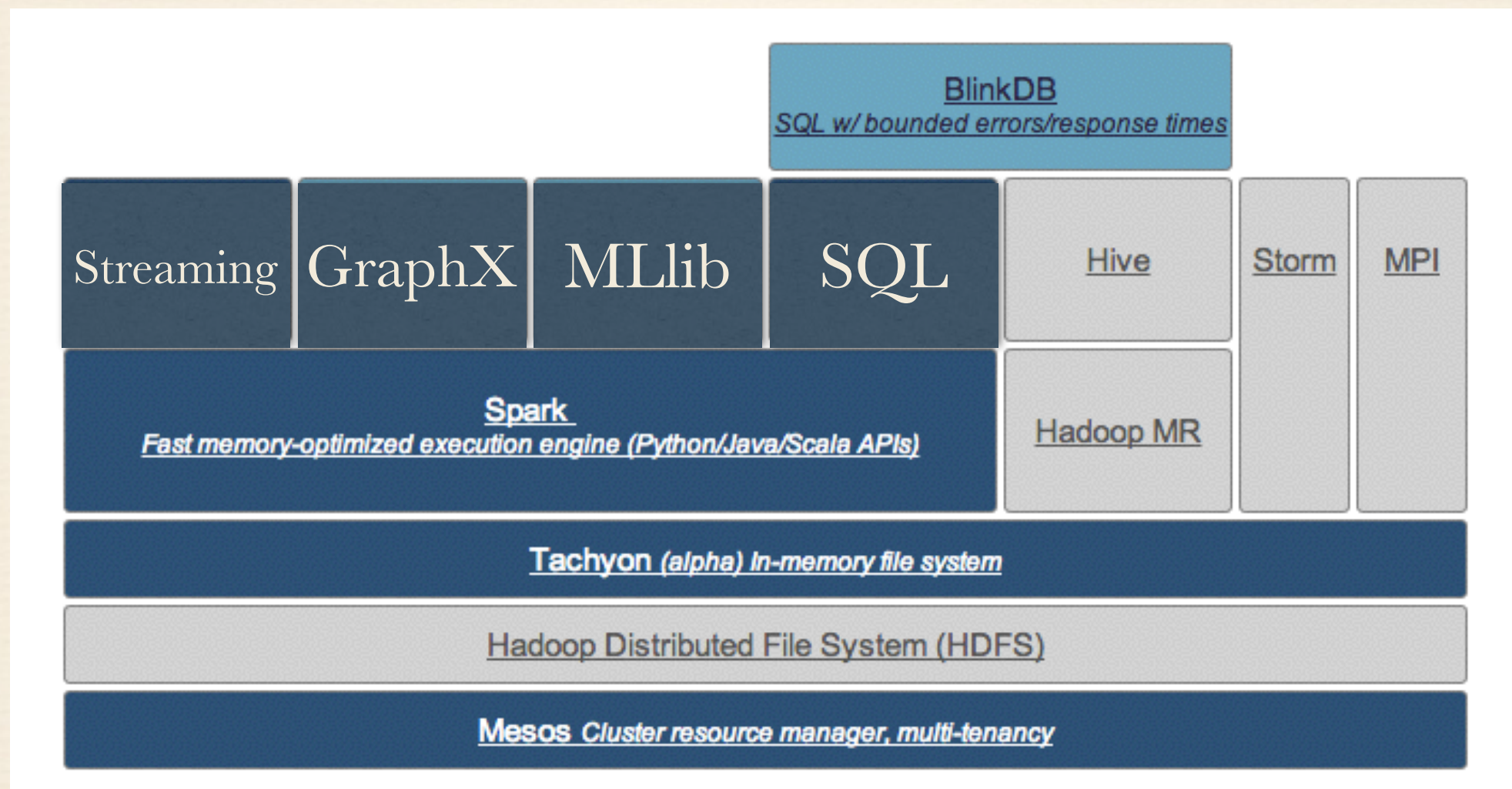# Spark Ecosystem
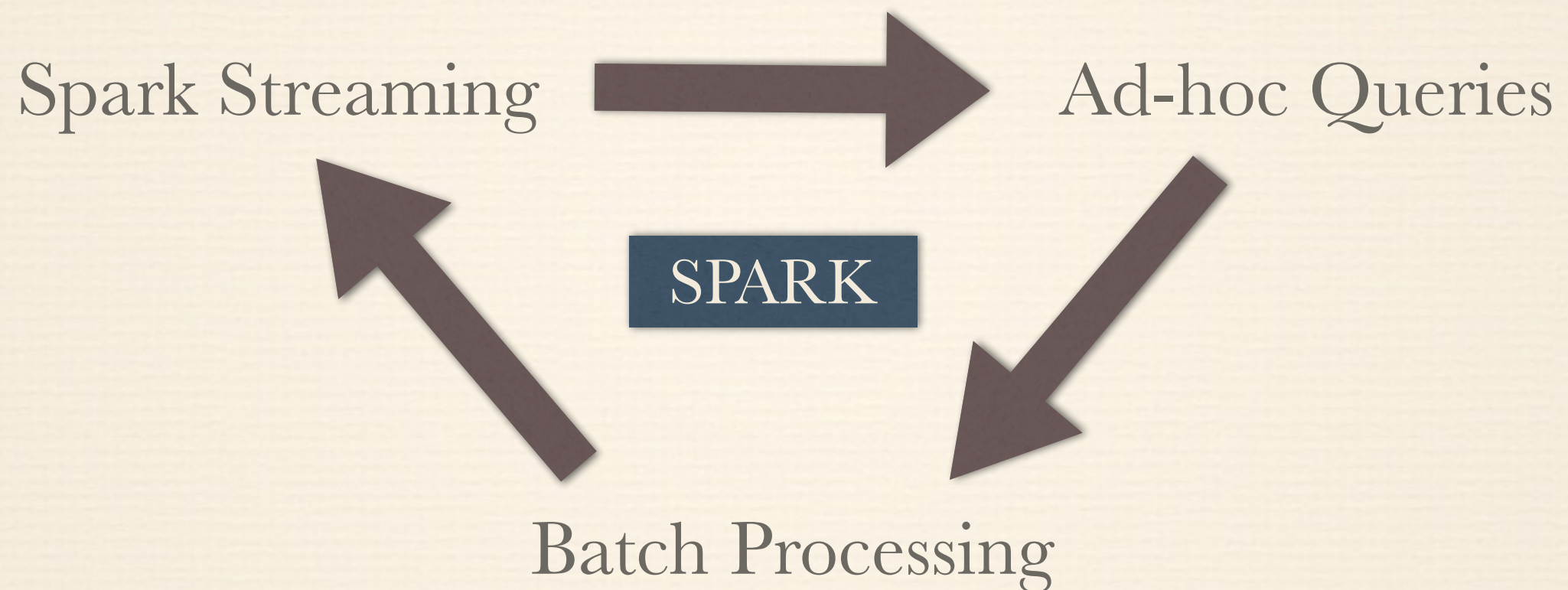
陈超 @CrazyJvm
*Spark Meetup @Hangzhou*
*2014.08.31*

# What is Spark

- Apache Spark is a fast and general engine for large-scale data processing.

- Speed

- Ease of Use

- Generality

- Integrated with Hadoop

# BDAS

# one stack to rule them all

Spark Streaming → Ad-hoc Queries

SPARK

Batch Processing

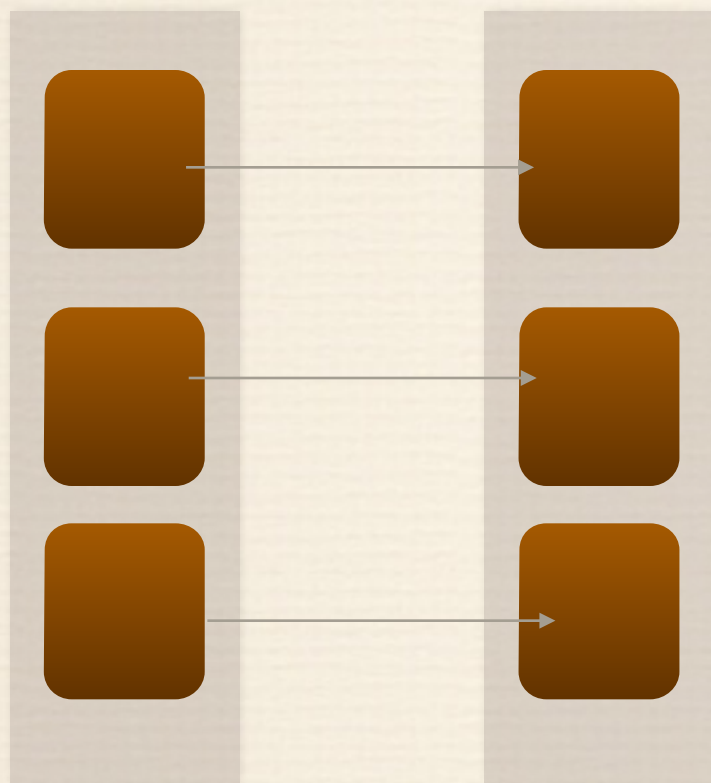# Key Concept-RDD

- A list of partitions

- A function for computing each split

- A list of dependencies on other RDDs

- Optionally, a Partitioner for key-value RDDs

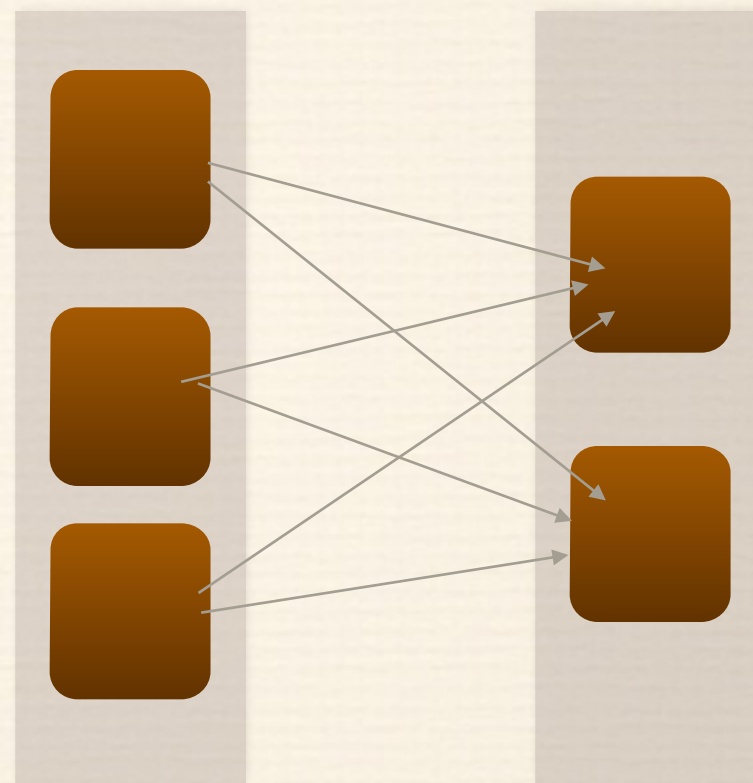- Optionally, a list of preferred locations to compute each split on

# Key Concept-Lineage

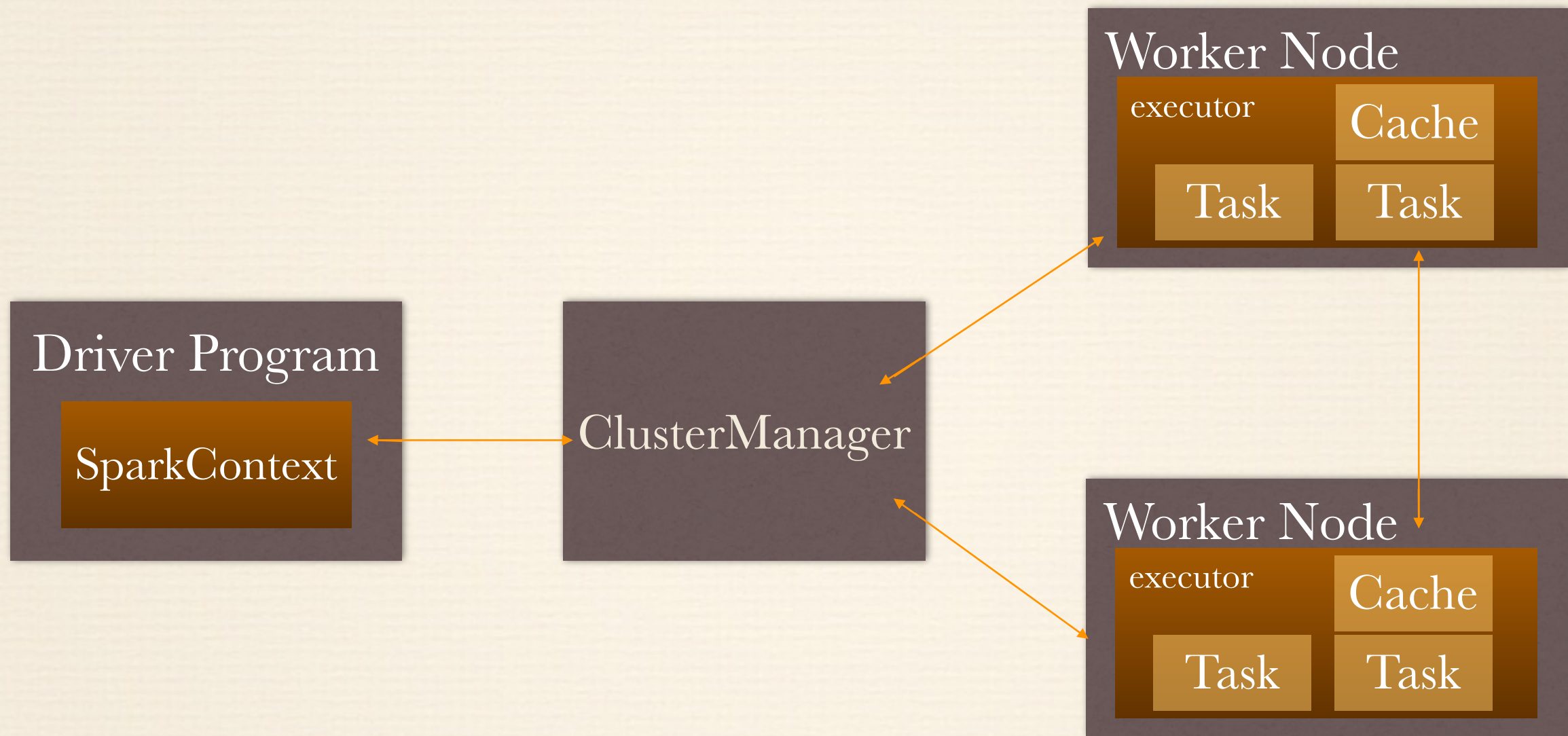# Key Concept-Dependency



Narrow Dependency

Wide Dependency

# Key Concept-ClusterManager
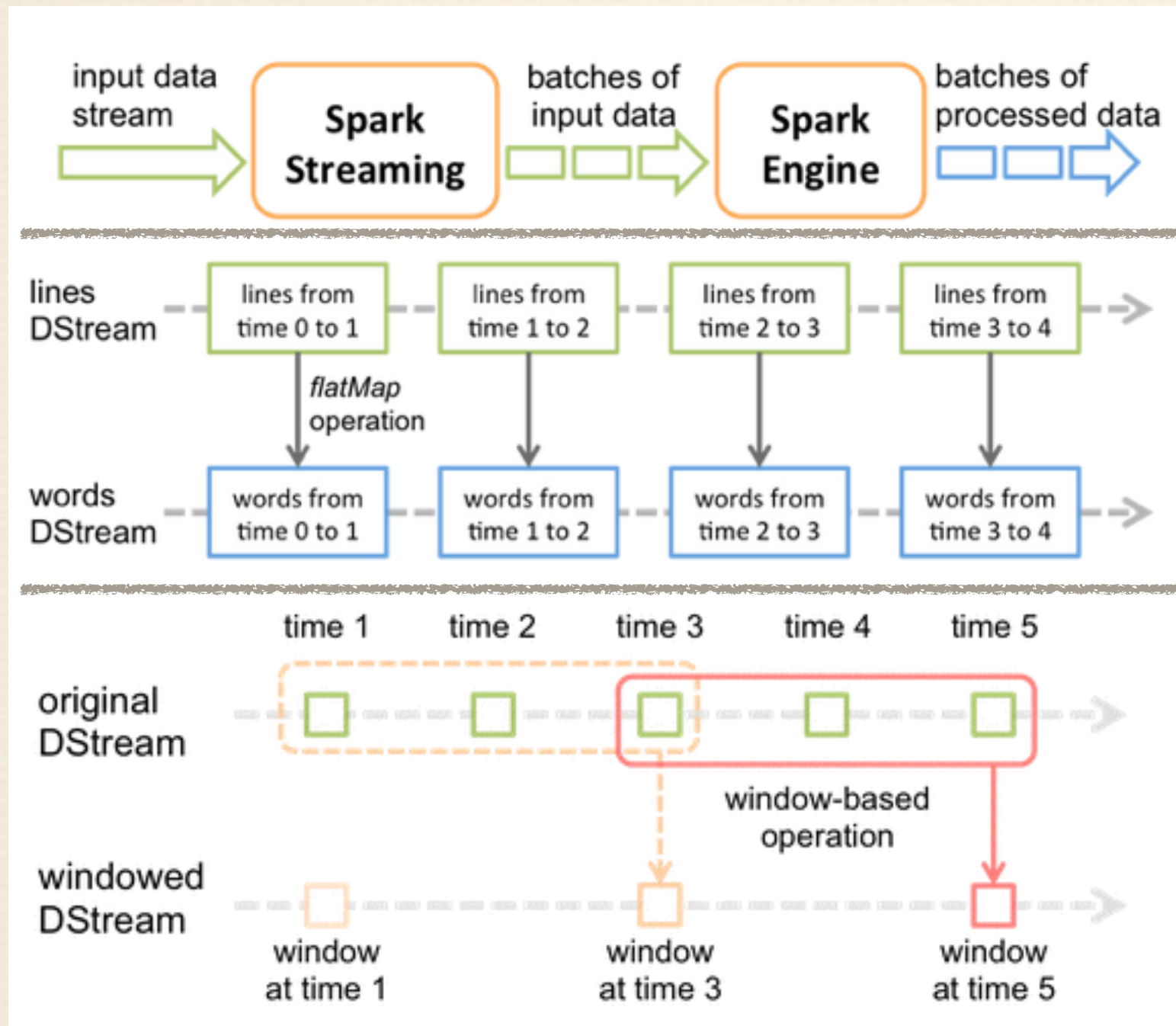
- Local

- Standalone

- Yarn

- Mesos

# Cluster Overview

**Driver Program**

SparkContext

**ClusterManager**

**Worker Node**

executor

Cache

Task | Task

**Worker Node**

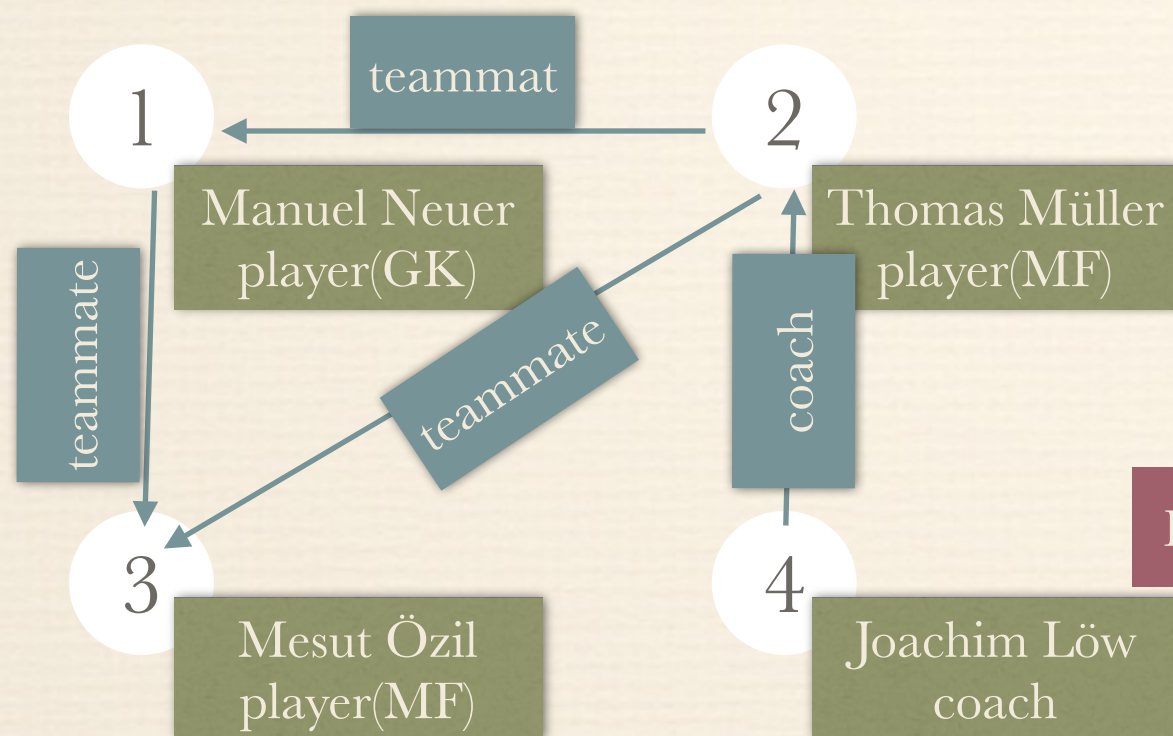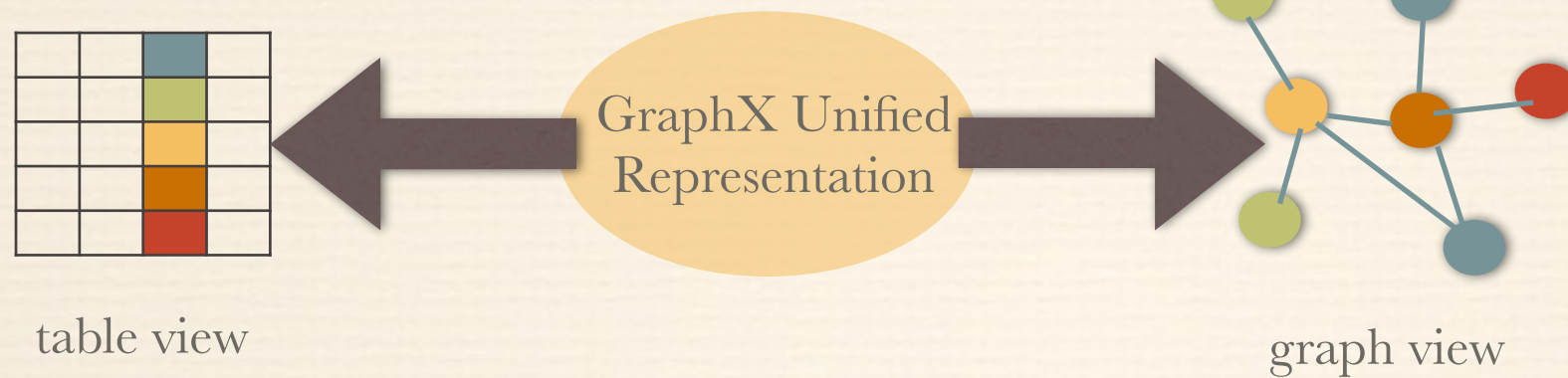executor

Cache

Task | Task

# Spark Streaming

- ❖ mini-batch

# MLlib

- Spark implementation of some common machine learning algorithms and utilities

- classification

- regression

- clustering

- collaborative filtering

- dimensionality reduction
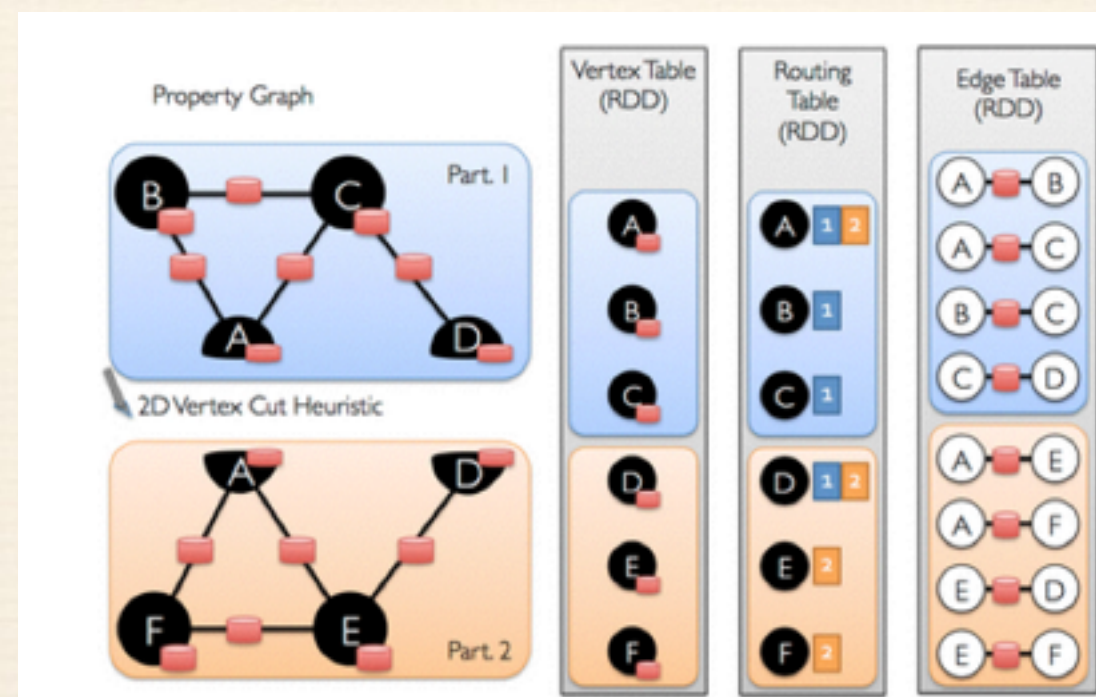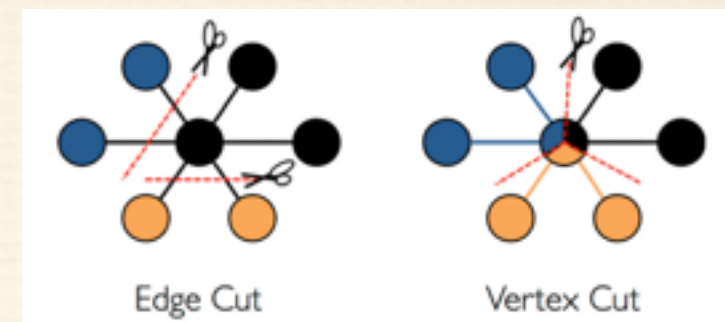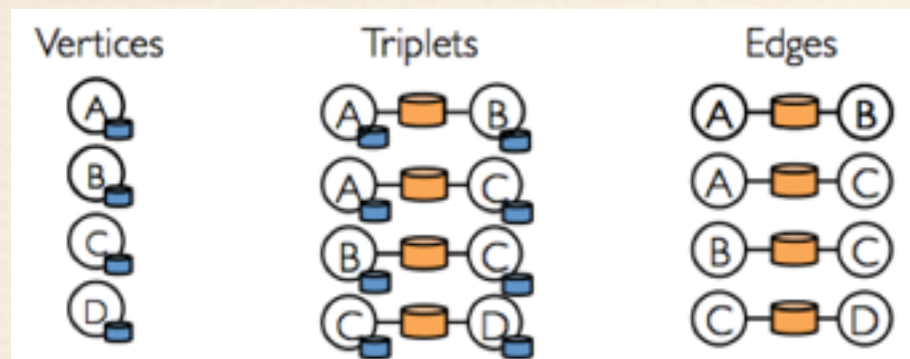
Sparse vector support

Evaluation support

# GraphX



table view

graph view

**1** — Manuel Neuer player(GK)

teammat

**2** — Thomas Müller player(MF)

teammate

coach

teammate

**3** — Mesut Özil player(MF)

**4** — Joachim Löw coach

**Vertex Table**

| id | Ver |
|----|-----|
| 1 | (Manuel Neuer,player) |
| 2 | (Thomas Müller,player) |
| 3 | (Mesut Özil,player) |
| 4 | (Joachim Löw,coach) |

**Edge Table**

| SrcId | DstId | Property(E) |
|-------|-------|-------------|
| 2 | 1 | teammate |
| 2 | 3 | teammate |
| 1 | 3 | teammate |
| 4 | 2 | coach |

GraphX Unified Representation

# GraphX

# Spark SQL

| Meta Store | HiveQL | UDFs | SerDes |
| --- | --- | --- | --- |

Spark SQL

Apache Spark

| BI Tools | ..... |
| --- | --- |

JDBC/ODBC

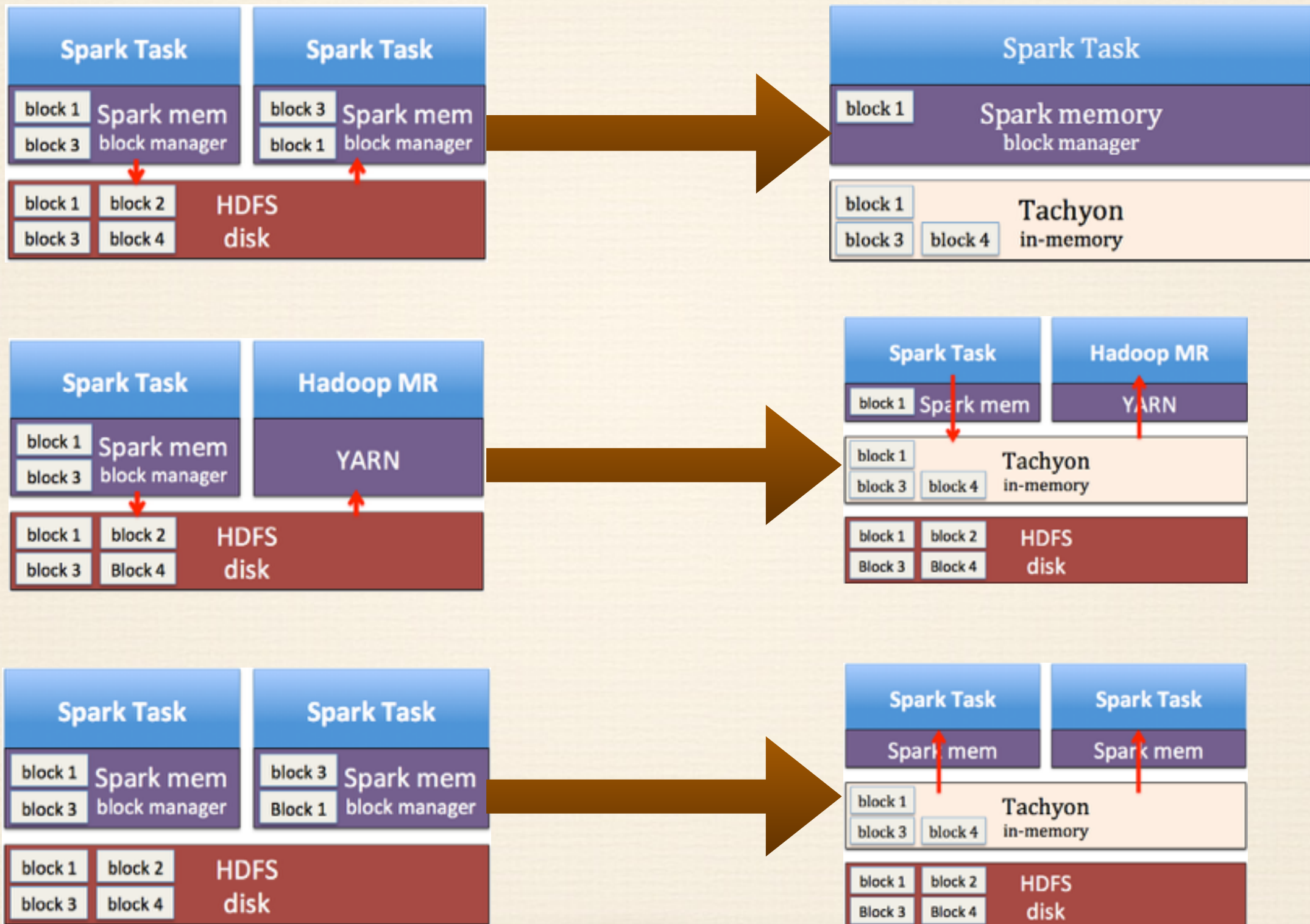Spark SQL

Apache Spark

# Spark SQL

- ❖ Data Sources
  - RDDs/Parquet Files/JSON Datasets/Hive Table
- ❖ DSL
- ❖ JDBC Server

# Shark

- ❖ Mission Completed!!!

# Tachyon

# Tachyon

| | | | | | | |
|---|---|---|---|---|---|---|
| MR | Spark | Tez | Shark | GraphX | Impala | …… |

**Tachyon**

| | | | | | | |
|---|---|---|---|---|---|---|
| HDFS | S3 | Localfs | Cluster fs | NFS | Ceph | …… |

# SparkR

❖ R + RDD = RRDD

RDDs as Distributed Lists

# QA & Thanks

weibo:@CrazyJvm

wechat public account : ChinaScala