

CSE 587 Spring 2023

Name: Pavana Lakshmi Venugopal

UBIT: pavanala

Analysis

Answer the following questions based on your WordCount application.

1. In the PySpark REPL, run your basic word count program on a single text file.
 - a. What are the 25 most common words? Include a screenshot of program output to back-up your claim.
 - b. How many stages is execution broken up into? Explain why. Include a screenshot of the DAG visualization from Spark's WebUI to back-up your claim.

Solution:

```
1 def countWords(sc):
2     import re, nltk
3
4     from nltk.corpus import stopwords
5
6     stopwords = stopwords.words('english')
7
8     file = 'book1.txt'
9
10    lines = sc.textFile(file)
11    counts = lines.flatMap(lambda x: x.split(' ')) \
12                    .map(lambda x: re.sub(r'^\w\s', '', x)) \
13                    .map(lambda x: (x.lower(),1)) \
14                    .filter(lambda x: x[0] not in stopwords + ['']) \
15                    .reduceByKey(lambda a,b: a+b) \
16                    .sortBy(lambda pair: pair[1], ascending=False)
17
18    # counts.saveAsTextFile("C:/DIC-Project/HW2/output")
19    return counts
20
21 if __name__ == '__main__':
22     from pyspark.context import SparkContext
23     sc = SparkContext('local', 'test')
24
25     results.saveAsTextFile("C:/DIC-Project/HW2/output")
```

a) The 25 most common words are shown below.

```
1 ('co', 3323)
2 ('biograph', 1949)
3 ('edison', 1827)
4 ('american', 1714)
5 ('mutoscope', 1641)
6 ('thomas', 1363)
7 ('scenes', 844)
8 ('mfg', 574)
9 ('14', 566)
10 ('lubin', 514)
11 ('lc', 483)
12 ('america', 438)
13 ('vitagraph', 429)
14 ('inc', 379)
15 ('new', 341)
16 ('york', 205)
17 ('méliès', 189)
18 ('see', 184)
19 ('love', 173)
20 ('parade', 155)
21 ('scene', 154)
22 ('fight', 148)
23 ('city', 146)
24 ('picture', 145)
25 ('war', 145)
```

b) The execution is broken into 2 stages.

Details for Job 1

Status: SUCCEEDED

Submitted: 2023/05/01 22:13:38

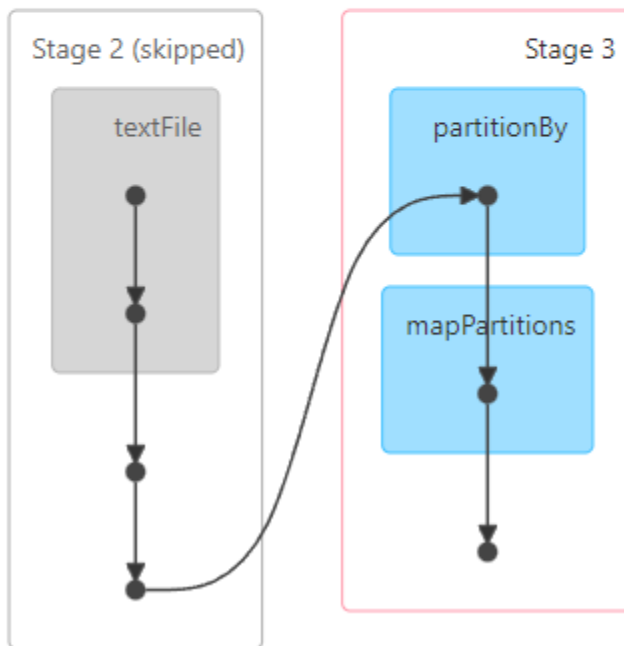
Duration: 2 s

Completed Stages: 1

Skipped Stages: 1

► Event Timeline

▼ DAG Visualization



Stage2 : Started with book1.txt. Spark uses in memory caching. It would have the output of those stages already, so it skipped it. Basically this is for reading the text file using `sc.textFile(file)` creates a RDD that needs to be stored in memory and partitioned for parallel processing, thereby creating the first stage of the execution.

Stage 3: It basically picks up where it left off. All this stage really does is shuffle the data around and remap the partitions. The transformations "flatMap," "map," and "filter" are applied to the RDD produced in the previous stage. This stage includes the `reduceByKey` and `sortBy` transformations and the take action on the RDD. The transformations are not executed immediately but are instead evaluated and executed together when they are needed. Then, using the 'sortBy' transformation to order the resulting word counts in descending order, see the top 25 frequently occurring words.

Details for Stage 3 (Attempt 0)

Resource Profile Id: 0

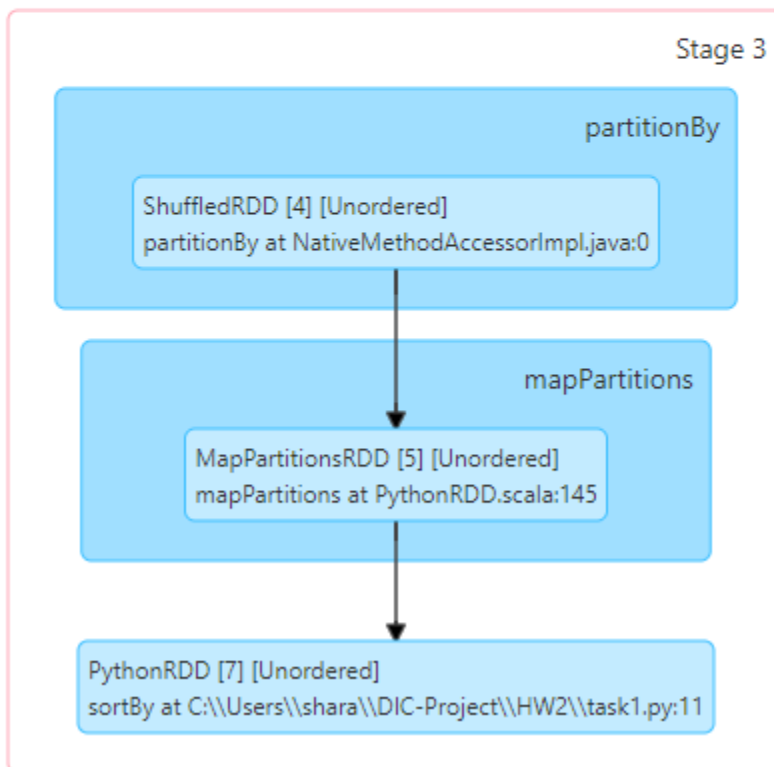
Total Time Across All Tasks: 3 s

Locality Level Summary: Node local: 2

Shuffle Read Size / Records: 170.1 KiB / 56

Associated Job Ids: 1

▼ DAG Visualization



2. In the PySpark REPL, run your extended word count program on all 10 text files.
 - a. What are the 25 most common words? Include a screenshot of program output to back-up your claim.
 - b. How many stages is execution broken up into? Explain why. Include a screenshot of the DAG visualization from Spark's WebUI to back-up your claim.

Solution:

a)

```
1 def task2(sc, files):
2     import re, nltk
3
4     from nltk.corpus import stopwords
5
6     stopwords = stopwords.words('english')
7
8     output = None
9
10    counts_list = []
11
12
13    count = [sc.textFile(file) \
14              .flatMap(lambda x: x.split(' ')) \
15              .map(lambda x: re.sub(r'^\w\s', '', x)) \
16              .map(lambda x: (x.lower(),1)) \
17              .filter(lambda x: x[0] not in stopwords + ['']) \
18              .reduceByKey(lambda a,b: a+b) for file in files]
19
20
21    counts = sc.union(count) \
22              .sortBy(lambda pair: pair[1], ascending=False)
23
24    # counts.saveAsTextFile("C:/DIC-Project/HW2/output")
25    return counts
26
27 if __name__ == '__main__':
28     from pyspark.context import SparkContext
29     sc = SparkContext('local', 'test')
30
31     results = task2(sc, ['book1.txt', 'book2.txt', 'book3.txt', 'book4.txt', 'book5.txt', 'book6.txt', 'book7.txt', 'book8.txt', 'book9.txt',
32                          'book10.txt'])
33
34     results.saveAsTextFile("C:/DIC-Project/HW2/output2")
```

The 25 most common words are as shown below.

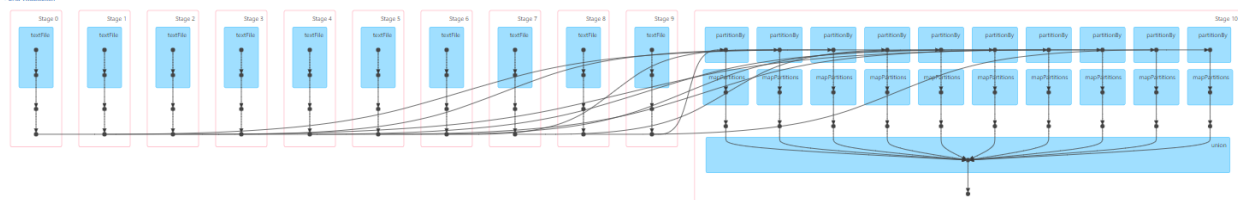
```
1 ('co', 3323)
2 ('biograph', 1949)
3 ('edison', 1827)
4 ('american', 1714)
5 ('mutoscope', 1641)
6 ('thomas', 1363)
7 ('mrs', 958)
8 ('scenes', 844)
9 ('mr', 805)
10 ('upon', 785)
11 ('fig', 754)
12 ('plant', 747)
13 ('plants', 732)
14 ('leaves', 675)
15 ('one', 640)
16 ('elizabeth', 603)
17 ('water', 597)
18 ('mfg', 574)
19 ('14', 566)
20 ('cells', 564)
21 ('would', 560)
22 ('could', 531)
23 ('said', 529)
24 ('illustration', 519)
25 ('would', 517)
```

b) The execution is broken down into 11 stages.

Details for Job 0

Status: SUCCEEDED
Submitted: 2022/05/01 22:36:27
Duration: 27 s
Completed Stages: 11

• Spark Timeline
• DAG Visualization



Stage0 to Stag9 : Reads all text files. Spark uses in memory caching. It would have the output of those stages already, so it skipped it. Basically this is for reading the text file using `sc.textFile(file)` creates a RDD that needs to be stored in memory and partitioned for parallel processing, thereby creating the first

stage of the execution. This transformation takes a line of text as input and separates it into individual words, resulting in a list of words. (flatmap)

Details for Stage 0 (Attempt 0)

Resource Profile Id: 0

Total Time Across All Tasks: 16 s

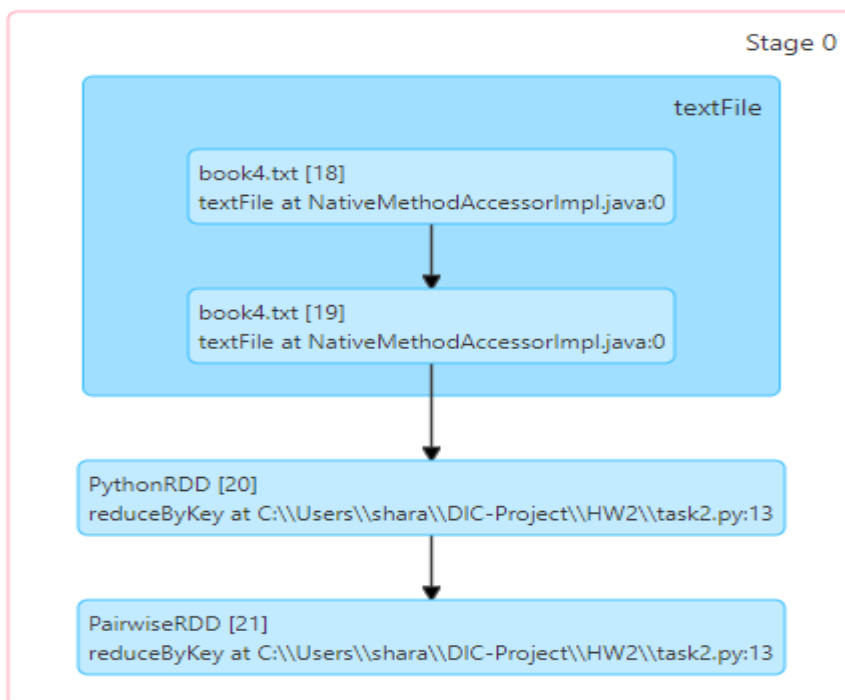
Locality Level Summary: Process local: 2

Input Size / Records: 522.9 KiB / 9100

Shuffle Write Size / Records: 90.5 KiB / 50

Associated Job Ids: 0

▼ DAG Visualization



Stage10: Each word is changed to lowercase and all punctuation is removed in this transition. (map). This transformation counts the number of occurrences of each word. (reduceByKey). It includes take action on the RDD.

Details for Stage 10 (Attempt 0)

Resource Profile Id: 0
Total Time Across All Tasks: 20 s
Shuffle Read Size / Records: 100.0 KiB / 10
Associated Job Ids: 0

▼ DAG Visualization



▼ Show Additional Details
 0 Task Failure
 0 Resource Manager Action 0 Unsuccessful Reads

3. Your WordCount application should compute the same results as your WordCount application from Homework #1. Answer the following based on your knowledge of both MapReduce and Spark:
- If you were running your WordCount programs in a large cluster or cloud environment, and one of the nodes you were running on died mid computation, how would your MapReduce and Spark programs handle this?
 - Explain one concrete benefit you experienced when writing the Spark version of WordCount compared to the MapReduce version.

Solution:

a) MapReduce offers linear scalability and fault tolerance for processing very large data sets. Spark maintains this revolutionary approach brought about by MapReduce.

In MapReduce, if a node fails during computation, the framework identifies the failed task and schedules its re-execution on another available node. The fact that the data is copied and stored on HDFS makes this possible. The intermediate results that were generated by the node that failed are recovered using the Hadoop Distributed File System (HDFS).

Spark's fault tolerance involves recovering partitions when there is a failure. In Spark, a node failure can be handled by its lineage mechanism. Each RDD's lineage, which shows the data transformations carried out on the RDD, is tracked by Spark. If a partition of a Resilient Distributed Dataset (RDD) is lost due to a node failure, the lineage of the RDD can be used to recompute the lost partition. Spark performs recovery automatically and ensures that the transformed RDDs are reprocessed only when needed, thereby minimizing the computation cost.

b) Spark transitions seamlessly between exploratory analytics and operational analytics.

Related to productivity, we saw wordcount in MapReduce had multiple java classes we had to define, we had a configuration to setup. Our mappers were very finely grained, explaining to read the file, split that into words go through each word, emit a key value pair and the reducer will take those key value pairs and it will reduce them into a single key value pair.


```

import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    // Mapper
    // takes the text, tokenize the words
    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 {
            // tokenize the lowercased text string
            String line = value.toString().toLowerCase(); // convert line to lowercase and to string variable
            StringTokenizer itr = new StringTokenizer(line, "!\"#$%&'()*+,-./:;<=>?@[^_`{|}~ "); // remove possible punctuations
            while (itr.hasMoreTokens()) {
                word.set(itr.nextToken());
                context.write(word, one);
            }
        }
    }
}

```

```

// Reducer
// counts the number of times the word is in data and writes the word and sum as key value pair
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();
    Job job = Job.getInstance(conf, "word count");
    job.setJar("WordCount.jar");
    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(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    job.waitForCompletion(true);
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}

```

In Spark its very streamlined into something that much more resembles something like a more basic script. There is again a series of operations like flat map, map and reduce by key. How those operate on data, how the data is communicated and how that data is distributed all that is handled by Spark for us.

```
def countWords(sc, file):
    import re, nltk

    nltk.download('stopwords')

    from nltk.corpus import stopwords

    stopwords = stopwords.words('english')

    lines = sc.textFile(file)
    counts = lines.flatMap(lambda x: x.split(' ')) \
        .map(lambda x: re.sub(r'^\w\s', '', x)) \
        .map(lambda x: (x.lower(),1)) \
        .filter(lambda x: x[0] not in stopwords + ['']) \
        .reduceByKey(lambda a,b: a+b) \
        .sortBy(lambda pair: pair[1], ascending=False) \
        .take(25)
    return counts

if __name__ == '__main__':
    from pyspark.context import SparkContext
    sc = SparkContext('local', 'test')

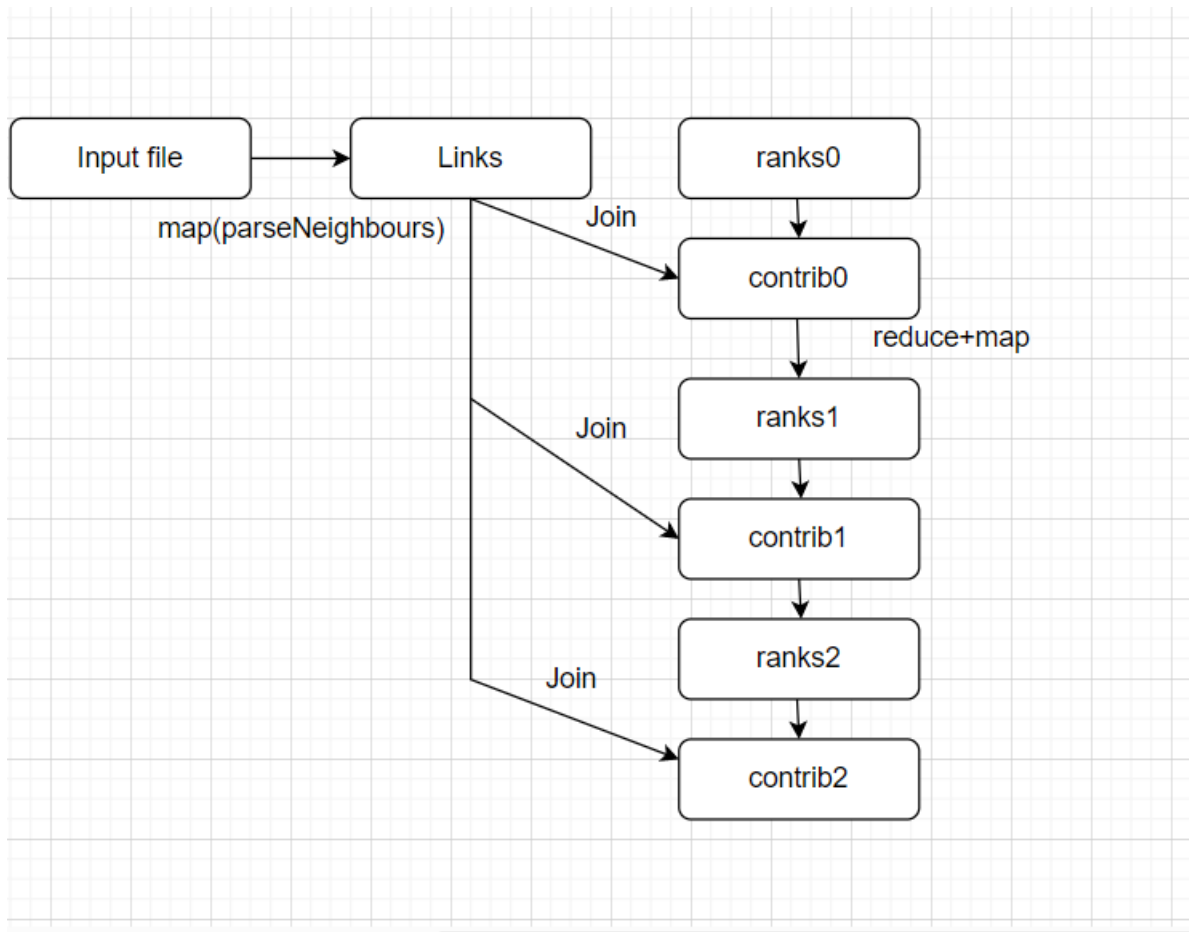
    counts = countWords(sc, "/*.book1.txt")

    for word, count in counts:
        print("{} : {}".format(word, count))
```

Also, the Spark implementation of WordCount is more brief and readable than the MapReduce implementation, while also providing better performance due to its ability to store intermediate data in memory and perform operations in parallel.

4. Given the above spark application, draw the lineage graph DAG for the RDD **ranks** on line 12 when the iteration variable **i** has a value of 2. Include nodes for all intermediate RDDs, even if they are unnamed.

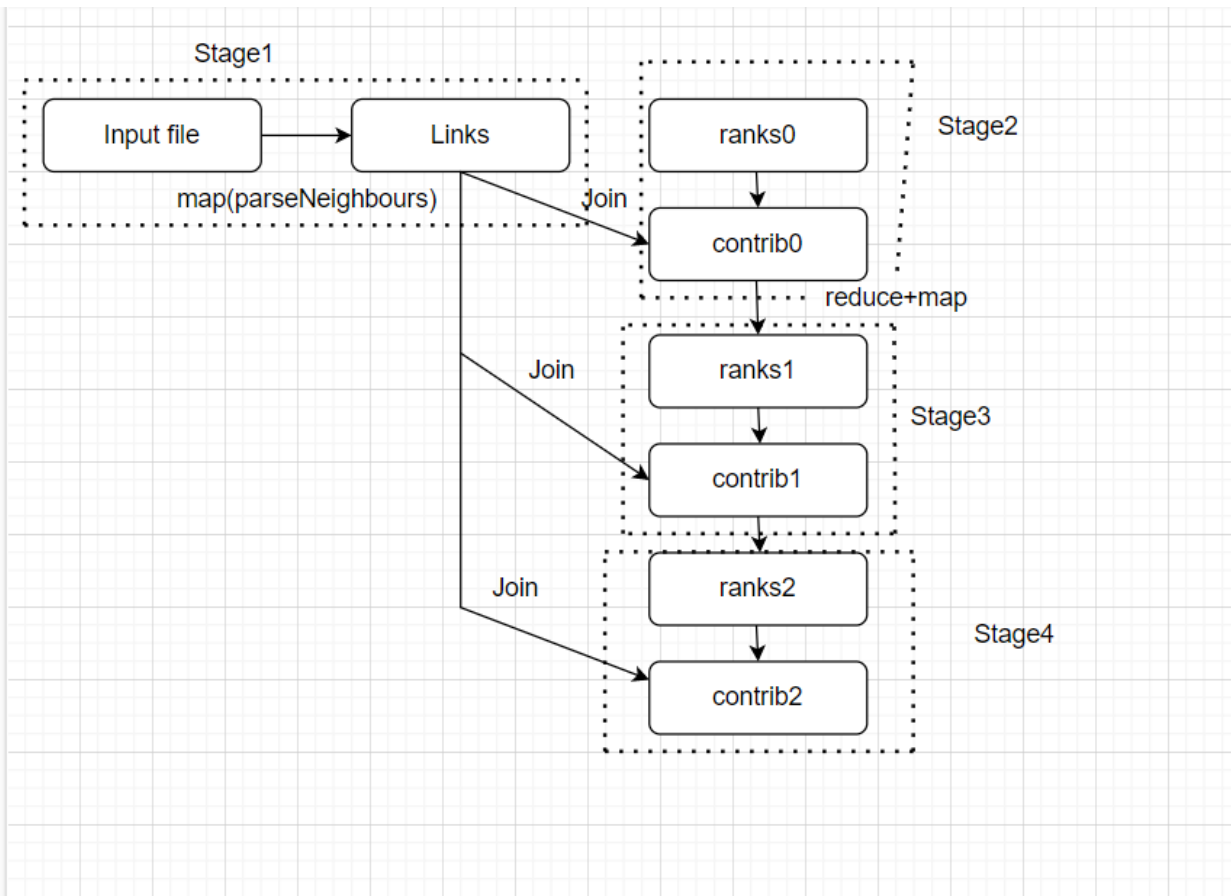
Solution:



5. How many stages will the above DAG be broken into? Give the number of stages AND draw stage boundaries on your diagram.

Solution:

Above DAG will be broken down into 4 stages.



6. Identify in the above code (by function name AND line number) one instance of:
- A transformation that results in a wide dependency
 - A transformation that results in a narrow dependency
 - A transformation that may result in a narrow dependency OR a wide dependency
 - An action

Solution:

a) `groupByKey()`

Line 3

b) `map(func)`

Line 2, Line 6

c) `join()`

Line 9

d) `count()`

Line 5

7. How many "jobs" will the above code run if **iters** has value 10?

Solution:

Number of jobs will still be the same that is one.

8. What algorithm is the above code an implementation of?

Solution:

The above code is an implementation of Page rank algorithm.