

6.14 [15/15/10/10] <6.2, 6.3> MapReduce enables large amounts of parallelism by having data-independent tasks run on multiple nodes, often using commodity hardware; however, there are limits to the level of parallelism. For example, for redundancy MapReduce will write data blocks to multiple nodes, consuming disk and, potentially, network bandwidth. Assume a total dataset size of 300 GB, a network bandwidth of 1 Gb/s, a 10 s/GB map rate, and a 20 s/GB reduce rate. Also assume that 30% of the data must be read from remote nodes, and each output file is written to two other nodes for redundancy. Use Figure 6.6 for all other parameters.

- a. [15] <6.2, 6.3> Assume that all nodes are in the same rack. What is the expected runtime with 5 nodes? 10 nodes? 100 nodes? 1000 nodes? Discuss the bottlenecks at each node size.
- b. [15] <6.2, 6.3> Assume that there are 40 nodes per rack and that any remote read/write has an equal chance of going to any node. What is the expected runtime at 100 nodes? 1000 nodes?
- c. [10] <6.2, 6.3> An important consideration is minimizing data movement as much as possible. Given the significant slowdown of going from local to rack to array accesses, software must be strongly optimized to maximize locality. Assume that there are 40 nodes per rack, and 1000 nodes are used in the MapReduce job. What is the runtime if remote accesses are within the same rack 20% of the time? 50% of the time? 80% of the time?
- d. [10] <6.2, 6.3> Given the simple MapReduce program in Section 6.2, discuss some possible optimizations to maximize the locality of the workload.

a. $300/5 = 60 \text{ GB} \rightarrow \text{data per node for 5 nodes}$

$60 \text{ GB} \cdot 30 = 180 \text{ GB} \rightarrow \text{read from remote}$

$180 \text{ GB} / 100 \text{ bandwidth (MB/s)} = 180 \text{ seconds} \rightarrow \text{remote access time}$

$60 \text{ GB} / 200 \text{ MB/s} = 300 \text{ seconds} \rightarrow \text{local disk access}$

$60 \text{ seconds} / \text{GB} = 600 \text{ seconds} \rightarrow \text{map time}$

$60 \text{ seconds} / \text{GB} = 1200 \text{ seconds} \rightarrow \text{reduce}$

$(180 \cdot 2) + (210 \cdot 2) + 600 + 1200$

$360 + 420 + 600 + 1200$

$= 2580 \text{ seconds total}$

the primary bottleneck is in the reduce stage with the next being in the data transfer stage.

$300/10 = 30 \text{ GB} \rightarrow \text{data per node for 10 nodes}$

$$30^{\text{data per node}} \cdot 30 = 9 \text{ GB} \rightarrow \text{read from remote}$$

$$9/100^{\text{bandwidth (MB/s)}} = 90 \text{ seconds} \rightarrow \text{remote access time}$$

$$30 - 9 / 200 (\text{MB/s}) = 21 / 200 = 105 \text{ seconds} \rightarrow \text{local disk access}$$

$$30 * 10 \text{ s/GB} = 300 \text{ seconds} \rightarrow \text{map time}$$

$$30 * 20 \text{ s/GB} = 600 \text{ seconds} \rightarrow \text{reduce}$$

$$90 + 105 + 300 + 600 = 1,095 \text{ seconds total}$$

the primary bottleneck is in the reduce stage with the next being in the data transfer stage.

$300/100 = 3 \text{ GB} \rightarrow \text{data per node for 100 nodes}$

$$3^{\text{data per node}} \cdot 30 = .9 \text{ GB} \rightarrow \text{read from remote}$$

$$.9/100^{\text{bandwidth (MB/s)}} = 9 \text{ seconds} \rightarrow \text{remote access time}$$

$$3 - .9 / 200 (\text{MB/s}) = 2.1 / 200 = 10.5 \text{ seconds} \rightarrow \text{local disk access}$$

$$3 * 10 \text{ s/GB} = 30 \text{ seconds} \rightarrow \text{map time}$$

$$3 * 20 \text{ s/GB} = 60 \text{ seconds} \rightarrow \text{reduce}$$

$$9 + 10.5 + 30 + 60 = 109.5 \text{ seconds total}$$

the primary bottleneck is in the reduce stage with the next being in the data transfer stage.

$$300/1000 = 0.3 \text{ GB} \rightarrow \text{data per node for } 1000 \text{ nodes}$$

$$0.3 \xrightarrow{\text{data per node}} .30 = .09 \text{ GB} \rightarrow \text{read from remote}$$

$$.09/100 \xrightarrow{\text{bandwidth (MB/s)}} = .9 \text{ seconds} \rightarrow \text{remote access time}$$

$$0.3 - .09/200 \text{ (MB/s)} = .21/200 = 1.05 \text{ seconds} \rightarrow \text{local disk access}$$

$$.3 * 10 \text{ s/GB} = 3 \text{ seconds} \rightarrow \text{map time}$$

$$.3 * 20 \text{ s/GB} = 6 \text{ seconds} \rightarrow \text{reduce}$$

$$.9 + 1.05 + 3 + 6 = 10.95 \text{ seconds total}$$

the primary bottleneck is STILL in the reduce stage with the next being in the data transfer stage.

b.

$$\frac{100}{10} = 2.5 \text{ racks}$$

$$100 \text{ nodes}$$

$$300/100 = 3 \text{ GB} \rightarrow \text{data per node}$$

$$3 \text{ GB} * 0.3 = 900 \text{ MB}$$

$$\hookrightarrow 900 \text{ MB} * \frac{2}{3} = 600 \text{ MB} \rightarrow \text{read}$$

$$600/10 \xrightarrow{\text{bandwidth (MB/s)}} = 60 \text{ seconds} \rightarrow \text{remote access}$$

$$300/100 = 3 \text{ seconds} \rightarrow \text{local}$$

$$60 + 30 + (60 * 2) + (3 * 2) = 216 \text{ seconds}$$

the greatest bottleneck is the data transfer.

$$1000 \text{ nodes} \quad 1000 / 10 = 100 \text{ racks}$$

$$300 / 1000 = .3 \text{ GB} \rightarrow \text{data per node}$$

$$.3 \text{ GB} * 0.3 = 90 \text{ MB}$$

$$\hookrightarrow 90 \text{ MB} * 2/3 = 60 \text{ MB} \rightarrow \text{read}$$

$$60 / 10 \xrightarrow{\text{bandwidth (MB/s)}} \frac{60}{6} \text{ seconds} \rightarrow \text{remote access}$$

$$30 / 100 = 0.3 \text{ seconds} \rightarrow \text{local}$$

$$6 + 3 + (6 * 2) + (.3 * 2) = 21.6 \text{ seconds}$$

the greatest bottleneck is the data transfer.

c.

$$20\% \text{ of the time} \rightarrow 900 * 20\% = 180 \text{ same rack, } 720 \text{ remote}$$

$$720 / 10 \text{ MB/s} = 72 \text{ seconds} \rightarrow \text{remote access}$$

$$180 / 100 \text{ MB/s} = 1.8 \text{ seconds} \rightarrow \text{local}$$

$$[(72 + 1.8) * 2] + 30 + 60 = 237.6 \text{ seconds total}$$

$$50\% \text{ of the time} \rightarrow 450 \text{ same, } 450 \text{ remote}$$

$$450 / 10 \text{ MB/s} = 45 \text{ seconds} \rightarrow \text{remote access}$$

$$450 / 100 \text{ MB/s} = 4.5 \text{ seconds} \rightarrow \text{local}$$

$$[(45 + 4.5) * 2] + 30 + 60 = 189 \text{ seconds total}$$

$$80\% \text{ of the time} \rightarrow 720 \text{ same, } 180 \text{ remote}$$

$$180 / 10 \text{ MB/s} = 18 \text{ seconds} \rightarrow \text{remote access}$$

$$720 / 100 \text{ MB/s} = 7.2 \text{ seconds} \rightarrow \text{local}$$

$$[(18 + 7.2) * 2] + 30 + 60 = 140.4 \text{ seconds total}$$

- d. Some possible optimizations that would maximize the locality of the workload would be decreasing duplicate words between the map phase and reduce phase which would allow more data to properly fit in node caches overall improving efficiency. The system could also calculate the number of bytes in the map phase as well as properly placing data for reduction which minimizes data transfer within racks.