

CS-350 - Fundamentals of Computing Systems

Homework Assignment #5 - EVAL

Due on October 24, 2024 — Late deadline: October 25, 2024 EoD at 11:59 pm

Prof. Renato Mancuso

Renato Mancuso

EVAL Problem 1

In this EVAL problem, we will study the behavior of a system with a configurable FIFO vs. SJN queue dispatch (a.k.a. scheduling) policy. The goal of this EVAL is to demonstrate that despite no changes are made to the workload, some key system metrics might change if the system is able to meaningfully reorder the pending workload for processing.

- a) First thing first, let us understand the overhead of using SJN as the queue scheduling policy. Measuring overheads is always hard, but luckily for you when things are implemented on a real machine like these assignments, measuring the true machine behavior becomes possible!

We are going to measure overheads in two situations: (1) at low utilization, and (2) at high utilization. In both cases, you will run the commands provided below 10 times and extract the average server runtime. We will then compare the average obtained under FIFO vs. those obtained under SJN.

(1) **Low Utilization.** In the low-utilization test, run 10 times the following command for the FIFO experiment:

```
/usr/bin/time -v ./server_pol -w 2 -q 100 -p FIFO 2222 & ./client -a 10 -s 20 -n 1500 2222
```

As you can see the `/usr/bin/time` utility is prefixed to the server to obtain the final runtime statistics. You should parse the 10 results from this command that are found in the “Elapsed (wall clock) time” statistic reported by the `/usr/bin/time` utility at the end of the run, and compute an average of these values.

Now, repeat the same thing by replacing `-p FIFO` with `-p SJN` in the command above. Collect the new 10 values of wall-clock times and compute the average. Finally, compare the average runtime computed for FIFO vs. SJN.

(2) **High Utilization.** In the high-utilization test, we follow an identical procedure as above, with the only difference that we change the amount of workload sent by the client. The new command template will be:

```
/usr/bin/time -v ./server_pol -w 2 -q 100 -p <policy> 2222 & ./client -a 40 -s 20 -n 1500 2222
```

where `<policy>` should be set to FIFO in the first 10 runs, and to SJN in the last 10 runs.

Compare the average runtime under FIFO and SJN at low and high utilization, reasoning on the extra per-request time spent by the server under the two policies.

- b) Okay, we have understood the cost of SJN with the measurements above. Now let us quantify the benefits. In a way similar to what we did in HW1, we can measure the average response time as a function of the server utilization. Recall that the utilization is computed as the sum of the length of all the processed requests over the total server runtime. With only one worker, that should be a number less than or equal to 1, with two workers it should be less than or equal to 2, and so on.

For this experiment, we will always set `-w 2` to spawn two workers. Next, do a series of experiments where you run the following template command 10 times:

```
./server_pol -w 2 -q 100 -p FIFO 2222 & ./client -a <arr_rate> -s 20 -n 1500 2222
```

Where the parameter `<arr_rate>` goes from a value of 22 (first run) to a value of 40 (10th run), thus increasing by 2 at every run. Use these experiments to construct a curve of the request response time as a function of the server utilization.

Now, do 10 more runs with the same strategy as described above, but where the server policy is set to SJN, thus the new command template will be:

```
./server_pol -w 2 -q 100 -p SJN 2222 & ./client -a <arr_rate> -s 20 -n 1500 2222
```

Plot the two curves (one from the FIFO experiments, and the other from the SJN experiments) in the same plot and try to quantify which queueing policy works better and by how much.

- c) In the engineering world, we say that there is no free lunch. If one of the two policies works better than the other, where is the catch? We are set to discover that with the following experiment. The response time average only tells part of the story, so let us study the full distribution of response times.

Focus on the outputs produced in the previous experiment by the 10th run (the one where you set `-a 40`) of each policy (FIFO and SJN). Produce TWO separate plots, each depicting the CDF of the response times observed in the two runs under analysis. Make sure that in both plots, the range of the x and y axis are the same, so that they are visually comparable.

In the CDF plots, add two vertical bars: (1) one marking the value of the average response time, and (2) another marking the 99th percentile response time value.

- d) Now focus on the differences that you see between the two CDF plots. Reason on the spread of the response time distribution and on the overall predictability of the system when evaluated by looking response time metric. *HINT*: see possible definitions of predictability in the CS350 book.

Finally, provide a motivated answer to the question: *which queuing policy leads to a more predictable system behavior?*