


TCSS 422: OPERATING SYSTEMS

**Common Scheduling Algorithms,
Multi-level Feedback Queue (MLFQ) Scheduler,
Proportional Share Schedulers**

Wes J. Lloyd
School of Engineering and Technology
University of Washington - Tacoma



April 11, 2024 TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

1

TEXT BOOK COUPON

- 10% off textbook code: **MYLIBRARY10** (through Friday Apr 12)
- <https://www.lulu.com/shop/andrea-arpaci-dusseau-and-remzi-arpaci-dusseau/operating-systems-three-easy-pieces-hardcover-version-110/hardcover/product-15gjeeqy.html?q=three+easy+pieces+operating+systems&page=1&pageSize=4>
- With coupon textbook is only \$35.78 + tax & shipping

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma L6.2

2

TCSS 422 – OFFICE HRS – SPRING 2024

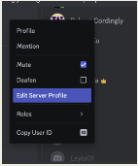
- ****Tuesdays after class until 7:00pm****
Hybrid (In-person/Zoom)
 - This session will be in person in CP 229.
 - Zoom will be monitored when no student is in CP 229.
- **Thursdays after class until 7:00pm – Hybrid (In-person/Zoom)**
 - Additional office time will be held on Thursdays after class when there is high demand indicated by a busy Tuesday office hour
 - When Thursday Office Hours are planned, Zoom links will be shared via Canvas
 - Questions after class on Thursdays are always entertained even when the formal office hour is not scheduled

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma L6.3

3

TCSS 422 DISCORD SERVER

- Please join the TCSS 422 A – Spring 2024 Discord Server
- <https://discord.gg/H7PPZ5ArFW>
- Under Edit Server Profile:
Please update your 'Server Nickname' to your real name or UW NET ID
THANK YOU



April 11, 2024 TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma L6.4

4

INDUSTRY GUEST SPEAKER RED HAT LINUX (IBM) APRIL 20

- Grad Certificate Soft Dev Eng. (GC-SDE) Spring Seminar. **open to TCSS 422 students**
- Damien Eversmann, RedHat Chief Architect for Education
- **Saturday, April 20** - 12:30 to 1:20 pm
- **Zoom Link:** <https://washington.zoom.us/j/96445774685>

Selected as one of the Industry Leaders of the Year in 2022 by EdScoop, Damien has over 25 years of experience as an IT professional. Having spent the bulk of his career working in or in support of the public sector, he is somewhat of an expert when it comes to IT in government and higher education. Throughout his working life, Damien has served as a Developer, System Administrator, Development Manager, Enterprise Architect and Technology Director. Living the life of an Academic and Research Administrator has also given Damien a vast knowledge of and a healthy respect for regulations and compliance. He has worked on projects running the gamut from desktop-based widgets to major, multi-tiered applications, from small, embedded systems to many-faceted infrastructures.

As Chief Architect for Education at Red Hat, Damien serves the role of bridging the gap between the mission and the business of education and the technologies and solutions that support it all. He has a penchant for teaching and demonstration and anything else that gets him in front of people to share the message of Continuous Learning, DevOps Culture, Innovation through Automation and IT Modernization.

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma L6.5

5

OBJECTIVES – 4/11

- **Questions from 4/11**
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 – Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma L6.6

6

ONLINE DAILY FEEDBACK SURVEY

- Daily Feedback Quiz in Canvas – Available After Each Class
- Extra credit available for completing surveys **ON TIME**
- Tuesday surveys: due by ~ Wed @ 11:59p
- Thursday surveys: due ~ Mon @ 11:59p

April 11, 2024 | TCSS422: Computer Operating Systems (Spring 2024) | L6.7

7

TCSS 422 - Online Daily Feedback Survey - 4/1

Quiz Instructions

Question 1 0.5 pts

On a scale of 1 to 10, please classify your perspective on material covered in today's class:

1 2 3 4 5 6 7 8 9 10

Mostly review to me Equal new and review Mostly new to me

Question 2 0.5 pts

Please rate the pace of today's class:

1 2 3 4 5 6 7 8 9 10

slow just right fast

April 11, 2024 | TCSS422: Computer Operating Systems (Spring 2024) | L6.8

8

MATERIAL / PACE

- Please classify your perspective on material covered in today's class (29 respondents):
- 1-mostly review, 5-equal new/review, 10-mostly new
- **Average – 7.00 (↑ - previous 6.44)**
- Please rate the pace of today's class:
- 1-slow, 5-just right, 10-fast
- **Average – 5.21 (no change - previous 5.21)**

April 11, 2024 | TCSS422: Computer Operating Systems (Spring 2024) | L6.9

9

FEEDBACK FROM 4/9

- **What does preemption mean?**
 - Preempt or preemption means to interrupt
 - In the context of scheduling for operating systems, preemption refers to the interruption and pausing of a job or task so that some other job or task is allowed to run
 - When a job is interrupt, it goes from RUNNING → READY
 - Why does the job not go from RUNNING → BLOCKED when preempted?
- **the timing of ABC, and how to determine the Initial speed**
 - For some problems, we will say that the job arrive in the sequence of 'A B C', but we do not provide distinctly different times. We say they all arrive at time t=0, but in the order of A B C
 - This ordering is required as schedulers like FIFO require the ability to infer the arrival ordering, but we do not distinguish distinct arrival times
 - Question-on-question: What is mean by 'the initial speed' ?

April 11, 2024 | TCSS422: Operating Systems (Spring 2024) | L6.10

10

FEEDBACK - 2

- **Some of the schedulers were a little confusing mainly the Faster First, random order.**
- Question-on-question:
 - What is "Faster First" ?
 - What is "random order" ?

April 11, 2024 | TCSS422: Operating Systems (Spring 2024) | L6.11

11

FEEDBACK - 3

- **The 'preemptive kernel' is a little unclear**
 - The primary idea here is that interrupt handlers (kernel functions that process interrupts) can now be interrupted in Linux starting with version >= 2.6
 - Non-Maskable-Interrupts have the highest priority and cannot be masked out by other interrupts. These are critical hardware events such as memory parity error or power loss.
 - Other interrupt handlers can be interrupted
 - Locks are added around kernel code that should not be interrupted. The locks increment a preemption counter and track the number or code sections running that can't be interrupted
 - Interrupts can only interrupt other interrupt handlers when counter is zero

April 11, 2024 | TCSS422: Operating Systems (Spring 2024) | L6.12

12

OBJECTIVES – 4/11

- Questions from 4/11
- **Assignment 0**
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 – Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.13

13

OBJECTIVES – 4/11

- Questions from 4/11
- Assignment 0
- **C Tutorial - Pointers, Strings, Exec in C**
- Quiz 1 – Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.14

14

OBJECTIVES – 4/11

- Questions from 4/11
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- **Quiz 1 – Active Reading Chapter 9**
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.15

15


OBJECTIVES – 4/11

- Questions from 4/11
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 – Active Reading Chapter 9
- **Chapter 7: Scheduling Introduction**
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.16

16

CHAPTER 7- SCHEDULING: INTRODUCTION



April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.17

17

CHAPTER 7

- Chapter 7: Scheduling Introduction
 - Scheduling metrics
 - Turnaround time, Jain's Fairness Index, Response time
 - FIFO, SJF, STCF, **RR** schedulers

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.18

18

SCHEDULING WITH I/O

- STCF scheduler
 - A: CPU=50ms, I/O=40ms, 10ms intervals
 - B: CPU=50ms, I/O=0ms
 - Consider A as 10ms subjobs (CPU, then I/O)
- Without considering I/O:

CPU utilization = $100/140 = 71\%$

Poor Use of Resources

April 11, 2024 TCCS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.19

19

SCHEDULING WITH I/O - 2

- When a job initiates an I/O request
 - A is blocked, waits for I/O to compute, frees CPU
 - STCF scheduler assigns B to CPU
- When I/O completes → raise interrupt
 - Unblock A, STCF goes back to executing A: (10ms sub-job)

Cpu utilization = $100/100 = 100\%$

Overlap Allows Better Use of Resources

April 11, 2024 TCCS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.20

20

Respond at poll.com/westloyd
 Text WESLLOYD to 22333 once to join, then 1, 2, 3, 4, 5...

Which scheduler, thus far, best address fairness and average response time of jobs?

First In - First Out (FIFO) 1
 Shortest Job First (SJF) 2
 Shortest Time to Completion First (STCF) 3
 Round Robin 4
 None of the Above 5
 All of the Above 6

Total Results: 0

Powered by Poll Everywhere
 Start the presentation to see live content. For screen share software, share the entire screen. Get help at poll.com/app

21

QUESTION: SCHEDULING FAIRNESS

- Which scheduler, this far, best addresses fairness and average response time of jobs?
 - First In - First Out (FIFO)
 - Shortest Job First (SJF)
 - Shortest Time to Completion First (STCF)
 - Round Robin (RR)
 - None of the Above
 - All of the Above

April 11, 2024 TCCS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.22

22

SCHEDULING METRICS

- Consider Three jobs (A, B, C) that require: $time_A = 400ms$, $time_B = 100ms$, and $time_C = 200ms$
- All jobs arrive at time=0 in the sequence of A B C.
- Draw a scheduling graph to help compute the **average response time (ART)** and **average turnaround time (ATT)** scheduling metrics for the FIFO scheduler.

Example:

April 11, 2024 TCCS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.23

23

What is the Average Response Time of the FIFO scheduler?

Powered by Poll Everywhere
 Start the presentation to see live content. For screen share software, share the entire screen. Get help at poll.com/app

24

What is the Average Turnaround Time of the FIFO scheduler?

Powered by Poll Everywhere

Start the presentation to see live content. For screen share software, share the entire screen. Get help at poller.com/app

25

SCHEDULING METRICS

- Consider Three jobs (A, B, C) that require: $time_A=400ms$, $time_B=100ms$, and $time_C=200ms$
- All jobs arrive at time=0 in the sequence of A B C.
- Draw a scheduling graph to help compute the **average response time (ART)** and **average turnaround time (ATT)** scheduling metrics for the SJF scheduler.

Example:

April 11, 2024 | TCSS422: Operating Systems (Spring 2024) School of Engineering and Technology, University of Washington - Tacoma | L6.26

26

What is the Average Response Time of the Shortest Job First Scheduler?

Powered by Poll Everywhere

Start the presentation to see live content. For screen share software, share the entire screen. Get help at poller.com/app

27

What is the Average Turnaround Time of the Shortest Job First Scheduler?

“ 7.75 milli ”

“ 2ms ”

“ Too long :(”

“ 1000 ”

Powered by Poll Everywhere

Start the presentation to see live content. For screen share software, share the entire screen. Get help at poller.com/app

28

OBJECTIVES – 4/11

- Questions from 4/11
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 – Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler**
 - Job Starvation
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 | TCSS422: Operating Systems (Spring 2024) School of Engineering and Technology, University of Washington - Tacoma | L6.29

29

CHAPTER 8 – MULTI-LEVEL FEEDBACK QUEUE (MLFQ) SCHEDULER

April 11, 2024 | TCSS422: Operating Systems (Spring 2024) School of Engineering and Technology, University of Washington - Tacoma | L6.30

30

MULTI-LEVEL FEEDBACK QUEUE

- Objectives:
 - Improve turnaround time:
Run shorter jobs first
 - Minimize response time:
Important for interactive jobs (UI)
- Achieve without a priori knowledge of job length

April 11, 2024
TCSS422: Operating Systems (Spring 2024)
School of Engineering and Technology, University of Washington - Tacoma
L6.31

31

MLFQ - 2

Round-Robin within a Queue

- Multiple job queues
- Adjust job priority based on observed behavior
- Interactive Jobs
 - Frequent I/O → keep priority high
 - Interactive jobs require fast response time (GUI/UI)
- Batch Jobs
 - Require long periods of CPU utilization
 - Keep priority low

[High Priority] Q8 → A → B

Q7

Q6

Q5

Q4 → C

Q3

Q2

[Low Priority] Q1 → D

April 11, 2024
TCSS422: Operating Systems (Spring 2024)
School of Engineering and Technology, University of Washington - Tacoma
L6.32

32

MLFQ: DETERMINING JOB PRIORITY

- New arriving jobs are placed into highest priority queue
- If a job uses its entire time slice, priority is reduced (↓)
 - Jobs appears CPU-bound ("batch" job), not interactive (GUI/UI)
- If a job relinquishes the CPU for I/O priority stays the same

MLFQ approximates SJF

April 11, 2024
TCSS422: Operating Systems (Spring 2024)
School of Engineering and Technology, University of Washington - Tacoma
L6.33

33

MLFQ: LONG RUNNING JOB

- Three-queue scheduler, time slice=10ms

Priority ↓

Long-running Job Over Time (msec)

April 11, 2024
TCSS422: Operating Systems (Spring 2024)
School of Engineering and Technology, University of Washington - Tacoma
L6.34

34

MLFQ: BATCH AND INTERACTIVE JOBS

- A_{arrival_time} = 0ms, A_{run_time} = 200ms,
- B_{run_time} = 20ms, B_{arrival_time} = 100ms

Priority ↓

Scheduling multiple jobs (ms)

April 11, 2024
TCSS422: Operating Systems (Spring 2024)
School of Engineering and Technology, University of Washington - Tacoma
L6.35

35

MLFQ: BATCH AND INTERACTIVE - 2

- Continuous interactive job (B) with long running batch job (A)
- Low response time is good for B
- A continues to make progress

The MLFQ approach keeps interactive job(s) at the highest priority

Q2

Q1

Q0

A Mixed I/O-intensive and CPU-intensive Workload (msec)

April 11, 2024
TCSS422: Operating Systems (Spring 2024)
School of Engineering and Technology, University of Washington - Tacoma
L6.36

36

WE WILL RETURN AT
4:55PM



April 11, 2024 TCSS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.37

37

OBJECTIVES - 4/11

- Questions from 4/11
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 - Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - **Job Starvation**
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCSS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.38

38

MLFQ: ISSUES

■ Starvation

[High Priority] Q8 → A → B → C → D → E → F

Q7

Q6

Q5

Q4

Q3

Q2

[Low Priority] Q1 → G → H *CPU bound batch job(s)*

April 11, 2024 TCSS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.39

39

OBJECTIVES - 4/11

- Questions from 4/11
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 - Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - **Gaming the Scheduler**
 - Examples
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCSS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.40

40

MLFQ: ISSUES - 2

- Gaming the scheduler
 - Issue I/O operation at 99% completion of the time slice
 - Keeps job priority fixed - never lowered
- Job behavioral change
 - CPU/batch process becomes an interactive process

Priority becomes stuck

[High Priority] Q8 → A → B → C → D → E → F

Q7

Q6

Q5

Q4

Q3

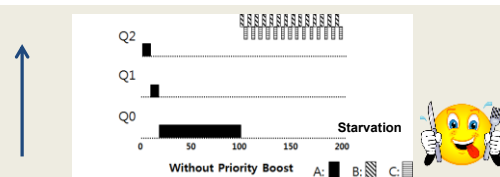
Q2

[Low Priority] Q1 → G → H *CPU bound batch job(s)*

April 11, 2024 TCSS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.41

41

RESPONDING TO BEHAVIOR CHANGE



Without Priority Boost

Starvation

- Priority Boost
 - Reset all jobs to topmost queue after some time interval S

April 11, 2024 TCSS422: Operating Systems (Spring 2024)
 School of Engineering and Technology, University of Washington - Tacoma L6.42

42

RESPONDING TO BEHAVIOR CHANGE - 2

- With priority boost
 - Prevents starvation

With Priority Boost

A: B: C:

April 11, 2024 TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.43

43

KEY TO UNDERSTANDING MLFQ – PB

- Without priority boost:
 - Rule 1:** If Priority(A) > Priority(B), A runs (B doesn't).
 - Rule 2:** If Priority(A) = Priority(B), A & B run in RR.
- KEY:** If time quantum of a higher queue is filled, then we don't run any jobs in lower priority queues!!!

April 11, 2024 TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.44

44

STARVATION EXAMPLE

- Consider 3 queues:
 - Q2 – HIGH PRIORITY – Time Quantum 10ms
 - Q1 – MEDIUM PRIORITY – Time Quantum 20 ms
 - Q0 – LOW PRIORITY – Time Quantum 40 ms
- Job A: 200ms no I/O
- Job B: 5ms then I/O
- Job C: 5ms then I/O
- Q2 fills up, starves Q1 & Q0
- A makes no progress

Without Priority Boost

A: B: C:

April 11, 2024 TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.45

45

PREVENTING GAMING

- Improved time accounting:
 - Track total job execution time in the queue
 - Each job receives a fixed time allotment
 - When allotment is exhausted, job priority is lowered

Without(Left) and With(Right) Gaming Tolerance

April 11, 2024 TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.46

46

MLFQ: TUNING

- Consider the tradeoffs:
 - How many queues?
 - What is a good time slice?
 - How often should we “Boost” priority of jobs?
 - What about different time slices to different queues?

Example) 10ms for the highest queue, 20ms for the middle, 40ms for the lowest

April 11, 2024 TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.47

47

PRACTICAL EXAMPLE

- Oracle Solaris MLFQ implementation
 - 60 Queues → w/ slowly increasing time slice (high to low priority)
 - Provides sys admins with set of editable table(s)
 - Supports adjusting time slices, boost intervals, priority changes, etc.
- Advice
 - Provide OS with hints about the process
 - Nice command → Linux

April 11, 2024 TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.48

48

MLFQ RULE SUMMARY

- The refined set of MLFQ rules:
- Rule 1:** If $Priority(A) > Priority(B)$, A runs (B doesn't).
- Rule 2:** If $Priority(A) = Priority(B)$, A & B run in RR.
- Rule 3:** When a job enters the system, it is placed at the highest priority.
- Rule 4:** Once a job uses up its time allotment at a given level (regardless of how many times it has given up the CPU), its priority is reduced (i.e., it moves down on queue).
- Rule 5:** After some time period S, move all the jobs in the system to the topmost queue.

April 11, 2024 TCCS422: Operating Systems (Spring 2024) School of Engineering and Technology, University of Washington - Tacoma L6.49

49

OBJECTIVES – 4/11

- Questions from 4/11
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 – Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - Gaming the Scheduler
 - Examples**
- Chapter 9: Proportional Share Schedulers

April 11, 2024 TCCS422: Operating Systems (Spring 2024) School of Engineering and Technology, University of Washington - Tacoma L6.50

50

Jackson deploys a 3-level MLFQ scheduler. The time slice is 1 for high priority jobs, 2 for medium priority, and 4 for low priority. This MLFQ scheduler performs a Priority Boost every 6 timer units. When the priority boost fires, the current job is preempted, and the next scheduled job is run in round-robin order.

| Job | Arrival Time | Job Length |
|-----|--------------|------------|
| A | T=0 | 4 |
| B | T=0 | 16 |
| C | T=0 | 8 |

(11 points) Show a scheduling graph for the MLFQ scheduler for the jobs above. Draw vertical lines for key events and be sure to label the X-axis times as in the example. Please draw clearly. An unreadable graph will lose points.

51

Jackson deploys a 3-level MLFQ scheduler. The time slice is 1 for high priority jobs, 2 for medium priority, and 4 for low priority. This MLFQ scheduler performs a Priority Boost every 6 timer units. When the priority boost fires, the current job is preempted, and the next scheduled job is run in round-robin order.

| Job | Arrival Time | Job Length |
|-----|--------------|------------|
| A | T=0 | 4 |
| B | T=0 | 16 |
| C | T=0 | 8 |

(11 points) Show a scheduling graph for the MLFQ scheduler for the jobs above. Draw vertical lines for key events and be sure to label the X-axis times as in the example. Please draw clearly. An unreadable graph will lose points.

52

Jackson deploys a 3-level MLFQ scheduler. The time slice is 1 for high priority jobs, 2 for medium priority, and 4 for low priority. This MLFQ scheduler performs a Priority Boost every 6 timer units. When the priority boost fires, the current job is preempted, and the next scheduled job is run in round-robin order.

| Job | Arrival Time | Job Length |
|-----|--------------|------------|
| A | T=0 | 4 |
| B | T=0 | 16 |
| C | T=0 | 8 |

(11 points) Show a scheduling graph for the MLFQ scheduler for the jobs above. Draw vertical lines for key events and be sure to label the X-axis times as in the example. Please draw clearly. An unreadable graph will lose points.

Handwritten notes: time slice is JOB time, Before C/S

53

EXAMPLE

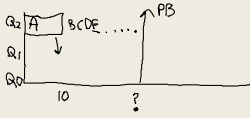
- Question:
- Given a system with a quantum length of 10 ms **for all jobs** in its highest queue, how often would you have to **boost jobs** back to the highest priority level to guarantee that a single long-running (and potentially starving) job gets at least 5% of the CPU?

April 11, 2024 TCCS422: Operating Systems (Spring 2024) School of Engineering and Technology, University of Washington - Tacoma L6.54

54

EXAMPLE

- Question:
- Given a system with a quantum length of 10 ms **for all jobs** in its highest queue, how often would you have to boost jobs back to the highest priority level to guarantee that a single long-running (and potentially starving) job gets at least 5% of the CPU?



$.05 PB = 10$
 $PB = \frac{10}{.05} = 200 \text{ ms}$

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.55

55

EXAMPLE

- Question:
- Given a system with a quantum length of 10 ms **for all jobs** in its highest queue, how often would you have to boost jobs back to the highest priority level to guarantee that a single long-running (and potentially starving) job gets at least 5% of the CPU?
- Some combination of n short jobs runs for a total of 10 ms per cycle without relinquishing the CPU
 - E.g. 2 jobs = 5 ms ea; 3 jobs = 3.33 ms ea, 10 jobs = 1 ms ea
 - n jobs always uses full time quantum in highest queue (10 ms)
 - Batch jobs starts, runs for full quantum of 10ms, pushed to lower queue
 - All other jobs run and context switch totaling the quantum per cycle
 - If 10ms is 5% of the CPU, when must the priority boost be ???

ANSWER → Priority boost should occur every 200ms

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.56

56


OBJECTIVES – 4/11

- Questions from 4/11
- Assignment 0
- C Tutorial - Pointers, Strings, Exec in C
- Quiz 1 – Active Reading Chapter 9
- Chapter 7: Scheduling Introduction
- Chapter 8: Multi-level Feedback Queue
 - MLFQ Scheduler
 - Job Starvation
 - Gaming the Scheduler
 - Examples
- Chapter 9: Proportional Share Schedulers**

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.57

57

CHAPTER 9 - PROPORTIONAL SHARE SCHEDULER



April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.58

58

OBJECTIVES – 4/11

- Chapter 9: Proportional Share Schedulers**
 - Lottery scheduler**
 - Ticket mechanisms
 - Stride scheduler
 - Linux Completely Fair Scheduler

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.59

59

PROPORTIONAL SHARE SCHEDULER

- Also called fair-share scheduler or lottery scheduler
- Guarantees each job receives some percentage of CPU time based on share of "tickets"
- Each job receives an allotment of tickets
- % of tickets corresponds to potential share of a resource
- Can conceptually schedule any resource this way
 - CPU, disk I/O, memory

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.60

60

LOTTERY SCHEDULER

- Simple implementation
 - Just need a random number generator
 - Picks the winning ticket
 - Maintain a data structure of jobs and tickets (list)
 - Traverse list to find the owner of the ticket
 - Consider sorting the list for speed

April 11, 2024
TCCS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L6.61

61

LOTTERY SCHEDULER IMPLEMENTATION

```

1 // counter: used to track if we've found the winner yet
2 int counter = 0;
3
4 // winner: use some call to a random number generator to
5 // get a value, between 0 and the total # of tickets
6 int winner = getrandom(0, totaltickets);
7
8 // current: use this to walk through the list of jobs
9 node_t *current = head;
10
11 // loop until the sum of ticket values is > the winner
12 while (current) {
13     counter = counter + current->tickets;
14     if (counter >= winner)
15         break; // found the winner
16     current = current->next;
17 }
18 // 'current' is the winner: schedule it...
```

April 11, 2024
TCCS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L6.62

62

OBJECTIVES – 4/11

- Chapter 9: Proportional Share Schedulers
 - Lottery scheduler
 - Ticket mechanisms
 - Stride scheduler
 - Linux Completely Fair Scheduler

April 11, 2024
TCCS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L6.63

63

TICKET MECHANISMS

- Ticket currency / exchange
 - User allocates tickets in any desired way
 - OS converts user currency into global currency
- Example:
 - There are 200 global tickets assigned by the OS

User A → 500 (A's currency) to A1 → 50 (global currency)
 → 500 (A's currency) to A2 → 50 (global currency)

User B → 10 (B's currency) to B1 → 100 (global currency)

April 11, 2024
TCCS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L6.64

64

TICKET MECHANISMS - 2

- Ticket transfer
 - Temporarily hand off tickets to another process
- Ticket inflation
 - Process can temporarily raise or lower the number of tickets it owns
 - If a process needs more CPU time, it can boost tickets.

April 11, 2024
TCCS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L6.65

65

LOTTERY SCHEDULING

- Scheduler picks a **winning** ticket
 - Load the job with the winning ticket and run it
- Example:
 - Given 100 tickets in the pool
 - Job A has 75 tickets: 0 - 74
 - Job B has 25 tickets: 75 - 99

Scheduler's winning tickets: 63 85 70 39 76 17 29 41 36 39 10 99 68 83 63

Scheduled job: A B A A B A A A A A A B A B A

- But what do we know about probability of a coin flip?

April 11, 2024
TCCS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L6.66

66

COIN FLIPPING

- Equality of distribution (fairness) requires a lot of flips!

Similarly, Lottery scheduling requires lots of "rounds" to achieve fairness.

Increasing number of coin tosses

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.67

67

LOTTERY FAIRNESS

- With two jobs
 - Each with the same number of tickets ($t=100$)

When the job length is not very long, average unfairness can be quite severe.

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.68

68

LOTTERY SCHEDULING CHALLENGES

- What is the best approach to assign tickets to jobs?
 - Typical approach is to assume users know best
 - Users are provided with tickets, which they allocate as desired
- How should the OS automatically distribute tickets upon job arrival?
 - What do we know about incoming jobs a priori ?
 - Ticket assignment is really an open problem...

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.69

69

OBJECTIVES – 4/11

- Chapter 9: Proportional Share Schedulers**
 - Lottery scheduler
 - Ticket mechanisms
 - Stride scheduler**
 - Linux Completely Fair Scheduler

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.70

70

STRIDE SCHEDULER

- Addresses statistical probability issues with lottery scheduling
- Instead of guessing a random number to select a job, simply count...

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.71

71

STRIDE SCHEDULER - 2

- Jobs have a "stride" value
 - A stride value describes the counter pace when the job should give up the CPU
 - Stride value is **Inverse In proportion** to the job's number of tickets (more tickets = smaller stride)
- Total system tickets = 10,000
 - Job A has 100 tickets $\rightarrow A_{stride} = 10000/100 = 100$ stride
 - Job B has 50 tickets $\rightarrow B_{stride} = 10000/50 = 200$ stride
 - Job C has 250 tickets $\rightarrow C_{stride} = 10000/250 = 40$ stride
- Stride scheduler tracks "pass" values for each job (A, B, C)

April 11, 2024 TCSS422: Operating Systems [Spring 2024]
 School of Engineering and Technology, University of Washington - Tacoma L6.72

72

STRIDE SCHEDULER - 3

- Basic algorithm:
 - Stride scheduler picks job with the lowest pass value
 - Scheduler increments job's pass value by its stride and starts running
 - Stride scheduler increments a counter
 - When counter exceeds pass value of current job, pick a new job (go to 1)
- KEY:** When the counter reaches a job's "PASS" value, the scheduler passes on to the next job...

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.73

73

STRIDE SCHEDULER - EXAMPLE

- Stride values
 - Tickets = priority to select job
 - Stride is inverse to tickets
 - Lower stride = more chances to run (higher priority)

Priority
 C stride = 40
 A stride = 100
 B stride = 200

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.74

74

STRIDE SCHEDULER EXAMPLE - 2

- Three-way tie:** randomly pick job A (all pass values=0)
- Set A's pass value to A's stride = 100
- Increment counter until > 100
- Pick a new job: two-way tie

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|-------------------------|-------------------------|------------------------|-----------|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | ... |

Tickets
 C = 250
 A = 100
 B = 50

← Initial job selection is random. All @ 0
 ← C has the most tickets and receives a lot of opportunities to run...

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.75

75

STRIDE SCHEDULER EXAMPLE - 3

- We set A's counter (pass value) to A's stride = 100
- Next scheduling decision between B (pass=0) and C (pass=0)
 - Randomly choose B
- C has the lowest counter for next 3 rounds

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|-------------------------|-------------------------|------------------------|-----------|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | ... |

Tickets
 C = 250
 A = 100
 B = 50

← C has the most tickets and is selected to run more often ...

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.76

76

STRIDE SCHEDULER EXAMPLE - 4

- Job counters support determining which job to run next
- Over time jobs are scheduled to run based on their priority represented as their share of tickets...
- Tickets are analogous to job priority**

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|-------------------------|-------------------------|------------------------|-----------|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | ... |

Tickets
 C = 250
 A = 100
 B = 50

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.77

77

OBJECTIVES - 4/11

- Chapter 9: Proportional Share Schedulers**
 - Lottery scheduler
 - Ticket mechanisms
 - Stride scheduler
 - Linux Completely Fair Scheduler**

April 11, 2024 TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma L6.78

78

LINUX: COMPLETELY FAIR SCHEDULER (CFS)

- Large Google datacenter study: "Profiling a Warehouse-scale Computer" (Kanev et al.)
- Monitored 20,000 servers over 3 years
- Found 20% of CPU time spent in the Linux kernel
- 5% of CPU time spent in the CPU scheduler!
- Study highlights importance for high performance OS kernels and CPU schedulers!

Figure 5: Kernel time, especially time spent in the scheduler, is a significant fraction of WSC cycles.

Ref: <https://dl.acm.org/cdoi/10.1145/3749448.3749449>

April 11, 2024 | TCS5422: Operating Systems (Spring 2024) | School of Engineering and Technology, University of Washington - Tacoma | L6.79

79

LINUX: COMPLETELY FAIR SCHEDULER (CFS)

- Loosely based on the stride scheduler
- CFS models system as a Perfect Multi-Tasking System
 - In perfect system every process of the same priority (class) receive exactly $1/n^{\text{th}}$ of the CPU time
- Each scheduling class has a runqueue
 - Groups process of same class
 - In class, scheduler picks task w/ lowest **vruntime** to run
 - Time slice varies based on how many jobs in shared runqueue
 - Minimum time slice prevents too many context switches (e.g. 3 ms)

April 11, 2024 | TCS5422: Operating Systems (Spring 2024) | School of Engineering and Technology, University of Washington - Tacoma | L6.80

80

COMPLETELY FAIR SCHEDULER - 2

- Every thread/process has a scheduling class (policy):
- Normal classes:** SCHED_OTHER (TS), SCHED_IDLE, SCHED_BATCH
 - TS = Time Sharing
- Real-time classes:** SCHED_FIFO (FF), SCHED_RR (RR)
- How to show scheduling class and priority:


```
#class
ps -elfc
```
- #priority (nice value)**

```
ps ax -o pid,ni,cls,pri,cmd
```

April 11, 2024 | TCS5422: Operating Systems (Spring 2024) | School of Engineering and Technology, University of Washington - Tacoma | L6.81

81

COMPLETELY FAIR SCHEDULER - 3

- Linux \geq 2.6.23: Completely Fair Scheduler (CFS)
- Linux $<$ 2.6.23: O(1) scheduler
- Linux maintains simple counter (vruntime) to track how long each thread/process has run
- CFS picks process with lowest vruntime to run next
- CFS adjusts timeslice based on # of proc waiting for the CPU
- Kernel parameters that specify CFS behavior:


```
$ sudo sysctl kernel.sched_latency_ns
kernel.sched_latency_ns = 24000000
$ sudo sysctl kernel.sched_min_granularity_ns
kernel.sched_min_granularity_ns = 3000000
$ sudo sysctl kernel.sched_wakeup_granularity_ns
kernel.sched_wakeup_granularity_ns = 4000000
```

April 11, 2024 | TCS5422: Operating Systems (Spring 2024) | School of Engineering and Technology, University of Washington - Tacoma | L6.82

82

COMPLETELY FAIR SCHEDULER - 4

- Sched_min_granularity_ns (3ms)**
 - Time slice for a process: busy system (w/ full runqueue)
 - If system has idle capacity, time slice exceed the min as long as difference in **vruntime** between running process and process with lowest **vruntime** is less than **sched_wakeup_granularity_ns** (4ms)
- Scheduling time period is: total cycle time for iterating through a set of processes where each is allowed to run (like round robin)
- Example:


```
sched_latency_ns (24ms)
if (proc in runqueue < sched_latency_ns / sched_min_granularity)
or
sched_min_granularity * number of processes in runqueue
```

Ref: https://www.systemd.io/sched_min_granularity_ns-sched_latency_ns-cfs-effect-timeline-processes/

April 11, 2024 | TCS5422: Operating Systems (Spring 2024) | School of Engineering and Technology, University of Washington - Tacoma | L6.83

83

CFS TRADEOFF

- HIGH**
 - sched_min_granularity_ns (timeslice)
 - sched_latency_ns
 - sched_wakeup_granularity_ns

reduced context switching → less overhead
 poor near-term fairness
- LOW**
 - sched_min_granularity_ns (timeslice)
 - sched_latency_ns
 - sched_wakeup_granularity_ns

increased context switching → more overhead
 better near-term fairness

April 11, 2024 | TCS5422: Operating Systems (Spring 2024) | School of Engineering and Technology, University of Washington - Tacoma | L6.84

84

COMPLETELY FAIR SCHEDULER - 5

- Runqueues are stored using a linux red-black tree
 - Self balancing binary tree - nodes indexed by **vruntime**
- Leftmost node has lowest **vruntime** (approx execution time)
- Walking tree to find left most node has very low big O complexity: $\sim O(\log N)$ for N nodes
- Completed processes removed

Nodes represent `sched_entity(s)` indexed by their virtual runtime

Virtual runtime

← Most need of CPU | Least need of CPU →

| | | |
|----------------|---|-------|
| April 11, 2024 | TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L6.85 |
|----------------|---|-------|

85

CFS: JOB PRIORITY

- Time slice: Linux "**Nice value**"
 - Nice predates the CFS scheduler
 - Top shows nice values
 - Process command (nice & priority):
`ps ax -o pid,ni,cmd,%cpu, pri`
- Nice Values: from -20 to 19
 - Lower is **higher** priority, default is 0
 - Vruntime is a weighted time measurement
 - Priority weights the calculation of vruntime within a runqueue to give high priority jobs a boost.
 - Influences job's position in rb-tree

```
static const int prio_to_weight[40] = {
/* -20 */ 88761, 71755, 56483, 46273, 36291,
/* -15 */ 29154, 23254, 18702, 14949, 11914,
/* -10 */ 9548, 7620, 6100, 4904, 3906,
/* -5 */ 3123, 2501, 1991, 1586, 1272,
/* 0 */ 1024, 820, 655, 526, 423,
/* 5 */ 335, 272, 219, 172, 137,
/* 10 */ 110, 87, 70, 56, 45,
/* 15 */ 36, 29, 23, 18, 15,
};
```

| | | |
|----------------|---|-------|
| April 11, 2024 | TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L6.86 |
|----------------|---|-------|

86

COMPLETELY FAIR SCHEDULER - 6

- CFS tracks cumulative job run time in **vruntime** variable
- The task on a given runqueue with the lowest **vruntime** is scheduled next
- `struct sched_entity` contains **vruntime** parameter
 - Describes process execution time in nanoseconds
 - Value is not pure runtime, is weighted based on job priority
 - Perfect scheduler → achieve equal vruntime for all processes of same priority
- Sleeping jobs: upon return reset vruntime to lowest value in system
 - Jobs with frequent short sleep **SUFFER !!**
- Key takeaway:
Identifying the next job to schedule is really fast!

| | | |
|----------------|---|-------|
| April 11, 2024 | TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L6.87 |
|----------------|---|-------|

87

COMPLETELY FAIR SCHEDULER - 7

- More information:
 - Man page: "man sched" : Describes Linux scheduling API
 - <http://manpages.ubuntu.com/manpages/bionic/man7/sched.7.html>
 - <https://www.kernel.org/doc/Documentation/scheduler/sched-design-CFS.txt>
 - https://en.wikipedia.org/wiki/Completely_Fair_Scheduler
 - See paper: The Linux Scheduler – a Decade of Wasted Cores
 - <http://www.ece.ubc.ca/~sasha/papers/eurosys16-final29.pdf>

| | | |
|----------------|---|-------|
| April 11, 2024 | TCCS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L6.88 |
|----------------|---|-------|

88

QUESTIONS

89