Université Grenoble Alpes & Grenoble INP M2 MoSIG Academic year –

Large-scale data management and distributed systems — Final Exam Data Management

1 Some Statements (2 points)

For each of the following statements, answer **True** or **False** and explain in a few words. Answers without explanations will be ignored.

- (a) Taking the topology of data centers (which nodes are in which racks) into account can help improving the fault tolerance of big data applications.
- (b) Spark is in general more efficient than Hadoop MapReduce because it moves the computation instead of the data.
- (c) To deal with large amount of data, HDFS is based on a scale-out approach.
- (d) The only advantage of NoSQL databases over relational databases is that they better support data partitioning.
- (e) Since network partitioning is impossible to avoid in a distributed system, it is impossible to build a distributed database that provides availability and consistency.

2 Processing data (3 points)

2.1 In the context of Map-Reduce frameworks, explain what "*Code once, benefit to all*" refers to. (A detailed explanation is expected)

2.2 We consider a CSV file including information about all students registered at UGA. All information about one student are stored in one line of the file.

The global statistics about UGA students say that about 50% of the students are not french¹. We would like to compute the number of foreign students in each university school (Science, Management, Architecture, etc.).

The algorithm presented in Figure 1 is an attempt to make the required computation. Here are some comments to help you analyze this code:

- The algorithm is written in Python.
- The algorithm is correct, in the sense that it does not crash.

¹This is not an official statistics, I just invented this number.

- The algorithm is based on RDDs. It does not use Spark dataframes.
- About lambda: The keyword lambda is used to define anonymous functions in Python.
- Semantic of the main RDD and Python functions used:
 - map(): Applies the defined function to each item in the RDD.
 - filter(): Only keeps the items for which the condition is true.
 - groupByKey(): Group the values for each key in the RDD into a single sequence.
 - len(): returns the number of items in a sequence

```
dataFile = sc.textFile("hdfs://.../data.csv")
1
     dataSplit = dataFile.map(lambda v: v.split(','))
\mathbf{2}
     # X is the column of the school -- Y is the column of the nationality
3
4
     data = dataSplit.map(lambda v:(v[X], v[Y]))
5
     d1 = data.filter(lambda v: v[1] != "French")
\mathbf{6}
7
     d2 = d1.groupByKey()
     d3 = d2.map(lambda v: (v[0], len(v[1])))
8
9
10
     print(d3.collect())
```

Figure 1: Algorithm for analyzing a CSV file

Answer the following questions:

- (a) Do you think that the algorithm presented in Figure 1 computes the number of foreign students per school? Answer YES or NO, and:
 - If your answer is yes, describe the logical steps of the algorithm.
 - If your answer is no, justify your answer.
- (b) Independently of whether it does the expected computation or not, propose 2 changes in the algorithm that would improve its performance and explain in a few words why each of the changes you propose can be beneficial.

2.3 A data engineer has to build a stream processing pipeline for a company. While discussing the design of his/her solution, he/she makes the following claim: "Since our number of data sources is a multiple of the number of workers on our stream processing engine, since all data belong to the same topic, and since we are only interested in computing trends (it is not a problem if we lose some data), I conclude that it is useless to put a message broker service between our data sources and the stream processing engine."

Assuming that the claims that the data engineer makes about the characteristics of the problem are correct, do you agree with his/her conclusion? (A detailed answer is expected)

3 Storing data (3 points)

3.1 A Bloom filter is a probabilistic data structure that can be used to test whether an element is in a set.

Provide a detailed answer to each of the following questions:

- (a) Explain the main principles of bloom filters, and what guarantees we have about the answers provided by a bloom filter.
- (b) Explain how and why bloom filters are used to improve the performance of LSM-Tree-based systems.

3.2 RocksDB is a very efficient key-value store targeting large scale distributed systems. RocksDB is based on an LSM-Tree data structure. One very attractive feature of RocksDB is the fact that it allows users to select between different compaction algorithms to be used internally for major compaction. Depending on the characteristics of the compaction algorithm that is selected, RocksDB can be more *read friendly* or *write friendly*.

- (a) Explain the purpose of major compaction in LSM-Tree algorithms.
- (b) Explain why the compaction algorithm can have an impact on the efficiency of RocksDB when handling read-heavy or write-heavy workloads.

3.3 A data engineer is in charge of selected the storage system that will be used for storing some data in a company. These data are stored in a table where each row has a large number of columns. The main access patterns for the data are as follows:

- Each write operation mostly update a few elements in a row
- Read operations access a large number of rows with filtering based on one column
- There are 90% of read operations and 10% of write operations

The data engineer claims that "because of the write access pattern, the only solution is to go with a solution based on row-oriented storage".

Do you agree with his/her claim? (a detailed answer is expected)