



MapReduce Programming Model Designed to operate on LARGE distributed input data sets stored e.g. in HDFS nodes Abstracts from parallelism, data distribution, load balancing, data transfer, fault tolerance Implemented in Hadoop and other frameworks Provides a high-level parallel programming construct (= a skeleton) called MapReduce A generalization of the data-parallel MapReduce skeleton of Lect. 1 Covers the following algorithmic design pattern: Map Phase Shuffle Phase Reduce Phase



R Reducer (+shuffle) tasks

M Mapper (+combiner) tasks









Mapper	
Applies a user-defined function to each element (i.e., key/value pair coming from the Record reader).	
• Examples:	
 Filter function – drop elements that do not fulfill a constraint 	
 Transformation function – calculation on each element 	
 Produces a list of zero or more new key/value pairs = intermediate elements 	
 Key in K₂: index for grouping of data 	

- Value in V₂: Data to be forwarded to reducer
- Buffered in memory

Combiner

- An optional local reducer run in the mapper task as postprocessor
- Applies a user-provided function to aggregate values in the intermediate elements of one mapper task

I.U LINKOP

I.U UNKÖPIN

- Reduction/aggregation could also be done by the reducer, but local reduction can improve performance considerably
 - Data locality key/value pairs still in cache resp. memory of same node
 - Data reduction aggregated information is often smaller
- Applicable if the user-defined Reduce function is commutative and associative
- Recommended if there is significant repetition of intermediate keys produced by each Mapper task

Partitioner

- Splits the intermediate elements from the mapper/combiner into shards (64MB blocks stored in local files)
 - one shard per reducer
 - Default: element to hashCode(element.key) modulo R for even (round-robin) distribution of elements

Usually good for load balancing

Writes the shards to the local file system

Shuffle-and-sort

- Downloads the needed files written by the partitioners to the node on which the reducer is running
- Sort the received (key,value) pairs by key into one list
 - · Pairs with equivalent keys will now be next to each other (groups)
 - To be handled by the reducer
- No customization here beyond how to sort and group by keys

Reducer

- Run a user-defined reduce function once per key grouping
 - Can aggregate, filter, and combine data
 - Output: 0 or more key/value pairs sent to output formatter.

Output Formatter Translates the final (key,value) pair from the reduce function and writes it to stdout -> to a file in HDFS Default formatting (key <TAB> value <NEWLINE>) can be customized C.Kester, DA, LinkOpings universe.











Further Examples for MapReduce MapReduce Implementation / Execution Flow Count URL frequencies (a variant of wordcount) User application calls MapReduce and waits. MapReduce library implementation **splits** the input data (if not already done) in *M* blocks (of e.g. 64MB) and **creates** (*P* MapReduce processes on different cluster nodes: 1 master and *P*-1 workers. Input: logs of web page requests < URL, 1> Reduce function adds together all values for same URL Master creates *M* mapper tasks and *R* reducer tasks, and dispatches them to idle workers (dynamic scheduling) Construct reverse web-link graph Input: <sourceURL, targetURL> pairs Worker executing a **Mapper** task reads its block of input, applies the *Map* (and local *Combine*) function, and buffers (key, value) pairs in memory. Buffered pairs are periodically written to local disk, locations of these files are sent to Master. • Mapper reverses: <targetURL, sourceURL> Shuffle-and-sort → <targetURL, list of all URLs pointing to targetURL> Worker executing a **Reducer** task is notified by Master about locations of intermediate data to shuffle+sort and **fetches** them by remote memory access request, then **sorts** them by key (K₂). It applies the *Reduce function* to the sorted data and appends its output to a local file. no reduction → Reduce function is identity function Indexing web documents Input: list of documents (e.g. web pages) Mapper parses documents and builds sequences <word, documentID> When all mapper and reducer tasks have completed, the master wakes up the user program and returns the locations of the *R* output files. Shuffle-and-sort produces for each word a list of all documentIDs when

word occurs (Reduce function is identity)





Run-time

schedule

...

Worker processes

Ş

I.U UNKÖPIN MapReduce Implementation: Granularity

Numbers M, R and work of tasks (block size) might be tuned

- Default: M = input file size / block size • User can set other value
- M, R should be >> P
 - · For flexibility in dynamic load balancing
 - Hadoop recommends ~10...100 mappers per cluster node, or more if lightweight
- Not too large, though...
 - ~ M+R scheduling decisions by master
 - Block size should be reasonably large (e.g. 64MB) to keep relative impact of communication and task overhead low

Questions for Reflection

- A MapReduce computation should process 12.8 TB of data in a distributed file with block (shard) size 64MB. How many mapper tasks will be created, by default? (Hint: 1 TB (Terabyte) = 10^{12} byte)
- Discuss the design decision to offer just one MapReduce construct that covers both mapping, shuffle+sort and reducing. Wouldn't it be easier to provide one separate construct for each phase? What would be the performance implications of such a design operating on distributed files?
- Reformulate the wordcount example program to use no Combiner.
- Consider the local reduction performed by a Combiner: Why should the user-defined Reduce function be associative and commutative? Give examples for reduce functions that are associative and commutative, and such that are not.
- Extend the wordcount program to discard words shorter than 4 characters.
- Write a wordcount program to only count *all* words of odd and of even length. There are several possibilities.
- Show how to calculate a database join with MapReduce.
- Sometimes, workers might be temporarily slowed down (e.g. repeated disk read errors) without being broken. Such workers could delay the completion of an entire MapReduce computation considerably. How could the master speed up the overall MapReduce processing if it observes that some worker is late?

References

- J. Dean, S. Ghemawat: MapReduce: Simplified Data Processing on Large Clusters. Proc. OSDI 2004. Also in: Communications of the ACM 51(1), 2008.
- D. Miner, A. Shook: MapReduce Design Patterns. O'Reilly, 2012.

Apache Hadoop: https://hadoop.apache.org

 $\langle \rangle$ Worker processes on different cluster node

Run-time

schedule

###