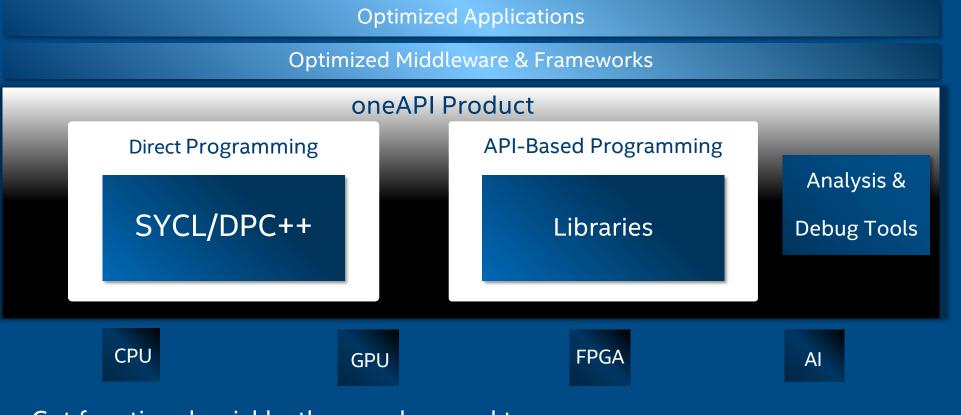
## Partition Communication and SYCL Bindings

Maria Garzaran

### A detour to SYCL

#### oneAPI for Cross-Architecture Performance



Get functional quickly, then analyze and tune

#### SYCL

- SYCL is an industry-driven Khronos standard that adds data parallelism to C++ for heterogeneous systems (https://www.khronos.org/sycl/)
- There are SYCL compilers that support different vendor's hardware:
- 1. Intel compiler generates code for Intel CPUs, GPUs, and FPGAs
- 2. ComputeCpp from CodePlay generates code for any CPU, Intel GPUs and FPGAs, AMD GPUs, ARM Mali, Nvidia GPUs
- 3. triSYCL from Xilinx generates code Xilinx FPGAs
- 4. hipSYCL from University of Heidelberg generates code for AMD GPUs and Nvidia GPUs.
- Libraries such as Kokkos and Raja have SYCL backends

#### Hello data-parallel SYCL code

```
#include <CL/sycl.hpp>
    #include <iostream>
    using namespace sycl;
    const std::string secret {
      "Ifmmp-!xpsme\"\012J(n!tpssz-!Ebwf/!"
      "J(n!bgsbje!J!dbo(u!ep!uibu/!.!IBM\01" };
    const auto sz = secret.size();
 9.
    int main() {
      queue Q;
11.
12.
      char *result = malloc shared<char>(sz, Q);
13.
14.
      std::memcpy(result,secret.data(),sz);
15.
      Q.parallel for(sz,[=](auto& i) { <
16.
        result[i] -= 1;
17.
        }).wait();
18.
19.
      std::cout << result << "\n":
20.
     return 0;
21.
22. }
```

- Work is submitted to a queue
- A queue is associated with a single device
- Work in the device is asynchronous
- Allocate memory on shared memory
- (both host and device can access)
- Submit work to the queue
- The work is expressed as a lambda function here
- Piece of work the kernel performs

#### Hello data-parallel SYCL code

```
#include <CL/sycl.hpp>
    #include <iostream>
                                                                        Where the code executes depends on
    using namespace sycl;
                                                                        where we setup the queue
 4.
    const std::string secret {
     "Ifmmp-!xpsme\"\012J(n!tpssz-!Ebwf/!"
     "J(n!bgsbje!J!dbo(u!ep!uibu/!.!IBM\01" };
                                                                        queue Q { default selector{} };
    const auto sz = secret.size();
                                                    Host code
                                                                        queue Q { host selector{} };
 9.
    int main() {
10.
                                                                        queue Q { cpu_selector{} };
     queue Q;
11.
                                                                        queue Q { gpu_selector{} };
12.
                                                                        queue Q { accelerator_selector{} };
     char *result = malloc shared<char>(sz, Q);
13.
                                                                        queue Q { INTEL::fpga_selector{} };
     std::memcpy(result,secret.data(),sz);
14.
15.
     Q.parallel_for(sz,[=](auto& i) {
16.
                                                    Device code
       result[i] -= 1;
17.
       }).wait();
18.
19.
     std::cout << result << "\n";
20.
                                                    Host code
21.
     return 0;
22. }
```

### Partition Communication

#### Example 4.2 MPI standard (OpenMP)

```
if (myrank == 1) {
if (myrank == 0) {
                                                           MPI_Precv_init(message, partitions, count, xfer_type
 MPI_Psend_init(message, partitions, count, xfer_type,
   dest, tag, info, MPI_COMM_WORLD, &request);
                                                             source, tag, info, MPI_COMM_WORLD, &request);
                                                            MPI Start(&request);
 MPI_Start(&request);
 #pragma omp parallel for shared(request)
num_threads(NUM_THREADS)
 for (int i = 0; i < partitions; i++) {
    /* compute and fill partition #i, then mark ready: */
    MPI_Pready(i, request);
 while (!flag) {
                                                           while (!flag) {
   /* Do useful work */
                                                            /* Do useful work */
   MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
                                                             MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
 MPI_Request_free(&request);
                                                           MPI_Request_free(&request);
```

#### Partition Communication Ex. 4.2 MPI Standard (SYCL)

```
<u>if (myrank == 1) {</u>
if (myrank == 0) {
                                                            MPI_Precv_init(message, partitions, count, xfer_type
 MPI_Psend_init(message, partitions, count, xfer_type,
   dest, tag, info, MPI_COMM_WORLD, &request);
                                                              source, tag, info, MPI_COMM_WORLD, &request);
                                                             MPI_Start(&request);
 MPI_Start(&request);
 mQueue.submit([&](sycl::handler &h) {
   h.parallel_for(partitions, [=](id<1> i) {
     MPI_Pready(i, request);
   });
 });
 mQueue.wait();
                                                           while (!flag) {
 while (!flag) {
                                                            /* Do useful work */
   /* Do useful work */
                                                              MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
   MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
                                                           MPI_Request_free(&request);
 MPI_Request_free(&request);
```

#### Partition Communication: MPI\_Pready executes from Device

```
if (myrank == 1) {
if (myrank == 0) {
 MPI Psend init(message, partitions, count, xfer type,
                                                                  MPI_Precv_init(message, partitions, count, xfer_type
   dest, tag, info, MPI_COMM_WORLD, &request);
                                                                    source, tag, info, MPI_COMM_WORLD, &request);
 MPI_Start(&request);
                                                                  MPI_Start(&request);
 MPI_Pbuf_prepare(&request);
                                                                  MPI_Pbuf_prepare(&request);
 mQueue.submit([&](sycl::handler &h) {
   h.parallel_for(partitions, [=](id<1> i) {
     MPI_Pready(i, request);
   });
 });
 mQueue.wait();
                                                                 while (!flag) {
 while (!flag) {
                                                                  /* Do useful work */
   /* Do useful work */
                                                                    MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
   MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
                                                                 MPI_Request_free(&request);
 MPI_Request_free(&request);
```

#### Partition Communication: MPI Prequest for device

```
if (myrank == 1) {
if (myrank == 0) {
 MPI Psend init(message, partitions, count, xfer type,
                                                                  MPI_Precv_init(message, partitions, count, xfer_type
   dest, tag, info, MPI_COMM_WORLD, &request);
                                                                    source, tag, info, MPI_COMM_WORLD, &request);
 MPI_Start(&request);
                                                                   MPI_Start(&request);
 MPI_Pbuf_prepare(&request);
                                                                   MPI_Pbuf_prepare(&request);
 // Create a request to pass to the kernel
 MPI_Prequest_create(request, info, &device_request);
 mQueue.submit([&](sycl::handler &h) {
   h.parallel_for(partitions, [=](id<1> i) {
     MPI_Pready(i, device_request);
   });
 });
mQueue.wait();
while (!flag) {
                                                                  while (!flag) {
   /* Do useful work */
                                                                  /* Do useful work */
   MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
                                                                    MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
MPI Request free(&request);
                                                                 MPI_Request_free(&request);
      1/19/22
```

#### Partition Communication: MPI\_Prequest for device (notes)

```
if (myrank == 0) {
 MPI Psend init(message, partitions, count, xfer type,
   dest, tag, info, MPI_COMM_WORLD, &request);
 MPI_Start(&request);
 MPI_Pbuf_prepare(&request);
 // Create a request to pass to the kernel
 MPI_Prequest_create(request, info, &device_request);
 mQueue.submit([&](sycl::handler &h) {
   h.parallel_for(partitions, [=](id<1> i) {
     MPI_Pready(i, device_request);
   });
 });
mQueue.wait();
while (!flag) {
   /* Do useful work */
   MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
MPI Request free(&request);
      1/19/22
```

```
if (myrank == 1) {
    MPI_Precv_init(message, partitions, count, xfer_type
        source, tag, info, MPI_COMM_WORLD, &request);
    MPI_Start(&request);
    MPI_Pbuf_prepare(&request);
```

- We are making the implicit assumption that the device request is allocated on the same device as the message
- We also assume that the message is only allocated in a single device
- Device information is not needed, because we assume the runtime is able to recover the device where message is allocated from the request created during the MPI\_Psend\_init() call
- An additional reason to have the device\_request is that in some devices (e.g. FPGAs), the kernel might not be able to access the host memory, so it would not be able to access the request in the host memory

```
MPI_Request_free(&request);
```

#### MPI\_Prequest

int MPI\_Prequest\_free(MPI\_Prequest \*preq)

Executes in the host

- MPI\_Prequest should be created on the same device as the message in MPI\_Send\_Init().
- The MPI implementation is free to allocate the MPI\_Prequest on USM shared, device or host.
- Each MPI implementation can determine the information that is needed on the MPI\_Prequest
- MPI\_Prequest objects are only valid for use in device functions

#### Device Ready

```
SYCL_EXTERNAL int MPI_Pready(MPI_Prequest prequest, int partition)

SYCL_EXTERNAL int MPI_Pready_range(int part_low, int part_high, MPI_Prequest prequest)

SYCL_EXTERNAL int MPI_Pready_list(int length, const int array_of_partitions[], MPI_Prequest prequest)
```

Execute in the device

#### Device Arrived

SYCL\_EXTERNAL int MPI\_Parrived(MPI\_Prequest preq, int partition, int \*flag)

Execute in the device

#### Others

For efficiency reasons we should define variants of MPI\_Ready()/MPI\_Parrived()
for work items (warp) and/or work groups (thread block)

#### Example 4.2 MPI standard (OpenMP offload)

```
if (myrank == 1) {
if (myrank == 0) {
                                                            MPI_Precv_init(message, partitions, count, xfer_type
 MPI_Psend_init(message, partitions, count, xfer_type,
   dest, tag, info, MPI_COMM_WORLD, &request);
                                                              source, tag, info, MPI_COMM_WORLD, &request);
                                                            MPI_Start(&request);
 MPI_Start(&request);
 #pragma omp target teams distribute parallel for
 for (int i = 0; i < partitions; i++) {
    /* compute and fill partition #i, then mark ready: */
    MPI_Pready(i, request);
 while (!flag) {
                                                           while (!flag) {
   /* Do useful work */
                                                            /* Do useful work */
   MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
                                                              MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
 MPI_Request_free(&request);
                                                           MPI_Request_free(&request);
```

#### OpenMP declare target directive

declare target directive is used to indicate that the corresponding call inside a target region is to a function that can execute on the default target device

```
#pragma omp declare target
int MPI_Pready(MPI_Request request, int partition) {
    ...
}
#pragma omp end declare target
```

#