## MapReduce, Spark and IO

Mentoring 11: November 25, 2019

## 1 Speaking Spark

1.1 In this question let's write Spark to find the mode of a list of values and how often it occurs. As a refresher, the mode is the number that appears most often. If there is a tie, select any of the options. Fill in the blanks for the Python code below. Use the following Spark Python functions when necessary: map, flatmap, reduce, reduceByKey.

## 2 More MapReduce

2.1 Imagine we're looking at Facebook's friendship graph, which we model as having a vertex for each user, and an undirected edge between friends. Facebook stores this graph as an adjacency list, with each vertex associated with the list of its neighbors, who are its friends. This representation can be viewed as a list of degree 1 friendships, since each user is associated with their direct friends. We're interested in finding the list of degree 2 friendships, that is, an association between each user and the friends of their direct friends.

You are given a list of associations of the form (user\_id, list(friend\_id)), where the user\_id is 1st degree friends with all the users in the list.

Your output should be another list of associations of the same form, where the first item of the pair is a user\_id, and the second item is a list of that user's 2nd degree friends. Note: a user is not their own 2nd degree friend, so the list of second degree friends must not include the user themselves.

Write pseudocode for the mapper and reducer to get the desired output from the input. Assume you have a set data structure. You can call set(list) to create a set, and use .remove(value), and .union(set) methods.

## 3 Some Miscellaneous Quest(IO)ns

An important advantage of interrupts over polling is the ability of the processor to perform other tasks while waiting for communication from an I/O device. Suppose that a 1 GHz processor needs to read 1000 bytes of data from a particular I/O device. The I/O device supplies 1 byte of data every 0.02 ms. The time to process the data and store it in a buffer is negligible.

- 3.1 Assume a polling iteration takes 60 cycles. If the processor detects that a byte of data is ready through polling:
  - (a) How many cycles does it take for the I/O device to supply 1 byte of data?
  - (b) How many polling iterations does it take to read 1 byte of data? (round up to an integer)
  - (c) How many cycles does it take to read the 1000 bytes of data?
- 3.2 If instead, the processor is interrupted when a byte is ready, and the processor spends the time between interrupts on another task, how many cycles of this other task can the processor complete while the I/O communication is taking place? The overhead for handling an interrupt is 2000 cycles.
- 3.3 The advantage of polling however arises when data rates become very large so that the interrupt overhead becomes substantial and at some point the system simply can't keep up. What is the data arrival time (in ms) at which point an interrupt-driven I/O scheme on this computer can't keep up with the data coming in? The overhead for handling an interrupt is 2000 cycles.

- 4 MapReduce, Spark and IO
- 3.4 Solve for the maximum controller overhead to meet the following specifications: We need disk latency under 18 ms while reading 800 B of data. The hard drive spins at 6000 rev/min with a seek time of 2.5 ms and transfer rate of 80 KB/s (SI prefix). Don't forget units!
- 3.5 To support interrupts, the CPU should be able to save and restore the current state. Which of the following should be saved before handling interrupts to ensure correct execution?
  - a. Program Counter b. User Registers c. TLB d. Caches
- 3.6 Consider the following three devices. For which device is Direct Memory Access (DMA) most beneficial?

Device	Data Rate	Transfer Block Size
Α	80 B/s	4 B
В	400 MB/s	4 B
С	400 MB/s	2 KB