# Project 3: Hadoop Blast
# Cloud Computing
# Spring 2017

Professor Judy Qiu

## Goal

By this point you should have gone over the sections concerning Hadoop Setup and a few Hadoop programs. Now you are going to blend these applications by implementing a parallel version of BLAST (Basic Local Alignment Search Tool: `http://blast.ncbi.nlm.nih.gov/Blast.cgi`) using the programming interfaces of the Hadoop MapReduce framework. Note that this application is written in "Map-Only" fashion, which means no reduce code is necessary.

## Deliverables

You are required to turn in the following items in a zip file (username_HadoopBlast.zip)

- The source code of Hadoop Blast you implemented.

- Technical report (username_HadoopBlast_report.docx) that answers the following questions.

    - What is Hadoop Distributed Cache and how is it used in this program?
    - Write the two lines that put and get values from Distributed cache. Also include the method and class information.
    - In previous projects we used Hadoop's TextInputFormat to feed in the file splits line by line to map tasks. In this program, however, we want to feed in a whole file to a single map task. What is the technique used to achieve this? Also, briefly explain what are the key and value pairs you receive as input to a map task and what methods are responsible for producing these pairs?
    - Do you think this particular implementation will work if the input files are larger than the default HDFS block size? Briefly explain why. [Hint: you can test what will happen by concatenating the same input file multiple times to create a larger input file in the resources/blast_input folder]
    - If you wanted to extend this program such that all output files will be concatenated into a single file, what key and value pairs would you need to emit from the map task? Also, how would you use these in the reduce that you would need to add?

- The 4 output FASTA files: celllines_1.fa to celllines_4.fa.

## Evaluation

The point total for this project is 3, where the distribution is as follows:

- Completeness of your code and output (1 points)

- Correctness of written report (2 points)

# Introduction

Hadoop-Blast is an advanced Hadoop program which helps BLAST, a bioinformatics application, to utilize the computing capability of Hadoop. This exercise shows the details of its implementation, and provides an example of how to handle similar approaches in other applications.

BLAST is one of the most widely used bioinformatics applications written in C++. The version we are using is v2.2.23, which houses new features and better performance. The database used in the following settings is a subset of a full 8.5GB(nr)database; its full name is Non-redundant protein sequence database. Optionally, for more details on how to run the BLAST binary, please see Big Data for Science tutorial page for Blast Installation [NOT required for the assignment].

In this project, we have provided a sketch code which contains just one java class for you to implement:

- RunnerMap.java: The pleasingly-parallel/map-only Map class which takes the prepackaged Blast (v2.2.23) Binary Program and optimized database from Hadoop's Distributed Cache, then executes BLAST binary as java external process with the assigned FASTA file. These are passed as key-value pairs of **(filename, filepath on HDFS)** handled by a provided customized Hadoop MapReduce InputFormat DataFileInputFormat.java.

The detail dataflow can be seen in Figure 1. You will implement the RunnerMap.java, which copies the distributed cache and assigned FASTA file to local, then run the BLAST binary with correct parameters.
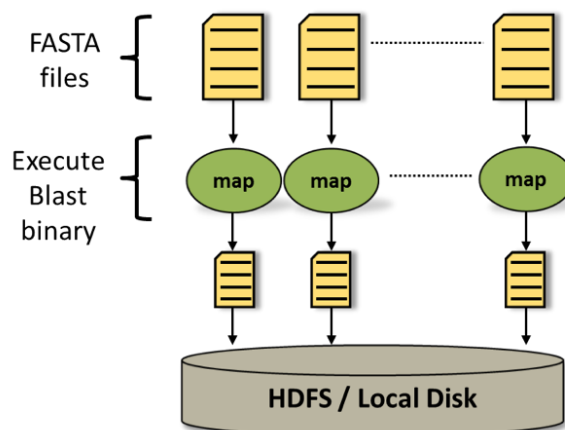


Figure 1: Hadoop Blast dataflow

Normally, for any Hadoop MapReduce program, input data is uploaded and stored in the Hadoop Distributed File System (HDFS) before computation in order to generate **(key, value)** pairs to the mapper. Initially, the BLAST input data is a set of FASTA files located in the local file system. Then it will be uploaded to the HDFS and distributed across the compute nodes. Hadoop framework reads the application records from HDFS with the InputFormat interface and generates **(key, value)** pair input streams; here, we use a provided customized Hadoop MapReduce InputFormat DataFileInputFormat.java to generate key-value pairs of **(filename, filepath on HDFS)**. For this Hadoop Blast program, the map function initially sets up the distributed cache and generates the two absolute location filepaths for Blast binary and Blast Database. Afterwards it copies the assigned FASTA file to local disk by looking up the file from HDFS and generating an absolute filepath. Once this is accomplished and file dependencies are stored in the local disk, we call an external java process and execute the Blast binary with the correct parameters. Finally, the output FASTA file of Blast binary will be uploaded back to HDFS.

# Sketch for Hadoop Blast

You need to complete one file in the provided pacakge inside "cgl/hadoop/apps/runner": RunnerMap.java. Code snapshots are shown below.

```
1  /*file: RunnerMap.java*/
2  package cgl.hadoop.apps.runner;
3
4  import java.io.File;
5  import java.io.IOException;
6  import java.net.URI;
7  import java.net.URISyntaxException;
8
9  import org.apache.hadoop.conf.Configuration;
10 import org.apache.hadoop.filecache.DistributedCache;
11 import org.apache.hadoop.fs.FileSystem;
12 import org.apache.hadoop.fs.Path;
13 import org.apache.hadoop.io.IntWritable;
14 import org.apache.hadoop.io.Text;
15 import org.apache.hadoop.mapreduce.Mapper;
16
17 /**
18  * @author Thilina Gunarathne (tgunarat@cs.indiana.edu)
19  *
20  * @editor Stephen, TAK-LON WU (taklwu@indiana.edu)
21  */
22
23 public class RunnerMap extends Mapper<String, String, IntWritable, Text> {
24
25   private String localDB = "";
26   private String localBlastProgram = "";
27
28
29   @Override
30   public void setup(Context context) throws IOException{
31     Configuration conf = context.getConfiguration();
32     Path[] local = DistributedCache.getLocalCacheArchives(conf);
33
34     /* Write your code here
35           get two absolute filepath for localDB and localBlastBinary
36       */
37
38   }
39
40
41   public void map(String key, String value, Context context) throws IOException,
42     InterruptedException {
43
44     long startTime = System.currentTimeMillis();
45     String endTime = "";
46
47     Configuration conf = context.getConfiguration();
48     String programDir = conf.get(DataAnalysis.PROGRAM_DIR);
49     String execName = conf.get(DataAnalysis.EXECUTABLE);
50     String cmdArgs = conf.get(DataAnalysis.PARAMETERS);
51     String outputDir = conf.get(DataAnalysis.OUTPUT_DIR);
52     String workingDir = conf.get(DataAnalysis.WORKING_DIR);
53
54       String localInputFile = null;
55       String outFile = null;
56       String stdOutFile = null;
57       String stdErrFile = null;
58
59     System.out.println("the map key : " + key);
60     System.out.println("the value path : " + value.toString());
61     System.out.println("Local DB : " + this.localDB);
62
63         /*
64           Write your code to get localInputFile, outFile,
65           stdOutFile and stdErrFile
66         */
67
```

```
68
69     // Prepare the arguments to the executable
70     String execCommand = cmdArgs.replaceAll("#_INPUTFILE_#", localInputFile);
71     if (cmdArgs.indexOf("#_OUTPUTFILE_#") > -1) {
72       execCommand = execCommand.replaceAll("#_OUTPUTFILE_#", outFile);
73     } else {
74       outFile = stdOutFile;
75     }
76
77     endTime = Double.toString(((System.currentTimeMillis() - startTime) / 1000.0));
78     System.out.println("Before running the executable Finished in " + endTime + " seconds");
79
80     execCommand = this.localBlastProgram + File.separator + execName
81     + " " + execCommand + " -db " + this.localDB;
82     //Create the external process
83
84     startTime = System.currentTimeMillis();
85
86     Process p = Runtime.getRuntime().exec(execCommand);
87
88     OutputHandler inputStream = new OutputHandler(p.getInputStream(), "INPUT", stdOutFile);
89     OutputHandler errorStream = new OutputHandler(p.getErrorStream(), "ERROR", stdErrFile);
90
91     // start the stream threads.
92     inputStream.start();
93     errorStream.start();
94
95     p.waitFor();
96     //end time of this procress
97     endTime = Double.toString(((System.currentTimeMillis() - startTime) / 1000.0));
98     System.out.println("Program Finished in " + endTime + " seconds");
99
100     //Upload the results to HDFS
101     startTime = System.currentTimeMillis();
102
103     Path outputDirPath = new Path(outputDir);
104     Path outputFileName = new Path(outputDirPath, fileNameOnly);
105     fs.copyFromLocalFile(new Path(outFile), outputFileName);
106
107     endTime = Double.toString(((System.currentTimeMillis() - startTime) / 1000.0));
108     System.out.println("Upload Result Finished in " + endTime + " seconds");
109
110   }
111 }
```

In addition, if you need to understand the dataflow and main program, please look into the DataAnalysis.java.

```
1  /*file: DataAnalysis.java*/
2  package cgl.hadoop.apps.runner;
3
4  import org.apache.hadoop.conf.Configuration;
5  import org.apache.hadoop.conf.Configured;
6  import org.apache.hadoop.filecache.DistributedCache;
7  import org.apache.hadoop.fs.FileSystem;
8  import org.apache.hadoop.fs.Path;
9  import org.apache.hadoop.io.IntWritable;
10 import org.apache.hadoop.io.Text;
11 import org.apache.hadoop.mapreduce.Job;
12 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
14 import org.apache.hadoop.mapreduce.lib.output.SequenceFileOutputFormat;
15 import org.apache.hadoop.util.Tool;
16 import org.apache.hadoop.util.ToolRunner;
17
18 import cgl.hadoop.apps.DataFileInputFormat;
19 import java.net.URI;
20
```

```java
public class DataAnalysis extends Configured implements Tool {

   public static String WORKING_DIR = "working_dir";
   public static String OUTPUT_DIR = "out_dir";
   public static String EXECUTABLE = "exec_name";
   public static String PROGRAM_DIR = "exec_dir";
   public static String PARAMETERS = "params";
   public static String DB_NAME = "nr";
   public static String DB_ARCHIVE = "BlastDB.tar.gz";

/**
 * Launch the MapReduce computation.
 * This method first, remove any previous working directories and create a new one
 * Then the data (file names) is copied to this new directory and launch the
 * MapReduce (map-only though) computation.
 * @param numMapTasks - Number of map tasks.
 * @param numReduceTasks - Number of reduce tasks =0.
 * @param programDir - The directory where the Cap3 program is.
 * @param execName - Name of the executable.
 * @param dataDir - Directory where the data is located.
 * @param outputDir - Output directory to place the output.
 * @param cmdArgs - These are the command line arguments to the Cap3 program.
 * @throws Exception - Throws any exception occurs in this program.
 */
   void launch(int numReduceTasks, String programDir,
       String execName, String workingDir, String databaseArchive, String databaseName,
       String dataDir, String outputDir, String cmdArgs) throws Exception {

     Configuration conf = new Configuration();
     Job job = new Job(conf, execName);

     // First get the file system handler, delete any previous files, add the
     // files and write the data to it, then pass its name as a parameter to
     // job
     Path hdMainDir = new Path(outputDir);
     FileSystem fs = FileSystem.get(conf);
     fs.delete(hdMainDir, true);

     Path hdOutDir = new Path(hdMainDir, "out");

     // Starting the data analysis.
     Configuration jc = job.getConfiguration();

     jc.set(WORKING_DIR, workingDir);
     jc.set(EXECUTABLE, execName);
     jc.set(PROGRAM_DIR, programDir); // this the name of the executable archive
     jc.set(DB_ARCHIVE, databaseArchive);
     jc.set(DB_NAME, databaseName);
     jc.set(PARAMETERS, cmdArgs);
     jc.set(OUTPUT_DIR, outputDir);


     long startTime = System.currentTimeMillis();
     DistributedCache.addCacheArchive(new URI(programDir), jc);
     System.out.println("Add Distributed Cache in "
         + (System.currentTimeMillis() - startTime) / 1000.0
         + " seconds");


     FileInputFormat.setInputPaths(job, dataDir);
     FileOutputFormat.setOutputPath(job, hdOutDir);

     job.setJarByClass(DataAnalysis.class);
     job.setMapperClass(RunnerMap.class);
     job.setOutputKeyClass(IntWritable.class);
     job.setOutputValueClass(Text.class);
```

```java
88        job.setInputFormatClass(DataFileInputFormat.class);
89        job.setOutputFormatClass(SequenceFileOutputFormat.class);
90        job.setNumReduceTasks(numReduceTasks);
91
92        startTime = System.currentTimeMillis();
93
94        int exitStatus = job.waitForCompletion(true) ? 0 : 1;
95        System.out.println("Job Finished in "
96            + (System.currentTimeMillis() - startTime) / 1000.0
97            + " seconds");
98        //clean the cache
99
100
101        System.exit(exitStatus);
102    }
103
104    public int run(String[] args) throws Exception {
105        if (args.length < 8) {
106            System.err.println("Usage: DataAnalysis <Executable and Database Archive on HDFS>
107             <Executable> <Working_Dir> <Database dir under archive> <Database name>
108             <HDFS_Input_dir> <HDFS_Output_dir> <Cmd_args>");
109            ToolRunner.printGenericCommandUsage(System.err);
110            return -1;
111        }
112        String programDir = args[0];
113        String execName = args[1];
114        String workingDir = args[2];
115        String databaseArchive = args[3];
116        String databaseName = args[4];
117        String inputDir = args[5];
118        String outputDir = args[6];
119        //"#_INPUTFILE_# -p 95 -o 49 -t 100"
120        String cmdArgs = args[7] ;
121
122        int numReduceTasks = 0;// We don't need reduce here.
123
124        launch(numReduceTasks, programDir, execName, workingDir ,
125            databaseArchive, databaseName,inputDir, outputDir, cmdArgs);
126        return 0;
127    }
128
129    public static void main(String[] argv) throws Exception {
130        int res = ToolRunner.run(new Configuration(), new DataAnalysis(), argv);
131        System.exit(res);
132    }
133 }
```

# Edit

The sketch code is stored within the provided VirtualBox image. Use linux text editor vi/vim to add your code.

```
1 $ cd /root/MoocHomeworks/HadoopBlast/
2 $ vim src/cgl/hadoop/apps/runner/RunnerMap.java
```

## Compile and run your code

Use the same one-click script compileAndExecHadoopBlast.sh as in prior homework. Standard error messages such as compile errors, execution errors, etc. will be redirected on the screen. Follow the same debugging format.

```
1 $ cd /root/MoocHomeworks/HadoopBlast/
2 $ ./compileAndExecHadoopBlast.sh
```

# View the result

The result is generated at /root/MoocHomeworks/HadoopBlast/output/HDFS_blast_output . There should be 4 output FASTA files with .fa extension

```
$ cd /root/MoocHomeworks/HadoopBlast/output/HDFS_blast_output
$ ls
```