

# Introduction to Data Science

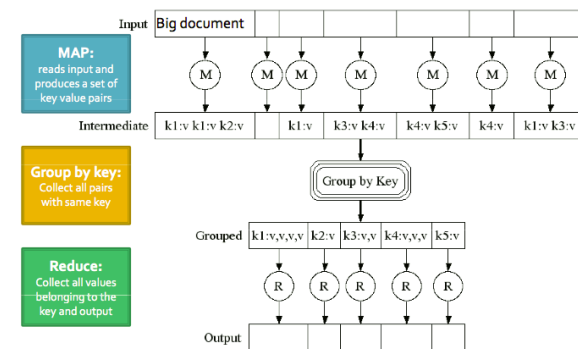
CS 5665  
Utah State University  
Department of Computer Science  
Instructor: Prof. Kyumin Lee

## MapReduce

### What do you do?

- Define two functions:
  - map  $(k, v) \rightarrow \langle k', v' \rangle^*$
  - reduce  $(k', \langle v' \rangle) \rightarrow \langle k', v'' \rangle^*$
- All  $v'$  with the same  $k'$  are reduced together and processed in  $v'$  order

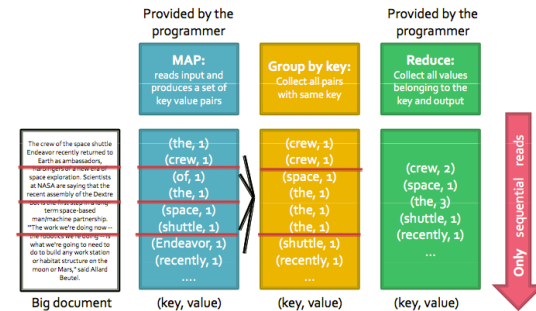
### MapReduce: A diagram



## MapReduce: Word counting

- Program specifies two primary methods
  - $\text{Map}(k,v) \rightarrow \langle k', v' \rangle^*$
  - $\text{Reduce}(k', \langle v' \rangle^*) \rightarrow \langle k', v'' \rangle^*$
- Ex: Two documents - (d1, "the crew") and (d2, "of the")
  - $\text{Map}(d1, \text{"the crew"}) \Rightarrow [(\text{the}, 1), (\text{crew}, 1)]$
  - $\text{Map}(d2, \text{"of the"}) \Rightarrow [(\text{of}, 1), (\text{the}, 1)]$
  - MapReduce runs its grouper module and calls reduce for every key
  - $\text{Reduce}(\text{the}, [1, 1]) \Rightarrow (\text{the}, 2)$
  - $\text{Reduce}(\text{crew}, [1]) \Rightarrow (\text{crew}, 1)$
  - $\text{Reduce}(\text{of}, [1]) \Rightarrow (\text{of}, 1)$

## MapReduce: Word counting



## Coordination

- Master data structures
  - Task status: (idle, in-progress, completed)
  - Idle tasks get scheduled as workers become available
  - When a map task completes, it sends the master the location and sizes of its R intermediate files, one for each reducer
  - Master pushes this info to reducers
- Master pings workers periodically to detect failures

## Failures

- Map worker failure
  - Map tasks completed or in-progress at worker are reset to idle
  - Reduce workers are notified when task is rescheduled on another worker
- Reduce worker failure
  - Only in-progress tasks are reset to idle
- Master failure
  - MapReduce task is aborted and client is notified

## Example 1: Host size

- Suppose we have a large web corpus
- Let's look at the metadata file
  - Lines of the form (URL, size, date, ...)
- For each host, find the total number of bytes
  - i.e., the sum of the page sizes for all URLs from that host
- Solution:
  - Map(line\_no, (URL, size, date, ...)) => [(URL\_host, size)]
  - Reduce(URL\_host, [size1, size2,...]) => (URL\_host, sum([size1, size2,...]))