

Praktikum im WS 2023

Evaluierung moderner HPC-Architekturen und -Beschleuniger
(LMU)

Evaluation of Modern Architectures and Accelerators (TUM)

Sergej Breiter, MSc., Minh Thanh Chung, MSc., MSc., Amir Raoofy, MSc., Bengisu Elis, MSc.,
Dr. Karl Furlinger, Dr. Josef Weidendorfer

Assignment 05 – Due: 23.11.2023

In Assignment #2, you worked with GPUs in BEAST and specifically focused on the performance characteristics (peak performance, memory bandwidth) of the Vector Triad. In this assignment, we focus on Matrix Multiplication (MxM) on GPUs. For this purpose, you will adapt the code to use OpenMP offloading directives and to utilize GPU resources in BEAST, i.e., AMD MI200 and Nvidia A100.

Before going ahead with tests and analysis on the microbenchmarks, it is important to keep in mind that we investigate Matrix Multiplication as a compute-bound benchmark. This information may be helpful when explaining the observed results from the MxM microbenchmark.

Offloading Matrix Multiplication (MM) with OpenMP

For this set of tasks, please use the provided matrix-matrix multiplication (MM) implementation called `assignment5.cpp`. This implementation already offloads MM to GPU by OpenMP directives.

1. **Offloading with Optimal Data Transfer Policy:** Make sure that the MM computation is offloaded to GPU. Take a look back at Assignment #2, choose the best data initialization policy and make sure this policy is used to initialize and migrate your matrices to GPU.
2. **Optimizations:** Unlike Assignment #4, where we used loop interchange (reordering the loops) as the first step, here we use another technique for optimization. You need to introduce two variants of matrix multiplication benchmark, both with the same “*ijk*” loop ordering.
 - **Variant 1 :** Store the elements of both multiplicand **B** and multiplier **C** in row-major layout.
 - **Variant 2 :** Use row-major layout for multiplicand **B** and column-major layout for multiplier **C**.

For both variants, complete the following tasks:

- (a) Apply cache blocking (similar to Assignment #4)

- (b) Use appropriate loop scheduling policies and vectorization directives based on your experience.
 - (c) Run experiments to measure FLOP rates and memory bandwidth utilization of your code on GPU.
 - (d) Explain your results and compare them to the results you achieved in Assignment #4, for matrix multiplication on CPUs.
3. **Execution Configuration on Target Device:** As you learned how to use OpenMP clauses that configure the execution of teams and threads on target devices, please apply them for the following subtasks:
- (a) Make sure parallelism is achieved by using OpenMP clauses on the GPUs by setting the number of teams = 100 and the number of threads = 128.
 - (b) Use the league and team size combinations from the following set: $\{(t, T)\} = \{8, 16, 32, 40, 48, 64, 72, 80, 88, 96, 104, 112, 120, 128, 256, 512, 1024\} \otimes \{32, 64, 80, 128, 256, 512, 1024\}$ and plot the flop rate vs. league/team size combination plot in 3D (or heatmap). Here, “t” denotes the number of threads per team, “T” denotes the number of teams. Perform this task only for a single matrix size for each thread and team size combination, but make sure the matrix size is large enough. Additionally, if you have better ideas about the above sets of (t, T), please just suggest a combination. Explain your observations for each system on BEAST with GPUs and compare your results with the corresponding results in Assignment #2 about the optimal combination of league and team size.
 - (c) Find the optimum team number vs thread number configuration at which you get the maximum performance and explain the reason for it. You can use this optimal combination for the following task.

4. Power Measurements:

To take power measurements, use The Data Center Data Base (DCDB) system monitoring framework which was introduced in the previous assignments.

- To recall using DCDB, you should run the following command on the gateway:
`source /opt/dcdb/dcdb.bash.`
- You can list available AC power sensors by the following command:
`dcdbconfig sensor list|grep power_ac.`
- Sensor values and timestamps can be read by the following command:
`dcdbquery /beast/<node name>/<sensor name> now-5m now`
This command outputs the average power measurements from 5 minutes before until now, reported by DCDB every 5 seconds.

Power Distribution Unit (PDU) is a device with multiple power outlets designed to distribute electric power, especially to racks of computers and networking equipment located within a data center. Power consumption from each outlet can be measured, for example, the power measured for by outlet #5 of PDU #3 is measured and stored in DCDDDB as sensor “(/beast-/pdu3/power5)”. Note that each BEAST node has 2 PDUs for redundancy, and the total power consumption of a node is the sum of measurements taken from both PDUs. The *power_ac* virtual sensor output performs this summation of two PDUs. In addition, the power measurements reported by DCDB are obtained by averaging PDU-internal power measurements that are taken every 30 msecs.

Also note that for correct and interference-free measurements, you need to determine and impose sleep time before and after your benchmark runs.

For CPU only and best-performing GPU versions of your code:

- (a) Present your results for speed up and power consumption, in a table.
- (b) Plot Mflops per Watt (Mflops/Watt) vs data set size results, similar to previous assignments' Mflops/sec vs data set size plots.

References

- [1] David B. Kirk and Wen-mei W. Hwu. 2010. Programming Massively Parallel Processors: A Hands-on Approach (1st. ed.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.