

Scheduling: metrics / FCFS+RR / SJF+SRTF

last time

finish pipe / read / write

read: wait till something available, then copy what's available

multithreaded process

xv6: only single-threaded processes

thread: registers, program counter, stack

process: open files, memory contents (heap, etc.), process id, etc.

xv6 scheduler organization

separate scheduler thread

switch to scheduler thread, scheduler thread finds next process

scheduler thread iterates through available processes

xv6 scheduler locking

don't look at thread state (running/waiting/etc.) while it might be changing

don't change thread state while something else might be looking at it

goal/policy/mechanism

end of last time — used terms *policy* and *mechanism* imprecisely

goal: properties we want schedule to satisfy

policy: how we choose to schedule to do it

mechanism: how we switch processes, do bookkeeping, etc.

the scheduling policy problem

what RUNNABLE program should we run?

xv6 answer: whatever's next in list

best answer?

well, what should we care about?

some simplifying assumptions

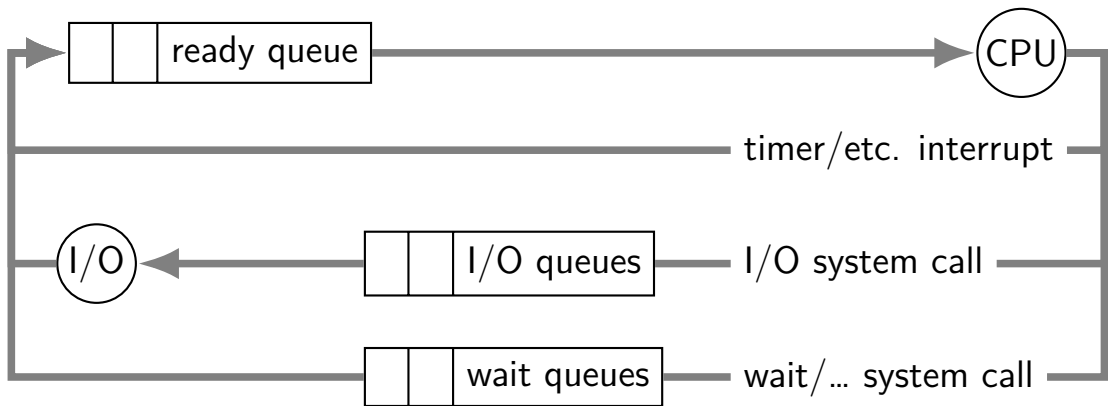
welcome to 1970:

one program per user

one thread per program

programs are independent

recall: scheduling queues



CPU and I/O bursts

...

compute

start read

(from file/keyboard/...)

wait for I/O

compute on read data

start read

wait for I/O

compute on read data

start write

wait for I/O

...

program alternates between computing and waiting for I/O

examples:

shell: wait for keypresses

drawing program: wait for mouse presses/etc.

web browser: wait for remote web server

...

CPU bursts and interactivity (one c. 1966 shared system)

shows compute time
from command entered
until next command prompt

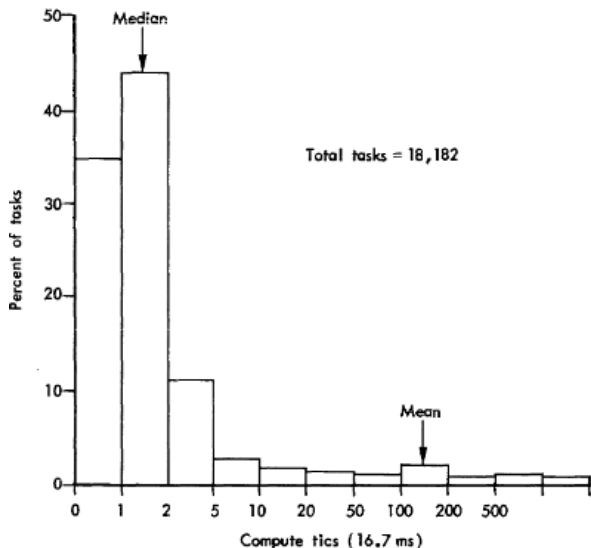
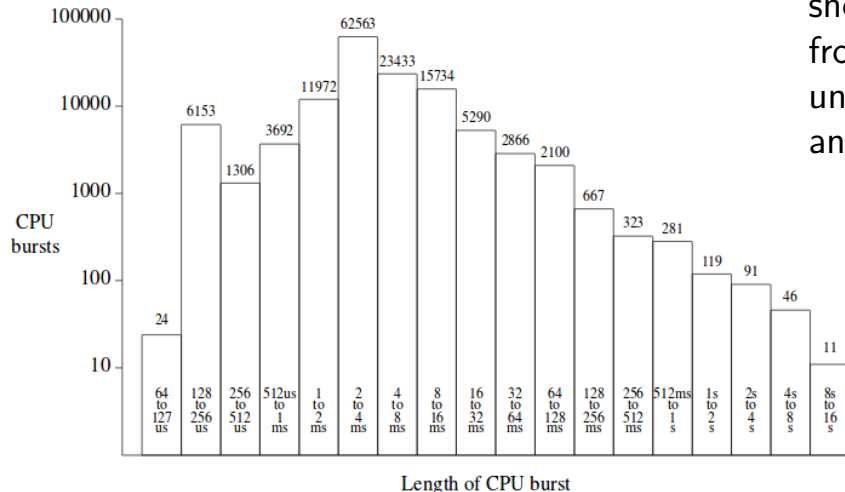


Figure 11—Compute time per task

CPU bursts and interactivity (one c. 1990 desktop)



shows CPU time from RUNNING until not RUNNABLE anymore

CPU bursts

observation: applications alternate between I/O and CPU

especially interactive applications

but also, e.g., reading and writing from disk

typically short “CPU bursts” (milliseconds) followed by short “IO bursts” (milliseconds)

scheduling CPU bursts

our typical view: ready queue, bunch of CPU bursts to run

to start: just look at running what's currently in ready queue best
same problem as 'run bunch of programs to completion'?

later: account for I/O after CPU burst

an historical note

historically applications were less likely to keep all data in memory

historically computers shared between more users

meant *more* applications alternating I/O and CPU

context many scheduling policies were developed in

scheduling metrics

response time (Anderson-Dahlin) AKA **turnaround time**
(Arpaci-Dusseau) (want *low*)

(what Arpaci-Dusseau calls response time is slightly different — more later)

what user sees: from *keypress* to *character on screen*
(submission until job finished)

throughput (want *high*)

total work per second

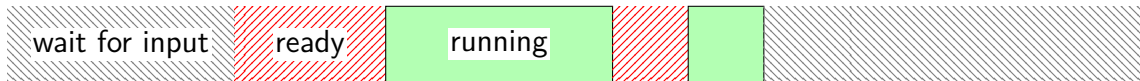
problem: overhead (e.g. from context switching)

fairness

many definitions

all **conflict** with best average throughput/turnaround time

turnaround and wait time

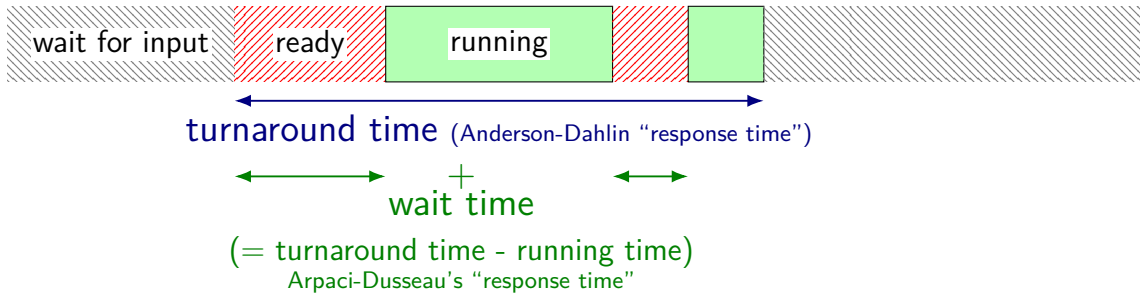


turnaround time (Anderson-Dahlin "response time")

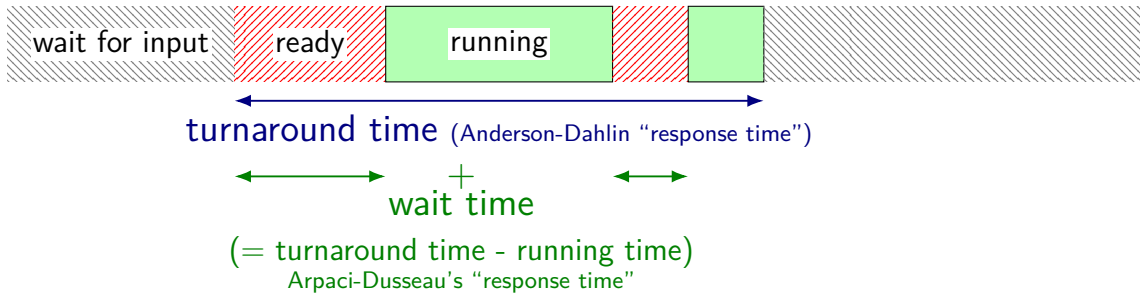
+
wait time

(= turnaround time - running time)
Arpaci-Dusseau's "response time"

turnaround and wait time

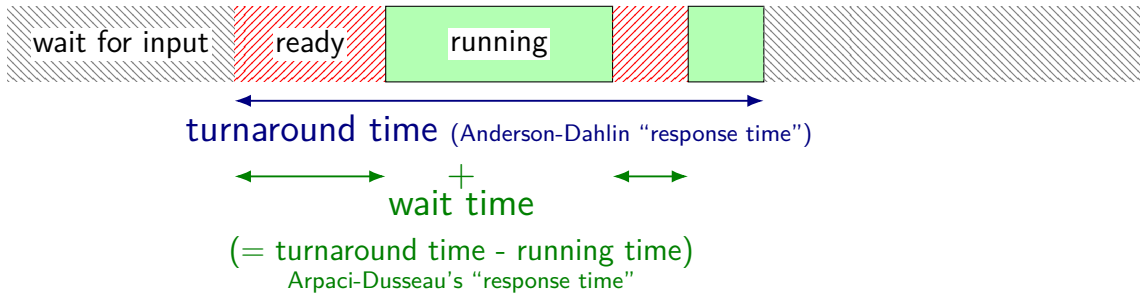


turnaround and wait time



common measure: *mean* turnaround time or *total* turnaround time

turnaround and wait time



common measure: *mean* turnaround time or *total* turnaround time

same as optimizing mean/total waiting time

turnaround time and I/O

scheduling CPU bursts? (what we'll mostly deal with)

turnaround time \approx time to start next I/O

important for fully utilizing I/O devices

closed loop: faster turnaround time \rightarrow program requests CPU sooner

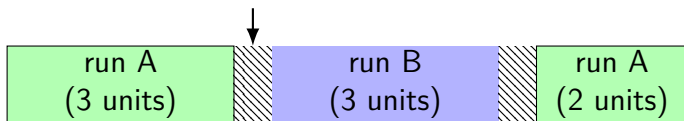
scheduling batch program on cluster?

turnaround time \approx how long does user wait

once program done with CPU, it's probably done

throughput

context switch(each .5 units)



throughput: **useful work** done per unit time

$$\text{non-context switch CPU utilization} = \frac{3 + 3 + 2}{3 + .5 + 3 + .5 + 2} = 88\%$$

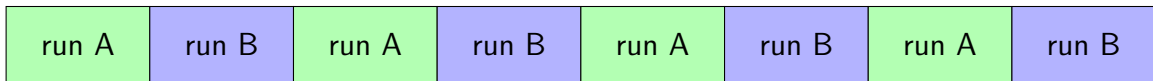
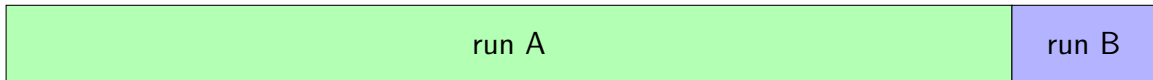
also other considerations:

- time lost due to cold caches

- time lost not starting I/O early as possible

- ...

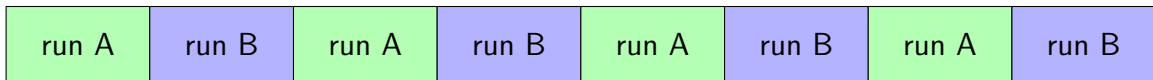
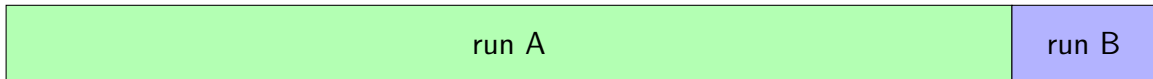
fairness



assumption: one program per user

two timelines above; which is fairer?

fairness



assumption: one program per user

two timelines above; which is fairer?

easy to answer — but formal definition?

two trivial scheduling algorithms

first-come first served (FCFS)

round robin (RR)

scheduling example assumptions

multiple programs become ready at almost the same time
alternately: became ready while previous program was running

...but in some order that we'll use
e.g. our ready queue looks like a linked list

two trivial scheduling algorithms

first-come first served (FCFS)

round robin (RR)

first-come, first-served

simplest(?) scheduling algorithm

no preemption — run program until it can't

suitable in cases where no context switch

e.g. not enough memory for two active programs

first-come, first-served (FCFS)

(AKA “first in, first out” (FIFO))

process	CPU time needed
A	24
B	4
C	3

first-come, first-served (FCFS)

(AKA “first in, first out” (FIFO))

process	CPU time needed	
A	24	} A ~ CPU-bound } B, C ~ I/O bound or interactive
B	4	
C	3	

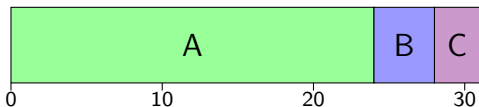
first-come, first-served (FCFS)

(AKA “first in, first out” (FIFO))

process	CPU time needed
A	24
B	4
C	3

} A ~ CPU-bound
B, C ~ I/O bound or interactive

arrival order: **A**, **B**, **C**



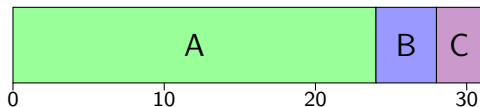
first-come, first-served (FCFS)

(AKA “first in, first out” (FIFO))

process	CPU time needed
A	24
B	4
C	3

} A ~ CPU-bound
B, C ~ I/O bound or interactive

arrival order: **A**, **B**, **C**



waiting times: (mean=17.3)

0 (**A**), 24 (**B**), 28 (**C**)

turnaround times: (mean=27.7)

24 (**A**), 28 (**B**), 31 (**C**)

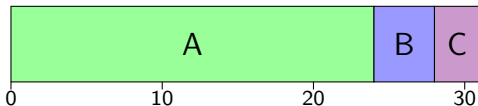
first-come, first-served (FCFS)

(AKA “first in, first out” (FIFO))

process	CPU time needed
A	24
B	4
C	3

} A ~ CPU-bound
B, C ~ I/O bound or interactive

arrival order: **A**, **B**, **C**



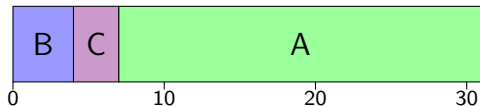
waiting times: (mean=17.3)

0 (**A**), 24 (**B**), 28 (**C**)

turnaround times: (mean=27.7)

24 (**A**), 28 (**B**), 31 (**C**)

arrival order: **B**, **C**, **A**



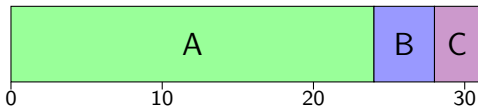
first-come, first-served (FCFS)

(AKA “first in, first out” (FIFO))

process	CPU time needed
A	24
B	4
C	3

A ~ CPU-bound
B, C ~ I/O bound or interactive

arrival order: **A**, **B**, **C**



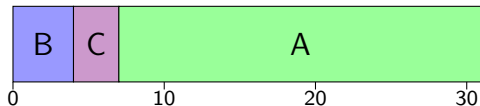
waiting times: (mean=17.3)

0 (**A**), 24 (**B**), 28 (**C**)

turnaround times: (mean=27.7)

24 (**A**), 28 (**B**), 31 (**C**)

arrival order: **B**, **C**, **A**



waiting times: (mean=3.7)

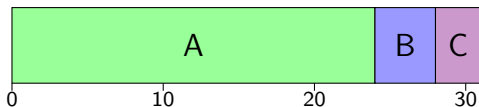
7 (**A**), 0 (**B**), 4 (**C**)

turnaround times: (mean=14)

31 (**A**), 4 (**B**), 7 (**C**)

FCFS orders

arrival order: **A**, **B**, **C**



waiting times: (mean=17.3)

0 (**A**), 24 (**B**), 28 (**C**)

turnaround times: (mean=27.7)

24 (**A**), 28 (**B**), 31 (**C**)

arrival order: **B**, **C**, **A**



waiting times: (mean=3.7)

7 (**A**), 0 (**B**), 4 (**C**)

turnaround times: (mean=14)

31 (**A**), 3 (**B**), 7 (**C**)

“convoy effect”

two trivial scheduling algorithms

first-come first served (FCFS)

round robin (RR)

round-robin

simplest(?) preemptive scheduling algorithm

run program until either

- it can't run anymore, or

- it runs for too long (exceeds "time quantum")

requires good way of interrupting programs

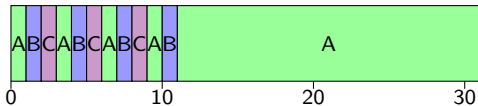
- like xv6's timer interrupt

requires good way of stopping programs whenever

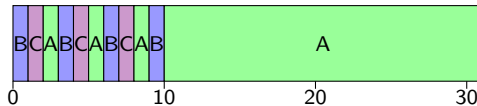
- like xv6's context switches

round robin (RR) (varying order)

time quantum = 1,
order **A**, **B**, **C**

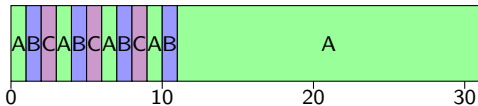


time quantum = 1,
order **B**, **C**, **A**



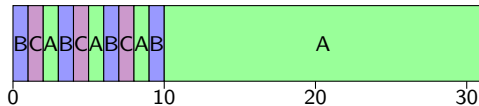
round robin (RR) (varying order)

time quantum = 1,
order **A**, **B**, **C**



waiting times: (mean=6.7)
7 (**A**), 7 (**B**), 6 (**C**)
turnaround times: (mean=17)
31 (**A**), 11 (**B**), 9 (**C**)

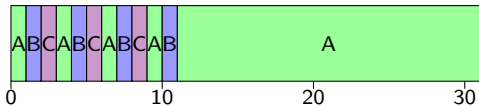
time quantum = 1,
order **B**, **C**, **A**



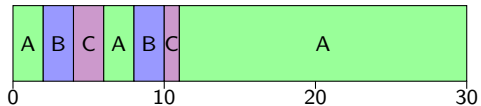
waiting times: (mean=6)
7 (**A**), 6 (**B**), 5 (**C**)
turnaround times: (mean=16.3)
31 (**A**), 10 (**B**), 8 (**C**)

round robin (RR) (varying time quantum)

time quantum = 1,
order **A**, **B**, **C**

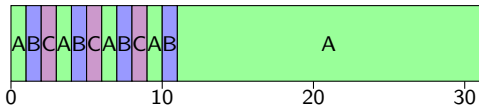


time quantum = 2,
order **A**, **B**, **C**



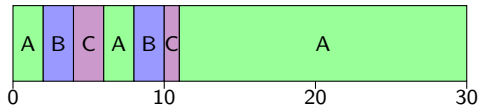
round robin (RR) (varying time quantum)

time quantum = 1,
order **A**, **B**, **C**



waiting times: (mean=6.7)
7 (**A**), 7 (**B**), 6 (**C**)
turnaround times: (mean=17)
31 (**A**), 11 (**B**), 9 (**C**)

time quantum = 2,
order **A**, **B**, **C**



waiting times: (mean=7)
7 (**A**), 6 (**B**), 8 (**C**)
turnaround times: (mean=17.3)
31 (**A**), 10 (**B**), 11 (**C**)

round robin idea

choose fixed time quantum Q

unanswered question: what to choose

switch to next process in ready queue after time quantum expires

this policy is what xv6 scheduler does

scheduler runs from timer interrupt (or if process not runnable)

finds next runnable process in process table

round robin and time quantum

many context switches
(lower throughput)

few context switches
(higher throughput)

order doesn't matter
(more fair)

first program favored
(less fair)

RR with
short quantum



FCFS

smaller quantum: more fair, worse throughput

round robin and time quantum

many context switches
(lower throughput)

few context switches
(higher throughput)

order doesn't matter
(more fair)

first program favored
(less fair)

RR with
short quantum



FCFS

smaller quantum: more fair, worse throughput

FCFS = RR with infinite quantum

more fair: at most $(N - 1)Q$ time until scheduled if N total processes

aside: context switch overhead

typical context switch: ~ 0.01 ms to 0.1 ms

but tricky: lot of indirect cost (cache misses)

(above numbers try to include likely indirect costs)

choose time quantum to manage this overhead

current Linux default: between ~ 0.75 ms and ~ 6 ms

varied based on number of active programs

Linux's scheduler is more complicated than RR

historically common: 1 ms to 100 ms

round robin and time quantum

many context switches
(lower throughput)

few context switches
(higher throughput)

order doesn't matter
(more fair)

first program favored
(less fair)

RR with
short quantum



FCFS

smaller quantum: more fair, worse throughput

FCFS = RR with infinite quantum

more fair: at most $(N - 1)Q$ time until scheduled if N total processes

but what about **turnaround/waiting time?**

exercise: round robin quantum

if there were no context switch overhead, *decreasing* the time quantum (for round robin) would cause average turnaround time to _____.

- A. always decrease or stay the same
- B. always increase or stay the same
- C. increase or decrease or stay the same
- D. something else?

increase turnaround time

A: 1 unit CPU burst

B: 1 unit

$Q = 1$



mean turnaround time =
 $(1 + 2) \div 2 = 1.5$

$Q = 1/2$



mean turnaround time =
 $(1.5 + 2) \div 2 = 1.75$

decrease turnaround time

A: 10 unit CPU burst

B: 1 unit



mean turnaround time =
 $(10 + 11) \div 2 = 10.5$



mean turnaround time =
 $(6 + 11) \div 2 = 8.5$

stay the same

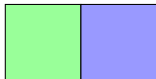
A: 1 unit CPU burst

B: 1 unit

$Q = 10$



$Q = 1$



FCFS and order

earlier we saw that with FCFS, arrival order mattered

big changes in turnaround/waiting time

let's use that insight to see how to optimize *mean* turnaround times

FCFS orders

arrival order: **A**, **B**, **C**



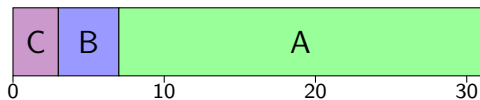
waiting times: (mean=17.3)

0 (**A**), 24 (**B**), 28 (**C**)

turnaround times: (mean=27.7)

24 (**A**), 28 (**B**), 31 (**C**)

arrival order: **B**, **C**, **A**



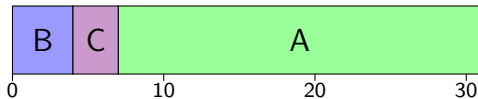
waiting times: (mean=3.3)

7 (**A**), 3 (**B**), 0 (**C**)

turnaround times: (mean=13.7)

31 (**A**), 7 (**B**), 3 (**C**)

arrival order: **B**, **C**, **A**



waiting times: (mean=3.7)

7 (**A**), 0 (**B**), 4 (**C**)

turnaround times: (mean=14)

31 (**A**), 4 (**B**), 7 (**C**)

order and turnaround time

best turnaround time = run shortest CPU burst first

worst turnaround time = run longest CPU burst first

intuition: “race to go to sleep”

diversion: some users are more equal

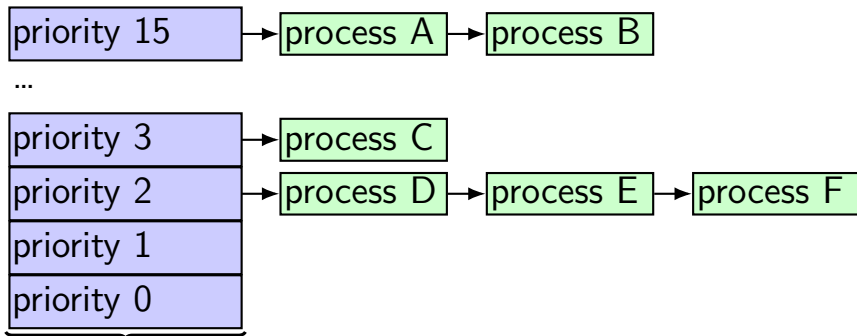
shells more important than big computation?

i.e. programs with short CPU bursts

faculty more important than students?

scheduling algorithm: schedule shells/faculty programs first

priority scheduling



ready queues for each priority level

choose process from **ready queue for highest priority**

within each priority, use some other scheduling (e.g. round-robin)

could have each process have unique priority

priority scheduling and preemption

priority scheduling can be preemptive

i.e. higher priority program comes along — stop whatever else was running

exercise: priority scheduling (1)

Suppose there are two processes:

process A

highest priority

repeat forever: 1 unit of I/O, then 10 units of CPU, ...

process Z

lowest priority

4000 units of CPU (and no I/O)

How long will it take process Z complete?

exercise: priority scheduling (2)

Suppose there are three processes:

process A

highest priority

repeat forever: 1 unit of I/O, then 10 units of CPU, ...

process B

second-highest priority

repeat forever: 1 unit of I/O, then 10 units of CPU, ...

process Z

lowest priority

4000 units of CPU (and no I/O)

How long will it take process Z complete?

starvation

programs can get “starved” of resources

never get those resources because of higher priority

big reason to have a ‘fairness’ metric

minimizing turnaround time

recall: first-come, first-served best order:
had shortest CPU bursts first

→ scheduling algorithm: 'shortest job first' (SJF)

= same as priority where CPU burst length determines priority

...but without preemption for now

priority = job length doesn't quite work with preemption
(preview: need priority = remaining time)

a practical problem

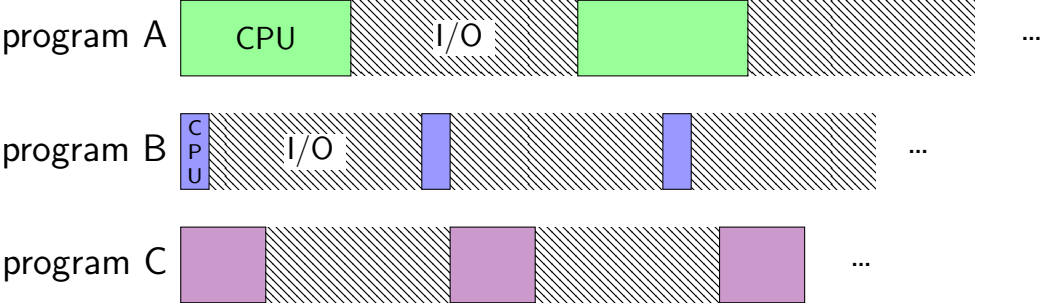
so we want to run the shortest CPU burst first

how do I tell which thread that is?

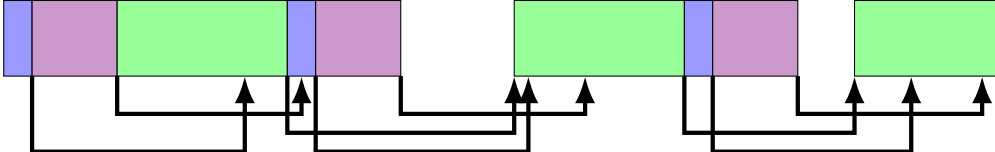
we'll deal with this problem later

...kinda

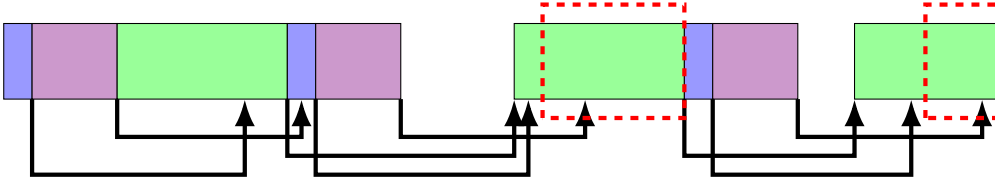
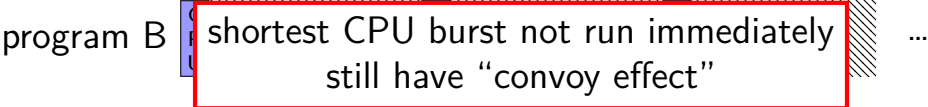
alternating I/O and CPU: SJF



alternating I/O and CPU: SJF



alternating I/O and CPU: SJF



minimizing turnaround time

recall: first-come, first-served best order:
had shortest CPU bursts first

→ scheduling algorithm: 'shortest job first' (SJF)

= same as priority where CPU burst length determines priority

...but without preemption for now

priority = job length doesn't quite work with preemption
(preview: need priority = remaining time)

a practical problem

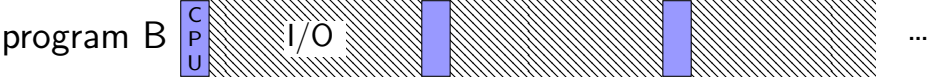
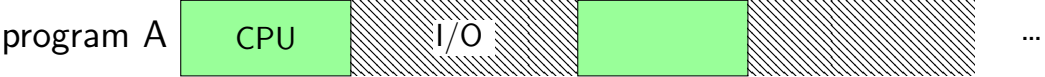
so we want to run the shortest CPU burst first

how do I tell which thread that is?

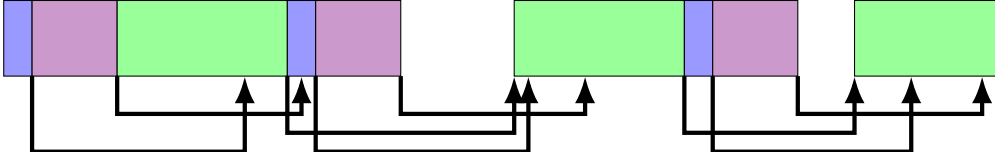
we'll deal with this problem later

...kinda

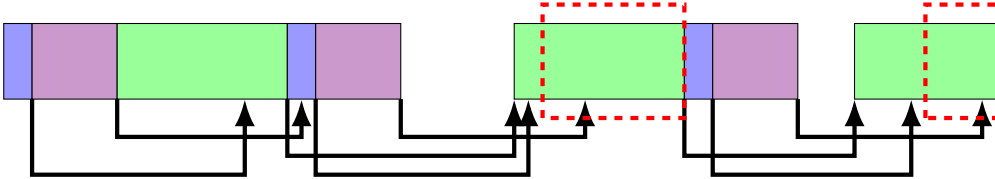
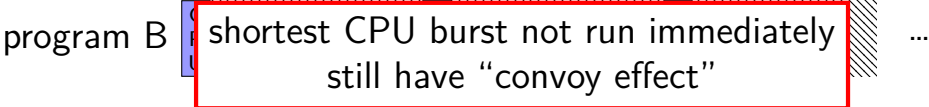
alternating I/O and CPU: SJF



alternating I/O and CPU: SJF



alternating I/O and CPU: SJF



adding preemption (1)

what if a long job is running, then a short job interrupts it?
short job will wait for too long

solution is preemption — reschedule when new job arrives
new job is shorter — run now!

adding preemption (2)

what if a long job is *almost done* running, then a medium job interrupts it?

recall: priority = job length

long job waits for medium job

...for longer than it would take to finish

worse than letting long job finish

adding preemption (2)

what if a long job is *almost done* running, then a medium job interrupts it?

recall: priority = job length

long job waits for medium job

...for longer than it would take to finish

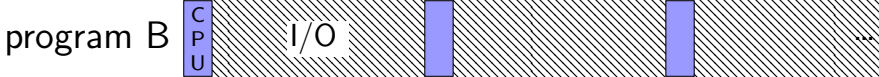
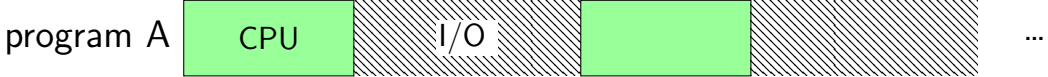
worse than letting long job finish

solution: priority = **remaining time**

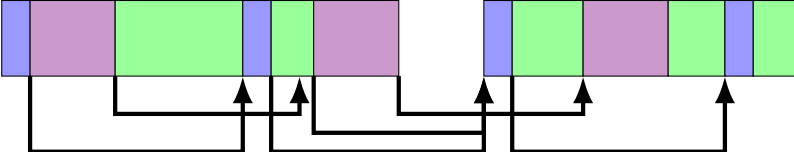
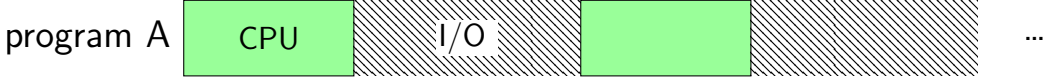
called shortest *remaining time* first (SRTF)

prioritize by what's left, not the total

alternating I/O and CPU: SRTF



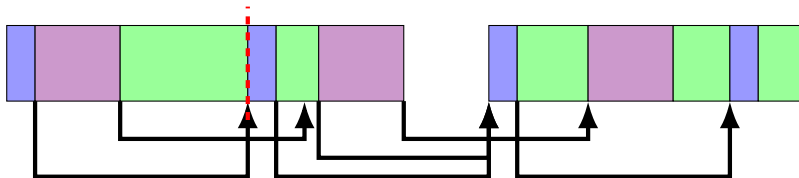
alternating I/O and CPU: SRTF



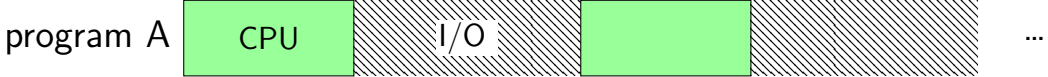
alternating I/O and CPU: SRTF



program **B** preempts **A** because it has less time left
(that is, **B** is shorter than the time **A** has left)

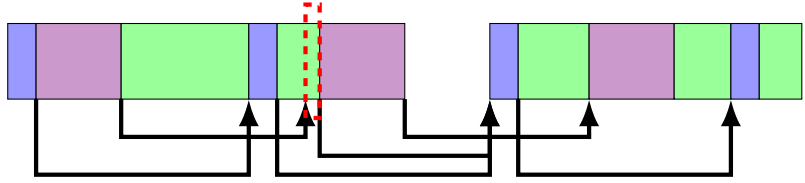


alternating I/O and CPU: SRTF



program B

C does not preempt A
because finishing A is faster than running C



SRTF, SJF are optimal (for response time)

SJF minimizes response time/waiting time

...if you disallow preemption/leaving CPU deliberately idle

SRTF minimizes response time/waiting time

...if you ignore context switch costs

aside on names

we'll use:

SRTF for preemptive algorithm with remaining time

SJF for non-preemptive with total time=remaining time

might see different naming elsewhere/in books, sorry...

knowing job lengths

seems hard

sometimes you can ask

common in batch job scheduling systems

and maybe you'll get accurate answers, even