Definition of real-time	Real-time task model 0000	Example 00000	Architectural aspects	Classification 000000000000000

# Real-time Systems Introduction

#### Tullio Facchinetti <tullio.facchinetti@unipv.it>

20<sup>th</sup> November, 2024

http://robot.unipv.it/toolleeo



which, among the following tasks, can be defined as a real-time task?

- the periodic sampling of a sensor every 30 minutes
- sending/receiving a TCP/IP packet every  $10\mu s$  (in average)
- the framerate of a videogame, which is 60 FPS
  - what would you say if it is 12 FPS?
- the video stream of a video-conference is 25 FPS, but sometimes some frames are skipping/freezing
- the update of financial information have frequency of 20ms
  - what would you say if it is 20 seconds?
  - and if it is 20 minutes?
- opening a program (software) takes 10 seconds from the mouse click, on average



the correct behavior of a real-time system depends not only on the logic correctness of the result produced by the calculation, but also from the time at which the result is made available

- in general, there are temporal constraints that must be satisfied
- temporal constraints depend on the interaction between the digital system and the process
- a result that is numerically correct but that is provided too late may cause problems as big as a wrong result



- the concept of "fast" is relative to the process
- the same system may be fast in some environments while being slow in others
- when (actually, always) more tasks must be managed by the same computing system, real-time systems aim at achieving the temporal constraints of each one
- fast systems usually aim to provide a low average response time on the set of tasks (all computer benchmarks are based on average values)
- guarantees on average values can not be trusted when individual temporal constraints must be achieved





- The cache memory is a hardware technology implemented in most of computers
- It improves the **average** access time to the information stored in memory, to obatin a faster execution of programs
- It works by maintaining a copy of information (caching) in a memory located between the processor and the RAM





- Cache memory uses a significantly more expensive technology (Static RAM) than RAM (Dynamic RAM)
- The cost depends on the density of circuits, and thus on the area of the chip
- More expensive technology implies that fetching information from the cache memory is much faster than from RAM
- The size of cache memory is kept much smaller than RAM to limit the cost of processors

Definition of real-time	Real-time task model	Example 00000	Architectural aspects	Classification 00000000000000
RAM vs cache				



	RAM	cache
Size (Intel Core i7-5775C)	up to 16 GB	6 MB
Speed [ns]	60 - 120	1 - 1.5
Cost [\$ / mm <sup>2</sup> ]	30 - 50	500 - 1,500

- Several levels of cache are usually implemented to improve the trade-off among speed, size and cost
- E.g., Intel Core i7-5775C: L1: 4x64KB, L2: 4x256KB, L3: 6MB





- The information stored in RAM can not fit within the cache memory
- The cache controller decides the information to keep within the cache (many sophisticated policies have been developed)

The data are exchanged when necessary between cache and RAM in order to keep data sync'ed.





When the CPU looks for some data in memory:

- cache hit: the information can be fetched from the cache memory
- cache miss: the information is not cached, and must be retrieved from RAM

The efficiency of the cache memory technology is based on the **data locality principle**: many subsequent memory accesses address the same variables.



- tram : time to access data stored in RAM
- t<sub>c</sub> : time to access data stored in cache memory
- for n >> 1 memory accesses there are
  - nh cache hits
  - n<sub>m</sub> cache misses
- total memory accesses:  $n = n_h + n_m$

Cache is faster than RAM:

 $t_c << t_{ram}$ 

Principle of locality:

 $n_h >> n_m$ 



- $n_h$ ,  $n_m$  and n: number of hits, misses, and total accesses
- $t_{ram}$  and  $t_c$ : time required to access RAM and cache

$$\begin{array}{c|c} & t_{avg} & t_{max} \\ \hline \text{cache} & (1/n)(t_c n_h + (t_{ram} + t_c)n_m) & t_{ram} + t_c \\ \text{no cache} & t_{ram} & t_{ram} \end{array}$$

- the cache memory achieves  $t_{avg} \rightarrow t_c \; (<< t_{ram}) \; \text{for} \; n_m \rightarrow 0$ (best average performance)
- the access time of the cache in the worst conditions (t<sub>max</sub>) is higher than without cache

The performance may decrease in the worst case w.r.t. the absence of cache.

Definition of real-time	Real-time task model ●000	Example 00000	Architectural aspects	Classification 00000000000000
Real-time task m	odel			



#### absolute parameters

- $r_i$  release time (or arrival time  $a_i$ )
- $s_i$  start time
- $d_i$  absolute deadline
- $f_i$  finishing time

#### relative parameters

- C<sub>i</sub> worst-case execution time (WCET)
- $D_i$  relative deadline

Definition of real-time	Real-time task model 0●00	Example 00000	Architectural aspects	Classification
Criticality				

# tasks can be classified as:

- hard: the violation of a temporal constraint may have "catastrophic" consequences on the system
- soft: a limited number of violations of temporal constraints bring to a tolerable degradation of the performance

# examples of (typical) hard/soft tasks

- hard: sensor sampling, actuator operation, control loops, communication
- soft: user I/O (keyboard input, messaging, ecc.), multimedia streaming

a computing system able to deal with hard tasks is said hard real-time



tasks can be associated to different kind of constraints

#### temporal constraints

• activation instant, finishing time, jitter, ...

#### precedence constraints

• refers to the temporal ordering between task

#### resource constraints

• mutual exclusion, synchronization

Definition of real-time	Real-time task model 000●	Example 00000	Architectural aspects	Classification 00000000000000
Temporal constra	ints			

### explicit constraints

- they are explicitly stated by the problem formulation
- e.g.: "the sensor must be sampled every 20ms"

# implicit constraints

- they must be inferred from the problem parameters, where they do not appear explicitly
- e.g.: "the sensor shall detect an object having length L = 10cm travelling at max speed of  $V_{max} = 1$ m/s"

Definition of real-time	Real-time task model	Example ●0000	Architectural aspects	Classification 00000000000000
- · · ·				

#### From physical to temporal parameters



- S<sub>a</sub> : minimum braking distance to stop safely
- S<sub>d</sub> : maximum sensing range of the sensor

Definition of real-time		ie	Real-time task model		Example 00000	Architectur 00000	Architectural aspects	Classification 00000000000000		
_										

#### From physical to temporal parameters



- r<sub>i</sub> : time when the sensor detects the obstacle
- d<sub>i</sub>: latest instant to start braking to stop without hitting the obstacle (D<sub>i</sub> = d<sub>i</sub> r<sub>i</sub>)
- C<sub>i</sub> : duration of the task controlling the braking system

Definition of real-time	Real-time task model 0000	Example 00●00	Architectural aspects	Classification 000000000000000

# From physical to temporal parameters

#### dynamic friction force:

$$F = -\mu mg$$

where

- $\mu$  : friction constant
- m : the vehicle mass

deceleration (Second Newton's Law):

$$a = F/m = -\mu g$$

law of the uniform acceleration (kept it simple):

$$egin{aligned} \mathsf{x}(t) &= \mathsf{v}t - rac{1}{2}\mu gt^2 \ \mathsf{v}(t) &= \mathsf{v} - \mu gt \end{aligned}$$

Definition of real-time	Real-time task model	Example 00000	Architectural aspects	Classification 00000000000000
From physical to t	temporal paramete	ers		

the system must be evaluated in worst conditions, i.e., when the speed equals the maximum speed  $v_{max}$ 

at time  $t_a$  the vehicle stops, i.e.,  $v(t_a) = 0$ ; therefore

$$t_a = rac{v_{max}}{\mu g}$$

by substitution within the first equation:

$$x(t_a) = S_a = \frac{v_{max}^2}{2\mu g}$$

 Definition of real-time
 Real-time task model
 Example
 Architectural aspects
 Classification

 0000000000
 0000
 0000
 00000
 00000
 00000000000

#### From physical to temporal parameters

#### known quantities:

- S<sub>a</sub> : just calculated
- $S_d$  : depends from the adopted sensor
- C<sub>i</sub> : known, since the implemented task has a given worst-case duration under the considered computing platform

the maximum time to travel the distance  $S_d - S_a$  in the worst case (assuming  $S_d > S_a$ , otherwise.... CRASH!!):

$$D_i = \frac{S_d - S_a}{v_{max}}$$

it is possible to state whether the timing required for the computation is suitable for the dynamics (timing constants) of the system, or tailor the computing parameters accordingly

Definition of real-time	Real-time task model 0000	Example 00000	Architectural aspects •0000	Classification
Sources of tempor	al non-determinis	m		

#### computing architecture

• cache, pipelining, interrupts, DMA

#### operating system

• scheduling, synchronization, communication

#### programming languages

• in most languages there is no explicit support for the specification of temporal constraints

# design methodologies

• lack of methods and tools for analysis and verification

Definition of real-time	Real-time task model	Example	Architectural aspects	Classification
	0000	00000	○●○○○	00000000000000
Solutions and dra	wbacks			

#### common solutions:

- low level programming (assembler)
- timing imposed by explicit manipulation of hardware/software timers
- processing within device drivers (faster response time)
- explicit manipulation of task priorities

# drawbacks:

- time-consuming development, strongly dependent on programmer's skills
- poor clarity of code
- reduced portability
- prone to errors



a sequence of assignments of tasks to the processor



formally, given a task set  $\Gamma = \{\tau_1, \ldots, \tau_n\}$ , a schedule is a function  $\sigma : \mathbb{R}^+ \to \mathbb{N}$  such that, for each time interval  $[t_1, t_2)$ , it holds

$$\sigma(t) = \begin{cases} k & \text{if } \tau_k \text{ is running at time } t \in [t_1, t_2) \\ 0 & \text{if the processor is idle} \end{cases}$$



- the task set is composed by 4 tasks
- $\sigma = 0$  when no tasks are running



Definition of real-time	Real-time task model	Example	Architectural aspects	Classification
	0000	00000	0000●	00000000000000
The ready queue				



- the ready queue contains tasks ready to execute
- ordered in descending priority order, such that the highest priority task will be the first to be fetched
- the scheduler is the OS component that inserts a task into the ready queue using some given criteria (called scheduling algorithm)

Definition of real-time	Real-time task model 0000	Example 00000	Architectural aspects	Classification •000000000000000000000000000000000000
T	d. Para alternitidades			

Faxonomy o	f sched	luling a	gorithms
------------	---------	----------	----------

preemptive	VS.	non-preemptive
static	VS.	dynamic
on-line	VS.	off-line
optimal	VS.	best effort
time-triggered	VS.	event-triggered

for each class, an algorithm is classified in either one or the other category

#### example

• EDF is a preemptive, dynamic, on-line and optimal algorithm

Definition of real-time	Real-time task model	Example 00000	Architectural aspects	Classification 000000000000000000000000000000000000
Preemptive vs. no	on-preemptive			

preemptivity refers to the possibility (or not) to temporary stop the running task

#### preemptive

• the kernel can stop the running task to execute one having higher priority

#### non-preemptive

• once the execution of a task has started, it can not be stopped by the kernel until its execution has ended (or self-suspended)



algorithms have different approaches regarding the assignment of priorities

#### static

- the priority of a task is based on static parameters
- the value of a static parameter does never change

# dynamic

- the priority of tasks is based on dynamic parameters
- the value of a dynamic parameter continuously changes over the time



# they differ regarding **when** scheduling decisions are taken

#### on-line

- scheduling decisions are made at run-time
- they are based on changing parameters, such as task activation time and temporal constraints

# off-line

- scheduling decisions are taken before starting the system
- they are suitably stored and applied at run-time



they optimize (or not) some cost function

#### optimal (regarding schedulability)

• if a feasible schedule exists, the scheduling algorithm will find it

# best effort

• a suitable heuristic method is used to obtain the best possible result



they have different policies regarding the activation of tasks

## time-triggered

- a timer triggers the activation of task on a regular basis
- periodic tasks

#### event-triggered

- a task is triggered by an external event
- aperiodic/sporadic tasks

examples of tasks are:

- time-triggered: periodic sampling of sensors
- event-triggered: interrupt management (e.g., keyboard input)

Definition of real-time	Real-time task model 0000	Example 00000	Architectural aspects	Classification 00000000000000
Task scheduling				

#### feasible schedule

the schedule achieves the temporal constraints of all tasks

#### examples:

- every task terminates within its deadline
- the precedence order of tasks is guaranteed

given a task set, if a feasible schedule exists, then the task set is said **feasible** 

NOTE: feasibility is a **property of the schedule** only; it is independent from the adopted schedulig algorithm

Definition	of	rea	l-ti	me
0000000	00	00		

Real-time task model

Example 00000 Architectural aspects

Classification 00000000000000

#### Schedulability test

Given a scheduling algorithm, the schedulability test is a formula or algorithm that allows to check the guarantee of timing constraints of each task.

- The schedulability test **depends from the algorithm**, i.e., different algorithms have different schedulability tests.
- It can be applied **before** running the scheduling algorithm to know whether the task set is schedulable or not.

 Definition of real-time
 Real-time task model
 Example
 Architectural aspects
 Classification

 0000000000
 00000
 00000
 00000
 00000
 000000

#### The general problem of scheduling

given a **task set**  $\Gamma$  and a **set of resources** R(processors, variables, peripherals, etc.), the problem of finding a feasible schedule for  $\Gamma$  is said general problem of scheduling

- the general problem of scheduling has NP-hard complexity
- in presence of simplifying assumptions, it is possible to fomulate algorithms having polynomial complexity

Definition of real-time	Real-time task model 0000	Example 00000	Architectural aspects	Classification 000000000000000000000000000000000000
Notes on the com	plexity			

The complexity of an algorithm is related to the number of elementary steps that are required to complete.

- The number of elementary steps determines the time required to complete.
- The complexity is expressed as a function of the problem parameters.



The complexity is typically expressed using the so-called Big O notation, indicated as O(.).

Mathematically speaking, "f(n) is O(g(n))" iff for some constants c and  $N_0$ , it holds  $f(N) \le cg(N)$  for all  $N > N_0$ .

Informally:

- The Big O notation indicates how fast a function increases.
- For this reason, the Big O notation contains the term of the original function that makes it to increase faster (without constants).
- Example: the function  $n^2 + 3n + 1$  is  $O(n^2)$ .





**Big-O Complexity Chart** 

Elements

From https://www.freecodecamp.org/news/ big-o-notation-why-it-matters-and-why-it-doesnt-1674cfa8a23c/



in real-time systems, the typical parameter is the number of tasks n

- O(1): does not depend on n (i.e., it is constant)
- O(n): linearly depends on n
- $O(n^4)$ : polynomially depends on n
- $O(4^n)$ : exponentially depends on n

let's assume n = 50 and a duration of the elementary step of 1ms

- O(1): does not depend on *n*, e.g., it is equal to 20ms
- O(n): milliseconds (e.g., 50ms)
- $O(n^4)$ : hours (e.g., 1 hour, 44 min e 10 sec)
- $O(4^n)$ : billions of billions of years (e.g., 40 b.b.y.)

Definition of real-time	Real-time task model 0000	Example 00000	Architectural aspects	Classification 0000000000000
Common simplify	ing assumptions			

- Only one processor
- est of tasks having common features
- I full preemptive systems
- simultaneous activations
- o precedence constraints
- no shared resources
- o no overhead (context switch, cache issues, etc.)