

# Промышленное машинное обучение на Spark

Лекция 1: Docker

21 октября 2023 г.

01

Оборачиваем модель в сервис. Первая часть: Docker. Flask

02

Оборачиваем модель в сервис. Вторая часть: Requests. REST API

03

Распределенные вычисления. HDFS. MapReduce. Spark DataFrame

04

Погружение в среду Spark. RDD, SQL, Pandas API.

05

Генерация признаков. Spark feature engineering

06

Распределенное обучение моделей. Spark ML

07

Обработка и хранение текстовых данных и картинок.  
Spark image processing. Spark NLP.

08

Обработка потоковых данных. Spark Streaming.

# План лекции

## Что такое большие данные и откуда они берутся?

- Сферы производящие большие данные. Data explosion.
- Большие данные - где начало. 4 основных принципа.

## Как хранить большие данные?

- Устройство файловой системы Linux.
- Отказоустойчивость.
- HDFS. Ее устройство и основные свойства.

## Как обрабатывать большие данные?

- Предпосылки к созданию MapReduce.
- Компоненты системы
- Задача подсчета слов. Map. Shuffle. Reduce.

## Какие есть технологии для работы с большими данными?

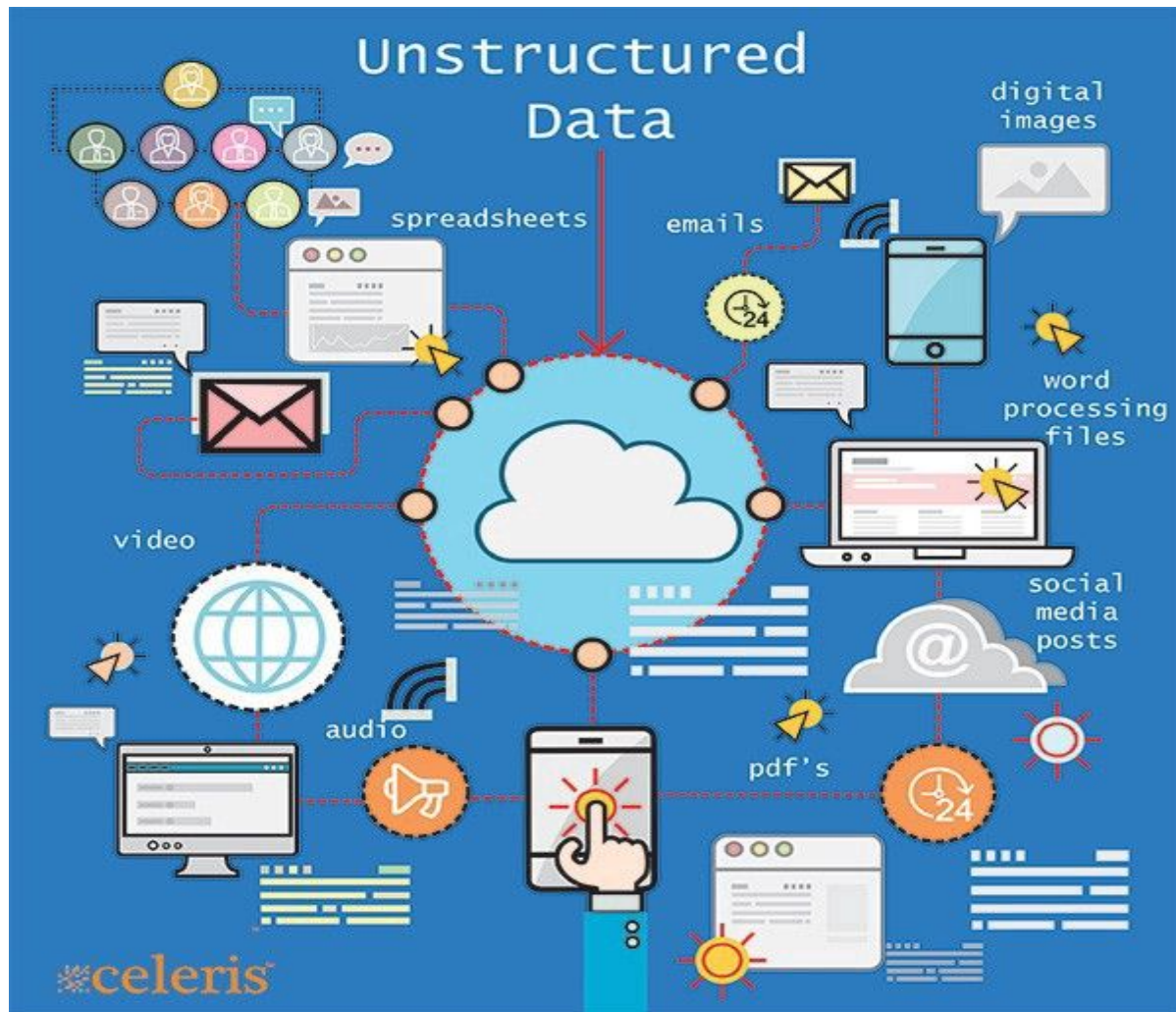
- В чём преимущество технологии Spark.
- Основные компоненты системы, RDD. DAG.
- Spark DataFrame. Поддерживаемые методы.

Сферы производящие большие данные.  
Data explosion.

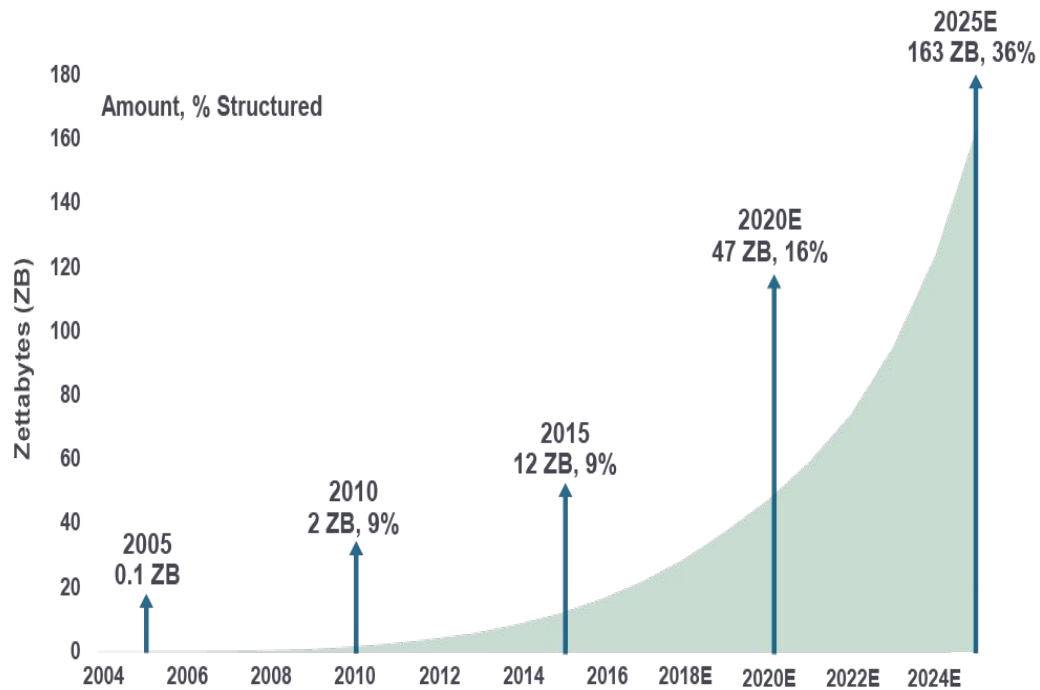
# Откуда пришла Big Data

## Сферы:

- Телеком
- Банки
- Социальные сети
- Медиа
- Промышленность
- Биоинформатика
- Интернет вещей



# Взрыв данных



Чему равен зеттабайт? - триллиону гигабайт

**1 kilobyte      1 000**

**1 megabyte    1 000 000**

**1 gigabyte     1 000 000 000**

**1 terabyte     1 000 000 000 000**

**1 petabyte     1 000 000 000 000 000**

**1 exabyte      1 000 000 000 000 000 000**

**1 zettabyte     1 000 000 000 000 000 000 000**

# Количество обрабатываемых данных



Обрабатывает 40 000PB в день



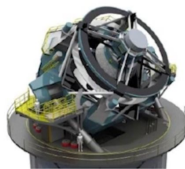
Содержит 300Pb пользовательских данных



Хранит 160Zb



Wayback Machine: 240B веб-страниц в архиве, 5 PB (1/2013)



LSST: 6-10 PB в год (~2015)

LHC: ~15 PB в год



Большие данные, где начало. 4 основных принципа.



# 4'V of big data

## Veracity

## Volume

- Петабайты
- Записи
- Транзакции
- Таблицы, файлы

## Velocity

- Твиты
- Посты в facebook
- Датчики устройств
- Скорость генерации данных
- Генерируемые в реал-тайм
- Онлайн и офлайн данные
- По стримам, батчам или битам

## Variety

- Видео
- Аудио
- Текстовые данные
- Временные ряды
- Структурированные
- Неструктурированные
- Полуструктурированные

Как компании справляются с большими данными?

Первый ключевой вопрос: как хранить большие данные?

# Файловая система

Основными функциями файловой системы являются:

- размещение и упорядочивание на носителе данных в виде файлов;
- определение максимально поддерживаемого объема данных на носителе информации;
- создание, чтение и удаление файлов;
- назначение и изменение атрибутов файлов (размер, время создания и изменения, владелец и создатель файла, доступен только для чтения, скрытый файл, временный файл, архивный, исполняемый, максимальная длина имени файла и т. п.);
- определение структуры файла;
- поиск файлов;
- организация каталогов для логической организации файлов;
- защита файлов при системном сбое;
- защита файлов от несанкционированного доступа и изменения их содержимого

# Linux File System Directories

/bin: Where Linux core commands reside like ls, mv.

/boot: Where boot loader and boot files are located.

/dev: Where all physical drives are mounted like USBs DVDs.

/etc: Contains configurations for the installed packages.

/home: Where every user will have a personal folder to put his folders with his name like /home/likegeeks.

/lib: Where the libraries of the installed packages located since libraries shared among all packages, unlike Windows, you may find duplicates in different folders.

/media: Here are the external devices like DVDs and USB sticks that are mounted, and you can access their files from here.

/mnt: Where you mount other things Network locations and some distros, you may find your mounted USB or DVD.

/opt: Some optional packages are located here and managed by the package manager.

/proc: Because everything on Linux is a file, this folder for processes running on the system, and you can access them and see much info about the current processes.

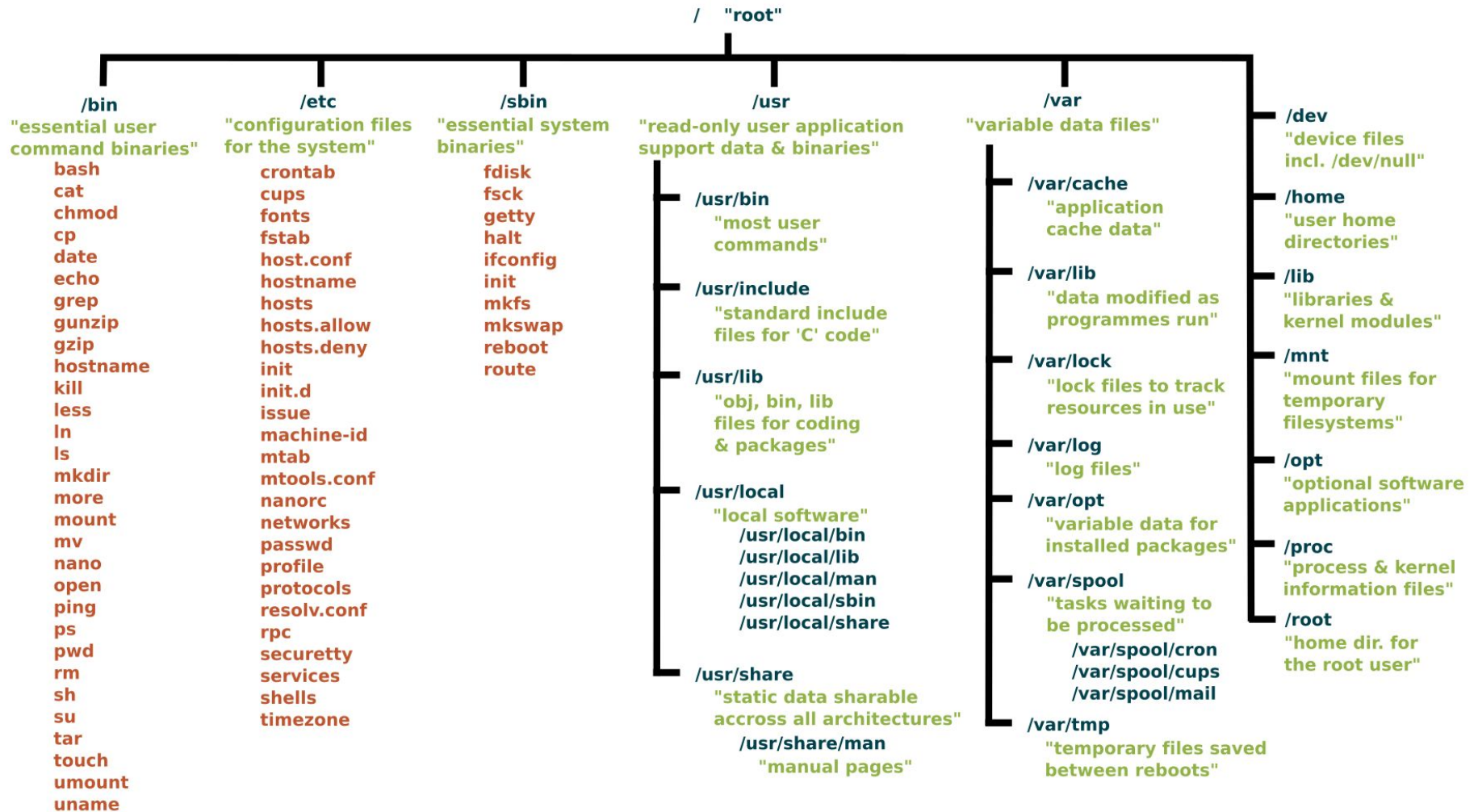
/root: The home folder for the root user.

/sbin: Like /bin, but binaries here are for root user only.

/tmp: Contains the temporary files.

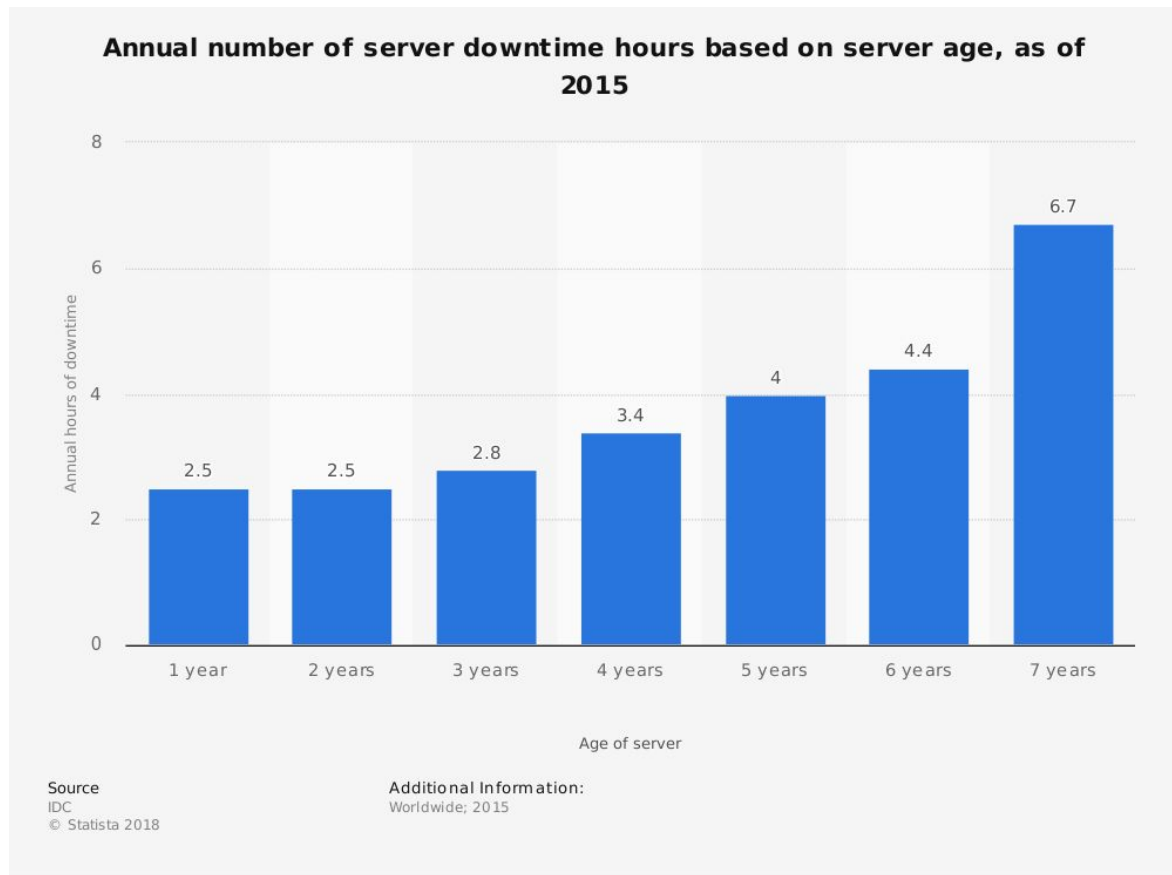
/usr: Where the utilities and files shared between [users on Linux](#).

/var: Contains system logs and other variable data.



Отказоустойчивость хранения данных и вычислений.

## Насколько надёжно хранить данные на сервере



Вероятность, что в следующий час случится поломка:

$$P = 2.5 / (24 * 365) \\ = 0.00028$$

$$P(\text{не выйдет из строя}) \\ = (1 - P) = 0.9997$$

Если в кластере 1000 машин, вероятность, что один из серверов будет недоступен в ближайший час:

$$1 - 0.9997^{1000} = 0.25$$

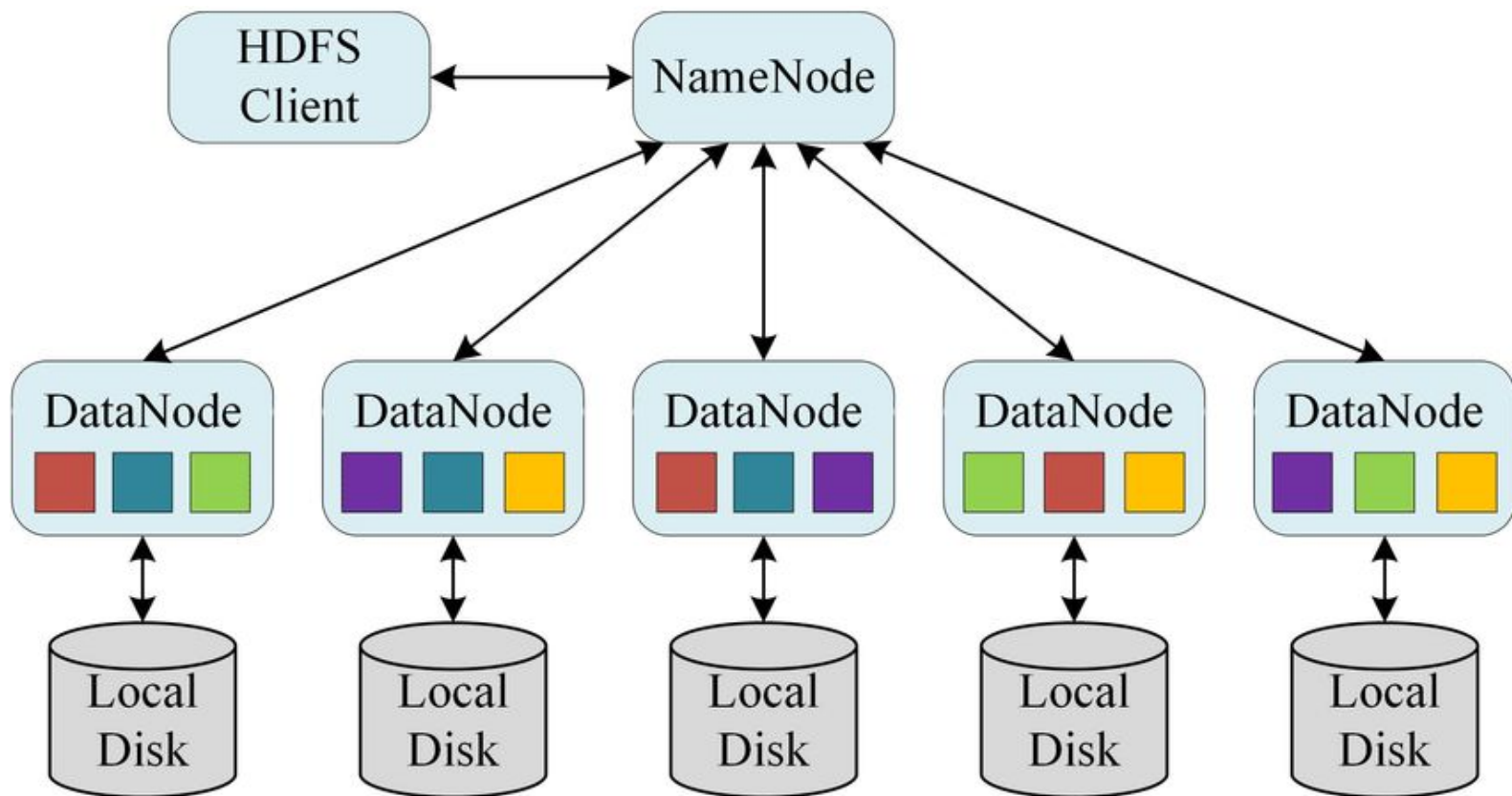


# HDFS - Hadoop Distributed File System

## Основные свойства:

- Доступность данных, за счёт механизмов отказоустойчивость
- Приложения, работающие в HDFS, имеют большие наборы данных. Типичный файл в HDFS имеет размер от гигабайтов до терабайт.
- Эффективно работает в парадигме Write once Read many
- Не требует дорогостоящего оборудования и может быть запущена на маломощных машинах

# Как устроена HDFS

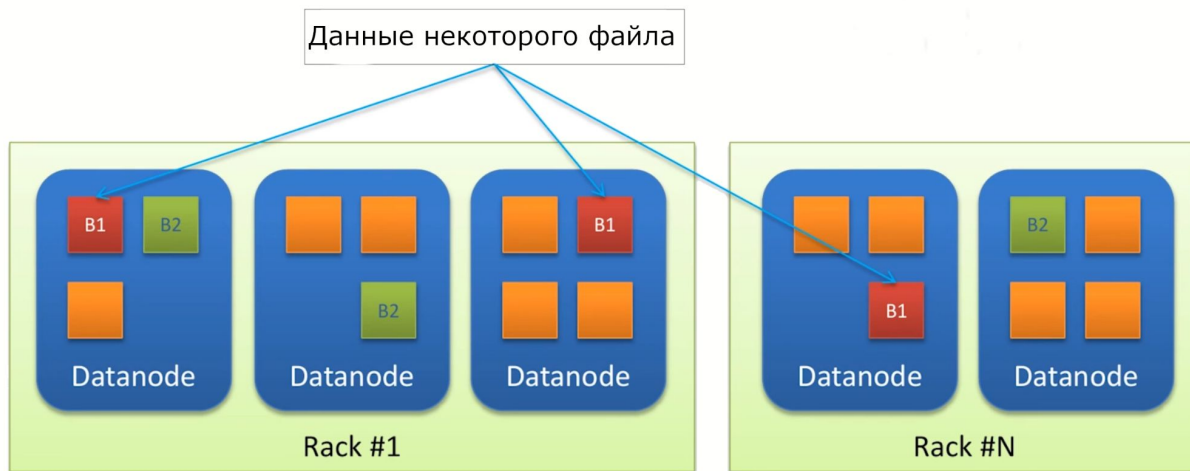


# Из чего состоит HDFS

- ***Классер*** - набор вычислительных узлов (серверов), соединённых сетью;
- ***HDFS client*** - аппаратный или программный компонент системы, откуда поступают команды на произведение некоторых действий в распределённом хранилище данных;
- ***NameNode*** - специально выделенный узел, на котором содержится мета-информация о том, на каком узле хранятся данные для того или иного файла
- ***DataNode*** - один из множества серверов, на котором непосредственно располагаются данные

# Отказоустойчивость в HDFS

Основной принцип, благодаря которому достигается надёжное хранение данных в HDFS заключается в том, что копии одних и тех же данных располагаются на нескольких отдельных друг от друга серверах.



# Фактор репликации

Число  $N$  равное количеству копий одних и тех же данных в системе называется *фактором репликации*

- При выборе  $N$  необходимо подбирать баланс между объёмом занимаемой памяти и уровнем надёжности. Чем больше фактор репликации тем на бóльшем количестве узлов будут содержаться копии данных, но при этом возрастает объём требуемой памяти для хранения
- Как правило на практике зачастую используют  $N = 3$

Второй ключевой вопрос: как обрабатывать большие данные?

# Модель вычислений Map-Reduce

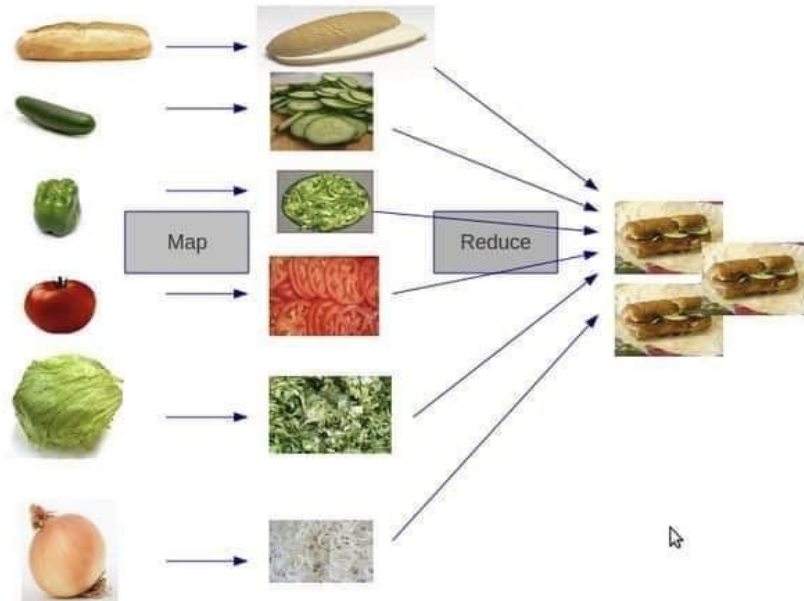
Для эффективной обработки больших объёмов данных применяется двухшаговая модель вычислений Map-Reduce.

- Шаг Map: на данном шаге к каждой порции данных применяется некоторая функция  $F$ .

$$\text{map}(F, [x_1, \dots, x_N]) \rightarrow [F(x_1), \dots, F(x_N)]$$

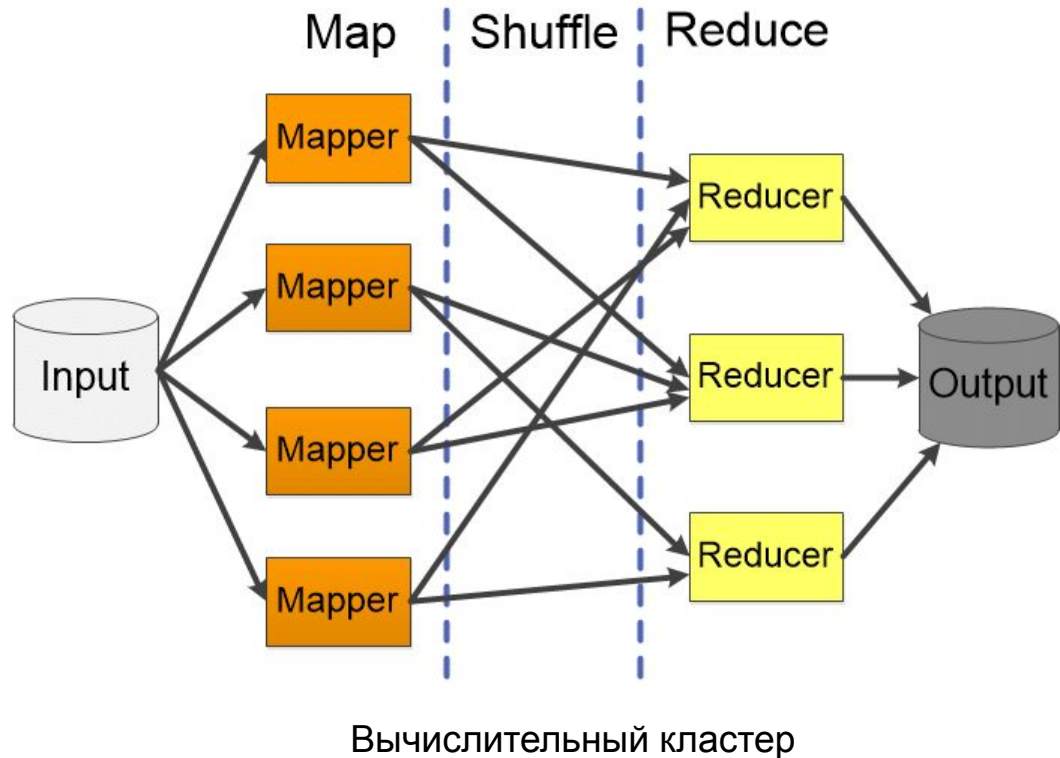
- Шаг Reduce: на данном шаге происходит объединение результатов из предыдущего шага Map и получение итогового результата

$$\text{Reduce}([F(x_1), \dots, F(x_N)]) \rightarrow Y$$



# Архитектура Map-Reduce

- Вычислительные узел, на котором происходит вычисление шага Map называются Mapper
- А вычислительный узел, на котором происходит шаг Reduce, называется Reducer





# Задача подсчета слов

Кошка Мышь Собака

Собака Собака Кошка

Собака Кошка Утка

Кошка Мышь  
Собака

1

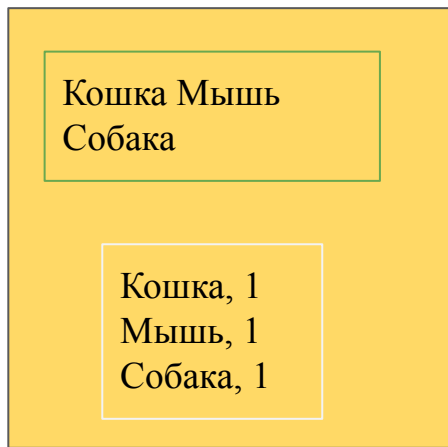
Собака Собака  
Кошка

2

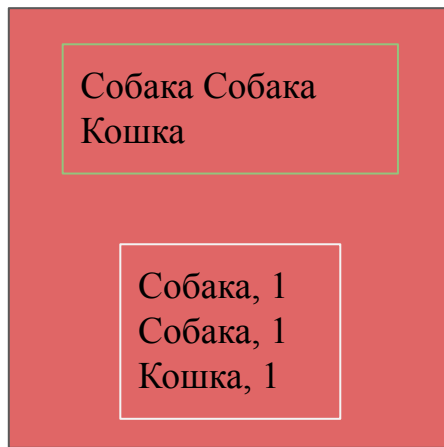
Собака Кошка  
Утка

3

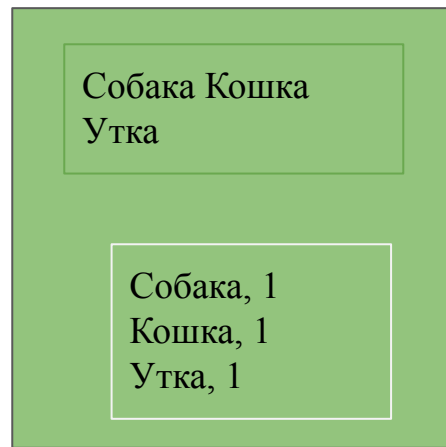
# Шаг Map



1



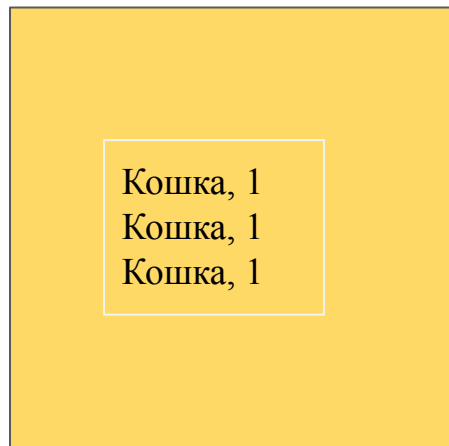
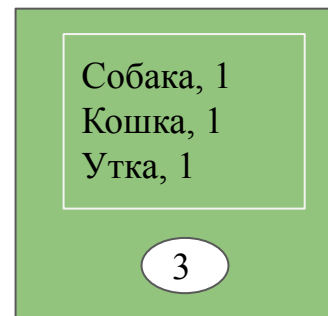
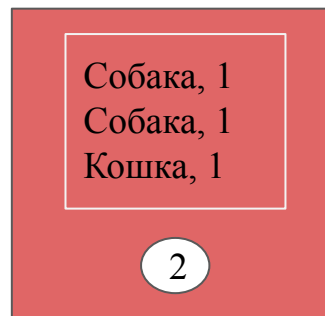
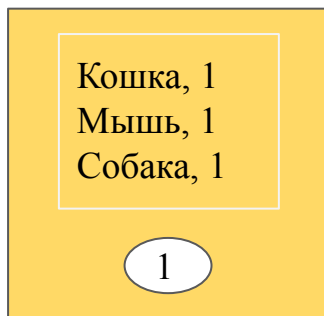
2



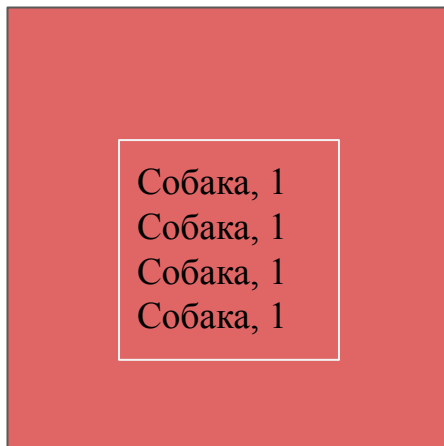
3

# Шаг Shuffle

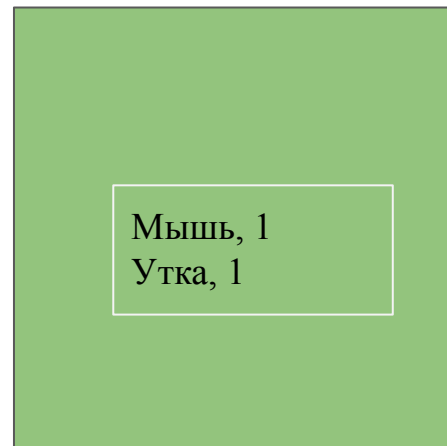
Цель данной процедуры  
отправить правильным  
образом данные на каждый  
Reducer



1

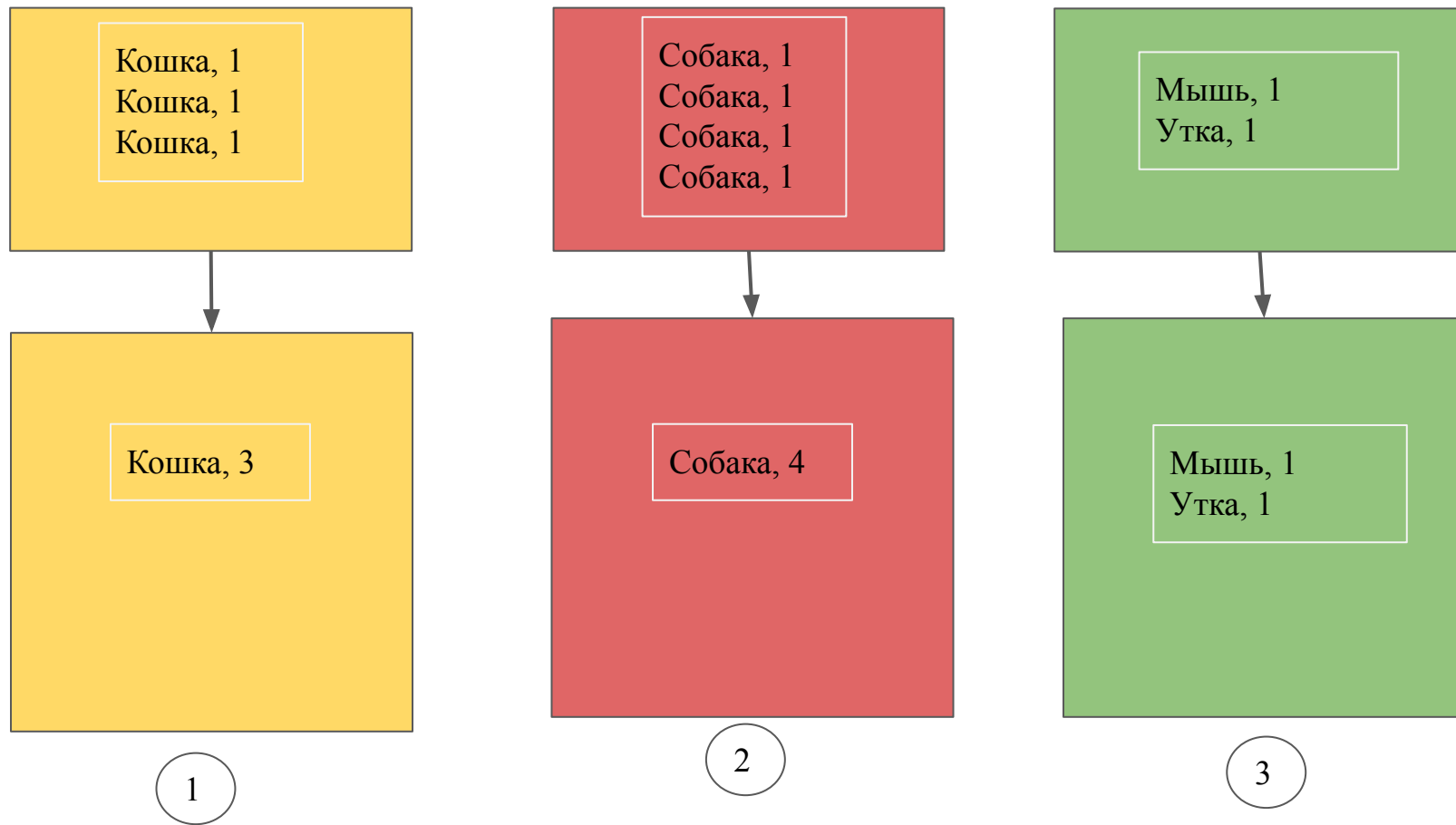


2



3

# Шаг Reduce



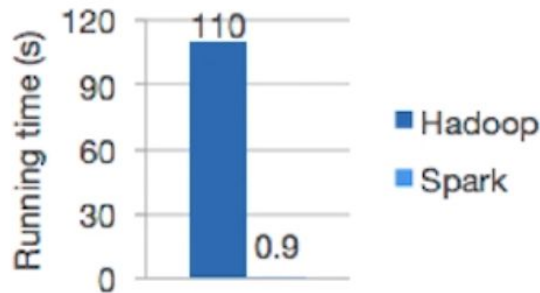
Как со всем этим работают на практике?

# Инструменты для работы с большими данными

Запускать вычисления на распределённом вычислительном кластере крайне непростая задача. Необходимо решить множество аспектов, связанных с как с настройкой самого кластера так и с написанием корректных и эффективных программ на нём.

Для упрощения работы с большими данными была разработана технология Spark, позволяющая разработчикам:

1. Абстрагироваться от технических деталей
2. Писать более эффективные программы за счёт оптимизации вычислений



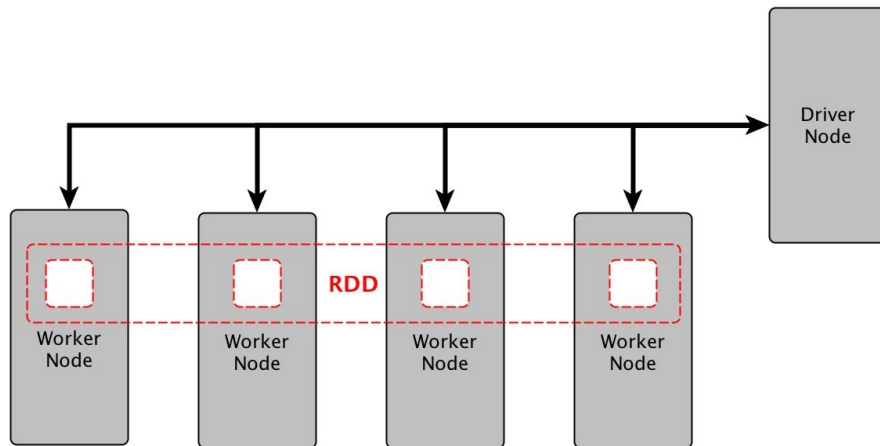
Logistic regression in Hadoop and Spark

# Основные преимущества Spark

1. Высокая скорость работы
2. Удобные абстракции над данными
3. Поддержка нескольких языков программирования: Java, Python, Scala
4. Поддержка разных источников данных
  - 3.1. Файловые системы: HDFS, Linux File System
  - 3.2. Базы данных: HBase, Cassandra, S3
5. Является Open-Source проектом
6. Имеет большое сообщество специалистов вокруг себя

# Основные компоненты. RDD.

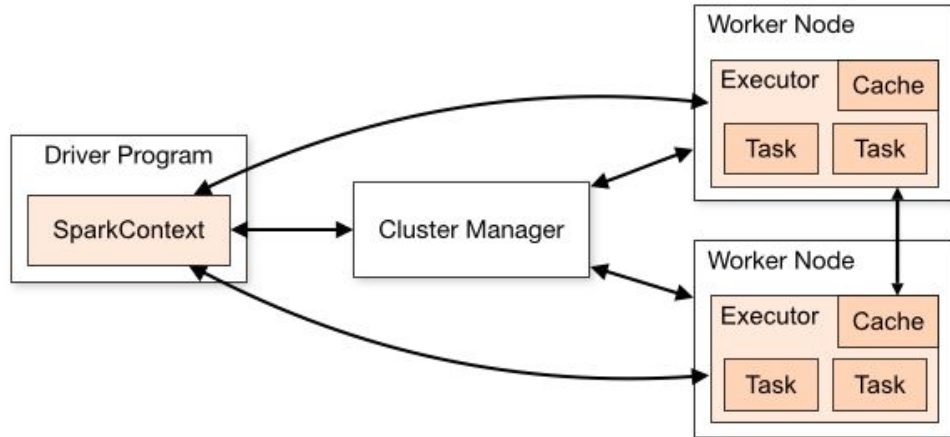
**RDD (Resilient Distributed Dataset)** — это фундаментальная структура данных Spark, которая представляет собой неизменяемый набор данных, который вычисляются и располагается на разных узлах кластера. Каждый набор данных в Spark RDD логически разделен на множество серверов, чтобы их можно было вычислить на разных узлах кластера.





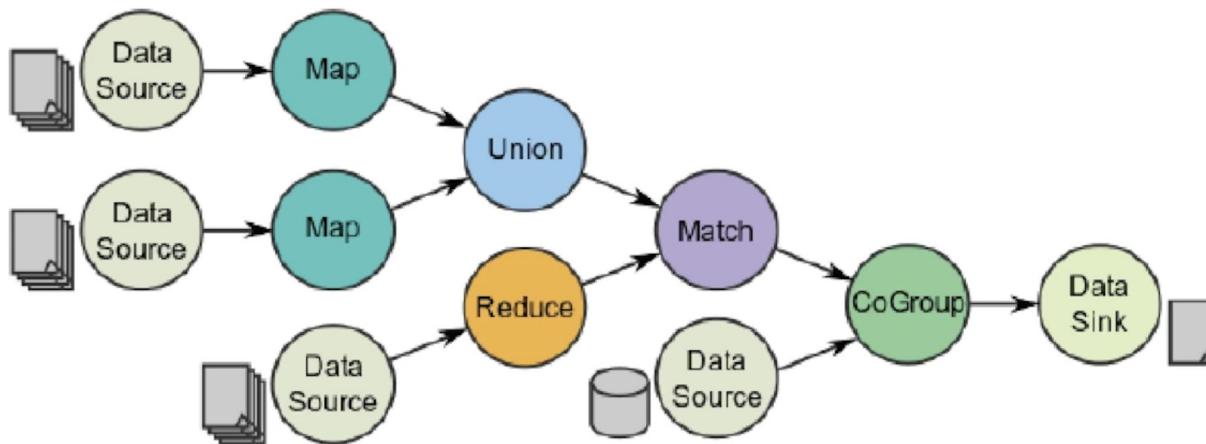
# Основные компоненты. Spark Context

**SparkContext** — это точка входа для всех операций Spark и средство, с помощью которого приложение подключается к ресурсам кластера Spark. Через этот специальный объект Spark происходит конфигурирование ресурсов кластера, обращение к данным и постановка различных задач на исполнение.



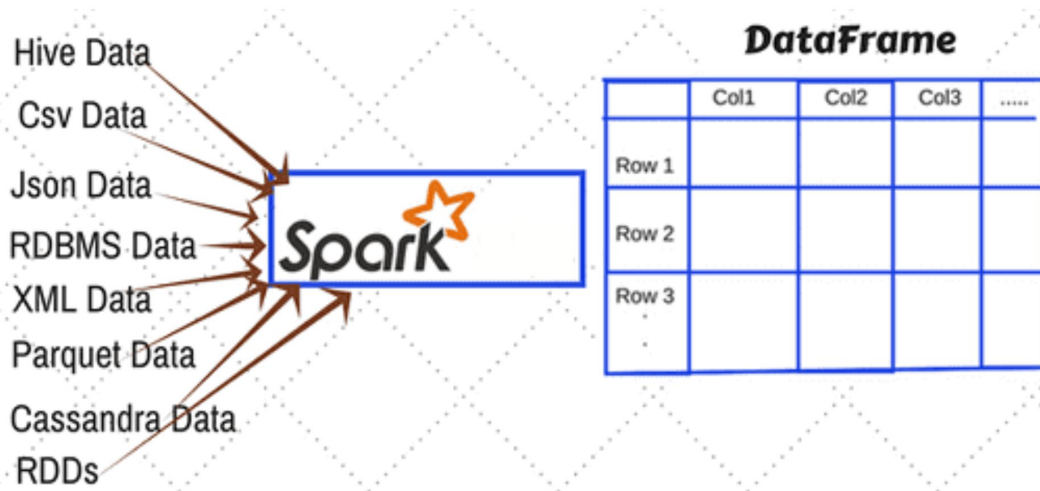
# Основные компоненты. DAG.

**DAG (Направленный ациклический граф)** — это набор вершин и ребер, где вершины представляют RDD, а ребра представляют операцию, которая будет применяться к RDD. В Spark DAG каждое ребро направляет от более раннего к более позднему в последовательности.



# Основные компоненты. Spark DataFrame.

**Spark DataFrame** — это набор данных, организованный в виде таблицы. Концептуально он эквивалентен таблице базы данных или фрейму данных в Python, но с более обширной внутренней оптимизацией. Spark DataFrames могут быть созданы из различных источников, таких как файлы структурированных данных, таблицы в Hive, внешние базы данных или существующие RDD.



# Операции над данными

В Spark поддерживается два вида операторов, которые работают над RDD:

- Transformations - операторы, которые создают новые RDD на основе имеющихся, при этом самого исполнения операции не происходит
- Actions - операторы, которые запускают вычисления и возвращают результат работы в программу или сохраняют его на диск

# Transformations

<b>map</b> ( <i>func</i> )	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .	<b>groupByKey</b> ([ <i>numPartitions</i> ])	When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable<V>) pairs. <b>Note:</b> If you are grouping in order to perform an aggregation (such as a sum or average) over each key, using <code>reduceByKey</code> or <code>aggregateByKey</code> will yield much better performance. <b>Note:</b> By default, the level of parallelism in the output depends on the number of partitions of the parent RDD. You can pass an optional <code>numPartitions</code> argument to set a different number of tasks.
<b>filter</b> ( <i>func</i> )	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.		
<b>mapPartitions</b> ( <i>func</i> )	Similar to <code>map</code> , but runs separately on each partition (block) of the RDD, so <i>func</i> must be of type <code>Iterator&lt;T&gt; =&gt; Iterator&lt;U&gt;</code> when running on an RDD of type T.		
<b>sortByKey</b> ([ <i>ascending</i> ], [ <i>numPartitions</i> ])	When called on a dataset of (K, V) pairs where K implements <code>Ordered</code> , returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean <i>ascending</i> argument.		
<b>repartition</b> ( <i>numPartitions</i> )	Reshuffle the data in the RDD randomly to create either more or fewer partitions and balance it across them. This always shuffles all data over the network.	<b>reduceByKey</b> ( <i>func</i> , [ <i>numPartitions</i> ])	When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function <i>func</i> , which must be of type <code>(V,V) =&gt; V</code> . Like in <code>groupByKey</code> , the number of reduce tasks is configurable through an optional second argument.

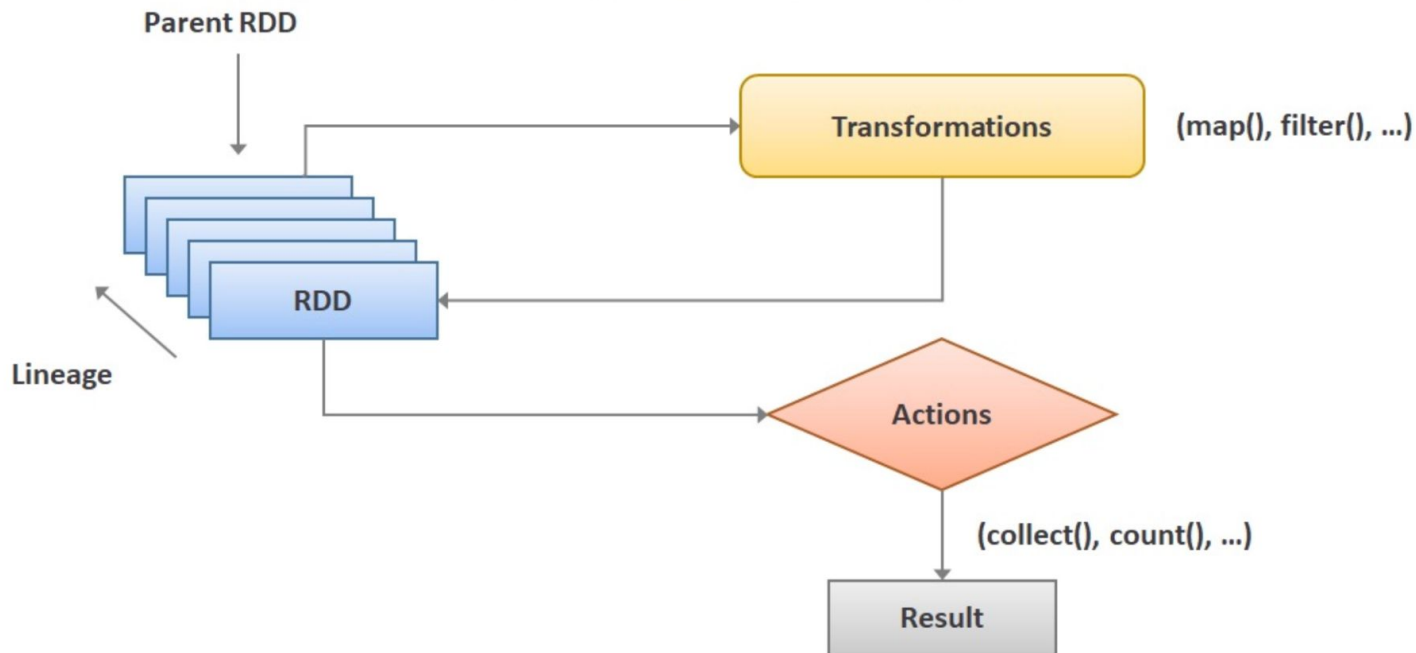
# Actions

<b>reduce</b> ( <i>func</i> )	Aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
<b>collect</b> ()	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
<b>count</b> ()	Return the number of elements in the dataset.
<b>first</b> ()	Return the first element of the dataset (similar to <code>take(1)</code> ).
<b>take</b> ( <i>n</i> )	Return an array with the first <i>n</i> elements of the dataset.
<b>countByKey</b> ()	Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.

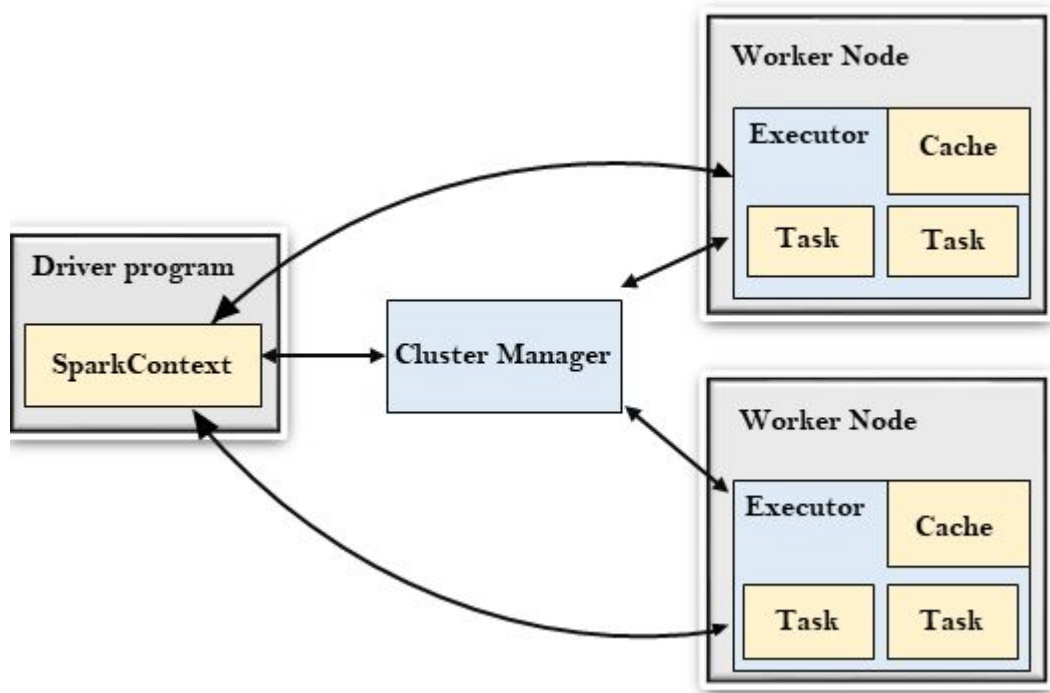
# Жизненный цикл RDD



*Spark RDD (Unstructured) Operations*



# Устройство Spark кластера



- ***Worker Node*** - компоненты системы Spark, которые выполняют фактическое действие над данными;
- ***Spark Driver*** - компонент, ответственный за планирование вычислений и подготовку задач для Worker'ов;
- ***Cluster Manager*** - компонент, который координирует обмен данными между узлами;



# Data locality

В Spark реализованы механизмы, поддерживающие локальность данных. Под локальностью данных понимается то, что вычисления производятся на том узле, где располагаются необходимые данные. Это позволяет уменьшить нагрузку на сеть и увеличить производительность системы.

Достигается это через передачу Spark-кода на узлы, на которых располагаются данные, а не наоборот.