

Due Wednesday, October 14th, 2020 at 11:59 PM PST, on Gradescope.

## 1 Scheduling Simulator Implementation

For both of these problems, we intend that you implement the schedulers in the IPython notebook and then run the simulation to obtain the CPU log output. We will not collect your code from the notebook, however, so you are free to execute the scheduling algorithm manually to obtain the log output if you wish.

- Implement the SRTF scheduler. To make the existing tests pass, you should break ties according to FIFO. Run it on `workload3` in the IPython notebook, with a quantum of 2, and produce the CPU log output here. Cell 21 sets this up for you.
- Implement the MLFQ scheduler. Use two queues, one high-priority queue for interactive tasks and one low-priority queue for CPU-bound tasks, each serviced in a round-robin fashion. Each CPU burst begins in the high-priority queue; if its quantum expires before the CPU burst ends, it is demoted to the low-priority queue. Run it on `workload3` in the IPython notebook, using a quantum of 4 in the low-priority queue and a quantum of 2 in the high-priority queue, and produce the CPU log output here. Cell 27 sets this up for you.

## 2 Approaching 100% Utilization

Consider a sequence  $B_i$  of CPU bursts, where (for simplicity) each CPU burst has a fixed length  $T_S$ . The first burst,  $B_0$ , arrives at time  $t = 0$  (i.e.,  $\text{ArrivalTime}(B_0) = 0$ ). For  $i \geq 1$ , the arrival time of  $B_i$  is given by  $\text{ArrivalTime}(B_i) = \text{ArrivalTime}(B_{i-1}) + X_i$ , where the  $X_i$  are i.i.d. exponentially distributed random variables with parameter  $\lambda$ .

To allow the CPU bursts to execute concurrently, we model each CPU burst as executing in its own task.

- Is this an open system or a closed system? Explain.
- What value of  $\lambda$  should we choose, such that the mean time between arrivals is equal to  $T_S$ ? Explain.
- What value of  $\lambda$  should we choose, such that the system runs at 50% utilization on average? Explain.

Now, set up a simulation in the IPython notebook to model the system with a large number of CPU bursts. Fix some value for  $T_S$ . Vary the arrival rate (i.e.,  $\lambda$ ) starting at a small number, increasing it to approach the value you determined in Part b. Run the simulation at various points along the way, with multiple trials for each point, **making sure to choose some points where  $\lambda$  is very close to the value you determined in Part b.**

- As you vary  $\lambda$  as described as above, what happens to CPU utilization? Show a line plot with the arrival rate on the  $x$  axis and CPU utilization on the  $y$  axis, depicting the results.
- As you vary  $\lambda$  as described as above, what happens to the response time for each CPU burst? Show a line plot with the arrival rate on the  $x$  axis and response time on the  $y$  axis, depicting the results. Be sure to consider both the median response time and the 95th percentile.
- Does using a different scheduler affect your answer to Part e? (Hint: If you choose to look at SRTF, think carefully about how your implementation breaks ties and how that might affect the median response time.)
- Qualitatively explain why running a system at close to 100% throughput results in poor latency.

### 3 Fairness for CPU Bursts

Consider two periodic tasks,  $S$  and  $T$ . Each task consists of a series of CPU bursts; each task yields to the scheduler between CPU bursts but (for simplicity) has zero I/O time between CPU bursts. Let  $S_i$  (respectively,  $T_i$ ) be the random variable that denotes the length of the  $i$ th CPU burst. **Assume that the CPU burst lengths (i.e.,  $S_i$  and  $T_i$ ) are i.i.d.**

Sam, a student in CS 162, wishes to run  $S$  and  $T$  concurrently on a single CPU. Noticing that  $S$  and  $T$  use the same burst length distribution, he reasons that the CPU will be shared fairly between the two tasks, even if the scheduler does not enforce fairness. Thus, he chooses a simple FCFS scheduler.

- (a) Explain why the length of the FCFS queue never exceeds 2.
- (b) What is  $\Pr[S_1 < T_1]$ ?
- (c) Suppose  $S$  has run for  $m$  CPU bursts, where  $m$  is large. Using the Central Limit Theorem, characterize  $\text{CPUTime}(S)$ , the total CPU time spent on  $S$ , as a normal distribution parameterized by  $\mathbb{E}[S_i]$ ,  $\text{Var}(S_i)$ , and  $m$ .
- (d) After the scheduler has run a large number of CPU bursts, what is  $\Pr[n \cdot \text{CPUTime}(S) < \text{CPUTime}(T)]$ ? Use your approximation from Part c and write your answer in terms of  $\Phi(x)$ , the CDF of the Standard Normal Distribution.

Part d quantifies unfairness. For example, the probability that  $T$  receives at least 5% more CPU time than  $S$  corresponds to  $n = 1.05$ .

For simplicity, assume now that  $\mathbb{E}[S_i] = \sqrt{\text{Var}(S_i)}$ . This is satisfied by, for example, the exponential distribution.

- (e) Using software of your choice, (e.g. Maple, WolframAlpha, Python, graphing calculator), calculate the probability that one task receives at least 10% more CPU time than the other, for  $m = 100$ . Be sure to consider both tails of the distribution (i.e.,  $T$  may get more CPU time than  $S$  OR  $S$  may get more CPU time than  $T$ ). How does this change if you use  $m = 10000$ ? Was Sam's reasoning that the CPU allocation will be fair with FCFS correct?
- (f) Run a simulation in the IPython notebook to confirm your result from Part e. Describe the simulation you performed and produce a graph supporting your conclusion.
- (g) (Optional Stretch) How would the result be different if Sam were to use a preemptive round-robin scheduler with a small quantum?