

## 深入理解Spark 2.1 Core （十一）： Shuffle Reduce 端的原理与源码分析

■ 版权声明：本文为博主原创文章，转载请附上原文地址。 <http://blog.csdn.net/u011239443/article/details/56843264>

<http://blog.csdn.net/u011239443/article/details/56843264>

在《深入理解Spark 2.1 Core （九）： 迭代计算和Shuffle的原理与源码分析 》我们讲解了，以传统Hadoop MapReduce类似的从HDFS中读取数据，再到 `rdd.HadoopRDD.compute` 便可以调用函数 `f`，即 `map` 中的函数的过程。在《深入理解Spark 2.1 Core （十）： Shuffle map端的原理与源码分析》 我们深入讲解了 `sorter.insertAll(records)`，即如何对数据进行排序并写入内存缓冲区。

我们曾经在《深入理解Spark 2.1 Core （一）： RDD的原理与源码分析 》讲解过：

为了有效地实现容错，RDD提供了一种高度受限的共享内存，即RDD是只读的，并且只能通过其他RDD上的批量操作来创建（注：还可以由外部存储系数数据集创建，如HDFS）

可知，我们在第九，第十篇博文所讲的是传统Hadoop MapReduce类似的，在最初从HDFS中读取数据生成 `HadoopRDD` 的过程。而RDD可以通过其他RDD上的批量操作来创建，所以这里的 `HadoopRDD` 对于下一个生成的 `ShuffledRDD` 可以视为 `Map` 端，当然下一个生成的 `ShuffledRDD` 可以被下下个 `ShuffledRDD` 视为 `Map` 端。反过来说，下一个 `ShuffledRDD` 可以被 `HadoopRDD` 视作 `Reduce` 端。

这篇博文，我们就来讲下 `Shuffle` 的 `Reduce` 端。其实在 `RDD` 迭代部分和第九篇博文类似，不同的是，这里调用的是 `rdd.ShuffledRDD.compute`：

```
1  override def compute(split: Partition, context: TaskContext): Iterator[(K, C)] = {
2  // 得到依赖
3      val dep = dependencies.head.asInstanceOf[ShuffleDependency[K, V, C]]
4      // 调用getReader, 传入dep.shuffleHandle 分区 上下文
5      // 得到Reader, 调用read()
6      // 得到迭代器
7      SparkEnv.get.shuffleManager.getReader(dep.shuffleHandle, split.index, split.index + 1, context)
8          .read()
9          .asInstanceOf[Iterator[(K, C)]]
10 }
```

这里调用的是 `shuffle.sort.SortShuffleManager` 的 `getReader`：

```
1  override def getReader[K, C](
2      handle: ShuffleHandle,
3      startPartition: Int,
4      endPartition: Int,
5      context: TaskContext): ShuffleReader[K, C] = {
6      // 生成返回 BlockStoreShuffleReader
7      new BlockStoreShuffleReader(
8          handle.asInstanceOf[BaseShuffleHandle[K, _, C]], startPartition, endPartition, context)
9  }
```

`shuffle.BlockStoreShuffleReader.read`：

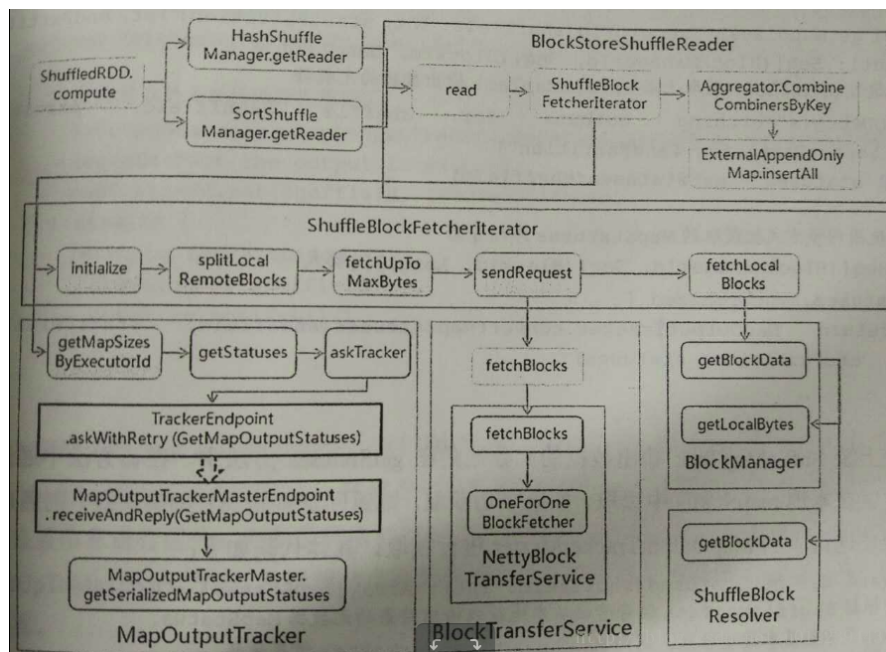
```
1  override def read(): Iterator[Product2[K, C]] = {
2  // 实例化ShuffleBlockFetcherIterator
3      val blockFetcherItr = new ShuffleBlockFetcherIterator(
4          context,
5          blockManager.shuffleClient,
6          blockManager,
7          // 通过消息发送获取 ShuffleMapTask 存储数据位置的元数据
8          mapOutputTracker.getMapSizesByExecutorId(handle.shuffleId, startPartition, endPartition),
9          // 设置每次传输的大小
10         SparkEnv.get.conf.getSizeAsMb("spark.reducer.maxSizeInFlight", "48m") * 1024 * 1024,
11         // // 设置Int的大小
12         SparkEnv.get.conf.getInt("spark.reducer.maxReqsInFlight", Int.MaxValue))
```

```

13
14 // 基于配置的压缩和加密来包装流
15 val wrappedStreams = blockFetcherItr.map { case (blockId, inputStream) =>
16     serializerManager.wrapStream(blockId, inputStream)
17 }
18
19 val serializerInstance = dep.serializer.newInstance()
20
21 // 对每个流生成 k/v 迭代器
22 val recordIter = wrappedStreams.flatMap { wrappedStream =>
23     serializerInstance.deserializeStream(wrappedStream).asKeyValueIterator
24 }
25
26
27 // 每条记录读取后更新任务度量
28 val readMetrics = context.taskMetrics.createTempShuffleReadMetrics()
29 // 生成完整的迭代器
30 val metricIter = CompletionIterator[(Any, Any), Iterator[(Any, Any)]](
31     recordIter.map { record =>
32         readMetrics.incRecordsRead(1)
33         record
34     },
35     context.taskMetrics().mergeShuffleReadMetrics())
36
37 // 传入metricIter到可中断的迭代器
38 // 为了能取消迭代
39 val interruptibleIter = new InterruptibleIterator[(Any, Any)](context, metricIter)
40
41 val aggregatedIter: Iterator[Product2[K, C]] = if (dep.aggregator.isDefined) {
42     // 若需要对数据进行聚合
43     if (dep.mapSideCombine) {
44         // 若需要进行Map端 (对于下一个Shuffle来说) 的合并
45         val combinedKeyValuesIterator = interruptibleIter.asInstanceOf[Iterator[(K, C)]]
46         dep.aggregator.get.combineCombinersByKey(combinedKeyValuesIterator, context)
47         // 若只需要进行Reduce端 (对于下一个Shuffle来说) 的合并
48     } else {
49         val keyValuesIterator = interruptibleIter.asInstanceOf[Iterator[(K, Nothing)]]
50         dep.aggregator.get.combineValuesByKey(keyValuesIterator, context)
51     }
52 } else {
53     require(!dep.mapSideCombine, "Map-side combine without Aggregator specified!")
54     interruptibleIter.asInstanceOf[Iterator[Product2[K, C]]]
55 }
56
57
58 dep.keyOrdering match {
59     case Some(keyOrd: Ordering[K]) =>
60         // 若需要排序
61         // 若spark.shuffle.spill设置为否的话
62         // 将不会spill到磁盘
63         val sorter =
64             new ExternalSorter[K, C, C](context, ordering = Some(keyOrd), serializer = dep.serializer)
65         sorter.insertAll(aggregatedIter)
66         context.taskMetrics().incMemoryBytesSpilled(sorter.memoryBytesSpilled)
67         context.taskMetrics().incDiskBytesSpilled(sorter.diskBytesSpilled)
68         context.taskMetrics().incPeakExecutionMemory(sorter.peakMemoryUsedBytes)
69         CompletionIterator[Product2[K, C], Iterator[Product2[K, C]]](sorter.iterator, sorter.stop())
70     case None =>
71         aggregatedIter
72 }
73 }

```

类调用关系图:



下面我们来深入讲解下实例化 `ShuffleBlockFetcherIterator` 的过程：

```

1 // 实例化ShuffleBlockFetcherIterator
2 val blockFetcherItr = new ShuffleBlockFetcherIterator(
3     context,
4     blockManager.shuffleClient,
5     blockManager,
6     // 通过消息发送获取 ShuffleMapTask 存储数据位置的元数据
7     mapOutputTracker.getMapSizesByExecutorId(handle.shuffleId, startPartition, endPartition),
8     // 设置每次传输的大小
9     SparkEnv.get.conf.getSizeAsMb("spark.reducer.maxSizeInFlight", "48m") * 1024 * 1024,
10    // // 设置Int的大小
11    SparkEnv.get.conf.getInt("spark.reducer.maxReqsInFlight", Int.MaxValue))

```

## 获取元数据

### mapOutputTracker.getMapSizesByExecutorId

首先我们会调用 `mapOutputTracker.getMapSizesByExecutorId`：

```

1 def getMapSizesByExecutorId(shuffleId: Int, startPartition: Int, endPartition: Int)
2   : Seq[(BlockManagerId, Seq[(BlockId, Long)])] = {
3   logDebug(s"Fetching outputs for shuffle $shuffleId, partitions $startPartition-$endPartition")
4   // 得到元数据
5   val statuses = getStatuses(shuffleId)
6   // 返回格式为:
7   // Seq[BlockManagerId,Seq[(shuffle block id, shuffle block size)]]
8   statuses.synchronized {
9       return MapOutputTracker.convertMapStatuses(shuffleId, startPartition, endPartition, statuses)
10  }
11 }

```

### mapOutputTracker.getStatuses

```

1 private def getStatuses(shuffleId: Int): Array[MapStatus] = {
2 // 尝试从本地获取数据
3 val statuses = mapStatuses.get(shuffleId).orNull
4 if (statuses == null) {
5 // 若本地无数据

```

```

6      logInfo("Don't have map outputs for shuffle " + shuffleId + ", fetching them")
7      val startTime = System.currentTimeMillis
8      var fetchedStatuses: Array[MapStatus] = null
9      fetching.synchronized {
10         // 若以及有其他人也准备远程获取这数据的话
11         // 则等待
12         while (fetching.contains(shuffleId)) {
13             try {
14                 fetching.wait()
15             } catch {
16                 case e: InterruptedException =>
17             }
18         }
19
20         // 尝试直接获取数据
21         fetchedStatuses = mapStatuses.get(shuffleId).OrNull
22         if (fetchedStatuses == null) {
23             // 若还是不得不远程获取,
24             // 则将shuffleId加入fetching
25             fetching += shuffleId
26         }
27     }
28
29     if (fetchedStatuses == null) {
30         logInfo("Doing the fetch; tracker endpoint = " + trackerEndpoint)
31         try {
32             // 远程获取
33             val fetchedBytes = askTracker(Array[Byte])(GetMapOutputStatuses(shuffleId))
34             // 反序列化
35             fetchedStatuses = MapOutputTracker.deserializeMapStatuses(fetchedBytes)
36             logInfo("Got the output locations")
37             // 将数据加入mapStatuses
38             mapStatuses.put(shuffleId, fetchedStatuses)
39         } finally {
40             fetching.synchronized {
41                 fetching -= shuffleId
42                 fetching.notifyAll()
43             }
44         }
45     }
46     logDebug(s"Fetching map output statuses for shuffle $shuffleId took " +
47         s"${System.currentTimeMillis - startTime} ms")
48
49     if (fetchedStatuses != null) {
50         // 若直接获取, 则直接返回
51         return fetchedStatuses
52     } else {
53         logError("Missing all output locations for shuffle " + shuffleId)
54         throw new MetadataFetchFailedException(
55             shuffleId, -1, "Missing all output locations for shuffle " + shuffleId)
56     }
57 } else {
58     // 若直接获取, 则直接返回
59     return statuses
60 }
61 }

```

## mapOutputTracker.askTracker

向 `trackerEndpoint` 发送消息 `GetMapOutputStatuses(shuffleId)`

```

1      protected def askTracker[T: ClassTag](message: Any): T = {
2          try {
3              trackerEndpoint.askWithRetry[T](message)

```

```

4      } catch {
5          case e: Exception =>
6              logError("Error communicating with MapOutputTracker", e)
7              throw new SparkException("Error communicating with MapOutputTracker", e)
8      }
9  }

```

## MapOutputTrackerMasterEndpoint.receiveAndReply

```

1      case GetMapOutputStatuses(shuffleId: Int) =>
2          val hostPort = context.senderAddress.hostPort
3          logInfo("Asked to send map output locations for shuffle " + shuffleId + " to " + hostPort)
4          val mapOutputStatuses = tracker.post(new GetMapOutputMessage(shuffleId, context))

```

可以看到，这里并不是直接返回消息，而是调用 `tracker.post`：

```

1      def post(message: GetMapOutputMessage): Unit = {
2          mapOutputRequests.offer(message)
3      }

```

向 `mapOutputRequests` 加入 `GetMapOutputMessage(shuffleId, context)` 消息。这里的 `mapOutputRequests` 是链式阻塞队列。

```

1      private val mapOutputRequests = new LinkedBlockingQueue[GetMapOutputMessage]

```

## MapOutputTrackerMaster.MessageLoop.run

`MessageLoop` 启一个线程不断的参数从 `mapOutputRequests` 读取数据：

```

1      private class MessageLoop extends Runnable {
2          override def run(): Unit = {
3              try {
4                  while (true) {
5                      try {
6                          val data = mapOutputRequests.take()
7                          if (data == PoisonPill) {
8                              mapOutputRequests.offer(PoisonPill)
9                              return
10                     }
11                     val context = data.context
12                     val shuffleId = data.shuffleId
13                     val hostPort = context.senderAddress.hostPort
14                     logDebug("Handling request to send map output locations for shuffle " + shuffleId +
15                             " to " + hostPort)
16                     // 若读到数据
17                     // 则序列化
18                     val mapOutputStatuses = getSerializedMapOutputStatuses(shuffleId)
19                     // 返回数据
20                     context.reply(mapOutputStatuses)
21                 } catch {
22                     case NonFatal(e) => logError(e.getMessage, e)
23                 }
24             }
25             } catch {
26                 case ie: InterruptedException => // exit
27             }
28         }
29     }

```

## MapOutputTracker.convertMapStatuses

我们回到 `mapOutputTracker.getMapSizesByExecutorId` 中返回的 `MapOutputTracker.convertMapStatuses`：

```

1 private def convertMapStatuses(
2     shuffleId: Int,
3     startPartition: Int,
4     endPartition: Int,
5     statuses: Array[MapStatus]): Seq[(BlockManagerId, Seq[(BlockId, Long)])] = {
6     assert(statuses != null)
7     val splitsByAddress = new HashMap[BlockManagerId, ArrayBuffer[(BlockId, Long)]]
8     for ((status, mapId) <- statuses.zipWithIndex) {
9         if (status == null) {
10             val errorMessage = s"Missing an output location for shuffle $shuffleId"
11             logError(errorMessage)
12             throw new MetadataFetchFailedException(shuffleId, startPartition, errorMessage)
13         } else {
14             for (part <- startPartition until endPartition) {
15                 // 返回的Seq中的结构是status.location, Seq[ShuffleBlockId,SizeForBlock]
16                 splitsByAddress.getOrElseUpdate(status.location, ArrayBuffer()) +=
17                     ((ShuffleBlockId(shuffleId, mapId, part), status.getSizeForBlock(part)))
18             }
19         }
20     }
21     // 对Seq根据status.location进行排序
22     splitsByAddress.toSeq
23 }

```

## 划分本地和远程Block

让我回到 `new ShuffleBlockFetcherIterator`

### storage.ShuffleBlockFetcherIterator.initialize

当我们实例化 `ShuffleBlockFetcherIterator` 时, 会调用 `initialize`:

```

1 private[this] def initialize(): Unit = {
2     context.addTaskCompletionListener(_ => cleanup())
3
4     // 划分本地和远程的blocks
5     val remoteRequests = splitLocalRemoteBlocks()
6     // 把远程请求随机的添加到队列中
7     fetchRequests += Utils.randomize(remoteRequests)
8     assert ((0 == reqsInFlight) == (0 == bytesInFlight),
9         "expected reqsInFlight = 0 but found reqsInFlight = " + reqsInFlight +
10         ", expected bytesInFlight = 0 but found bytesInFlight = " + bytesInFlight)
11
12     // 发送远程请求获取blocks
13     fetchUpToMaxBytes()
14
15     val numFetches = remoteRequests.size - fetchRequests.size
16     logInfo("Started " + numFetches + " remote fetches in" + Utils.getUsedTimeMs(startTime))
17
18     // 获取本地的Blocks
19     fetchLocalBlocks()
20     logDebug("Got local blocks in " + Utils.getUsedTimeMs(startTime))
21 }

```

### storage.ShuffleBlockFetcherIterator.splitLocalRemoteBlocks

```

1 private[this] def splitLocalRemoteBlocks(): ArrayBuffer[FetchRequest] = {
2     // 是的远程请求最大长度为 maxBytesInFlight / 5
3     // maxBytesInFlight: 为单次航班请求的最大字节数
4     // 航班: 一批请求
5     // 1/5 : 是为了提高请求批发度, 允许5个请求分别从5个节点获取数据
6     val targetRequestSize = math.max(maxBytesInFlight / 5, 1L)

```

```

7      logDebug("maxBytesInFlight: " + maxBytesInFlight + ", targetRequestSize: " + targetRequestSize)
8
9      // 缓存需要远程请求的FetchRequest对象
10     val remoteRequests = new ArrayBuffer[FetchRequest]
11
12     // 总共 blocks 的数量
13     var totalBlocks = 0
14     // 我们从上文可知blocksByAddress是根据status.location进行排序的
15     for ((address, blockInfos) <- blocksByAddress) {
16         totalBlocks += blockInfos.size
17         if (address.executorId == blockManager.blockManagerId.executorId) {
18             // 若 executorId 相同 与本 blockManagerId.executorId,
19             // 则从本地获取
20             localBlocks += blockInfos.filter(_._2 != 0).map(_._1)
21             numBlocksToFetch += localBlocks.size
22         } else {
23             // 否则 远程请求
24             // 得到迭代器
25             val iterator = blockInfos.iterator
26             // 当前累计块的大小
27             var curRequestSize = 0L
28             // 当前累加块
29             // 累加: 若向一个节点频繁的请求字节很少的Block,
30             // 那么会造成网络阻塞
31             var curBlocks = new ArrayBuffer[(BlockId, Long)]
32             // iterator 中的block 都是同一节点的
33             while (iterator.hasNext) {
34                 val (blockId, size) = iterator.next()
35                 if (size > 0) {
36                     curBlocks += ((blockId, size))
37                     remoteBlocks += blockId
38                     numBlocksToFetch += 1
39                     curRequestSize += size
40                 } else if (size < 0) {
41                     throw new BlockException(blockId, "Negative block size " + size)
42                 }
43                 if (curRequestSize >= targetRequestSize) {
44                     // 若累加到大于远程请求的尺寸
45                     // 往remoteRequests加入FetchRequest
46                     remoteRequests += new FetchRequest(address, curBlocks)
47                     curBlocks = new ArrayBuffer[(BlockId, Long)]
48                     logDebug(s"Creating fetch request of $curRequestSize at $address")
49                     curRequestSize = 0
50                 }
51             }
52             // 增加最后的请求
53             if (curBlocks.nonEmpty) {
54                 remoteRequests += new FetchRequest(address, curBlocks)
55             }
56         }
57     }
58     logInfo(s"Getting $numBlocksToFetch non-empty blocks out of $totalBlocks blocks")
59     remoteRequests
60 }

```

## 获取Block

### storage.ShuffleBlockFetcherIterator.fetchUpToMaxBytes

我们回到 `storage.ShuffleBlockFetcherIterator.initialize` 的 `fetchUpToMaxBytes()` 来深入讲解下如何获取远程的 `Block` :

```

1  private def fetchUpToMaxBytes(): Unit = {
2      // Send fetch requests up to maxBytesInFlight

```



```

3 // 单次航班请求数要小于最大航班请求数
4 // 单次航班字节数要小于最大航班字节数
5 while (fetchRequests.nonEmpty &&
6     (bytesInFlight == 0 ||
7     (reqsInFlight + 1 <= maxReqsInFlight &&
8     bytesInFlight + fetchRequests.front.size <= maxBytesInFlight))) {
9     sendRequest(fetchRequests.dequeue())
10 }
11 }

```

## storage.ShuffleBlockFetcherIterator.sendRequest

```

1 private[this] def sendRequest(req: FetchRequest) {
2     logDebug("Sending request for %d blocks (%s) from %s".format(
3         req.blocks.size, Utils.bytesToString(req.size), req.address.hostPort))
4     bytesInFlight += req.size
5     reqsInFlight += 1
6
7     // 可根据blockID查询block大小
8     val sizeMap = req.blocks.map { case (blockId, size) => (blockId.toString, size) }.toMap
9     val remainingBlocks = new HashSet[String]() ++= sizeMap.keys
10    val blockIds = req.blocks.map(_._1.toString)
11
12    val address = req.address
13    // 关于shuffleClient.fetchBlocks我们会在之后的博文讲解
14    shuffleClient.fetchBlocks(address.host, address.port, address.executorId, blockIds.toArray,
15        new BlockFetchingListener {
16            // 请求成功
17            override def onBlockFetchSuccess(blockId: String, buf: ManagedBuffer): Unit = {
18                ShuffleBlockFetcherIterator.this.synchronized {
19                    if (!isZombie) {
20                        buf.retain()
21                        remainingBlocks -= blockId
22                        results.put(new SuccessFetchResult(BlockId(blockId), address, sizeMap(blockId), buf,
23                            remainingBlocks.isEmpty))
24                        logDebug("remainingBlocks: " + remainingBlocks)
25                    }
26                }
27                logTrace("Got remote block " + blockId + " after " + Utils.getUsedTimeMs(startTime))
28            }
29
30            // 请求失败
31            override def onBlockFetchFailure(blockId: String, e: Throwable): Unit = {
32                logError(s"Failed to get block(s) from ${req.address.host}:${req.address.port}", e)
33                results.put(new FailureFetchResult(BlockId(blockId), address, e))
34            }
35        }
36    )
37 }

```

## storage.ShuffleBlockFetcherIterator.fetchLocalBlocks

我们再回过头来看获取本地blocks:

```

1 private[this] def fetchLocalBlocks() {
2     // 获取迭代器
3     val iter = localBlocks.iterator
4     while (iter.hasNext) {
5         val blockId = iter.next()
6         try {
7             // 遍历获取数据
8             // blockManager.getBlockData 会在后续博文讲解
9             val buf = blockManager.getBlockData(blockId)

```



```
10 shuffleMetrics.incLocalBlocksFetched(1)
11 shuffleMetrics.incLocalBytesRead(buf.size)
12 buf.retain()
13 results.put(new SuccessFetchResult(blockId, blockManager.blockManagerId, 0, buf, false))
14 } catch {
15   case e: Exception =>
16     logError(s"Error occurred while fetching local blocks", e)
17     results.put(new FailureFetchResult(blockId, blockManager.blockManagerId, e))
18     return
19   }
20 }
21 }
```