# Технологии обработки больших данных

**«**38.04.05 – Бизнес-информатика направленность интеллектуальное управление цифровым предприятием**»** 

http//vikchas.ru

https://www.famous-scientists.ru/3653/

# Лекция 2 «Архитектура системы обработки Больших данных»

# Часовских Виктор Петрович

д-р техн. наук, профессор кафедры ШИиКМ

ФГБОУ ВО «Уральский государственный экономический университет»

Екатеринбург 2022

Для работы с Большими данными используются сложные системы, в которых можно выделить несколько компонентов или слоёв (Layers).

Обычно выделяют четыре уровня компонентов таких систем:

- > прием,
- > сбор,
- > анализ данных
- > представление результатов

Это деление является в значительной мере условным так как, с одной стороны, каждый компонент в свою очередь может быть разделен на подкомпоненты, а с другой некоторые функции компонентов могут перераспределяться в зависимости от решаемой задачи и используемого программного обеспечения, например, выделяют хранение данных в отдельный слой.

Для работы с Большими данными разработчиками систем создаются модели данных, содержательно связанные с реальным миром. Разработка адекватных моделей данных представляет собой сложную аналитическую задачу, выполняемую системными архитекторами и аналитиками.

Модель данных позволяет создать математическую взаимодействий объектов реального мира и включает в себя описание структуры данных, методы манипуляции данными и аспекты сохранения целостности данных.

Для хранения данных используются распределенные системы различных типов. Это могут быть файловые системы, базы данных, журналы, механизмы доступа к общей виртуальной памяти.

Большинство систем хранения ориентированы исключительно на работу с Большими данными, они имеют крайне ограниченное число (например, может отсутствовать возможность модификации, но и удаления поступивших данных) что объясняется внутренней сложностью создания высокоэффективных распределенных систем.

# Прием данных (Data Ingestion)

Источники данных имеют различные параметры, такие как частоту поступления данных из источника, объём порции данных, скорость передачи данных, тип поступающих данных и их достоверность.

Для эффективного сбора данных необходимо установить источники данных.

Это могут быть хранилища данных, поставщики агрегированных данных, АРІ каких-либо датчиков, системные журналы, сгенерированный человеком контент в социальных сетях, в корпоративных информационных системах, геофизическая информация, научная информация, унаследованные данные из других систем. Источники данных определяют исходный формат данных

Например, мы можем самостоятельно проводить погодные на территории аэропорта, использовать данные, поступающие с взлетающих и садящихся самолетов, закупить данные со спутников, пролетающих над аэропортом и у местной метеослужбы, а также найти их где-то в сети в другом месте.

В общем случае для каждого источника необходимо создавать собственный сборщик (Data Crawler для сбора информации в сети и Data Acquisition для проведения измерений).

Прием данных заключается в начальной подготовке данных от источников с целью приведения данных к общему формату представления данных.

Этот единый формат выбирается в соответствии с принятой моделью данных. Выполняются преобразования систем измерения, типов (типизация), верификация. Обработка данных содержательно не затрагивает имеющуюся в данных информацию, но может изменять ее представление (например, приводить координаты к единой системе координат, а значения к единой размерности).

# Сбор данных (Data Staging)

Этап сбора данных характеризуется непосредственным взаимодействием с системами хранения данных.

Устанавливается точка сбора, в которой собранные данные снабжаются локальными метаданными и помещаются в хранилище либо передаются для последующей обработки.

Данные, по каким-либо причинам не прошедшие точку сбора, игнорируются. Для структурированных данных проводится преобразование из исходного формата по заранее заданным алгоритмам. Это наиболее эффективная процедура в случае, если структура данных известна. Однако если данные представлены в двоичном виде, структура и связи между данными утеряны, то разработка алгоритмов и основанного на них программного обеспечения для обработки данных может оказаться крайне затруднительной

Для полуструктурированных данных требуется интерпретация поступающих данных и использование программного обеспечения, умеющего работать с используемым языком описания данных.

Существенным плюсом полуструктурированных данных является то, что в них зачастую содержатся не только сами данные, но метаданные в виде информации о связях между данными и способах их получения.

Разработка программного обеспечения для обработки полуструктурированных данных представляет собой достаточно сложную задачу. Однако имеется значительное количество готовых конвертеров, которые могут, например, извлечь данные из формата XML в сформированное табличное представление. Наибольшего объема работ требует обработка неструктурированных данных. Для их перевода к заданному формату может потребоваться создание специального ПО, сложная ручная обработка, распознавание и выборочный ручной контроль

На этапе сбора проводится контроль типов данных и может выполняться базовый контроль достоверности данных.

Например, координаты молекул газа, содержащихся в какой-либо области, не могут лежать за пределами этой области, а скорости — существенно превышать скорость звука.

Для того, чтобы избежать ошибок типизации, необходимо проверять правильно ли заданы единицы измерения.

Например, в одном наборе данных высота может измеряться в километрах, а в другом — в футах. В этом случае необходимо произвести преобразование высоты в те единицы измерения, которые приняты в используемой модели.

При сборе данные систематизируются и снабжаются метаданными, хранимыми в связанных метаданных. При наличии большого количества источников данных может потребоваться управление сбором данных для того, чтобы сбалансировать объемы информации, поступающие из различных источников. Собранные данные либо сохраняются в системах хранения, либо (в особенности, для потоковых данных) передаются для анализа в реальном времени

# Анализ данных (Analysis Layer)

Анализ данных, в отличии от сбора данных, использует информацию, содержащуюся в самих данных.

Анализ может проводиться как в реальном времени, так и в пакетном режиме.

Анализ данных составляет основную по трудоемкости задачу при работе с Большими данными.

Существует множество методик обработки данных: предиктивный анализ, запросы и отчетность, реконструкция по математической модели, трансляция, аналитическая обработка и другие.

Методики используют специфические алгоритмы в зависимости от поставленных целей.

Например, аналитическая обработка может являться анализом изображений, социальных сетей, географического местоположения, распознавания по признакам, текстовым анализом, статистической обработкой, анализом голоса, транскрибированием.

Алгоритмы анализа данных также, как и алгоритмы обработки данных, опираются на модель данных.

При этом при анализе может быть использовано несколько моделей, задающих общий формат данных, но по-разному моделирующие содержательные процессы, данные о которых мы обрабатываем.

При использовании при анализе методов искусственного интеллекта, в частности нейронных сетей, производится динамическое обучение моделей на различных наборах данных.

При анализе данных производится идентификация сущностей, описываемых данными на основании имеющейся в данных информации и используемых моделей.

Сущностью анализа является аналитических механизм, использующий аналитические алгоритмы, управление моделями и идентификацию сущностей для получения новой содержательной информации, являющейся результатом анализа. Для анализа данных также используются методы искусственного интеллекта на нейронных сетях.

# Представление результатов (Consumption Layer)

Результаты анализа данных предоставляются на уровне потребления. Имеется несколько механизмов, позволяющих использовать результаты анализа больших данных.

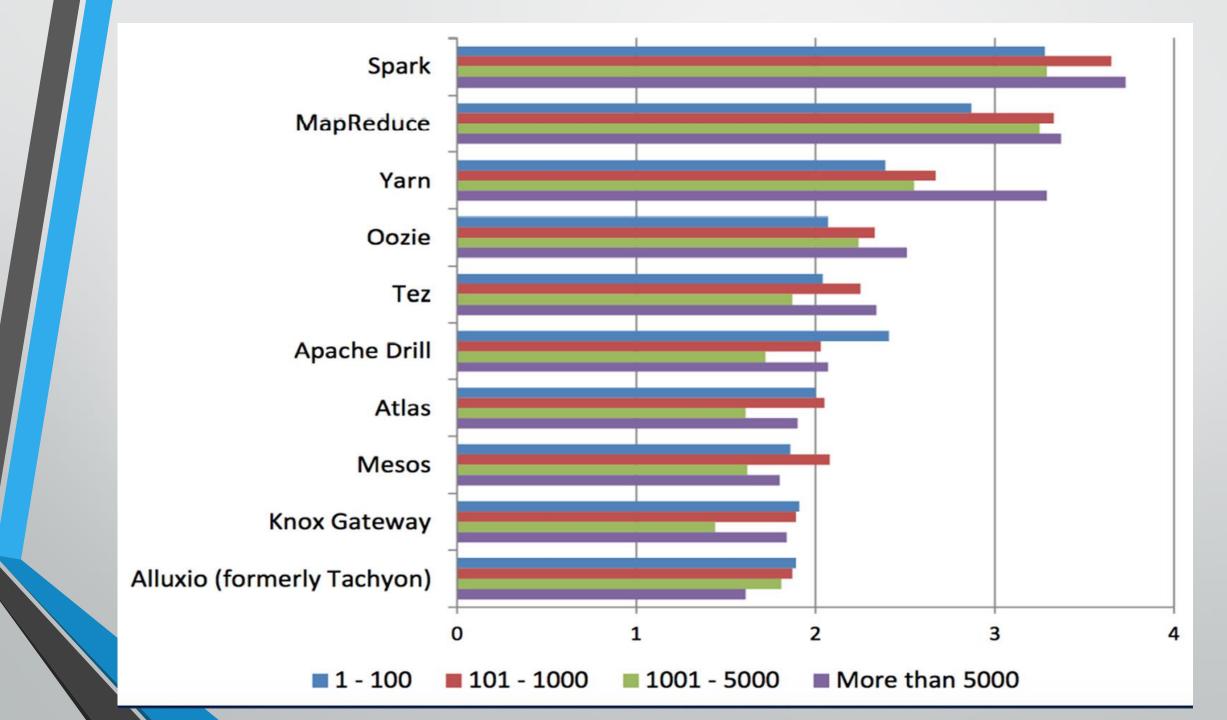
- ▶ Мониторинг метаинформации. Подсистема отображения в реальном времени существенных параметров работы системы, загруженности вычислителей, распределение задач в кластере, распределение информации в хранилищах, наличие свободного места в хранилищах, поступление данных от источников, активности пользователей, отказов оборудования и т.д.
- Мониторинг данных. Подсистема отображения в реальном времени процессов приема, сбора и анализа данных, навигация по данным.
- > Генерация отчетов, запросы к данным, представление данных в виде визуализации на дэшбордах(*информационная панель*) (Dashboard), в формате PDF, инфографике, сводных таблицах и кратких справках
- Преобразование данных и экспорт в другие системы, интерфейс 
   системами.

# Программные платформы и системы для Больших данных

В настоящее время используется значительное количество платформ и систем Больших данных.

Системы обработки больших данных являются фреймворками, то есть каркасами, для использования которых необходимо состыковать их с другими фреймворками, прикладным программным обеспечением пользователя и системой хранения данных.

В аналитическом отчете Big Data Analytics Market Study Edition приводится следующая диаграмма инфраструктур Больших данных, внедренных на предприятиях, представленная в разрезе размеров предприятий.



# Системы управления потоками данных

Flume Система управления потоками данных.

https://flume.apache.org/

Разработана в 2017 году.

Apache Kafka Масштабируемый отказоустойчивый журнал коммитов.

https://kafka.apache.org/

Разработан в Linked In в 2011 году.

Niagara Files (NiFi) Система управления потоками данных

https://nifi.apache.org/

Разработан в NSA в 2014.

### Системы хранения Больших данных

**HDFS** (Hadoop Distributed File System)
https://hadoop.apache.org/docs/r1.2.1/hdfs\_design.html
Файловая система, входящая в проект Hadoop.

#### **OpenStack Swift**

https://docs.openstack.org/swift/latest/ Платформа хранения, входящая в проект OpenStack.

#### Cassandra

http://cassandra.apache.org/ Табличная СУБД, написана на языке Java.

#### **HBase**

https://hbase.apache.org/ Табличная СУБД, написана на языке Java.

#### **Apache Drill**

https://drill.apache.org/

SQL-интерфейс к NoSQL базам данным (HBase, MongoDB, MapR-DB, HDFS, MapR-FS, Amazon S3, Azure Blob Storage, Google Cloud Storage, Swift, NAS and local files).

### Платформы Больших данных

### **Hadoop**

http://hadoop.apache.org/

Предоставляет интерфейс к Java (а через коннекторы и к другим языкам), свободно распространяется под лицензиями Apache License 2.0 и GNU GPL пакет программного обеспечения, состоящий из управляющего модуля Hadoop Common, распределенной файловой системы HDFS, планировщика заданий YARN и вычислительной платформы Hadoop MapReduce. Развивается с 2005 года.

### **Spark**

https://spark.apache.org/

Предоставляет интерфейсы к Scala, Java, Python и R, распространяется под лицензией Apache License 2.0. Вычислительная платформа, развивающаяся с 2014 года

#### Elasticsearch

https://www.elastic.co/products/elasticsearch

Совместно с системой сбора Logstash и платформой аналитики Kibana составляют интегрированную систему сбора, хранения, поиска и аналитики данных.

#### Solr

http://lucene.apache.org/solr/

Еще одна система поиска и анализа в базах Больших данных.

### **Hortonworks Data Platform (HDP)**

https://hortonworks.com/products/data-platforms/hdp/

Платформа управления данными, включающая HDFS, Hadoop, HBase, HCatalog, Pig, Hive, Oozie, Zookeper, Ambari, WebHDFS, TalentOS, Sqoop, Flume и Mahout.

### **Windows Azure HDInsight**

https://azure.microsoft.com/en-ca/services/hdinsight/ Система от Microsoft для развёртывания hadoop на Azure.

# Аналитические платформы

## **RapidMiner**

https://rapidminer.com/

Система прогнозной аналитики, поддерживает глубинный анализ, проверку, оптимизацию и визуализацию, имеет графический интерфейс программирования.

#### **IBM SPSS Modeler**

https://www.ibm.com/products/spss-modeler Коммерческая аналитическая система автоматизированного моделирования,

геопространственной аналитики, анализа текстовой информации. Плохо

подходит для больших объёмов информации.

#### KNIME

https://www.knime.com/

Бесплатная система анализа данных, имеющая глубинный анализ, веб-анализ, обработку изображений, анализ социальных сетей, обработку текстов

# **Qlik Analytics Platform**

https://www.qlik.com/us/products/qlik-analytics-platform Система визуальной аналитики, предоставляет доступ к ассоциативной машине индексации данных QIX Engine.

### **STATISTICA Data Miner**

http://statsoft.ru/products/STATISTICA\_Data\_Miner/data-mining-tools.php

Система от российского производителя

# **IBM Watson Analytics**

https://www.ibm.com/watson-analytics Мощная облачная система от IBM (применяемая в том числе twitter'ом)

# **Dell EMC Analytic Insights Module**

https://www.dellemc.com/en-us/big-data/index.htm Многокомпонентная система от Dell EMC

# **SAP Predictive Analytics**

https://www.sap.com/products/predictive-analytics.html Система от мирового лидера ERP, интегрируется с SAP HANA

# **Oracle Big Data Preparation**

https://cloud.oracle.com/bigdatapreparation Облачное решение от мирового лидера баз данных

# Оборудование для обработки Больших данных

Комплект оборудования для обработки Больших данных монтируется в ЦОД. Основными компонентами системы являются система управления, вычислительные ресурсы, система хранения данных, локальная сеть.

Электропитание, мониторинг, доступ к интернету и другие внешние ресурсы предоставляют центры обработки данных (ЦОД).

Система управления предназначена для общего управления системой, внешнего доступа, обеспечения аутентификации, авторизации и представления результатов пользователям. Выполняется на базе обычной серверной платформы, специфических требований не имеет <sup>21</sup>



` Вычислительные ресурсы кластера состоят из узлов, основными параметрами которых является объем оперативной памяти с контролем четности (ECC) и максимальное количество ядер CPU.

Расчётным параметрами является размер оперативной памяти на ядро и скорость работы процессора.

Высокая отказоустойчивость узлов желательна, но в целом не требуется, некоторое количество отказов является нормальным и компенсируется при помощи программного обеспечения — отказавший узел автоматически выводится из эксплуатации, и его работа перераспределяется на другие узлы.

В некоторых случаях для обработки Больших данных, особенно при использовании нейронных сетей, используют вычислители на базе GPU, аналогичные используемым для HPC.

Распределённая система хранения данных состоит из дисковых полок, обеспечивающих максимально быстрый доступ к данным. Зачастую также используются компьютеры, в которых подключено большое количество локальных дисков.



Обработка Больших данных является сложным технологическим процессом, требующим глубоких программно инженерных знаний для разработки модели данных, выбора соответствующих программно-аппаратных средств и оценки совокупной стоимости управления данными.

Во многих случаях обработка данных может быть проведена достаточно скромными средствами, при помощи аренды систем хранения и обработки в облачной среде, в других случаях требуется аренда или даже строительство собственного ЦОД и установка собственного оборудования, в третьих случаях — стоимость работы с данными может превысить доход от их обработки и обработка данных своими силами нецелесообразна, однако может быть выполнена при помощи подрядчика.