

Analyzing Cache Behaviour and System Performance

Objectives.

- Study and run a program that experimentally calculates the L1 data cache size
- Case Study #1. Analyze cache behaviour and performance using Timers and Valgrind tool
- Case Study #2. Analyze cache behaviour and performance using Timers and Valgrind tool
- To introduce the roofline model

Aim

The **aim** of this lab session is to understand the behaviour of cache memories and how memory hierarchy affects system performance. Developing software that utilizes memory hierarchy is of critical importance in High Performance Computing, whether we are using CPUs or GPUs.

Section 1 – Experimental procedure that calculates the L1 data cache size

In `cache_benchmark.c`, I have developed a program that stores an array into memory many times. In `cache_benchmark()` routine, the $X[N]$ array is stored 'TIMES' times. I have included a 'weird' code in the beginning, which maps/allocates the running thread to core number zero (first CPU core only); modern multi-core processors have more than one CPU cores. You do not have to understand how this code works. So, the running process which consists of a single thread will run on one CPU core only. If we omit this piece of code, then, the Operating System (OS) will toggle our process amongst different CPU cores. The OS does this to reduce the overall heat dissipation of the CPU chip.

Task1: Compile and run the program `cache_benchmark.c` for $N=1000,2000,4000,8000,16000$ and 32000 . Measure the execution time for each case. Draw a graph, where execution time is the y-axis and N is the x-axis. To compile this program you will need extra gcc options, which they are given in the first line of `cache_benchmark.c`. **Make sure the execution time lasts a few seconds, otherwise the execution time measured is not accurate.** This is because other processes run too; if the execution time of this process is much larger than the others', then the value measured is more accurate. To this end, you must appropriately specify the 'TIMES' value.

You will realize that the execution time is not linear to the input size. Why?

Answer: As long as $X[]$ is smaller than L1 data cache size, then it is stored 'TIMES' times to L1 data cache memory, which is very fast. If $X[]$ no longer fits in L1 data cache, then it is written 'TIMES' times to L2 cache, which is slower.

Section 2 – Analysing Cache Behaviour and Performance using Valgrind

Case study #1. Array Initialization

Cache Analysis using Cachegrind tool

Task1: Download *'init_array.c'* file. This program initializes a 2-d array in a column wise manner (N=1000). We have discussed in the lecture that accessing an array in column-wise order is not efficient. The following example will verify the theory. Compile it and use Valgrind Cachegrind to extract the number of L1 data cache accesses and misses. Then, apply loop interchange and apply the above process again. Compare the results. I have included specific instructions about how to do that hereafter.

Compile the program using:

```
gcc init_array.c -o p -O2 -g
```

Use valgrind to simulate the cache by typing:

```
valgrind --tool=cachegrind ./p
```

Use the annotate command to get a more detailed picture about L1 data cache accesses and misses by typing:

```
cg_annotate cachegrind.out.24486
```

```
or cg_annotate cachegrind.out.24486 > column.txt
```

Then, apply loop interchange and repeat the above steps.

```

31
32 -- Auto-annotated source: /home/user01/Desktop/comp3001/my/Labs/cache/init_arrays.c
33
34 Ir          IImr      IImr      Dr          DImr      DLmr      Dw          DImw      DLmw
35
36 -- line 4 -----
37 .           .         .         .           .         .         .         .         .
38 .           .         .         .           .         .         .         .         .
39 declared    .         .         .           .         .         .         .         .
40 .           .         .         .           .         .         .         .         .
41 .           .         .         .           .         .         .         .         .
42 .           .         .         .           .         .         .         .         .
43 .           .         .         .           .         .         .         .         .
44 .           .         .         .           .         .         .         .         .
45 2,001 ( 0.04%) 0       0         0         0         0         0         0         0
46 .           .         .         .           .         .         .         .         .
47 .           .         .         .           .         .         .         .         .
48 .           .         .         .           .         .         .         .         .
49 .           .         .         .           .         .         .         .         .
50 zero, we mean that the program ended successfully.
51 2 ( 0.00%) 0       0         1 ( 0.00%) 1 ( 0.07%) 0       0         0         0
52 .           .         .         .           .         .         .         .         .
53 .           .         .         .           .         .         .         .         .
54 .           .         .         .           .         .         .         .         .
55 .           .         .         .           .         .         .         .         .
56 .           .         .         .           .         .         .         .         .
57 6,002 ( 0.12%) 1 ( 0.12%) 1 ( 0.13%) 0       0         0         0         0
58 2,000,000 (39.00%) 0       0         0         0         0         0         0
59 3,000,000 (58.51%) 0       0         0         0         1,000,000 (99.05%) 1,000,000 (99.95%) 62,500 (99.18%)
60 .           .         .         .           .         .         .         .         .
61 .           .         .         .           .         .         .         .         .
62 .           .         .         .           .         .         .         .         .
63 .           .         .         .           .         .         .         .         .
64 .           .         .         .           .         .         .         .         .
65 .           .         .         .           .         .         .         .         .
66 .           .         .         .           .         .         .         .         .
67 .           .         .         .           .         .         .         .         .
68 .           .         .         .           .         .         .         .         .
69
70 -- line 36 -----
71
72 -----
73 The following files chosen for auto-annotation could not be found:
74 -----

```

Fig.1 Valgrind output for init_arrays.c (column-wise)

```

32 -- Auto-annotated source: /home/user01/Desktop/comp3001/my/Labs/cache/init_arrays.c
33
34 Ir          IImr      IImr      Dr          DImr      DLmr      Dw          DImw      DLmw
35
36 -- line 4 -----
37 .           .         .         .           .         .         .         .         .
38 .           .         .         .           .         .         .         .         .
39 declared    .         .         .           .         .         .         .         .
40 .           .         .         .           .         .         .         .         .
41 .           .         .         .           .         .         .         .         .
42 .           .         .         .           .         .         .         .         .
43 .           .         .         .           .         .         .         .         .
44 .           .         .         .           .         .         .         .         .
45 2,001 ( 0.04%) 0       0         0         0         0         0         0         0
46 .           .         .         .           .         .         .         .         .
47 .           .         .         .           .         .         .         .         .
48 .           .         .         .           .         .         .         .         .
49 .           .         .         .           .         .         .         .         .
50 zero, we mean that the program ended successfully.
51 2 ( 0.00%) 0       0         1 ( 0.00%) 1 ( 0.07%) 0       0         0         0
52 .           .         .         .           .         .         .         .         .
53 .           .         .         .           .         .         .         .         .
54 .           .         .         .           .         .         .         .         .
55 .           .         .         .           .         .         .         .         .
56 .           .         .         .           .         .         .         .         .
57 .           .         .         .           .         .         .         .         .
58 4,002 ( 0.08%) 0       0         0         0         0         0         0         0
59 2,000,000 (39.02%) 0       0         0         0         0         0         0         0
60 3,000,000 (58.53%) 0       0         0         0         1,000,000 (99.05%) 62,500 (99.13%) 62,500 (99.18%)
61 .           .         .         .           .         .         .         .         .
62 .           .         .         .           .         .         .         .         .
63 .           .         .         .           .         .         .         .         .
64 .           .         .         .           .         .         .         .         .
65 .           .         .         .           .         .         .         .         .
66 .           .         .         .           .         .         .         .         .
67 .           .         .         .           .         .         .         .         .
68 .           .         .         .           .         .         .         .         .
69
70 -- line 36 -----
71
72 -----
73 The following files chosen for auto-annotation could not be found:
74 -----
75 /build/glibc-0TsELS/glibc-2.27/elf/./sysdeps/x86_64/dl-machine.h
76 /build/glibc-0TsELS/glibc-2.27/elf/dl-lookup.c

```

Fig.2 Valgrind output for init_arrays.c (row-wise)

The Cachegrind output is shown in Fig.1 and Fig2. The memory statistics look very similar, however, if you look closer, there is a big difference in the 'D1mw' column (L1 data cache write misses). Although, both programs write 1,000,000 elements into L1 data cache (Dw column), in the first case they are written in column-wise order, while in the second case they are written row-wise.

In the column-wise case, to write just $A[0][0]$ into memory, an entire L1 data cache line is loaded to dL1 memory (it will contain 16 elements of A, i.e., $A[0][0:15]$) and then just $A[0][0]$ is stored into it. Then, to write $A[1][0]$, another cache line is loaded (16 elements) and just $A[1][0]$ is stored etc. Thus, there is one dL1 miss for every store, as $A[0][0]$ is stored in DDR, not dL1. On the contrary, in the row-wise case, a dL1 miss is followed by 15 subsequent dL1 hits. This is evidenced by the 'D1mw' column (L1 data miss writes). In the column-wise case, there are 1,000,000 dL1 write misses, while in the row-wise case there are just 62,500 dL1 write misses. There are 16 times less misses as each dL1 cache line contains 16 elements. Keep in mind that $62,500 \times 16 = 1,000,000$.

Comparison using timers

Cache misses introduce significant performance penalties, and this is reflected on the program's execution time.

Task3: Run the two different versions of the *'init_arrays.c'* program discussed above and compare their execution times. You will find out that initializing an array of integers in a row-wise manner is about 16 times faster than the column-wise. Recall that to get an accurate execution time value, **the execution time must be at least a few seconds, otherwise is not accurate.**

Generating the assembly

You can generate the assembly of the above two program versions and compare them. You can generate the assembly code by using the following command

```
gcc -S init_arrays.c -o assembly.s -O2
```

Case study #2. Matrix-Matrix Multiplication

Task1. Download 2mmm.c file. The routine we are interested in is *void mmm()*. Study the program and make sure you understand what it does. Use Cachegrind and pg_annotate to simulate the behaviour of dL1. Measure the dL1 accesses and misses. You should see something like Fig.3.


```

for (j=0;j<N;j++){
  c=C[i][j];
  for (k=0;k<N;k++) {
    c+=A[i][k]*B[k][j];
  }
  C[i][j]=c;
}

for (i=0;i<N;i++)
  for (j=0;j<N;j++){
    e=E[i][j];
    for (k=0;k<N;k++){
      e+=C[i][k]*B[k][j];
    }
    E[i][j]=e;
  }
}

```

Thus, the binary is closer to the code above rather than the code we developed. As you can observe from Fig.3, in `mmm()` routine there are 500.000 writes in total (see Dw column). This value was expected as two 2d arrays of size 500x500 are stored into memory.

The Roofline Model

The roofline model [3] provides an easy way to get performance bounds for compute-bound and memory-bound loop kernels. It allows us to know how far the achieved performance is from the optimum. It is based on the concept of computational intensity, sometimes also called arithmetic or operational intensity. The arithmetic intensity is given by the following formula: '*FP.arithmetical.instructions / number.of.bytes.loaded.stored*'. This model has several limitations [3], e.g., does not consider all features of modern processors and ignores integer computations.

Roofline Model for MMM algorithm: MMM has N^3 iterations and each iteration contains 4 Floating Point (FP) L/S operations and 2 FP arithmetical operations. So, the arithmetical intensity of MMM (the ratio between FP arithmetical operations and number of bytes loaded/stored), when the arrays are of type 'float', is $2/(4*4\text{bytes})=1/8$.

MMM loop kernel contains integer arithmetical instructions too (we could generate the assembly and check them out by typing `gcc source.c -S assembly.s -O2`). The integer arithmetical instructions are responsible for a) computing the L/S memory addresses (e.g., `...=A[i][j]` will be broken down to multiple assembly instructions), b) controlling the iterators (increment i, compare i to N, branch back); however, these integer operations take 1 CPU cycle each and they are performed in parallel with the FP ones. So, for this loop kernel we could assume that performance is not affected by the integer operations. Furthermore, the roofline model does not consider integer operations, and this is a serious limitation of this model.

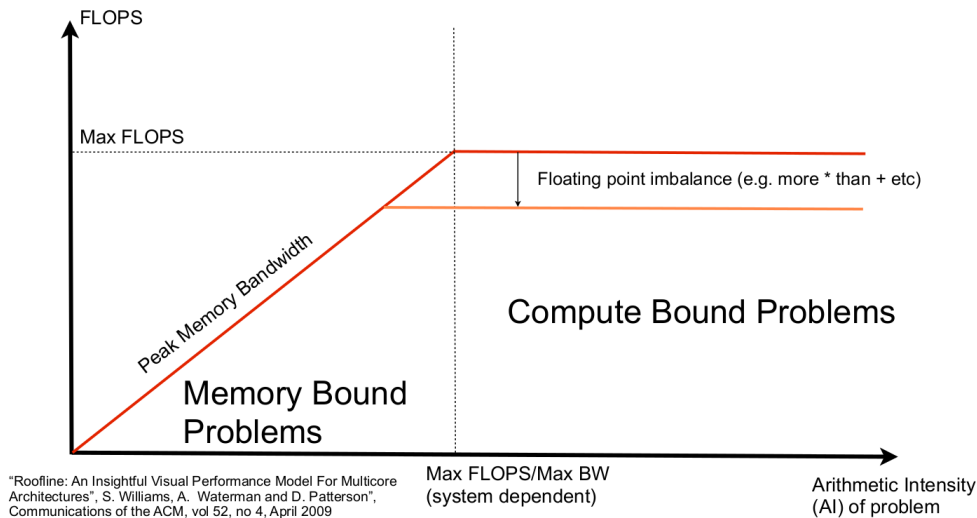


Fig.4. The Roofline Model [3]

Algorithms that have a low arithmetical intensity are **memory-bound**, while algorithms that have a high arithmetical intensity are **compute-bound**. Memory-bound means that their performance is bounded on the memory latency and bandwidth values; in simple words, no matter how high the CPU frequency is, or no matter how many cores the CPU supports, performance depends on how fast the data are loaded/stored from/to memory.

$$\text{Attainable FLOPS} = \min(\text{Peak Floating Point Performance}, \text{Peak Memory Bandwidth} \times \text{Arithmetical Intensity})$$

The peak FP performance is the maximum we can get, and it refers to compute-bound loop kernels with a perfect balance between simple and complex FP operations. The latter is an advanced topic and it is not further explained here; for those who want to go deeper and learn more about it, ask the module leader. The peak FP performance is CPU dependent.

Tip. each instruction has a latency and a throughput value, where the latter one is always larger or equal to the first; thus, to achieve the optimum performance of a bunch of instructions, e.g., multiply instructions, multiple multiply instructions must be feed the instruction pipeline one after another.

The peak memory bandwidth depends on the DDR memory and memory controller hardware characteristics. Furthermore, if the data fit and remain in the cache, the peak memory bandwidth is the cache bandwidth. As a reminder, a DDR memory access takes about 100-200 CPU cycles, an L3 memory accesses about 40-70 CPU cycles, an L2 memory access 6-14 CPU cycles, while an L1 memory access takes about 1-4 CPU cycles.

So, if the peak memory bandwidth is 21GBytes/sec, then the maximum MMM performance will be 21GB/sec x (1/8)=2.65gigaflops. If the arrays fit in the precious cache memories, then the memory bandwidth is higher and thus performance is increased.

How can we improve the performance of memory-bound algorithms? The main strategies are as follows. Reducing the number of memory accesses through the whole memory hierarchy; this relates to

reducing the number of cache-misses too. Another strategy is to use software prefetching. The above can be achieved by using code optimizations such as loop tiling, register blocking, array copying, loop merge/distribution etc. We will study those next week.

Further Reading:

1. The Valgrind Quick Start Guide, available at <https://www.valgrind.org/docs/manual/quick-start.html#quick-start.intro>
2. Cachegrind: a cache and branch-prediction, available at <https://valgrind.org/docs/manual/cg-manual.html>
3. Samuel Williams, Andrew Waterman, and David Patterson. 2009. Roofline: an insightful visual performance model for multicore architectures. Commun. ACM 52, 4 (April 2009), 65-76. DOI=10.1145/1498765.1498785, available at <https://people.eecs.berkeley.edu/~kubitron/cs252/handouts/papers/RooflineVyNoYellow.pdf>