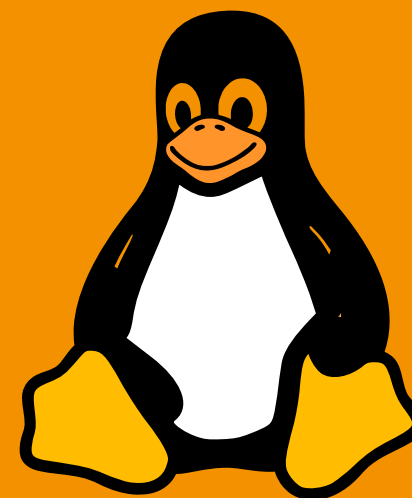




Prof. Florian Künzner

Technical University of Applied Sciences Rosenheim, Computer Science

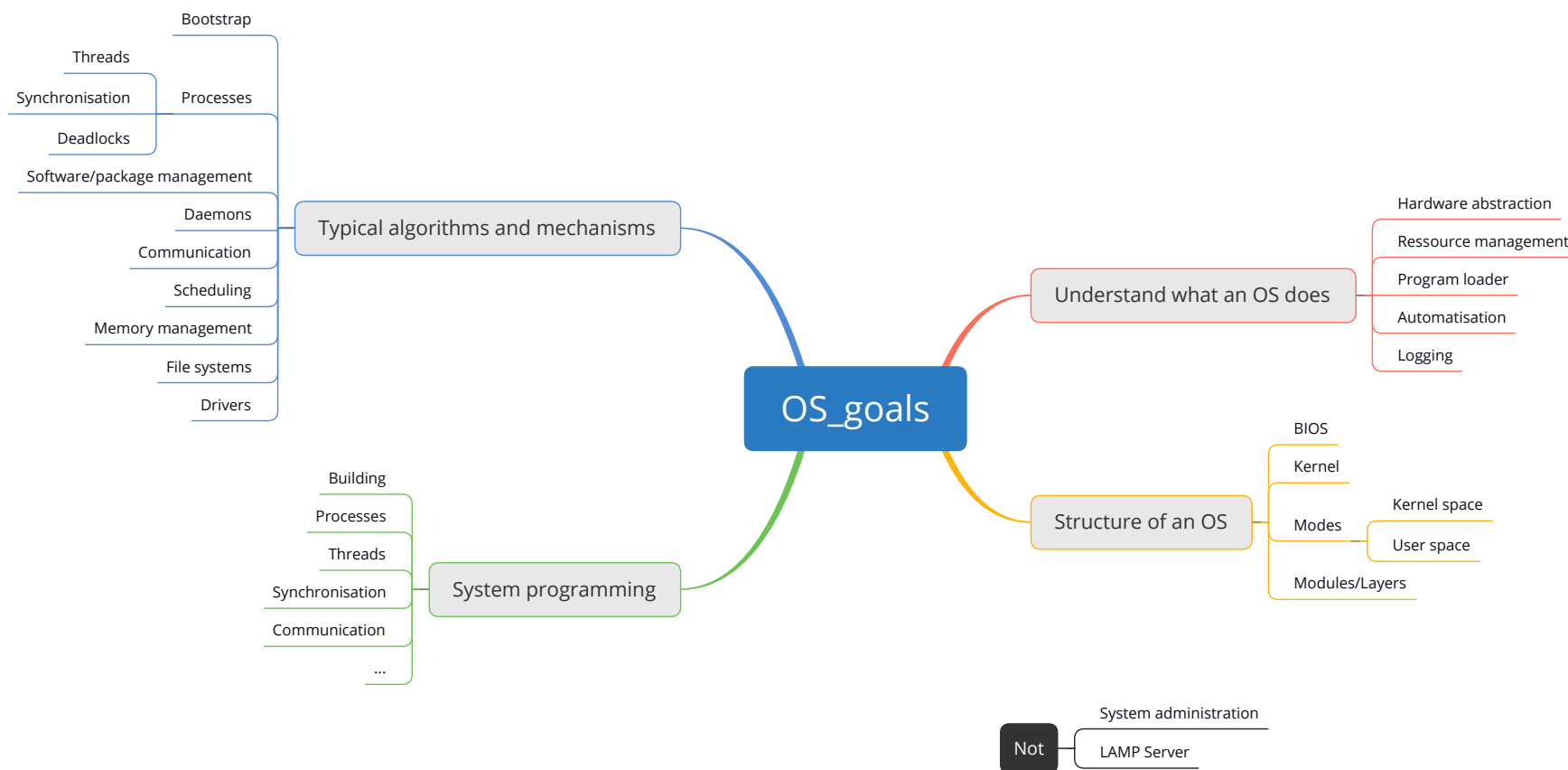
OS 11 – Scheduling



source: [iconspng.com](https://www.iconspng.com)

The lecture is based on the work and the documents of Prof. Dr. Ludwig Frank

Goal



Goal

OS::Scheduling

- Scheduling theory
- Scheduling strategies
- Scheduling on Linux

Scheduling

Scheduling is a technique to **distribute computing resources** like processor time, bandwidth, memory, or device I/O **to the processes** on a computer system.

Scheduler

The **scheduler** is such a **program** that **distributes** the **resources** to the processes.

Scheduling time frame

Long term scheduling

Determines which processes (jobs) are admitted to the system for processing.

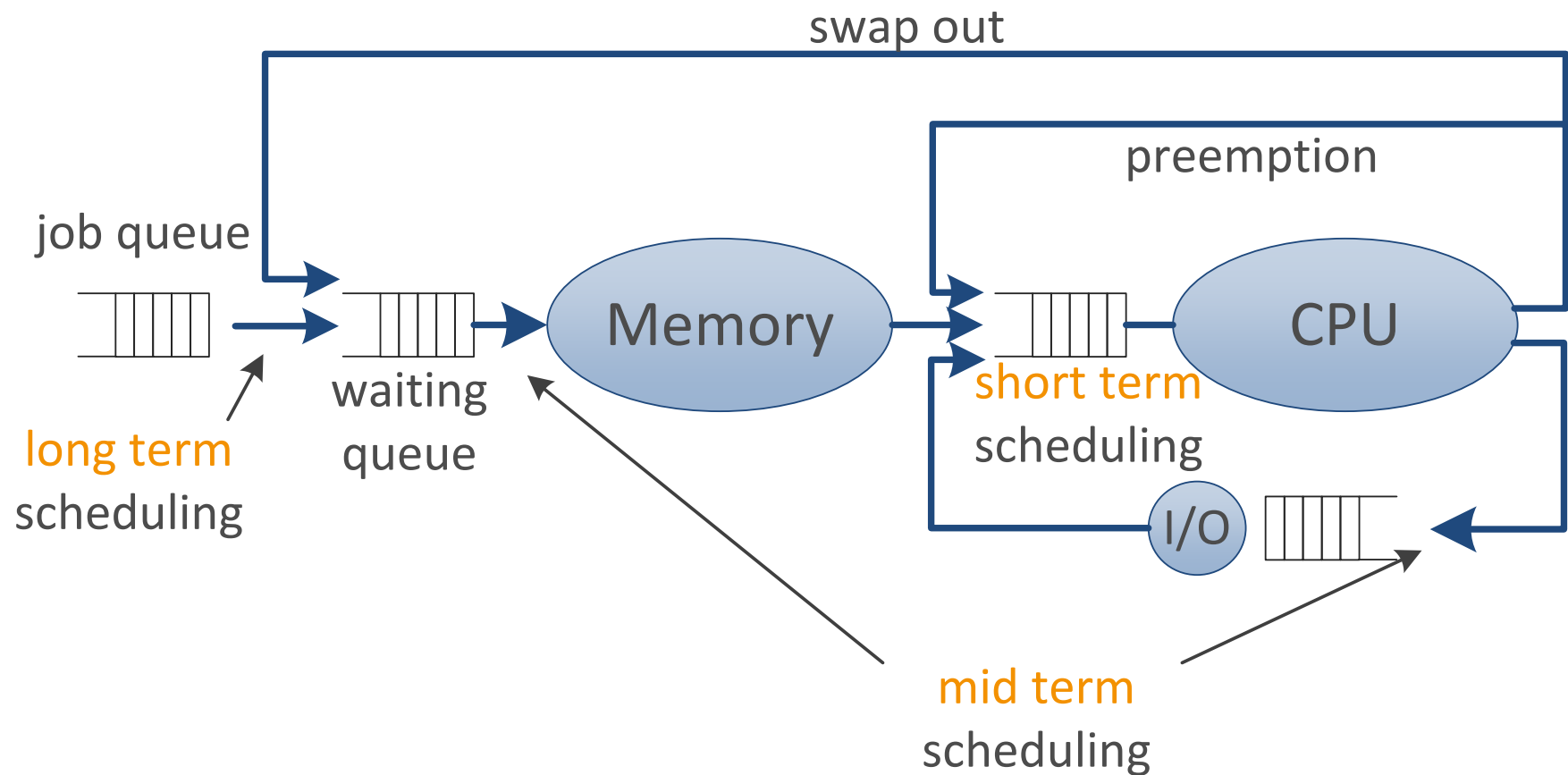
Mid term scheduling

Allocation of bandwidth, memory, or device I/O to a process/thread.

Short term scheduling

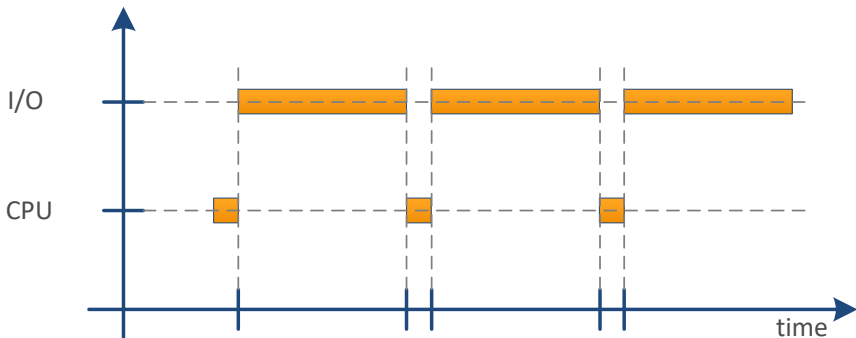
Allocation of a process/thread to a cpu core.

Scheduling time frame overview

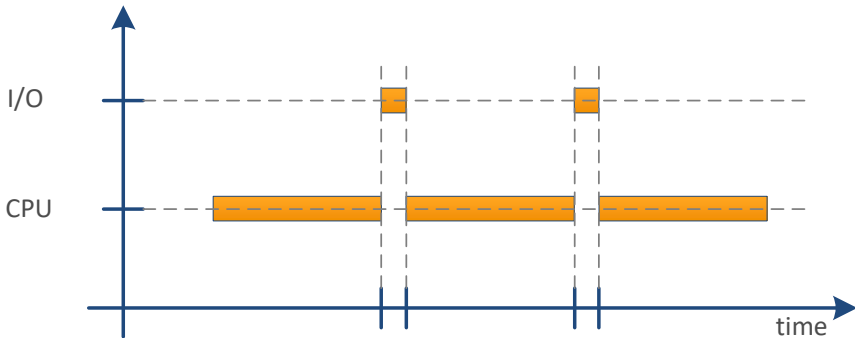


I/O-bound vs. CPU-bound

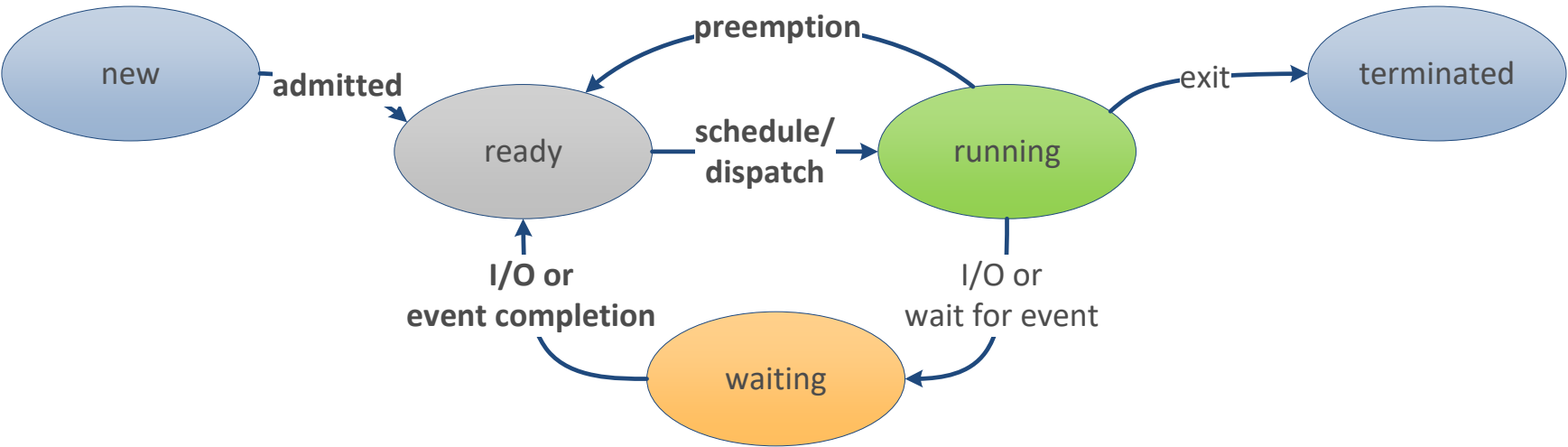
I/O-bound



CPU-bound



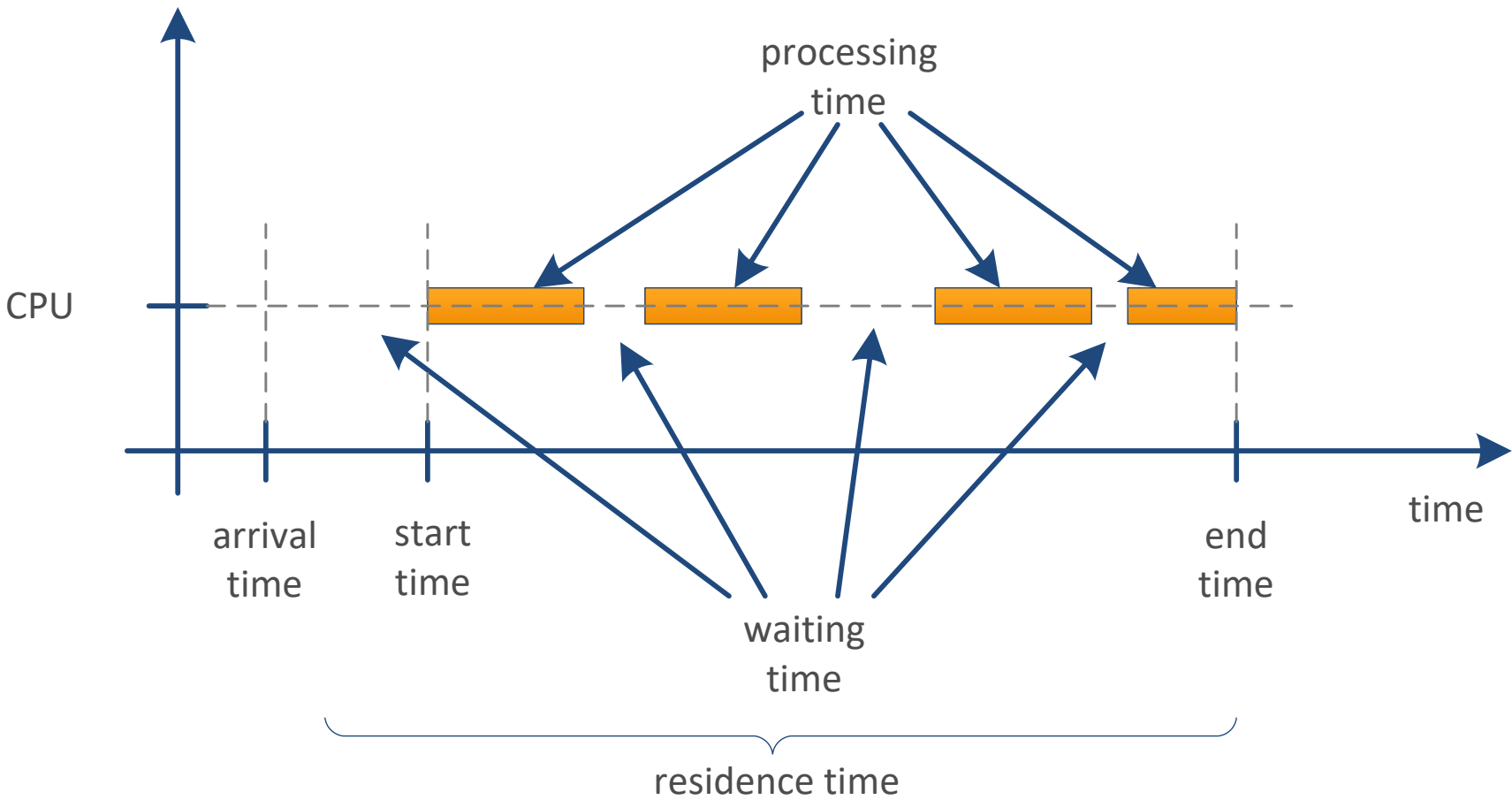
Process states



Scheduling terms

Term	Description
CPU usage	CPU usage up to 100% if possible.
throughput	The number of completed processes in a time frame. Should be as high as possible.
arrival time	The point in time at which a process arrives for execution in the system.
processing time	The time a process takes to run on the CPU.
waiting time	The time a process has to wait until it can run.
residence time	The total time a process takes to finish (= processing time + waiting time).

Time aspects



System categories

System type

Description

Job processing system

Typically users submit jobs (programs with its parameters) to a system. The job scheduler decides when the job starts.

Interactive system

With an interactive system, the users work directly: a PC (terminal, desktop), smartphone, ...

Real-time system

A real-time system usually observes and controls a physical process in the real world. It must guarantee that it reacts fast enough.

Scheduling Goals

All systems

- fairness - Every process can run on the CPU.
- policy enforcement - The system's policy is enforced.
- balance - All parts of the system are busy.

Job processing system

- throughput - Maximize the number of jobs in a time frame.
- residence time - Minimize the residence time for each job.
- CPU usage - The CPU is constantly used as long as there are jobs in the queue.

Interactive system

- response time - Respond quickly to requests.
- proportionality - Consider the requirements for all users.

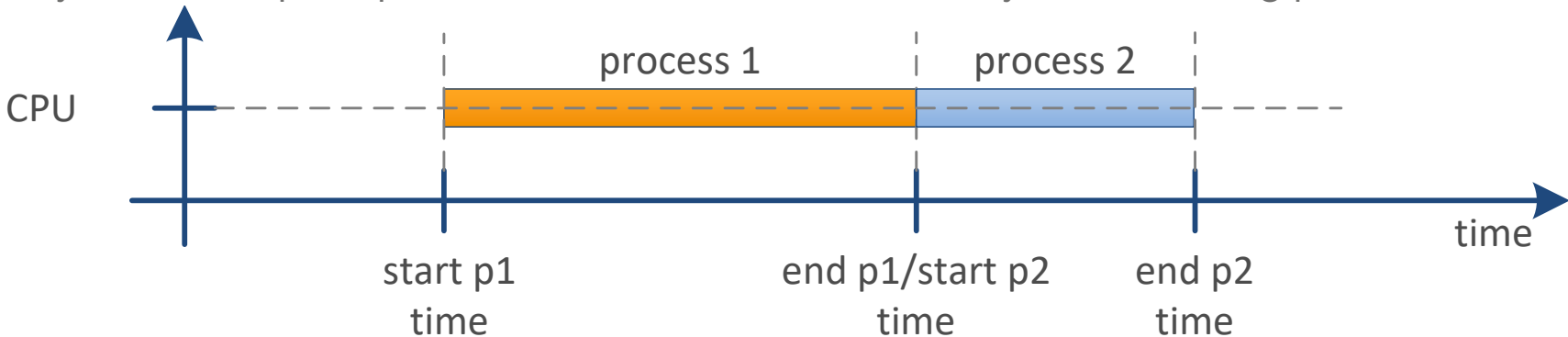
Real-time system

- meet deadline - Meet always the deadline of all processes.
- predictability - Always guarantee the same periodic execution (small jitter).

Preemption

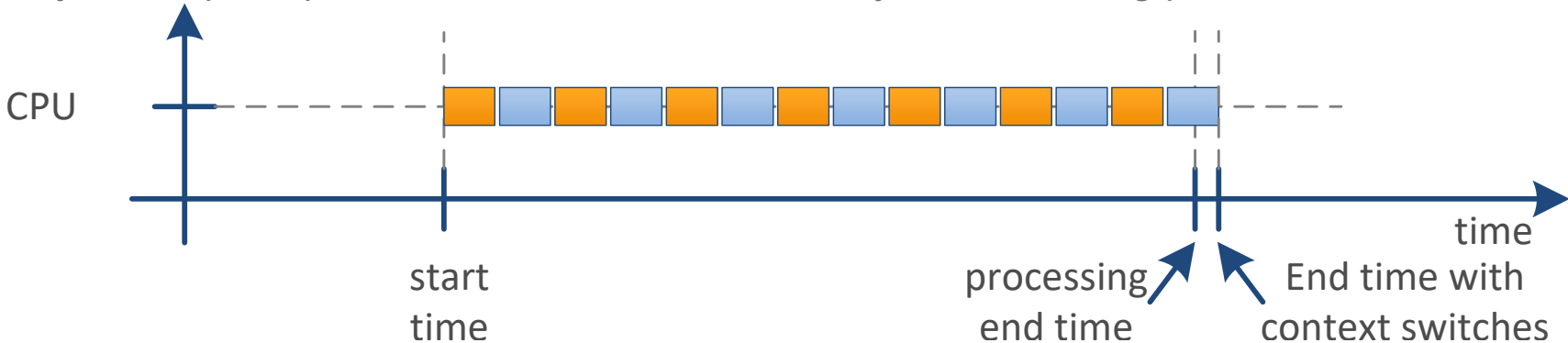
Non-preemptive system

A system is non-preemptive if the CPU cannot be taken away from a running process.



Preemptive system

A system is preemptive if the CPU can be taken away from a running process.



Context switch

A context switch **changes the active process or thread** on the CPU. This may be expensive (takes some time).

Procedure:

- Save the register content of the currently running process into its PCB.
- Select a new process to run.
- Load/restore memory information into the CPU (MMU) from its PCB.
- Load/restore the register content of the new process from its PCB.

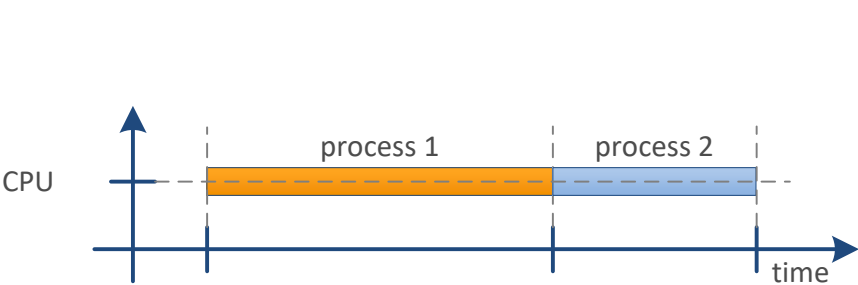
Time for context switches:

- Time to switch processes: $\approx 3600ns$ per context switch (Intel E5440 CPU)
- Time to switch threads: $\approx 1300ns$ per context switch (Intel E5440 CPU)

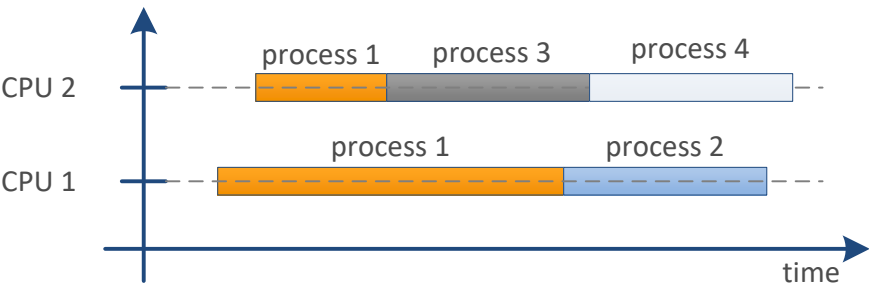
Source for measured times: <https://blog.tsunanet.net/2010/11/how-long-does-it-take-to-make-context.html>

Scheduling: single vs. multi core CPU

Single core CPU



Multi core CPU



Scheduling strategies

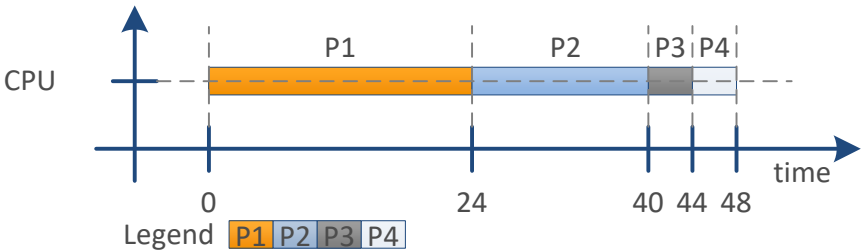
How would you schedule processes on your machine?

FCFS - first come first served

The arrival order (time) in the waiting queue is the scheduling order.

Or.	ar. time	Process	proc. time	res. time
1	0	P1	24	24
2	0	P2	16	40
3	0	P3	4	44
4	0	P4	4	48

Mean res. size = $(24 + 40 + 44 + 48)/4 = 39$



Properties

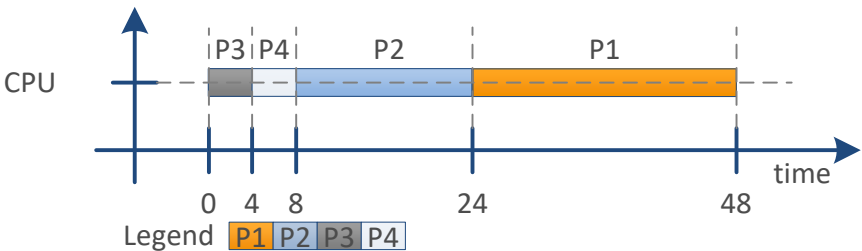
- A faire order (arrival time).
- Small jobs may wait long.
- Not good for interactive system: does not guarantee a good response time.

SJF - shortest jobs first

The job with the smallest processing time is scheduled first.

Or.	ar. time	Process	proc. time	res. time
1	0	P1	24	48
2	0	P2	16	24
3	0	P3	4	4
4	0	P4	4	8

Mean res. size = $(48 + 24 + 4 + 8)/4 = 21$



Properties

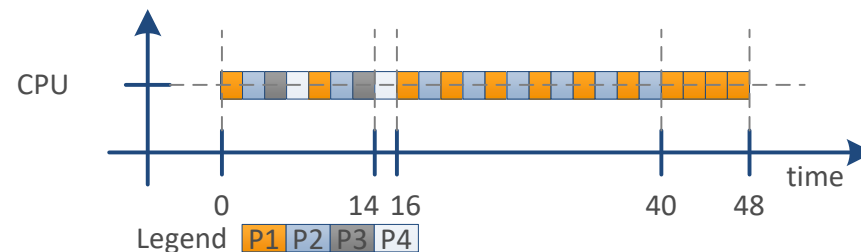
- The processing time is often not available (prediction also hard).
- Processing time is predicted: by user, automatically?
- SJF optimises the throughput.
- Jobs with a long processing time may not get scheduled (starvation).

RR - round robin

Every process in the queue get the required resource for a limited amount of time. Then the preempted process is placed to the last position in the waiting queue.

Or.	ar. time	Process	proc. time	res. time
1	0	P1	24	48
2	0	P2	16	40
3	0	P3	4	14
4	0	P4	4	16

Mean res. size = $(48 + 40 + 14 + 16)/4 = 29.5$



Properties

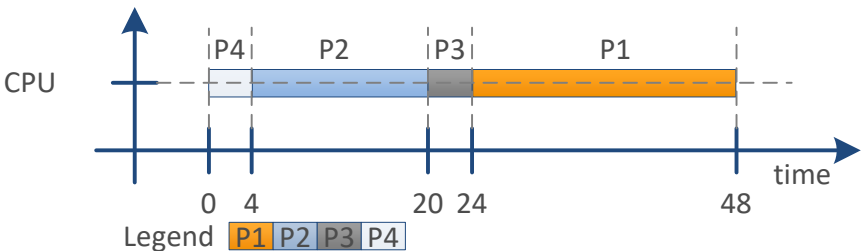
- Good for interactive systems.
- Too many process switches causes context switches—and that are expensive.
- The average waiting time is often longer than with other scheduling strategies.

EDF - earliest deadline first

The process whose deadline ends first is processed first.

Or.	ar. time	Process	proc. time	deadl.	res. time
1	0	P1	24	60	48
2	0	P2	16	20	20
3	0	P3	4	28	24
4	0	P4	4	4	4

Mean res. size = $(48 + 20 + 24 + 4) / 4 = 24$



Properties

- Used in real-time systems.
- Is not always optimal on multi-core CPUs.

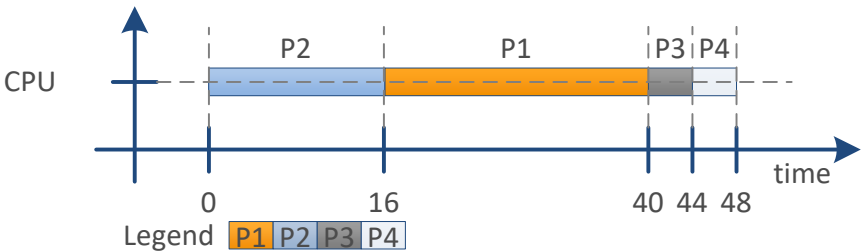
Priority based

External priority

The process with the highest priority is scheduled first. The user can define the priority on startup or change it during execution.

Or.	ar. time	Process	proc. time	prio	res. time
1	0	P1	24	2	40
2	0	P2	16	1	16
3	0	P3	4	3	44
4	0	P4	4	4	48

Mean res. size = $(40 + 16 + 44 + 48)/4 = 37$



Properties

- The important processes are scheduled first.
- Long-running processes with a high priority can cause low priority processes to be kept away from the CPU for a long time.

Priority based

Internal priority

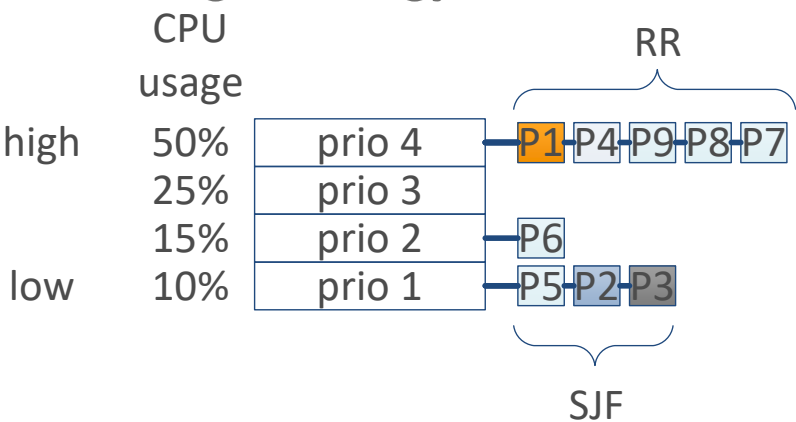
After a process runs for a while its priority is automatically lowered.
After some long waiting time, the priority can automatically be increased again.

Properties

- Solves the issue of the external priority.
- Improves the response time.

Multilevel queue scheduling

There exists multiple queues, whereas each can have its own scheduling strategy.



Properties

- Combination of different strategies.
- Each queue can have its own CPU usage.
- Improves response time while taking priorities into account.

Linux priorities

Supports 140 different priority classes (0...139)
The lower the value, the higher the priority.

Normal processes: 100...139 [0...39]

Parameter	Definition	Description
<i>NI</i>	$NI = -20 \dots +19$	Nice value ranges from $-20 \dots +19$. Users can only increase the nice value, but not lower it.
<i>PR</i>	$PR = 20 + NI$	The priority <i>PR</i> ranges from 0...39. Default user processes usually submitted with $NI = 0 \Rightarrow PR = 20$.

Real-time processes: 0...99 [-100... -1]

Parameter	Definition	Description
<i>RT</i>	$RT = 1 \dots 99$	The real-time priority <i>RT</i> ranges from 1...99
<i>PR</i>	$PR = -1 - RT$	The priority <i>PR</i> ranges from $-100 \dots -1$.

Linux commands

Command	Description
<code>top</code>	Shows processes in live view.
<code>htop</code>	Shows processes in live view.
<code>ps ax -o pid,rtprio,pri,ni,cmd</code>	Shows processes with pid, rtprio, pri, ni, and its cmd.
<code>nice -n 15 make -j</code>	Starts a parallel build with $NI = 15$
<code>nice -n -5 make -j</code>	Starts a parallel build with $NI = -5$
<code>renice -n 15 -p 1000</code>	Change $NI = 15$ of process with PID 1000.
<code>renice -n -5 -p 1000</code>	Change $NI = -5$ of process with PID 1000.

Summary and outlook

Summary

- Scheduling theory
- Scheduling strategies
- Scheduling on Linux

Outlook

- Memory management