

HDFS Problem

Consider a data storage for University students. Each student data, stuData which is in a file of size less than 64 MB (1 MB = 220B). A data block stores the full file data for a student of stuData_idN, where N = 1 to 500.

(i) How the files of each student will be distributed at a Hadoop cluster? How many student data can be stored at one cluster? Assume that each rack has two DataNodes for processing each of 64 GB(1 GB = 230B) memory. Assume that cluster consists of 120 racks, and thus 240 DataNodes.

(ii) What is the total memory capacity of the cluster in TB ((1 TB = 240B) and DataNodes in each rack?

(iii) Show the distributed blocks for students with ID= 96 and 1025. Assume default replication in the DataNodes = 3.

(iv) What shall be the changes when a stuData file size ≤ 128 MB?

Solution :

(i) Data block default size is 64 MB. Each students file size is less than 64MB. Therefore, for each student file one data block suffices. A data block is in a DataNode. Assume, for simplicity, each rack has two nodes each of memory capacity = 64 GB. Each node can thus store $64 \text{ GB}/64\text{MB} = 1024$ data blocks = 1024 student files. Each rack can thus store $2 \times 64 \text{ GB}/64\text{MB} = 2048$ data blocks = 2048 student files. Each data block default replicates three times in the DataNodes. Therefore, the number of students whose data can be stored in the cluster = number of racks multiplied by number of files divided by 3 = $120 \times 2048/3 = 81920$. Therefore, the maximum number of 81920 stuData_IDN files can be distributed per cluster, with N = 1 to 81920.

(ii) Total memory capacity of the cluster = $120 \times 128 \text{ MB} = 15360 \text{ GB} = 15 \text{ TB}$. Total memory capacity of each DataNode in each rack = $1024 \times 64 \text{ MB} = 64 \text{ GB}$.

(iii) Figure 2.3 shows a Hadoop cluster example, and the replication of data blocks in racks for two students of IDs 96 and 1025. Each stuData file stores at two data blocks, of capacity 64 MB each. (Refer the diagram in ppt)

(iv) Changes will be that each node will have half the number of data blocks.(Don't confuse this as we discussed in class as twice, we said twice because the current 64 needs to be made 128 however in terms of node size it should be half)

