



## Exam

DAT277/DIT053: Energy-Aware Computing

DAT278/DIT054: Sustainable Computing

Thursday 20 August 2020, 08:30-12:30 (+10min tolerance)

---

### Examiner and contact person during exam:

Pedro Petersen Moura Trancoso (contact through canvas email)

### Grading intervals:

The maximum grade for this exam is 100 points

DAT277/278 (Chalmers): U < 40p ≤ Grade 3 < 60p ≤ Grade 4 < 80p ≤ Grade 5

DIT053/054 (GU): U < 40p ≤ Grade G < 75p ≤ Grade VG

### Examination solution and review:

Exam solution will be posted in canvas within 24h after the exam has finished. A canvas announcement regarding the exam review.

### Instructions:

- **ANSWERS:** All answers need to be typed in the space provided. Different questions need to start in a new page. Additional material may be added to your typed answers as pictures inserted in the same space with the typed text. (pictures can be used only for diagrams not for text)
- **LANGUAGE:** All answers should be written in English.
- **SUBMISSION:** Submit your answers in a SINGLE PDF file. Name your file "<course code>-<first name>-<last name>.pdf" (e.g. for John Smith taking DAT278 would be "DAT278-John-Smith.pdf"). You may submit many times during the exam time but only your LAST submission will be considered.
- **DEADLINE:** The deadline is 12:30, there is a 10 min tolerance but then the assignment closes at 12:40, no more submissions will be accepted after that. Submit before the deadline to avoid problems.
- **RESOURCES:** All available material is allowed. For any of your answers, if you use information from the internet, include the source (the web address).
- **PLAGIARISM:** Copying of material from any source (unless requested in the question) is FORBIDDEN (you need to write the answers in your own words). Communication and exchange of information with other people during the exam period is FORBIDDEN. Exam answers will be tested for plagiarism.
- **COMMUNICATION WITH EXAMINER:** Before communicating with the examiner test the Frequently Asked Questions (FAQ) page which will contain the answers given to other students. This page will be continuously updated. If you still need to communicate with the examiner, do this through canvas email.
- **CONTINGENCY PLAN:** In case there is a failure in canvas during submission you will have to submit your file with the answers as an attachment to an email to the examiner ppedro@chalmers.se. Use the same email in case the canvas site is down during the exam period and you need to communicate with the examiner.

### Question 1: Sustainability and Sustainable Computing (20 points)

- (a) (10 points) Find in the internet a definition for Sustainable Computing, copy-paste this definition to your answer along with the web address (URL) of where you got it from. Now explain it in your own words.

From [https://computing.fs.cornell.edu/Sustainable/fsit\\_definition.cfm](https://computing.fs.cornell.edu/Sustainable/fsit_definition.cfm)

*“Sustainable Computing is a principle that embraces a range of policies, procedures, programs, and attitudes that run the length and breadth of any use of information technologies. It is a holistic approach that stretches from power to waste to purchasing to education and is a life-cycle management approach to the deployment of IT across an organization. The concept of Sustainable Computing considers total cost of ownership, the total impact, and the total benefit of technology systems.”*

*So sustainable computing is a holistic concept that considers the whole life of the computer systems and its contribution to the environment. So, this considers manufacturing, shipping, running, and end-of-life phases. In each phase we should consider policies to reduce the carbon footprint for the use of the system.*

- (b) (10 points) Find in the internet one (1) example of Computing for Sustainability. Describe them in your own words and show the web address (URL) of where you got this example from. Note: Using an example that has been given in the solutions of previous exams solutions is not a valid answer.

*The article (“Computational Sustainability: Computing for a Better World and a Sustainable Future”) discussed in class includes several examples. We could consider using computers for controlling better the water distribution and consumption in a country as to save fresh water, improving the overall sustainability for the planet.*

## Question 2: Energy Efficiency (20 points)

- (a) (10 points) Check the Green 500 list of June 2020 in <https://www.top500.org/lists/green500/>. Explain in your own words what this ranking means, i.e. what makes a system be on the top of this list and what makes a system be on the bottom of this list? In the list it shows the ranking of each system on the Green500 and Top500 lists – Check the A64FX system which ranks high on the Green500 but low on the Top500 and explain why this happens?

*The Green 500 list ranks the most energy-efficient systems. A system at the top of the list is able to do the more work with the less amount of energy. The A64FX system ranks 4th in the June 2020 Green 500 list but it ranks 204th in the Top 500 list! What is happening is that the system is built for efficiency, not performance. So, it uses low power components that make its work consuming low energy such as the processor which is ARM based like the processors in a mobile phone. This is a very efficient processor, but it is not able to solve the problems at the same speed as other systems with more powerful processors.*

- (b) (10 points) The system on the 3<sup>rd</sup> place in the June 2020 list (NA-1) is equipped with the PEZY-SC2 processor. Find out information about this processor and explain in your own words which are the characteristics of this processor that help this system achieve such a high ranking?

*From <https://en.wikichip.org/wiki/pezy/pezy-scx/pezy-sc2>:  
“Introduced by PEZY along with their second-generation ZettaScaler-2.0 supercomputer series, the SC2 incorporates 2,048 cores along with 8-way SMT support for a total of 16,384 threads, twice as many cores as its predecessor. The PEZY-SC2 powers many of the top Green500 most efficient supercomputers with upward of 14 GFLOPS/watt in performance. Operating at 1 GHz, the PEZY-SC2 has a peak performance of 8.192 TFLOPS (single-precision) and 4.1 TFLOPS (double-precision) while consuming around 180 Watts. The PEZY-SC2 is designed using over 2.4 billion gates and is manufactured on TSMC's 16FF+ process. In attempt to increase adaptability in the field of deep learning and AI as well as to increase throughput for specialized workloads, the PEZY-SC2 introduced support for 16-bit half precision floating point support. At 1 GHz, the SC2 can peak at 16.4 TFLOPS for half precision.”  
This system has many small processors with support for many threads each. It is also a simple architecture and thus consumes low power compared with traditional processors. This looks like a GPU with many processing units for many parallel threads. Thus, it achieves very high efficiency and thus is highly ranked in Green 500 list.*

### Question 3: Metrics and Techniques (20 points)

- (a) (10 points) Assume that a processor P1 has a 4MB Last-Level Cache (LLC) (composed of two parts of 2MB) which dissipates in average 40W. When application A runs on this processor it takes 20s to complete and during the execution, we record 1000 LLC misses. Which is the “Energy-per-miss” for this setup? Then we scale the cache to half which means that it will dissipate approximately 20W. The execution now is 50s and the misses increase to 3000. What is the “Energy-per-miss” now? Which setup is more efficient? In a regular simulation, which tool or tools do you need to use to get a good estimate of the values for the “Energy-per-miss” metric?

Execution with 4MB LLC:

*Energy-for-execution = 40W x 20s = 800Ws = 800J*

*Energy-per-miss = 800 / 1000 = 0.8 J/miss*

Execution with 2MB LLC:

*Energy-for-execution = 20W x 50s = 1000Ws = 1000J*

*Energy-per-miss = 1000 / 3000 = 0.3 J/miss*

*According to the Energy-per-miss metric, the 2MB setup is more efficient.*

*In order to get these values from a simulation we need to use a tool to measure the power of the devices, we could use the CACTI tool. This, together with an execution driven simulator like SimpleScalar or GEM5 can be used to determine the energy consumption. The same execution driven simulator can be used to extract the number of misses giving all the required data to determine the energy-per-miss metric.*

- (b) (10 points) Explain what “power gating” is and how it may help improve the efficiency of a processor. Find a research paper that presents some work on a processor using power gating to improve its efficiency. Include the title and the authors, as well as the site where you found this work, in your answer. Briefly describe in your own words the contributions of that work (read the abstract to get the idea of the work).

*Power gating is a technique where we control the power delivery to the hardware components, thus when a component is not in use we are able to turn its power off and thus reduce the power dissipation and even not consume any static power since the component has no power at all.*

*In dl.acm.org if you search for power gating you can find many articles. We can select the following:*

*Po-Han Wang, Chia-Lin Yang, Yen-Ming Chen, and Yu-Jung Cheng. 2011. Power gating strategies on GPUs. ACM Trans. Archit. Code Optim. 8, 3, Article 13 (October 2011), 25 pages. DOI:<https://doi.org/10.1145/2019608.2019612>*

*From the abstract of this paper it is possible to understand that the authors propose different techniques to improve the power gating control for internal components of the GPUs. One of these techniques, the one which seems to have the most impact, is called simple time-out power gating which is able to reduce by 83.3% the leakage (static-power) in all units but the shaders. It is also important to note that the proposed techniques do not result in significant performance overhead (only 1%).*

### Question 4: Dynamic Power (20 points)

- (a) (10 points) AVX is a set of instruction extensions for the Intel processors. Find information on the internet about AVX (include the links for the sites where you found the info in your answer). Describe in your own words what AVX is and how it can help improve the performance and efficiency of the execution of an application. What do we need to do to use the AVX instructions in the execution? Which applications can benefit the most with the AVX instructions?

*From [https://en.wikipedia.org/wiki/Advanced\\_Vector\\_Extensions](https://en.wikipedia.org/wiki/Advanced_Vector_Extensions):*

*AVX instructions are SIMD (Single Instruction Multiple Data) instructions where in a single cycle, multiple operations of the same type can be done at small data types in parallel. The latest extension is the AVX512 where each instruction can for example operate simultaneously on eight 32bit single precision FP numbers. You can either add the instructions in the assembly code directly or use a compiler with the proper directives and command-line options to use these AVX extensions. Obviously, the processor needs to support these instruction extensions.*

*Usually scientific applications with many loops where the same operation is done on many data elements are good candidates for using these instruction extensions. For example, Machine Learning codes may benefit from these instructions.*

- (b) (10 points) Describe in your own words what “Cache compression” is. You may find information in the internet about this topic and include any sites in your answer. How can “cache compression” improve the energy efficiency of a system? Explain the advantages and disadvantages of “cache compression”.

*Cache compression is a technique where it is possible to compress the data in the cache, i.e. to somehow make the data be represented in a more compressed format and thus occupy less space than if uncompressed. While there are different compression methods, for cache compression in most cases we should use lossless compression methods so that no data information is lost with the compression.*

*The energy efficiency can be improved with cache compression as for example the same hardware cache is able to store more data and thus, less misses will happen and fewer costly requests to memory if this applied to the LLC. So, the advantages are that it is possible to store more data in the same space and thus improve the cache performance. As disadvantage it is the fact that now if the data is compressed, before it is utilized it needs to be decompressed and thus there is an overhead to the access latency to decompress and only after the data is usable for the computations.*

### Question 5: Heterogeneity and Data Centers (20 points)

- (a) (10 points) Find in the internet information on “Domain-Specific Architectures” (or DSAs) and include all links in your answer. Describe very briefly in your own words what Domain-Specific Architectures are. How can these architectures improve the energy efficiency in a system? Give a real-life example (a commercial product) of a DSA. Justify your answer with a link to the product.

*A DSA is an architecture that is designed to solve a set of problems from the same domain. So, this is a more specialized architecture than a common CPU or even a GPU. As such, since these architectures are design for specific tasks they are able to do those tasks more efficiently than a general-purpose CPU. The DSA term was presented in class with the article “A New Golden Age for Computer Architecture”. One commercial example of a DSA is the Google TPU. A commercial TPU edge product is sold by Coral and information for one of their products can be found in <https://coral.ai/products/dev-board>.*

- (b) (10 points) Find in the internet the definition for “Power Usage Effectiveness (PUE)”, include any links for the definition that you found. Explain it briefly in your own words. Find examples of PUE values reported by Data Center operators, include the links from the sites where you found the information. Explain at least one technique that Data center operators can use to improve the PUE.

*From [https://en.wikipedia.org/wiki/Power\\_usage\\_effectiveness](https://en.wikipedia.org/wiki/Power_usage_effectiveness):  
“Power usage effectiveness (PUE) is a ratio that describes how efficiently a computer data center uses energy; specifically, how much energy is used by the computing equipment (in contrast to cooling and other overhead)... An ideal PUE is 1.0. Anything that isn't considered a computing device in a data center (i.e. lighting, cooling, etc.) falls into the category of facility energy consumption.”  
So PUE accounts for the total energy consumed for a data center to do its job. The more “overhead” energy, the “worse” (larger) the PUE.  
One simple technique that has been employed by many Data Center operators is to equip the data center with solar panels thus reducing the need for energy from the outside grid, thus reducing the PUE. Another way is to make the systems more efficient so that energy is reduced in the cooling of the systems such as the Facebook Data Center in Luleå which uses open windows to cool thus reducing the need of energy for cooling.*