# 作业3

姓名： 吴双

学号： 10164102141

# 1 软件体系结构

WordSort

要求

将之前的WordCount改为WordSort排序

思路

同样，使用一个list来记录当前位置，同时使用二分查找找到当前位置并插入

top-k

要求

单词的top-k：求最频繁的k个word 要考虑代码的性能

思路

首先仍然与WordCount一样，需要一个 hashmap 来存当前word对应的count数量 考虑到性能原因，只需要维护前k个word即可 需要设置最小count，表示前k个中最小的count数，如果当前有word的count数大于这个数，则将其加入TopK数组，然后剔除

# 2 详细设计

## Wordsort:

SentenceSpout

```
1  public class SentenceSpout extends BaseRichSpout {
2
3      private SpoutOutputCollector spoutOutputCollector;
4      private String[] sentences = {"the cow jumped over the moon", "an apple a day
   keeps the doctor away",
```

```
  5              "four score and seven years ago", "snow white and the seven dwarfs", "i
     am at two with nature"};
  6
  7      public void open(Map map, TopologyContext topologycontext, SpoutOutputCollector
     spoutoutputcollector) {
  8          this.spoutOutputCollector = spoutoutputcollector;
  9      }
 10
 11      public void nextTuple() {
 12          for (String sentence : sentences) {
 13              Values values = new Values(sentence);
 14              UUID msgId = UUID.randomUUID();
 15              this.spoutOutputCollector.emit(values, msgId);
 16          }
 17          Utils.sleep(1000);
 18      }
 19
 20      public void declareOutputFields(OutputFieldsDeclarer outputfieldsdeclarer) {
 21          outputfieldsdeclarer.declare(new Fields("sentence"));
 22      }
 23
 24  }
```

SplitSentenceBolt

```
  1  public class SplitSentenceBolt extends BaseBasicBolt {
  2      public void execute(Tuple tuple, BasicOutputCollector collector) {
  3          String sentence = tuple.getStringByField("sentence");
  4          String[] words = sentence.split(" ");
  5          for (String word : words) {
  6              collector.emit(new Values(word));
  7          }
  8      }
  9
 10      public void declareOutputFields(OutputFieldsDeclarer outputfieldsdeclarer) {
 11          outputfieldsdeclarer.declare(new Fields("word"));
 12      }
 13
 14  }
```

WordSortBolt

```
  1  public class WordSortBolt extends BaseBasicBolt {
  2      List<String> wordList = new ArrayList<String>();
  3
  4      public int arrayIndexOf(String key) {
  5          int min, max, mid;
  6          min = 0;
  7          max = wordList.size() - 1;
```

```
 8
 9          while (min <= max) {
10              mid = (min + max) >> 1;
11              String tmp = wordList.get(mid);
12              if (key.compareTo(tmp) > 0) {
13                  min = mid + 1;
14              } else if (key.compareTo(tmp) < 0) {
15                  max = mid - 1;
16              } else {
17                  return mid;
18              }
19          }
20          return min;
21      }
22
23
24      public void execute(Tuple tuple, BasicOutputCollector collector) {
25          String word = tuple.getString(0);
26          if (wordList == null) {
27              wordList.add(word);
28              System.out.println(word);
29          } else {
30              int addIndex = arrayIndexOf(word);
31              wordList.add(addIndex,word);
32              for (String tmp:wordList
33                   ) {
34                  System.out.println(tmp);
35              }
36          }
37          collector.emit(new Values(word));
38      }
39
40      public void declareOutputFields(OutputFieldsDeclarer outputfieldsdeclarer) {
41          outputfieldsdeclarer.declare(new Fields("word"));
42      }
43  }
```

WordSortTopology

```
 1  public class WordSortTopology {
 2
 3      public static void main(String[] args) throws Exception {
 4
 5          SentenceSpout sentenceSpout = new SentenceSpout();
 6          SplitSentenceBolt splitSentenceBolt = new SplitSentenceBolt();
 7          WordSortBolt wordSortBolt = new WordSortBolt();
 8
 9          TopologyBuilder builder = new TopologyBuilder();
10          builder.setSpout("sentenceSpout-1", sentenceSpout);
11          builder.setBolt("splitSentenceBolt-1",
    splitSentenceBolt).shuffleGrouping("sentenceSpout-1");
```

```
12          builder.setBolt("wordSortBolt-1",
    wordSortBolt).fieldsGrouping("splitSentenceBolt-1", new Fields("word"));
13
14          Config config = new Config();
15          LocalCluster cluster = new LocalCluster();
16
17          cluster.submitTopology("wordSortTopology-1", config,
    builder.createTopology());
18          Thread.sleep(999999999);
19          cluster.shutdown();
20      }
21
22  }
```

将程序打包成JAR包，并shell input:

```
1  storm jar WordSort.jar WordSortTopology ws
```

output:



## K-top:

SentenceSpout

```
1  public class SentenceSpout extends BaseRichSpout {
2
3      private SpoutOutputCollector spoutOutputCollector;
```

```
4    private String[] sentences = {"the cow jumped over the moon", "an apple a day
     keeps the doctor away",
5              "four score and seven years ago", "snow white and the seven dwarfs", "i
     am at two with nature"};
6
7    public void open(Map map, TopologyContext topologycontext, SpoutOutputCollector
     spoutoutputcollector) {
8        this.spoutOutputCollector = spoutoutputcollector;
9    }
10
11   public void nextTuple() {
12       for (String sentence : sentences) {
13           Values values = new Values(sentence);
14           UUID msgId = UUID.randomUUID();
15           this.spoutOutputCollector.emit(values, msgId);
16       }
17       Utils.sleep(1000);
18   }
19
20   public void declareOutputFields(OutputFieldsDeclarer outputfieldsdeclarer) {
21       outputfieldsdeclarer.declare(new Fields("sentence"));
22   }
23
24 }
```

SplitSentenceBolt

```
1  public class SplitSentenceBolt extends BaseBasicBolt {
2      public void execute(Tuple tuple, BasicOutputCollector collector) {
3          String sentence = tuple.getStringByField("sentence");
4          String[] words = sentence.split(" ");
5          for (String word : words) {
6              collector.emit(new Values(word));
7          }
8      }
9
10     public void declareOutputFields(OutputFieldsDeclarer outputfieldsdeclarer) {
11         outputfieldsdeclarer.declare(new Fields("word"));
12     }
13
14 }
```

pair

```
1  public class pair {
2      public final String content;
3      public final Integer count;
4      public pair(String content, Integer count) {
5          this.content = content;
```

```
 6          this.count = count;
 7      }
 8
 9      public int compareCount(pair other) {
10          return this.count - other.count;
11      }
12
13      public int compareWord(pair other) {
14          return this.content.compareTo(other.content);
15      }
16  }
```

TopK 使用HashMap来保存当前所有word的count数量，使用TopList保存前k个word

```
 1  public class TopK extends BaseBasicBolt {
 2      private HashMap<String, Integer> counts;
 3      private ArrayList<pair> TopList;
 4      public int K;
 5      public int minCount;
 6
 7      TopK(int k) {
 8          this.K = k;
 9          this.counts = new HashMap<>();
10          this.TopList = new ArrayList<>();
11          this.minCount = 0;
12      }
13
14
15      public void insertWord(pair word) {
16          int max = TopList.size() - 1;
17
18  //          if the same word, the new.count > old.count
19          for (int i = max; i >= 0; i--) {
20              pair tmp = TopList.get(i);
21              // find the same word,replace older one
22              if (word.compareWord(tmp) == 0) {
23                  TopList.set(i, word);
24                  return;
25              }
26              if (word.compareCount(tmp) <= 0) {
27                  TopList.add(i + 1, word);
28                  return;
29              }
30          }
31          TopList.add(0, word);
32      }
33
34
35      public void execute(Tuple tuple, BasicOutputCollector collector) {
36          String word = tuple.getStringByField("word");
37          Integer count = counts.get(word);
```

```
38          if (count == null) {
39              count = 0;
40          }
41          count++;
42          this.counts.put(word, count);
43          if (count > minCount || TopList.size() < K) {
44              insertWord(new pair(word, count));
45              if (TopList.size() > K) {
46                  TopList.remove(TopList.size() - 1);
47                  minCount = TopList.get(TopList.size() - 1).count;
48              }
49              for (pair tmpWord : TopList)
50                  System.out.println(tmpWord.content + " " + tmpWord.count.toString());
51              collector.emit(new Values(word, count));
52          }
53
54
55      }
56
57      public void declareOutputFields(OutputFieldsDeclarer outputfieldsdeclarer) {
58          outputfieldsdeclarer.declare(new Fields("word", "count"));
59      }
60 }
```
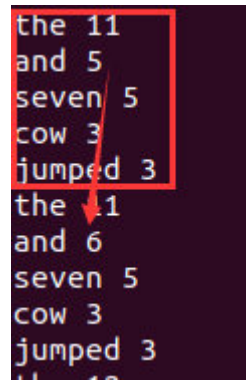
WordTopKTopology

```
1  public class WordTopKTopology {
2      public static void main(String[] args) throws Exception {
3
4          SentenceSpout sentenceSpout = new SentenceSpout();
5          SplitSentenceBolt splitSentenceBolt = new SplitSentenceBolt();
6          TopK wordTopKBolt = new TopK(5);
7
8          TopologyBuilder builder = new TopologyBuilder();
9          builder.setSpout("sentenceSpout-1", sentenceSpout);
10         builder.setBolt("splitSentenceBolt-1",
   splitSentenceBolt).shuffleGrouping("sentenceSpout-1");
11         builder.setBolt("wordTopKBolt-1",
   wordTopKBolt).shuffleGrouping("splitSentenceBolt-1");
12
13         Config config = new Config();
14         LocalCluster cluster = new LocalCluster();
15
16         cluster.submitTopology("wordTopKTopology-1", config,
   builder.createTopology());
17         Thread.sleep(999999999);
18         cluster.shutdown();
19     }
20 }
```

shell input:

```
storm jar TopK.jar WordTopKTopology tk
```

output: