

# Spark分布式计算



美柚@郭厦同

[alien.xt.xm@gmail.com](mailto:alien.xt.xm@gmail.com)

<http://github.com/alienxt>

# 目录

- 1 Spark介绍
- 2 Spark核心及基础概念
- 3 Spark运行架构
- 4 Spark SQL核心与实践
- 5 Spark Streaming核心与实践



# Spark是什么？

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

Apache Spark is an open source cluster computing system that aims to make data analytics fast.

# 为什么使用Spark ?



## 快速的

处理速度快于Hadoop ,  
10x - 100x倍



## 多语言支持

可以各种语言开发 , 支持Java,  
Scala, Python, R



## 通用性

结合了sql , mllib , graphx ,  
streaming库



## 分布式可水平扩展

master/slave结构 , master管  
理worker , master支持ha ,  
worker可扩展



## Runs Everywhere

数据源丰富  
随地可部署



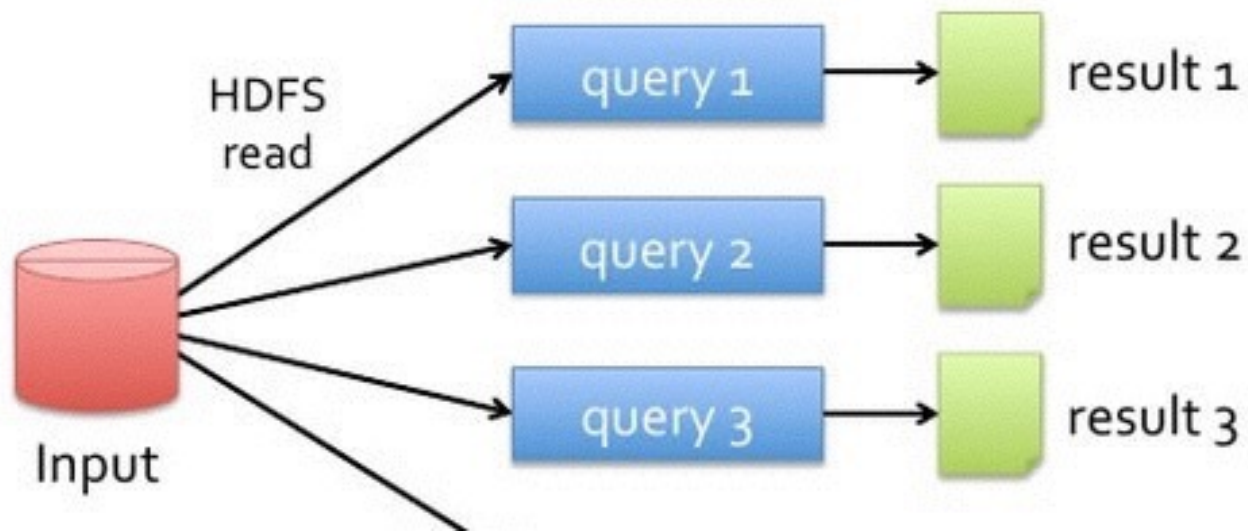
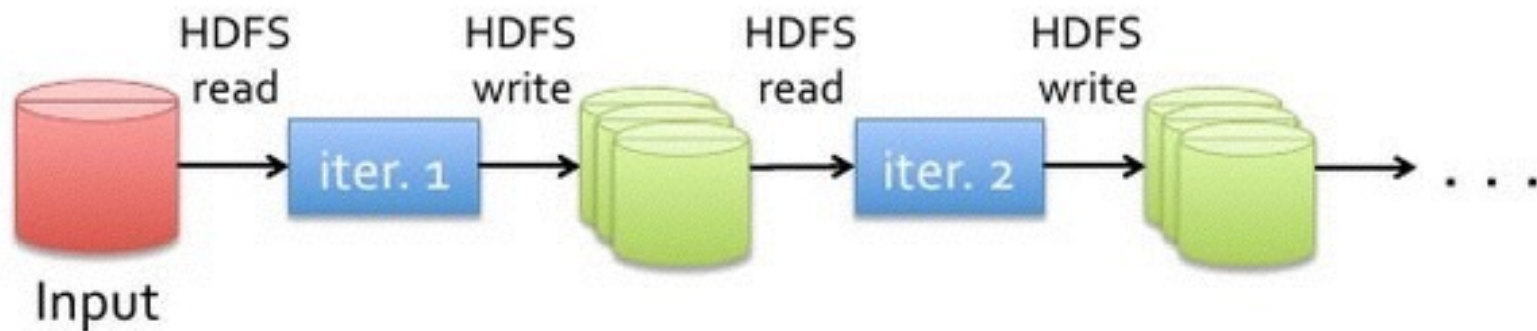
## 本地模式

Spark有一个本地模式 , 可以  
在开发和测试中 , 完全模拟集  
群

# Spark vs Hadoop

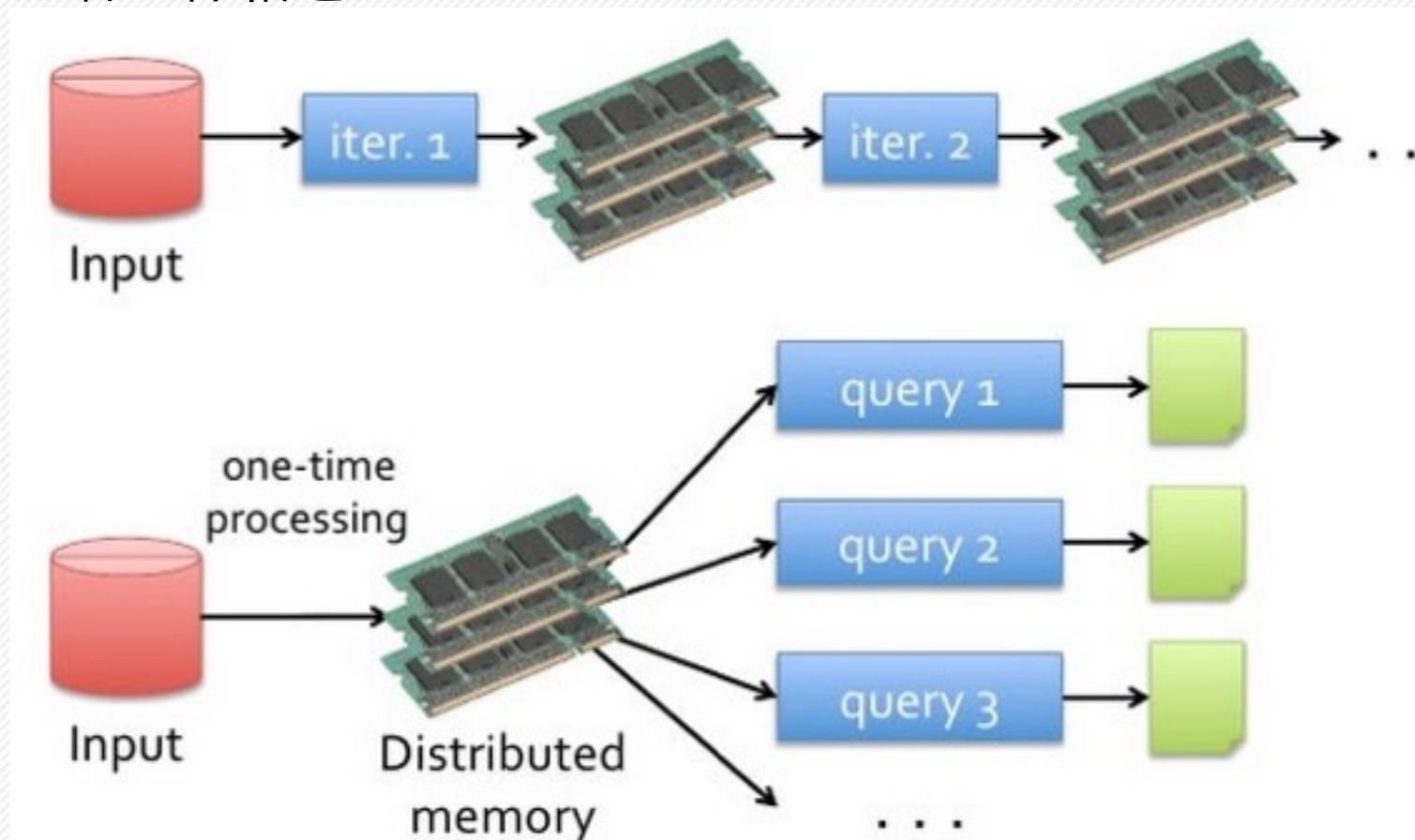
## Hadoop缺点

- 计算的中间结果，保存在磁盘，产生大量的io，降低计算速度



# Spark为什么快？

- 减少磁盘io
- DAG（统筹性运算）
- 各组件相通





# Spark核心及基础概念

1. RDD
2. 基础概念
3. 实践



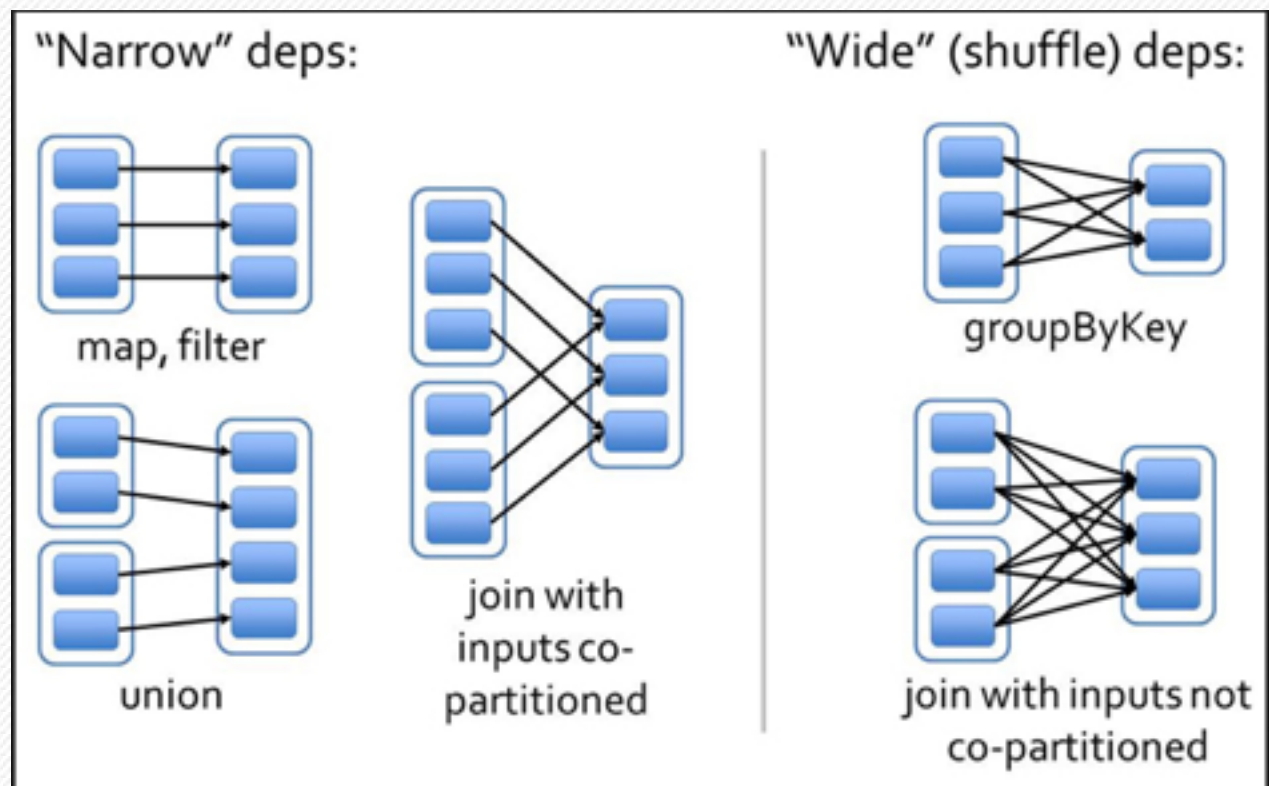


# RDD是什么？

RDD：Spark的核心概念是RDD (resilient distributed dataset)，指的是一个只读的，可分区的分布式数据集，这个数据集的全部或部分可以缓存在内存中，在多次计算间重用。

特性：

- 分区性
- 只读性
- 函数性
- 容错性



<b><i>map(func)</i></b>	通过将源数据中的每个元素传递给 <b><i>func</i></b> 函数来返回一个新的 <b><i>RDD</i></b>
<b><i>filter(func)</i></b>	从源数据中选出那些传递给 <b><i>func</i></b> 函数后返回 <b><i>true</i></b> 的元素构成的新 <b><i>RDD</i></b>
<b><i>flatMap(func)</i></b>	与 <b><i>map</i></b> 类似，但是每个输入元素可以被映射为 <b><i>0</i></b> 或更多的输出元素（所以 <b><i>func</i></b> 应该返回一个序列而不是单个元素）
<b><i>mapPartitions(func)</i></b>	与 <b><i>map</i></b> 类似, 但是在 <b><i>RDD</i></b> 的每个分区上分开运算, <b><i>mapPartitions()</i></b> 的输入函数是应用于每个分区
<b><i>mapPartitionsWithIndex(func)</i></b>	与 <b><i>mapPartitions</i></b> 类似, 但是提供了一个整型参数给函数 <b><i>func</i></b> , 代表了分区的索引。所以函数 <b><i>func</i></b> 的类型必须为 <b><i>(Int, Iterator&lt;T&gt;) =&gt; Iterator&lt;U&gt;</i></b> ，当 <b><i>RDD</i></b> 的类型是 <b><i>RDD[T]</i></b> .
<b><i>sample(withReplacement, fraction, seed)</i></b>	对数据的一部分进行采样, 可选的替换参数, 使用一个给定的随机数发生器种子。
<b><i>union(otherDataset)</i></b>	返回一个新的数据集，包含了参数与源数据的并集。
<b><i>intersection(otherDataset)</i></b>	返回一个新的 <b><i>RDD</i></b> 与参数的交集。
<b><i>distinct([numTasks])</i></b>	返回一个新的数据集包含了与源数据集的不同元素
<b><i>groupByKey([numTasks])</i></b>	
<b><i>reduceByKey(func, [numTasks])</i></b>	
<b><i>sortByKey([ascending], [numTasks])</i></b>	
<b><i>join(otherDataset, [numTasks])</i></b>	

<b><i>reduce(func)</i></b>	聚集，传入的函数是两个参数输入返回一个值，这个函数必须是满足交换律和结合律的
<b><i>collect()</i></b>	一般在 <b><i>filter</i></b> 或者足够小的结果的时候，再用 <b><i>collect</i></b> 封装返回一个数组
<b><i>count()</i></b>	返回的是 <b><i>dataset</i></b> 中的 <b><i>element</i></b> 的个数
<b><i>first()</i></b>	返回的是 <b><i>dataset</i></b> 中的第一个元素
<b><i>take(n)</i></b>	返回前 <b><i>n</i></b> 个 <b><i>elements</i></b> ，这个在 <b><i>driverprogram</i></b> 返回的
<b><i>takeSample(withReplacement, num, seed)</i></b>	抽样返回一个 <b><i>dataset</i></b> 中的 <b><i>num</i></b> 个元素，随机种子 <b><i>seed</i></b>
<b><i>saveAsTextFile (path)</i></b>	把 <b><i>dataset</i></b> 写到一个 <b><i>textfile</i></b> 中，或者 <b><i>hdfs</i></b> ，或者 <b><i>hdfs</i></b> 支持的文件系统中， <b><i>Spark</i></b> 把每条记录都转换为一行记录，然后写到 <b><i>file</i></b> 中
<b><i>saveAsSequenceFile(path)</i></b>	只能用在 <b><i>key-value</i></b> 对上，然后生成 <b><i>SequenceFile</i></b> 写到本地或者 <b><i>hadoop</i></b> 文件系统
<b><i>countByKey()</i></b>	返回的是 <b><i>key</i></b> 对应的个数的一个 <b><i>map</i></b> ，作用于一个 <b><i>RDD</i></b>
<b><i>foreach(func)</i></b>	对 <b><i>dataset</i></b> 中的每个元素都使用 <b><i>func</i></b>

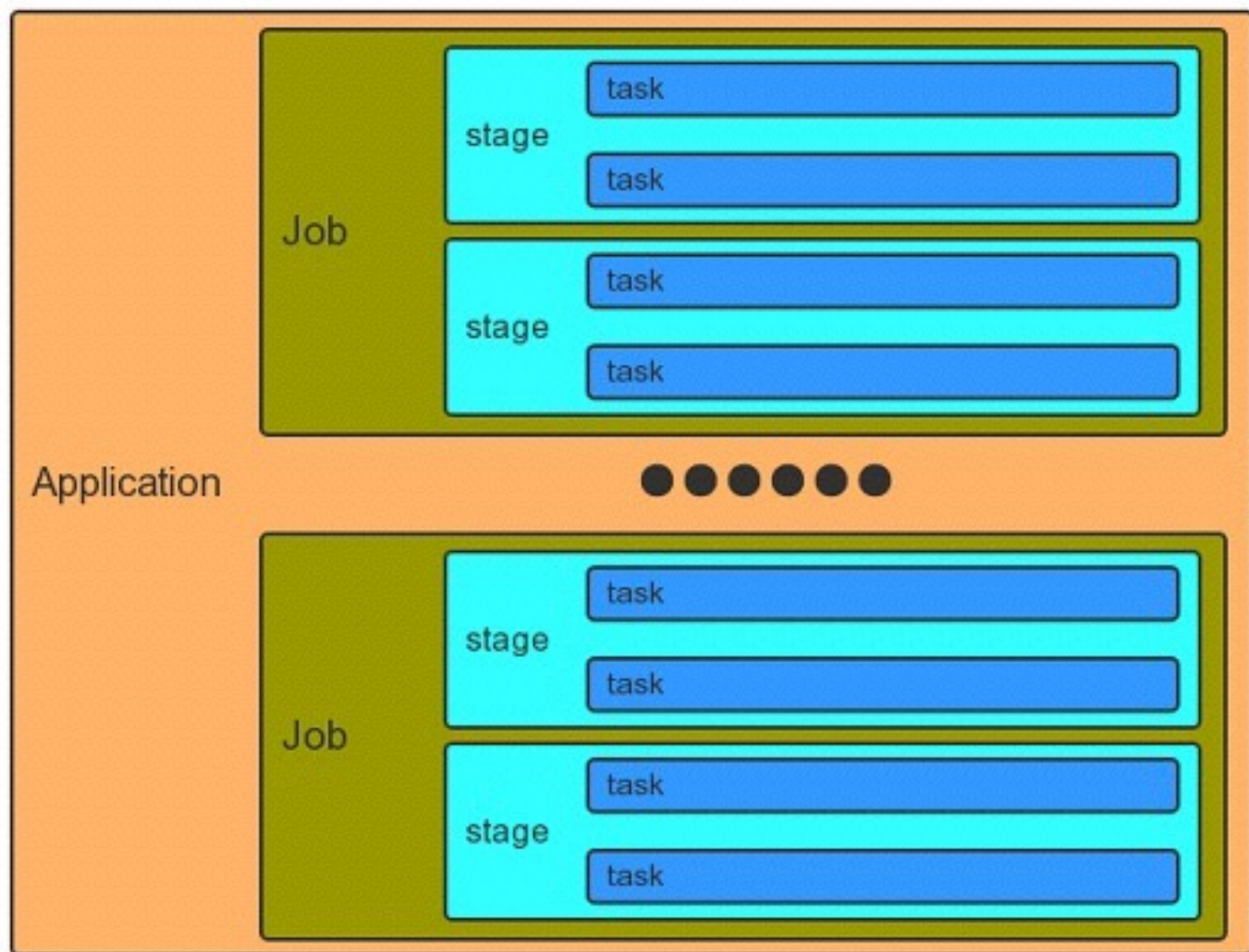
<b><i>persist()</i></b>	<b><i>persist</i></b> 方法手工设定 <i>StorageLevel</i> 来满足工程需要的存储级别
<b><i>cache()</i></b>	<i>RDD</i> 的 <i>cache()</i> 方法其实调用的就是 <b><i>persist</i></b> 方法，缓存策略均为 <b><i>MEMORY_ONLY</i></b> ;
<b><i>unpersist()</i></b>	清空缓存

StorageLevel

<https://github.com/apache/spark/blob/v1.6.1/core/src/main/scala/org/apache/spark/storage/StorageLevel.scala>

# Spark的基本概念

- **Application:** 基于Spark的用户程序，包含了一个driver program和多个executor
- **Driver Program:** 运行application的main函数，并且创建SparkContext
- **Executor:** application中运行在worker node上的一个进程，该进程负责运行task，并且将数据存在内存或磁盘中，executor是独立的进程
- **Cluster Manager:** 在集群上获取资源的外部服务（如standalone, mesos, yarn, cloud）
- **Worker node:** 集群中负责运行运行executor
- **Task:** 在executor中的工作单元，线程
- **Job:** 包含多个task组成的计算，由action操作催生
- **Stage:** 每个job会拆分成多组的task，每组task称为一个stage，也称为taskset
- **RDD:** 基本计算单元，可以通过一系列的算子进行操作

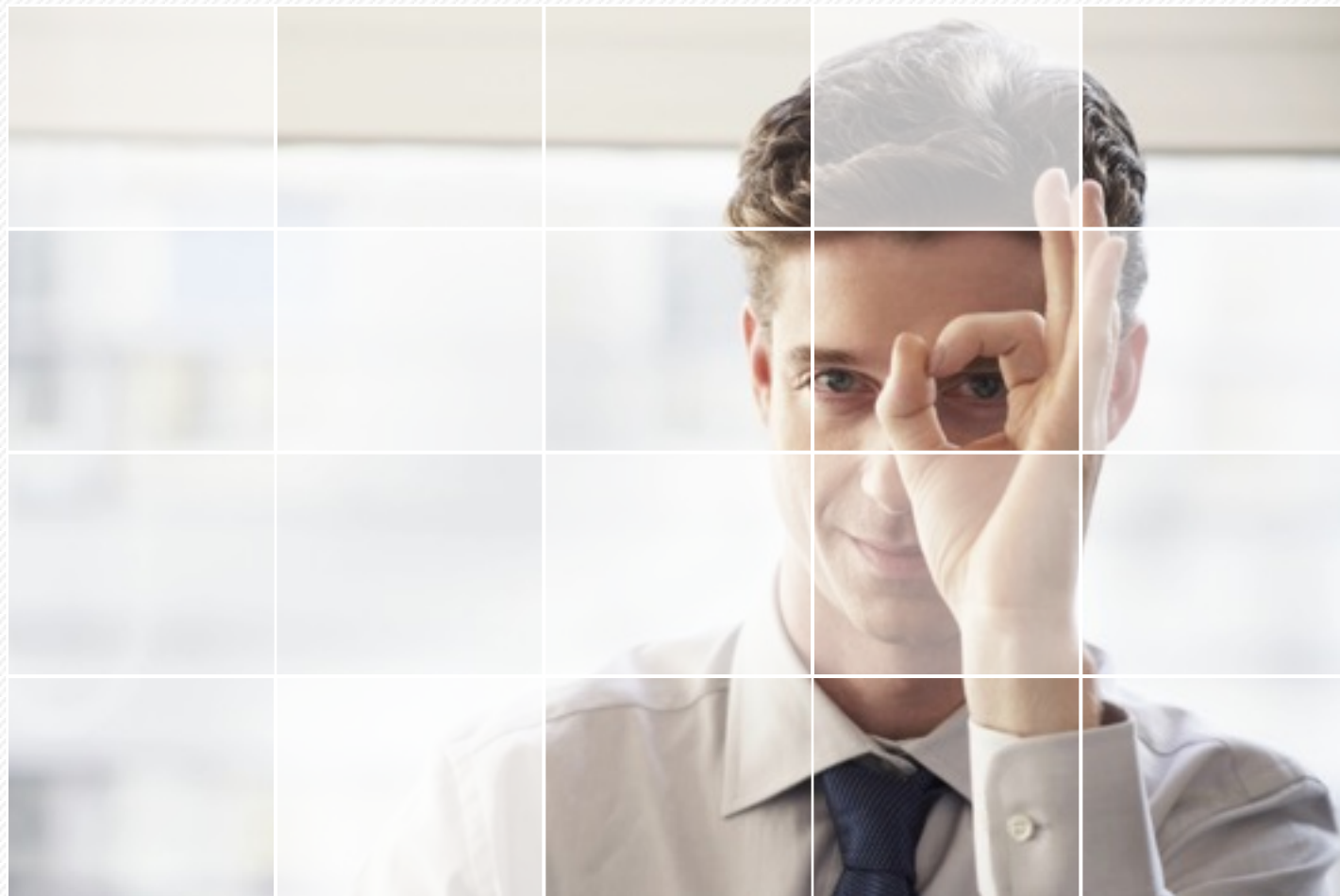


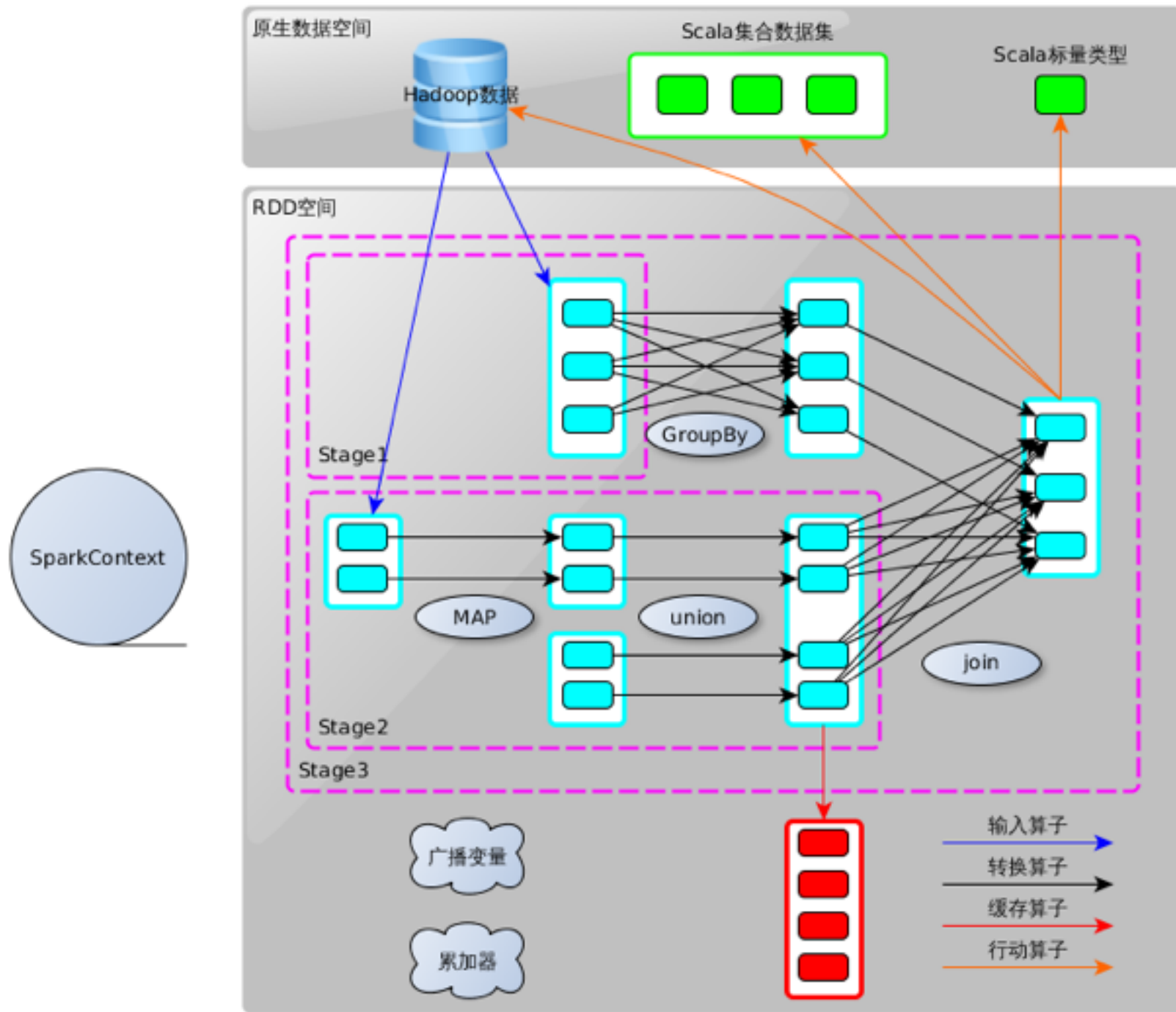
*Let's try it.*

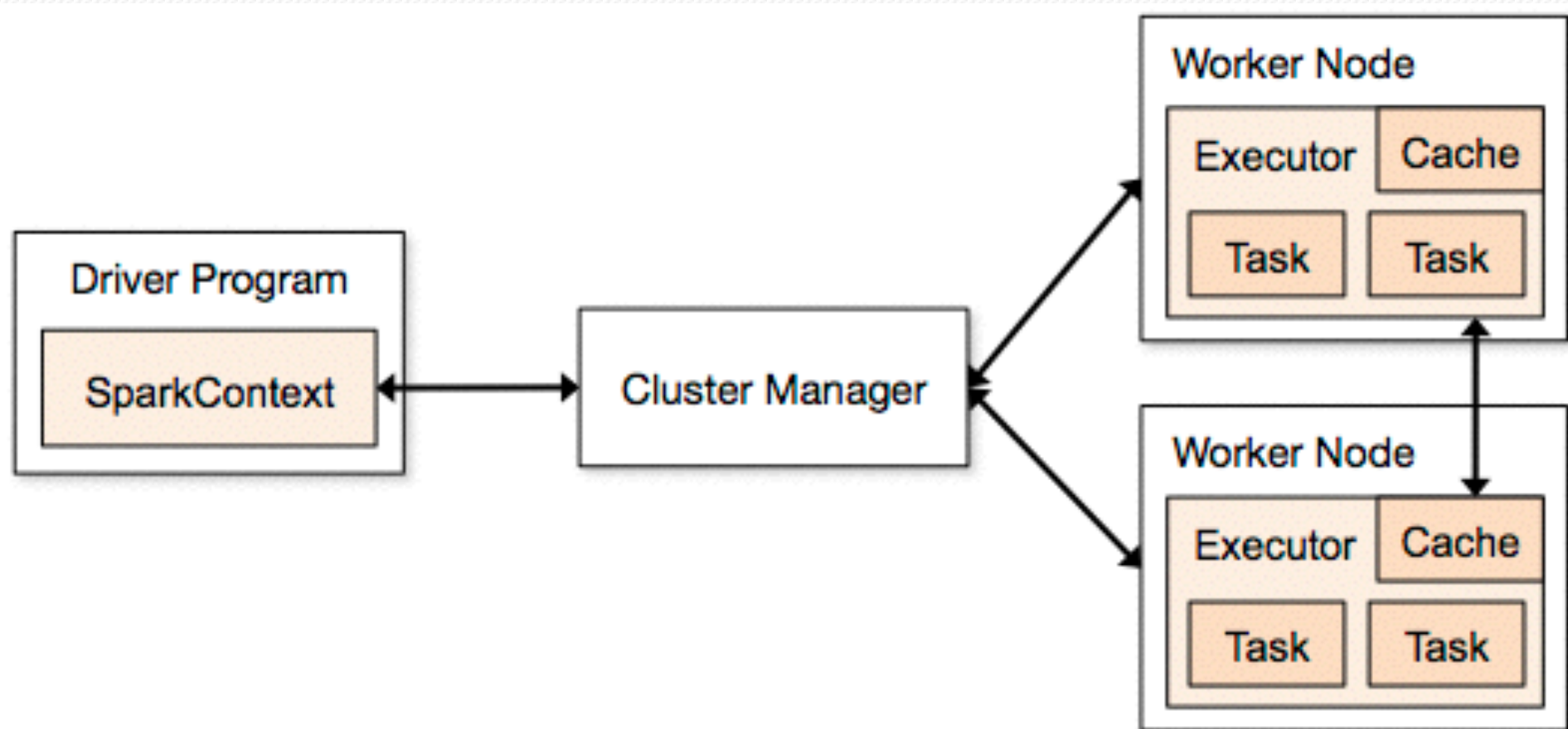


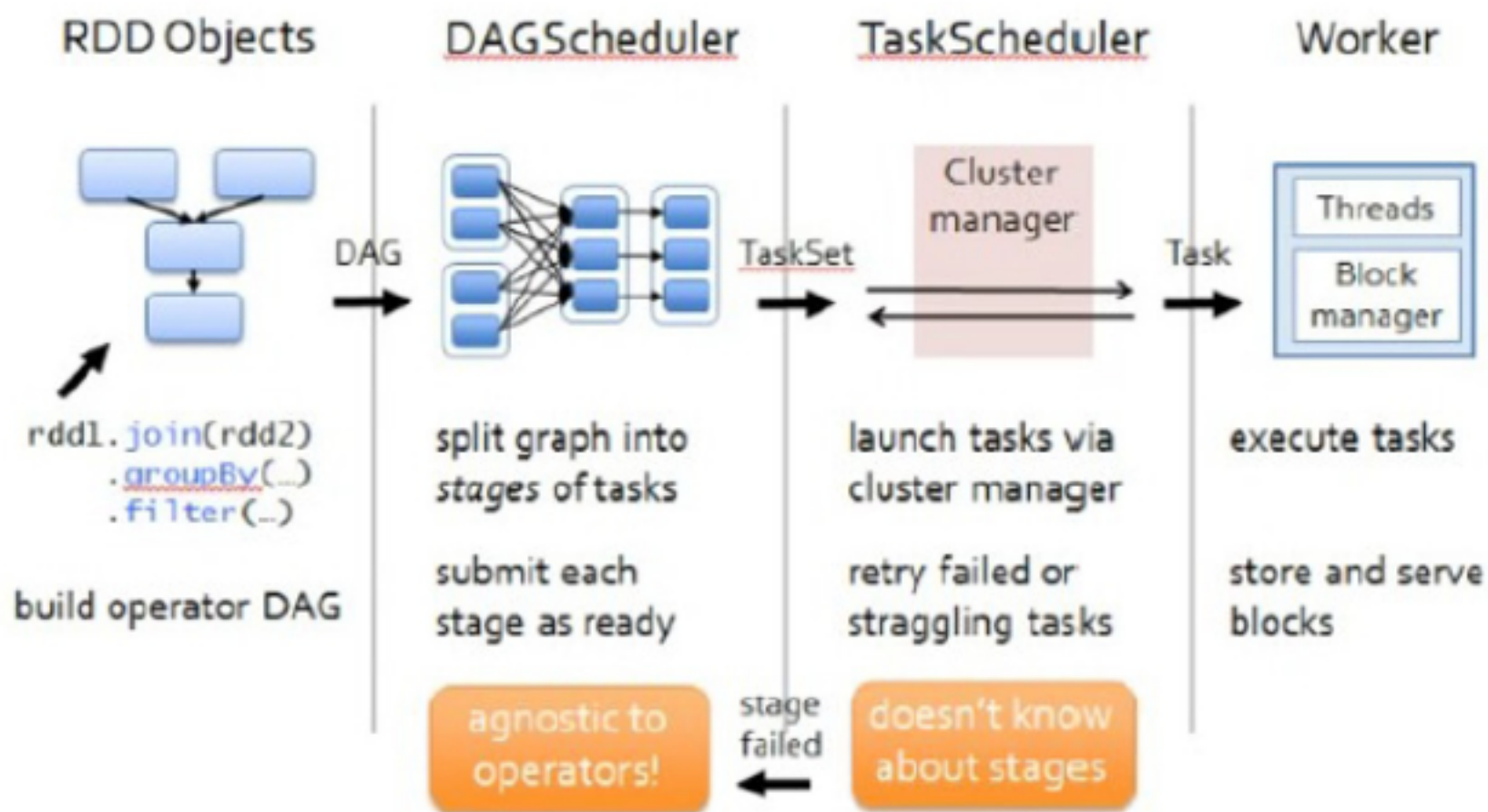


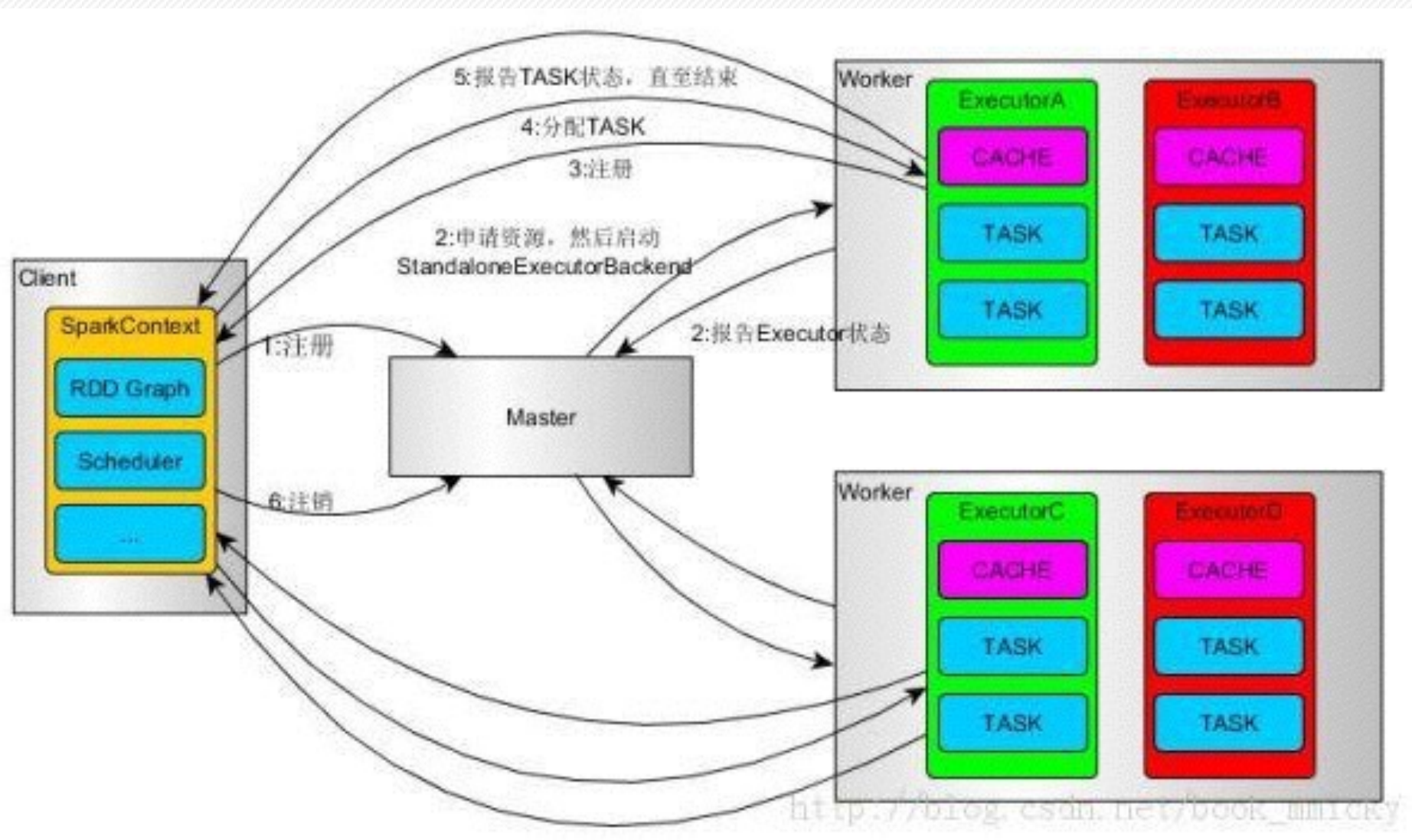
# Spark运行架构











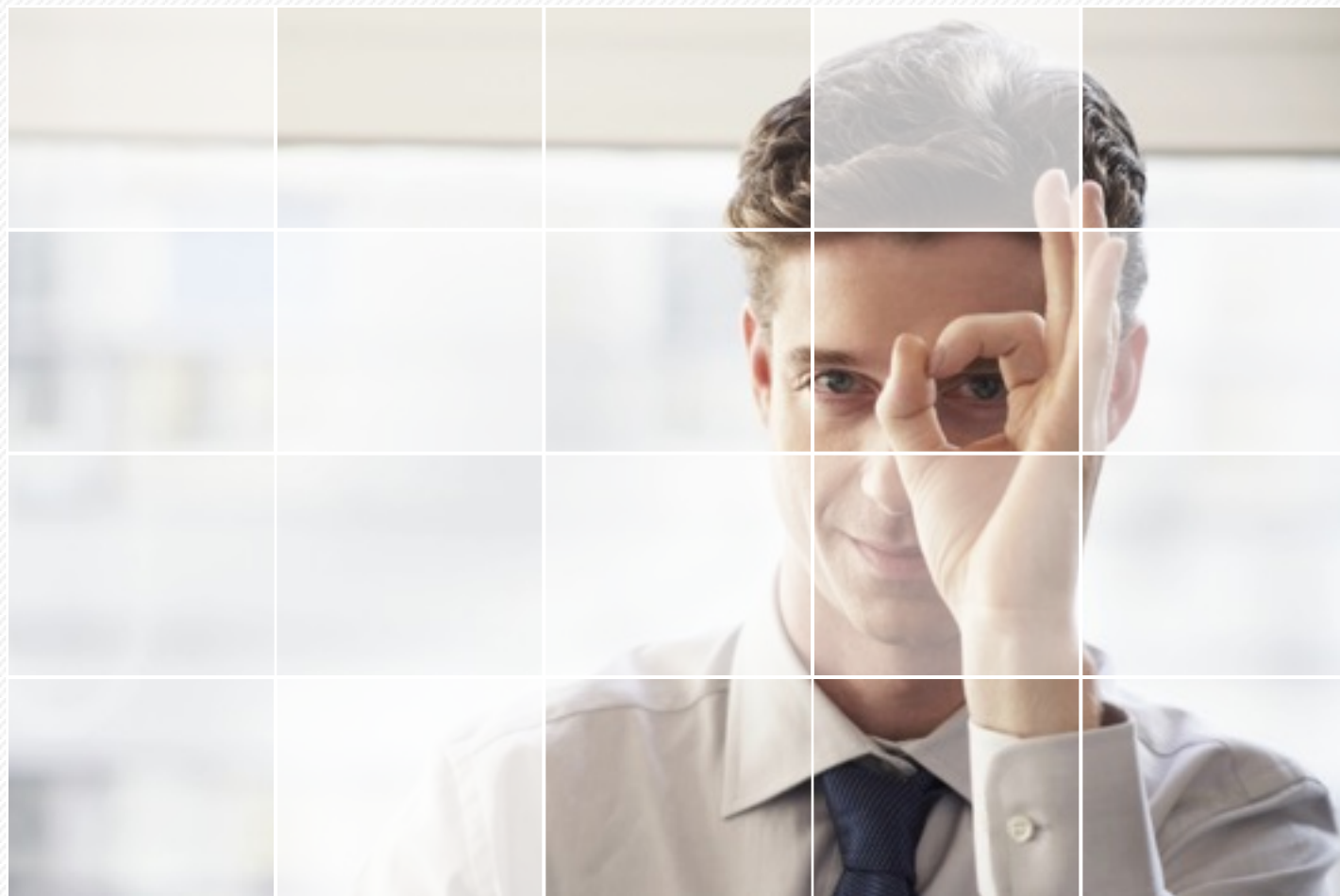




# Spark SQL核心与实践

1. 编程模型DataFrame&DataSets

2. 实践



# DataFrame && Datasets

Spark SQL 是 Spark 生态系统里用于处理结构化大数据的模块，它的前身是 Shark。随着 Spark 自身的发展，Spark 团队开始试图放弃 Shark 这个对于 Hive 有太多依赖 (查询优化，语法解析) 的东西，于是才有了 Spark SQL 这个全新的模块

该模块里最重要的概念就是 DataFrame, DataFrame 是基于早期版本中的 SchemaRDD，所以很自然的使用分布式大数据处理的场景。Spark DataFrame 以 RDD 为基础，但是带有 Schema 信息，它类似于传统数据库中的二维表格。

Spark SQL 模块目前支持将多种外部数据源的数据转化为 DataFrame，并像操作 RDD 或者将其注册为临时表的方式处理和分析这些数据。当前支持的数据源有：

- Json
- 文本文件
- RDD
- 关系数据库
- Hive
- Parquet

一旦将 DataFrame 注册成临时表，我们就可以使用类 SQL 的方式操作这些数据，我们将在下文的案例中详细向读者展示如何使用 Spark SQL/DataFrame 提供的 API 完成数据读取，临时表注册，统计分析等步骤。

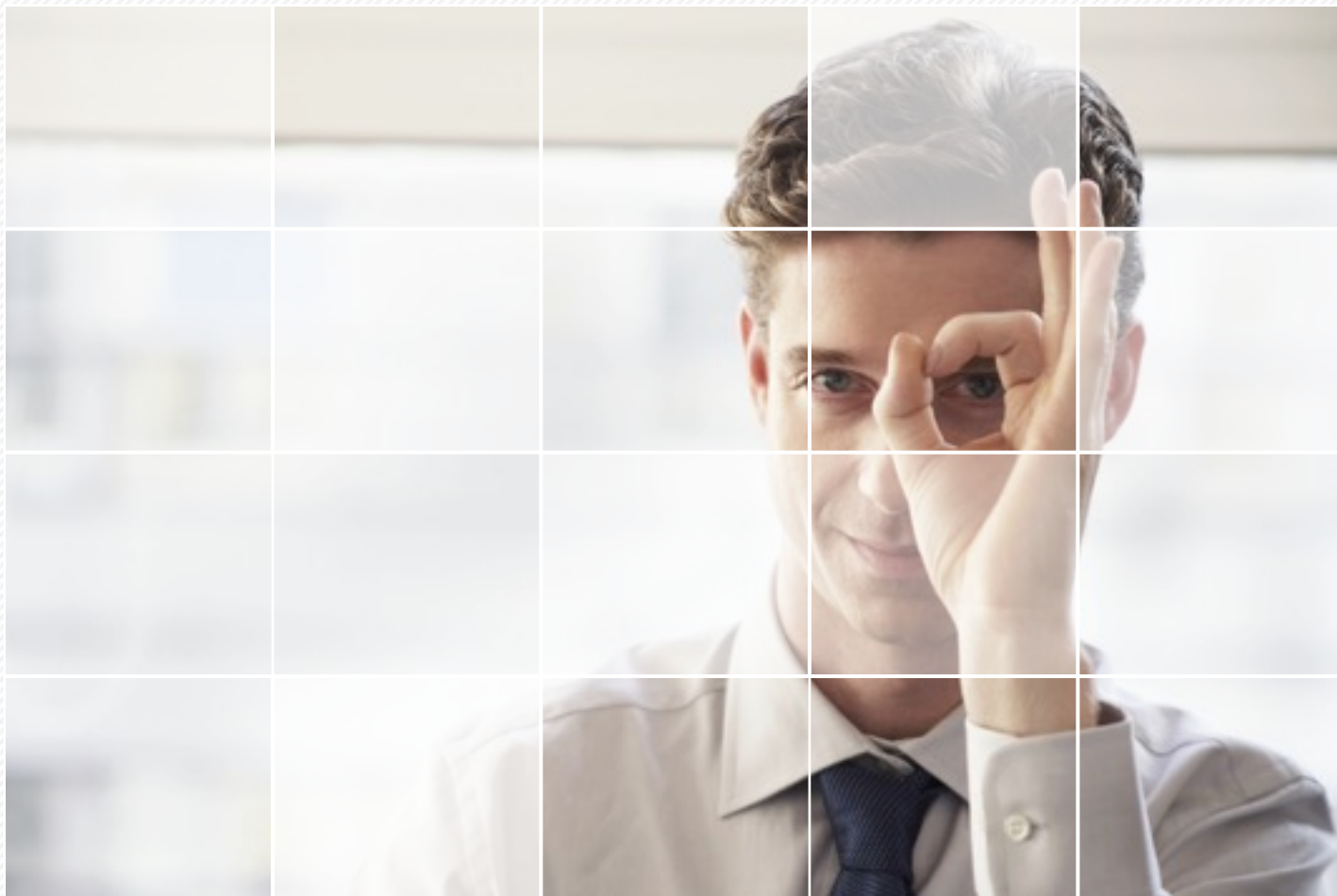


*Let's try it.*



# Streaming核心与实践

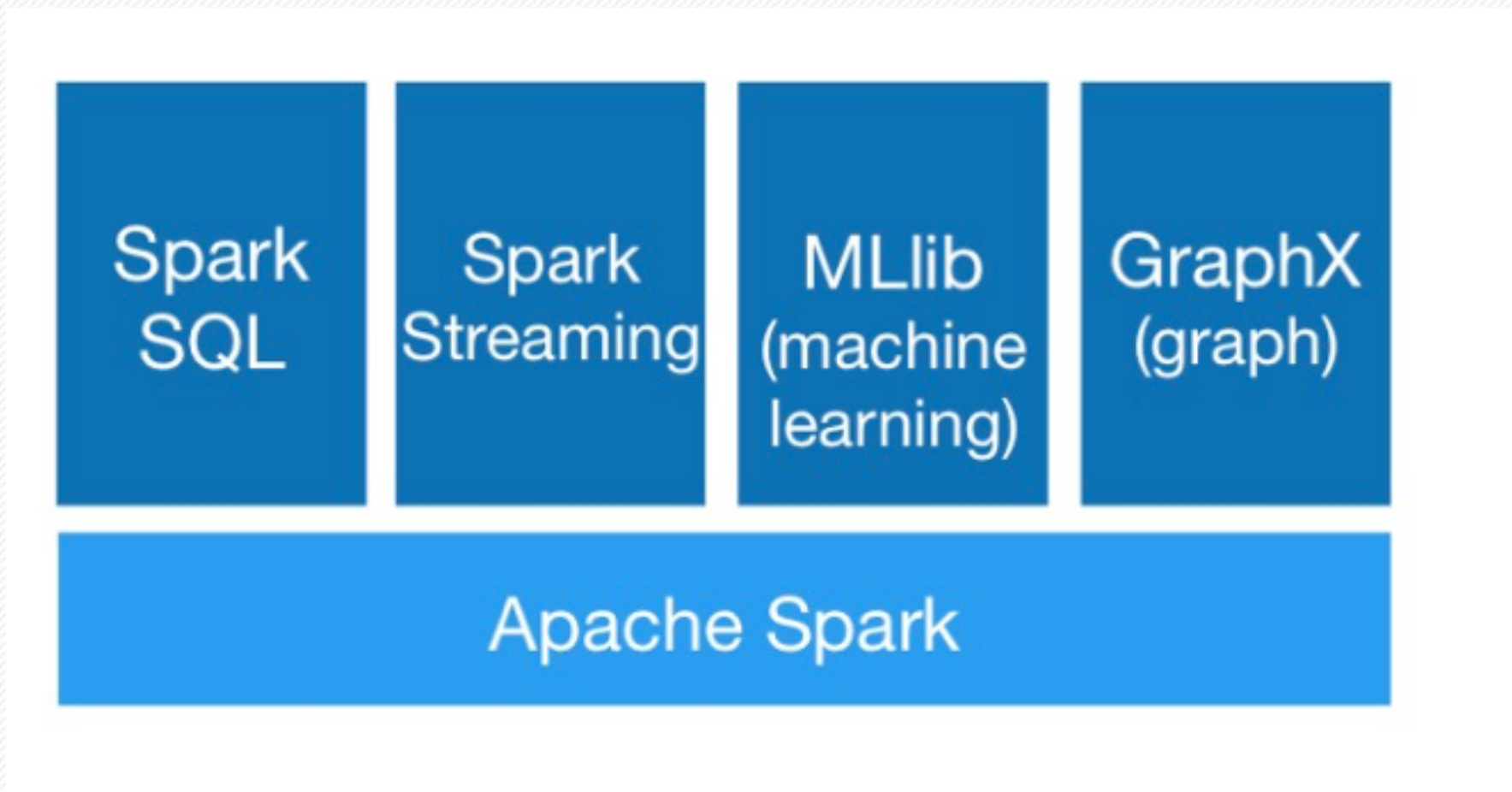
1. 运行原理
2. 运用场景
3. 编程模型DStream
4. 持久化和Checkpointing机制
5. 实践



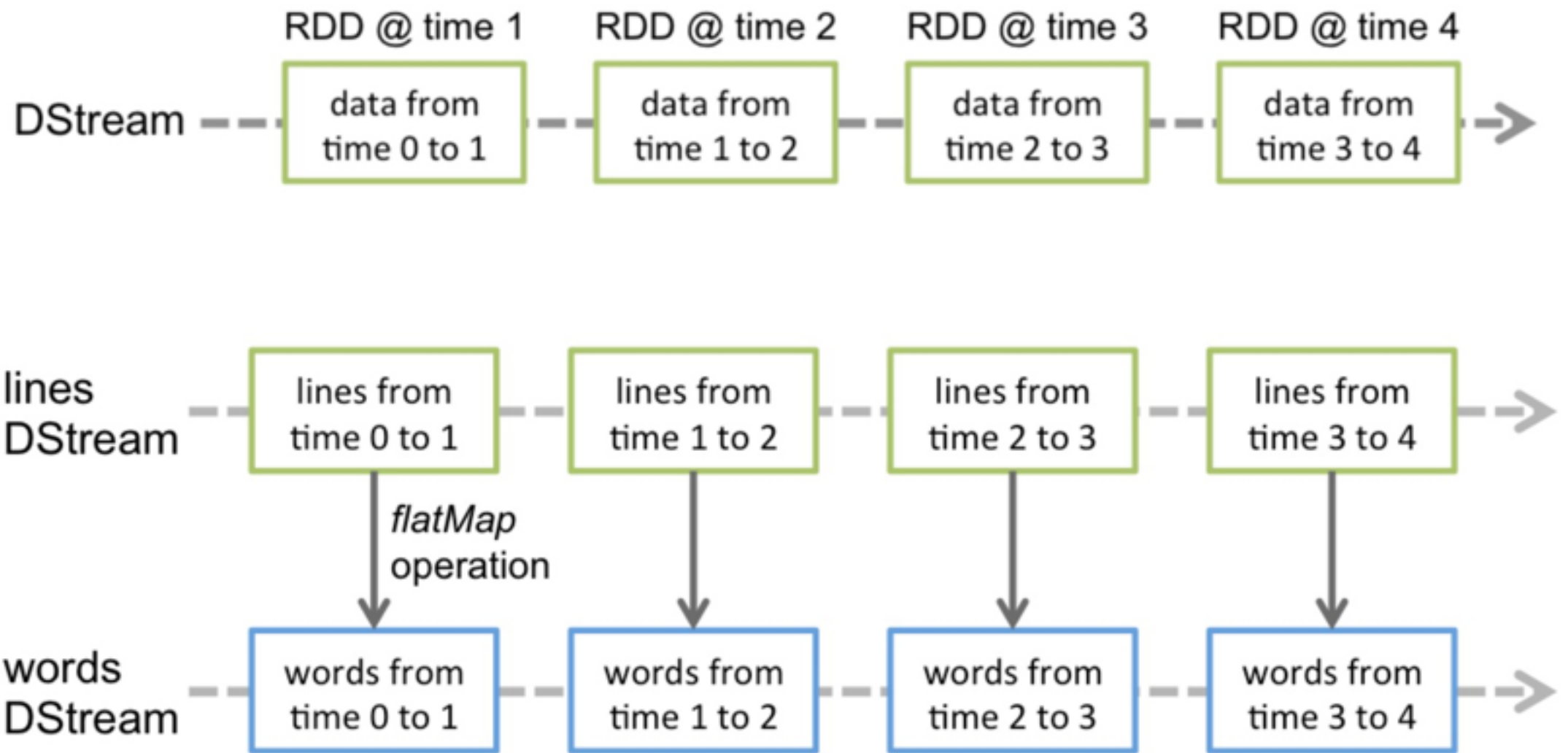
# Spark Streaming运行原理

是什么？

- 低延迟
- 横向扩展
- 容错性
- 可结合其他组件







# Streaming的运行场景

- 无状态操作
- 状态操作
  - UpdateStateByKey
- window操作

# 编程模型DStream

- Discretized Stream是Streaming的基础抽象
- DStream是RDD的系列
- DStream拥有大部分的RDD操作
- 从流输入源创建，从DStream之间transformation而来



# Input Sources

- 外部组件
  - Kafka, Flume, ZeroMQ, Akka, Akka Actors, TCP Sockets, HDFS
- 自定义接收器
- 生成RDD系列，作为流

Source	Artifact
Kafka	spark-streaming-kafka_2.10
Flume	spark-streaming-flume_2.10
Kinesis	spark-streaming-kinesis-asl_2.10 [Amazon Software License]
Twitter	spark-streaming-twitter_2.10
ZeroMQ	spark-streaming-zeromq_2.10
MQTT	spark-streaming-mqtt_2.10

# Output Operations on DStreams

<b><i>print()</i></b>	打印到控制台
<b><i>foreachRDD(func)</i></b>	对 <i>Dstream</i> 里面的每个 <i>RDD</i> 执行 <i>func</i> ，保存到外部系统
<b><i>saveAsObjectFiles(prefix, [suffix])</i></b>	保存流的计算结果为 <i>SequenceFile</i> , 文件名 : " <i>prefix-TIME_IN_MS[.suffix]</i> ".
<b><i>saveAsTextFiles(prefix, [suffix])</i></b>	保存流的计算结果为文本文件, 文件名 : " <i>prefix-TIME_IN_MS[.suffix]</i> ".
<b><i>saveAsHadoopFiles(prefix, [suffix])</i></b>	保存流的计算结果为 <i>hadoop</i> 文件, 文件名 : " <i>prefix-TIME_IN_MS[.suffix]</i> ".

# Transformations on DStreams

<i>map(func)</i>	对每一个元素执行 <i>func</i> 方法
<i>flatMap(func)</i>	类似 <i>map</i> 函数，但是可以 <i>map</i> 到0+个输出
<i>filter(func)</i>	过滤
<i>repartition(numPartitions)</i>	增加分区，提高并行度
<i>union(otherStream)</i>	合并两个流
<i>count()</i>	统计元素的个数
<i>reduce(func)</i>	对 <i>RDDs</i> 里面的元素进行聚合操作，2个输入参数，1个输出参数
<i>countByKey()</i>	针对类型统计，当一个 <i>Dstream</i> 的元素的类型是 <i>K</i> 的时候，调用它会返回一个新的 <i>Dstream</i> ，包含< <i>K</i> , <i>Long</i> >键值对， <i>Long</i> 是每个 <i>K</i> 出现的频率。
<i>reduceByKey(func, [numTasks])</i>	对于一个( <i>K</i> , <i>V</i> )类型的 <i>Dstream</i> ，为每个 <i>key</i> ，执行 <i>func</i> 函数，默认是 <i>local</i> 是2个线程， <i>cluster</i> 是8个线程，也可以指定 <i>numTasks</i>
<i>join(otherStream, [numTasks])</i>	把( <i>K</i> , <i>V</i> )和( <i>K</i> , <i>W</i> )的 <i>Dstream</i> 连接成一个( <i>K</i> , ( <i>V</i> , <i>W</i> ))的新 <i>Dstream</i>
<i>transform(func)</i>	转换操作，把原来的 <i>RDD</i> 通过 <i>func</i> 转换成一个新的 <i>RDD</i>
<i>updateStateByKey(func)</i>	针对 <i>key</i> 使用 <i>func</i> 来更新状态和值，可以将 <i>state</i> 该为任何值

# Window Operations

<b><i>window(windowLength, slideInterval)</i></b>	返回一个根据制定的窗口大小和滑动大小的 <b><i>dstream</i></b>
<b><i>countByWindow(windowLength, slideInterval)</i></b>	返回流中窗口滑动的大小技术
<b><i>reduceByWindow(func, windowLength, slideInterval)</i></b>	对指定的窗口进行 <b><i>func</i></b> 的操作
<b><i>reduceByKeyAndWindow(func, windowLength, slideInterval, [numTasks])</i></b>	对流中 <b><i>kv</i></b> 类型的窗口，进行 <b><i>reduceByKey</i></b> ，可以设置 <b><i>task</i></b> 个数
<b><i>reduceByKeyAndWindow(func, invFunc, windowLength, slideInterval, [numTasks])</i></b>	...
<b><i>countByValueAndWindow(windowLength, slideInterval, [numTasks])</i></b>	对流中 <b><i>kv</i></b> 类型的窗口，进行 <b><i>countByValue</i></b> ，可以设置 <b><i>task</i></b> 个数

# 持久化 & Checkpointing

- 持久化

- 通过persist()来持久化
- 默认的持久化级别是MEMORY\_ONLY\_SER
- 对于来源于网络的数据源采用MEMORY\_AND\_DISK\_SER\_2

- Checkpointing

- window和stateful操作必须checkpoint
- ssc.checkpoint指定目录，dstrem.checkpoint指定时间间隔
- 间隔必须是窗口slide interval的倍数
- 合理的设置checkponit时间

*Let's try it.*



# THANKS!



[alien.xt.xm@gmail.com](mailto:alien.xt.xm@gmail.com)

<http://github.com/alienxt>