
Solution for Project 1

HPC Lab — Submission Instructions

(Please, notice that following instructions are mandatory:
submissions that don't comply with, won't be considered)

- Assignments must be submitted to iCorsi (i.e. in electronic format).
- Provide both executable package and sources (e.g. C/C++ files, Matlab). If you are using libraries, please add them in the file. Sources must be organized in directories called:
Project_number_lastname_firstname
and the file must be called:
project_number_lastname_firstname.zip
project_number_lastname_firstname.pdf
- The TAs will grade your project by reviewing your project write-up, and looking at the implementation you attempted, and benchmarking your code's performance.
- You are allowed to discuss all questions with anyone you like; however: (i) your submission must list anyone you discussed problems with and (ii) you must write up your submission independently.

In this project you will practice memory access optimization, performance-oriented programming, and OpenMP parallelization on the Rosa Cluster .

1. Rosa Warm-Up

(5 Points)

1. Modules concept is used to allow the user to manage between different versions of software packages without needing to configure the corresponding environment variables that could affect every user of the cluster. For example, by calling *module load modulename/module-version*, the user can load the desired version of the software if available. To check availability is as simple as calling *module avail*. In addition to this commands, it is possible to check the loaded modules *module list*, unload all loaded modules *module purge*, among others. ¹
2. Slurm is a highly scalable cluster management and job scheduling system for Linux clusters (workload manager). As a manager it assigns access between nodes and users for certain time interval given a task (allocation). Furthermore, it provides a framework for starting, executing, and monitoring work (normally a parallel job) on a set of allocated nodes and it does prioritization of the tasks according the queued jobs. ²
3. The code is provided under the folder "exercise1" with the name "printhost" with the cited sources used during the development inside the file.

¹<https://www.ci.inf.usi.ch/research/resources/>

²<https://slurm.schedmd.com/overview.html>

```
[bellaj@icslogin01 ~]$ srun --nodes=1 --time=00:05:00 --pty bash -i
[bellaj@icsnode17 ~]$ ls
avx_example.cpp  ciao  loop_reader_example.cpp  main.exe  read_environment.cpp  tp1
[bellaj@icsnode17 ~]$ cd tp1
[bellaj@icsnode17 tp1]$ ls
exe1
[bellaj@icsnode17 tp1]$ cd exe1
[bellaj@icsnode17 exe1]$ ls
printhost.cpp
[bellaj@icsnode17 exe1]$ gcc printhost.cpp -o printhost
[bellaj@icsnode17 exe1]$ ls
printhost  printhost.cpp
[bellaj@icsnode17 exe1]$ ./printhost
Hostname: icsnode17
[bellaj@icsnode17 exe1]$
```

Figure 1: Terminal of Cluster print while testing the code

- Here i implemented the batch command in the file *printhostbatch.sh* and the output is saved in the file *printhostbatch.txt*. Notice that in the instructions is recommended to use the `-constraint` option to specify an available feature, however the available features are not set "(null)". Therefore, i use the `-partition` option to specify the partition to run the job as an example.

```
[bellaj@icslogin01 exe1]$ sbatch printhostbatch.sh
Submitted batch job 9272
[bellaj@icslogin01 exe1]$ ls
helloworld1      output_9260.txt  printhostbatch.sh  slurm-9258.out  slurm-9262.out
helloworld1.txt  printhost        printhost.cpp      slurm-9259.out  slurm-9270.out
[bellaj@icslogin01 exe1]$ cat helloworld1
Loading gcc/13.2.0-gcc-8.5.0-5hqhkwo
Loading requirement: gcc-runtime/8.5.0-gcc-8.5.0-7fyorqa
gmp/6.2.1-gcc-8.5.0-lrpcvy5 mpfr/4.2.1-gcc-8.5.0-ybeybcx
mpc/1.3.1-gcc-8.5.0-cv2gjfw zlib-ng/2.1.6-gcc-8.5.0-ztbc5xt
zstd/1.5.5-gcc-8.5.0-azepnn7
Hostname: icsnode07
```

Figure 2: Terminal of Cluster print while testing the code

- The file *printhostbatch2node.sh* contain the implementation of the batch script that runs the program on two nodes as requested:

```
[bellaj@icslogin01 exe1]$ sbatch printhostbatch2node.sh
Submitted batch job 9288
[bellaj@icslogin01 exe1]$ cat helloworld-2nodes.txt
Loading gcc/13.2.0-gcc-8.5.0-5hqhkwo
Loading requirement: gcc-runtime/8.5.0-gcc-8.5.0-7fyorqa
gmp/6.2.1-gcc-8.5.0-lrpcvy5 mpfr/4.2.1-gcc-8.5.0-ybeybcx
mpc/1.3.1-gcc-8.5.0-cv2gjfw zlib-ng/2.1.6-gcc-8.5.0-ztbc5xt
zstd/1.5.5-gcc-8.5.0-azepnn7
Hostname: icsnode19
Hostname: icsnode18
```

Figure 3: Terminal of Cluster print while testing the code

As in previous exercises, the sources utilized are cited in the file. ³

2. Performance Characteristics

(30 Points)

2.1. Peak performance

Following the example, I visit the USI Rosa cluster specifications website ⁴. It is possible to find that each node has 2 Intel Xeon E5-2650 v3 (25M Cache, 2.3 GHz).

³At the moment I started the assignment the batch files were not included, reason why the naming is different to the provided later on and i added links that i used to understand how to construct such a file to the .sh document

⁴<https://www.ci.inf.usi.ch/research/resources/>

Knowing this, it was possible to find the processor specification in intel's website ⁵, where it is shown that the processor has 10 cores, a base frequency of 2.3 GHz ($f = 2.30$) and supports AVX2 SIMD instructions, as specified on "Instruction Set Extensions". Since AVX2 uses 256-bit wide vector registers, which can process 4 double-precision (64-bit) numbers at once, then $n_{SIMD} = 4$.

With respect to FMA instructions, on Chapter 15 - section 15 of Intel's Optimization Reference Manual Volume 1, it is mentioned that Haswell microarchitectures implements FMA instructions with execution units on port 0 and 1 and 256 - bit data paths. "Dot product, matrix multiplication and polynomial evaluations are expected to benefit from the use of FMA, 256 - bit data path and the independent executions on two ports. The peak throughput of FMA from each processor core are 16 single-precision and 8 double-precision results each cycle". Therefore, $n_{FMA} = 2$. In addition, it also tells us that since "...the independent executions on two ports", we can expect that the number of FP operations the FPU of a core can execute in parallel within a single clock cycle is 2 (n_{super}).⁶

With this information we can compute the peak performance of a single core:

$$\begin{aligned} P_{core} &= f \times n_{SIMD} \times n_{FMA} \times n_{super} \\ &= 2.3 \times 4 \times 2 \times 2 \\ &= 36.8 \text{ GFLOP/s} \end{aligned}$$

As i mentioned before, each processor has 10 cores, therefore each CPU has a peak performance of 368 GFLOP/s.

$$\begin{aligned} P_{CPU} &= n_{cores} \times P_{core} \\ &= 10 \times 36.8 \text{ GFLOP/s} \\ &= 368 \text{ GFLOP/s} \end{aligned}$$

Since in the cluster website it is specified that there are two sockets per node, then the peak performance of a node is 736 GFLOP/s:

$$\begin{aligned} P_{node} &= n_{sockets} \times P_{CPU} \\ &= 2 \times 368.0 \text{ GFLOP/s} \\ &= 736 \text{ GFLOP/s} \end{aligned}$$

Finally, it is also specified that there are 42 compute nodes in total in the cluster:

$$\begin{aligned} P_{cluster} &= n_{nodes} \times P_{node} \\ &= 42 \times 736.0 \text{ GFLOP/s} \\ &= 30.912 \text{ TFLOP/s} \end{aligned}$$

Which means that the peak performance of the cluster is 30.912 TFLOP/s.

2.2. Memory Hierarchies

I present each terminal output after following the instructions:

⁵<https://ark.intel.com/content/www/us/en/ark/products/81705/intel-xeon-processor-e5-2650-v3-25m-cache-2-30-ghz.html>

⁶It is also possible to check on <https://uops.info/table.html>, however the results are quite extensive even after filtering by Haswell architecture and instruction set categories.

```

[bellaj@icslogin01 ~]$ lscpu
Architecture: x86_64
CPU op-mode(s): 32-bit, 64-bit
Byte Order: Little Endian
CPU(s): 20
On-Line CPU(s) list: 0-19
Thread(s) per core: 1
Core(s) per socket: 10
Socket(s): 2
NUMA node(s): 2
Vendor ID: GenuineIntel
CPU family: 6
Model: 63
Model name: Intel(R) Xeon(R) CPU E5-2650 v3 @ 2.30GHz
Stepping: 2
CPU MHz: 3800.000
CPU max MHz: 3800.0000
CPU min MHz: 1200.0000
bogomips: 4559.95
L1d cache: 32K
L1i cache: 32K
L2 cache: 256K
NUMA node0 CPU(s): 0-9
NUMA node1 CPU(s): 10-19
Flags: fpu_sme de pse tsc msr pae mce cx8 apic sep mtrr pge mca cmov pat pse36 clflush dts acpi mmx fxsr sse sse2 ss ht tm pbe syscall nx pdpe1gb rdtscp lm constant_tsc arch_perfmon pebs bts rep_good nopl xtopology nonstop_tsc cpuid aperfperf pni pclmulqdq dtcs64 mon
itor ds_cpl smx est tm2 ssse3 sdbg fma cx16 xtpr pdcm pcid dca sse4_1 sse4_2 xzavic novbe popcnt tsc_deadline_timer aes xsave avx f16c rdrand lahf_lm abm cquid_fault epb invpcid_single intel_pspin ssbd ibrs ibpb stibp fsgsbase tsc_adjust bmi1 avx2 smep bmi2 erms invpcid cqm xsaveopt
cqm_llc cqm_occup_llc dtherm ida arat pin pmu no_clear_flush tid

```

Figure 4: CPU model and specifications

Here we can observe the CPU model and specifications of the cluster. Which confirms the information we could obtain from the cluster website and the Intel's website.

We proceed to obtain the main memory information as indicated in the commands:

```

[bellaj@icslogin01 ~]$ cat /proc/meminfo
MemTotal: 57263108 kB
MemFree: 19679164 kB
MemAvailable: 53880264 kB
Buffers: 485856 kB
Cached: 33913188 kB
SwapCached: 164 kB
Active: 6693988 kB
Inactive: 29388532 kB
Active(anon): 5032 kB
Inactive(anon): 1888704 kB
Active(file): 6688956 kB
Inactive(file): 27499828 kB
Unevictable: 12388 kB
Mlocked: 12388 kB
SwapTotal: 33554428 kB
SwapFree: 33543616 kB
Dirty: 68 kB
Writeback: 0 kB
AnonPages: 1639692 kB
Mapped: 633736 kB
Shmem: 199840 kB
KReclaimable: 750428 kB
Slab: 968820 kB
SReclaimable: 750428 kB
SUnreclaim: 218392 kB
KernelStack: 14352 kB
PageTables: 48852 kB
NFS_Unstable: 0 kB
Bounce: 0 kB
WritebackTmp: 0 kB
CommitLimit: 62185980 kB
Committed_AS: 4269472 kB
VmallocTotal: 34359738367 kB
VmallocUsed: 168740 kB
VmallocChunk: 0 kB
Percpu: 85056 kB
HardwareCorrupted: 0 kB
AnonHugePages: 436224 kB
ShmemHugePages: 0 kB
ShmemPmdMapped: 0 kB
FileHugePages: 1251328 kB
FilePmdMapped: 145408 kB
CmaTotal: 0 kB
CmaFree: 0 kB
HugePages_Total: 0
HugePages_Free: 0
HugePages_Rsvd: 0
HugePages_Surp: 0
Hugepagesize: 2048 kB
Hugetlb: 0 kB
DirectMap4k: 1076396 kB
DirectMap2M: 47042560 kB
DirectMap1G: 12582912 kB

```

Figure 5: Main memory information

Now we proceed with the cache information:

```
[beta@icsslogin0 ~]$ module load hwloc
Loading hwloc/2.9.1-gcc-8.5.0-adhoag
Loading requirement: gcc-runtime/8.5.0-gcc-8.5.0-77forqa libpciaccess/0.17-gcc-8.5.0-ng764nd libiconv/1.17-gcc-8.5.0-rznb3y xz/5.4.6-gcc-8.5.0-lmjijpf zlib-ng/2.1.6-gcc-8.5.0-ztbc5xt libxml2/2.10.3-gcc-8.5.0-7kur2mf ncurses/6.4-gcc-8.5.0-njp6ntc
[beta@icsslogin0 ~]$ hwloc-ls
Machine (55GB total)
Package L#0
  NUMANode L#0 P#0 (23GB)
    L3 L#0 (25MB)
      L2 L#0 (256KB) + L1d L#0 (32KB) + L1i L#0 (32KB) + Core L#0 + PU L#0 (P#0)
      L2 L#1 (256KB) + L1d L#1 (32KB) + L1i L#1 (32KB) + Core L#1 + PU L#1 (P#1)
      L2 L#2 (256KB) + L1d L#2 (32KB) + L1i L#2 (32KB) + Core L#2 + PU L#2 (P#2)
      L2 L#3 (256KB) + L1d L#3 (32KB) + L1i L#3 (32KB) + Core L#3 + PU L#3 (P#3)
      L2 L#4 (256KB) + L1d L#4 (32KB) + L1i L#4 (32KB) + Core L#4 + PU L#4 (P#4)
      L2 L#5 (256KB) + L1d L#5 (32KB) + L1i L#5 (32KB) + Core L#5 + PU L#5 (P#5)
      L2 L#6 (256KB) + L1d L#6 (32KB) + L1i L#6 (32KB) + Core L#6 + PU L#6 (P#6)
      L2 L#7 (256KB) + L1d L#7 (32KB) + L1i L#7 (32KB) + Core L#7 + PU L#7 (P#7)
      L2 L#8 (256KB) + L1d L#8 (32KB) + L1i L#8 (32KB) + Core L#8 + PU L#8 (P#8)
      L2 L#9 (256KB) + L1d L#9 (32KB) + L1i L#9 (32KB) + Core L#9 + PU L#9 (P#9)
HostBridge
  PCI 01:00.0 (Ethernet)
    Net "eno1"
  PCI 01:00.1 (Ethernet)
    Net "eno2"
  PCI 02:00.0 (InfiniBand)
    Net "ibp"
  OpenFabrics "qib0"
  PCI 00:11.4 (SATA)
  PCI 07:00.0 (VGA)
  PCI 00:1f.2 (SATA)
  Block (disk) "sda"
Package L#1
  NUMANode L#1 P#1 (31GB)
    L3 L#1 (25MB)
      L2 L#10 (256KB) + L1d L#10 (32KB) + L1i L#10 (32KB) + Core L#10 + PU L#10 (P#10)
      L2 L#11 (256KB) + L1d L#11 (32KB) + L1i L#11 (32KB) + Core L#11 + PU L#11 (P#11)
      L2 L#12 (256KB) + L1d L#12 (32KB) + L1i L#12 (32KB) + Core L#12 + PU L#12 (P#12)
      L2 L#13 (256KB) + L1d L#13 (32KB) + L1i L#13 (32KB) + Core L#13 + PU L#13 (P#13)
      L2 L#14 (256KB) + L1d L#14 (32KB) + L1i L#14 (32KB) + Core L#14 + PU L#14 (P#14)
      L2 L#15 (256KB) + L1d L#15 (32KB) + L1i L#15 (32KB) + Core L#15 + PU L#15 (P#15)
      L2 L#16 (256KB) + L1d L#16 (32KB) + L1i L#16 (32KB) + Core L#16 + PU L#16 (P#16)
      L2 L#17 (256KB) + L1d L#17 (32KB) + L1i L#17 (32KB) + Core L#17 + PU L#17 (P#17)
      L2 L#18 (256KB) + L1d L#18 (32KB) + L1i L#18 (32KB) + Core L#18 + PU L#18 (P#18)
      L2 L#19 (256KB) + L1d L#19 (32KB) + L1i L#19 (32KB) + Core L#19 + PU L#19 (P#19)
```

Figure 6: Cache memory information

Where it is possible to see that, since we are reporting about the login node, it has 55 GB of total memory, and it is composed of 2 "NUMA" nodes with 10 cores each. Its total memory is distributed such as 23 GB of RAM memory belonging to the first package and 31 to the second one. Then, each package has an L3 cache memory of 25 MB each, and it is shared among all the cores each of the packages. The L2 cache memory has 256 KB and L1 cache memory has 32 KB (for each L1i - instruction and L1d - data) on each package too.

This can be visualized in the downloaded figure as required:⁷

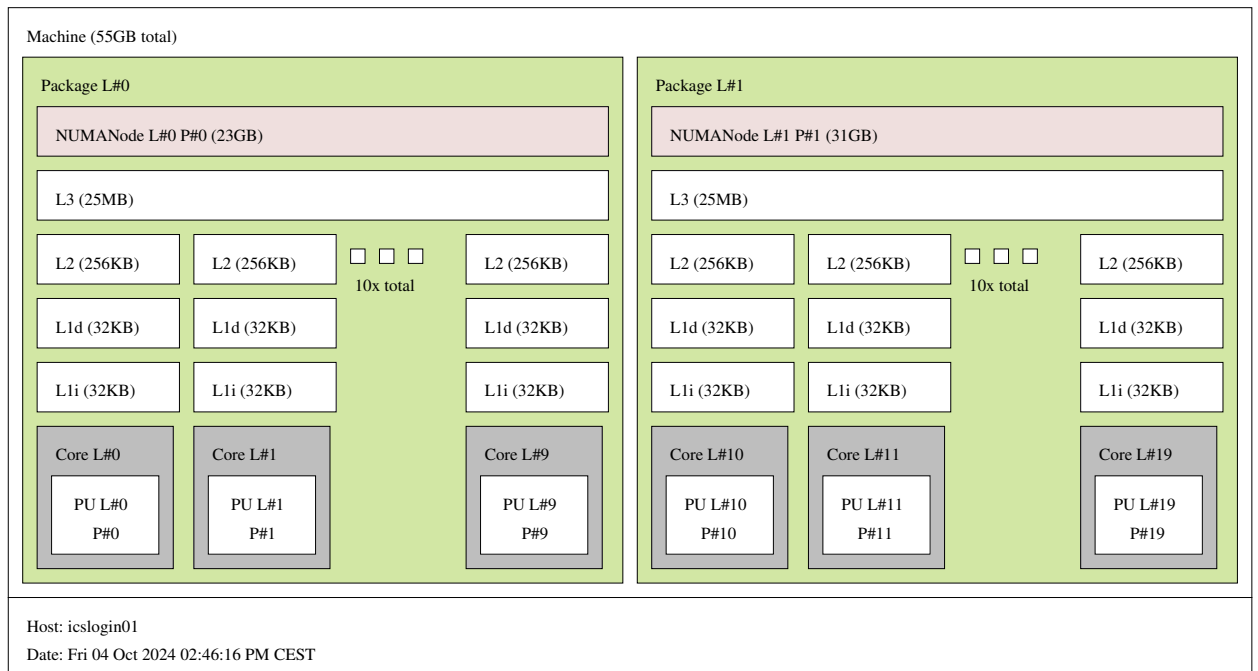


Figure 7: Resulting PDF figure from ROSA Cluster - memory Hierarchies

2.3. Bandwidth: STREAM benchmark

After downloading the proper files (stream.c and mysecond.c), I proceed with the calculations required to run properly the file. In particular, as it is mentioned in the instructions, each array must be at least four times the size of the L3 cache memory. Therefore, since each L3 cache memory is 25 MB, then the size of each array must be at least 100 MB. In line 176 of the stream.c file, it is

⁷I had troubles downloading it through the terminal with the provided command given permission issues, but I could download it through VS Code.

defined that the stream type is a double and since a double in C is 8 bytes (64 bits), then 100 MB / 8 = 12.5 M doubles (12500000). To approximate it into a power of two number, we can choose 12800000. Notice that this number is one magnitude smaller than the provided in the example, but it can be seen that both provide similar results, since, as calculated, the size of the array is at least four times the size of the L3 cache memory in both cases.

Now, I run the commands with the specified changes:

```
[bellaj@icslogin01 exe2]$ module load gcc
[bellaj@icslogin01 exe2]$ gcc -O3 -march=native -DSTREAM_TYPE=double -DSTREAM_ARRAY_SIZE=128000000 -DNTIMES=20 stream.c -o stream_c.exe
[bellaj@icslogin01 exe2]$ ls
mysecond.c stream.c stream_c.exe stream.f
[bellaj@icslogin01 exe2]$ sbatch --mem-per-cpu=4G --wrap "./stream_c.exe"
Submitted batch job 9844
[bellaj@icslogin01 exe2]$ ls
mysecond.c slurm-9844.out stream.c stream_c.exe stream.f
```

Figure 8: Terminal commands on ROSA cluster to obtain the STREAM benchmark - same I did with 12800000

With the resulting benchmarks as shown in the file *slurm9844.out* and *slurm9846.out*:

```
=====
STREAM version $Revision: 5.10 $
=====
This system uses 8 bytes per array element.
=====
Array size = 128000000 (elements), Offset = 0 (elements)
Memory per array = 976.6 MiB (= 1.0 GiB).
Total memory required = 2929.7 MiB (= 2.9 GiB).
Each kernel will be executed 20 times.
The *best* time for each kernel (excluding the first iteration)
will be used to compute the reported bandwidth.
=====
Your clock granularity/precision appears to be 1 microseconds.
Each test below will take on the order of 117053 microseconds.
(= 117053 clock ticks)
Increase the size of the arrays if this shows that
you are not getting at least 20 clock ticks per test.
=====
WARNING -- The above is only a rough guideline.
For best results, please be sure you know the
precision of your system timer.
=====
```

Function	Best Rate MB/s	Avg time	Min time	Max time
Copy:	19178.7	0.106838	0.106785	0.107000
Scale:	11218.7	0.182619	0.182553	0.182755
Add:	12294.0	0.249960	0.249878	0.250121
Triad:	12288.2	0.250086	0.249996	0.250206

```
=====
Solution Validates: avg error less than 1.000000e-13 on all three arrays
=====
```

Figure 9: STREAM benchmark results

```

=====
STREAM version $Revision: 5.10 $
=====

This system uses 8 bytes per array element.
=====

Array size = 12800000 (elements), Offset = 0 (elements)
Memory per array = 97.7 MiB (= 0.1 GiB).
Total memory required = 293.0 MiB (= 0.3 GiB).
Each kernel will be executed 20 times.
The *best* time for each kernel (excluding the first iteration)
will be used to compute the reported bandwidth.
=====

Your clock granularity/precision appears to be 1 microseconds.
Each test below will take on the order of 11704 microseconds.
(= 11704 clock ticks)
Increase the size of the arrays if this shows that
you are not getting at least 20 clock ticks per test.
=====

WARNING -- The above is only a rough guideline.
For best results, please be sure you know the
precision of your system timer.
=====

Function      Best Rate MB/s  Avg time     Min time     Max time
Copy:          19133.0    0.010742    0.010704    0.010901
Scale:         11280.7    0.018181    0.018155    0.018200
Add:           12282.6    0.025058    0.025011    0.025237
Triad:         12310.6    0.025001    0.024954    0.025033
=====
+=====
Solution Validates: avg error less than 1.000000e-13 on all three arrays
=====

```

Figure 10: STREAM benchmark results - with 12800000 array size as calculated

2.4. Performance model: A simple roofline model

Let's put together all the information we gathered to build a simple roofline model as in the provided example. On one side, we know that the y-axis of the roofline model measures the attainable floating-point performance, as gigaflops per second (Giga Floating-Point Operations per Second - GFLOP/s). On the other side, the x-axis measures the operational intensity as flops per byte of DRAM traffic ⁸ (the number of floating-point operations per byte of data moved from memory to the processor - Operational Intensity / [Flops/byte]). Therefore, low operational intensity means that the system is moving a lot of data but doing little computation - memory bound. whereas high operational intensity means that the system is performing a lot of computations with the data it moves.

As we computed in 2.1, we know that the theoretical peak performance of a core with the described characteristics in these exercises is 36.8 GFLOP/s. Whereas, using the same approach that in the provided example of taking the STREAM benchmark results of Scale, Add, Triad (treating Copy as an exception), we can average the results of these three:

$$\begin{aligned}
 \text{Average maximum bandwidth} &= \frac{11280.7 + 12282.6 + 12310.6}{3} \\
 &= 11957.97 \text{ MB/s} \\
 &= 11.96 \text{ GB/s}
 \end{aligned}$$

⁸Based on Samuel Williams, Andrew Waterman, and David Patterson. Roofline: an insightful visual performance model for multicore architectures. Communications of the ACM, 52(4):65, April 2009. <https://dl.acm.org/doi/10.1145/1498765.1498785> and class slide.

Following the assignment explanation and the mentioned paper, the sloped line in the graph represents the memory bound region and by consequence, the performance there is restricted by how fast the system can move data between memory and the processor (bandwidth - slope). Then, the machine balance is the operational intensity at which the memory bandwidth is saturated. This is the point where the performance is limited by the memory bandwidth and the peak performance of the processor is reached, the memory bound region meets the peak performance line. To put this into math terms since it makes it easier to understand:

$$\text{Attainable GFlops/sec} = \min(\text{Peak Floating-Point Performance}, \\ \text{Peak Memory Bandwidth} \times \text{Operational Intensity})$$

By consequence, on the left side of the graph, the expression summarize that the attainable performance is proportional to how fast it can transfer data from memory (bandwidth) and the FLOPs it can perform per byte of data (Operational Intensity). Another way of seeing this is that for every additional FLOP per byte of data transfer, the performance increases by a factor equal to the memory bandwidth.

Now, in the compute-bound region, the performance is capped by the peak core performance (independently of the data transfer)

$$\text{Attainable performane (GFLOP/s)} = \text{Peak performane (GFLOP/s)}$$

Then, the ridge, as defined in the assignment, is the point where the two regions meet:

$$\text{Peak performance (GFLOP/s)} = \text{Peak Memory Bandwidth (GB/s)} \times \text{OI (FLOP/Byte)}$$

$$36.8 \text{ GFLOP/s} = 11.96 \text{ GB/s} \times \text{Operational Intensity (FLOP/Byte)}$$

$$\begin{aligned} \text{Operational Intensity (FLOP/Byte)} &= \frac{36.8 \text{ GFLOP/s}}{11.96 \text{ GB/s}} \\ &= 3.08 \text{ FLOP/Byte} \end{aligned}$$

Finally, the ridge point is at 3.08 FLOP/Byte. The following graph summarize the calculations as required: ⁹

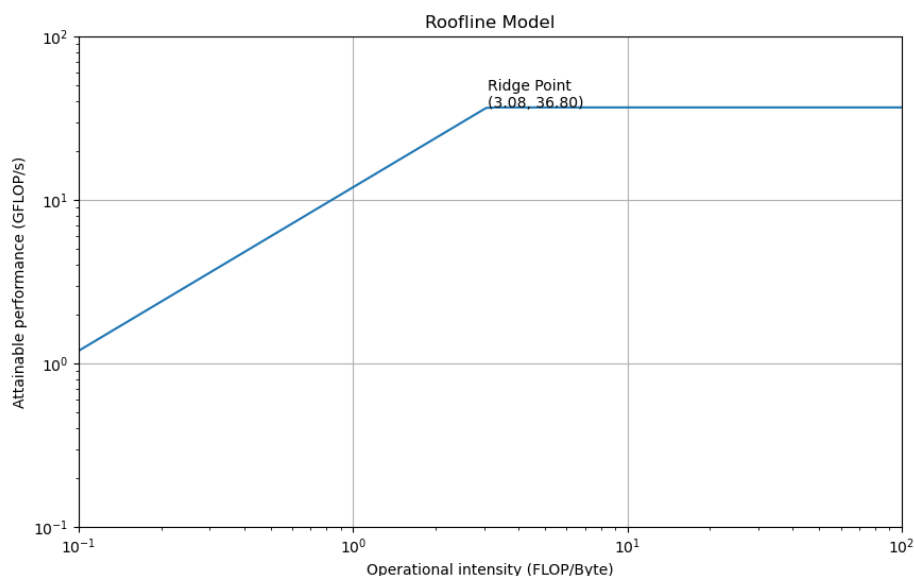


Figure 11: Resulting Roofline plot

⁹The code of the plot is named *plot.py* and the resulting plot is named *roofline-plot.png*

In conclusion, given the data, kernels with operational intensity with less 3.08 FLOP/Byte will benefit the most from optimizations that improve memory access patterns or reduce memory traffic by restructuring loops for unit stride accesses, ensuring memory affinity, and using software prefetching.¹¹ Whereas, kernels with operational intensity with more than 3.08 FLOP/Byte could improve its computational efficiency by improving instruction-level parallelism (ILP) and applying SIMD.

3. Optimize Square Matrix-Matrix Multiplication (50 Points)

3.1. Blocked matrix multiplication - Optimization

To implement blocked matrix multiplication I started by computing the maximum size of each block such that it fits on the L3 cache memory.¹²

Since the L3 cache memory is 25 MB which is equivalent to 26214400 bytes¹³

Then the maximum size of each block can be calculated as $\sqrt{\frac{26214400}{3 \times 8 \text{ bytes}}}$ since we are working with 3 matrices (A, B, C) that are symmetric and that contains double values (8 bytes). Therefore, the $(block_size)^2 = (\frac{26214400}{3 \times 8 \text{ bytes}})$.

After this calculation, on the main function *square_dgemm*, I implement the block as follows:

First, I check if the block size is bigger than the matrix size, since in that case then we should store the matrix completely in L3 and avoid indexing errors. Then, each of the three loops increments by calculated block size, taking into account the edge case that emerge when one of the loops block jumps can surpass the length of the matrix. If that occurs, the faction that computes the actual multiplication will receive the appropriate size of the corresponding block to not surpass the matrix size when indexing.

Most importantly, the order of indexing follows the column-major order as specified in the assignment. Citing from the main class book¹⁴ "If an inner loop variable is used as an index to a multidimensional array, it should be the index that ensures stride-one access".

¹⁰I used python to generate the graph since, according to what I could found about graph generation in C++, the most common libraries used as matplotlib-cpp calls for python during compilation anyway. <https://stackoverflow.com/questions/63667255/plotting-graphs-in-c>

¹¹Based on mentioned paper.

¹²I wrote the majority of my code based on the pseudocode from the last slide of HPC - "Blocked matrix multiplication", Chapter 3 of Georg Hager and Gerhard Wellein. Introduction to high performance computing for scientists and engi- neers; and the example code provided about AVX on previous to last class.

¹³<https://www.gbmb.org/mb-to-bytes>

¹⁴Introduction to High Performance Computing for Scientists and Engineers by Georg Hager and Gerhard Wellein, chapter 3, page 70

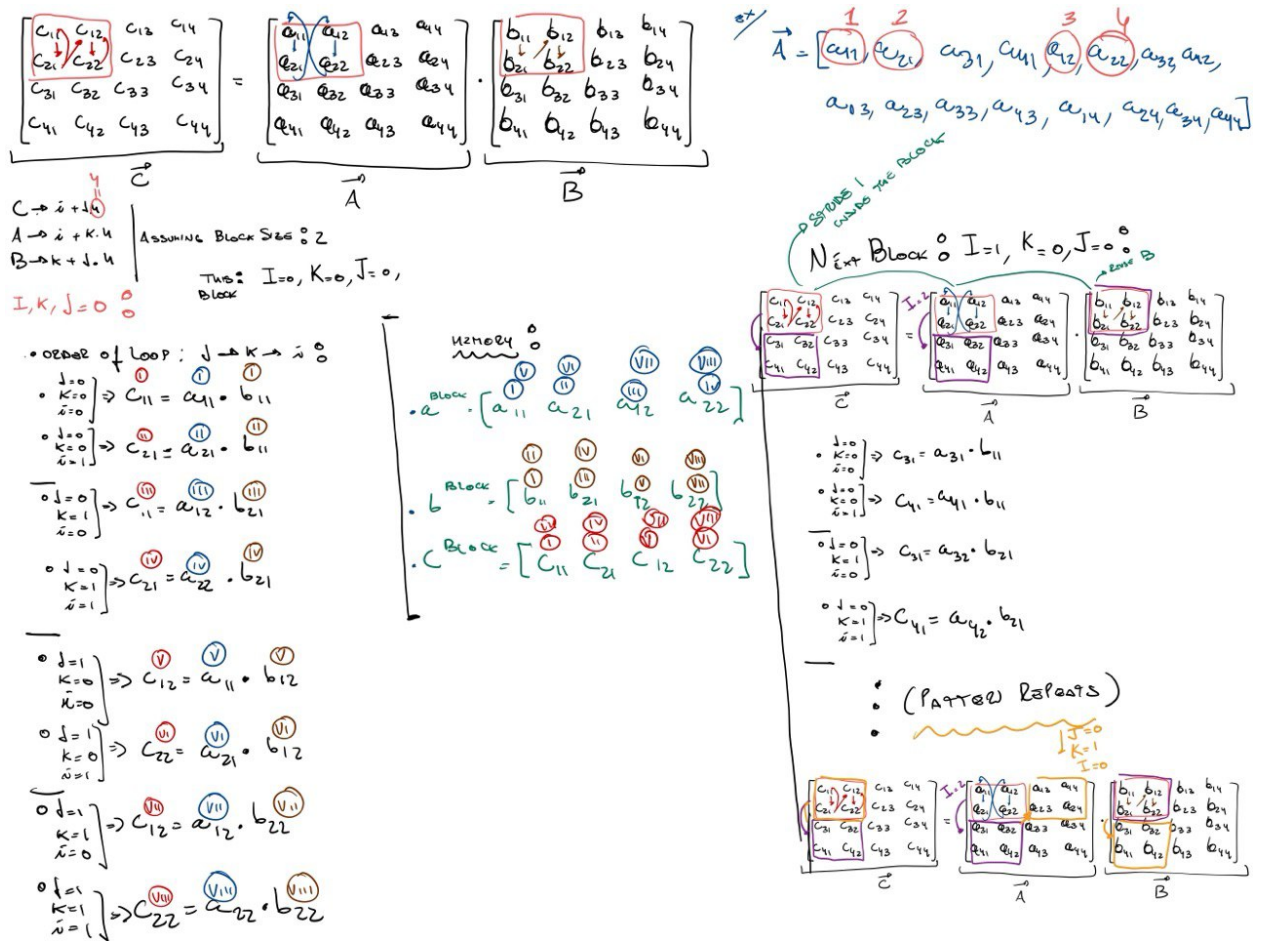


Figure 12: An illustration and summary of how I thought the implementation

Since we pass a pointer with the address of the beginning of each matrix (double* A, as we saw in previous assignment) then, we need to move in each block iteration the pointer to the right position of the block. This means, as an example, to take the first block, process it in the way I will describe below and then move the pointer to the beginning of the next block.

The block multiplication function call *matrix_multiply_block* receives the matrix size, the block size of each matrix (by debugging the edge cases), and the matrices. Following again the column major order as described in the book, which means that if a matrix, for example, of 3x3 is store in memory as:

$[a_{11}, a_{21}, a_{31}, a_{12}, a_{22}, a_{32}, a_{13}, a_{23}, a_{33}]$

To index for a particular position, we can think of "jumping" to the right column $j \cdot n$ and then move I step into the column row that is expected. However, we need to do this in a way that we can access the memory efficiently. Remembering how our matrices are stored, then we should make operations such that each movement is done as going down in rows into a column and then moving to the next one such that I access the C and A consecutive elements in column (stride one access) while keeping B elements constant (same element - repeatedly) until that element will not be necessary anymore for operations.

As we can see in my handwritten illustration (Figure 12), which shows the process that my implementation follows and its benefits over the naive implementation, we can see that takes into advantage two main things:

Loading a block into cache once and reusing it multiple times - temporal locality. As it is the case with the B matrix. And stride one access to the elements in the corresponding block of A and C - spatial locality. Iterating on each block in the column-major order.

This provided a substantial increase over the naive implementation but still far away with the provided by blas library.

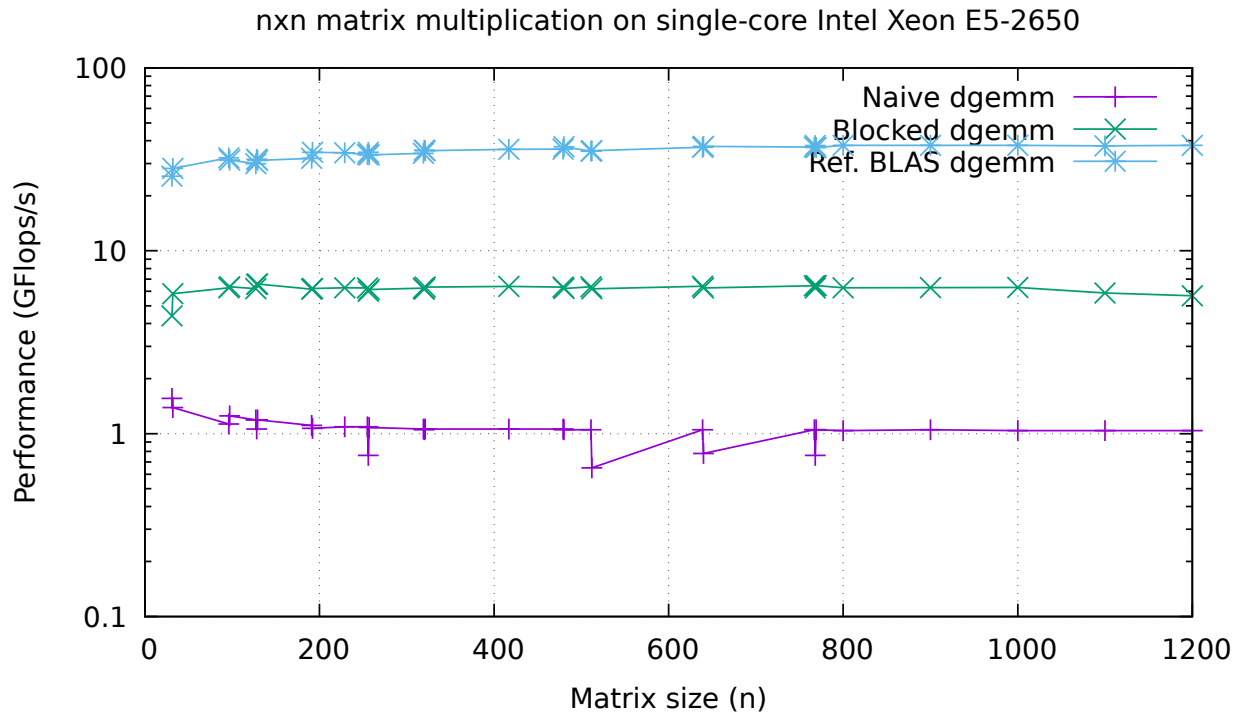


Figure 13: Performance of the different implementations - Average percentage of peak performance = 16.8045

However, there is also an additional slight improvement exist if we choose to keep fixed C block instead of B as we saw in class. This means that the loop order of the blocks is changed to J - I - K, where K is the innermost one. Given that in my code K index for the columns of A and J for the rows of B (as we saw in last class example, slide 44 - "*blocked_matrix_multiply*"). Notice that for the multiplication itself, the loop order does not change since we need to keep stride one access for column-major.

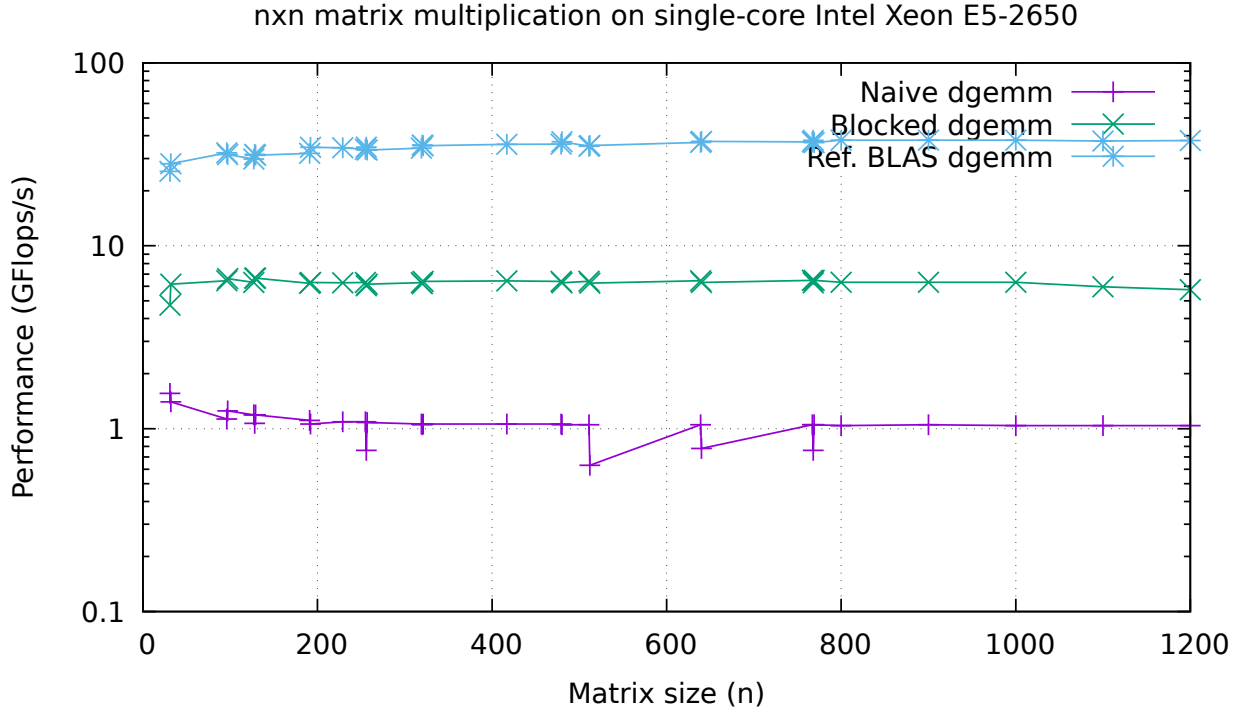


Figure 14: Performance of the different implementations - Average percentage of peak performance = 17.0158

The corresponding codes can be found in the folder *exercise3*, where, in *block* folder it is the B fixed block implementation of the illustration whereas in *blockCfixed* can be found the C fixed implementation.

3.2. AVX optimization

Considering these results, I tried to improve my performance by using the AVX instructions.¹⁵

Notice that the same illustration (Figure 12) allow us to visualize how to construct this. As i mentioned before, i constructed *square_dgemm* such that it is not necessary to modify it later on, since the later optimization can be implemented directly on the multiplication related loops, not on the blocks.

With respect to *matrix_multiply_block*, the loops are almost unchanged, however, there are two main differences. Since we can process 4 doubles at a time with AVX (256-bits), we need to take this into account. First, as i noted when i built the previous one, B (or C - with loop order change as in previous implementation) matrix element will be repeated until passing to the next j or k index. This pattern emerges since the independence of i which is the inner loop. Then, we can broadcast the value of B (or C) into 4 and keeping it constant through the loop over C and A sub-blocks respectively.

The second main difference, is that we need to account for 4 loads at once in the internal loop. Which means that we loop over the index on jumps of 4. Which also means that we should account for the edge case where the block is not divisible by 4. This explains why I define *i* outside I loop, such that when the "expected" case is finish, it keeps the value where it finishes and pass to the next loop that will manage that fragment of maximum 3 elements each in the "normal" way. Notice that the loading of the C and A value follows the same that in the class example of AVX. The difference here is on *_mm256_fmadd_pd*, since we need to multiply A and B sub-blocks and then add the result to the C sub-block. The function that does this, thanks to ETH slide set¹⁶, is the

¹⁵For this implementation i based on the following documents: Class code example on avx, <https://acl.inf.ethz.ch/teaching/fastcode/2021/slides/07-simd-avx.pdf>, and <http://const.me/articles/simd/simd.pdf>

¹⁶<https://acl.inf.ethz.ch/teaching/fastcode/2021/slides/07-simd-avx.pdf>

previously mentioned one. Finally, I store the results in the same manner as class example.

It is finally not that exiting to report that the result was not as expected, since the implementation provided the same CPU usage as the previous one. I am still thinking why this could be. I suspect that in part I could be using the wrong compiling configuration:

Listing 1: Makefile example

```
CC = gcc                                # compiler (& linker)
OPT = -O3 -march=core-avx2 -mavx \
      -funroll-loops \
      -ftree-vectorize                  # optimization flags (you may add more)
CFLAGS = -Wall -std=gnu99 $(OPT) # standard compiler flags
```

I tried using the commented Intel config for this, however, I keep getting an error related to not being able to load icx or icc. I tried to compensate this by adding the flag -mavx as it was in the AVX class example but still I do not see any improvement, which could also mean that my AVX implementation is not correct.¹⁷

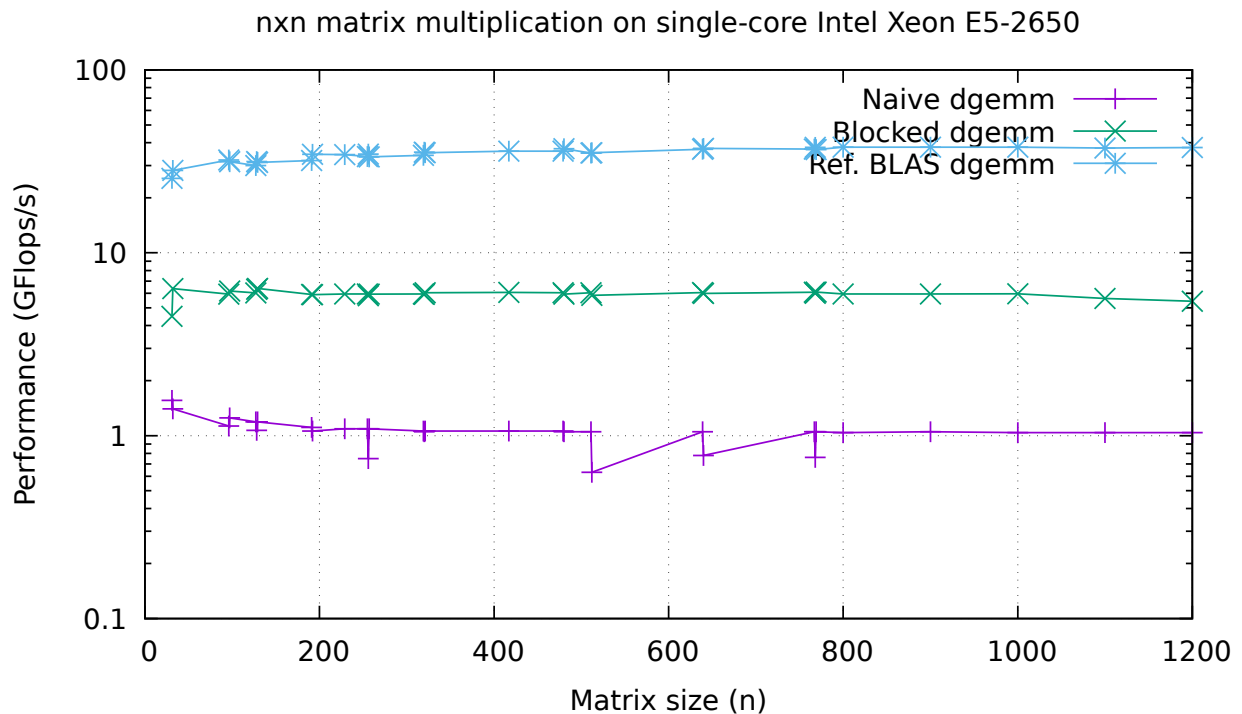


Figure 15: Performance of the different implementations with AVX - Average percentage of peak performance = 16.1474

The corresponding code can be found in the folder *exercise3*, where, in *avx* subfolder you will find the AVX related code whereas in *block* the original blocked matrix multiplication implementation. The *"provided"* subfolder contains the rest of the code provided by the assignment.

4. Quality of the Report

(15 Points)

¹⁷To run this code, I worked in a separate file call *dgemmblockedavx*, however, for running I copy and paste it inside the *dgemmblocked* file to not modify the *run_matrixmult.sh*.