Міністерство освіти і науки України
Національний університет водного господарства та природокористування
Навчально-науковий інститут кібернетики, інформаційних технологій та інженерії
Кафедра обчислювальної техніки

04-04-278M

# METHODICAL GUIDELINES
for laboratory classes in the course
**«Parallel and Distributed Calculations»**
for higher education students of the first (Bachelor's) level
in the educational degree programme «Computer Engineering»
of the field of study 123 «Computer Engineering»
of full-time and part-time forms of study

# МЕТОДИЧНІ РЕКОМЕНДАЦІЇ
до лабораторних робіт з навчальної дисципліни
**«Паралельні та розподілені обчислення»**
для здобувачів вищої освіти першого (бакалаврського) рівня
за освітньо-професійною програмою «Комп'ютерна інженерія»
спеціальності 123 «Комп'ютерна інженерія»
денної та заочної форм навчання

Рекомендовано
науково-методичною радою
з якості ННІКІТІ
Протокол № 7 від 17.06.2024 р.

Rivne – 2024

Compilers:

Boichura M. V., PhD, Associate Professor of the Department of Computer Engineering;
Shatnyi S. V., PhD, Associate Professor of the Department of Computer Engineering;
Shatna A. V., Senior Lecturer of the Department of Computer Engineering.

Укладачі:

Бойчура М. В., к.т.н., доцент кафедри обчислювальної техніки;
Шатний С. В., к.т.н., доцент кафедри обчислювальної техніки;
Шатна А. В., старший викладач кафедри обчислювальної техніки.

Responsible for the publication: Sydor A. I., PhD, Associate Professor, Acting Head of the Department of Computer Engineering.

Відповідальний за випуск: Сидор А. І., к.т.н., доцент, в.о. завідувача кафедри обчислювальної техніки.

Head of the support group
in the field of study
123 «Computer Engineering»                          Sydor A. I.

Керівник групи забезпечення
спеціальності
123 «Комп'ютерна інженерія»                        Сидор А. І.

# **CONTENT**

# Introduction

In the field of computer engineering, it is crucial to be able to build applications that can operate with maximum capacity of computer hardware. This can be achieved through the ability to operate with parallel threads, which allows for the use of multiple cores of the central processor during the program's operation. Unfortunately, many students learn how to compose program code without taking into consideration the possibility of executing it on parallel threads.

However, the advantages of parallel computing go beyond just using the central processor's multiple cores. Modern computers possess powerful computing devices like graphics cards, which can often execute code hundreds of thousands of times faster than on dozens of powerful CPU cores. By acquiring knowledge and skills in parallel computing, one can build effective programs that extract the maximum potential of the hardware.

The discipline «Parallel and Distributed Computing» is an optional component of the educational program 123 "Computer Engineering" at the bachelor's level of education. The subject provides basic training for a specialist in parallel programming and distributed computing, setting the trajectory for further application of the acquired skills in other aspects of the specialty defined in the Bachelor of Computer Engineering course.

Students studying this subject will acquire theoretical knowledge and basic practical skills in creating parallel programs using common specialized libraries. In addition to the three laboratory works, it is possible to complete other tasks or use other high-performance computing technologies on CPUs/graphics cards or arbitrary programming languages with the consent of the lecturer.

Presenting reports is not required, however, to obtain at least a minimal mark, students should be able to explain each

line of code. Only such an approach can help you understand when a given application requires parallelization and what technology should be used.

Finally, completing the course will enable graduates to design and develop parallel and distributed programs. This will prove valuable when planning and executing scientific research, creating games, cryptocurrency mining, or implementing real-time information systems.

As a result of completing all laboratory works of the first module, students should:

know:

− advantages and disadvantages of using parallel computing;

− cases when it is advisable to replace non-parallel code with parallel code;

− basic syntax of CUDA, OpenCL, MPI, OpenMP technologies.

to be able to:

− install libraries for working with CUDA, OpenCL, MPI, OpenMP on various platforms, update drivers, and customize development environments;

− modify non-parallel code into parallel code;

− operate with video card memory and RAM;

− develop parallel programs in teams;

− analyze the characteristics of computing devices according to the criteria of the possibility and feasibility of parallelization;

− successfully combine several high-performance computing technologies;

− programing in accordance with the principles of Coding Conventions.

# Вступ

У галузі комп'ютерної інженерії дуже важливо вміти створювати програми, які здатні використовувати максимальні потужності комп'ютерного обладнання. Цього можна досягти завдяки використанню паралельних потоків, що дозволяє задіювати декілька ядер центрального процесора під час роботи програми. На жаль, зазвичай студентів навчають писати програмний код без врахування можливості його виконання на паралельних потоках.

Однак переваги паралельних обчислень виходять за рамки простого використання декількох ядер центрального процесора. Сучасні комп'ютери мають потужні обчислювальні пристрої, такі як відеокарти, які часто можуть виконувати код у сотні тисяч разів швидше, ніж на десятках потужних ядер центрального процесора. Здобувши знання та навички паралельних обчислень, студенти зможуть створювати ефективний програмний код, який враховує максимальний потенціал апаратного забезпечення.

Освітня компонента «Паралельні та розподілені обчислення» є вибірковою складовою освітньо-професійної програми 123 «Комп'ютерна інженерія» бакалаврського рівня вищої освіти. Навчальна дисципліна забезпечує базову підготовку фахівця з паралельного програмування та розподілених обчислень, задаючи траєкторію для подальшого застосування набутих навичок в інших аспектах спеціальності, визначених у курсі підготовки бакалавра з комп'ютерної інженерії.

Студенти, які вивчають цю дисципліну, отримають теоретичні знання та базові практичні навички створення паралельних програм з використанням поширених спеціалізованих бібліотек. В межах лабораторних робіт, за згодою викладача, можна виконувати й інші, окрім

запропонованих викладачем, завдання або використовувати інші мови програмування чи технології високопродуктивних обчислень на процесорах/відеокартах.

Представлення звітів не вимагається, однак для отримання хоча б мінімальної оцінки з навчальної дисципліни, студенти повинні вміти пояснити кожен рядок коду. Лише таких підхід сприятиме розумінню того, в яких випадках конкретна програма вимагає розпаралелювання обчислень і яку технологію варто використати.

Виконання усіх наведених завдань дозволить випускникам проектувати та розробляти паралельні і розподілені програми. Це стане в нагоді при плануванні та виконанні наукових досліджень, створенні ігор, майнінгу криптовалют чи впровадженні інформаційних систем реального часу.

У результаті виконання всіх лабораторних робіт першого модуля студенти повинні знати переваги та недоліки використання паралельних обчислень, випадки, коли доцільно замінити звичайний код на паралельний, базовий синтаксис технологій CUDA, OpenCL, MPI, OpenMP; вміти: встановлювати бібліотеки для роботи з CUDA, OpenCL, MPI, OpenMP на різних платформах, оновлювати драйвери та налаштовувати середовища розробки, модифікувати звичайний код у паралельний, працювати з пам'яттю відеокарти та оперативною пам'яттю, командно розробляти програми з реалізацією розпаралелення, аналізувати характеристики обчислювальних пристроїв за критеріями можливості та доцільності розпаралелювання, вдало поєднують декілька високопродуктивних обчислювальних технологій, дотримуватись принципів Coding Conventions.

# Laboratory work 1
## High-performance computing using video cards

### Goals

1. Learn how to investigate hardware to check its compatibility with various HPC technologies on graphics cards.

2. Learn to code relatively simple parallel programs using high-performance computing technologies on video cards.

## 1.1. Introduction to OpenCL and CUDA technologies

CUDA (Compute Unified Device Architecture) is a software and hardware architecture for parallel computing that enhances computing performance through the use of NVIDIA graphics processors. The CUDA SDK provides the possibility to include calls to subprograms executed on NVIDIA GPUs to the code of C programs. This is achieved through commands written in a special dialect of C. The CUDA architecture allows developers to organize access to the instruction set of the graphics accelerator and manage its memory at the programmer's discretion.

OpenCL (Open Computing Language) is a framework that allows for parallelization of computing on GPUs and CPUs, as well as on FPGAs. The framework includes a programming language based on the C99 standard and an application programming interface (API). OpenCL provides instruction-level and data-level parallelism, complementing the capabilities of the GPU.

The goal of OpenCL is to complement OpenGL and OpenAL, which are industry standards for three-dimensional computer graphics and sound, by leveraging the capabilities of the GPU. It is a fully open standard available under free licenses

and is developed and supported by the Khronos Group, a consortium of large companies.

For viewing the video characteristics of your computer, you can use the free GPU Caps Viewer software, available for download at the following link: https://www.geeks3d.com/dl/show/710.

## 1.2. Provide OpenCL programming capabilities in Microsoft Visual Studio with an existing NVIDIA graphics card

### Installing OpenCL

To ensure your driver supports OpenCL technology, you can utilize the GPU Caps Viewer program (Fig. 1.1) or inspect the driver settings (Fig. 1.2). If OpenCL support is not present, consider updating to a more modern driver version. If the latest driver does not support OpenCL, you might need to upgrade to a newer graphics card or execute tasks on a different computer or in the cloud.

If your existing driver version does support OpenCL, the next step is to install the CUDA-Toolkit from the official NVIDIA website. However, it is essential to first determine which versions of the CUDA-Toolkit are compatible with your driver. To do this, follow these steps:
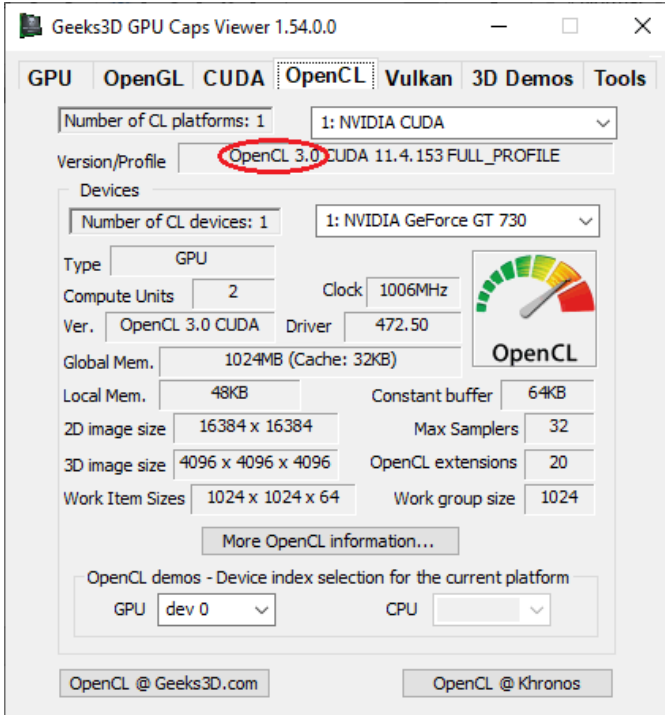
Fig. 1.1. An example of the GPU Caps Viewer program's tray in case of OpenCL technology support (the presence of the OpenCL version indicates support)

1. Check the version of your graphics card driver using a resource like: https://www.nvidia.com/download/index.aspx.

2. Visit https://docs.nvidia.com/cuda/cuda-toolkit-release-notes/index.html and identify the suitable version of the CUDA-Toolkit (see Table 2 or Table 3 on the web page).

3. Proceed to install the CUDA-Toolkit from https://developer.nvidia.com/cuda-toolkit-archive.

4. After installation, ensure to restart your computer to save the new system paths settings.
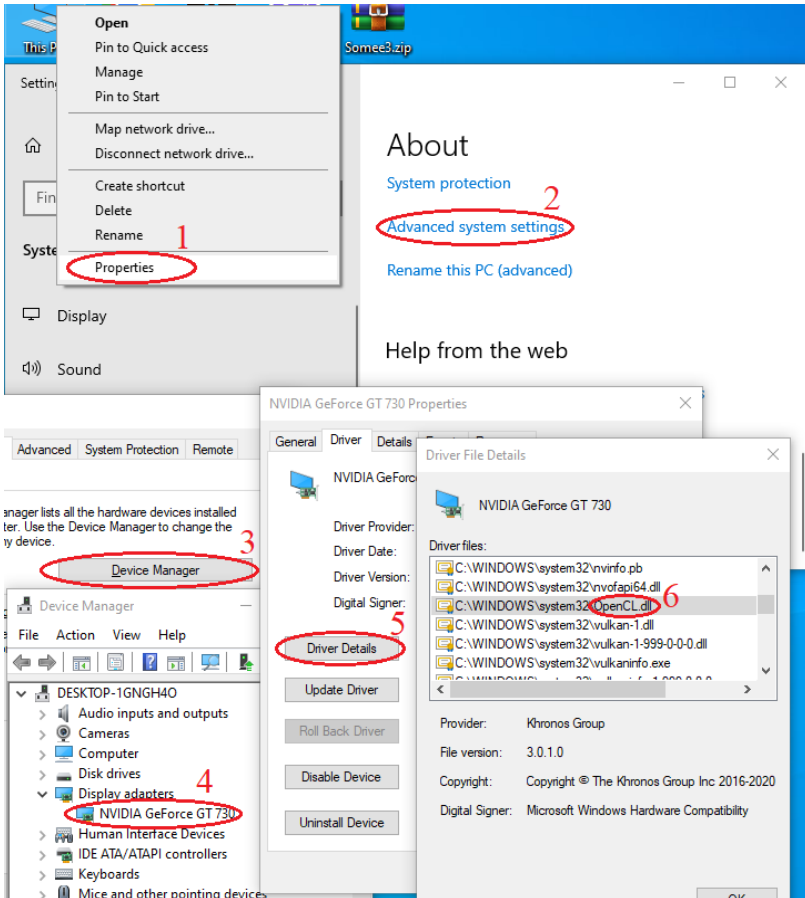
Fig. 1.2. A typical sequence of steps to check OpenCL technology support (support is indicated by the presence of the OpenCL.dll file)

In scenarios where older driver versions necessitate additional component installations apart from the CUDA-Toolkit, it is advisable to seek guidance from your lecturer.

# Creating your first OpenCL project

It is convenient to develop parallel programs using Microsoft Visual Studio. In order to create and run an OpenCL project for NVIDIA graphics cards under a 64-bit Windows operating system, you only need to create an empty C++ console project. To create your first OpenCL project, follow these steps to ensure smooth development:

1. Create an empty C++ console project to get started.
2. Include the path of header files (e.g., *C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.5\include*) in the properties of the project at *C/C++ → General → Additional Include Directories* (Fig. 1.3). This directory is created upon installing the CUDA-Toolkit.
3. Add a required path like *C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.5\lib\x64* to the location *Linker → General → Additional Library Directories* (Fig. 1.4).
4. Specify the file name *OpenCL.lib* in the location *Linker → Input → Additional Dependencies* to complete the setup (refer to Fig. 1.5 for illustrations).
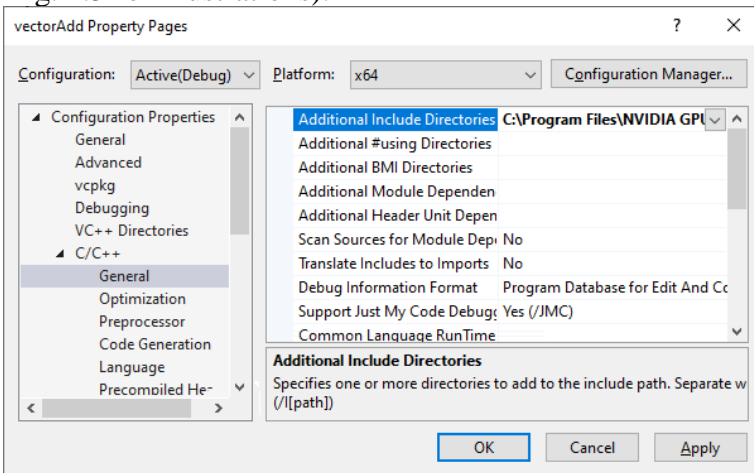


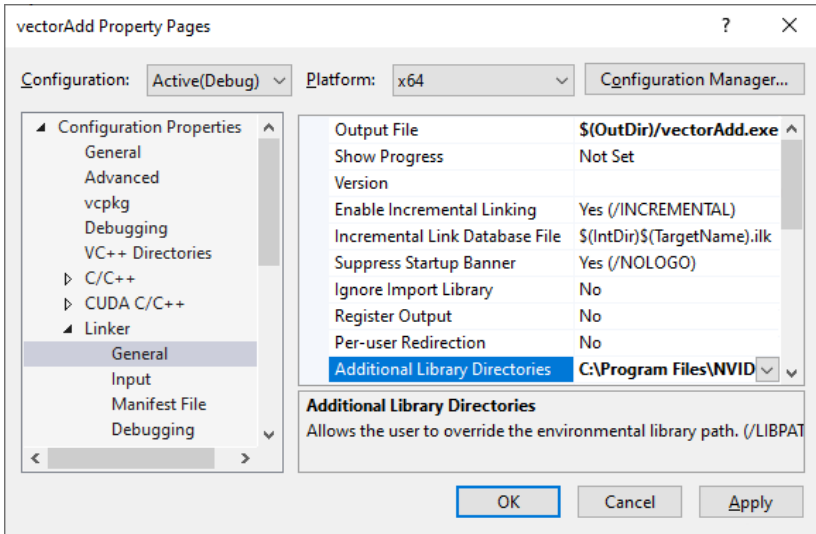Fig. 1.3. Example of adding the path of the OpenCL header files location

Fig. 1.4. An example of adding the path to the location of OpenCL library files
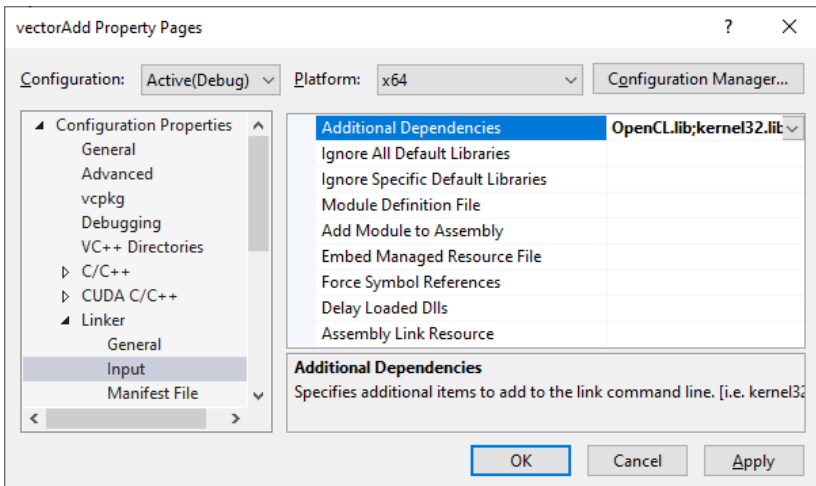


Fig. 1.5. An example of adding the OpenCL.lib library as an additional dependency

13

There are only 2 files we just need to create:

1. Kernel file: vector_add_kernel.cl with code:

```
1.  __kernel void vector_add(__global const int* A,
        __global const int* B, __global int* C)
1.  {
2.    //Get the index of the current item for further
        processing
3.    int i = get_global_id(0);
4.    // Adding vectors
5.    C[i] = A[i] + B[i];
6.  }
```

2. A file of the form *.cpp:

```
1.  #define _CRT_SECURE_NO_WARNINGS
2.  #include <stdio.h>
3.  #include <stdlib.h>
4.
5.  #pragma comment(lib, "OpenCL.lib");
6.
7.  #ifdef __APPLE__
8.    #include <OpenCL/opencl.h>
9.  #else
10.   #include <CL/cl.h>
11. #endif
12. #define MAX_SOURCE_SIZE (0x100000)
13.
14. int main(void)
15. {
16.   // creating two vectors
17.   int i;
18.   const int LIST_SIZE = 1024;
19.   int* A = (int*)malloc(sizeof(int) * LIST_SIZE);
20.   int* B = (int*)malloc(sizeof(int) * LIST_SIZE);
21.   for (i = 0; i < LIST_SIZE; i++)
22.   {
23.     A[i] = i;
24.     B[i] = LIST_SIZE - i;
25.   }
26.   // reading kernel source code from
        vector_add_kernel.cl
27.   FILE* fp;
```

```c
28.    char* source_str;
29.    size_t source_size;
30.    fp = fopen("vector_add_kernel.cl", "r");
31.    if (!fp)
32.    {
33.      fprintf(stderr, "Failed to load kernel.\n");
34.      exit(1);
35.    }
36.    source_str = (char*)malloc(MAX_SOURCE_SIZE);
37.    source_size = fread(source_str, 1,
          MAX_SOURCE_SIZE, fp);
38.    fclose(fp);
39.    // Get information about platforms and devices
40.    cl_platform_id platform_id = NULL;
41.    cl_device_id device_id = NULL;
42.    cl_uint ret_num_devices;
43.    cl_uint ret_num_platforms;
44.    cl_int ret = clGetPlatformIDs(1, &platform_id,
          &ret_num_platforms);
45.    ret = clGetDeviceIDs(platform_id,
          CL_DEVICE_TYPE_GPU, 1, &device_id,
          &ret_num_devices);
46.    // Creating an OpenCL context
47.    cl_context context = clCreateContext(NULL, 1,
          &device_id, NULL, NULL, &ret);
48.    // Creating a command queue
49.    cl_command_queue command_queue =
          clCreateCommandQueue(context, device_id, 0,
          &ret);
50.    // Creating memory buffers on the device for each
          vector
51.    cl_mem a_mem_obj = clCreateBuffer(context,
52.      CL_MEM_READ_ONLY, LIST_SIZE * sizeof(int),
          NULL, &ret);
53.    cl_mem b_mem_obj = clCreateBuffer(context,
54.      CL_MEM_READ_ONLY, LIST_SIZE * sizeof(int),
          NULL, &ret);
55.    cl_mem c_mem_obj = clCreateBuffer(context,
56.      CL_MEM_WRITE_ONLY, LIST_SIZE * sizeof(int),
          NULL, &ret);
57.    // Copying vectors to memory buffers
58.    ret = clEnqueueWriteBuffer(command_queue,
```

```
59.    a_mem_obj, CL_TRUE, 0, LIST_SIZE * sizeof(int),
       A, 0, NULL, NULL);
60.  ret = clEnqueueWriteBuffer(command_queue,
61.    b_mem_obj, CL_TRUE, 0, LIST_SIZE * sizeof(int),
       B, 0, NULL, NULL);
62.  // Creating a program from the kernel source code
63.  cl_program program =
       clCreateProgramWithSource(context, 1,
64.    (const char**)&source_str, (const
       size_t*)&source_size, &ret);
65.  // Creating an executable file
66.  ret = clBuildProgram(program, 1, &device_id,
       NULL, NULL, NULL);
67.  // Creating an OpenCL kernel
68.  cl_kernel kernel = clCreateKernel(program,
       "vector_add", &ret);
69.  // Setting kernel arguments
70.  ret = clSetKernelArg(kernel, 0, sizeof(cl_mem),
       (void*)&a_mem_obj);
71.  ret = clSetKernelArg(kernel, 1, sizeof(cl_mem),
       (void*)&b_mem_obj);
72.  ret = clSetKernelArg(kernel, 2, sizeof(cl_mem),
       (void*)&c_mem_obj);
73.  // Executing the kernel
74.  size_t global_item_size = LIST_SIZE;
75.  size_t local_item_size = 64;
76.  ret = clEnqueueNDRangeKernel(command_queue,
77.    kernel, 1, NULL, &global_item_size,
78.    &local_item_size, 0, NULL, NULL);
79.  // Reading the result from the device to the
       local list C
80.  int* C = (int*)malloc(sizeof(int) * LIST_SIZE);
81.  ret = clEnqueueReadBuffer(command_queue,
82.    c_mem_obj, CL_TRUE, 0, LIST_SIZE * sizeof(int),
       C, 0, NULL, NULL);
83.  ret = clFlush(command_queue);
84.  ret = clFinish(command_queue);
85.  ret = clReleaseKernel(kernel);
86.  ret = clReleaseProgram(program);
87.  ret = clReleaseMemObject(a_mem_obj);
88.  ret = clReleaseMemObject(b_mem_obj);
89.  ret = clReleaseMemObject(c_mem_obj);
```

```
90.    ret = clReleaseCommandQueue(command_queue);
91.    ret = clReleaseContext(context);
92.    free(A);
93.    free(B);
94.    free(C);
95.    return 0;
96. }
```

If we include the print function, the last lines of the program result may look like the one shown in Fig. 1.6.
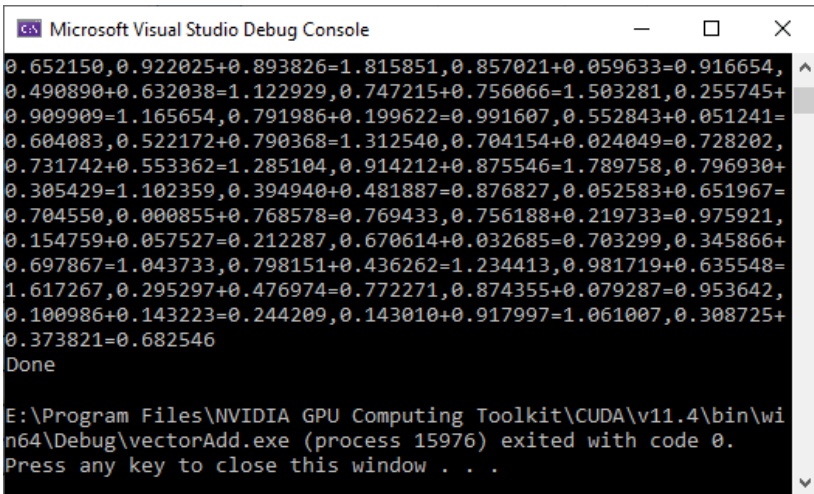


Fig. 1.6. An example of displaying the last lines in the console window

### Additionally

There may be cases when the code works correctly, but there is no syntax highlighting for OpenCL. This problem can be solved as follows:

1. Choose the *Tools → Options* menu item.
2. Find the *Text Editor → File Extension* in the list of options.

17

3. Type *cl* in the *Extension* field; *Microsoft Visual C++* in the *Editor* field.
4. Click *Add*.
5. Click *OK*.
6. Restart the Microsoft Visual Studio environment.

## 1.3. Provide CUDA programming capabilities in Microsoft Visual Studio with an existing NVIDIA graphics card

### Installing CUDA

To verify CUDA technology support, utilize the GPU Caps Viewer program (Fig. 1.7) or inspect the driver settings (Fig. 1.8). If the driver does not support CUDA, consider installing a more recent version or upgrading to a more modern graphics card. Alternatively, you can perform tasks on another computer or in the cloud.
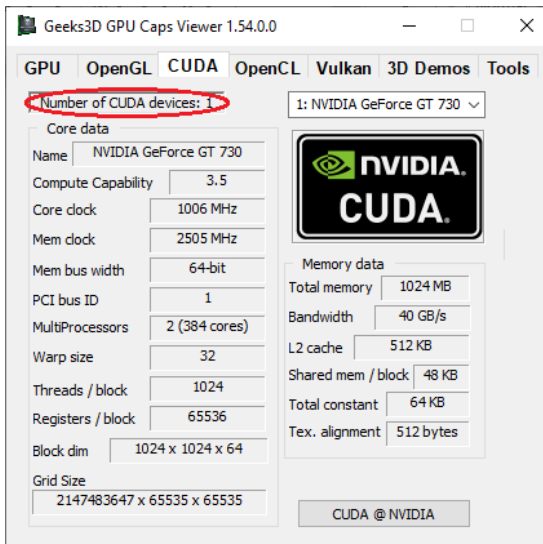


Fig. 1.7. An example of the GPU Caps Viewer program's wreath in the case of CUDA technology support (the number of devices is greater than 0, which indicates CUDA support)
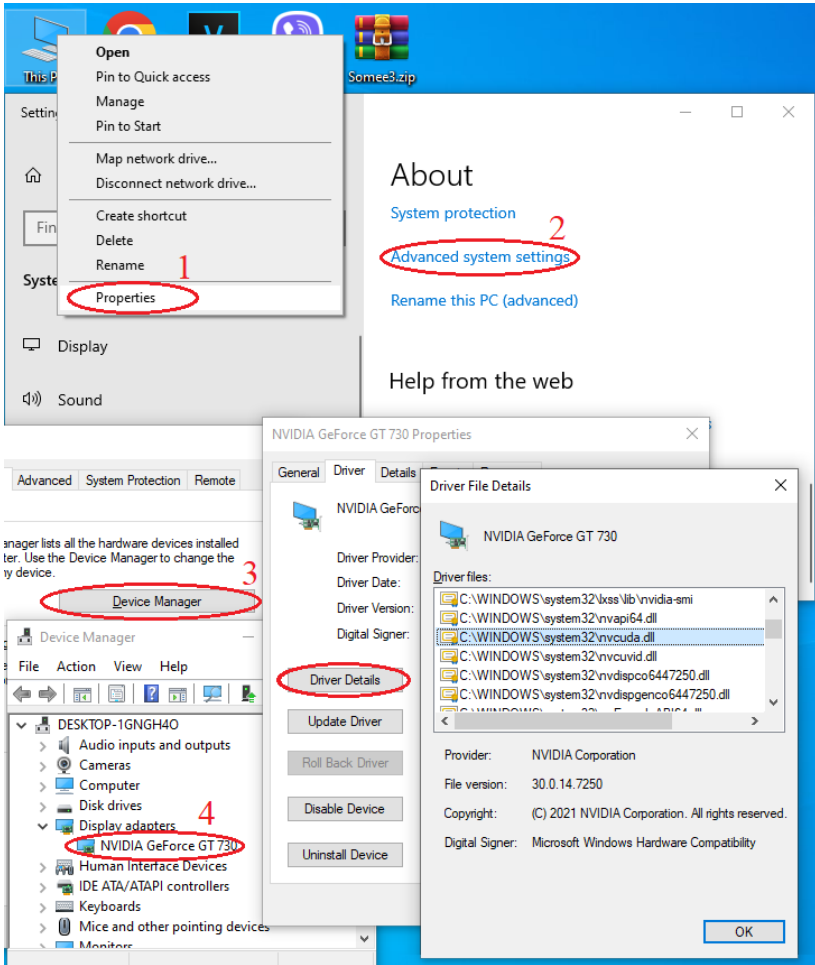
Fig. 1.8. A typical sequence of steps to check for CUDA support (the presence of a text like in the Driver File Details window is a sign of support)

Follow these steps to ensure the possibility of programming in CUDA:

1. Use the following link: https://www.nvidia.com/download/index.aspx, select the current driver version for your graphics card, and install it (see Fig. 1.9).

2. Compare the driver version with Table 2 or 3 at https://docs.nvidia.com/cuda/cuda-toolkit-release-notes/index.html to determine the compatible CUDA Toolkit version.

3. Install the identified driver and CUDA-Toolkit from https://developer.nvidia.com/cuda-toolkit-archive.

4. After installation, restart your computer.



Fig. 1.9. A window with the latest version of some NVIDIA driver

**Creating your first CUDA project**

To create your first CUDA project on a 64-bit Windows operating system in Microsoft Visual Studio, follow these steps:

1. Start Microsoft Visual Studio and search among the project templates for one containing the word "CUDA". Create a project using this template.

2. If the CUDA template is not available, create an empty C++ console project instead.

3. Access the properties of the created project and add the header files' address located at C:\Programme Files\NVIDIA GPU Computing Toolkit\CUDA\v11.5\include\ to the location *C/C++ → General → Additional Include Directories* as shown in Fig. 1.10. This directory is created during CUDA-Toolkit installation.

4. Next, add a path such as C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.5\lib\x64 to the location *Linker → General → Additional Library Directories* as shown in Fig. 1.11.

5. Finally, locate the *Linker → Input → Additional Dependencies* location and add a file similar to *CUDA.lib*, as displayed in Fig. 1.12 to complete the setup.

Fig. 1.10. Example of adding the address of the location of CUDA header files



Fig. 1.11. Example of adding a path to the location of CUDA library files

Fig. 1.12. Example of adding the CUDA.lib library as an additional dependency

It remains to create 1 file:

```
1.  __global__ void vectorAdd(const float *A, const
        float *B, float *C, int numElements)
2.  {
3.    int i = blockDim.x * blockIdx.x + threadIdx.x;
4.    if (i < numElements)
5.    {
6.      C[i] = A[i] + B[i];
7.    }
8.  }
9.  int main()
10. {
11.   cudaError_t err = cudaSuccess;
12.   int numElements = 50000000;
13.   size_t size = numElements * sizeof(float);
14.   printf("[Vector addition of %d elements]\n",
        numElements);
15.
16.   float *h_A = (float *)malloc(size);
17.   float *h_B = (float *)malloc(size);
```

```
18.    float *h_C = (float *)malloc(size);
19.
20.    if (h_A == NULL || h_B == NULL || h_C == NULL)
21.      exit(EXIT_FAILURE);
22.    for (int i = 0; i < numElements; ++i)
23.    {
24.      h_A[i] = rand()/(float)RAND_MAX;
25.      h_B[i] = rand()/(float)RAND_MAX;
26.    }
27.
28.    float *d_A = NULL;
29.    float *d_B = NULL;
30.    float *d_C = NULL;
31.
32.    err = cudaMalloc((void **)&d_A, size);
33.    if (err != cudaSuccess)
34.      exit(EXIT_FAILURE);
35.    err = cudaMalloc((void **)&d_B, size);
36.    if (err != cudaSuccess)
37.      exit(EXIT_FAILURE);
38.    err = cudaMalloc((void **)&d_C, size);
39.    if (err != cudaSuccess)
40.      exit(EXIT_FAILURE);
41.
42.    printf("Copy input data from the host memory to
          the CUDA device\n");
43.    err = cudaMemcpy(d_A, h_A, size,
          cudaMemcpyHostToDevice);
44.    if (err != cudaSuccess)
45.      exit(EXIT_FAILURE);
46.    err = cudaMemcpy(d_B, h_B, size,
          cudaMemcpyHostToDevice);
47.    if (err != cudaSuccess)
48.      exit(EXIT_FAILURE);
49.
50.    int threadsPerBlock = 256;
51.    int blocksPerGrid = (numElements +
          threadsPerBlock - 1) / threadsPerBlock;
52.    printf("CUDA kernel launch with %d blocks of %d
          threads\n", blocksPerGrid, threadsPerBlock);
53.    vectorAdd<<<blocksPerGrid,
          threadsPerBlock>>>(d_A, d_B, d_C, numElements);
```

```
54.    err = cudaGetLastError();
55.    if (err != cudaSuccess)
56.      exit(EXIT_FAILURE);
57.
58.    printf("Copy output data from the CUDA device to
            the host memory\n");
59.    err = cudaMemcpy(h_C, d_C, size,
            cudaMemcpyDeviceToHost);
60.    if (err != cudaSuccess)
61.      exit(EXIT_FAILURE);
62.
63.    err = cudaFree(d_A);
64.    if (err != cudaSuccess)
65.      exit(EXIT_FAILURE);
66.    err = cudaFree(d_B);
67.    if (err != cudaSuccess)
68.      exit(EXIT_FAILURE);
69.    err = cudaFree(d_C);
70.    if (err != cudaSuccess)
71.      exit(EXIT_FAILURE);
72.
73.    free(h_A);
74.    free(h_B);
75.    free(h_C);
76.
77.    printf("Done\n");
78.    return 0;
79. }
```

If you also add output functions, the last lines of the program result may look like the one shown in Fig. 1.6.

### **Additionally**

There may be cases when the code works correctly, but there is no syntax highlighting for CUDA. This problem can be solved as follows:

1. Select the *Tools → Options* menu item.

2. Search for *Text Editor → File Extension* in the list of options

3. Type *cu* in the *Extension* field; select *Microsoft Visual C++* in the *Editor* field.

4. Click Add.

5. Click OK.

6. Restart the *Microsoft Visual Studio* environment.

## 1.4. Provide OpenCL programming capabilities in Microsoft Visual Studio with an existing AMD graphics card

### Installing OpenCL

To determine if your driver supports OpenCL technology, check with the GPU Caps Viewer program (Fig. 1.1) or refer to the driver settings. An indication of support is the existence of a file like OpenCL.dll (Fig. 1.2). If this specific file is not present, it is recommended to update to a more recent version of the driver. In cases where the latest driver version does not offer OpenCL support, consider upgrading to a more modern graphics card or completing tasks on a different computer or in the cloud.

If the current driver version does support OpenCL, the next step involves downloading the SDK through one of the two methods provided:

1. Utilizing vcpkg (this's preferred method).

1.1. Install Git for Windows from https://gitforwindows.org/ and then restart your computer.

1.2. Use a command prompt in administrator mode to execute the following steps sequentially:

- cd c:\Program Files
- git clone https://github.com/Microsoft/vcpkg.git
- .\vcpkg\bootstrap-vcpkg.bat
- .\vcpkg install opencl:x64-windows
- .\vcpkg integrate install

1.3. Restart your computer to apply the changes.

2. Access the open source project on GitHub by visiting: https://github.com/GPUOpen-LibrariesAndSDKs/OCL-SDK/releases.

**Creating your first OpenCL project**

If you want to create and execute your first OpenCL project for AMD graphics card on a 64-bit Windows operating system, it is beneficial to use Microsoft Visual Studio. Utilize an empty C++ console project as the basis for development. Should any obstacles arise when attempting to access libraries, follow the listed steps:

1. Create an empty C++ console project.
2. In the project's properties, go to *C/C++ → General → Additional Include Directories* and add the path of header files, such as C:\Users\Roma\Downloads\lightOCLSDK\include, from the unzipped SDK downloaded on GitHub.
3. Add a new path, such as *C:\Users\Roma\Downloads\lightOCLSDK\lib\x86_64*, to the project's properties location *Linker → General → Additional Library Directories*.
4. Add a filename like *OpenCL.lib* to the *Linker → Input → Additional Dependencies* location in the project's properties.

It remains to create 2 files:

1. The kernel file: vector_add_kernel.cl with the code

```
1. __kernel void vector_add(__global const int* A,
       __global const int* B, __global int* C)
2. {
3.   //Get the index of the current item to process
4.   int i = get_global_id(0);
5.   // Adding vectors
6.   C[i] = A[i] + B[i];
7. }
```

2. A file of the form *.cpp:

```
1. #define _CRT_SECURE_NO_WARNINGS
2. #include <stdio.h>
3. #include <stdlib.h>
```

```
4.
5.  #pragma comment(lib, "OpenCL.lib");
6.
7.  #ifdef __APPLE__
8.    #include <OpenCL/opencl.h>
9.  #else
10.   #include <CL/cl.h>
11. #endif
12. #define MAX_SOURCE_SIZE (0x100000)
13.
14. int main(void)
15. {
16.   // creating two vectors
17.   int i;
18.   const int LIST_SIZE = 100000000 / 2;
19.   int* A = (int*)malloc(sizeof(int) *LIST_SIZE  );
20.   int* B = (int*)malloc(sizeof(int) *LIST_SIZE  );
21.   for (i = 0; i < LIST_SIZE; i++)
22.   {
23.     A[i] = i;
24.     B[i] = LIST_SIZE - i;
25.   }
26.   // reading kernel source code from
          vector_add_kernel.cl
27.   FILE* fp;
28.   char* source_str;
29.   size_t source_size;
30.   fp = fopen("vector_add_kernel.cl", "r");
31.   if (!fp)
32.   {
33.     fprintf(stderr, "Failed to load kernel.\n");
34.     exit(1);
35.   }
36.   source_str = (char*)malloc(MAX_SOURCE_SIZE);
37.   source_size = fread(source_str, 1,
          MAX_SOURCE_SIZE, fp);
38.   fclose(fp);
39.   // Get information about platforms and devices
40.   cl_platform_id platform_id = NULL;
41.   cl_device_id device_id = NULL;
42.   cl_uint ret_num_devices;
43.   cl_uint ret_num_platforms;
```

```
44.  cl_int ret = clGetPlatformIDs(1, &platform_id,
        &ret_num_platforms);
45.  ret = clGetDeviceIDs(platform_id,
        CL_DEVICE_TYPE_GPU, 1, &device_id,
        &ret_num_devices);
46.  // Creating an OpenCL context
47.  cl_context context = clCreateContext(NULL, 1,
        &device_id, NULL, NULL, &ret);
48.  // Creating a command queue
49.  cl_command_queue command_queue =
        clCreateCommandQueue(context, device_id, 0,
        &ret);
50.  // Creating memory buffers on the device for each
        vector
51.  cl_mem a_mem_obj = clCreateBuffer(context,
52.   CL_MEM_READ_ONLY, LIST_SIZE * sizeof(int),
        NULL, &ret);
53.  cl_mem b_mem_obj = clCreateBuffer(context,
54.   CL_MEM_READ_ONLY, LIST_SIZE * sizeof(int),
        NULL, &ret);
55.  cl_mem c_mem_obj = clCreateBuffer(context,
56.   CL_MEM_WRITE_ONLY, LIST_SIZE * sizeof(int),
        NULL, &ret);
57.  // Copying vectors to memory buffers
58.  ret = clEnqueueWriteBuffer(command_queue,
59.   a_mem_obj, CL_TRUE, 0, LIST_SIZE * sizeof(int),
        A, 0, NULL, NULL);
60.  ret = clEnqueueWriteBuffer(command_queue,
61.   b_mem_obj, CL_TRUE, 0, LIST_SIZE * sizeof(int),
        B, 0, NULL, NULL);
62.  // Creating a program from the kernel source code
63.  cl_program program =
        clCreateProgramWithSource(context, 1,
64.  (const char**)&source_str, (const
        size_t*)&source_size, &ret);
65.  // Creating an executable file
66.  ret = clBuildProgram(program, 1, &device_id,
        NULL, NULL, NULL);
67.  // Creating an OpenCL kernel
68.  cl_kernel kernel = clCreateKernel(program,
        "vector_add", &ret);
69.  // Setting kernel arguments
```
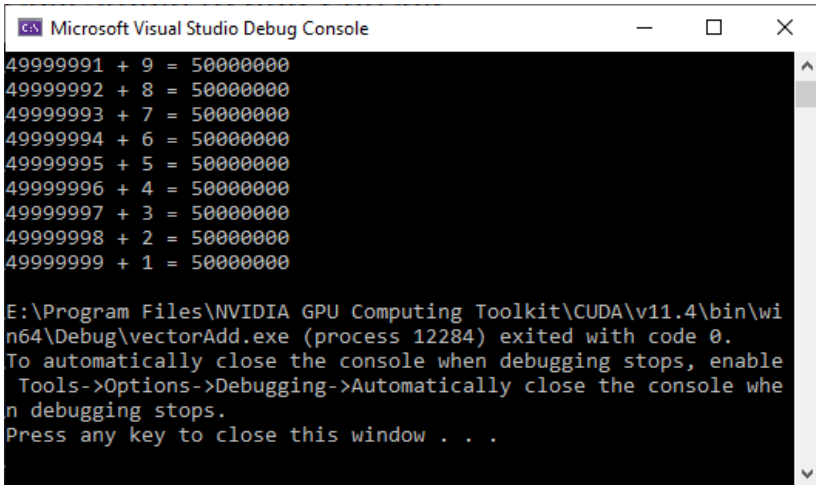
```
70.   ret = clSetKernelArg(kernel, 0, sizeof(cl_mem),
        (void*)&a_mem_obj);
71.   ret = clSetKernelArg(kernel, 1, sizeof(cl_mem),
        (void*)&b_mem_obj);
72.   ret = clSetKernelArg(kernel, 2, sizeof(cl_mem),
        (void*)&c_mem_obj);
73.   // Executing the kernel
74.   size_t global_item_size = LIST_SIZE;
75.   size_t local_item_size = 64;
76.   ret = clEnqueueNDRangeKernel(command_queue,
77.     kernel, 1, NULL, &global_item_size,
78.     &local_item_size, 0, NULL, NULL);
79.   // Reading the result from the device to the
        local list C
80.   int* C = (int*)malloc(sizeof(int) * LIST_SIZE);
81.   ret = clEnqueueReadBuffer(command_queue,
82.     c_mem_obj, CL_TRUE, 0, LIST_SIZE * sizeof(int),
        C, 0, NULL, NULL);
83.   ret = clFlush(command_queue);
84.   ret = clFinish(command_queue);
85.   ret = clReleaseKernel(kernel);
86.   ret = clReleaseProgram(program);
87.   ret = clReleaseMemObject(a_mem_obj);
88.   ret = clReleaseMemObject(b_mem_obj);
89.   ret = clReleaseMemObject(c_mem_obj);
90.   ret = clReleaseCommandQueue(command_queue);
91.   ret = clReleaseContext(context);
92.   free(A);
93.   free(B);
94.   free(C);
95.   return 0;
96. }
```

If you also add output functions, the last lines of the program result can look like the one shown in Fig. 1.13.

Fig. 1.13. Example of the last lines of the output window that
displays the result of adding two vectors

**Additionally**

There may be cases when the code works correctly, but
there is no syntax highlighting for OpenCL. This problem can
be solved as follows:

1. Select the Tools → Options menu item.
2. In the list of options, find Text Editor → File Extension.
3. Enter *cl* in the Extension field; select *Microsoft Visual
C++* in the Editor field.
4. Click Add.
5. Click OK.
6. Restart the Microsoft Visual Studio environment.

## 1.5. Statements

1. Configure the working capability of OpenCL or CUDA on
your computer or laptop (*20% of points*).

2. Create two separate kernel functions, each dedicated to its unique operation according to the chosen option. If OpenCL technology is being used, it is recommended to create two corresponding files with the extension *.cl (*10% of points*).

3. Determine the number of array elements to check the maximum capabilities of your computing device (*10% of points*).

4. Implement the execution of both cores separately and in parallel according to the chosen option (*50% of points*).

5. Programmatically compare the amount of time required to perform calculations in parallel versus sequentially (*10% of points*).

6. Analyze and explain the calculated results.

## Options for accomplishing the task

The following options are available for selecting the operation to be performed:

Option 1.   Addition and subtraction.
Option 2.   Addition and multiplication.
Option 3.   Addition and division.
Option 4.   Addition and floor division.
Option 5.   Addition and division with remainder.
Option 6.   Subtraction and multiplication.
Option 7.   Subtraction and division.
Option 8.   Subtraction and floor division.
Option 9.   Subtraction and division with remainder.
Option 10. Multiplication and division.
Option 11. Multiplication and floor division.
Option 12. Multiplication and division with remainder.
Option 13. Division and floor division.
Option 14. Division and division with remainder.
Option 15. Floor division and division with remainder.

**Additional materials**

1. An introductory video on CUDA basics and some examples in Python: https://www.youtube.com/watch?v=r9IqwpMR9TE&list=PLq XS1b2lRpYTUHPp2MYkgXS7v6qA-JsF.

2. Lessons about CUDA: https://www.youtube.com/watch?v=m0nhePeHwFs&list=PLK K11Ligqititws0ZOoGk3SW-TZCar4dK&index=1

3. An introductory video about OpenCL: https://www.youtube.com/watch?v=V4RfPfHQPC8&list=PLi wt1iVUib9s6vyEqdpcgAq7NBRlp9mAY&index=1

4. Lectures about OpenCL: https://www.youtube.com/watch?v=8D6yhpiQVVI&list=PLDi vN33Lbf6cLqZ5_i_k0KeMaQKEctMgS

**Questions for independent training**

1. What is the `__global__` qualifier used for?

2. What is the `blockdim.x * blockidx.x + threadidx.x` expression for?

3. Why is it important to specify an if statement `if (i < numElements)` in the kernel?

4. What serious limitation has the `new` operator compared to the `malloc` function?

5. Why is it advisable to use the `if` statement directly after using the `malloc` function?

6. What is the difference between the `malloc` and `cudamalloc` functions?

7. How does the `cudamemcpy` function perform?

8. What are the essence of the `threadsperblock` and `blockspergrid` variables?

9. What do angle brackets mean in the expression `func<<<blockspergrid, threadsperblock>>>(d_a, d_b, d_c, numElements)`?

10. What is the difference between the `cudafree` and `free` functions?

11. What is the difference between a device and a host in CUDA terminology?

12. What is the `__kernel` keyword used for?

13. What is the purpose of the `get_global_id(0)` function?

14. What is the difference between the `clGetPlatformIDs` and `clGetDeviceIDs` functions?

15. What is the purpose of the `clCreateContext` function?

16. What is the purpose of the `clCreateCommandQueue` function?

17. How do the `clEnqueueWriteBuffer` and `clEnqueueReadBuffer` functions performs?

18. What is the purpose of the `clBuildProgram` function?

19. What is the purpose of the `clCreateKernel` function?

20. What is the purpose of the `clEnqueueNDRangeKernel` function?

# Laboratory work 2
## High-performance computing on processors

### Goals
1. Learn how to create relatively simple programs using high-performance computing technologies on CPU cores.
2. Learn how to analyze and compare different HPC technologies.

## 2.1. Introduction to MPI technology

Message Passing Interface (MPI) serves as a critical software interface (API) for facilitating information transfer by enabling the communication of messages among processes striving for a shared objective. Essentially, the concept behind MPI involves adapting the technology of message transmission between pairs or groups of subscribers for the realm of parallel computing.

The utilization of MPI stands as the predominant standard for data exchange interfaces in the realm of parallel programming, with implementations available across a plethora of computer platforms. This technology is particularly instrumental in crafting programs designed for clusters and supercomputers, with the primary mode of communication among MPI processes involving the exchange of messages with one another.

## 2.2. Provide MPI programming capabilities in Microsoft Visual Studio with an existing AMD graphics card

### Installing MPI
A guide for Installing MPI on a Windows operating system and creating the first MPI project is as below.

Implementing MPI Technology:

1. Begin by installing Microsoft MPI from https://www.microsoft.com/en-us/download/details.aspx?id=100593. Following the installation, restart the computer.

2. Proceed by installing Git for Windows, which can be obtained from https://gitforwindows.org, and then restart the computer.

3. Using a command prompt that runs in administrator mode, execute the following steps in sequence:

 3.1. Navigate to the directory by typing "cd c:\Program Files".

 3.2. Clone the vcpkg repository by running the command "git clone https://github.com/Microsoft/vcpkg.git".

 3.3. Execute the bootstrap-vcpkg.bat file by typing ".\vcpkg\bootstrap-vcpkg.bat".

 3.4. Install the msmpi package for x64 Windows using ".\vcpkg install msmpi:x64-windows".

 3.5. Integrate the installed by running ".\vcpkg integrate install". After completing these steps, restart the computer.

## Creating the first MPI project

Microsoft Visual Studio provides a convenient environment for developing parallel programs. To set up an MPI project for a 64-bit Windows operating system, simply create an empty C++ console project and input the desired MPI code. Should issues relating to library access arise, you will need to follow these steps:

1. Start by creating an empty C++ console project.

2. In the project's properties, navigate to *C/C++ → General → Additional Include Directories* and include the path of the header files, such as *C:/Users/I-HAVE-NO-TIME-for/vcpkg/packages/msmpi_x64-windows/include/*. This

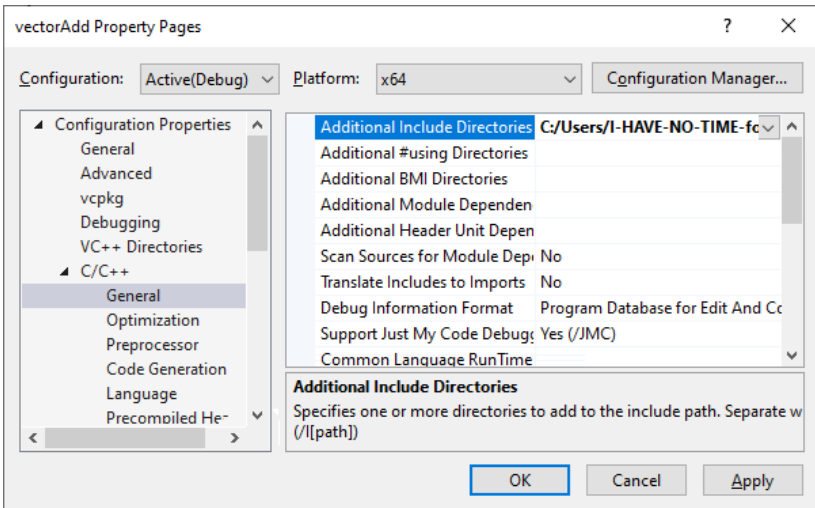directory becomes accessible after installing the msmpi package (refer to Fig. 2.1 for guidance).



Fig. 2.1. Example of insert the address of MPI header files location

We need to place the following code in a *.c or *.cpp file:

```
1.   #include "mpi.h"
2.   #include <stdio.h>
3.   #include <stdlib.h>
4.
5.   #define MASTER 0 //Create a synonym for the rank
         of the main process
6.   #define ARRAY_SIZE 1024 //Set the sizes for
         arrays.
7.
8.   int main(int argc, char* argv[]) //Entry point to
         the parallel program
9.   {
10.
11.    // Create pointers to the addresses where the
         arrays will be located
12.    int* a = NULL;
```

```
13.    int* b = NULL;
14.    int* c = NULL;
15.
16.    int total_proc;   // total number of processes
17.    int rank;         // the rank of each process
18.    int n_per_proc;   // elements allocated to the
          process
19.    int n = ARRAY_SIZE;
20.    int i;
21.
22.    MPI_Status status; // used to check for MPI
          errors.
23.
24.    // 1. Initialize the MPI environment
25.    MPI_Init(&argc, &argv);
26.    MPI_Comm_size(MPI_COMM_WORLD, &total_proc);
27.    // 2. Assign the total number of processes to
          the total_proc variable
28.    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
29.    // 3. Assign the number of the current process
          to the rank variable
30.
31.    // Pointers to a number of short arrays that
          will be used for parallel additions
32.    int* ap;
33.    int* bp;
34.    int* cp;
35.
36.    // 4. The main process initializes the arrays
          that will be used for addition.
37.    if (rank == MASTER)
38.    {
39.      a = (int*)malloc(sizeof(int) * n);
40.      b = (int*)malloc(sizeof(int) * n);
41.      c = (int*)malloc(sizeof(int) * n);
42.
43.      for (i = 0; i < n; i++)
44.      {
45.        a[i] = i + 1;
46.      }
47.      for (i = 0; i < n; i++)
48.      {
```

```
49.        b[i] = i * 2;
50.      }
51.    }
52.
53.    // set the number of pairs of elements that will
          be added in each process
54.    n_per_proc = n / total_proc;
55.
56.    // 5. Memory allocation for shortened arrays
57.    ap = (int*)malloc(sizeof(int) * n_per_proc);
58.    bp = (int*)malloc(sizeof(int) * n_per_proc);
59.    cp = (int*)malloc(sizeof(int) * n_per_proc);
60.
61.    // 6. Uniform distribution of array a into
          shortened arrays ap
62.    MPI_Scatter(a, n_per_proc, MPI_INT, ap,
          n_per_proc, MPI_INT, MASTER, MPI_COMM_WORLD);
63.    // Uniform distribution of array b into
          shortened arrays bp
64.    MPI_Scatter(b, n_per_proc, MPI_INT, bp,
          n_per_proc, MPI_INT, MASTER, MPI_COMM_WORLD);
65.
66.    // 7. Pairwise addition of elements of shortened
          arrays
67.    for (i = 0; i < n_per_proc; i++)
68.    {
69.      cp[i] = ap[i] + bp[i];
70.    }
71.
72.    // 8. Collecting shortened arrays cp into a
          single array c
73.    MPI_Gather(cp, n_per_proc, MPI_INT, c,
          n_per_proc, MPI_INT, MASTER, MPI_COMM_WORLD);
74.
75.    // Outputting results in the main process
76.    if (rank == MASTER)
77.    {
78.      for (i = 0; i < n; i++)
79.      {
80.        printf("%lld + %lld = %lld\n", a[i], b[i],
          c[i]);
81.      }
```
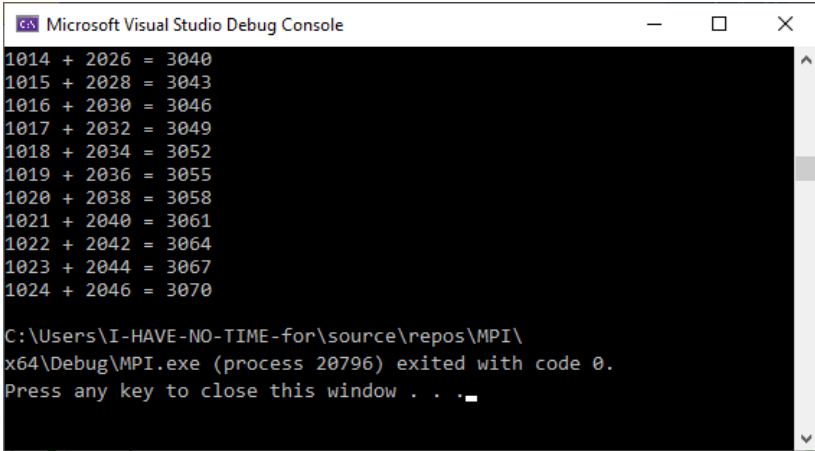
```
82.    }
83.
84.    // Freeing up memory
85.    if (rank == MASTER)
86.    {
87.      free(a);
88.      free(b);
89.      free(c);
90.    }
91.    free(ap);
92.    free(bp);
93.    free(cp);
94.
95.    // 9. Terminating MPI processes and environment
96.    MPI_Finalize();
97.
98.    return 0;
99. }
```

An example of outputting the last lines of the program result is shown in Fig. 2.2.
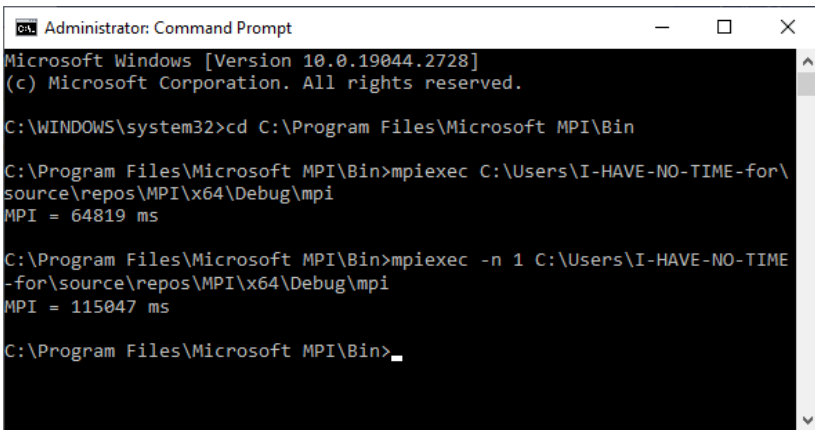
If you are executing the provided code from the Microsoft Visual Studio environment, it is important to note that the computations will be carried out on a single processor core. For actual parallelization, it is recommended to utilize the command line of your operating system, preferably in administrator mode. By making a few adjustments to the source code and assuming the mpiexec file is situated in the folder C:\Program Files\Microsoft MPI\Bin, while the compiled C++ project executable file is located in C:\Users\I-HAVE-NO-TIME-for\source\repos\MPI\x64\Debug\mpi, the command line prompt may resemble the example depicted in Fig. 2.3.

Fig. 2.2. Console window of outputting the result of the
program for adding two vectors



Fig. 2.3. Examples of executing a parallel program on MPI
technology using a different number of threads

**Additionally**

It is recommended to install the Microsoft MPI program in
the default folder. Otherwise, you may have problems with the
installation of the msmpi package.

41

## 2.3. Statements

1. Configure MPI technology on your computer or laptop to enable its functionality (*20% of points*).

2. Adjust the existing code based on the chosen option from the first lab (*50% of points*).

3. Efficiently distribute the workload among processor cores to ensure optimal performance of the operations in accordance the chosen option (*10% of points*).

4. Determine the appropriate number of processes to check the maximum capabilities of your computing device (*10% of points*).

5. Perform a programmatically time-based comparison of the calculation performance across MPI, CUDA(OpenCL), and sequential technologies (*10% of points*).

6. Perform a thorough analysis of the results obtained and provide insightful explanations.

### Additional materials

1. Fundamentals of MPI:
https://www.youtube.com/watch?v=c0C9mQaxsD4

2. More complex cases of MPI usage:
https://www.youtube.com/watch?v=q9OfXis50Rg

### Questions for independent training

1. What is the purpose of the `MPI_Init` function?

2. What is the difference between the `MPI_Comm_size` and `MPI_Comm_rank` functions?

3. What is a process rank?

4. What does the `MPI_Scatter` function do?

5. What does the `MPI_Gather` function do?

6. What is the purpose of the `MPI_Finalize` function?

7. Which rank is usually used as the main process?

8.  What does `MPI_COMM_WORLD` mean?

9.  What types of data are available in MPI?

10. What function is used to free previously allocated memory in MPI technology?

11. What of the well-known things does MPI technology imitate?

12. What physical device is designed for MPI computing?

13. What is the usual limitation when running an MPI program in Microsoft Visual Studio?

14. What command should be entered to calculate on exactly two threads?

15. How many threads/cores are used by default to perform calculations when running an MPI program from the command line?

16. Which of the following functions does not exist in MPI: `MPI_Init, MPI_Send, MPI_Free`?

17. What is `MPI_Status`?

18. What arguments does the `MPI_Init` function receive?

19. What does the `MPI_Recv` function do?

20. What does the `MPI_Reduce` function do?

# Laboratory work 3
## Compatible high-performance computing on graphics cards and central processors

### Goals
1. Learn how to use a fast method of parallelizing computations on processors.
2. Learn how to combine the computing resources of HPC of video cards and processors.
3. Learn ways to analyse different computing technologies.

## 3.1. Introduction to OpenMP technology

OpenMP, which stands for Open Multi-Processing, is an open standard for parallelizing programs between processor cores. It employs preprocessor directives, or "pragmas", library procedures, and environment variables to facilitate the creation of multi-threaded applications intended for use on multiprocessor systems with shared memory.

Key elements of the standard include structures for creating threads (*parallel* directive), work distribution between threads (*DO/for* and *section* directives), data management (use of the *shared* and *private* memory classes to define variable memory), thread synchronization (*critical*, *atomic*, and *barrier* directives), and runtime support library procedures (such as *omp_get_thread_num*) and environment variables (e.g. *OMP_NUM_THREADS*).

## 3.2. Creating the first OpenMP project in Microsoft Visual Studio

Developing parallel programs with Microsoft Visual Studio is a convenient option. To create a relatively simplistic

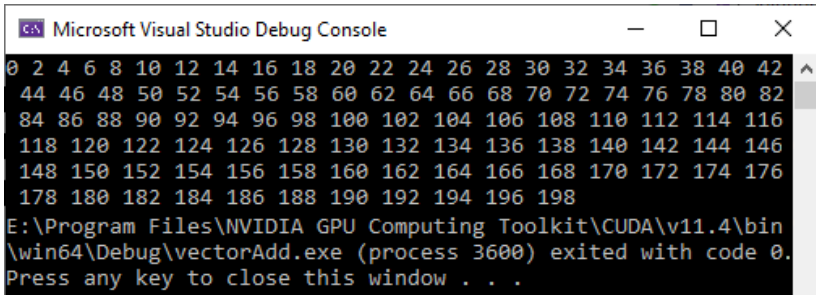program, such as one that adds two vectors with OpenMP parallelization, follow the below-mentioned steps:

1. Create an empty C++ console project.
2. Place the following code in a *.c or *.cpp file:

```cpp
1.  #include <omp.h>
2.  #include <iostream>
3.
4.  int main()
5.  {
6.    const long n = 100;
7.    int* a = (int*)malloc(sizeof(int) * n);
8.    int* b = (int*)malloc(sizeof(int) * n);
9.    int* c = (int*)malloc(sizeof(int) * n);
10.   for (int i = 0; i < n; i++)
11.   {
12.     a[i] = i;
13.     b[i] = i;
14.   }
15.
16.   #pragma omp parallel for
17.   for (int i = 0; i < n; i++)
18.   {
19.     c[i] = a[i] + b[i];
20.   }
21.
22.   for (int i = 0; i < n; i++)
23.   {
24.     printf("%i ", c[i]);
25.   }
26. }
```

An example of the output is shown in Fig. 3.1.

### 3.3.Statements

The tasks involve the integration of CUDA/OpenCL and OpenMP technologies to enhance performance, compare calculation times across multiple technologies, and analyze the obtained results.

Fig. 3.1. An example of executing a program code using OpenMP technology

1. Adjust the existing code based on your chosen option from the first lab (10*% of points*).

2. Integrate the capabilities of CUDA/OpenCL and OpenMP technologies effectively to combine their power and efficiency (*50% of points*).

3. Compare the time performance of computations using MPI, OpenMP, CUDA (OpenCL) and sequential technologies programmatically (*20% of points*).

4. Provide a comprehensive analysis of the results and explain them in detail to demonstrate your understanding (*20% of points*).

### Additional materials

1. Video tutorials about OpenMP:
https://www.youtube.com/watch?v=SLnh7yS52-I&list=PL3xCBlatwrsWhsdHq3JFJiuQ60gHzYtFU&index=1.

2. Fundamentals of OpenMP:
https://www.youtube.com/watch?v=YdHv_2AT4GI&list=PLQsgurGJM5piDxUkHX15V6xrAFHzzWfMh&index=2.

## Questions for independent training

1. What is the purpose of the `omp.h` library?

2. How can the `#pragma omp parallel` fragment affect on program execution?

3. How does the `omp_get_thread_num()` function work?

4. What is the difference between `shared` and `private` variables?

5. What does `#pragma omp barrier` mean?

6. What does `#pragma omp critical { ... }` mean?

7. How does the code controlled by the `#pragma omp for { ... }` directive work?

8. What does the `#pragma omp parallel private`(partial_Sum) shared(total_Sum) mean?

9. What is the complexity of combining several HPC technologies?

10. What might be the compiler's reaction if you use `#pragma omp` without connecting the `omp.h` library?

11. What construction is used to create threads?

12. What structures exist for distributing work between threads?

13. What constructions are used for determining the memory class of variables?

14. What are the ways to synchronize threads?

15. What is *OMP_NUM_THREADS*?

16. What technology is more efficient: MPI or OpenMP?

17. What does the `#pragma omp master` string mean?

18. What steps should be performed to ensure the ability to use OpenMP?

19. What is the difference between the `#pragma omp parallel` and `#pragma omp` expressions in the case of adding elements of two arrays?

20. How does the compiler react on the `#pragma parallel omp parallel` entry?

# Bibliography

## Primary

1. Бойчура М. В., Шатний С. В., Шатна А. В. Методичні вказівки до лабораторних робіт з навчальної дисципліни «Паралельні та розподілені обчислення» (частина 1) для здобувачів вищої освіти першого (бакалаврського) рівня за освітньо-професійною програмою «Комп'ютерна інженерія» спеціальності 123 «Комп'ютерна інженерія» денної і заочної форм навчання [Електронне видання]. Рівне : НУВГП, 2023. 49 с.

2. Пасічник В. В., Лупенко С. А., Луців А. М. Паралельні та розподілені обчислення. Львів : Магнолія 2006, 2024. 664 с.

3. Hwu W.-M.W., Kirk D.B., Hajj I.E. Programming Massively Parallel Processors: A Hands-on Approach. 4th ed. Cambridge : Morgan Kaufmann, 2022. 580 p.

4. Deakin T., Mattson T.G. Programming Your GPU with OpenMP: Performance Portability for GPUs. Cambridge : The MIT Press, 2023. 336 p.

5. Barlas G. Multicore and GPU Programming: An Integrated Approach. 2nd ed. Burlington : Morgan Kaufmann, 2022. 1024 p.

6. Chopp D.L. Introduction to High Performance Scientific Computing. Philadelphia : SIAM - Society for Industrial and Applied Mathematics, 2019. 453 p.

7. Dash N. Ultimate Parallel and Distributed Computing with Julia For Data Science. Delhi : Orange Education Pvt Ltd, 2024. 484 p.

8. Peters T. Parallel Python with Dask: Perform distributed computing, concurrent programming and manage large dataset. Delhi : GitforGits, 2023. 172 p.

9. Nelli F. Parallel and High Performance Programming with Python. Delhi : AVA, 2023. 392 p.

10. Scott L.R., Clark T., Bagheri B. Scientific Parallel Computing. Princeton : Princeton University Press, 2021. 392 p.

**Secondary**

11. Bomba A. Ya., Kuzlo M. T., Michuta O. R., Boichura M.V. On a method of image reconstruction of anisotropic media using applied quasipotential tomographic data. *Mathematical Modeling and Computing*. 2019. Vol. 6 (2). P. 211–219.

12. Vlasyuk A., Zhukovskyy V., Zhukovska N., Shatnyi S. Parallel Computing optimization of Two- Dimensional Mathematical Modeling of Contaminant Migration in Catalytic Porous Media. *ACIT'2020:* Proceedings of 2020 10th International Conference on Advanced Computer Information Technologies (Deggendorf, Sep. 16-18, 2020). Deggendorf, 2020. P. 23–28.

13. Overview - CUDA Python 12.5.0 documentation. URL: https://nvidia.github.io/cuda-python/overview.html (Last accessed: 01.06.2024).

14. CUDA C++ Best Practices Guide. Release 12.5. URL: https://docs.nvidia.com/cuda/pdf/CUDA_C_Best_Practices_Guide.pdf (Last accessed: 01.06.2024).

15. The programming guide to the CUDA model and interface. URL: https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html (Last accessed: 01.06.2024).

16. The OpenCL™ C Specification. URL: https://registry.khronos.org/OpenCL/specs/3.0-unified/pdf/OpenCL_C.pdf (Last accessed: 01.06.2024).

17. MPI: A Message-Passing Interface Standard. Version 4.1. URL: https://www.mpi-forum.org/docs/mpi-4.1/mpi41-report.pdf (Last accessed: 01.06.2024).

18. Using MPI with C - RC University of Colorado Boulder documentation. URL: https://curc.readthedocs.io/en/latest/programming/MPI-C.html (Last accessed: 01.06.2024).

19. Home - OpenMP. URL: http://www.openmp.org (Last accessed: 01.06.2024).

20. OpenMP | LLNL HPC Tutorials. URL: https://hpc-tutorials.llnl.gov/openmp/ (Last accessed: 01.06.2024).

21. OpenMP Application Programming Interface. URL: https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5-2.pdf (Last accessed: 01.06.2024).

22. Using OpenMP with C - RCU of Colorado Boulder documentation. URL: https://curc.readthedocs.io/en/latest/programming/OpenMP-C.html (Last accessed: 01.06.2024).

**Useful**

23. MPI Forum. URL: https://www.mpi-forum.org/ (Last accessed: 01.06.2024).

24. Ukraine - Distributed Computing Team. URL: http://distributed.org.ua (Last accessed: 01.06.2024).