



黑客派

# MapReduce 概述

作者: [ibywind](#)

原文链接: <https://hacpai.com/article/1580301635611>

来源网站: [黑客派](#)

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

## 一、MapReduce 概述

`<script async src="https://pagead2.googlesyndication.com/pagead/js/adsbygoogle.js"></script>`

`<!-- 黑客派PC帖子内嵌-展示 -->`

`<ins class="adsbygoogle" style="display:block" data-ad-client="ca-pub-5357405790190342" data-ad-slot="8316640078" data-ad-format="auto" data-full-width-responsive="true"></in>`

`<script>`

`(adsbygoogle = window.adsbygoogle || []).push({};`

`</script>`

Hadoop MapReduce 是一个分布式计算框架，用于编写批处理应用程序。编写好的程序可以提到 Hadoop 集群上用于并行处理大规模的数据集。


MapReduce 作业通过将输入的数据集拆分为独立的块，这些块由 `map` 以并的方式处理，框架对 `map` 的输出进行排序，然后输入到 `reduce` 中。MapReduce 框架专门用于 `<code>&lt;key, value&gt;` 键值对处理，它将作业的输视为一组 `<code>&lt;key, value&gt;` 对，并生成一组 `<code>&lt;key, value&gt;` 对作为输出。输出和输出的 `key` 和 `value` 都必须实现 [Writable](https://link.hacpai.com/forward?goto=http%3A%2F%2Fhadoop.apache.org%2Fdocs%2Fstable%2Fapi%2Forg%2Fapache%2Fhadoop%2Fio%2FWritable.html) 接口。

```
(input) &lt;k1, v1&gt; -&gt; map -&gt; &lt;k2, v2&gt; &gt; combine -&gt; &lt;k2, v2&gt; -&gt; reduce -&gt; &lt;k3, v3&gt; (output)
```

`</code>`

## 二、MapReduce 编程模型简述

这里以词频统计为例进行说明，MapReduce 处理的流程如下：

<https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2FmapreduceProcess.png>  `</a></p>`

`<ol>`

`<li><strong>input</strong> : 读取文本文件; </li>`

`<li><strong>splitting</strong> : 将文件按照行进行拆分，此时得到的 K1 行，V1 表示对应行的文本内容; </li>`

`<li><strong>mapping</strong> : 并行将每一行按照空格进行拆分，拆分得到的 List(K2, 2)，其中 K2 代表每一个单词，由于是做词频统计，所以 V2 的值为 1，代表出现 1 次; </li>`


`<li><strong>shuffling</strong> : 由于 Mapping 操作可能是在不同的机器上行处理的，所以需要通过 shuffling 将相同 key 值的数据分发到一个节点上去合并，这样才能统计出最终的结果，此时得到 K2 为每一个单词，List(V2) 为可迭代集合，V2 就是 Mapping 中的 V2; </li>`

`<li><strong>Reducing</strong> : 这里的案例是统计单词出现的总次数，所以 Reducing 对 List(V2) 进行归约求和操作，最终输出。 </li>`

`</ol>`

MapReduce 编程模型中 `splitting` 和 `shuffling` 操作都是由架实现的，需要我们自己编程实现的只有 `mapping` 和 `reducing`，这也就是 MapReduce 这个称呼的来源。

## 三、combiner & partitioner

<https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2FDetailed-Hadoop-MapReduce-Data-Flow-14.png>  `</a></p>`

`<h3 id="-"></h3>`

<p></p>

<p>[<a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fnotes%2FHadoop-MapReduce.md%2331-inputformat--recordreaders%293.1" target="\_blank" rel="nofollow ugc">https://github.com/heibaiying/BigData-Notes/blob/master/notes/Hadoop-MapReduce.md#31-inputformat--recordreaders)3.1</a> InputFormat & RecordReaders</p>

<p><code>InputFormat</code> 将输出文件拆分为多个 <code>InputSplit</code>, 并由 <code>RecordReaders</code> 将 <code>InputSplit</code> 转换为标准的 &lt;key, value&gt; 键对, 作为 map 的输出。这一步的意义在于只有先进行逻辑拆分并转为标准的键值对格式后, 才能为一个 <code>map</code> 提供输入, 以便进行并行处理。</p>

<h3 id="3-2-Combiner">3.2 Combiner</h3>

<p><code>combiner</code> 是 <code>map</code> 运算后的可选操作, 它实际上是一个本化的 <code>reduce</code> 操作, 它主要是在 <code>map</code> 计算出中间文件后做一个简单的合并重复 <code>key</code> 值的操作。这里以词频统计为例: </p>

<script async src="https://pagead2.googlesyndication.com/pagead/js/adsbygoogle.js"></script>

<!-- 黑客派PC帖子内嵌-展示 -->

<ins class="adsbygoogle" style="display:block" data-ad-client="ca-pub-5357405790190342" data-ad-slot="8316640078" data-ad-format="auto" data-full-width-responsive="true"></ins>

<script>

(adsbygoogle = window.adsbygoogle || []).push({});

</script>

<p><code>map</code> 在遇到一个 hadoop 的单词时就会记录为 1, 但是这篇文章里 hadoop 能会出现 n 多次, 那么 <code>map</code> 输出文件冗余就会很多, 因此在 <code>reduce</code> 计算前对相同的 key 做一个合并操作, 那么需要传输的数据量就会减少, 传输效率就可以得到提高。</p>

<p>但并非所有场景都适合使用 <code>combiner</code>, 使用它的原则是 <code>combiner</code> 的输出不会影响到 <code>reduce</code> 计算的最终输入, 例如: 求总数, 最大值, 最小时都可以使用 <code>combiner</code>, 但是做平均值计算则不能使用 <code>combiner</code>。</p>

<p>不使用 combiner 的情况: </p>

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fmapreduce-without-combiners.png" target="\_blank" rel="nofollow ugc"></a></p>

<p>使用 combiner 的情况: </p>

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fmapreduce-with-combiners.png" target="\_blank" rel="nofollow ugc"></a></p>

<p>可以看到使用 combiner 的时候, 需要传输到 reducer 中的数据由 12keys, 降低到 10keys。低的幅度取决于你 keys 的重复率, 下文词频统计案例会演示用 combiner 降低数百倍的传输量。</p>

<h3 id="3-3-Partitioner">3.3 Partitioner</h3>

<p><code>partitioner</code> 可以理解成分类器, 将 <code>map</code> 的输出按照 key 的不同分别分给对应的 <code>reducer</code>, 支持自定义实现, 下文案例会给出演示。</p>

<h2 id="四-MapReduce-词频统计案例">四、MapReduce 词频统计案例</h2>

<h3 id="4-1-项目简介">4.1 项目简介</h3>

<p>这里给出一个经典的词频统计的案例: 统计如下样本数据中每个单词出现的次数。</p>

```
Spark HBase  
Hive Flink Storm Hadoop HBase Spark
```

```
Flink
HBase Storm
HBase Hadoop Hive Flink
HBase Flink Hive Storm
Hive Flink Hadoop
HBase Hive
Hadoop Spark HBase Storm
HBase Hadoop Hive Flink
HBase Flink Hive Storm
Hive Flink Hadoop
HBase Hive</pre>
```

<p>为方便大家开发，我在项目源码中放置了一个工具类 `WordCountDataUtils` 用于模拟产生词频统计的样本，生成的文件支持输出到本地或者直接写到 HDFS 上。</p>

### 4.2 项目依赖</h3>

<p>想要进行 MapReduce 编程，需要导入 `hadoop-client` 依赖：</p>

```
<pre>&lt;<span>dependency</span>&gt;
    &lt;<span>groupId</span>&gt;org.apache.hadoopgroupId&gt;
    &lt;<span>artifactId</span>&gt;hadoop-clientartifactId&gt;
    &lt;<span>version</span>&gt;${hadoop.version}version&gt;
dependency&gt;</pre>
```

### 4.3 WordCountMapper</h3>

<p>将每行数据按照指定分隔符进行拆分。这里需要注意在 MapReduce 中必须使用 Hadoop 定义类型，因为 Hadoop 预定义的类型都是可序列化，可比较的，所有类型均实现了 `WritableComparable` 接口。</p>

```
<pre><span>public</span> <span>class</span> <span>WordCountMapper</span> <span>extends</span> <span>Mapper<LongWritable, Text, Text, IntWritable>&gt;</span> {
```

```
</pre>
```

```
<script async src="https://pagead2.googlesyndication.com/pagead/js/adsbygoogle.js"></script>
```

```
<!-- 黑客派PC帖子内嵌-展示 -->
```

```
<ins class="adsbygoogle" style="display:block" data-ad-client="ca-pub-5357405790190342" data-ad-slot="8316640078" data-ad-format="auto" data-full-width-responsive="true"></ins>
```

```
<script>
    (adsbygoogle = window.adsbygoogle || []).push({});
</script>
```

<p><code>WordCountMapper</code> 对应下图的 Mapping 操作：</p>

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fhadoop-code-mapping.png" target="\_blank" rel="nofollow ugc"></a></p>

<p><code>WordCountMapper</code> 继承自 `Mapper` 类，这是一个泛型类定义如下：</p>

```
<pre><span>WordCountMapper</span> extends <span>Mapper<LongWritable, Text, Text, IntWritable>&gt;
```

```
<span>public</span> <span>class</span> <span>Mapper</span> {  
<span>.....</span>  
}</pre>
```

<ul>

<li><strong>KEYIN</strong> : <code>mapping</code> 输入 key 的类型, 即每行的偏移量 (行第一个字符在整个文本中的位置), <code>Long</code> 类型, 对应 Hadoop 中的 <code>LongWritable</code> 类型; </li>

<li><strong>VALUEIN</strong> : <code>mapping</code> 输入 value 的类型, 即每行数据 <code>String</code> 类型, 对应 Hadoop 中 <code>Text</code> 类型; </li>

<li><strong>KEYOUT</strong> : <code>mapping</code> 输出的 key 的类型, 即每个单词 <code>String</code> 类型, 对应 Hadoop 中 <code>Text</code> 类型; </li>

<li><strong>VALUEOUT</strong> : <code>mapping</code> 输出 value 的类型, 即每个单词出现的次数; 这里用 <code>int</code> 类型, 对应 <code>IntWritable</code> 类型。 </li>  
</ul>

### 4.4 WordCountReducer</h3>

<p>在 Reduce 中进行单词出现次数的统计: </p>

```
<pre> <span>public</span> <span>class</span> <span>WordCountReducer</span> <span>extends</span> <span>Reducer</span>&lt;<span>Text</span>, <span>IntWritable</span>, <span>Text</span>, <span>IntWritable</span>&gt;</pre> {
```

```
</pre>
```

<p>如下图, <code>shuffling</code> 的输出是 reduce 的输入。这里的 key 是每个单词, values 是一个可迭代的数据类型, 类似 <code>(1,1,1,...)</code>。 </p>

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fhadoop-code-reducer.png" target="\_blank" rel="nofollow ugc"></a></p>

### 4.4 WordCountApp</h3>

<p>组装 MapReduce 作业, 并提交到服务器运行, 代码如下: </p>

```
<pre> <span>/**</span></pre>
```

```
<span> * 组装作业 并提交到集群运行</span>
```

```
<span>*/</span></pre>
```

```
<span>public</span> <span>class</span> <span>WordCountApp</span> {
```

```
</pre>
```

```
<script async src="https://pagead2.googlesyndication.com/pagead/js/adsbygoogle.js"></script>
```

```
<!-- 黑客派PC帖子内嵌-展示 -->
```

```
<ins class="adsbygoogle" style="display:block" data-ad-client="ca-pub-5357405790190342" data-ad-slot="8316640078" data-ad-format="auto" data-full-width-responsive="true"></ins>
```

```
<script>
```

```
  (adsbygoogle = window.adsbygoogle || []).push({});
```

```
</script>
```



<p>需要注意的是：如果不设置 <code>Mapper</code> 操作的输出类型，则程序默认它和 <code>Reducer</code> 操作输出的类型相同。</p>

### 

<p>在实际开发中，可以在本机配置 hadoop 开发环境，直接在 IDE 中启动进行测试。这里主要介绍一下打包提交到服务器运行。由于本项目没有使用除 Hadoop 外的第三方依赖，直接打包即可：</p>

```
<pre><span><span>#</span> mvn clean package</span></pre>
```

<p>使用以下命令提交作业：</p>

```
<pre>hadoop jar /usr/appjar/hadoop-word-count-1.0.jar \
com.heibaiying.WordCountApp \
/wordcount/input.txt /wordcount/output/WordCountApp</pre>
```

<p>作业完成后查看 HDFS 上生成目录：</p>

```
<pre><span><span>#</span> 查看目录</span>
hadoop fs -ls /wordcount/output/WordCountApp
```

<span><span>#</span> 查看统计结果</span>

```
hadoop fs -cat /wordcount/output/WordCountApp/part-r-00000</pre>
```

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fhadoop-wordcountapp.png" target="\_blank" rel="nofollow ugc"></a></p>

## 

### 

<p>想要使用 <code>combiner</code> 功能只要在组装作业时，添加下面一行代码即可：</p>

```
<pre><span><span>//</span> 设置 Combiner</span>
job.</span>setCombinerClass(<span>WordCountReducer</span><span>.</span></pre>
```

### 

<p>加入 <code>combiner</code> 后统计结果是不会有变化的，但是可以从打印的日志看出 <code>combiner</code> 的效果：</p>

<p>没有加入 <code>combiner</code> 的打印日志：</p>

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fhadoop-no-combiner.png" target="blank" rel="nofollow ugc"></a></p>

<p>加入 <code>combiner</code> 后的打印日志如下：</p>

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fhadoop-combiner.png" target="blank" rel="nofollow ugc"></a></p>

<p>这里我们只有一个输入文件并且小于 128M，所以只有一个 Map 进行处理。可以看到经过 combiner 后，records 由 <code>3519</code> 降低为 <code>6</code> (样本中单词种类就只有 6 种，在这个用例中 combiner 就能极大地降低需要传输的数据量。</p>

## 

### 

<p>这里假设有个需求：将不同单词的统计结果输出到不同文件。这种需求实际上比较常见，比如统计产品的销量时，需要将结果按照产品种类进行拆分。要实现这个功能，就需要用到自定义 <code>Partitioner</code>。</p>

<p>这里先介绍下 MapReduce 默认的分类规则：在构建 job 时候，如果不指定，默认使用的是 <

ode>HashPartitioner</code>：对 key 值进行哈希散列并对 <code>numReduceTasks</code> 取余。其实现如下：</p></div>

```
<pre> <span>public</span> <span>class</span> <span>HashPartitioner</span> <span>extends</span> <span>Partitioner<&lt;span>K</span>, <span>V</span>&gt;</span> {
```

```
<span>public</span> <span>int</span> <span>getPartition</span>(<span>K</span> <span>key</span>, <span>V</span> <span>value</span>,
```

```
<span>int</span> <span>numReduceTasks</span>) {
```

```
<span>return</span> (key.<span>.hashCode() <span>not found render function for node [type=NodeHTMLEntity, Tokens=&]not found render function for node [type=NodeHTMLEntity, Tokens=&]</span> <span>Integer</span><span><span>/span>MAX_VALUE</span>) <span>%</span> numReduceTasks;
```

```
}
```

```
</pre>
```

### ``` <script async src="https://pagead2.googlesyndication.com/pagead/js/adsbygoogle.js"></script> ``` ``` <!-- 黑客派PC帖子内嵌-展示 --> ``` ``` <ins class="adsbygoogle" style="display:block" data-ad-client="ca-pub-5357405790190342" data-ad-slot="8316640078" data-ad-format="auto" data-full-width-responsive="true"></ins> ``` ``` <script> ``` ``` (adsbygoogle = window.adsbygoogle || []).push({}); ``` ``` </script> ``` <p>这里我们继承 <code>Partitioner</code> 自定义分类规则，这里按照单词进行分类：</p></div> ``` <pre> <span>public</span> <span>class</span> <span>CustomPartitioner</span> <span>extends</span> <span>Partitioner<&lt;span>Text</span>, <span>IntWritable</span>&gt;</span> { ``` ``` </pre> ``` <p>在构建 <code>job</code> 时候指定使用我们自己的分类规则，并设置 <code>reduce</code> 的个数：</p></div> ``` <pre> <span><span>//</span> 设置自定义分区规则</span> ``` ``` job.<span>.setPartitionerClass(<span>CustomPartitioner</span><span>.class</span>); ``` ``` <span><span>//</span> 设置 reduce 个数</span> ``` ``` job.<span>.setNumReduceTasks(<span>WordCountDataUtils.<span><span>WORD_LIST</span><span>.size()</span>);</span> ``` <p>执行结果如下，分别生成 6 个文件，每个文件中为对应单词的统计结果：</p></div> <p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fgithub.com%2Fheibaiying%2FBigData-Notes%2Fblob%2Fmaster%2Fpictures%2Fhadoop-wordcountcombinerpartition.png" target="\_blank" rel="nofollow ugc"></a></p></div> 原文链接：MapReduce 概述