

Lecture 4: MapReduce (continue)

COSC 526: Introduction to Data Mining



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE
KNOXVILLE

BIG ORANGE. BIG IDEAS.®

Today Outline

- More about MR: portable and efficient implementations
- Text manipulation
- Introduction in XSEDE Jetstream
- Overview of Spark core concepts
- Create SparkContext and Resilient distributed datasets (RDDs)
- Run parallel operations on RDDs
- Live chat

**One problems,
many solutions**

WordCount++: A solution

```
1: class MAPPER
2:     method MAP(string  $t$ , integer  $r$ )
3:         EMIT(string  $t$ , integer  $r$ )

1: class REDUCER
2:     method REDUCE(string  $t$ , integers  $[r_1, r_2, \dots]$ )
3:          $sum \leftarrow 0$ 
4:          $cnt \leftarrow 0$ 
5:         for all integer  $r \in$  integers  $[r_1, r_2, \dots]$  do
6:              $sum \leftarrow sum + r$ 
7:              $cnt \leftarrow cnt + 1$ 
8:          $r_{avg} \leftarrow sum / cnt$ 
9:         EMIT(string  $t$ , integer  $r_{avg}$ )
```

Is this an efficient solution?

```
1: class MAPPER
2:     method MAP(string  $t$ , integer  $r$ )
3:         EMIT(string  $t$ , integer  $r$ )

1: class REDUCER
2:     method REDUCE(string  $t$ , integers  $[r_1, r_2, \dots]$ )
3:          $sum \leftarrow 0$ 
4:          $cnt \leftarrow 0$ 
5:         for all integer  $r \in$  integers  $[r_1, r_2, \dots]$  do
6:              $sum \leftarrow sum + r$ 
7:              $cnt \leftarrow cnt + 1$ 
8:          $r_{avg} \leftarrow sum / cnt$ 
9:         EMIT(string  $t$ , integer  $r_{avg}$ )
```

A more efficient solution?

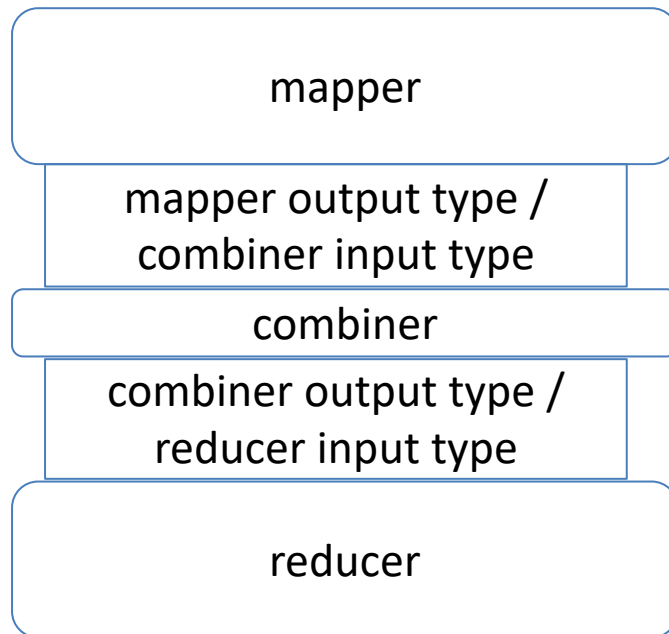
```
1: class MAPPER
2:     method MAP(string  $t$ , integer  $r$ )
3:         EMIT(string  $t$ , integer  $r$ )

1: class COMBINER
2:     method COMBINE(string  $t$ , integers  $[r_1, r_2, \dots]$ )
3:          $sum \leftarrow 0$ 
4:          $cnt \leftarrow 0$ 
5:         for all integer  $r \in$  integers  $[r_1, r_2, \dots]$  do
6:              $sum \leftarrow sum + r$ 
7:              $cnt \leftarrow cnt + 1$ 
8:         EMIT(string  $t$ , pair ( $sum, cnt$ ))

1: class REDUCER
2:     method REDUCE(string  $t$ , pairs  $[(s_1, c_1), (s_2, c_2) \dots]$ )
3:          $sum \leftarrow 0$ 
4:          $cnt \leftarrow 0$ 
5:         for all pair  $(s, c) \in$  pairs  $[(s_1, c_1), (s_2, c_2) \dots]$  do
6:              $sum \leftarrow sum + s$ 
7:              $cnt \leftarrow cnt + c$ 
8:          $r_{avg} \leftarrow sum / cnt$ 
9:         EMIT(string  $t$ , integer  $r_{avg}$ )
```

Combiners Constraints and Portability

- Combiners in, for example, Apache Hadoop cannot change the correctness of the MapReduce algorithm
- Combiners must have same input and output key-value types



mapper output type
==
combiner input type
==
combiner output type
==
reducer input type

Does the solution meet the constraints?

```
1: class MAPPER
2:   method MAP(string  $t$ , integer  $r$ )
3:     EMIT(string  $t$ , integer  $r$ )

1: class COMBINER
2:   method COMBINE(string  $t$ , integers  $[r_1, r_2, \dots]$ )
3:      $sum \leftarrow 0$ 
4:      $cnt \leftarrow 0$ 
5:     for all integer  $r \in$  integers  $[r_1, r_2, \dots]$  do
6:        $sum \leftarrow sum + r$ 
7:        $cnt \leftarrow cnt + 1$ 
8:     EMIT(string  $t$ , pair ( $sum$ ,  $cnt$ ))

1: class REDUCER
2:   method REDUCE(string  $t$ , pairs  $[(s_1, c_1), (s_2, c_2) \dots]$ )
3:      $sum \leftarrow 0$ 
4:      $cnt \leftarrow 0$ 
5:     for all pair  $(s, c) \in$  pairs  $[(s_1, c_1), (s_2, c_2) \dots]$  do
6:        $sum \leftarrow sum + s$ 
7:        $cnt \leftarrow cnt + c$ 
8:      $r_{avg} \leftarrow sum / cnt$ 
9:     EMIT(string  $t$ , integer  $r_{avg}$ )
```


Does the solution meet the constraints?

```
1: class MAPPER
2:   method MAP(string t, integer r)
3:     EMIT(string t, integer r)

1: class COMBINER
2:   method COMBINE(string t, integers [r1, r2, ...])
3:     sum ← 0
4:     cnt ← 0
5:     for all integer r ∈ integers [r1, r2, ...] do
6:       sum ← sum + r
7:       cnt ← cnt + 1
8:     EMIT(string t, pair (sum, cnt))

1: class REDUCER
2:   method REDUCE(string t, pairs [(s1, c1), (s2, c2) ...])
```

Mismatch between combiner input key-value type and output key-value type violates the MapReduce programming model!!!

A portable and efficient solution

```
1: class MAPPER
2:     method MAP(string  $t$ , integer  $r$ )
3:         EMIT(string  $t$ , pair ( $r$ , 1))

1: class COMBINER
2:     method COMBINE(string  $t$ , pairs  $[(s_1, c_1), (s_2, c_2) \dots]$ )
3:          $sum \leftarrow 0$ 
4:          $cnt \leftarrow 0$ 
5:         for all pair  $(s, c) \in$  pairs  $[(s_1, c_1), (s_2, c_2) \dots]$  do
6:              $sum \leftarrow sum + s$ 
7:              $cnt \leftarrow cnt + c$ 
8:         EMIT(string  $t$ , pair ( $sum$ ,  $cnt$ ))

1: class REDUCER
2:     method REDUCE(string  $t$ , pairs  $[(s_1, c_1), (s_2, c_2) \dots]$ )
3:          $sum \leftarrow 0$ 
4:          $cnt \leftarrow 0$ 
5:         for all pair  $(s, c) \in$  pairs  $[(s_1, c_1), (s_2, c_2) \dots]$  do
6:              $sum \leftarrow sum + s$ 
7:              $cnt \leftarrow cnt + c$ 
8:          $r_{avg} \leftarrow sum / cnt$ 
9:         EMIT(string  $t$ , integer  $r_{avg}$ )
```

Assignment 4

Assignment 4

- Goal: Continue building our expertise with Jupyter and Python
 - Sequential manipulation of a classical in literature
 - Visualization of statistics
- Deadline: **Friday, Feb 19 - 8AM ET**

Assignment 4: Text manipulation

- Given a literature classic such as the “The Count of Monte Cristo”
- Problem 1: Analyze the text for word length frequency
- Problem 2: Analyze the text for letter frequency
- Problem 3: Count the positional frequencies of each letter (first, interior, and last)
- Problem 4: Visualize your findings in histograms (one for each one of Problems 1-3)
 - One code is give to you, write the other two codes
- Repeat with a different manuscript
 - Run your code with e.g., a manuscript written in a different language or an old manuscript

*Deadline: **Feb 19 - 8AM ET***



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

BIG ORANGE. BIG IDEAS.®