

CPU VIRTUALIZATION: SCHEDULING

Questions answered in this lecture:

- What are different scheduling policies, such as: FCFS, SJF, STCF, RR and MLFQ?
- What type of workload performs well with each scheduler?

ANNOUNCEMENTS

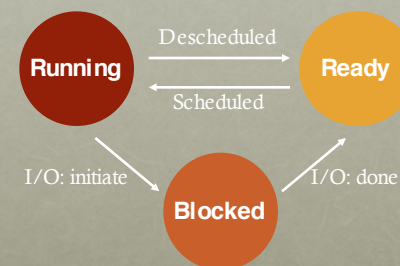
- Reading:
 - Today cover Chapters 7-9
- Project 1: Sorting and System Calls
 - Sorting : Warm-up with using C
 - Finish Part A this week
 - Competition:
 - Free text book or t-shirt to fastest (average) sort in each discussion section
 - Handin directories not yet available
 - Goal is for everyone to learn material
 - Do not copy code from others!

CPU VIRTUALIZATION: TWO COMPONENTS

Dispatcher (Previous lecture)

- Low-level mechanism
- Performs context-switch
 - Switch from user mode to kernel mode
 - Save execution state (registers) of old process in PCB
 - Insert PCB in ready queue
 - Load state of next process from PCB to registers
 - Switch from kernel to user mode
 - Jump to instruction in new user process
- Scheduler (Today)
 - Policy to determine which process gets CPU when

REVIEW: STATE TRANSITIONS



How to transition? (“mechanism”)
When to transition? (“policy”)

VOCABULARY

Workload: set of **job** descriptions (arrival time, run_time)

- Job: View as current CPU burst of a process
- Process alternates between CPU and I/O
process moves between ready and blocked queues

Scheduler: logic that decides which ready job to run

Metric: measurement of scheduling quality

SCHEDULING PERFORMANCE METRICS

Minimize turnaround time

- Do not want to wait long for job to complete
- $\text{Completion_time} - \text{arrival_time}$

Minimize response time

- Schedule interactive jobs promptly so users see output quickly
- $\text{Initial_schedule_time} - \text{arrival_time}$

Minimize waiting time

- Do not want to spend much time in Ready queue

Maximize throughput

- Want many jobs to complete per unit of time

Maximize resource utilization

- Keep expensive devices busy

Minimize overhead

- Reduce number of context switches

Maximize fairness

- All jobs get same amount of CPU over some time interval

WORKLOAD ASSUMPTIONS

1. Each job runs for the same amount of time
2. All jobs arrive at the same time
3. All jobs only use the CPU (no I/O)
4. Run-time of each job is known

SCHEDULING BASICS

Workloads:

arrival_time
run_time

Schedulers:

FIFO
SJF
STCF
RR

Metrics:

turnaround_time
response_time

EXAMPLE: WORKLOAD, SCHEDULER, METRIC

JOB	arrival_time (s)	run_time (s)
A	~0	10
B	~0	10
C	~0	10

FIFO: First In, First Out

- also called FCFS (first come first served)
- run jobs in *arrival_time* order

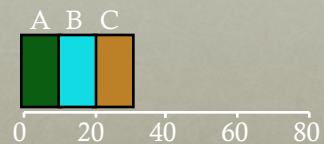
What is our **turnaround?**: $completion_time - arrival_time$

FIFO: EVENT TRACE

JOB	arrival_time (s)	run_time (s)	Time	Event
A	~0	10	0	A arrives
B	~0	10	0	B arrives
C	~0	10	0	C arrives
			0	run A
			10	complete A
			10	run B
			20	complete B
			20	run C
			30	complete C

FIFO (IDENTICAL JOBS)

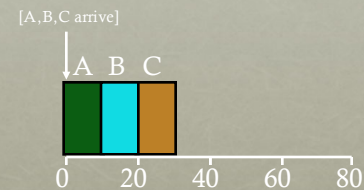
JOB	arrival_time (s)	run_time (s)
A	~0	10
B	~0	10
C	~0	10



Gantt chart:

Illustrates how jobs are scheduled over time on a CPU

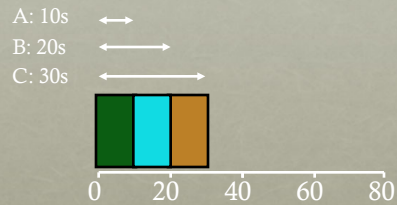
FIFO (IDENTICAL JOBS)



What is the average turnaround time?

Def: $turnaround_time = completion_time - arrival_time$

FIFO (IDENTICAL JOBS)



What is the average turnaround time?
Def: $turnaround_time = completion_time - arrival_time$
 $(10 + 20 + 30) / 3 = 20s$

SCHEDULING BASICS

Workloads:
arrival_time
run_time

Schedulers:
FIFO
SJF
STCF
RR

Metrics:
turnaround_time
response_time

WORKLOAD ASSUMPTIONS

1. Each job runs for the same amount of time
2. All jobs arrive at the same time
3. All jobs only use the CPU (no I/O)
4. The run-time of each job is known

ANY PROBLEMATIC WORKLOADS FOR FIFO?

Workload: ?

Scheduler: FIFO

Metric: turnaround is high

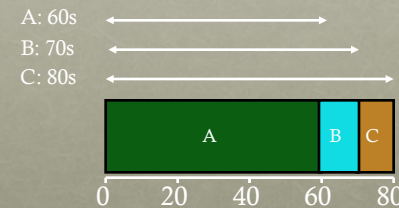
EXAMPLE: BIG FIRST JOB

JOB	arrival_time (s)	run_time (s)
A	~0	60
B	~0	10
C	~0	10

Draw Gantt chart for this workload and policy...
What is the average turnaround time?

EXAMPLE: BIG FIRST JOB

JOB	arrival_time (s)	run_time (s)
A	~0	60
B	~0	10
C	~0	10



Average turnaround time: **70s**

CONVOY EFFECT



PASSING THE TRACTOR

Problem with Previous Scheduler:

FIFO: Turnaround time can suffer when short jobs must wait for long jobs

New scheduler:

SJF (Shortest Job First)

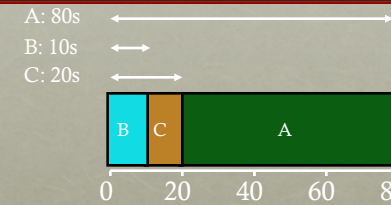
Choose job with smallest *run_time*

SHORTEST JOB FIRST

JOB	arrival_time (s)	run_time (s)
A	~0	60
B	~0	10
C	~0	10

What is the average turnaround time with SJF?

SJF TURNAROUND TIME



What is the average turnaround time with SJF?

$$(80 + 10 + 20) / 3 = \sim 36.7s \quad \text{Average turnaround with FIFO: 70s}$$

For minimizing average turnaround time (with no preemption):
SJF is provably optimal

Moving shorter job before longer job improves turnaround time of short job more than it harms turnaround time of long job

SCHEDULING BASICS

Workloads:	Schedulers:	Metrics:
arrival_time	FIFO	turnaround_time
run_time	SJF	response_time
	STCF	
	RR	

WORKLOAD ASSUMPTIONS

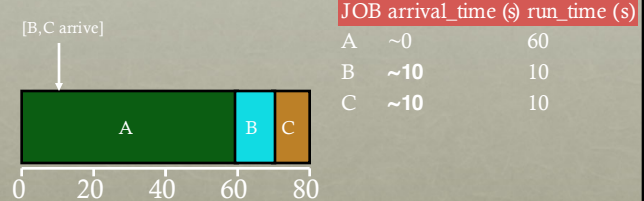
1. Each job runs for the same amount of time
2. All jobs arrive at the same time
3. All jobs only use the CPU (no I/O)
4. The run-time of each job is known

SHORTEST JOB FIRST (ARRIVAL TIME)

JOB	arrival_time (s)	run_time (s)
A	~0	60
B	~10	10
C	~10	10

What is the average turnaround time with SJF?

STUCK BEHIND A TRACTOR AGAIN



What is the average turnaround time?

$$(60 + (70 - 10) + (80 - 10)) / 3 = 63.3s$$

PREEMPTIVE SCHEDULING

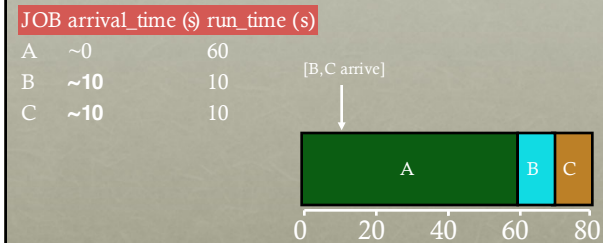
Prev schedulers:

- FIFO and SJF are non-preemptive
- Only schedule new job when previous job voluntarily relinquishes CPU (performs I/O or exits)

New scheduler:

- Preemptive: Potentially schedule different job at any point by taking CPU away from running job
- STCF (Shortest Time-to-Completion First)
- Always run job that will complete the quickest

NON-PREEMPTIVE: SJF



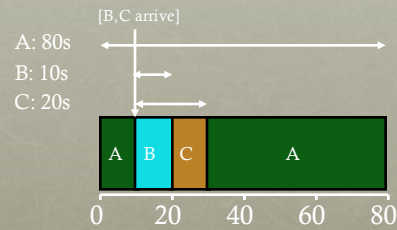
Average turnaround time:

$$(60 + (70 - 10) + (80 - 10)) / 3 = 63.3s$$

PREEMPTIVE: STCF

JOB arrival_time (s) run_time (s)

A ~0 60
B ~10 10
C ~10 10



Average turnaround time with STCF?

36.6

Average turnaround time with SJF: **63.3s**

SCHEDULING BASICS

Workloads:

arrival_time
run_time

Schedulers:

FIFO
SJF
STCF
RR

Metrics:

turnaround_time
response_time

RESPONSE TIME

Sometimes care about when job starts instead of when it finishes

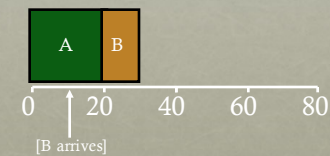
New metric:

$$response_time = first_run_time - arrival_time$$

RESPONSE VS. TURNAROUND

B's turnaround: 20s ←→

B's response: 10s ←→



ROUND-ROBIN SCHEDULER

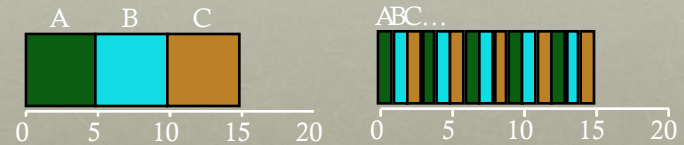
Prev schedulers:

FIFO, SJF, and STCF can have poor response time

New scheduler: RR (Round Robin)

Alternate ready processes every fixed-length time-slice

FIFO VS RR



Avg Response Time?
 $(0+5+10)/3 = 5$

Avg Response Time?
 $(0+1+2)/3 = 1$

In what way is RR worse?

Ave. turn-around time with equal job lengths is horrible

Other reasons why RR could be better?

If don't know run-time of each job, gives short jobs a chance to run and finish fast

SCHEDULING BASICS

Workloads:

arrival_time
run_time

Schedulers:

FIFO
SJF
STCF
RR

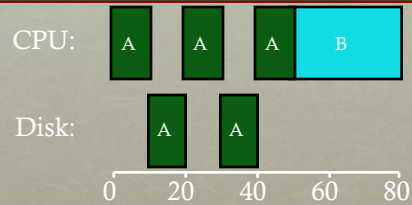
Metrics:

turnaround_time
response_time

WORKLOAD ASSUMPTIONS

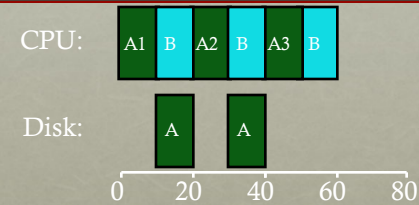
1. Each job runs for the same amount of time
2. All jobs arrive at the same time
3. All jobs only use the CPU (no I/O)
4. The run-time of each job is known

NOT I/O AWARE



Don't let Job A hold on to CPU while blocked waiting for disk

I/O AWARE (OVERLAP)



Treat Job A as 3 separate CPU bursts
When Job A completes I/O, another Job A is ready
Each CPU burst is shorter than Job B, so with SCTF,
Job A preempts Job B

WORKLOAD ASSUMPTIONS

1. Each job runs for the same amount of time
2. All jobs arrive at the same time
3. All jobs only use the CPU (no I/O)
4. The run time of each job is known
(need smarter, fancier scheduler)

MLFQ (MULTI-LEVEL FEEDBACK QUEUE)

Goal: general-purpose scheduling

Must support two job types with distinct goals

- "interactive" programs care about response time
- "batch" programs care about turnaround time

Approach: multiple levels of round-robin;
each level has higher priority than lower levels and preempts them

PRIORITIES

Rule 1: If $\text{priority}(A) > \text{Priority}(B)$, A runs

Rule 2: If $\text{priority}(A) == \text{Priority}(B)$, A & B run in RR

Q3 → A

“Multi-level”

Q2 → B

How to know how to set priority?

Q1

Approach 1: nice

Q0 → C → D

Approach 2: history “feedback”

HISTORY

- Use past behavior of process to predict future behavior
 - Common technique in systems
- Processes alternate between I/O and CPU work
- Guess how CPU burst (job) will behave based on past CPU bursts (jobs) of this process

MORE MLFQ RULES

Rule 1: If $\text{priority}(A) > \text{Priority}(B)$, A runs

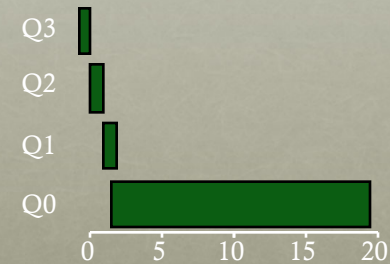
Rule 2: If $\text{priority}(A) == \text{Priority}(B)$, A & B run in RR

More rules:

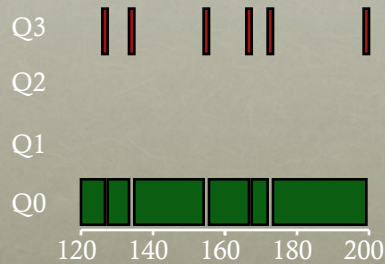
Rule 3: Processes start at top priority

Rule 4: If job uses whole slice, demote process
(longer time slices at lower priorities)

ONE LONG JOB (EXAMPLE)

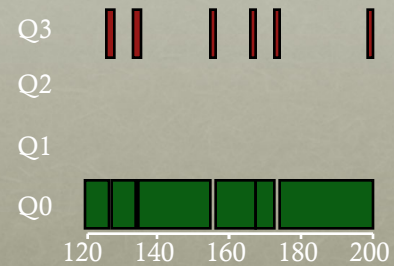


AN INTERACTIVE PROCESS JOINS



Interactive process never uses entire time slice, so never demoted

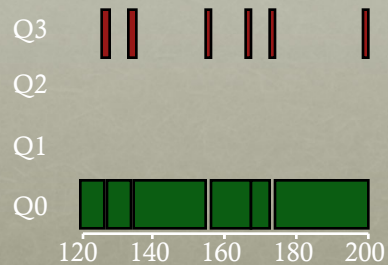
PROBLEMS WITH MLFQ?



Problems

- unforgiving + starvation
- gaming the system

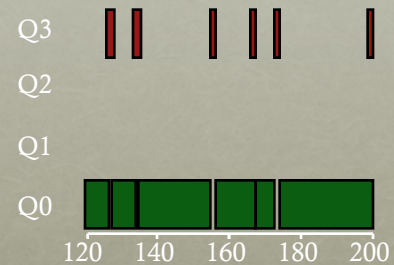
PREVENT STARVATION



Problem: Low priority job may never get scheduled

Periodically boost priority of all jobs (or all jobs that haven't been scheduled)

PREVENT GAMING



Problem: High priority job could trick scheduler and get more CPU by performing I/O right before time-slice ends

Fix: Account for job's total run time at priority level (instead of just this time slice); downgrade when exceed threshold

LOTTERY SCHEDULING

Goal: proportional (fair) share

Approach:

- give processes lottery tickets
- whoever wins runs
- higher priority => more tickets

Amazingly simple to implement

LOTTERY CODE

```
int counter = 0;
int winner = getrandom(0, totaltickets);
node_t *current = head;
while (current) {
    counter += current->tickets;
    if (counter > winner) break;
    current = current->next;
}
// current is the winner
```

LOTTERY EXAMPLE

```
int counter = 0;
int winner = getrandom(0, totaltickets);
node_t *current = head;
while(current) {
    counter += current->tickets;
    if (counter > winner) break;
    current = current->next;
}
// current gets to run
```

Who runs if winner is:
50
350
0



OTHER LOTTERY IDEAS

Ticket Transfers

Ticket Currencies

Ticket Inflation

(read more in OSTEP)

SUMMARY

Understand goals (metrics) and workload, then design scheduler around that

General purpose schedulers need to support processes with different goals

Past behavior is good predictor of future behavior

Random algorithms (lottery scheduling) can be simple to implement, and avoid corner cases.