

Apache Spark and Scala

Module 2: Spark and Big Data

Module 1

Getting Started /
Introduction to Scala

Module 2

Spark and Big Data

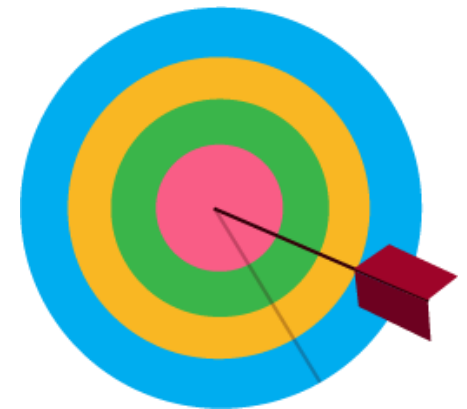
Module 3

Understanding RDDs

Session Objectives

In this session, you will understand:

- Analyze Batch Processing and Real-time Processing
- Understand Spark Ecosystem
- Analyze MapReduce Limitations
- Go through Spark History
- Analyze Spark Architecture
- Understand Spark and Hadoop Advantages
- Analyze benefits of Spark and Hadoop combined
- Install Spark



Bombay Stock Exchange – Big Data Case Study

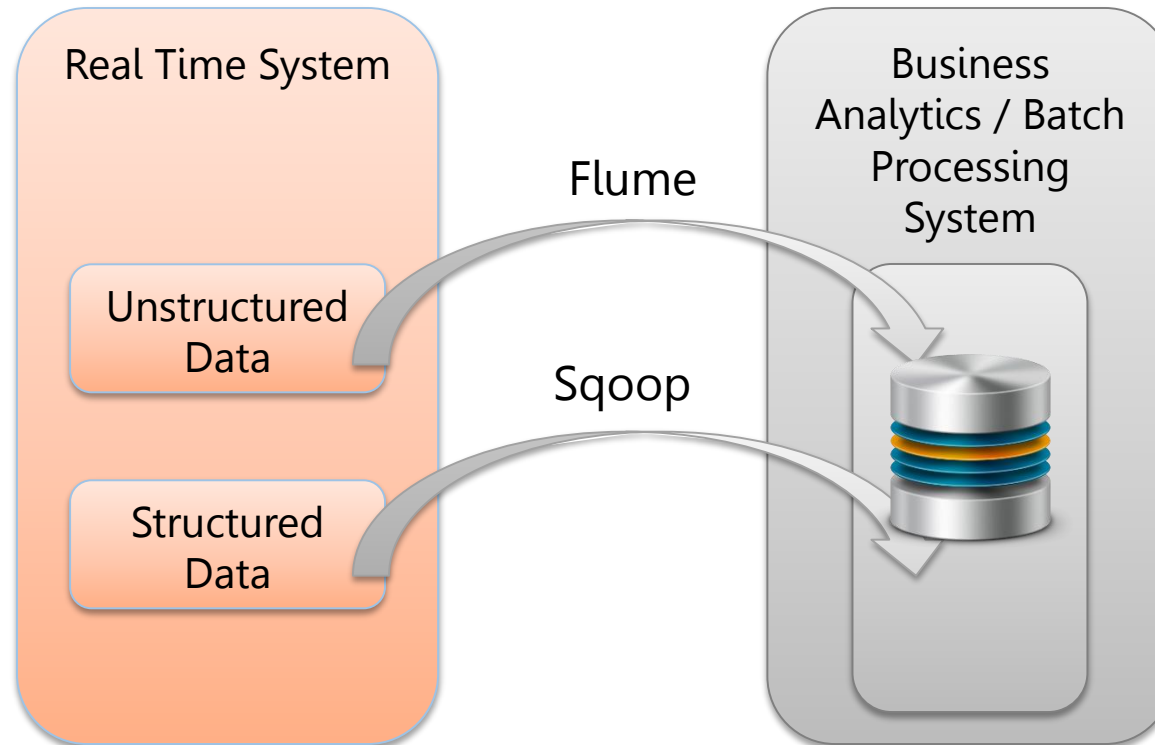
- When Bombay Stock Exchange (the seventh largest stock exchange in the world, in terms of market capitalization) wanted to ramp up / scale up its operations, the company faced major challenges
- These challenges were in terms of exponential growth of data (read big data), need for complex analytics and managing information that was scattered across multiple and monolithic system
- DataMetica (a Mumbai / Pune based big data organization) suggested a 3 phased solution to BSE:
 - In the first phase, they created a POC which demonstrated how a Hadoop based Big data implementation can work for BSE
 - In the second phase, they worked with BSE to pick up the most critical business use cases (which had the maximum ROI for BSE) and implemented them
 - Finally in the third phase they delivered the complete solution in a multi-faced manner for a full fledged implementation
- That's how Hadoop got implemented at BSE in a cost effective and scalable fashion

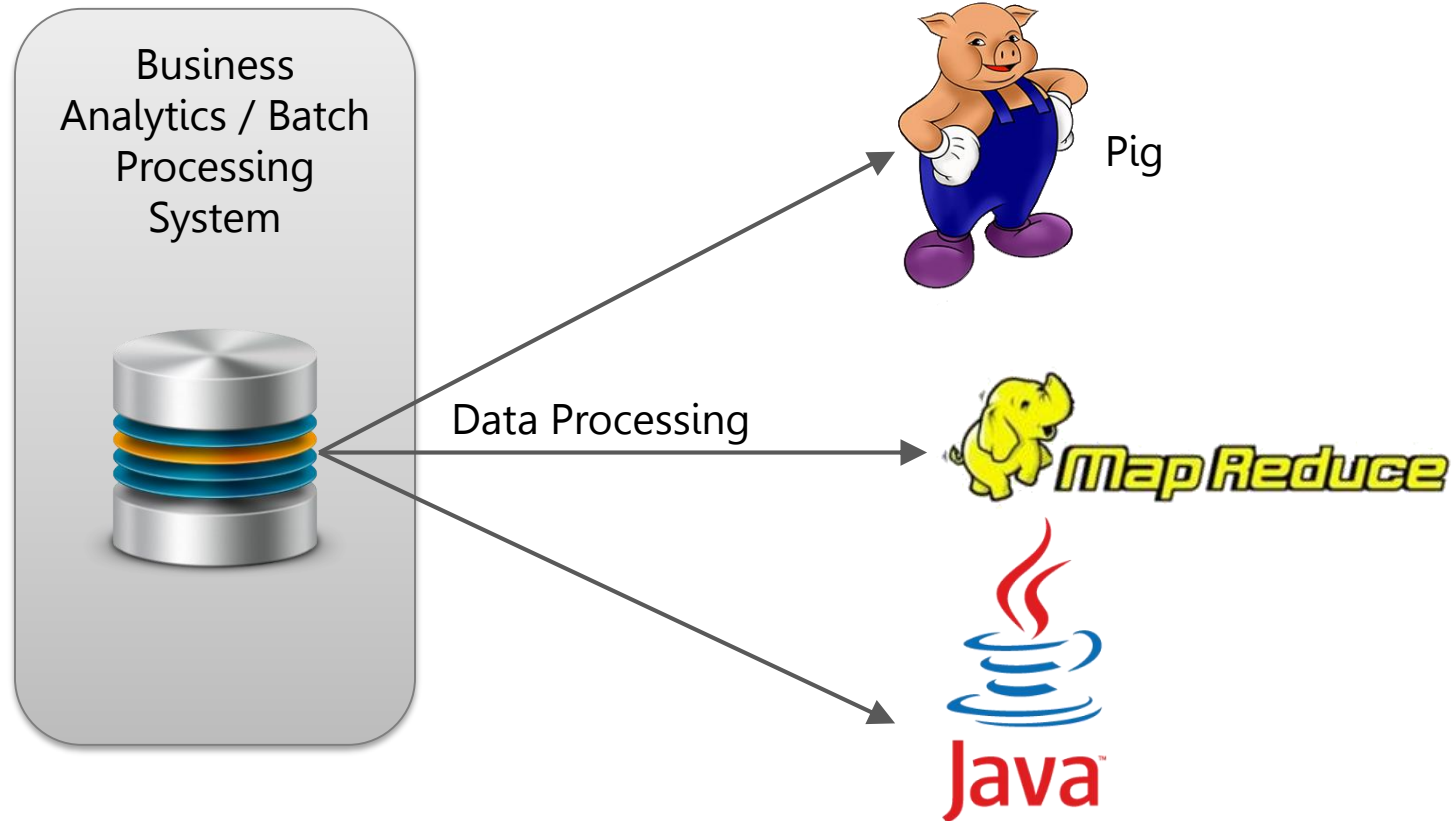


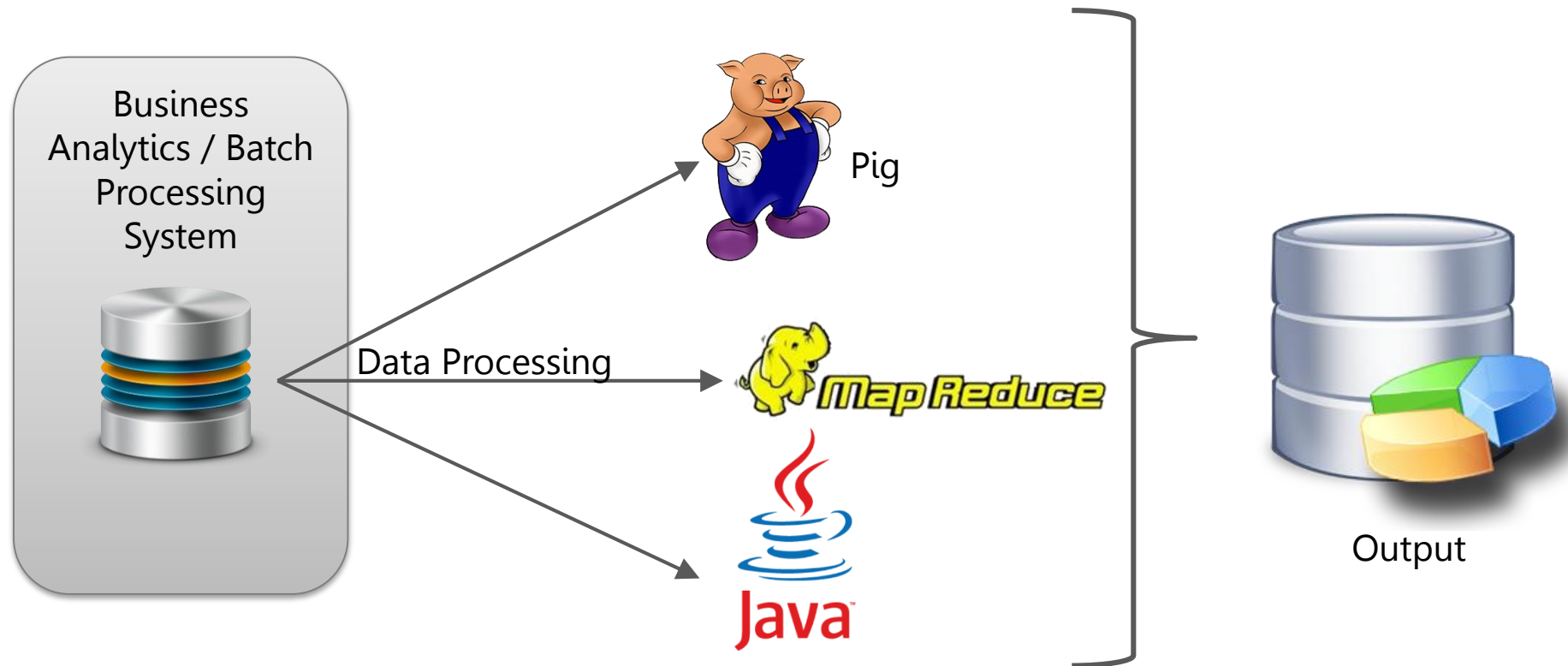
Batch Processing Phase/Life Cycle

- Processing transactions in a group or batch
- Following three phases are common to batch processing or business analytics project, irrespective of the type of data (structured or unstructured)









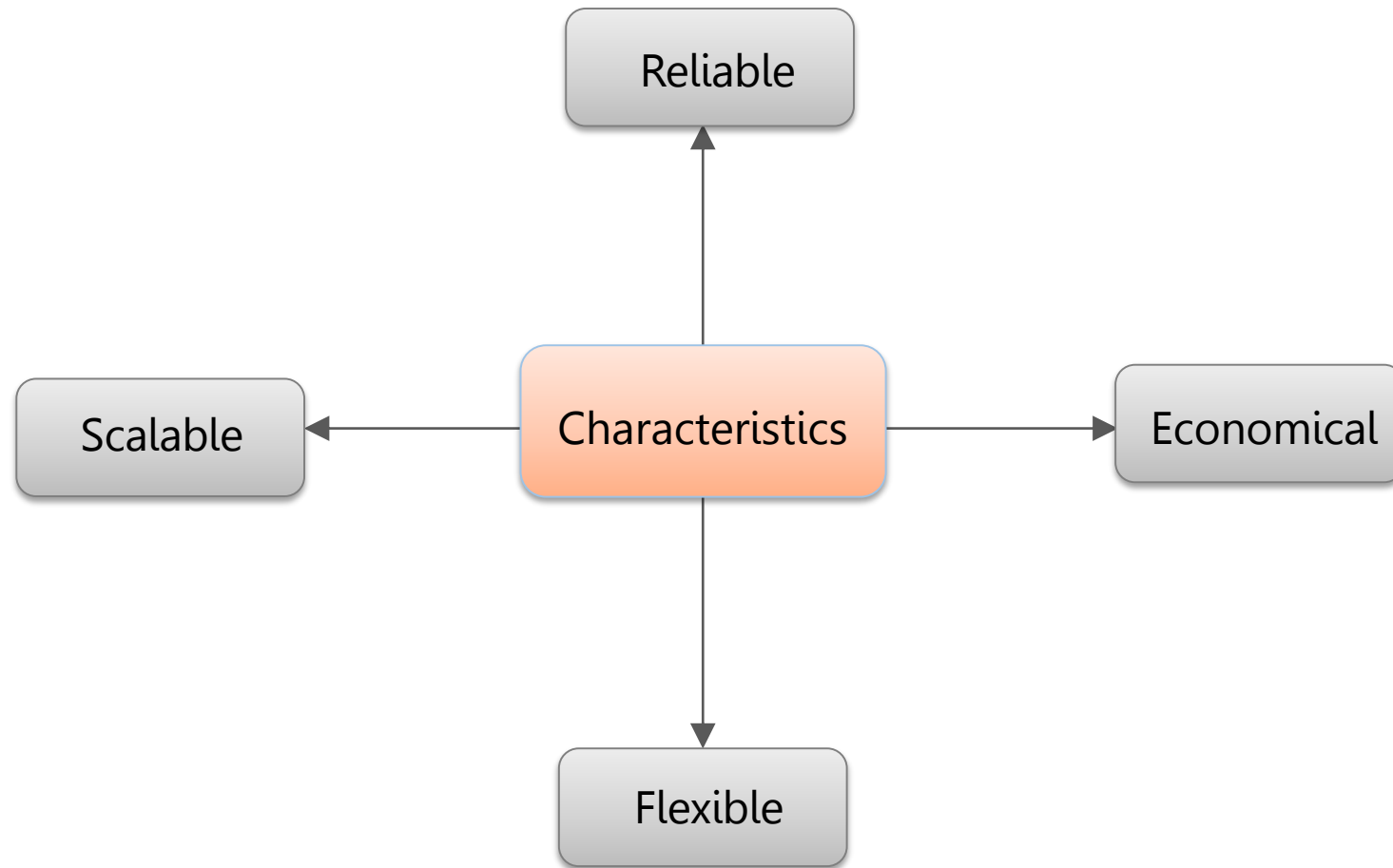
What is Hadoop?

Apache Hadoop is a **framework** that allows the distributed processing of large data sets across clusters of commodity computers using a simple programming mode



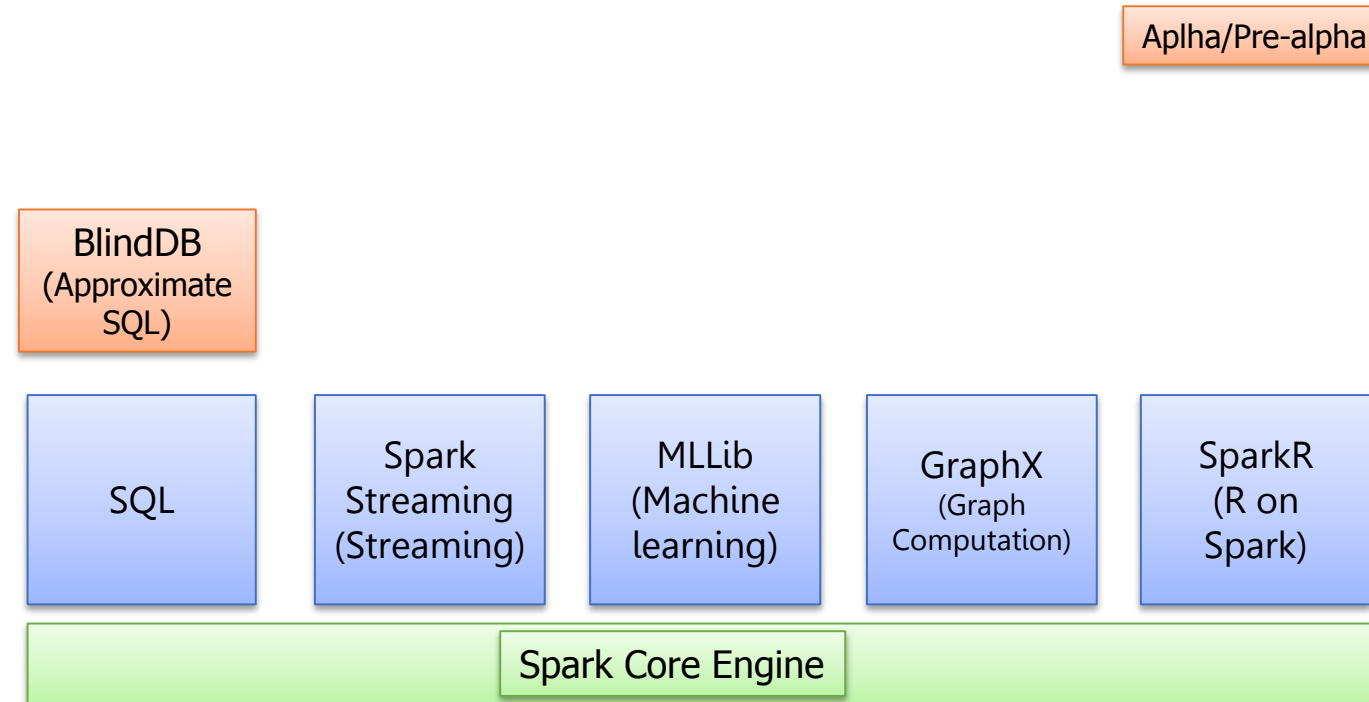
It is an **Open-source Data Management** with scale-out storage and distributed processing

Hadoop Key Characteristics

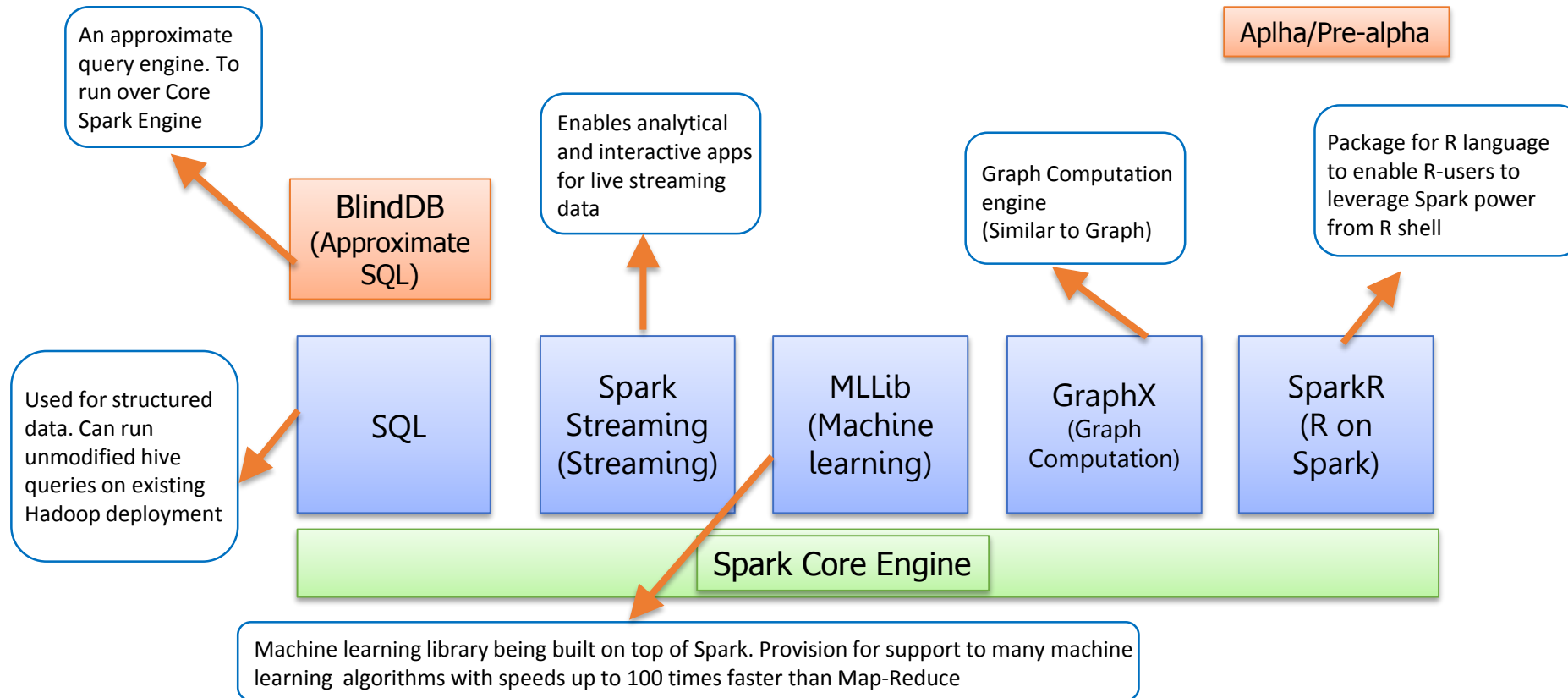


What is Spark?

- **Apache Spark** is a fast and general engine for large-scale data processing
- Apache Spark is a general-purpose cluster in-memory computing system
- It is used for fast data analytics
- It abstracts APIs in Java, Scala and Python, and provides an optimized engine that supports general execution graphs
- Provides various high level tools like Spark SQL for structured data processing, Mlib for Machine Learning and more



Spark Ecosystem (Cont'd)



Spark Ecosystem (Cont'd)

- Spark Core Engine
 - The core engine for entire Spark framework
 - Provides utilities and architecture for other components
- Spark SQL
 - Spark SQL is the newest component of Spark and provides a SQL like interface
 - Used for Structured data
 - Can expose many datasets as tables
 - Spark SQL is tightly integrated with the various spark programming languages like hive
- Spark Streaming
 - Spark Streaming is an extension of the core Spark API that enables scalable, high-throughput, fault-tolerant stream processing of live data streams
 - A good alternative of Storm

Spark Ecosystem (Cont'd)

▶ BlinkDB

- An approximate query engine. To run over Core Spark Engine
- Accuracy trade-off for response time

▶ MLlib

- Machine learning library being built on top of Spark
- Provision for support to many machine learning algorithms with speeds up to 100 times faster than Map-Reduce
- Mahout is also being migrated to MLlib

▶ GraphX

- Graph Computation engine (Similar to Giraph)
- Combines data-parallel and graph-parallel concepts

▶ SparkR

- Package for R language to enable R-users to leverage Spark power from R shell

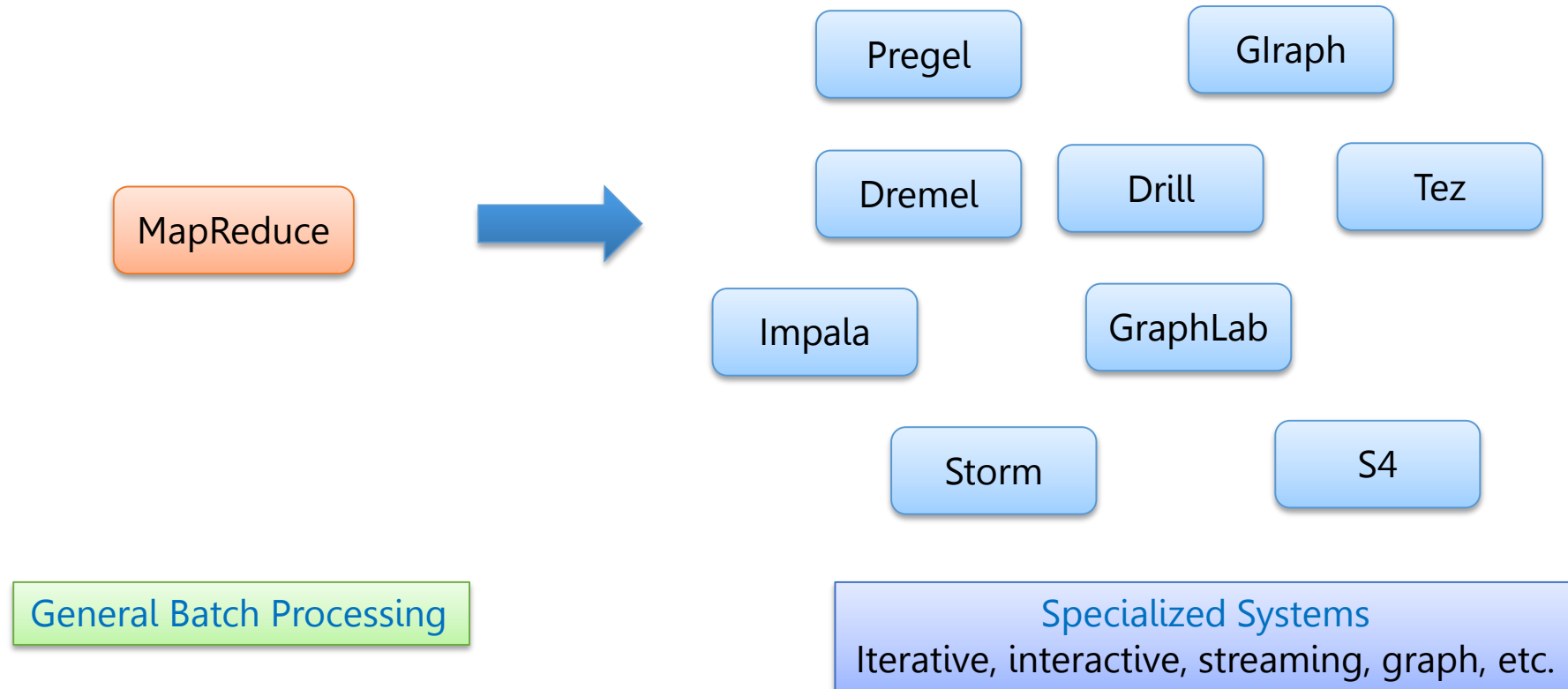
Why Spark?

- Spark exposes a simple programming layer which provides powerful caching and disk persistence capabilities
- The Spark framework can be deployed through Apache Mesos, Apache Hadoop via Yarn, or Spark's own cluster manager
- Spark framework is polyglot – Can be programmed in several programming languages (Currently Scala, Java and Python supported)
- Has super active community
- Spark fits well with existing Hadoop ecosystem
 - Can be launched in existing **YARN** Cluster
 - Can fetch the data from Hadoop 1.0
 - Can be integrated with Hive

Brief History: M/R Limitations

- Map Reduce is a very powerful programming paradigm, but it has some limitations:
 - Difficult to Program an algorithm directly in Native Map Reduce
 - Performance bottlenecks, specifically for small batch not fitting the use cases
 - Many categories of algorithms not supported (e.g. iterative algorithms, asynchronous algorithms etc.)
- In short, MR doesn't compose well for large applications
- We are forced to take **hybrid** approaches many times
- Therefore, many specialized systems evolved over a period of time as workarounds

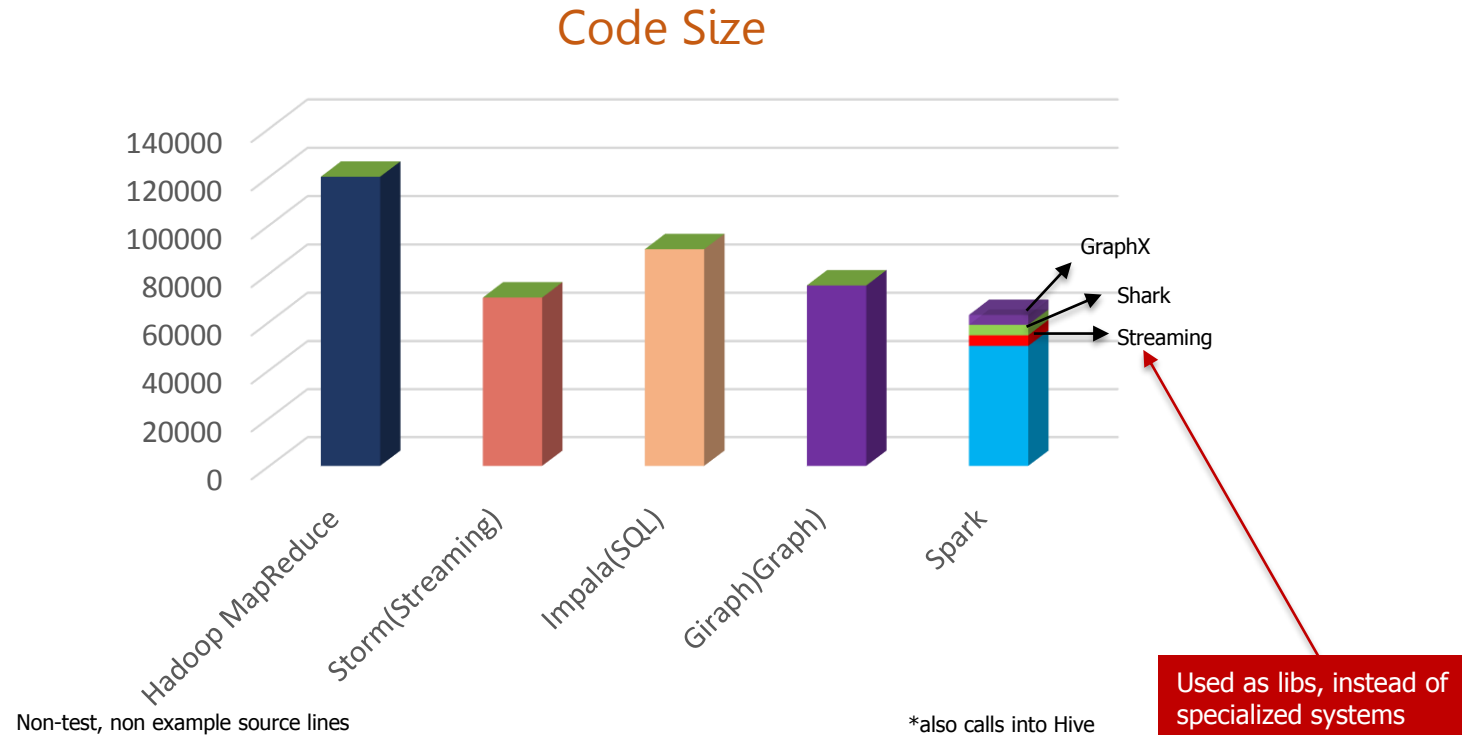
Brief History: Evolution of Specialized Systems



Brief History: Spark

- Unlike other evolved specialized systems, Spark's design goal is to generalize Map Reduce concept to support new apps within same engine
- Two reasonably small additions are enough to express the previous models:
 - Fast data sharing (For Faster Processing)
 - General DAGs (For Lazy Processing)
- This allows for an approach which is more efficient for the engine, and much simpler for the end users

Brief History: Spark Key Points



The State of Spark, and where we're going next
Matei Zaharia
Spark Summit(2013)
you.be/nU6v02EJAb4

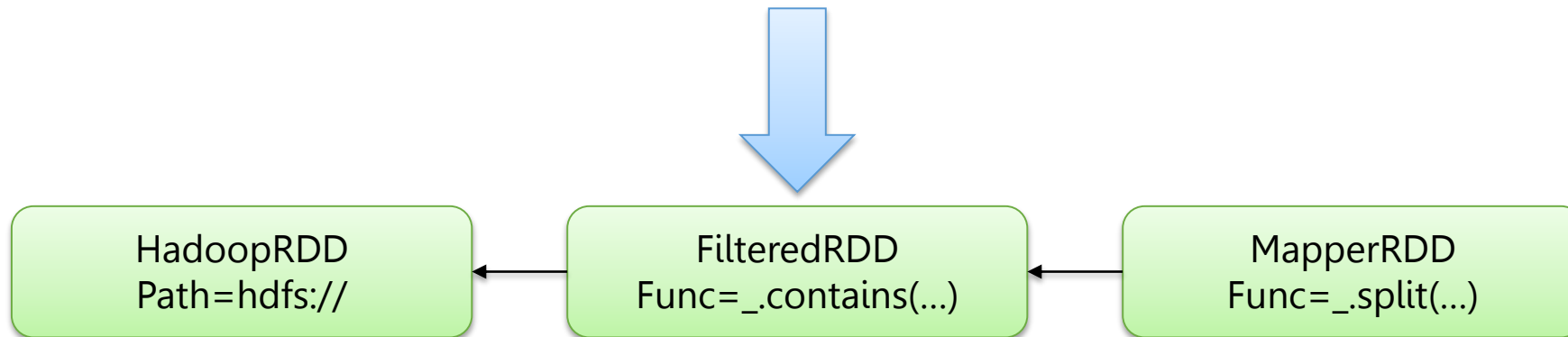
Brief History: Spark Key Points (Cont'd)

RDD Fault Tolerance

RDDs track the series of transformation used to build them (their lineage) to recompute lost data

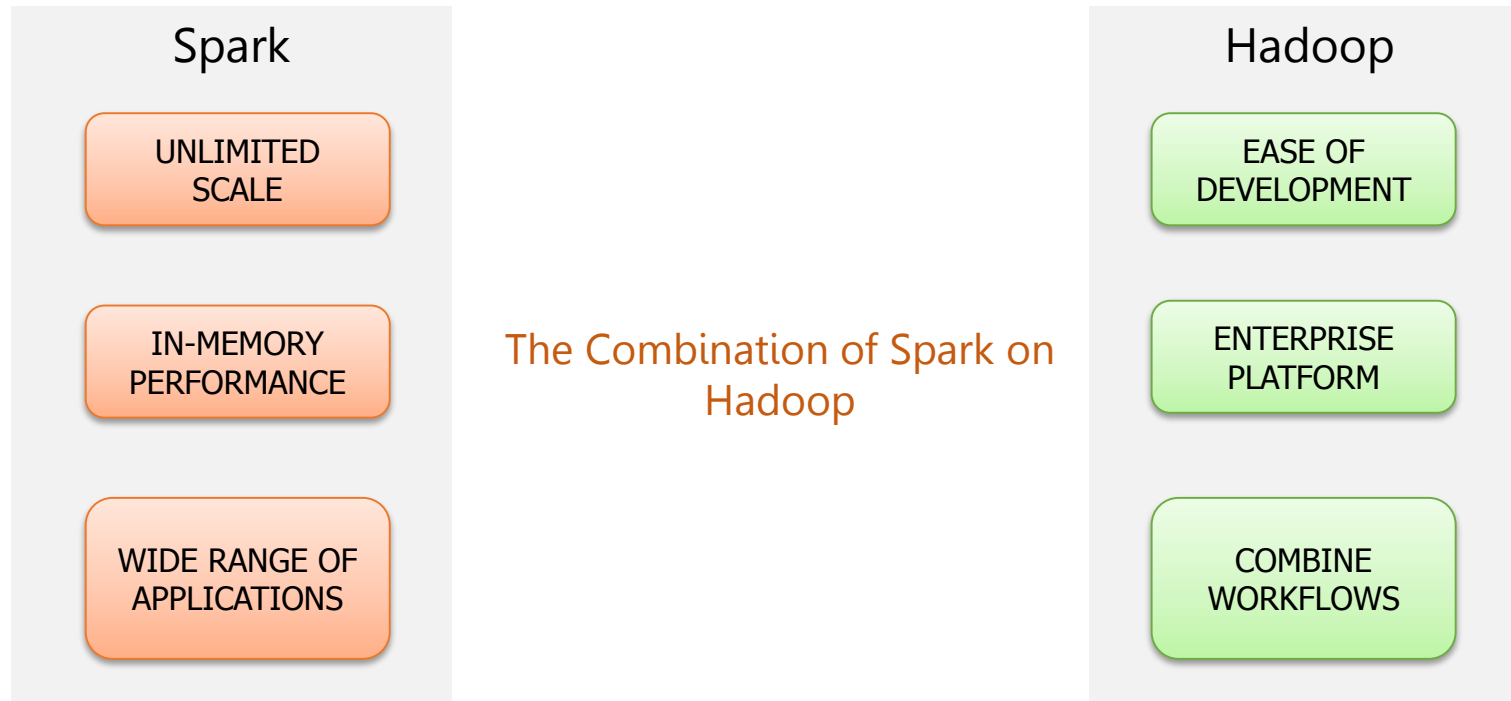
Example:

```
messages=textFile(...).filter(_.contains("error"))  
                        .map(_.split('\t')(2))
```





Operational Applications Augmented by In-Memory Performance



Simple Spark Apps: Word Count

This simple program provides a good test case for parallel processing, since it:

- Requires a minimal amount of code
- Demonstrates use of both symbolic and numeric values
- Isn't many steps away from search indexing

```
val f = sc.textFile("README.md")  
  
val wc = f.flatMap(l => l.split(" ")).map(word => (word, 1)).reduceByKey(_ + _)  
  
wc.saveAsTextFile("wc_out")
```


Using Hadoop as Storage

▶ Spark can use Hadoop as Storage

- Spark is NOT limited to HDFS only for its storage needs
- HDFS provides distributed storage of large datasets
- High Availability is assured natively through HDFS
- No extra software installation is required
- Compatible with Hadoop 1.x also. Using HDFS as storage doesn't require Hadoop 2.x
- Data Loss during computation is handled by HDFS itself

Using Hadoop as Execution Engine

- ▶ Spark can use Hadoop as execution engine
 - Spark can be integrated with Yarn for it's execution
 - Spark can be used with other engines (like Mesos, Spark Cluster manager) also
 - Yarn integration automatically provides processing scalability to Spark
 - Spark needs Hadoop 2.0+ versions in order to use it for execution
 - Every node in Hadoop cluster need Spark also to be installed
 - Using Hadoop cluster for Spark processes, requires RAM upgrading of data nodes
 - The integration distribution of Spark is quite new and still in the process of stabilization



thank
you!