

Exercises for High Performance Computing (MA-INF 1108) WS 2023/2024

E. Suarez, M. Wolter and B. Kostrzewa
Tutors: O. Vrapcani and N. Pillath

6 Memory Bandwidth

ATTENTION!: eCampus can become unavailable without previous announcement due to an urgent maintenance. Students are responsible for submitting their checklists enough ahead of time of the deadline. Submissions via Email will not be accepted unless tutors explicitly authorized it beforehand.

This exercise will be performed on the **JURECA-DC** system.

Important recommendations for benchmarking runs (exercises 6.2 and 6.3):

- Try first with just 4-5 different vector lengths until you are sure all your scripts are doing what they should.
- For the final run (after debugging), increase in vector length in factors of 2, starting with a size for which the full problem fits in the L2\$, and ending with a size that exceeds the size of the L3\$.
- Include in the beginning of your output files a header containing in a few lines all important information of the run conditions, e.g.: name and version of source code, compiler version and flags used, node number on which it run, number of threads and/or processes, etc.
- Include in your output plots title, axis labels (with units) and legend, so that it is clear what is being plot and how the results have been obtained. See listings example.

```
1 =====
2 DAXPY daxpy.c version: 2023
3 gcc -O1 -fopenmp -info
4 gcc (GCC) 11.3.0
5 Scalar Version (no OpenMP or MPI)
6 Run: Node: jrc0736 ; ntasks: 1
7 =====
8      Size      BW[GB/s]      Runtime[ms]
9      2048      39.961      0.001
10     4096      40.289      0.002
11 [ ... ]
```

1: Theoretical Peak Performance

Estimate the maximum peak performance of a standard compute node in the JURECA-DC.

- a) [1pt] Use any of the tools presented in the exercise sheet #3 (*Ex03 Hardware characteristics of the processor*) and/or the JSC-system documentation to find out following parameters:
- Vendor
 - CPU model
 - Cores per socket
 - Threads per core
 - Basis clock frequency
 - Turbo clock frequency
 - Capacity of L1D\$ per core
 - Capacity of L2\$ per core
 - Capacity of L3\$ per core
 - Capacity of main memory
- b) [1pt] Find out through the chip specifications (e.g. via the vendor website or reliable web sources as `wikichip`):
- Amount of FMA units
 - Maximum SIMD vector length supported
 - Type of DDR memory
 - Maximum memory frequency
 - Number of memory channels
 - Maximum memory bandwidth
- c) [1pt] Estimate the maximum performance per core.
- d) [1pt] Estimate the maximum performance for the whole socket.
- e) [1pt] Most HPC applications are memory bound. Traditionally, the ideal ratio between memory bandwidth and performance was defined to be 1 Byte/FLOP. Estimate the value for our JURECA-DC CPU (socket).

Solution:

- a) [1pt] JURECA-DC standard compute node, most data can be obtained from `lstopo`.
- Vendor: AMD
 - CPU model: AMD EPYC 7742
 - Cores per socket: 64

- Threads per core: 2
- Basis clock frequency: 2.25 GHz
- Turbo clock frequency: 3.4 GHz
- Capacity of L1D\$ per core: 32 KB
- Capacity of L2\$ per core: 512 KB
- Capacity of L3\$ per core: 16 MB
- Capacity of main memory: 512 GB (lstopo shows 503 GB free)

b) [1pt] Additional documentation

(<https://en.wikichip.org/wiki/amd/epyc/7742>,
https://en.wikichip.org/wiki/amd/microarchitectures/zen_2),
<https://www.amd.com/en/products/cpu/amd-epyc-7742>:

- Amount of FMA units: 2 (from e.g. https://en.wikichip.org/wiki/amd/microarchitectures/zen_2)
supported : 256 – bit(AVX – 256)(fromwikichip)
- Type of DDR memory: DDR4
- Maximum memory frequency: 3200 MT/s (from AMD web)
- Number of memory channels: 8 (from AMD web)
- Maximum memory bandwidth: 204.8 GB/s (from AMD web)

c) [1pt] Maximum performance per core can be estimated with:

$$P_{core} = f_{turbo} \cdot n_{super}^{FP} \cdot n_{FMA} \cdot n_{SIMD} \quad (1)$$

where, f_{turbo} is the turbo frequency, n_{FMA} is the number of fuse multiply add execution units, and n_{SIMD} is the maximum SIMD length, measured on how many double precision numbers (DP) fit into the SIMD unit. n_{super}^{FP} is the amount of floating point operations that come out of the superscalar pipeline. The execution unit that delivers the most operations per cycle is FMA, which does simultaneously an addition and a multiplication.

For the AMD EPYC 7742 processor, using the data collected before:

- $f_{turbo} = 3.4$ GHz
- $n_{super}^{FP} = FLOP$, since the FMA unit can do 2 operations simultaneously.
- $n_{FMA} = 2$, there are two FMA units (see e.g. https://www.nas.nasa.gov/hecc/support/kb/amd-rome-processors_658.html)
- $n_{SIMD} = 4$, since AMD supports AVX-256, which fits 4 DP numbers.

Resolving:

$$P_{core} = 3.4\text{GHz} \cdot 2 \text{ FLOP} \cdot 2 \cdot 4 = 54.4 \text{ GFLOP/s} \quad (2)$$

- d) [1pt] Maximum performance for the whole socket, can be calculated with

$$P_{socket} = f_{allcores} \cdot n_{super}^{FP} \cdot n_{FMA} \cdot n_{SIMD} \cdot n_{cores} \quad (3)$$

where the number of cores in our case is $n_{cores} = 64$. Notice that the frequency is changed by the value when using all cores, which falls down to the basis frequency: $f_{allcores} = 2.25$ GHz.

Resolving:

$$P_{socket} = 2.25 \text{ GHz} \cdot 2 \text{ FLOP} \cdot 2 \cdot 4 \cdot 64 = 2304 \text{ GFLOP/s} = 2.3 \text{ TFLOP/s} \quad (4)$$

- e) [1pt] Maximum Byte/FLOP ratio:

$$\text{Byte per FLOP} = \frac{\text{Memory Bandwidth}}{\text{Peak Performance}} = \frac{204.8 \text{ GB/s}}{2304 \text{ GFLOP/s}} = 0.09 \text{ B/FLOP} \quad (5)$$

This CPU (and most CPUs nowadays) is therefore far away from the ideal ratio of 1 B/FLOP.

2: Implement daxpy

You will implement and benchmark `daxpy`, a function that performs following operation on vectors:

$$\vec{z} = a \cdot \vec{x} + \vec{y}$$

where a is a scalar, and \vec{z} , \vec{x} , and \vec{y} are vectors of a given length N containing double precision numbers.

- [2pt] Solve the TODOs in the source code `daxpy.c`.
- [2pt] Compile the source code into an executable. Create a script (`run_daxpy.sh`) to run the executable (on the `dc-cpu-devel` partition of the JURECA-DC cluster) for increasing vector lengths, starting by a value for which all vectors fit in the L1D\$, and ending by a value in which they do not fit into the L3\$. The output of this script should be saved into the file (`daxpy.txt`) in your solutions folder.
- [1pt] Write a script (e.g. `plot_daxpy.py` in python or matplotlib) that plots the memory bandwidth usage in GB/s vs. the memory footprint of `daxpy`. Hint: the memory footprint is the amount of memory used by the function, which you can calculate from the vector size and amount of reads or writes performed per element. You can neglect the impact of the scalar. Store the plot as `daxpy.png`, in your solutions folder.
- [1pt] Add to the plot (`daxpy.png`, in your solutions folder) vertical lines showing the capacity of the L1D\$, L2\$, and L3\$, plus an horizontal line for the memory bandwidth of this CPU according to its specifications (values from exercise 6.1). What do you observe?

- e) [4pt] Use the pragma `pragma omp parallel for` before your `for` loops, save the source file as `daxpy_omp.c` and compile it. Write a script (`run_daxpy.sh`) to run `daxpy` for increasing number of OpenMP threads, starting by 1 and ending by the maximum number of hardware threads supported by the CPU, using its SMT capabilities. Store the resulting text files on your solutions folder following the naming convention `daxpy_omp<ThreadNum>.txt` (e.g. `daxpy_omp1.txt`, `daxpy_omp2.txt`, etc.). Plot the curves in (`daxpy_omp.png`) in your solutions folder). What do you observe?

Solution:

See solution folder in GitHub

3: The Stream benchmark

With this exercise you will learn to run the `stream` benchmark, developed by John D. McCalpin. This benchmark suite has been designed to measure the memory bandwidth on almost any computer platform. It calculates three different vector operations:

Copy: $\vec{y} = \vec{x}$

Scale: $\vec{y} = c \cdot \vec{x}$

Add: $\vec{y} = \vec{x} + \vec{z}$

Triad: $\vec{y} = \vec{x} + c \cdot \vec{z}$

where \vec{x} , \vec{y} , and \vec{z} are vectors in double precision of a given length, and c is a scalar.

First of all: Download the stream benchmark from its original website (<https://www.cs.virginia.edu/stream/>) or the git repository:

<https://github.com/jeffhammond/STREAM.git>

- a) [1pt] Read the provided `README` file and adapt the source file `stream.c` to make sure that you run the benchmark under its defined specifications. Compile it also according to the description.
- b) 1pt Write a script (`run_stream.sh`) to run the `stream` benchmark on a **standard compute node of JURECA-DC** for growing vector sizes (vector length increasing in factors of 2), from a size fitting in the L1D\$ up to sizes that do not fit in the L3\$. Follow the recommendations at the top of the exercise sheet.
- c) [1pt] Write a script that extracts the following from the standard output files of `stream`: the vector size, and the **best rate** for the three benchmarks contained in the suite. Print these four values into four columns in a text file (`stream_bw.txt`) in your solutions folder.
- d) [1pt] Write a script (e.g. `plot_stream.py` in python or matplotlib) that plots (`stream.png`, in your solutions folder) the memory bandwidth (in GB/s) measured by the four `stream` benchmarks vs. their memory footprint. Include in the plot vertical lines marking the physical capacity of the L1D\$, L2\$, and L3\$ caches (values from exercise 6.1), and an horizontal line for the memory bandwidth as given in the chip specifications.

- e) [1pt] What are the bandwidths measured at the L1D\$, L2\$, L3\$, and main memory? Do the latter match with the specifications of memory bandwidth that you found out in exercise 6.1?
- f) [3pt] Run now the **stream** benchmarks using OpenMP, following the indications given in the source file. Redo the previous steps, saving the **best rate** for increasing number of OpenMP threads in text files (**stream_omp1.txt**, **stream_omp2.txt**, etc.). Plot the copy, scale, add, and triad results for the run with 64 OpenMP threads (stream_omp64.png), including vertical lines for the L1\$, L2\$, and L3\$ capacities, and an horizontal line for the memory bandwidth given by the chip specs. Compare your result with the plot you previously obtained with the sequential version.
- g) [1pt] Produce a plot (**triad_omp.png**) of bandwidth vs. memory footprint with the stream triad results alone, but for increasing number of OpenMP threads.

Solution:

See solution folder in GitHub

Commit your solutions to the GitHub Classroom.

If you have used Jupyter, close your Jupyter session and stop JupyterLabs.