

TDDD93/TEN2 – Large-scale distributed systems and networks

Final Examination: 8:00-12:00, Thursday, Aug. 17, 2017

Time: 240 minutes

Total Marks: 40

Grade Requirements: Three (20/40); four (28/40); and five (36/40).

Assistance: None (closed book, closed notes, and no electronics)

Instructor: Niklas Carlsson

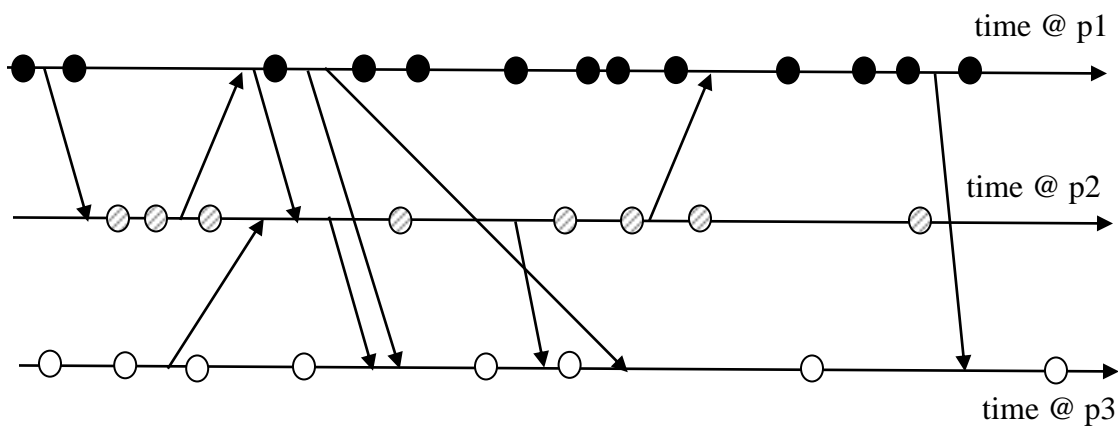
**Instructions:**

- Read all instructions carefully (including these)!!! Some questions have multiple tasks/parts. Please make sure to address *all* of these.
- The total possible marks granted for each question are given in parentheses. The entire test will be graded out of 40. This gives you 10 marks per hour, or six minutes per mark, plan your time accordingly.
- This examination consists of a total of 12 questions. Check to ensure that this exam is complete.
- When applicable, please explain how you derived your answers. Your final answers should be clearly stated.
- Write answers legibly; no marks will be given for answers that cannot be read easily.
- Where a discourse or discussion is called for, be concise and precise.
- If necessary, state any assumptions you made in answering a question. However, remember to read the instructions for each question carefully and answer the questions as precisely as possible. Solving the *wrong* question may result in deductions! It is better to solve the *right* question incorrectly, than the *wrong* question correctly.
- Please write your AID number, exam code, page numbers (even if the questions indicate numbers as well), etc. at the top/header of each page. (This ensures that marks always can be accredited to the correct individual, while ensuring that the exam is anonymous.)
- Please answer in English and utilize figures and tables to the largest extent.
- If needed, feel free to bring a dictionary from an official publisher. Hardcopy, not electronic!! Also, your dictionary is not allowed to contain any notes; only the printed text by the publisher.
- Good luck with the exam.

## Part A: Distributed Systems

### Question 1 (4 points)

Assume that you have three processes p1, p2, and p3 which are implementing Lamport's clocks. There are many events that take place at these processes, including some messages being sent between the processes. In the figure below we use circles and arrows to specify in-processor events and messages being sent between processes, respectively. Please provide the logical timestamps associated with each event. You can assume that all three clocks start at zero, at the left-most point in time. (Also, explain how the processes would adjust their clocks if using Lamport's logical clocks.)



### Question 3 (4 points)

Mutual exclusion. Consider a simple scenario in which there are five nodes: A, B, C, D, and E. Use a sequence of figures to illustrate and explain the message sequences and coordination between these nodes when node A acts as a central coordinator for a shared memory resources (that all five nodes can use) and both nodes B and C almost at the same time decides that they want to write to the resource. You can assume that each write access (to memory) takes 1 second and that there is 10ms between the times when B's and C's decisions, and that the round trip times between the nodes are random in the approximate range 60-120ms.

### Question 3 (2 points)

Transparency plays a central role in some distributed systems. Consider a simple multi-tier system with three levels: a user interface, an application server, and two replicated database servers. Within this context and example scenario, please give two concrete examples of two different types of transparency and explain how transparency is used here to provide improve service for the end users.

## Part B: Methodology

### Question 4 (4 points)

When designing experiments, it is important to carefully identify the most appropriate factors, levels, and metrics to consider. Consider a researcher wanting to assess the performance of a webserver. The researcher has identified three factors of interest: (i) the request rate, (ii) the job size, and (iii) the processor speed. For each of these factors, the researcher has identified 9, 8 and 7 levels of interest, respectively, including identified a default request rate, job size, and processor speed. Let us call the request rate levels  $R_1, R_2, \dots, R_9$ ; the job size levels  $S_1, S_2, \dots, S_8$ ; and the processor speed levels  $P_1, P_2, \dots, P_7$ . Please estimate the number of experiments that the researcher would need to perform if performing (a) one factor experiments with the default scenario as baseline, (b) two factor experiments with the default scenario as baseline, and (c) full factor experiments. Also, please explain which experiments would be performed in each case.

### Question 5 (3 points)

Assume that you are interested in understanding the Autonomous Systems (ASes) that your internet packets take between two end points. For example, let us assume that you have access to PlanetLab and wants to estimate the path between a node located in California (USA) and a node in Italy. For this scenario, please describe a methodology for how you can estimate the AS path packets takes as they go from the PlanetLab node in California to that in Italy. Note that direction is important here. Please provide a figure of your experimental setup, describe the tools you would use, and how you would use them to solve the task at hand. (Also, as a semi-bonus questions, please give “guestimates” of the number of ASes and routers along such an example path. Note that in practice, having reasonable “guestimates” help quickly sanity check your results.)

### Question 6 (3 points)

Consider a system with two states: “on” and “off”. Assume that the system is “on” whenever there are jobs to serve and the system instantaneously can go between the “on” and “off” states whenever a new job arrive to an empty system or when the system is done serving all jobs, respectively. Furthermore, assume that the system only can serve one job at a time (as with any G/G/1 queue system), on average 100 jobs/second arrive to the system, each job on average takes 5ms to serve, and each job stays in the system for on average 25ms.

- What is the system utilization?
- How many jobs are on average in the system?
- Assuming that the “on” state consumes 100 Watts and the “off” state 10 Watts. What is the average power consumption of the system, given the described workload and system characteristics?

## Part C: Multicore and Parallel Programming

### Question 7 (4 points)

Questions on parallel computer architecture.

- a) Sketch the network topology (graph structure) for the 2D-torus network. (Hint: annotated drawing). Analyze how its (A) average communication distance, (B) overall network cost (#links) and (C) node degree grow with the number  $N$  of nodes connected by the network (give commented formulas in  $N$ ). (Hint: 3 commented formulas in total. Assume here for simplicity that  $N$  is a square number and the torus is organized accordingly.) (2p)

If you do not recall the 2D-torus network, you may instead do it for any other interconnection network discussed in the lecture, with half the amount of points.

- b) What is a "heterogeneous (multicore) system"? What are its strengths over homogeneous multicore CPUs, and what additional challenges does it bring for the software? (1.5p)
- c) Explain how "hardware multithreading" can improve the utilization of a CPU core. (0.5p)

### Question 8 (3 points)

Questions on MPI/algorithm design. You have a cluster computer running  $P$  ( $P > 1$ ) MPI processes. Process 0 has an array  $A$  of  $N$  elements, e.g., float values. Each process also allocates in its main memory an array  $B$  of  $N$  elements.

- i) Design a parallel message-passing program (using all  $P$  nodes), using MPI-like explicit *send()* and *receive()* operations only, that broadcasts the array  $A$  to all  $P$  processes so that afterwards each process  $i$ ,  $0 \leq i \leq P-1$ , shall hold in its array  $B$  a copy of array  $A$ . (MPI or pseudocode is fine, explain your code.) Hint: Keep it simple, no optimality is required here. (1p)
- ii) Assume that the time for sending and receiving a block of  $K$  elements is  $a \times K + b$  for constant system parameters  $a > 1$  and  $b > 1$ , that each node can only send to or receive from one node at a time, and that there are direct network links between all nodes. Assume further that copying an element and other local operations on an element cost 1 unit of time. Derive the asymptotic time complexity for completing the entire broadcast operation according to your algorithm of (i) (i.e., derive a formula in  $N$ ,  $P$ ,  $a$  and  $b$ ; use big-O notation where appropriate, and explain your calculation). (1p)
- iii) How does your algorithm perform if  $N = 1$  while  $P$  grows very large? Describe the main idea of an alternative algorithm that is asymptotically faster, and explain why. (1p)

### Question 9 (3 points)

Question on theory.

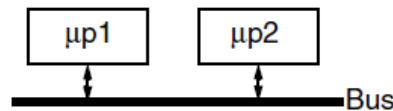
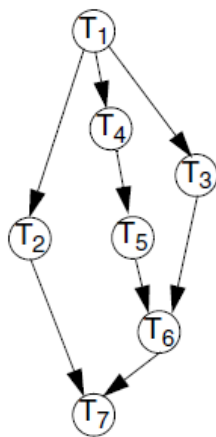
- a) Explain Amdahl's Law (including its assumptions) and give its proof, and explain its implications for where to focus one's efforts when parallelizing a sequential program. (2p)
- b) Give an example of a (parallel) speedup anomaly. (1p)

### Part D: Embedded Systems

Consider an application modelled as the task graph below. Each task, when activated, consumes one message on each input edge and generates, at termination, one message on each output edge. The task graph is executed on the architecture shown in the figure. Execution times of the tasks, when executed on the corresponding processor, are shown in the table. All messages transmitted over the bus, between tasks mapped on different processors, consume 2 time units to reach the destination. Communication between tasks mapped to the same processor is considered to not consume any time.

Propose an efficient task mapping (indicate on which processor each task is executed) and a corresponding static (nonpreemptive) schedule for the application. Illustrate your schedule as a Gantt chart (similar to the way we captured schedules in Lecture 1&2).

Try to achieve a maximum delay (the time interval between the start of T1 and the finishing of T7) of 46 time units!



Task	WCET	
	μp1	μp2
T <sub>1</sub>	5	6
T <sub>2</sub>	12	15
T <sub>3</sub>	10	11
T <sub>4</sub>	5	6
T <sub>5</sub>	3	4
T <sub>6</sub>	17	21
T <sub>7</sub>	10	14

#### Question 11 (3 points)

In the lectures we have particularly emphasized three design steps: architecture selection, task mapping, elaboration of a schedule. Explain, in short, what each step is doing. Illustrate the three steps by a small example.

#### Question 12 (3 points)

Think at the sources of power dissipation as we discussed at the lectures. What are main opportunities to reduce power consumption?

*Good luck!!*