

Map Reduce & Hadoop



Julian Kunkel

Learning Objectives



- Describe the architecture and features of Apache Hadoop
- Formulate processing problems in the MapReduce programming model
- Justify architectural decisions made in Apache Hadoop
- Sketch the execution phases of MapReduce and describe their behavior
- Describe limitations of Hadoop1 and the benefit of Hadoop2 with TEZ
- Sketch the parallel file access performed by MapReduce jobs on files

Iulian M. Kunkel CS3DP, 2019 2/45

Outline



- 1 Hadoop
- 2 Map Reduce
- 3 Hadoop 2
- 4 TEZ Execution Engine
- 5 Development
- 6 Summary

 Julian M. Kunkel
 CS3DP, 2019
 3/4

Hadoop

•00000000000000

- Apache Hadoop: Framework for scalable processing of data
 - Based on Google's MapReduce paper
- Consists of:
 - Hadoop distributed file system (HDFS)
 - MapReduce execution engine: schedules tasks on HDFS
- Why should we combine storage and execution paradigms?
 - Execution exploits data locality to avoid network data transfer
 - Ship compute to data and not (big) data to compute
- A complete ecosystem has been layered on top of MapReduce

Iulian M. Kunkel CS3DP, 2019 4/45

Hadoop Distributed File System (HDFS)



- Goal: Reliable storage on commodity-of-the-shelf hardware
- Implemented in Java
- Provides single-writer, multiple-reader concurrency model
- Has demonstrated scalability to 200 PB of storage and 4500 servers [12]

Features

Hadoop

00000000000000

- Hiearchical namespace (with UNIX/ACL permissions)
- High availability and automatic recovery
- Replication of data (pipelined write)
- Rack-awareness (for performance and high availability)
- Parallel file access

Julian M. Kunkel CS3DP, 2019 5/45

Hadoop File System Shell



Overview

00000000000000

Hadoop

- Invoke via: hadoop fs <command> <args>
 - Example: hadoop fs -ls hdfs://serverName/

HDFS command overview

- Read files: cat, tail, get, getmerge (useful!)
- Write files: put, appendToFile, moveFromLocal
- Permissions: chmod, chgrp, ..., getfacl
- Management: Is, rm, rmdir, mkdir, df, du, find, cp, mv, stat, touchz

Special commands

- distcp: map-reduce parallelized copy command between clusters
- checksum
- expunge (clear trash)
- setrep (change replication factor)

Iulian M. Kunkel CS3DP, 2019 6/45 Hadoop

00000000000000



- Namenode: Central manager for the file system namespace
 - ▶ Filenames, permissions
 - ▶ Information about file block (location)
 - ► For HA, a secondary NameNode backups data
- DataNode: Provide storage for objects

Map Reduce

- ▶ Directly communicates with other DataNodes for replication
- TaskTracker: accept and runs map, reduce and shuffle
 - Provides a number of slots for tasks (logical CPUs)
 - ▶ A **task** is tried to be scheduled on a slot of the machine hosting data
 - ▶ If all slots are occupied, run the task on the same rack
- JobTracker: Central manager for running MapReduce jobs
 - ► For HA, a secondary JobTracker backups data
- Tools to access and manage the file system (e.g., rebalancing)

Julian M. Kunkel CS3DP, 2019 7/45

High-Level Perspective

Hadoop

00000000000000



Hadoop Server Roles

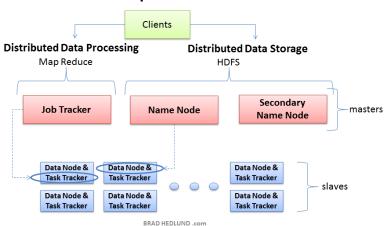


Figure: Source: B. Hedlund. [15]

 Julian M. Kunkel
 CS3DP, 2019
 8/45

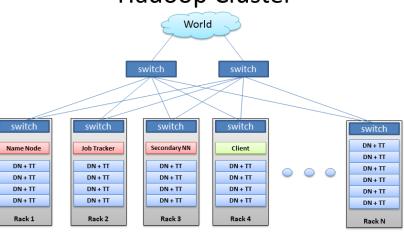
System-Level Perspective of Hadoop Clusters

Hadoop

00000000000000



Hadoop Cluster



BRAD HEDLUND .com

Figure: Source: B. Hedlund. [15]

 Julian M. Kunkel
 CS3DP, 2019
 9/45

Mapping of In-Memory Data Structures to Streams



(De)Serialization

Hadoop

00000000000000

- Data structure \Rightarrow byte stream \Rightarrow data structure
- Serialization is the process creating a byte stream from a data structure
- De-serialization creates a data structure in memory from the byte stream
- Byte streams can be transferred via network or stored on block storage

Serialization frameworks

- Provide serialization code for basic types
- Support writing of datatype-specific serializers
- Examples:
 - ▶ Java: Apache Avro¹, Kryo [40]
 - Python: Pickle
 - ► R: serialize(), unserialize() (functions for objects)
 - ▶ Apache Thrift supports multiple languages
- Requirements: Performance, platform independence

Julian M. Kunkel CS3DP, 2019 10/45

Excerpt of File Formats Supported by MapReduce



Text files

00000000000000

Hadoop

- Delimiters can be choosen
- Splittable at newlines (only decompressed files)

This is a simple file.\n With three lines – \n this is the end.

Comma-separated values (CSV)

- No header supported but JSON records are supported
- Does not support block compression

'max', 12.4, 10 \n
'john', 10.0, 2 \n

Sequence files

- Flat binary file for key/value pairs
- Supports splitting in HDFS using a synchronization marker
- Optional block compression for values (and keys)
- Widely used within Hadoop as internal structure

Julian M. Kunkel CS3DP, 2019 11/45

File Formats (2)



MapFile [21]

- Extends the sequence file
- Provides an index for keys

Avro

Hadoop

000000000000000

- Apache Avro's serialization system format
- Self-describing data format², allows inference of schema
 - ▶ Schema can also be changed upon read
- Enables exchange of data types between tools
- ⇒ Popular file format for Hadoop ecosystem

Mapping to Storage: Files are split into **blocks**

- A typical block size is 64 MiB
- Blocks are distributed across nodes
- Blocks may be compressed individually

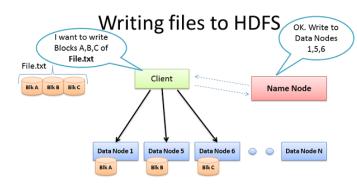
Julian M. Kunkel CS3DP, 2019 12/45

The HDFS I/O Path

Hadoop

00000000000000





- Client consults Name Node
- Client writes block directly to one Data Node
- Data Nodes replicates block
- Cycle repeats for next block

Figure: Source: B. Hedlund. [15]

Iulian M. Kunkel CS3DP, 2019 13/45

The HDFS Write Path

Hadoop



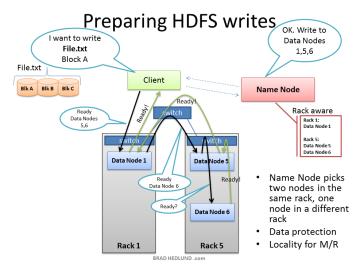


Figure: Source: B. Hedlund [15]

 Julian M. Kunkel
 CS3DP, 2019
 14/45

The HDFS Write Path

Hadoop

000000000000000



Multi-block Replication Pipeline

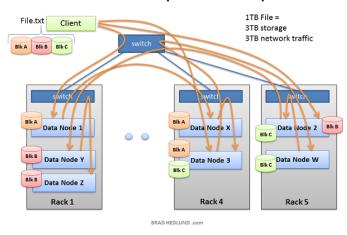


Figure: Source: B. Hedlund [15]

 Julian M. Kunkel
 CS3DP, 2019
 15/45

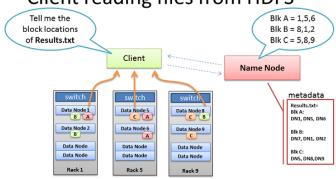
The HDFS Read Path

Hadoop

000000000000000



Client reading files from HDFS



- Client receives Data Node list for each block
- Client picks first Data Node for each block
- Client reads blocks sequentially

 BRAD HEDLUND COM

 BRAD HEDL

Figure: Source: B. Hedlund [15]

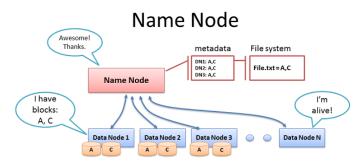
 Julian M. Kunkel
 CS3DP, 2019
 16/45

Name Node and High Availability

Hadoop

000000000000000





- Data Node sends Heartbeats
- Every 10th heartbeat is a Block report
- Name Node builds metadata from Block reports
- TCP every 3 seconds
- If Name Node is down, HDFS is down

BRAD HEDLUND .com

Figure: Source: B. Hedlund. [15]

Julian M. Kunkel CS3DP, 2019 17/45

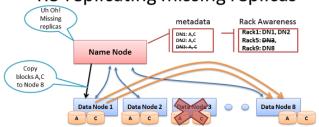
Name Node and High Availability

Hadoop

00000000000000



Re-replicating missing replicas



- Missing Heartbeats signify lost Nodes
- Name Node consults metadata, finds affected data
- Name Node consults Rack Awareness script
- Name Node tells a Data Node to re-replicate

BRAD HEDLUND .com

Figure: Source: B. Hedlund. [15]

Julian M. Kunkel CS3DP, 2019 18/45

Outline

000000000000000



- 2 Map Reduce
 - Overview
 - Execution on Hadoop

Map Reduce

•0000000000

Iulian M. Kunkel CS3DP, 2019 19/45

Map Reduce Execution Paradigm



Idea: Appy a processing pipeline consisting of map and reduce operations

- 1. Map: filter and convert input records (pos, data) to tuples (key, value)
- 2. Reduce: receives all tuples with the same key (key, list<value>)
- Hadoop takes care of reading input, distributing (key,value) to reduce
- Types for key, value & format, records depend on the configuration

Example: WordCount [10]: Count word frequency in large texts

```
map(key, text): # input: key=position, text=line
for each word in text:
    Emit(word,1) # outputs: key/value

reduce(key, list of values): # input: key == word, our mapper output
    count = 0
for each v in values:
    count += v
Emit(key, count) # it is possible to emit multiple (key, value) pairs here
```

Julian M. Kunkel CS3DP, 2019 20/45

Map Reduce Execution: Aggregation of Tables



Example from [17]

Goal: aggregate a CSV file by grouping certain entries

```
Country State City Population
USA,
        CA,
               Su, 12
                                         Country State
                                                          Population
USA.
        CA.
               SA. 42
                                \Rightarrow
                                         USA
                                                 CA
                                                          54
USA.
        MO.
               XY. 23
                                         USA
                                                 MO
                                                          33
USA.
        MO.
               AB. 10
```

Algorithm

```
map(kev. line):
  (county, state, city, population) = line.split(',')
  EmitIntermediate( (country, state), population )
reduce(key, values): # key=(country.state) values=list of populations
  count = 0
  for each v in values:
    count += v
  Emit(kev. count)
```

Iulian M. Kunkel CS3DP, 2019 21/45

Phases of MapReduce Execution

Map Reduce

00000000000



Phases of MapReduce Execution

- Distribute code (IAR files)
- Determine files to read, blocks and file splits, assign mappers to splits and slots
- Map: Invoke (local) map functions
- 4 Shuffle: Sort by the key, exchange data
- 5 Combine: Perform a local reduction by the key
- 6 Partition: Partition key space among reducers (typically via hashing)
- Reduce: Invoke reducers
- Write output, each reducer writes to its own file³

Iulian M. Kunkel CS3DP, 2019 22/45

Parallel Access to Files



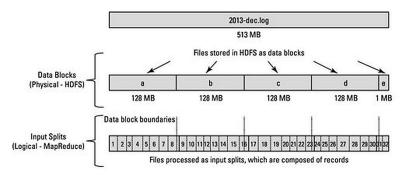
- A MapReduce job processes all files in a directory
 - Provides parallelism on the file level, each file is read independently
- MapReduce jobs process records that are grouped in input splits
 - ► Input splits == logical organization of blocks
 - ► Each input split is processed by one **mapper** (local processing preferred)
 - Processing for records spanning blocks
 - Skip partial records at the beginning of a split
 - · For truncated records, read data from a remote
 - ▶ Input splitting (intelligence) depends on the file format
- File formats that are not splittable must be avoided
 - e.g., XML, JSON Files, compressed text files
 - ▶ They enforce sequential read by one mapper
- Usage of file formats depends on the tools to query data

 Julian M. Kunkel
 CS3DP, 2019
 23/45

Hadoop

Mapping of Data Blocks to Input Splits [23]





map: (K1, V1) → list(K2, V2)

reduce: (K2, list(V2)) -> list(K3, V3)

Figure: Source: [23]

Julian M. Kunkel CS3DP, 2019 24/45

Execution of MapReduce – the Big Picture



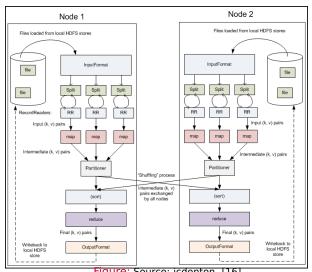


Figure: Source: jcdenton. [16]

Execution of MapReduce on HDFS – the Combiner



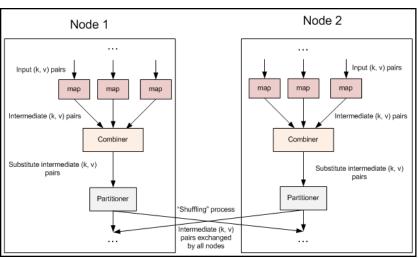


Figure: Source: jcdenton. [16]

 Julian M. Kunkel
 CS3DP, 2019
 26/45

Execution of MapReduce Tasks on Hadoop [14]



Steps in the execution of tasks

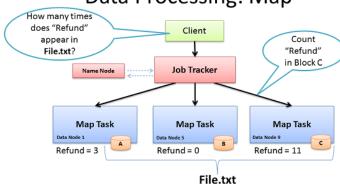
- Client submits a job to the JobTracker
- 2 JobTracker identifies the location of data via the NameNode
- JobTracker locates TaskTracker nodes with free slots close to the data
- 4 JobTracker starts tasks on the TaskTracker nodes
- Monitoring of TaskTrack nodes
 - ▶ If heartbeat signals are missed, work is rescheduled on another TaskTracker
 - ▶ A TaskTracker will notify the JobTracker when a task fails
- The JobTracker constantly updates its status
 - Clients can query this information

Julian M. Kunkel CS3DP, 2019 27/45

Execution of MapReduce



Data Processing: Map



- Map: "Run this computation on your local data"
- Job Tracker delivers Java code to Nodes with local data

BRAD HEDLUND .com

Figure: Source: B. Hedlund. [15]

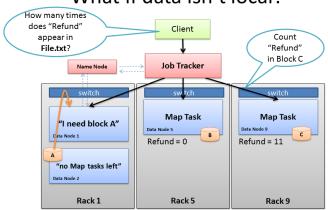
 Julian M. Kunkel
 C53DP, 2019
 28/45

Execution of MapReduce



Development

What if data isn't local?



- Job Tracker tries to select Node in same rack as data
- Name Node rack awareness

Figure: Source: B. Hedlund. [15]

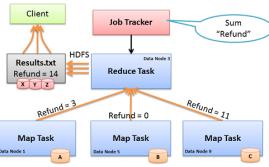
 Julian M. Kunkel
 C53DP, 2019
 29/45

Execution of MapReduce

Hadoop



Data Processing: Reduce



- **Reduce:** "Run this computation across Map results"
- Map Tasks send output data to Reducer over the network
- Reduce Task data output written to and read from HDFS

BRAD HEDLUND .com

Figure: Source: B. Hedlund. [15]

Iulian M. Kunkel CS3DP, 2019 30/45

Outline



- 3 Hadoop 2
 - Overview
 - System Architecture

Iulian M. Kunkel CS3DP, 2019 31/45

University of 💎 Reading

Hadoop 2, the Next Generation [12]

- Goal: real-time and interactive processing of events
- Introduction of YARN: Yet Another Resource Negotiator
- Supports of classical MapReduce and, via TEZ, DAG of tasks
- Support for NFS access to HDFS data
- Compatability to Hadoop v1

Hadoop

High-availability, federation and snapshots for HDFS

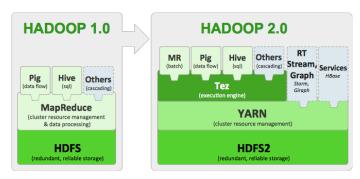


Figure: Source: Apache Hadoop 2 is now GA. Hortonworks. [12]

Iulian M. Kunkel CS3DP, 2019 32/45

System Architecture



Yarn modularizes JobTracker functionality

- Resource management
- 2 Job scheduling/execution inclusive monitoring

Data computation framework

- Applications are executed in containers
- ResourceManager component (global daemon)
 - Partitiones resources and schedules applications
 - Scheduler: distributes resources among applications
 - ApplicationsManager: accepts jobs, negotiates execution of AppMaster
- Per-node NodeManager: manages and monitors local resources
- Per-application ApplicationMaster
 - Framework-specific library
 - Negotiates container resources with ResourceManager
 - Works with Scheduler/NodeManager to execute and monitor tasks

Iulian M. Kunkel CS3DP, 2019 33/45

YARN System Architecture

000000000000000



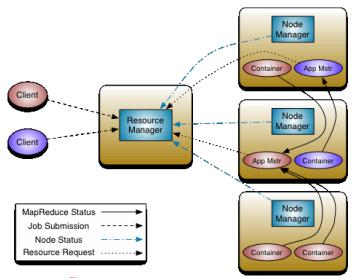


Figure: Source: Apache Hadoop NextGen [18]

Iulian M. Kunkel CS3DP, 2019 34/45

Outline

00000000000000



- 4 TEZ Execution Engine

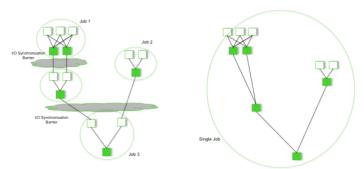
Iulian M. Kunkel CS3DP, 2019 35/45

university of 💎 Reading

TEZ Execution Engine

Hadoop

- TEZ: Hindi for "speed"
- Allow modelling and execution of data processing logic
 - Directed acyclic graph (DAG) of tasks
 - Vertex with input (dependencies) and output edges
- VertexManager defines parallelism and resolves dependencies



Pia/Hive - MR Pig/Hive - Tez Figure: Source: Introducing... Tez. Hortonworks [19]

Iulian M. Kunkel CS3DP, 2019 36/45

TEZ Example DAG [20]



```
1 // Define DAG
2 DAG dag = new DAG();
3 // Define Vertex, which class to execute
4 Vertex Map1 = new Vertex(Processor.class);
5 // Define Edge
  Edge edge = Edge(Map1, Reduce2,
    SCATTER_GATHER, // Distribution of data from
         \hookrightarrow source to target(s)
    PERSISTED, // Persistency of data
    SEOUENTIAL. // Scheduling: either concurrent
         \hookrightarrow or sequential execution
    Output.class, Input.class);
10
11 // Connect edges with vertex
dag.addVertex(Map1).addEdge(edge)...
```

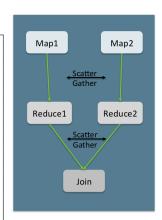


Figure: Source: Apache Tez. H.
Shah [20]

Julian M. Kunkel CS3DP, 2019 37/45

TEZ DAG API

Hadoop



Edge properties define the connection between producer and consumer tasks in the DAG

- Data movement Defines routing of data between tasks
 - **One-To-One**: Data from the ith producer task routes to the ith consumer task.
 - **Broadcast**: Data from a producer task routes to all consumer tasks.
 - Scatter-Gather: Producer tasks scatter data into shards and consumer tasks gather the data. The ith shard from all producer tasks routes to the ith consumer task.
- Scheduling Defines when a consumer task is scheduled
 - Sequential: Consumer task may be scheduled after a producer task completes.
 - **Concurrent**: Consumer task must be co-scheduled with a producer task.
- Data source Defines the lifetime/reliability of a task output
 - Persisted: Output will be available after the task exits. Output may be lost later on.
 - Persisted-Reliable: Output is reliably stored and will always be available
 - Ephemeral: Output is available only while the producer task is running

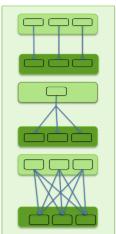


Figure: Source: Apache Tez. H. Shah [20]

Julian M. Kunkel CS3DP, 2019 38/45

TEZ Dynamic Graph Reconfiguration



- Reconfigure dataflow graph based on data sizes and target load
- Controlled by vertex management modules
 - ▶ State changes of the DAG invoke plugins on the vertices
 - ▶ Plugins monitor runtime information and provide hints to TEZ

Example: Adaption of the number of reducers

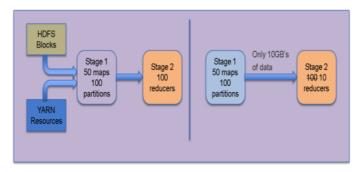


Figure: Source: Introducing... Tez. Hortonworks [19]

 Julian M. Kunkel
 CS3DP, 2019
 39/45

TEZ Resource Management

Hadoop



- Task and resource aware scheduling
- Pre-launch and re-use containers and intermediate results (caching)

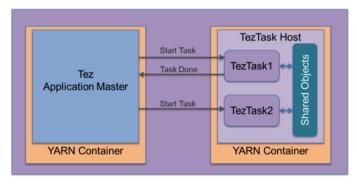


Figure: Source: Introducing... Tez. Hortonworks [19]

Julian M. Kunkel C53DP, 2019 40/45

Outline

00000000000000



- 5 Development ■ Coding

Iulian M. Kunkel CS3DP, 2019 41/45

Coding



- Programming Map-Reduce can be done in various languages
 - Java (low-level)
 - Python
 - **...**
- Process:
 - ▶ Implement map/reduce functions
 - Main method controls process:
 - Define mapper/reducer/combiner
 - · Define input/output formats and files
 - · Run the job
- Programming of TEZ in Java
- Command line tools to run the "application"
- There are some tools for debugging / performance analysis

 Julian M. Kunkel
 CS3DP, 2019
 42/45

Coding: Wordcount, Mapper & Reducer



Goal: Count the frequency of each word in a text

```
package org.myorg;
   import java.io.IOException; import java.util.*; import org.apache.hadoop.fs.Path; import org.apache.hadoop.conf.*;
   import org.apache.hadoop.io.*: import org.apache.hadoop.mapred.*: import org.apache.hadoop.util.*:
   public class WordCount {
     public static class Map extends MapReduceBase implements Mapper<LonoWritable. Text. Text. IntWritable> {
        private final static IntWritable one = new IntWritable(1); // for small optimization of object cleaning, reuse object
 8
 9
        // Mapper splits sentence and creates the tuple (word, 1) for each word
10
        public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
11
          String line = value.toString():
12
          Text word = new Text();
13
         StringTokenizer tokenizer = new StringTokenizer(line):
         while (tokenizer.hasMoreTokens()) {
14
15
            word.set(tokenizer.nextToken());
16
            output.collect(word, one):
17
18
       }}
19
20
     // Reducer accumulates tuples with the same key by summing their frequency
     public static class Reduce extends MapReduceBase implements Reducer<Text. IntWritable. Text. IntWritable> {
21
22
        public void reduce(Text key. Iterator<IntWritable> values. OutputCollector<Text. IntWritable> output. Reporter reporter) throws IOException
               \hookrightarrow 4
23
          int sum = \theta:
24
          while (values.hasNext()) {
25
            sum += values.next().get();
26
27
          output.collect(key, new IntWritable(sum));
28
29
      } // Continued => see the next slide
```

Julian M. Kunkel CS3DP, 2019 43/45

Development 0000

Coding: Wordcount, Main Method



The main method configures the Hadoop Job⁴

```
public static void main(String[] args) throws Exception {
       JobConf conf = new JobConf(WordCount.class):
       conf.set.lobName("wordcount"):
       // Set data types of output
       conf.setOutputKeyClass(Text.class);
       conf.setOutputValueClass(IntWritable.class):
 8
 9
       // Set classes for map, reduce and combiner
10
       conf.setMapperClass(Map.class):
11
       conf.setReducerClass(Reduce.class);
12
       conf.setCombinerClass(Reduce.class):
13
14
       // Set file input and output format
15
       conf.setInputFormat(TextInputFormat.class):
16
       conf.setOutputFormat(TextOutputFormat.class);
17
18
       // Configure input and output paths
19
       FileInputFormat.setInputPaths(conf, new Path(args[0]));
20
       FileOutputFormat.setOutputPath(conf. new Path(args[1])):
21
22
       JobClient.runJob(conf);
23
24
```

See https://github.com/apache/tez/tree/master/tez-examples/src/main/java/ org/apache/tez/examples for examples with TEZ

Iulian M. Kunkel CS3DP, 2019 44/45



- Hadoop provides the file system HDFS and concepts for processing
- HDFS
 - Single writer, multiple reader concurrency
 - ► Robust and high availability
- MapReduce: fixed function pipeline, reliable execution
- Hadoop2 with YARN: refined architecture for resource management
- TEZ: Execution of DAGs with various configurations

Bibliography



- Book: Lillian Pierson. Data Science for Dummies. John Wiley & Sons
- 10 Wikipedia
- Hortonworks http://hortonworks.com/
- B. Ramamurthy. Hadoop File System. http://www.cse.buffalo.edu/faculty/bina/MapReduce/HDFS.ppt
- Hadoop Wiki. https://wiki.apache.org/hadoop/
- B. Hedlund. Understanding Hadoop Clusters and the Network. http://bradhedlund.com/2011/09/10/understanding-hadoop-clusters-and-the-network/
- 16 jcdenton. H-stack Introduction Guide. https://github.com/jcdenton/hadoop-quide/blob/master/hadoop.md
- http://tutorials.techmytalk.com/2014/11/14/mapreduce-composite-key-operation-part2/
- http://hadoop.apache.org/docs/
- http://hortonworks.com/blog/introducing-tez-faster-hadoop-processing/
- Presentation: H. Shah. Apache Tez. Hortonworks.
- Hadoop I/O http://blog.cloudera.com/blog/2011/01/hadoop-io-sequence-map-set-array-bloommap-files/
- http://hadoop.apache.org/docs/current/hadoop-streaming/HadoopStreaming.html
- 23 http://www.dummies.com/how-to/content/input-splits-in-hadoops-mapreduce.html

Iulian M. Kunkel CS3DP, 2019 46/45

Appendix

Compilation



Here we compile manually and are not using ant or maven:

- Prepare the class path for dependencies (may be complex)
- Compile each Java file
- Create a JAR package

Example

Execution



Syntax: [hadoop|yarn] jar FILE.jar ClassWithMain Arguments

Example

```
> hadoop jar averagePerformance.jar de.wr.AveragePerformance data-energy-efficiency.csv summary
   STARTUP: Computing average ## NOTE: This is output of the main() method
   15/10/15 13:49:24 INFO impl.TimelineClientImpl: Timeline service address: http://abu3.cluster:8188/ws/v1/timeline/
   15/10/15 13:49:25 INFO client.RMProxy: Connecting to ResourceManager at abu3.cluster/10.0.0.65:8050
   15/10/15 13:49:25 INFO impl.TimelineClientImpl: Timeline service address: http://abu3.cluster:8188/ws/v1/timeline/
   15/10/15 13:49:25 INFO client.RMProxy: Connecting to ResourceManager at abu3.cluster/10.0.0.65:8050
   15/10/15 13:49:26 INFO mapred.FileInputFormat: Total input paths to process: 1
   15/10/15 13:49:26 INFO mapreduce.JobSubmitter: number of splits:8
   15/10/15 13:49:26 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1444759114226_0016
   15/10/15 13:49:27 INFO impl.YarnClientImpl: Submitted application application_1444759114226_0016
   15/10/15 13:49:27 INFO mapreduce.Job: The url to track the job: http://abu3.cluster:8088/proxy/application_1444759114226_0016/
   15/10/15 13:49:27 INFO mapreduce.Job: Running job: job_1444759114226_0016
   15/10/15 13:49:37 INFO mapreduce.Job: Job job_1444759114226_0016 running in uber mode : false
   15/10/15 13:49:37 INFO mapreduce.Job: map 0% reduce 0%
   15/10/15 13:49:54 INFO mapreduce.Job: map 11% reduce 0%
   15/10/15 13:50:02 INFO mapreduce.Job: map 100% reduce 100%
   15/10/15 13:50:02 INFO mapreduce. Job: Job job_1444759114226_0016 completed successfully
   15/10/15 13:50:02 INFO mapreduce.Job: Counters: 50
19
     File System Counters
20
       FILE: Number of bytes read=768338
21
       FILE: Number of bytes written=2679321
22
       FILE: Number of read operations=0
23
       FILE: Number of large read operations=0
24
       FILE: Number of write operations=0
25
       HDFS: Number of bytes read=1007776309
26
       HDFS: Number of bytes written=1483856
27
       HDFS: Number of read operations=27
28
       HDFS: Number of large read operations=0
29
       HDFS: Number of write operations=2
30
     Job Counters
31
       Launched map tasks=8
32
       Launched reduce tasks=1
```

Retrieving Output Data



The output is a directory containing one file per reducer

```
# Retrieve the summary directory
2 $ hadoop fs -get summary
3 $ ls -lah summary/
4 -rw-r--r-- 1 kunkel wr 1.5M Okt 15 14:45 part-00000
  -rw-r--r-- 1 kunkel wr 0 0kt 15 14:45 _SUCCESS
6 $ head summary/part-00000
7 ESM_example_ESM_example_ESM_example_ESM_example 4397 112.69512266727315

→ 186388.93997432772 ....

8 EXX_example_EXX_example_EXX_example_EXX_example 4511 118.44219725094219
       \hookrightarrow 251865.2199417397 ...
  . . .
10
11 # A merged file can be retrieved via getmerge
12 hadoop fs -getmerge summary summary.csv
```

Using Arbitrary Tools/Languages via Streaming



- Hadoop Streaming [22] allows to pipe data through arbitrary tools
- This allows easy integration of Python code, e.g.

```
yarn jar /usr/hdp/current/hadoop-mapreduce/hadoop-streaming.jar \
-Dmapred.map.tasks=11 -mapper $PWD/mein-map.py \
-Dmapred.reduce.tasks=1 -reducer $PWD/mein-reduce.py \
-input <input> -output <output-directory>
```

- Map/reduce apps receive lines with key value pairs and emit them
 - ► ANY other (disturbing) output must be avoided to avoid errors
- Trivial mapper:

```
#!/usr/bin/python3
import sys

for line in sys.stdin:
    print("\t".join(line.split(","))) # Split CSV into key (first word) and values
```

Easy testing on the shell:

```
cat Input.csv | ./mein-map.py | sort | ./mein-reduce.py
```

Using Arbitrary Tools/Languages via Streaming



■ We can use the streaming also to integrate Rscripts

```
1 #!/usr/bin/env Rscript
 3 # WordCount Example
 4 | # Discard error messages for loading libraries (if needed) as this would be seen as a "tuple"
 5 sink(file=NULL, type="message")
 6 library('stringi')
 7 # Remove redirection
 8 sink(type="message")
10 stdin=file('stdin', open='r')
11
12 # Batch processing of multiple lines, here 100 elements
13 while(length( lines=readLines(con=stdin, n=100L) ) > 0){
   # paste concatenates all lines (the array) together
    # stri_extract_all_words() returns an 2D array of lines with words
    # Instead of paste, we could use unlist() to take care of multiple lines and returns a single array
     # table() counts number of occurences of factor levels (that are strings)
     tblWithCounts = table(stri_extract_all_words(paste(lines, collapse=" ")))
     words = names(tblWithCounts)
     counts = as.vector(tblWithCounts)
21
     cat(stri_paste(words, counts, sep="\t"), sep="\n")
22 1
```

Still: easy testing on the shell, similar execution with streaming

```
cat Input.csv | ./mein-map.R | sort | ./mein-reduce.py
```

Julian M. Kunkel CS3DP, 2019 52/45

Debugging of MapReduce and YARN Applications



Runtime information

- Call: yarn logs -applicationId < ID >
 - ▶ The ID is provided upon startup of the job
- Provides for each phase of the execution
 - Log4j output
 - Node information (logfiles)
 - Container information
 - Stdout, stderr of your application
- Increase log verbosity

```
1 export YARN_ROOT_LOGGER=DEBUG,console
2 or
3 run yarn --loglevel DEBUG ...
```

- ▶ Properties: mapreduce.map.log.level, mapreduce.reduce.log.level
- Dump the current configuration of (X) by adding the argument:
 - ▶ Parent class: hadoop org.apache.hadoop.conf.Configuration
 - ▶ Yarn: hadoop org.apache.hadoop.yarn.conf.YarnConfiguration
 - ▶ MapReduce: hadoop org.apache.hadoop.mapred.JobConf

Example Logfile Output



```
> yarn logs -applicationId application_1444759114226_0016
   Container: container_1444759114226_0016_01_000005 on abu3.cluster_45454
   LogType:stderr
   Log Upload Time: Thu Oct 15 13:50:09 +0200 2015
   LogLength: 243
   Log Contents:
   log4j:WARN No appenders could be found for logger (org.apache.hadoop.metrics2.impl.MetricsSystemImpl).
   log4i:WARN Please initialize the log4i system properly.
   log4j:WARN See http://logging.apache.org/log4j/1.2/fag.html#noconfig for more info.
12
   End of LogType:stderr
13
   LogType:stdout
   Log Upload Time: Thu Oct 15 13:50:09 +0200 2015
   LogLength:751944
   Log Contents:
18
   KEY: 134195662 word cpu_idl_idl_idl
   ACCEPTING LINE
   KEY: 134204510 word cpu_idl_idl_idl
   ACCEPTING LINE
   KEY: 134213460 word cpu_idl_idl_idl
   ACCEPTING LINE
   End of LogType:stdout
26
```

Job Information via Web Interface



- The task tracker keeps detailed information about job execution
- Access via an internal web-server on Port 8088 and 19888
- An internal web-server on each node provides node information

Example

```
# Output when submitting the job:
```

2 15/10/15 13:49:27 INFO mapreduce.Job: The url to track the job: http://abu3.cluster:8088/proxy/application_1444759114226_0016/

4 # After SSH forward visit localhost:8088, you may need to change the hostname from abu3.cluster to localhost again

Job Status



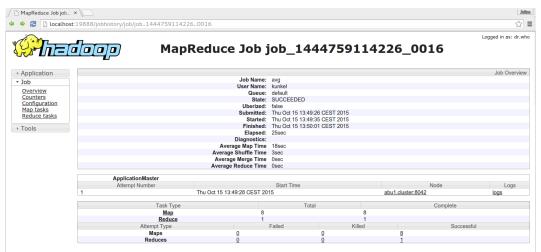
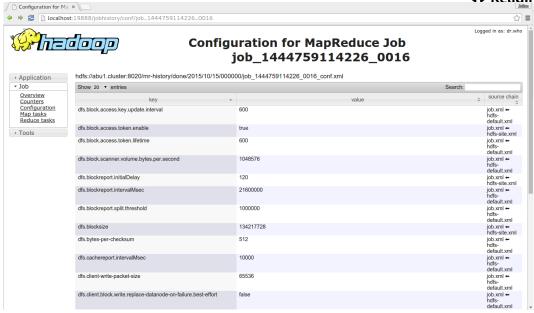


Figure: Overview, when using the tracking url

Julian M. Kunkel CS3DP, 2019 56/45

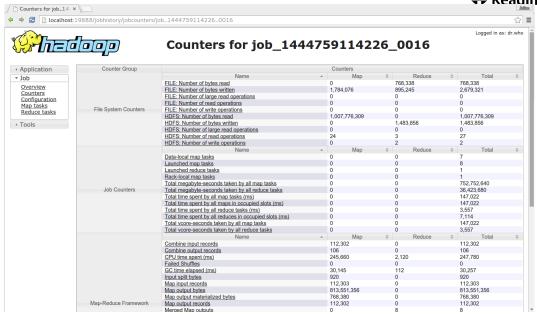
Job Configuration





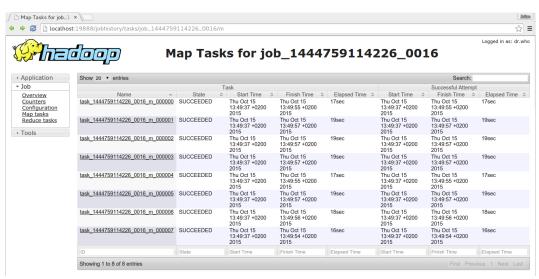
Performance Counters





Information About Map Tasks





Logfile



localhost:19888/job × h localhost:19888/iobhistory/logs/abu5.cluster:45454/container_1444759114226_0016_01_000010/attempt_1444759114226_0016_r_000000_0/kunkel Log Type: stderr Application Log Upload Time: Thu Oct 15 13:50:09 +0200 2015 About Log Length: 243 Jobs log4::WARN No appenders could be found for logger (org.apache.hadoop.metrics2.impl.MetricsSystemImpl). log4j:WARN Please initialize the log4j system properly. loq4j:WARN See http://loqqinq.apache.org/loq4j/1.2/faq.html#noconfig for more info. → Tools Log Type: stdout Log Upload Time: Thu Oct 15 13:50:09 +0200 2015 Log Length: 0 Log Type: syslog Log Upload Time: Thu Oct 15 13:50:09 +0200 2015 Log Length: 2269 2015-10-15 13:49:58,877 WARN [main] org.apache.hadoop.metrics2.impl.MetricsConfig: Cannot locate configuration: tried hadoop-metrics2-reducetask.properties,hadoop-metrics2.propertie 2015-10-15 13:49:58,972 INFO [main] org.apache.hadoop.metrics2.impl.MetricsSystemImpl: Scheduled snapshot period at 10 second(s). 2015-10-15 13:49:58,973 INFO [main] org.apache.hadoop.metrics2.impl.MetricsSystemImpl: ReduceTask metrics system started 2015-10-15 13:49:58,984 INFO [main] org.apache.hadoop.mapred.YarnChild: Executing with tokens: 2015-10-15 13:49:58.984 INFO [main] org.apache.hadoop.mapred.YarnChild: Kind: mapreduce.job. Service: job 1444759114226 0016. Ident: (org.apache.hadoop.mapreduce.security.token.JobT 2015-10-15 13:49:59.061 INFO [main] org.apache.hadoop.mapred.YarnChild: Sleeping for 0ms before retrying again. Got null now. 2015-10-15 13:49:59,334 INFO [main] org.apache.hadoop.mapred.YarnChild: mapreduce.cluster.local.dir for child: /tmp/hadoop/yarn/local/usercache/kunkel/appcache/application 144475911 2015-10-15 13:49:59,599 INFO [main] org.apache.hadoop.conf.Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session.id 2015-10-15 13:50:00,182 INFO [main] org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter: File Output Committer Algorithm version is 1 2015-10-15 13:50:00.195 INFO [main] org.apache.hadoop.mapred.Task: Using ResourceCalculatorProcessTree : [] 2015-10-15 13:50:00.264 INFO [main] org.apache.hadoop.mapred.ReduceTask; Using ShuffleConsumerPlugin; org.apache.hadoop.mapreduce.task.reduce.Shuffle@180da663 2015-10-15 13:50:00.899 INFO [main] org.apache.hadoop.mapred.Task: Task:attempt 1444759114226 0016 r 000000 0 is done. And is in the process of committing 2015-10-15 13:50:00 964 INFO [main] org.apache.hadoop.mapred.Task: Task attempt 1444759114226 0016 r 0000000 0 is allowed to commit now 2015-10-15 13:50:00,974 INFO [main] org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter: Saved output of task 'attempt 1444759114226 0016 r 000000 0' to hdfs://abul.cluster:8 2015-10-15 13:50:01,021 INFO [main] org.apache.hadoop.mapred.Task: Task 'attempt 1444759114226 0016 r 000000 0' done. Log Type: syslog.shuffle Log Upload Time: Thu Oct 15 13:50:09 +0200 2015 Log Length: 8352 Showing 4096 bytes of 8352 total. Click here for the full log. mpt 1444759114226 8816 m 888886 8 2015-10-15 13:50:00,491 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl: closeInMemoryFile -> map-output of size: 101418, inMemoryMapOutputs.size() -> 4, c 2015-10-15 13:50:00.492 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.Fetcher: fetcher#9 about to shuffle output of map attempt 1444759114226 0016 m 000002 0 decomp: 1014 2015-10-15 13:50:00.492 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput: Read 101418 bytes from map-output for attempt 1444759114226 0016 m 0000002 0 2015-10-15 13:50:00.492 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl; closeInMemoryFile -> map-output of size; 101418, inMemoryMapOutputs.size() -> 5. c 2015-10-15 13:50:00,493 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.Fetcher# about to shuffle output of map attempt 1444759114226 0016 m 000001 0 decomp: 1451 2015-10-15 13:50:00,494 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput: Read 145170 bytes from map-output for attempt 1444759114226 0016 m 0000010 2015-10-15 13:50:00,494 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl: closeInMemoryFile -> map-output of size: 145170, inMemoryMapOutputs.size() -> 6, c 2015-10-15 13:50:00.494 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.Fetcher: fetcher#9 about to shuffle output of map attempt 1444759114226 0016 m 000003 0 decomp: 9417 2015-10-15 13:50:00.495 INFO [fetcher#9] org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput: Read 94174 bytes from map-output for attempt 1444759114226 0016 m 0000030 0 2015-10-15 13:50:00.495 INFO | fetcher#9| org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl; closeInMemoryFile -> map-output of size; 94174, inMemoryMapOutputs.size() -> 7, co v

Node Manager



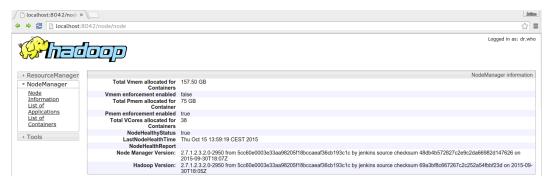


Figure: The Node Manager provides information about a particular node