

Life is short, you need Spark!



从零开始

不需要任何基础，带领您无痛入门 Spark

云计算分布式大数据 Spark 实战高手之路

王家林著

Spark 亚太研究院系列丛书 版权所有

伴随着大数据相关技术和产业的逐步成熟，继 Hadoop 之后，Spark 技术以其无可比拟的优势，发展迅速，将成为替代 Hadoop 的下一代云计算、大数据核心技术。

本书特点

- ▶ 云计算分布式大数据 Spark 实战高手之路三部曲之第一部
- ▶ 网络发布版为图文并茂方式，边学习，边演练
- ▶ 不需要任何前置知识，从零开始，循序渐进

本书作者



Spark 亚太研究院院长和首席专家，中国目前唯一的移动互联网和云计算大数据集大成者。在 Spark、Hadoop、Android 等方面有丰富的源码、实务和性能优化经验。彻底研究了 Spark 从 0.5.0 到 0.9.1 共 13 个版本的 Spark 源码，并已完成 2014 年 5 月 31 日发布的 Spark1.0 源码研究。

Hadoop 源码级专家，曾负责某知名公司的类 Hadoop 框架开发工作，专注于 Hadoop 一站式解决方案的提供，同时也是云计算分布式大数据处理的最早实践者之一。

Android 架构师、高级工程师、咨询顾问、培训专家。

通晓 Spark、Hadoop、Android、HTML5，迷恋英语播音和健美。

“真相会使你获得自由。”

—— 耶稣《圣经》约翰 8:32KJV

“所有人类的幸福都来源于不能直面事实。”

—— 释迦摩尼

“道法自然”

—— 老子《道德经》第 25 章

《云计算分布式大数据 Spark 实战高手之路》

系列丛书三部曲

《云计算分布式大数据 Spark 实战高手之路---从零开始》：

不需要任何基础，带领您无痛入门 Spark 并能够轻松处理 Spark 工程师的日常编程工作，内容包括 Spark 集群的构建、Spark 架构设计、RDD、Shark/SparkSQL、机器学习、图计算、实时流处理、Spark on Yarn、JobServer、Spark 测试、Spark 优化等。

《云计算分布式大数据 Spark 实战高手之路---高手崛起》：

大话 Spark 源码，全世界最有情趣的源码解析，过程中伴随诸多实验，解析 Spark 1.0 的任何一句源码！更重要的是，思考源码背后的问题场景和解决问题的设计哲学和实现招式。

《云计算分布式大数据 Spark 实战高手之路---高手之巅》：

通过当今主流的 Spark 商业使用方法和最成功的 Hadoop 大型案例让您直达高手之巅，从此一览众山小。



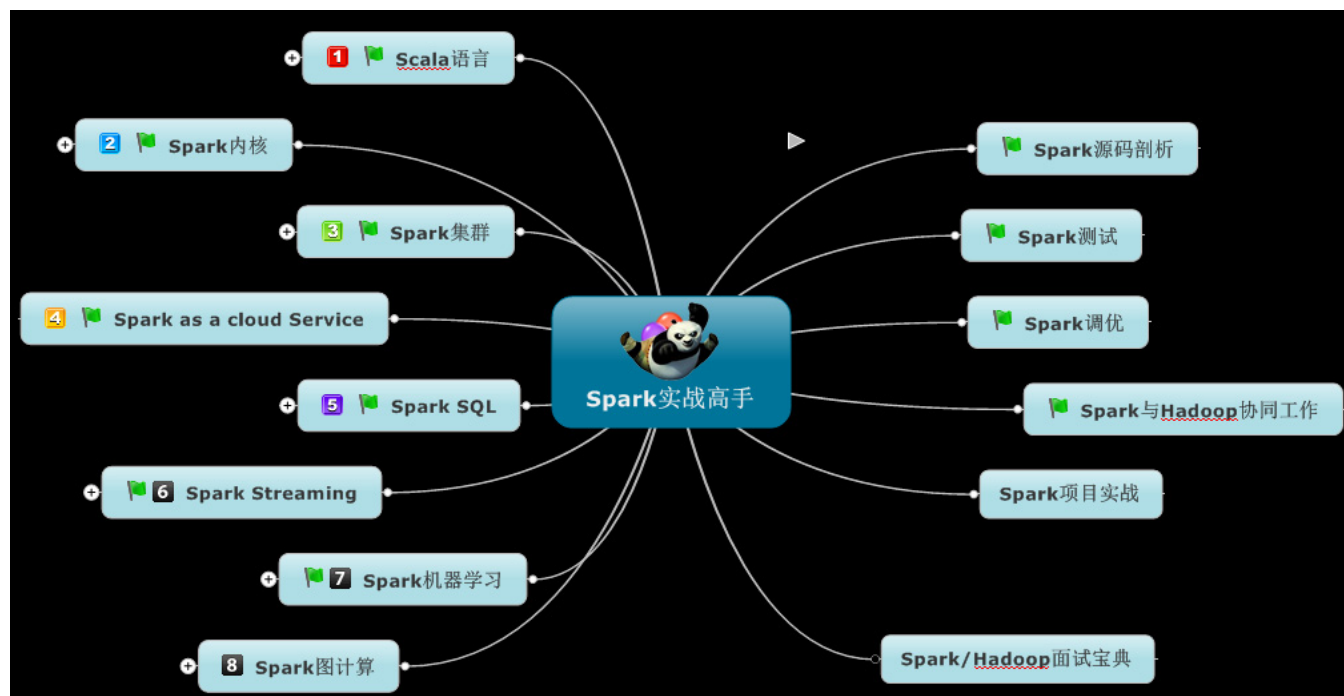
《前言》

Spark采用一个统一的技术堆栈解决了云计算大数据的如流处理、图技术、机器学习、NoSQL查询等方面的所有核心问题，具有完善的生态系统，这直接奠定了其统一云计算大数据领域的霸主地位；

要想成为Spark高手，需要经历六大阶段



Spark 实战高手之核心技能点



第一阶段：熟练的掌握Scala语言

1. Spark 框架是采用 Scala 语言编写的，精致而优雅。要想成为 Spark 高手，你就必须阅读 Spark 的源代码，就必须掌握 Scala；
 2. 虽然说现在的 Spark 可以采用多语言 Java、Python 等进行应用程序开发，但是最快速的和支持最好的开发 API 依然并将永远是 Scala 方式的 API，所以你必须掌握 Scala 来编写复杂的和高性能的 Spark 分布式程序；
 3. 尤其要熟练掌握 Scala 的 trait、apply、函数式编程、泛型、逆变与协变等；
- 推荐课程：“精通Spark的开发语言：Scala最佳实践”

第二阶段：精通Spark平台本身提供给开发者API

1. 掌握 Spark 中面向 RDD 的开发模式 掌握各种 transformation 和 action 函数的使用；
2. 掌握 Spark 中的宽依赖和窄依赖以及 lineage 机制；
3. 掌握 RDD 的计算流程，例如 Stage 的划分、Spark 应用程序提交给集群的基本过程和 Worker 节点基础的工作原理等

推荐课程：“18 小时内掌握Spark：把云计算大数据速度提高 100 倍以上!”

第三阶段：深入Spark内核

此阶段主要是通过 Spark 框架的源码研读来深入 Spark 内核部分：

1. 通过源码掌握 Spark 的任务提交过程；
2. 通过源码掌握 Spark 集群的任务调度；
3. 尤其要精通 DAGScheduler、TaskScheduler 和 Worker 节点内部的工作的每一步的细节；

推荐课程：[“Spark 1.0.0 企业级开发动手：实战世界上第一个Spark 1.0.0 课程，涵盖Spark 1.0.0 所有的企业级开发技术”](#)

第四阶段:掌握基于Spark上的核心框架的使用

Spark 作为云计算大数据时代的集大成者，在实时流处理、图技术、机器学习、NoSQL 查询等方面具有显著的优势，我们使用 Spark 的时候大部分时间都是在使用其上的框架例如 Shark、Spark Streaming 等：

1. Spark Streaming 是非常出色的实时流处理框架，要掌握其 DStream、transformation 和 checkpoint 等；
2. Spark 的离线统计分析功能，Spark 1.0.0 版本在 Shark 的基础上推出了 Spark SQL，离线统计分析的功能的效率有显著的提升，需要重点掌握；
3. 对于 Spark 的机器学习和 GraphX 等要掌握其原理和用法；

推荐课程：[“Spark企业级开发最佳实践”](#)

第五阶段:做商业级别的Spark项目

通过一个完整的具有代表性的 Spark 项目来贯穿 Spark 的方方面面，包括项目的架构设计、用到的技术的剖析、开发实现、运维等，完整掌握其中的每一个阶段和细节，这样就可以让您以后可以从容面对绝大多数 Spark 项目。

推荐课程：[“Spark架构案例鉴赏：Conviva、Yahoo！、优酷土豆、网易、腾讯、淘宝等公司的实际Spark案例”](#)

第六阶段：提供Spark解决方案

1. 彻底掌握 Spark 框架源码的每一个细节；
2. 根据不同的业务场景的需要提供 Spark 在不同场景的下的解决方案；
3. 根据实际需要，在 Spark 框架基础上进行二次开发，打造自己的 Spark 框架；

推荐课程：[“精通Spark：Spark内核剖析、源码解读、性能优化和商业案例实战”](#)

《第四章：Spark 内核揭秘》

掌握 Spark 内核是精通 Spark 的关键，也是驾驭 Spark 的精髓所在。

基于 Spark 内核，Spark 构建起了一体化多元化的大数据处理流水线，在一个技术堆栈中即可以同时完成批处理、实时流处理、交互式查询、机器学习、图计算以及这些子框架之间数据和 RDD 算子的无缝共享与互操作。

可以说，Spark 内核是每个想彻底掌握 Spark 的人员的必修课，通过对内核的探索，我们对整个 Spark 的运行机制会了如指掌，这对 Spark 的大规模应用、性能优化、系统自定义开发 Spark 系统都是至关重要的。

本章首先带你初探 Spark 内核，接着通过源码游离 Spark 内核中，然后通过源码细致解析 Spark 作业的全生命周期，最后分享 Spark 性能优化，内容组织循序渐进，希望助力诸位 Spark 爱好者能够掌握 Spark 内核。

Spark 内核揭秘共分六个部分：

- 第一部分：Spark 内核初探
- 第二部分：Spark 内核核心源码解析
- 第三部分：Job 全生命周期源码解读
- 第四部分：解读 RDD 源码
- 第五部分：Driver、Master、Worker
- 第六部分：Spark 性能优化

本讲是 Spark 内核揭秘的第四部分：解读 RDD 源码，具体内容如下所示：

- RDD 源码研读

不需任何前置知识，从零开始，循序渐进，成为 Spark 高手！



目录

RDD源码研读.....8

RDD 源码研读

RDD 的官方解释是：

A Resilient Distributed Dataset (RDD), the basic abstraction in Spark. Represents an immutable, partitioned collection of elements that can be operated on in parallel. This class contains the basic operations available on all RDDs, such as `map`, `filter`, and `persist`. In addition, [org.apache.spark.rdd.PairRDDFunctions](#) contains operations available only on RDDs of key-value pairs, such as `groupByKey` and `join`; [org.apache.spark.rdd.DoubleRDDFunctions](#) contains operations available only on RDDs of Doubles; and [org.apache.spark.rdd.SequenceFileRDDFunctions](#) contains operations available on RDDs that can be saved as SequenceFiles. These operations are automatically available on any RDD of the right type (e.g. `RDD[(Int, Int)]`) through implicit conversions when you `import org.apache.spark.SparkContext._`.

Internally, each RDD is characterized by five main properties:

- A list of partitions
- A function for computing each split
- A list of dependencies on other RDDs
- Optionally, a Partitioner for key-value RDDs (e.g. to say that the RDD is hash-partitioned)
- Optionally, a list of preferred locations to compute each split on (e.g. block locations for an HDFS file)

All of the scheduling and execution in Spark is done based on these methods, allowing each RDD to implement its own way of computing itself. Indeed, users can implement custom RDDs (e.g. for reading data from a new storage system) by overriding these functions. Please refer to the [Spark paper](#) for more details on RDD internals.

每个 RDD 的核心方法如下所示:

```
/**
 * :: DeveloperApi ::
 * Implemented by subclasses to compute a given partition.
 */
@DeveloperApi
def compute(split: Partition, context: TaskContext): Iterator[T]

/**
 * Implemented by subclasses to return the set of partitions in this RDD. This method will only
 * be called once, so it is safe to implement a time-consuming computation in it.
 */
protected def getPartitions: Array[Partition]

/**
 * Implemented by subclasses to return how this RDD depends on parent RDDs. This method will only
 * be called once, so it is safe to implement a time-consuming computation in it.
 */
protected def getDependencies: Seq[Dependency[_]] = deps

/**
 * Optionally overridden by subclasses to specify placement preferences.
 */
protected def getPreferredLocations(split: Partition): Seq[String] = Nil
```


RDD 主要分为两种：

Transformations	<code>map(f: T ⇒ U)</code>	: <code>RDD[T] ⇒ RDD[U]</code>
	<code>filter(f: T ⇒ Bool)</code>	: <code>RDD[T] ⇒ RDD[T]</code>
	<code>flatMap(f: T ⇒ Seq[U])</code>	: <code>RDD[T] ⇒ RDD[U]</code>
	<code>sample(fraction: Float)</code>	: <code>RDD[T] ⇒ RDD[T]</code> (Deterministic sampling)
	<code>groupByKey()</code>	: <code>RDD[(K, V)] ⇒ RDD[(K, Seq[V])]</code>
	<code>reduceByKey(f: (V, V) ⇒ V)</code>	: <code>RDD[(K, V)] ⇒ RDD[(K, V)]</code>
	<code>union()</code>	: <code>(RDD[T], RDD[T]) ⇒ RDD[T]</code>
	<code>join()</code>	: <code>(RDD[(K, V)], RDD[(K, W)]) ⇒ RDD[(K, (V, W))]</code>
	<code>cogroup()</code>	: <code>(RDD[(K, V)], RDD[(K, W)]) ⇒ RDD[(K, (Seq[V], Seq[W]))]</code>
	<code>crossProduct()</code>	: <code>(RDD[T], RDD[U]) ⇒ RDD[(T, U)]</code>
	<code>mapValues(f: V ⇒ W)</code>	: <code>RDD[(K, V)] ⇒ RDD[(K, W)]</code> (Preserves partitioning)
	<code>sort(c: Comparator[K])</code>	: <code>RDD[(K, V)] ⇒ RDD[(K, V)]</code>
	<code>partitionBy(p: Partitioner[K])</code>	: <code>RDD[(K, V)] ⇒ RDD[(K, V)]</code>
Actions	<code>count()</code>	: <code>RDD[T] ⇒ Long</code>
	<code>collect()</code>	: <code>RDD[T] ⇒ Seq[T]</code>
	<code>reduce(f: (T, T) ⇒ T)</code>	: <code>RDD[T] ⇒ T</code>
	<code>lookup(k: K)</code>	: <code>RDD[(K, V)] ⇒ Seq[V]</code> (On hash/range partitioned RDDs)
	<code>save(path: String)</code>	: Outputs RDD to a storage system, e.g., HDFS

Table 2: Transformations and actions available on RDDs in Spark. Seq[T] denotes a sequence of elements of type T.

下面我们初步分析一下 RDD 的四个核心方法，首先看一下 getPartitions 方法的源码：

```
/**
 * Implemented by subclasses to return the set of partitions in this RDD. This method will only
 * be called once, so it is safe to implement a time-consuming computation in it.
 */
protected def getPartitions: Array[Partition]
```

getPartitions 返回的是一系列 partitions 的集合，即一个 Partition 类型的数组；

而 getDependencies 表达式 RDD 之间的依赖关系，如下所示：

```
/**
 * Implemented by subclasses to return how this RDD depends on parent RDDs. This method will only
 * be called once, so it is safe to implement a time-consuming computation in it.
 */
protected def getDependencies: Seq[Dependency[_]] = deps
```

getDependencies 返回的是依赖关系的一个 Seq 集合，里面的 Dependency 数组中的下划线是类型的 Placeholder；

每个 RDD 都会具有计算的函数，如下所示：

```
/**
 * :: DeveloperApi ::
 * Implemented by subclasses to compute a given partition.
 */
@DeveloperApi
def compute(split: Partition, context: TaskContext): Iterator[T]
```

Compute 方法是针对 RDD 的每个 Partition 进行计算的，其 TaskContext 参数的源码如下：

```

/**
 * :: DeveloperApi ::
 * Contextual information about a task which can be read or mutated during execution.
 *
 * @param stageId stage id
 * @param partitionId index of the partition
 * @param attemptId the number of attempts to execute this task
 * @param runningLocally whether the task is running locally in the driver JVM
 * @param taskMetrics performance metrics of the task
 */
@DeveloperApi
class TaskContext(
  val stageId: Int,
  val partitionId: Int,
  val attemptId: Long,
  val runningLocally: Boolean = false,
  private[spark] val taskMetrics: TaskMetrics = TaskMetrics.empty)
  extends Serializable {

  @deprecated("use partitionId", "0.8.1")
  def splitId = partitionId

  // List of callback functions to execute when the task completes.
  @transient private val onCompleteCallbacks = new ArrayBuffer[TaskCompletionListener]

```

getPreferredLocations 是寻找 Partition 的首选位置：

```

/**
 * Optionally overridden by subclasses to specify placement preferences.
 */
protected def getPreferredLocations(split: Partition): Seq[String] = Nil

```

其实 RDD 还有一个可选的分区策略：

```

/** Optionally overridden by subclasses to specify how they are partitioned. */
@transient val partitioner: Option[Partitioner] = None

```

Partitioner 的源码如下：

```

/**
 * An object that defines how the elements in a key-value pair RDD are partitioned by key.
 * Maps each key to a partition ID, from 0 to `numPartitions - 1`.
 */
abstract class Partitioner extends Serializable {
  def numPartitions: Int
  def getPartition(key: Any): Int
}

object Partitioner {
  /**
   * Choose a partitioner to use for a cogroup-like operation between a number of RDDs.
   *
   * If any of the RDDs already has a partitioner, choose that one.
   *
   * Otherwise, we use a default HashPartitioner. For the number of partitions, if
   * spark.default.parallelism is set, then we'll use the value from SparkContext
   * defaultParallelism, otherwise we'll use the max number of upstream partitions.
   *
   * Unless spark.default.parallelism is set, the number of partitions will be the
   * same as the number of partitions in the largest upstream RDD, as this should
   * be least likely to cause out-of-memory errors.
   */
  /**
   * We use two method parameters (rdd, others) to enforce callers passing at least 1 RDD.
   */
  def defaultPartitioner(rdd: RDD[_], others: RDD[_]*): Partitioner = {
    val bySize = (Seq(rdd) ++ others).sortBy(_.partitions.size).reverse
    for (r <- bySize if r.partitioner.isDefined) {
      return r.partitioner.get
    }
    if (rdd.context.conf.contains("spark.default.parallelism")) {
      new HashPartitioner(rdd.context.defaultParallelism)
    } else {
      new HashPartitioner(bySize.head.partitions.size)
    }
  }
}

/**
 * A [[org.apache.spark.Partitioner]] that implements hash-based partitioning using
 * Java's `Object.hashCode`.
 *
 * Java arrays have hashCodes that are based on the arrays' identities rather than their contents,
 * so attempting to partition an RDD[Array[_]] or RDD[(Array[_], _)] using a HashPartitioner will
 * produce an unexpected or incorrect result.
 */
class HashPartitioner(partitions: Int) extends Partitioner {

```

```
def numPartitions = partitions

def getPartition(key: Any): Int = key match {
  case null => 0
  case _ => Utils.nonNegativeMod(key.hashCode, numPartitions)
}

override def equals(other: Any): Boolean = other match {
  case h: HashPartitioner =>
    h.numPartitions == numPartitions
  case _ =>
    false
}

override def hashCode: Int = numPartitions
}
```

可以看出默认使用的是 HashPartitioner，要注意 key 为 Array 的情况；

■ Spark 亚太研究院

Spark 亚太研究院，提供 Spark、Hadoop、Android、Html5、云计算和移动互联网一站式解决方案。以帮助企业规划、部署、开发、培训和使用为核心，并规划和实施人才培训完整路径，提供源码研究和应用技术训练。

■ 近期活动及相关课程

1、决战云计算大数据时代 Spark 亚太研究院 100 期公益大奖堂

每周四晚上 20:00—21:00

课程介绍：http://edu.51cto.com/course/course_id-1659.html#showDesc

报名参与：http://ke.qq.com/cgi-bin/courseDetail?course_id=6167

2、大数据 Spark 实战高手之路—熟练掌握 Scala 语言视频课程



国内第一个 Scala 视频学习课程！

成为 Spark 高手必备技能，必修课程！

现在购买，即可享受套餐优惠！

课程地址：<http://edu.51cto.com/pack/view/id-124.html>

■ 近期公开课：

《决胜大数据时代：Hadoop、Yarn、Spark 企业级最佳实践》

集大数据领域最核心三大技术：Hadoop 方向 50%：掌握生产环境下、源码级别下的 Hadoop 经验，解决性能、集群难点问题；Yarn 方向 20%：掌握最佳的分布式集群资源管理框架，能够轻松使用 Yarn 管理 Hadoop、Spark 等；Spark 方向 30%：未来统一的大数据框架平台，剖析 Spark 架构、内核等核心技术，对未来转向 SPARK 技术，做好技术储备。课程内容落地性强，即解决当下问题，又有助于驾驭未来。

开课时间：9 月 26—28 日 上海、10 月 26—28 日 北京、11 月 1—3 日 深圳

咨询电话：4006-998-758

QQ 交流群：1 群：317540673（已满）
2 群 297931500



微信公众号：spark-china