OPENACC DATA MANAGEMENT





CPU AND GPU MEMORIES



CPU + GPU Physical Diagram

- CPU memory is larger, GPU memory has more bandwidth
- CPU and GPU memory are usually separate, connected by an I/O bus (traditionally PCI-e)
- Any data transferred between the CPU and GPU will be handled by the I/O Bus
- The I/O Bus is relatively slow compared to memory bandwidth
- The GPU cannot perform computation until the data is within its memory





CUDA UNIFIED MEMORY



CUDA UNIFIED MEMORY

Simplified Developer Effort

Without Managed Memory

Commonly referred to as *"managed memory."*





CPU and GPU memories are combined into a single, shared pool

Managed Memory

This material is released by NVIDIA Corporation under the Creative Commons Attribution 4.0 International (CC BY 4.0)

CUDA MANAGED MEMORY Usefulness

- Handling explicit data transfers between the host and device (CPU and GPU) can be difficult
- The NVIDIA HPC compiler can utilize CUDA Managed Memory to defer data management
- This allows the developer to concentrate on parallelism and think about data movement as an optimization

nvc -fast -acc -ta=tesla:managed -Minfo=accel main.c

\$

\$

nvfortran -fast -acc -ta=tesla:managed -Minfo=accel main.f90



MANAGED MEMORY Limitations

- The programmer will almost always be able to get better performance by manually handling data transfers
- Memory allocation/deallocation takes longer with managed memory
- Cannot transfer data asynchronously
- Currently only available on NVIDIA GPUs with NVIDIA HPC SDK.





Managed Memory



DATA SHAPING



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.

Principal use: For many important data structures in your code, this is a logical default to input, modify and return the data.

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

Principal use: Think of this like an array that you would use as just an input to a subroutine.

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

Principal use: A result that isn't overwriting the input data structure.

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

OpenACC

Principal use: Temporary arrays.

ARRAY SHAPING

- Sometimes the compiler needs help understanding the shape of an array
- The first number is the start index of the array
- In C/C++, the second number is how much data is to be transferred
- In Fortran, the second number is the ending index





OPTIMIZED DATA MOVEMENT

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

```
#pragma acc parallel loop reduction(max:err) copyin(A[0:n*m]) copy(Anew[0:n*m])
       for( int j = 1; j < n-1; j++) {
         for(int i = 1; i < m-1; i++) {
                                                                   Data clauses
           Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                                A[j-1][i] + A[j+1][i]);
                                                               provide necessary
                                                                  "shape" to the
           err = max(err, abs(Anew[j][i] - A[j][i]));
                                                                       arrays.
       }
     #pragma acc parallel loop copyin(Anew[0:n*m]) copyout(A[0:n*m])
       for( int j = 1; j < n-1; j++) {
         for( int i = 1; i < m-1; i++ ) {</pre>
           A[j][i] = Anew[j][i];
       iter++;
         💿 nvidia,
OpenACC
```

OPENACC SPEED-UP SLOWDOWN

45.00X			
40.00X			
2 5.00X ——			
20.00X			
0.00X	1.00X		0.33X
			V100 (DATA CLAUSES)

This material is released by NVIDIA Corporation under the Creative Commons Attribution 4.0 International (CC BY 4.0)



RUNTIME BREAKDOWN



OPTIMIZED DATA MOVEMENT

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

```
#pragma acc parallel loop reduction(max:err) copyin(A[0:n*m]) copy(Anew[0:n*m])
       for( int j = 1; j < n-1; j++) {
         for(int i = 1; i < m-1; i++) {
                                                                Currently we're
          Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                               A[j-1][i] + A[j+1][i]);
                                                              copying to/from the
                                                              GPU for each loop,
          err = max(err, abs(Anew[j][i] - A[j][i]));
         }
                                                                can we reuse it?
       }
     #pragma acc parallel loop copyin(Anew[0:n*m]) copyout(A[0:n*m])
       for( int j = 1; j < n-1; j++) {
         for( int i = 1; i < m-1; i++ ) {</pre>
          A[j][i] = Anew[j][i];
       iter++;
OpenACC
        💿 nvidia,
```

OPTIMIZE DATA MOVEMENT



OPENACC DATA DIRECTIVE Definition

- The data directive defines a lifetime for data on the device beyond individual loops
- During the region data is essentially "owned by" the accelerator
- Data clauses express shape and data movement for the region

#pragma acc data <i>clauses</i> {				
< Sequential and/or Parallel code >				
}				
!\$acc data <i>cLauses</i>				

Sequential and/or Parallel code >

!\$acc end data



OPTIMIZED DATA MOVEMENT

```
#pragma acc data copy(A[:n*m]) copyin(Anew[:n*m])
while ( err > tol && iter < iter_max ) {
    err=0.0;</pre>
```

```
#pragma acc parallel loop reduction(max:err) copyin(A[0:n*m]
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] + A[j+1][i]);
```

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

}



Copy A to/from the accelerator only when needed.

Copy initial condition of Anew, but not final value

REBUILD THE CODE

nvc -fast -ta=tesla -Minfo=accel laplace2d_uvm.c
main:

- 60, Generating copy(A[:m*n])
 Generating copyin(Anew[:m*n])
- 64, Accelerator kernel generated Generating Tesla code
 - 64, Generating reduction(max:error)
 - 65, #pragma acc loop gang /* blockIdx.x */
 - 67, #pragma acc loop vector(128) /* threadIdx.x */
- 67, Loop is parallelizable
- 75, Accelerator kernel generated

Generating Tesla code

- 76, #pragma acc loop gang /* blockIdx.x */
- 78, #pragma acc loop vector(128) /* threadIdx.x */
- 78, Loop is parallelizable



Now data movement only happens at our data region.



WHAT WE'VE LEARNED SO FAR

- CUDA Unified (Managed) Memory is a powerful porting tool
- GPU programming without managed memory often requires data shaping
- Moving data at each loop is often inefficient
- The OpenACC Data region can decouple data movement and computation



DATA SYNCHRONIZATION



OPENACC UPDATE DIRECTIVE

update: Explicitly transfers data between the host and the device

Useful when you want to synchronize data in the middle of a data region Clauses:

self: makes host data agree with device data

device: makes device data agree with host data

<pre>#pragma acc update #pragma acc update</pre>	<pre>self(x[0:count]) device(x[0:count])</pre>			
	C/C++			
<pre>!\$acc update self(x(1:end_index)) !\$acc update device(x(1:end_index)) Fortran</pre>				



OPENACC RESOURCES

Guides • Talks • Tutorials • Videos • Books • Spec • Code Samples • Teaching Materials • Events • Success Stories • Courses • Slack • Stack Overflow



💿 nvidia.

OpenACC

Resources

Success Stories

About Tools News Events Re





<section-header><section-header><section-header><section-header><section-header><section-header>