# *Real-Time Compute Virtualization*

Luca Abeni

`luca.abeni@santannapisa.it`

January 13, 2022

# Real-Time Applications

- Real-Time Application: The time when a result is produced matters
    - A correct result produced *too late* is equivalent to a wrong result (or to no result)
- What does "*too late*" mean, here?
    - Applications characterised by temporal constraints that have to be respected!
- Examples:
    - Control applications, autonomous driving, ...
    - But also infotainment, gaming, telecommunications, ...!!!

# Temporal Constraints

- Temporal constraints are modelled through *deadlines*

    - Finish some activity before a time (deadline)
    - Generate some data before a deadline
    - Terminate some process/thread before a deadline

    - ...

- What happens if a constraint is not respected?

    - Simple: the application fails!
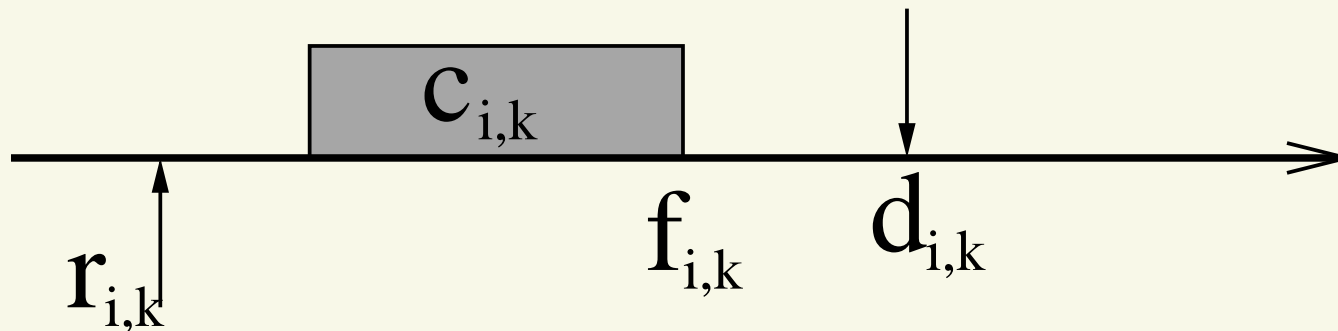
# Processes, Threads, and Tasks

- Algorithm → logical procedure used to solve a problem
- Program → formal description of an algorithm, using a *programming language*
- Process → instance of a program (program in execution)

  - Program: static entity
  - Process: dynamic entity

- The term *task* is used to indicate a schedulable entity (either a process or a thread)

  - Thread → flow of execution
  - Process → flow of execution + private resources (address space, file table, etc...)

# Real-Time Tasks

- A task can be seen as a sequence of actions ...
- ... and a deadline must be associated to each one of them!
  - Some kind of formal model is needed to identify these "actions" and associate deadlines to them

# Mathematical Model of a Task - 1

- Real-Time task $\tau_i$: stream of jobs (or instances) $J_{i,k}$
- Each job $J_{i,k} = (r_{i,k}, c_{i,k}, d_{i,k})$:
  - Arrives at time $r_{i,k}$ (activation time)
  - Executes for a time $c_{i,k}$
  - Finishes at time $f_{i,k}$
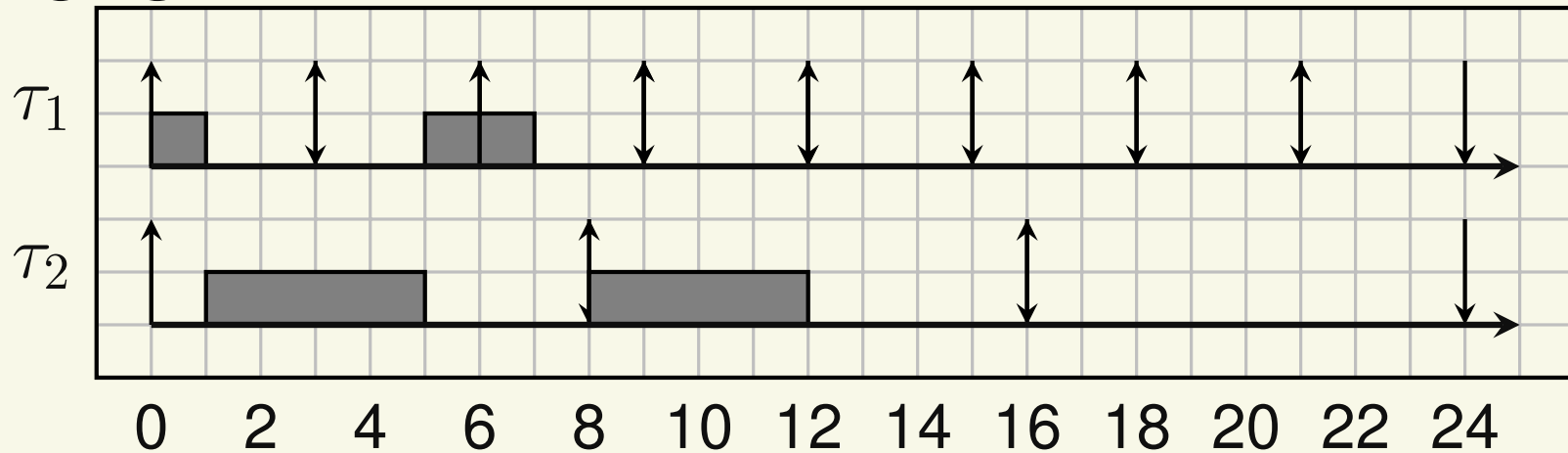  - Should finish within an absolute deadline $d_{i,k}$
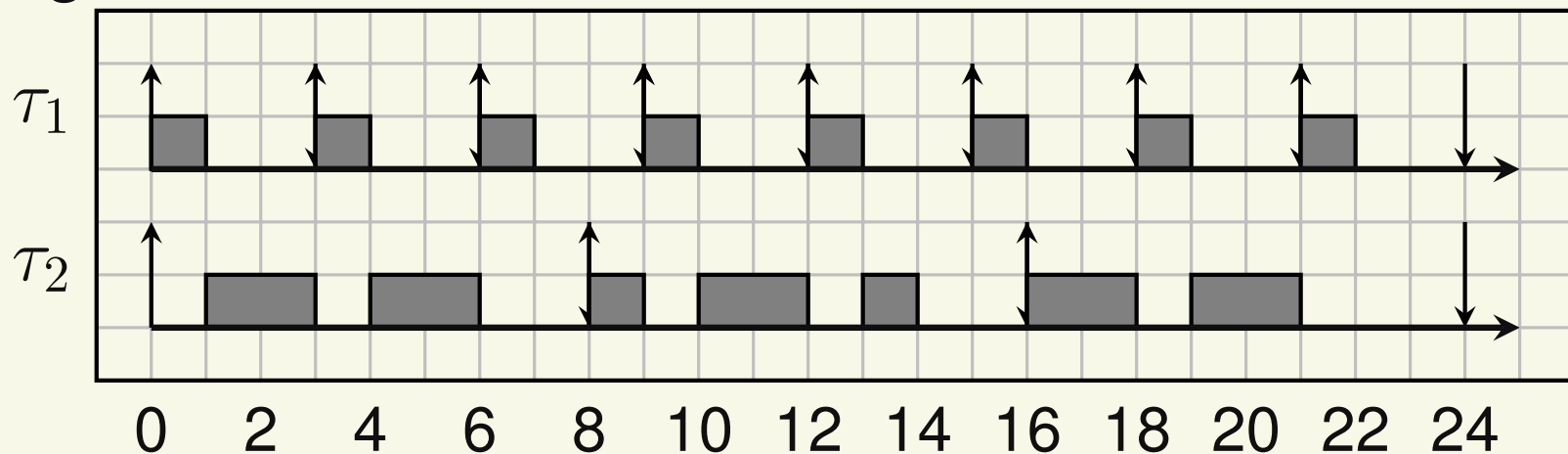
# Mathematical Model of a Task - 2

- Job: abstraction used to associate deadlines (temporal constraints) to activities
  - $r_{i,k}$: time when job $J_{i,k}$ is *activated* (by an external event, a timer, an explicit activation, etc...)
  - $c_{i,k}$: computation time needed by job $J_{i,k}$ to complete
  - $d_{i,k}$: absolute time instant by which job $J_{i,k}$ must complete
    - job $J_{i,k}$ respects its deadline if $f_{i,k} \leq d_{i,k}$
- Response time of job $J_{i,k}$: $\rho_{i,k} = f_{i,k} - r_{i,k}$

# RT Scheduling: Why?

- The task set $\mathcal{T} = \{(1,3), (4,8)\}$ is not schedulable by FCFS



- $\mathcal{T} = \{(1,3), (4,8)\}$ is schedulable with other algorithms
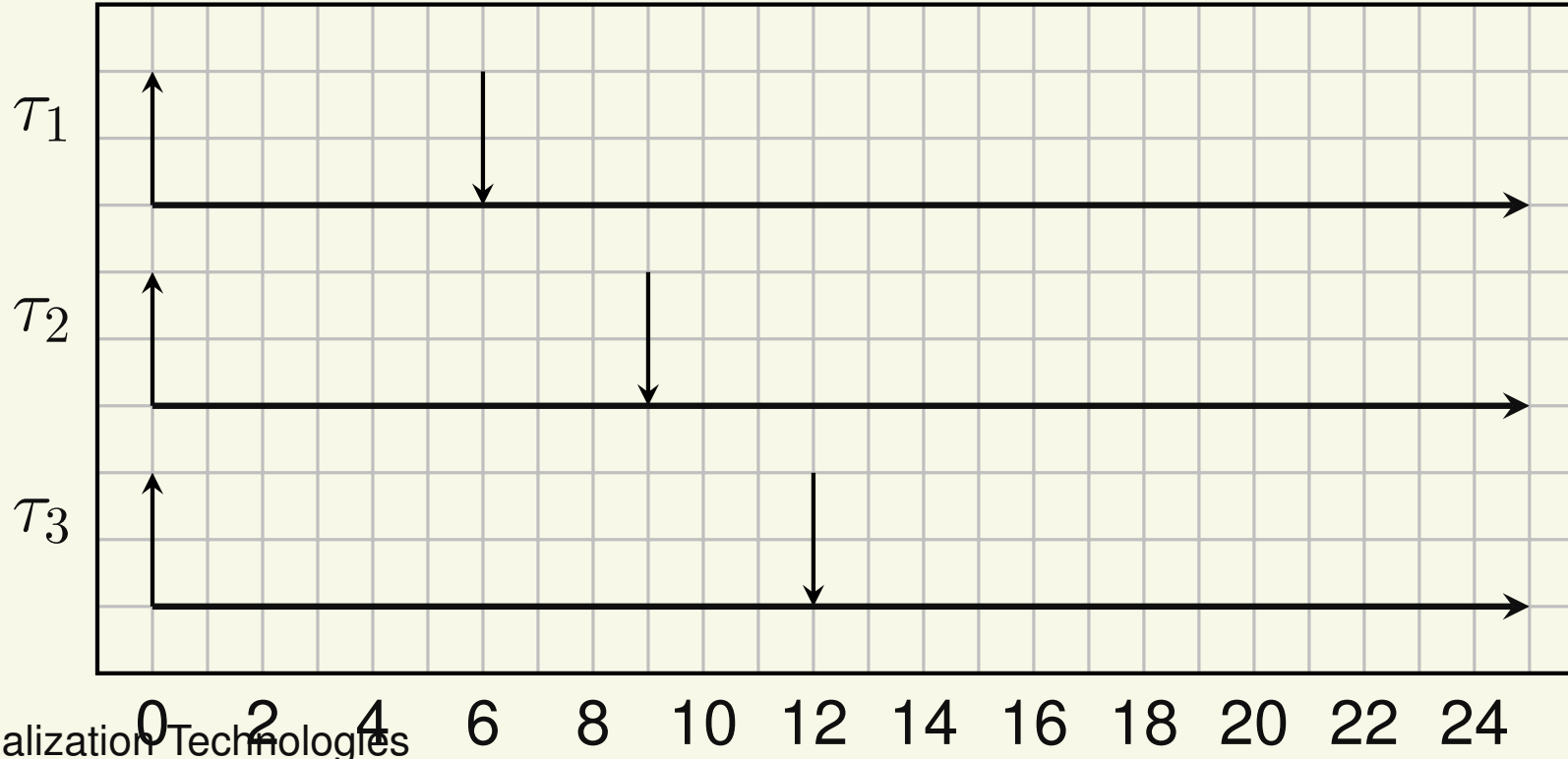
# The Scheduling Problem

- A real-time task $\tau_i$ is properly served if all jobs respect their deadline...
- ...Appropriate scheduling is important!
  - The CPU scheduler must somehow know the temporal constaints of the tasks...
  - ...To schedule them so that such temporal constraints are respected
- How to schedule real-time tasks? (scheduling algorithm)
- Is it possible to respect all the deadlines?
- Do commonly used OSs provide appropriate scheduling algorithms?

# Fixed Priority Scheduling

- Very simple *preemptive* scheduling algorithm
  - Every task $\tau_i$ is assigned a fixed priority $p_i$
  - The active task with the highest priority is scheduled
- Priorities are integer numbers: the higher the number, the higher the priority
  - In the research literature, sometimes authors use the opposite convention: the lowest the number, the highest the priority
- In the following we show some examples, considering periodic tasks, constant execution times, and deadlines equal to the period
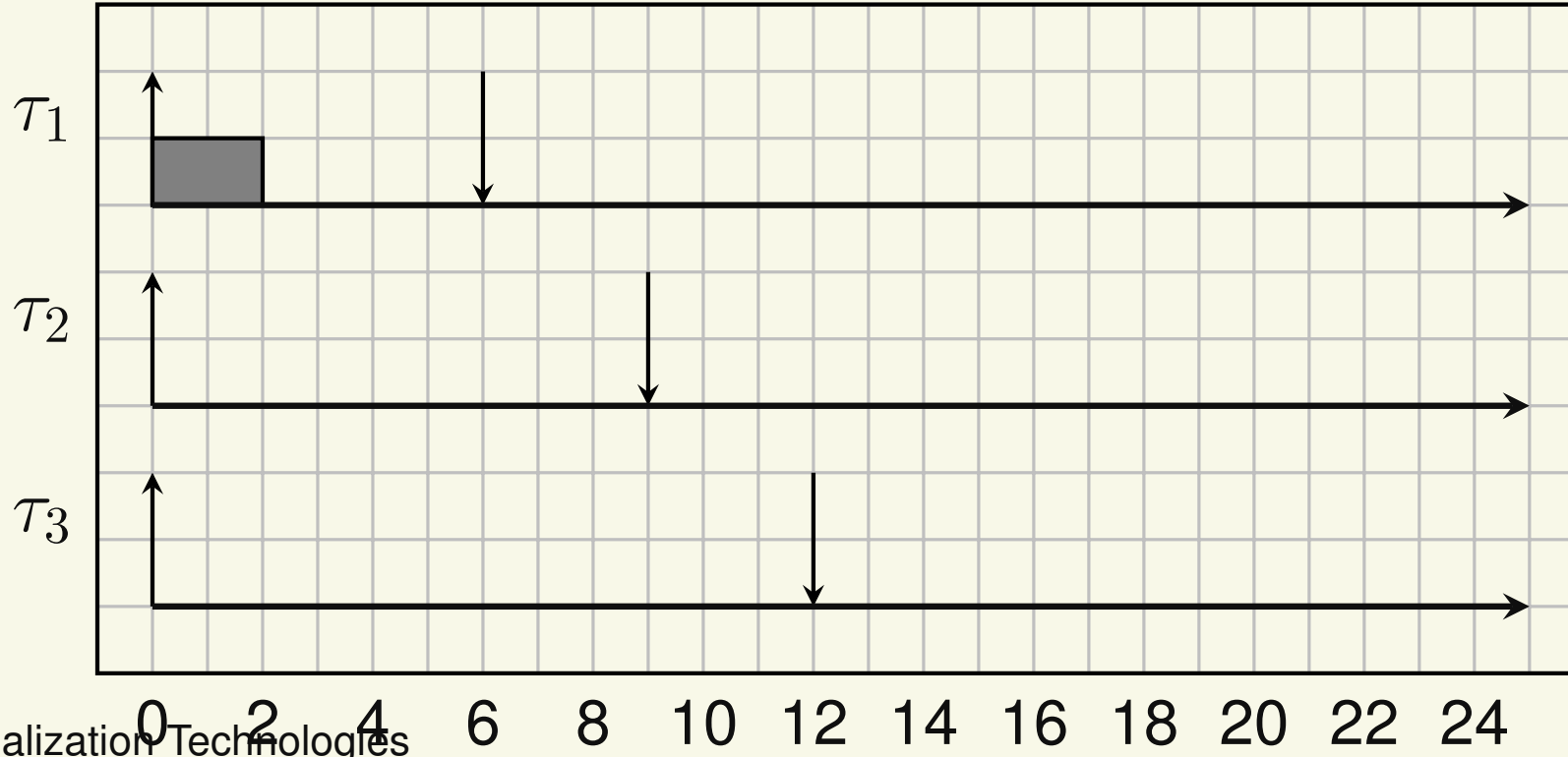
# Example of Schedule

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
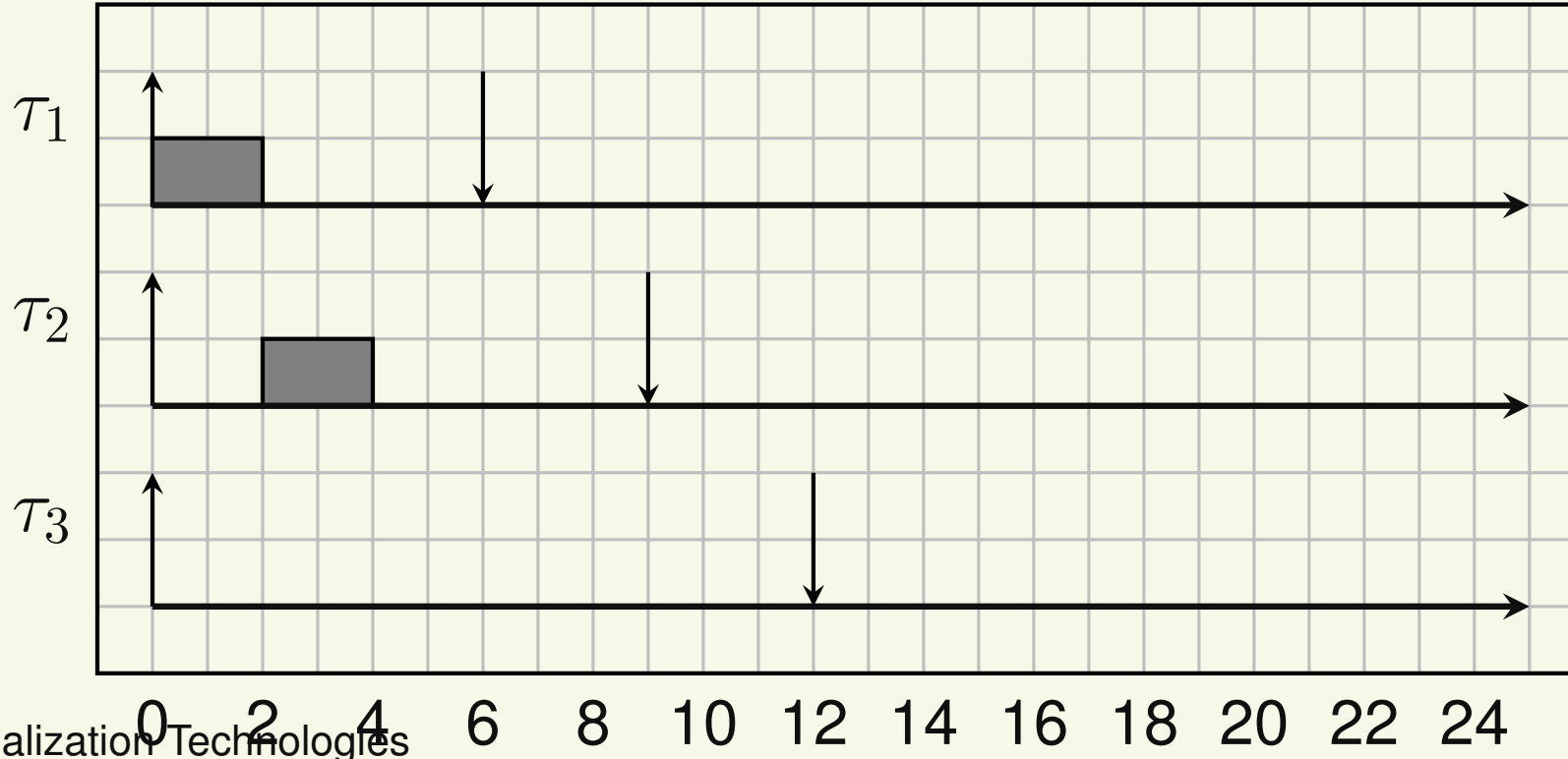
# Example of Schedule

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)

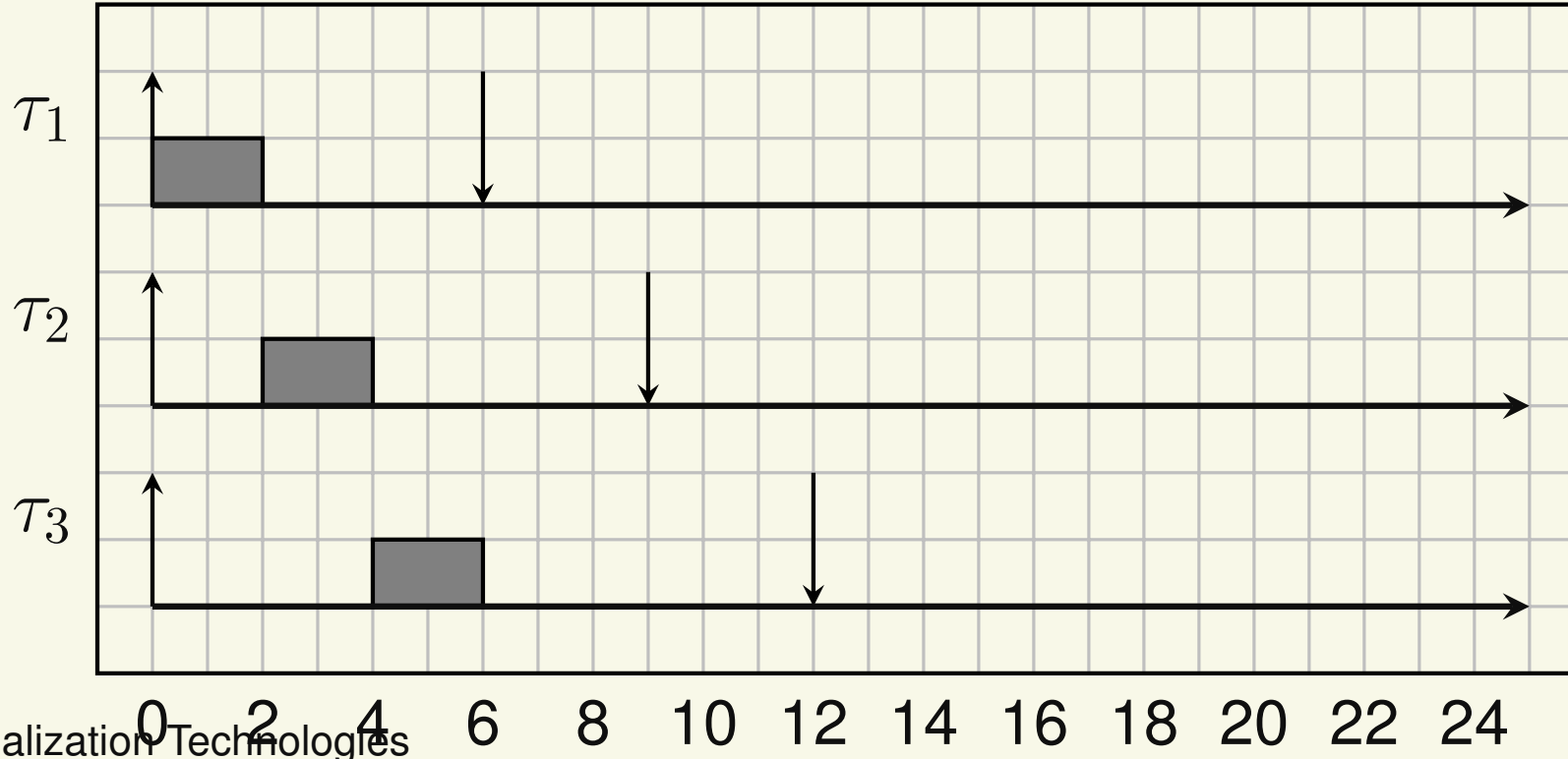- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
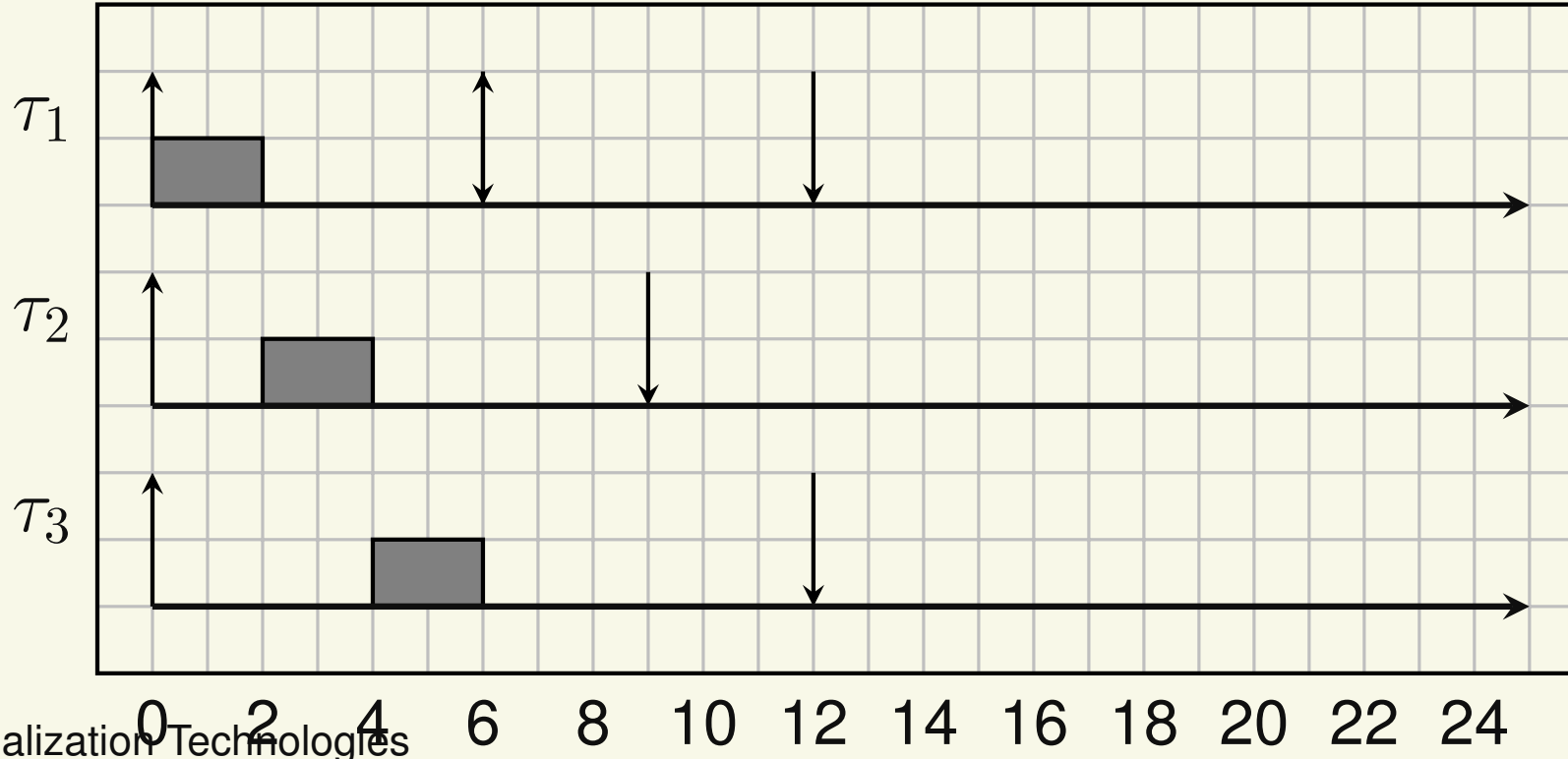
# Example of Schedule

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)

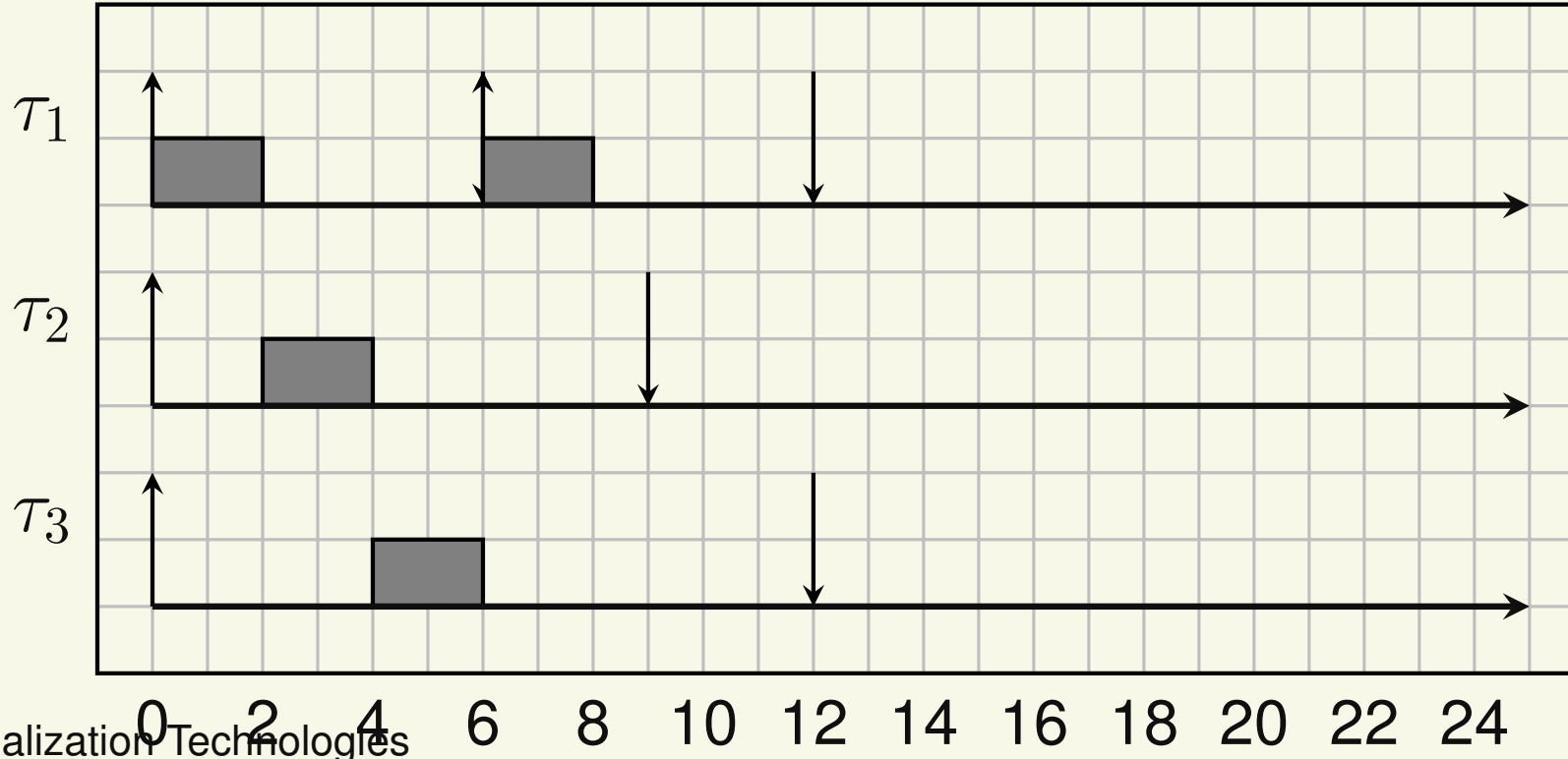# Example of Schedule

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
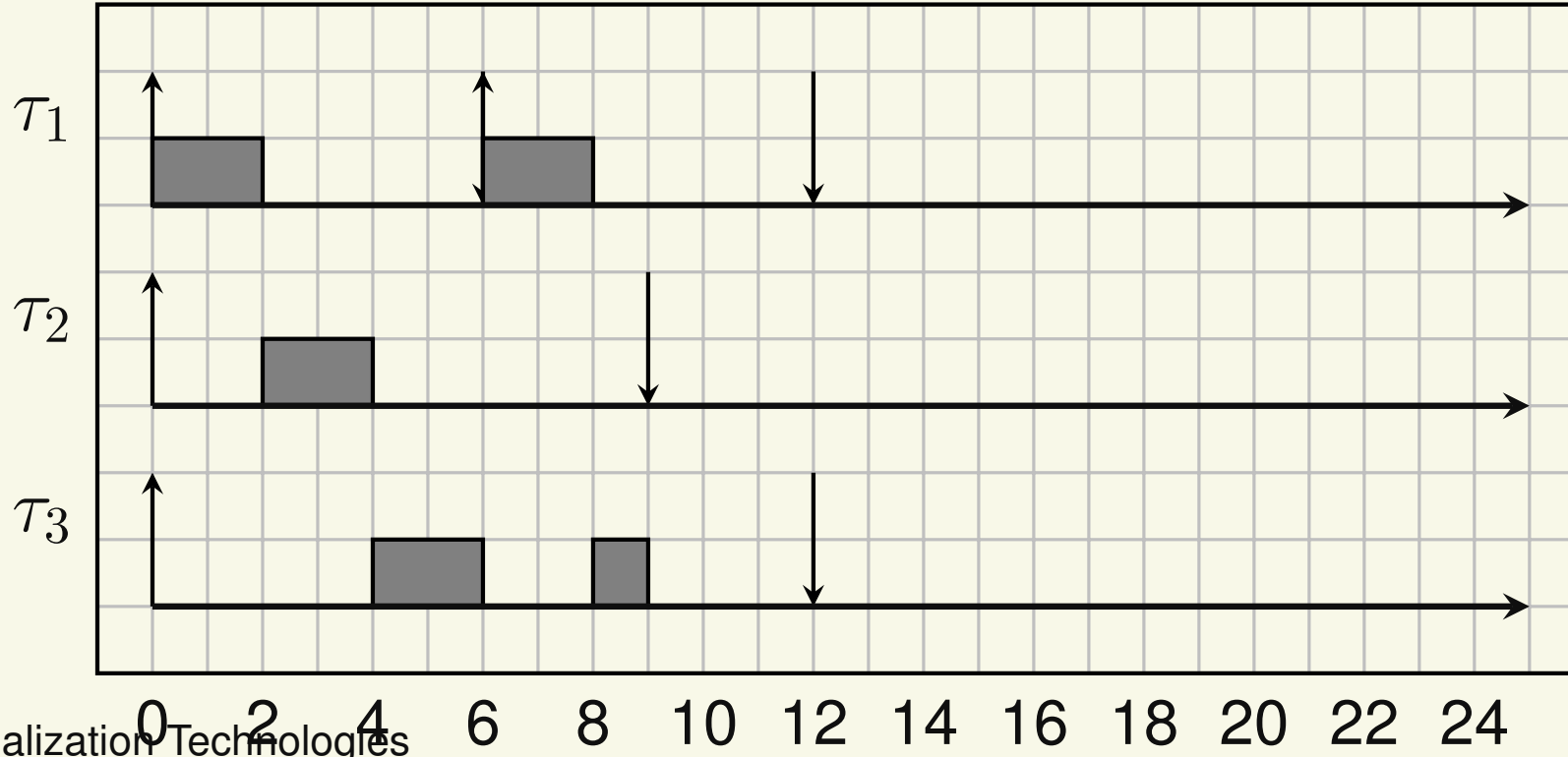
# Example of Schedule

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
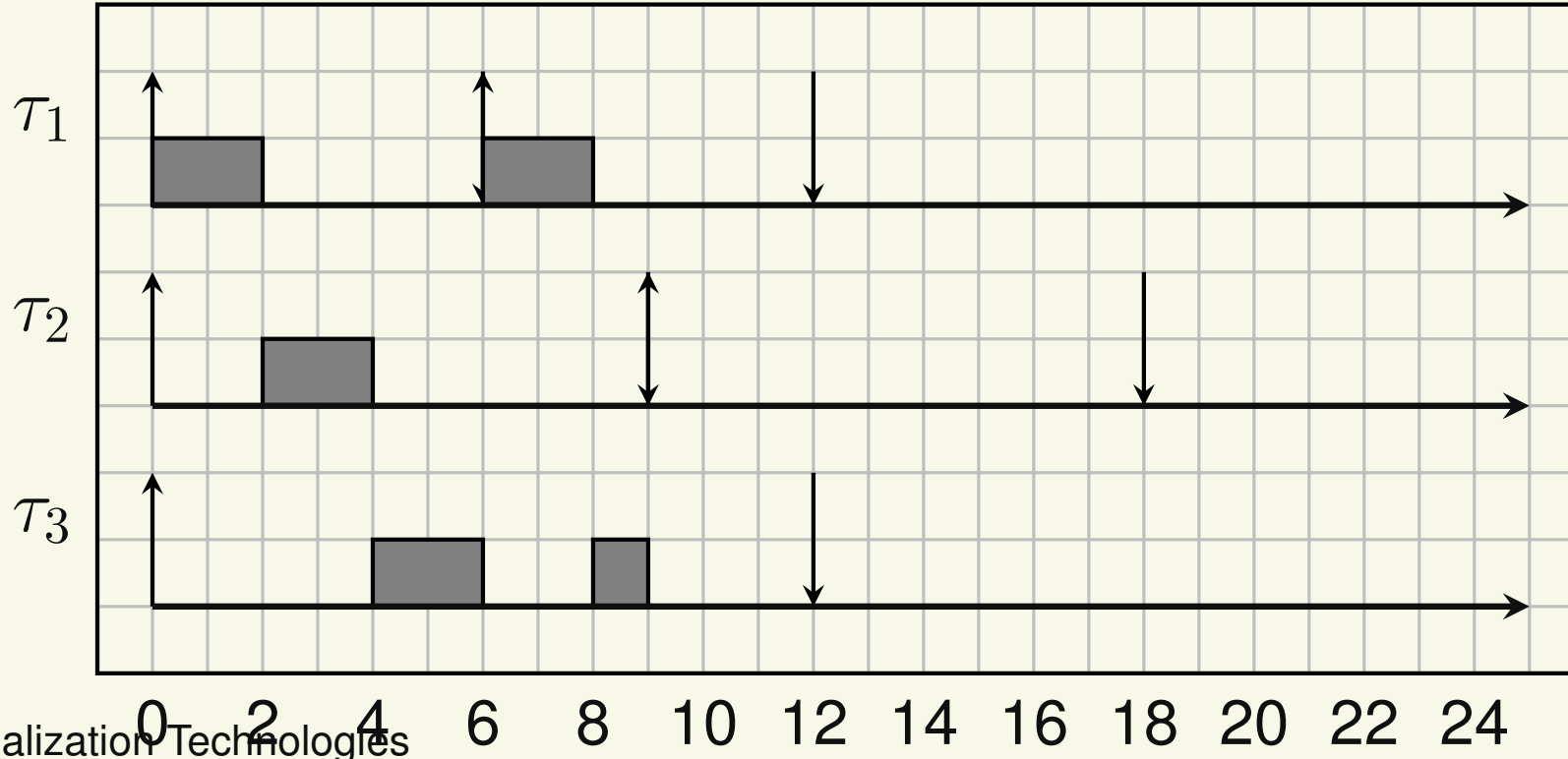
- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
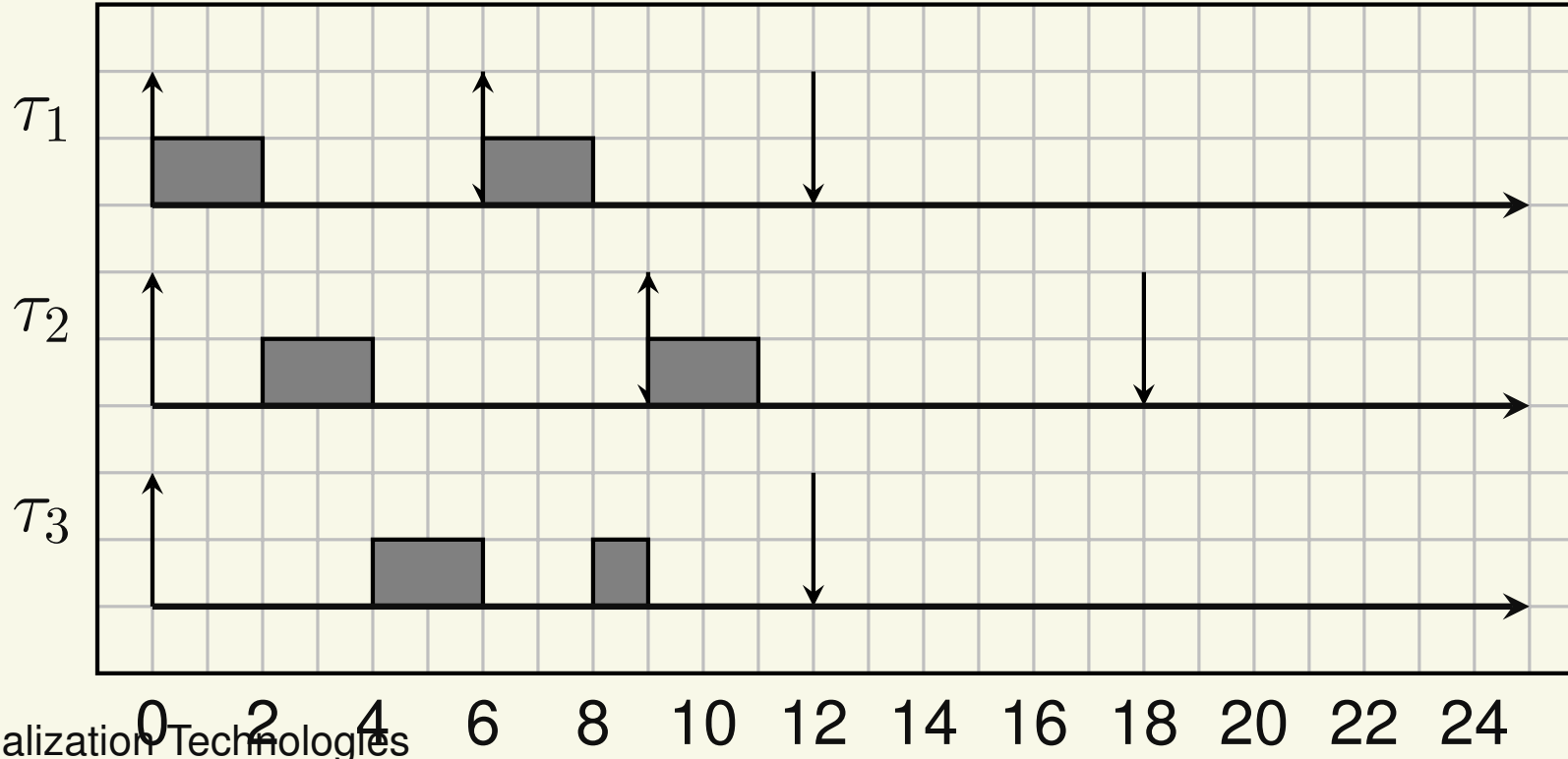
- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)

# Example of Schedule

- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
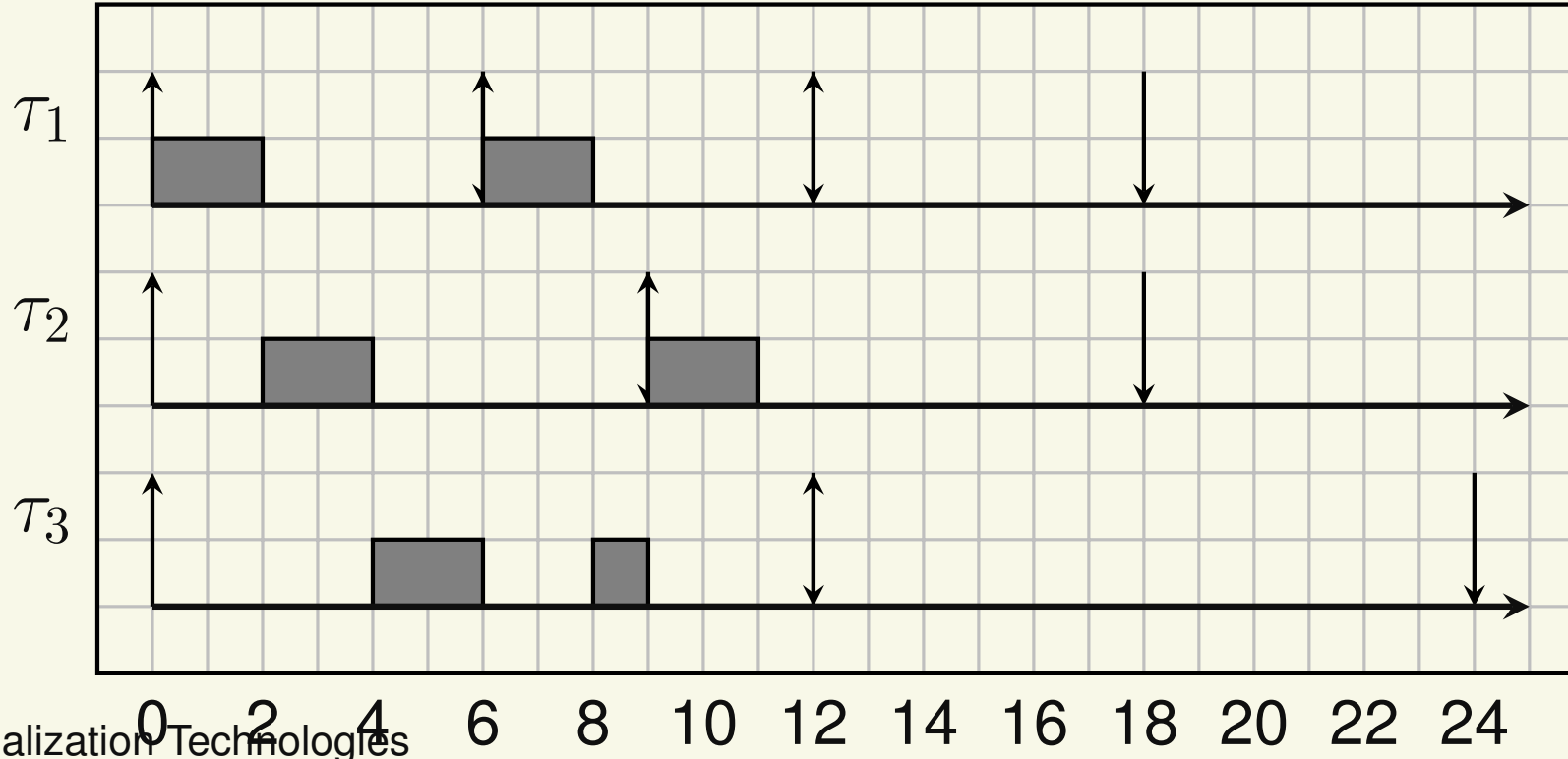
- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
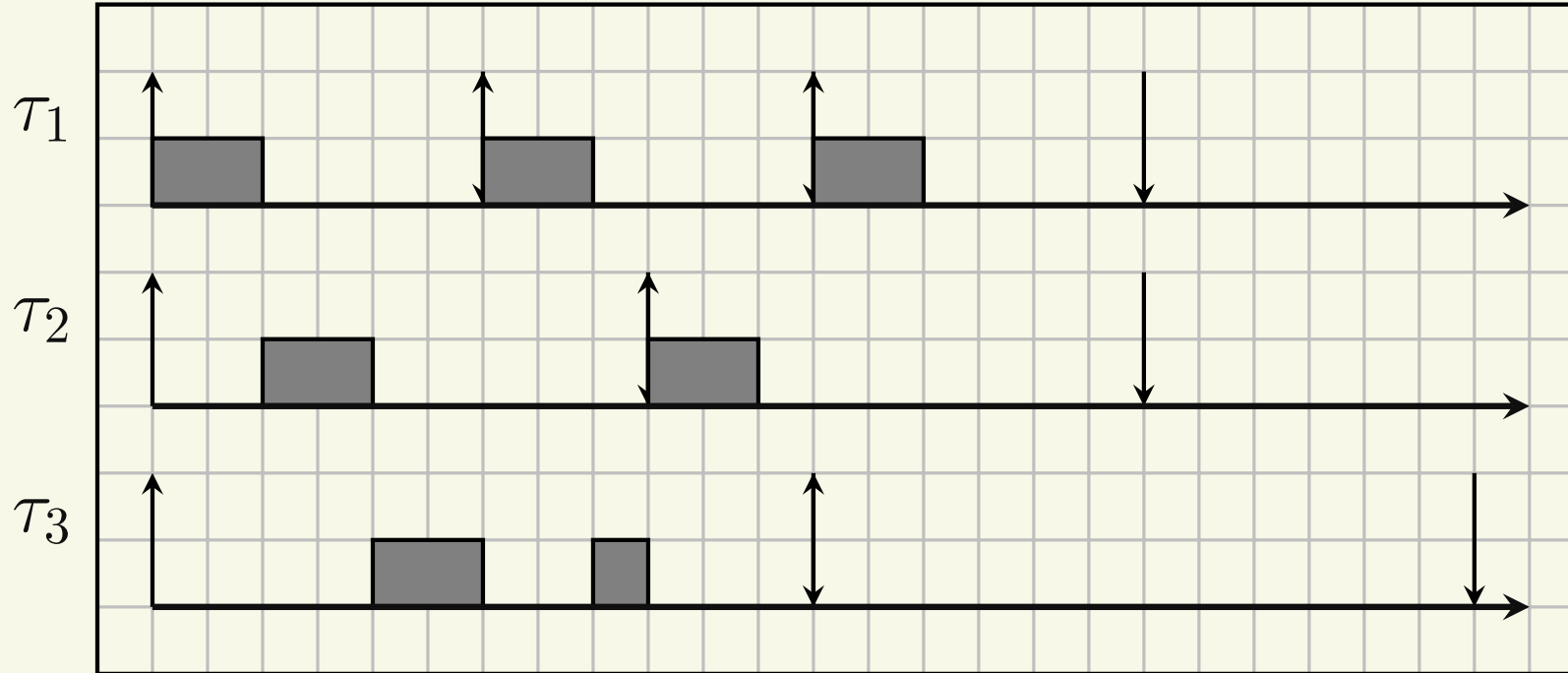
- Consider the following task set: $\tau_1 = (2, 6, 6)$, $\tau_2 = (2, 9, 9)$, $\tau_3 = (3, 12, 12)$. Task $\tau_1$ has priority $p_1 = 3$ (highest), task $\tau_2$ has priority $p_2 = 2$, task $\tau_3$ has priority $p_3 = 1$ (lowest)
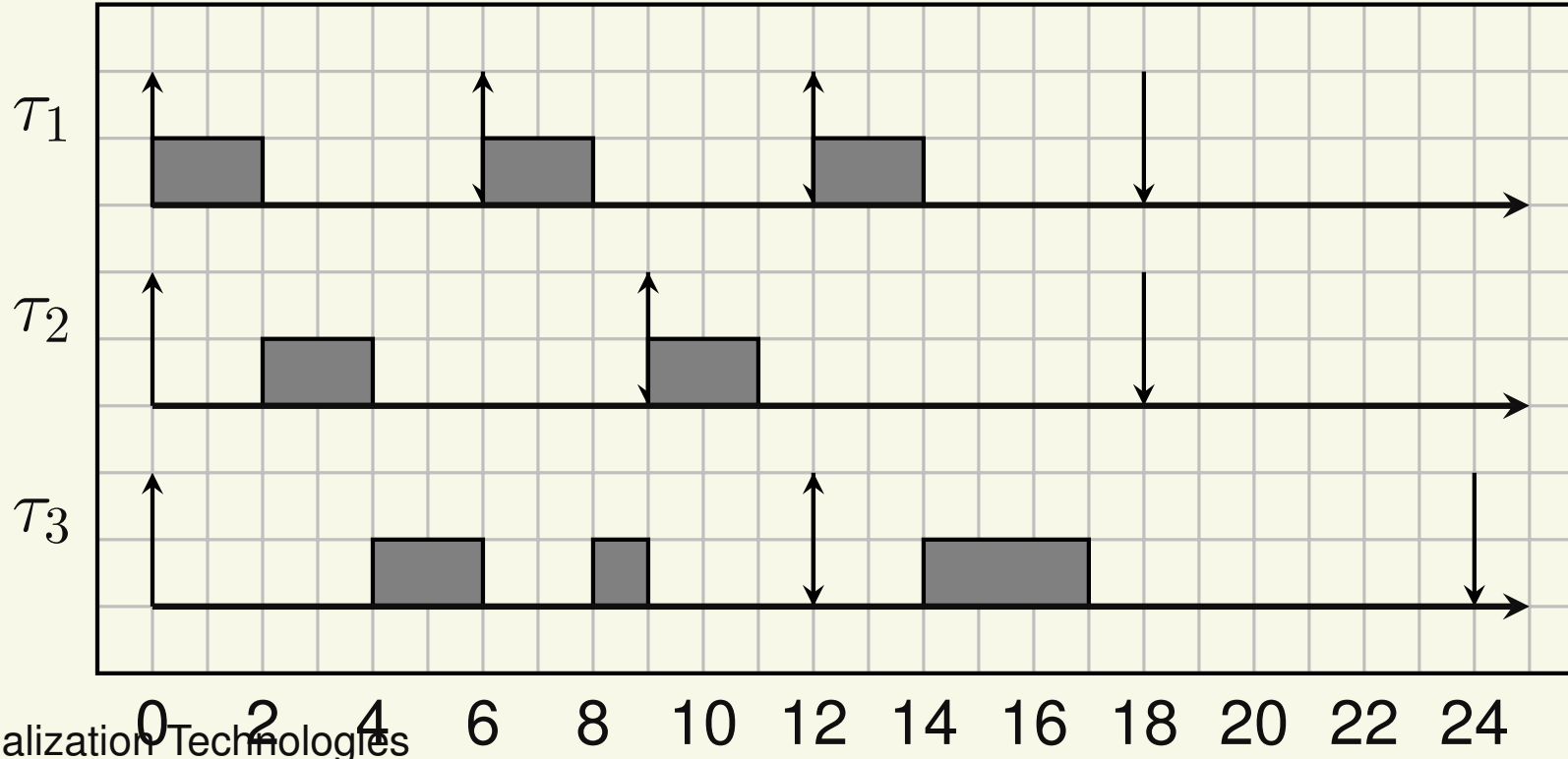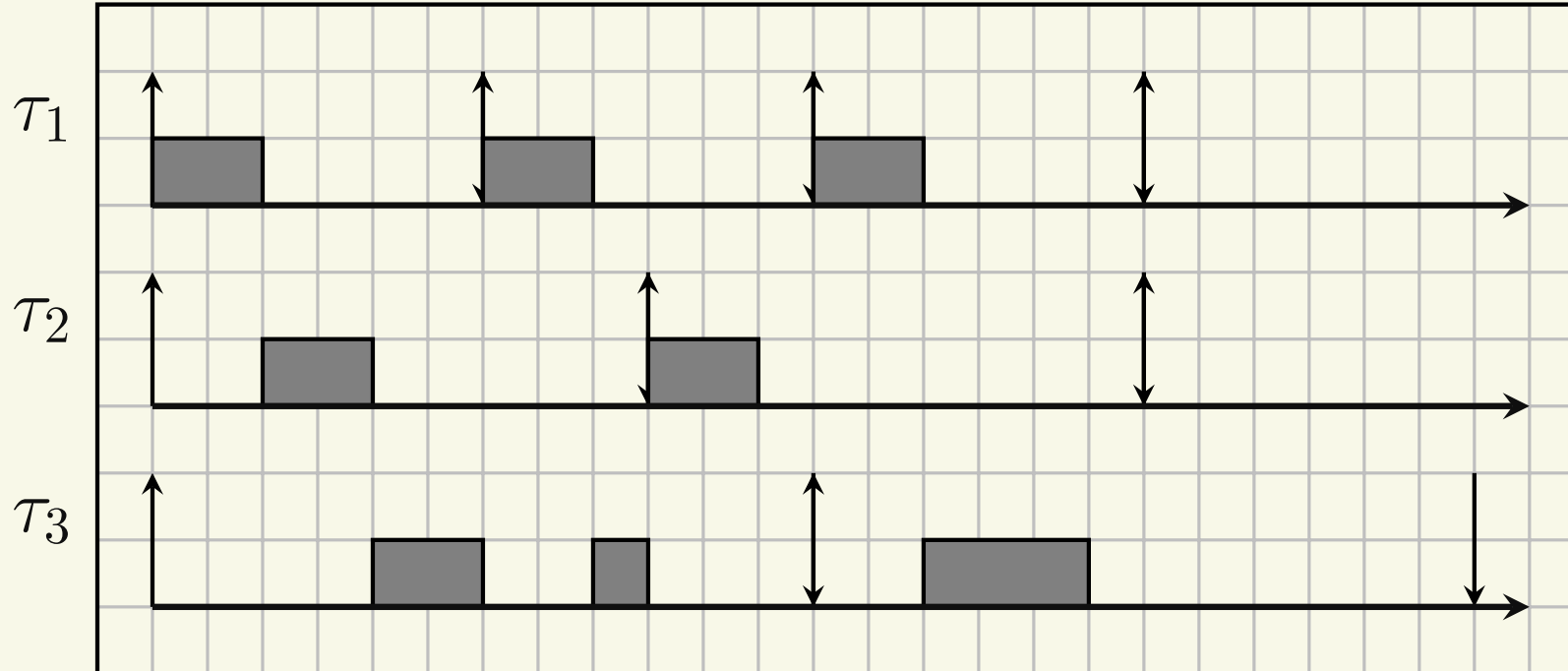
- Consider the following task set: $\tau_1 = (3, 6, 6)$, $p_1 = 3$, $\tau_2 = (2, 4, 8)$, $p_2 = 2$, $\tau_3 = (2, 12, 12)$, $p_3 = 1$



In this case, task $\tau_2$ misses its deadline!

# Another Example (non-schedulable)

- Consider the following task set: $\tau_1 = (3, 6, 6)$, $p_1 = 3$, $\tau_2 = (2, 4, 8)$, $p_2 = 2$, $\tau_3 = (2, 12, 12)$, $p_3 = 1$



In this case, task $\tau_2$ misses its deadline!

- Consider the following task set: $\tau_1 = (3, 6, 6)$, $p_1 = 3$, $\tau_2 = (2, 4, 8)$, $p_2 = 2$, $\tau_3 = (2, 12, 12)$, $p_3 = 1$



In this case, task $\tau_2$ misses its deadline!

# Notes about Priority Scheduling

- Some considerations about the schedule shown before:
    - The response time of the task with the highest priority is minimum and equal to its WCET
    - The response time of the other tasks depends on the *interference* of the higher priority tasks
    - The priority assignment may influence the schedulability of a task set
        - Problem: how to assign tasks' priorities so that a task set is schedulable?

# What About Multiple Cores?

- How to schedule tasks on multiple CPUs / cores?
  - First idea: partitioned scheduling
- Statically assign tasks to CPU cores
- Reduce the problem of scheduling on $M$ cores to $M$ instances of uniprocessor scheduling

# Or...

- One single task queue, shared by $M$ CPU cores
  - The first $M$ ready tasks are selected
  - What happens using fixed priorities?
  - Tasks are not bound to specific CPUs
  - Tasks can often migrate between different CPUs

- Problem: UP schedulers do not work well!

# Using Fixed Priorities in Linux

- `SCHED_FIFO` and `SCHED_RR` use fixed priorities

  - They can be used for real-time tasks, to implement RM and DM
  - Real-time tasks have priority over non real-time (`SCHED_OTHER`) tasks

- The difference between the two policies is visible when more tasks have the same priority

  - In real-time applications, try to avoid multiple tasks with the same priority

# Setting the Scheduling Policy

```
int sched_get_priority_max(int policy);
int sched_get_priority_min(int policy);

int sched_setscheduler(pid_t pid, int policy,
                       const struct sched_param *param);
int sched_setparam(pid_t pid,
                   const struct sched_param *param);
```

- If `pid == 0`, then the parameters of the running task are changed
- The only meaningful field of `struct sched_param` is `sched_priority`

# Problems with Real-Time Priorities

- In general, "regular" (`SCHED_OTHER`) tasks are scheduled in background respect to real-time ones
- Real-time tasks can / starve other applications
- Example: the following task scheduled at high priority can make a CPU / core unusable

```c
void bad_bad_task()
{
    while(1);
}
```

- Real-time computation have to be limited (use real-time priorities only when **really needed**!)
- Using real-time priorities requires root privileges (or part of them!)

# Real-Time Throttling

- A "bad" rt task can make a CPU / core unusable...
- ...Linux provides the *real-time throttling* mechanism
    - How does real-time throttling interfere with real-time guarantees?
    - Given a priority assignment, a taskset is guaranteed all the deadlines if no throttling mechanism is used...
    - ...But, what happens in case of throttling?
- Very useful idea, but something more "theoretically founded" might be needed...

# What About EDF?

- Can EDF (or similar) be supported in Linux?
- Problem: the kernel is not aware of tasks deadlines...
- ...But deadlines are needed to schedule the tasks
  - EDF schedules tasks based on absolute deadlines
- So, a more advanced API is needed...

# EDF on a real OS

- More advanced API:

  - Assign relative deadlines $D_i$ to the tasks...
  - A *runtime* and a *period* are also needed

- Moreover, $d_{i,j} = r_{i,j} + D_i$...

  - ...However, how can the scheduler know $r_{i,j}$?
  - The scheduler is not aware of jobs...

- To use EDF, the scheduler must know when a job starts / finishes

  - Modify applications, or guess...

# Tasks and Jobs... And Scheduling Deadlines!

- **Applications must be modified** to signal the beginning / end of a job (some kind of `startjob()` / `endjob()` system call)...
- ...Or the scheduler can assume that **a new job arrives each time a task wakes up**!
- Alternative:assign dynamic *scheduling deadlines*

    - Scheduling deadline $d_i^s$: **assigned by the kernel**
    - If the scheduling deadline $d_i^s$ matches the absolute deadline $d_{i,j}$ of a job, then the scheduler can respect $d_{i,j}$!!!

# Real-Time in VMs???

- Running real-time applications on an RTOS is not a problem...
- ...But, can real-time applications run in virtual machines?
    - Real-Time in Virtual Machines??? But... Why?
- Component-Based Development
    - Complex applications: sets of smaller components
    - Both functional and temporal interfaces
- Security (isolate real-time applications in a VM)
- Easy deployment; Time-sensitive clouds

# Real-Time in VMs

- Real-Time applications running in a VM?
  - As for OSs, two different aspects

# Real-Time in VMs

- Real-Time applications running in a VM?
  - As for OSs, two different aspects
    - Resource allocation/management (scheduling)

  - CPU allocation/scheduling: lot of work in literature

# Real-Time in VMs

- Real-Time applications running in a VM?
  - As for OSs, two different aspects

    - Latency (host and guest)

  - Latencies not investigated too much (yet!)

# Real-Time in VMs

- Real-Time applications running in a VM?
  - As for OSs, two different aspects
    - Resource allocation/management (scheduling)
    - Latency (host and guest)
  - CPU allocation/scheduling: lot of work in literature
  - Latencies not investigated too much (yet!)
- Virtualization: full hw or OS-level
  - OS-Level virtualization: real-time performance of the host kernel
  - Hw virtualization: hypervisors (example: KVM or Xen) can introduce latencies!

# Latency

- Latency: measure of the difference between the <span style="color:blue">theoretical</span> and <span style="color:red">actual</span> schedule

  - Task $\tau$ <span style="color:blue">expects</span> to be scheduled at time $t$ ...
  - ... but <span style="color:red">is actually scheduled</span> at time $t'$
  - $\Rightarrow$ Latency $L = t' - t$

- The latency $L$ can be accounted for in schedulability analysis

  - Similar to what is done for shared resources, etc...
  - Strange "shared resource": the OS kernel (or the hypervisor)

# Example: Periodic Task

- Consider a periodic task

```
/* ... */
while(1) {
    /* Job body */
    clock_nanosleep(CLOCK_REALTIME,
                    TIMER_ABSTIME, &r, NULL);
    timespec_add_us(&r, period);
}
```

- The task expects to be executed at time $r$ $(= r_0 + jT)$...

- ...But is sometimes delayed to $r_0 + jT + \delta$

# Theoretical Schedule

# Actual Schedule



- What happens if the $2^{nd}$ job of $\tau_1$ arrives a little bit later???

  - The $2^{nd}$ job of $\tau_2$ misses a deadline!!!

# Effects of the Latency

- Upper bound for $L$? If not known, no schedulability analysis!!!

  - The latency must be *bounded*: $\exists L^{max} : L < L^{max}$

- If $L^{max}$ is too high, only few task sets result to be schedulable

  - The worst-case latency $L^{max}$ cannot be too high

- Task: stream of jobs (activations) arriving at time $r_j$
- Task scheduled at time $t' > r_j \to$ Delay $t' - r_j$ caused by:

1. Job arrival (task activation) signaled at time $r_j + L^1$
2. Event served at time $r_j + L^1 + L^2$
3. Task actually scheduled at $r_{i,j} + L^1 + L^2 + I$

# Sources of Latency — 2

- $L = L^1 + L^2 + I$
- $I$: interference from higher priority tasks

    - Not really a latency!!!

- $L^2$: *non-preemptable section latency* $L^{np}$

    - Due to non-preemptable sections in the kernel (or hypervisor!) or to deferred interrupt processing

- $L^1$: delayed interrupt generation

    - Generally small
    - Hardware (or virtualized) timer interrupt: *timer resolution latency* $L^{timer}$

# Latency in Linux

- Tool (`cyclictest`) to measure the latency
  - Periodic task scheduled at the highest priority
  - Response time equal to execution time (almost $0$)
- Vanilla kernel: depends on the configuration
  - Can be tens of milliseconds
- Preempt-RT patchset (`https://wiki.linuxfoundation.org/realtime`): reduce latency to less than $100$ microseconds
  - Tens of microseconds on well-tuned systems!
- So, real-time on Linux is not an issue
  - Is this valid for hypervisors/VMs too?

# What About VM Latencies?

- Hypervisor: software component responsible for executing multiple OSs on the same physical node
  - Can introduce latencies too!
- Different kinds of hypervisors:
  - Xen: bare-metal hypervisor (*below* the Linux kernel)
    - Common idea: the hypervisor is small/simple, so it causes small latencies
  - KVM: hosted hypervisor (Linux kernel module)
    - Latencies reduced by using Preempt-RT
    - Linux developers already did lot of work!!!

# Hypervisor Latency

- Same strategy/tools used for measuring kernel latency
- Idea: run `cyclictest` in a VM
  - `cyclictest` process ran in the guest OS...
  - ...instead of host OS
- `cyclictest` period: $50\mu s$
- "Kernel stress" to trigger high latencies
  - Non-real-time processes performing lot of syscalls or triggering lots of interrupts
  - Executed in the host OS (for KVM) or in Dom0 (for Xen)
- Experiments on multiple x86-based systems

# Hypervisor Latencies



Intel Core Duo

Intel Core i7

| | |
|---|---|
| kvm, RT host, RT guest | ⊡ |
| kvm, RT host, NRT guest | ○ |
| kvm, NRT host, RT guest | △ |
| kvm, NRT host, NRT guest | ▽ |
| Xen, RT Dom0, RT DomU | ■ |
| Xen, RT Dom0, NRT DomU | ● |
| Xen, NRT Dom0, RT DomU | ▲ |
| Xen, NRT Dom0, NRT DomU | ▼ |

# Worst Cases

| Kernels | Core Duo | | Core i7 | |
|---------|----------|----------|---------|----------|
| | Xen | KVM | Xen | KVM |
| NRT/NRT | $3216\mu s$ | $851\mu s$ | $785\mu s$ | $275\mu s$ |
| NRT/RT | <span style="color:red">$4152\mu s$</span> | $463\mu s$ | <span style="color:red">$1589\mu s$</span> | $243\mu s$ |
| RT/NRT | $3232\mu s$ | $233\mu s$ | $791\mu s$ | $99\mu s$ |
| RT/RT | <span style="color:red">$3956\mu s$</span> | <span style="color:blue">$71\mu s$</span> | <span style="color:red">$1541\mu s$</span> | <span style="color:blue">$72\mu s$</span> |

- Preempt-RT helps a lot with KVM

  - <span style="color:blue">Good worst-case values (less than $100\mu s$)</span>

- Preempt-RT in the guest is dangerous for Xen

  - <span style="color:red">Worst-case values stay high</span>

# Hypervisor vs Kernel



- **Worst Cases:**
  - Host: $29\mu s$
  - Dom0: $201\mu s$ with Preempt-RT, $630\mu s$ with NRT

# Investigating Xen Latencies

- KVM: usable for real-time workloads
- Xen: strange results
    - Larger latencies in general
    - Using Preempt-RT in the guest increases the latencies?
- Xen latencies are not due to the hypervisor's scheduler
    - Repeating the experiments with the null scheduler did not decrease the experienced latencies

# Impact of the Kernel Stress

- Experiments repeated without "Kernel Stress" on Dom0
  - This time, using Preempt-RT in the guest reduces latencies!
  - Strange result: Dom0 load *should not* affect the guest latencies...

| Kernels | Core Duo | | Core i7 | |
|---|---|---|---|---|
| | Stress | No Stress | Stress | No Stress |
| NRT/NRT | $3216\mu s$ | $3179\mu s$ | $785\mu s$ | $1607\mu s$ |
| NRT/RT | $4152\mu s$ | $1083\mu s$ | $1589\mu s$ | $787\mu s$ |
| RT/NRT | $3232\mu s$ | $3359\mu s$ | $791\mu s$ | $1523\mu s$ |
| RT/RT | $3956\mu s$ | $960\mu s$ | $1541\mu s$ | $795\mu s$ |

# Virtualization Mechanisms

- Xen virtualization: PV, HVM, PVH, ...

  - PV: everything is para-virtualized
  - HVM: full hardware emulation (through qemu) for devices (some para-virtualized devices, too); use CPU virtualization extensions (Intel VT-x, etc...)
  - PVH: hardware virtualization for the CPU + para-virtualized devices (trade-off between the two)

- Dom0 kernel does not affect results; focus on guest kernel

| Guest Kernel | PV | PVH | HVM |
|---|---|---|---|
| NRT | $661\mu s$ | $1276\mu s$ | $1187\mu s$ |
| RT | $178\mu s$ | $216\mu s$ | $4470\mu s$ |

# What's up with HVM?

- HVM uses qemu as *Device Model* (DM)

  - Qemu instance running in Dom0
  - Used for boot and emulating some devices...
  - ...But somehow involved in the strange latencies!!!

- Scheduling all qemu threads with priority 99, the worst-case latencies are comparable with PV / PVH!!!

  - High HVM latencies due to the Kernel Stress workload preempting qemu...

- Summing up: for good real-time performance, use PV or PVH!

# Cyclictest Period

- Most of the latencies larger than cyclictest period...
- Are hypervisor's timers able to respect that period?

  - Example of timer resolution latency...

- So, let's try a larger period!

  - $500\mu s$ and $1ms$ instead of $50\mu s$
  - Measure timer resolution latency $\rightarrow$ no kernel stress

- Results are much better!

  - $P = 500\mu s$: worst-case latency $112\mu s$ (HVM), $82\mu s$ (PVH) or $101\mu s$ (PV)
  - $P = 1000\mu s$: worst-case latency $129\mu s$ (HVM), $124\mu s$ (PVH) or $113\mu s$ (PV)

# Further Analysis

- Xen latencies seem to be mainly due to timer resolution latency
    - Turned out to be an issue in the Linux code handling Xen's para-virtualized timers
        - Linux jargon: "clockevent device"
    - Does not activate a timer at less than $100\mu s$ from current time (`TIMER_SLOP`)
- After reducing the timer slop, average latency smaller than $50\mu s$ even for cyclictest with period $50\mu s$
    - Still larger than KVM latencies (probably due to non-preemptable sections?)

# Final Results

- Xen with a properly configured `TIMER_SLOP`:

  - Timer resolution latency reduced to almost $0$
  - Non-preemptable section latency dependent on the virtualization technology
  - Worst-case latencies higly dependent on the hardware

    - Example: some old CPUs need to (trap and) emulate `rdtsc` $\Rightarrow 15\mu s$ additional latency

- Xeon CPU: $28\mu s$ with PVH, $72\mu s$ for PV (KVM is $44\mu s$)
- Core 2 CPU: $88\mu s$ for PV, $182\mu s$ for PVH (KVM is $71\mu s$)

# Reproducible Results

- Results can be reproduced on your test machine

  - You just need some manual installation of KVM, Xen, etc...

  `http://retis.santannapisa.it/luca/VMLatencies`

- Scripts to reproduce the previous experiments

  - Numbers depend on the hw, but the obtained figures are consistent with the previous results

- Other figures can be easily obtained by modifying scripts / configuration files
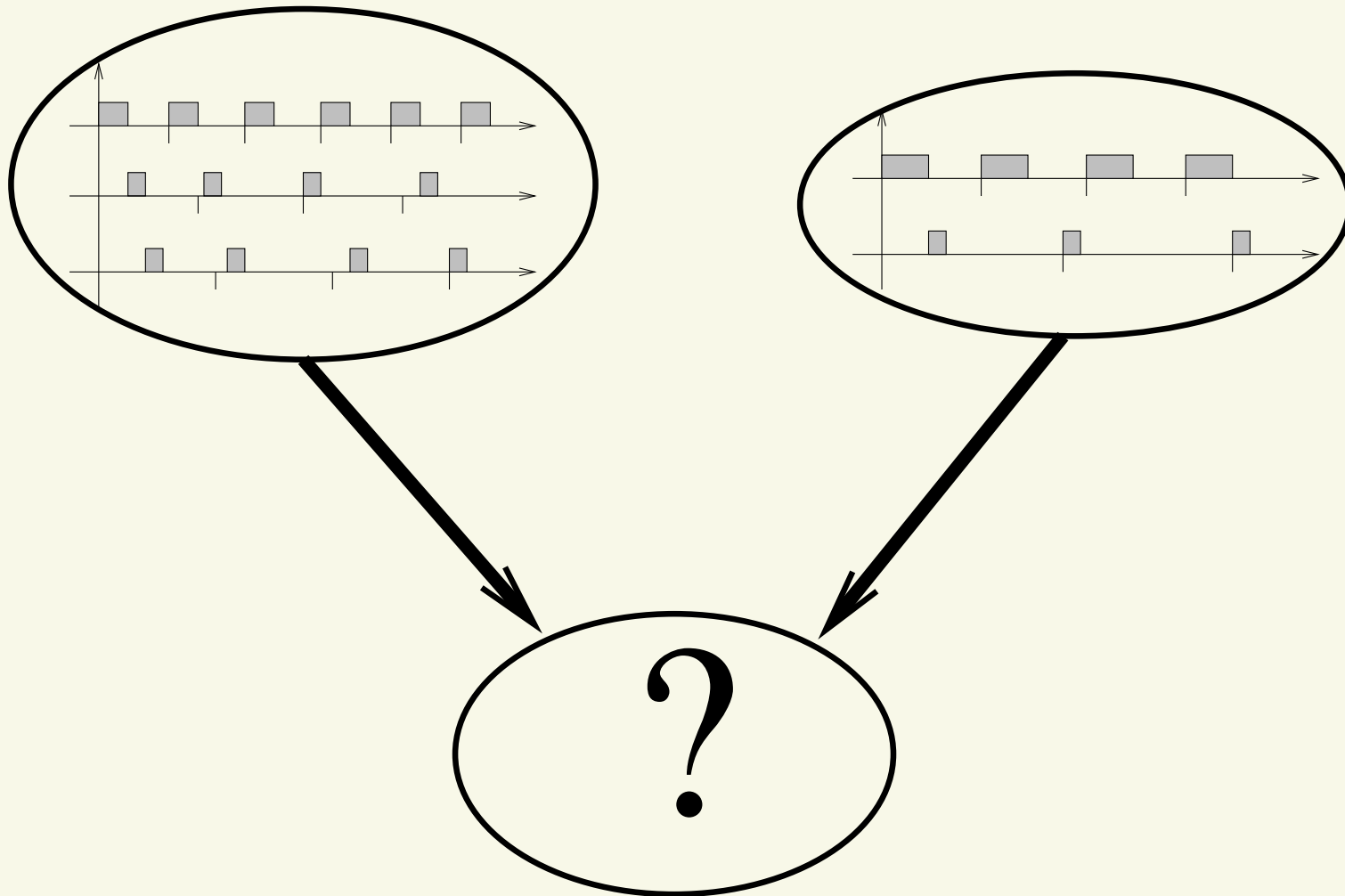
# Summing Up

- Latencies experienced in a VM (`cyclictest`)

  - KVM: Preempt-RT allows to achieve low latencies → usable for real-time
  - Xen: high latencies, Preempt-RT does not help, strange impact of the Dom0 load

- Xen behaves better when PV or PVH is used

  - Part of the latencies due to the DM (qemu running in Dom0)?

- Xen experiences a large timer resolution latency

  - Fixable by modifying the guest kernel

# Latencies and Scheduling

- Most of the industrial work on real-time virtualization <span style="color:red">focused on latency reduction</span>

  - Example: real-time KVM industrial solution based on vCPU pinning — No scheduling!!!

- Scheduling VMs is still needed to share hardware resources...

  - Bounded latencies are needed to have precise and accurate vCPU scheduling...

  - ...But <span style="color:blue">appropriate scheduling algorithms are still needed</span>!!!

- Advanced scheduling algoritms are useless if latencies are not bounded, and bounded latencies are useless if appropriate scheduling is not used!

# Combining Real-Time Guarantees



- Schedulability analysis in each VM...
- What about the resulting system?

# Real-Time Applications Inside VMs

- VM $\mathcal{C}^i$ contains $n^i$ tasks
- How to analyze its schedulability?

  - We only know how to schedule single tasks...
  - And we need to somehow "summarise" the requirements of a VM!

$$\mathcal{C}^i = \{(C_0^i, D_0^i, T_0^i), (C_1^i, D_1^i, T_1^i), \ldots, (C_{n^i}^i, D_{n^i}^i, T_{n^i}^i)\}$$

- So, $2$ main issues:

  1. Describe the temporal requirements of a VM in a simple way
  2. Schedule the VMs, and somehow "combine" their temporal guarantees

# The "not so smart" Solution
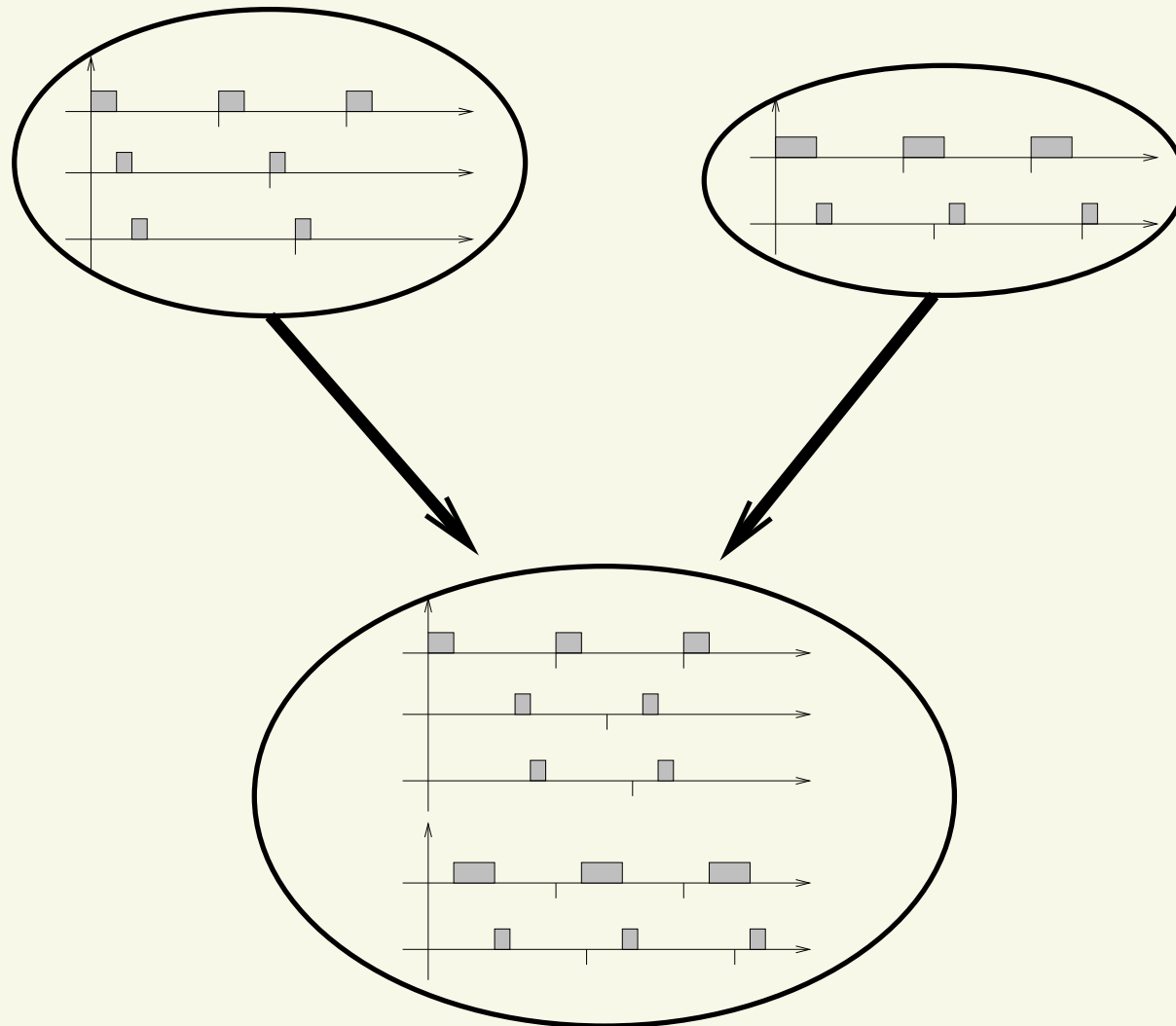
- Each VM is a set of real-time tasks:

$$\mathcal{C}^i = \{(C_j^i, D_j^i, T_j^i)\}$$

- Build the "global taskset" composed by all the tasks from all the VMs

$$\Gamma = \bigcup_i \mathcal{C}^i$$

- ...And use some known real-time scheduler (RM, EDF, ...) on $\Gamma$!

# Flattened Scheduling



- One single "flattened" scheduler seeing all the tasks
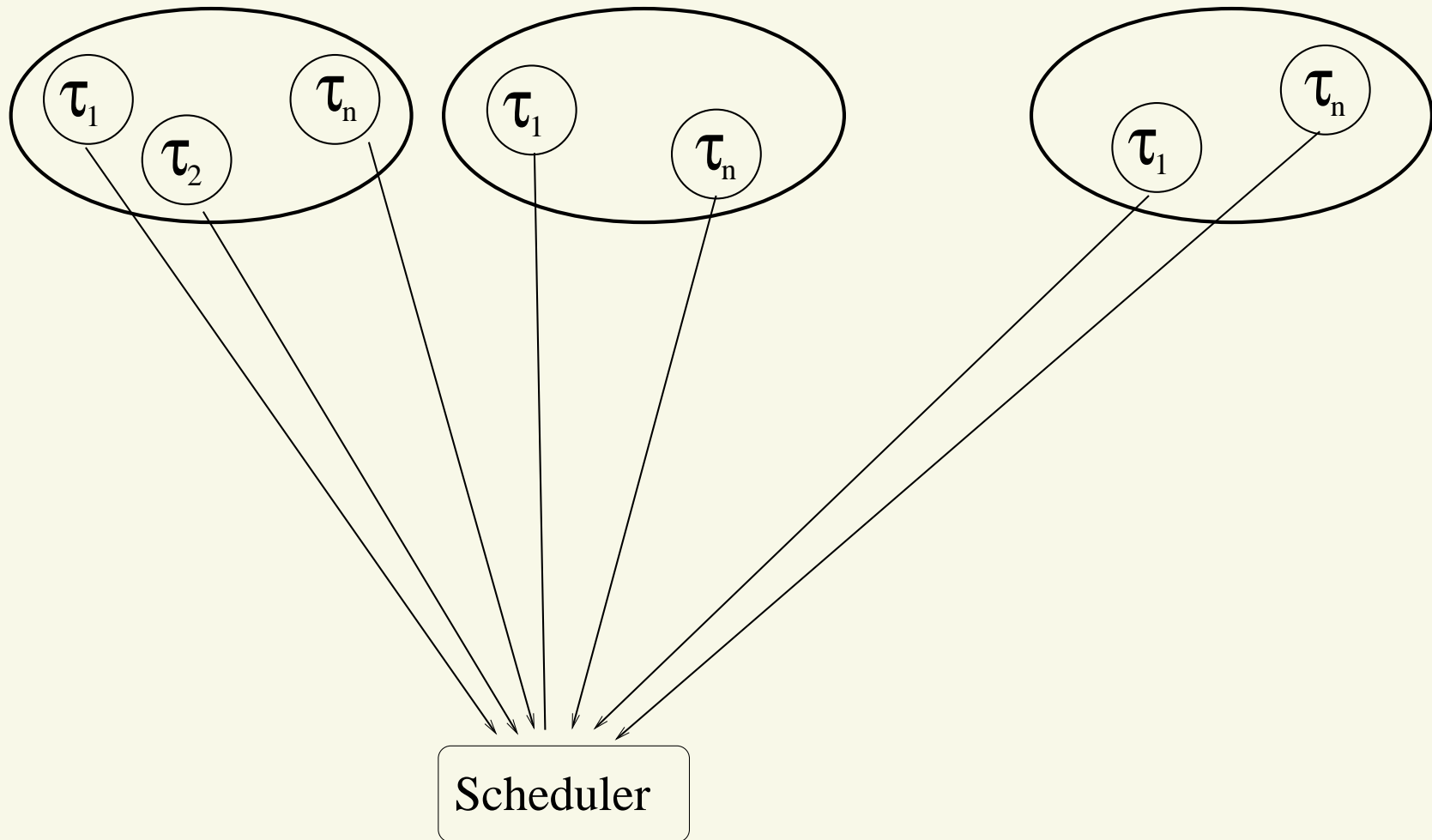
# Why it is "not so smart"

- One single scheduler, that must "see" all the tasks of all the VMs

    - Internals of the VMs have to be exposed!
    - VMs cannot run their own "local" schedulers
    - Misbehaving tasks in a VM can affect other VMs

        - No isolation!!!

- Using fixed priorities might be "not so simple"

    - Think about RM: priorities in a VM might depend on other VMs...

# Practical Issues

- The host/hypervisor scheduler only sees a VMs, but <span style="color:red">cannot see the tasks inside</span> it
- Para-virtualization (of the OS scheduler) could be used to address this issue, but it is not so simple...
- ...And requires huge modifications to host, guest, and applications!
- So, how to schedule VMs?
- Two-level hierarchical scheduling system

  - Host (global / root) scheduler, scheduling VMs
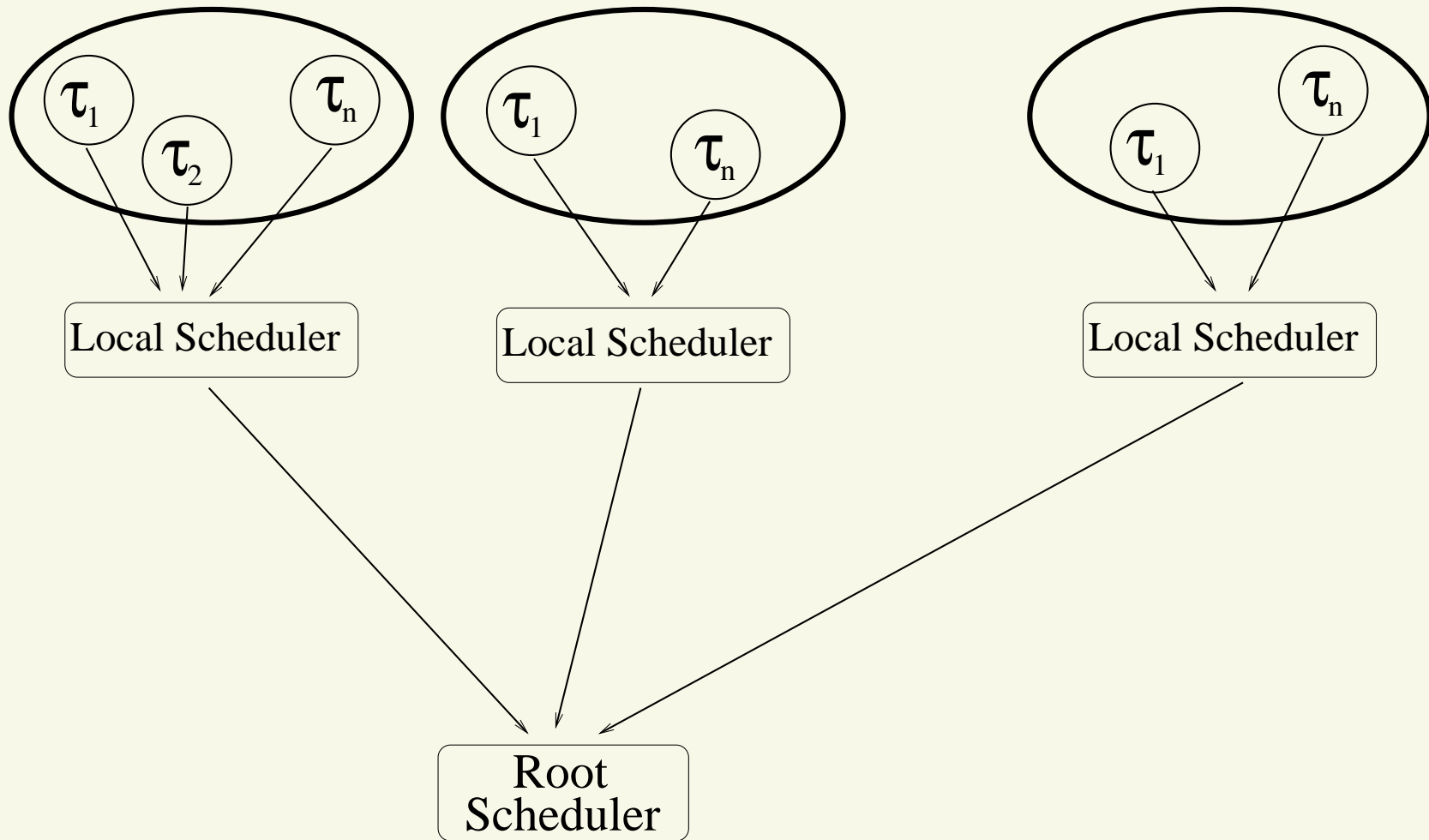  - Each VM contains its (local / 2nd level) scheduler

- Scheduler assigns CPU to tasks "inside the VMs"

# ...To a 2-Levels Hierarchy



- Host Scheduler assigns CPU to VMs
- Local Schedulers assign CPU to single tasks

# Hierarchical Scheduling

- The root scheduler does not see the tasks
- The OSs inside VMs are free to define their own (fixed priorities, EDF, whatever) schedulers

  - No problems in assigning fixed priorities to tasks!

- Root scheduler: host / hypervisor scheduler
- Local scheduler: guest scheduler
- Problem: what to use as a root scheduler?

  - We must have a model for it
  - Must allow to compose the "local guarantees"

- Before going on, summary of RT definitions and concepts

# Real-Time Guarantees in a Component

- First requirement: analyse the schedulability of a component independently from other components
  - This means that the root scheduler must provide some kind of temporal protection between components
- Various possibilities
  - Resource Reservations / server-based approach
  - Static time partitioning
  - ...
- In any case, the root scheduler must guarantee that each VM receives a minimum amount of resources in a time interval

# Schedulability Analysis: the Basic Idea

- (Over?)Simplifying things a little bit...
- ...Suppose to know the amount of time needed by a component to respect its temporal constraints and the amount of time provided by the root scheduler
- A component is "schedulable" if

$$\text{demanded time} \leq \text{supplied time}$$

- "demanded time": amount of time (in a time interval) needed by a component
- "supplied time": amount of time (in a time interval) given by the root scheduler to a component

- Of course the devil is in the details

# Demanded Time

- Amount of time needed by a component to respect its temporal constraints
  - Depends on the time interval we are considering
  - Depends on the component's local scheduler
    - EDF $\rightarrow dbf(t) = \sum_j \max\{0, \left\lfloor \frac{t+T_j-D_j}{T_j} \right\rfloor\}C_j$
    - RM: $\rightarrow$ workload $W(t) = C_i + \sum_{j<i} \left\lceil \frac{t}{T_i} \right\rceil C_j$
  - Note: $W(t)$ is very pessimistic, $dbf(t)$ is not
- This is the description of the temporal requirements of a component we were searching for...
- And what about the supplied time?

# Supplied Time

- Description of the root scheduler temporal behaviour
- More formally:

  - Depends on the time interval $t$ we are considering
  - Depends on the root scheduler $\mathcal{A}$

- Minimum amount of time given by $\mathcal{A}$ to a VM in a time interval of size $s$

  - Given all the time interval $(t_0, t_1) : t_1 - t_0 = s$...
  - ...Compute the size of the sub-interval in which $\sigma(t) = VM$...
  - ...And then find the minimum!

# Supplied Time Bound Function

- Even more formally:

  - Define $s(t) = \begin{cases} 1 & \text{if } \alpha(t) = VM \\ 0 & \text{otherwise} \end{cases}$

  - Time for VM in $(t_0, t_0 + s)$: $\int_{t_0}^{t_0+s} s(t)dt$

  - Then, compute the minimum over $t_0$
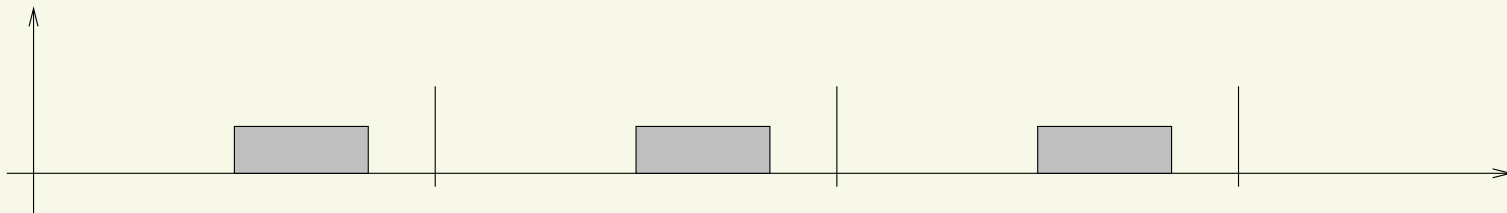
- $sbf(t) = \min_{t_0} \int_{t_0}^{t_0+t} s(x)dx$

# Example: Static Time Partitioning

- First (very simple) example of VM scheduling: static time partitioning
  - Static schedule describing when time is assigned to each VM
  - Pre-computed $\sigma(t)$
- Generally, periodic!
  - Otherwise, need to store an infinite schedule...
  - ...Might be problematic!
- Example: VM$_\mathcal{A}$ is scheduled in $(3, 4)$, $(9, 10)$, $(15, 16)$, ...
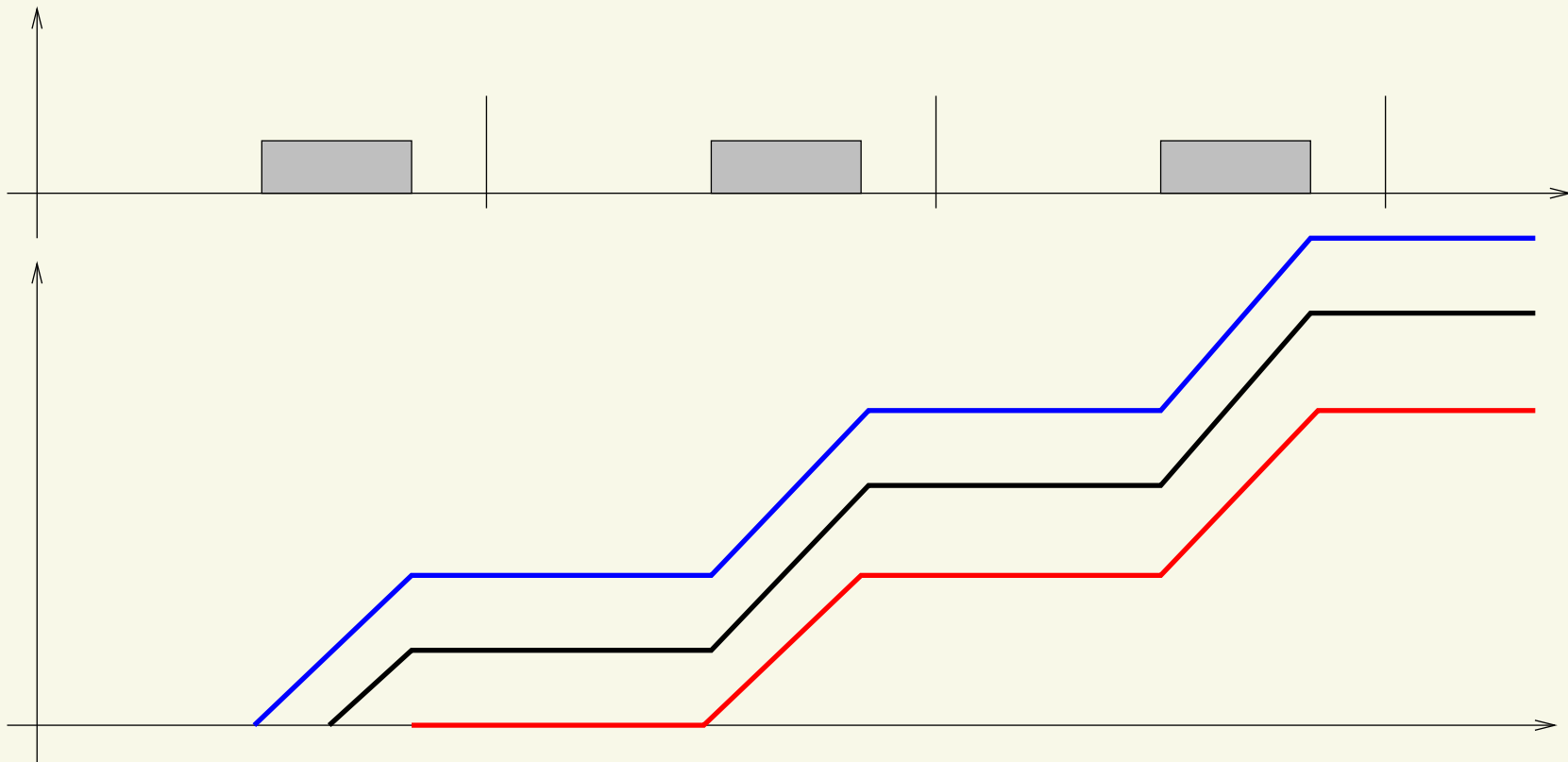  - More formally: $s(t) = 1$ if $6k + 3 \leq t \leq 6k + 4$, $s(t) = 0$ otherwise

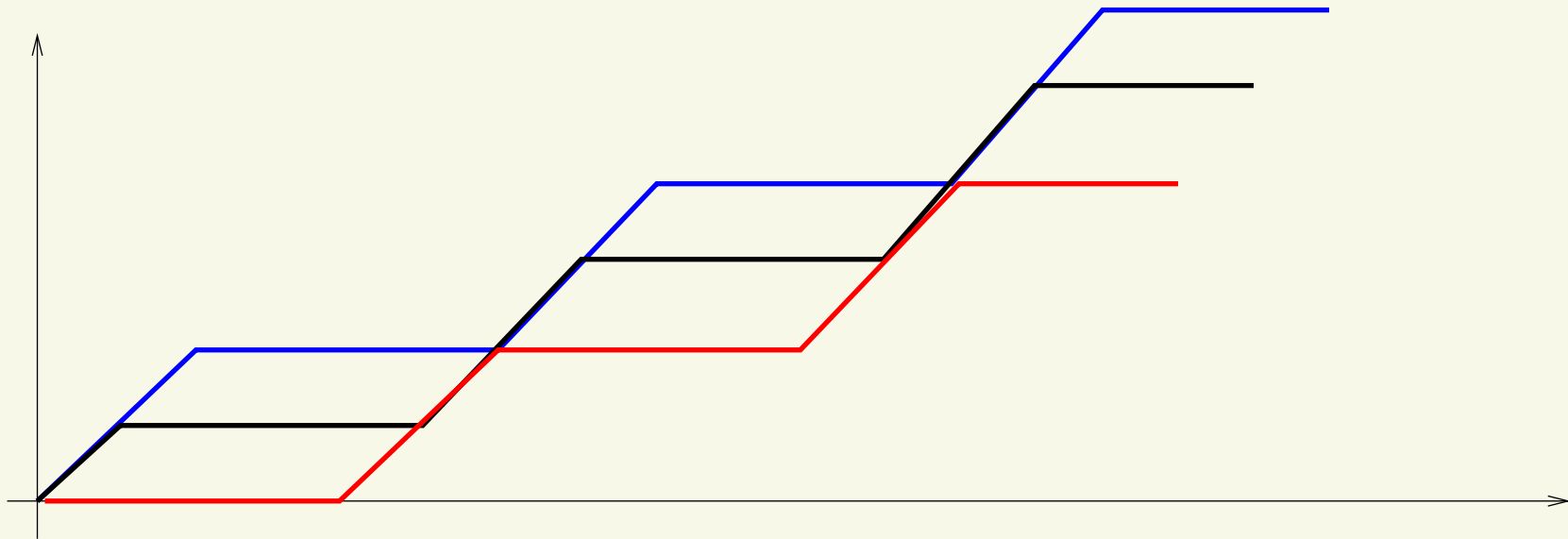$$s(t) = \begin{cases} 1 & \text{if } 6k + 3 \leq t \leq 6k + 4 \\ 0 & \text{otherwise} \end{cases}$$



- What is the supply bound function $sbf(t)$ in this case?
- Let's try different supply functions compatibe with this schedule...
- ...And see what is the worst case!
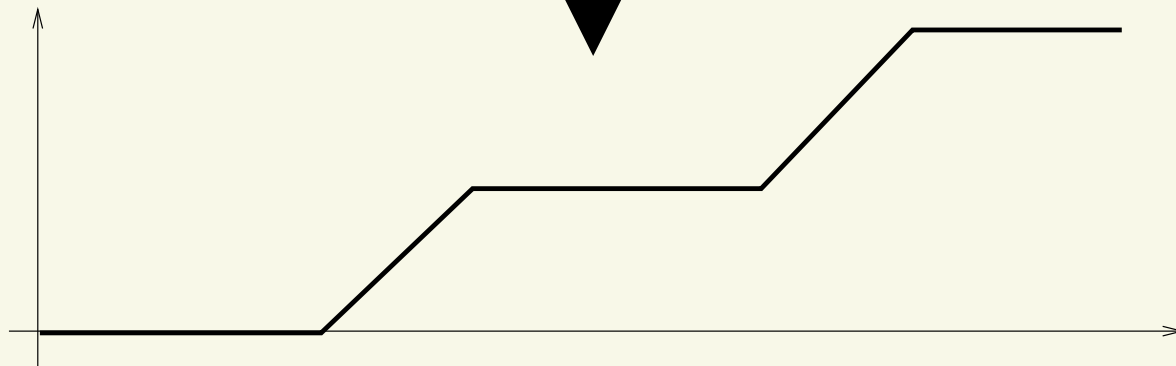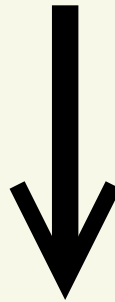  - Intervals of size $t$ starting at different times...

- Different supply functions depending on when the considered interval begins
- Which one is the worst case (supply bound function)?
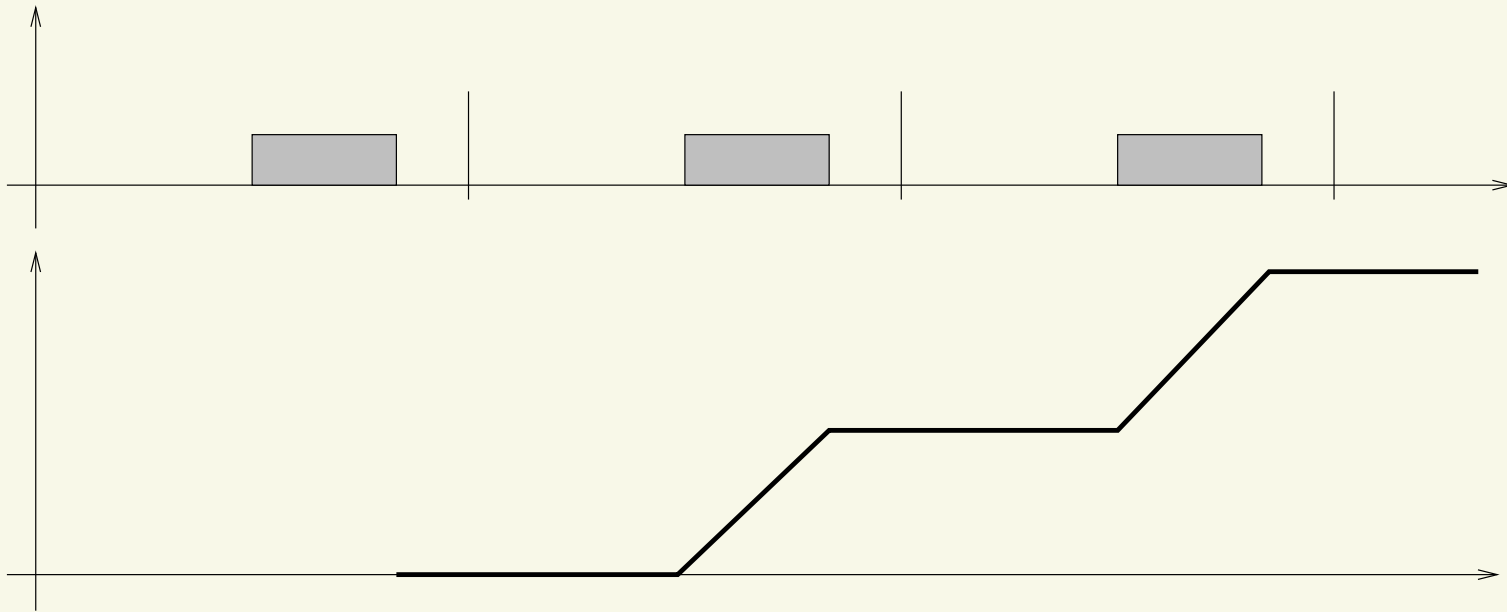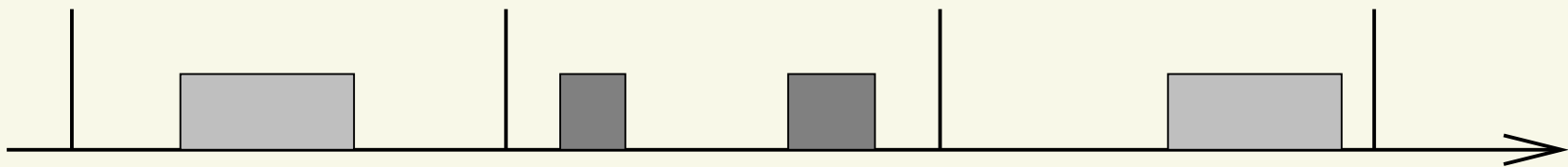
- Different supply functions depending on when the considered interval begins
- Which one is the worst case (supply bound function)?
  - The red one!

# Periodic Servers

- Periodic Server $\mathcal{S} = (Q, P)$: guarantees $Q$ units of time every period $P$
  - Can be implemented in different ways (example: CBS)
- Different from static allocation: we do not know where in the period the $Q$ time units are allocated
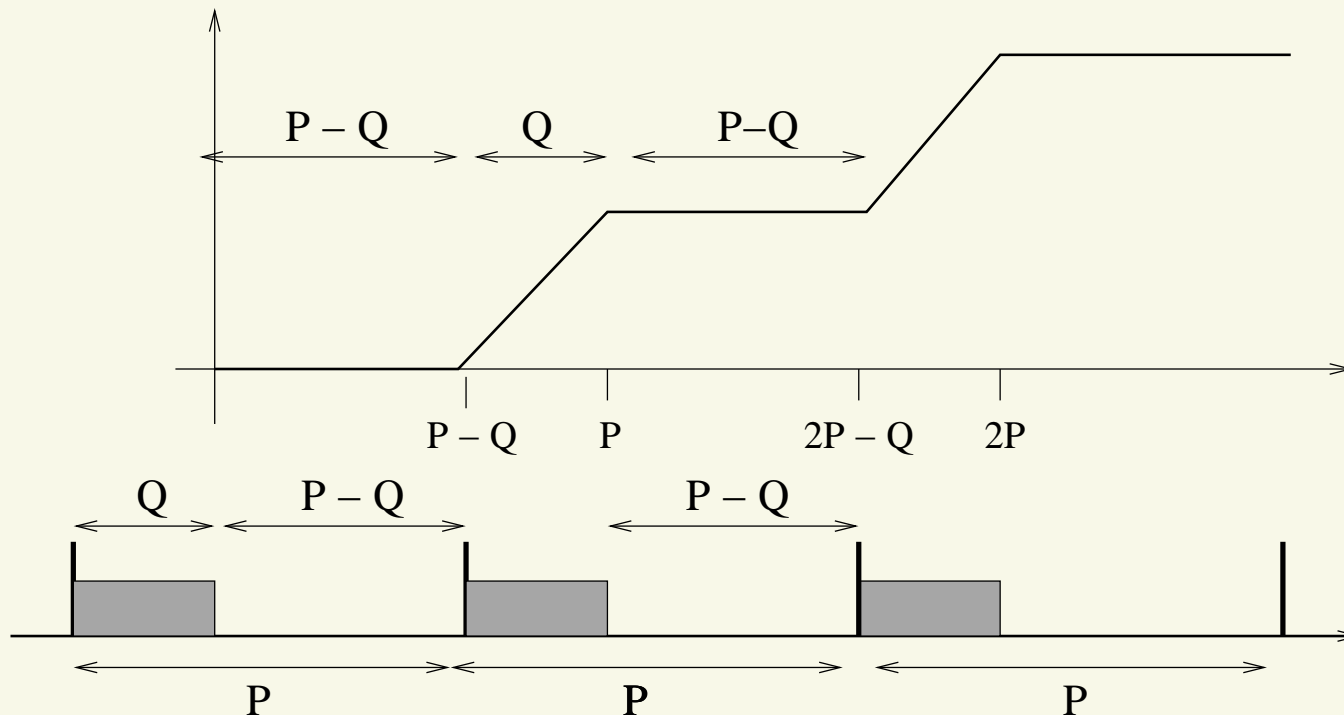  - Execution inside a period can even be preempted!

# Periodic Servers — Supplied Time

- $sbf(t)$: minimum amount of time that a VM is guaranteed to receive in a time interval of size $t$
  - Consider all the possible intervals of size $t$...
    - As already seen for static time partitioning
  - ...And all the possible "legal CPU allocations" generated by the periodic server!
- Big difference with static time partitioning: consider all the possible allocations of $Q$ in the period

# The Wrong Solution

- Immagine $Q$ is allocated at the beginning of the period
  - Worst case allocation: $t0$ immediately after $Q$
  - The time interval starts when the root scheduler deschedules the component
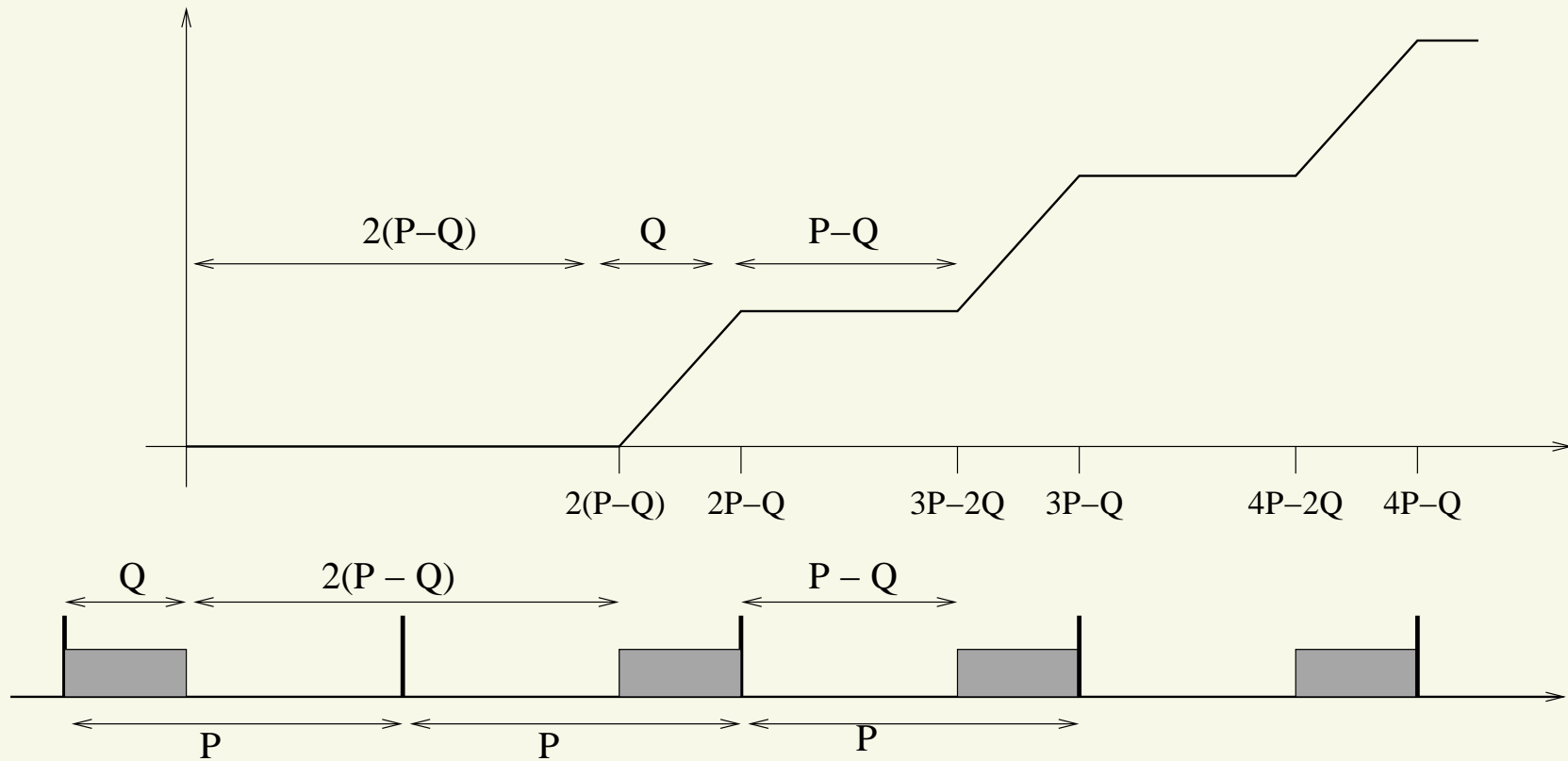
- Supplied time: $0$ until $P - Q$...
- ...Then increases with slope $1$ until $P$...
- ...Then flat again until $2P - Q$...
- ...

$$sbf(t) = \begin{cases} 0 & \text{if } t < (P - Q) \\ (n - 1)Q & \text{if } (n - 1)P \leq t < nP - Q \\ t + nQ - (n - 1)P & \text{if } nP - Q \leq t < nP \end{cases}$$

# Why Wrong?

- The previous computation assumed $Q$ always at the beginning of a period...
- ...But this is not the worst case!
  - Think about the second period...
  - ...What happens if the root scheduler delays the allocation?
  - The initial "$0$ allocation period" increases!!!
- Worst-case schedule: $Q$ at the beginning of the first period and at the end of the second one
  - See the difference with static time partitioning?

# Considering the Worst-Case Situation



$$sbf(t) = \begin{cases} 0 & \text{if } t < 2(P - Q) \\ (n - 1)Q & \text{if } nP - Q \leq t < (n + 1)P - 2Q \\ t - (n + 1)(P - Q) & \text{if } (n + 1)P - 2Q \leq t < (n + 1)P - Q \end{cases}$$
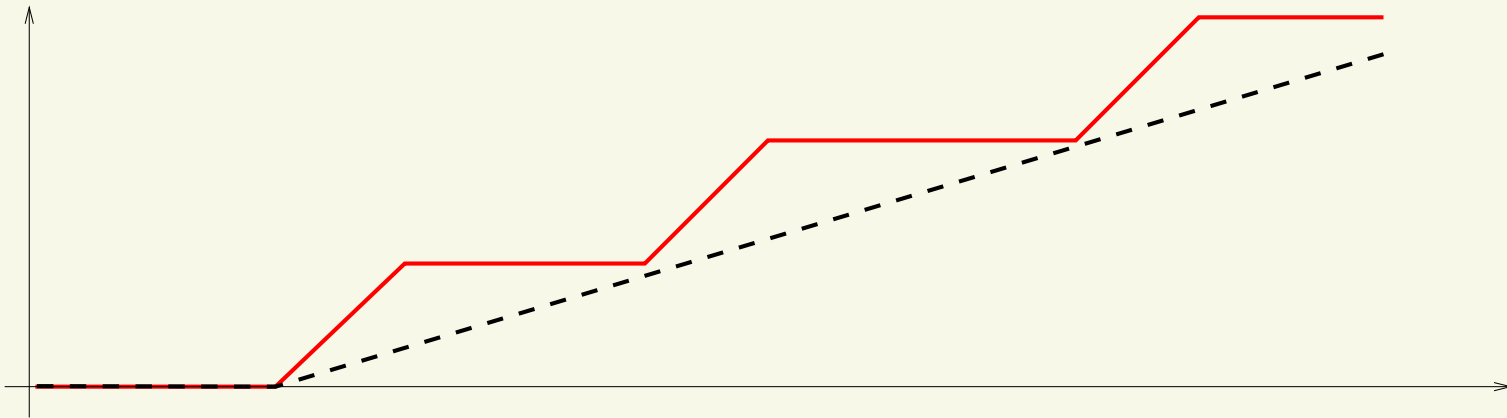
# Understanding the Supplied Bound Function

- Supplied bound function $sbf(t)$: minimum amount of time that a VM is guaranteed to receive in a time interval of size $t$

  - Considers all the possible intervals of size $t$...

- Strange looking function!

  - Flat for large intervals of time...
  - $\frac{\delta sbf(t)}{\delta t} = 1$ in the other intervals

- Can we "summarise" it with something simpler?
- What about a line ($y = ax + b$)?

  - $sbf(t) < 0$ makes no sense...
  - So, better $sbf(t) = max\{0, at + b\}$

- $sbf(t) = max\{0, at + b\}$... $at + b$ is below $0$ for $t < -b/a$

  - Let's rewrite the equation... $at + b = a(t - \Delta)$ with $\Delta = -b/a$

$$sbf(t) = \begin{cases} 0 & \text{if } t < \Delta \\ a(t - \Delta) & \text{otherwise} \end{cases}$$

# Interpreting the Linear Approximation

- $t < \Delta \Rightarrow sbf(t) = 0$: $\Delta$ is the *allocation delay* for the VM

  - Worst-case delay between the VM becoming active and the root scheduler scheduling it
  - How much time should I wait before the root scheduler starts giving the CPU to my VM?

- $a$ (sometimes referred as $\alpha$) is the *bandwidth* of the VM

  - Minimum fraction of CPU time reserved for the VM <span style="color:red">after the initial delay</span>

- Of course, $(a, \Delta)$ should be so that $a(t - \Delta)$ is below the real $sbf()$

# Periodic Servers Revisited

- How to compute $(\alpha, \Delta)$ for a periodic server $(Q^s, T^s)$?
  - $\alpha = \frac{Q^s}{T^s}$, $\Delta = 2(T^s - Q^s)$
- So, after the initial delay $2(T^s - Q^s)$ the VM is really receiving the expected fraction of CPU time $(Q^s/T^s)$
  - If we reduce $T^s$ (keeping $Q^s/T^s$ unchanged)...
  - $...sbf(t)$ tends to the "fluid allocation"!
- Why not using very very small server periods?
  - Of course there is a reason...

## The Design Problem

- Given a component (set of tasks and a local scheduler)...
  - Described by a time demand function (workload for fixed priorities)
- ...Find a root scheduler (and scheduling parameters) able to respect the components' temporal constraints
- Problem reduced to solving "$sbf(t) \geq dbf(t)$" for a set of points
  - Must be verified for all the points in case of EDF
  - Must be verified for at least one point in case of fixed priorities

# Simplified Design

- $sbf(t) \geq dbf(t)$
- Using $sbf(t) = a(t - \Delta)$...

$$a(t - \Delta) \geq dbf(t) \Rightarrow \Delta \leq t - \frac{dbf(t)}{a}$$

- Solve this for every $(t, dbf(t))$, and plot the solution on a $a - \Delta$ plane...
- ...Then compute the intersection (for EDF) or union (for fixed priorities)