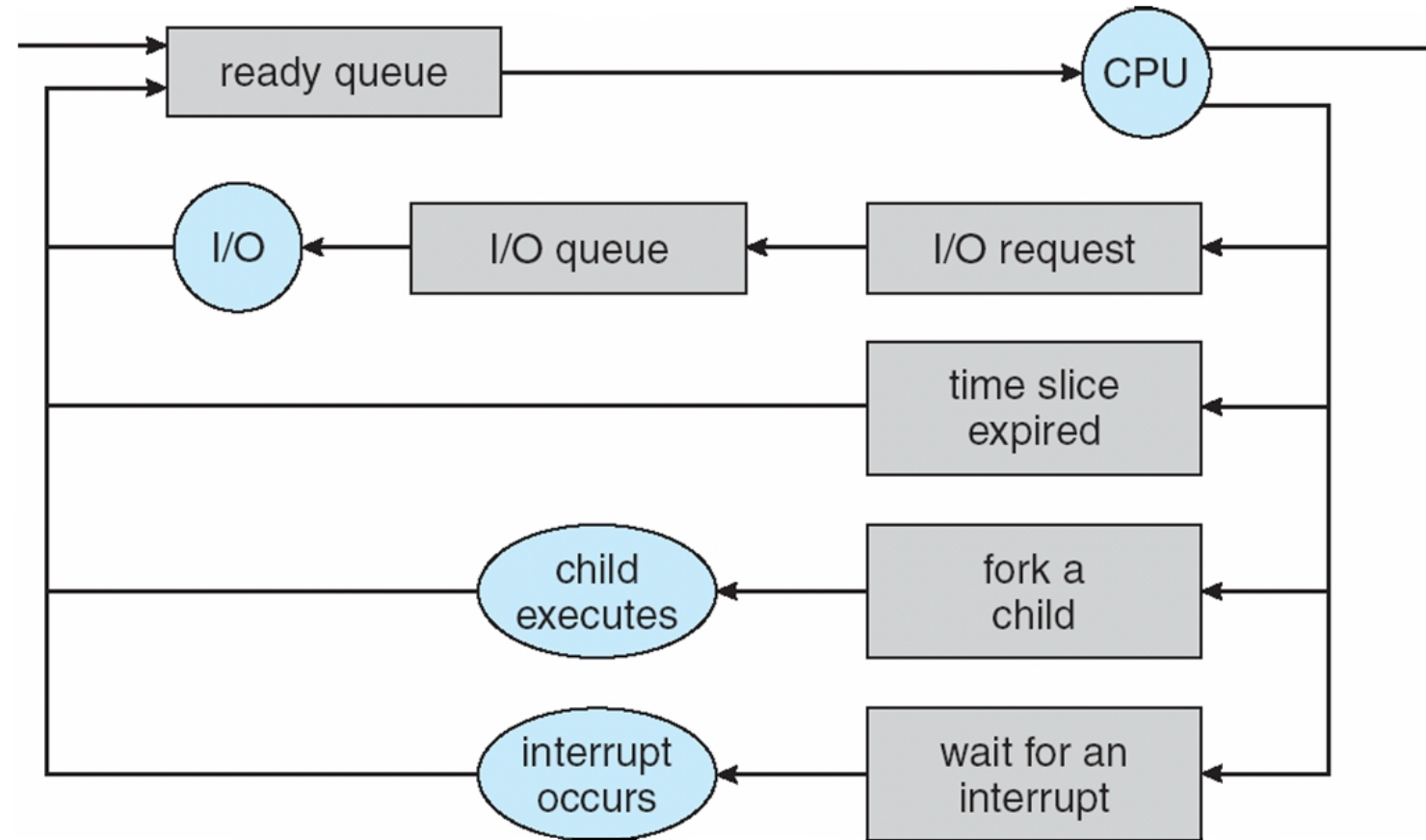


Operating Systems

Process Scheduling Algorithms

CPU Scheduling

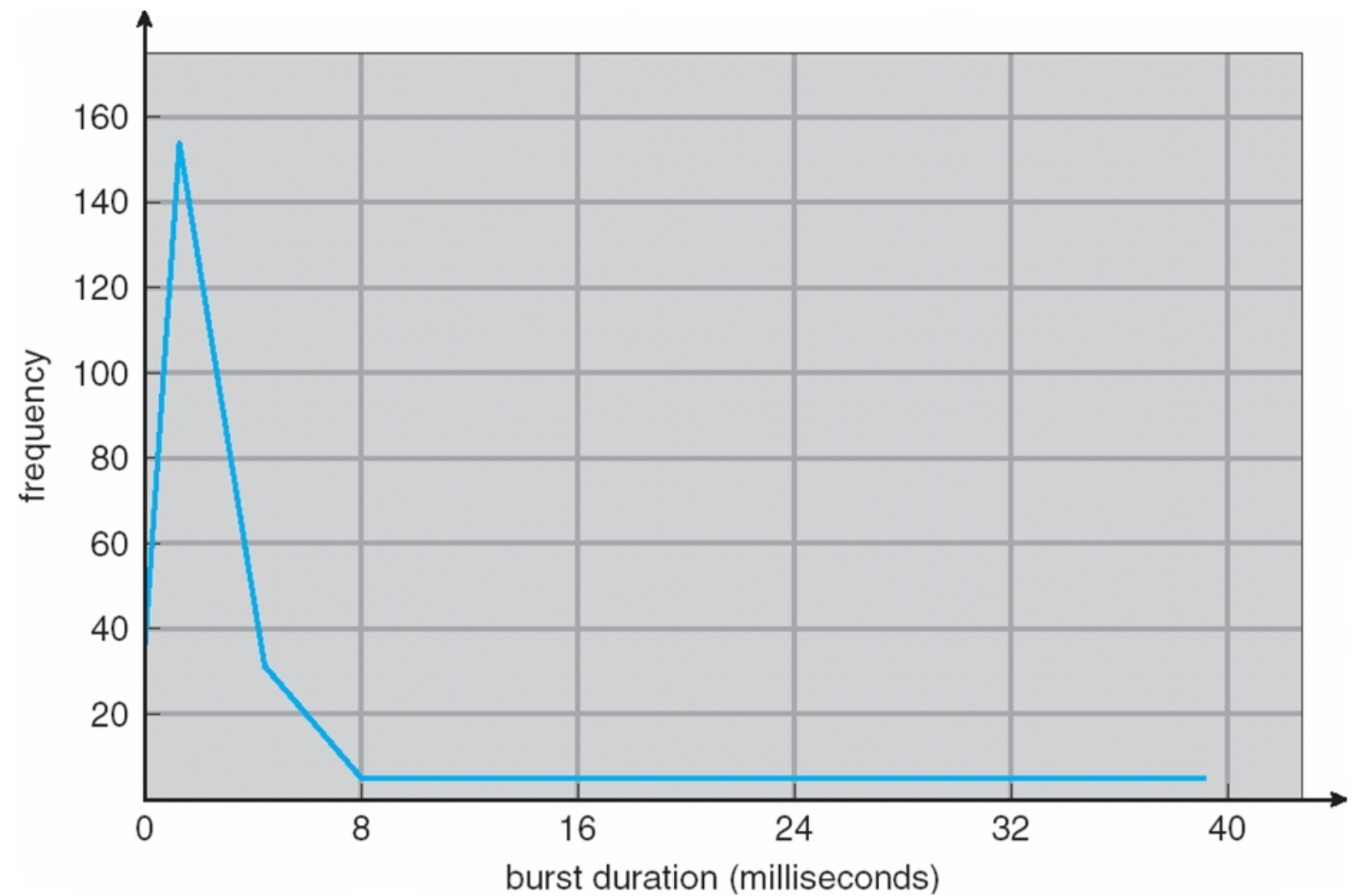
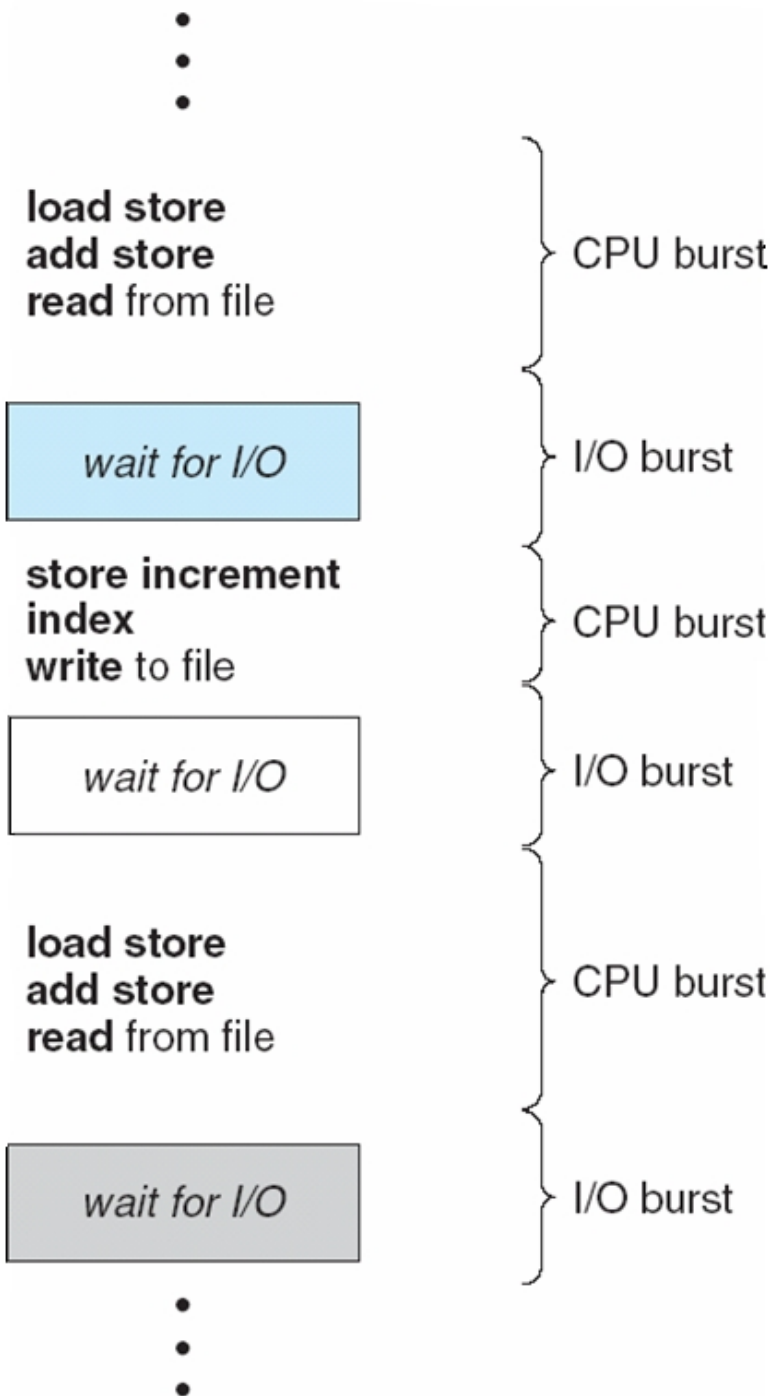


- **How is the OS to decide which of several tasks to take off a queue?**
- **Scheduling: deciding which threads are given access to resources from moment to moment.**

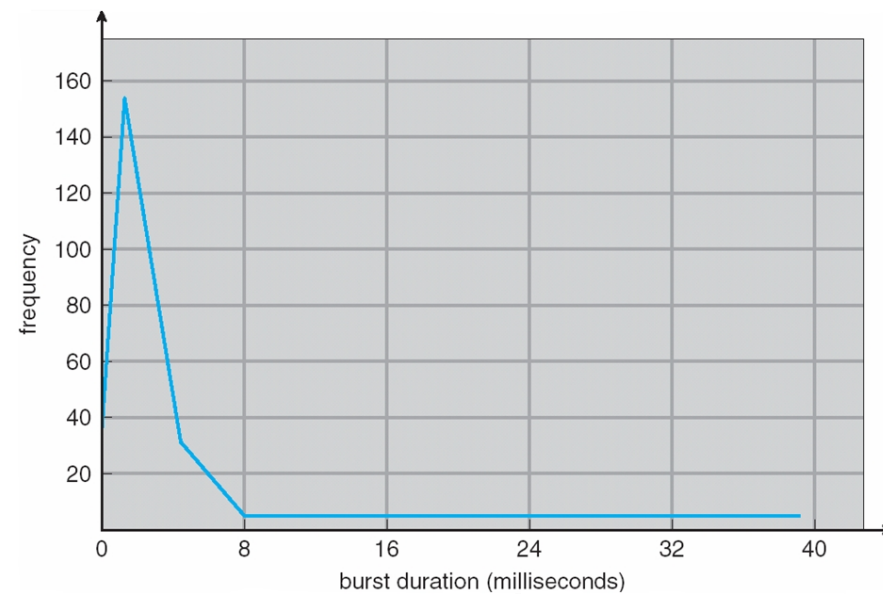
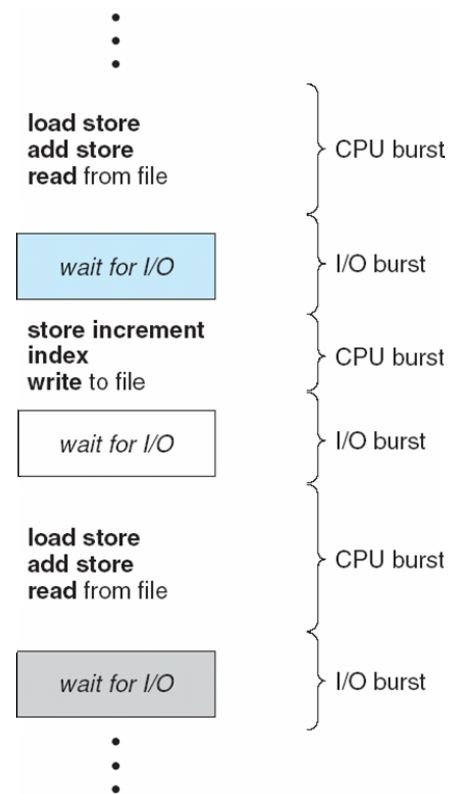
Assumptions about Scheduling

- **CPU scheduling big area of research in early '70s**
- **Many implicit assumptions for CPU scheduling:**
 - One program per user
 - One thread per program
 - Programs are independent
- **These are unrealistic but simplify the problem**
- **Does “fair” mean fairness among users or programs?**
 - If I run one compilation job and you run five, do you get five times as much CPU?
 - Often times, yes!
- **Goal: dole out CPU time to optimize some desired parameters of the system.**
 - What parameters?

Assumption: CPU Bursts



Assumption: CPU Bursts



- **Execution model: programs alternate between bursts of CPU and I/O**
 - Program typically uses the CPU for some period of time, then does I/O, then uses CPU again
 - Each scheduling decision is about which job to give to the CPU for use by its next CPU burst
 - With timeslicing, thread may be forced to give up CPU before finishing current CPU burst.

Process Execution Phases

```
move(A, B)
add(R1, 3)
move(R1, A)
load(File_1)
```

CPU-intensive

I/O-intensive

```
copy(R1, R3)
sub(R3, 5)
move(R3, B)
store(File_1)
```

CPU-intensive

I/O-intensive

```
div(B, 2)
move(B, A)
```

CPU-intensive

What is Important in a Scheduling Algorithm?



Performance Evaluation Criteria

- **fairness**
 - often: CPU time per process
- **Efficiency/CPU utilization**
 - resource utilization; often CPU
- **response time** (interactive users)
 - start of the process/burst until the first response to the user
- **turnaround time** (batch users)
 - start of the process/burst until the end of the process/burst
 - The interval from the time of submission of a process to the time of completion is the *turnaround time*. Waiting to get into memory + waiting in the ready queue + executing on the CPU + doing I/O

Evaluation Criteria (cont.)

- **throughput**
 - number of finished processes per time unit
- **waiting time**
 - time spent in the *ready* queue
- **context switches**
 - indication for the amount of overhead
- **complexity of the scheduling algorithm**
 - indication of the time needed to select the next process

What is Important in a Scheduling Algorithm?

- **Minimize Response Time**
 - Elapsed time to do an operation (job)
 - Response time is what the user sees
 - Time to echo keystroke in editor
 - Time to compile a program
 - Real-time Tasks: Must meet deadlines imposed by World
- **Maximize Throughput**
 - Jobs per second
 - Throughput related to response time, but not identical
 - Minimizing response time will lead to more context switching than if you maximized only throughput
 - Minimize overhead (context switch time) as well as efficient use of resources (CPU, disk, memory, etc.)
- **Fairness**
 - Share CPU among users in some equitable way
 - Not just minimizing average response time

Scheduling

- **deals with the allocation of resources to processes**
 - CPU scheduling
 - disk scheduling
- **important aspects**
 - when does a process get a resource
 - how long may the process keep the resource
 - how are conflicting requests resolved

Scheduling of Processes

- **long-term scheduling (job scheduling)**
 - a process is added to the pool of processes
 - medium-term scheduling
 - a process is swapped in or out of main memory
- **short-term scheduling**
 - which process will get the CPU
 - invoked in the following situations
 - a process is done with its task or CPU burst
 - a process has an I/O request
 - the time slice of a process is over
 - a new process with a higher priority arrives
 - possibly after an interrupt

Factors in Scheduling

- **is the process CPU- or I/O-bound**
- **interactive or batch**
- **process priority**
- **execution time used so far**
- **execution time required to complete**
- **preemption frequency**
- **page fault frequency**

CPU Scheduler

- **component of the operating system**
 - is itself a process
 - uses resources of the computer system
 - in particular CPU time, memory
 - should not consume too much CPU time
 - otherwise the overhead is too high

CPU Scheduler (cont.)

- **responsible for the selection of the next process to be executed**
- **the selection is done according to a scheduling algorithm**
- **short-term scheduling of processes**
 - **job scheduler and mid-term scheduler also schedule processes, but with a longer-term perspective**
 - **basic principles are similar for all three process schedulers**

Preemptive & non preemptive Scheduling

- **CPU-scheduling decisions may take place under the following four circumstances:**
- **When a process switches from the running state to the waiting state (for example, as the result of an I/O request or an invocation of wait for the termination of one of the child processes)**
- **When a process switches from the running state to the ready state (for example, when an interrupt occurs)**
- **When a process switches from the waiting state to the ready state (for example, at completion of I/O)**
- **When a process terminates**

Scheduling Algorithm

- **prescribes the way the next process is selected**
- **determines the order in which processes are executed**
- **defines the actions for the scheduler**

Scheduling Algorithms

- **types of scheduling algorithms**
 - preemptive, non-preemptive
- **instances of scheduling algorithms**
 - FCFS, SJF, SRTF, priority-based, round robin
 - multilevel, multilevel feedback

Non-Preemptive Scheduling

- **a process stays on the CPU until it voluntarily releases the CPU**
 - long waiting and response times
 - may lead to starvation
- **simple, easy to implement**
- **not suited for multi-user systems**
- **euphemism: “cooperative multitasking”**

Preemptive Scheduling

- **the execution of a process may be interrupted by the operating system at any time**
 - interrupt
 - higher priority process
 - arrival of a new process, change of status
 - time limit
- **prevents a process from using the CPU for too long**
- **may lead to race conditions**
 - can be solved by using process synchronization

Context Switch

- **removal of the “old” process from the CPU**
 - saves the process block in main memory
- **selection of the next process**
 - see scheduling
- **installation of the “new” process in the CPU**
 - read the process block from main memory
 - set the program counter, registers, etc.
 - (re-)start the process

Dispatcher

- **transfer of control to the newly selected process**
 - context switch
 - switch to user mode
 - continue execution of the process at the position indicated by the program counter
- **dispatch latency**
 - the time it takes to stop one process and start the execution of another one

First-Come, First-Served (FCFS)

- **principle**
 - processes are served in the order in which they arrive
 - even if they have the same arrival time, an order of arrival is distinguishable, or a process is chosen randomly

FCFS Example

- **five processes arrive at time 0 in the order P1, P2, P3, P4, P5**
- **processes have different (expected) burst lengths**
- **all processes have the same priority**
- **Gantt chart illustrates the execution of processes on the CPU**
- **comparison criteria with other algorithms**
 - waiting time
 - response time
 - turnaround time
 - number of context switches

Terminology

- **waiting time**
 - time a process spends in the ready queue
 - **not** in the waiting queue
 - may consist of several separate periods
- **response time**
 - time from the arrival of the process until its first activity
- **turnaround time**
 - time from the arrival until the termination of the process
 - usually only one single CPU burst is considered
- **number of context switches**
 - simplified assumption: for each period of CPU activity, a process requires one context switch
 - roughly one half context switch to be loaded into the CPU, and another half to be removed

Scheduling Algorithms: First-Come, First-Served (FCFS)

- **“Run until Done:” FIFO algorithm**
- **In the beginning, this meant one program runs non-preemptively until it is finished (including any blocking for I/O operations)**
- **Now, FCFS means that a process keeps the CPU until one or more threads block**
- **Example: Three processes arrive in order P1, P2, P3.**
 - **P1 burst time: 24**
 - **P2 burst time: 3**
 - **P3 burst time: 3**
- **Draw the Gantt Chart and compute Average Waiting Time and Average Completion Time.**

Scheduling Algorithms: First-Come, First-Served (FCFS)

- **Example:** Three processes arrive in order P1, P2, P3.

- P1 burst time: 24
- P2 burst time: 3
- P3 burst time: 3



- **Waiting Time**

- P1: 0
- P2: 24
- P3: 27

- **Completion Time:**

- P1: 24
- P2: 27
- P3: 30

- **Average Waiting Time:** $(0+24+27)/3 = 17$
- **Average Completion Time:** $(24+27+30)/3 = 27$

Scheduling Algorithms: First-Come, First-Served (FCFS)

- **What if their order had been P2, P3, P1?**
 - **P1 burst time: 24**
 - **P2 burst time: 3**
 - **P3 burst time: 3**

Scheduling Algorithms: First-Come, First-Served (FCFS)

- What if their order had been P2, P3, P1?

- P1 burst time: 24
- P2 burst time: 3
- P3 burst time: 3



- Waiting Time

- P1: 0
- P2: 3
- P3: 6

- Completion Time:

- P1: 3
- P2: 6
- P3: 30

- Average Waiting Time: $(0+3+6)/3 = 3$ (compared to 17)
- Average Completion Time: $(3+6+30)/3 = 13$ (compared to 27)

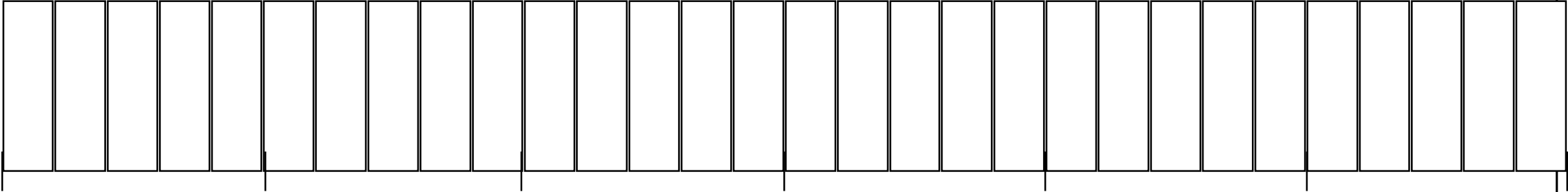
Scheduling Algorithms: First-Come, First-Served (FCFS)

- **Average Waiting Time: $(0+3+6)/3 = 3$ (compared to 17)**
- **Average Completion Time: $(3+6+30)/3 = 13$ (compared to 27)**
- **FIFO Pros and Cons:**
 - Simple (+)
 - Short jobs get stuck behind long ones (-)
 - If all you're buying is milk, doesn't it always seem like you are stuck behind a cart full of many items
 - Performance is highly dependent on the order in which jobs arrive (-)
 - Convoy effect(favors long processes compared to shorter ones)

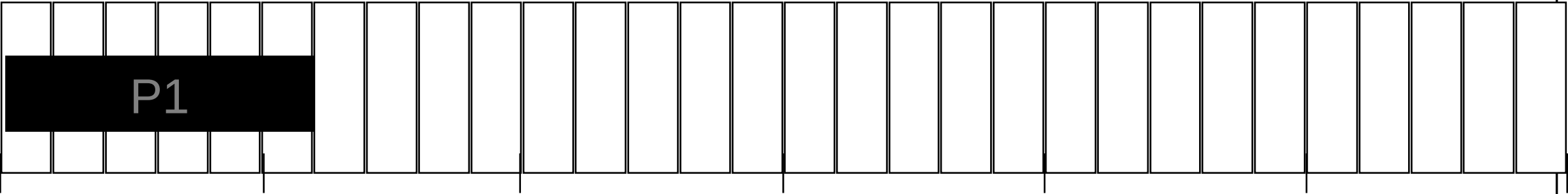
Convoy Effect

- **A convoy effect happens when a set of processes need to use a resource for a short time, and one process holds the resource for a long time, blocking all of the other processes. Causes poor utilization of the other resources in the system**

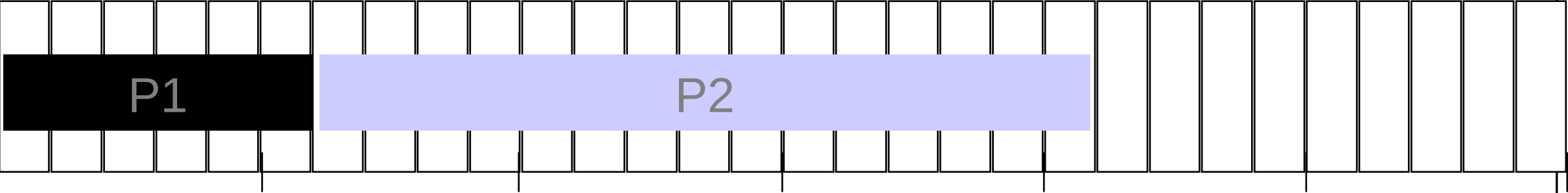
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



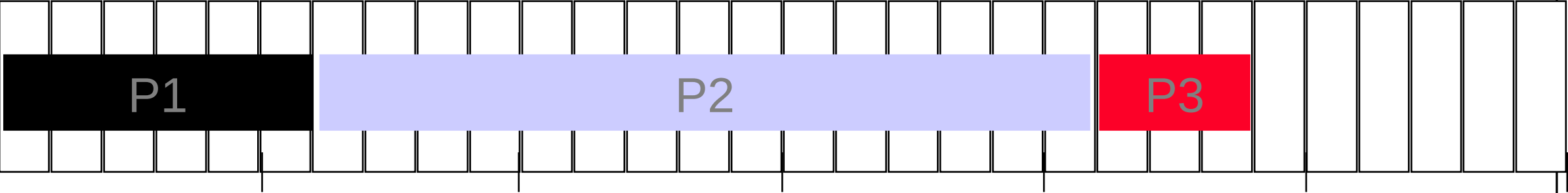
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



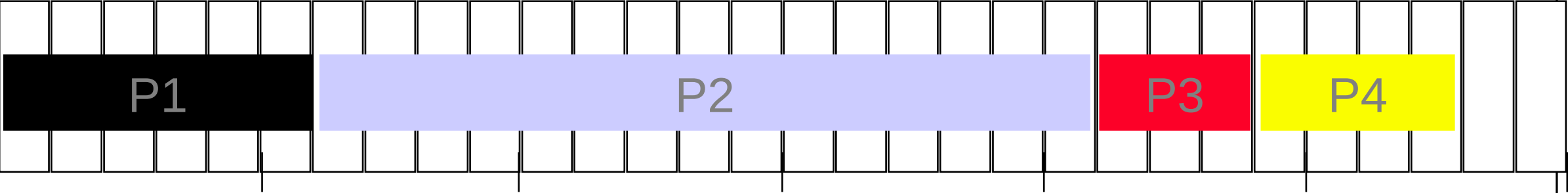
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



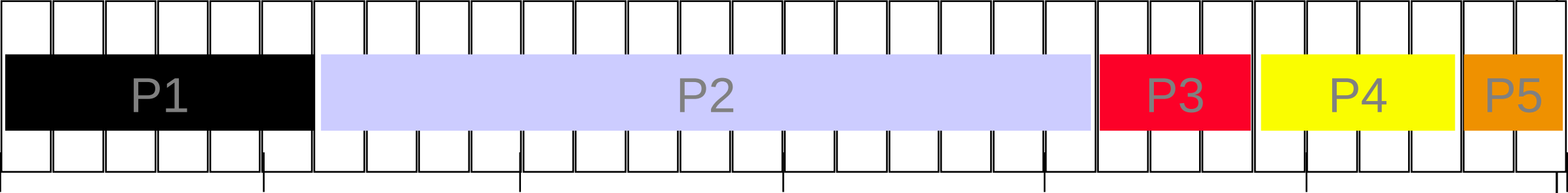
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



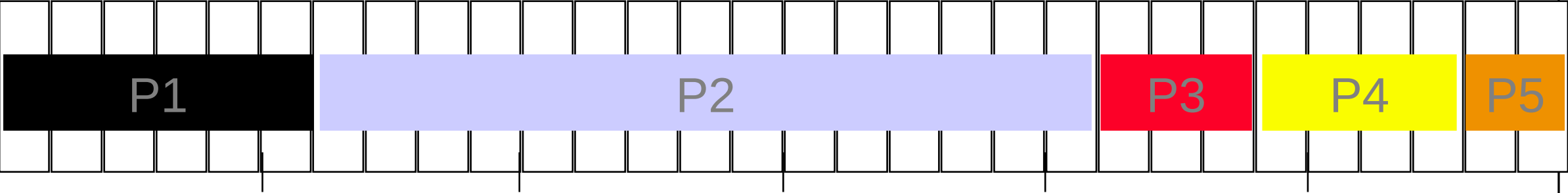
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



| <i>Process #</i> | <i>Waiting Time</i> | <i>Response Time</i> | <i>Turnaround Time</i> | <i>#of Context Switches</i> |
|------------------|---------------------|----------------------|------------------------|-----------------------------|
| P1 | 0 | 0 | 6 | 1 |
| P2 | 6 | 6 | 21 | 1 |
| P3 | 21 | 21 | 24 | 1 |
| P4 | 24 | 24 | 28 | 1 |
| P5 | 28 | 28 | 30 | 1 |
| Average | $79/5 = 15.8$ | $79/5 = 15.8$ | 21.8 | 1 |

How Can We Improve on This?



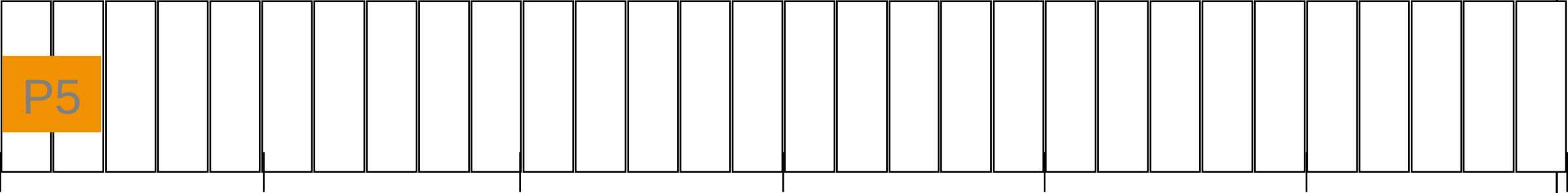
Shortest Job First (SJF)

- **principle:**
the process with the shortest CPU burst length is selected
 - obvious improvement from FCFS

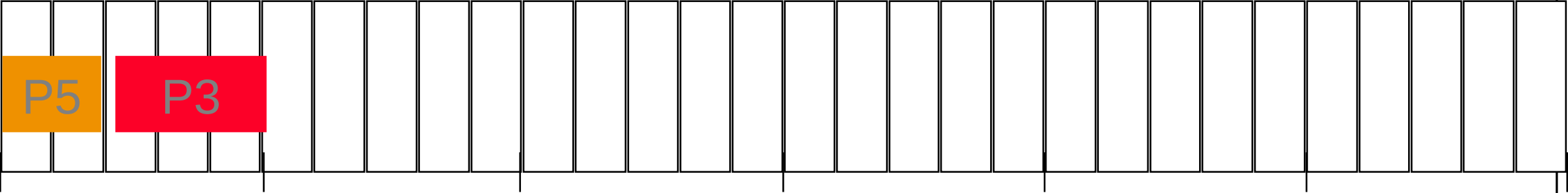
SJF Example

- **same conditions as in the FCFS example**
 - arrival time 0 for all processes
 - same priorities
 - different (expected) burst lengths

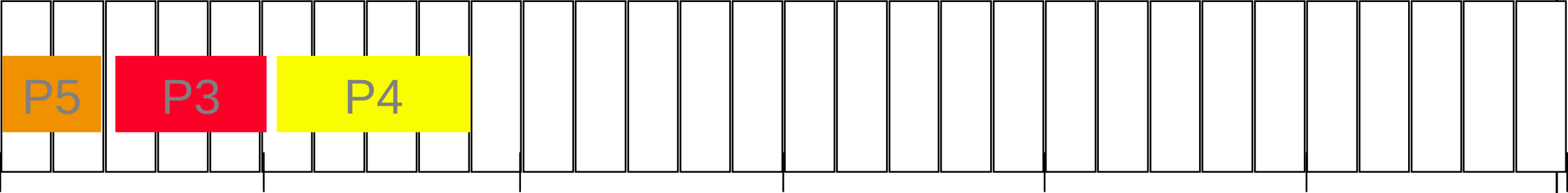
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



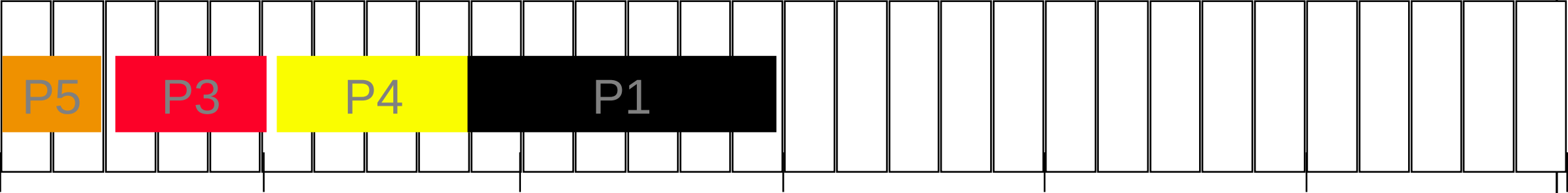
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



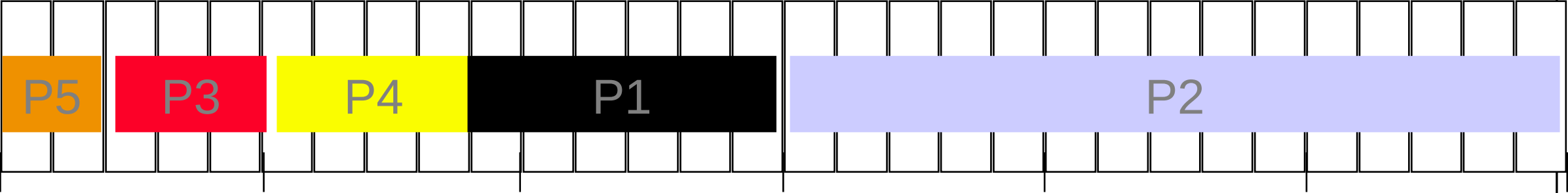
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



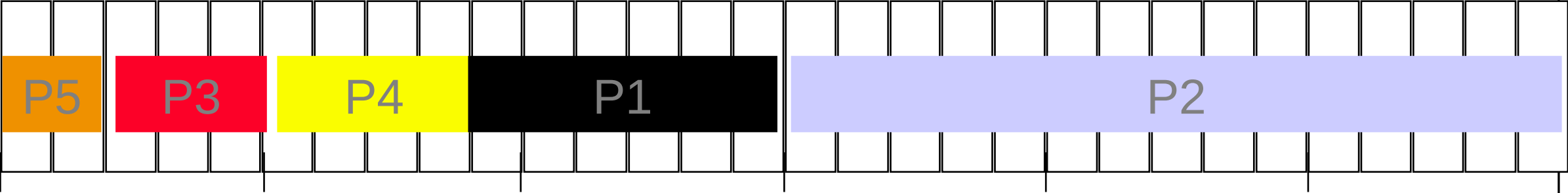
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 0 | 15 | 1 |
| P3 | 0 | 3 | 1 |
| P4 | 0 | 4 | 1 |
| P5 | 0 | 2 | 1 |



| <i>Process #</i> | <i>Waiting Time</i> | <i>Response Time</i> | <i>Turnaround Time</i> | <i>#of Context Switches</i> |
|------------------|---------------------|----------------------|------------------------|-----------------------------|
| P1 | 9 | 9 | 15 | 1 |
| P2 | 15 | 15 | 30 | 1 |
| P3 | 2 | 2 | 6 | 1 |
| P4 | 5 | 5 | 9 | 1 |
| P5 | 0 | 0 | 2 | 1 |
| Average | $31/5 = 6.2$ | $31/5 = 6.2$ | $62/5 = 12.4$ | 1 |

SJF Example 2

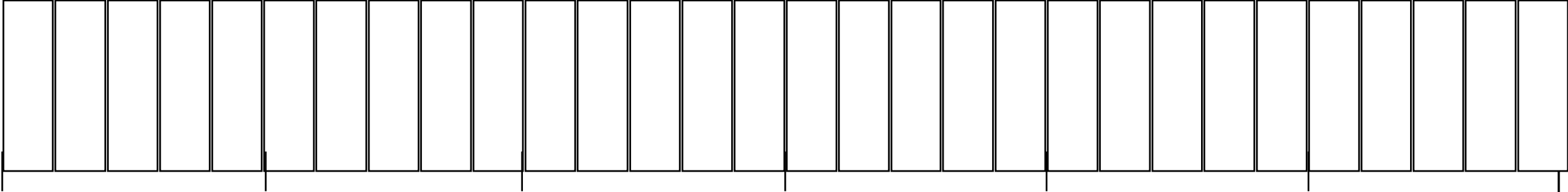
- **slight modification**
 - different arrival times for the processes
 - same priorities
 - different (expected) burst lengths
 - the processes waiting in the ready queue are added to the diagram

SJF Example 2

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 5 | 3 | 1 |
| P4 | 8 | 4 | 1 |
| P5 | 14 | 2 | 1 |

P1: 6

← ready queue at time 0



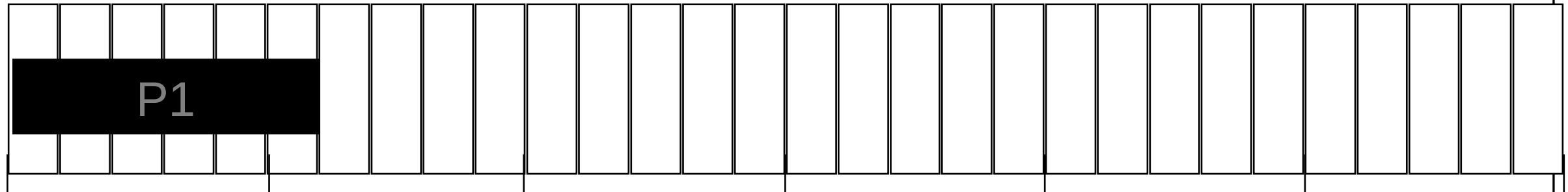
SJF Example 2

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 5 | 3 | 1 |
| P4 | 8 | 4 | 1 |
| P5 | 14 | 2 | 1 |

P1: 6

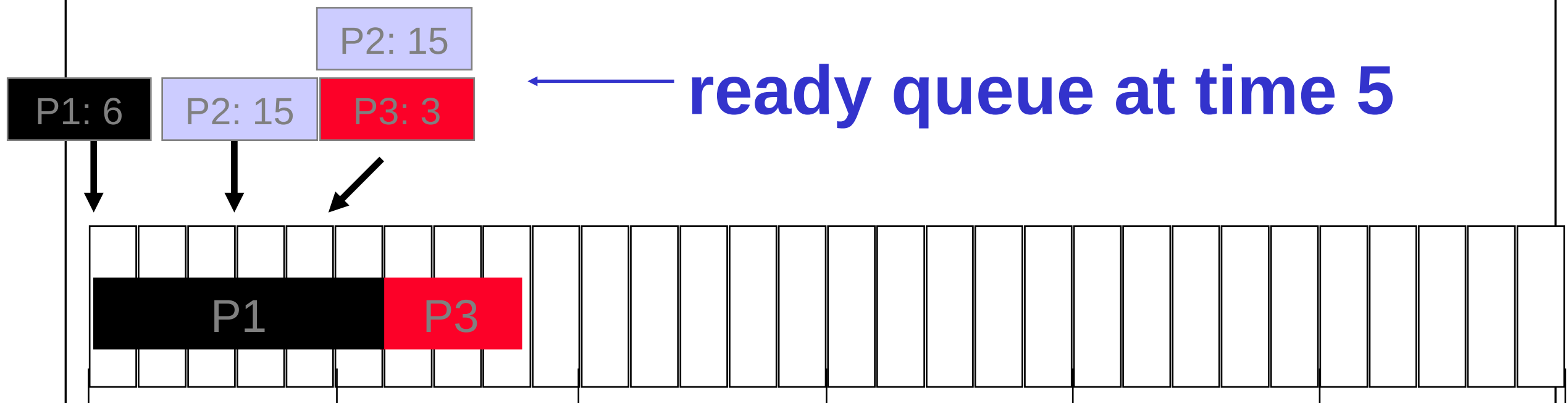
P2: 15

← ready queue at time 3



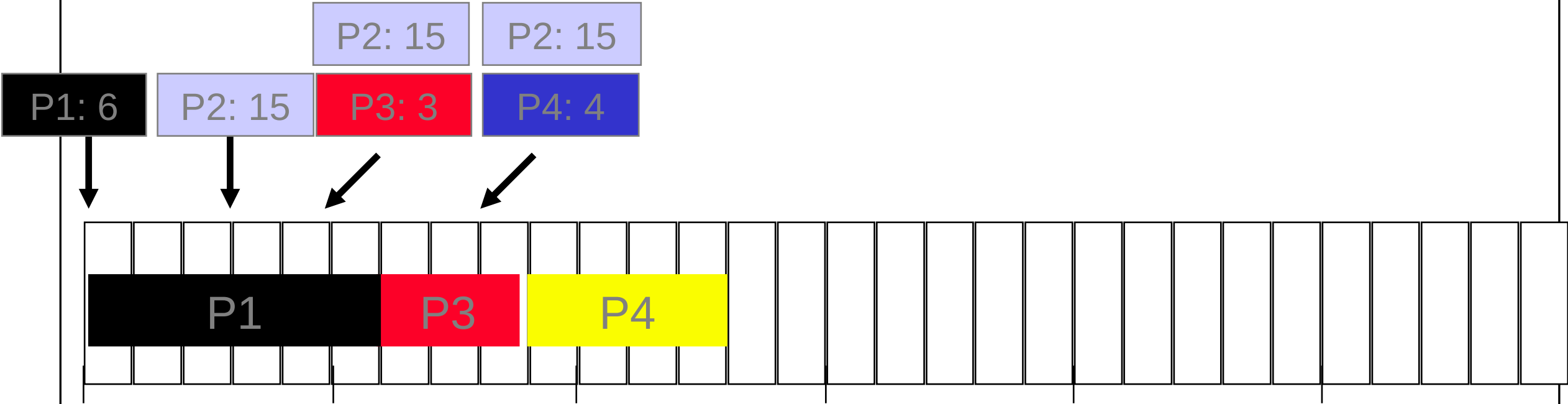
SJF Example 2

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 5 | 3 | 1 |
| P4 | 8 | 4 | 1 |
| P5 | 14 | 2 | 1 |



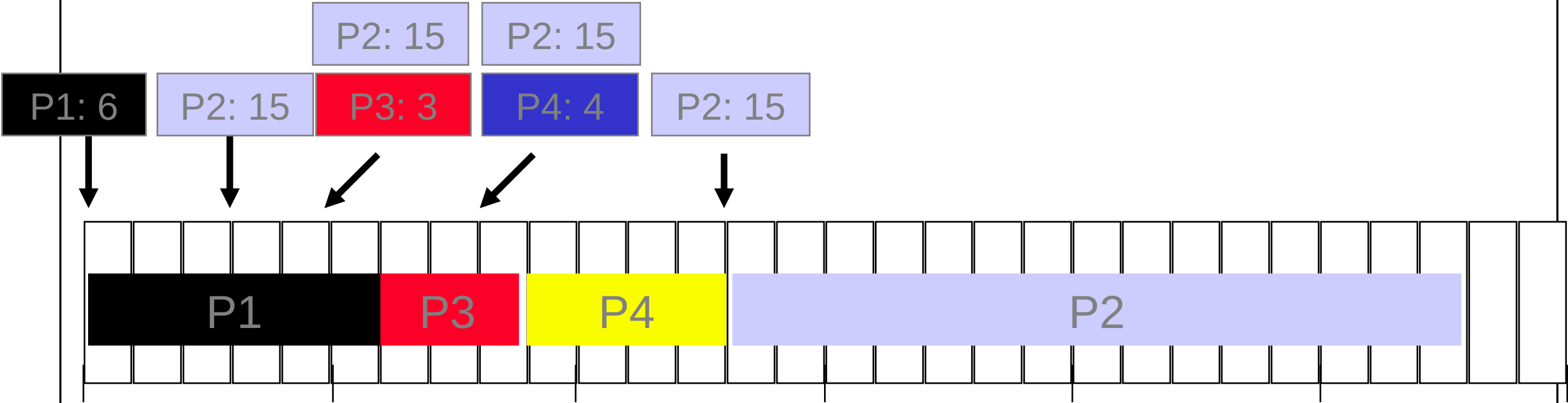
SJF Example 2

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 5 | 3 | 1 |
| P4 | 8 | 4 | 1 |
| P5 | 14 | 2 | 1 |



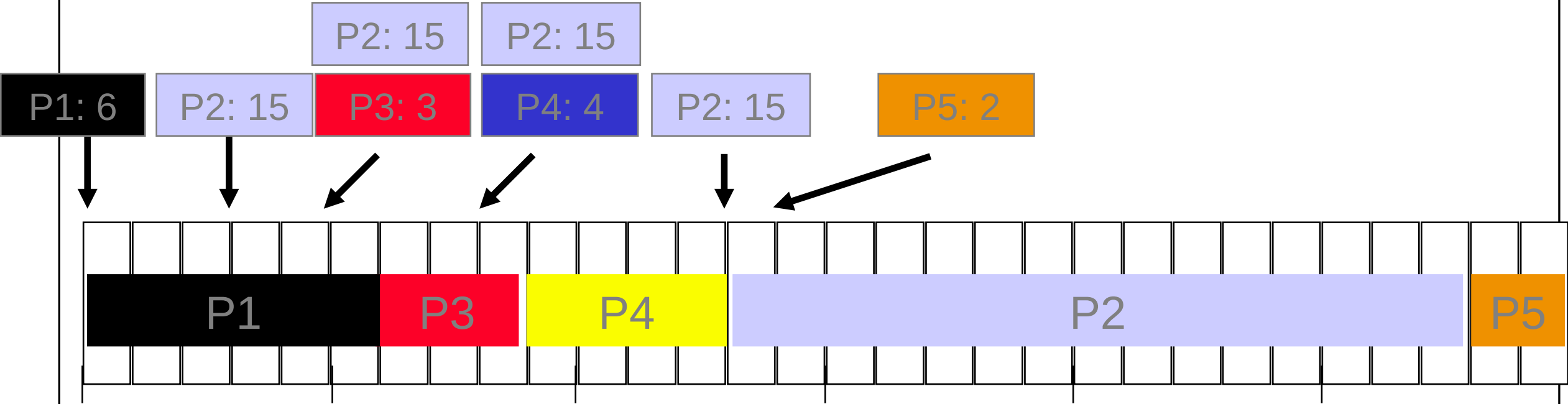
SJF Example 2

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 5 | 3 | 1 |
| P4 | 8 | 4 | 1 |
| P5 | 14 | 2 | 1 |



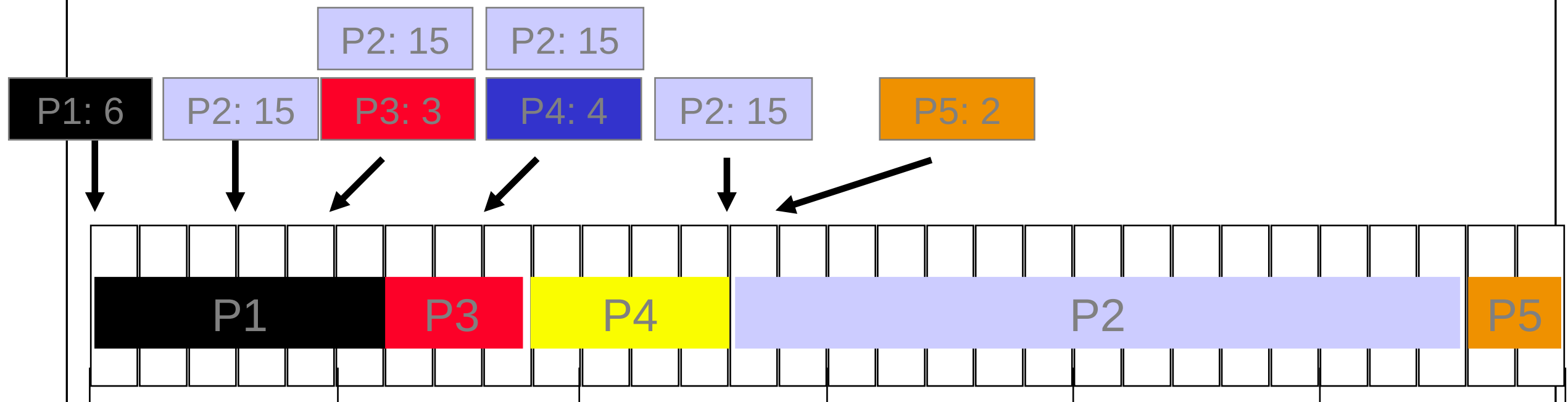
SJF Example 2

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 5 | 3 | 1 |
| P4 | 8 | 4 | 1 |
| P5 | 14 | 2 | 1 |



SJF Example 2

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 5 | 3 | 1 |
| P4 | 8 | 4 | 1 |
| P5 | 14 | 2 | 1 |



| <i>Process #</i> | <i>Waiting Time</i> | <i>Response Time</i> | <i>Turnaround Time</i> | <i>#of Context Switches</i> |
|------------------|---------------------|----------------------|------------------------|-----------------------------|
| P1 | 0 | 0 | 6 | 1 |
| P2 | $13-3 = 10$ | $13-3 = 10$ | $28-3 = 25$ | 1 |
| P3 | $6-5 = 1$ | $6-5 = 1$ | $9-5 = 4$ | 1 |
| P4 | $9-8 = 1$ | $9-8 = 1$ | $13-8 = 5$ | 1 |
| P5 | $28-14 = 14$ | $28-14 = 14$ | $30-14 = 16$ | 1 |
| Average | $26/5 = 5.2$ | $26/5 = 5.2$ | $50/5 = 10$ | 1 |

SJF properties

- **much better average waiting time than FCFS**
 - provably optimal with respect to the average waiting time
- **non-preemptive**
- **relies on knowing the length of the CPU bursts**
 - in general difficult to impossible
- **selection is more complex than FCFS**
 - linear w.r.t. number of processes in the ready queue
- **starvation is possible**
 - if new, short processes keep on arriving, old, long processes may never be served

Intermezzo: Burst Length

- **burst length prediction**
- **burst length estimation**
- **burst length calculation**
- **example**

Burst Length Prediction

- **in practice, the length of the next CPU burst of a process is not known**
 - as a consequence, algorithms such as SJF, SRTF in their pure form can't be used in practical systems
- **the CPU burst length can be estimated based on previous CPU bursts of the process**
- **estimation by analysis**
 - requires analysis of the code to be executed while the scheduling decision is made
 - in general intractable

Burst Length Estimation

- the length of the next CPU burst is estimated from the lengths of previous bursts
- frequently, recent bursts are given more importance than older bursts
- additional overhead during the scheduling decision
 - time to calculate the estimates
 - memory space to keep values of recent CPU burst lengths

Burst Length Calculation

- generic formula to estimate the length of the next CPU burst from the lengths of previous ones

$$E_i = a * B_{i-1} + (1-a) * E_{i-1}$$

- E_i estimate at time i
- B_{i-1} actual burst length at time $i-1$
- a factor to balance the importance of recent and not so recent bursts

Burst Length Example

- at time T_i we estimate the length of the next CPU burst based on information we have about previous bursts according to
$$E_i = a * B_{i-1} + (1-a) * E_{i-1}$$
 - we have to select a value for a (here 0.75)
 - for the very first burst B_0 , we have to guess, after that we use the measured time for the previous burst
 - the measured burst time may be significantly different from the estimate

Burst Length Example (cont.)

Time

Estimate

Previous

Burst

| | | | | | | |
|-------|-------------|-----|------------|-----|-----|-----|
| T_0 | $0.75 * 10$ | $+$ | $0.25 * 0$ | $=$ | 7.5 | --- |
|-------|-------------|-----|------------|-----|-----|-----|

Burst Length Example (cont.)

Time

Estimate

Previous

Burst

$$T_0 \quad 0.75 * 10 \quad + \quad 0.25 * 0 \quad = \quad 7.5 \quad \quad \quad 6$$

$$T_1 \quad 0.75 * 6 \quad + \quad 0.25 * 7.5 \quad = \quad 6.375$$

Burst Length Example (cont.)

Time

Estimate

Previous

Burst

$$T_0 \quad 0.75 * 10 \quad + \quad 0.25 * 0 \quad = 7.5 \quad 6$$

$$T_1 \quad 0.75 * 6 + 0.25 * 7.5 \quad = 6.375 \quad \mathbf{3}$$

$$T_2 \quad 0.75 * 3 + 0.25 * 6.375 \quad = \mathbf{3.85}$$

Burst Length Example (cont.)

Time

Estimate

Previous

Burst

$$T_0 \quad 0.75 * 10 \quad + \quad 0.25 * 0 \quad = 7.5 \quad 6$$

$$T_1 \quad 0.75 * 6 + 0.25 * 7.5 \quad = 6.375 \quad 3$$

$$T_2 \quad 0.75 * 3 + 0.25 * 6.375 \quad = 3.85 \quad 7$$

$$T_3 \quad 0.75 * 7 + 0.25 * 3.85 \quad = 6.2$$

Burst Length Example (cont.)

| <i>Time</i> | <i>Estimate</i> | <i>Previous</i> | <i>Burst</i> |
|-------------|-----------------|-----------------|--------------|
|-------------|-----------------|-----------------|--------------|

| | | | |
|-------|------------------------------|--|---|
| T_0 | $0.75 * 10 + 0.25 * 0 = 7.5$ | | 6 |
|-------|------------------------------|--|---|

| | | | |
|-------|---------------------------------|--|---|
| T_1 | $0.75 * 6 + 0.25 * 7.5 = 6.375$ | | 3 |
|-------|---------------------------------|--|---|

| | | | |
|-------|----------------------------------|--|---|
| T_2 | $0.75 * 3 + 0.25 * 6.375 = 3.85$ | | 7 |
|-------|----------------------------------|--|---|

| | | | |
|-------|--------------------------------|--|----|
| T_3 | $0.75 * 7 + 0.25 * 3.85 = 6.2$ | | 12 |
|-------|--------------------------------|--|----|

| | | | |
|-------|---------------------------------|--|--|
| T_4 | $0.75 * 12 + 0.25 * 6.2 = 9.55$ | | |
|-------|---------------------------------|--|--|

Shortest Remaining Time First (SRTF)

- **principle:**
the process with the shortest remaining time is selected
 - remaining time is the CPU burst length minus the time the CPU has already spent serving the process
- if a process with a shorter CPU burst length than the remaining time of the current process arrives, the current process is preempted

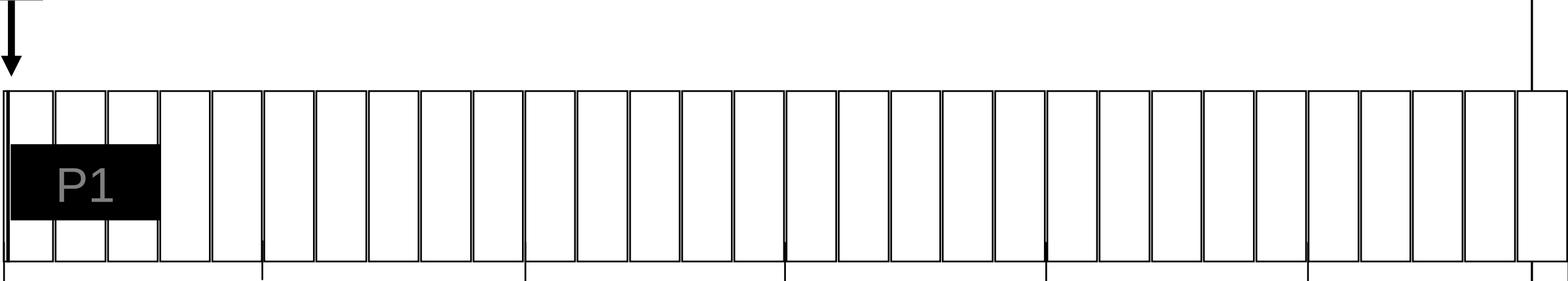
SRTF Example

- **arrival time of new processes is important**
- **it is useful to keep track of the processes currently in the ready queue**
 - **policy used here: preempted processes go to the end of the ready queue**
 - **may depend on the actual implementation**
- **a scheduling decision must be made when**
 - **a process is done with its CPU burst**
 - **a new process arrives in the ready queue**

SRTF Example

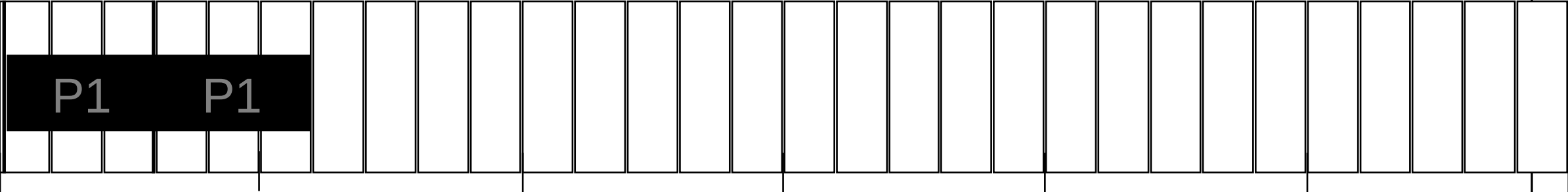
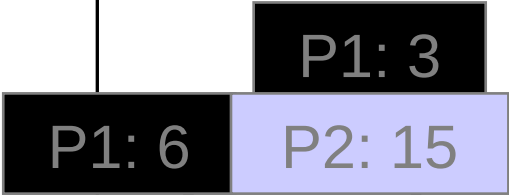
| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |

P1: 6



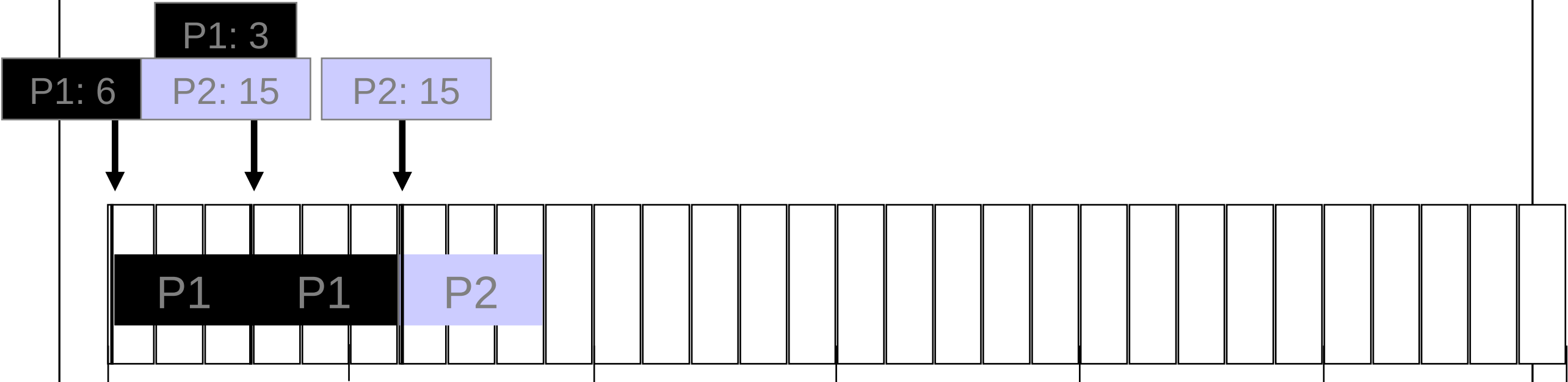
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



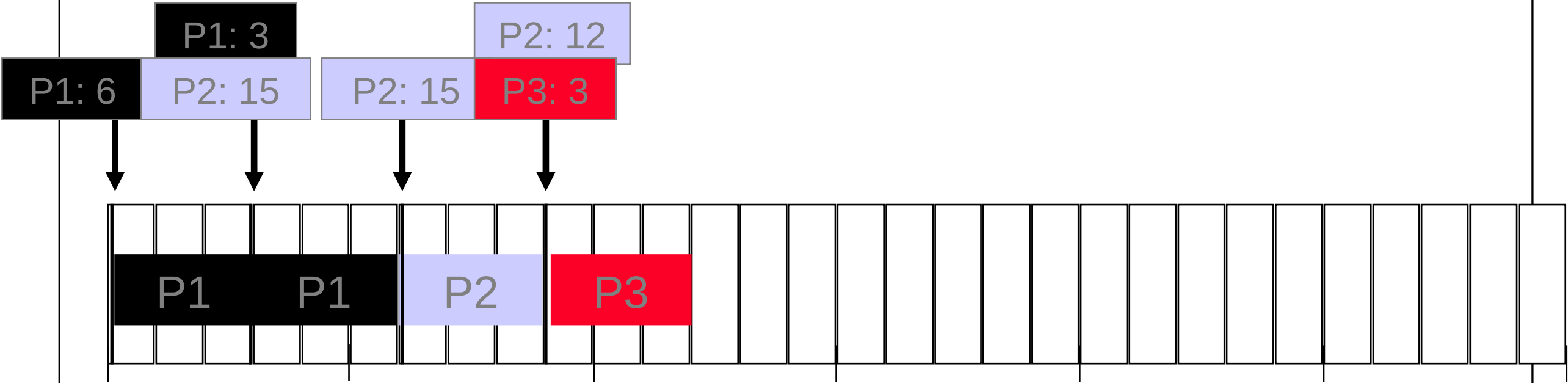
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



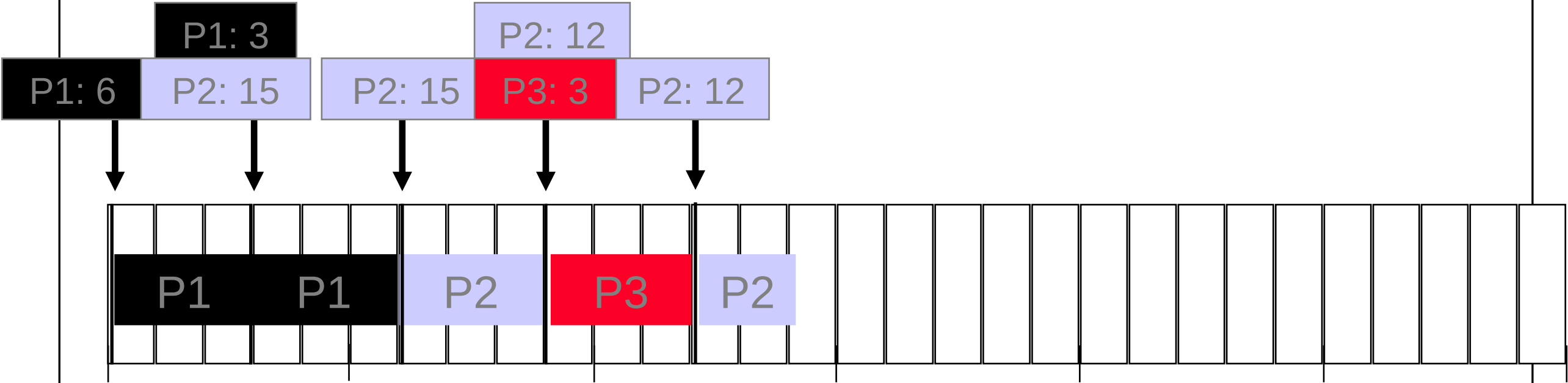
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



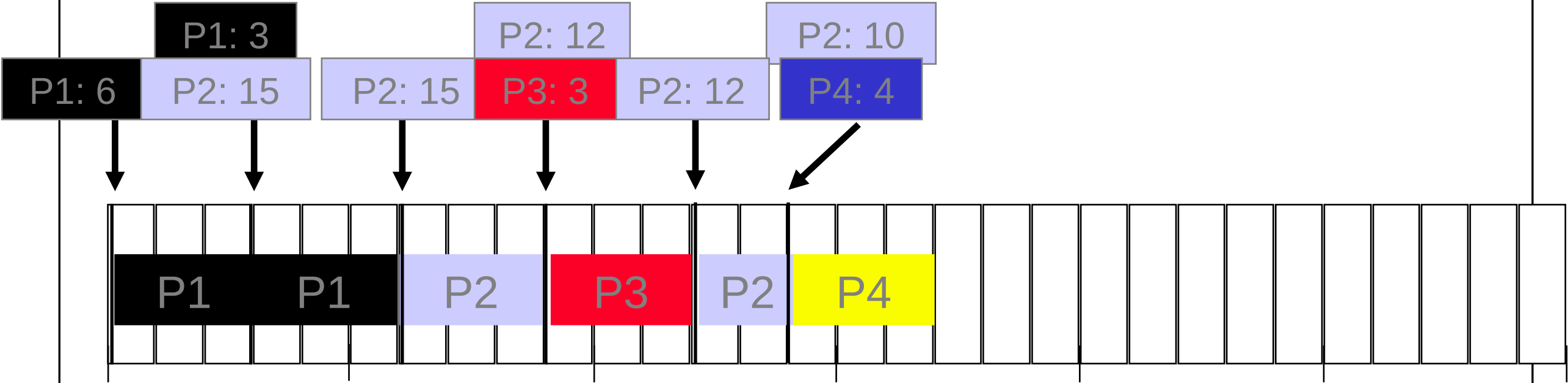
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



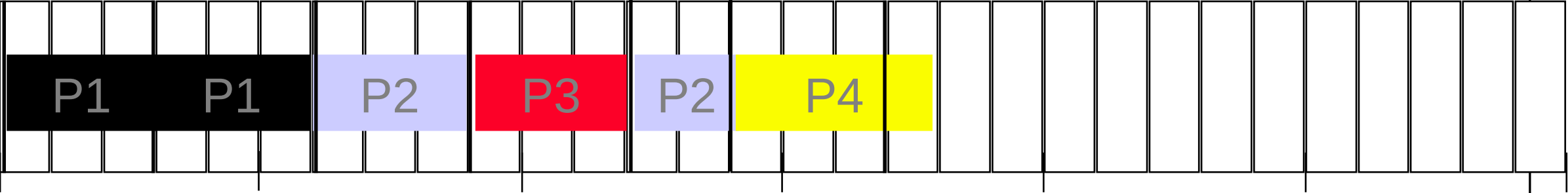
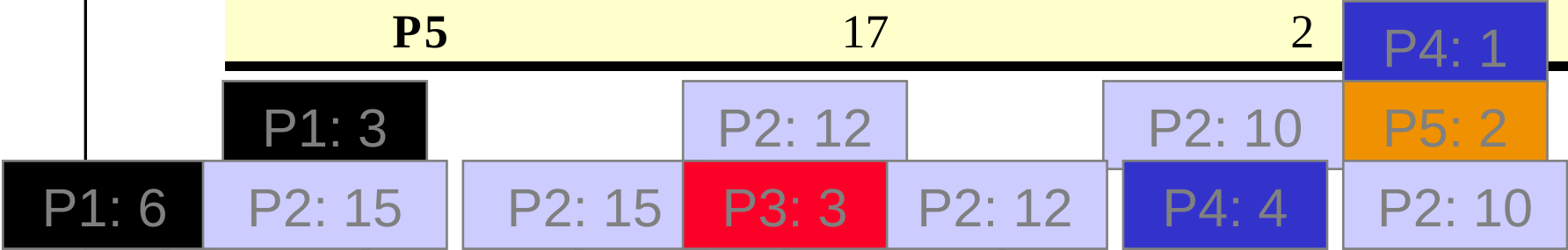
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



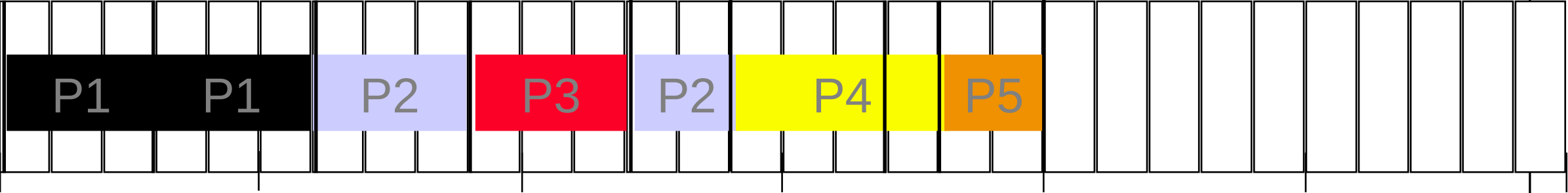
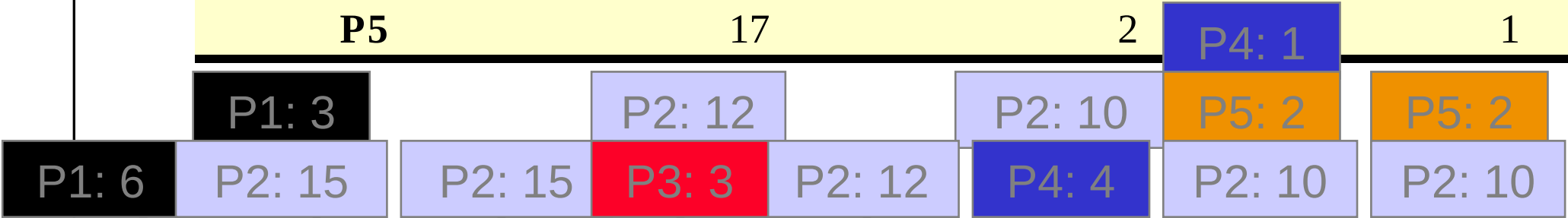
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



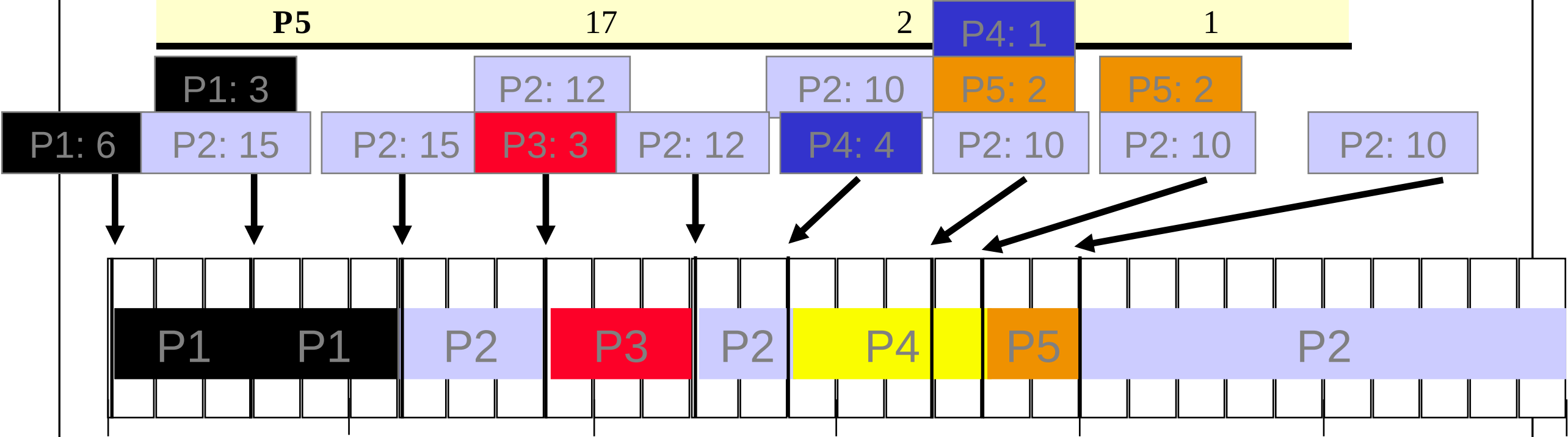
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



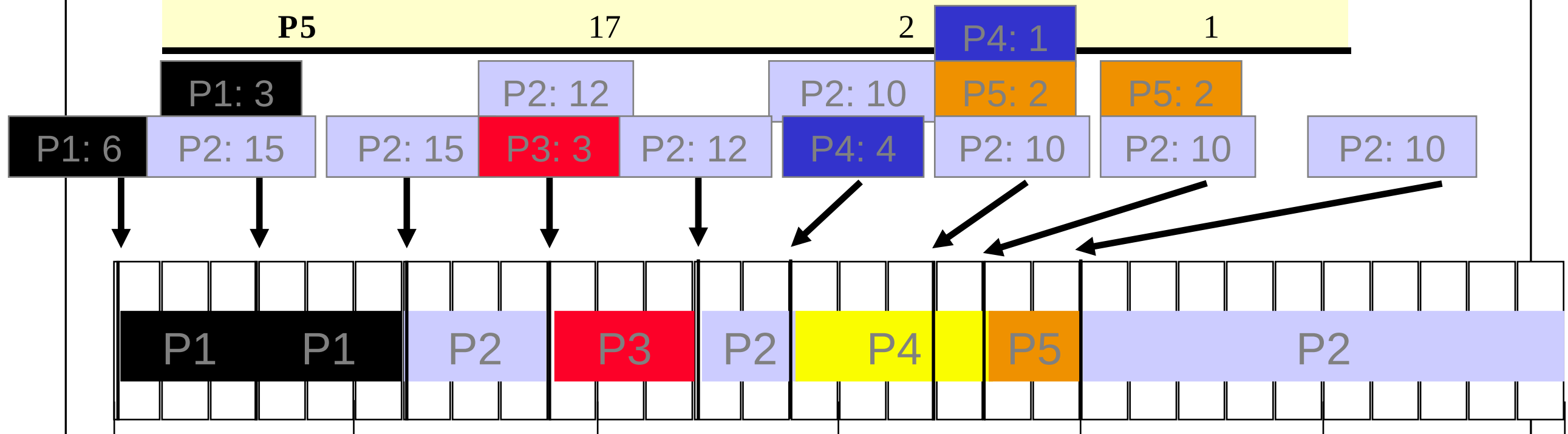
SRTF Example

| <i>Process #</i> | <i>Arrival Time</i> | <i>Burst Length</i> | <i>Priority</i> |
|------------------|---------------------|---------------------|-----------------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



SRTF Example

| Process # | Arrival Time | Burst Length | Priority |
|-----------|--------------|--------------|----------|
| P1 | 0 | 6 | 1 |
| P2 | 3 | 15 | 1 |
| P3 | 9 | 3 | 1 |
| P4 | 14 | 4 | 1 |
| P5 | 17 | 2 | 1 |



| Process # | Waiting Time | Response Time | Turnaround Time | #of Context Switches |
|-----------|---------------------------------|---------------|-----------------|----------------------|
| P1 | 0 | 0 | 6 | 1 (2) |
| P2 | $(6-3) + (12-9) + (20-14) = 12$ | $(6-3) = 12$ | $(30 - 3) = 27$ | 3 |
| P3 | 0 | 0 | $(12-9) = 3$ | 1 |
| P4 | 0 | 0 | $(18-14) = 3$ | 1 (2) |
| P5 | $(18-17) = 1$ | $(18-17) = 1$ | $(20-17) = 3$ | 1 |
| Average | $13/5 = 2.6$ | $13/5 = 2.6$ | $36/5 = 7.2$ | 7 |

SRTF Properties

- **good response time for short processes**
 - attractive for multi-user systems
- **preemptive version of the SJF algorithm**
 - higher overhead (context switches, elapsed time table)
- **starvation possible**
- **impractical due to burst length prediction problem**

Highest Response Ratio Next

- **principle**
 - **priority of a process is a function of its execution time and the time it has been waiting for service**
 - **priority = (time waiting + execution time) / execution time**

HRNN Properties

- **non-preemptive in its basic form**
 - preemptive variations exist
- **preference for shorter processes over longer ones**
- **aging ensures that long processes will eventually get the CPU**
- **still requires estimation of the execution time (burst length)**
- **minimum priority = 1.0**

Round Robin (RR) Scheduling


- **FCFS Scheme: Potentially bad for short jobs!**
 - Depends on submit order
 - If you are first in line at the supermarket with milk, you don't care who is behind you; on the other hand...
- **Round Robin Scheme**
 - Each process gets a small unit of CPU time (time quantum)
 - Usually 10-100 ms
 - After quantum expires, the process is preempted and added to the end of the ready queue
 - Suppose N processes in ready queue and time quantum is Q ms:
 - Each process gets $1/N$ of the CPU time
 - In chunks of at most Q ms
 - What is the maximum wait time for each process?



Round Robin (RR) Scheduling

- **FCFS Scheme: Potentially bad for short jobs!**
 - Depends on submit order
 - If you are first in line at the supermarket with milk, you don't care who is behind you; on the other hand...
- **Round Robin Scheme**
 - Each process gets a small unit of CPU time (time quantum)
 - Usually 10-100 ms
 - After quantum expires, the process is preempted and added to the end of the ready queue
 - Suppose N processes in ready queue and time quantum is Q ms:
 - Each process gets $1/N$ of the CPU time
 - In chunks of at most Q ms
 - What is the maximum wait time for each process?
 - No process waits more than $(n-1)q$ time units

Round Robin (RR) Scheduling

- **Round Robin Scheme**
 - Each process gets a small unit of CPU time (time quantum)
 - Usually 10-100 ms
 - After quantum expires, the process is preempted and added to the end of the ready queue
 - Suppose N processes in ready queue and time quantum is Q ms:
 - Each process gets $1/N$ of the CPU time
 - In chunks of at most Q ms
 - What is the maximum wait time for each process?
 - No process waits more than $(n-1)q$ time units
- **Performance Depends on Size of Q**
 - Small $Q \Rightarrow$ interleaved
 - Large Q is like... 
 - Q must be large w.r.t. to context switch time, otherwise overhead is too high (spending most of your time context switching!)

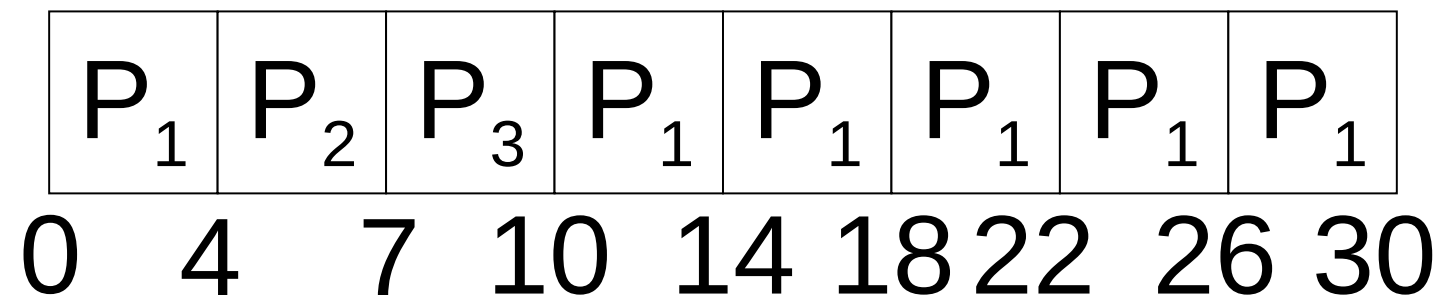
Round Robin (RR) Scheduling

- **Round Robin Scheme**
 - Each process gets a small unit of CPU time (time quantum)
 - Usually 10-100 ms
 - After quantum expires, the process is preempted and added to the end of the ready queue
 - Suppose N processes in ready queue and time quantum is Q ms:
 - Each process gets $1/N$ of the CPU time
 - In chunks of at most Q ms
 - What is the maximum wait time for each process?
 - No process waits more than $(n-1)q$ time units
- **Performance Depends on Size of Q**
 - Small $Q \Rightarrow$ interleaved, Q must be greater than the context switch time, otherwise overhead is too high.
 - Large Q is like FCFS
 - Q must be large with respect to context switch time, otherwise overhead is too high (spending most of your time context switching!)

Example of RR with Time Quantum = 4

| <u>Process</u> | <u>Burst Time</u> |
|----------------|-------------------|
| P_1 | 24 |
| P_2 | 3 |
| P_3 | 3 |

- The Gantt chart is:



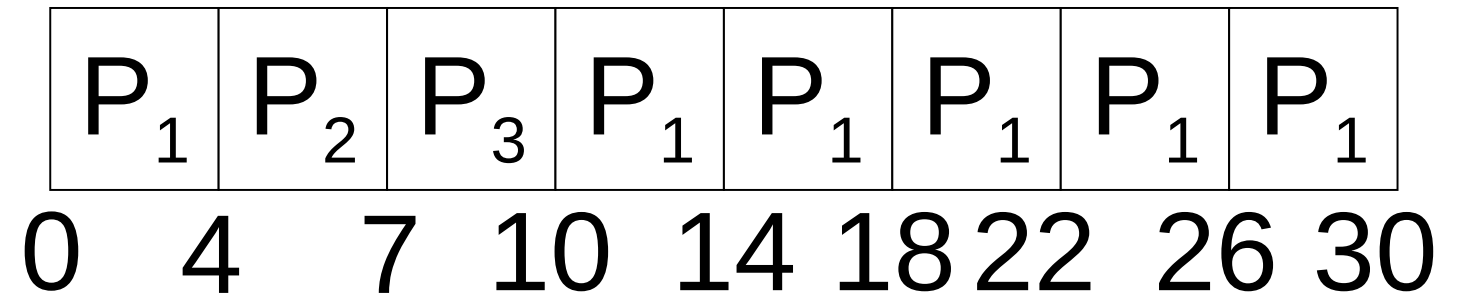
Example of RR with Time Quantum = 4

Process Burst Time

P_1 24

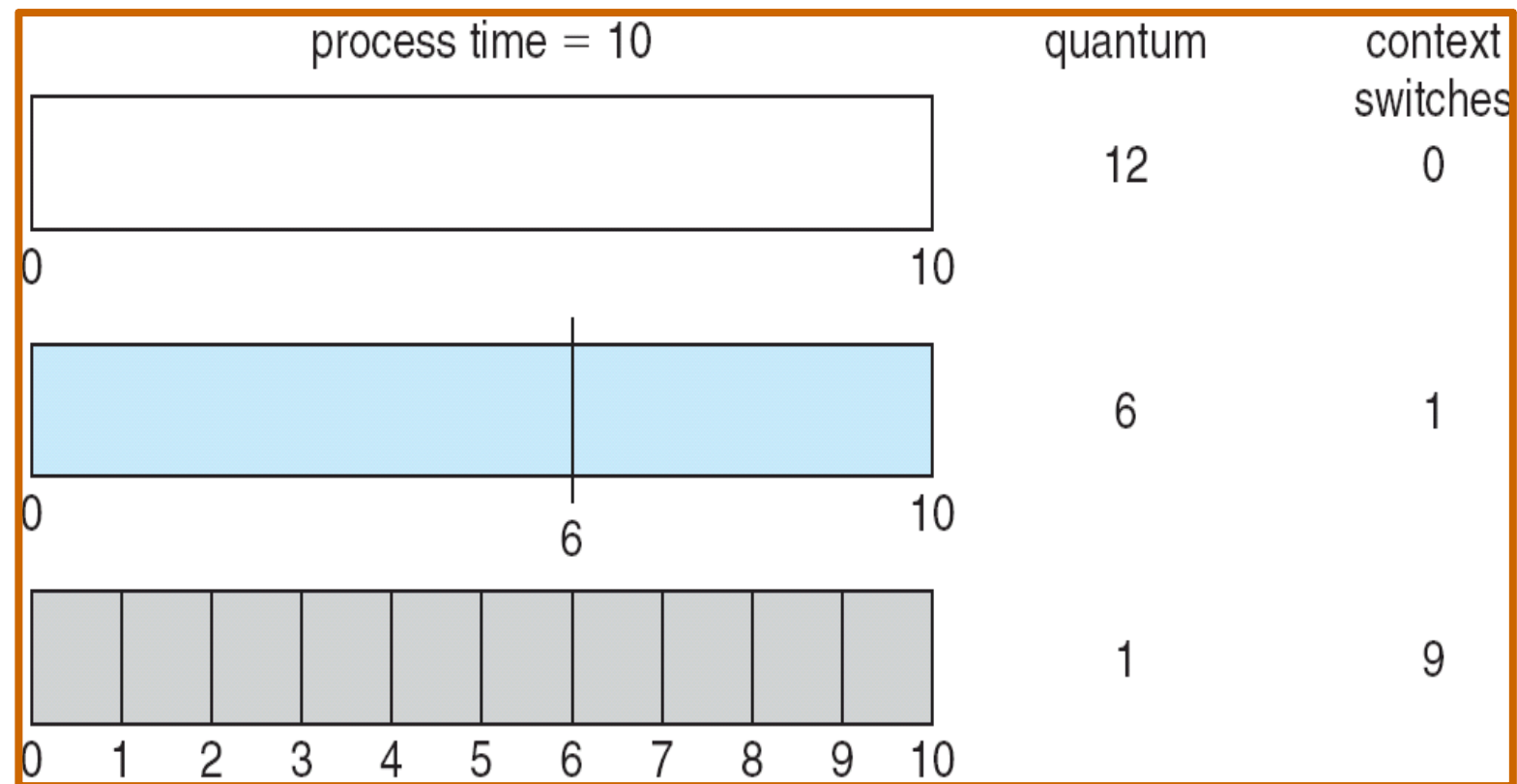
P_2 3

P_3 3

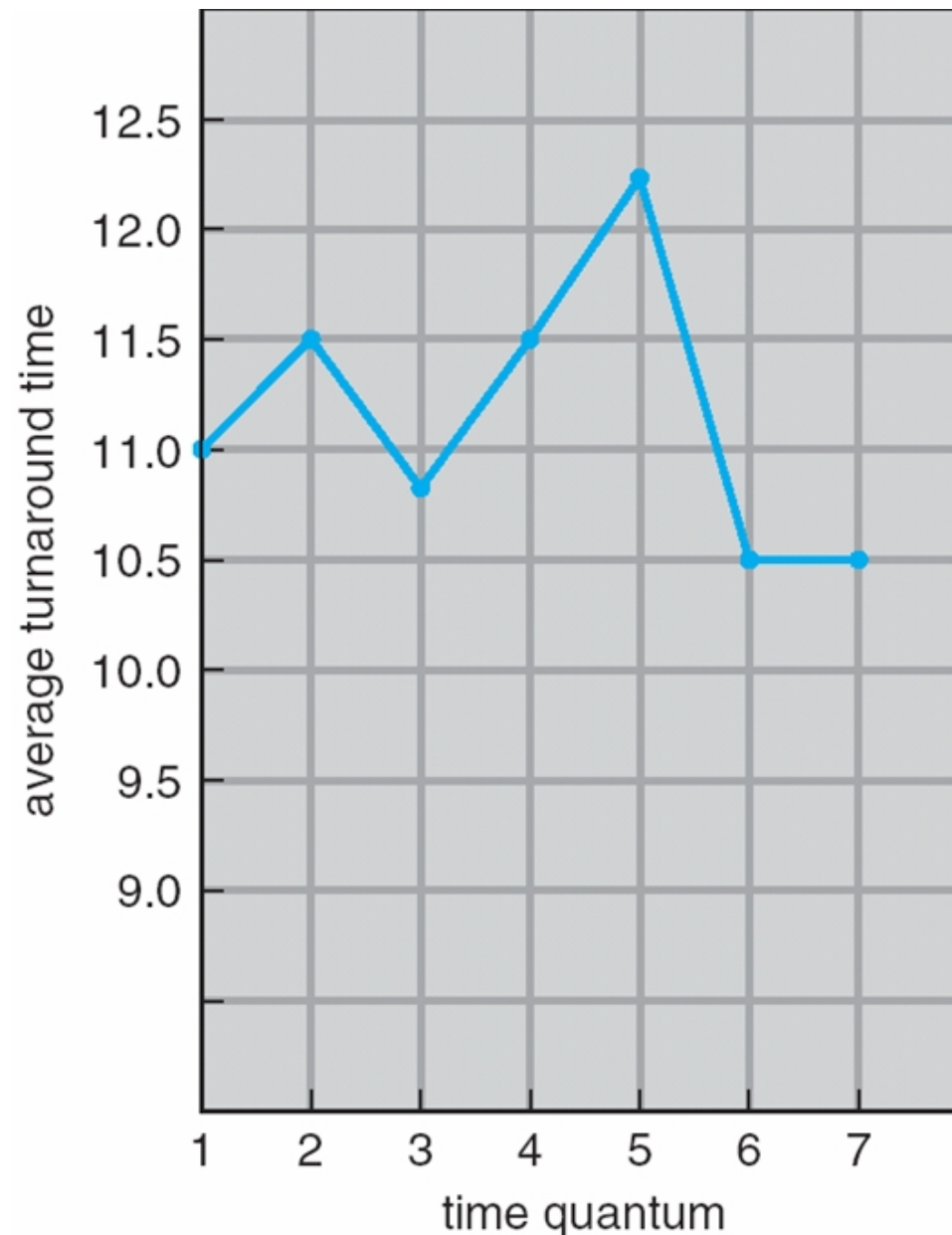


- **Waiting Time:**
 - P_1 : $(10-4) = 6$
 - P_2 : $(4-0) = 4$
 - P_3 : $(7-0) = 7$
- **Completion Time:**
 - P_1 : 30
 - P_2 : 7
 - P_3 : 10
- **Average Waiting Time:** $(6 + 4 + 7)/3 = 5.67$
- **Average Completion Time:** $(30+7+10)/3=15.67$

Time Quantum & Context Switches



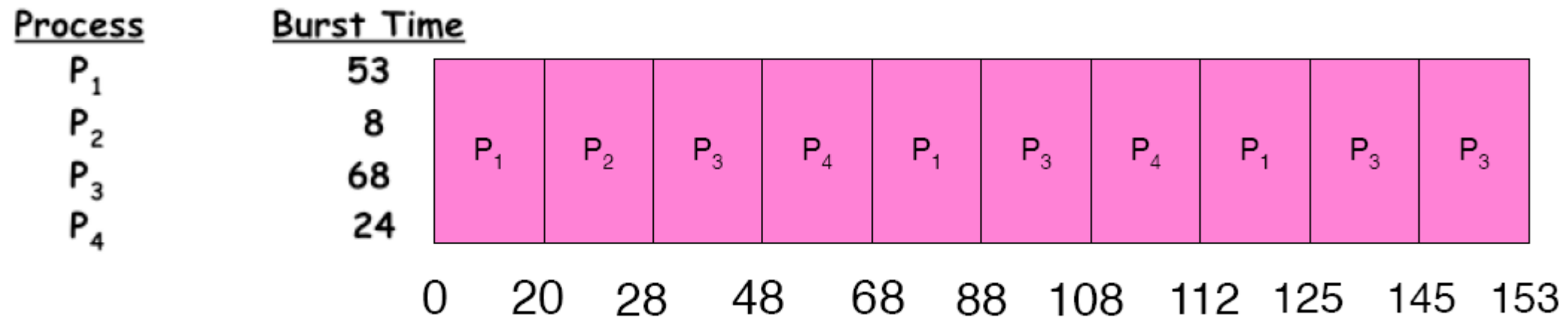
Turnaround Time Varies With The Time Quantum



| process | time |
|---------|------|
| P_1 | 6 |
| P_2 | 3 |
| P_3 | 1 |
| P_4 | 7 |

As can be seen from this graph, the average turnaround time of a set of processes does not necessarily improve as the time quantum size increases. In general, the average turnaround time can be improved if most processes finish their next CPU burst in a single time quantum.

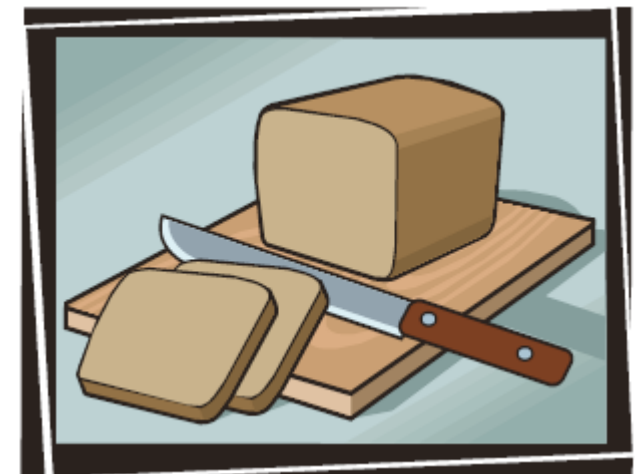
Example of RR with Time Quantum = 20



- **Waiting Time:** A process can finish before the time quantum expires, and release the CPU.
 - P1: $(68-20)+(112-88) = 72$
 - P2: $(20-0) = 20$
 - P3: $(28-0)+(88-48)+(125-108) = 85$
 - P4: $(48-0)+(108-68) = 88$
- **Completion Time:**
 - P1: 125
 - P2: 28
 - P3: 153
 - P4: 112
- **Average Waiting Time:** $(72+20+85+88)/4 = 66.25$
- **Average Completion Time:** $(125+28+153+112)/4 = 104.5$

RR Summary

- **Pros and Cons:**
 - Better for short jobs (+)
 - Fair (+)
 - Context-switching time adds up for long jobs (-)
 - The previous examples assumed no additional time was needed for context switching – in reality, this would add to wait and completion time without actually progressing a process towards completion.
 - Remember: the OS consumes resources, too!
- **If the chosen quantum is**
 - too large, response time suffers
 - infinite, performance is the same as FIFO
 - too small, throughput suffers and percentage overhead grows
- **Actual choices of timeslice:**
 - UNIX: initially 1 second:
 - Worked when only 1-2 users
 - If there were 3 compilations going on, it took 3 seconds to echo each keystroke!
 - In practice, need to balance short-job performance and long-job throughput:
 - Typical timeslice **10ms-100ms**
 - Typical context-switch overhead **0.1ms – 1ms (about 1%)**

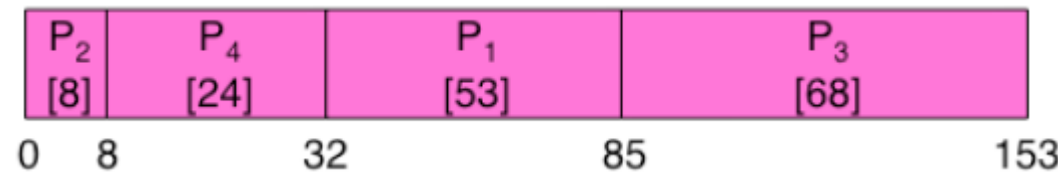


Comparing FCFS and RR

- Assuming zero-cost context switching time, is RR always better than FCFS?
- Assume 10 jobs, all start at the same time, and each require 100 seconds of CPU time
- RR scheduler quantum of 1 second
- Completion Times (CT)
 - Both FCFS and RR finish at the same time
 - But average response time is much worse under RR!
 - Bad when all jobs are same length
- Also: cache state must be shared between all jobs with RR but can be devoted to each job with FIFO
 - Total time for RR longer even for zero-cost context switch!

| Job # | FCFS CT | RR CT |
|-------|---------|-------|
| 1 | 100 | 991 |
| 2 | 200 | 992 |
| ... | ... | ... |
| 9 | 900 | 999 |
| 10 | 1000 | 1000 |

Comparing FCFS and RR



| | Quantum | P_1 | P_2 | P_3 | P_4 | Average |
|-----------------|------------|-------|-------|-------|-------|------------------|
| Wait Time | Best FCFS | 32 | 0 | 85 | 8 | $31\frac{1}{4}$ |
| | Q = 1 | 84 | 22 | 85 | 57 | 62 |
| | Q = 5 | 82 | 20 | 85 | 58 | $61\frac{1}{4}$ |
| | Q = 8 | 80 | 8 | 85 | 56 | $57\frac{1}{4}$ |
| | Q = 10 | 82 | 10 | 85 | 68 | $61\frac{1}{4}$ |
| | Q = 20 | 72 | 20 | 85 | 88 | $66\frac{1}{4}$ |
| | Worst FCFS | 68 | 145 | 0 | 121 | $83\frac{1}{2}$ |
| Completion Time | Best FCFS | 85 | 8 | 153 | 32 | $69\frac{1}{2}$ |
| | Q = 1 | 137 | 30 | 153 | 81 | $100\frac{1}{2}$ |
| | Q = 5 | 135 | 28 | 153 | 82 | $99\frac{1}{2}$ |
| | Q = 8 | 133 | 16 | 153 | 80 | $95\frac{1}{2}$ |
| | Q = 10 | 135 | 18 | 153 | 92 | $99\frac{1}{2}$ |
| | Q = 20 | 125 | 28 | 153 | 112 | $104\frac{1}{2}$ |
| | Worst FCFS | 121 | 153 | 68 | 145 | $121\frac{3}{4}$ |

Scheduling

- The performance we get is somewhat dependent on what “kind” of jobs we are running (short jobs, long jobs, etc.)
- If we could “see the future,” we could mirror best FCFS
- Shortest Job First (SJF) a.k.a. Shortest Time to Completion First (STCF):
 - Run whatever job has the least amount of computation to do
- Shortest Remaining Time First (SRTF) a.k.a. Shortest Remaining Time to Completion First (SRTCF):
 - Preemptive version of SJF: if a job arrives and has a shorter time to completion than the remaining time on the current job, immediately preempt CPU
- These can be applied either to a whole program or the current CPU burst of each program
 - Idea: get short jobs out of the system
 - Big effect on short jobs, only small effect on long ones
 - Result: better average response time

Scheduling

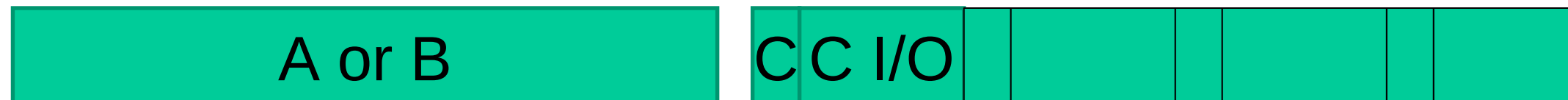
- But, this is hard to estimate
- We could get feedback from the program or the user, but they have incentive to lie!
- SJF/SRTF are the best you can do at minimizing average response time
 - Provably optimal (SJF among non-preemptive, SRTF among preemptive)
 - Since SRTF is always at least as good as SJF, focus on SRTF
- Comparison of SRTF with FCFS and RR
 - What if all jobs are the same length?
 - What if all jobs have varying length?



Scheduling

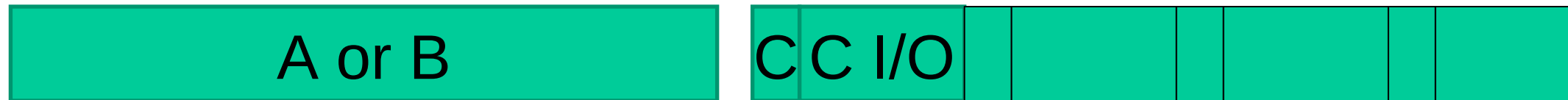
- But, this is hard to estimate
- We could get feedback from the program or the user, but they have incentive to lie!
- SJF/SRTF are the best you can do at minimizing average response time
 - Provably optimal (SJF among non-preemptive, SRTF among preemptive)
 - Since SRTF is always at least as good as SJF, focus on SRTF
- Comparison of SRTF with FCFS and RR
 - What if all jobs are the same length?
 - SRTF becomes the same as FCFS (i.e. FCFS is the best we can do)
 - What if all jobs have varying length?
 - SRTF (and RR): short jobs are not stuck behind long ones

Example: SRTF

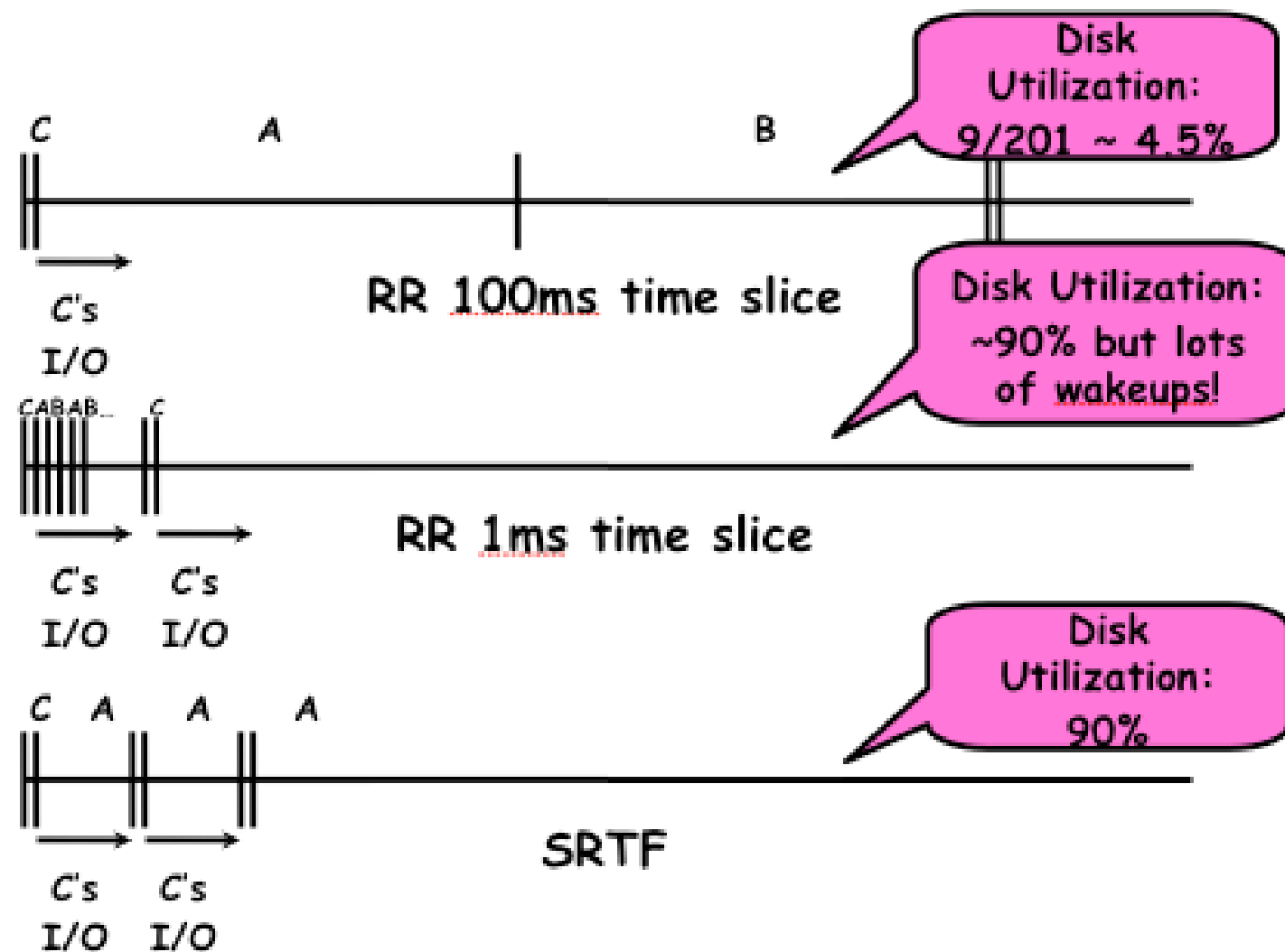


- **A,B: both CPU bound, run for a week**
- **C: I/O bound, loop 1ms CPU, 9ms disk I/O**
- **If only one at a time, C uses 90% of the disk, A or B could use 100% of the CPU**
- **With FIFO: Once A and B get in, the CPU is held for two weeks**
- **What about RR or SRTF?**
 - **Easier to see with a timeline**

Example: SRTF



- **A,B: both CPU bound, run for a week**
- **C: I/O bound, loop 1ms CPU, 9ms disk I/O**



Last Word on SRTF

- **Starvation**
 - SRTF can lead to starvation if many small jobs!
 - Large jobs never get to run
- **Somehow need to predict future**
 - How can we do this?
 - Some systems ask the user
 - When you submit a job, you have to say how long it will take
 - To stop cheating, system kills job if it takes too long
 - But even non-malicious users have trouble predicting runtime of their jobs
- **Bottom line, can't really tell how long job will take**
 - However, can use SRTF as a yardstick for measuring other policies, since it is optimal
- **SRTF Pros and Cons**
 - Optimal (average response time) (+)
 - Hard to predict future (-)
 - Unfair, even though we minimized average response time! (-)

Predicting the Future

- **Back to predicting the future... perhaps we can predict the next CPU burst length?**
- **Iff programs are generally repetitive, then they may be predictable**
- **Create an adaptive policy that changes based on past behavior**
 - CPU scheduling, virtual memory, file systems, etc.
 - If program was I/O bound in the past, likely in the future
- **Example: SRTF with estimated burst length**
 - Use an estimator function on previous bursts
 - Let $T(n-1)$, $T(n-2)$, $T(n-3)$, ..., be previous burst lengths. Estimate next burst $T(n) = f(T(n-1), T(n-2), T(n-3), \dots)$
 - Function f can be one of many different time series estimation schemes (Kalman filters, etc.)

Determining Length of Next CPU Burst

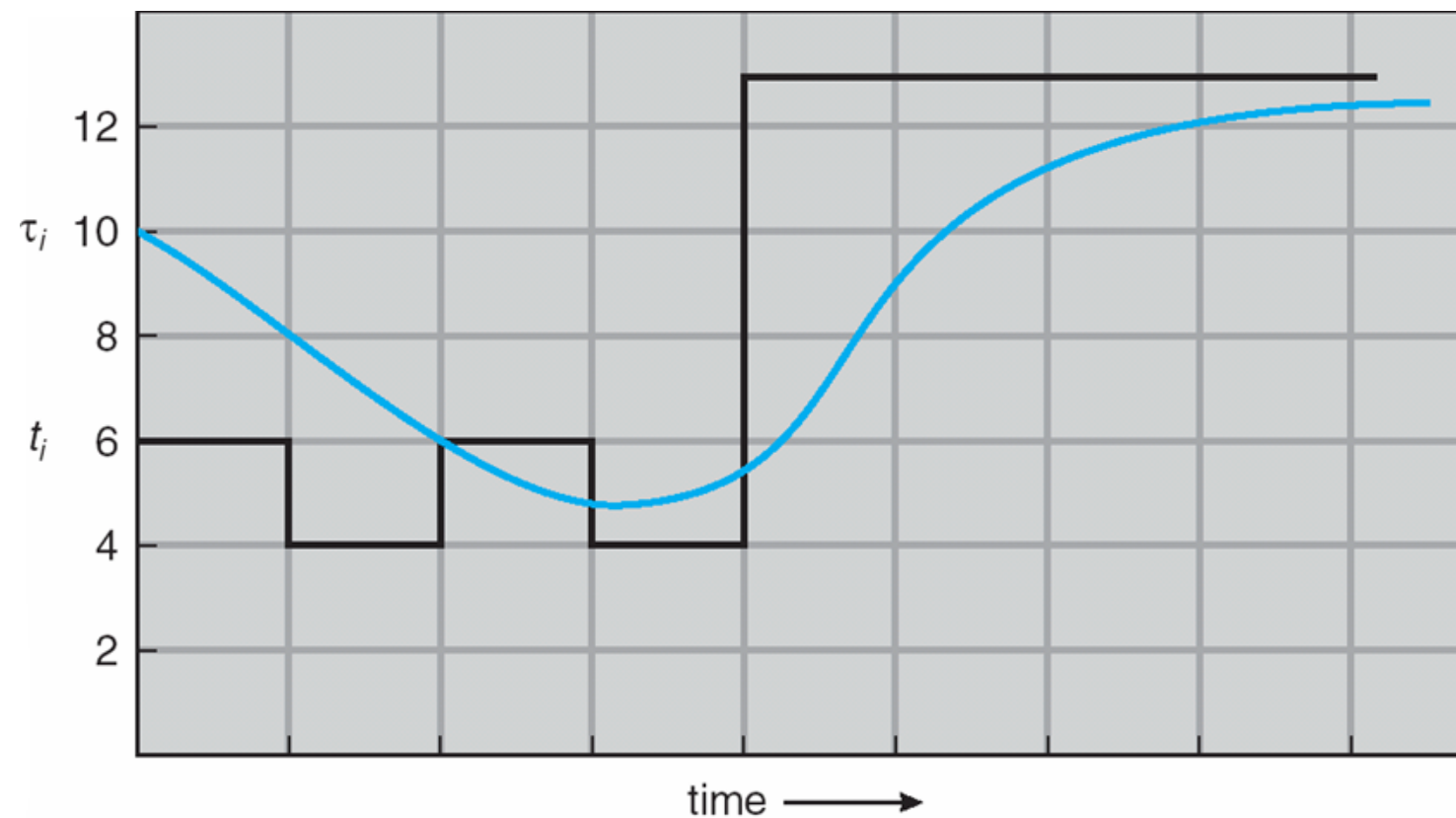
- Can only estimate the length
- Can be done by using the length of previous CPU bursts, using exponential averaging

$$\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n.$$

1. t_n = actual length of n^{th} CPU burst
2. τ_{n+1} = predicted value for the next CPU burst
3. $\alpha, 0 \leq \alpha \leq 1$
4. Define :

Predicting the Future

$$\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n.$$



| | | | | | | | | |
|----------------------|----|---|---|---|----|----|----|-----|
| CPU burst (t_i) | 6 | 4 | 6 | 4 | 13 | 13 | 13 | ... |
| "guess" (τ_i) | 10 | 8 | 6 | 5 | 9 | 11 | 12 | ... |

Examples of Exponential Averaging

∀ $\alpha = 0$

$$\tau_{n+1} = \tau_n$$

- Recent history does not count

∀ $\alpha = 1$

- $\tau_{n+1} = \alpha t_n$

- Only the actual last CPU burst counts

- If we expand the formula, we get:

$$\begin{aligned}\tau_{n+1} = & \alpha t_n + (1 - \alpha)\alpha t_{n-1} + \dots \\ & + (1 - \alpha)^j \alpha t_{n-j} + \dots \\ & + (1 - \alpha)^{n+1} \tau_0\end{aligned}$$

- Since both α and $(1 - \alpha)$ are less than or equal to 1, each successive term has less weight than its predecessor

Priority Scheduling

- A priority number (integer) is associated with each process.
- The SJF algorithm is a special case of the general priority scheduling algorithm.
- The CPU is allocated to the process with the highest priority (smallest integer \equiv highest priority)
 - Preemptive (if a higher priority process enters, it receives the CPU immediately)
 - Nonpreemptive (higher priority processes must wait until the current process finishes; then, the highest priority ready process is selected)
- SJF is a priority scheduling where priority is the predicted next CPU burst time
- Problem \equiv Starvation – low priority processes may never execute
- Solution \equiv Aging – as time progresses increase the priority of the process

Priority Inversion

- Consider a scenario in which there are three processes, one with high priority (H), one with medium priority (M), and one with low priority (L).
- Process L is running and successfully acquires a resource, such as a lock or semaphore.
- Process H begins; since we are using a preemptive priority scheduler, process L is preempted for process H.
- Process H tries to acquire L's resource, and blocks (because it is held by L).
- Process M begins running, and, since it has a higher priority than L, it is the highest priority ready process. It preempts L and runs, thus starving high priority process H.
- This is known as priority inversion.
- What can we do?



Priority Inversion

- **Process L should, in fact, be temporarily of “higher priority” than process M, on behalf of process H.**
- **Process H can donate its priority to process L, which, in this case, would make it higher priority than process M.**
- **This enables process L to preempt process M and run.**
- **When process L is finished, process H becomes unblocked.**
- **Process H, now being the highest priority ready process, runs, and process M must wait until it is finished.**
- **Note that if process M’s priority is actually higher than process H, priority donation won’t be sufficient to increase process L’s priority above process M. This is expected behavior (after all, process M would be “more important” in this case than process H).**

Multi-level Feedback Scheduling

- **Another method for exploiting past behavior**
 - **Multiple queues, each with different priority**
 - Higher priority queues often considered “foreground” tasks
 - **Each queue has its own scheduling algorithm**
 - E.g. foreground → RR, background → FCFS
 - Sometimes multiple RR priorities with quantum increasing exponentially (highest queue: 1ms, next: 2ms, next: 4ms, etc.)
 - **Adjust each job’s priority as follows (details vary)**
 - Job starts in highest priority queue
 - If entire CPU time quantum expires, drop one level
 - If CPU is yielded during the quantum, push up one level (or to top)

Scheduling Details

- **Result approximates SRTF**
 - CPU bound jobs drop rapidly to lower queues
 - Short-running I/O bound jobs stay near the top
- **Scheduling must be done between the queues**
 - Fixed priority scheduling: serve all from the highest priority, then the next priority, etc.
 - Time slice: each queue gets a certain amount of CPU time (e.g., 70% to the highest, 20% next, 10% lowest)
- **Countermeasure: user action that can foil intent of the OS designer**
 - For multilevel feedback, put in a bunch of meaningless I/O to keep job's priority high
 - But if everyone does this, it won't work!
 - Consider an Othello program, playing against a competitor. Key was to compute at a higher priority than the competitors.
 - Put in printf's, run much faster!

Scheduling Details

- It is apparent that scheduling is facilitated by having a “good mix” of I/O bound and CPU bound programs, so that there are long and short CPU bursts to prioritize around.
- There is typically a long-term and a short-term scheduler in the OS.
- We have been discussing the design of the short-term scheduler.
- The long-term scheduler decides what processes should be put into the ready queue in the first place for the short-term scheduler, so that the short-term scheduler can make fast decisions on a good mix of a subset of ready processes.
- The rest are held in memory or disk
 - Why else is this helpful?



Scheduling Details

- It is apparent that scheduling is facilitated by having a “good mix” of I/O bound and CPU bound programs, so that there are long and short CPU bursts to prioritize around.
- There is typically a long-term and a short-term scheduler in the OS.
- We have been discussing the design of the short-term scheduler.
- The long-term scheduler decides what processes should be put into the ready queue in the first place for the short-term scheduler, so that the short-term scheduler can make fast decisions on a good mix of a subset of ready processes.
- The rest are held in memory or disk
 - This also provides more free memory for the subset of ready processes given to the short-term scheduler.

Fairness

- **What about fairness?**
 - Strict fixed-policy scheduling between queues is unfair (run highest, then next, etc.)
 - Long running jobs may never get the CPU
 - In Multics, admins shut down the machine and found a 10-year-old job
 - Must give long-running jobs a fraction of the CPU even when there are shorter jobs to run
 - Tradeoff: fairness gained by hurting average response time!
- **How to implement fairness?**
 - Could give each queue some fraction of the CPU
 - i.e., for one long-running job and 100 short-running ones?
 - Like express lanes in a supermarket – sometimes express lanes get so long, one gets better service by going into one of the regular lines
 - Could increase priority of jobs that don't get service (as seen in the multilevel feedback example)
 - This was done in UNIX
 - Ad hoc – with what rate should priorities be increased?
 - As system gets overloaded, no job gets CPU time, so everyone increases in priority
 - Interactive processes suffer

Lottery Scheduling

- **Yet another alternative: Lottery Scheduling**
 - Give each job some number of lottery tickets
 - On each time slice, randomly pick a winning ticket
 - On average, CPU time is proportional to number of tickets given to each job over time
- **How to assign tickets?**
 - To approximate SRTF, short-running jobs get more, long running jobs get fewer
 - To avoid starvation, every job gets at least one ticket (everyone makes progress)
- **Advantage over strict priority scheduling: behaves gracefully as load changes**
 - Adding or deleting a job affects all jobs proportionally, independent of how many tickets each job possesses

Example: Lottery Scheduling

- Assume short jobs get 10 tickets, long jobs get 1 ticket
- What percentage of time does each long job get? Each short job?



- What if there are too many short jobs to give reasonable response time
 - In UNIX, if load average is 100%, it's hard to make progress
 - Log a user out or swap a process out of the ready queue (long term scheduler)

Example: Lottery Scheduling

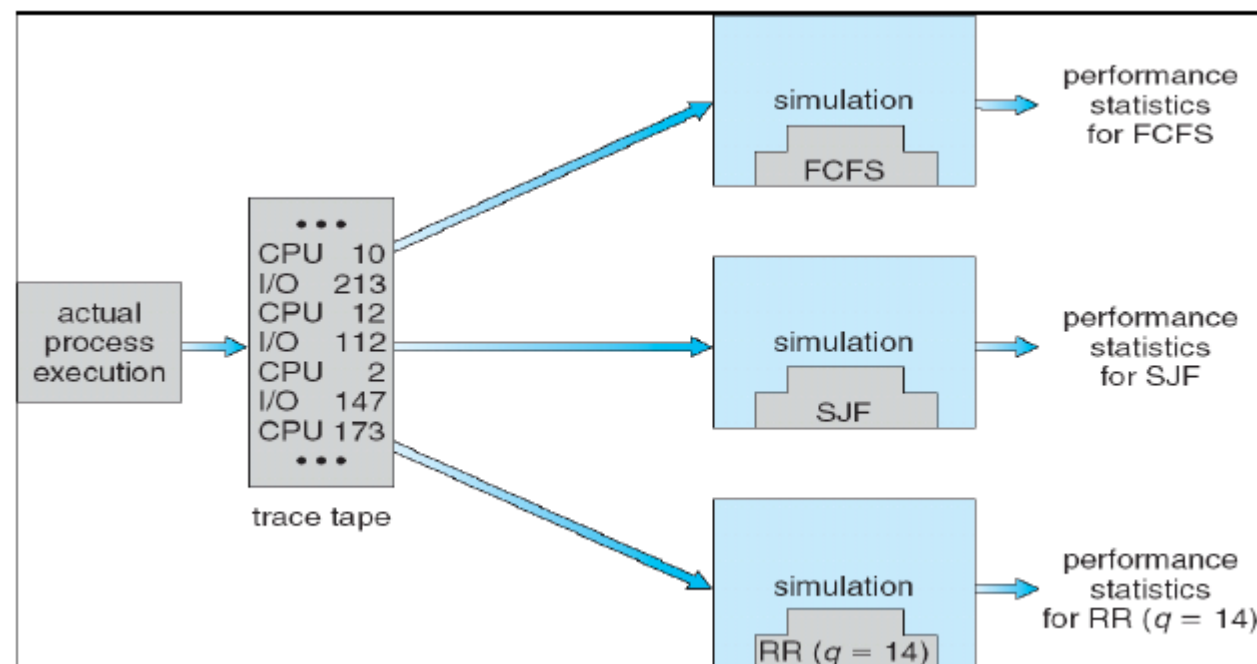
- Assume short jobs get 10 tickets, long jobs get 1 ticket

| # short jobs / # long jobs | % of CPU each short job gets | % of CPU each long job gets |
|-------------------------------|---------------------------------|--------------------------------|
| 1/1 | 91% | 9% |
| 0/2 | N/A | 50% |
| 2/0 | 50% | N/A |
| 10/1 | 9.9% | 0.99% |
| 1/10 | 50% | 5% |

- What if there are too many short jobs to give reasonable response time
 - In UNIX, if load average is 100%, it's hard to make progress
 - Log a user out or swap a process out of the ready queue (long term scheduler)

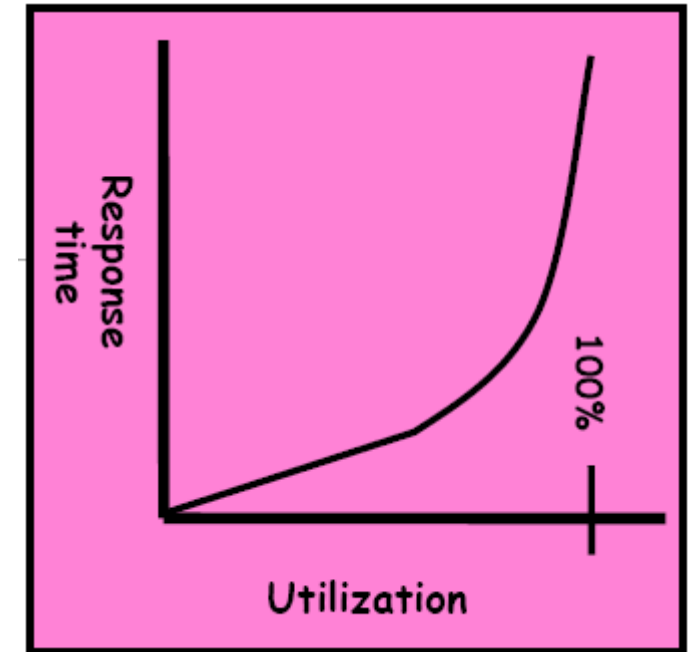
Scheduling Algorithm Evaluation

- **Deterministic Modeling**
 - Takes a predetermined workload and compute the performance of each algorithm for that workload
- **Queuing Models**
 - Mathematical Approach for handling stochastic workloads
- **Implementation / Simulation**
 - Build system which allows actual algorithms to be run against actual data. Most flexible / general.



Conclusion

- **Scheduling:** selecting a waiting process from the ready queue and allocating the CPU to it
- **When do the details of the scheduling policy and fairness really matter?**
 - When there aren't enough resources to go around
- **When should you simply buy a faster computer?**
 - Or network link, expanded highway, etc.
 - One approach: buy it when it will pay for itself in improved response time
 - Assuming you're paying for worse response in reduced productivity, customer angst, etc.
 - Might think that you should buy a faster X when X is utilized 100%, but usually, response time goes to infinite as utilization goes to 100%
 - Most scheduling algorithms work fine in the “linear” portion of the load curve, and fail otherwise
 - Argues for buying a faster X when utilization is at the “knee” of the curve



- **FCFS scheduling, FIFO Run Until Done:**
 - Simple, but short jobs get stuck behind long ones
- **RR scheduling:**
 - Give each thread a small amount of CPU time when it executes, and cycle between all ready threads
 - Better for short jobs, but poor when jobs are the same length
- **SJF/SRTF:**
 - Run whatever job has the least amount of computation to do / least amount of remaining computation to do
 - Optimal (average response time), but unfair; hard to predict the future
- **Multi-Level Feedback Scheduling:**
 - Multiple queues of different priorities
 - Automatic promotion/demotion of process priority to approximate SJF/SRTF
- **Lottery Scheduling:**
 - Give each thread a number of tickets (short tasks get more)
 - Every thread gets tickets to ensure forward progress / fairness
- **Priority Scheduling:**
 - Preemptive or Nonpreemptive
 - Priority Inversion