



# Spark源码分析之-Storage模块

[architecture \(1\) \(/categories.html#architecture-ref\)](/categories.html#architecture-ref)

[cloud <sup>7</sup> \(/tags.html#cloud-ref\)](/tags.html#cloud-ref)

[spark <sup>8</sup> \(/tags.html#spark-ref\)](/tags.html#spark-ref)

08 October 2013

## Background

前段时间琐事颇多，一直没有时间整理自己的博客，Spark源码分析写到一半也搁置了。之前介绍了[deploy](http://jerryshao.me/architecture/2013/04/21/Spark%E6%BA%90%E7%A0%81%E5%88%86%E6%9E%90%E4%B9%8B-scheduler%E6%A8%A1%E5%9D%97/) (<http://jerryshao.me/architecture/2013/04/21/Spark%E6%BA%90%E7%A0%81%E5%88%86%E6%9E%90%E4%B9%8B-scheduler%E6%A8%A1%E5%9D%97/>)和[scheduler](http://jerryshao.me/architecture/2013/04/21/Spark%E6%BA%90%E7%A0%81%E5%88%86%E6%9E%90%E4%B9%8B-scheduler%E6%A8%A1%E5%9D%97/) (<http://jerryshao.me/architecture/2013/04/21/Spark%E6%BA%90%E7%A0%81%E5%88%86%E6%9E%90%E4%B9%8B-scheduler%E6%A8%A1%E5%9D%97/>)两大模块，这次介绍Spark中的另一大模块 - **storage**模块。

在写Spark程序的时候我们常常和**RDD** (*Resilient Distributed Dataset*) 打交道，通过RDD为我们提供的各种**transformation**和**action**接口实现我们的应用，RDD的引入提高了抽象层次，在接口和实现上进行有效地隔离，使用户无需关心底层的实现。但是RDD提供给我们的仅仅是一个“形”，我们所操作的数据究竟放在哪里，如何存取？它的“体”是怎样的？这是由**storage**模块来实现和管理的，接下来我们就要剖析一下**storage**模块。

## Storage模块整体架构

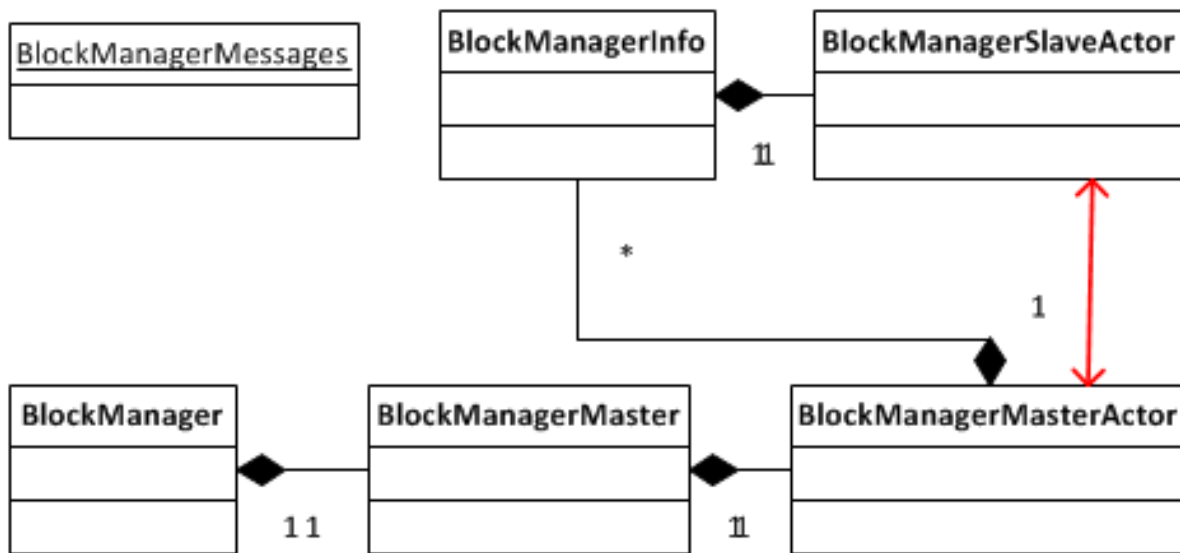
Storage模块主要分为两层：

1. 通信层：**storage**模块采用的是**master-slave**结构来实现通信层，**master**和**slave**之间传输控制信息、状态信息，这些都是通过通信层来实现的。
2. 存储层：**storage**模块需要把数据存储到**disk**或是**memory**上面，有可能还需**replicate**到远端，这都是由存储层来实现和提供相应接口。

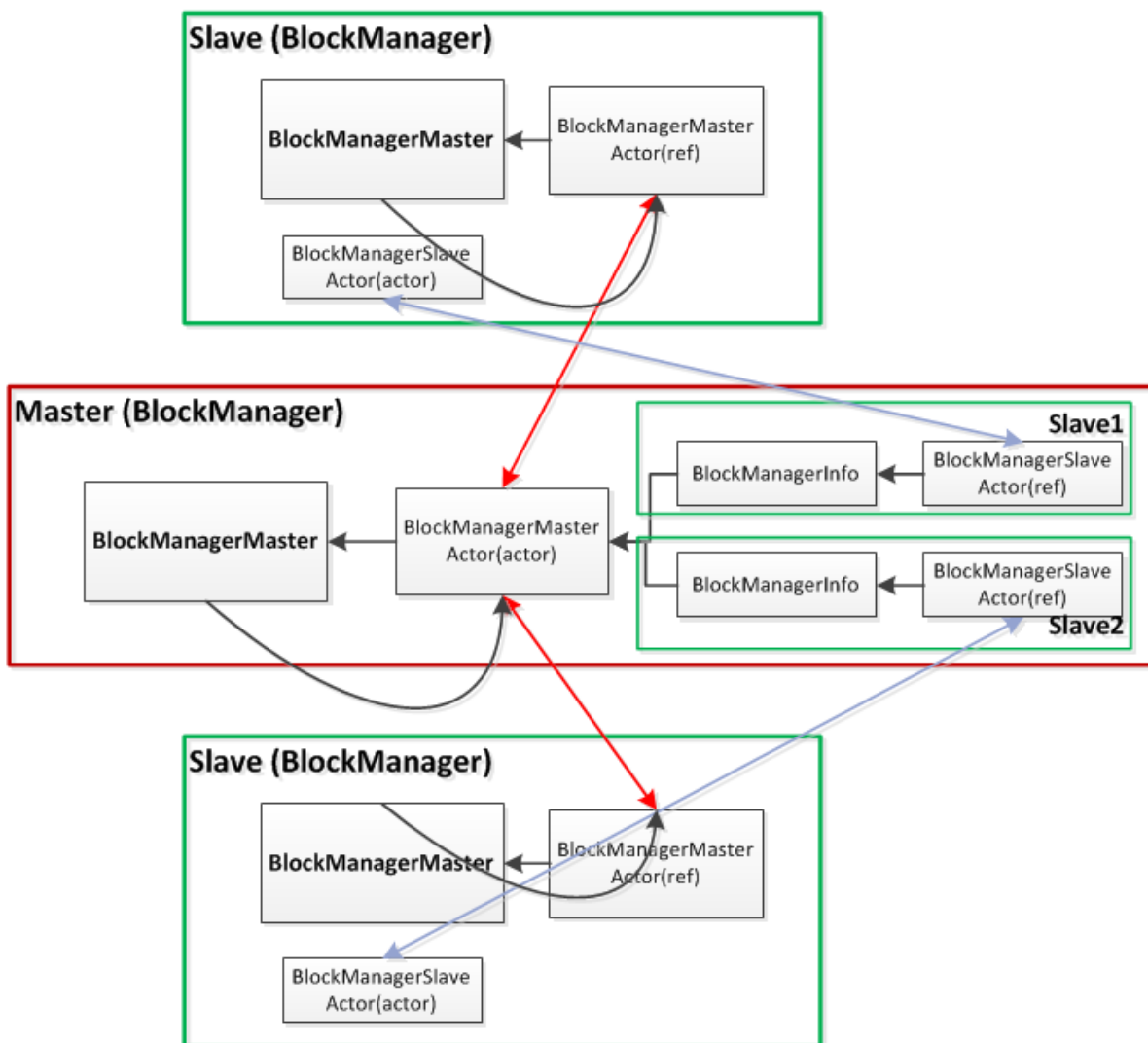
而其他模块若要和**storage**模块进行交互，**storage**模块提供了统一的操作类 **BlockManager**，外部类与**storage**模块打交道都需要通过调用 **BlockManager** 相应接口来实现。

## Storage模块通信层

首先来看一下通信层的UML类图：



其次我们来看看各个类在master和slave上所扮演的不同角色：



对于master和slave，`BlockManager` 的创建有所不同：

- Master (client driver)

`BlockManagerMaster` 拥有 `BlockManagerMasterActor` 的 *actor* 和所有 `BlockManagerSlaveActor` 的 *ref*。

- Slave (executor)

对于slave，`BlockManagerMaster` 则拥有 `BlockManagerMasterActor` 的 *ref* 和自身 `BlockManagerSlaveActor` 的 *actor*。

`BlockManagerMasterActor` 在`ref`和`actor`之间进行通信；`BlockManagerSlaveActor` 在`ref`和`actor`之间通信。

`actor`和`ref`。

`actor`和`ref`是Akka (<http://akka.io/>)中的两个不同的actor reference，分别由`actorOf` 和 `actorFor` 所创建。`actor`类似于网络服务中的server端，它保存所有的状态信息，接收client端请求并返回给客户端；`ref`类似于网络服务中的client端，通过向server端发起请求获取结果。

`BlockManager` wrap 了 `BlockManagerMaster`，通过 `BlockManagerMaster` 进行通信。Spark 会在 client driver 和 executor 端创建各自的 `BlockManager`，通过 `BlockManager` 对 `storage` 模块进行操作。

`BlockManager` 对象在 `SparkEnv` 中被创建，创建的过程如下所示：

```
1. def registerOrLookup(name: String, newActor: => Actor): ActorRef = {
2.   if (isDriver) {
3.     logInfo("Registering " + name)
4.     actorSystem.actorOf(Props(newActor), name = name)
5.   } else {
6.     val driverHost: String = System.getProperty("spark.driver.host", "localhost")
7.     val driverPort: Int = System.getProperty("spark.driver.port", "7077").toInt
8.     Utils.checkHost(driverHost, "Expected hostname")
9.     val url = "akka://spark@%s:%s/user/%s".format(driverHost, driverPort, name)
10.    logInfo("Connecting to " + name + ": " + url)
11.    actorSystem.actorFor(url)
12.  }
13. }
14.
15. val blockManagerMaster = new BlockManagerMaster(registerOrLookup(
16.  "BlockManagerMaster",
17.  new BlockManagerMasterActor(isLocal)))
18. val blockManager = new BlockManager(executorId, actorSystem, blockManagerMaster, serializer)
```

可以看到对于 client driver 和 executor，Spark 分别创建了 `BlockManagerMasterActor` `actor` 和 `ref`，并被 wrap 到 `BlockManager` 中。

## 通信层传递的消息

### • `BlockManagerMasterActor`

#### \* `executor to client driver`

1. `RegisterBlockManager` (executor 创建 `BlockManager` 以后向 client driver 发送请求注册自身)
2. `HeartBeat`
3. `UpdateBlockInfo` (更新 block 信息)
4. `GetPeers` (请求获得其他 `BlockManager` 的 id)
5. `GetLocations` (获取 block 所在的 `BlockManager` 的 id)
6. `GetLocationsMultipleBlockIds` (获取一组 block 所在的 `BlockManager` id)

#### \* `client driver to client driver`

1. `GetLocations` (获取 block 所在的 `BlockManager` 的 id)
2. `GetLocationsMultipleBlockIds` (获取一组 block 所在的 `BlockManager` id)
3. `RemoveExecutor` (删除所保存的已经死亡的 executor 上的 `BlockManager`)
4. `StopBlockManagerMaster` (停止 client driver 上的 `BlockManagerMasterActor`)

有些消息例如 `GetLocations` 在 `executor` 端和 `client driver` 端都会向 `actor` 请求，而其他的消息比如 `RegisterBlockManager` 只会由 `executor` 端的 `ref` 向 `client driver` 端的 `actor` 发送，于此同时例如 `RemoveExecutor` 则只会由 `client driver` 端的 `ref` 向 `client driver` 端的 `actor` 发送。具体消息是从哪里发送，哪里接收和处理请看代码细节，在这里就不再赘述了。

- **BlockManagerSlaveActor**

\* **client driver to executor**

1. **RemoveBlock** (删除block)
2. **RemoveRdd** (删除RDD)

通信层中涉及许多控制消息和状态消息的传递以及处理，这些细节可以直接查看源码，这里就不在一一罗列。下面就只简单介绍一下 `exeuctor` 端的 `BlockManager` 是如何启动以及向 `client driver` 发送注册请求完成注册。

## Register BlockManager

前面已经介绍了 `BlockManager` 对象是如何被创建出来的，当 `BlockManager` 被创建出来以后需要向 `client driver` 注册自己，下面我们来看一下这个流程：

首先 `BlockManager` 会调用 `initialize()` 初始化自己

```
1. private def initialize() {
2.   master.registerBlockManager(blockManagerId, maxMemory, slaveActor)
3.   ...
4.   if (!BlockManager.getDisableHeartBeatsForTesting) {
5.     heartBeatTask = actorSystem.scheduler.schedule(0.seconds, heartBeatFrequency.milliseconds) {
6.       heartBeat()
7.     }
8.   }
9. }
```

在 `initialized()` 函数中首先调用 `BlockManagerMaster` 向 `client driver` 注册自己，同时设置 `heartbeat` 定时器，定时发送 `heartbeat` 报文。可以看到在注册自身的时候向 `client driver` 传递了自身的 `slaveActor`，`client driver` 收到 `slaveActor` 以后会将其与之对应的 `BlockManagerInfo` 存储到 `hash map` 中，以便后续通过 `slaveActor` 向 `executor` 发送命令。

`BlockManagerMaster` 会将注册请求包装成 `RegisterBlockManager` 报文发送给 `client driver` 的 `BlockManagerMasterActor`，`BlockManagerMasterActor` 调用 `register()` 函数注册 `BlockManager`：

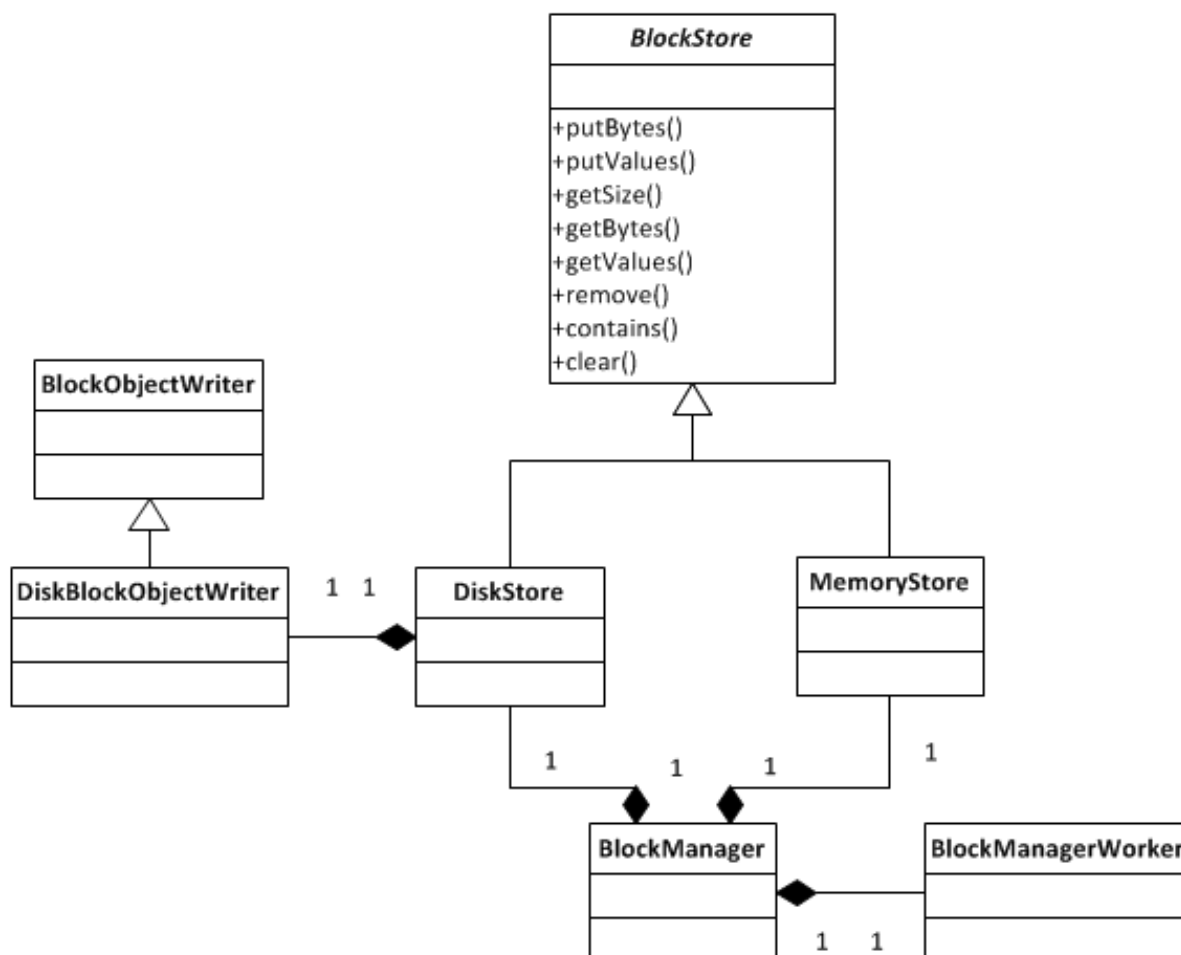
```
1. private def register(id: BlockManagerId, maxMemSize: Long, slaveActor: ActorRef) {
2.   if (id.executorId == "<driver>" && !isLocal) {
3.     // Got a register message from the master node; don't register it
4.   } else if (!blockManagerInfo.contains(id)) {
5.     blockManagerIdByExecutor.get(id.executorId) match {
6.       case Some(manager) =>
7.         // A block manager of the same executor already exists.
8.         // This should never happen. Let's just quit.
9.         logError("Got two different block manager registrations on " + id.executorId)
10.        System.exit(1)
11.       case None =>
12.         blockManagerIdByExecutor(id.executorId) = id
13.     }
14.     blockManagerInfo(id) = new BlockManagerMasterActor.BlockManagerInfo(
15.       id, System.currentTimeMillis(), maxMemSize, slaveActor)
16.   }
17. }
```

需要注意的是在client driver端也会执行上述过程，只是在最后注册的时候如果判断是 "<driver>" 就不进行任何操作。可以看到对应的 BlockManagerInfo 对象被创建并保存在hash map中。

## Storage模块存储层

在RDD层面上我们了解到RDD是由不同的partition组成的，我们所进行的transformation和action是在partition上面进行的；而在storage模块内部，RDD又被视为由不同的block组成，对于RDD的存取是以block为单位进行的，本质上partition和block是等价的，只是看待的角度不同。在Spark storage模块中存取数据的最小单位是block，所有的操作都是以block为单位进行的。

首先我们来看一下存储层的UML类图：



`BlockManager` 对象被创建的时候会创建出 `MemoryStore` 和 `DiskStore` 对象用以存取block，同时在 `initialize()` 函数中创建 `BlockManagerWorker` 对象用以监听远程的block存取请求来进行相应处理。

```

1. private[storage] val memoryStore: BlockStore = new MemoryStore(this, maxMemory)
2. private[storage] val diskStore: DiskStore =
3.   new DiskStore(this, System.getProperty("spark.local.dir", System.getProperty("java.io.tmpdir")))
4.
5. private def initialize() {
6.   ...
7.   BlockManagerWorker.startBlockManagerWorker(this)
8.   ...
9. }
  
```

下面就具体介绍一下对于 `DiskStore` 和 `MemoryStore`，block的存取操作是怎样进行的。

## DiskStore如何存取block

`DiskStore` 可以配置多个folder，Spark会在不同的folder下面创建Spark文件夹，文件夹的命名方式为(spark-local-yyyyMMddHHmmss-xxx, xxx是一个随机数)，所有的block都会存储在所创建的folder里面。`DiskStore` 会在对象被创建时调用 `createLocalDirs()` 来创建文件夹：

```

1. private def createLocalDirs(): Array[File] = {
2.     logDebug("Creating local directories at root dirs '" + rootDirs + "'")
3.     val dateFormat = new SimpleDateFormat("yyyyMMddHHmmss")
4.     rootDirs.split(",").map { rootDir =>
5.         var foundLocalDir = false
6.         var localDir: File = null
7.         var localDirId: String = null
8.         var tries = 0
9.         val rand = new Random()
10.        while (!foundLocalDir && tries < MAX_DIR_CREATION_ATTEMPTS) {
11.            tries += 1
12.            try {
13.                localDirId = "%s-%04x".format(dateFormat.format(new Date), rand.nextInt(65536))
14.                localDir = new File(rootDir, "spark-local-" + localDirId)
15.                if (!localDir.exists) {
16.                    foundLocalDir = localDir.mkdirs()
17.                }
18.            } catch {
19.                case e: Exception =>
20.                    logWarning("Attempt " + tries + " to create local dir " + localDir + " failed", e)
21.            }
22.        }
23.        if (!foundLocalDir) {
24.            logError("Failed " + MAX_DIR_CREATION_ATTEMPTS +
25.                " attempts to create local dir in " + rootDir)
26.            System.exit(ExecutorExitCode.DISK_STORE_FAILED_TO_CREATE_DIR)
27.        }
28.        logInfo("Created local directory at " + localDir)
29.        localDir
30.    }
31. }

```

在 `DiskStore` 里面，每一个 `block` 都被存储为一个 `file`，通过计算 `block id` 的 `hash` 值将 `block` 映射到文件中，`block id` 与文件路径的映射关系如下所示：

```

1. private def getFile(blockId: String): File = {
2.     logDebug("Getting file for block " + blockId)
3.
4.     // Figure out which local directory it hashes to, and which subdirectory in that
5.     val hash = Utils.nonNegativeHash(blockId)
6.     val dirId = hash % localDirs.length
7.     val subDirId = (hash / localDirs.length) % subDirsPerLocalDir
8.
9.     // Create the subdirectory if it doesn't already exist
10.    var subDir = subDirs(dirId)(subDirId)
11.    if (subDir == null) {
12.        subDir = subDirs(dirId).synchronized {
13.            val old = subDirs(dirId)(subDirId)
14.            if (old != null) {
15.                old
16.            } else {
17.                val newDir = new File(localDirs(dirId), "%02x".format(subDirId))
18.                newDir.mkdir()
19.                subDirs(dirId)(subDirId) = newDir
20.                newDir
21.            }
22.        }
23.    }
24.
25.    new File(subDir, blockId)
26. }

```

根据block id计算出hash值，将hash取模获得 dirId 和 subDirId，在 subDirs 中找出相应的 subDir，若没有则新建一个 subDir，最后以 subDir 为路径、block id为文件名创建file handler，DiskStore 使用此file handler将block写入文件内，代码如下所示：

```

1. override def putBytes(blockId: String, _bytes: ByteBuffer, level: StorageLevel) {
2.     // So that we do not modify the input offsets !
3.     // duplicate does not copy buffer, so inexpensive
4.     val bytes = _bytes.duplicate()
5.     logDebug("Attempting to put block " + blockId)
6.     val startTime = System.currentTimeMillis
7.     val file = createFile(blockId)
8.     val channel = new RandomAccessFile(file, "rw").getChannel()
9.     while (bytes.remaining > 0) {
10.        channel.write(bytes)
11.    }
12.    channel.close()
13.    val finishTime = System.currentTimeMillis
14.    logDebug("Block %s stored as %s file on disk in %d ms".format(
15.        blockId, Utils.bytesToString(bytes.limit), (finishTime - startTime)))
16. }

```

而获取block则非常简单，找到相应的文件并读取出来即可：



```
1. override def getBytes(blockId: String): Option[ByteBuffer] = {
2.   val file = getFile(blockId)
3.   val bytes = getFileBytes(file)
4.   Some(bytes)
5. }
```

因此在 `DiskStore` 中存取block首先是要将block id映射成相应的文件路径，接着存取文件就可以了。

## MemoryStore如何存取block

相对于 `DiskStore` 需要根据block id hash计算出文件路径并将block存放到对应的文件里面，`MemoryStore` 管理block就显得非常简单：`MemoryStore` 内部维护了一个hash map来管理所有的block，以block id为key将block存放到hash map中。

```
1. case class Entry(value: Any, size: Long, deserialized: Boolean)
2.
3. private val entries = new LinkedHashMap[String, Entry](32, 0.75f, true)
```

在 `MemoryStore` 中存放block必须确保内存足够容纳下该block，若内存不足则会将block写到文件中，具体的代码如下所示：

```
1. override def putBytes(blockId: String, _bytes: ByteBuffer, level: StorageLevel) {
2.   // Work on a duplicate - since the original input might be used elsewhere.
3.   val bytes = _bytes.duplicate()
4.   bytes.rewind()
5.   if (level.deserialized) {
6.     val values = blockManager.dataDeserialize(blockId, bytes)
7.     val elements = new ArrayBuffer[Any]
8.     elements += values
9.     val sizeEstimate = SizeEstimator.estimate(elements.asInstanceOf[AnyRef])
10.    tryToPut(blockId, elements, sizeEstimate, true)
11.  } else {
12.    tryToPut(blockId, bytes, bytes.limit, false)
13.  }
14. }
```

在 `tryToPut()` 中，首先调用 `ensureFreeSpace()` 确保空闲内存是否足以容纳block，若可以则将该block放入hash map中进行管理；若不足以容纳则通过调用 `dropFromMemory()` 将block写入文件。

```

1. private def tryToPut(blockId: String, value: Any, size: Long, deserialized: Boolean): Boolean = {
2.     // TODO: Its possible to optimize the locking by locking entries only when selecting blocks
3.     // to be dropped. Once the to-be-dropped blocks have been selected, and lock on entries has been
4.     // released, it must be ensured that those to-be-dropped blocks are not double counted for
5.     // freeing up more space for another block that needs to be put. Only then the actually dropping
6.     // of blocks (and writing to disk if necessary) can proceed in parallel.
7.     putLock.synchronized {
8.         if (ensureFreeSpace(blockId, size)) {
9.             val entry = new Entry(value, size, deserialized)
10.            entries.synchronized {
11.                entries.put(blockId, entry)
12.                currentMemory += size
13.            }
14.            if (deserialized) {
15.                logInfo("Block %s stored as values to memory (estimated size %s, free %s)".format(
16.                    blockId, Utils.bytesToString(size), Utils.bytesToString(freeMemory)))
17.            } else {
18.                logInfo("Block %s stored as bytes to memory (size %s, free %s)".format(
19.                    blockId, Utils.bytesToString(size), Utils.bytesToString(freeMemory)))
20.            }
21.            true
22.        } else {
23.            // Tell the block manager that we couldn't put it in memory so that it can drop it to
24.            // disk if the block allows disk storage.
25.            val data = if (deserialized) {
26.                Left(value.asInstanceOf[ArrayBuffer[Any]])
27.            } else {
28.                Right(value.asInstanceOf[ByteBuffer].duplicate())
29.            }
30.            blockManager.dropFromMemory(blockId, data)
31.            false
32.        }
33.    }
34. }

```

而从 `MemoryStore` 中取得block则非常简单，只需从hash map中取出block id对应的value即可。

```

1. override def getValues(blockId: String): Option[Iterator[Any]] = {
2.     val entry = entries.synchronized {
3.         entries.get(blockId)
4.     }
5.     if (entry == null) {
6.         None
7.     } else if (entry.deserialized) {
8.         Some(entry.value.asInstanceOf[ArrayBuffer[Any]].iterator)
9.     } else {
10.        val buffer = entry.value.asInstanceOf[ByteBuffer].duplicate() // Doesn't actually copy data
11.        Some(blockManager.dataDeserialize(blockId, buffer))
12.    }
13. }

```

## Put or Get block through BlockManager

上面介绍了 `DiskStore` 和 `MemoryStore` 对于 `block` 的存取操作，那么我们是要直接与它们交互存取数据吗，还是封装了更抽象的接口使我们无需关心底层？

`BlockManager` 为我们提供了 `put()` 和 `get()` 函数，用户可以使用这两个函数对 `block` 进行存取而无需关心底层实现。

首先我们来看一下 `put()` 函数的实现：

```
1. def put(blockId: String, values: ArrayBuffer[Any], level: StorageLevel,
2.   tellMaster: Boolean = true) : Long = {
3.
4.   ...
5.
6.   // Remember the block's storage level so that we can correctly drop it to disk if it needs
7.   // to be dropped right after it got put into memory. Note, however, that other threads will
8.   // not be able to get() this block until we call markReady on its BlockInfo.
9.   val myInfo = {
10.    val tinfo = new BlockInfo(level, tellMaster)
11.    // Do atomically !
12.    val oldBlockOpt = blockInfo.putIfAbsent(blockId, tinfo)
13.
14.    if (oldBlockOpt.isDefined) {
15.      if (oldBlockOpt.get.waitForReady()) {
16.        logWarning("Block " + blockId + " already exists on this machine; not re-adding it")
17.        return oldBlockOpt.get.size
18.      }
19.
20.      // TODO: So the block info exists - but previous attempt to load it (?) failed. What do we d
o now ? Retry on it ?
21.      oldBlockOpt.get
22.    } else {
23.      tinfo
24.    }
25.  }
26.
27.  val startTimeMs = System.currentTimeMillis
28.
29.  // If we need to replicate the data, we'll want access to the values, but because our
30.  // put will read the whole iterator, there will be no values left. For the case where
31.  // the put serializes data, we'll remember the bytes, above; but for the case where it
32.  // doesn't, such as deserialized storage, let's rely on the put returning an Iterator.
33.  var valuesAfterPut: Iterator[Any] = null
34.
35.  // Ditto for the bytes after the put
36.  var bytesAfterPut: ByteBuffer = null
37.
38.  // Size of the block in bytes (to return to caller)
39.  var size = 0L
40.
41.  myInfo.synchronized {
42.    logTrace("Put for block " + blockId + " took " + Utils.getUsedTimeMs(startTimeMs)
43.      + " to get into synchronized block")
44.
45.    var marked = false
46.    try {
47.      if (level.useMemory) {
48.        // Save it just to memory first, even if it also has useDisk set to true; we will later
```

```

49.         // drop it to disk if the memory store can't hold it.
50.         val res = memoryStore.putValues(blockId, values, level, true)
51.         size = res.size
52.         res.data match {
53.             case Right(newBytes) => bytesAfterPut = newBytes
54.             case Left(newIterator) => valuesAfterPut = newIterator
55.         }
56.     } else {
57.         // Save directly to disk.
58.         // Don't get back the bytes unless we replicate them.
59.         val askForBytes = level.replication > 1
60.         val res = diskStore.putValues(blockId, values, level, askForBytes)
61.         size = res.size
62.         res.data match {
63.             case Right(newBytes) => bytesAfterPut = newBytes
64.             case _ =>
65.         }
66.     }
67.
68.     // Now that the block is in either the memory or disk store, let other threads read it,
69.     // and tell the master about it.
70.     marked = true
71.     myInfo.markReady(size)
72.     if (tellMaster) {
73.         reportBlockStatus(blockId, myInfo)
74.     }
75. } finally {
76.     // If we failed at putting the block to memory/disk, notify other possible readers
77.     // that it has failed, and then remove it from the block info map.
78.     if (! marked) {
79.         // Note that the remove must happen before markFailure otherwise another thread
80.         // could've inserted a new BlockInfo before we remove it.
81.         blockInfo.remove(blockId)
82.         myInfo.markFailure()
83.         logWarning("Putting block " + blockId + " failed")
84.     }
85. }
86. }
87. logDebug("Put block " + blockId + " locally took " + Utils.getUsedTimeMs(startTimeMs))
88.
89. // Replicate block if required
90. if (level.replication > 1) {
91.     val remoteStartTime = System.currentTimeMillis
92.     // Serialize the block if not already done
93.     if (bytesAfterPut == null) {
94.         if (valuesAfterPut == null) {
95.             throw new SparkException(
96.                 "Underlying put returned neither an Iterator nor bytes! This shouldn't happen.")
97.         }
98.         bytesAfterPut = dataSerialize(blockId, valuesAfterPut)
99.     }
100.    replicate(blockId, bytesAfterPut, level)
101.    logDebug("Put block " + blockId + " remotely took " + Utils.getUsedTimeMs(remoteStartTime))
102. }
103. BlockManager.dispose(bytesAfterPut)

```

```
004.  
005.     return size  
006. }
```

对于 `put()` 操作，主要分为以下3个步骤：

1. 为block创建 `BlockInfo` 结构体存储block相关信息，同时将其加锁使其不能被访问。
2. 根据block的storage level将block存储到memory或是disk上，同时解锁标识该block已经ready，可被访问。
3. 根据block的replication数决定是否将该block replicate到远端。

接着我们来看一下 `get()` 函数的实现：

```
1. def get(blockId: String): Option[Iterator[Any]] = {  
2.     val local = getLocal(blockId)  
3.     if (local.isDefined) {  
4.         logInfo("Found block %s locally".format(blockId))  
5.         return local  
6.     }  
7.     val remote = getRemote(blockId)  
8.     if (remote.isDefined) {  
9.         logInfo("Found block %s remotely".format(blockId))  
10.        return remote  
11.    }  
12.    None  
13. }
```

`get()` 首先会从local的 `BlockManager` 中查找block，如果找到则返回相应的block，若local没有找到该block，则发起请求从其他的executor上的 `BlockManager` 中查找block。在通常情况下Spark任务的分配是根据block的分布决定的，任务往往会被分配到拥有block的节点上，因此 `getLocal()` 就能找到所需的block；但是在资源有限的情况下，Spark会将任务调度到与block不同的节点上，这样就必须通过 `getRemote()` 来获得block。

我们先来看一下 `getLocal()`：

```
1. def getLocal(blockId: String): Option[Iterator[Any]] = {  
2.     logDebug("Getting local block " + blockId)  
3.     val info = blockInfo.get(blockId).orNull  
4.     if (info != null) {  
5.         info.synchronized {  
6.  
7.             // In the another thread is writing the block, wait for it to become ready.  
8.             if (!info.waitForReady()) {  
9.                 // If we get here, the block write failed.  
10.                logWarning("Block " + blockId + " was marked as failure.")  
11.                return None  
12.            }  
13.  
14.            val level = info.level  
15.            logDebug("Level for block " + blockId + " is " + level)  
16.  
17.            // Look for the block in memory  
18.            if (level.useMemory) {  
19.                logDebug("Getting block " + blockId + " from memory")  
20.                memoryStore.getValues(blockId) match {  
21.                    case Some(iterator) =>  
22.                        return Some(iterator)  
23.                }  
24.            }  
25.        }  
26.    }  
27.    None  
28. }
```

```

23.         case None =>
24.             logDebug("Block " + blockId + " not found in memory")
25.         }
26.     }
27.
28.     // Look for block on disk, potentially loading it back into memory if required
29.     if (level.useDisk) {
30.         logDebug("Getting block " + blockId + " from disk")
31.         if (level.useMemory && level.deserialized) {
32.             diskStore.getValues(blockId) match {
33.                 case Some(iterator) =>
34.                     // Put the block back in memory before returning it
35.                     // TODO: Consider creating a putValues that also takes in a iterator ?
36.                     val elements = new ArrayBuffer[Any]
37.                     elements += iterator
38.                     memoryStore.putValues(blockId, elements, level, true).data match {
39.                         case Left(iterator2) =>
40.                             return Some(iterator2)
41.                         case _ =>
42.                             throw new Exception("Memory store did not return back an iterator")
43.                     }
44.                 case None =>
45.                     throw new Exception("Block " + blockId + " not found on disk, though it should be")
46.             }
47.         } else if (level.useMemory && !level.deserialized) {
48.             // Read it as a byte buffer into memory first, then return it
49.             diskStore.getBytes(blockId) match {
50.                 case Some(bytes) =>
51.                     // Put a copy of the block back in memory before returning it. Note that we can't
52.                     // put the ByteBuffer returned by the disk store as that's a memory-mapped file.
53.                     // The use of rewind assumes this.
54.                     assert (0 == bytes.position())
55.                     val copyForMemory = ByteBuffer.allocate(bytes.limit)
56.                     copyForMemory.put(bytes)
57.                     memoryStore.putBytes(blockId, copyForMemory, level)
58.                     bytes.rewind()
59.                     return Some(dataDeserialize(blockId, bytes))
60.                 case None =>
61.                     throw new Exception("Block " + blockId + " not found on disk, though it should be")
62.             }
63.         } else {
64.             diskStore.getValues(blockId) match {
65.                 case Some(iterator) =>
66.                     return Some(iterator)
67.                 case None =>
68.                     throw new Exception("Block " + blockId + " not found on disk, though it should be")
69.             }
70.         }
71.     }
72. }
73. } else {
74.     logDebug("Block " + blockId + " not registered locally")
75. }
76. return None
77. }

```

`getLocal()` 首先会根据block id获得相应的 `BlockInfo` 并从中取出该block的storage level，根据storage level的不同 `getLocal()` 又进入以下不同分支：

1. `level.useMemory == true`: 从memory中取出block并返回，若没有取到则进入分支2。
2. `level.useDisk == true`:
  - `level.useMemory == true`: 将block从disk中读出并写入内存以便下次使用时直接从内存中获得，同时返回该block。
  - `level.useMemory == false`: 将block从disk中读出并返回
3. `level.useDisk == false`: 没有在本地找到block，返回None。

接下来我们来看一下 `getRemote()`：

```
1. def getRemote(blockId: String): Option[Iterator[Any]] = {
2.   if (blockId == null) {
3.     throw new IllegalArgumentException("Block Id is null")
4.   }
5.   logDebug("Getting remote block " + blockId)
6.   // Get locations of block
7.   val locations = master.getLocations(blockId)
8.
9.   // Get block from remote locations
10.  for (loc <- locations) {
11.    logDebug("Getting remote block " + blockId + " from " + loc)
12.    val data = BlockManagerWorker.syncGetBlock(
13.      GetBlock(blockId), ConnectionManagerId(loc.host, loc.port))
14.    if (data != null) {
15.      return Some(dataDeserialize(blockId, data))
16.    }
17.    logDebug("The value of block " + blockId + " is null")
18.  }
19.  logDebug("Block " + blockId + " not found")
20.  return None
21. }
```

`getRemote()` 首先取得该block的所有location信息，然后根据location向远端发送请求获取block，只要有一个远端返回block该函数就返回而不继续发送请求。

至此我们简单介绍了 `BlockManager` 类中的 `get()` 和 `put()` 函数，使用这两个函数外部类可以轻易地存取block数据。

## Partition如何转化为Block

在storage模块里面所有的操作都是和block相关的，但是在RDD里面所有的运算都是基于partition的，那么partition是如何与block对应上的呢？

RDD计算的核心函数是 `iterator()` 函数：

```
1. final def iterator(split: Partition, context: TaskContext): Iterator[T] = {
2.   if (storageLevel != StorageLevel.NONE) {
3.     SparkEnv.get.cacheManager.getOrCompute(this, split, context, storageLevel)
4.   } else {
5.     computeOrReadCheckpoint(split, context)
6.   }
7. }
```

如果当前RDD的storage level不是NONE的话，表示该RDD在BlockManager中有存储，那么调用CacheManager中的getOrCompute()函数计算RDD，在这个函数中partition和block发生了关系：

首先根据RDD id和partition index构造出block id (rdd\_xx\_xx)，接着从BlockManager中取出相应的block。

- 如果该block存在，表示此RDD在之前已经被计算过和存储在BlockManager中，因此取出即可，无需再重新计算。
- 如果该block不存在则需要调用RDD的computeOrReadCheckpoint()函数计算出新的block，并将其存储到BlockManager中。

需要注意的是block的计算和存储是阻塞的，若另一线程也需要用到此block则需等到该线程block的loading结束。



```

1. def getOrCompute[T](rdd: RDD[T], split: Partition, context: TaskContext, storageLevel: StorageLevel)
2.   : Iterator[T] = {
3.     val key = "rdd_%d_%d".format(rdd.id, split.index)
4.     logDebug("Looking for partition " + key)
5.     blockManager.get(key) match {
6.       case Some(values) =>
7.         // Partition is already materialized, so just return its values
8.         return values.asInstanceOf[Iterator[T]]
9.
10.    case None =>
11.      // Mark the split as loading (unless someone else marks it first)
12.      loading.synchronized {
13.        if (loading.contains(key)) {
14.          logInfo("Another thread is loading %s, waiting for it to finish...".format (key))
15.          while (loading.contains(key)) {
16.            try {loading.wait()} catch {case _ : Throwable =>}
17.          }
18.          logInfo("Finished waiting for %s".format(key))
19.          // See whether someone else has successfully loaded it. The main way this would fail
20.          // is for the RDD-level cache eviction policy if someone else has loaded the same RDD
21.          // partition but we didn't want to make space for it. However, that case is unlikely
22.          // because it's unlikely that two threads would work on the same RDD partition. One
23.          // downside of the current code is that threads wait serially if this does happen.
24.          blockManager.get(key) match {
25.            case Some(values) =>
26.              return values.asInstanceOf[Iterator[T]]
27.            case None =>
28.              logInfo("Whoever was loading %s failed; we'll try it ourselves".format (key))
29.              loading.add(key)
30.          }
31.        } else {
32.          loading.add(key)
33.        }
34.      }
35.      try {
36.        // If we got here, we have to load the split
37.        logInfo("Partition %s not found, computing it".format(key))
38.        val computedValues = rdd.computeOrReadCheckpoint(split, context)
39.        // Persist the result, so long as the task is not running locally
40.        if (context.runningLocally) { return computedValues }
41.        val elements = new ArrayBuffer[Any]
42.        elements += computedValues
43.        blockManager.put(key, elements, storageLevel, true)
44.        return elements.iterator.asInstanceOf[Iterator[T]]
45.      } finally {
46.        loading.synchronized {
47.          loading.remove(key)
48.          loading.notifyAll()
49.        }
50.      }
51.    }
52.  }

```

这样RDD的transformation、action就和block数据建立了联系，虽然抽象上我们的操作是在partition层面上进行的，但是partition最终还是被映射成为block，因此实际上我们的所有操作都是对block的处理和存取。

# End

本文就storage模块的两个层面进行了介绍-通信层和存储层。通信层中简单介绍了类结构和组成以及类在通信层中所扮演的不同角色，还有不同角色之间通信的报文，同时简单介绍了通信层的启动和注册细节。存储层中分别介绍了DiskStore和MemoryStore中对于block的存和取的实现代码，同时分析了BlockManager中put()和get()接口，最后简单介绍了Spark RDD中的partition与BlockManager中的block之间的关系，以及如何交互存取block的。

本文从整体上分析了storage模块的实现，并未就具体实现做非常细节的分析，相信在看完本文对storage模块有一个整体的印象以后再去分析细节的实现会有事半功倍的效果。

[← Previous \(/functional%20programming/2013/08/30/Functional-Abstraction-of-Sequence-Spark-API-Design\)](#)

[Archive \(/archive.html\)](#) [Next →](#)

1 comment



Join the discussion...

Best ▾ Community

Share

 Login ▾

Wangda Tan • 14 days ago

楼主太强大了，这些正是最近想学的，谢谢楼主的分享！！

| • Reply • Share ›

ALSO ON JERRY SHAO'S HOMEPAGE

WHAT'S THIS?

Spark Overview

3 comments • 7 months ago

vincent\_hv —

Spark源码分析之-scheduler模块

19 comments • 7 months ago

jerryshao —

Spark源码分析之-deploy模块

3 comments • 8 months ago

Huangdong Meng — 明白～ 期待大神的更多大作哈～ 比如shark～ RDD的部分复杂的operator的实现 等data processing层的讲解哈～



## CATEGORIES

lessons (1) (/categories.html#lessons-ref)

test (1) (/categories.html#test-ref)

architecture (6) (/categories.html#architecture-ref)

functional programming (1) (/categories.html#functional programming-ref)

arhitecture (1) (/categories.html#arhitecture-ref)

## LINKS

阮一峰的网络日志 (<http://www.ruanyifeng.com/blog/>)

刘未鹏 (<http://mindhacks.cn/>)

酷壳 (<http://coolshell.cn/>)

BeiYuu.com (<http://beiyuu.com/>)

## MY FAVORITES