# Laboratory Assignment 1:
# Exploring the Impact of Cache Hierarchy on Processor Performance

### Contributors: Mehrzad Nejat, Mohammad Waqar Azhar, Per Stenstrom

**Learning outcome:** The objective of this laboratory assignment is to understand the impact of cache design on overall processor performance. You will use the SimpleScalar toolset [1] and the MiBench benchmark suite [2] to perform the experiments. At the end of this assignment, you should know how to:

- measure and report the performance of a cache model

- compare the relative performance of various cache models

- evaluate the impact of cache parameters on cache performance

**Suggested reading:** It is strongly recommended that you go through the following study material ([1] and [2] are available in the Canvas Repository) before you start this assignment:

[1] D. Burger and T. Austin, "*The SimpleScalar Tool Set Version 2.0*", University of Wisconsin-Madison, Computer Sciences Department, Technical Report 1342, 1997.

[2] M. R. Guthaus, J.S. Ringenberg, D. Ernst, T.M. Austin, T. Mudge and R.B. Brown, "MiBench: A free, commercially representative embedded benchmark suite", In *Proceedings of IEEE International Workshop on the Workload Characterization (WWC-4)*, pages 3 – 14, 2001.

**Evaluation:** You must write a report that includes: A brief description of the problem and the method you used, important assumptions you made (if any), simulation results and observations, your design choices and the reason behind them. Detailed report guidelines will be available on Canvas.
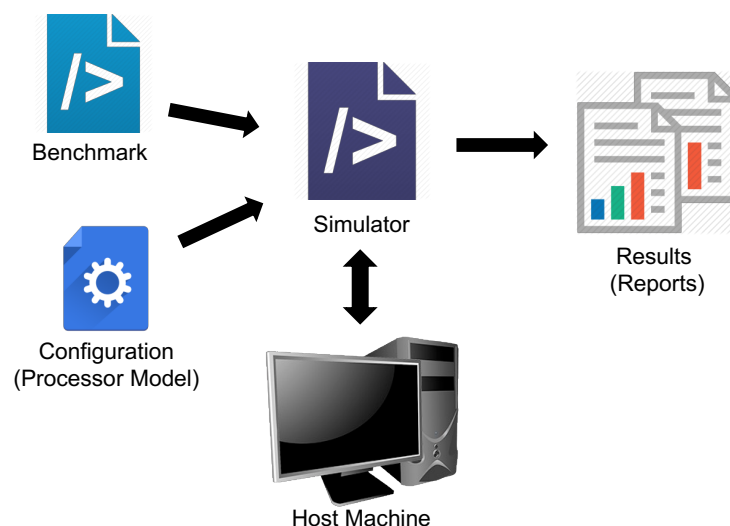
## Simulation Environment

As depicted in Figure 1, the simulation environment basically consists of the following sections:

- **Benchmark:** This is a test program that is executed on the simulated processor. It can have its own inputs and outputs. You can run any program on the simulated processor; but in order to have a standard base for comparison, people usually use programs and input sets from well-known benchmark suites.

- **Configuration File:** the architectural simulator uses this file to set the configuration and parameters of the processor and cache model in simulation. This is not the only method for modifying the processor and cache model. You can use parameter-specific command-line options or change the simulator code. But we will not use these methods in this lab. We provide a base

configuration file. You can create different processor and cache models by modifying this file.

- **Simulator:** This is the simulator program. It is executed on a real computer referred to as the *host machine*. The host machine can be your laptop, a local computer in the Lab, or a server. The simulator is a software model of a computer system and runs as software on the host machine. It includes probes for measurements while simulating the execution of benchmark on the computer system model being simulated with the specified configuration. The number of executed instructions, the number of loads and stores, and the average cycles per instruction are just a few examples of many parameters that the simulator can measure.

- **Results:** The simulator eventually reports the measurement results in some output files. After simulating different processor models, you can compare these results to evaluate the performance of the models when running a specific benchmark program with a specific input set.

In summary, you must remember that you are running a benchmark program on a simulated computer system model with the specified configuration on a real computer. Therefore, for example you have two different run times: program runtime on the simulated processor, and runtime of the simulation – the simulation time – on the host machine. Simulation time is usually several orders of magnitude larger than program run time.



***Figure 1: Simulation Environment***

## Running a simulation

You must run the simulator in a Linux system. You can use computers in the lab or connect to Chalmers server with a Secure Shell (SSH) connection[1].
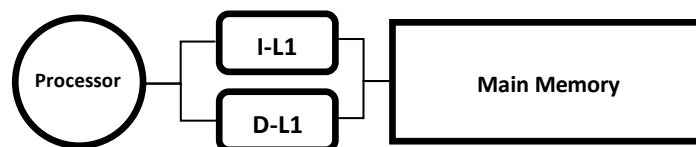
The required files can be found in the documents section of the course page on Ping-Pong. You must download '*simhome.tar.gz*'. After uncompressing, you can see three folders: '*apps*', '*bin*', and '*configs*'. The first one contains 5 benchmark applications from MiBench. We are using the small input sets to reduce the simulation time. The second folder contains the executable files of two SimpleScalar simulators. '*sim-profile*' is used for fast functional simulation of the benchmark applications to provide statistics on instruction mix. We will not use this simulator in this lab; but you can try it yourself. '*sim-outorder*' performs detailed architectural simulation of an out of order processor. This is the main simulator that you will use. '*configs*' folder is where you place different processor models and simulation results are stored in the corresponding sub-folders.

Besides these folders, you can find the base configuration file '*base1.txt*' and two scripts for running out of order simulator and profiling simulator: '*runsim_sim*' and '*runsim_profile*' respectively. Open these scripts and see how each simulation is initiated[2]. You can pass the name of the configuration as a command line argument to the script. You can deactivate any benchmark application by simply changing the corresponding lines to comment.

Remember, after each simulation, two different sets of outputs can be generated: Application outputs in the same location of the application executables, and simulation outputs in the same location as the configuration files. These locations are set in the script.

## Task 1: Evaluating the base configuration

In this assignment you will study a very simple system presented in Figure 2. Usually, there are two or three levels of cache between the processor and the main memory. One reason is that bigger caches come with longer access time and each level tries to hide the access time of the next level. However, we want to start with a comprehensive yet interesting system. Hence, here we only have one level of cache with one cycle access time. But, we keep the instruction and data caches separated as they are in most of processors for the first-level cache (denoted I-L1 and D-L1).



---

[1] You can connect to the server even from a windows operating system

[2] If you are not familiar with Linux scripts, you should check at least some basic structures and simple examples on internet. It will not take a lot of your time, but it can be a powerful tool.

We know that the execution time of a code is modeled as:

$$Execution\_Time = IC \times CPI \times T_{cycle}$$
( 1 )

Where IC is the instruction count, CPI is Cycles Per Instruction and $T_{cycle}$ is cycle time[1]. Because the processor frequency is constant in our case, we can remove cycle time from both sides of the equation and express all the times in processor cycles. Furthermore, since we run the same program on different computer system models with the same ISA, IC will be the same. Therefore, we can divide both sides by IC and consider only CPI – Cycles Per Instruction. This is especially useful when comparing the performance of multiple applications with different number of instructions. This means CPI could be used as an important parameter for measuring the performance of a processor.

CPI can be divided into two components[2]:
$$CPI = CPI_0 + MPI \times MP$$
( 2 )

Here MPI and MP stand for number of Misses Per Instruction and Miss Penalty in cycles, respectively. This is a simple model that applies to the system in Figure 2. The miss penalty component in ( 2 ) corresponds to the extra waiting time when a data or an instruction is not found in the data and instruction cache, respectively, and the processor needs to bring it from main memory. If we remove this component from CPI, what remains (we call it $CPI_0$) is the number of cycles per instruction related to processor operations and even if we had a perfect cache with no misses, it would not be affected.

Now, let's start by simulating the base model and measuring its performance. **Note:** It takes several minutes to run a benchmark. Do not sit and wait for an experiment to finish. Go ahead and plan for the next experiment.
- Place the provided "base1.txt" configuration file in the corresponding folder in '*simhome/configs*'. Open it and see how different processor configurations are set in this file. Especially, check the configurations for caches. Use the SimpleScalar user guide [1] to understand the parameters. You will change these parameters to create new processor models.
- Execute '*runsim_sim*' for the base model[3].
- Read the results from txt files starting with "*Stats_*" and fill out the first five rows of Table 1.

---

[1] The inverse of processor frequency

[2] Usually in a multi-level cache and when cache access time is longer than one cycle, the second component of this equation is expressed in more details that include the multiplication of accesses per instruction to each cache or memory by the corresponding access time.

[3] Hint: All the files that need to be executed, should have executable permissions.

***We need to establish CPI₀ for each benchmark. How would you configure the simulator to get CPI₀ for each benchmark?***

- Simulate the ideal model and collect the CPI results as an estimation of $CPI_0$ and fill out the corresponding row in Table 1.
- Using $CPI_{base}$ and $CPI_0$, calculate the upper bound on speedup ($SP_{ideal}$). Use the following formula for calculating the speed up of any configuration x:

$$SP_x = \frac{CPI_{base}}{CPI_x} \qquad\qquad (3)$$

- Choose the two applications that suffer the most from accessing memory. You will focus on these applications in the next task[1].

*Table 1: Simulation results for base configuration*

| Application | dijkstra | qsort | stringsearch | gsm-untoast | jpeg-cjpeg |
|---|---|---|---|---|---|
| Instruction Count* | | | | | |
| Execution Time in cycles | | | | | |
| $CPI_{base}$ | | | | | |
| MPI  I-L1 | | | | | |
| MPI  D-L1 | | | | | |
| $CPI_0$ | | | | | |
| $SP_{ideal}$ | | | | | |

\* use number of committed instructions

Elaborate on these questions:
- ➢ How much is the memory subsystem responsible for each application's execution time?
- ➢ If you could only optimize one cache to improve the performance, which cache would you choose?
- ➢ Is the effect of the memory subsystem on execution time similar for different applications?

## Task 2: Optimization of the Caches

You have studied a base processor model in Task 1. In this task, the goal is to find a cache model that improves performance for the benchmarks you have chosen.

To do this, you need to devise a methodology to understand the impact each cache organization parameter has on the performance of the processor. These parameters are cache size, associativity and block size. It is important that you change only one parameter at a time to isolate the effect of other parameters. For example, if you

---

[1] Hint: You can deactivate the simulations of each benchmark by changing the corresponding lines in the simulation script into comment. This way, you can easily reactivate it later.

change associativity and cache size[1] together and observe some performance variations, you cannot determine how much each parameter in isolation impacts on the result.

Every time you want to move to the next step in the simulations, create a new folder and configuration file in "*configs*" folder. You can start by copying '*base1.txt*' and then modifying it. In each step, alter only one parameter of one cache and derive MPI and CPI from the simulation result. Then calculate the speedup relative to base configuration and compare it with the upper bound you found in the previous task for an ideal memory subsystem. You can also elaborate the results by generating curves and graphs. Perform the simulations for the two applications you chose in the previous task. Whenever you make a decision, you should consider both applications.

### 2.1. Optimize D-Cache

***2.1.1. Cache Size:*** The objective here is to find the optimal size of the D-Cache in the range 4 – 32 KB. You do that by freezing other design parameters. Use a configuration with a 2-way associative cache with a block size of 16 Bytes.

**Step 1:** Devise a methodology using the simulator for how to determine the optimal size of the D-Cache.

**Step 2:** Determine the optimal size and the performance of the two selected applications using that size.

***2.1.2. Associativity:*** The objective here is to find the optimal associativity of the D-Cache in the range 1 – 8 ways. You do that by freezing other design parameters. Use a configuration with the optimal cache size found in 2.1.1 with a block size of 16 Bytes.

**Step 1:** Devise a methodology using the simulator for how to determine the optimal associativity of the D-Cache.

**Step 2:** Determine the optimal associativity and the performance of the two selected applications using that associativity.

***2.1.3. Block size:*** The objective here is to find the block size of the D-Cache in the range 16 – 128 B. You do that by freezing other design parameters. Use a configuration with the optimal cache size found in 2.1.1 and the optimal associativity found in 2.1.2.

**Step 1:** Devise a methodology using the simulator for how to determine the optimal block size of the D-Cache.

**Step 2:** Determine the optimal block size and the performance of the two selected applications using that block size.

---

[1] Hint: In configuration file, cache size is not directly available. It is the multiplication of number of sets, associativity and block size. Therefore, you must be careful when you want to change each parameter. Check simulator user guide.

## 2.2. Optimize the I-Cache

You have found an optimum configuration for the D-Cache so far. Start from the base configuration and perform the same steps as in Task 2.1 for the I-Cache.

## 2.3. Compare the results

By now, you should have found an optimum configuration for each cache. Put it together and make one final configuration named OPT1. But, choose one block size for both caches. Even though it is not impossible to have different block sizes, it complicates the system since these blocks of data move through different components of the memory subsystem. Perform the simulations for the rest of the applications with OPT1, calculate speedup compared to base and fill out Table 2. Compare the geometric mean of the speed ups.

*Table 2: Comparing the speedup of OPT1 model with an Ideal memory system*

| Application | dijkstra | qsort | stringsearch | gsm-untoast | jpeg-cjpeg | GM |
|-------------|----------|-------|--------------|-------------|------------|----|
| $SP_{OPT1}$ | | | | | | |
| $SP_{ideal}$ | | | | | | |

Elaborate on the following questions:
- ➢ Which cache parameters have the largest effect on performance? Or are the effects in a similar range?
- ➢ How close did you get to an Ideal memory system? What are preventing you from reaching the maximum speed up?

# Task 3: A more realistic system model

So far, we have assumed that changing the cache size and associativity doesn't have any cost. But, this is not true in a real system. Other than area cost, larger caches with higher associativity comes with longer access time and higher power consumption. This is one of the main reasons of having multiple levels of a cache in real computers. We do not want to complicate the problem by looking into area or power costs. But, we can have a simple case with a more realistic assumption on performance cost of increasing size and associativity. Use the values[1] in Table 3 and update the cache access time in cycles in OPT1 and name the updated model OPT2.

---

[1] Note that these values are hypothetical and not based on real measurements or data. They are only useful for this experiment.

Find now the best cache configuration using the methodology you used in Task 2. Calculate the speedups and fill out Table 4. Note that the base model remains unchanged by applying the values in Table 3.

Think about these questions:
- ➢ What is the effect of considering latency costs of your proposed cache architecture?
- ➢ Is it still worth making those changes? If yes, how much further increase in the latency costs do you think would destroy the improvement you have achieved?

*Table 3: Cache access latency in cycles for different sizes and associativities.*

|  |  | Size | | | |
|---|---|---|---|---|---|
|  |  | 4 KB | 8 KB | 16 KB | 32 KB |
| Associativity | 1 | 1 | 1 | 2 | 2 |
|  | 2 | 1 | 2 | 2 | 3 |
|  | 4 | 2 | 2 | 3 | 3 |
|  | 8 | 2 | 3 | 3 | 4 |

*Table 4: Comparing the speedups of OPT1, OPT2 and Ideal memory system*

| Application | dijkstra | qsort | stringsearch | gsm-untoast | jpeg-cjpeg | GM |
|---|---|---|---|---|---|---|
| $SP_{OPT1}$ |  |  |  |  |  |  |
| $SP_{OPT2}$ |  |  |  |  |  |  |
| $SP_{ideal}$ |  |  |  |  |  |  |