

# Introduction to OpenACC: Multi-platform Directive-based Parallelism

Yiru Sun, Chen Chen, Hengjiu Kang

9 June 2016\*

## Abstract

This is an OpenACC tutorial for beginners who want to write parallel programs with OpenACC directives. This tutorial is written under the class framework of ECS158 Parallel Programming taught by Professor Matloff in University of California, Davis.

---

\*First drafted on 29 May 2016

# Contents

<b>Contents</b>	<b>2</b>
<b>List of Figures</b>	<b>3</b>
<b>1 Introduction</b>	<b>4</b>
1.1 What Exactly is OpenACC? . . . . .	4
1.2 Motivations and Significance . . . . .	4
1.3 What Will Be Covered and What Will Not . . . . .	6
1.4 Setting up Environment Variables for Omni ompcc Compiler . . . . .	6
<b>2 Matrix Multiplication: Connecting C, OpenMP and CUDA to OpenACC</b>	<b>7</b>
2.1 Why Matrix Multiplication . . . . .	7
2.2 Matrix Multiplication in C . . . . .	7
2.3 Matrix Multiplication in OpenMP . . . . .	8
2.4 Matrix Multiplication in CUDA . . . . .	9
<b>3 Introduction to OpenACC</b>	<b>12</b>
3.1 Available Compilers . . . . .	12
3.2 Basic Directives . . . . .	12
3.3 Matrix Multiplication in OpenACC . . . . .	14
<b>4 OpenACC Example 1: Path Finder</b>	<b>18</b>
<b>5 OpenACC Example 2: Bright Spot Counter</b>	<b>23</b>

<i>List of Figures</i>	3
------------------------	---

<b>Bibliography</b>	<b>28</b>
---------------------	-----------

# List of Figures

# Chapter 1

## Introduction

### 1.1 What Exactly is OpenACC?

OpenACC is a directive-based programming interface that brings portability and ease of use to parallel programming. To understand how it works, it is important to have some prior knowledge about OpenMP and CUDA. Long story short, OpenACC is similar to OpenMP in a sense that both use directives to signal compilers what to do with the code. On top of that, OpenACC has added portability so that programs written with OpenACC can be easily transported to work with multiple platforms such as CUDA. This introduction will explain the significance of OpenACC and offer code examples to show how to start writing parallel programs with OpenACC directives.

### 1.2 Motivations and Significance

Modern image processing and machine learning tasks require more and more computational resources. Computer scientists around the world have been trying to address this ever-growing demand by turning to parallel computing. However, shifting to parallel programming is by no means apparent. While proper implementations of parallel algorithms offer potential performance improvements in

tasks such as matrix operations, non-optimized ones could even hinder the performance when compared with their serial counterparts. In order to hide programmers from the hassles of tricky low-level implementations, many parallel libraries and extensions have been developed to help programmers write parallel softwares more easily and promise to deliver reasonable results. Among those libraries, OpenMP is deemed by many as "*the **de facto** standard for shared-memory programming*"[1]. It is a software implementation that can be deployed on most of the multi-core CPU systems. OpenMP uses pragma directives that are relatively easy to program and can be ignored by compilers that don't understand them. Great convenience is probably one of the many reasons why OpenMP has been so successful, and the use of directives went on to inspire the interface of OpenACC as we will discuss later in detail.

Besides OpenMP, another very important event that has vastly shaped parallel programming and inspired OpenACC is the advent of general-purpose graphics processing unit (GPGPU). GPGPUs are very good at running parallel programs because they are packed with so many processing units in them. Over the years, the companies that design and manufacture GPGPUs have offered programmers with some interfaces and standards to interact with GPGPUs. One example is CUDA which was introduced by NVIDIA in order to interact with CUDA-enabled NVIDIA graphic cards. Despite the great potential of CUDA in performance improvement, many programmers are deterred by the horrendous details and efforts that go into writing a robust and efficient CUDA program. Some called for simpler interfaces which still reap the benefits of the massive parallelism offered by GPGPUs.

In a nutshell, OpenACC was born to inherit the ease of use of directive-based interfaces such as OpenMP, and was designed to be very portable across different platforms such as NVIDIA GPU, AMD Radeon, x86/ARM CPU and Intel Xeon Phi. In the following chapters, we will explore OpenACC directives in greater details by using an introductory matrix multiplication example and two homework problems adapted from ECS158 coursework. It is worth noting that the code examples in this introduction are in C style.

### 1.3 What Will Be Covered and What Will Not

This introduction will focus on how to use OpenACC directives to parallelize portions of the code (usually for loops), and how to optimize data locality to minimize unnecessary data transfer overheads. This introduction will not cover hardware-specific optimization, interoperability and other advanced features.

### 1.4 Setting up Environment Variables for Omni ompcc Compiler

The compiler of choice of this tutorial is Omni ompcc on UC Davis CSIF machines. In order for ompcc to run and use backends properly, we need to configure some environment variables first. The following scripts can be written into a .login file to run automatically when logging in, or they can be run manually.

```

1 set path = ($path /home/matloff/Pub)
2 set path = ($path /home/matloff/Pub/omni/bin)
3 set path = ($path /usr/lib64/mpich/bin)
4 set path = ($path /usr/local/cuda-5.5/bin)
5 set path = ($path /usr/local/cuda-7.5/bin)
6 setenv LD_LIBRARY_PATH
7 setenv LD_LIBRARY_PATH /usr/local/cuda-5.5/targets/x86_64-linux/
   lib:$LD_LIBRARY_PATH
8 setenv LD_LIBRARY_PATH /usr/local/cuda-7.5/targets/x86_64-linux/
   lib:$LD_LIBRARY_PATH

```

Listing 1.1: mmul\_acc

To compile OpenACC GPU code with ompcc:  
 ompcc -o <outputfile name> -acc <input c source code>

To compile OpenACC CPU code with ompcc:  
 ompcc -o <outputfile name> <input c source code>

## Chapter 2

# Matrix Multiplication: Connecting C, OpenMP and CUDA to OpenACC

### 2.1 Why Matrix Multiplication

In this chapter, we will show you different versions of a simple matrix multiplication example that are written in C, OpenMP and CUDA. Then we will try to link different aspects of C, OpenMP and CUDA to the design of OpenACC. We choose to start with matrix multiplication for it is one of the most important operations in science and engineering, which has sprung many modern applications such as image processing and voice recognition. The examples in this chapter should be relative easy to understand and mainly serve to refresh your memory. The detailed implementations can be skipped.

### 2.2 Matrix Multiplication in C

We started by writing a matrix multiplication function in C.

```
1 /*  
2  Introduction to OpenACC  
3  mmul.c Matrix Multiplication in C
```

```

4      Not optimized
5      Not claimed to be efficient
6  */
7  #include "stdio.h"
8  #include "stdlib.h"
9
10 void mmul_c(float* A, float* B, int row, int col, float* out){
11     // This code assume that out has been malloc-ed
12     int i=0;
13     int j=0;
14     int inner=0;
15     for (i = 0; i < row; i++) {
16         for (j = 0; j < col; j++) {
17             out[i*col+j] = 0.0;
18             for (inner = 0; inner < col; inner++) {
19                 out[i*col+j] += A[i*col+inner] * B[inner*col+j
20             ]*1.0;
21         }
22     }
23 }
24
25 void printMatrix(float* A, int row, int col){
26     int i=0;
27     int j=0;
28     for (i=0;i<row;i++){
29         for (j=0;j<col;j++){
30             printf("%.2f\t", A[i*col+j]);
31         }
32         printf("\n");
33     }
34 }
35
36 int main(){
37     // crate 3*3 matrix
38     float A[]={1,2,3,4,5,6,7,8,9};
39     float B[]={2,2,2,2,2,2,2,2,2};
40     float out[9];
41     printf("A is:\n");
42     printMatrix(A, 3, 3);
43     printf("B is:\n");
44     printMatrix(B, 3, 3);
45     mmul_c(A, B, 3, 3, out);
46     printf("out is:\n");
47     printMatrix(out, 3, 3);
48     return 0;
49 }

```

Listing 2.1: mmul\_c

## 2.3 Matrix Multiplication in OpenMP

Now we can turn the mmul function written in C into a parallel version of itself by using OpenMP directives.



```

1 void mmulOmp(float *A, float *B, int row, int col, float *out, )
2 {
3     #pragma omp parallel
4     {
5         int nth = omp_get_thread_num();
6         int total = omp_get_num_threads();
7         int division = row/total;
8         int i=0;
9         int j=0;
10        int inner = 0;
11        #pragma omp for
12        for (i = (division*nth); i < (nth==(total-1) ? row : (
13            division*nth+division)); i++) {
14            for (j = 0; col < col; j++) {
15                out[i*col+j] = 0;
16                for (inner = 0; inner < col; inner++) {
17                    out[i*col+j] += A[i*col+inner] * B[inner*col+
18                j]*1.0;
19            }
20        }
21    }
22 }

```

Listing 2.2: mmul\_omp

The OpenMP example above shows how directives can be used to start parallelism. Each thread starts to execute on its own starting from the **pragma omp parallel** directive, and the for loop is parallelized by using the **pragma omp for** directive. The designers of OpenACC adopted the concept of directives and they came up with their own versions of the similar directives in OpenACC.

## 2.4 Matrix Multiplication in CUDA

CUDA is a language extension that interfaces with NVIDIA GPUs. Matrix Multiplication can also be performed on a GPU as shown below.

```

1
2 --global-- void kernel(float *A, float *B, int row, int col,
3     float *out){
4     int nth = threadIdx.x;
5     int division = r1/total;
6     int i=0;
7     int j=0;
8     int inner=0;
9     for (i = (division*nth); i < (nth==(total-1) ? row : (
10        division*nth+division)); i++) {
11         for (j = 0; j < col; j++) {
12             out[i*col+j] = 0;

```

```

11         for (inner = 0; inner < col; inner++) {
12             out[i*col+j] += a[i*col+inner] * b[inner*col+j]
13             } * 1.0;
14         }
15     }
16 }
17
18 void mmul_cuda(float *A, float *B, int row, int col, float *out){
19     float *dA, *dB, *dout;
20     int msize = row*col*sizeof(float);
21     cudaMalloc((void**)&dA, msize);
22     cudaMalloc((void**)&dB, msize);
23     cudaMalloc((void**)&dm, msize);
24     cudaMemcpy(dA, A, msize, cudaMemcpyHostToDevice);
25     cudaMemcpy(dB, B, msize, cudaMemcpyHostToDevice);
26
27     int blockSize = 4; // use 4 blocks
28     dim3 dimGrid(blockSize, 1);
29     dim3 dimBlock(1,1,1);
30
31     kernel<<<dimGrid, dimBlock>>>(dA, dB, row, col, dm);
32
33     cudaMemcpy(out, dm, msize, cudaMemcpyDeviceToHost);
34     cudaFree(dA);
35     cudaFree(dB);
36     cudaFree(dm);
37
38     return;
39 }
40 }

```

Listing 2.3: mmul\_omp

Two key features of CUDA which have inspired OpenACC design are the explicit allocation of memory and transfer of data. With CUDA, programmers have decent control over data flow and memory allocation, which enables them to optimize data locality. OpenACC designers also want programmers to have some degree of control over data and memory, so they incorporate some special directives and structures dedicated to the purpose of data locality optimization. By reducing the overheads of data movement, a parallel program can see great performance improvement.

Another important design detail of CUDA is the separation between host memory space and device memory space. Device memory space can be thought of as a programmer-managed cache which is crucial to the performance of a GPU program. The unnecessary transfer of data between host and device often causes huge performance problem and OpenACC offers some directives to

control and manage the transfer of data between them.

Up until this point, chapter 2 has established connections between C, OpenMP, CUDA and OpenACC. Chapter 3 will introduce OpenACC in greater details.

## Chapter 3

# Introduction to OpenACC

### 3.1 Available Compilers

OpenACC offers a set of directives to suggest compilers what to do with the code, but the backend implementation depends on the compiler and flags being used. The list below shows several OpenACC compatible compilers.

Name	Developer	Commercial	Backend	Platform
Omni	Omni Compiler Project	NO	CUDA	X86 Only
accULL	High Performance Computing Group	NO	CUDA/OpenCL	ARM/X86
GNU-Openacc	GNU project	NO	CUDA/OpenCL	Not complete
PGI	PGI	YES	CUDA/OpenCL	X86 Only

Table 3.1: OpenACC compilers

As mentioned in the introductory chapter, we are using Omni omppc compiler in this tutorial.

### 3.2 Basic Directives

Just like OpenMP, OpenACC uses directives to communicate with compilers. Below is a list of basic directives that are useful in starting off with OpenACC.

1. `#pragma acc kernels`

**kernels** construct is probably the easiest one to work with. It can be placed before a for loop to signal the compiler that the for loop immediately after **kernels** directive might be parallelized. That being said, it is up to the compiler to decide whether the parallelism is going to happen and the outcome is not guaranteed. **kernels** is suitable for inexperienced programmers and cases where the loops are too complicated to analyze.

2. `#pragma acc parallel`

**parallel** construct is similar to `#pragma acc kernels`, but it goes further and tells the compiler explicitly that the for loop immediately below can be parallelized. The programmer is responsible for making sure that the loop is safe to parallelize. It is worth noting that **parallel** needs to be used with another directive to properly function.

3. `#pragma acc loop [options]`

**loop** directive is usually used together with **parallel** directive in the form of `#pragma acc parallel loop` in order to signal the compiler to parallelize the loop immediately after the resultant directive.

4. `#pragma acc data [options]`

**data** construct of OpenACC precedes a special set of clauses for data transfer and memory allocation. The options of **data** directive are similar to a combination of CUDA functions working collectively. For example, the **copy** clause can be used to form `#pragma acc data copy(A)` and the resultant directive is similar to the following CUDA pseudocode:

```

1 cudaMalloc(A_device) // Allocate Memory on Device
2 cudaMemcpy(A_host, A_device, cudaMemcpyHostToDevice) //
  Initialize A with the Data from Host
3 cudaMemcpy(A_host, A_device, cudaMemcpyDeviceToHost) //
  Copy the Resultant Data back to Host
4 cudaFree(A_device) // Free Device Memory
```

Listing 3.1: data\_copy\_CUDA

### 3.3 Matrix Multiplication in OpenACC

We can now implement matrix multiplication in OpenACC. The following code is based on the C version of matrix multiplication in Chapter 2.

```

1 void mmul_acc(float* A, float* B, int row, int col, float* out){
2     // This code assume that out has been malloc-ed
3     int i=0;
4     int j=0;
5     int inner=0;
6     float temp_sum=0;
7     #pragma acc data create(A[0:(row*col)],B[0:(row*col)]) copy(
8         out[0:row*col])
9     {
10         #pragma acc parallel loop
11         for (i = 0; i < row; i++) {
12             #pragma acc loop
13             for (j = 0; j < col; j++) {
14                 out[i*col+j] = 0.0;
15                 #pragma acc loop reduction(+:temp_sum)
16                 for (inner = 0; inner < col; inner++) {
17                     temp_sum = A[i*col+inner] * B[inner*col+j]
18                 }
19                 out[i*col+j]=temp_sum;
20                 temp_sum=0;
21             }
22         }
23     }

```

Listing 3.2: mmul\_acc

In this example, we use the following directives to implement parallelism:

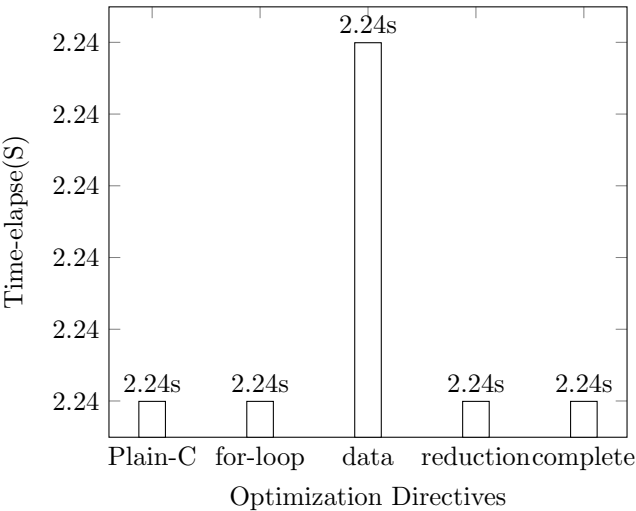
- `#pragma acc data create(A, B)`  
A and B are two input matrices and will not be retrieved from the device after the computations are finished. So it is safe to use **data create** construct instead of **data copy** to create dedicated memory space for A and B in the device but do not initialize or copy back to host after the job has been done. This move is usually very effective in improving data locality and reduce communication overhead.
- `#pragma acc parallel loop`  
This directive enables loop optimization and parallelism as we mentioned in point 3 of 3.2.

- `#pragma acc loop`  
Inside a nested for-loop, we need to use **`#pragma acc parallel loop`** first right before the outermost for-loop, then use **`#pragma acc loop`** in the inner for-loop to signal further parallelization.
- `reduction(+:temp_sum)`  
**reduction** is a special directive in OpenACC. It can monitor the value of a certain variable (in this case: `temp_sum`), perform a certain function (in this case: `+`) to the variable and reduce the results of all the operations into one single result. One of the following basic functions can act as a reduction function: `+`, `*`, **`min`**, **`max`**.

### Performance comparison in directives

Also, we need to be aware that, different directives have different optimization performance, and the performance depends on each application. For example, we want **`for loop`** optimization and **`reduce`** optimization for dot multiplication summation. In other application, like path-finding algorithm, we may want to create on-device memory space for faster variable access, because path-finding is memory accessing massive algorithm.

In this section, we compare the improvement of each directive to original program step by step. When we are using GCC to compile plain-C code without any flag, the execution time is about 12 seconds.

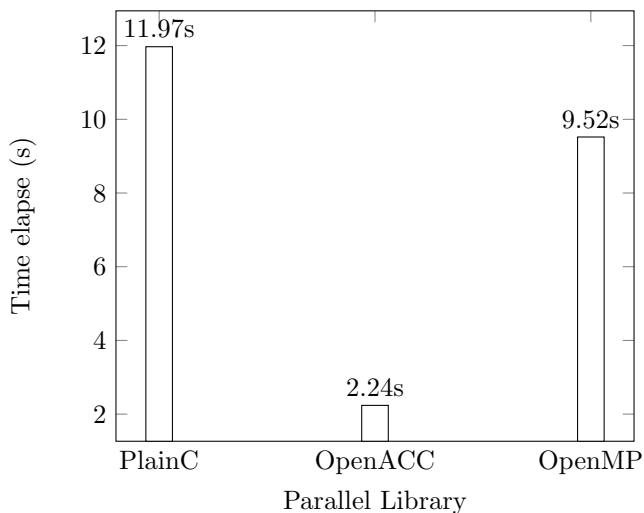


We found that the differences between directives are not significant, which is out of expectation. An possible explanation is that this Matrix multiplication is an typical application that ompcc can easily maximize the optimization. So it is reasonable that ompcc compiled plain-C code is already 5 times faster than gcc compiler, and user added customized directives will be ignored or have the same effect.

### Performance Comparison between Difference Libraries

Our test program performs a 50\*50 matrix multiplication for 30000 times, and the following time elapse data were recorded from running the command: `time mmul_x.out`, where `x` refers to `c`, `acc` or `omp`





It is easily observed that the performance improvement by using OpenACC directives is enormous in the case of matrix multiplication, partly because it is embarrassingly parallel and therefore benefit largely from the use of parallel computing.

Chapter 3 has covered most of the OpenACC directives in detail. The following two chapters will introduce some code examples to show how OpenACC can be applied in various applications.

## Chapter 4

# OpenACC Example 1: Path Finder

This chapter is all about OpenACC in action and a new **data update** directive. We use OpenACC to parallelize a path finder program originally written in C.

The path finder program takes an adjacency matrix and find all possible paths of length **k** that are connected, starting from a given vertex. Repeated vertices are allowed.

For example, if **k** is 5 and row **i** of **paths** is 99 2 0 3 8 88  
This is the path  $99 \rightarrow 2 \rightarrow 0 \rightarrow 3 \rightarrow 8 \rightarrow 88$   
Loops are allowed, e.g.  $5 \rightarrow 2 \rightarrow 0 \rightarrow 2 \rightarrow 0 \rightarrow 88$

Here is a list of arguments taken by the path finder function:  
**adjm**: Adjacency matrix of a directed graph.  
**n**: Number of vertices in the graph.  
**k**: Length of each path.  
**paths**: Output matrix of paths.  
**numpaths**: Number of paths.

The output matrix: **paths**, will have **numpaths** rows and **k+1**

columns, and likely to have unused rows at the bottom. The space will be pre-allocated in the main function.

Here is the C code to implement the path finder:

```

1 #include <stdio.h>
2 #include <time.h>
3 #include <sys/time.h>
4 #include <stdlib.h>
5
6 double get_wall_time(){
7     struct timeval time;
8     if (gettimeofday(&time,NULL)) return 0;
9     return (double)time.tv_sec + (double)time.tv_usec * .000001;
10 }
11
12 void findpaths(int *adjm, int n, int k, int *paths, int *numpaths)
13 {
14     int count = 0;
15     int times = n;
16     for (int i = 0; i < k; ++i) {times = times * n;} //
17     // Dimension of time = n^k
18
19     int* arr = (int*) malloc (k+1); // Create an array arr that
20     // is indexed from 0 to k
21
22     for (int i = 0; i < k+1; ++i) {arr[i] = 0;} // Initialize all
23     // elements in arr to be 0
24     for (int i = 0; i < times; ++i)
25     {
26         int temp = i, index = k, check = 1, r;
27         while(temp!=0)
28         {
29             r=(temp%n);
30             temp/=n;
31             arr[index] = r;
32             index--;
33         }
34
35         for (int i = 0; i < k; ++i)
36         {
37             check *= adjm[arr[i]*n+arr[i+1]];
38         }
39
40         if(check)
41         {
42             for (int i = 0; i < k+1; ++i)
43             {
44                 paths[count*(k+1)+i] = arr[i];
45             }
46             count++;
47         }
48     }
49     *numpaths = count;
50 }
51
52 int main(int argc, char *argv[])
53 {
54     int n = 5, k = 4;
55     int path[2000000];
56     int numpaths;
57     int loopitr;

```

```

54
55 static int A[400] = {
56     0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0,
57     1, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1,
58     1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0,
59     1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0,
60     1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0,
61     1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1,
62     1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
63     0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0,
64     0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1,
65     1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1,
66     0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0,
67     1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 1, 1,
68     0, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0,
69     0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0,
70     0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0,
71     1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1,
72     0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,
73     0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1,
74     1, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0,
75     0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 0, 0};
76
77 double time_begin=get_wall_time();
78 for (loopitr = 0; loopitr < 30000; loopitr++) {
79     findpaths(A, n, k, path, &numpaths);
80 }
81 double time_end=get_wall_time();
82 printf("Total elapsed time is: %.4lf\n", (time_end - time_begin));
83
84 printf("%d\n", numpaths);
85 for (int i = 0; i < numpaths; ++i)
86 {
87     for (int j = 0; j < k+1; ++j)
88     {
89         printf("%d->", path[i*(k+1)+j]);
90     }
91     printf("\n");
92 }
93
94 return 0;
95 }

```

Listing 4.1: path

To look for paths with length of 4, each path would have 5 vertices. Since loop is allowed and we have a total of  $n$  vertexes, we would check all possible paths from  $0 \rightarrow 0 \rightarrow 0 \rightarrow 0 \rightarrow 0$  to  $n \rightarrow n \rightarrow n \rightarrow n \rightarrow n$ .

The total run times would be  $n$  to the power of 5 in line 20. Let's give an order to these vertexes:  $i \rightarrow j \rightarrow k \rightarrow l \rightarrow m$ . We need to check each path in the matrix to make sure the whole path

is correct. If  $m_{i,j}$  equals 1, there is a path from  $i$  to  $j$ . Thus, we would check  $m_{i,j}$ ,  $m_{j,k}$ ,  $m_{k,l}$  and  $m_{l,m}$  in line 31. If the path exists, we count it and put it in to the output path matrix.

```

1 #include <stdio>
2 #include <ctime>
3 #include <omp.h>
4 #include <sys/time.h>
5
6 using namespace std;
7
8 double get_wall_time(){
9     struct timeval time;
10    if (gettimeofday(&time,NULL)) return 0;
11    return (double)time.tv_sec + (double)time.tv_usec * .000001;
12 }
13
14 void findpaths(int *adjm, int n, int k, int *paths, int *numpaths)
15 {
16     int count = 0;
17     int times = n;
18     for (int i = 0; i < k; ++i) {times = times * n;} //
19     // Dimension of time = n^k
20     int *arr = new int [k+1]; // Create an array arr that is
21     // indexed from 0 to k
22     for (int i = 0; i < k+1; ++i) {arr[i] = 0;} // Initialize all
23     // elements in arr to be 0
24     #pragma acc parallel loop
25     for (int i = 0; i < times; ++i)
26     {
27         int temp = i, index = k, check = 1, r;
28         while(temp!=0)
29         {
30             r=(temp%n);
31             temp/=n;
32             arr[index] = r;
33             index--;
34         }
35         for (int i = 0; i < k; ++i)
36         {
37             check *= adjm[arr[i]*n+arr[i+1]];
38         }
39         if (check)
40         {
41             for (int i = 0; i < k+1; ++i)
42             {
43                 paths[count*(k+1)+i] = arr[i];
44             }
45             #pragma acc atomic update
46             count++;
47         }
48     }
49
50     *numpaths = count;
51 }
52
53 int main(int argc, char const *argv[])

```

```

54 {
55     int n = 5, k = 4;
56     int path[2000000], numpaths;
57     static int A[400] = {
58         0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0,
59         1, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1,
60         1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0,
61         1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0,
62         1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0,
63         1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1,
64         1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
65         0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0,
66         0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1,
67         1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1,
68         0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0,
69         1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1,
70         0, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0,
71         0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0,
72         0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0,
73         1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1,
74         0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
75         0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1,
76         1, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0,
77         0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 0, 0};
78
79     double time_begin=get_wall_time();
80     findpaths(A, n, k, path, &numpaths);
81     double time_end=get_wall_time();
82     printf("Total elapsed time is: %.4lf\n", (time_end - time_begin));
83
84     printf("%d\n", numpaths);
85     for (int i = 0; i < numpaths; ++i)
86     {
87         for (int j = 0; j < k+1; ++j)
88         {
89             printf("%d->", path[i*(k+1)+j]);
90         }
91         printf("\n");
92     }
93
94
95     return 0;
96 }

```

Listing 4.2: path\_acc

## Chapter 5

# OpenACC Example 2: Bright Spot Counter

Consider an  $n \times n$  matrix of image pixels, with brightness values in  $[0,1]$ . Let's define a bright spot of size  $k$  and threshold  $b$  to be a  $k \times k$  subimage, with the pixels being contiguous and with each one having brightness at least  $b$ . Here is a C++ program that returns a count of all bright spots.

$n$  is matrix dimension

$k$  is subimage dimension

start and end are for parallelism

$\text{pix}$  is the pointer to the matrix.

```
1 #include <stdlib.h>
2 #include <stdio.h>
3 #include <time.h>
4 #define MATRIX_MAX 1
5 #define MAX_SIZE 1400
6 #define N_ITERS 1000
7
8
9 void seed(int r, int c, float *out) { //generate a random
    matrix
10     time_t t;
11     srand((unsigned) time(&t));
12     for (int i=0; i<r; i++){
13         for (int j=0; j<c; j++){
14             out[i*c+j]=(float)(MATRIX_MAX*((int)(rand()%100)
15                 /100.0));
16         }
17     }
```

```

17 }
18
19
20 // n = matrix dimension; k = subimage dimension; start and
21 // end are for parallelism; pix is the pointer to the matrix.
22 int brightssplit (float *pix, int n, int k, float thresh, int
23 // start, int end) {
24 // float count = 0, check[n];
25 // for (int i = 0; i < n; ++i)
26 // { // iterate throught ALL n rows
27 // // Clear book keeping variables for a new round
28 // // float minus, sum = 0;
29 // // Count the number of bright pixel in a single row of 1 x k
30 // // i * n locate row, + j locate column
31 // // for (int j = start; j < start+k; ++j) if (pix[i*n+j] >
32 // // thresh) {sum++;}
33 // // end - k + 1 basically recovers the boundary by the 1st
34 // // column index of each window
35 // // which could be understood as the window column index
36 // // instead of the last column index which is end itself
37 // // Now we run from the starting window column index to the
38 // // end WCI
39 // // Visually, we are moving a 1 x k block across row i to
40 // // track bright block
41 // // for (int j = start; j < (end-k+1); ++j)
42 // // {
43 // // // minus:
44 // // // This is a smart way of moving the 1 x k block to the right
45 // // //
46 // // // Instead of using excessive for loops and repetitive
47 // // // checking,
48 // // // we simply look at the left most pixel of that block,
49 // // // and later we can see what we could do with sum.
50 // // // if (pix[i*n+j] > thresh) {minus = 1;} else {minus = 0;}
51 // // // Check List for keeping track of the number of 1 x k bright
52 // // // unit
53 // // // in terms of window column index
54 // // // Initialized when going through the first row
55 // // //
56 // // // if (i == 0) {if (sum >= k) {check[j] = 1;} else {check[j]
57 // // // = 0;}}
58 // // // Now we are out of the first row, that means we get to +/-
59 // // // to the list,
60 // // // and hopefully hit a bring window when count >= k!
61 // // // Keep Count++ when check[j] >= k
62 // // // Clear it when we hit a 1 x k block with non-bright pixels.
63 // // // else {
64 // // // if (sum >= k) {
65 // // // check[j] ++;
66 // // // if (check[j] >= k) {count ++;}}
67 // // // else {check[j] = 0;}}
68 // // // Now we are about to move the 1 x k block to the right. Rmb
69 // // // the minus?
70 // // // We can now use it to change sum, together with the
71 // // // information of
72 // // // the one new pixel the 1 x k block now includes.
73 // // // Again, elegant way of avoiding repetitive checking!
74 // // //
75 // // // if(minus == 0 && pix[i*n+j+k] > thresh) {sum++;}
76 // // // else if(minus == 1 && pix[i*n+j+k] <= thresh) {sum--;}
77 // // // }
78 // // }
79 }

```



```

66
67     return count;
68 }
69
70 int brights (float *pix, int n, int k, float thresh) {
71     brightssplit (pix, n, k, thresh, 0, n);
72 }
73
74
75 int main (int argc, char const *argv[])
76 {
77
78     float A[MAX_SIZE*MAX_SIZE]; // use one dimensional array to
79                                // represent the matrix
80     seed (MAX_SIZE,MAX_SIZE,A); // generate random matrix
81     int count = 0;
82     int i=0;
83     for (i=0;i<N_ITSERS;i++){
84         count += brights (A, 1000, 5, 0.2);
85     }
86     printf("There are average %d starts in %d iters.\n", count/
87           N_ITSERS, N_ITSERS);
88     return 0;

```

Listing 5.1: Brightcounter

We run the function `brights()` 1000 times to get a large run time. In the function `brightssplit()`, we search the matrix row by row with the for loop in line 23. Then, we check if there is `k` numbers of continuous bright pixels from line 29 and store the information in the first row of `pix` matrix. When we checking each row, we add the information to the first row of `pix` matrix. Finally, we calculate the data in the first row of `pix` matrix and give out the number of bright spots.

```

1 #include <stdlib.h>
2 #include <stdio.h>
3 #include <time.h>
4 #define MATRIX_MAX 1
5 #define MAX_SIZE 1400
6 #define N_ITSERS 1000
7
8
9 void seed(int r, int c, float *out) { //generate a random
10     matrix
11     time_t t;
12     srand((unsigned) time(&t));
13     for (int i=0;i<r;i++){
14         for (int j=0;j<c;j++){
15             out[i*c+j]=(float)(MATRIX_MAX*((int)(rand()%100)
16             /100.0));
17         }
18     }

```

```

17 }
18
19
20 #pragma acc routine seq
21 int brightssplit (float *pix, int n, int k, float thresh, int
22   start, int end) {
23   float count = 0, check[n];
24   #pragma acc data copy(pix[0:n*n], check[0:n])
25   {
26       #pragma acc loop
27       for (int i = 0; i < n; ++i)
28       {
29           float minus, sum = 0;
30
31           for (int j = start; j < start+k; ++j){
32               if (pix[i*n+j] > thresh){
33                   sum++;
34               }
35           }
36           #pragma acc loop
37           for (int j = start; j < (end-k+1); ++j)
38           {
39               if (pix[i*n+j] > thresh){
40                   minus = 1;
41               } else {
42                   minus = 0;
43               }
44
45               if (i == 0) {
46                   if (sum >= k) {
47                       check[j] = 1;
48                   } else {
49                       check[j] = 0;
50                   }
51               }
52               else {
53                   if (sum >= k) {
54                       check[j] ++;
55                       if (check[j] >= k) {
56                           count ++;
57                       }
58                   }
59                   else {
60                       check[j] = 0;
61                   }
62               }
63
64               if(minus == 0 && pix[i*n+j+k] > thresh) {sum++;}
65               else if(minus == 1 && pix[i*n+j+k] <= thresh) {
66                   sum--;}
67           }
68       }
69
70       return count;
71   }
72
73 int brights (float *pix, int n, int k, float thresh) {
74     int i=0;
75     int blocks = n/k-1;
76     int count = 0;

```

```
77 #pragma acc data create(pix)
78 #pragma acc parallel loop reduction(+:count)
79 for (i=0;i<blocks;i++){
80     count += brightssplit (pix, n, k, thresh, k*i, k*(i+1)+k
81         -1);
82 }
83 count += brightssplit(pix, n, k, thresh, k*i, n);
84 return count;
85 }
86
87
88 int main (int argc, char const *argv[])
89 {
90
91     float A[MAX_SIZE*MAX_SIZE];
92     seed(MAX_SIZE,MAX_SIZE,A);
93
94     int count = 0;
95     int i=0;
96     for (i=0;i<N_ITERS;i++){
97         count += brights (A, 1000, 5, 0.2);
98     }
99
100     printf("There are average %d starts in %d iters.\n", count/
101         N_ITERS, N_ITERS);
102     return 0;
103 }
```

Listing 5.2: Brightcounter\_acc

# Bibliography

- [1] Norm Matloff. *Programming on Parallel Machines*. University of California, Davis, 1993.
- [2] Jeff Larkin (2015) *COMPARING OPENACC AND OPENMP PERFORMANCE AND PROGRAMMABILITY*. <http://on-demand.gputechconf.com/gtc/2015/presentation/S5196-Jeff-Larkin.pdf>
- [3] Omni Compiler, <http://omni-compiler.org/>
- [4] OpenACC official site, <http://www.openacc.org/>
- [5] OpenACC at Nvidia: <https://developer.nvidia.com/openacc>