

Data-Intensive Distributed Computing

CS 431/631 451/651 (Winter 2019)

Part 2: From MapReduce to Spark (1/2) January 22, 2019

Adam Roegiest Kira Systems

These slides are available at http://roegiest.com/bigdata-2019w/



This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States See http://creativecommons.org/licenses/by-nc-sa/3.0/us/ for details



Debugging at Scale

Works on small datasets, won't scale... why? Memory management issues (buffering and object creation) Too much intermediate data Mangled input records

Real-world data is messy!

There's no such thing as "consistent data" Watch out for corner cases Isolate unexpected behavior, bring local

The datacenter *is* the computer! What's the instruction set?

So you like programming in assembly?

202222222

Source: Wikipedia (ENIAC)

i Filler of

Hadoop is great, but it's really waaaaaay too low level! (circa 2007)

Source: Wikipedia (DeLorean time machine)

What's the solution?

Design a higher-level language Write a compiler Hadoop is great, but it's really waaaaaay too low level! (circa 2007)



What we really need is SQL!



What we really need is a scripting language!















Both open-source projects today!

Story for another day....

facebook.

Jeff Hammerbacher, Information Platforms and the Rise of the Data Scientist. In, *Beautiful Data*, O'Reilly, 2009.

> "On the first day of logging the Facebook clickstream, more than 400 gigabytes of data was collected. The load, index, and aggregation processes for this data set really taxed the Oracle data warehouse. Even after significant tuning, we were unable to aggregate a day of clickstream data in less than 24 hours."



Pig: Example

Task: Find the top 10 most visited pages in each category

Visits

URL Info

User	Url	Time	Url	Category	PageRank
Amy	cnn.com	8:00	cnn.com	News	0.9
Amy	bbc.com	10:00	bbc.com	News	0.8
Amy	flickr.com	10:05	flickr.com	Photos	0.7
Fred	cnn.com	12:00	espn.com	Sports	0.9
	•			•	

Pig: Example Script

```
visits = load '/data/visits' as (user, url, time);
```

```
gVisits = group visits by url;
```

```
visitCounts = foreach gVisits generate url, count(visits);
```

```
urlInfo = load '/data/urlInfo' as (url, category, pRank);
```

```
visitCounts = join visitCounts by url, urlInfo by url;
```

```
gCategories = group visitCounts by category;
```

```
topUrls = foreach gCategories generate top(visitCounts,10);
```

```
store topUrls into '/data/topUrls';
```

Pig Slides adapted from Olston et al. (SIGMOD 2008)

Pig Query Plan



Pig Slides adapted from Olston et al. (SIGMOD 2008)

Pig: MapReduce Execution



Pig Slides adapted from Olston et al. (SIGMOD 2008)

visits = load '/data/visits' as (user, url, time);

gVisits = group visits by url;

visitCounts = foreach gVisits generate url, count(visits);

urlinfo = load '/data/urlinfo' as (url, category, pRank);

visitCounts = join visitCounts by url, urlinfo by url;

gCategories = group visitCounts by category;

topUrls = foreach gCategories generate top(visitCounts,10);

This

store topUrls into '/data/topUrls';



ln.setflutnutKevClass(Text.class): lp.setOutputValueClass(Text.class); ln.setMannerClass(LoadPages.class); FileInputFormat.addInputPath(lp, new Path("/ user/gates/pages")); FileOutputFormat.setOutputPath(1p, new Path("/user/gates/tmp/indexed_pages")); lp.setNumReduceTasks(0); Job loadPages = new Job(1p); JobConf lfu = new JobConf(NRExample.class); lfu.s etJobName("Load and Filter Users"); lfu.setInputFormat(TextInputFormat.class); lfu.setUutputKeyClass(Text.class); lfu.setOutputValueClass(Text.class); lfu.setMapperClass(LoadAndFilterUsers.class); FileInputFormat.add InputPath(lfu, new Path("/user/mates/users")): FileOutputFormat.setOutputPath(lfu, new Path("/user/gates/tmp/filtered users")); lfu.setNumReduceTasks(0). Job loadUsers = new Job(lfu); JohConf join = new JohConf(MRExample.class); join.setJobName("Join Users and Pages"); join.setInputFormat(KeyValueTextInputFormat.class); join.setOutputKeyClass(Text.class); join.setOutputValueClass(Text.class); join.setMapperClass(IdentityMap per.class); join setReducerClass(Join class); FileInputFormat.addInputPath(join, new Path("/user/gates/tmp/indexed_pages")); FileInputFormat.addInputPath(join, new Path("/user/gates/tmp/filtered_users")); FileOutputFormat.se tOutputPath(join, new Path("/user/gates/tmp/joined")); join.setNumReduceTasks(50); Job joinJob = new Job(join); joinJob.addDependingJob(loadPages); joinJob.addDependingJob(loadUsers); JobConf group = new JobConf(MRE group.setJobName("Group URLs"); xample.class); group.setInputFormat(KeyValueTextInputFormat.class); group.setOutputKevClass(Text.class) group.setOutputValueClass(LongWritable.class); group.setOutputFormat(SequenceFi group.setMapperClass(LoadJoined.class); leOutputFormat.class); group.setCombinerClass(ReduceUrls.class); group.setReducerClass(ReduceUrls.class); FileInputFormat.addInputPath(group, new Path("/user/gates/tmp/joined")); FileOutputFormat.setOutputPath(group, new Path("/user/gates/tmp/grouped")); group.setNumReduceTasks(50); Job groupJob = new Job(group) groupJob.addDependingJob(joinJob); JobConf top100 = new JobConf(MRExample.class); top100.setJobName("Top 100 sites"); top100.setInputFormat(SequenceFileInputFormat.class); top100.setOutputKevClass(LongWritable.class); top100.setOutputValueClass(Text.class); top100.setOutputFormat(SequenceFileOutputF top100.setMapperClass(LoadClicks.class); top100.setCombinerClass(LimitClicks.class); ormat.class); top100.setReducerClass(LimitClicks.class); FileInputFormat.addInputPath(top100, new Path("/user/gates/tmp/grouped")); FileOutputFormat.setOutputPath(top100, new Path("/user/gates/top100sitesforusers18to25")); top100.setNumReduceTasks(1); Job limit = new Job(top100); limit.addDependingJob(groupJob);

But isn't Pig slower? Sure, but c can be slower than assembly too...



Pig: Basics

Sequence of statements manipulating relations (aliases)

Data model atoms tuples bags maps json

Pig: Common Operations

LOAD: load data (from HDFS) FOREACH ... GENERATE: per tuple processing FILTER: discard unwanted tuples "map" GROUP/COGROUP: group tuples "reduce" JOIN: relational join STORE: store data (to HDFS)

Pig: GROUPing

A = LOAD 'myfile.txt' AS (f1: int, f2: int, f3: int);

(1, 2, 3)
(4, 2, 1)
(8, 3, 4)
(4, 3, 3)
(7, 2, 5)
(8, 4, 3)

X = GROUP A BY f1;

(1, {(1, 2, 3)})
(4, {(4, 2, 1), (4, 3, 3)})
(7, {(7, 2, 5)})
(8, {(8, 3, 4), (8, 4, 3)})

Pig: COGROUPing

A:	В:
(1, 2, 3)	(2, 4)
(4, 2, 1)	(8 <i>,</i> 9)
(8, 3, 4)	(1, 3)
(4, 3, 3)	(2, 7)
(7, 2, 5)	(2, 9)
(8, 4, 3)	(4, 6)
	(4, 9)

X = COGROUP A BY\$0, B BY \$0;

(1, {(1, 2, 3)}, {(1, 3)})
(2, {}, {(2, 4), (2, 7), (2, 9)})
(4, {(4, 2, 1), (4, 3, 3)}, {(4, 6), (4, 9)})
(7, {(7, 2, 5)}, {})
(8, {(8, 3, 4), (8, 4, 3)}, {(8, 9)})

Pig: JOINing

A:	В:
(1, 2, 3)	(2, 4)
(4, 2, 1)	(8, 9)
(8, 3, 4)	(1, 3)
(4, 3, 3)	(2, 7)
(7, 2, 5)	(2, 9)
(8, 4, 3)	(4, 6)
	(4, 9)

X = JOIN A BY \$0, B BY \$0;

(1,2,3,1,3)
(4,2,1,4,6)
(4,3,3,4,6)
(4,2,1,4,9)
(4,3,3,4,9)
(8,3,4,8,9)
(8,4,3,8,9)

Pig UDFs

User-defined functions: Java Python JavaScript Ruby

...

UDFs make Pig arbitrarily extensible Express "core" computations in UDFs Take advantage of Pig as glue code for scale-out plumbing

The datacenter *is* the computer!

What's the instruction set? Okay, let's fix this!

Analogy: NAND Gates are universal



Let's design a data processing language "from scratch"!

What ops do you need?

(Why is MapReduce the way it is?)

Data-Parallel Dataflow Languages

We have a collection of records, want to apply a bunch of operations to compute some result

Assumption: static collection of records (what's the limitation here?)

We need per-record processing



Remarks: Easy to parallelize maps, record to "mapper" assignment is an implementation detail

Map alone isn't enough

(If we want more than embarrassingly parallel processing)

Where do intermediate results go? We need an addressing mechanism! What's the semantics of the group by?

Once we resolve the addressing, apply another computation That's what we call reduce! (What's with the sorting then?)

MapReduce



MapReduce is the minimally "interesting" dataflow!

MapReduce



(note we're abstracting the "data-parallel" part)

MapReduce Workflows



What's wrong?

Want MM?



Want MRR?



The datacenter *is* the computer!

Let's enrich the instruction set!

Source: Google

Dryad: Graph Operators



Dryad: Architecture



The Dryad system organization. The job manager (JM) consults the name server (NS) to discover the list of available computers. It maintains the job graph and schedules running vertices (V) as computers become available using the daemon (D) as a proxy. Vertices exchange data through files, TCP pipes, or shared-memory channels. The shaded bar indicates the vertices in the job that are currently running.

Dryad: Cool Tricks

Channel: abstraction for vertex-to-vertex communication File TCP pipe Shared memory

Runtime graph refinement

Size of input is not known until runtime Automatically rewrite graph based on invariant properties

Dryad: Sample Program



DryadLINQ

LINQ = Language INtegrated Query

.NET constructs for combining imperative and declarative programming

Developers write in DryadLINQ

Program compiled into computations that run on Dryad

Sound familiar?

Source: Yu et al. (2008) DryadLINQ: A System for General-Purpose Distributed Data-Parallel Computing Using a High-Level Language. OSDI.

What's the solution?

Design a higher-level language Write a compiler

DryadLINQ: Word Count

```
PartitionedTable<LineRecord> inputTable =
  PartitionedTable.Get<LineRecord>(uri);
```

```
IQueryable<string> words = inputTable.SelectMany(x => x.line.Split(' '));
IQueryable<IGrouping<string, string>> groups = words.GroupBy(x => x);
IQueryable<Pair> counts = groups.Select(x => new Pair(x.Key, x.Count()));
IQueryable<Pair> ordered = counts.OrderByDescending(x => x.Count);
IQueryable<Pair> top = ordered.Take(k);
```

Compare:

a = load 'file.txt' as (text: chararray);

b = foreach a generate flatten(TOKENIZE(text)) as term;

c = group b by term;

d = foreach c generate group as term, COUNT(b) as count;

store d into 'cnt';

Compare and contrast...

What happened to Dryad?



The Dryad system organization. The job manager (JM) consults the name server (NS) to discover the list of available computers. It maintains the job graph and schedules running vertices (V) as computers become available using the daemon (D) as a proxy. Vertices exchange data through files, TCP pipes, or shared-memory channels. The shaded bar indicates the vertices in the job that are currently running.

Data-Parallel Dataflow Languages

We have a collection of records, want to apply a bunch of operations to compute some result

What are the dataflow operators?

Spark

Answer to "What's beyond MapReduce?"

Brief history: Developed at UC Berkeley AMPLab in 2009 Open-sourced in 2010 Became top-level Apache project in February 2014 Commercial support provided by DataBricks

Spark vs. Hadoop



Google Trends

Source: Datanami (2014): http://www.datanami.com/2014/11/21/spark-just-passed-hadoop-popularity-web-heres/

What's an RDD? Resilient Distributed Dataset (RDD)

Much more next session...

MapReduce



Map-like Operations



Reduce-like Operations



Sort Operations



Join-like Operations



Join-like Operations



Set-ish Operations



Set-ish Operations



MapReduce in Spark?



Not quite...

MapReduce in Spark?



Still not quite...

Spark Word Count

Aside: Scala tuple access notation, e.g.,a._1

Don't focus on Java verbosity!

```
val textFile = sc.textFile(args.input())
```

```
textFile
.map(object mapper {
    def map(key: Long, value: Text) =
        tokenize(value).foreach(word => write(word, 1))
    })
.reduce(object reducer {
    def reduce(key: Text, values: Iterable[Int]) = {
        var sum = 0
        for (value <- values) sum += value
        write(key, sum)
    })
.saveAsTextFile(args.output())</pre>
```

Next Time...

What's an RDD? How does Spark actually work? Algorithm design: redux

Meanwhile, at 1600 Amphitheatre Parkway...

Sawzall – circa 2003 Lumberjack – circa ?? Flume(Java) – circa 2009 Cloud Dataflow (Flume + MillWheel) – circa 2014

Flume(Java)

Core data types

PCollection<T> - a (possibly huge) immutable bag of elements of type T PTable<K, V> - a (possibly huge) immutable bag of key-value pairs

Hmm... sounds suspiciously familiar...

Flume(Java)

Primitive operations



Hmm... looks suspiciously familiar...

Flume(Java) **Primitive operations**

PTable<K, V> PTable<URL,DocInfo> backlinks =



Hmm... looks suspiciously familiar...

Flume(Java) Primitive operations



Hmm... looks suspiciously familiar...

Data-Parallel Dataflow Languages

We have a collection of records, want to apply a bunch of operations to compute some result

Pig, Dryad(LINQ), Flume(Java), Spark are all variations on a theme!

Assumption: static collection of records What if this assumption is violated?

Remember: CS 451/651 Assignment 1 due 2:30pm Thursday, Jan 24 You must tell us if you wish to take the late penalty.

12 . 1

Source: Wikipedia (The Scream)