

An Introduction to OpenACC

William Killian

Department of Computer and Information Science
University of Delaware



CISC 879 — Advanced Parallel Programming

Outline

- 1 Introduction
 - Accelerating Applications on Various Architectures
 - About OpenACC
- 2 The OpenACC API
 - Directives
 - Clauses
 - Runtime Library Routines
- 3 Compiling
 - Compiling
- 4 Examples
 - Vector-Vector Addition
 - Matrix-Matrix Multiplication
- 5 Summary

Methods of Accelerating Applications

- Targeting Multi-core CPUs
- Targeting Multi-node CPUs
- Targeting GPUs

Methods of Accelerating Applications

- Targeting Multi-core CPUs
 - OpenMP
 - PThreads, QThreads, etc ...
- Targeting Multi-node CPUs
- Targeting GPUs

Methods of Accelerating Applications

- Targeting Multi-core CPUs
 - OpenMP
 - PThreads, QThreads, etc ...
- Targeting Multi-node CPUs
 - MPI
- Targeting GPUs

Methods of Accelerating Applications

- Targeting Multi-core CPUs
 - OpenMP
 - PThreads, QThreads, etc ...
- Targeting Multi-node CPUs
 - MPI
- Targeting GPUs
 - CUDA

Methods of Accelerating Applications

- Targeting Multi-core CPUs
 - OpenMP
 - PThreads, QThreads, etc ...
 - OpenCL
- Targeting Multi-node CPUs
 - MPI
- Targeting GPUs
 - CUDA
 - OpenCL

Methods of Accelerating Applications

- Targeting Multi-core CPUs
 - OpenMP
 - PThreads, QThreads, etc ...
 - OpenCL
- Targeting Multi-node CPUs
 - MPI
- Targeting GPUs
 - CUDA
 - OpenCL

What's left? ...**OpenACC!**

Methods of Accelerating Applications

- Targeting Multi-core CPUs
 - OpenMP
 - PThreads, QThreads, etc ...
 - OpenCL
- Targeting Multi-node CPUs
 - MPI
- Targeting GPUs
 - CUDA
 - OpenCL
 - OpenACC

What's left? ...**OpenACC!**

What is OpenACC?

- is an API using compiler directives that
- allows for small segments of code, called kernels, to be run on the GPU and
- requires little to no modifications to the original program
- is compatible with C/C++ and Fortran

Reason Behind Formation

OpenACC was formed to help create and foster a cross platform API that would allow any scientist or programmer to easily accelerate their application on modern many-core and multi-core processors using directives.

History of OpenACC

- Initially collaboration between CAPS Enterprise, Cray Inc., The Portland Group (PGI), and NVIDIA
- Built from OpenMP-style directives
 - `#pragma omp parallel` vs. `#pragma acc parallel`
 - Creators of OpenACC are all members of the OpenMP Working Group on accelerators
- Standardized in November 2011 at SuperComputing 2011
- Compilers available from Cray, CAPS, and PGI
- Potential API merge with OpenMP in the future?

Directives Overview

Format

```
#pragma acc directive-name [clause [[, clause] ...]
```

Possible directives are:

<code>parallel</code>	starts parallel execution on the accelerator
<code>kernels</code>	defines a region that should be converted to a kernel
<code>data</code>	defines contiguous data to be allocated on the accelerator
<code>host_data</code>	makes the address of accelerator data available on the host
<code>loop</code>	defines type of parallelism to apply to proceeding loop

Directives Overview (continued)

<code>cache</code>	defines elements or subarrays that should be fetched into cache
<code>declare</code>	defines that a variable should be allocated in accelerator memory
<code>update</code>	update all or part of host memory from device memory, or vice versa
<code>wait</code>	forces program to wait for completion of asynchronous activity

Clauses Overview

Each directive can have zero (or more) clauses associated.

Example clauses are:

`if (e)` condition used to determine if command should be executed (data transfer, accelerator computation, etc)

`async [(n)]` tells the current command to be executed asynchronously. Used with `wait` for synchronization.

Clauses found in either kernels or parallel directives:

`reduction (op:list)`

`private (list)`

`firstprivate (list)`

Clauses Overview

Each directive can have zero (or more) clauses associated.

Example clauses are:

`if (e)` condition used to determine if command should be executed (data transfer, accelerator computation, etc)

`async [(n)]` tells the current command to be executed asynchronously. Used with `wait` for synchronization.

Clauses found in either kernels or parallel directives:

`reduction (op:list)`

`private (list)`

similar to OpenMP

`firstprivate (list)`

Clauses – parallel and loop

Clauses – parallel directive

`num_gangs(e)` specify the number of gangs to execute in the region

`num_workers(e)` specify number of workers to launch in each gang

`vector_length(e)` define vector length to use

Clauses – loop directive

`collapse(n)` specifies # of loops associated

`gang(e)` distribute across gang

`worker(e)` distribute across worker (within gang)

`vector(e)` operate in SIMD (within gang or worker)

`seq` execute sequentially on the accelerator

`independent` tell the compiler loops are data-independent

Clauses — Data Operations (optional with most directives)

<code>copy(list)</code>	transfer to/from device
<code>copyin(list)</code>	transfer to device
<code>copyout(list)</code>	transfer from device
<code>create(list)</code>	allocate on device
<code>present(list)</code>	data which is already on the device
<code>deviceptr(list)</code>	used to inform which variables are device pointers (opposed to host)
<code>device_resident(list)</code>	allocate on device instead of host

Clauses — Data Operations (optional with most directives)

<code>pcopy(list)</code>	transfer to/from device
<code>pcopyin(list)</code>	transfer to device
<code>pcopyout(list)</code>	transfer from device
<code>pcreate(list)</code>	allocate on device
<code>present(list)</code>	data which is already on the device
<code>deviceptr(list)</code>	used to inform which variables are device pointers (opposed to host)
<code>device_resident(list)</code>	allocate on device instead of host

Checks for presence before issuing data command

Clauses — host_data and update

Clauses — use with host_data directive

`use_device (list)` make the device address data available
in host code

Clauses — use with update directive

`host (list)` variables to copy from device to host

`device (list)` variables to copy from host to device

Data Clauses

We mentioned data clauses such as `copy` and `create` but never went over what we do when the memory was dynamically allocated (using `malloc`).

What should we do?

We can specify the size of the data!

`A` is our array of size n but we need to provide a hint to OpenACC

Solution:

```
#pragma acc kernels copyin(A[0:n])
```

Note

This will also work for 2-dimensional arrays i.e. `A[0:m*n]`

Combining Clauses

Observation

Similar to OpenMP, we can combine directives

- `#pragma acc parallel loop [clause [[,] clause]...]`
- `#pragma acc kernels loop [clause [[,] clause]...]`

Notice

A loop must directly follow, similar to `parallel for` in OpenMP

Runtime Routines

Runtime calls allow the programmer to obtain information about the host and accelerators during runtime, instead of compile time.

List of library routines:

```
int acc_get_num_devices (acc_device_t);  
    gets number of devices of passed type  
int acc_set_device_type (acc_device_t);  
    sets device type to use  
acc_device_t acc_get_device_type ();  
    gets current device type  
void acc_set_device_num (int, acc_device_t);  
    sets device based on index and type  
void acc_get_device_num (acc_device_t);  
    gets current device number
```

Runtime Routines — Synchronization

<code>int acc_async_test</code> <code>(int);</code>	tests to see if a specified async. tasks are completed
<code>int acc_async_test_all</code> <code>();</code>	tests to see if all async. tasks are completed
<code>void acc_async_wait</code> <code>(int);</code>	waits until specified async. task is completed
<code>void acc_async_wait_all</code> <code>();</code>	waits until all async. tasks are completed

Runtime Routines — Setup and Teardown

<code>void acc_init (acc_device_t);</code>	initialize OpenACC runtime for passed device type
<code>void acc_shutdown (acc_device_t);</code>	shut down connection to passed device type
<code>int acc_on_device (acc_device_t);</code>	tells the program whether it's executing on passed device type
<code>void* acc_malloc (size_t);</code>	allocates memory on the device
<code>void acc_free (void*);</code>	frees memory on the device

Compiling OpenACC

There are a few different compilers available for OpenACC.

We will be using HMPP Workbench 3.2.1 by CAPS Enterprise

In addition to OpenACC, HMPP Workbench also supports another directive-based accelerator language, HMPP (and consequently OpenHMPP).

Compiling using Built-in (supplied) Makefiles

- Obtain HMPP Workbench 3.2.1 (see website for details)
- Extract the tarball
- Extract the OpenACC Labs tarball
- Navigate to the OpenACC_Labs/CUDA/C/ directory
- Copy an existing lab
- Edit the code
- Invoke “make”

Compiling using the command line

Sample Invocation

```
hmpc --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

Compiling using the command line

Sample Invocation

```
hmpp --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

- **hmpp compiler**

Compiling using the command line

Sample Invocation

```
hmpc --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

- hmpc compiler
- specify OpenACC codelet target (CUDA or OPENCL)

Compiling using the command line

Sample Invocation

```
hmpc --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

- hmpc compiler
- specify OpenACC codelet target (CUDA or OPENCL)
- **force proper codelet creation for compilation**

Compiling using the command line

Sample Invocation

```
hmpc --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

- hmpc compiler
- specify OpenACC codelet target (CUDA or OPENCL)
- force proper codelet creation for compilation
- **compiler to use for host code**

Compiling using the command line

Sample Invocation

```
hmpp --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

- hmpp compiler
- specify OpenACC codelet target (CUDA or OPENCL)
- force proper codelet creation for compilation
- compiler to use for host code
- flags for host compiler

Compiling using the command line

Sample Invocation

```
hmpp --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

- hmpp compiler
- specify OpenACC codelet target (CUDA or OPENCL)
- force proper codelet creation for compilation
- compiler to use for host code
- flags for host compiler
- specify output file

Compiling using the command line

Sample Invocation

```
hmpc --openacc-target=CUDA --codelet-required  
gcc -O2 -o mvmult mvmult.c
```

- hmpc compiler
- specify OpenACC codelet target (CUDA or OPENCL)
- force proper codelet creation for compilation
- compiler to use for host code
- flags for host compiler
- specify output file
- **source file(s)**

Problem Overview

Problem

Given two vectors, A and B , each of size n , we wish to compute the per-component addition and store the result into C .

Pseudocode

```
int i;
for (i = 0; i < N; ++i) {
  C [i] = A [i] + B [i];
}
```

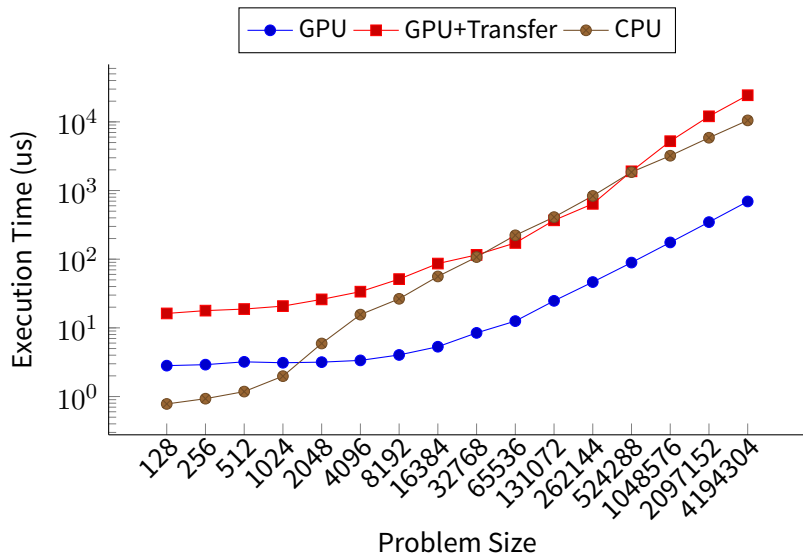
C Implementation

```
1  const int N = 1000;
2  float A [N];
3  float B [N];
4  float C [N];
5  int i;
6
7  // Initialization Loop
8  for (i = 0; i < N; ++i) {
9      A [i] = i;
10     B [i] = 2* i - 1;
11 }
12
13 // Computation Loop
14
15 for (i = 0; i < N; ++i) {
16     C [i] = A [i] + B [i];
17 }
```

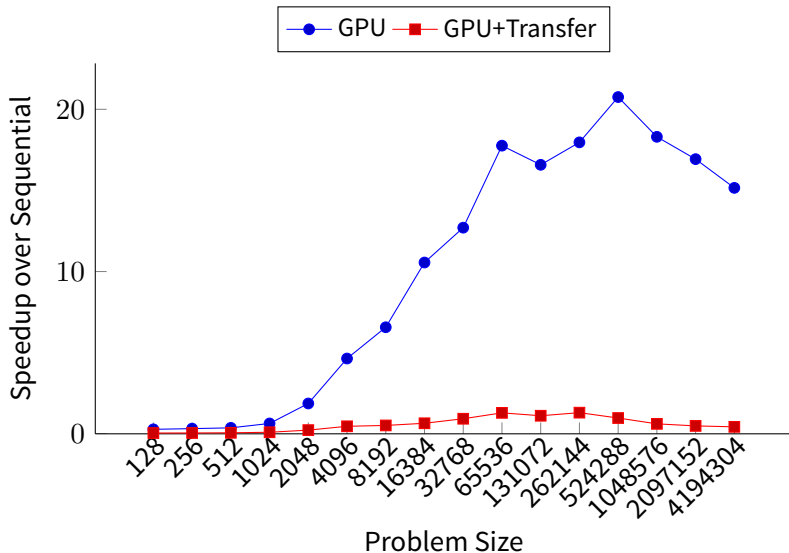
OpenACC Implementation

```
1  const int N = 1000;
2  float A [N];
3  float B [N];
4  float C [N];
5  int i;
6
7  // Initialization Loop
8  for (i = 0; i < N; ++i) {
9      A [i] = i;
10     B [i] = 2* i - 1;
11 }
12
13 // Computation Loop
14 #pragma acc kernels loop independent copyin(A,B), copyout(C)
15 for (i = 0; i < N; ++i) {
16     C [i] = A [i] + B [i];
17 }
```

Execution time of Vector-Vector Addition



Speedup of Vector-Vector Addition



Problem Overview

Problem

Given two matrices, A and B , with A having dimensions $m \times p$ and B having dimensions $p \times n$, we wish to compute the row-column inner product into C , a matrix with dimensions $m \times n$.

Pseudocode

```
int i, j, k;
for (i = 0; i < M; ++i)
  for (j = 0; j < N; ++j) {
    C [i][j] = 0;
    for (k = 0; k < P; ++k)
      C [i][j] += A [i][k] * B [k][j];
  }
```


C Implementation (Initialization)

```
5 #define INDEX(M,N,i,j) (i + j * M)
6
7 int main() {
8
9     float* A; float* B; float* C;
10    int i, j, k;
11
12    A = (float*) malloc (M * P * sizeof (float));
13    B = (float*) malloc (P * N * sizeof (float));
14    C = (float*) malloc (M * N * sizeof (float));
15
16    for (i = 0; i < P; ++i) {
17        for (j = 0; j < M; ++j)
18            A [INDEX(M,P,i,j)] = (float) rand () / RAND_MAX;
19        for (k = 0; k < N; ++k)
20            B [INDEX(P,N,k,i)] = (float) rand () / RAND_MAX;
21    }
```

C Implementation (Computation)

```
22  for (i = 0; i < M; ++i) {
23      for (j = 0; j < N; ++j) {
24          float sum = 0.0f;
25          for (k = 0; k < P; ++k) {
26              sum += A [INDEX(M,P,i,k)] * B [INDEX(P,N,k,j)];
27          }
28          C [INDEX(M,N,i,j)] = sum;
29      }
30  }
```

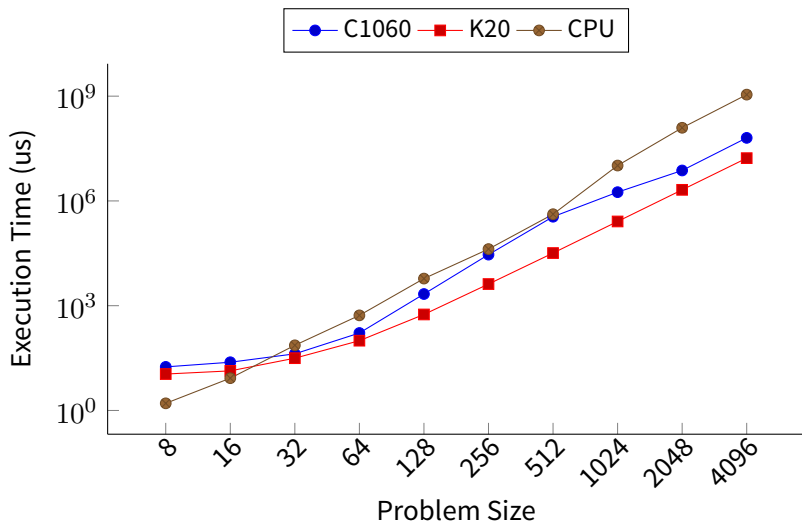
OpenACC Implementation (Computation)

```

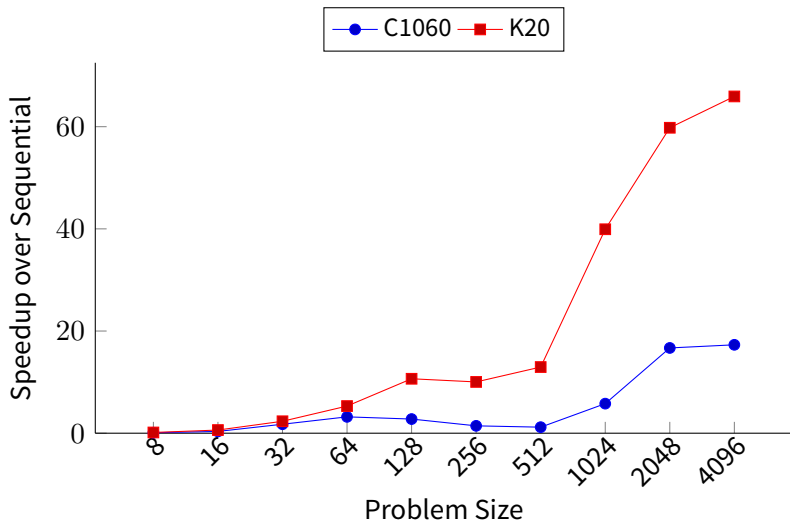
23  int m, n, p;
24  m = M; n = N; p = P;
25  // computation
26  #pragma acc kernels copyin(A[0:m*p],B[0:p*n]), copyout(C[0:m*n])
27  {
28  #pragma acc loop independent
29      for (i = 0; i < M; ++i) {
30  #pragma acc loop independent
31      for (j = 0; j < N; ++j) {
32          float sum = 0.0f;
33          for (k = 0; k < P; ++k) {
34              sum += A [INDEX(M,P,i,k)] * B [INDEX(P,N,k,j)];
35          }
36          C [INDEX(M,N,i,j)] = sum;
37      }
38  }
39  }
40

```

Execution time of Matrix-Matrix Multiplication



Speedup of Vector-Vector Addition



Summary

- OpenACC makes targeting **accelerators much easier**
- Designed for use by **scientists** to make GPGPU much easier
- Syntax similar to **OpenMP**
- Compiling with **HMPP Workbench 3.2.1** can target **CUDA** or **OpenCL**

OpenACC Reference API

http://www.openacc.org/sites/default/files/OpenACC.1.0_0.pdf