



**CHALMERS**  
UNIVERSITY OF TECHNOLOGY



UNIVERSITY OF GOTHENBURG

# Exam

DAT278 / DIT055: Sustainable Computing

Saturday 15 January 2022, 14:00-18:00

---

**Examiner:**

Pedro Petersen Moura Trancoso

**Contact person during exam:**

Pedro Petersen Moura Trancoso, tel. 7726319

**Supporting Materials/tools:**

Chalmers approved calculator.

**Instructions:**

All answers should be written in English.

**Grading intervals:**

The maximum grade for this exam is 100 points

DAT278 (Chalmers):

Fail < 40p ≤ Grade 3 < 60p ≤ Grade 4 < 80p ≤ Grade 5

DIT055 (University of Gothenburg):

Fail < 40p ≤ Grade G < 75p ≤ Grade VG

**Examination solution:**

Will be posted in canvas within 24h after the exam has finished.

**Examination review session:**

A canvas announcement with the date and location will be sent to all students.

---

**Q1 Sustainability ((2+2)+6 points)**

---

- (a) Explain in your own words the terms “Sustainability” and “Sustainable computing” and
- (b) Discuss, using real examples, how the material in this course addresses these issues.

- (a) *In terms of “Sustainability” we understand as the approaches and processes that use the available resources with a focus on reducing the impact on the environment so that the availability and quality of the environment are assured for future generations. In terms of “Sustainable computing” we understand as the methods and technologies used to develop and/or use computer systems in a more efficient way thus again reducing the impact on the environment and available resources.*
- (b) *Relevant topics that were discussed in class are for example the topic of energy-efficient data center cooling, awareness of design for end-of-life, hardware technologies to improve the energy efficiency for large memories, and more.*

---

**Q2 Metrics ((3+3)+4 points)**

---

In the course we discussed several metrics that are relevant to analysing the efficiency such as power, energy, and different combined efficiency metrics like the energy-delay product (EDP).

- (a) Select from those two different metrics and justify for which different devices of the compute continuum (from the edge to the cloud) they apply to.
- (b) Explain how you would measure those two metrics selected in (a).

- (a) *Energy for a mobile device as for these devices one of the most important efficiency factors is the battery lifetime, and energy is the determinant factor for that. Energy-delay-product (EDP) for a system on the cloud as for such systems, even though energy is very relevant for their consumption and cooling, at the same time those systems need to sustain a required performance and thus the combined metric works best for those.*
- (b) *Energy can be measured with a device that is plugged to the wall power plug. Alternatively, it can be estimated from models and measurement of relevant performance counters – some software tools handle these measurements and calculations on real hardware devices. EDP requires Energy that can be done as mentioned above and time which can be measured with a time function introduced in the program being executed or the call of time from the shell when execution on a terminal or a Linux shell for example. The EDP is then obtained from the multiplication of the Energy and the Time.*

---

**Q3 Technology/Circuits (5+5 points)**

---

- (a) Explain in your own words what Dennard scaling is and how it relates to the sustainable computing issues discussed in this course.
- (b) In this course, among many other, we presented the following two techniques to improve efficiency: Clock gating and Power gating. Which of these two techniques is better in reducing the static energy? Justify your answer, briefly describe the implementation for this technique and discuss the advantages and disadvantages of the technique.

- (a) Dennard scaling refers to the fact that the power density maintains constant as we scale the chip technology. This is not valid any longer with the latest reductions in the technology point.*
- (b) Power gating as it "cuts" the power from the circuit thus reducing completely the power for those components (static and dynamic). While it is very effective, it is difficult to implement in hardware and also it is more time consuming to enable and disable so it is important to be sure that the time the circuit can be off is enough to compensate for all overheads.*

---

**Q4 Dynamic Voltage-Frequency Scaling (DVFS) (2+8 points)**

---

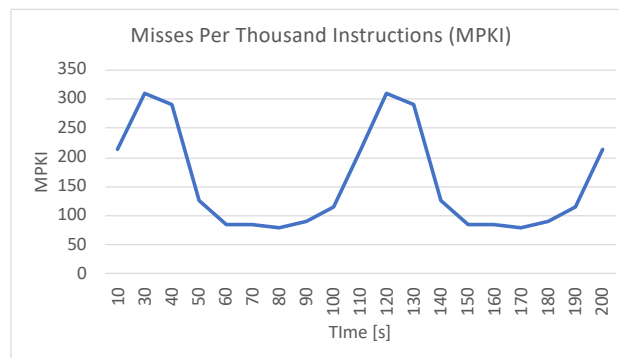
- (a) Briefly explain what the DVFS technique is.
- (b) Assume now that you have a multithreaded program and focus on a parallel section of this program where several threads are executing different iterations of a loop. Assume also that different iterations have different execution time.
- Discuss how DVFS implemented at the software level (programmer is given an API – set of functions – to control directly the voltage-frequency) could improve the efficiency of the parallel execution.
  - Describe in your own words how this can work - use your own small example and show how DVFS works for your example, which support is needed, which analysis needs to be done and which actions need to be taken at the software level.

- (a) DVFS is a technique where we change the voltage and frequency. Voltage and frequency need to be changed together and high frequencies require higher voltages and lower frequencies can be operated under lower voltages.*
- (b) Depending on the criticality of a thread the programmer could introduce class to have the non-critical threads run at slower voltage-frequency states. If the performance of the different threads is known at compile time, the code for controlling the DVFS can be introduced statically by the programmer and/or the compiler. If on the other hand the performance of the threads is only determined or is dependent of factors that are only happening during runtime, that will mean that the programmer will have to include in the code calls to monitor the performance and then calls to react based on that to select the appropriate DVFS state dynamically during runtime.*

### Q5 Dynamic Power (3+5+2 points)

Consider the figure below as the value of MPKI for the L2 cache for the execution of a certain application over time. As a reference, the L2 cache is 128KB of size, 32-way set-associative.

- Start by briefly describing what you observe and give an idea of what could be the reason for the high and low values observed in the MPKI.
- Assume that for an L2 of half its size (64KB), the high MPKI values get higher, but the low MPKI values stay the same. Given the techniques described in this course, what could you suggest to improve the efficiency?
- How does your answer change if the horizontal scale was in millisecond instead of second?



- The high values in the curves could represent the application reading in new data that is not present in the caches while the low points could represent the application reading data that is present in the cache, i.e. data being reused (read during the high peaks and then used again and again).*
- If the smaller cache shows the same performance as the bigger cache for the low points than that could mean that the smaller cache is large enough to hold the data during that portion of the execution. So, then it would be ok to reduce the original 128KB cache to half by disabling half of the ways, thus reducing the power consumption but hopefully keeping the same performance and consequently increasing the efficiency.*
- If the scale now is in milliseconds it means that the periods of time that we could have the cache reduced to half are much shorter so it is more difficult to make all the operations to reduce the cache to half and still overcome the overheads of the change in size. So it may not "pay off" to change if the program behavior changes so fast!*

---

**Q6 Memory (3+5+2 points)**

---

Recently announced dedicated accelerators such as the Imsys processor presented in the invited guest lecture or the Google TPU, among others, use scratchpad memories in their designs.

- (a) Briefly explain what a scratchpad memory is.
- (b) In your own words justify why the designers of these processors use scratchpads instead of cache memories, what the benefits are and why it is possible to achieve them for these processors.
- (c) Explain in your own words which is/are the main reason/s why general-purpose CPU designers do not choose scratchpad memories for their designs.

- (a) A scratchpad memory is a memory module that is controlled by software, as opposed as a cache memory which is a memory module controlled by hardware.*
- (b) The applications being executed on these processors have very regular and predictable memory access patterns, so it is easy to program the loading and managing of the scratchpad.*
- (c) In general-purpose CPUs the applications are much more irregular and there is a huge variety of applications that can run at a time and also at the same time so the management of the scratchpad would be extremely difficult. To leave the burden to the programmers would be asking too much from the programmers, and to keep it as part of the runtime would be to make it too complex and again hard to achieve good performance.*

---

**Q7 Approximate computing (8+2 points)**

---

Approximate computing is a technique that was discussed in class as a way to improve the efficiency. This technique can be implemented at the hardware, memory, and software.

- (a) Select one of these three levels, explain the technique and describe an example for that level and how we could achieve the efficiency benefits from approximate computing with that example.
- (b) Relate your answer to one of the papers that have been assigned for reading in the course.

- (a) In class we discussed a technique that could avoid the energy overhead of refreshing the memory for part of the data in an application/system, i.e. for the non-critical data. After analyzing a program and identifying which data is critical and non-critical, the allocation of the data could be assigned to the reliable or the approximate memory area (i.e. the one having the regular refresh and the one not being refreshed). If the non-refresh area is large enough the gains could be considerable without major impacts on the application (if the application is tolerant to the errors introduced by approximation).*
- (b) The paper that was mentioning this technique was "Shimmer: Implementing a Heterogenous-Reliability DRAM Framework on a Commodity Server"*

---

**Q8 Reconfigurable computing (3+4+3 points)**

---

Field-Programmable Gate Array (FPGA) is a type of hardware that is composed of many small logic, memory, and connection modules.

- (a) In your own words describe what an FPGA is.
- (b) Why is it that FPGA can be more efficient than other processors? Give a real example of an application or part of an application that could benefit from FPGAs to improve its efficiency.
- (c) One use-case for FPGAs mentioned during the class discussion was the hardware in a 5G wireless router device. Was an FPGA used in that case because it is more efficient or for any other reason? Justify your answer.

- (a) FPGA is a hardware that can be programmed/modified to perform a specific pre-defined function.*
- (b) FPGA implementations can be more efficient because they are dedicated to a single function, they directly implement in hardware a specific solution, they have the capability to exploit a large degree of parallelism, as long as resources are available, and the data is usually shipped directly from the output of one operation to the input of the next operation in a dataflow way avoiding temporary storage of intermediate results. This is all achieved even though the frequencies for FPGAs are usually much lower than the ones found in common general-purpose CPUs.*
- (c) In this case that we discussed in class, the point was using FPGAs in a product to be launched before the 5G standard had been finalized and thus any necessary changes could be made after the deployment of the product instead of waiting the production until all had been finalized.*

---

**Q9 Data centers (4+3+3 points)**

---

In class we have discussed the project Natick in which the data center servers are closed inside a cylinder which is placed underwater.

- (a) Describe in your own words what you remember from that discussion, and which were some of the advantages and disadvantages of this approach as compared to more traditional data centers.
- (b) Discuss if you think this will be a reality any time soon?
- (c) Think about going even beyond this setup and have the data center be in space! How would this compare with the underwater data center proposed by Microsoft?

- (a) Some advantages include the controlled environment (no dust or variable humidity) and more controlled temperature (fewer temperature variations) Another advantage was the fact that the data centers could be deployed closer to the end users. As disadvantages one of them is the obvious fact that it is impossible to maintain or upgrade the equipment and another is that it was unclear what the final impact on the environment this was causing (specially it was deployed in large scales). Cost may also be a disadvantage.*
- (b) This could be a reality if there is the available investment for the infrastructure required and also at the end of the day if the total cost is lower than more traditional on-shore approaches.*
- (c) In space a key advantage would be lower operational temperature. But the disadvantages are many including the fact that there is no available power network so it would have to be powered by solar panels and that connectivity is also inexistant and even if available it would have very long latency! In addition, the deployment is a complex, costly and risky operation!*

---

**Q10 Article review (10 points)**

---

You have read several articles as part of the weekly reading. Select one of the papers you have read, briefly describe the material in the paper, highlight the contributions of the work and explain in your own words how it relates to the different topics discussed in this course?

*This answer depends on the paper selected.*