"If you torture the data long enough, it will confess"

Ronald Coase, Professor at UChicago, 1981

Lecture C1 Batch Data Processing

CS205: Computing Foundations for Computational Science Dr. Ignacio M. Llorente Spring Term 2020



HARVARD

School of Engineering and Applied Sciences



INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE AT HARVARD UNIVERSITY

Before We Start

Where We Are

Computing Foundations for Computational and Data Science

How to use modern computing platforms in solving scientific problems

Intro: Large-Scale Computational and Data Science

- A. Parallel Processing Fundamentals
- B. Parallel Computing
- C. Parallel Data Processing
 - **C1. Batch Data Processing**
 - C2. Dataflow Processing
 - C3. Stream Data Processing

Wrap-Up: Advanced Topics



CS205: Contents

APPLICATION SOFTWARE





INSTITUTE FOR APPLIED LE COMPUTATIONAL SCIENCE CS

Before We Start

Where We Are

Week 9: Batch Data Processing => MapReduce

3/23	3/24 <u>Lecture C1</u>	3/25 Lab I8	3/26 <u>Hands-on H4</u>	3/27
	Batch Data Processing	MapReduce Hadoop Cluster	MapReduce	
	(Quiz & Reading)			

Week 10: Dataflow Processing => Spark

3/30	3/31 Lecture C2	4/1 Lab I9	4/2 <u>Hands-on H5</u>	4/3
	Dataflow Processing	Spark Single Node	Spark	
	(Quiz & Reading)			



Context

Foundations of Data Processing

"Big" Compute







Data-intensive

Bringing compute to data



1001 (001 (031 HARVARD

IACS INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE School of Engineering AT HARVARD UNIVERSITY and Applied Sciences

Lecture C1. Batch Data Processing CS205: Computing Foundations for Computational Science

Context Foundations of Data Processing

Paradigm	Independent parallel tasks that are performed simultaneously to address a particular part of the problem
Challenge	Decompose the application into tasks and define their communication and synchronization
Bottleneck	CPU
Input data	Gigabyte-scale to describe initial conditions
Programming	OpenMP, OpenACC and MPI







Lecture C1. Batch Data Processing

CS205: Computing Foundations for Computational Science

Context Foundations of Data Processing

Paradigm	Same task is applied to large volumes of data
Challenge	Partition the data into multiple segments and the subsequent combination of the intermediate results in multiple stages
Bottleneck	Storage
Input data	Far beyond gigabyte-scale: datasets are commonly on the order of tens, hundreds, or thousands of terabytes
Programming	MapReduce, Spark





IACS INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE AT HARVARD UNIVERSITY Lecture C1. Batch Data Processing

CS205: Computing Foundations for Computational Science

Hands-on Examples Requirements

- 1. Unix-like shell (Linux, Mac OS or Windows/Cygwin)
- 2. Python installed
- 3. Download example python codes

https://harvard-iacs.github.io/2020-CS205/pages/materials.html





Why Is Big Data Processing Different? The MapReduce Programming Model The Hadoop Processing Framework



The WordCount HelloWorld Example



counter.py

#!/usr/bin/python

```
import sys
import re
```

 $sums = \{\}$

```
for line in sys.stdin:
   line = re.sub( r'^W+W+, '', line )
   words = re.split(r'\\+', line)
```

INSTITUTE FOR APPLIED

AT HARVARD UNIVERSITY

```
for word in words:
   word = word.lower()
    sums[word] = sums.get( word, 0 ) + 1
```

Implementation

 Centralized key-value data structure, hash table (dictionary sums) to keep track of counts

Scalability Limitations

- Compute-bound: Limited by the speed of the system
- Memory-bound: Limited by the memory size of the system

print sums

Why Is Big Data Processing Different? The WordCount Example



Is the counter application limited by the CPU?

How would you develop a parallel version of the counter application?



The WordCount Example





The WordCount Example





Implementation

- Each node processes a subset in parallel
- Each node executes a sequential instance of the counter code and keeps its own local data structure
- Big final reduction operation for the complete data structure

Scalability Limitations

- Memory-bound: Limited by the memory size of each node
- Communication-bound: Cost of final aggregation with reduction of all the data structure



Data-Intensive Applications: Bring Compute to the Data

We want to avoid

- Centralized resources that are likely bottlenecks
- Replication of data structures across nodes
- Communication of too much intermediate data

We need a programming model with data locality

- Same computation to be applied to large volumes of data
- Assign tasks to machines that already have the input data
- Efficient combination of intermediate results from multiple processors
- Highly distributed and scale-out







Lecture C1. Batch Data Processing CS205: Computing Foundations for Computational Science

The MapReduce Programming Model Core Idea and Benefits

MapReduce is a programming model for processing <u>big data sets</u> with a parallel, distributed algorithm on a cluster

The core idea behind MapReduce is **mapping** your data set into a collection of <key, value> pairs, and then **reducing** over all pairs with the same key

The concept is quite **powerful** because almost all data can be mapped into <key, value> pairs somehow, and keys and values can be of any type (strings, integers, user-defined...)

The concept is very **simple** because developers are required to only write <u>simple map and reduce functions</u>, while distribution and parallelism are handled by the MapReduce framework

The concept is very **efficient** because computation operations are performed on data local to the computing node, data transfer over the network is reduced to a minimum



The MapReduce Programming Model Not So New

MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

Google, Inc.

Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value pairs, and a *reduce* function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code to deal with these issues.

OSDI'04, San Francisco, CA, December, 2004 https://www.usenix.org/legacy/event/osdi04/tech/full_papers/dean/dean.p

df



The MapReduce Programming Model Assign Compute to Machines that already Have the Data

The programmer essentially only specifies two (sequential) functions

STEP 1. MAP: $map(k1,v1) \rightarrow list(k2,v2)$

- Inputs each record consisting of key of type k1 and value of type v1
- Outputs a set of intermediate key-value pairs, each of type k2 and v2
- Types can be simple or complex user-defined objects
- Each map call is independent

STEP 2. SUFFLING: Internal grouping of all intermediate pairs with same key together and passes them to the workers executing reduce

STEP 3. REDUCE: $reduce(k2, list(v2)) \rightarrow list(k3, v3)$

• Combines information across records that share this same **intermediate key**







WordCount Example on a Single System

STEP 1. MAP: $map(k1,v1) \rightarrow list(k2,v2)$

```
$ mapper.py < input.txt</pre>
```

```
...
email 1
newsletter 1
to 1
hear 1
about 1
new 1
```

WordCount Example on a Single System

STEP 2. SUFFLING

<pre>\$ mapper.py < input.txt</pre>	sort
•••	
zodiac 1	
zodiac 1	
zogranda 1	
zone 1	
zone 1	
zone 1	
zone 1	
zone 1	
zoned 1	
zoned 1	
zones 1	
zones 1	
zones 1	
zoology 1	
zoology 1	
zoroaster 1	
HARVARD School of Engineering and Applied Sciences	C1. Batch Data Computing Fo

WordCount Example on a Single System

STEP 3. REDUCE: $reduce(k2, list(v2)) \rightarrow list(k3, v3)$

reducer.py Count the number of times each key #!/usr/bin/python occurs by summing values as long as import sys they have the same key previous = None Publish the result once the key changes sum = 0 for line in sys.stdin: key, value = line.split('\t') if key != previous: if previous is not None: print str(sum) + '\t' + previous previous = kevsum = 0 sum = sum + int(value) print str(sum) + '\t' + previous mapper.py < input.txt | sort |</pre> reducer.py

•••

3 zones

- 2 zoology
- 1 zoroaster

Prototyping and Debugging - Hadoop Streaming

Both the mapper and the reducer should be python executable scripts that read the input from stdin (line by line) and emit the output to stdout

\$ cat files | mapper.py| sort | reducer.py

1. Copy files to HDFS

bin/hadoop dfs -copyFromLocal /tmp/gutenberg /user/hduser/gutenberg

2. Execute Hadoop command

\$ bin/hadoop jar contrib/streaming/hadoop-*streaming*.jar \ -file
/home/hduser/mapper.py -mapper /home/hduser/mapper.py \ -file
/home/hduser/reducer.py -reducer /home/hduser/reducer.py \ -input
/user/hduser/input/* -output /user/hduser/g-output

3. Read all output files (one per reducer)





Programmer focus on the algorithm while the framework takes care of:

- Parallelizing program execution
- Partitioning input data
- Delivering data chunks to the different worker machines
- Scheduling the map/reduce tasks for execution on the worker machines
- Handling machine failures and slow responses





The MapReduce Programming Model WordCount Example on a Parallel System









Apache Hadoop and Alternatives











Lecture C1. Batch Data Processing CS205: Computing Foundations for Computational Science

The Hadoop Processing Framework **Bare-metal Deployment**





Elastic Map Reduce - AWS



108; 100; 103;

HARVARD

The Hadoop Processing Framework Scale Horizontally!







More, smaller servers







Elastic Map Reduce - AWS







Clustered Architectures





IACS INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE AT HARVARD UNIVERSITY Lecture C1. Batch Data Processing

CS205: Computing Foundations for Computational Science

Next Steps

- Get ready for next lab: I8. Hadoop Cluster on AWS (Wednesday 3/25)
- Get ready for next hands-on: H4. MapReduce Design Patterns (Thursday 3/26)

Questions Batch Data Processing

http://piazza.com/harvard/spring2020/cs205/home



Backup

Ecosystem: Pig Latin and HiveQL

Key Components in the Hadoop Ecosystem

 Both ease the complexity of writing complex java MapReduce programs

PIG	HIVE
Procedural Data Flow Language	Declarative SQLish Language
For programming	For creating reports
Used by Researchers and Programmers	Used by Data Analysts
Operates on the client side of a cluster	Operates on the server side of a cluster.
Pig is SQL like but varies to a great extent.	Directly leverages SQL and is easy to learn for database experts.



Clustered Architecture





INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE AT HARVARD UNIVERSITY

The Hadoop Processing Framework Clustered Architecture





INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE AT HARVARD UNIVERSITY

The Hadoop Processing Framework Storage Architecture - HDFS





IACS INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE at HARVARD UNIVERSITY

Storage Architecture - HDFS

Scalable Fault-Tolerant Storage of Distributed Commodity Hardware

- Designed with hardware failure in mind
- •Built for large datasets, with a default block size of 128 MB
- Optimized for sequential operations
- Rack-aware
- Cross-platform and supports heterogeneous clusters

Block-oriented Storage

- Data in a Hadoop cluster is broken down into smaller units (called blocks) and distributed throughout the cluster
- Each block is duplicated twice (for a total of three copies), with the two replicas stored on two nodes in a rack somewhere else in the cluster
- If a copy is lost (because of machine failure, for example), HDFS will automatically re-replicate it elsewhere in the cluster, ensuring that the threefold replication factor is maintained



The Hadoop Processing Framework File Access - HDFS



NameNode

- Manages the file system namespace and associated metadata (file-to-block maps), and acts as the master and brokers access to files by clients (though once brokered, clients communicate directly with **DataNodes**).
- Operates entirely in memory, persisting its state to disk. It represents a single point of failure for a Hadoop cluster that is not running in high-availability mode.
- To mitigate against this, production clusters typically persist state to two local disks or install in high-availability mode with a standby **NameNode** to guard against failures



The Hadoop Processing Framework Storage Architecture - HDFS

http://ec2-54-175-139-110.compute-1.amazonaws.com:50070

Hadoop	Overview	Datanodes	Datanode Volume Failures	Snapshot	Startup Progress	Utilities -

Overview 'ip-10-2-1-142.ec2.internal:8020' (active)

Started:	Wed Mar 21 16:18:29 -0400 2018
Version:	2.8.3-amzn-0, rcfe28705e7dfdec92539cc7b24fc97936c259a05
Compiled:	Thu Feb 01 20:50:00 -0500 2018 by ec2-user from (HEAD detached at cfe28705e7)
Cluster ID:	CID-568f181d-feec-4d2c-895e-00804903fab0
Block Pool ID:	BP-1084522759-10.2.1.142-1521663504725

Summary

Security is off.

Safemode is off.

1.121 files and directories, 1.074 blocks = 2.195 total filesystem object(s).

Heap Memory used 214.45 MB of 490.5 MB Heap Memory. Max Heap Memory is 1.51 GB

Non Heap Memory used 65.85 MB of 67.25 MB Committed Non Heap Memory. Max Non Heap Memory is <u href="https://www.ubev.com">www.ubev.com</u> Non Heap Memory used 65.85 MB of 67.25 MB Committed Non Heap Memory. Max Non Heap Memory is <u href="https://www.ubev.com">www.ubev.com</u> Non Heap Memory. Max Non Heap Memory is <u href="https://www.ubev.com">www.ubev.com</u> Non Heap Memory. Max Non Heap Memory is <u href="https://www.ubev.com">www.ubev.com</u> Non Heap Memory. Max Non Heap Memory is <u href="https://www.ubev.com">www.ubev.com</u> Non Heap Memory. Max Non Heap Memory is <u href="https://www.ubev.com">www.ubev.com</u> Non Heap Memory. Max Non Heap Memory is </u>

Configured Capacity:	137.9 GB
DFS Used:	784.98 MB (0.56%)
Non DFS Used:	0 B
DFS Remaining:	137.13 GB (99.44%)
Block Pool Used:	784.98 MB (0.56%)
DataNodes usages% (Min/Median/Max/stdDev):	0.51% / 0.60% / 0.60% / 0.04%
Live Nodes	2 (Decommissioned: 0)
Dead Nodes	0 (Decommissioned: 0)
Decommissioning Nodes	0
Total Datanode Volume Failures	0 (0 B)
Number of Under-Replicated Blocks	0
Number of Blocks Pending Deletion	0
Block Deletion Start Time	Wed Mar 21 16:18:29 -0400 2018
Last Checkpoint Time	Wed Mar 21 16:18:24 -0400 2018



Lecture C1. Batch Data Processing CS205: Computing Foundations for Computational Science

Job Scheduling and Resource Manager - YARN

- ResourceManager (one per cluster): Takes inventory of available resources and runs the **Scheduler** which allocates resource containers to running applications
- NodeManager (one per node): Oversees resource containers and monitors their resource usage, periodically communicating with ResourceManager
- ApplicationMaster (one per app/job): Framework specific, runs in one of the containers, oversees application execution, monitors their resource usage periodically, and potentially communicates with **ResourceManager** to request additional containers



Sources: https://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/YARN.html



Application Execution - YARN



- 1. Client program submits the MapReduce application to the **ResourceManager**, along with information to launch the application-specific **ApplicationMaster**.
- 2. ResourceManager negotiates a container for the ApplicationMaster and launches the ApplicationMaster.
- **3. ApplicationMaster** boots and registers with the **ResourceManager**, allowing the original calling client to interface directly with the **ApplicationMaster**.
- 4. ApplicationMaster negotiates resources (resource containers) for client application.
- **5. ApplicationMaster** gives the container launch specification to the **NodeManager**, which launches a container for the application.
- 6. During execution, client polls **ApplicationMaster** for application status and progress.
- 7. Upon completion, **ApplicationMaster** deregisters with the **ResourceManager** and shuts down, returning its containers to the resource pool.





Application Execution - YARN

http://ec2-54-175-139-110.compute-1.amazonaws.com:8088



