# Module 21 Lab

# OpenACC CUDA Vector Add

*GPU Teaching Kit – Accelerated Computing*

## OBJECTIVE

Implement a vector addition using OpenACC directives.

## LOCAL SETUP INSTRUCTIONS

The most recent version of source code for this lab along with the build-scripts can be found on the Bitbucket repository. A description on how to use the CMake tool in along with how to build the labs for local development found in the README document in the root of the repository.

The executable generated as a result of compiling the lab can be run using the following command:

```
./OpenAccVectorAdd_Template -e <expected.raw> \
  -i <input0.raw>,<input1.raw> -o <output.raw> -t vector
```

where <expected.raw> is the expected output, <input0.raw>,<input1.raw> is the input dataset, and <output.raw> is an optional path to store the results. The datasets can be generated using the dataset generator built as part of the compilation process.

## LOCAL DEVELOPMENT & OBTAINING A PGI COM-PILER LICENSE

The usage of OpenACC directives requires access to the PGI OpenACC compiler. Please follow the instructions on Bitbucket repository to download the tools, generate the license file and install the license.

## QUESTIONS

(1) Compare the implementation difficulty of this kernel compared to the CUDA implementation. What difficulties did you have with this implementation?

ANSWER: **CUDA programming requires a lot of boilerplate code and runtime calls that are totally removed when using OpenACC.**

(2) What features are available to you in CUDA or OpenCL that made programming in either simpler?

ANSWER: **Answers may vary, but for vector add there is no debate that OpenACC is simpler than CUDA or OpenCL.**

## CODE TEMPLATE

The following code is suggested as a starting point for students. The code handles the import and export as well as the checking of the solution. Students are expected to insert their code is the sections demarcated with //@@. Students expected the other code unchanged. The tutorial page describes the functionality of the wb∗ methods.

```
1   #include <math.h>
2   #include <stdio.h>
3   #include <stdlib.h>
4   #include <wb.h>
5
6   int main(int argc, char *argv[]) {
7     wbArg_t args;
8     float *__restrict__ input1;
9     float *__restrict__ input2;
10    float *__restrict__ output;
11    int inputLength;
12
13    args = wbArg_read(argc, argv);
14
15    wbTime_start(Generic, "Importing data and creating memory on host");
16    input1 = (float *)wbImport(wbArg_getInputFile(args, 0), &inputLength);
17    input2 = (float *)wbImport(wbArg_getInputFile(args, 1), &inputLength);
18    output = (float *)malloc(inputLength * sizeof(float));
19    wbTime_stop(Generic, "Importing data and creating memory on host");
20
21    //@@ Insert vector addition code here.
22    wbTime_start(GPU, "Copy to GPU, compute, and copy back to host.");
23    wbTime_stop(GPU, "Copy to GPU, compute, and copy back to host.");
24
25    wbSolution(args, output, inputLength);
26
27    // Release memory
28    free(input1);
29    free(input2);
```

```
30      free(output);
31
32      return 0;
33    }
```

## CODE SOLUTION

The following is a possible implementation of the lab. This solution is intended for use only by the teaching staff and should not be distributed to students.

```
1    #include <math.h>
2    #include <stdio.h>
3    #include <stdlib.h>
4    #include <wb.h>
5
6    int main(int argc, char *argv[]) {
7      wbArg_t args;
8      float *__restrict__ input1;
9      float *__restrict__ input2;
10     float *__restrict__ output;
11     int inputLength;
12
13     args = wbArg_read(argc, argv);
14
15     wbTime_start(Generic, "Importing data and creating memory on host");
16     input1 = (float *)wbImport(wbArg_getInputFile(args, 0), &inputLength);
17     input2 = (float *)wbImport(wbArg_getInputFile(args, 1), &inputLength);
18     output = (float *)malloc(inputLength * sizeof(float));
19     wbTime_stop(Generic, "Importing data and creating memory on host");
20
21     //@@ Insert vector addition code here.
22     wbTime_start(GPU, "Copy to GPU, compute, and copy back to host.");
23     int i;
24   // sum component wise and save result into vector c
25   #pragma acc kernels copyin(input1[0 : inputLength],                    \
26                                     input2[0 : inputLength]),             \
27                                     copyout(output[0 : inputLength])
28     for (i = 0; i < inputLength; ++i) {
29       output[i] = input1[i] + input2[i];
30     }
31     wbTime_stop(GPU, "Copy to GPU, compute, and copy back to host.");
32
33     wbSolution(args, output, inputLength);
34
35     // Release memory
36     free(input1);
37     free(input2);
38     free(output);
39
40     return 0;
41    }
```