# Introduction to Big Data

## ECAP456

**Dr. Rajni Bhalla**

**Associate Professor**
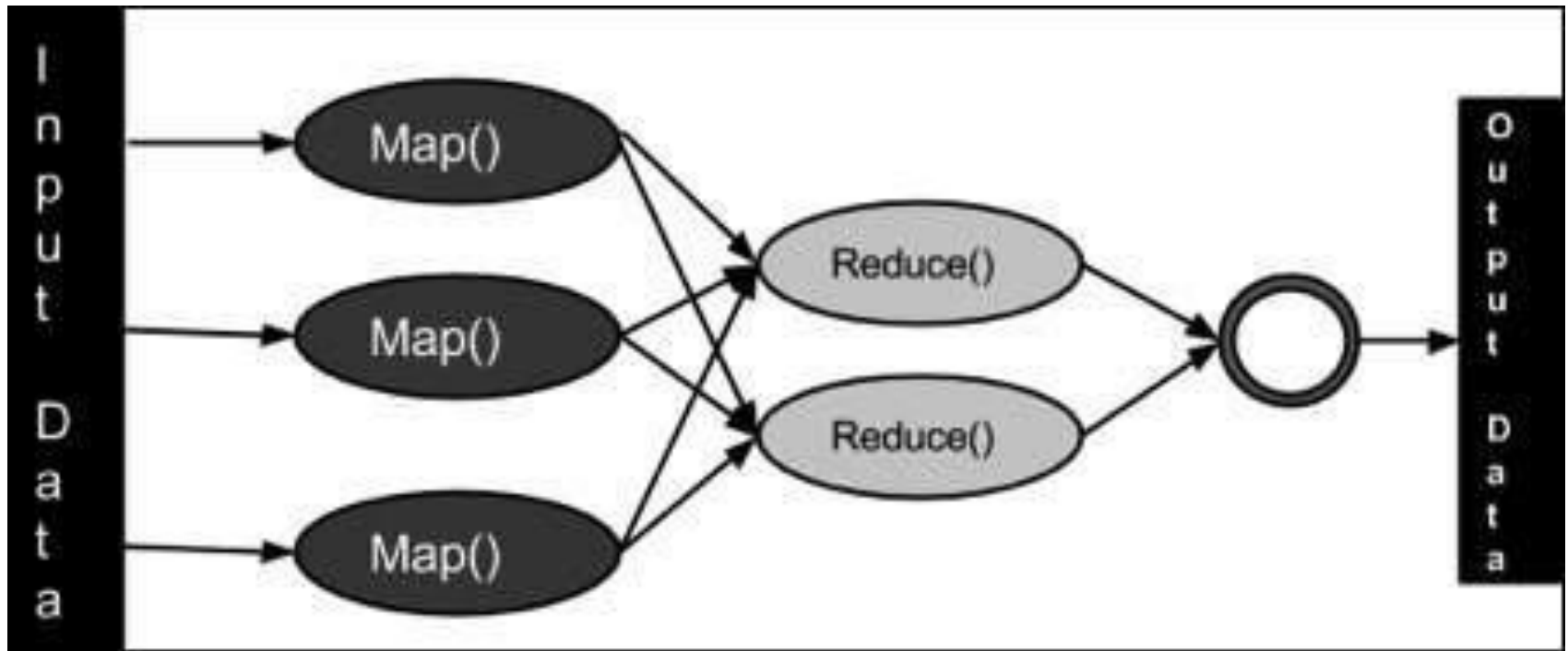
# Learning Outcomes

After this lecture, you will be able to
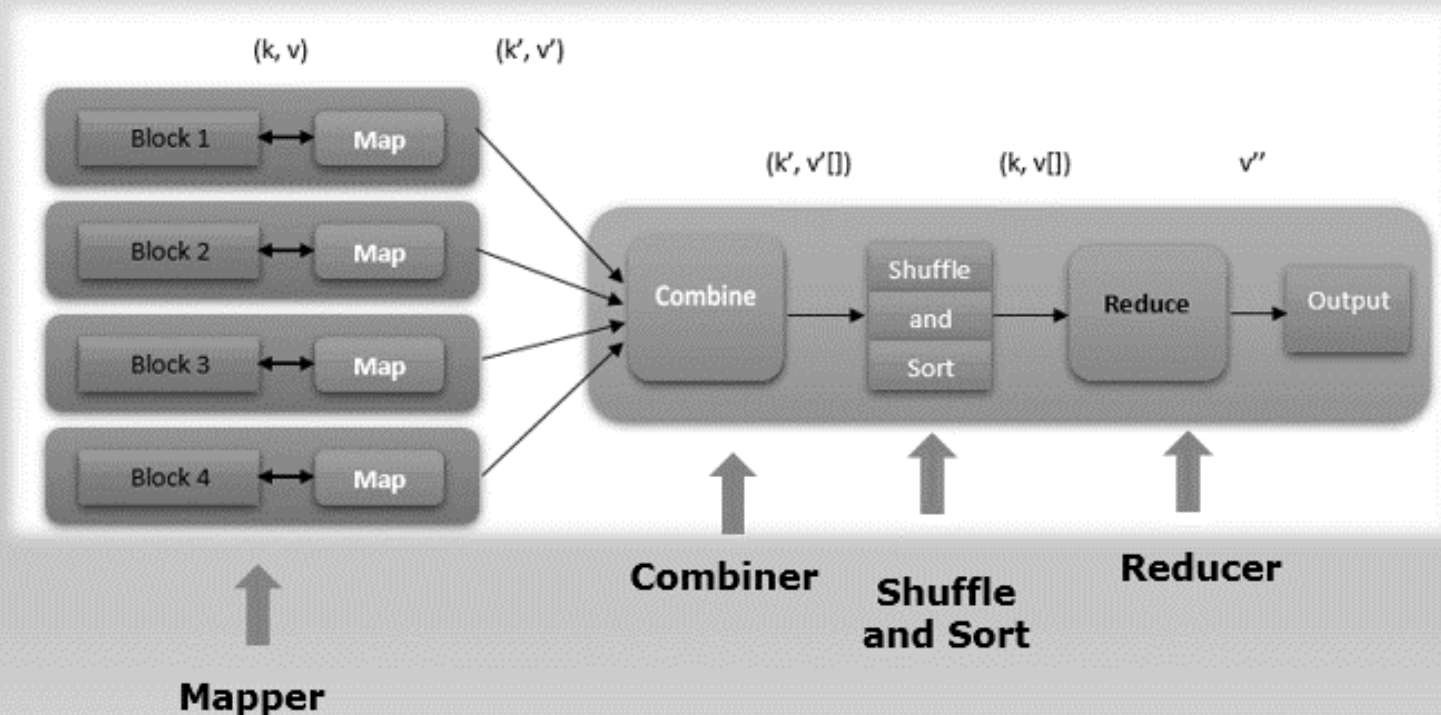
- learn anatomy of a MapReduce Job Run

# Introduction

# Introduction

# Introduction



Job Execution and job Completion

# Introduction

```
// Learning MapReduce by Nitesh Jain
import org.apache.hadoop.fs.Path;

public class WordCountWithCombiner extends Configured
implements Tool{

    @Override
    public int run(String[] args) throws Exception {
        Configuration conf = getConf();

        Job job = new Job(conf, "MyJob");

        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count With Combiners");

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setCombinerClass(WordCountReducer.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        return job.waitForCompletion(true) ? 0 : 1;

    }

    public static void main(String[] args) throws Exception {
        int exitCode = ToolRunner.run(new Configuration(), new WordCountWithCombiner(), args);
        System.exit(exitCode);
    }
}
```

**Last statement in run method**

# Introduction



```java
WordCountWithCombiner.java        WordCount.java ☒

  // Learning MapReduce by Nitesh Jain
 import org.apache.hadoop.fs.Path;

  public class WordCount {

    public static void main(String[] args) throws Exception {
      if (args.length != 2) {
        System.err.println("Usage: WordCount <input path> <output path>");
        System.exit(-1);
      }

      Job job = new Job();
      job.setJarByClass(WordCount.class);
      job.setJobName("Word Count");

      FileInputFormat.addInputPath(job, new Path(args[0]));
      FileOutputFormat.setOutputPath(job, new Path(args[1]));

      job.setMapperClass(WordCountMapper.class);
      job.setReducerClass(WordCountReducer.class);

      job.setOutputKeyClass(Text.class);
      job.setOutputValueClass(IntWritable.class);

      System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
  }
```

Last line in main method

# What happened behind the scene?

- **Jobs gets submitted by WaitforCompletion().**

- Properties that decide the job run:

  - "mapred.job.tracker" (for 0.22 release)

    - This is set in mapred-site.xml

    - Default is local.

    - Set as host:part

  - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

    - Can be set to local,classic or yarn

# What happened behind the scene?

- Jobs gets submitted by WaitforCompletion().

- Properties that decide the job run:

  - "mapred.job.tracker" (for 0.22 release)

    - This is set in mapred-site.xml

    - Default is local.

    - Set as host:part

  - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

    - Can be set to local,classic or yarn

# What happened behind the scene?

- Jobs gets submitted by WaitforCompletion().

- Properties that decide the job run:

  - "mapred.job.tracker" (for 0.22 release)

    - This is set in mapred-site.xml

    - Default is local.

    - Set as host:part

  - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

    - Can be set to local,classic or yarn

# What happened behind the scene?

- Jobs gets submitted by WaitforCompletion().

- Properties that decide the job run:

  - "mapred.job.tracker" (for 0.22 release)

    - This is set in mapred-site.xml

    - Default is local.

    - Set as host:part

  - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

    - Can be set to local,classic or yarn

# What happened behind the scene?

- Jobs gets submitted by WaitforCompletion().

- Properties that decide the job run:

  - "mapred.job.tracker" (for 0.22 release)

    - This is set in mapred-site.xml

    - Default is local.

    - Set as host:part

  - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

    - Can be set to local,classic or yarn

# What happened behind the scene?

- Jobs gets submitted by WaitforCompletion().

- Properties that decide the job run:

  - "mapred.job.tracker" (for 0.22 release)

    - This is set in mapred-site.xml

    - Default is local.

    - Set as host:part

  - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

    - Can be set to local,classic or yarn

# What happened behind the scene?

- Jobs gets submitted by WaitforCompletion().

- Properties that decide the job run:

    - "mapred.job.tracker" (for 0.22 release)

        - This is set in mapred-site.xml

        - Default is local.

        - Set as host:part

    - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

        - Can be set to local,classic or yarn

# What happened behind the scene?

- Jobs gets submitted by WaitforCompletion().

- Properties that decide the job run:

  - "mapred.job.tracker" (for 0.22 release)

    - This is set in mapred-site.xml

    - Default is local.

    - Set as host:part

  - "mapreduce.framework.name" (for .23 and later(eg 1.X.X)):

    - Can be set to local,classic or yarn

# Classic MapReduce – Overview

# Classic MapReduce – Overview

**Job Tracker**

Job Tracker node

# Classic MapReduce – Overview

# Classic MapReduce – Overview



1. Client submits the MapReduce job.

# Classic MapReduce – Overview



1. **Client submits the MapReduce job.**

2. Job tracker co-ordinate the job run

# Classic MapReduce – Overview



1. **Client submits the MapReduce job.**

2. Job tracker co-ordinate the job run

3. Task tracker run the map and reduce tasks.

# Classic MapReduce – Overview



1. **Client submits the MapReduce job.**

2. Job tracker co-ordinate the job run

3. Task tracker run the map and reduce tasks.

4. Task tracker sends progress status

# Classic MapReduce – Overview



**1. Client submits the MapReduce job.**

2. Job tracker co-ordinate the job run

3. Task tracker run the map and reduce tasks.

4. Task tracker sends progress status

5. Job Tracker sends progress status and completion message to client.

# Components of MapReduce

1. Client

2. Yarn node manager

3. Yarn resource manager

4. MapReduce application master

5. Distributed Filesystem

# Components of MapReduce

1. Client  2. Yarn node manager  3. Yarn resource manager

4. MapReduce application master 5. Distributed Filesystem

# How to submit Job?

- **use the submit()**

- submitJobInternal()

- waitForCompletion()

-  reports the progress to the console.

- job counters are displayed

# How to submit Job?

- use the submit()

- submitJobInternal()

- waitForCompletion()

-  reports the progress to the console.

- job counters are displayed

# How to submit Job?

- use the submit()

- submitJobInternal()

- waitForCompletion()

- reports the progress to the console.

- job counters are displayed

# How to submit Job?

- use the submit()

- submitJobInternal()

- waitForCompletion()

- reports the progress to the console.

- job counters are displayed

# How to submit Job?

- use the submit()

- submitJobInternal()

- waitForCompletion()

- reports the progress to the console.

- job counters are displayed

Processes implemented by **JobSubmitter** for submitting the Job :

1. The resource manager askes for a new application ID that is used for MapReduce Job ID.

2. Output specification of the job is checked.

3. If the splits cannot be computed, it computes the input splits for the job

4. Resources needed to run the job is copied

Processes implemented by JobSubmitter for submitting the Job :

1. The resource manager askes for a new application ID that is used for MapReduce Job ID.

2. Output specification of the job is checked.

3. If the splits cannot be computed, it computes the input splits for the job

4. Resources needed to run the job is copied

Processes implemented by **JobSubmitter** for submitting the Job :

1. The resource manager askes for a new application ID that is used for MapReduce Job ID.

2. Output specification of the job is checked.

3. If the splits cannot be computed, it computes the input splits for the job

4. Resources needed to run the job is copied

# Processes implemented by JobSubmitter for submitting the Job :

Processes implemented by **JobSubmitter** for submitting the Job :

1.  The resource manager askes for a new application ID that is used for MapReduce Job ID.

2.  Output specification of the job is checked.

3.  If the splits cannot be computed, it computes the input splits for the job

4.  Resources needed to run the job is copied

Processes implemented by **JobSubmitter** for submitting the Job :

5. It copies job **JAR** with a high replication factor, which is controlled by **mapreduce.client.submit.file.replication property.**

6. By calling submitApplication(), submits the job on the resource manager.

Processes implemented by **JobSubmitter** for submitting the Job :

5. It copies job **JAR** with a high replication factor, which is controlled by **mapreduce.client.submit.file.replication property.**

6. By calling **submitApplication(),** submits the job on the resource manager.

That's all for now...