



Computing Infrastructures

January 13, 2025

Course Section: Prof. Ardagna Prof. Palermo Prof. Roveri

Student ID (Codice Persona):

Last Name:
(LAST NAME IN CAPITAL LETTERS)

First Name:
(FIRST NAME IN CAPITAL LETTERS)

Exam Duration: 1hour and 30min

Students are not permitted to use mobile phones and similar connected devices. Course materials and programmable devices (e.g. programmable calculators) cannot be used as well. **Any violation of the rules is considered a cheating action.**

Answers must be given on the Answer Sheets and in English. Any box filled or answer provided on the other sheets will be ignored. Students must use a pen (black or blue) to mark the answers (no pencil).

Write the LAST and FIRST name in CAPITAL LETTER, and in this order, in all places where requested. **Where it is requested only the STUDENT ID (Codice Persona), do not write your name.**

Check that the first number of the code for the Answer Sheet is the same as for the other sheets. The code can be found in the top-right corner of each page in the form +NN/KK/XX+. The parts that should correspond is ONLY the first digit NN.

Mark clearly the box corresponding to your answers, without overlapping on other boxes. If you make a mistake on them, circle the word *Question* together with the related number, and write the correct letter to its side.

Numerical exercises require writing the formulas and procedure used to solve the problem just after the question in the left space. Exercises without the procedure used to reach the result will not be considered for the evaluation. Only the numeric answer and its unit should be reported on the corresponding dotted line in the Answer Sheet.

The answers to the *Open Questions* should be written using ONLY the space available on in the boxes within the Answer Sheets. The answers should be readable by the professor. Unreadable answers will not be considered for the evaluation.

Scores: correct answers take positive points, unanswered questions take 0 points, **wrong answers can have negative points.** An indication of the points is available at the beginning of each section. The final score can be re-modulated at the end of the evaluation.

**True false questions**

Correct answer: +1, No answer: 0, Wrong Answer -0.5

Answers must be given on the ANSWER SHEETS. Any box filled here will be ignored. Pay attention to the position (A or B) of the True/False answers, since they are not always in the same position.

Question 1 SANs are primarily used for block-level access to data, while NAS devices provide file-level access.

A False

B True

Question 2 RAID can be used with both hardware and software implementations.

A False

B True

Question 3 GPUs (Graphics Processing Units) are primarily used in datacenters for accelerating graphical applications and gaming.

A False

B True

Question 4 Cloud computing applications do not require any server to run.

A True

B False

Question 5 Availability zones within a Compute Region can have a roundtrip time larger than few milli seconds.

A False

B True

Question 6 Edge computing is a distributed computing paradigm that brings computation and data storage closer to the location where it is needed.

A False

B True

Question 7 The power consumption of an average datacenter is in the order of a few GigaWatts.

A False

B True

Question 8 Leaf-spine topologies are more suitable for small-scale datacenter deployments with limited server connections.

A False

B True

Question 9 Warehouse-scale computers are typically composed of a large set of homogeneous servers.

A False

B True

Question 10 Virtualization can help simplify IT management by reducing the number of physical machines that need to be managed.

A True

B False



Exercises

Correct answer: +2, No answer: 0.

The formulas and procedures used to solve the exercises should be included here close to the question. The numeric answer, and only that, must be given on the ANSWER SHEETS. Any number written only here will be ignored. The correct number is ONLY a necessary condition for a correct answer. If the formulas are not available after each exercise, they will be considered as not answered.

Question 11

A scientific computation that needs to be carried out within the PoliMi data center uses a server composed of 2 CPUs and 4 GPUs. Knowing that:

- The computation takes 8 days to complete,
- The computation requires at least one CPU and all GPUs within the server to be operational to complete successfully;
- $MTTF_{CPU} = 180$ days and $MTTF_{GPU} = 120$ days.

How many parallel instances of the computation must be launched to ensure a probability higher than 98% that at least one computation produces results successfully? Notes: (i) Use at least 4 decimal places for all intermediate calculations. (ii) All other components of the server can be considered ideal.

$$\begin{aligned}
 R_{CPU}(8) &= e^{-\frac{8}{180}} = 0.9565 \\
 R_{GPU}(8) &= e^{-\frac{8}{120}} = 0.9355 \\
 R_{\text{instance}} &= [1 - (1 - 0.9565)^2] \times 0.9355^4 = 0.7645 \\
 1 - (1 - 0.7645)^n &= 98\% \rightarrow n = 3 \\
 0.2355
 \end{aligned}$$

Question 12

If the time required to recover and replace a component of the server described in the previous exercise (whether it is a CPU or a GPU) is equal to 24 days, what is the total availability of the server? Notes: (i) Use at least 5 decimal places for all intermediate calculations. (ii) All other components of the server can be considered ideal.

$$\begin{aligned}
 A_{CPU} &= \frac{180}{180+24} = 0.88235 \\
 A_{GPU} &= \frac{120}{120+24} = 0.83333 \\
 [1 - (1 - 0.88235)^2] \times 0.83333^4 &= 0.4756
 \end{aligned}$$

**Question 13**

A company is planning to use a RAID 5 array composed of 8 disks for critical data storage. The desired Mean Time to Failure for the entire RAID system ($MTTF_{RAID5}$) is 16 years. Each disk in the array has a $MTTF_{disk}$ equal to 500 days. What should be the MTTR required *in hours* to meet the target RAID $MTTF_{RAID5}$.

Question 14

Consider a system composed of three stations: the CPU that is characterized by $V_{CPU} = 100$ visits and an average service time of $S_{CPU} = 10\text{ms}$; the disk, characterized by a throughput of 12 IOPS, and a demand of $D_{DISK}=150\text{ms}$; and the GPU whose demand is $D_{GPU}=40\text{s}$ and the number of visits $V_{GPU}=10$. Finally, the system throughput is $X=20\text{ jobs/min}$ while the response time when there are $N = 20$ end-users in the system is $R = 55\text{s}$.

Compute the CPU demand and the GPU throughput.

Write in the answer sheet: $D_{CPU} = \dots; X_{GPU} = \dots$

$$D_{CPU} = V_{CPU} \cdot S_{CPU} = 100 \times 10 = 1\text{s}$$

$$D_{GPU} = V_{GPU} \cdot S_{GPU} \rightarrow S_{GPU} = \frac{D_{GPU}}{V_{GPU}} = \frac{40}{10} = 4\text{s}$$

$$X_{GPU} = X \cdot V_{GPU} = 20 \times 10 = 200 / 60 = 3.33\text{s}$$

**Question 15**

Considering the system described in Question 14, what is the users' think time Z?

$$D_{CPU} = V_{CPU} \cdot S_{CPU} = 100 \times 10 = 1s$$

$$D_{GPU} = V_{GPU} \cdot S_{GPU} \rightarrow S_{GPU} = \frac{D_{GPU}}{V_{GPU}} = \frac{40}{10} = 4s$$

$$X_{GPU} = X \cdot V_{GPU} = 20 \times 10 = 200 / 60 = 3.33s$$

$$R = \frac{N}{X} - z \quad 55 = \frac{\cancel{20}/\cancel{60} 20}{20/60} - z \Rightarrow z = 5s$$

Question 16

Based on the system in Question 14, the number of end users is predicted to reach 40 in one month. Considering the response time lower bound, which option is better?

- Upgrade the system by adding one more GPU (you can assume the new GPU is equal to the one initially available and to balance evenly the GPU processing).
- Replace the GPU with one 2.5 times faster than the original one.

Provide the estimated bounds in the two scenarios to motivate your answer.

Write in the answer sheet: A or B ; $R_{LOW}^A = \dots$; $R_{LOW}^B = \dots$

Case A: $40s \rightarrow D_{GPU}' = 20s$

$$D_{GPU} = 40s \rightarrow D_{GPU} = 20s$$

$$R_A = N \cdot D_{GPU}' - z$$

$$\therefore = 40 \times 20 - 5 = 795s$$

Case B: $D_{GPU}' = \frac{D_{GPU}}{2.5} = 16s$

$$R_B = N \cdot D_{GPU}' - z = (6 \times 40) - 5$$

$$= 635s$$

$\approx 5s$

So Case B is better.



Open Questions

Correct answer: +5, No answer: 0. Points are modulated considering the written text

Write the answer using ONLY the space available in the boxes on the ANSWER SHEETS. The answers should be readable by the professor. Unreadable answers will be considered wrong.

Question 17

⇒ What are the advantages and disadvantages of using Type 1 hypervisors versus Type 2 hypervisors in a virtualized environment?

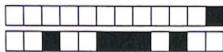
Question 18

⇒ What are the implications of hardware heterogeneity on software stack development in datacenters?

!!!ANY ANSWER PROVIDED ON THIS PAGE WILL BE IGNORED!!!

If needed, you can use the space hereafter to organize your answer.

Type 1 : Runs on bare metal machine Advantages : 1. Runs on bare metal machine , which means higher performance 2. The security is strong 3. Use case = enterprise , production 4. Setup complexity : complex
Type 2 : Runs on host operating system Advantages : 1. Easy to install and setup



Computing Infrastructures - January 13, 2025

Answer Sheets (Page 3)

Student ID (Codice Persona):

True/False Questions

Question 01 : A BQuestion 02 : A BQuestion 03 : A BQuestion 04 : A BQuestion 05 : A BQuestion 06 : A BQuestion 07 : A BQuestion 08 : A BQuestion 09 : A BQuestion 10 : A B

Exercises

3

Question 11 :

0.4756 σ 0.8558

Question 12 :

18,35 h (σ 0.7644 giorni)Question 14 : $D_{CPU} = 1 \frac{sec}{sec}$ $X_{GPU} = 200 \frac{sec}{mm} \sigma 3.33 \frac{sec}{sec}$ Question 15 : $Z = 5 \frac{sec}{sec}$ Question 16 : B ; $R_{MIN}^A = 785 \frac{sec}{sec}$; $R_{MIN}^B = 635 \frac{sec}{sec}$