

“Civilization advances by extending the number of important operations which we can perform without thinking about them”

**Alfred North Whitehead, Professor at Harvard,
1910s**

Hands-on H4

MapReduce Design Patterns

CS205: Computing Foundations for Computational Science
Dr. David Sondak
Spring Term 2021



**INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE**
AT HARVARD UNIVERSITY



HARVARD
School of Engineering
and Applied Sciences

Lectures adapted from Ignacio M. Llorente

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

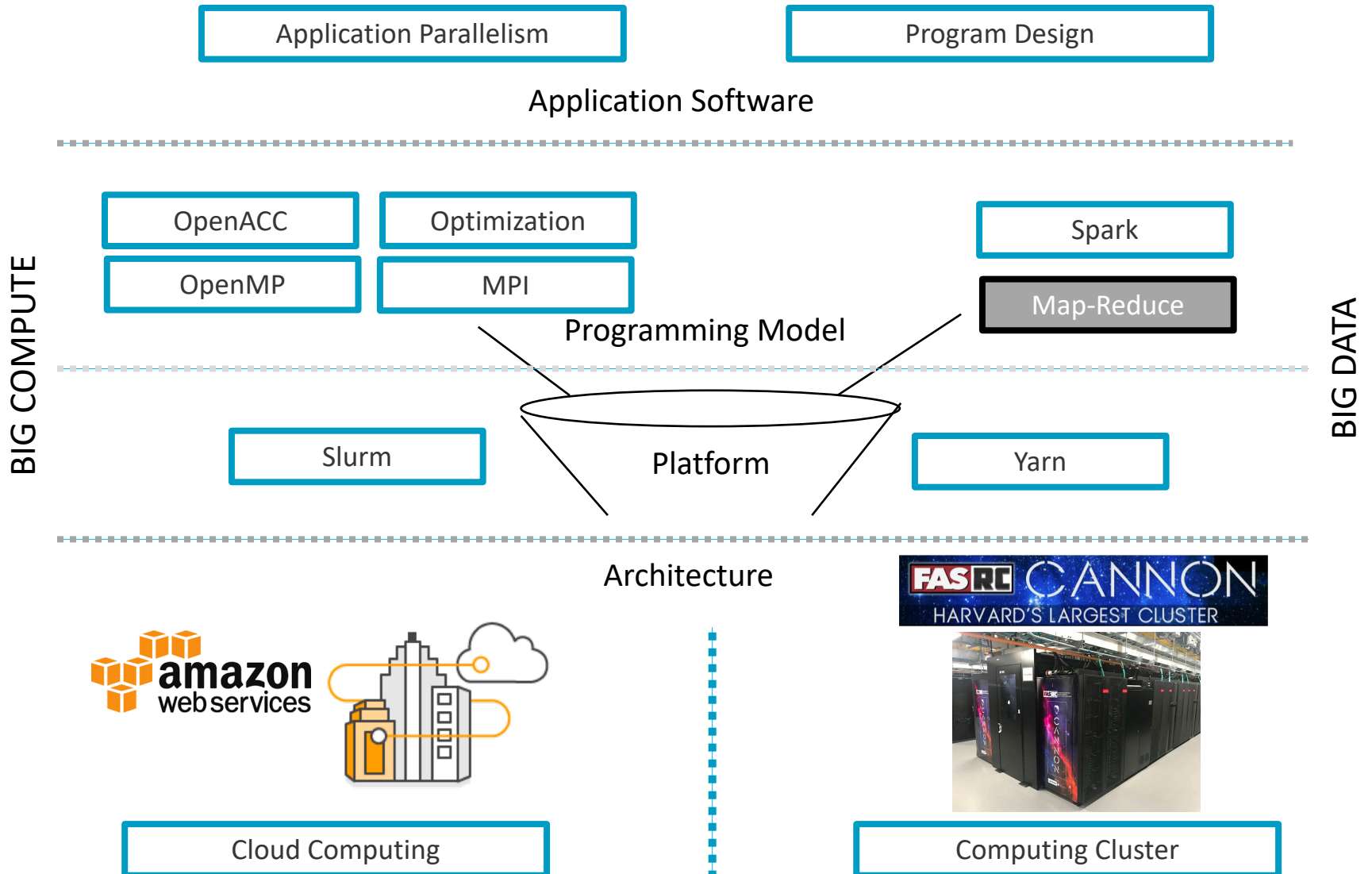
C2. Dataflow Processing

C3. Stream Data Processing

Wrap-Up: Advanced Topics

CS205: Contents

APPLICATION SOFTWARE



Before We Start

Where We Are



Week 9: Batch Data Processing => MapReduce

3/22	3/23 <u>Hands-on H4</u> MapReduce Programming	<u>Lab I8</u> Hadoop	3/25 <u>Lecture C2</u> Dataflow Processing (Quiz & Reading)	3/26
-------------	---	-------------------------	--	-------------

Week 10: Dataflow Processing => Spark

3/29	3/30 <u>Hands-on H5</u> Spark Programming	<u>Lab I9</u> Spark Single Node	4/1 <u>Lecture C3</u> Stream Data Processing	4/2
-------------	---	---------------------------------------	--	------------

Context

MapReduce Programming Model

The programmer essentially only specifies two (sequential) functions

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

- Inputs data record and 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 **fully independent (no execution ordering, sync or comm)**

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

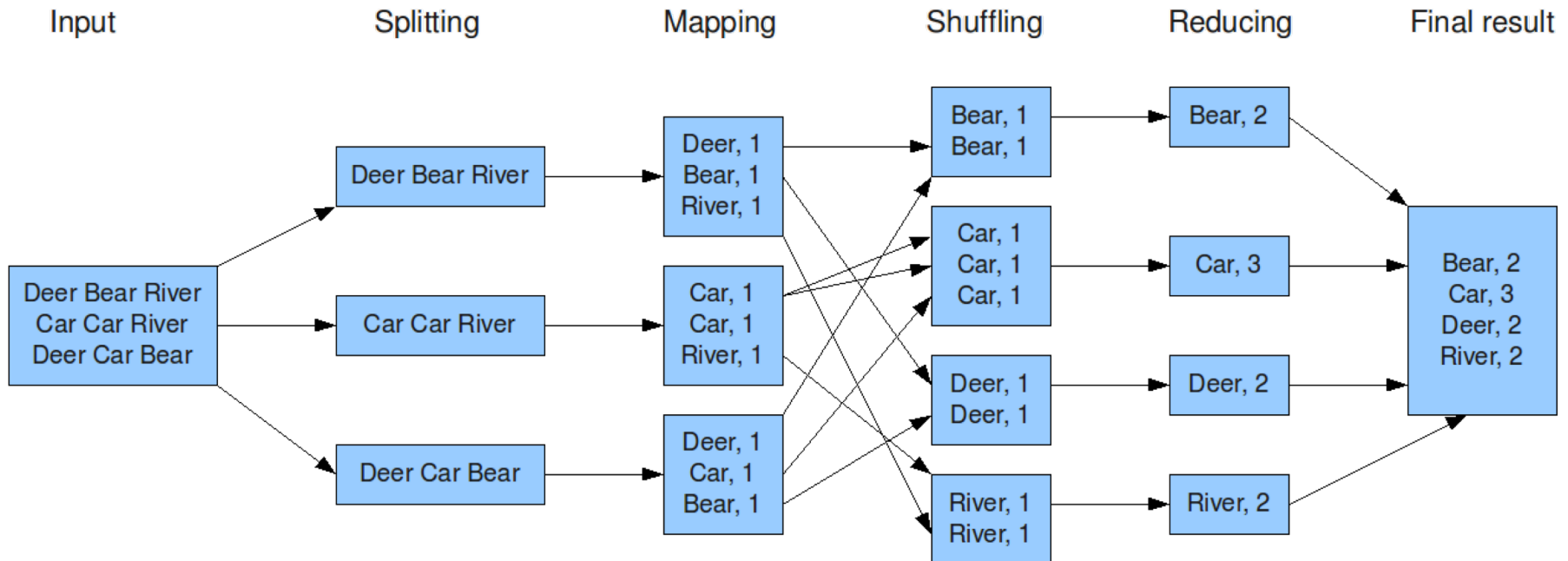
STEP 3. REDUCE: $reduce(k2, list(v2)) \rightarrow list(k3, v3)$

- Combines information across records that share this same intermediate key
- Each reduce call is **fully independent**

Context

MapReduce Programming Model

The overall MapReduce word count process



Hands-on Examples

Requirements



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. Unix-like shell (Linux, Mac OS or Windows/Cygwin)
2. Python installed

Roadmap

MapReduce Design Patterns

Design Patterns

Summarization

Inverted Index

Filtering

Other Patterns

Design Patterns

What Are Design Patterns?

- ✓ Reusable solutions to problems (HWC!)
- ✓ Domain independent
- ✓ Not a cookbook
- ✓ Not a guide
- ✓ Not a finished solution

Design Patterns

Why Design Patterns?

- ✓ Makes the intent of model and platform easier to understand
- ✓ Provides a common language for solutions
- ✓ Be able to reuse code
- ✓ Describes known performance profiles and limitations of solutions

Design Patterns

When Should I Use MapReduce?

Query

- Index and Search: inverted index
- Filtering
- Classification

Analytics

- Summarization and statistics
- Sorting and merging
- Frequency distribution
- SQL-based queries: group-by, having, etc.
- Generation of graphics: histograms, scatter plots.

... **large datasets in off-line mode** for boosting other on-line processes

Design Patterns

Main Functions and Patterns

Main Patterns

1. Summarization
2. Inverted Index
3. Filtering

Summarization

Calculating Aggregate Statistical Values

Description

- A general pattern for calculating aggregate statistical values over your data

Intent

- Group records together by a key field and calculate a numerical aggregate per group to get a top-level view of the larger data set

Examples

1. Word count
2. Record count
3. Min/Max/Count
4. Average/Median/Standard deviation
5. ...

Summarization

Word Count

Find the frequency of each word in text files

- Map: Process lines and generate as output $\langle \text{word}, 1 \rangle$
- Reduce: Add all values for the same *word*

map

input: [*line* of text file]

for each *word*

output: $\langle \text{word}, 1 \rangle$



reduce

input: [$\langle \text{word}, 1 \rangle$]

count for same *word*

output: $\langle \text{word}, \text{sum} \rangle$

Summarization

Word Count

mapper.py

```
#!/usr/bin/python
import sys
import re

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

    for word in words:
        print( word.lower() + "\t1" )
```

reducer.py

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

previous = None
summ = 0

for line in sys.stdin:
    key, value = line.split( '\t' )

    if key != previous:
        if previous is not None:
            print(str( summ ) + '\t' + previous)
        previous = key
        summ = 0

    summ = summ + int( value )

print(str( summ ) + '\t' + previous)
```


Summarization

Record Count

Find the frequency of each URL in web logs

- Map: Process web page access logs and generate $\langle URL, 1 \rangle$ as output
- Reduce: Add all values for the same URL

```
64.242.88.10 - - [07/Mar/2004:16:37:27 -0800] "GET /twiki/bin/view/TWiki/DontNotify HTTP/1.1" 200 4140
64.242.88.10 - - [07/Mar/2004:16:39:24 -0800] "GET /twiki/bin/view/Main/TokyoOffice HTTP/1.1" 200 3853
...
```

map

input: [*line* of log file]
for each line with a URL
output: $\langle URL, 1 \rangle$



reduce

input: $\langle URL, 1 \rangle$
Count for same URL
output: $\langle URL, \# \rangle$

Summarization

Max-Min

Given a list of tweets determine first and last time a user commented and the number of times.

- Data is a set of lines $\langle \textit{username}, \textit{date}, \textit{text} \rangle$

```
Peter [07/Mar/2020:16:39:24 -0800] "Stay at home"
```

```
John [07/Mar/2020:16:39:25 -0800] "Me too"
```

...

map

input: $\langle \textit{username}, \textit{date}, \textit{text} \rangle$

for each line

output: $\langle \textit{username}, \textit{date}, 1 \rangle$



reduce

input: $\langle \textit{username}, \textit{date}, 1 \rangle$

First, Last and Count for same *username*

output: $\langle \textit{username}, \textit{first_date}, \textit{last_date} \rangle$

Summarization

Average

Find average daily gains in stock for each company

- Data is a set of lines $\langle date, company, start_price, end_price \rangle$
- This example is for company from 1/1/2000 – 12/31/2015

Date , Company , Open , Close

2009-01-02 , Alphabet , 153 . 302917 , 159 . 870193

...

map

input: $\langle date, company, start_price, end_price \rangle$
output: $\langle company, end_price - start_price \rangle$



reduce

input: $\langle company, end_price - start_price \rangle$
Average for same *company*
output: $\langle company, average \rangle$

Inverted Index

Mapping Content to Location

Description

- A general pattern for mapping content, such as words or numbers, to its locations in a database file or in a document or a set of documents

Intent

- Most of the text searching systems rely on inverted index to search for documents that contain a given word or a term

Inverted Index

Word to Documents

Find what documents contain a specific word

- Map: Parse document and generate $\langle word, doc_id \rangle$ pairs
- Reduce: For each word, sort the corresponding document IDs



```
all id_432, id_76
also id_432
...
```

map

input: [*line* from document
doc_id]

for each *word*

output: $\langle word, doc_id \rangle$



reduce

input: $\langle word, doc_id \rangle$

concatenate for same *word*

output: $\langle word, [doc_ids] \rangle$

Hands-on

Word to Documents – Inverted Index

- ✓ Before beginning, run `file_name_ii.py` and write the output to a file
 - ✓ The result is a file that has the file name as the first item in each line
- ✓ Implement word to documents
 - ✓ Adapt `mapper` and `reducer` from `wordcount`
 - ✓ Run it with the output file that you generated in the first step

Inverted Index

Reverse Web-link Graph

Find where page links come from

- Map: Output $\langle \text{target}, \text{source} \rangle$ for each link to target in a page source
- Reduce: Concatenate the list of all source URLs associated with a target



map

input: [*line of HTML file*
URL_source]

for each *URL_target*

output: $\langle \text{URL_target},$
URL_source \rangle



reduce

input: [$\langle \text{URL_target},$
URL_source \rangle]

concatenate for same
URL_target

output: $\langle \text{URL_target},$
[*URL_sources*] \rangle

Filtering

Filtering Out Records

Description

- It evaluates each record separately and decides, based on some condition, whether it should stay or go

Intent

- Filter out records that are not of interest and keep ones that are.

Examples

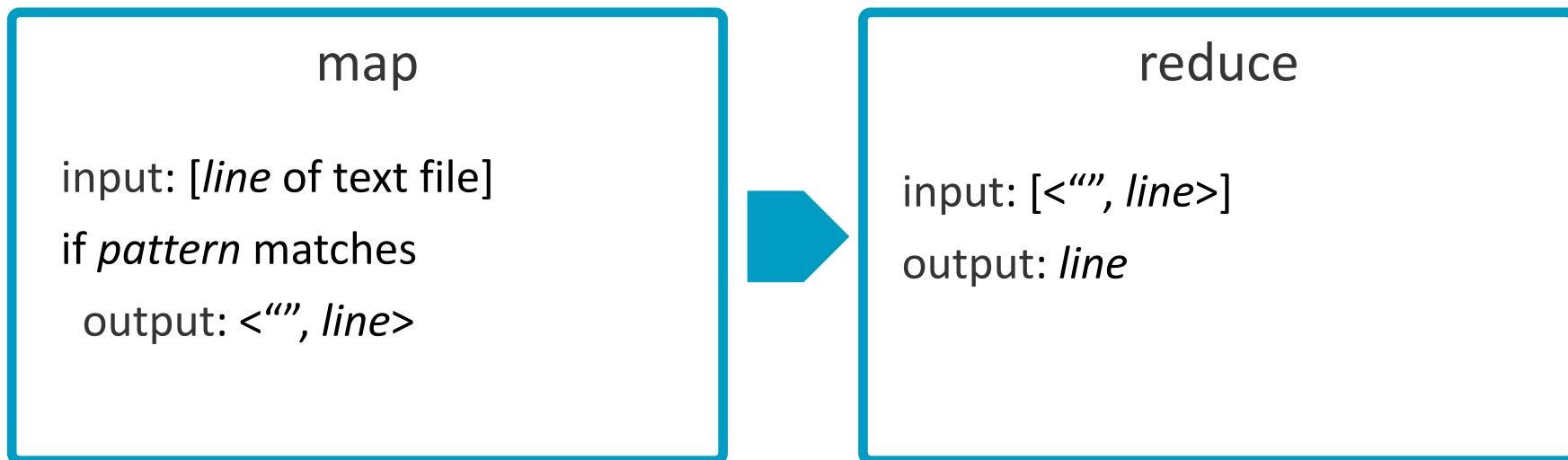
1. Closer view of dataset
2. Data cleansing
3. Tracking a thread of events
4. Simple random sampling
5. Distributed Grep
6. Removing low scoring dataset
7. Log Analysis
8. Data Querying and Validation
- 9....

Filtering

Distributed Grep

Search for words in a document

- Map: Generate a line if it matches a given *pattern*
- Reduce: Just copy the intermediate data to the output



Other Patterns

Organization, Join and Input/Output



- ✓ Summarization patterns: Get a top-level view by summarizing and grouping data
- ✓ Filtering patterns: View data subsets such as records generated from one user
- ✓ Data organization patterns: Reorganize data to work with other systems, or to make MapReduce analysis easier
- ✓ Join patterns: Analyze different datasets together to discover interesting relationships
- ✓ Metapatterns: Piece together several patterns to solve multi-stage problems, or to perform several analytics in the same job
- ✓ Input and output patterns: Customize the way you use Hadoop to load or store data

Next Steps

- Get ready for next lecture:
 C2. Dataflow Processing (Thursday 3/25)
- Project proposal presentations in two weeks!

Questions

MapReduce Design Patterns

