

《并行计算》上机报告

姓名:	魏钊	学号:	PB18111699	日期:	2021/6/2					
上机题目:	Hadoop 编程实验									
实验环境: CPU: i7-8750H; 内存: 16GB;操作系统 Ubuntu20.04;软件平台:Visual Studio 2017										
<p>一、算法设计与分析:</p> <p>题目一: 按照 Hadoop 安装运行说明文档中的指导自己搭建伪分布式 Hadoop 环境,熟悉 HDFS 的常用操作(参考 Hadoop 实战 第 31-36 页),运行 WordCount 程序,得到统计结果。</p> <p>题目二: 实现一个统计输入文件中各个长度的单词出现频次的程序。</p> <p>二、核心代码:</p> <p>题目一: 略</p> <p>题目二: 修改 MAP 即可,将映射由单词内容改为长度</p> <pre> public void map(Object key, Text value, Context context) throws IOException, InterruptedException { StringTokenizer itr = new StringTokenizer(value.toString()); while (itr.hasMoreTokens()) { word.set(itr.nextToken()); Text truth = new Text(Integer.toString(word.getLength())); context.write(truth, one); } } </pre> <p>三、结果与分析:</p> <p>题目一: 按照说明配置完成:</p>										

```
vmware@ubuntu: ~/wordcount
vmware@ubuntu:~/wordcount$ start-all.sh
This script is Deprecated. Instead use start-dfs.sh and start-yarn.sh
Starting namenodes on [localhost]
localhost: starting namenode, logging to /usr/local/hadoop/logs/hadoop-vmware-na
menode-ubuntu.out
localhost: starting datanode, logging to /usr/local/hadoop/logs/hadoop-vmware-da
tanode-ubuntu.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /usr/local/hadoop/logs/hadoop-vm
ware-secondarynamenode-ubuntu.out
starting yarn daemons
starting resourceemanager, logging to /usr/local/hadoop/logs/yarn-vmware-resource
manager-ubuntu.out
localhost: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-vmware-n
odemanager-ubuntu.out
vmware@ubuntu:~/wordcount$ jps
2994 NameNode
3699 NodeManager
3399 SecondaryNameNode
3545 ResourceManager
4043 Jps
3181 DataNode
vmware@ubuntu:~/wordcount$
```

文本编辑器 6月2日 20:18

part-r-00000
~/wordcount

1	can	2
2	die	2
3	no	4
4	up	2
5	you	4
6	zuo	2

题目二:

```
vmware@ubuntu: ~/lengthcount
vmware@ubuntu:~/lengthcount$ hadoop fs -mkdir /user/vmware/lengthcount
vmware@ubuntu:~/lengthcount$ hadoop fs -mkdir /user/vmware/lengthcount/input
vmware@ubuntu:~/lengthcount$ touch input1.txt
vmware@ubuntu:~/lengthcount$ touch input2.txt
vmware@ubuntu:~/lengthcount$ hadoop fs -put ./input1.txt /user/vmware/lengthcoun
t/input
vmware@ubuntu:~/lengthcount$ hadoop fs -put ./input2.txt /user/vmware/lengthcoun
t/input
vmware@ubuntu:~/lengthcount$
```

<div><div>vmware@ubuntu: ~/lengthcount</div><div><div>Spilled Records=10</div><div>Shuffled Maps =2</div><div>Failed Shuffles=0</div><div>Merged Map outputs=2</div><div>GC time elapsed (ms)=59</div><div>Total committed heap usage (bytes)=973602816</div><div>Shuffle Errors</div><div>BAD_ID=0</div><div>CONNECTION=0</div><div>IO_ERROR=0</div><div>WRONG_LENGTH=0</div><div>WRONG_MAP=0</div><div>WRONG_REDUCE=0</div><div>File Input Format Counters</div><div>Bytes Read=44</div><div>File Output Format Counters</div><div>Bytes Written=12</div><div>vmware@ubuntu:~/lengthcount\$ hadoop fs -ls /user/vmware/lengthcount/output</div><div>Found 2 items</div><div>-rw-r--r-- 1 vmware supergroup 0 2021-06-02 20:24 /user/vmware/lengthcount/output/_SUCCESS</div><div>-rw-r--r-- 1 vmware supergroup 12 2021-06-02 20:24 /user/vmware/lengthcount/output/part-r-00000</div><div>vmware@ubuntu:~/lengthcount\$</div></div></div>	
<div><div>打开(O)</div><div>part-r-00000</div><div>~/lengthcount</div><div><div>116</div><div>228</div><div>332</div></div></div>	
<div>四、备注 (* 可选):</div> <div>有可能影响结论的因素:</div>	
<div>总结:</div> <div>通过分布式系统实现并行</div>	
附录 (源代码)	<div>算法源代码 (C/C++/JAVA 描述)</div> <div>import java.io.IOException;</div> <div>import java.util.StringTokenizer;</div> <div>import org.apache.hadoop.conf.Configuration;</div> <div>import org.apache.hadoop.fs.Path;</div> <div>import org.apache.hadoop.io.IntWritable;</div> <div>import org.apache.hadoop.io.Text;</div> <div>import org.apache.hadoop.mapreduce.Job;</div> <div>import org.apache.hadoop.mapreduce.Mapper;</div> <div>import org.apache.hadoop.mapreduce.Reducer;</div> <div>import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;</div> <div>import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;</div>

```
import org.apache.hadoop.util.GenericOptionsParser;

public class WordCount {

    public static class TokenizerMapper
        extends Mapper<Object, Text, Text, IntWritable>{

        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(Object key, Text value, Context context
            ) throws IOException, InterruptedException {
            StringTokenizer itr = new StringTokenizer(value.toString());
            while (itr.hasMoreTokens()) {
                word.set(itr.nextToken());
                Text truth = new Text(Integer.toString(word.getLength()));
                context.write(truth, one);
            }
        }
    }

    public static class IntSumReducer
        extends Reducer<Text, IntWritable, Text, IntWritable> {
        private IntWritable result = new IntWritable();

        public void reduce(Text key, Iterable<IntWritable> values,
            Context context
            ) throws IOException, InterruptedException {

            int sum = 0;
            for (IntWritable val : values) {
                sum += val.get();
            }
            result.set(sum);
            context.write(key, result);
        }
    }

    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        String[] otherArgs = new GenericOptionsParser(conf,
args).getRemainingArgs();
        if (otherArgs.length != 2) {
            System.err.println("Usage: wordcount <in> <out>");
            System.exit(2);
        }
    }
}
```

	<pre>} Job job = new Job(conf, "word count"); job.setJarByClass(WordCount.class); job.setMapperClass(TokenizerMapper.class); job.setCombinerClass(IntSumReducer.class); job.setReducerClass(IntSumReducer.class); job.setOutputKeyClass(Text.class); job.setOutputValueClass(IntWritable.class); FileInputFormat.addInputPath(job, new Path(otherArgs[0])); FileOutputFormat.setOutputPath(job, new Path(otherArgs[1])); System.exit(job.waitForCompletion(true) ? 0 : 1); } }</pre>
--	--