

INTRODUCTION TO ACCELERATED COMPUTING WITH OPENACC

Jeff Larkin, NVIDIA Developer Technologies

AGENDA

- NVIDIA Introduction
- Accelerated Computing Basics
- > What are Compiler Directives?
- Accelerating Applications with OpenACC
 - Identifying Available Parallelism
 - Exposing Parallelism
 - > Optimizing Data Locality
- Next Steps

NVIDIA Introduction



ACCELERATED COMPUTING 10X PERFORMANCE & 5X ENERGY EFFICIENCY



THE INTERNATIONAL WEEKLY JOURNAL DF SCIENCE

THE HIV-1

CAPSID

Atomic structure of the AIDS pathogen's protein coat PAGE 643

ACCELERATING DISCOVERIES

USING A SUPERCOMPUTER POWERED BY 3,000 TESLA PROCESSORS, UNIVERSITY OF ILLINOIS SCIENTISTS PERFORMED THE FIRST ALL-ATOM SIMULATION OF THE HIV VIRUS AND DISCOVERED THE CHEMICAL STRUCTURE OF ITS CAPSID — "THE PERFECT TARGET FOR FIGHTING THE INFECTION."

WITHOUT GPU, THE SUPERCOMPUTER WOULD NEED TO BE 5X LARGER FOR SIMILAR PERFORMANCE.

TESLA ACCELERATED COMPUTING PLATFORM

Data Center Infrastructure





Enterprise Services Support & Maintenance

ACCELERATED COMPUTING BASICS

WHAT IS ACCELERATED COMPUTING?

Application Execution



SIMPLICITY & PERFORMANCE

Simplicity

Accelerated Libraries

- Little or no code change for standard libraries; high performance
- Limited by what libraries are available

Compiler Directives

- High Level: Based on existing languages; simple and familiar
- High Level: Performance may not be optimal
- Parallel Language Extensions
 - Expose low-level details for maximum performance
 - > Often more difficult to learn and more time consuming to implement

Performance

CODE FOR SIMPLICITY & PERFORMANCE

Libraries

• Implement as much as possible using portable libraries.

Directives

• Use directives to rapidly accelerate your code.

Languages

• Use lower level languages for important kernels.

WHAT ARE COMPILER DIRECTIVES?

WHAT ARE COMPILER DIRECTIVES?



Your original Fortran, C, or C++ code

- Insert portable compiler directives
- Compiler parallelizes code and manages data movement
- Programmer optimizes incrementally
- Designed for multi-core CPUs, GPUs & many-core Accelerators

OPENACC: THE STANDARD FOR GPU DIRECTIVES

Simple: Easy path to accelerate compute intensive applications

Open: Open standard that can be implemented anywhere

Portable: Represents parallelism at a high level making it portable to any architecture



OPENACC MEMBERS AND PARTNERS



ACCELERATING APPLICATIONS WITH OPENACC



EXAMPLE: JACOBI ITERATION

Iteratively converges to correct value (e.g. Temperature), by computing new values at each point from the average of neighboring points.



JACOBI ITERATION: C CODE

while (err > tol && iter < iter max) {</pre> err=0.0;

```
for( int j = 1; j < n-1; j++) {
  for(int i = 1; i < m-1; i++) {</pre>
```

```
Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                     A[j-1][i] + A[j+1][i]);
```

err = max(err, abs(Anew[j][i] - A[j][i]));

```
for( int j = 1; j < n-1; j++) {</pre>
  for( int i = 1; i < m-1; i++ ) {</pre>
    A[j][i] = Anew[j][i];
```

iter++;



Compute max error for convergence

Swap input/output arrays



Calculate new value from neighbors



Iterate across matrix

elements

19



IDENTIFY PARALLELISM

while (err > tol && iter < iter max) {</pre> err=0.0;

```
for( int j = 1; j < n-1; j++) {
  for(int i = 1; i < m-1; i++) {</pre>
```

```
Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                     A[j-1][i] + A[j+1][i]);
```

```
err = max(err, abs(Anew[j][i] - A[j][i]));
```

Data dependency between iterations.



Independent loop iterations





```
for( int j = 1; j < n-1; j++) {</pre>
  for( int i = 1; i < m-1; i++ ) {</pre>
    A[j][i] = Anew[j][i];
```

iter++;



OPENACC DIRECTIVE SYNTAX

⊳ C/C++

#pragma acc directive [clause [,] clause] ...]
...often followed by a structured code block

Fortran

!\$acc directive [clause [,] clause] ...]

...often paired with a matching end directive surrounding a structured code block:

!\$acc end directive



OPENACC PARALLEL LOOP DIRECTIVE

parallel - Programmer identifies a block of code containing parallelism. Compiler generates a *kernel*.

1000 - Programmer identifies a loop that can be parallelized within the kernel.

NOTE: parallel & loop are often placed together

```
#pragma acc parallel loop
```

```
for(int i=0; i<N; i++)</pre>
```

{

```
y[i] = a*x[i]+y[i];
```

Kernel: A function that runs in parallel on the GPU

PARALLELIZE WITH OPENACC

while (err > tol && iter < iter_max) {
 err=0.0;</pre>

```
#pragma acc parallel loop reduction(max:err)
for( int j = 1; j < n-1; j++) {
   for(int i = 1; i < m-1; i++) {</pre>
```

```
Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] + A[j-1][i] + A[j+1][i]);
```

```
err = max(err, abs(Anew[j][i] - A[j][i]));
```



Parallelize loop on accelerator

Parallelize loop on

accelerator

```
#pragma acc parallel loop
for( int j = 1; j < n-1; j++) {
   for( int i = 1; i < m-1; i++ ) {
      A[j][i] = Anew[j][i];
   }
}</pre>
```

iter++;

* A *reduction* means that all of the N*M values for err will be reduced to just one, the max.

BUILDING THE CODE

\$ pgcc -fast -acc -ta=tesla -Minfo=all laplace2d.c
main:

- 40, Loop not fused: function call before adjacent loop Generated vector sse code for the loop
- 51, Loop not vectorized/parallelized: potential early exits
- 55, Accelerator kernel generated
 - 55, Max reduction generated for error
 - 56, #pragma acc loop gang /* blockIdx.x */
 - 58, #pragma acc loop vector(256) /* threadIdx.x */
- 55, Generating copyout (Anew[1:4094][1:4094]) Generating copyin(A[:][:]) Generating Tesla code
- 58, Loop is parallelizable
- 66, Accelerator kernel generated
 - 67, #pragma acc loop gang /* blockIdx.x */
 - 69, #pragma acc loop vector(256) /* threadIdx.x */
- 66, Generating copyin(Anew[1:4094][1:4094]) Generating copyout(A[1:4094][1:4094]) Generating Tesla code
- 69, Loop is parallelizable

OPENACC KERNELS DIRECTIVE

The kernels construct expresses that a region *may contain parallelism* and *the compiler determines* what can safely be parallelized.



PARALLELIZE WITH OPENACC KERNELS

while (err > tol && iter < iter_max) {
 err=0.0;</pre>

```
#pragma acc kernels
 for( int j = 1; j < n-1; j++) {
    for(int i = 1; i < m-1; i++) {</pre>
     Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                            A[j-1][i] + A[j+1][i]);
      err = max(err, abs(Anew[j][i] - A[j][i]));
  for( int j = 1; j < n-1; j++) {
    for( int i = 1; i < m-1; i++ ) {</pre>
     A[j][i] = Anew[j][i];
    }
  iter++;
```



Look for parallelism within this region.

BUILDING THE CODE



OPENACC PARALLEL LOOP VS. KERNELS

PARALLEL LOOP

- Requires analysis by programmer to ensure safe parallelism
- Will parallelize what a compiler may miss
- Straightforward path from OpenMP

KERNELS

- Compiler performs parallel analysis and parallelizes what it believes safe
- Can cover larger area of code with single directive
- Gives compiler additional leeway to optimize.

Both approaches are equally valid and can perform equally well.

🔍 NVIDIA Visual Profiler				_ D X	
File View Run Help					
🔍 *NewSession1 🖾					
	s	25 s 50 s 75 s	i	100 s	$J = \chi = J$
Process "a.out" (11805)					$\wedge \forall h$
Thread 3991209728					× X
L Driver API					X = -7 X
Profiling Overhead					
🖻 [0] Tesla K20c	JII.				
Context 1 (CUDA)					X.
🗆 🍸 MemCpy (HtoD)					
🗆 🍸 MemCpy (DtoH)					
Compute					
└ 🍸 59.5% main_61_gpu					
└ 🍸 35.3% main_72_gpu					
L 🍸 5.2% main_65_gpu_red				Ver	/ low
Streams				Compute	Momeny
analysis 🕸 🗔 Details 🗳 Console 🗔 Settings					
E Export PDF Report Results ratio					tio
1. CUDA Application Analysis		Low Compute / Memcpy Efficiency [5.073 s / 62.219 s = 0.082]			
2. Check Overall GPU Usage					
		тетсру.		More	
indicate potential problems in how your application is taking advantage of the GPU's available compute and data movement capabilities. You should examine the information provided with each result to		▲ Low Memcpy/Compute Overlap [0 ns / 5.073 s = 0%]			
	Ξ	The percentage of time when memcpy is being performed in parallel	mpute		5.0s
		Low Kernel Concurrency [0 ns / 5.073 s = 0%]	mony Co		67 75
		The percentage of time when two kernels are being executed in paral		Ру	02.25
determine if you can make changes		Low Memcpy Throughput [83.659 MB/s avg, for memcpys accounting for 0% of all memcpy time]			
to your application to increase GPU	<u> </u>	- Low memory initiagiput [05.055 mb/s avg, for memorys accounting to	i ovo or an memo	by time 1	

EXCESSIVE DATA TRANSFERS

IDENTIFYING DATA LOCALITY

```
while ( err > tol && iter < iter_max ) {
    err=0.0;</pre>
```

```
#pragma acc parallel loop
for( int j = 1; j < n-1; j++) {
   for( int i = 1; i < m-1; i++ ) {
      A[j][i] = Anew[j][i];
   }
}</pre>
```

Does the CPU need the data between iterations of the convergence loop?

```
iter++;
```


DEFINING DATA REGIONS

The data construct defines a region of code in which GPU arrays remain on the GPU and are shared among all kernels in that region.

#pragma acc data

#pragma acc parallel loop

#pragma acc parallel loop

Data Region

Arrays used within the data region will remain on the GPU until the end of the data region.

DATA CLAUSES

COPY (**list**) Allocates memory on GPU and copies data from host to GPU when entering region and copies data to the host when exiting region.

copyin (list) Allocates memory on GPU and copies data from host to GPU when entering region.

copyout (list) Allocates memory on GPU and copies data to the host when exiting region.

create (**list**) Allocates memory on GPU but does not copy.

present (*list*) Data is already present on GPU from another containing data region.

and present_or_copy[in|out], present_or_create, deviceptr.

ARRAY SHAPING

Compiler sometimes cannot determine size of arrays

Must specify explicitly using data clauses and array "shape"

<u>C/C++</u>

#pragma acc data copyin(a[0:size]) copyout(b[s/4:3*s/4])

Fortran

!\$acc data copyin(a(1:end)) copyout(b(s/4:3*s/4))

Note: data clauses can be used on data, parallel, or kernels

OPTIMIZE DATA LOCALITY

```
#pragma acc data copy(A) create(Anew)
while ( err > tol && iter < iter_max ) {
    err=0.0;
#pragma acc parallel loop
    for( int j = 1; j < n-1; j++) {
        for(int i = 1; i < m-1; i++) {
    }
}</pre>
```

```
Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] + A[j-1][i] + A[j+1][i]);
```

```
err = max(err, abs(Anew[j][i] - A[j][i]));
}
```

Copy A to/from the accelerator only when needed.

Create Anew as a device temporary.

```
#pragma acc parallel loop
  for( int j = 1; j < n-1; j++) {
    for( int i = 1; i < m-1; i++ ) {
        A[j][i] = Anew[j][i];
        }
    }
    iter++;
}</pre>
```

REBUILDING THE CODE

\$ pgcc -fast -acc -ta=tesla -Minfo=all laplace2d.c main: 40, Loop not fused: function call before adjacent loop Generated vector sse code for the loop 51, Generating copy(A[:][:]) Generating create (Anew[:][:]) Loop not vectorized/parallelized: potential early exits 56, Accelerator kernel generated 56, Max reduction generated for error 57, #pragma acc loop gang /* blockIdx.x */ 59, #pragma acc loop vector(256) /* threadIdx.x */ 56, Generating Tesla code 59, Loop is parallelizable 67, Accelerator kernel generated 68, #pragma acc loop gang /* blockIdx.x */ 70, #pragma acc loop vector(256) /* threadIdx.x */ 67, Generating Tesla code 70, Loop is parallelizable

VISUAL PROFILER: DATA REGION

Speed-Up (Higher is Better)

OPENACC PRESENT CLAUSE

It's sometimes necessary for a data region to be in a different scope than the compute region.

When this occurs, the **present** clause can be used to tell the compiler data is already on the device.

Since the declaration of A is now in a higher scope, it's necessary to shape A in the present clause.

High-level data regions and the present clause are often critical to good performance.


```
function laplace2D(double[N][M] A,n,m)
{
    #pragma acc data present(A[n][m]) create(Anew)
    while ( err > tol && iter < iter_max ) {
        err=0.0;
        ...
    }
}</pre>
```


NEXT STEPS

ACCELERATING WITH OPENACC

- 1. Identify Available Parallelism
 - What important parts of the code have available parallelism?
- 2. Parallelize Loops
 - Express as much parallelism as possible and ensure you still get correct results.
 - Because the compiler *must* be cautious about data movement, the code will generally slow down.
- 3. Optimize Data Locality
 - The programmer will *always* know better than the compiler what data movement is unnecessary.
- 4. Optimize Loop Performance
 - Don't try to optimize a kernel that runs in a few us or ms until you've eliminated the excess data motion that is taking many seconds.

TYPICAL PORTING EXPERIENCE WITH OPENACC DIRECTIVES

Application Speed-up

FOR MORE INFORMATION

- Check out <u>http://openacc.org/</u>
- Watch tutorials at <u>http://www.gputechconf.com/</u>
- Share your successes at WACCPD at SC15. <u>http://www.openacc.org/content/Events/waccpd_2015</u>
- Take a self-paced lab at <u>https://nvidia.qwiklab.com/</u>