

CSE 486/586 Distributed Systems Data Analytics

Steve Ko
Computer Sciences and Engineering
University at Buffalo

CSE 486/586

Recap

- RPC enables programmers to call functions in remote processes.
- IDL (Interface Definition Language) allows programmers to define remote procedure calls.
- Stubs are used to make it appear that the call is local.
- Semantics
 - Cannot provide exactly once
 - At least once
 - At most once
 - Depends on the application requirements

CSE 486/586

2

Two Questions We'll Answer

- What is **data analytics**?
- What are the **programming paradigms** for it?

CSE 486/586

3

Example 1: Scientific Data

- CERN (European Organization for Nuclear Research) @ Geneva: Large Hadron Collider (LHC) Experiment
 - 300 GB of data per second
 - “15 petabytes (15 million gigabytes) of data annually – enough to fill more than 1.7 million dual-layer DVDs a year”



Example 2: Web Data

- Google
 - 20+ billion web pages
 - » ~20KB each = 400 TB
 - ~ 4 months to read the web
 - And growing...
 - » 1999 vs. 2009: ~ 100X
- Yahoo!
 - US Library of Congress every day (20TB/day)
 - 2 billion photos
 - 2 billion mail + messenger sent per day
 - And growing...

Google

YAHOO!

CSE 486/586

5

Data Analytics

- Computations on very large data sets
 - How large? TBs to PBs
 - Much time is spent on data moving/reading/writing
- Shift of focus
 - Used to be: **computation** (think supercomputers)
 - Now: **data**



CSE 486/586

6

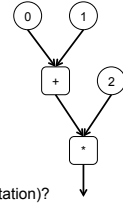
Popular Environment

- Environment for storing TBs ~ PBs of data
- **Cluster** of cheap **commodity PCs**
 - As we have been discussing in class...
 - 1000s of servers
 - Data stored as plain files on file systems
 - Data scattered over the servers
 - Failure is the norm
- How do you process all this data?



Turn to History

- Dataflow programming
 - Data sources and operations
 - Data items go through a series of transformations using operations.
 - Very popular concept
- Many examples
 - Even CPU designs back in 80's and 90's
 - SQL, data streaming, etc.
- Challenges
 - How to efficiently fetch data?
 - When and how to schedule different operations?
 - What if there's a failure (both for data and computation)?



CSE 486/586

8

Dataflow Programming

- This style of programming is now very popular with large clusters.
- Many examples
 - MapReduce, Pig, Hive, Dryad, Spark, etc.
- Two examples we'll look at
 - MapReduce and Pig



CSE 486/586

9

What is MapReduce?



- A system for processing large amounts of data
- Introduced by Google in 2004
- Inspired by **map & reduce in Lisp**
- OpenSource implementation: Hadoop by Yahoo!
- Used by many, many companies
 - A9.com, AOL, Facebook, The New York Times, Last.fm, Baidu.com, Joost, Veoh, etc.

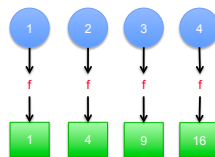
CSE 486/586

10

Background: Map & Reduce in Lisp

- Sum of squares of a list (in Lisp)
- **(map square '(1 2 3 4))**
 - Output: (1 4 9 16)

[processes each record individually]



CSE 486/586

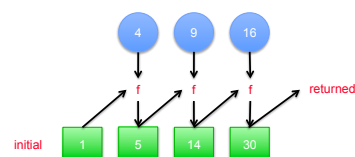
11

Background: Map & Reduce in Lisp

- Sum of squares of a list (in Lisp)
- **(reduce + '(1 4 9 16))**
 - (+ 16 (+ 9 (+ 4 1)))

– Output: 30

[processes set of all records in a batch]



CSE 486/586

12

Background: Map & Reduce in Lisp

- Map
 - processes each record individually
- Reduce
 - processes (combines) set of all records in a batch

CSE 486/586

13

What Google People Have Noticed

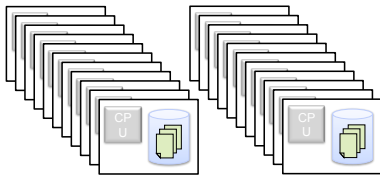
- Keyword search
 - Map** – Find a keyword in each web page **individually**, and if it is found, return the URL of the web page
 - Reduce** – **Combine** all results (URLs) and return it
- Count of the # of occurrences of each word
 - Map** – Count the # of occurrences in each web page **individually**, and return the list of <word, #>
 - Reduce** – For each word, **sum up (combine)** the count
- Notice the similarities?

CSE 486/586

14

What Google People Have Noticed

- Lots of storage + compute cycles nearby
- Opportunity
 - Files are **distributed already!** (GFS)
 - A machine can process its own web pages (**map**)



CSE 486/586

15

Google MapReduce

- Took the **concept** from Lisp, and applied to large-scale data-processing
- Takes two functions from a programmer (**map** and **reduce**), and performs three steps
- **Map**
 - Runs **map** for each file **individually in parallel**
- **Shuffle**
 - Collects the output from all **map** executions
 - **Transforms the map output into the reduce input**
 - Divides the **map** output into chunks
- **Reduce**
 - Runs **reduce (using a map output chunk as the input)** in parallel

CSE 486/586

16

Programmer's Point of View

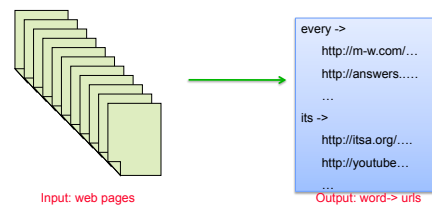
- Programmer writes two functions – **map()** and **reduce()**
- The programming interface is fixed
 - map (in_key, in_value) -> list of (out_key, intermediate_value)
 - reduce (out_key, list of intermediate_value) -> (out_key, out_value)
- **Caution: not exactly the same as Lisp**

CSE 486/586

17

Inverted Indexing Example

- Word -> list of web pages containing the word

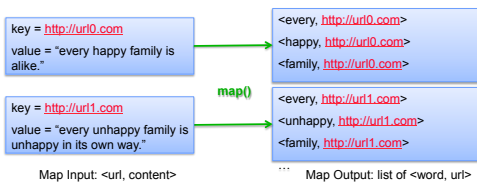


CSE 486/586

18

Map

- Interface
 - Input: $\langle \text{in_key}, \text{in_value} \rangle$ pair \Rightarrow $\langle \text{url}, \text{content} \rangle$
 - Output: list of intermediate $\langle \text{key}, \text{value} \rangle$ pairs \Rightarrow list of $\langle \text{word}, \text{url} \rangle$



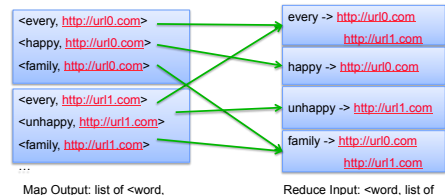
Map Input: $\langle \text{url}, \text{content} \rangle$ Map Output: list of $\langle \text{word}, \text{url} \rangle$

CSE 486/586

19

Shuffle

- MapReduce system
 - Collects outputs from all *map* executions
 - Groups all intermediate values by the same key



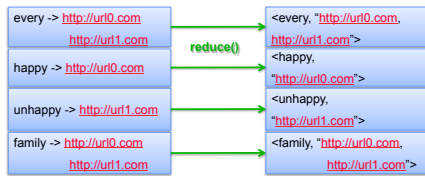
Map Output: list of $\langle \text{word}, \text{url} \rangle$ Reduce Input: $\langle \text{word}, \text{list of urls} \rangle$

CSE 486/586

20

Reduce

- Interface
 - Input: $\langle \text{out_key}, \text{list of intermediate_value} \rangle$
 - Output: $\langle \text{out_key}, \text{out_value} \rangle$

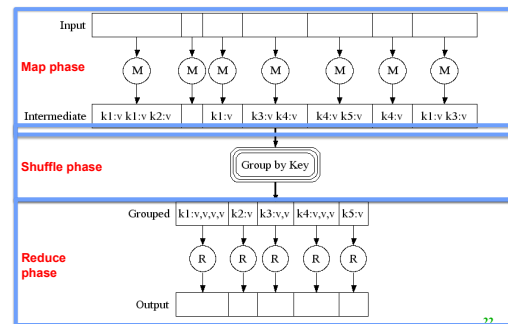


Reduce Input: $\langle \text{word}, \text{list of urls} \rangle$ Reduce Output: $\langle \text{word}, \text{string of urls} \rangle$

CSE 486/586

21

Execution Overview



22

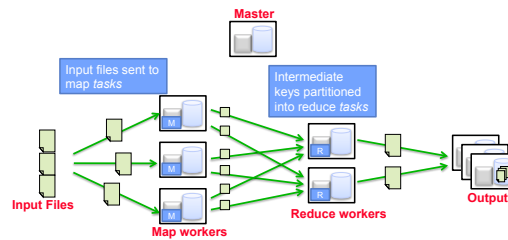
Implementing MapReduce

- Externally for user
 - Write a map function, and a reduce function
 - Submit a job; wait for result
 - No need to know anything about the environment (Google: 4000 servers + 48000 disks, many failures)
- Internally for MapReduce system designer
 - Run map in parallel
 - Shuffle: combine map results to produce reduce input
 - Run reduce in parallel
 - Deal with failures

CSE 486/586

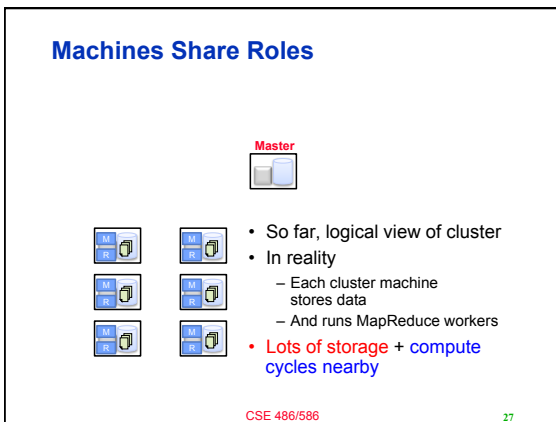
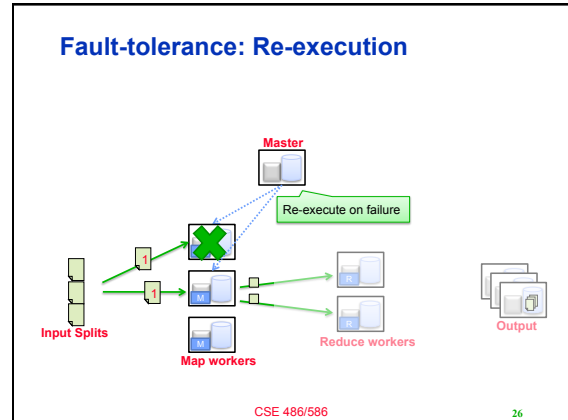
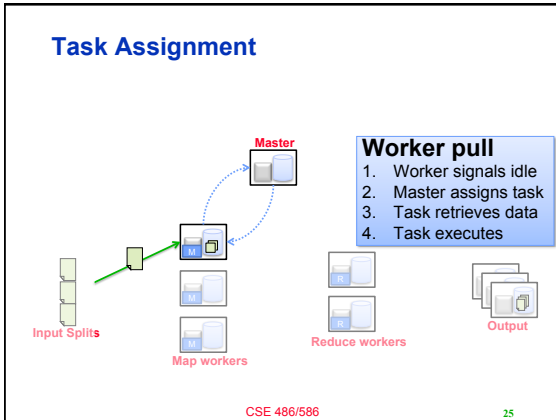
23

Execution Overview




CSE 486/586

24



- ### Problems of MapReduce
- Any you can think of?
 - There's only two functions you can work with (not expressive enough sometimes.)
 - Functional-style (a barrier for some people)
 - Turing completeness (or computationally universal)
 - If it can simulate a single-taped Turing machine.
 - Most general languages (C/C++, Java, Lisp, etc.) are.
 - SQL is.
 - MapReduce is not.
- CSE 486/586 28

- ### Pig
- 
- Why Pig?
 - MapReduce has limitations: only two functions
 - Many tasks require more than one MapReduce
 - Functional thinking: barrier for some
 - Pig
 - Defines a set of high-level simple "commands"
 - Compiles the commands and generates multiple MapReduce jobs
 - Runs them in parallel
- CSE 486/586 29

Pig Example

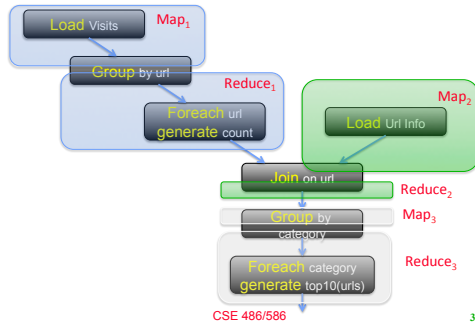
```
load '/data/visits';
group visits by url;
foreach gVisits generate url, count(visits);

load '/data/urlInfo';
join visitCounts by url, urlInfo by url;

group visitCounts by category;
foreach gCategories generate top(visitCounts,10);
```

CSE 486/586 30

Pig Example



Summary

- Data analytics shifts the focus from computation to data.
- Many programming paradigms are emerging.
 - MapReduce
 - Pig
 - Many others

CSE 486/586

32

More Details

- Papers
 - J. Dean and S. Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters," OSDI 2004
 - C. Olston, B. Reed, U. Srivastava, R. Kumar, and A. Tomkins, "Pig Latin: A Not-So-Foreign Language For Data Processing," SIGMOD 2008
- URLs
 - <http://hadoop.apache.org/core/>
 - <http://wiki.apache.org/hadoop/>
 - <http://hadoop.apache.org/pig/>
 - <http://wiki.apache.org/pig/>
- Slides
 - <http://labs.google.com/papers/mapreduce-osdi04-slides/index.html>
 - <http://www.systems.ethz.ch/education/past-courses/hs08/map-reduce/slides/intro.pdf>
 - http://www.cs.uiuc.edu/class/sp09/cs525/L4tmp_B.ppt
 - <http://infolab.stanford.edu/~usriv/talks/sigmod08-pig-latin.ppt>

CSE 486/586

33

Acknowledgements

- These slides contain material developed and copyrighted by Indranil Gupta (UIUC).

CSE 486/586

34