

CPU Virtualization: Scheduling (continued...)

Questions answered in this lecture:

Linux CFS scheduler

Announcements

- Reading: Today cover Chapters 12-16
- Project 2: User-level thread library and scheduling
 - To be posted soon, Due ~March 10th.
- Midterm - 03/05 – In Class
- Please attend recitations, TA and office hours. Don't wait until deadlines
- When emailing, CC' all TAs and instructor for faster response
 - Also add [CS 416] or [416] in your subject

Review- 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

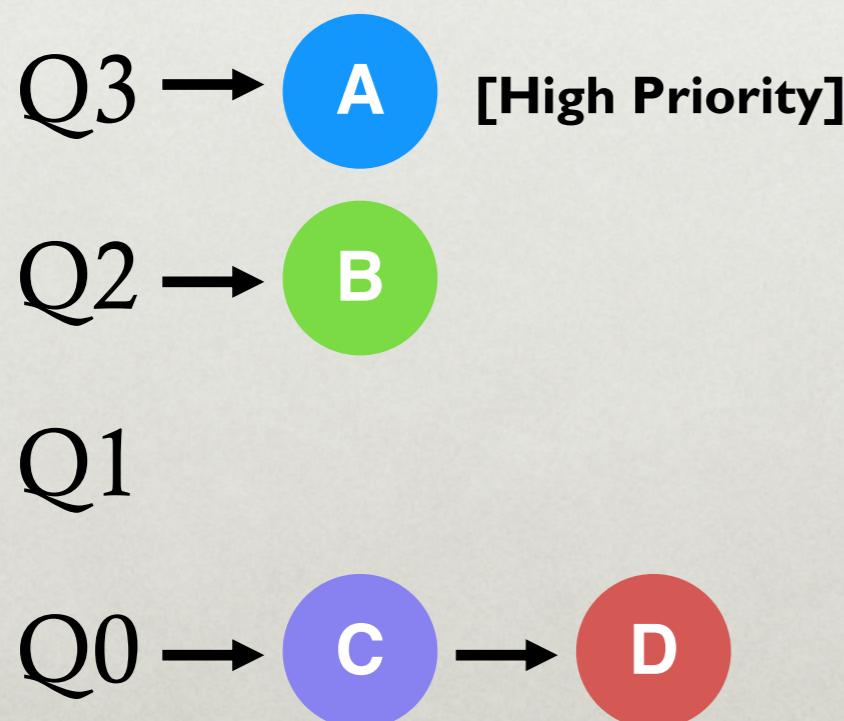
(Multi-Level Feedback Queue)

- 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
- MLFQ has a number of distinct queues.
- Each queue is assigned a different priority level.

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



“Multi-level”

How to know how to set priority?

Approach 1: nice

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

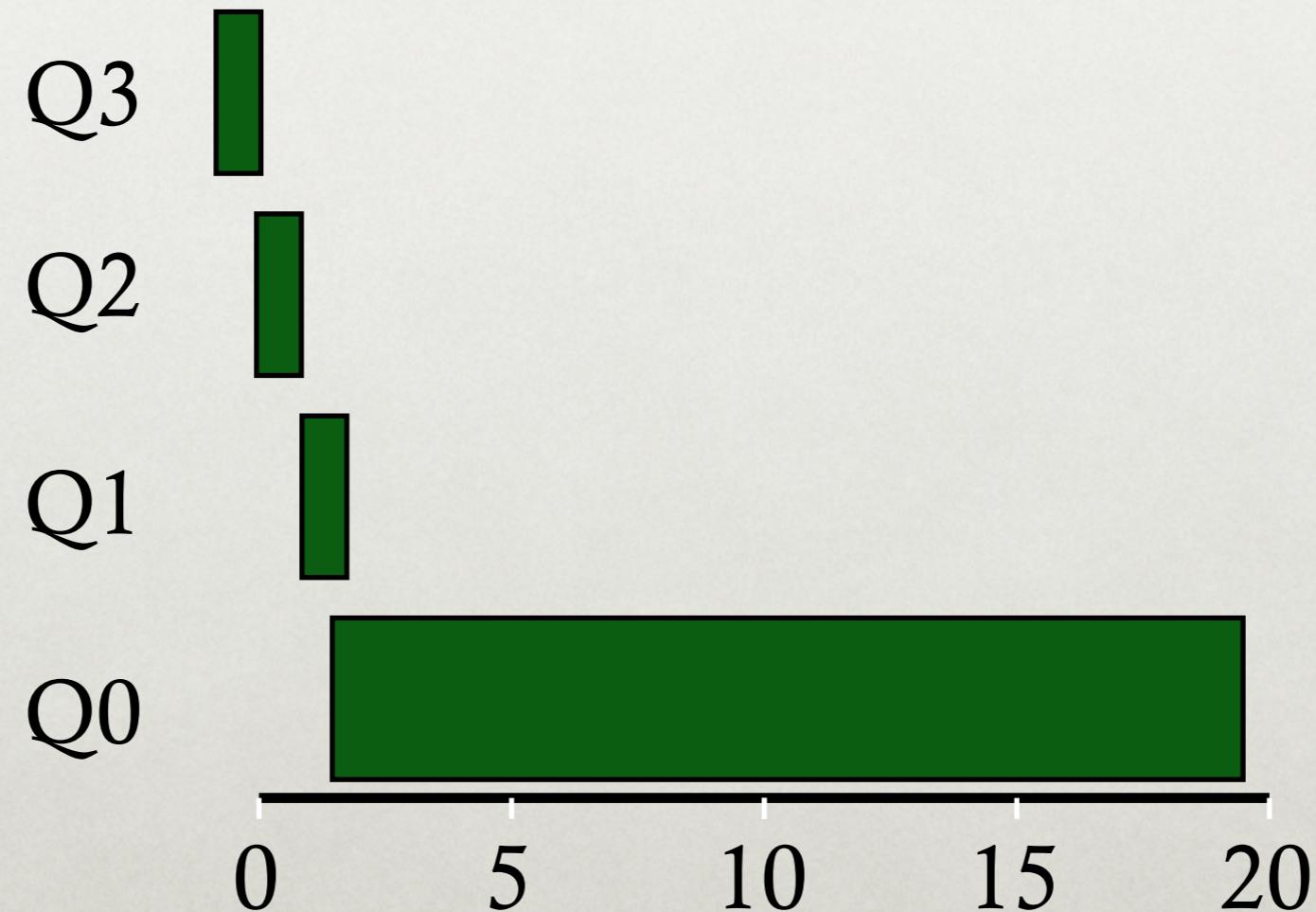
Rule 2: If $\text{priority}(A) == \text{Priority}(B)$, A & B run in RR

More rules:

Rule 3: Processes start at top priority

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

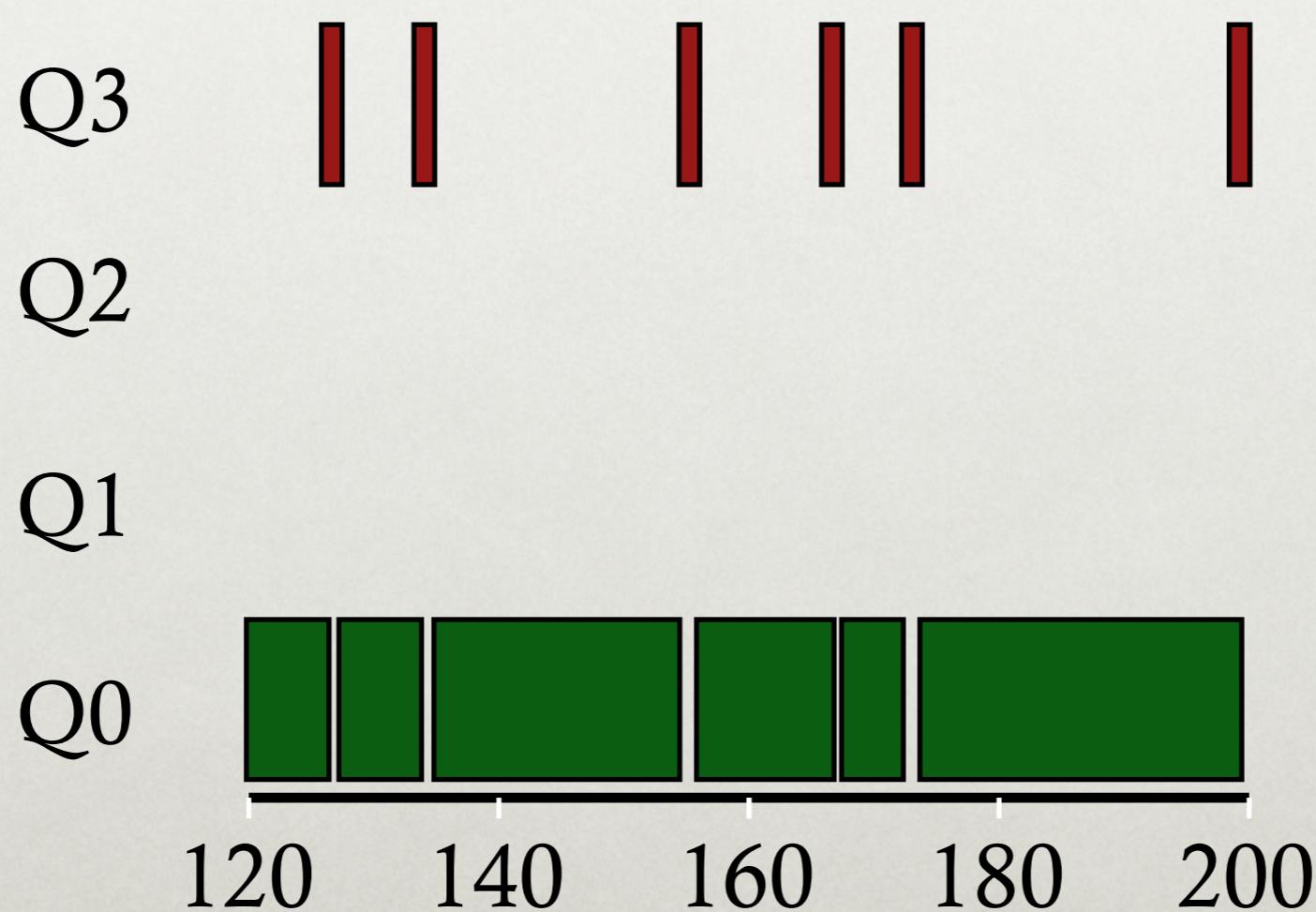
One Long Job (Example)



A four-queue scheduler with time slice 10ms

Long batch job – DNA analysis

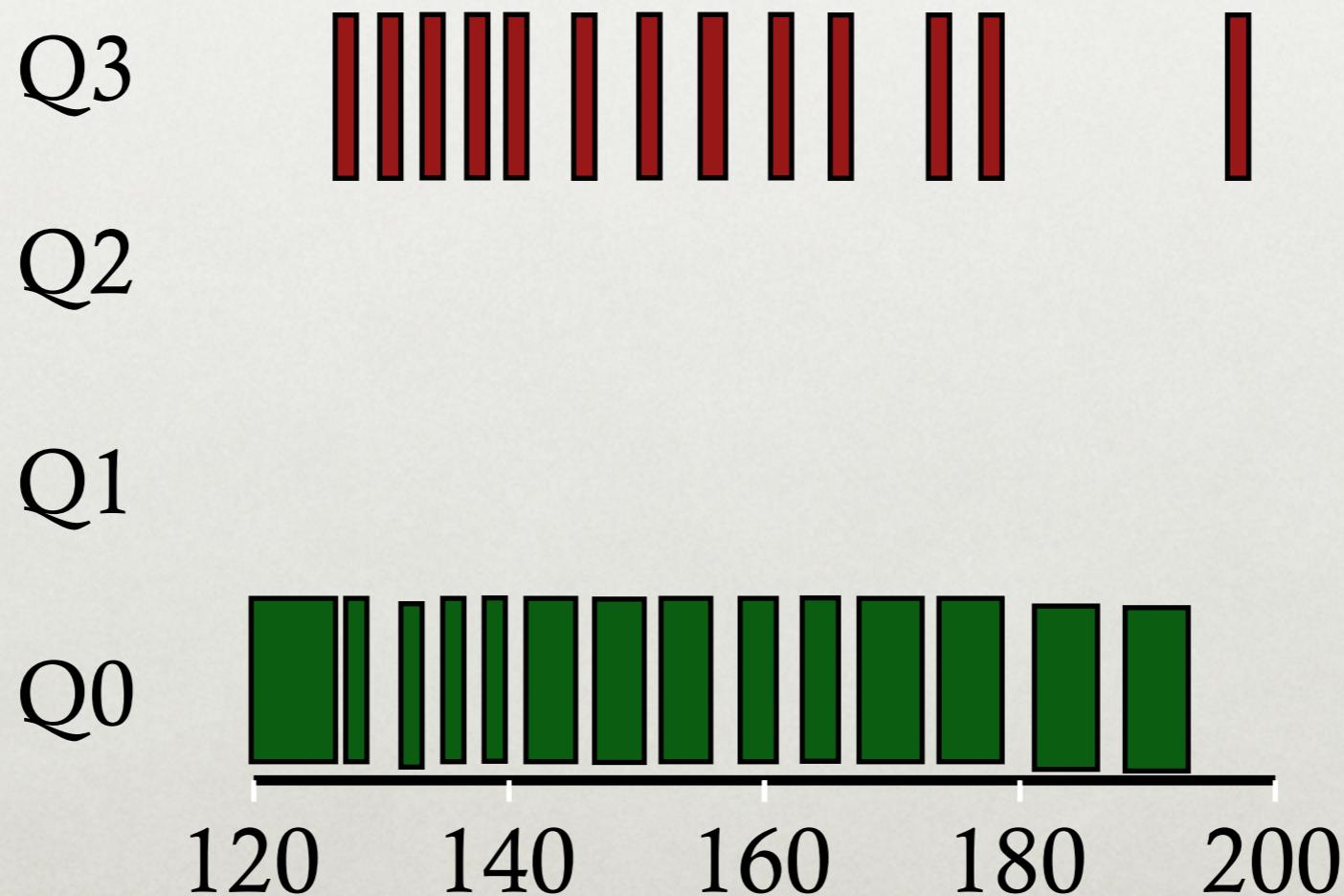
An Interactive Process Joins



Interactive job performs quick operation and does an I/O

Interactive process never uses entire time slice, so never demoted

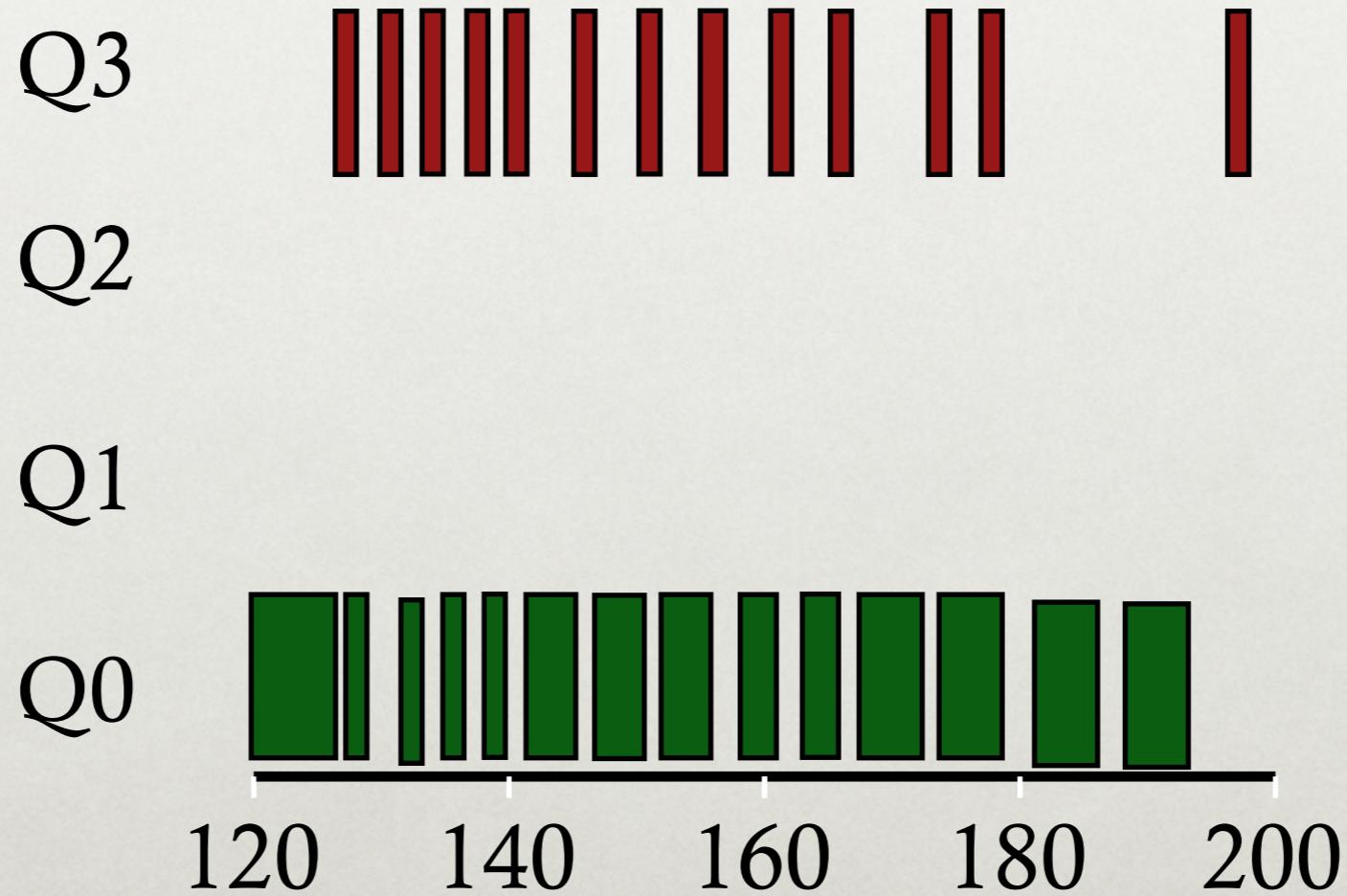
Problems with MLFQ?



Problems

- unforgiving + starvation
- gaming the system

Problems with MLFQ?



Problem: Low priority job may never get scheduled

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

Lottery Scheduling

Goal: proportional (fair) share

Fair-share scheduler

- Guarantee that each job obtain *a certain percentage* of CPU time.
- Not optimized for turnaround or response time

Approach:

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

Amazingly simple to implement

Lottery Scheduling

- Tickets
 - Represent the share of a resource that a process should receive
 - Percent of tickets represents its share of the system resource in question.
- Example
 - There are two processes, A and B.
 - Process A has 75 tickets → receive 75% of the CPU
 - Process B has 25 tickets → receive 25% of the CPU

Lottery Scheduling

- The scheduler picks a winning ticket.
 - Load the state of that *winning process* and runs it.
- Example
 - There are 100 tickets
 - Process A has 75 tickets: 0 ~ 74
 - Process 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
Resulting scheduler: A B A A B A A A A A A B A B A

Intuition:

The longer these two jobs compete,

The more likely they are to achieve the desired percentages.

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

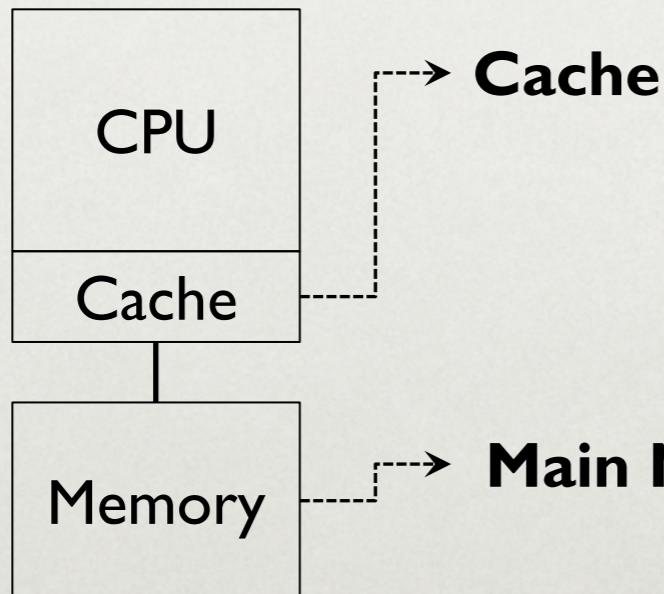
(read more in OSTEP)

Multiprocessor Scheduling

- The rise of the multicore processor is the source of multiprocessor-scheduling proliferation.
- **Multicore:** Multiple CPU cores are packed onto a single chip.
- Adding more CPUs does not make that single application run faster.
→ Rewrite application to run in parallel, using **threads**.

How to schedule jobs on **Multiple CPUs?**

Single CPU with cache



- Small, fast memories
- Hold copies of popular data that is found in the main memory.
- Utilize *temporal* and *spatial* locality

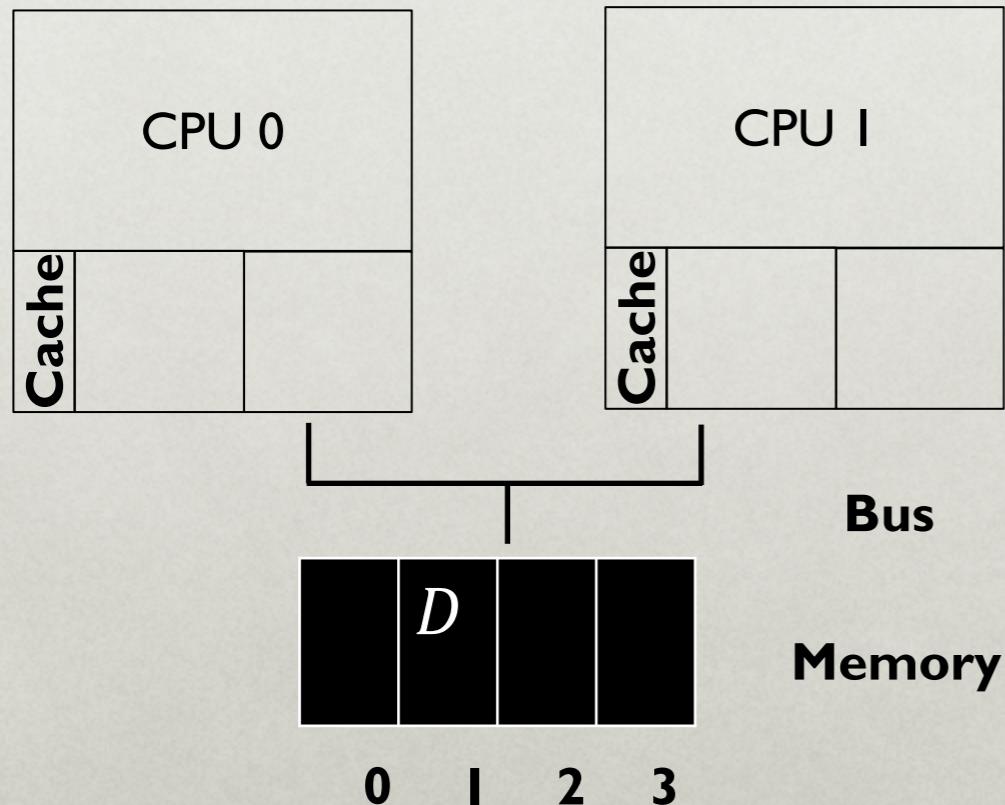
- Holds all of the data
- Access to main memory is slower than cache.

By keeping data in cache, the system can make slow memory appear to be a fast one

Cache Coherence

- Consistency of shared resource data stored in multiple caches.

0. Two CPUs with caches sharing memory

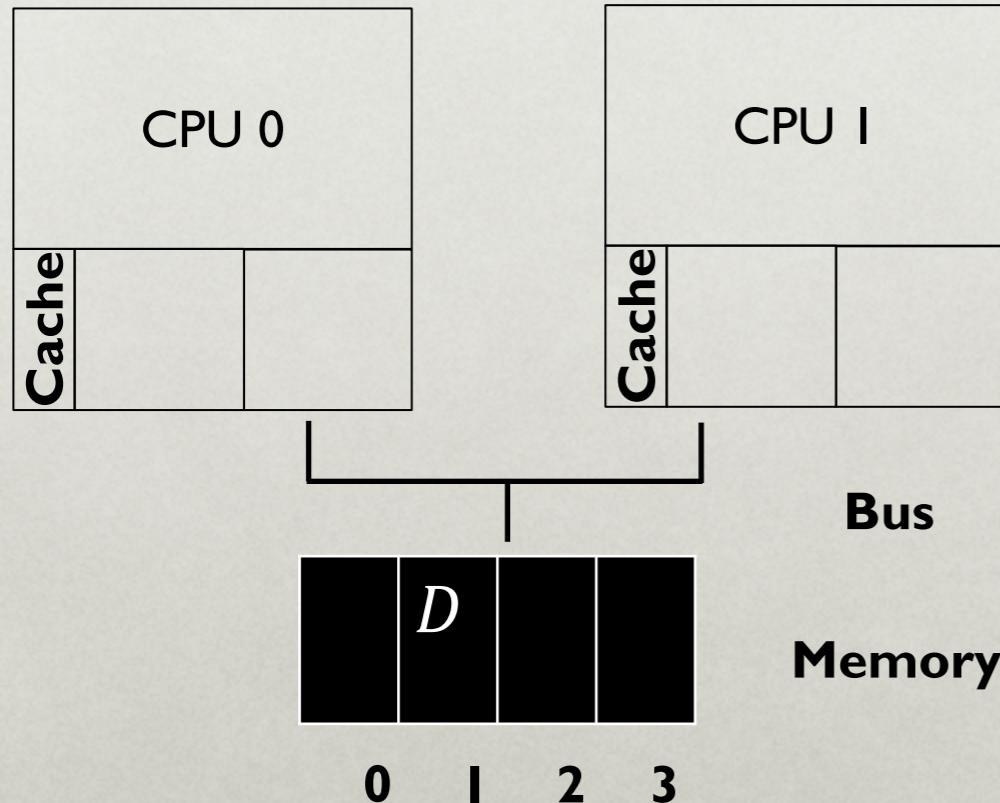


Initial D value is 0

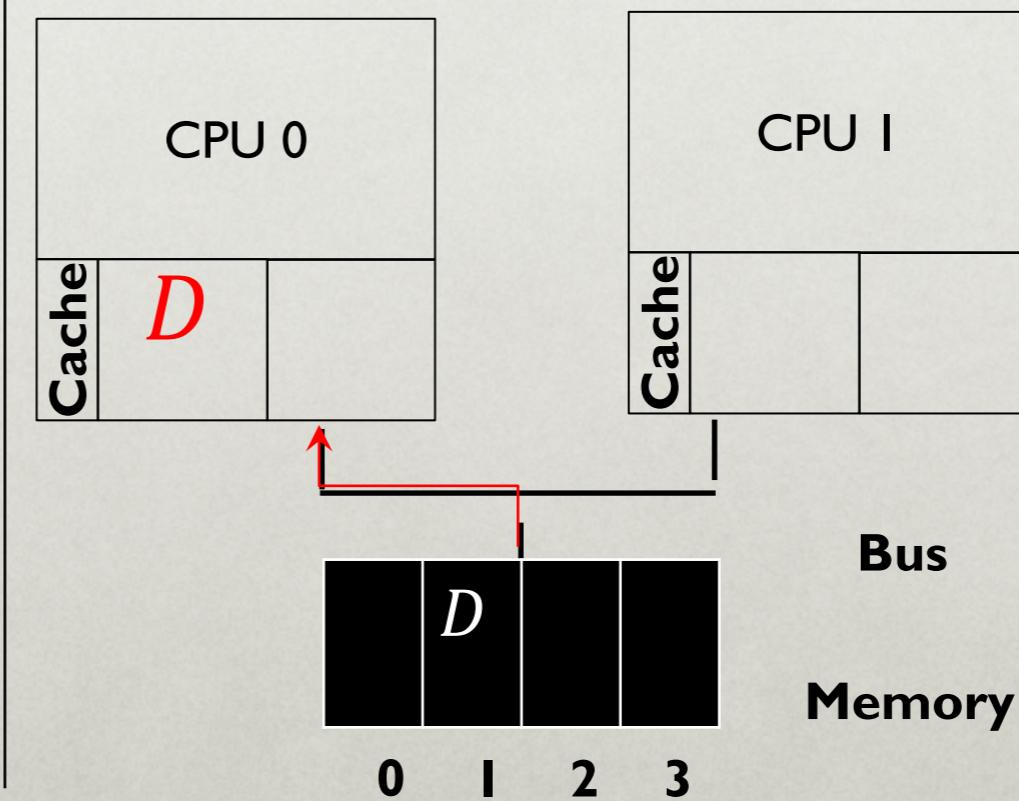
Cache Coherence

- Consistency of shared resource data stored in multiple caches.

0. Two CPUs with caches sharing memory



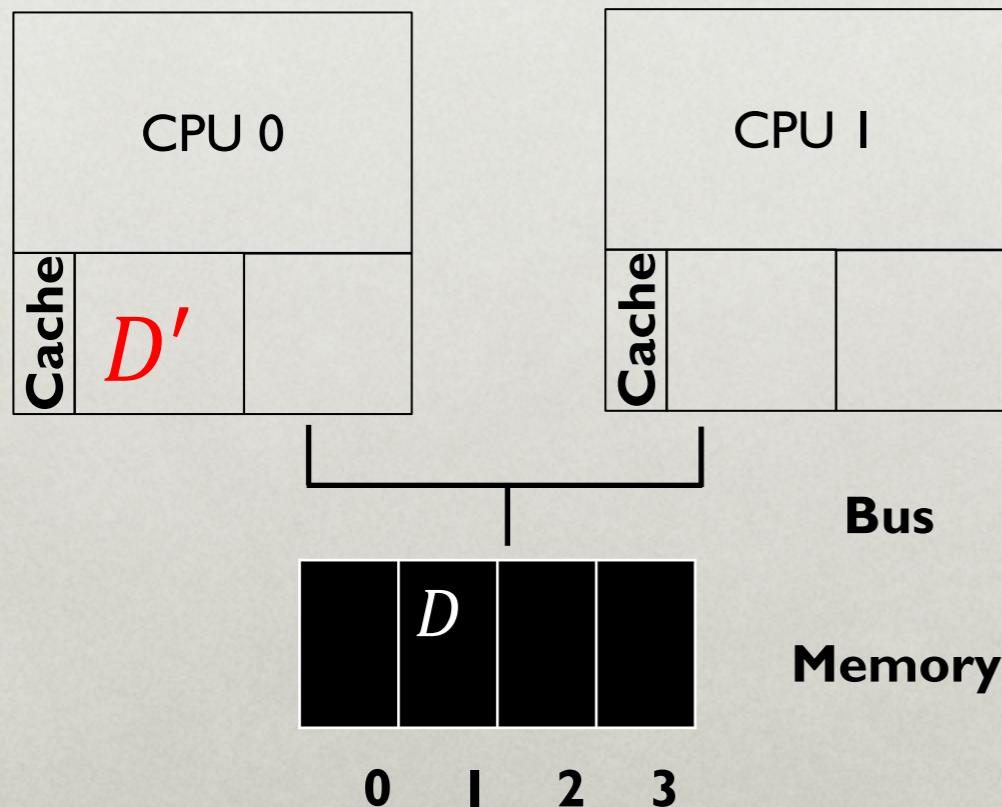
I. CPU0 reads a data at address 1.



Cache Coherence

- Consistency of shared resource data stored in multiple caches.

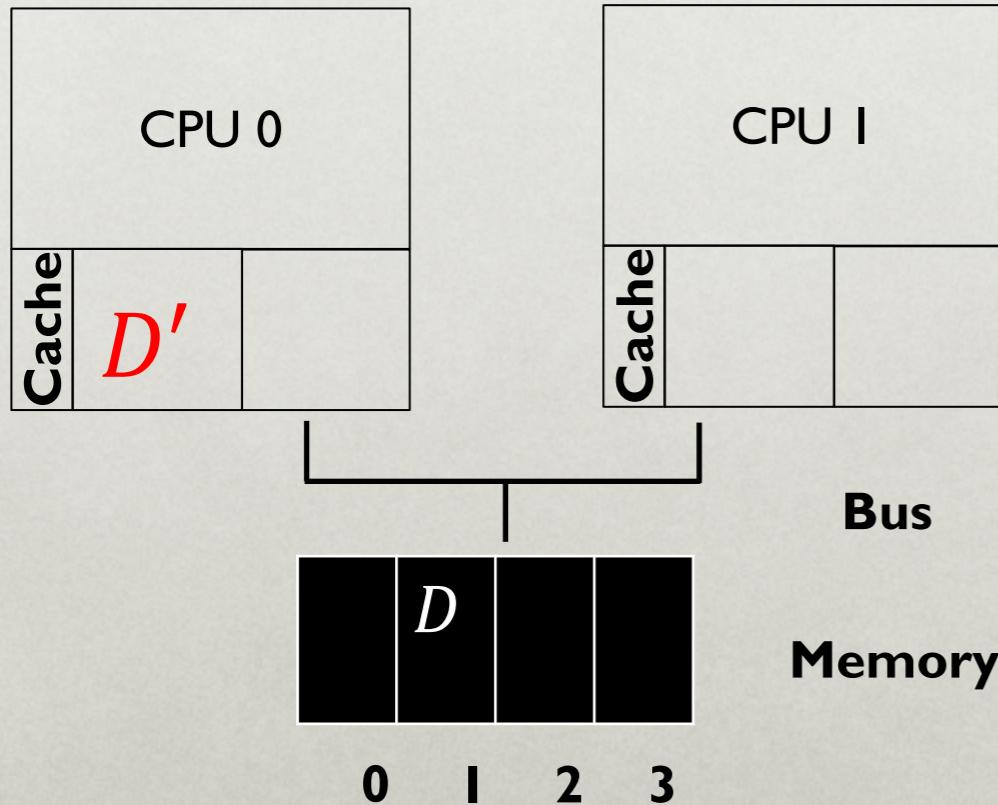
2. D is updated and CPU I is scheduled.



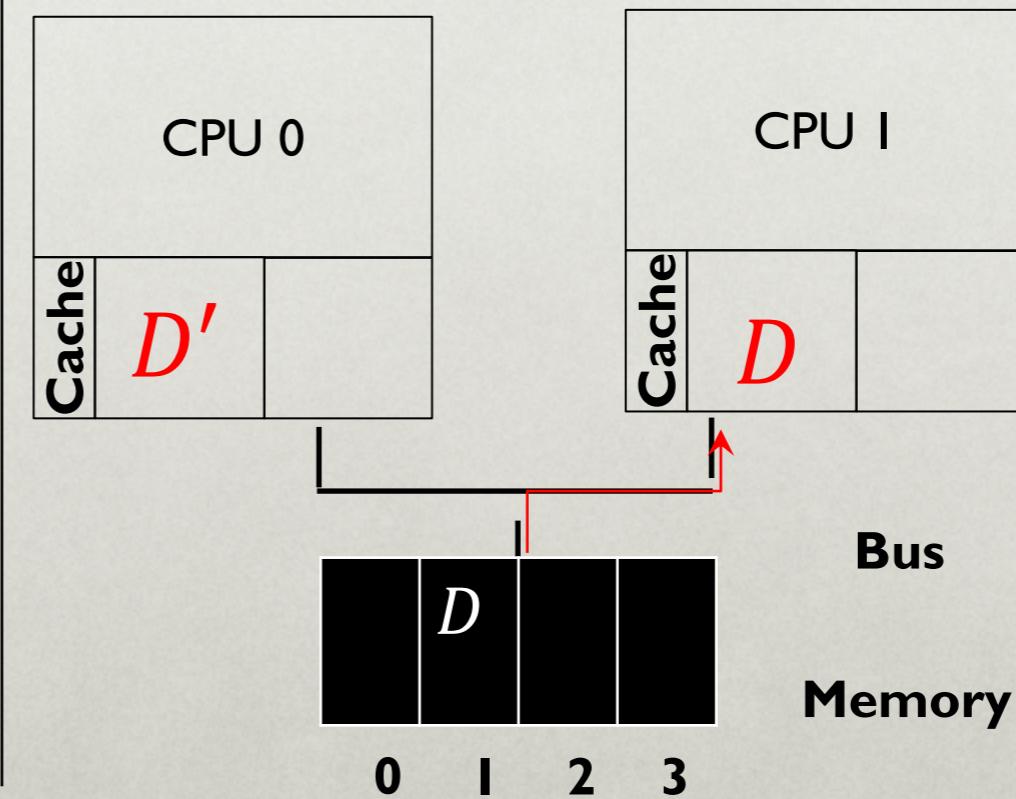
Cache Coherence

- Consistency of shared resource data stored in multiple caches.

2. D is updated and CPU I is scheduled.



3. CPU I re-reads the value at address A



Data Inconsistency Problem!

Cache Coherence Solution

- Bus snooping
 - Each cache pays attention to memory updates by **observing the bus**.
 - When a CPU sees an update for a data item it holds in its cache, it will notice the change and either invalidate its copy or update it.
- When accessing shared data across CPUs, mutual exclusion primitives should likely be used to guarantee correctness

Don't forget synchronization

```
1  typedef struct __Node_t {
2      int value;
3      struct __Node_t *next;
4  } Node_t;
5
6  int List_Pop() {
7      Node_t *tmp = head;    // remember old head ...
8      int value = head->value; // ... and its value
9      head = head->next;    // advance head to next pointer
10     free(tmp);           // free old head
11     return value;         // return value at head
12 }
```

Don't forget synchronization

```
1  typedef struct __Node_t {
2      int value;
3      struct __Node_t *next;
4  } Node_t;
5
6  int List_Pop() {
7      lock (&m)
8      Node_t *tmp = head;    // remember old head ...
9      int value = head->value; // ... and its value
10     head = head->next;    // advance head to next pointer
11     free(tmp);            // free old head
12     unlock (&m)
13     return value;          // return value at head
14 }
```

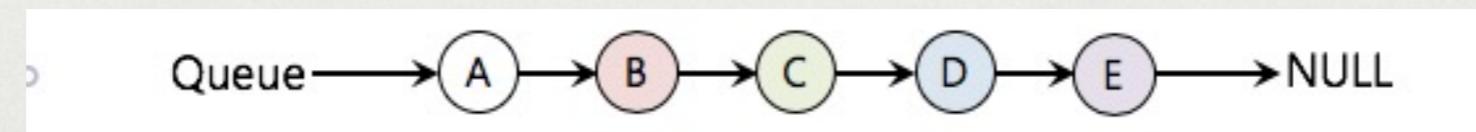
Cache Affinity

- Keep a process on the same CPU if at all possible
 - A process builds up a fair bit of state in the cache of a CPU.
 - The next time the process runs, it will run faster if some of its state is *already present* in the cache on that CPU.

A multiprocessor scheduler should consider **cache affinity** when making its scheduling decision.

Cache Affinity

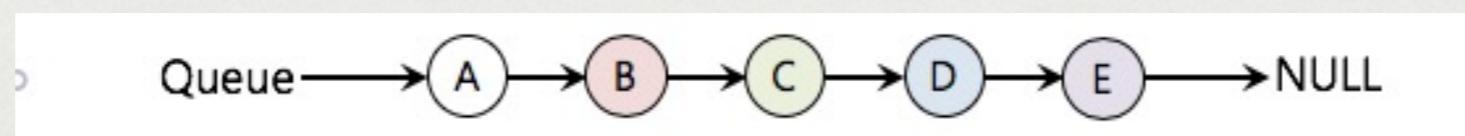
- Put all jobs that need to be scheduled into a single queue.



- Each CPU simply picks the next job from the globally shared queue.
- Cons:
 - Some form of **locking** have to be inserted → **Lack of scalability**
 - **Cache affinity**

Cache Affinity

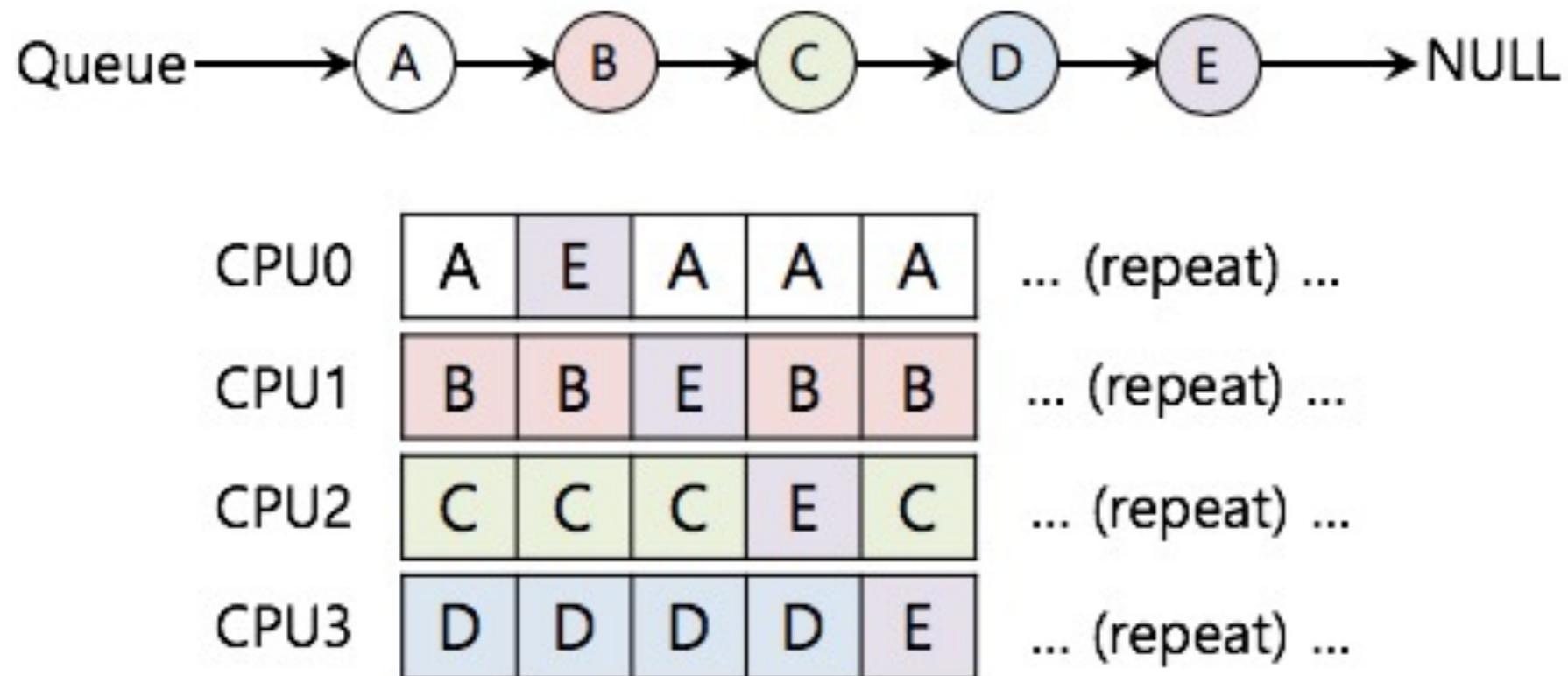
- Put all jobs that need to be scheduled into a single queue.



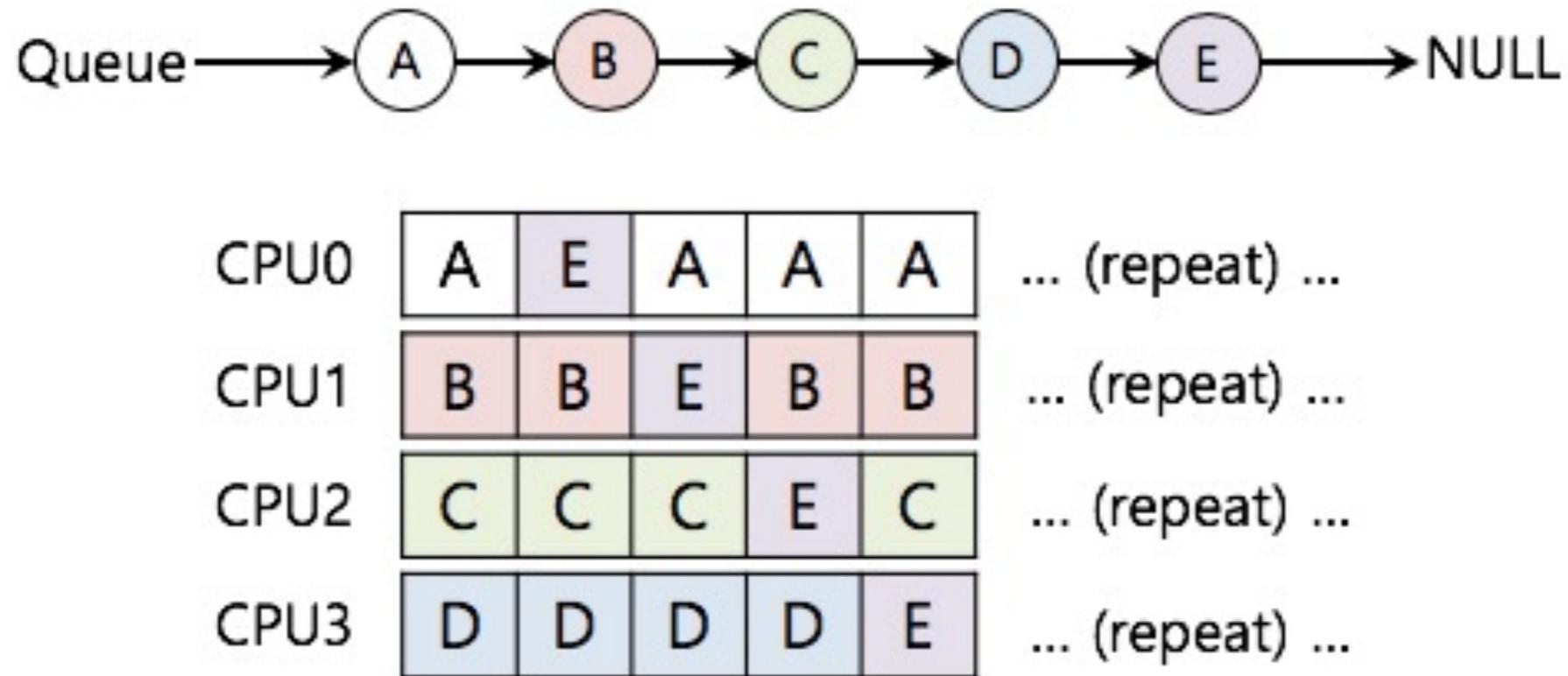
- Each CPU simply picks the next job from the globally shared queue.
- Cons:
 - Some form of **locking** have to be inserted → **Lack of scalability**
 - Cache affinity?**

CPU0	A	E	D	C	B	... (repeat) ...
CPU1	B	A	E	D	C	... (repeat) ...
CPU2	C	B	A	E	D	... (repeat) ...
CPU3	D	C	B	A	E	... (repeat) ...

Scheduling Cache Affinity



Scheduling Cache Affinity



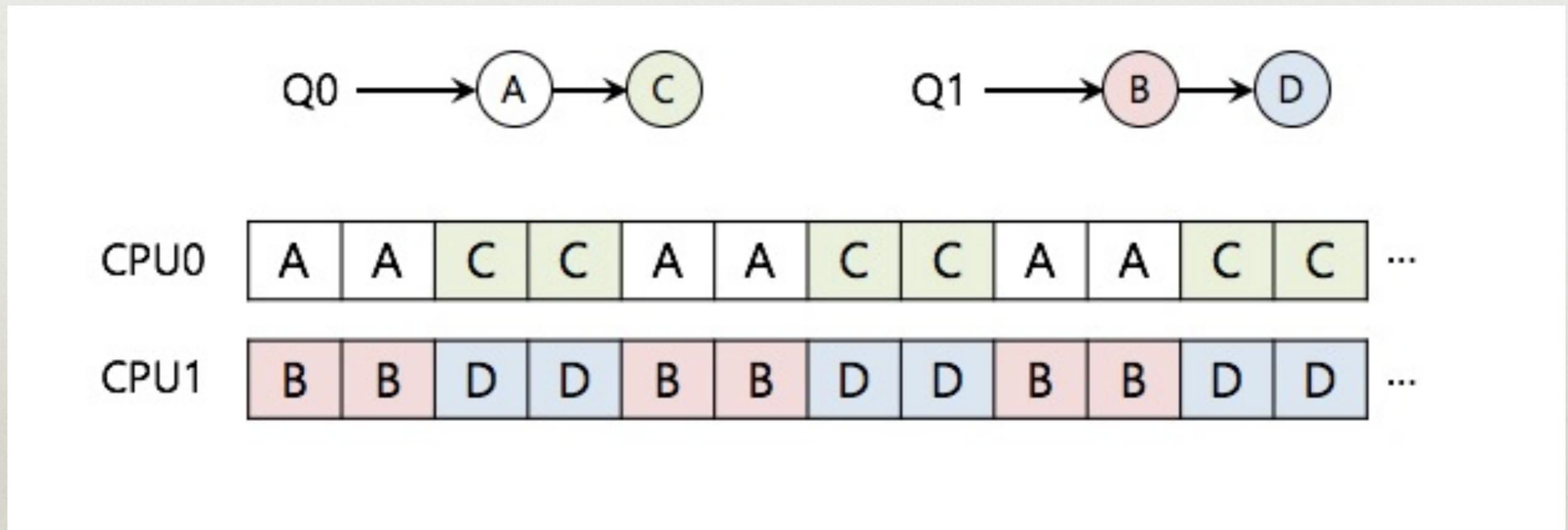
- Preserving affinity for most
 - Jobs A through D are not moved across processors.
 - Only job **E** Migrating from CPU to CPU.
- Implementing such a scheme can be **complex**.

Multi-queue Multiprocessor Scheduling (MQMS)

- MQMS consists of multiple scheduling queues.
 - Each queue will follow a particular scheduling discipline.
 - When a job enters the system, it is placed on **exactly one** scheduling queue.
 - Avoid the problems of information sharing and synchronization.

Multi-queue Multiprocessor Scheduling (MQMS)

- With **round robin**, the system might produce a schedule that looks like this:



MQMS provides more scalability and cache affinity.

Completely Fair Scheduling (CFS)

CFS changes or removes time-slice allotment

- Do away with timeslices completely
- Assign each process a proportion of the processor
- Yields constant fairness but a variable switching rate

CFS is based on a simple concept

Model scheduling as if the system had an ideal, perfectly multitasking processor

- Each process receives $1/n$ of the processor's time, where n is the number of runnable processes
- Would be scheduled for infinitely small durations
- In any measurable period, all n processes would run for the same amount of time

Completely Fair Scheduling (CFS)

CFS changes or removes time-slice allotment

- Do away with timeslices completely
- Assign each process a proportion of the processor
- Yields constant fairness but a variable switching rate

Assume there are two processes

- Unix model - run one process for 5ms and another for 5ms
- Each process receives 100% of the processor
- IDEAL System: Perfect multitasking processor
 - Run both processes simultaneously for 10 milliseconds
 - Each at 50% power - called perfect multitasking

Ideal Multiprocessor possible?

Why not?

Completely Fair Scheduling (CFS)

Rank processes based on their worth and need for processor time

Processes with a higher priority run before those with a lower priority

Linux has two priority ranges

- Nice value: ranges from -20 to +19 (default is 0)
- High values of nice means lower priority
- Real-time priority: ranges from 0 to 99
- Higher values mean higher priority
- Real-time processes always executes before standard (nice) processes

```
ps ax -eo pid,ni,rt,prio,cmd
```

Linux CFS

Rank processes based on their worth and need for processor time

Processes with a higher priority run before those with a lower priority

Linux has two priority ranges

- Nice value: ranges from -20 to +19 (default is 0)
- High values of nice means lower priority
- Real-time priority: ranges from 0 to 99
- Higher values mean higher priority
- Real-time processes always executes before standard (nice) processes

```
ps ax -eo pid,ni,rt,prio,cmd
```

Linux CFS

Assume the targeted latency is 20 milliseconds

If there are two runnable tasks at the same priority

Each will run for 10 milliseconds before preempting in favor of the other

If we have four tasks at the same priority

Each will run for 5 milliseconds

If there are 20 tasks, each will run for 1 millisecond

As the number of runnable tasks approaches infinity, the proportion of allotted processor and the assigned timeslice approaches zero

Problem? (Solution: Run any task for a minimum of 1 millisecond)

Linux CFS

target latency (TL)

- minimum amount of time—idealized to an infinitely small duration—required for every runnable task to get at least one turn on the processor
- Window where every process gets some CPU

minimum granularity (MG)

- imposes a floor on the timeslice assigned to each process
- For example: run for 1 ms to ensure there is a ceiling on the incurred switching costs
- not perfectly fair when the number of processes grows so large

Linux CFS

No default timeslice, CFS calculates how long a process should run according to the number of runnable processes

- Dynamic timeslice is weighted by the process priority (nice)
$$\text{Weight} = 1024 / (1.25 ^ \text{nice_value})$$
- To calculate the actual timeslice, CFS sets a targeted latency
- Targeted latency: period during which all runnable processes should be scheduled at least once
- Minimum granularity: floor at 1 ms (default)

For queuing tasks, CFS maintains a Red-Black tree instead of a linked list. The tasks are sorted based on their current runtime

Linux CFS

Given for instance a target latency of 20 milliseconds

Example 1: two runnable processes of equal niceness, then both processes will run for 10 milliseconds each before being pre-empted in favour of the other process.

Example 2: If there are 10 processes of equal niceness, each runs for 2 milliseconds each.

Example 3: one with a niceness of 0 (the default), the other with a niceness of 5?

proportional difference - 1/3

higher priority process – 15ms

lowest priority process – 5ms

Summary

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

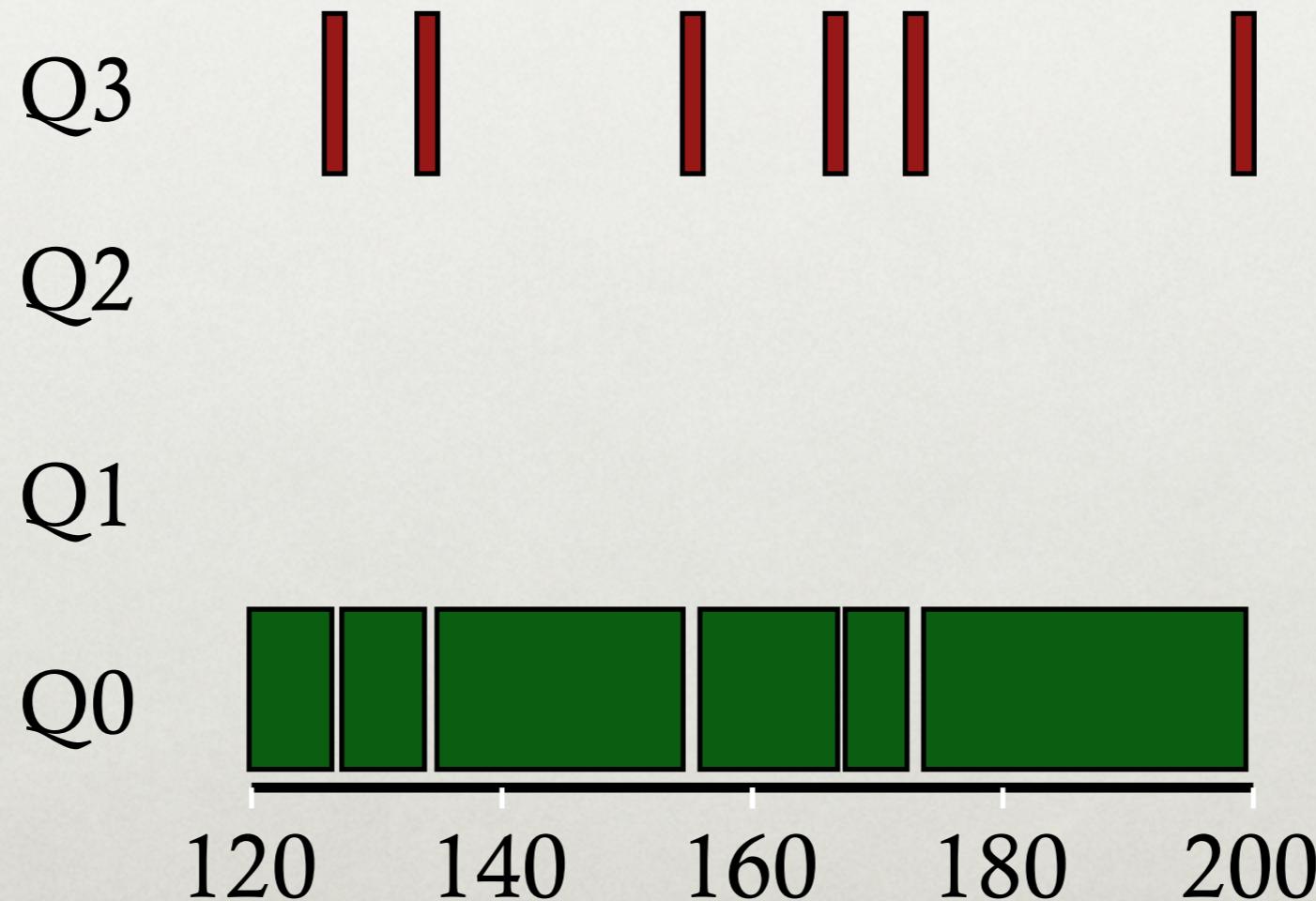
General purpose schedulers need to support processes with different goals

Past behavior is good predictor of future behavior

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

Important to consider caches in multiprocessor scheduling

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