# The Pilot Of Spark

2017.5    XenRon
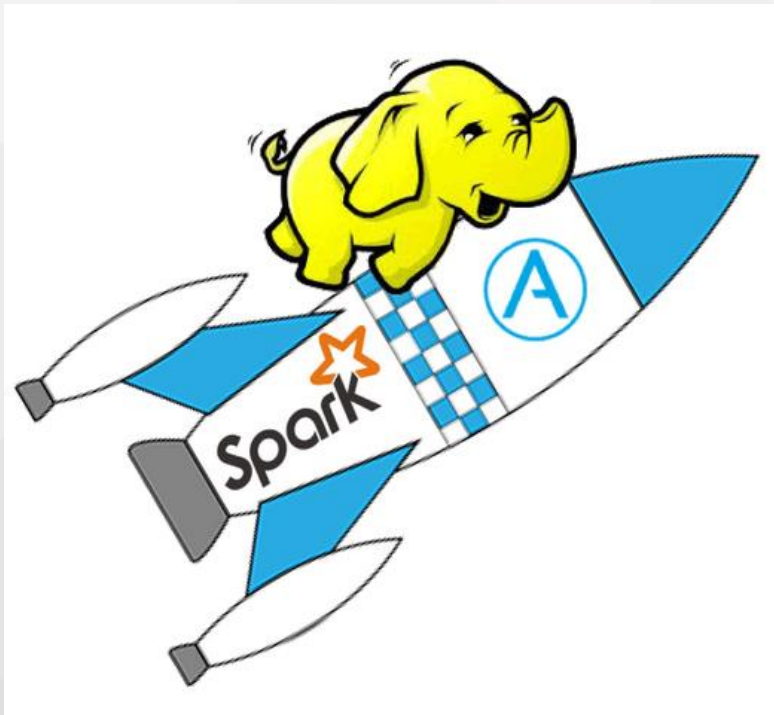
http://spark.apache.org/docs/latest/

# 📊 CONTENTS 📊
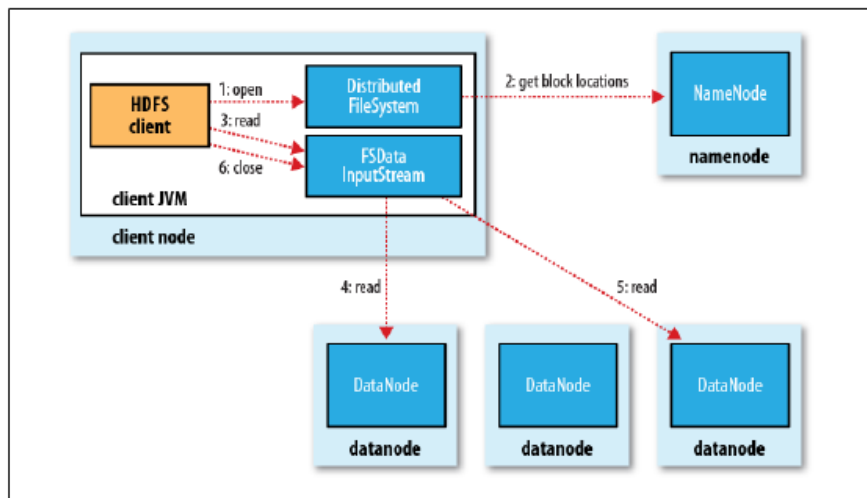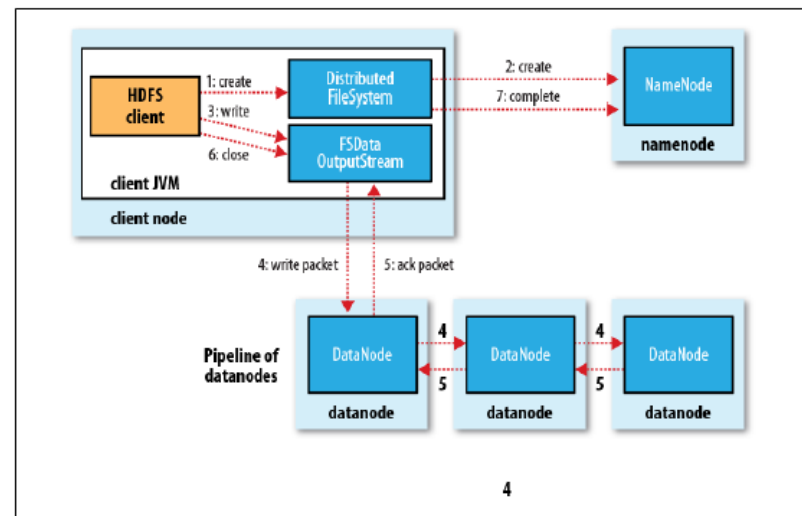
Review

# Map Reduce

Figure 3-2. A client reading data from HDFS
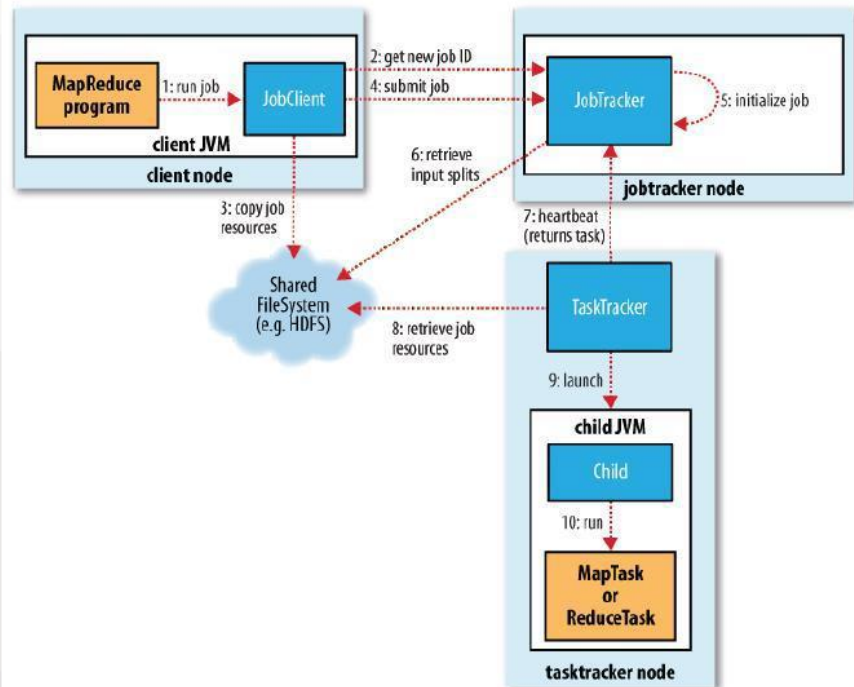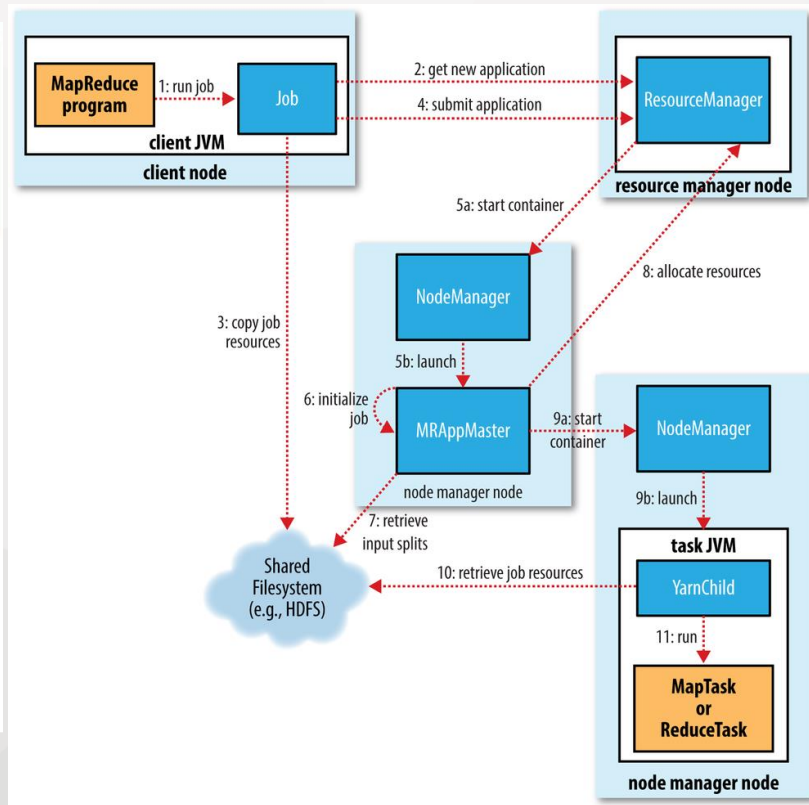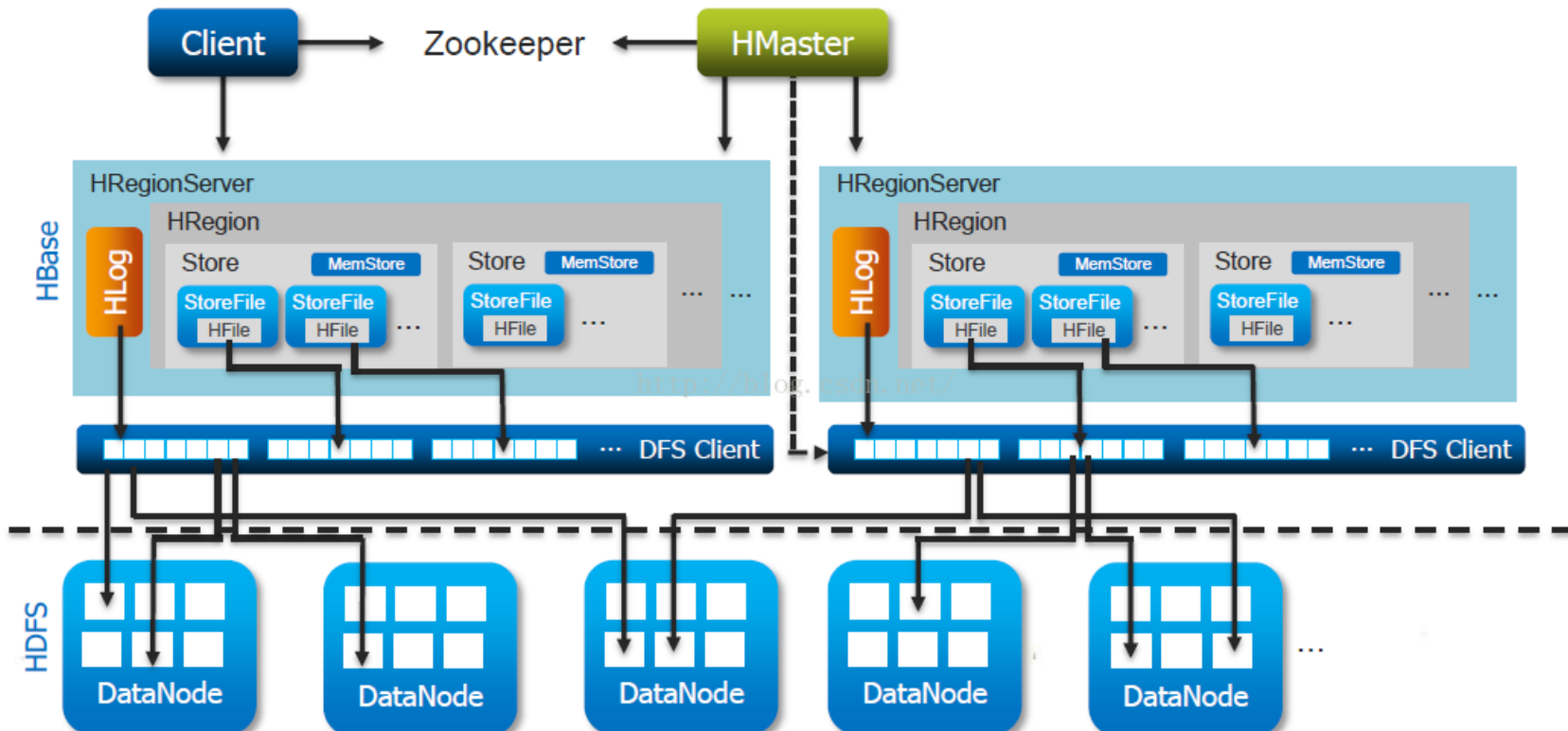


Figure 3-4. A client writing data to HDFS

# Map Reduce

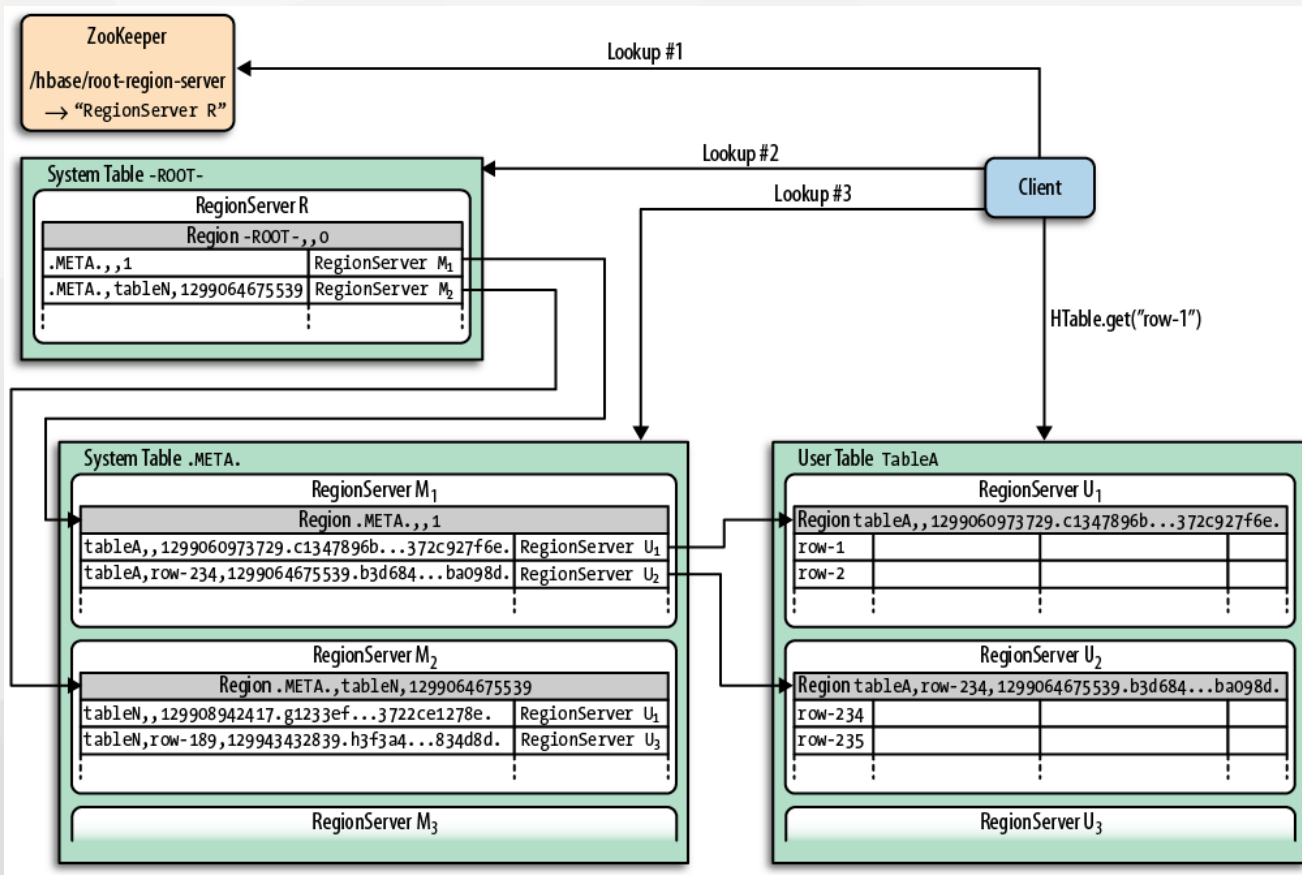Hadoop 1.x

Hadoop 2.x

# PART1

**Preliminary Topics**
**事前準備**

JVM Languages: "The Big Three" + JRuby

(Bar chart showing Groovy ~1150, Scala ~1310, Clojure ~1230, JRuby ~480)

Java
Scala
JRuby
Groovy

Compiler

Class file

cafe babe
0000 0032
0017 0100

JVM

- Runs
- Interprets
- Translates bytecode into Native Machine Code

# BUILD TOOL – SBT

http://www.scala-sbt.org/

```scala
name := "hello world"

version := "0.0.1"

scalaVersion := "2.11.1"

resolvers ++= Seq (
Resolver.mavenLocal,
Resolver.sonatypeRepo ("releases"),
Resolver.typesafeRepo ("releases")
)

libraryDependencies ++=
Seq ("org.scala-lang" % "scala-compiler" % "2.11.1")

addSbtPlugin ("com.typesafe.play" % "sbt-plugin" % "2.3.1")
```

```xml
1  <project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
2            xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org/maven-v4_0_0.xsd">
3      <modelVersion>4.0.0</modelVersion>
4      <groupId>info.solidsoft.rnd</groupId>
5      <artifactId>spock-10-groovy-24-gradle-maven</artifactId>
6      <version>0.0.1-SNAPSHOT</version>
7      <properties>
8          <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
9          <surefire.version>2.18.1</surefire.version>
10     </properties>
11     <build>
12         <plugins>
13             <plugin>
14                 <groupId>org.codehaus.gmavenplus</groupId>
15                 <artifactId>gmavenplus-plugin</artifactId>
16                 <version>1.4</version>
17                 <executions>
18                     <execution>
19                         <goals>
20                             <goal>compile</goal>
21                             <goal>testCompile</goal>
22                         </goals>
23                     </execution>
24                 </executions>
25             </plugin>
26             <plugin>
27                 <artifactId>maven-surefire-plugin</artifactId>
28                 <version>${surefire.version}</version>
29                 <configuration>
30                     <includes>
31                         <include>**/*Spec.java</include> <!-- Yes, .java extension -->
32                         <include>**/*Test.java</include> <!-- Just in case having "normal" JUnit tests -->
33                     </includes>
34                 </configuration>
35             </plugin>
36         </plugins>
37     </build>
38     <dependencies>
39         <dependency>
40             <groupId>org.codehaus.groovy</groupId>
41             <artifactId>groovy-all</artifactId>
42             <version>2.4.1</version>
43         </dependency>
44         <dependency>
45             <groupId>org.spockframework</groupId>
46             <artifactId>spock-core</artifactId>
47             <version>1.0-groovy-2.4</version>
48             <scope>test</scope>
49         </dependency>
50     </dependencies>
51 </project>
52
```

pom.xml

```groovy
1  apply plugin: 'groovy'
2
3  group = "info.solidsoft.rnd"
4  version = "0.0.1-SNAPSHOT"
5
6  repositories {
7      mavenCentral()
8  }
9
10 dependencies {
11     compile 'org.codehaus.groovy:groovy-all:2.4.1'
12
13     testCompile 'org.spockframework:spock-core:1.0-groovy-2.4'
14 }
15
```

**build.gradle**

```
1  rootProject.name = 'spock-10-groovy-24-gradle-maven'
2
```
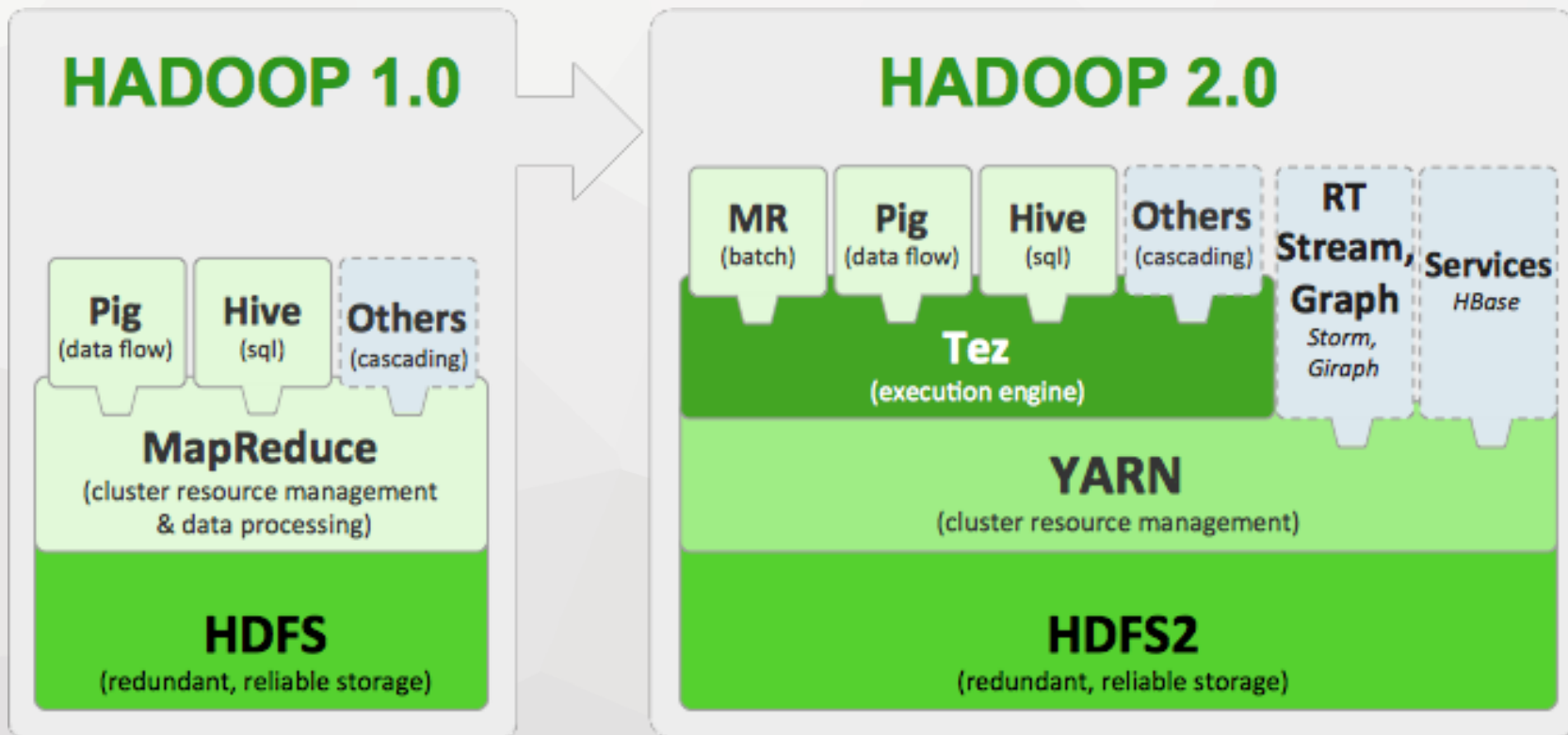
**settings.xml**

maven

gradle

## Applications Run Natively IN Hadoop

| Pig | Hive | HBase | Accumulo | Storm | Solr | Spark | Cascading | Others |
|-----|------|-------|----------|-------|------|-------|-----------|--------|
| Script | SQL | NoSQL | NoSQL | Stream | Search | In-Memory | Java | ISV Engines |

### YARN: Data Operating System

### HDFS
(Hadoop Distributed File System)

Node Manager
- Container
- App Mstr
  - 5
  - 6

Node Manager
- App Mstr
- Container
  - 2

Node Manager
- Container
- Container
  - 6
  - 6
  - 5
  - 5

7

Client

Client

1

Resource Manager

3, 4, 8

Legend:
- MapReduce Status ⟶
- Job Submission ----⟶
- Node Status -·-·-⟶
- Resource Request ·····⟶

https://hadoop.apache.org/docs/stable/hadoop-yarn/hadoop-yarn-site/WritingYarnApplications.html

PART2

Spark Environment

# Spark Source Compile

PART3

Spark Architecture

RDD Objects

DAGScheduler

TaskScheduler

Worker

DAG

TaskSet

Cluster manager

Task

Threads

Block manager

```
rdd1.join(rdd2)
    .groupBy(...)
    .filter(...)
```

build operator DAG

split graph into *stages* of tasks

submit each stage as ready

launch tasks via cluster manager

retry failed or straggling tasks

execute tasks

store and serve blocks

agnostic to operators!

stage failed

doesn't know about stages

# What is an RDD?



**Some RDD Characteristics**

- Hold references to Partition objects

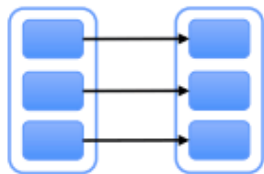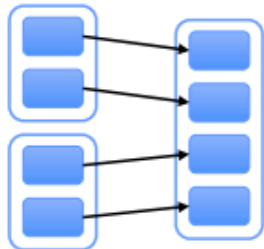- Each Partition object references a subset of your data

- Partitions are assigned to nodes on your cluster

- Each partition/split will be in RAM (by default)

Dependency Types

"Narrow" (pipeline-able)

map, filter

union

join with inputs co-partitioned

"Wide" (shuffle)

groupByKey on non-partitioned data

join with inputs not co-partitioned

# Example

31

# Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
messages = errors.map(_.split('\t')(2))
cachedMsgs = messages.cache()


cachedMsgs.filter(_.contains("foo")).count
cachedMsgs.filter(_.contains("bar")).count

. . .
```

**Result:** scaled to 1 TB data in 5-7 sec
(vs 170 sec for on-disk data)

# Driver
# Executor
# Application

Stage

ComplexJob
including map(), partitionBy(), union(), and join()

# Spark Program Flow by RDD

**Create input RDD**
- spark context
- create from external
- create from SEQ

Example:
myRDD = sc.parallelize([1,2,3,4])
myRDD = sc.textFile("hdfs://tmp/shakepeare.txt")

**Transform RDD**
- Using map, flatmap, distinct, filter
- Wide/Narrow

**Persist intermediate RDD**
- Cache
- Persist

**Action on RDD**
- Produce result

# PART4

Spark SQL

# Hive

- Data warehousing package built on top of Hadoop
- Bringing structure to unstructured data
- Query petabytes of data with HiveQL
- Schema on read



Metastore

Stores schema information
Provides a structure to data stored

Hive Engine

Runs query processing, compiler, optimizer and executer (using MapReduce)

MapReduce

HDFS

# Hadoop MapReduce Vs Pig Vs Hive

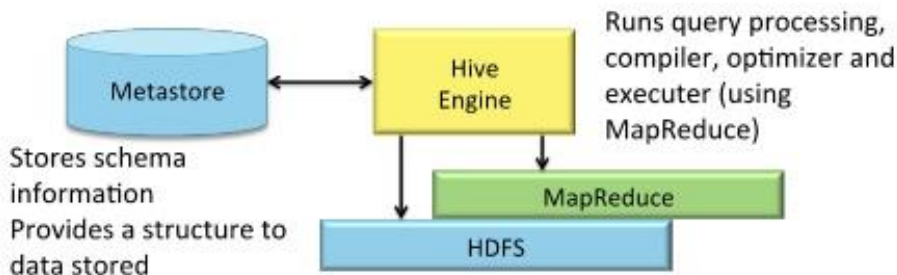| Hadoop MapReduce | Pig | Hive |
|---|---|---|
| Compiled Language | Scripting Language | SQL like query Language |
| Lower Level of Abstraction | Higher Level of Abstraction | Higher Level of Abstraction |
| More lines of Code | Comparatively less lines of Code than MapReduce | Comparatively less lines of Code than MapReduce and Apache Pig |
| More Development Effort is involved | Development Effort is less Code Efficiency is relatively less | Development Effort is less Code Efficiency is relatively less |
| Code Efficiency is high when compared to Pig and Hive | Code Efficiency is relatively less | Code Efficiency is relatively less |

DeZyre

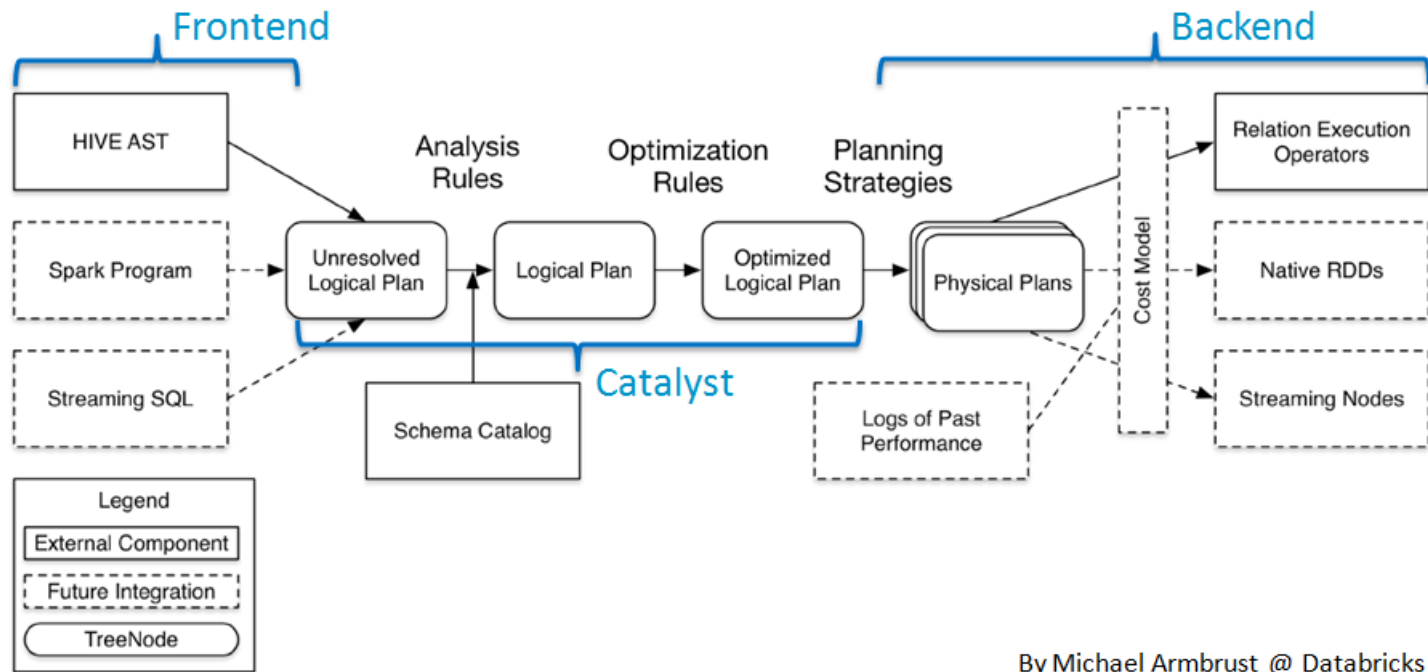# The Right SQL Engine for the Use Case

| BI and SQL Analytics | Batch Processing | Spark Developers |
|---|---|---|

cloudera

# Spark SQL Architecture

Frontend

Backend

HIVE AST

Spark Program

Streaming SQL

Unresolved Logical Plan

Analysis Rules

Schema Catalog

Logical Plan

Optimization Rules

Optimized Logical Plan

Catalyst

Planning Strategies

Physical Plans

Logs of Past Performance

Cost Model

Relation Execution Operators

Native RDDs

Streaming Nodes

Legend

External Component

Future Integration

TreeNode

By Michael Armbrust @ Databricks

Software and Services

PART5

Spark Streaming

**Resource Management**

| Standalone | YARN | Mesos |

**Spark Ecosystems**

| Spark SQL | Spark Streaming | BlinkDB |
| Spark Machine Learning | GraphX | Tachyon |

**Spark Core**

BACK TO BAZICS

**Spark DataFrame API**

| Java | Scala | Python | R |

**Spark Core**

# Kafka

**HTTP Proxy**

**Event Producer**

**Control Plane**

**Fronting Kafka**

**Router**

**S3**

**EMR**

**ELASTIC SEARCH**

**Consumer Kafka**

**Stream Consumers (Spark, Mantis, Custom Apps)**

# ELK (ElasticSearch LogStash Kibana)

华为|                                                  搜索

华为p9                                      约39个商品
华为手机                                    约431个商品
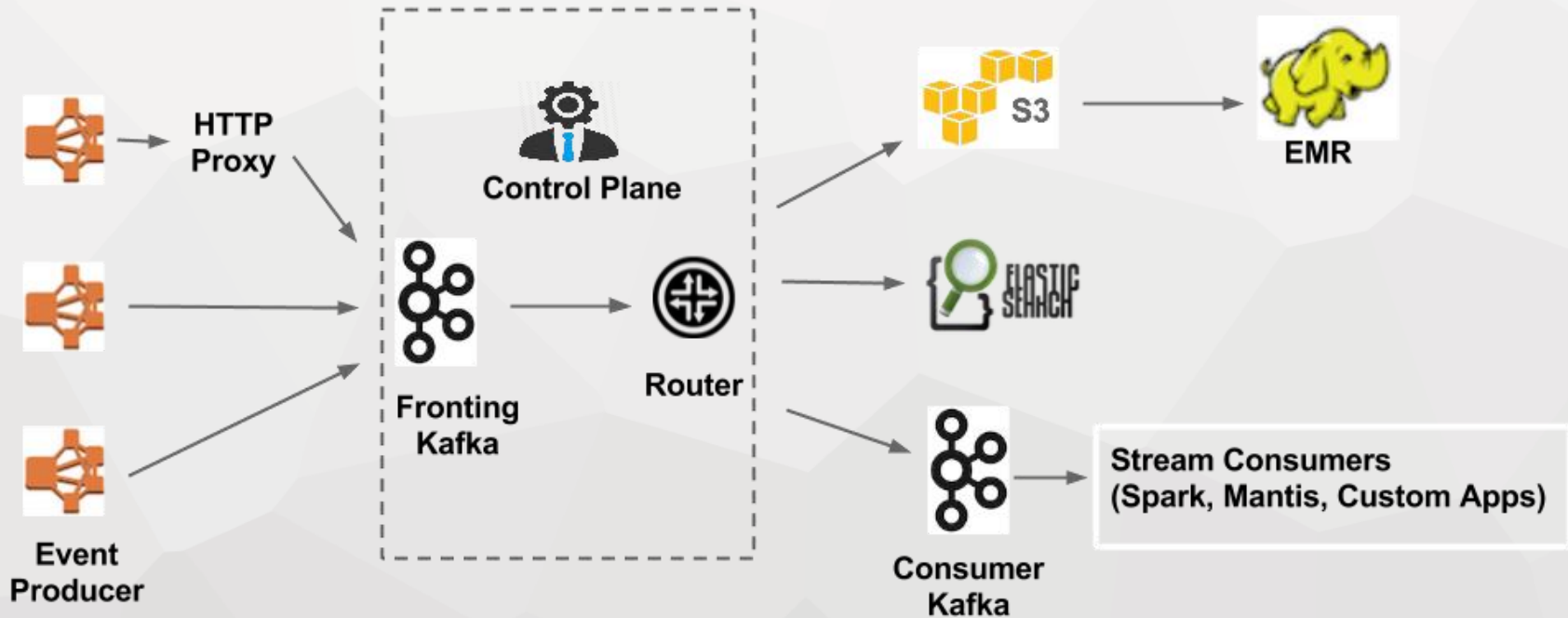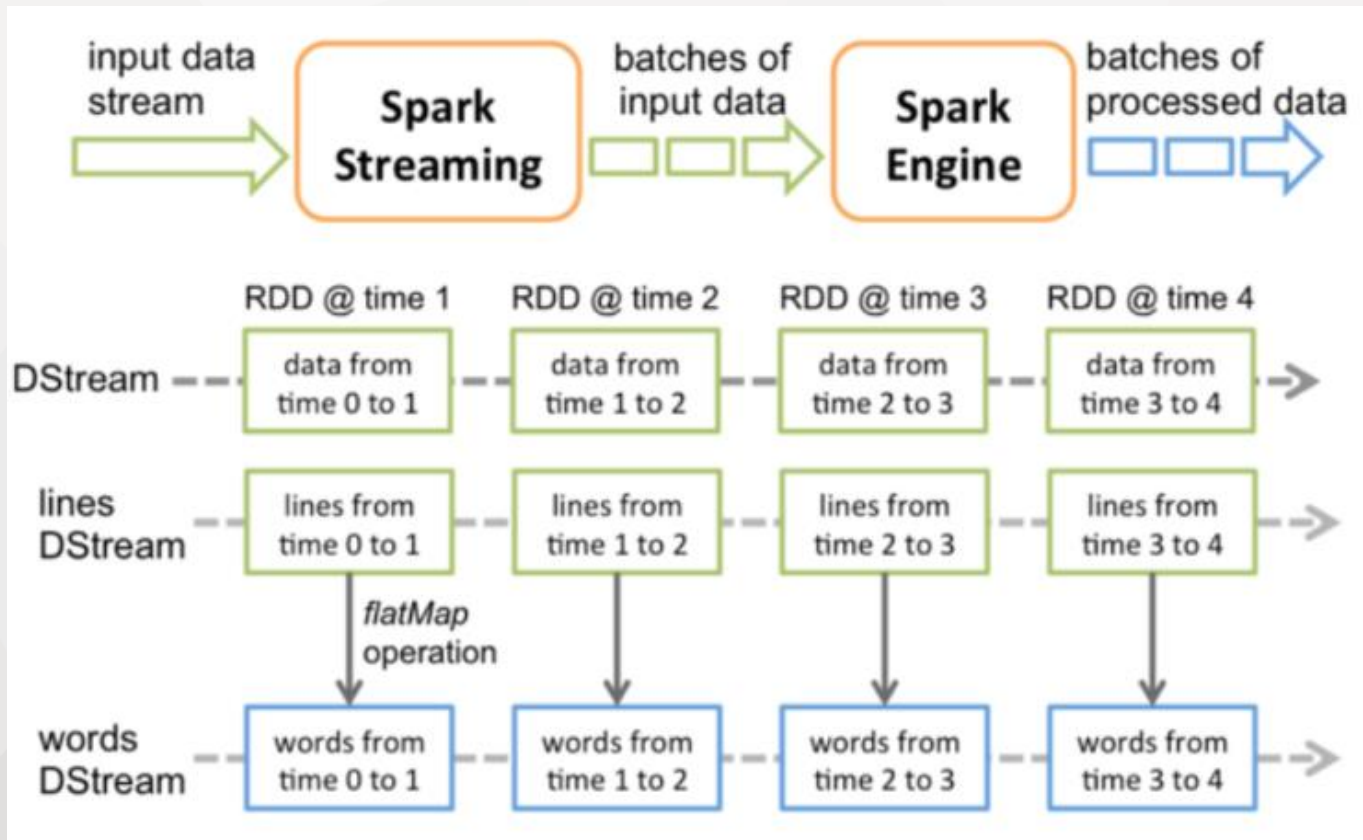华为v8                                    约1273个商品      金融
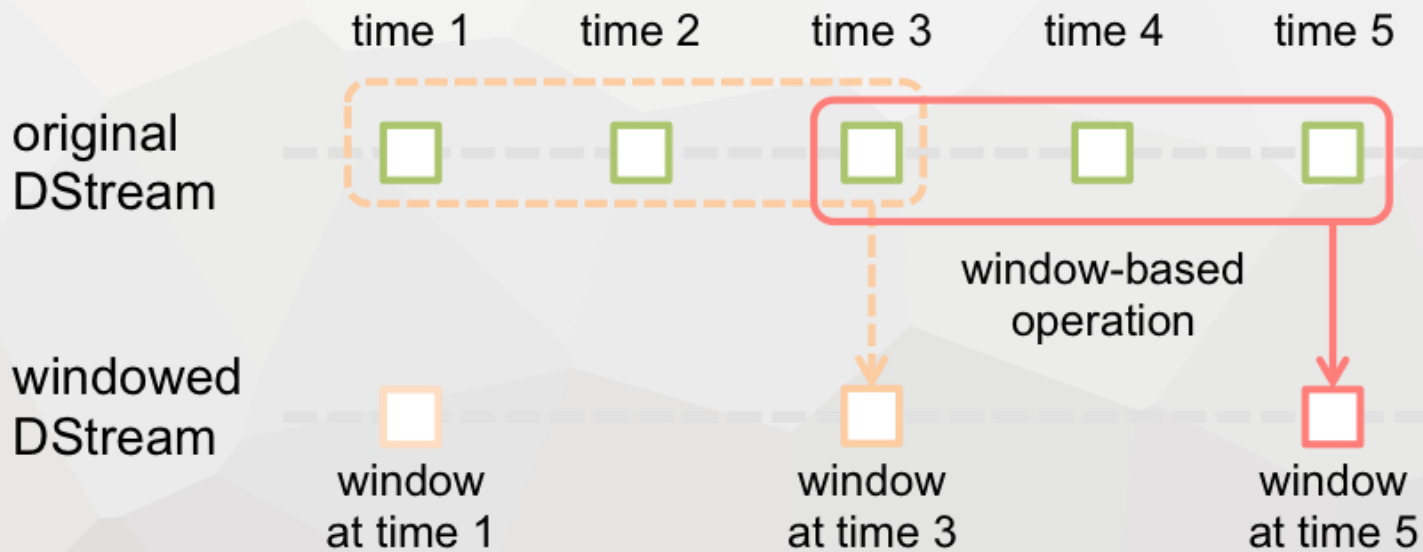华为p8                                      约38个商品
华为p7                                      约41个商品
华为5X                                      约60个商品
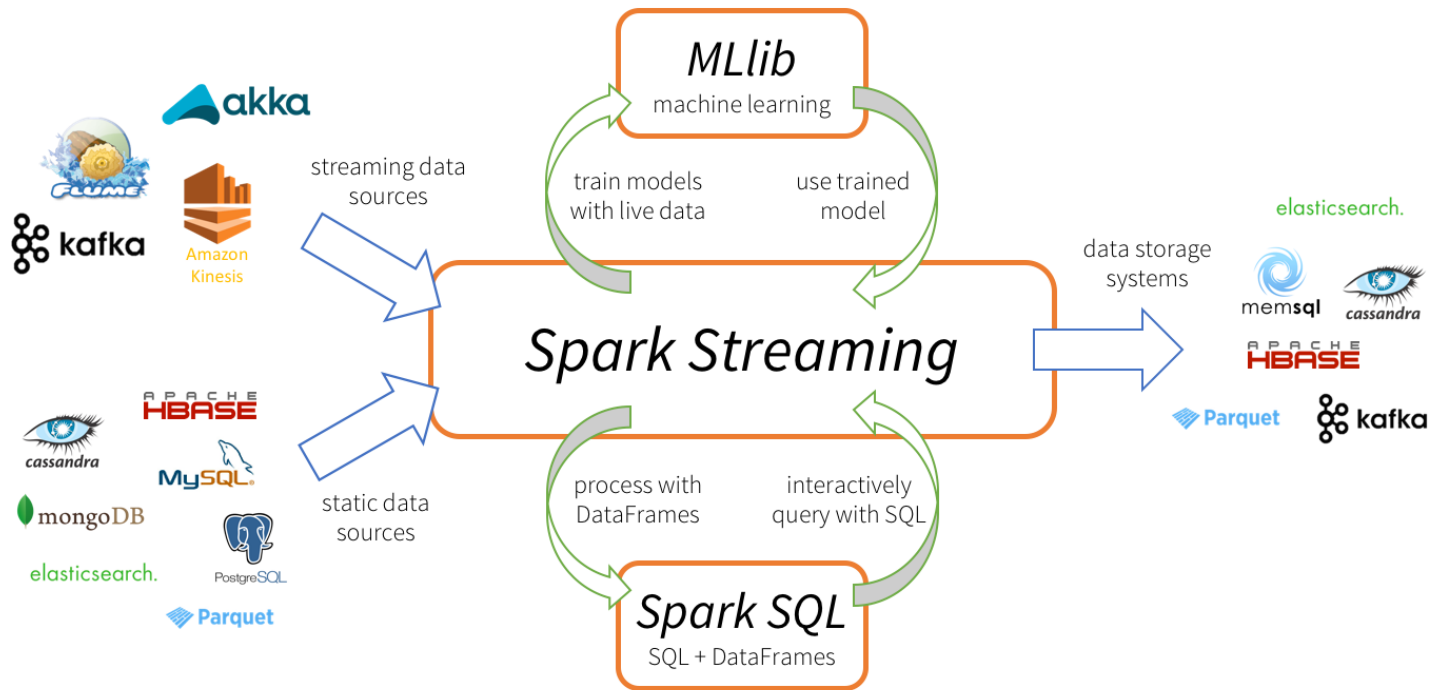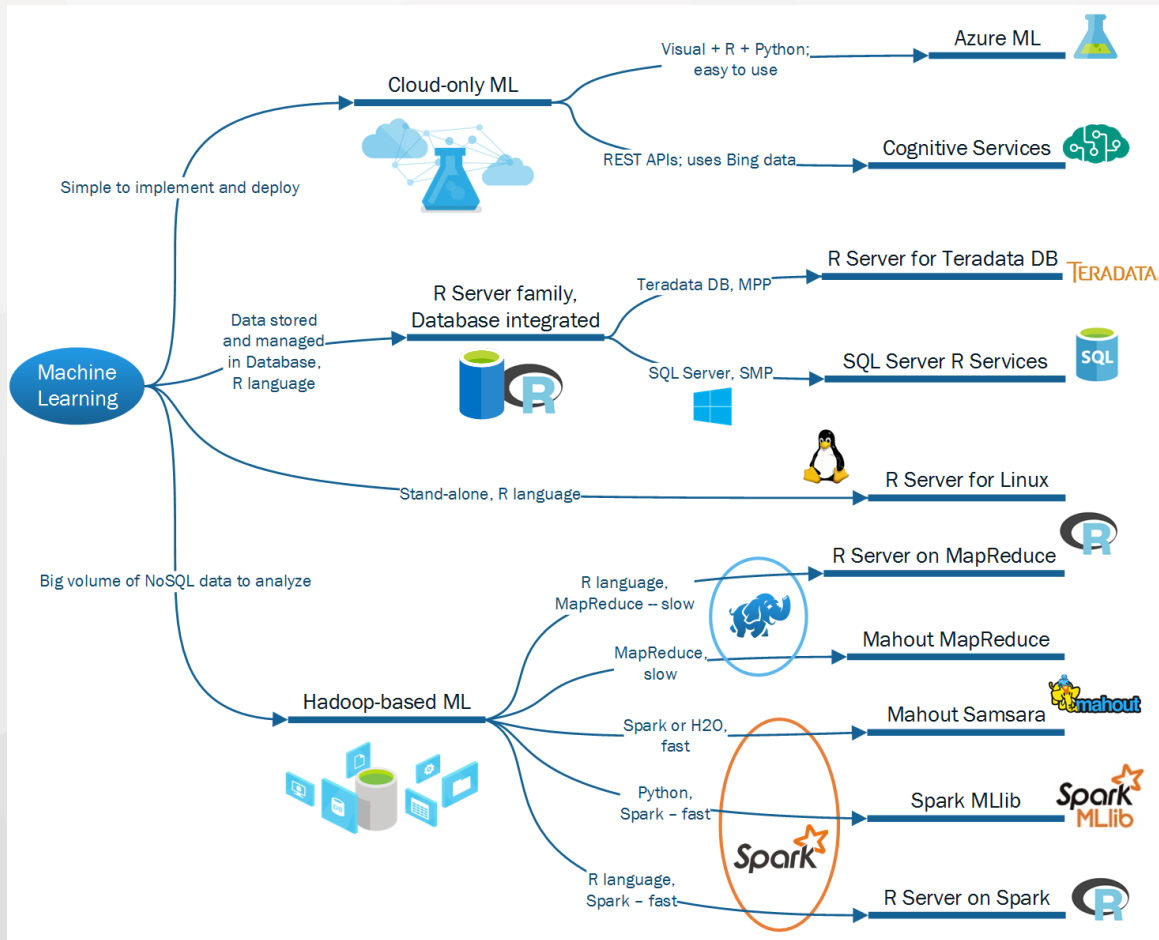华为G9                                     约810个商品
华为荣耀                                    约187个商品
华为5c                                     约778个商品
华为5s                                      约19个商品
                                              关闭
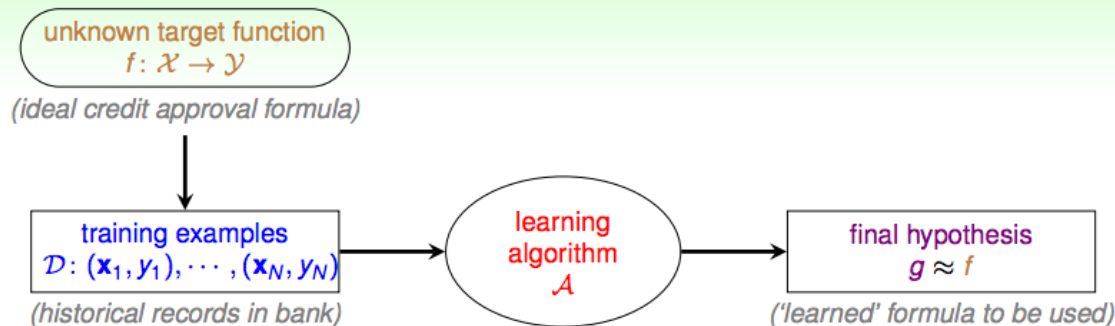
# Windows based operations

# PART6

**Spark Machine Learning**

# Use Case : Credit Approval



The Learning Problem          Components of Machine Learning

## Learning Flow for Credit Approval

unknown target function
$f: \mathcal{X} \rightarrow \mathcal{Y}$

*(ideal credit approval formula)*

training examples
$\mathcal{D}: (\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_N, y_N)$

*(historical records in bank)*

learning algorithm
$\mathcal{A}$

final hypothesis
$g \approx f$

*('learned' formula to be used)*

- target $f$ **unknown**
  (i.e. no programmable definition)
- hypothesis $g$ hopefully $\approx f$
  but possibly **different** from $f$
  (perfection 'impossible' when $f$ unknown)

What does $g$ look like?

# Use Case : Hatsune Miku Upscale

PART7

Spark GraphX
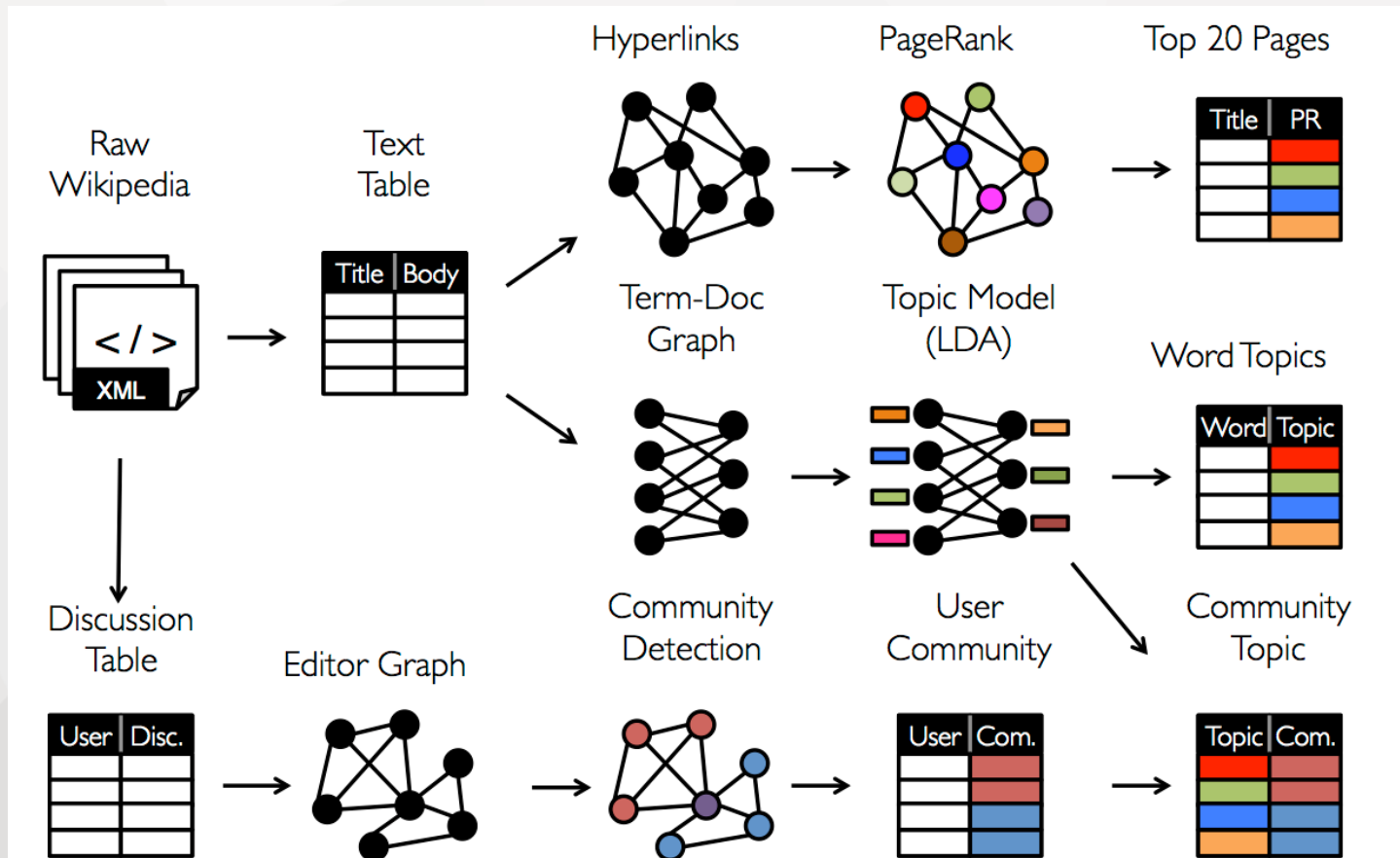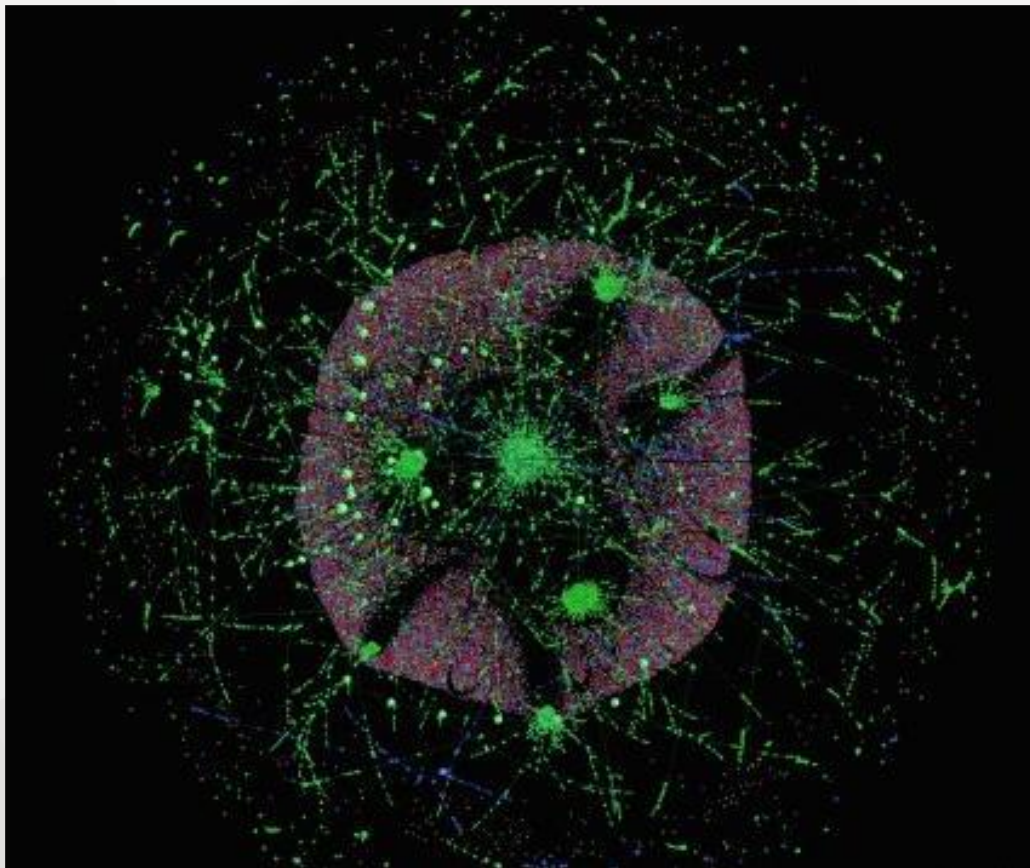
# GraphX
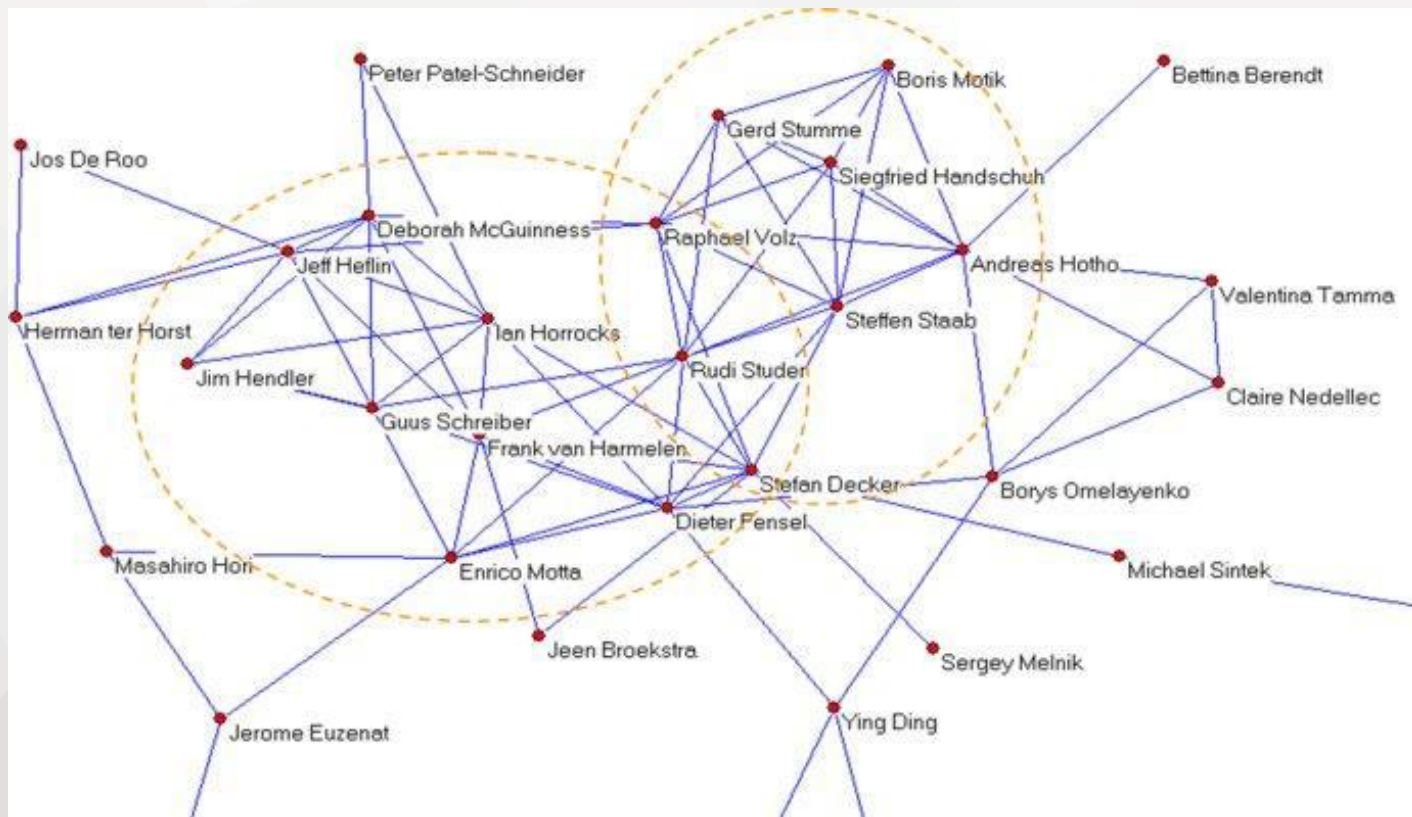
The End