

Running Programs on a System

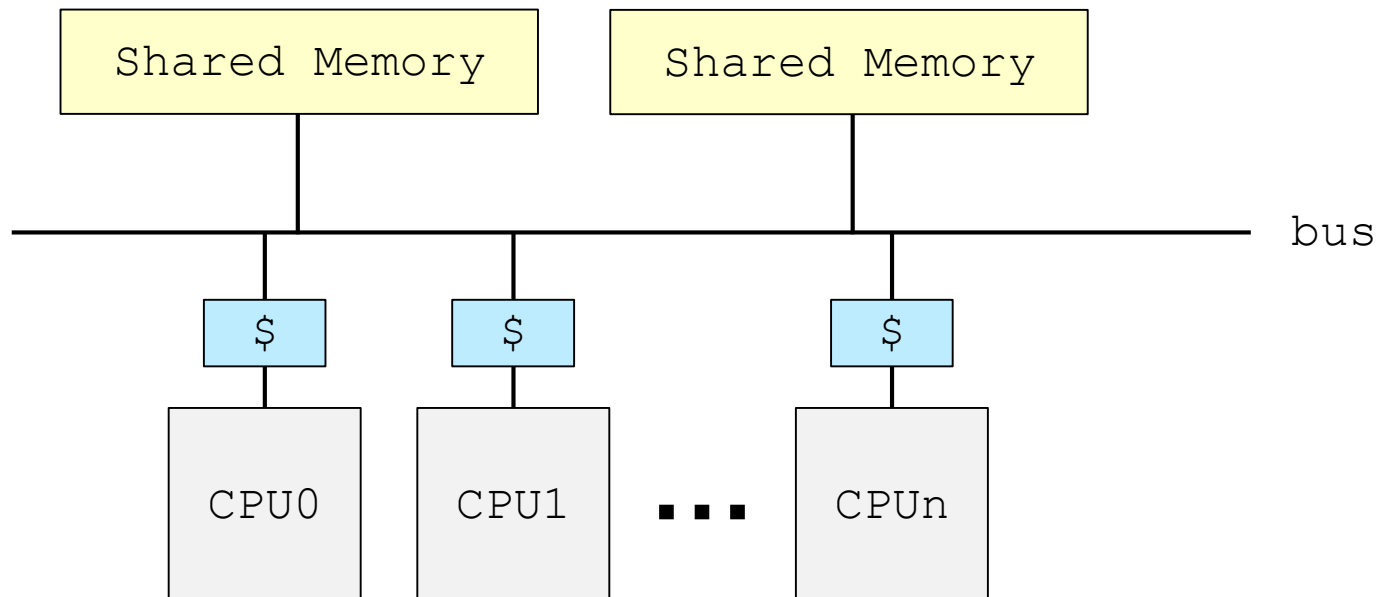
Advanced Process Scheduling



SMP Scheduling

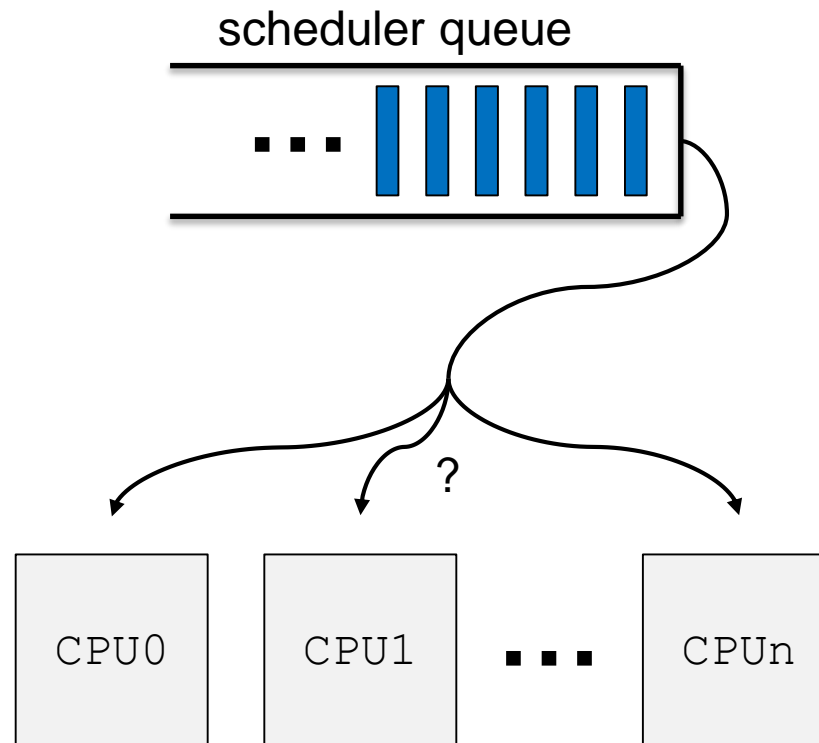
Symmetric Multiprocessor (SMP)

- Multiprocessor with shared memory
 - limited number of homogeneous CPUs with private cache
 - equal/similar access time to shared memory



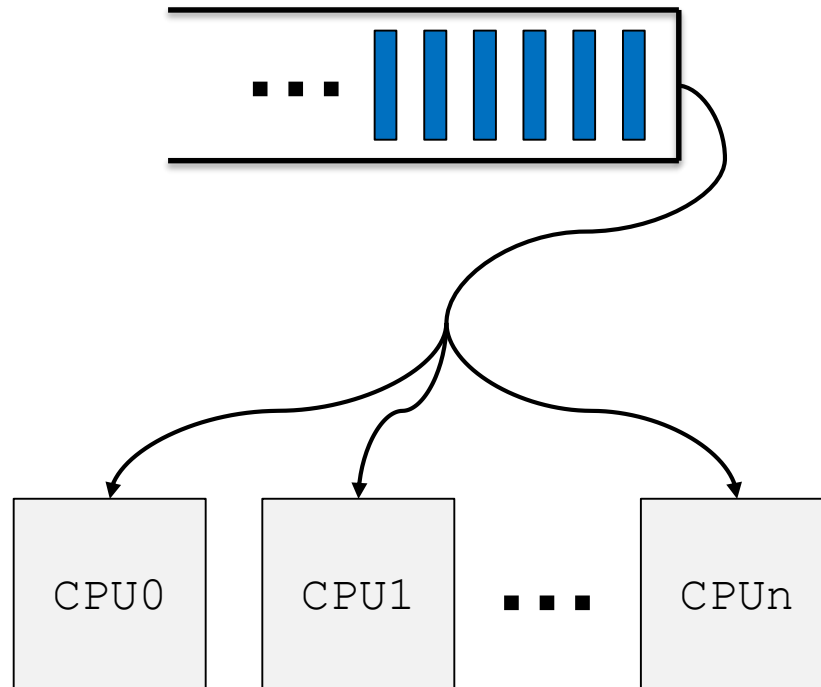
Multiprocessor Scheduling

- How should tasks get distributed?



Global Queue

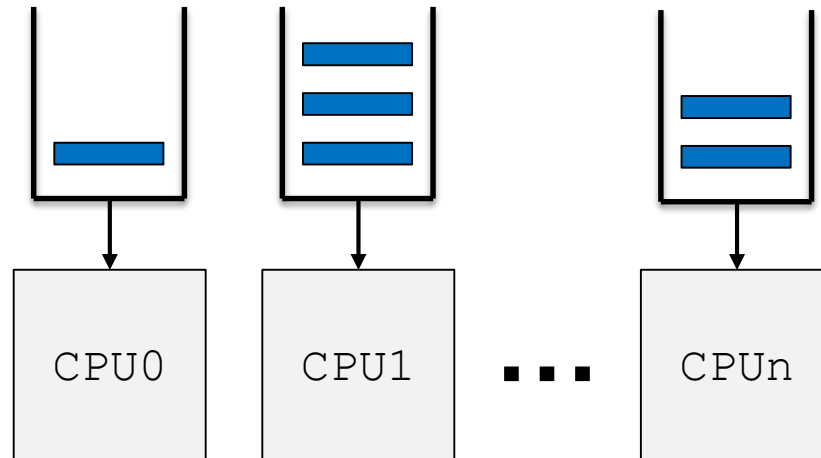
- each core runs his own scheduler
- global queue shared data structure



- + fairness
- + CPU utilization
- scalability
- cache locality

Local (per-CPU) Queue

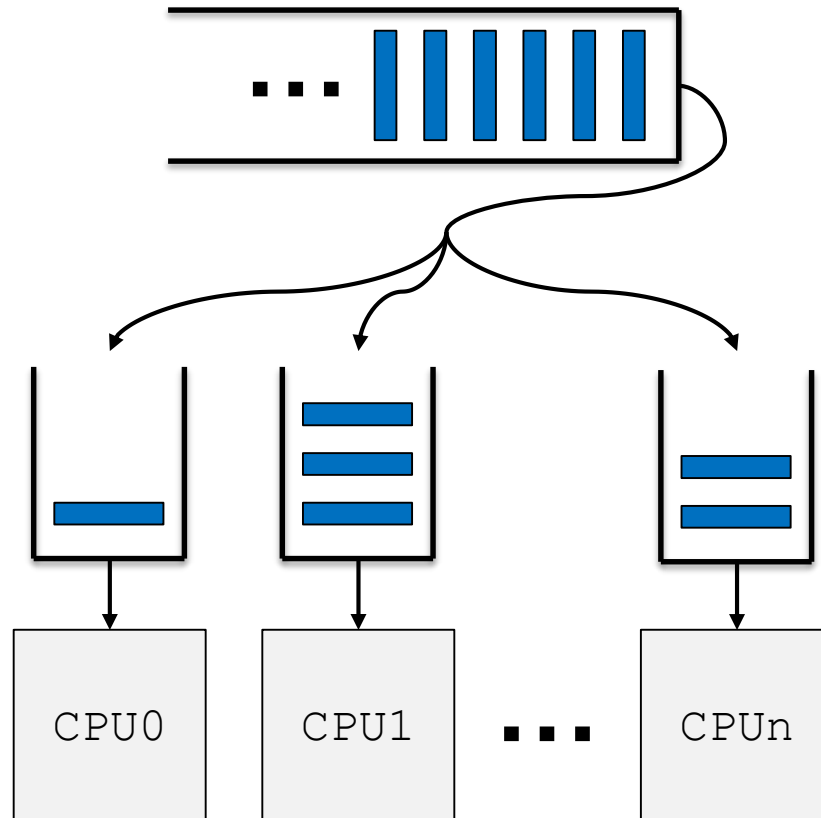
- each core runs his own scheduler
- local queue private data structure
- static assignment of processes to CPUs



- + scalability
- + cache locality
- + simplicity
- load balancing

Local (per-CPU) Queue with Load Balancing

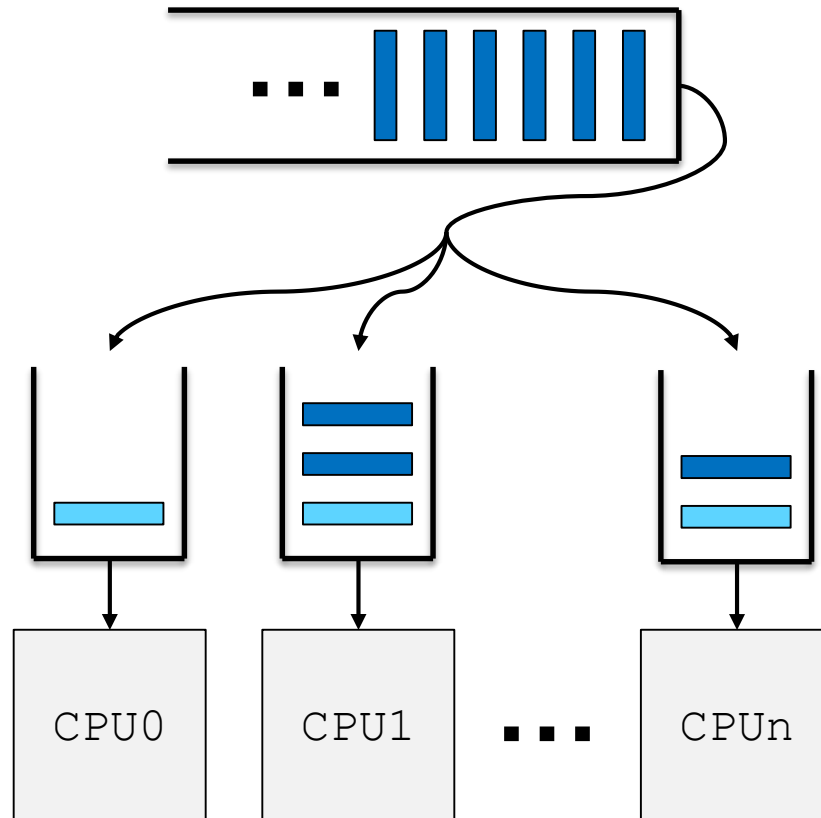
- hybrid approach between global & local queues
- schedule locally, balance globally
- processor affinity



- + best of both worlds
- complexity

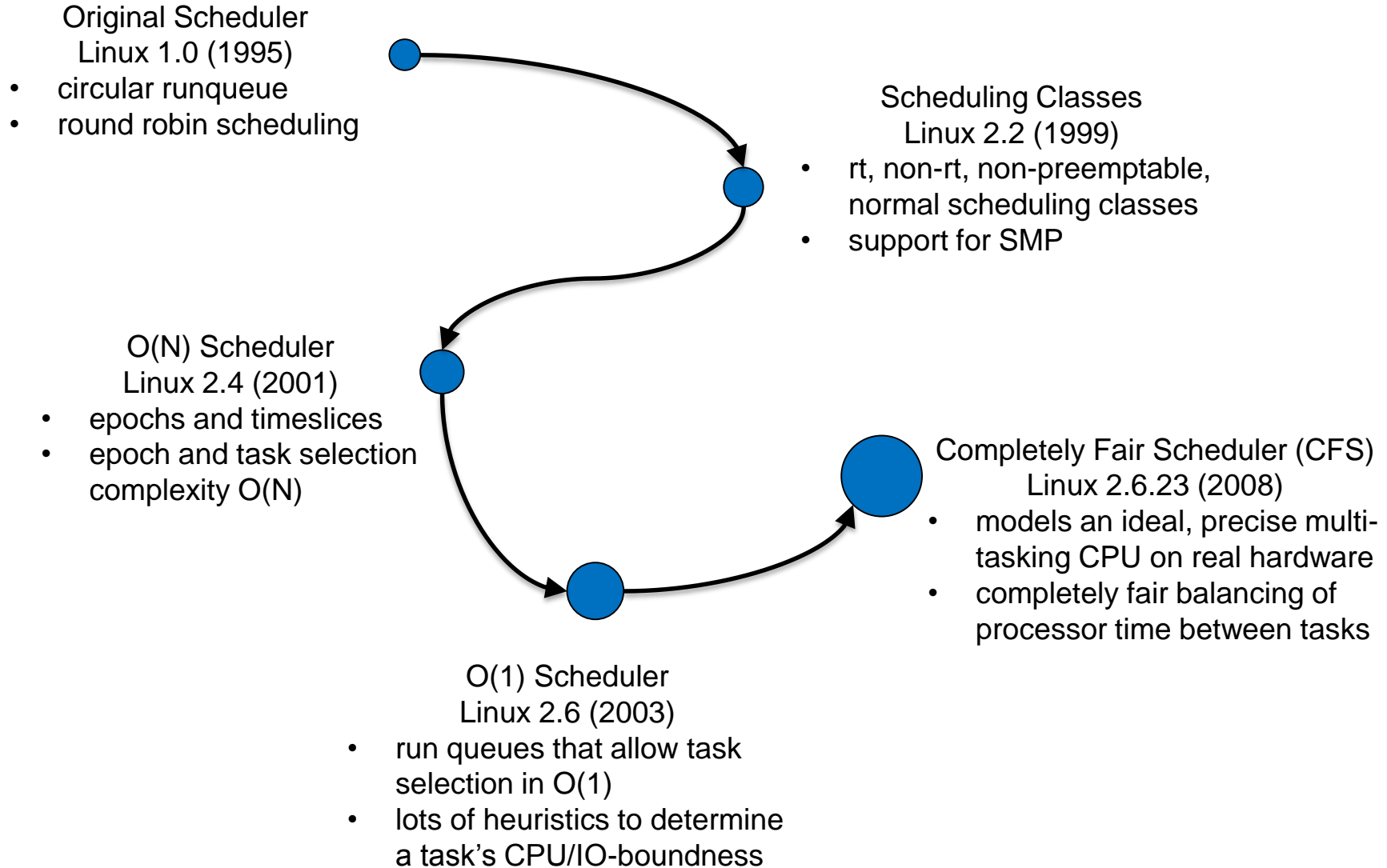
Gang Scheduling

- parallel tasks may need coordination and should be scheduled simultaneously
- local scheduler queue does not act independently
- global context switch



History of the Linux Scheduler

History of the Linux Scheduler



The $O(N)$ Scheduler

Linux 2.4 (2001)

- **timeslice**: assigned to each task at the beginning of an epoch
 - depends on priority and unused timeslice in last epoch
 - `fork()`: parent/child share remaining timeslice (each gets half)
- **epoch**: scheduling “generation”
 - initialization: sum of all timeslices
 - end of epoch: no ready-to-run tasks has timeslice > 0
- **goodness()**: selection of next task to run

```
max_goodness = 0
max_t = idle_task

foreach runnable task t
    if ((g = goodness(t)) > max_goodness)
        max_goodness = g
        max_t = t

dispatch(max_t)
```

The $O(N)$ Scheduler

Linux 2.4 (2001)

■ issues

- does not scale well
 - ▶ epoch calculation: $O(\# \text{ tasks})$
 - ▶ task selection: $O(\# \text{ runnable tasks})$
- favors I/O bound tasks
 - ▶ maybe okay for server, less so for desktops
- predefined minimal quantum is too long
 - ▶ long latency under high system load
- weak support for real-time processes

The O(1) Scheduler

Linux 2.6 (2003)

■ two main ideas

- *active* and *expired* runqueues
 - ▶ active runqueue: tasks with timeslice > 0
 - ▶ expired: tasks that have expired their timeslice
 - ▶ active \rightarrow expired: task moved when timeslice reaches 0
 - new timeslice is computed upon insertion into expired
 - ▶ active runqueue empty \rightarrow array swap
- pair of (active,expired) runqueues *per priority level*
 - ▶ scheduling: pick first task from first non-empty active runqueue of highest priority level $\rightarrow O(1)$

The O(1) Scheduler

Linux 2.6 (2003)

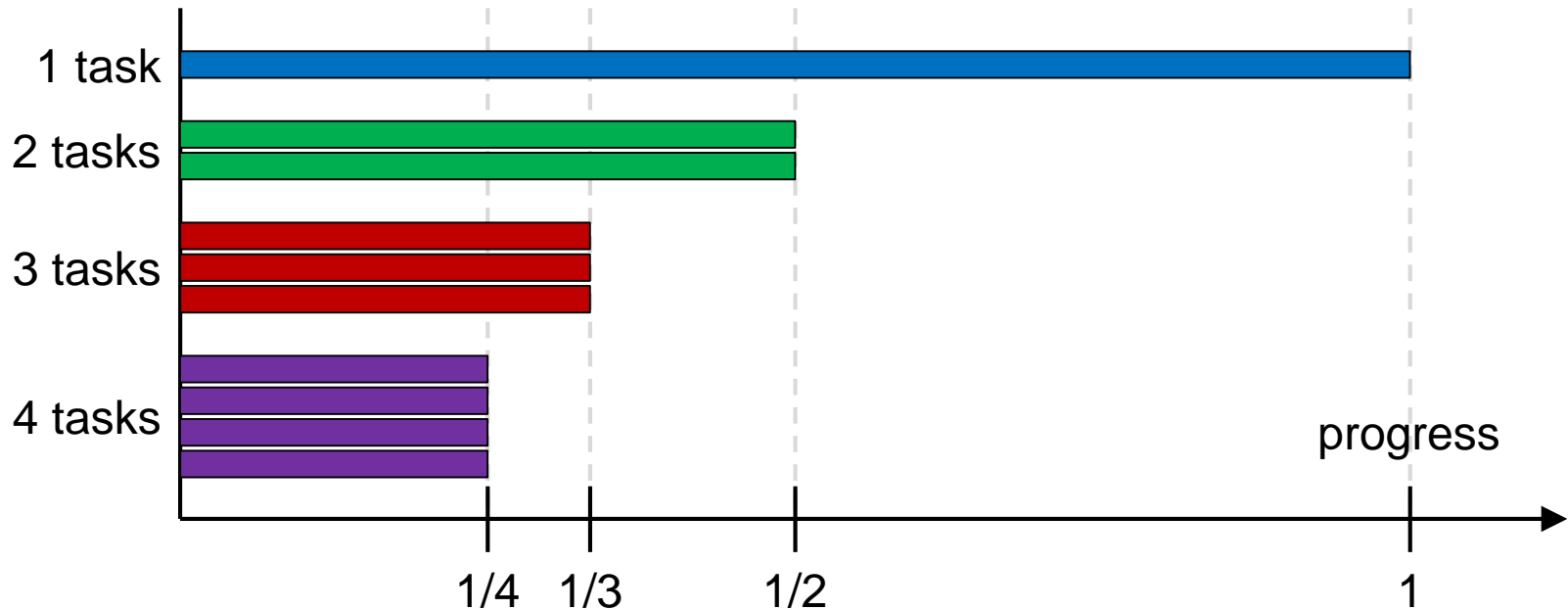
■ issues

- heuristics to determine CPU- resp. I/O-boundness of task
 - ▶ complex and error-prone
 - ▶ permitted attacks on the scheduler

The Completely Fair Scheduler

Linux 2.6.23 (2008)

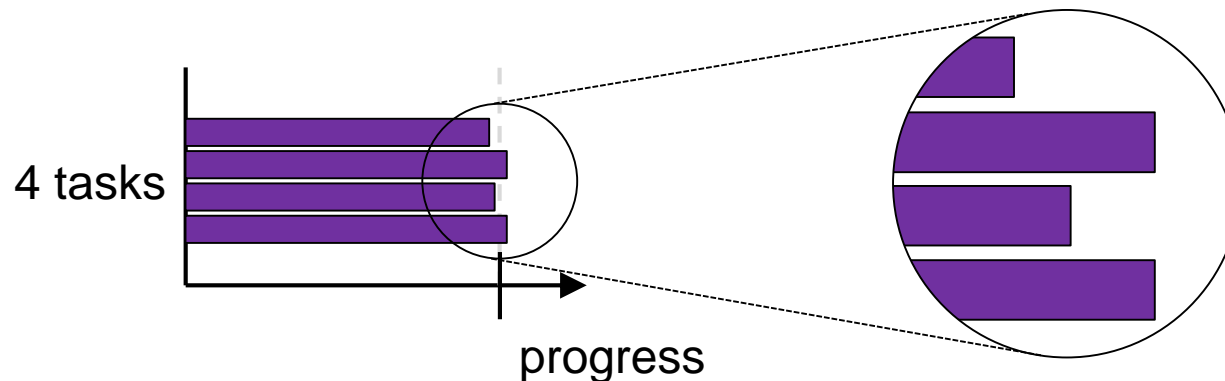
- Models an ideal, precise multitasking CPU
 - simultaneous process of tasks
infinitesimally small timeslices, no task switching overhead
 - n runnable tasks progress uniformly at $1/n^{\text{th}}$ of the CPU speed



The Completely Fair Scheduler

Linux 2.6.23 (2008)

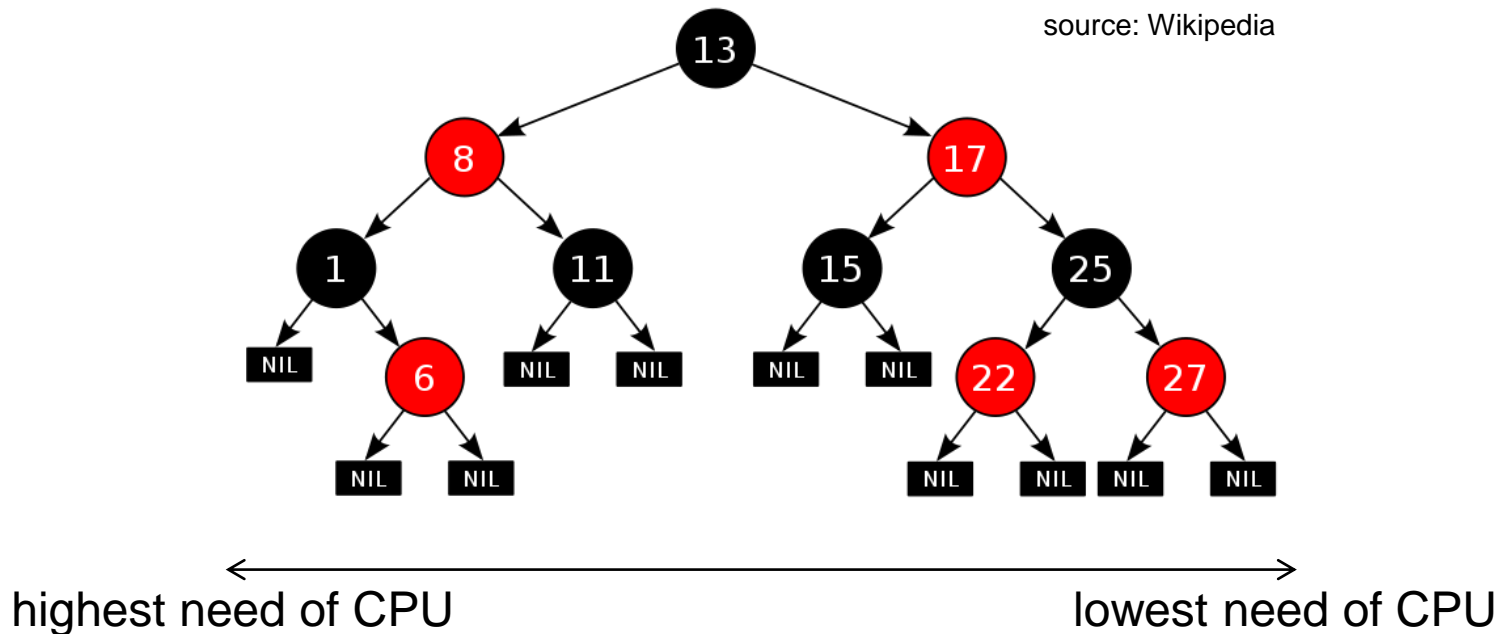
- Each task's (of equal priority) timeslice = $1/n$ of the default task latency
 - task latency = maximum wait time
 - default: 20ms
 - lower bound timeslice (default: 4ms) may force increasing the task latency
 - priorities modeled with proportional sharing
- Selection of next task to run = pick task with minimum runtime so far
 - runtime accounting: vruntime (nanosecond granularity)



The Completely Fair Scheduler

Linux 2.6.23 (2008)

- Efficiently finding the task with minimum runtime
 - red-black tree
 - insertion & deletion: $O(\log N)$, find minimum: $O(1)$
 - ordered by runtime
 - task with minimum runtime is always left-most node

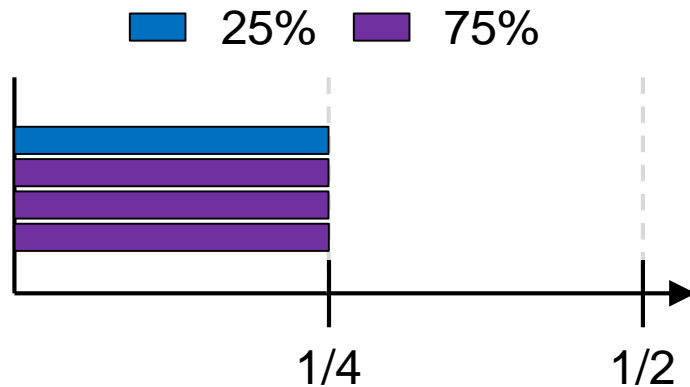


The Completely Fair Scheduler

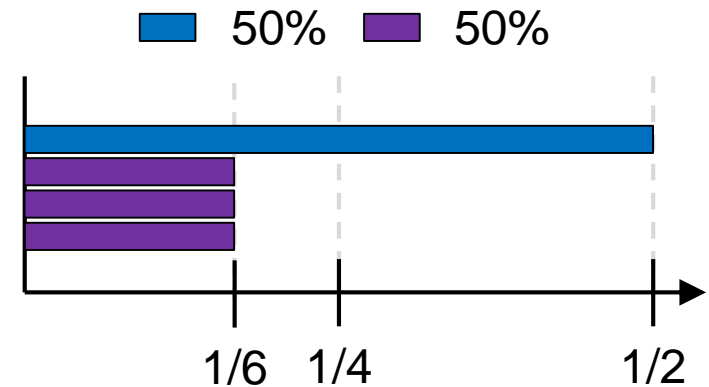
■ Group Scheduling

- divide CPU time equally between *groups* instead of *tasks*

■ user A ■ user B



no group scheduling



with group scheduling

The Completely Fair Scheduler

- Load balancing
 - active balancing
periodically pull tasks over from busiest CPU
 - idle balancing
as soon as there is no runnable task
 - ▶ migrated only if average idle time > migration cost