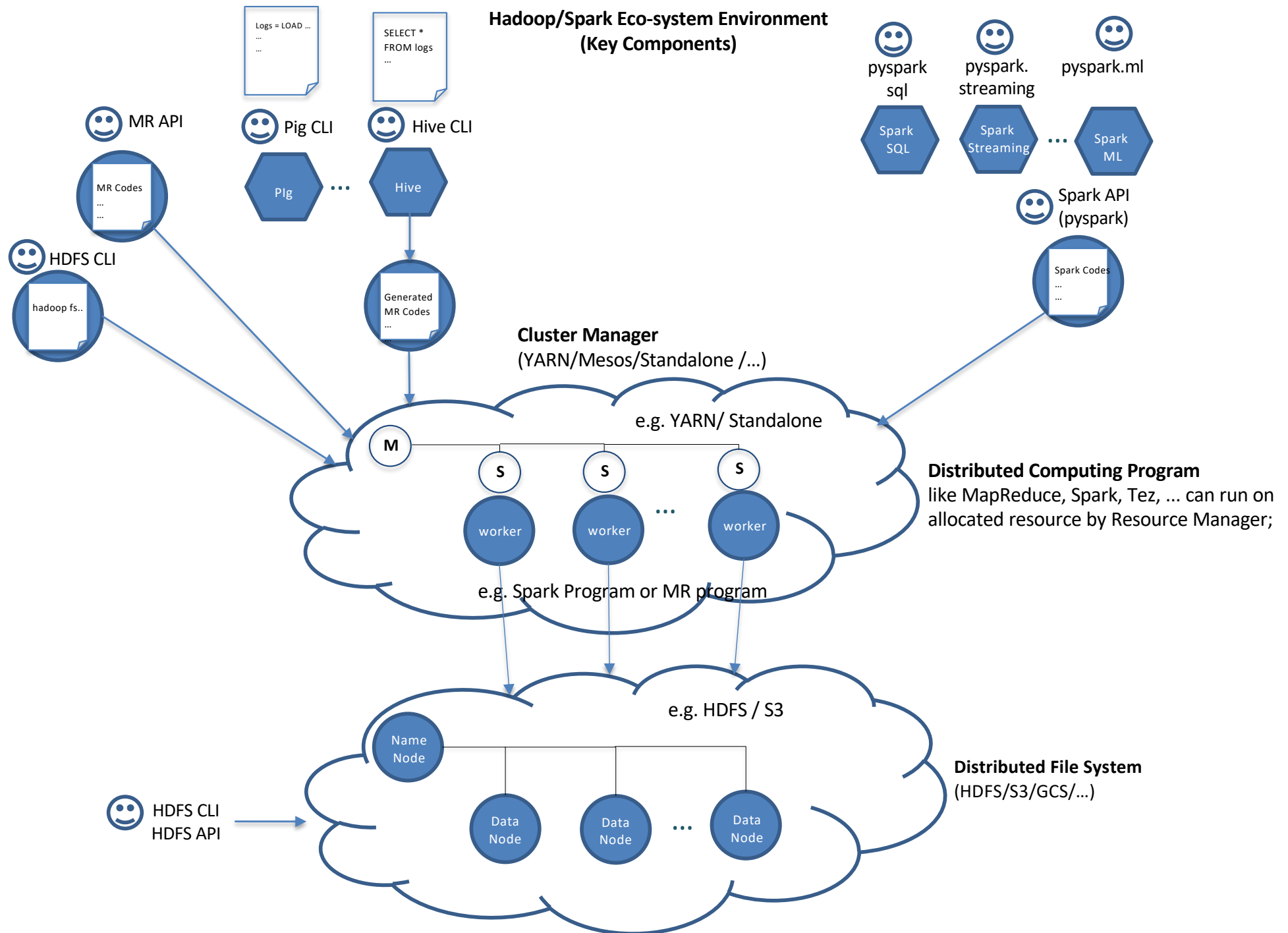
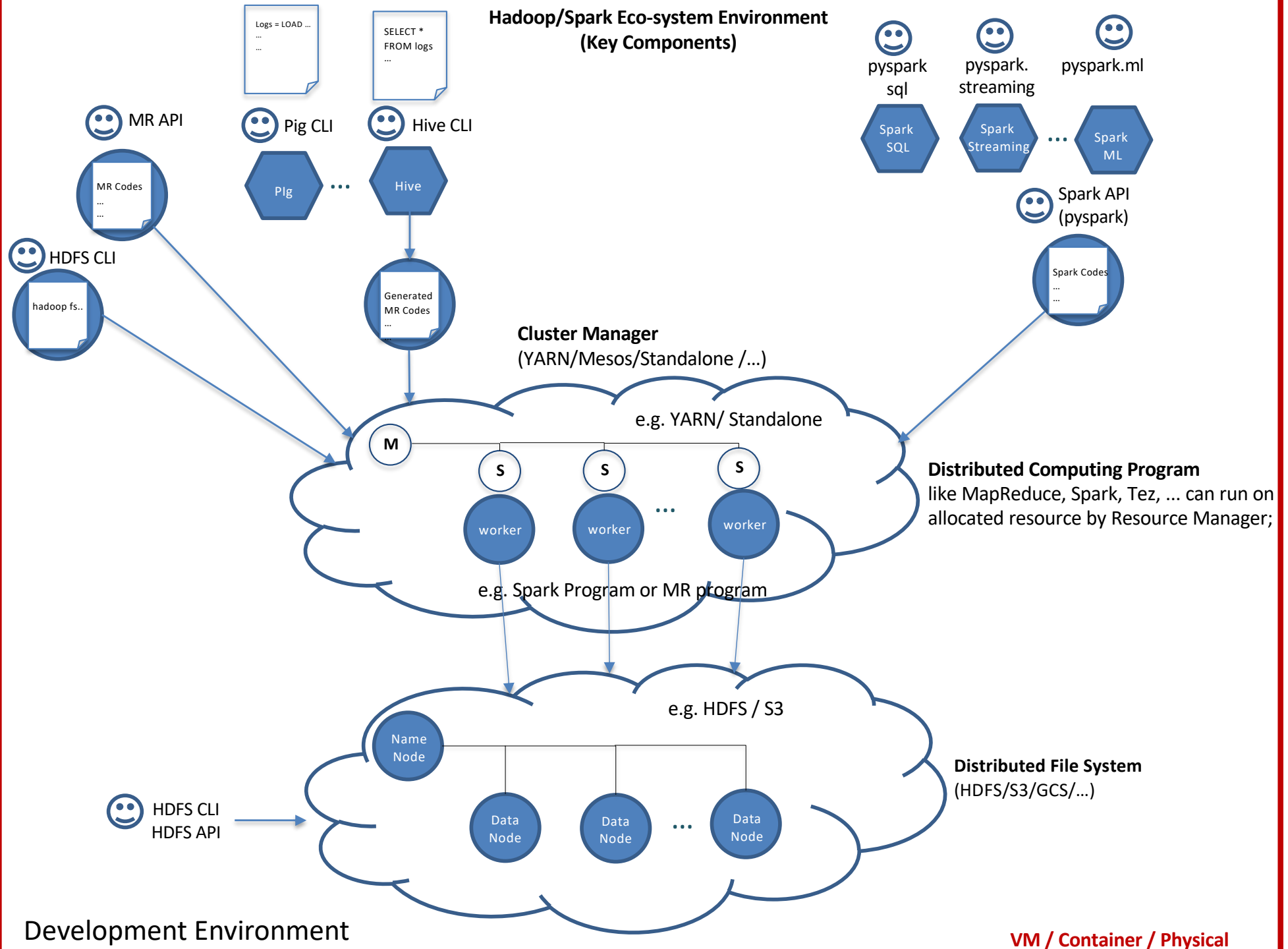


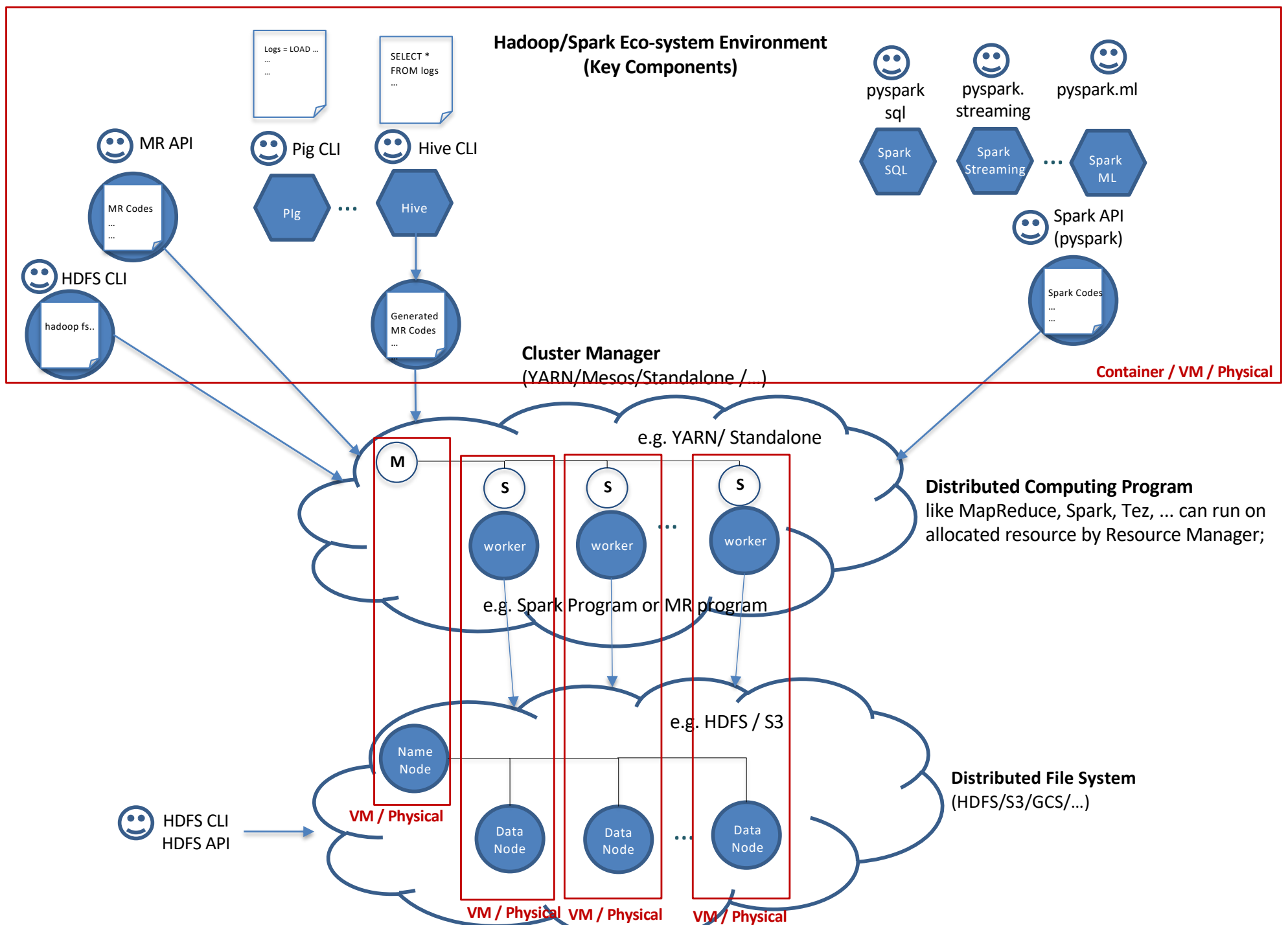
# Hadoop/Spark Eco-system Overview

## Hadoop/Spark Eco-system Environment (Key Components)



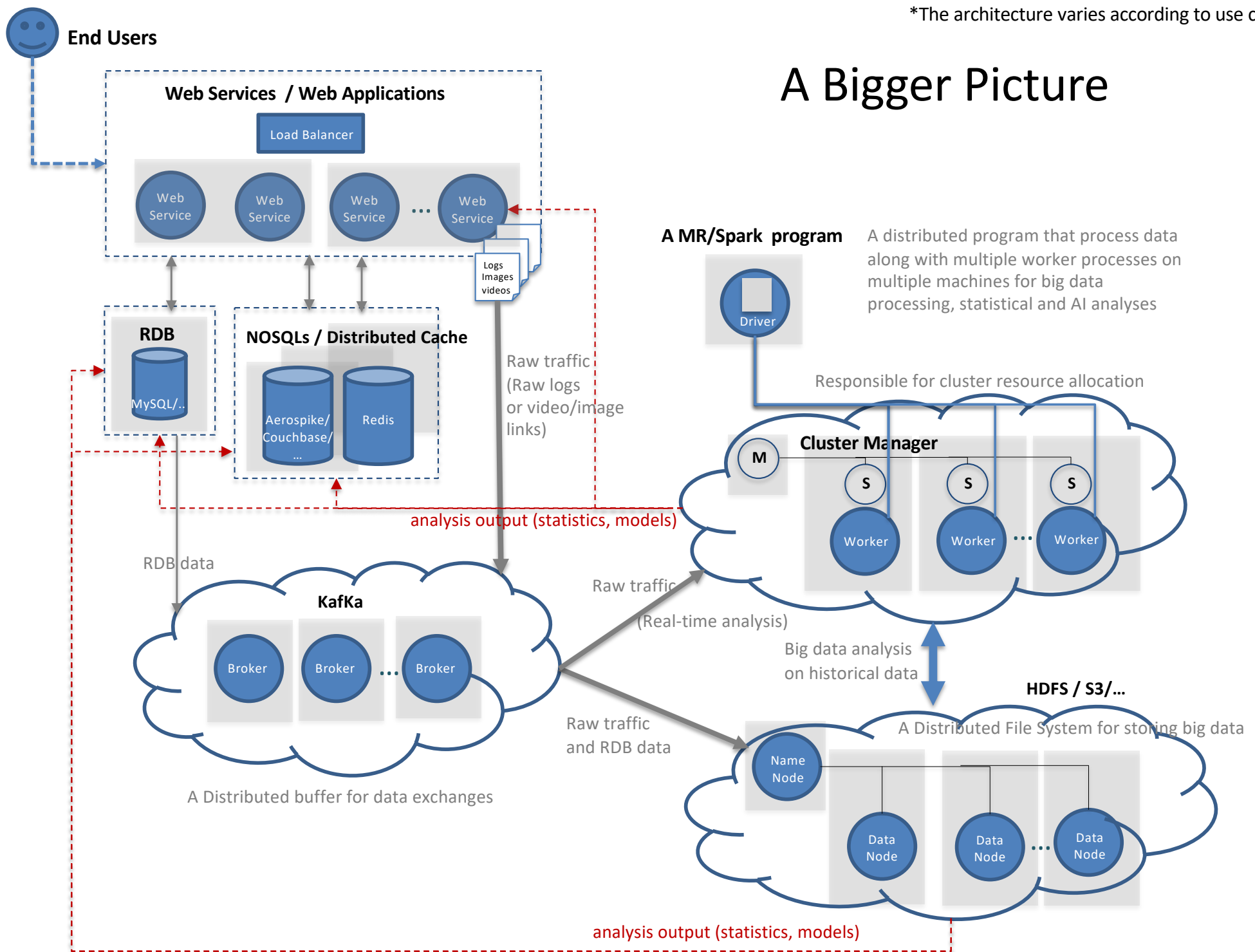
## Hadoop/Spark Eco-system Environment (Key Components)





Production Environment

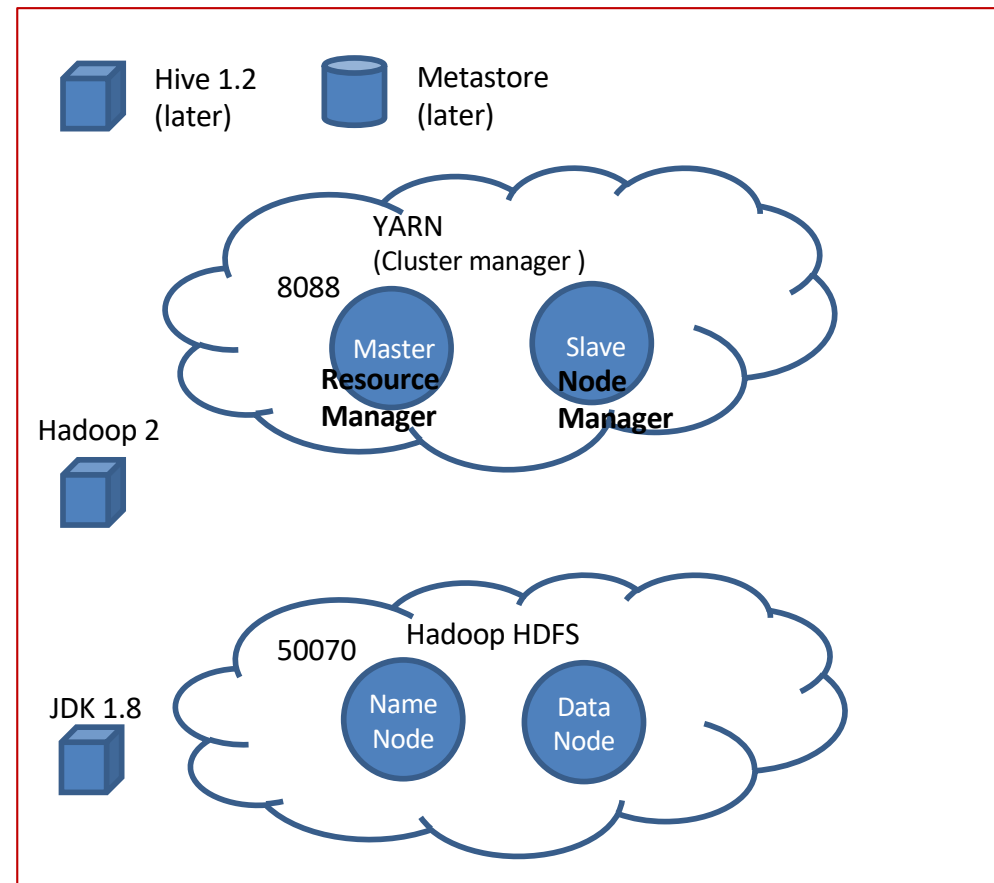
# A Bigger Picture



# Demo

## Development Environment Setup

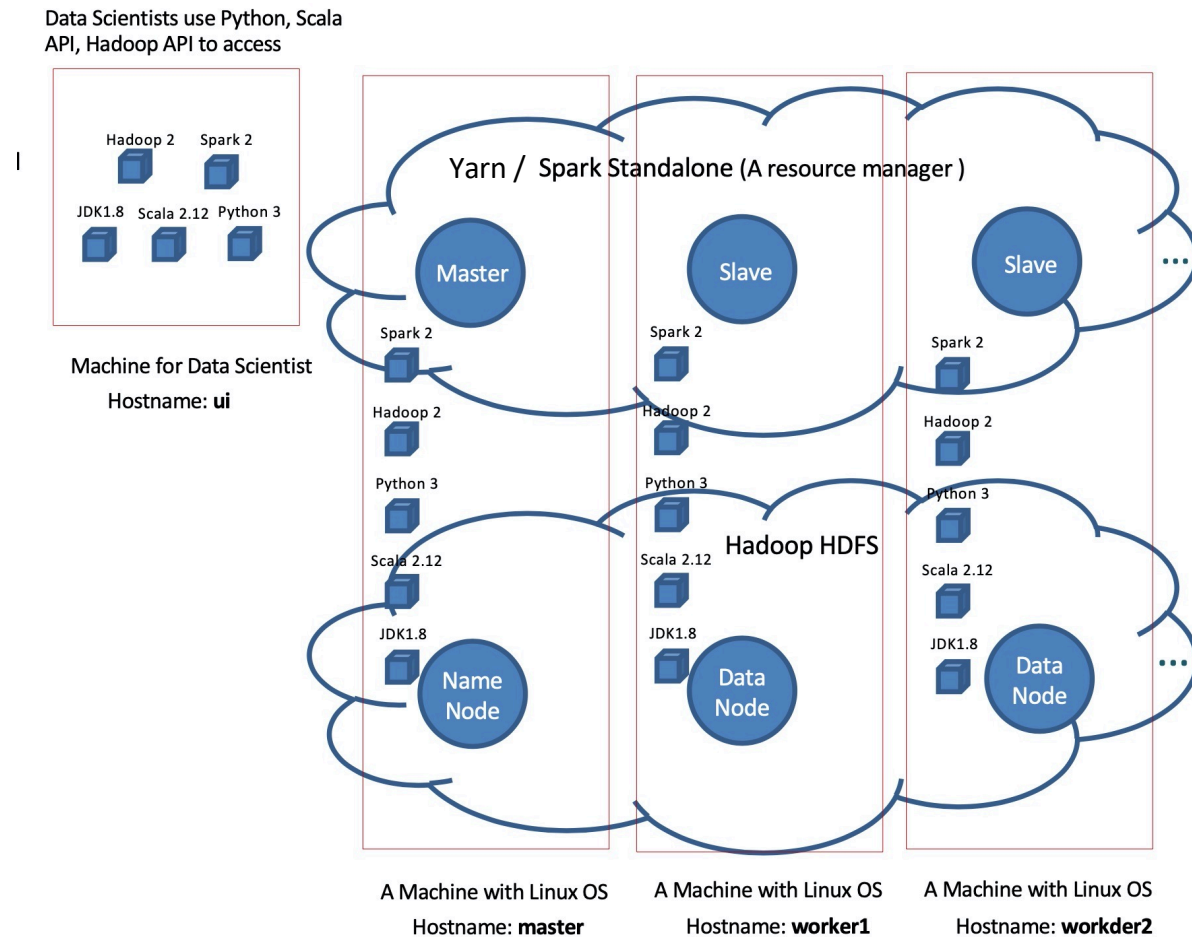
### Development Environment Setup



A Machine with Linux OS

# Demo (later on)

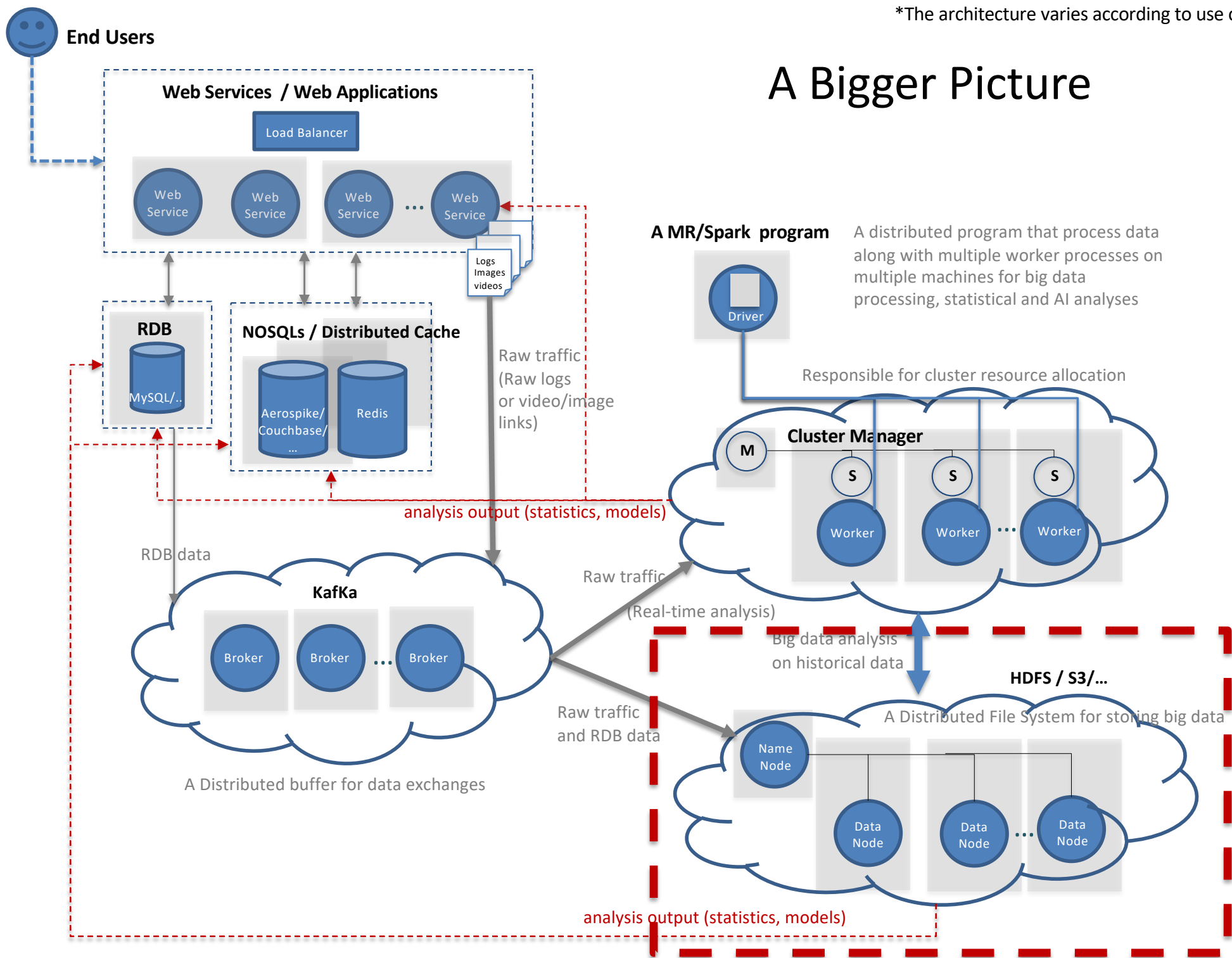
## Production Environment Setup



# Introduction to HDFS



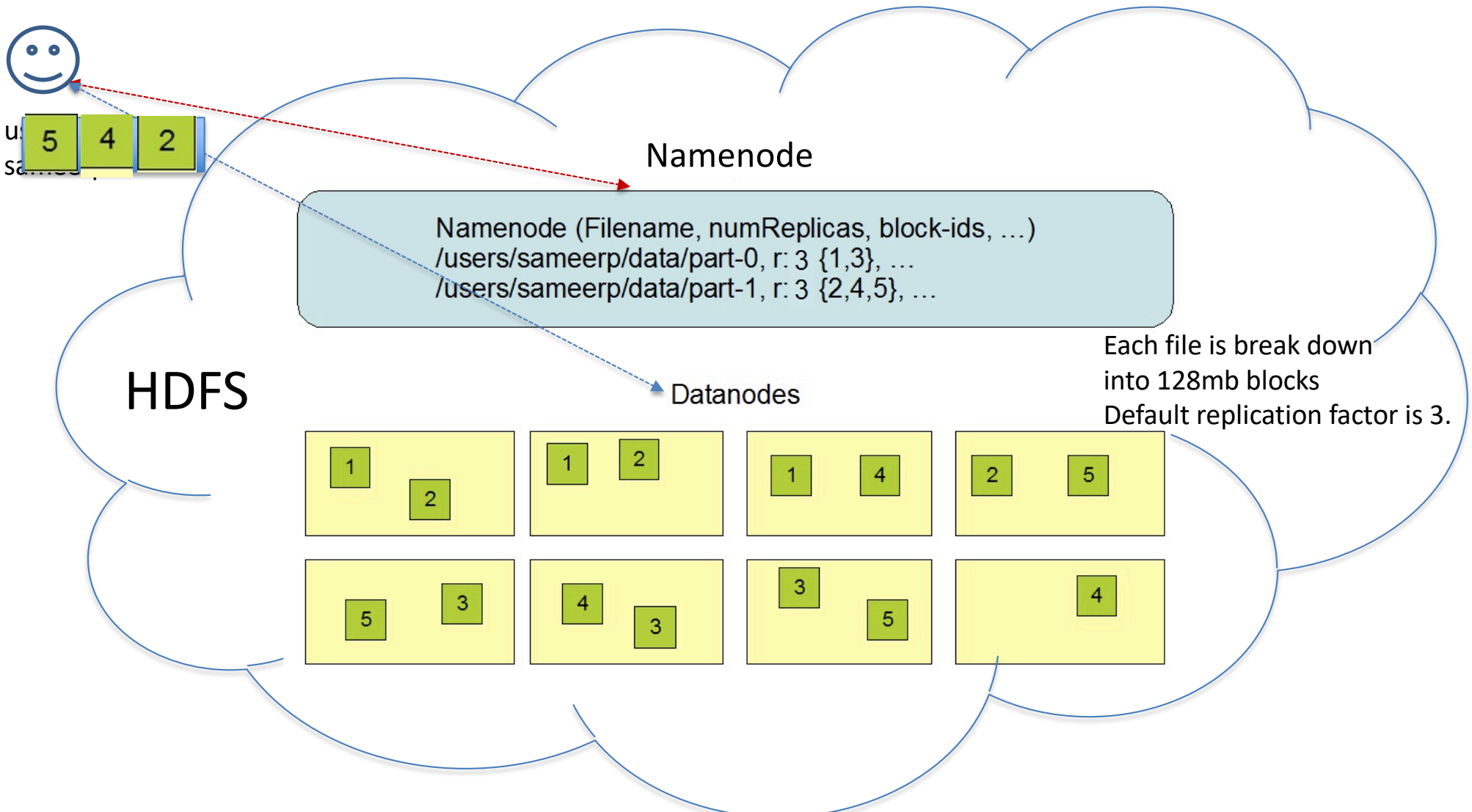
# A Bigger Picture



# HDFS/S3/GCS Architecture

```
hadoop fs -mkdir /user/sameerp/data  
hadoop fs -put part-0 /user/sameerp/data  
hadoop fs -put part-1 /user/sameerp/data
```

HDFS CLI

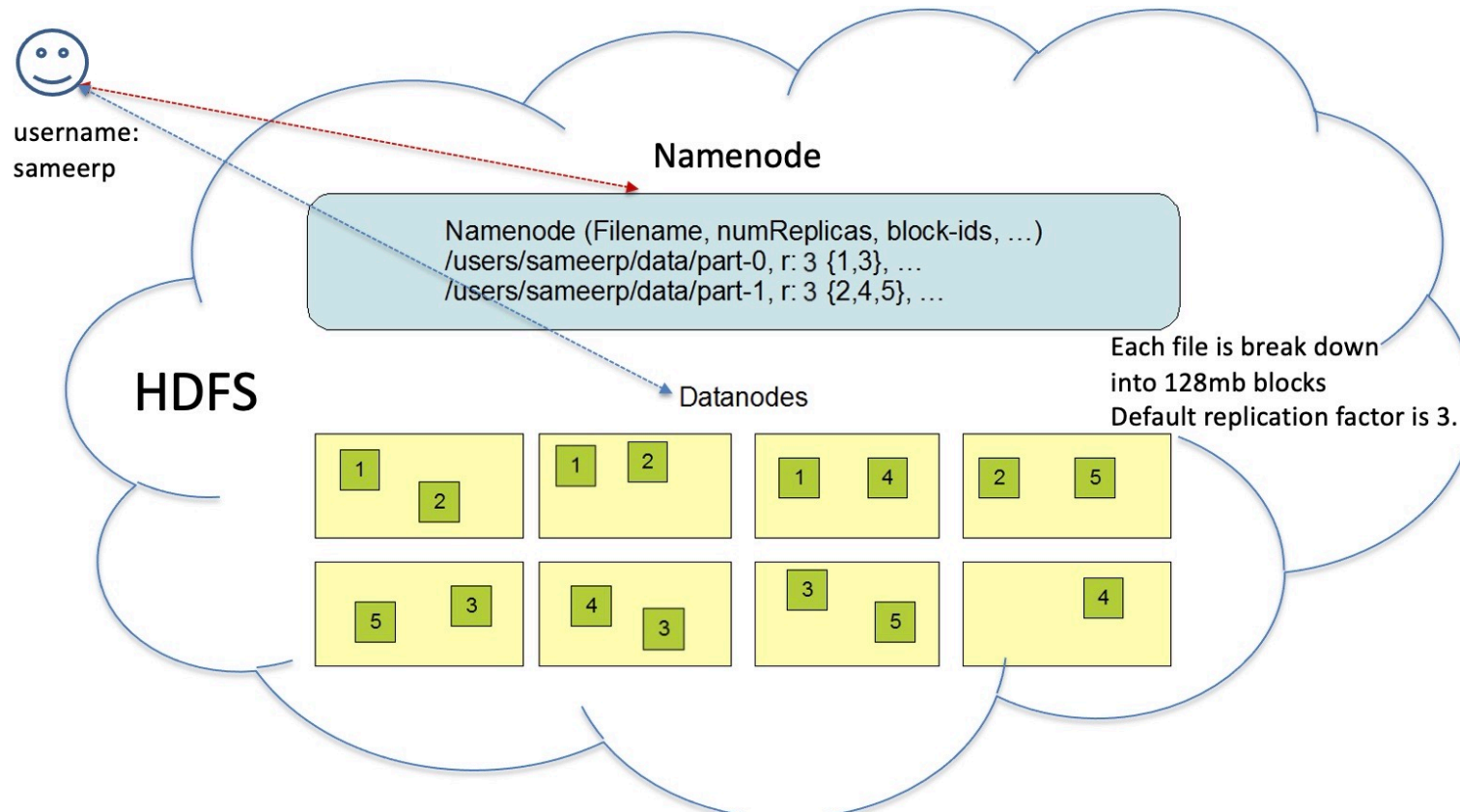


# HDFS Architecture

- Present a **single file system** for the end user though it consists of multiple machines under the hood (s3, gcs,..)
- Used to **store small to very large files**
- Files are broken up into blocks stored on Data Nodes
  - Typically, **128 MB** each block
  - Each block has **3 copies** by default
  - **Configurable by clients**
- **Commodity hardware**
  - Files are replicated to handle hardware failure
  - Detect failures and recovers from them

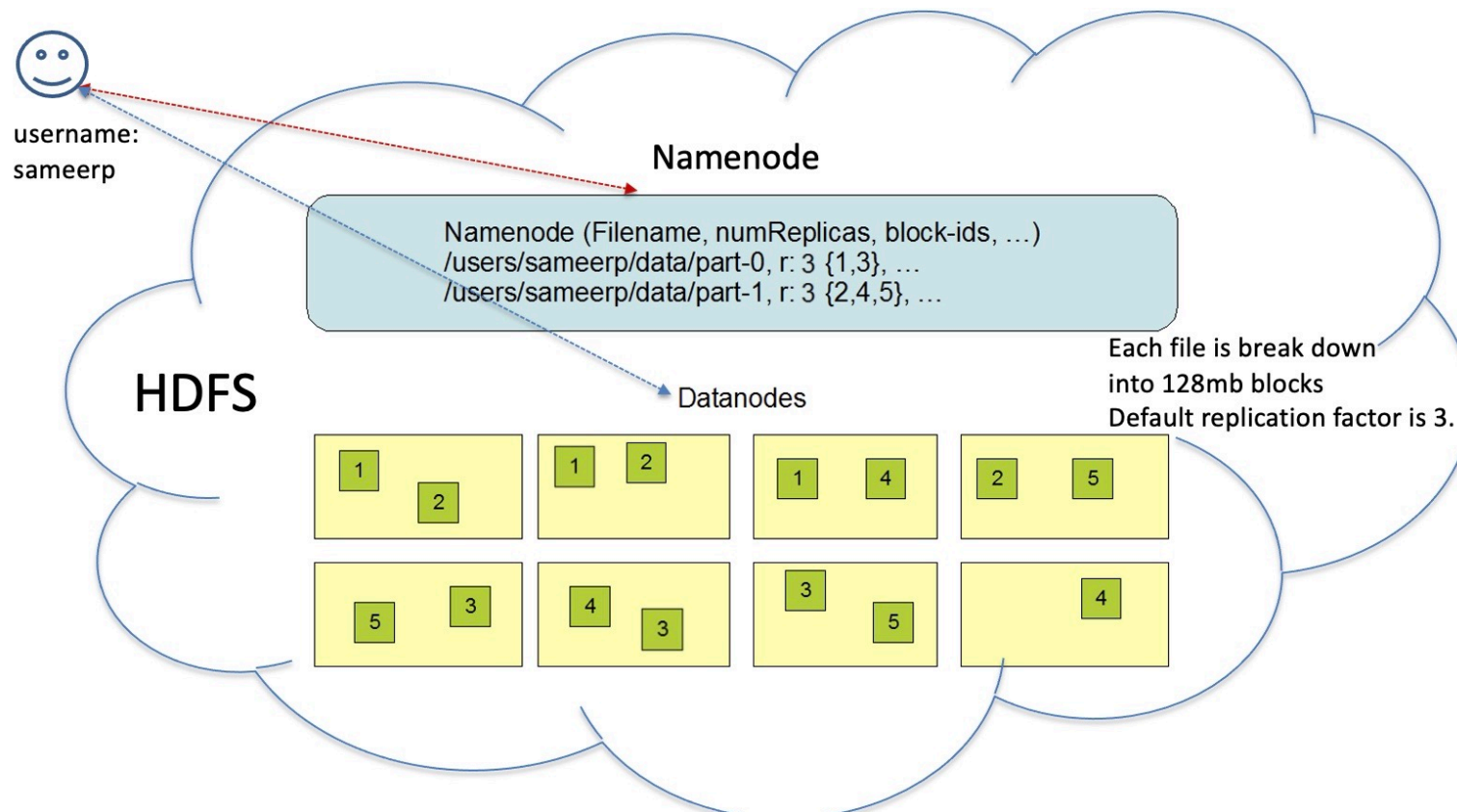
# NameNode

- Manages File System Namespace
  - Maps a file name to a set of blocks
  - Maps a block to the DataNodes where it resides



# Name Node Metadata

- Meta-data in Memory
  - The entire metadata is in main memory



# NameNode Metadata (cont.)

- A Transaction Log (called *EditLog*)
  - Records file creations, file deletions. Etc
- *FsImage*
  - The entire namespace, mapping of blocks to files and file system properties are stored in a file called *FsImage*.
  - NameNode can be configured to maintain multiple copies of *FsImage* and *EditLog*.

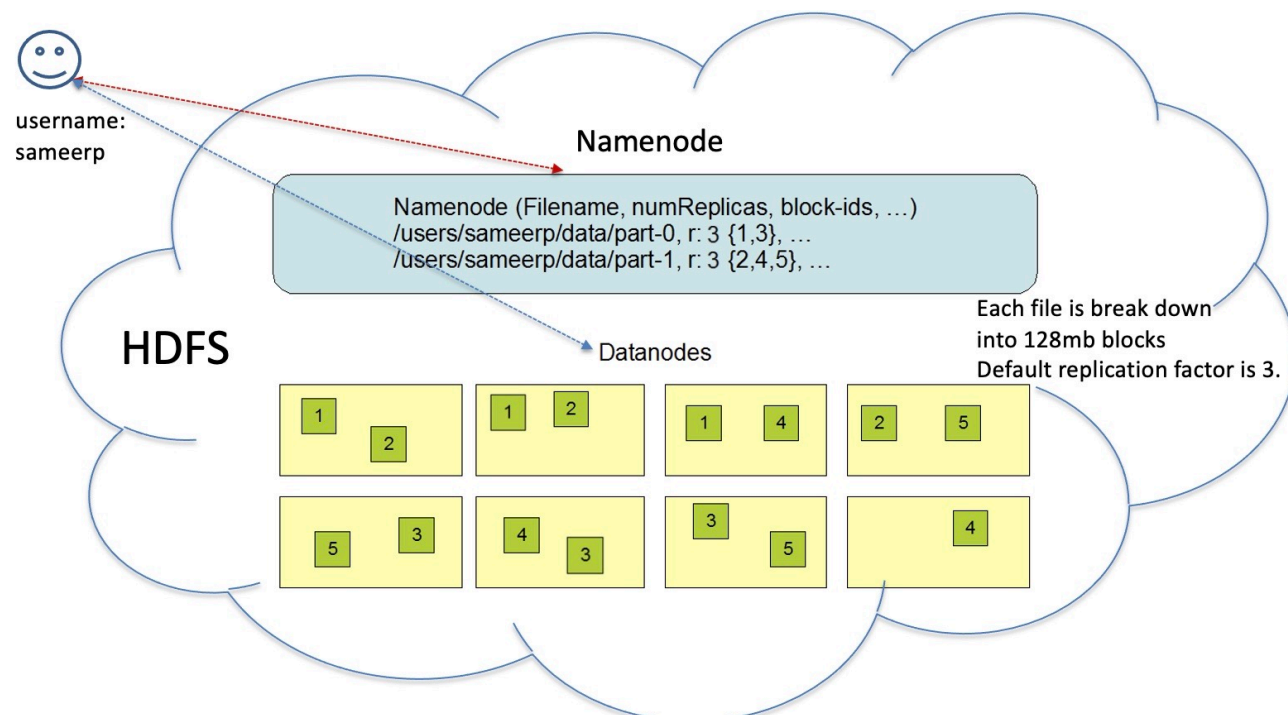
# Secondary NameNode

- Copies *FsImage* and *EditLog* from NameNode to a temporary directory
- Merges *FsImage* and *EditLog* into a new *FsImage* in temporary directory.
- Uploads new *FsImage* to the NameNode
  - Transaction Log on NameNode is purged



# DataNode

- A Block Server
  - Stores data in the local file system (e.g. ext3)
  - Serves data and meta-data to Clients
- Block Report
  - Periodically sends a report of all existing blocks to the NameNode





# HDFS User Interface

- The command line interface
  - `hadoop fs -mkdir /foodir`
  - `hadoop fs -cat /foodir/myfile.txt`
  - `hadoop fs -rm /foodir myfile.txt`

`hdfs dfs -mkdir /foodir`

- APIs in JAVA, Python, etc.

# HDFS User Interface

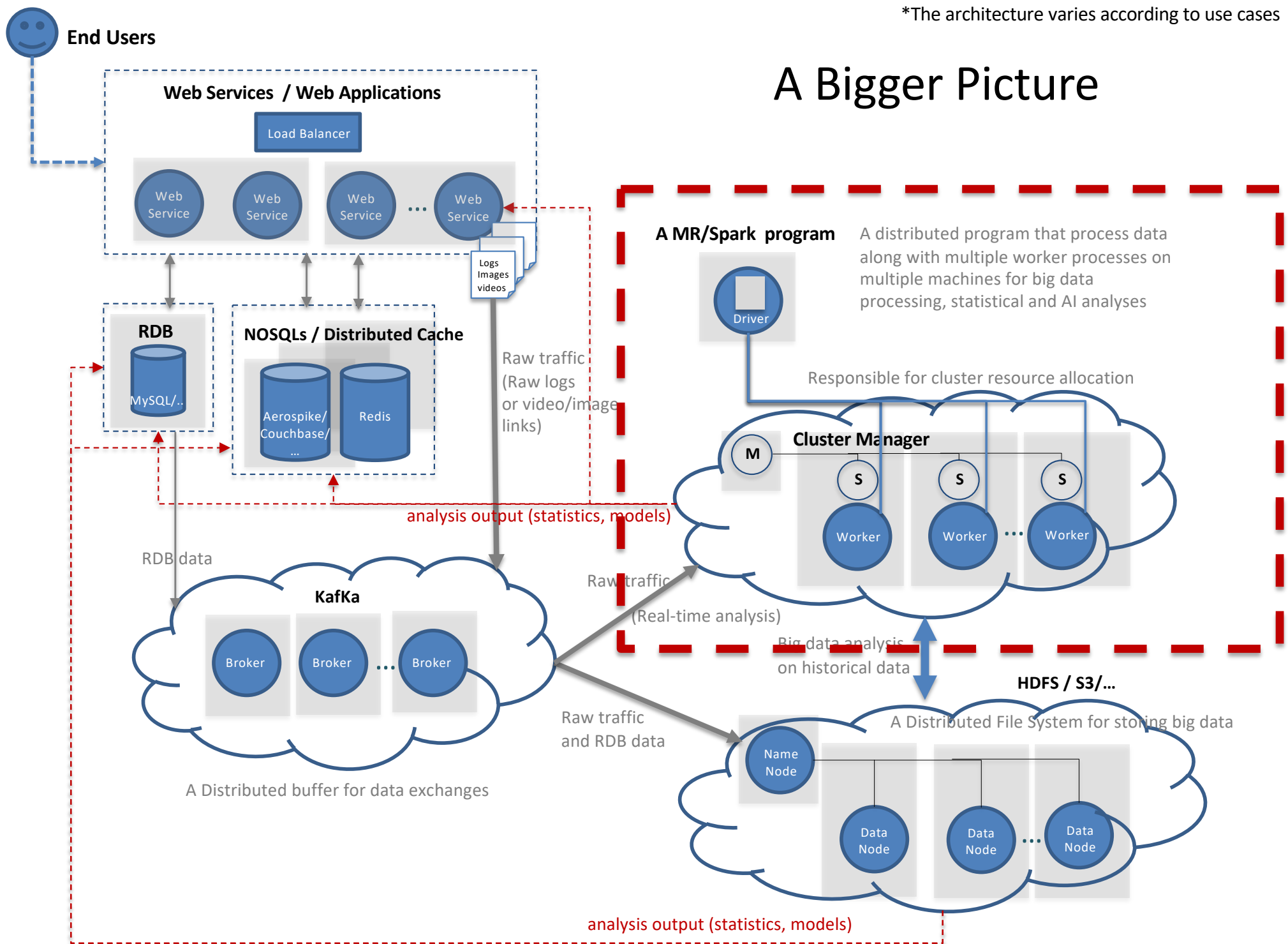
```
spark@devenv: ~$ hadoop fs
Usage: hadoop fs [generic options]
    [-appendToFile <localsrc> ... <dst>]
    [-cat [-ignoreCrc] <src> ...]
    [-checksum <src> ...]
    [-chgrp [-R] GROUP PATH...]
    [-chmod [-R] <MODE[,MODE]... | OCTALMODE> PATH...]
    [-chown [-R] [OWNER][:[GROUP]] PATH...]
    [-copyFromLocal [-f] [-p] [-l] [-d] <localsrc> ... <dst>]
    [-copyToLocal [-f] [-p] [-ignoreCrc] [-crc] <src> ... <localdst>]
    [-count [-q] [-h] [-v] [-t <storage type>]] [-u] [-x] <path> ...]
    [-cp [-f] [-p | -p[topax]] [-d] <src> ... <dst>]
    [-createSnapshot <snapshotDir> [<snapshotName>]]
    [-deleteSnapshot <snapshotDir> <snapshotName>]
    [-df [-h] [<path> ...]]
    [-du [-s] [-h] [-x] <path> ...]
    [-expunge]
    [-find <path> ... <expression> ...]
    [-get [-f] [-p] [-ignoreCrc] [-crc] <src> ... <localdst>]
    [-getfacl [-R] <path>]
    [-getfattr [-R] {-n name | -d} [-e en] <path>]
    [-getmerge [-nl] [-skip-empty-file] <src> <localdst>]
    [-help [cmd ...]]
    [-ls [-C] [-d] [-h] [-q] [-R] [-t] [-S] [-r] [-u] [<path> ...]]
    [-mkdir [-p] <path> ...]
    [-moveFromLocal <localsrc> ... <dst>]
    [-moveToLocal <src> <localdst>]
    [-mv <src> ... <dst>]
    [-put [-f] [-p] [-l] [-d] <localsrc> ... <dst>]
    [-renameSnapshot <snapshotDir> <oldName> <newName>]
    [-rm [-f] [-r|-R] [-skipTrash] [-safely] <src> ...]
    [-rmdir [--ignore-fail-on-non-empty] <dir> ...]
    [-setfacl [-R] [{-b|-k} {-m|-x <acl_spec>} <path>]|[--set <acl_spec> <path>]]
    [-setfattr {-n name [-v value] | -x name} <path>]
    [-setrep [-R] [-w] <rep> <path> ...]
    [-stat [format] <path> ...]
    [-tail [-f] <file>]
    [-test [-defsz] <path>]
    [-text [-ignoreCrc] <src> ...]
    [-touchz <path> ...]
    [-truncate [-w] <length> <path> ...]
    [-usage [cmd ...]]
```

# Demo

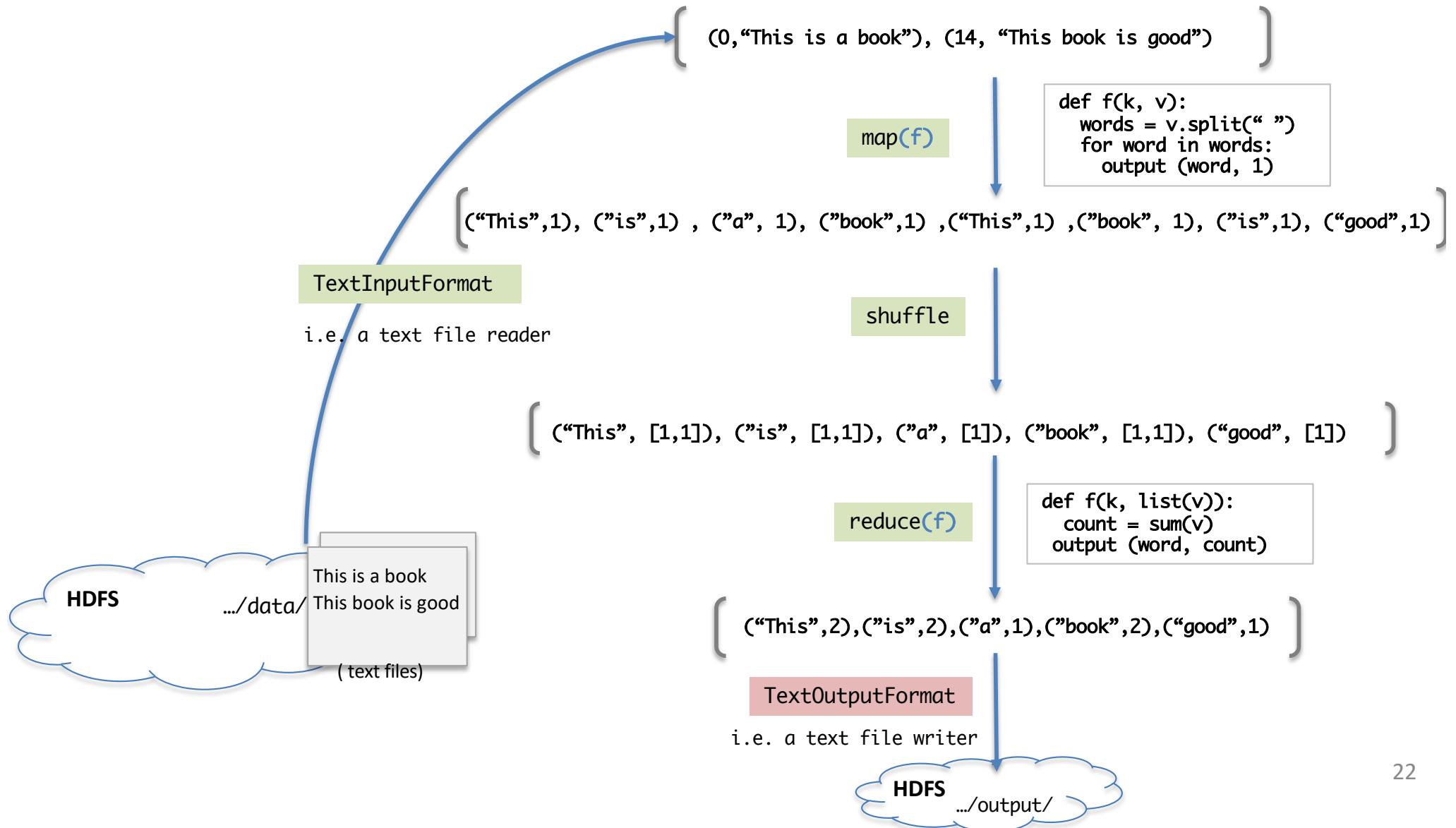
- Use HDFS CLI to handle files and folders

# Introduction to MapReduce

# A Bigger Picture



# MapReduce Programming Model (Concept)



# How a MR program is like ?

WordCount.java

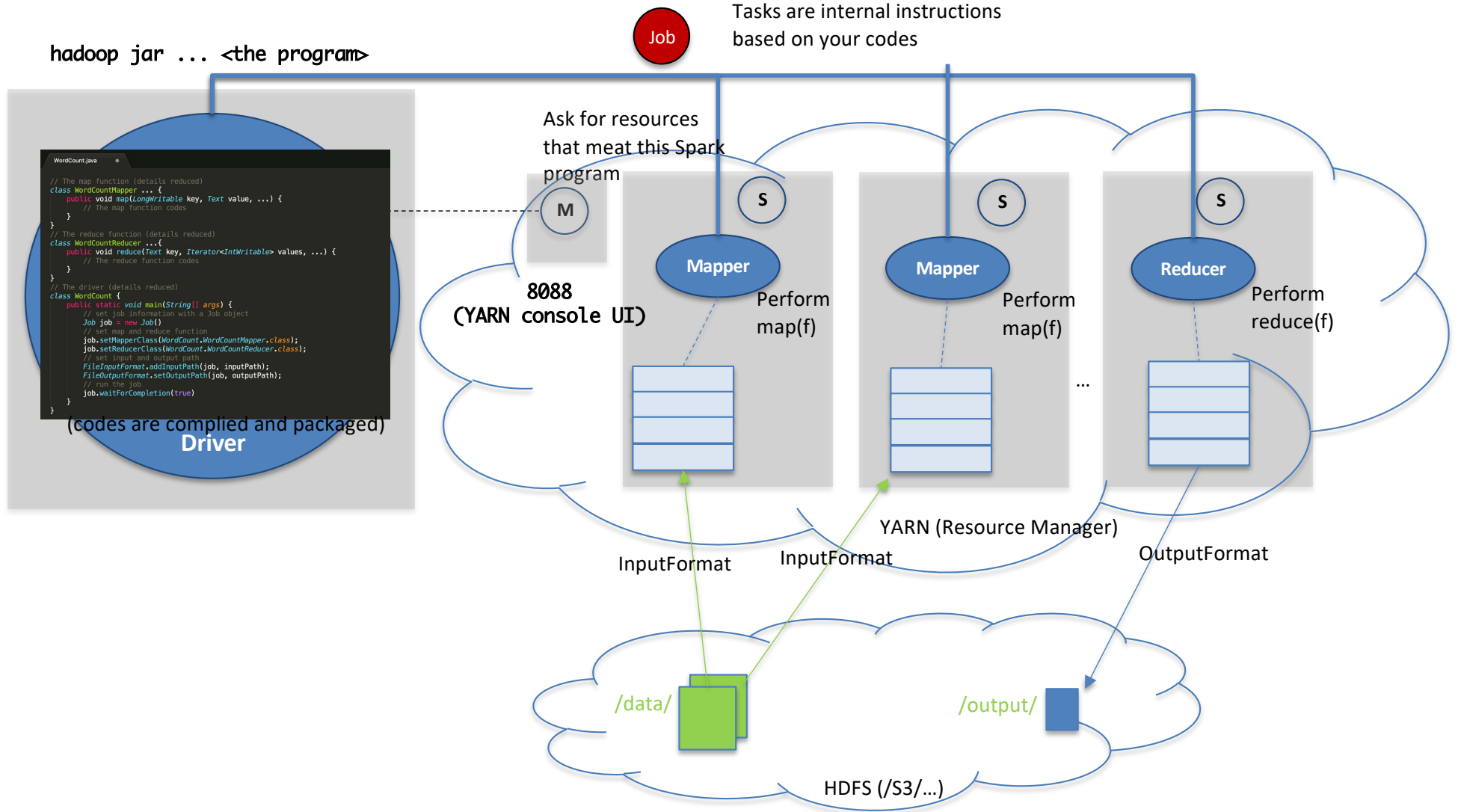
```
// The map function (details reduced)
class WordCountMapper ... {
    public void map(LongWritable key, Text value, ...) {
        // The map function codes
    }
}

// The reduce function (details reduced)
class WordCountReducer ...{
    public void reduce(Text key, Iterator<IntWritable> values, ...) {
        // The reduce function codes
    }
}

// The driver (details reduced)
class WordCount {
    public static void main(String[] args) {
        // set job information with a Job object
        Job job = new Job()
        // set map and reduce function
        job.setMapperClass(WordCount.WordCountMapper.class);
        job.setReducerClass(WordCount.WordCountReducer.class);
        // set input and output path
        FileInputFormat.addInputPath(job, inputPath);
        FileOutputFormat.setOutputPath(job, outputPath);
        // run the job
        job.waitForCompletion(true)
    }
}
```

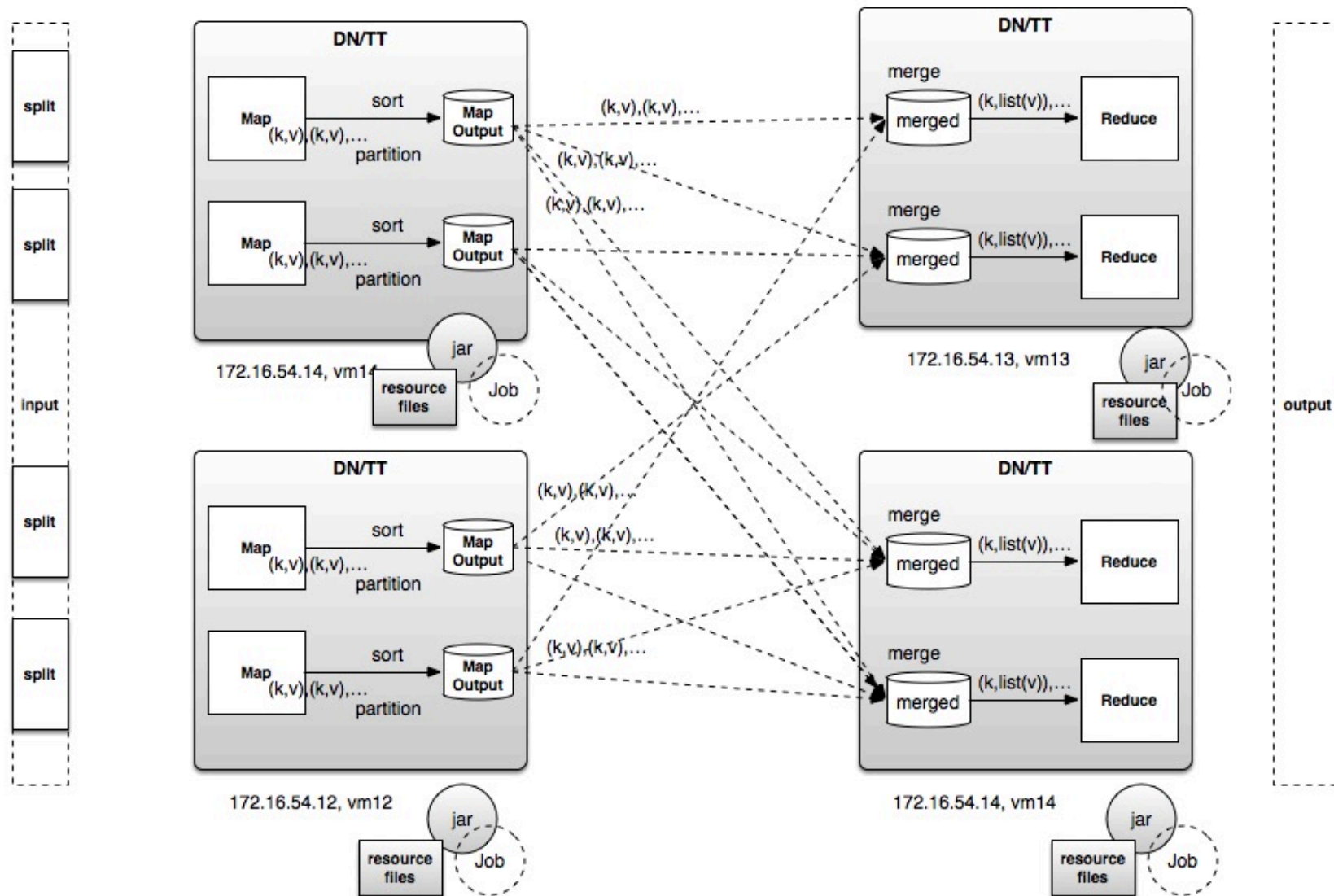
# MR Program Run-time Execution Flow

```
hadoop jar mr101.jar com.iii.mr101.WordCount hdfs://devenv/user/spark/mr101/wordcount/data hdfs://devenv/user/spark/mr101/wordcount/output
```





# MR Program Run-time Execution Flow (another view)



# Demo

- Install Java IDE
- Calculate word count example in MR
- Calculate average temperature in MR

(Just to have an idea how the MR program are like and run. Don't worry about JAVA 😊 )