CS535  Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

**CS535 BIG DATA**

**PART 0. INTRODUCTION**
**3. DATA MODEL FOR BIG DATA**
**: APACHE THRIFT**

Sangmi Lee Pallickara
Computer Science, Colorado State University
http://www.cs.colostate.edu/~cs535

---

## FAQs

• Programming Assignment 1
  • Due Sept. 28
  • Submission via Canvas
  • Please check the course Web Page at least twice a week

---

## Today's topics

• Apache Thrift continued
• Introduction to MapReduce

---

## Thrift Architecture



You can change the protocol and transport without regenerating code

Source:
http://jnb.ociweb.com/jnb/jnbJun2009.html

---

## Protocol

---

## Protocols                    (1/2)

• Describes "what" is transmitted

• Thrift supports both text and binary protocols
  • The binary protocol outperforms text protocol
  • The text protocol may be useful (such as in debugging)

• `TBinaryProtocol`
  • A straight-forward binary format encoding numeric values as binary, rather than converting to text.

• `TCompactProtocol`
  • Very efficient, dense encoding of data

CS535  Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

---

8/30/2016 · CS535 Big Data - Fall 2016 · **W2.A.6**

## Protocols (1/2)

- `TDenseProtocol`
  - Similar to `TCompactProtocol` but strips off the meta information from what is transmitted, and adds it back in at the receiver.

- `TJSONProtocol`
  - Uses JSON for encoding of data.

- `TSimpleJSONProtocol`
  - A write-only protocol using JSON. It cannot be parsed by Thrift. Suitable for parsing by scripting languages

- `TDebugProtocol`
  - Uses a human-readable text format to aid in debugging.

---

8/30/2016 · CS535 Big Data - Fall 2016 · **W2.A.7**

## TCompactProtocol

- The `TCompactProtocol` is the most-efficient method in the Java implementation of Thrift

  - Writes numeric tags for each piece of data
  - The recipient is expected to properly match these tags with the data
  - If the data is not present, there simply is no tag/data pair

  | Tag | DATA | Tag | DATA | Tag | DATA |

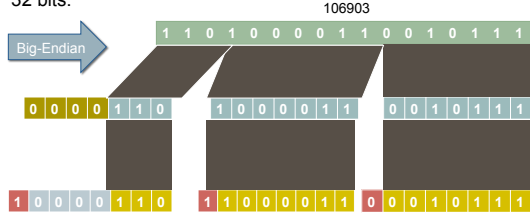---

8/30/2016 · CS535 Big Data - Fall 2016 · **W2.A.8**

## TCompactProtocol (1/2)

- For integers, the TCompactProtocol performs compression Variable-Length Quantity (VLQ) encoding

- VLQ uses 7 of 8 bits out of each byte for information
  - the 8th bit used as a continuation bit

---

8/30/2016 · CS535 Big Data - Fall 2016 · **W2.A.9**

## TCompactProtocol (2/2)

- Decimal value 106903 (0x1A197) saving 1 byte if it was stored in 32 bits:



VLQ's worst-case encoding
For a 32-bit int, the worst case is 5 bytes
For 64-bit ints, the worst case is 10 bytes

---

8/30/2016 · CS535 Big Data - Fall 2016 · **W2.A.10**

## Versioning

---

8/30/2016 · CS535 Big Data - Fall 2016 · **W2.A.11**

## Versioning

- **Applications** evolve over time
  - You added an extra field to your message format

- The system **must be able to support reading of old data** from log files
  - Process requests from out-of-date clients to new servers and vice versa

CS535  Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

---

**Slide W2.A.12**

8/30/2016 — CS535 Big Data - Fall 2016 — W2.A.12

## Versioning Field Identifiers

- Any new fields should be optional
  - Any messages serialized by code using the "old" message format **can be parsed** by the new code
    - Without any missing required fields (Parsing error)
  - Messages created by new code can be parsed by old code
    - Old binary will ignore the new field when parsing
    - Unknown fields are not discarded: it will be serialized along with other fields

- Non-required fields can be removed

- Changing a default value is allowed

- Mark Slee, Aditya Agarwal and Marc Kwiatkowski, Thrift: Scalable Cross-Language Services Implementation, https://thrift.apache.org/static/files/thrift-20070401.pdf

---

**Slide W2.A (title)**

**CS535 BIG DATA**

PART 1. BATCH COMPUTING MODELS FOR BIG DATA ANALYTICS
### 1. DISTRIBUTED MODEL FOR SCALABLE BATCH COMPUTING - MAPREDUCE

Sangmi Lee Pallickara
Computer Science, Colorado State University
http://www.cs.colostate.edu/~cs535

---

**Slide W2.A.14**

8/30/2016 — CS535 Big Data - Fall 2016 — W2.A.14

## This material is built based on

- Jeffrey Dean and Sanjay Ghemawat, "MapReduce:Simplified Data Processing on Large Clusters" In Proceeding OSDI'04 Proceedings of the 6th conference on Symposium on Opearting Systems Design & Implementation - Vol. 6

- Hadoop: The definitive Guide, Tom White, O'Reilly, 3rd Edition, 2014

- MapReduce Design Patterns, Donald Miner and Adam Shook, O'Reilly, 2013

- Anand Rajaraman, Jure Leskovec, and Jeffrey Ullman, "Mining of Massive Datasets", Cambridge University Press, 2012 -- Chapter 2

---

**Slide W2.A.15**

8/30/2016 — CS535 Big Data - Fall 2016 — W2.A.15

## Why MapReduce?

---

**Slide W2.A.16**

8/30/2016 — CS535 Big Data - Fall 2016 — W2.A.16

## What is MapReduce?

- MapReduce is inspired by the concepts of *map* and *reduce* in Lisp.

---

**Slide W2.A.17**

8/30/2016 — CS535 Big Data - Fall 2016 — W2.A.17

```
# Grammar for Python
# Note:  Changing the grammar specified in this file will most likely
#        require corresponding changes in the parser module
#        (../Modules/parsermodule.c).  If you can't make the changes to
#        that module yourself, please co-ordinate the required changes
#        with someone who can; ask around on python-dev for help.  Fred
#        Drake <fdrake@acm.org> will probably be listening there.

# NOTE WELL: You should also follow all the steps listed at
# https://docs.python.org/devguide/grammar.html

# Start symbols for the grammar:
#       single_input is a single interactive statement;
#       file_input is a module or sequence of commands read from an input file;
#       eval_input is the input for the eval() functions.
# NB: compound_stmt in single_input is followed by extra NEWLINE!
single_input: NEWLINE | simple_stmt | compound_stmt NEWLINE
file_input: (NEWLINE | stmt)* ENDMARKER
eval_input: testlist NEWLINE* ENDMARKER

stmt: simple_stmt | compound_stmt
simple_stmt: small_stmt (';' small_stmt)* [';'] NEWLINE
small_stmt: (expr_stmt | del_stmt | pass_stmt | flow_stmt |
             import_stmt | global_stmt | nonlocal_stmt | assert_stmt)
expr_stmt: testlist_star_expr (augassign (yield_expr|testlist) |
                     ('=' (yield_expr|testlist_star_expr))*)
testlist_star_expr: (test|star_expr) (',' (test|star_expr))* [',']
augassign: ('+=' | '-=' | '*=' | '@=' | '/=' | '%=' | '&=' | '|=' | '^=' |
            '<<=' | '>>=' | '**=' | '//=')
# For normal assignments, additional restrictions enforced by the
interpreter
del_stmt: 'del' exprlist
pass_stmt: 'pass'
flow_stmt: break_stmt | continue_stmt | return_stmt | raise_stmt |
yield_stmt
break_stmt: 'break'
continue_stmt: 'continue'
return_stmt: 'return' [testlist]
yield_stmt: yield_expr
raise_stmt: 'raise' [test ['from' test]]
import_stmt: import_name | import_from
import_name: 'import' dotted_as_names
# note below: the ('.' | '...') is necessary because '...' is tokenized as
ELLIPSIS
import_from: ('from' (('.' | '...')* dotted_name | ('.' | '...')+)
              'import' ('*' | '(' import_as_names ')' | import_as_names))
import_as_name: NAME ['as' NAME]
dotted_as_name: dotted_name ['as' NAME]
import_as_names: import_as_name (',' import_as_name)* [',']
dotted_as_names: dotted_as_name (',' dotted_as_name)*
dotted_name: NAME ('.' NAME)*
global_stmt: 'global' NAME (',' NAME)*
nonlocal_stmt: 'nonlocal' NAME (',' NAME)*
assert_stmt: 'assert' test [',' test]
```

**Python implementations**
- Cpython
- Brython
- CLPython
- HotPy
- IronPython
- Jython
- Pyjs
- PyPy
- SNAPpy …

Programming language specification: Python 3.5.2

CS535 Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

---

## What is MapReduce?

- Developed within Google as a mechanism for processing large amounts of raw data.
  - Crawled documents or web request logs
  - Distributes these data across thousands of machines
  - Same computations are performed on each CPU with a different portion of the dataset

---

## Why MapReduce?

- MapReduce provides an abstraction that allows engineers to perform simple computations while hiding the details of:
  - Parallelization
  - Data distribution
  - Load balancing
  - Fault tolerance

---

## Example

- Climate data from National Climatic Data Center (NCDC)

```
0057
332130 # USAF weather station identifier
99999 # WBAN weather station identifier
19500101 # observation date
0300 # observation time 4
+ 51317 # latitude (degrees x 1000)
+ 028783 # longitude (degrees x 1000)
FM-12
+ 0171 # elevation (meters)
99999
V020
320 # wind direction (degrees)
1 # quality code
```

---

## The first entries for 1990

```
% ls raw/ 1990 | head
010010-99999-1990.gz
010014-99999-1990.gz
010015-99999-1990.gz
010016-99999-1990.gz
010017-99999-1990.gz
010030-99999-1990.gz
010040-99999-1990.gz
010080-99999-1990.gz
010100-99999-1990.gz
010150-99999-1990.gz
```

---

## Analyzing the data with Unix Tools (1/2)

- A program for finding the **maximum recorded temperature by year** from NCDC weather records
  - e.g. Weather change for a century

```
#!/ usr/ bin/ env bash
for year in all/*
do
 echo -ne ` basename $ year .gz `"\ t"
 gunzip -c $ year | \
   awk '{
      temp = substr( $ 0, 88, 5) + 0;
      q = substr( $ 0, 93, 1);
      if (temp != 9999 && q ~ /[01459]/ && temp > max)
          max = temp
      }
      END { print max }'
Done
```

---

## Analyzing the data with Unix Tools (2/2)

- The script loops through the compressed year files
  - Printing the year
  - Processing each file using `awk`
    - Extracts two fields
    - Air temperature and the quality code
    - Check if it is greater than the maximum value seen so far

```
% ./ max_temperature.sh
1901 317
1902 244
1903 289
1904 256
1905 283
…
```

CS535 Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

---

## Results?

- The complete run for the century took 42 minutes
- To speed up the processing
  - We need to run parts of the program in parallel

  - Process different years in different processes
  - What will be the problems?

---

## Challenges

- **Dividing the work** into equal-size pieces
  - Data size per year?

- **Combining the results** from independent processes
  - Combining results and sorting by year?

---

## Map and Reduce

- MapReduce works by breaking the processing into two phases
  - The map phase
  - The reduce phase

- Each phase has key-value pairs as input and output

- Programmers should specify
  - Types of input/output key-values
  - The map function
  - The reduce function

---

## Visualizing the way MapReduce works　　　(1/4)

Sample lines of input data

```
0067011990999991950051507004... 9999999N9 + 00001 +99999999999...
0043011990999991950051512004... 9999999N9 + 00221 +99999999999...
0043011990999991950051518004... 9999999N9-00111 +99999999999...
0043012650999991949032412004... 0500001N9 + 01111 +99999999999...
0043012650999991949032418004... 0500001N9 + 00781 +99999999999…
```

These lines are presented to the map function as key-value pairs

```
(0, 0067011990999991950051507004...9999999N9 + 00001+
99999999999...)
(106, 0043011990999991950051512004...9999999N9 + 00221+
99999999999...)
(212, 0043011990999991950051518004...9999999N9−0011 1 +
99999999999...)
```
The keys are the line offsets within the file

---

## Visualizing the way MapReduce works　　　(2/4)

The **map function** extracts **the year and the air temperature** and emits them as its output

```
(1950, 0)
(1950, 22)
(1950, − 11)
(1949, 111)
(1949, 78)
```

This output key-value pairs will be **sorted and grouped by key**.
Our reduce function will see the following input:

```
(1949, [111, 78])
(1950, [0, 22, −11])
```
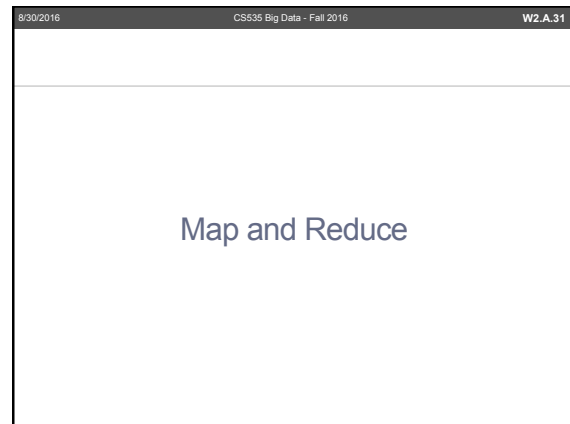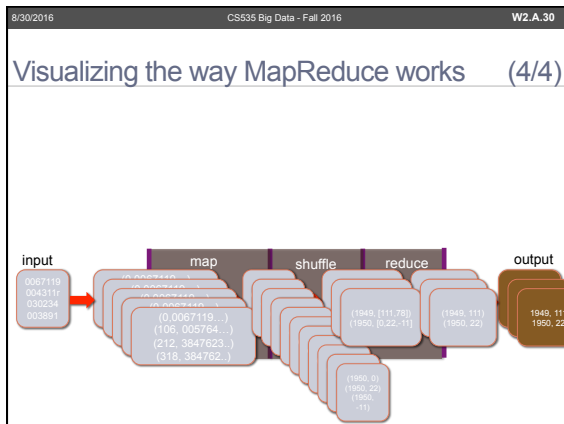
---

## Visualizing the way MapReduce works　　　(3/4)

Reduce function iterates through the list and picks the maximum reading

```
(1949, 111)
(1950, 22)
```

This is the final output



---

CS535 Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

---

8/30/2016     CS535 Big Data - Fall 2016     **W2.A.30**

## Visualizing the way MapReduce works    (4/4)



---

8/30/2016     CS535 Big Data - Fall 2016     **W2.A.31**

## Map and Reduce

---

8/30/2016     CS535 Big Data - Fall 2016     **W2.A.32**

## Map function

- Takes an input pair

- **Produces a set of intermediate key/value pairs**

- The MapReduce library groups together all intermediate values associated with the same intermediate key $i$ and passes them to the reduce function

---

8/30/2016     CS535 Big Data - Fall 2016     **W2.A.33**

## Reduce function

- **Accepts an intermediate key** $i$ and a set of values for that key

- **Merges** together these values
  - Forms a possibly smaller **set** of values

---

8/30/2016     CS535 Big Data - Fall 2016     **W2.A.34**

## Example: counting words

- Count the number of occurrences of each word in a large collection of documents:

```
map(String key, String value):
// key: document name
// value: document contents
for each word w in value:
    EmitIntermediate(w,"1");


reduce(String key, Iterator values):
//key: a word
//values: a list of counts
int result = 0;
for each v in values:
    result += ParseInt(v)
Emit(AsString(result));
```

---

8/30/2016     CS535 Big Data - Fall 2016     **W2.A.35**

## Comparison with other systems

- MPI vs. MapReduce
  - MapReduce tries to collocate the data with the compute node
  - Data access is fast
    - Data is local!

- Volunteer computing vs. MapReduce
  - SETI@home
    - Uses donated (volunteered) CPU time

تSorry, let me produce proper output.

CS535  Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

---

8/30/2016 · CS535 Big Data - Fall 2016 · W2.A.36

# MapReduce Data Flow

---

8/30/2016 · CS535 Big Data - Fall 2016 · W2.A.37



1. Shards the input files into M pieces
2. Starts up many copies of program.
3. Assigns work
8. Wake up the user program
7. Local write
(4) read
(5) Local write
6. Accesses the location notified by Master and perform reduce function

Input files · Map phase · Intermediate files (on local disks) · Reduce phase · Output files

4. Read contents of the corresponding input shard
Parses & passes the key-value pair to the Map function

5. Buffered pairs are written to local disk
Location is reported to the Master, which forwards them to appropriate reducer

---

8/30/2016 · CS535 Big Data - Fall 2016 · W2.A.38

# Job execution framework

- `mapred.job.tracker`
  - If the configuration property is set to `local`
    - Local job runner is used
    - Whole job will be in a single JVM
    - Designed for testing and running with small datasets
  - Using `classic`
    - For the classic MapReduce framework (MapReduce 1)
  - Using `yarn`
    - For the YARN framework

---

8/30/2016 · CS535 Big Data - Fall 2016 · W2.A.39

# Job submission

- Single method call
  - `submit()` on a `Job` object
    - You can also call `waitForCompletion()`
- Creates an internal `JobSubmitter` instance
  - Calls `submitJobInternal()` on it
- `waitForCompletion()`
  - Polls the job's progress once per second
  - Reports the progress to the console if it has changed since the last report
- Job completes successfully
  - The job counters are displayed
  - Otherwise, the error is logged to the console

---

8/30/2016 · CS535 Big Data - Fall 2016 · W2.A.40

# MapReduce job in the Classic framework



---

8/30/2016 · CS535 Big Data - Fall 2016 · W2.A.41

# Job submission process in Classic framework

- Asks the jobtracker for a new jobID
  - Step 2
- Checks the output specification of the job
- Computes the input splits for the job
- Copies the resources to the jobtracker's filesystem in a directory named after the jobID
  - The job JAR file, the configuration file, and the computed input splits
- Tells the jobtracker that the job is ready for execution by calling `sumitJob()` on `JobTracker`
  - Step 4

CS535  Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

---

## Job Initialization in the Classic mode

- JobTracker puts the job into an internal queue
  - Creating an object to represent the job being run
  - Bookkeeping information to keep track of the status and progress of its tasks

- Creating the list of tasks
  - The job scheduler retrieves the input splits
  - The job scheduler creates a map task for each split
  - The number of reduce tasks to create
    - mapred.reduce.tasks property in the Job
      - setNumReduceTasks()
  - A job setup task and a job cleanup task are created

---

## Task assignment

- Jobtracker select a job to run
  - Default algorithm
    - Based on the priority list of jobs

- Tasktrackers have a fixed number of slots for map tasks and reduce tasks
  - These are set independently
  - These are selected based on the number of cores and the memory
- The default scheduler fills empty map task slots first
  - Before it fills the reduce task slots
  - Scheduling the reduce task does not need to consider data locality

- Tasktrackers run a simple loop that periodically sends heartbeat method calls to the jobtracker

---

## Task execution

- Tasktracker
  - Copies the job Jar to the tasktracker's file system
    - Any files needed from the distributed cache to the local file system
  - Creates a local working directory
  - Creates instance of TaskRunner

- TaskRunner
  - Launches a new Java Virtual Machine to run each task
    - Any failed map or reduce does not affect the tasktracker

---

## Streaming and pipes

- Runs special map and reduce tasks
  - To launch the user supplied executable
  - To communicate with it

- Streaming task communicates with the process using standard input and output streams

- Pipes task listens on a socket and passes the C++ process a port number in its environment
  - On startup, the C++ process establish a persistent socket connection back to the parent Java Pipes task

---

## Progress and status updates                    [1/3]

- A Job and each of its tasks have a status
  - State of the job or task
    - E.g. running, successfully complete, failed
  - The progress of maps and reduces
  - The values of the job's counters
  - A status message or description

- Progress of a task
  - The proportion of the task completed
  - Map task
    - The proportion of the input that has been processed
  - Reduce task
    - Divides the total progress into 3 parts (copy/sort/reduce)
    - If the task has run the reducer on half its input
      - 1/3 (copy) + 1/3 (sort) + a half of 1/3(reduce phase) = 5/6

---

## Progress and status updates                    [2/3]

- Tasks
  - have a set of counters
  - Count various events at the task run
  - E.g. the number of map output records written
  - If a task reports progress
    - it sets a flag to indicate that the status change should be sent to the tasktracker
      - Checked every 3 seconds

- Tasktracker
  - Tasks notify the current task status to the tasktracker
    - if the flag is set
  - Tasktracker sends heartbeats to the jobtracker every 5 seconds (minimum)

---

CS535  Big Data
Fall 2016 Colorado State University
http://www.cs.colostate.edu/~cs535

8/30/2016
Sangmi Pallickara
Week 2 - A

## Progress and status updates        [2/3]

- Jobtracker
  - Combines updates to produce a global view

- Job
  - Receives the latest status by polling the jobtracker every second
  - Prints job statistics and counters to the console

## YARN (MapReduce 2)

- To provide the scalability to MapReduce
  - Splitting responsibility of the jobtracker
    - Scheduling
    - Task progress monitoring

- MapReduce is one type of YARN application