

深入理解Spark 2.1 Core（九）：迭代计算和Shuffle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/54981376>

在博文《深入理解Spark 2.1 Core（七）：任务执行的原理与源码分析》我们曾讲到过：

Task有两个子类，一个是非最后的Stage的Task，ShuffleMapTask；一个是最后的Stage的Task，ResultTask。它们都覆盖了Task的runTask方法。

我们来看一下ShuffleMapTask的runTask方法中的部分代码：

```
1    var writer: ShuffleWriter[Any, Any] = null
2    try {
3        //获取 shuffleManager
4        val manager = SparkEnv.get.shuffleManager
5        // 获取 writer
6        writer = manager.getWriter[Any, Any](dep.shuffleHandle, partitionId, context)
7        // 调用writer.write 开始计算RDD,
8        // 这部分 我们会在后续博文讲解
9        writer.write(rdd.iterator(partition, context).asInstanceOf[Iterator[_ <: Product2[Any, Any]]])
10       // 停止计算，并返回结果
11       writer.stop(success = true).get
12    }
```

这篇博文，我们就来深入这部分源码。

RDD迭代

调用栈如下：

- rdd.iterator
 - rdd.computeOrReadCheckpoint
 - rdd.MapPartitionsRDD.compute
 -
 - rdd.HadoopRDD.compute

rdd.RDD.iterator

我们先来看 `writer.write` 传入的参数：

```
1    rdd.iterator(partition, context).asInstanceOf[Iterator[_ <: Product2[Any, Any]]]
```

partition是该任务所在的分区，context为该任务的上下文。

rdd.iterator的方法如下：

```
1    final def iterator(split: Partition, context: TaskContext): Iterator[T] = {
2        // 此部分关于存储，会在后续讲解
3        if (storageLevel != StorageLevel.NONE) {
4            getOrCompute(split, context)
5        } else {
6            computeOrReadCheckpoint(split, context)
7        }
8    }
```

rdd.RDD.computeOrReadCheckpoint

我们来看下上述的computeOrReadCheckpoint方法：

```

1 private[spark] def computeOrReadCheckpoint(split: Partition, context: TaskContext): Iterator[T] =
2 {
3   // 若Checkpointed 获取结果
4   if (isCheckpointedAndMaterialized) {
5     firstParent[T].iterator(split, context)
6   } else {
7     // 否则计算
8     compute(split, context)
9   }
10 }

```

rdd.MapPartitionsRDD.compute

这里对 `compute` 实现的RDD是 `MapPartitionsRDD`：

```

1 override def compute(split: Partition, context: TaskContext): Iterator[U] =
2   f(context, split.index, firstParent[T].iterator(split, context))

```

我们可以看到，这里还是会调用 `firstParent[T].iterator`，这样父RDD继续调用 `MapPartitionsRDD.compute`，这样一层层的向上调用，直到最初的RDD。

rdd.HadoopRDD.compute

若是从HDFS读取生成的最初的RDD，则经过层层调用，会调用到 `HadoopRDD.compute`。下面我们来看下该方法：

```

1 override def compute(theSplit: Partition, context: TaskContext): InterruptibleIterator[(K, V)] = {
2   // iter 是NextIterator匿名类的一个对象
3   val iter = new NextIterator[(K, V)] {
4
5     //***** 以下为NextIterator匿名类内容 *****
6
7     private val split = theSplit.asInstanceOf[HadoopPartition]
8     logInfo("Input split: " + split.inputSplit)
9     // hadoop的配置
10    private val jobConf = getJobConf()
11
12    // 用于计算字节读取
13    private val inputMetrics = context.taskMetrics().inputMetrics
14    // 之前写入的值
15    private val existingBytesRead = inputMetrics.bytesRead
16
17    // 设置 文件名的 线程本地值
18    split.inputSplit.value match {
19      case fs: FileSplit => InputFileNameHolder.setInputFileName(fs.getPath.toString)
20      case _ => InputFileNameHolder.unsetInputFileName()
21    }
22
23
24    // 用于返回该线程从文件读取的字节数
25    // 需要在RecordReader创建前创建
26    // 因为RecordReader的构造函数可能需要读取一些字节
27    private val getBytesReadCallback: Option[() => Long] = split.inputSplit.value match {
28      case _: FileSplit | _: CombineFileSplit =>
29        SparkHadoopUtil.get.getFSBytesReadOnThreadCallback()
30      case _ => None
31    }
32
33    //对于 Hadoop 2.5以上的版本，我们从线程本地HDFS统计中得到输入的字节数。

```

```
34 // 如果我做一个合并操作的话,
35 // 我们需要在同一个任务且同一个线程里计算多个分区。
36 // 在这种情况下, 我们需要去避免覆盖之前分区中已经被写入的值
37 private def updateBytesRead(): Unit = {
38     getBytesReadCallback.foreach { getBytesRead =>
39         inputMetrics.setBytesRead(existingBytesRead + getBytesRead())
40     }
41 }
42
43 private var reader: RecordReader[K, V] = null
44 // 即 TextInputFormat
45 private val inputFormat = getInputFormat(jobConf)
46 // 添加hadoop相关任务配置
47 HadoopRDD.addLocalConfiguration(
48     new SimpleDateFormat("yyyyMMddHHmmss", Locale.US).format(createTime),
49     context.stageId, theSplit.index, context.attemptNumber, jobConf)
50 // 创建RecordReader
51 reader =
52     try {
53         inputFormat.getRecordReader(split.inputSplit.value, jobConf, Reporter.NULL)
54     } catch {
55         case e: IOException if ignoreCorruptFiles =>
56             logWarning(s"Skipped the rest content in the corrupted file: ${split.inputSplit}", e)
57             finished = true
58             null
59     }
60 // 注册任务完成回调来关闭输入流
61 context.addTaskCompletionListener{ context => closeIfNeeded() }
62 // key: LongWritable
63 private val key: K = if (reader == null) null.asInstanceOf[K] else reader.createKey()
64 // v: Text
65 private val value: V = if (reader == null) null.asInstanceOf[V] else reader.createValue()
66
67 // 对reader.next的代理
68 override def getNext(): (K, V) = {
69     try {
70         finished = !reader.next(key, value)
71     } catch {
72         case e: IOException if ignoreCorruptFiles =>
73             logWarning(s"Skipped the rest content in the corrupted file: ${split.inputSplit}", e)
74             finished = true
75     }
76     if (!finished) {
77         inputMetrics.incRecordsRead(1)
78     }
79     if (inputMetrics.recordsRead % SparkHadoopUtil.UPDATE_INPUT_METRICS_INTERVAL_RECORDS == 0) {
80         // 更新inputMetrics的BytesRead
81         updateBytesRead()
82     }
83     (key, value)
84 }
85
86 // 关闭
87 override def close() {
88     if (reader != null) {
89         InputFileNameHolder.unsetInputFileName()
90         try {
91             reader.close()
92         } catch {
93             case e: Exception =>
94                 if (!ShutdownHookManager.inShutdown()) {
95                     logWarning("Exception in RecordReader.close()", e)
96                 }
97         } finally {
```

```

98         reader = null
99     }
100     if (getBytesReadCallback.isDefined) {
101         updateBytesRead()
102     } else if (split.inputSplit.value.isInstanceOf[FileSplit] ||
103         split.inputSplit.value.isInstanceOf[CombineFileSplit]) {
104         try {
105             inputMetrics.incBytesRead(split.inputSplit.value.getLength)
106         } catch {
107             case e: java.io.IOException =>
108                 logWarning("Unable to get input size to set InputMetrics for task", e)
109         }
110     }
111 }
112 }
113 }
114 new InterruptibleIterator[(K, V)](context, iter)
115 }

```

InterruptibleIterator传入参数iter，可以看成是NextIterator类的代理：

```

1 class InterruptibleIterator[+T](val context: TaskContext, val delegate: Iterator[T])
2 extends Iterator[T] {
3
4     def hasNext: Boolean = {
5         if (context.isInterrupted) {
6             throw new TaskKilledException
7         } else {
8             delegate.hasNext
9         }
10    }
11
12    def next(): T = delegate.next()
13 }

```

迭代返回

当 `rdd.HadoopRDD.compute` 运算完毕后，生成的初始的RDD计算结果。退回到 `rdd.HadoopRDD.compute` 便可以调用函数 `f`：

```
1 f(context, split.index, firstParent[T].iterator(split, context))
```

`f` 计算出第二个的RDD计算结果，以此类推，一层层的返回。最终回到 `writer.write`：

```
1 writer.write(rdd.iterator(partition, context).asInstanceOf[Iterator[_ <: Product2[Any, Any]])
```

`ShuffleWriter` 是一个抽象类，它有子类 `SortShuffleWriter`。`SortShuffleWriter.write`：

```

1 override def write(records: Iterator[Product2[K, V]]): Unit = {
2     // 创建ExecutorSorter,
3     // 用于Shuffle Map Task 输出结果排序
4     sorter = if (dep.mapSideCombine) {
5         // 当计算结果需要combine,
6         // 则外部排序进行聚合
7         require(dep.aggregator.isDefined, "Map-side combine without Aggregator specified!")
8         new ExternalSorter[K, V, C](
9             context, dep.aggregator, Some(dep.partitionner), dep.keyOrdering, dep.serializer)
10    } else {
11        // 否则，外部排序不进行聚合
12        new ExternalSorter[K, V, C](
13            context, aggregator = None, Some(dep.partitionner), ordering = None, dep.serializer)
14    }
15    // 根据排序方式，对数据进行排序并写入内存缓冲区。

```

```

16 // 若排序中计算结果超出的阈值,
17 // 则将其溢写到磁盘数据文件
18 sorter.insertAll(records)
19
20 // 通过shuffle编号和map编号来获取该数据文件
21 val output = shuffleBlockResolver.getDataFile(dep.shuffleId, mapId)
22 val tmp = Utils.tempFileWith(output)
23 try {
24 // 通过shuffle编号和map编号来获取 ShuffleBlock 编号
25 val blockId = ShuffleBlockId(dep.shuffleId, mapId, IndexShuffleBlockResolver.NOOP_REDUCE_ID)
26 // 在外部排序中,
27 // 有部分结果可能在内存中
28 // 另外部分结果在一个或多个文件中
29 // 需要将它们merge成一个大文件
30 val partitionLengths = sorter.writePartitionedFile(blockId, tmp)
31 // 创建索引文件
32 // 将每个partition在数据文件中的起始与结束位置写入到索引文件
33 shuffleBlockResolver.writeIndexFileAndCommit(dep.shuffleId, mapId, partitionLengths, tmp)
34 // 将元数据写入mapStatus
35 // 后续任务通过该mapStatus得到处理结果信息
36 mapStatus = MapStatus(blockManager.shuffleServerId, partitionLengths)
37 } finally {
38   if (tmp.exists() && !tmp.delete()) {
39     logError(s"Error while deleting temp file ${tmp.getAbsolutePath}")
40   }
41 }
42 }

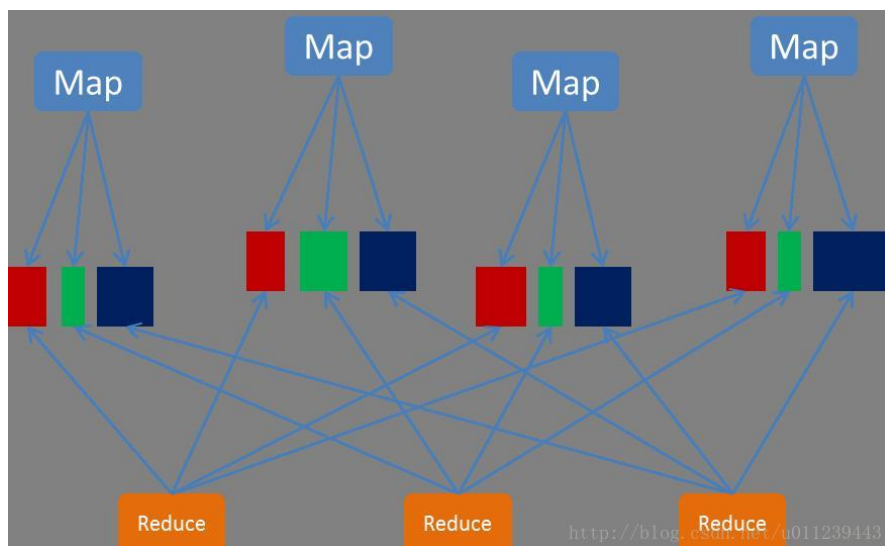
```

Shuffle原理概要

MapReduce Shuffle原理 与 Spark Shuffle 原理

MapReduce Shuffle原理 与 Spark Shuffle原理可以参阅腾讯的一篇博文 [《MapReduce Shuffle原理 与 Spark Shuffle原理》](#)

在这里我们重新讲解下早起Spark Shuffle的过程：



- map任务会给每个reduce任务分配一个bucket。假设有M个map任务，每个map任务有N个reduce任务，则map阶段一共会创建M×N个bucket。
- map任务会将产生的中间结果按照partition写入到不同的bucket中
- reduce任务从本地或者远端的map任务所在的BlockManager获取相应的bucket作为输入

MapReduce Shuffle 与 Spark Shuffle缺陷

MapReduce Shuffle缺陷

- map任务产生的结果排序后会写入磁盘，reduce获取map任务产生的结果会在磁盘上merge sort，产生很多磁盘I/O
- 当数量很小，但是map和reduce任务很多时，会产生很多网络I/O

Spark Shuffle缺陷

- map任务产生的结果先写入内存，当一个节点输出的结果集很大是，容易内存紧张
- map任务数量与reduce数量大了，bucket数量容易变得非常大，这就带来了两个问题：
 - 每打开一个文件（bucket为一个文件）都会暂用一定内存，容易内存紧张
 - 若bucket本身很小，而对于系统来说遍历多个文件是随机读取，那么磁盘I/O性能会变得非常差

Spark shuffle 的优化

- 把相同的partition的bucket放在一个文件中
- 使用缓存及聚合算法对map任务的输出结果进行聚合
- 使用缓存及聚合算法对reduce从map拉取的输出结果进行聚合
- 缓存超出阈值时，将数据写入磁盘
- reduce任务将同一BlockManager地址的Block累计，减少网络请求