

Hadoop 词频统计实验

一、实验准备：环境与文件创建

1. 环境确认
2. 创建词频文件

二、核心命令操作（对应 Q1-Q3）

1. Q1: 启动 Hadoop
2. Q2: 上传文件到 HDFS 的 input 文件夹
3. Q3: 检查并删除 HDFS 中的 output 目录

三、编写 MapReduce 程序（对应 Q4-Q8，基于 Eclipse）

1. 环境配置（无插件方案，通用）
2. 编写核心代码（WordCount.java）
3. 代码逻辑解释（对应 Q4-Q8）

四、编译打包与运行程序

1. 打包 JAR 文件
2. 运行 MapReduce 程序
3. 查看运行结果

五、实验收尾

一、Hadoop 词频统计实验

1.1 一、实验准备：环境与文件创建

1.1.1 1. 环境确认

确保已在 **Ubuntu 系统** 中安装 Hadoop，并配置好环境变量（可通过在终端输入 `hadoop version` 验证，能显示版本则配置正常）；Windows 系统主要用于 Eclipse 开发（若 Ubuntu 装了 Eclipse 也可直接用）。

1.1.2 2. 创建词频文件

在 Ubuntu 终端中，通过以下命令创建 `wordfile1.txt` 和 `wordfile2.txt`（文件内容需固定，用于后续词频统计）：

```
1  #创建并编辑wordfile1.txt，输入指定内容
2  vim wordfile1.txt
3  #按i进入编辑模式，输入：
4  Hello Hadoop
5  Hello Hdfs
6  Hello MapReduce
7  #按Esc，输入:wq保存退出
8  #同理创建wordfile2.txt
9  vim wordfile2.txt
10 #输入内容：
11 Hadoop is good
12 Hdfs is distributed
13 MapReduce is efficient
```




1.2 二、核心命令操作（对应 Q1-Q3）

所有命令需在 Ubuntu 终端中执行，且确保当前用户有 Hadoop 操作权限（建议切换到 Hadoop 安装用户，如 `su hadoop`）。

1.2.1 1. Q1: 启动 Hadoop

Hadoop 启动需先初始化（首次启动时），再启动集群，命令如下：

```
1 #启动Hadoop集群（包含NameNode、DataNode、ResourceManager等）
2 start-dfs.sh # 启动分布式文件系统（HDFS）
3 start-yarn.sh # 启动资源调度系统（YARN）
4 #验证是否启动成功：输入jps，能看到NameNode、DataNode、ResourceManager、
  NodeManager则正常
5 jps
```



1.2.2 2. Q2: 上传文件到 HDFS 的 input 文件夹

HDFS 中默认没有 `input` 文件夹，需先创建，再上传本地文件：


```
1 #1. 在HDFS根目录创建input文件夹
2 hdfs dfs -mkdir /input
3 #2. 上传本地的wordfile1.txt和wordfile2.txt到HDFS的/input目录
4 （注意：需替换命令中的“本地文件路径”为实际路径，如\~/wordfile1.txt，\~代表当前用
  户家目录）
5 hdfs dfs -put \~/wordfile1.txt /input
6 hdfs dfs -put \~/wordfile2.txt /input
7 #验证上传：查看HDFS input文件夹下的文件
8 hdfs dfs -ls /input
```



1.2.3 3. Q3: 检查并删除 HDFS 中的 output 目录

MapReduce 运行时要求输出目录（如 output）**不存在**，否则会报错，需先检查并删除：

```
1 #1. 检查HDFS根目录是否存在output目录
2 hdfs dfs -test -e /output && echo "output存在" || echo "output不存在"
3 #2. 若存在，删除output目录（-r表示递归删除，包含子目录和文件）
4 hdfs dfs -rm -r /output
5 #验证删除：再次查看是否存在
6 hdfs dfs -ls /
```



1.3 三、编写 MapReduce 程序（对应 Q4-Q8，基于 Eclipse）

1.3.1 1. 环境配置（无插件方案，通用）

- ① 在 Eclipse 中创建 **Java Project**（如命名为 WordCountProject）。
- ② 添加 Hadoop 依赖 JAR 包：

- 找到 Ubuntu 中 Hadoop 的 `share/hadoop` 目录 (如 `/usr/local/hadoop/share/hadoop`)，将以下目录下的所有 JAR 包导入项目：
 - `common`、`common/lib`、`hdfs`、`hdfs/lib`、`mapreduce`、`mapreduce/lib`、`yarn`、`yarn/lib`。
- 导入方式：右键项目→Build Path→Add External Archives→选中上述 JAR 包。

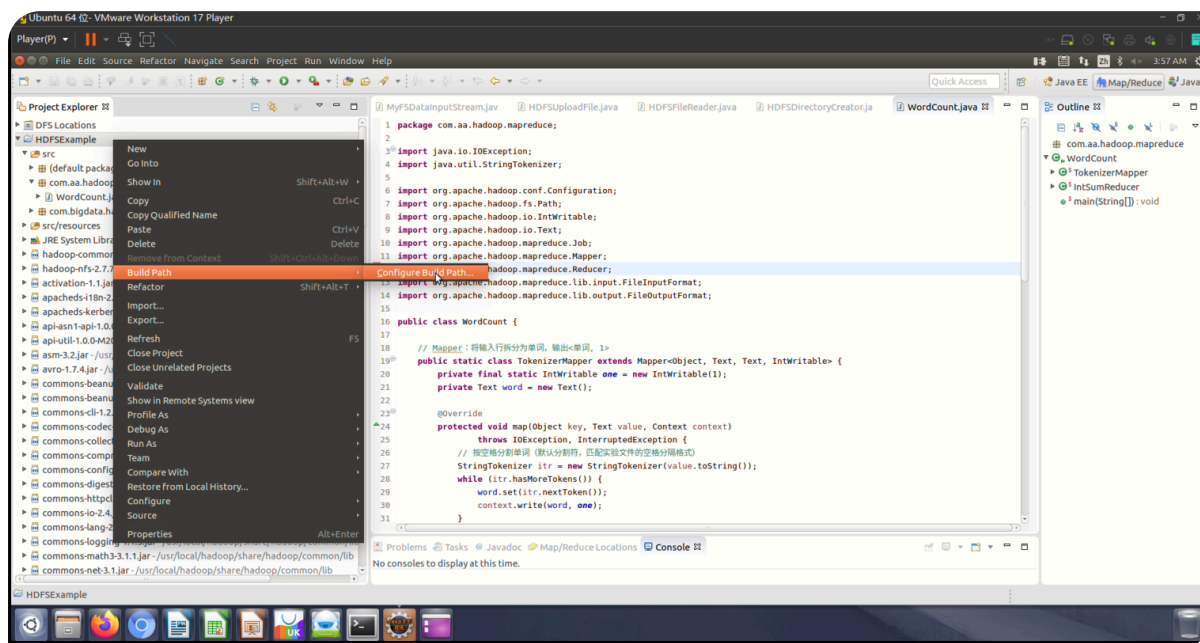


图 1.3.1-1 导入方式-1

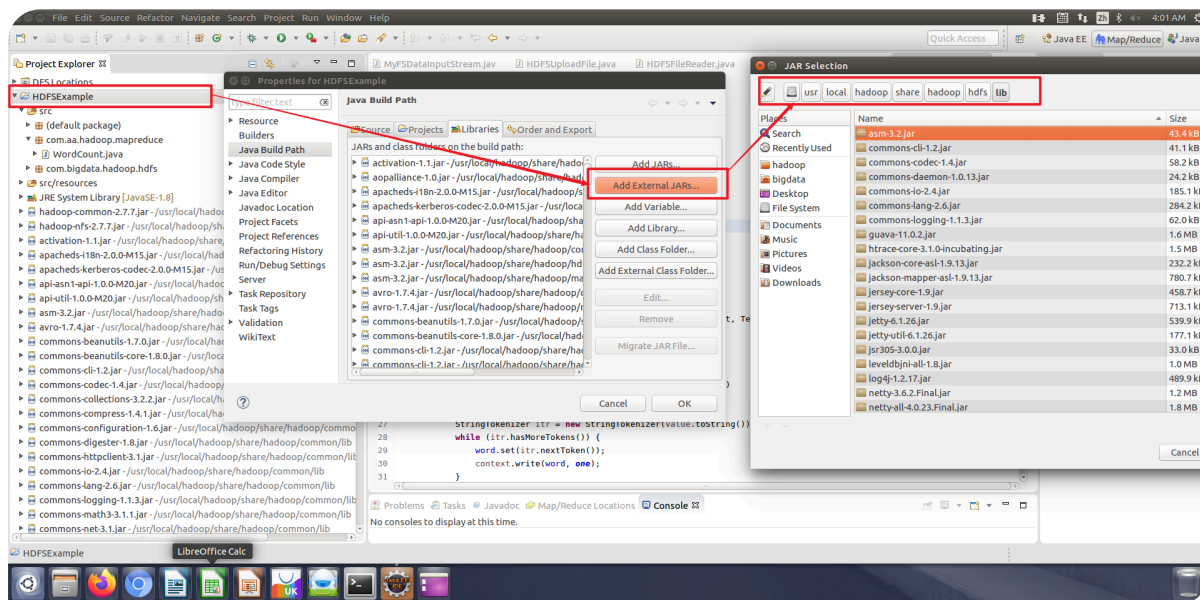


图 1.3.1-2 导入方式-2

1.3.2 2. 编写核心代码 (WordCount.java)

在项目的 `src` 目录下创建 `com.aa.mapreduce.wordcount` 包，新建 `WordCount.java` 文件，代码如下 (已包含 Map、Reduce、Main 逻辑)：


```

1  package com.aa.hadoop.mapreduce;
2
3  import java.io.IOException;
4  import java.util.StringTokenizer;
5
6  import org.apache.hadoop.conf.Configuration;
7  import org.apache.hadoop.fs.Path;
8  import org.apache.hadoop.io.IntWritable;
9  import org.apache.hadoop.io.Text;
10 import org.apache.hadoop.mapreduce.Job;
11 import org.apache.hadoop.mapreduce.Mapper;
12 import org.apache.hadoop.mapreduce.Reducer;
13 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
14 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
15
16 public class WordCount {
17
18     // Mapper: 将输入行拆分为单词, 输出<单词, 1>
19     public static class TokenizerMapper extends Mapper<Object, Text,
20 Text, IntWritable> {
21         private final static IntWritable one = new IntWritable(1);
22         private Text word = new Text();
23
24         @Override
25         protected void map(Object key, Text value, Context context)
26             throws IOException, InterruptedException {
27             // 按空格分割单词 (默认分割符, 匹配实验文件的空格分隔格式)
28             StringTokenizer itr = new
29 StringTokenizer(value.toString());
30             while (itr.hasMoreTokens()) {
31                 word.set(itr.nextToken());
32                 context.write(word, one);
33             }
34         }
35     }
36
37     // Reducer: 汇总同一单词的计数, 输出<单词, 总次数>
38     public static class IntSumReducer extends Reducer<Text,
39 IntWritable, Text, IntWritable> {
40         private IntWritable result = new IntWritable();
41
42         @Override
43         protected void reduce(Text key, Iterable<IntWritable> values,
44 Context context)
45             throws IOException, InterruptedException {
46             int sum = 0;
47             // 累加同一单词的所有计数
48             for (IntWritable val : values) {

```

```

45         sum += val.get();
46     }
47     result.set(sum);
48     context.write(key, result);
49 }
50 }
51
52 // 主方法：配置Job并提交任务
53 public static void main(String[] args) throws Exception {
54     Configuration conf = new Configuration();
55     // 显式指定HDFS地址（与你的Hadoop配置一致，避免路径歧义）
56     conf.set("fs.defaultFS", "hdfs://localhost:9000");
57
58     Job job = Job.getInstance(conf, "word count");
59     job.setJarByClass(WordCount.class);
60     // 配置Mapper、Reducer类
61     job.setMapperClass(TokenizerMapper.class);
62     job.setCombinerClass(IntSumReducer.class); // Combiner优化：本地
先汇总，减少网络传输
63     job.setReducerClass(IntSumReducer.class);
64     // 配置输出键值对类型（与Mapper输出、Reducer输出一致）
65     job.setOutputKeyClass(Text.class);
66     job.setOutputValueClass(IntWritable.class);
67
68     // 关键修改：仅添加两个目标文件作为输入（排除其他无关文件）
69     FileInputFormat.addInputPath(job, new
Path("/input/wordfile1.txt"));
70     FileInputFormat.addInputPath(job, new
Path("/input/wordfile2.txt"));
71     // 输出路径：必须是HDFS中不存在的目录（每次运行前需删除旧目录）
72     FileOutputFormat.setOutputPath(job, new
Path("/output_wordcount_final"));
73
74     // 提交Job并等待完成（0表示成功，1表示失败）
75     System.exit(job.waitForCompletion(true) ? 0 : 1);
76 }
77 }

```



1.3.3 3. 代码逻辑解释（对应 Q4-Q8）

问题	答案（基于上述代码）
Q4	key : 输入文本的“行偏移量”（即当前行在整个文件中的起始字节位置，如第 0 行偏移量为 0，第 1 行偏移量为 15 等，类型为 Object）； value : 输入文本的“一行内容”（如 <code>Hello Hadoop</code> ，类型为 Text）。

问题	答案（基于上述代码）
Q5	itr: <code>StringTokenizer</code> 对象，存储对“当前行文本（value）按空白字符拆分后的所有单词”（如对 <code>Hello Hadoop</code> 拆分后，itr 包含 <code>Hello</code> 和 <code>Hadoop</code> 两个元素）。
Q6	word: 存储“当前遍历到的单个单词”（如循环中每次从 itr 取出的 <code>Hello</code> 、 <code>Hadoop</code> 等，类型为 <code>Text</code> ）； one: 固定值 <code>1</code> （类型为 <code>IntWritable</code> ），代表当前单词在该行出现 1 次。
Q7	key: “去重后的单个单词”（Map 输出后会按 key 分组，如所有 <code>Hello</code> 会被分到同一组，key 为 <code>Hello</code> ，类型为 <code>Text</code> ）； values: “同一单词对应的所有 value 集合”（即多个 <code>1</code> ，类型为 <code>Iterable</code> ，如 <code>Hello</code> 出现 3 次，values 包含 3 个 <code>1</code> ）。
Q8	result: “同一单词的总出现次数”（将 values 中的所有 <code>1</code> 累加后的结果，如 <code>Hello</code> 累加后为 3，类型为 <code>IntWritable</code> ）。

1.4 四、编译打包与运行程序

1.4.1 1. 打包 JAR 文件

- ① 在 Eclipse 中右键项目→Export→Java→JAR file→Next。

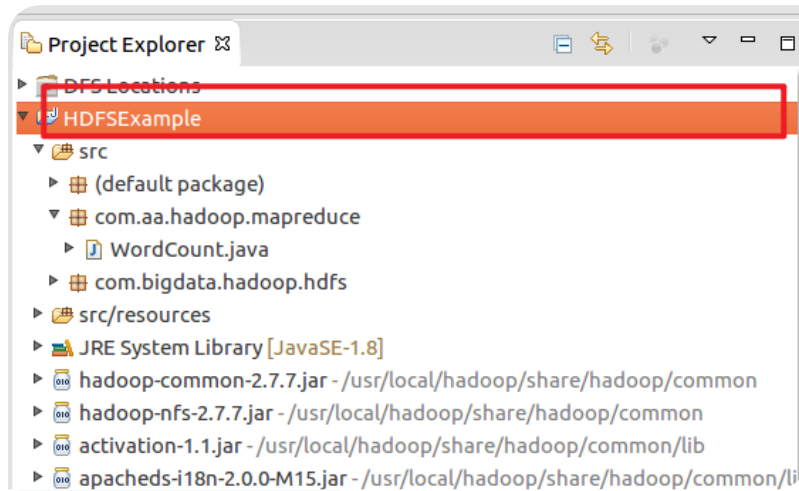


图 1.4.1-1 编译打包-1

- ② 选择要导出的 `WordCount.java` (确保勾选 `src` 目录下的包和类), 指定输出路径 (如 Ubuntu 的 `/usr/local/hadoop/myapp/WordCount.jar`), 点击 Finish。

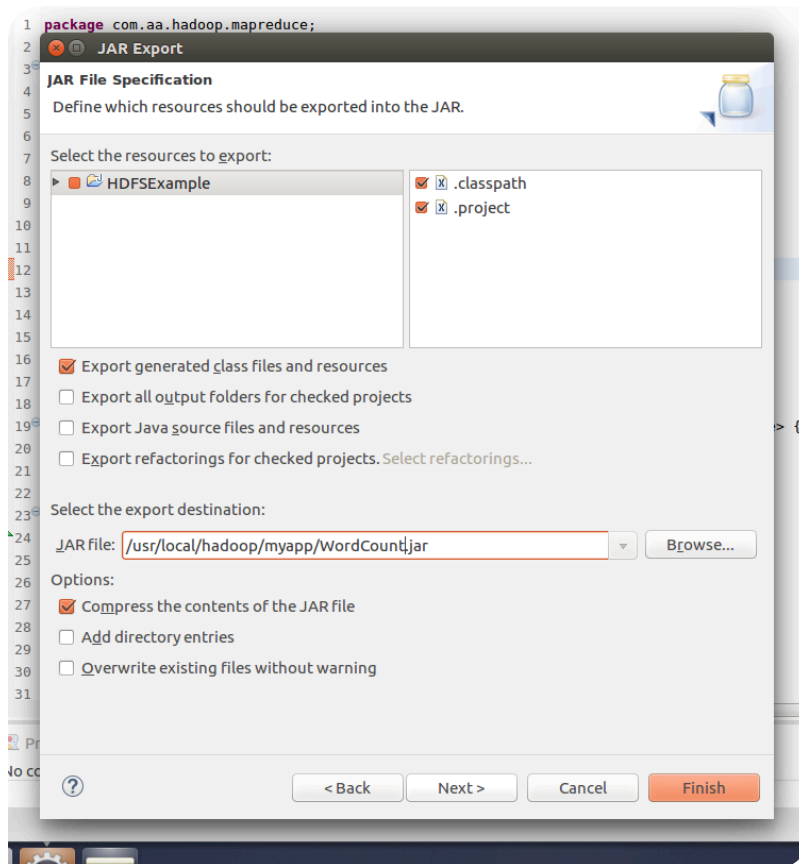


图 1.4.1-2 image-20251030190716028

- ③ 将打包好的 JAR 文件复制到 Ubuntu 的 `/usr/local/hadoop/myapp/` 目录 (若 `myapp` 不存在, 先通过 `mkdir /usr/local/hadoop/myapp` 创建)。

1.4.2 2. 运行 MapReduce 程序

在 Ubuntu 终端中执行以下命令，提交 Job 到 Hadoop 集群：


```
1  #切换到Hadoop安装目录（如/usr/local/hadoop）
2
3  cd /usr/local/hadoop
4
5  #运行JAR包，格式：hadoop jar JAR路径 主类全限定名 输入路径 输出路径
6
7  hadoop jar /usr/local/hadoop/myapp/WordCount.jar
   com.aa.hadoop.mapreduce.WordCount /input/wordfile1.txt
   /input/wordfile2.txt /output_wordcount_final
8
9
10
```



1.4.3 3. 查看运行结果


程序运行成功后，HDFS 会生成 `output` 目录，包含结果文件 `part-r-00000`，通过以下命令查看：

```
1  #查看output目录下的文件
2  hdfs dfs -ls /output
3
4  #执行下面的命令查看结果
5  hdfs dfs -cat /output_wordcount_final/part-r-00000
```



结果：

```
1  Hadoop  2
2  Hdfs    2
3  Hello   3
4  MapReduce  2
5  distribute  1
6  efficient  1
7  good     1
8  is       3
```



```
Terminal File Edit View Search Terminal Help
pInformation.java:1762)
    at org.apache.hadoop.mapreduce.Job.submit(Job.java:1287)
    at org.apache.hadoop.mapreduce.Job.waitForCompletion(Job.java:13
08)
    at com.aa.hadoop.mapreduce.WordCount.main(WordCount.java:75)
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAcces
sorImpl.java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMet
hodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.util.RunJar.run(RunJar.java:226)
    at org.apache.hadoop.util.RunJar.main(RunJar.java:141)
bigdata@ubuntu:/usr/local/hadoop$ hdfs dfs -cat /output_wordcount_final/
part-r-00000
Hadoop      2
Hdfs        2
Hello       3
MapReduce   2
distribute  1
efficient   1
good        1
is          3
bigdata@ubuntu:/usr/local/hadoop$
```

图 1.4.3-1 运行结果

1.5 五、实验收尾

① 运行完成后，可停止 Hadoop 集群：

```
1 stop-yarn.sh # 停止YARN
2 stop-dfs.sh  # 停止HDFS
```