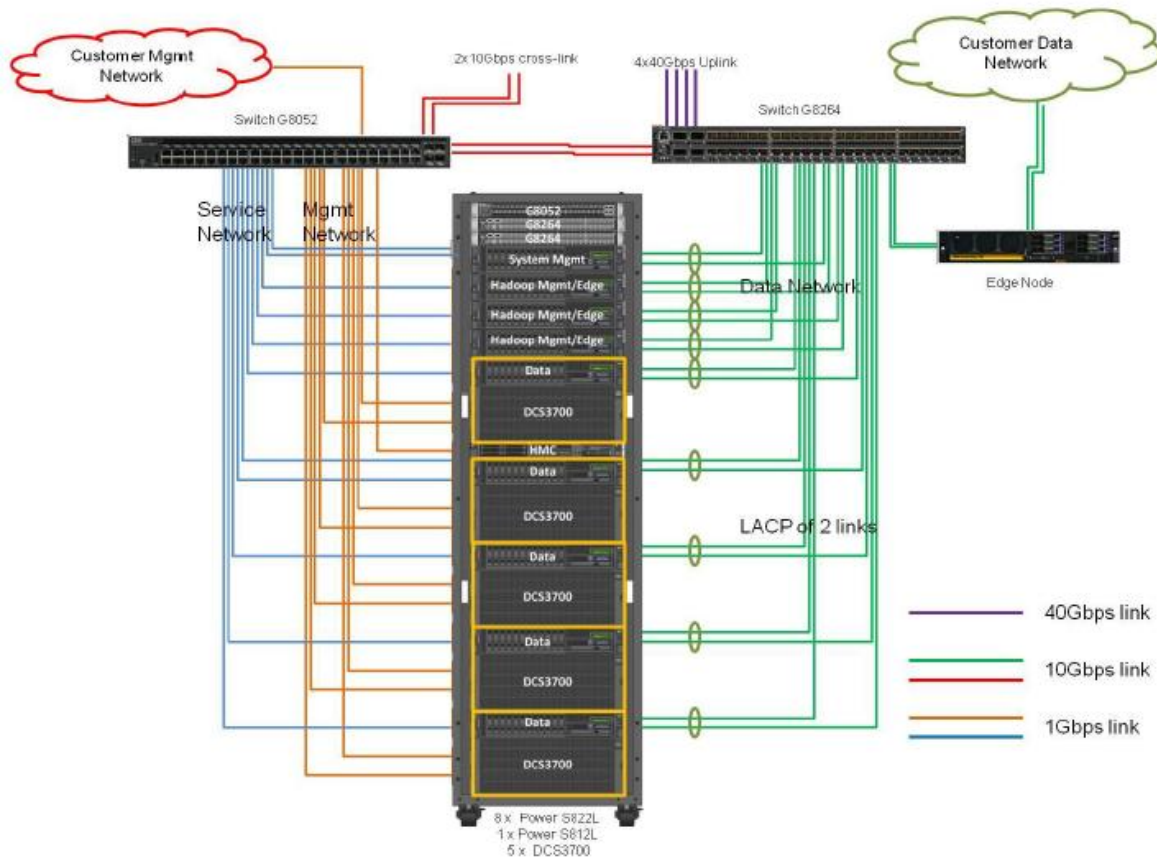
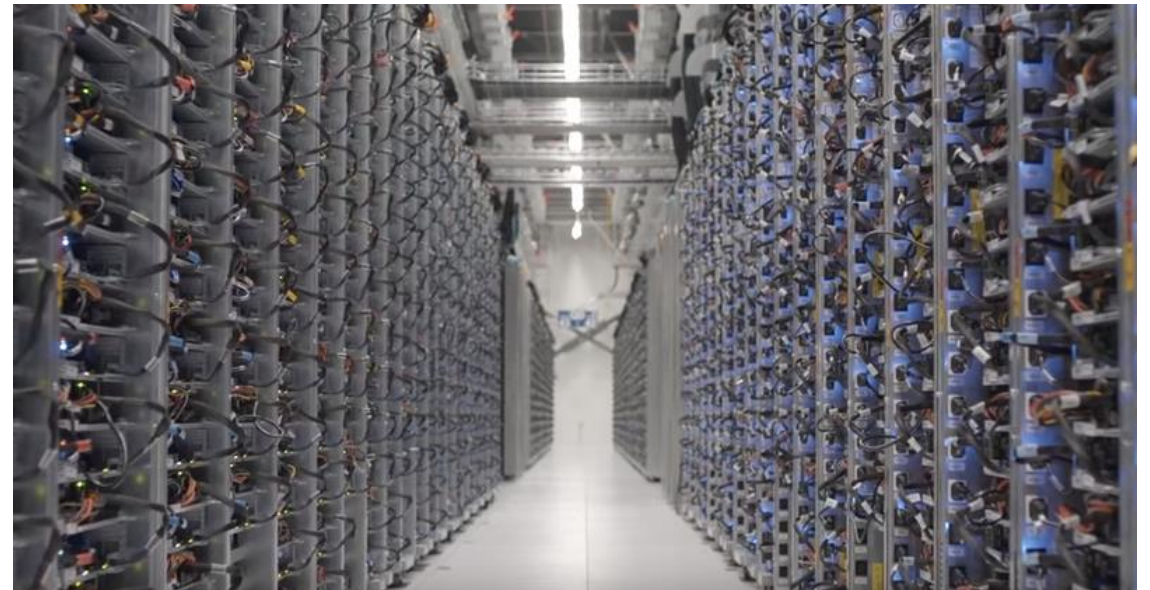
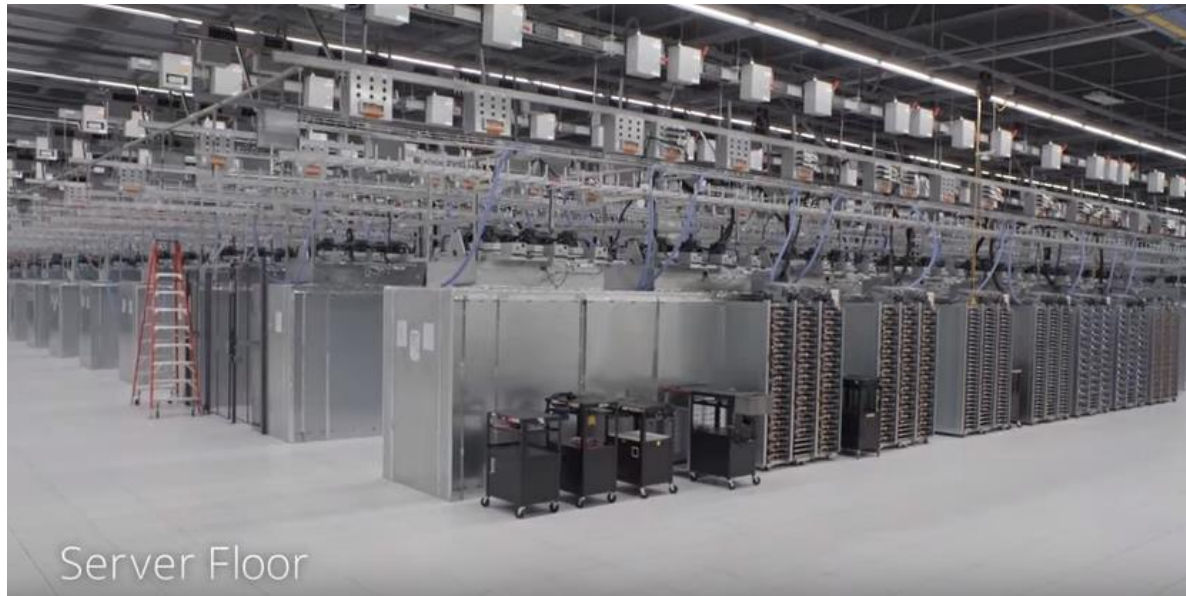
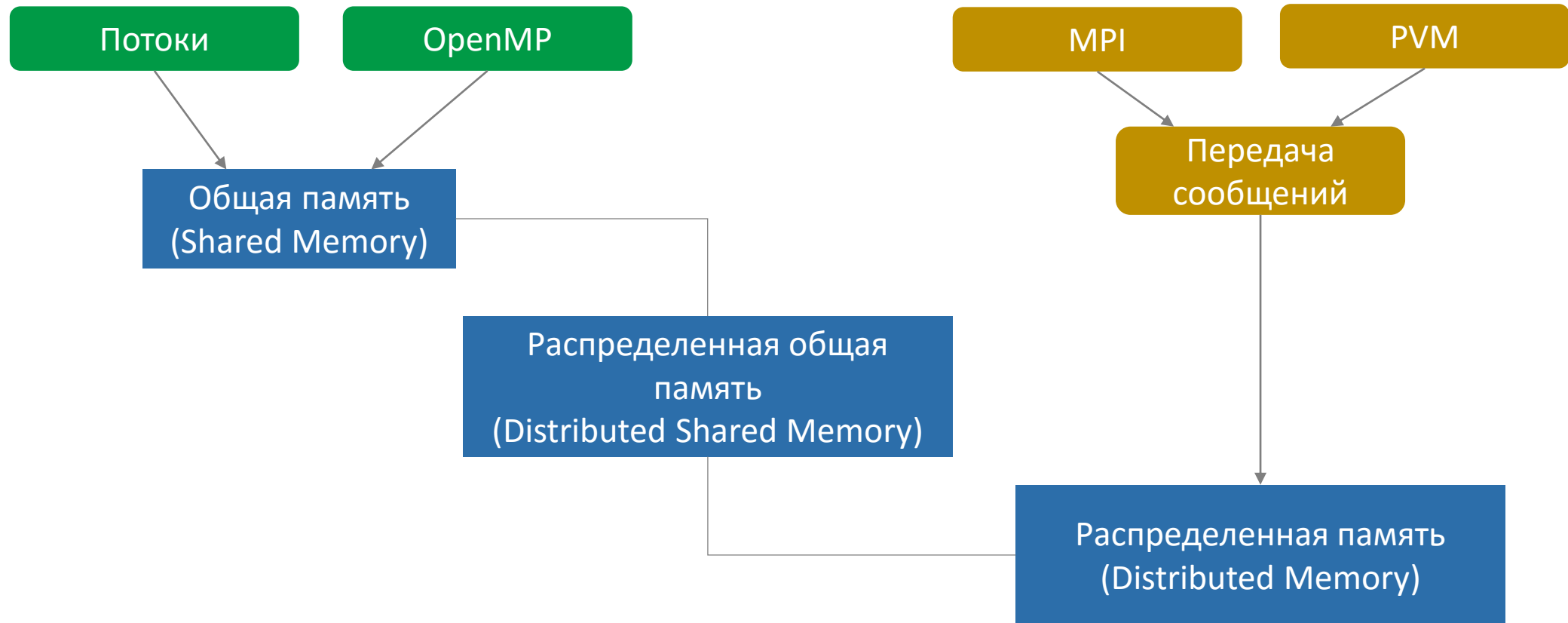


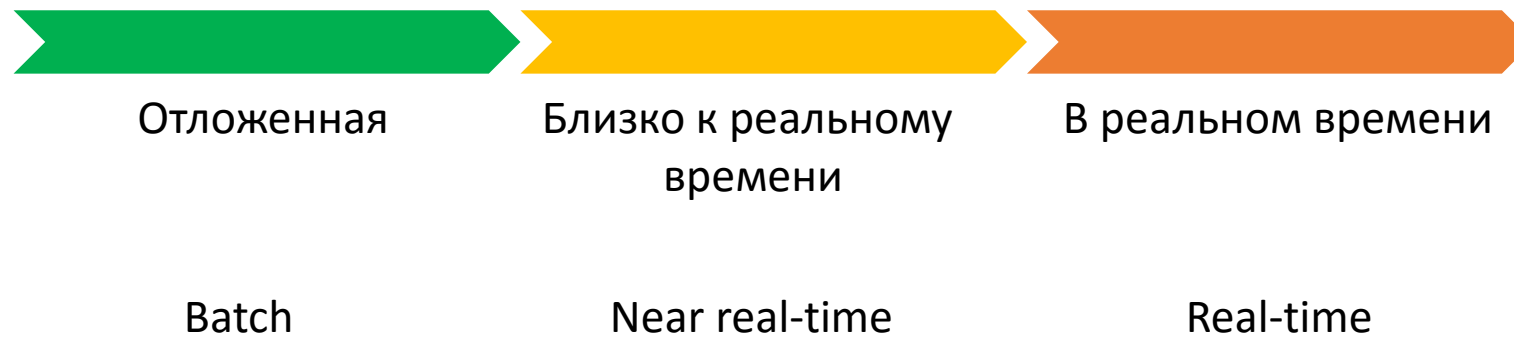
Figure 6: Cross-rack networking



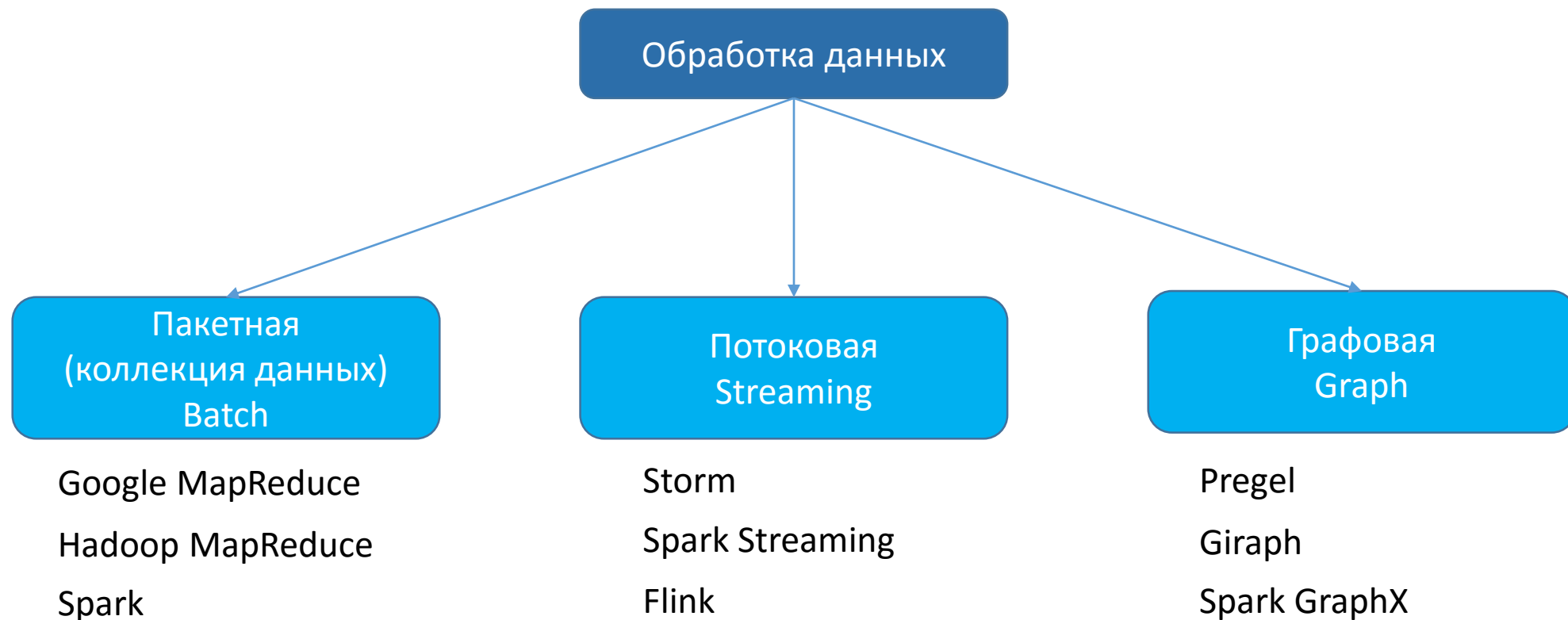
Параллельные и распределенные вычисления



Системы обработки и хранения больших данных



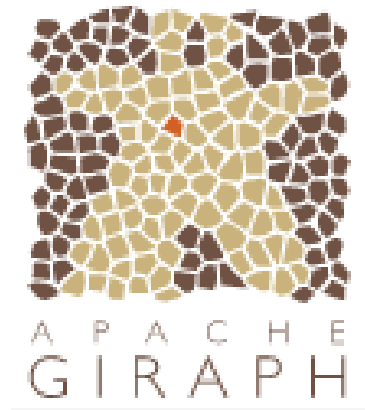
Принципы выполнения вычислений

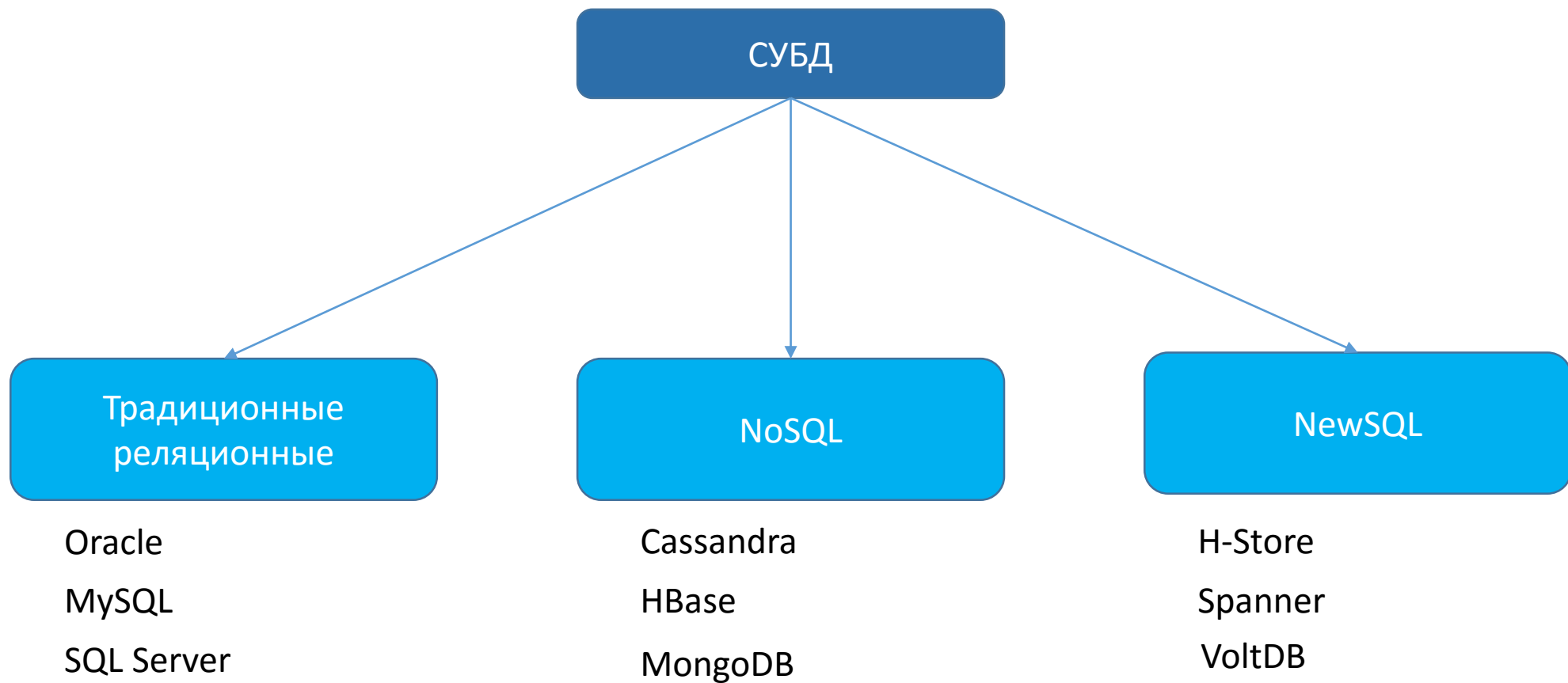


Обработка коллекций данных (batch processing)









ACID

Атомарность (Atomicity)

Согласованность (Consistency)

Изолированность (Isolation)

Долговечность (Durability)

BASE

Basic Availability

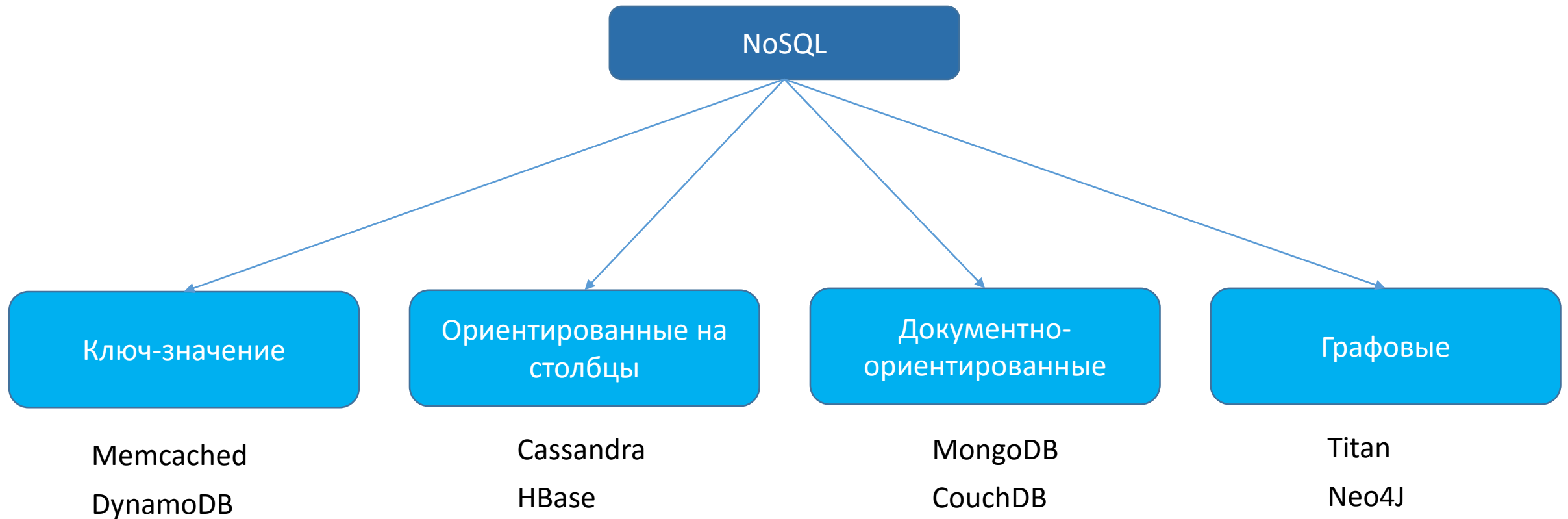
Soft-state

Eventual consistency

CAP



Классификация NoSQL СУБД



Примеры СУБД

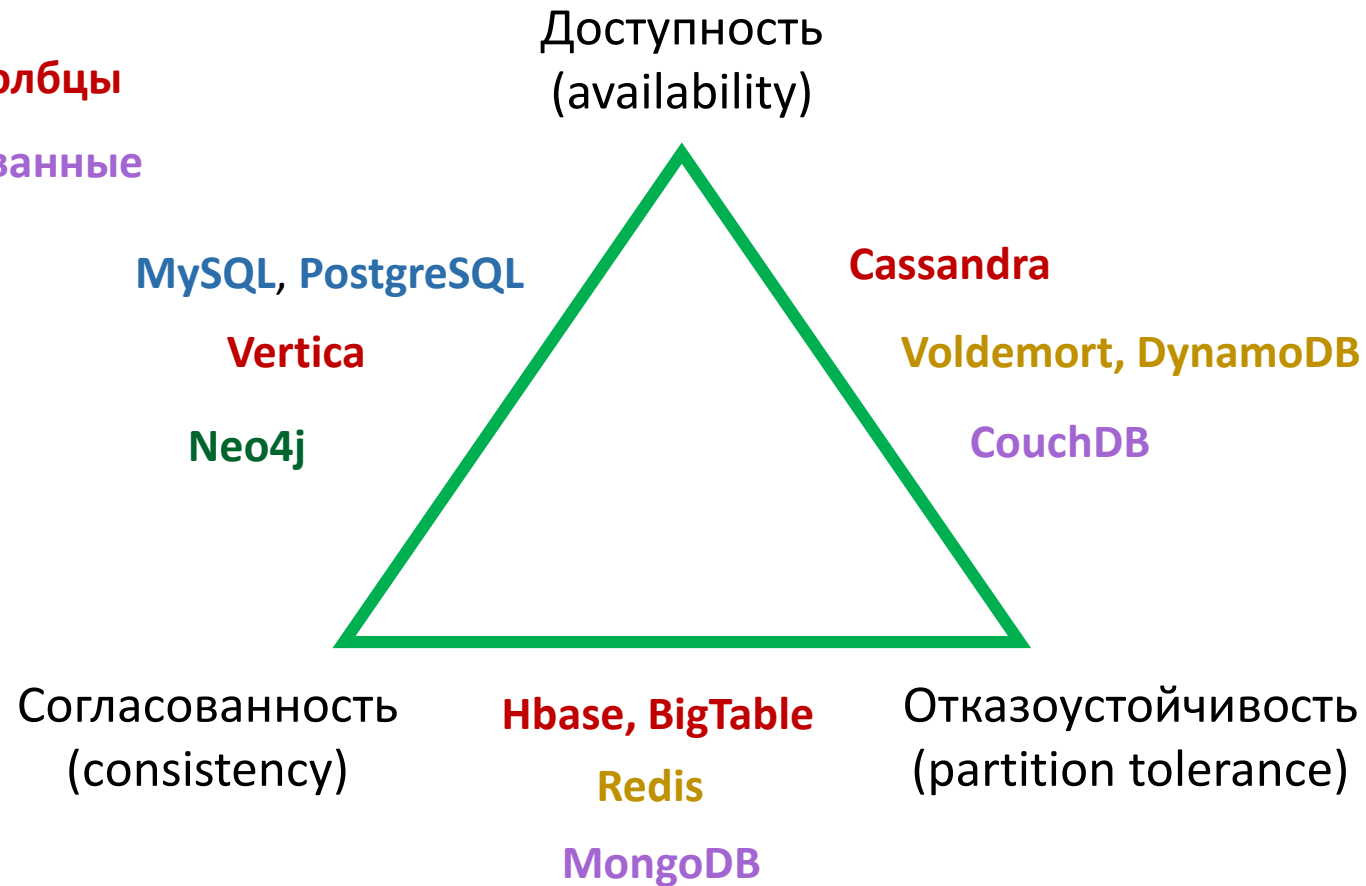
Реляционные

Ключ-значение

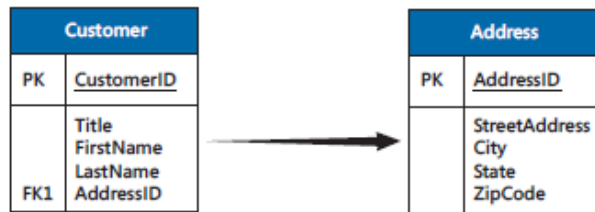
Ориентированные на столбцы

Документно-ориентированные

Графовые



Примеры СУБД



Customer Table

CustomerID	Title	FirstName	LastName	AddressID
1	Mr	Mark	Hanson	500
2	Ms	Lisa	Andrews	501
3	Mr	Walter	Harp	500

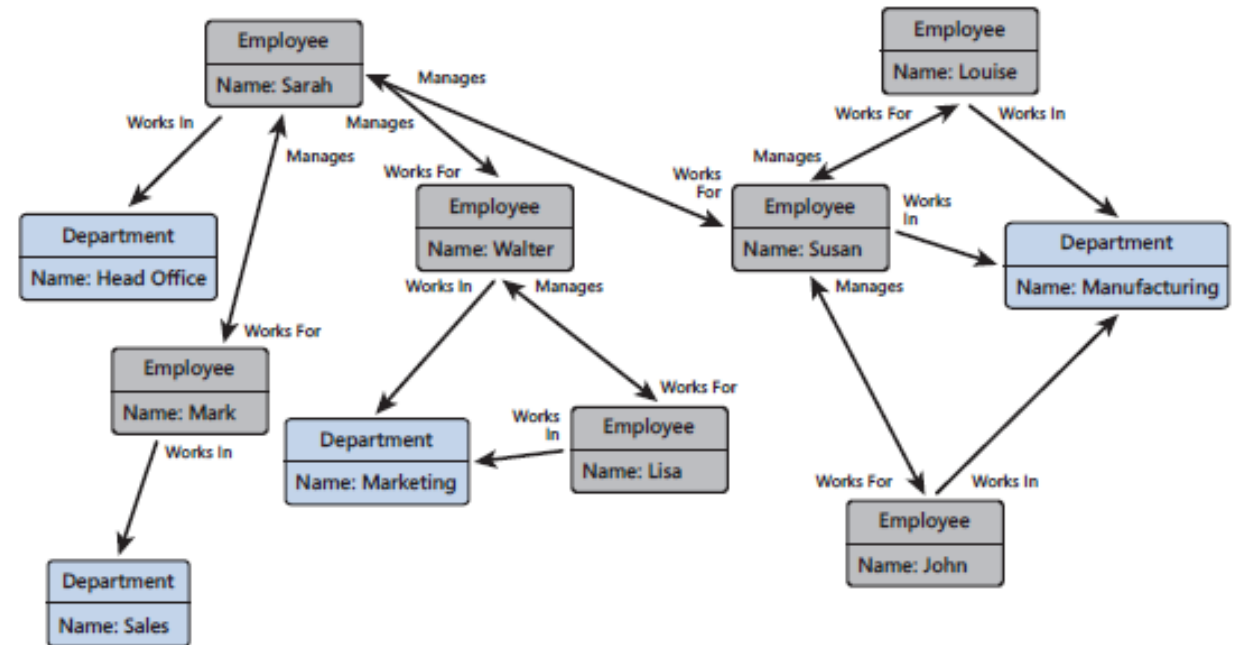
Address Table

AddressID	StreetAddress	City	State	ZipCode
500	999 500th Ave	Bellevue	WA	12345
501	888 W. Front St	Boise	ID	54321

Key	Value (blob)
AAAAA	110100100100111101001001001
AABAB	000110100111100100011110010
DFA766	01011001100100110011111001011
FABCC4	1111000011001010010110011001

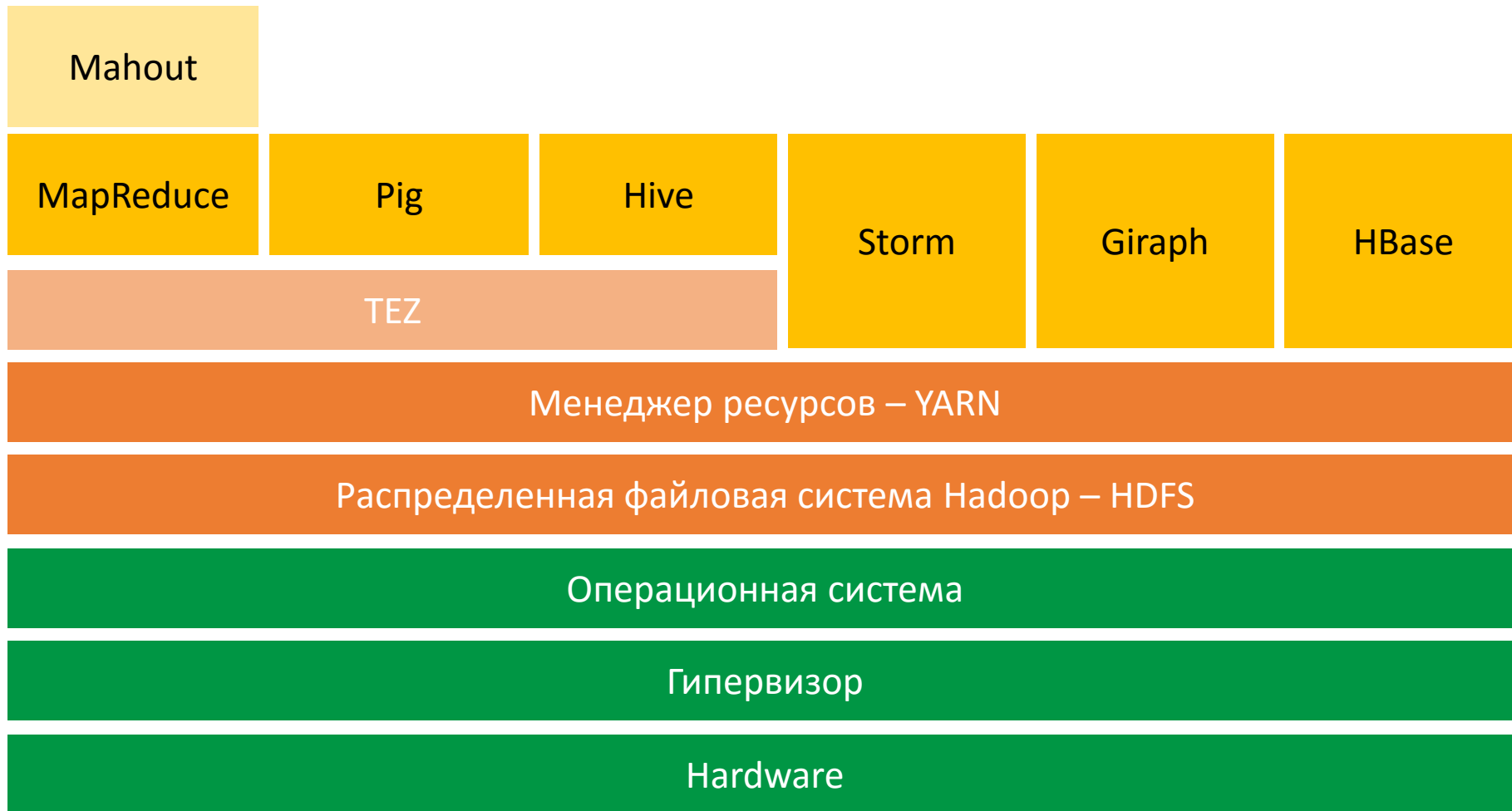
Row Key	Column Families			
CustomerID	CustomerInfo		AddressInfo	
1	CustomerInfo:Title	Mr	AddressInfo:StreetAddress	999 Thames St
	CustomerInfo:FirstName	Mark	AddressInfo:City	Reading
	CustomerInfo:LastName	Hanson	AddressInfo:County	Berkshire
			AddressInfo:PostCode	RG99 922
2	CustomerInfo:Title	Ms	AddressInfo:StreetAddress	888 W. Front St
	CustomerInfo:FirstName	Lisa	AddressInfo:City	Boise
	CustomerInfo:LastName	Andrews	AddressInfo:State	ID
			AddressInfo:ZipCode	54321
3	CustomerInfo:Title	Mr	AddressInfo:StreetAddress	999 500th Ave
	CustomerInfo:FirstName	Walter	AddressInfo:City	Bellevue
	CustomerInfo:LastName	Harp	AddressInfo:State	WA
			AddressInfo:ZipCode	12345

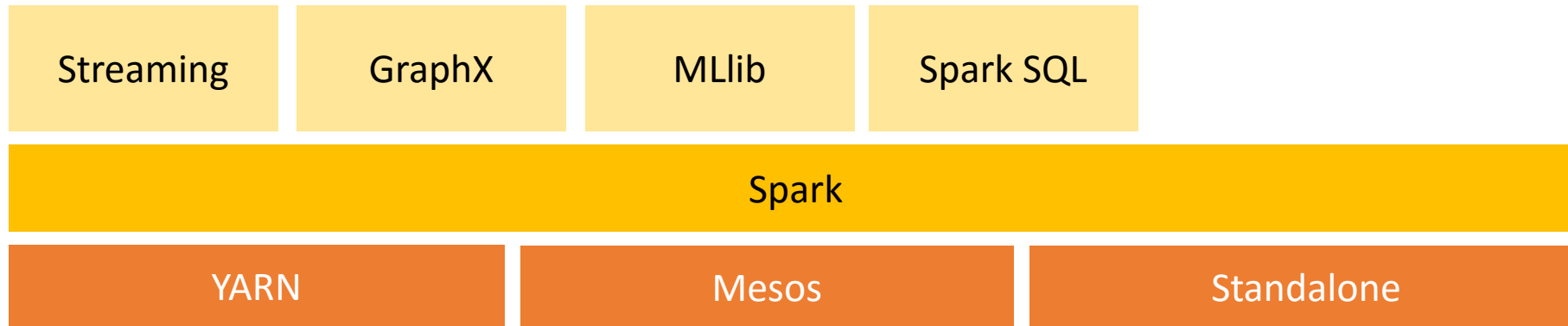
Row Key	Document
1001	<p>OrderDate: 06/06/2013 OrderItems: ProductID: 2010 Quantity: 2 Cost: 520</p> <p>ProductID: 4365 Quantity: 1 Cost: 18</p> <p>OrderTotal: 1058 Customer ID: 99 ShippingAddress: StreetAddress: 999 500th Ave City: Bellevue State: WA ZipCode: 12345</p>
1002	<p>OrderDate: 07/07/2013 OrderItems: ProductID: 1285 Quantity: 1 Cost: 120</p> <p>OrderTotal: 120 Customer ID: 220 ShippingAddress: StreetAddress: 888 W. Front St City: Boise State: ID ZipCode: 54321</p>



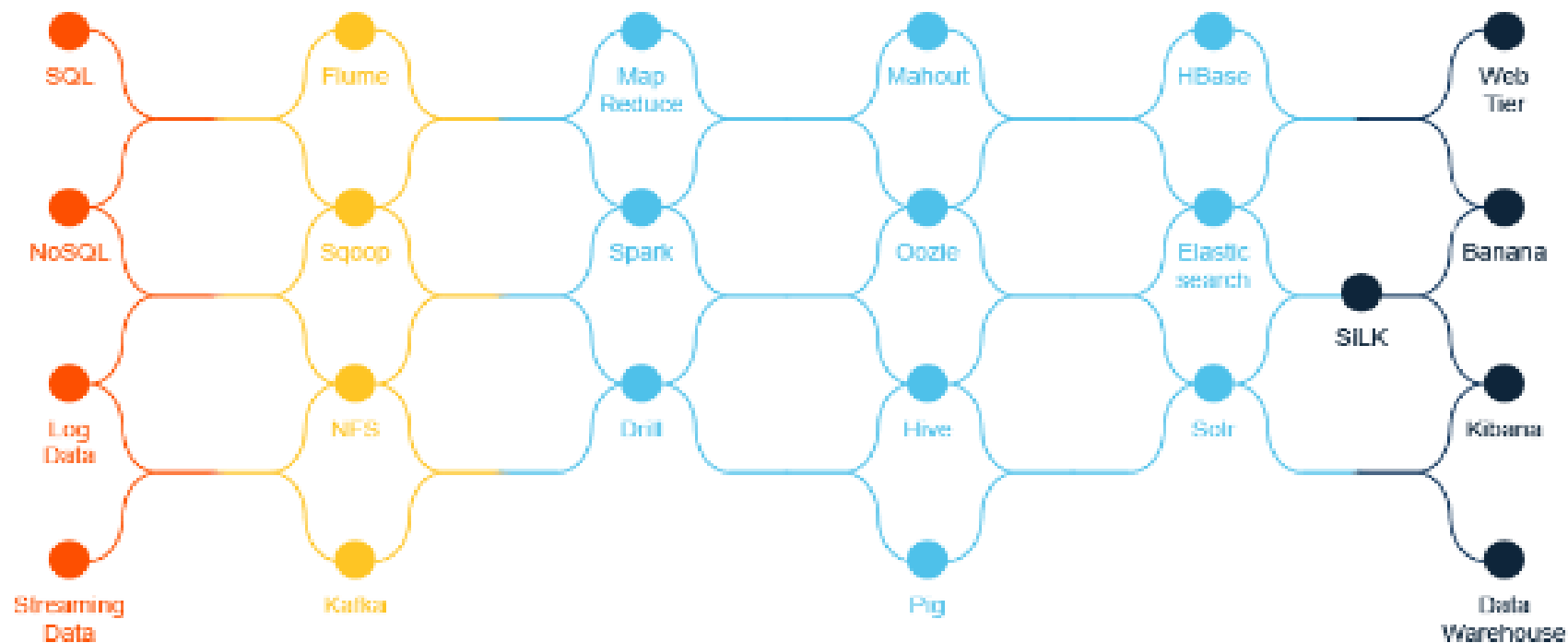
Стек технологий

Стек Hadoop



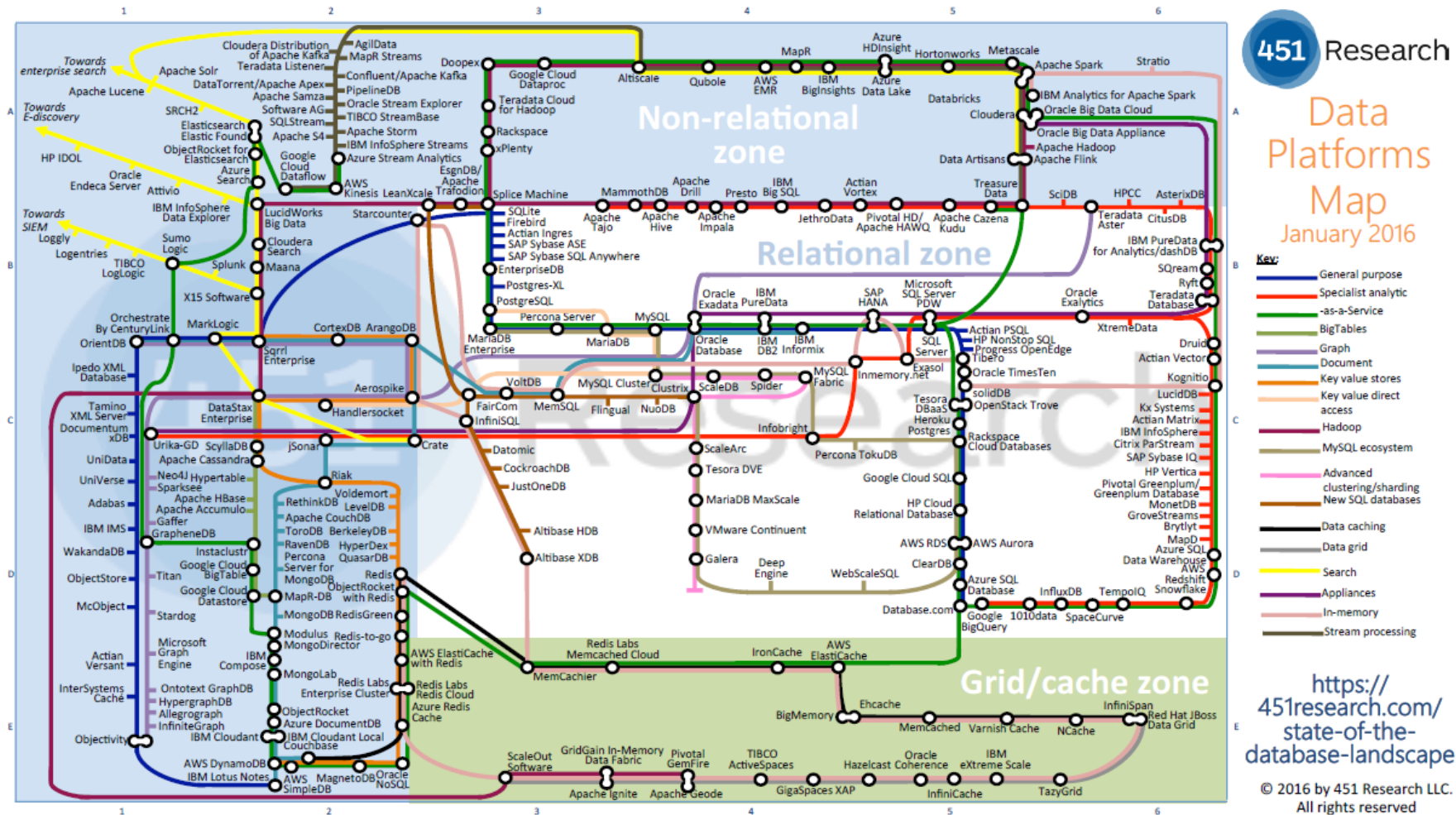


Стек технологий



MapR

Карта технологий



Облачные ресурсы

Облачные ресурсы





Инфраструктуры анализа больших данных

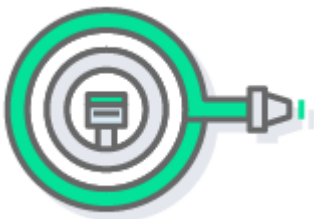


Hadoop и Spark
Amazon EMR



Elasticsearch
Сервис Amazon Elasticsearch

Анализ больших данных в режиме реального времени



Amazon Kinesis Firehose



Amazon Kinesis Streams



Amazon Kinesis Analytics

Хранилища и базы больших данных



Объектное хранилище
Amazon S3



NoSQL
Amazon DynamoDB



HBase в Amazon EMR



Реляционные базы данных
Amazon RDS



Графовые базы данных
Amazon DynamoDB для БД Titan