



# Exercise 5

Information Retrieval

## 7. Implementing IR-Systems II

### Index Construction

## Exercise 7.1

- Are the following statements true or false? Give reasons for your answer.
  - a) External sorting algorithms are used to sort lists which do not fit in main memory.
  - b) The term-at-a-time approach to query answering requires more run-time memory than the document-at-a-time approach.
  - c) It is impossible to calculate the top-k documents for a query without completely reading the postings list of each query term.
  - d) When distributing a large index across a cluster of computers, we can partition the index by term or by document
  - e) In IR, we do not consider the case of evolving indexes, since document collections (almost) never change.

# Efficient Top-k for One-Term-Queries

## Exercise 7.2

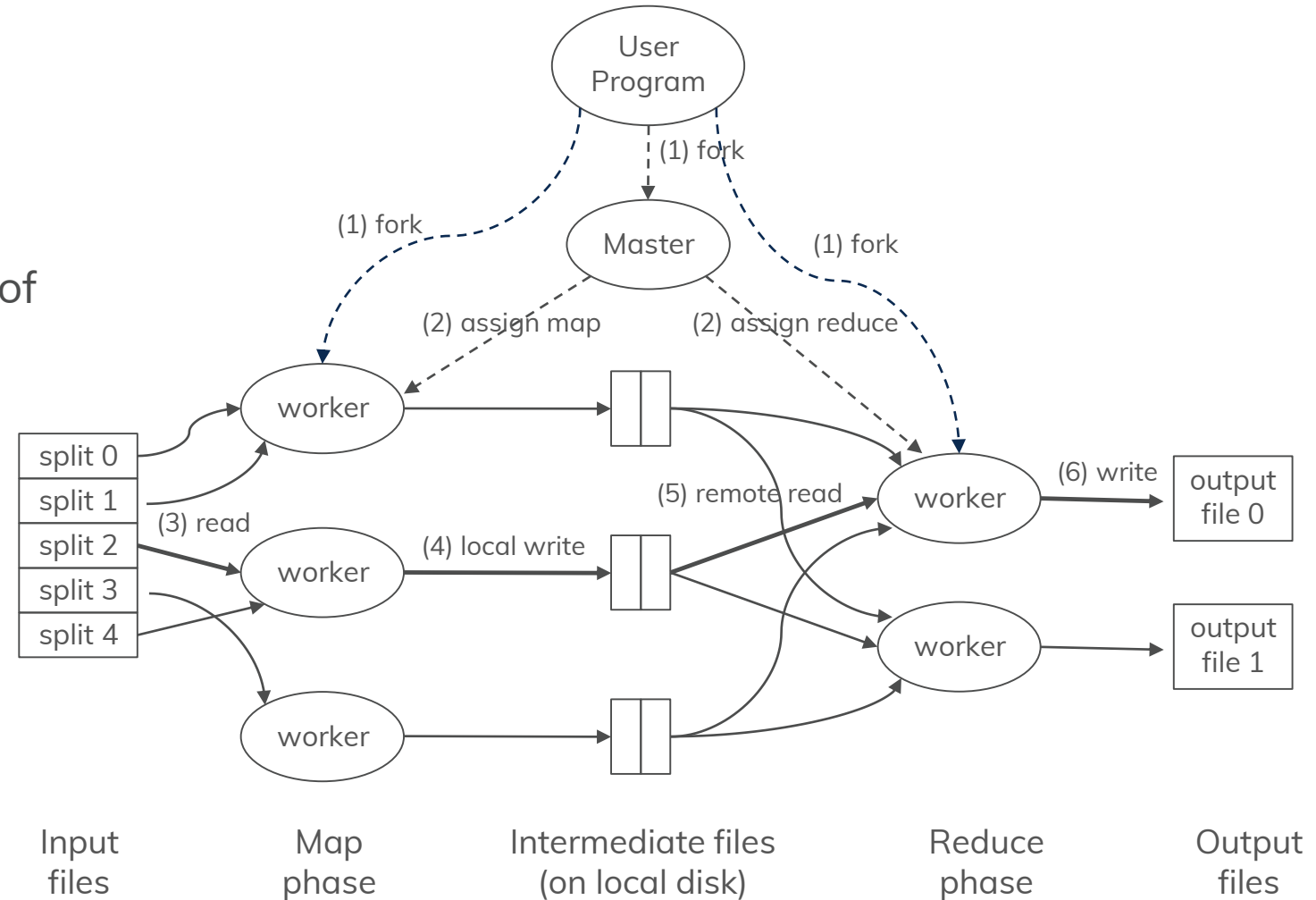
Exercise 0.99

- Assume
  - we only have one-term queries
  - we order postings lists according to weight, instead of  $\text{docId}$
  - we truncate every postings list at position  $k$
  - the weight  $w$  stored for doc.  $d$  in the postings list of  $t$  is the cosine-normalized weight of  $t$  for  $d$
- Explain why this approach suffices for identifying the  $k$  highest scoring documents for a query

# Distributed Indexing: MapReduce

## Key Facts

- MapReduce is a programming model
- Programmer implements a `map` function and a `reduce` function
- Underlying runtime system takes care of
  - Parallelization across large-scale cluster
  - Machine failures
  - Communication
  - Performance
  - etc.
- see lecture slides for details



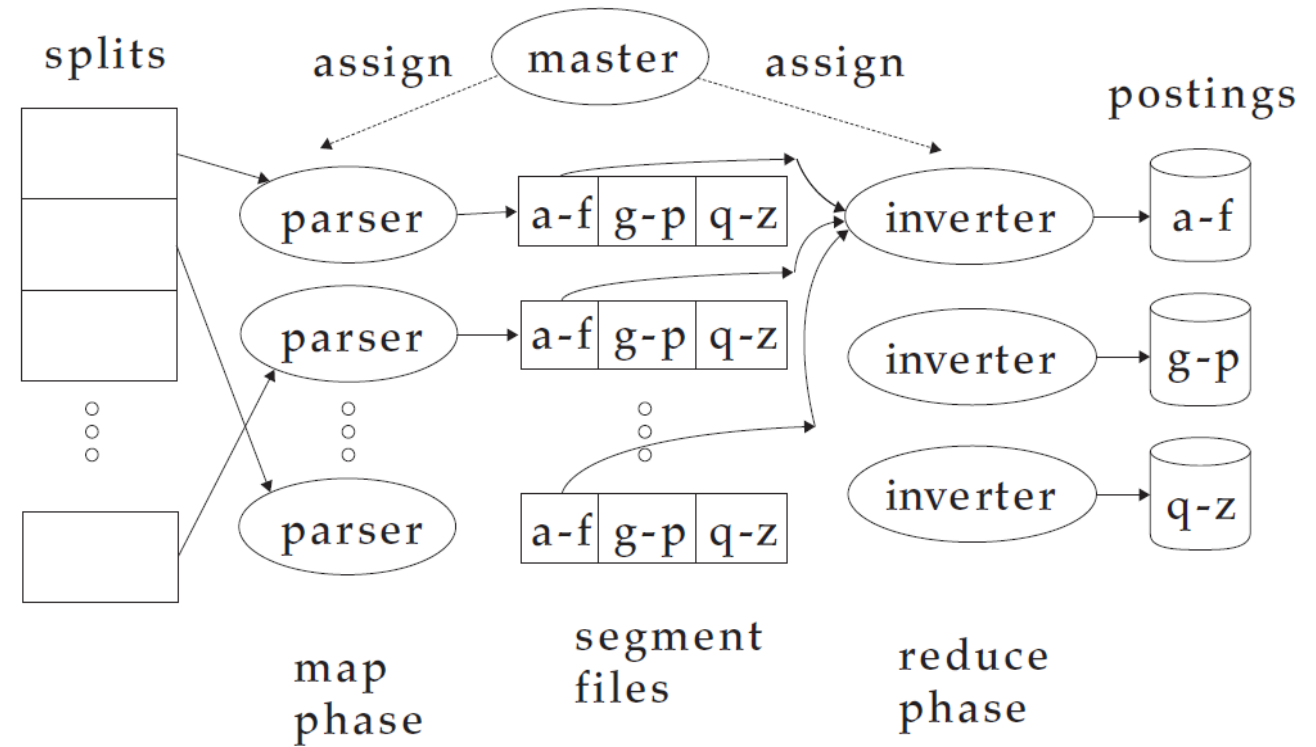
# Index Construction in MapReduce

## Parsers

- Master assigns a split to an idle parser machine
- Parser reads a document at a time and emits (term, docID)-pairs
- Parser writes pairs into  $j$  term-partitions
- Each for a range of terms' first letters
- E.g., a-f, g-p, q-z (here:  $j = 3$ )

## Inverters

- Inverter collects all (term, docID)-pairs for one term-partition (e.g., for a-f)
- Sorts and writes to postings lists



# map and reduce functions

## Schema of *map* and *reduce* functions

- *map*:  $\text{input} \rightarrow \text{list}(k, v)$
- *reduce*:  $\text{list}(k, \text{list}(v)) \rightarrow \text{output}$

## Example: Index Construction

- Parser (*map*):  $\text{web collection} \rightarrow \text{list}(\text{term}, \text{docID})$
- Inverter (*reduce*):  $\text{list}(\langle \text{term}, \text{list}(\text{docID}) \rangle) \rightarrow \text{postings\_list1}, \text{postings\_list2}, \dots$

## Exercise 7.3

- Apply MapReduce to the problem of counting how often each term occurs in a set of files by specifying *map* and *reduce* operations for this task

Exercise 0.53

Exercise 0.51

## Exercise 7.4

- Assume that machines in MapReduce have 100 GB of disk space each
- Assume further that the postings list of the term `the` has a size of 200 GB
- Why can the MapReduce algorithm as described in the lecture not be run in this case?
- How would you modify MapReduce so that it can handle this case?

Exercise 0.52

## Exercise 7.5

- For optimal load balancing, the inverters in MapReduce must get segmented postings files of similar sizes
- For a new collection, the distribution of key-value pairs may not be known a-priori
- How would you solve this problem?