

CS 6240 Study Guide for the Exam

You should be well-prepared since you studied the online modules and practiced your knowledge with the end-of-module quizzes, attended the lectures, did the homework where we practiced many of the concepts, and were reading along in the textbooks where needed. At this point, we recommend you browse through all modules again to see the “big picture” of what we covered and how material in different modules relates to each other. Then go over your own lecture notes and try to explain the most important material (ideally to somebody else). Expect questions of the type a more technical interviewer might ask you when you apply for a job in a team that develops MapReduce programs for large-scale data analytics:

1. What are Mapper, Reducer, Combiner, and Partitioner? Give examples for each.
2. Given a small MapReduce program and some data, explain how the data is processed and derive the final output. How many `map()` and `reduce()` calls will be made? Practice by picking some of the programs we discussed and running them manually (on paper) on some small input.
3. Most questions will be of the type “write a MapReduce program in pseudo-code for a given problem.” This is very important. Expect problems of the size and complexity like the ones we discussed in class. Practice writing pseudo-code (like in the modules), not detailed Java code. A good way to prepare is to take problems we discussed (or a new one) and try to write the MapReduce pseudo-code without looking at the solution.
4. Know how the important design patterns we discussed in class affect program behavior. E.g., how does a Combiner or in-mapper combining improve performance and what are its disadvantages? How do secondary sort and order inversion work, and what do they achieve?
5. Understand important non-trivial algorithms, e.g., sorting, equi-join, PageRank, matrix product. Walk through Map and Reduce calls for small data samples.
6. Know how choice of key and Partitioner affect load distribution and memory consumption (heap space use).
7. Given two different programs for a problem, compare their communication cost by analyzing total amount of data transferred into Map, from Map to Reduce, and out of Reduce.
8. Given a new problem scenario, argue why (not) MapReduce would be a good match for it.

We will not ask low-level details about HDFS, there will be no question about PigLatin, and you do not need to know details about data mining techniques beyond what we discussed in class.

Specific suggestions for writing pseudo-code

- Keep function definitions compact. E.g., instead of `public void map(IntWritable key, Text value, Context context)`, just write `map(matrix cell c)`.
- Don't worry about details like writing a complete string parser just to access the fields of a record. E.g., instead of `String[] index = value.toString().split(",")` etc., just refer to the different fields of a record by their name: `c.row`, `c.column`, `c.value`.

- Use *for all records r in file/list/etc* (or *for each*—same thing) when going through all records in a list or file. This is easier to read than statements like *while (tokenizer.hasMoreTokens())*.

Example question and answer

Problem:

We are given a large input file (word1, word2, word3,...) containing words. Write the best possible MapReduce program (pseudo-code) you can think of to find the word with the most occurrences of letter “y”. (If there are multiple such words, output at least one of them.) You do not need to write any text parsers. In a couple of sentences, explain the main idea and design decisions of your program.

Example solution:

Brief summary of solution idea: A Map task receives a split of the file, usually containing many words. We can find the word with most “y” in such a split by using in-mapper combining. All Mappers then send their “local winner” to a single Reduce call, which determines the global winner. Instead, we could also globally sort the words by their number of y’s, but this would move a lot more data from Mappers to Reducers.

Class Mapper {

 localWinner = (NULL, 0) // localWinner is a pair (word, count), keeping track of the word with most y’s in the task

```
Map( ..., word w) {
    if (localWinner = NULL or yCount(w) > localWinner.count)
        localWinner = (w, yCount(w))
}
```

```
Cleanup()
    emit( dummy, localWinner.word )
```

}

// There is only a single Reduce call: for key=dummy

```
Reduce( dummy, [ word1, word2,...] ) {
    globalWinner = (word1, yCount(word1))

    for each word w in input list do
        if (yCount(w) > globalWinner.count)
            globalWinner.set( w, yCount(w) )

    emit( globalWinner.word )
}
```

Possible improvement: We can use the secondary sort design pattern to avoid having to go through the entire Reduce input list. To implement it, we make `localWinner.count` part of the key in the Map output. (Notice that this increases the amount of data sent from Mappers to Reducers due to the additional count field, but overall the amount of data is still very small because only a single record per Mapper is emitted.) We also need to add a custom Partitioner that partitions on the dummy field only, ignoring the counter field of the key. Similarly, we have to add a grouping comparator that ignores the counter field of the key.

For the Reduce function to receive the global winner as the first input record, we need to define a regular key comparator that sorts on the dummy field first (not really needed here because all keys have the same dummy value), then breaks ties by sorting in decreasing order of the count field.

With secondary sort, the Reduce function simply emits the first input record.