

Using HIP and GPU Libraries with OpenMP

OLCF Preparing for Frontier Training Series

Reuben D. Budiardja
Swaroop Pophale
Wael Elwasif
Suzanne Parete-Koon

Oak Ridge Leadership Computing Facility
Oak Ridge National Laboratory

ORNL is managed by UT-Battelle LLC for the US Department of Energy



U.S. DEPARTMENT OF
ENERGY

The Question

GPU libraries and HIP kernels often expect **pointers to GPU memory addresses** as arguments.

How do we get those from OpenMP?

Introducing “use_device_addr” clause

Potentially common use case for this clause is to call GPU libraries / routines from host.

```
use_device_addr ( var1, [ var2, ... ] )
```

[...] references to the list item in [...] are converted into references to the corresponding storage [in the device data environment].

- item should already be *mapped*
 - if not mapped, assumed to be accessible from the device (i.e. GPU)
 - For C/C++, since dynamically allocated array is likely a pointer, use `use_device_ptr` instead
 - In OpenMP 4.5, only `use_device_ptr` exists
- use this with compiler supporting OpenMP 4.5 (e.g. IBM XL on Summit)

use_device_addr / use_device_ptr simple example

```
1  call print_addr ( A )  !-- print host addr
2
3  !$OMP target enter data map ( to: A )
4
5  !$OMP target data use_device_addr ( A )
6  call print_addr ( A )  !-- print device addr
7  !$OMP end target data
8
9  call print_addr ( A )  !-- print host addr
10 |
```

```
1  float *A = (float *) malloc ( sizeof ( float ) * 10 );
2
3  printf("%p\n", A); // print A device addr
4
5  #pragma omp target enter data map ( to: A[0:10] )
6
7  #pragma omp target data use_device_ptr ( A )
8  { printf("%p\n", A); } // print A device addr
9
10 printf("%p\n", A); // print A host addr
11
12
```

Repository containing examples in this tutorial:

<https://github.com/olcf/openmp-gpu-library>

Illustrative Fortran Example - Main program

```
1  program MatrixMultiply_Example
2
3  !-- Initialize matrices
4  allocate ( Matrix_A ( N_VALUES, N_VALUES ), &
5             Matrix_B ( N_VALUES, N_VALUES ), &
6             Matrix_C ( N_VALUES, N_VALUES ) )
7
8  call random_number ( Matrix_A )
9  call random_number ( Matrix_B )
10 Matrix_C = huge ( 1.0_real64 ) !-- make error more obvious
11
12 !$OMP target enter data map ( to: Matrix_A, Matrix_B, Matrix_C )
13
14 Matrix_C_Ref = matmul ( Matrix_A, Matrix_B )
15
16 call MatrixMultiply_OpenMP ( Matrix_A, Matrix_B, Matrix_C )
17 !$OMP target update from ( Matrix_C )
18
19 !-- [Verification and resetting of Matrix_C]
20
21
22 call MatrixMultiply_HIP ( Matrix_A, Matrix_B, Matrix_C )
23 !$OMP target update from ( Matrix_C )
24
25 !-- [Verification and resetting of Matrix_C]
26
27
28 call MatrixMultiply_GPU_Library ( Matrix_A, Matrix_B, Matrix_C )
29 !$OMP target update from ( Matrix_C )
30
31 !-- [Verification and resetting of Matrix_C]
32
33 end program MatrixMultiply_Example
```

Initializations on host

Mapping to device; OpenMP runtime creates corresponding data storage on device

Multiple ways to do matrix multiplication:

- Using Fortran intrinsic matmul() on host
- Using subroutine with OpenMP offload
- **Using subroutine with HIP kernel**
- **Using subroutine with GPU library (hipblas)**

OpenMP Offload Subroutine

```
1  module MM_OpenMP
2
3  use iso_fortran_env
4
5  implicit none
6  private
7
8  public :: &
9      MatrixMultiply_OpenMP
10
11 contains
12
13     subroutine MatrixMultiply_OpenMP ( A, B, C )
14
15         real ( real64 ), dimension ( :, : ), intent ( in ) :: &
16             A, B
17         real ( real64 ), dimension ( :, : ), intent ( out ) :: &
18             C
19
20         integer :: &
21             iV, jV, kV
22
23         !$OMP target teams distribute parallel do collapse ( 2 )
24         do iV = 1, size ( A, dim = 1 )
25             do jV = 1, size ( B, dim = 2 )
26                 C ( iV, jV ) = 0.0_real64
27                 do kV = 1, size ( A, dim = 2 )
28                     C ( iV, jV ) = C ( iV, jV ) + A ( iV, kV ) * B ( kV, jV )
29                 end do
30             end do
31         end do
32
33     end subroutine MatrixMultiply_OpenMP
34
35 end module MM_OpenMP
36
```

Computation if offloaded to GPU;
No data movement because of data mappings
are already done previously in the main
program.

Using HIP Kernel - Fortran Calling Subroutine

```
1  #ifdef POWER_XL
2  #define use_device_addr use_device_ptr
3  #endif
4
5  module MM_HIP
6
7      use iso_fortran_env
8      use iso_c_binding
9
10     implicit none
11
12     interface
13
14         subroutine MatrixMultiply_HIP_Launch ( A, B, C, nValues ) &
15             bind ( C, name = 'MatrixMultiply_HIP_Launch' )
16             use iso_c_binding
17             implicit none
18             type ( c_ptr ), value :: &
19                 A, B, C
20             integer ( c_int ), value :: &
21                 nValues
22         end subroutine MatrixMultiply_HIP_Launch
23
24     end interface
25
26 contains
27
28     subroutine MatrixMultiply_HIP ( A, B, C )
29
30         real ( real64 ), dimension ( :, : ), intent ( in ), target :: &
31             A, B
32         real ( real64 ), dimension ( :, : ), intent ( inout ), target :: &
33             C
34
35         !$OMP target data use_device_addr ( A, B, C )
36         call MatrixMultiply_HIP_Launch &
37             ( c_loc ( A ), c_loc ( B ), c_loc ( C ), size ( A, dim = 1 ) )
38         !$OMP end target data
39
40     end subroutine MatrixMultiply_HIP
41
42 end module MM_HIP
```

Preprocessor replaces “use_device_addr” to “use_device_ptr” for IBM XL compiler on Summit since it provides OpenMP 4.5.

Interface for the C helper subroutine (i.e. function returning void) using the Fortran - C interoperability standard.
c_ptr points to GPU memory addresses.

use_device_addr tells OpenMP to use the corresponding addresses on device. c_loc() returns the c_ptr type compatible to void * in C.

Using HIP Kernel - C & HIP Subroutines

```
1 // Include and headers ellided
2
3 #ifdef __cplusplus
4 extern "C"
5 {
6 #endif
7 void MatrixMultiply_HIP_Launch
8     ( double *d_A, double *d_B, double *d_C, int nV );
9 #ifdef __cplusplus
10 }
11 #endif
12
13 __global__ void matrix_multiply(double *a, double *b, double *c, int N)
14 {
15     int row    = blockDim.x * blockIdx.x + threadIdx.x;
16     int column = blockDim.y * blockIdx.y + threadIdx.y;
17     if (row < N && column < N)
18     {
19         double element = 0.0;
20         for(int i=0; i<N; i++){
21             element += a[row + N * i] * b[i + N * column];
22         }
23         c[row + N * column] = element;
24     }
25 }
26
27 void MatrixMultiply_HIP_Launch ( double *d_A, double *d_B, double *d_C, int nV )
28 {
29     dim3 threads_per_block ( 16, 16, 1 );
30     dim3 blocks_in_grid ( ceil( float(nV) / threads_per_block.x ),
31                           ceil( float(nV) / threads_per_block.y ), 1 );
32
33     // Launch kernel
34     double start = omp_get_wtime ( );
35     hipLaunchKernelGGL ( matrix_multiply, blocks_in_grid, threads_per_block, 0,
36                         0, d_A, d_B, d_C, nV );
37     hipDeviceSynchronize ( );
38 }
39
```

Create un-mangled C symbol available for linkage with the corresponding Fortran interface.

HIP / CUDA kernel for matrix multiplication, expecting the pointers to GPU memory addresses as argument.

C helper function to be called from Fortran.

Using GPU Libraries (hipblas / cublas) - Fortran Interfaces

```
1  #ifndef POWER_XL
2  #define use_device_addr use_device_ptr
3  #define __mm_lib_create__ 'cublasCreate_v2'
4  #define __mm_lib_dgemm__ 'cublasDgemm_v2'
5  #define __mm_lib_destroy__ 'cublasDestroy_v2'
6  ...
7  #endif
8
9  #ifdef Cray_CCE
10 #define __mm_lib_create__ 'hipblasCreate'
11 #define __mm_lib_dgemm__ 'hipblasDgemm'
12 #define __mm_lib_destroy__ 'hipblasDestroy'
13 ...
14 #endif
15
16 module MM_GPU_Library
17   use iso_fortran_env
18   use iso_c_binding
19
20   implicit none
21
22   enum, bind ( c ) !-- From cublas_api.h
23     enumerator :: BLAS_OP_N = __blas_op_n__
24     ...
25   end enum
26
27   interface
28
29     function MM_LibCreate ( Handle ) &
30       bind ( C, name = __mm_lib_create__ ) result ( Status )
31       use iso_c_binding
32       implicit none
33       type ( c_ptr ) :: &
34         Handle
35       ...
36     end function MM_LibCreate
37
38     function MM_LibDgemm ( Handle, Trans_A, Trans_B, M, N, K, Alpha, &
39       d_A, LDA, d_B, LDB, Beta, d_C, LDC ) &
40       bind ( c, name = __mm_lib_dgemm__ ) result ( Status )
41       use iso_c_binding
42       implicit none
43       type ( c_ptr ), value :: &
44         Handle
45       ...
46     end function MM_LibDgemm
47
```

Preprocessor to switch between HIP or CUDA version based on defined macros (set in Makefile)

Note: On Frontier / Crusher, these Fortran interfaces are provided by hipfort:

- module load hipfort on terminal
- use hipfort in Fortran

<https://github.com/ROCmSoftwarePlatform/hipfort>

Interfaces to the hipblas / cublas functions; name replacements are done via preprocessor above.

Using GPU Libraries (hipblas / cublas) - Fortran Call

```
1  subroutine MatrixMultiply_GPU_Library ( A, B, C )
2
3     real ( real64 ), dimension ( :, : ), intent ( in ), target :: &
4     A, B
5     real ( real64 ), dimension ( :, : ), intent ( inout ), target :: &
6     C
7
8     integer ( c_int ) :: &
9     nV, &
10    Status
11    real ( c_double ) :: &
12    Alpha = 1.0_c_double, &
13    Beta = 0.0_c_double
14    type ( c_ptr ) :: &
15    Handle
16
17    Status = MM_LibCreate ( Handle )
18
19    nV = size ( A, dim = 1 )
20    !$OMP target data use_device_addr ( A, B, C )
21    Status = MM_LibDgemm ( Handle, BLAS_OP_N, BLAS_OP_N, nV, nV, nV, &
22                        Alpha, c_loc ( A ), nV, c_loc ( B ), nV, &
23                        Beta, c_loc ( C ), nV )
24    !$OMP end target data
25
26    Status = MM_LibDestroy ( Handle )
27
28    end subroutine MatrixMultiply_GPU_Library
29
```

Calling hipblas / cublas

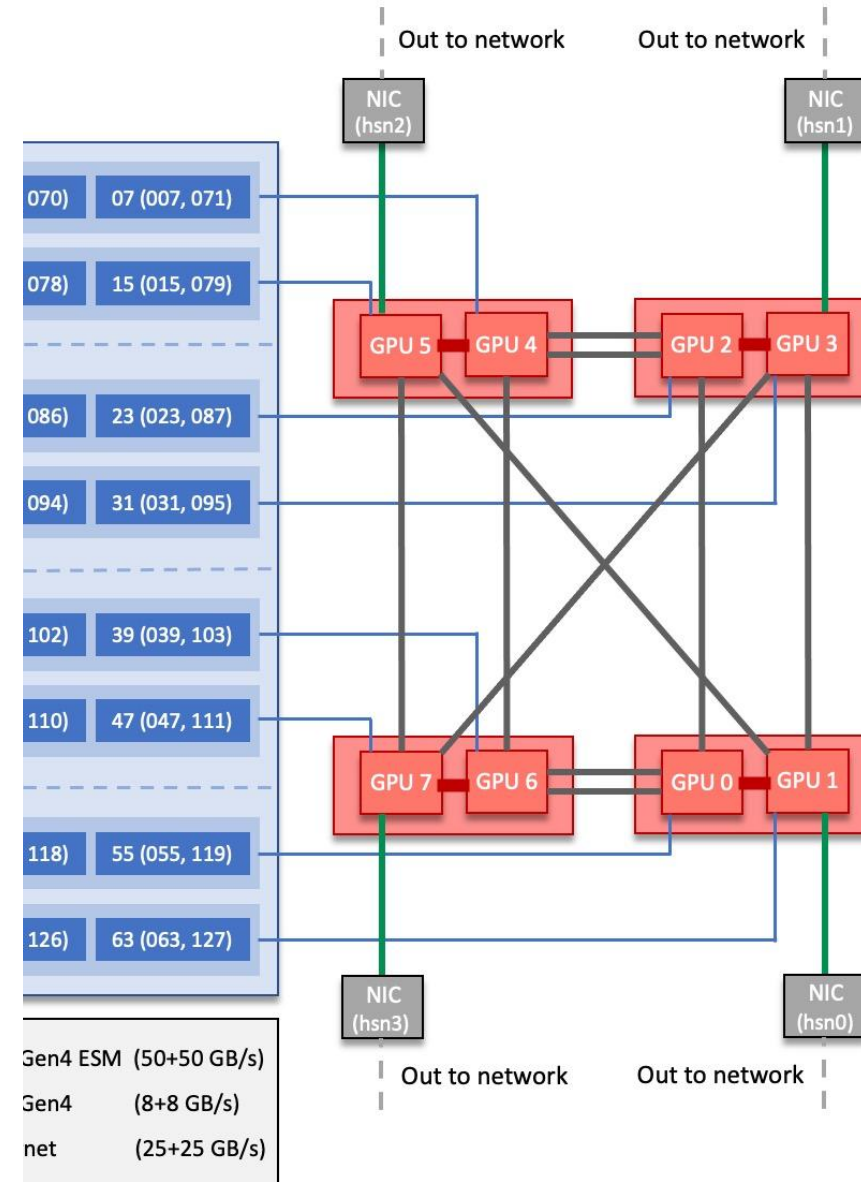
GPU-Aware MPI

GPU-aware MPI library allows GPU memory addresses to be used as arguments to MPI subroutine calls:

- avoids staging (i.e. copying) data to host first
- allows optimized / faster GPU-to-GPU communications

In the next few slides we look at examples on how to do this with OpenMP-managed data.

Frontier Node Diagram (partial)



Calling GPU-aware MPI routine

```
1  program MPI_Send_Recv_Example
2
3  use mpi
4
5  implicit none
6
7  real ( real64 ), allocatable, dimension ( : ) :: &
8  Buffer
9  ! [other type definitions]
10
11  ...
12
13  !$OMP target data use_device_addr ( Buffer )
14  call MPI_SEND &
15       ( Buffer, nV, MPI_DOUBLE_PRECISION, 1, 0, MPI_COMM_WORLD, &
16         Error )
17  !$OMP end target data
18
19  ...|
20
21  end program MPI_Send_Recv_Example
22
```

Bonus: “omp_target_associate_ptr” routine*

```
omp_target_associate_ptr(host_ptr, device_ptr, size,  
device_offset, device_num)
```

... maps a device pointer to a host pointer.

This routine tells OpenMP to use the `device_ptr` when `host_ptr` appear in target region and subsequent map clause.

- `device_ptr` may be returned from runtime routine (e.g. `omp_target_alloc()`, `hipMalloc()` / `cudaMalloc()`)
- OpenMP will not create its own device storage for the `host_ptr` (i.e. host variables)

**Note: `omp_target_associate_ptr()` is not needed for GPU-aware MPI, but useful to know.*

Calling GPU-aware MPI routine (2)

```
1 program MPI_Send_Recv_Example
2   use mpi
3   use omp_lib
4   use iso_c_binding
5   use iso_fortran_env
6   implicit none
7
8   real ( real64 ), allocatable, dimension ( : ), target :: &
9     Buffer
10  real ( real64 ), dimension ( : ), pointer :: &
11    F_Buffer
12  type ( c_ptr ) :: &
13    D_Buffer
14  !-- [other type definitions]
15
16  allocate ( Buffer ( nV ) )
17
18  !-- allocate memory on GPU, then associate it with host variable
19  associate ( Size => c_sizeof ( 1.0_c_double ) * nV )
20    D_Buffer = omp_target_alloc ( c_sizeof ( 1.0_c_double ) * nV, 0 )
21    Error = omp_target_associate_ptr ( c_loc ( Buffer ), D_Buffer, Size, &
22      0_c_size_t, 0 )
23  end associate
24
25  !-- [ MPI Initializations and set values]
26
27  if ( MPI_Rank == 0 ) then
28
29    !$OMP target
30    Buffer = 2.0_Real64 !-- Set value on device
31    !$OMP end target
32
33    call c_f_pointer ( D_Buffer, F_Buffer, [ nV ] )
34    call MPI_SEND &
35      ( F_Buffer, nV, MPI_DOUBLE_PRECISION, 1, 0, MPI_COMM_WORLD, &
36        Error )
37  end if
38
```

Buffer: host variable

D_Fuffer: device_ptr to a GPU address

F_Buffer: Fortran pointer on the host

Buffer is allocated; D_Buffer points to GPU-memory allocated to the same size. Both are associated with `omp_target_associate_ptr()`

Values are set on addressed pointed by D_Buffer

D_Buffer is dereferenced to F_Buffer, which is now a rank-2 real arrays backed by GPU-memory addressed. F_Buffer can be passed to GPU-aware MPI routines.

Building and Running Matrix Multiplication Examples

On Summit

```
module load xl
module load cuda
make clean
make MACHINE=POWER_XL
```

```
jsrun -n 1 -g 1 -c 1 \
./MatrixMultiply_POWER_XL
```

```
OpenMP Verification : PASSED
HIP Verification : PASSED
GPU Library Verification : PASSED
```

On Crusher / Frontier

```
module load PrgEnv-cray
module load craype-accel-amd-gfx90a
module load rocm
make clean
make MACHINE=Cray_CCE
```

```
srun -n1 ./MatrixMultiply_Cray_CCE
```

```
OpenMP Verification : PASSED
HIP Verification : PASSED
GPU Library Verification : PASSED
```

Repository containing examples in this tutorial:

<https://github.com/olcf/openmp-gpu-library>

Conclusion

- We learned how to use GPU libraries with OpenMP code
 - use_device_addr clause to get the corresponding device address of a host variable
 - omp_target_associate_ptr to map a generic device pointer to a host variable for OpenMP runtime
- Complete examples code with README to run on Summit and Crusher / Frontier is available at [REPO INFO]
- Questions & comments: reubendb@ornl.gov