# The Memory Allocation Kinds Side Document
## Version 1.0

Authors
Rohit Zambre, James Dinan, Maria Garzaran, Daniel Holmes, Edgar Gabriel,
Michael Klemm, and Pedram Alizadeh

# Contents

# Chapter 1

# Overview

Modern computing systems contain a variety of memory types, each closely associated with a distinct type of computing hardware. For example, compute accelerators such as GPUs typically feature their own memory that is distinct from the memory attached to the host processor. Additionally, GPUs from different vendors also differ in their memory types. The differences in memory types influence feature availability and performance behavior of an application running on such modern systems. Hence, MPI libraries need to be aware of and support additional memory types. For a given type of memory, MPI libraries need to know the associated memory allocator and the limitations on memory access. The different memory kinds capture the differentiating information needed by MPI libraries for different memory types.

This MPI side document defines the memory allocation kinds and their associated restrictors that users can use to query the support for different memory kinds provided by the MPI library. These definitions supplement those found in section 11.4.3 of the MPI standard, which also explains their usage model.

# Chapter 2

# Definitions

This section contains definitions of memory allocation kinds and their restrictors.

## 2.1   Kind: cuda

The cuda memory kind refers to the memory allocated by the CUDA runtime system [1].

Restrictors

- host: Support for memory allocations on the host system that are page-locked for direct access from the CUDA device (e.g., memory allocations from the `cudaHostAlloc()` function). These memory allocations are attributed with `cudaMemoryTypeHost`.

- device: Support for memory allocated on a CUDA device (e.g., memory allocations from the `cudaMalloc()` function). These memory allocations are attributed with `cudaMemoryTypeDevice`.

- managed: Support for memory that is managed by CUDA's Unified Memory system (e.g., memory allocations from the `cudaMallocManaged()` function). These memory allocations are attributed with `cudaMemoryTypeManaged`.

**Example 2.1** This example demonstrates the usage of the different kinds to perform communication in a manner that is supported by the underlying MPI library.

```c
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <mpi.h>

int main(int argc, char *argv[])
{
    int cuda_device_aware = 0;
    int cuda_managed_aware = 0;
    int len = 0, flag = 0;
    int *managed_buf = NULL;
    int *device_buf = NULL, *system_buf = NULL;
    MPI_Info info;
    MPI_Session session;
    MPI_Group wgroup;
```

```
1    MPI_Comm system_comm;
2    MPI_Comm cuda_managed_comm = MPI_COMM_NULL;
3    MPI_Comm cuda_device_comm = MPI_COMM_NULL;
4
5    MPI_Info_create(&info);
6    MPI_Info_set(info, "mpi_memory_alloc_kinds",
7                       "system,cuda:device,cuda:managed");
8    MPI_Session_init(info, MPI_ERRORS_ARE_FATAL, &session);
9    MPI_Info_free(&info);
10
11   MPI_Session_get_info(session, &info);
12   MPI_Info_get_string(info, "mpi_memory_alloc_kinds",
13                       &len, NULL, &flag);
14
15   if (flag) {
16       char *val, *valptr, *kind;
17
18       val = valptr = (char *) malloc(len);
19       if (NULL == val) return 1;
20
21       MPI_Info_get_string(info, "mpi_memory_alloc_kinds",
22                           &len, val, &flag);
23
24       while ((kind = strsep(&val, ",")) != NULL) {
25           if (strcasecmp(kind, "cuda:managed") == 0) {
26               cuda_managed_aware = 1;
27           }
28           else {
29               if (strcasecmp(kind, "cuda:device") == 0) {
30                   cuda_device_aware = 1;
31               }
32       }
33       free(valptr);
34   }
35
36   MPI_Info_free(&info);
37
38   MPI_Group_from_session_pset(session, "mpi://WORLD" , &wgroup);
39
40   // Create a communicator for operations on system memory
41   MPI_Info_create(&info);
42   MPI_Info_set(info, "mpi_assert_memory_alloc_kinds", "system");
43   MPI_Comm_create_from_group(wgroup,
44           "org.mpi-side-doc.mem-kind.example.system",
45           info, MPI_ERRORS_ABORT, &system_comm);
46
47   MPI_Info_free(&info);
48
49   // Check if all processes have CUDA managed support
50   MPI_Allreduce(MPI_IN_PLACE, &cuda_managed_aware, 1, MPI_INT,
51                 MPI_LAND, system_comm);
52
53   if (cuda_managed_aware) {
54       // Create a communicator for operations that use
```

```
1        // CUDA managed buffers.
2        MPI_Info_create (& info );
3        MPI_Info_set (info , " mpi_assert_memory_alloc_kinds ",
4                      "cuda:managed ");
5        MPI_Comm_create_from_group (wgroup ,
6            "org.mpi-side-doc.mem-kind.example.cuda.managed",
7            info , MPI_ERRORS_ABORT , & cuda_managed_comm );
8        MPI_Info_free (& info );
9    }
10    else {
11        // Check if all processes have CUDA device support
12        MPI_Allreduce (MPI_IN_PLACE , & cuda_device_aware , 1, MPI_INT ,
13                      MPI_LAND , system_comm );
14        if ( cuda_device_aware ) {
15            // Create a communicator for operations that use
16            // CUDA device buffers.
17            MPI_Info_create (& info );
18            MPI_Info_set (info , " mpi_assert_memory_alloc_kinds ",
19                          "cuda:device ");
20            MPI_Comm_create_from_group (wgroup ,
21                "org.mpi-side-doc.mem-kind.example.cuda.device",
22                info , MPI_ERRORS_ABORT , & cuda_device_comm );
23            MPI_Info_free (& info );
24        }
25        else {
26            printf (" Warning: cuda alloc kind not supported\n");
27        }
28    }
29
30    MPI_Group_free (& wgroup );
31
32    if ( cuda_managed_aware ) {
33        // Allocate managed buffer and initialize it
34        cudaMallocManaged (& managed_buf , sizeof (int ));
35        *managed_buf = 1;
36
37        // Perform communication using cuda_managed_comm
38        // if it's available.
39        MPI_Allreduce (MPI_IN_PLACE , managed_buf , 1, MPI_INT ,
40                      MPI_SUM , cuda_managed_comm );
41
42        cudaFree ( managed_buf );
43    }
44    else {
45        // Allocate system buffer and initialize it
46        system_buf = malloc ( sizeof (int ));
47        *system_buf = 1;
48
49        // Allocate CUDA device buffer and initialize it
50        cudaMalloc (& device_buf , sizeof (int ));
51        cudaMemcpy ( device_buf , system_buf , sizeof (int),
52                    cudaMemcpyHostToDevice );
53
54        if ( cuda_device_aware ) {
```

```
1          // Perform communication using cuda_comm
2          // if it's available.
3          MPI_Allreduce(MPI_IN_PLACE, device_buf, 1, MPI_INT,
4                        MPI_SUM, cuda_comm);
5      }
6      else {
7          // Otherwise, copy data to a system buffer,
8          // use system_comm, and copy data back to device buffer
9          cudaMemcpy(system_buf, device_buf, sizeof(int),
10                 cudaMemcpyDeviceToHost);
11         MPI_Allreduce(MPI_IN_PLACE, system_buf, 1, MPI_INT,
12                       MPI_SUM, system_comm);
13         cudaMemcpy(device_buf, system_buf, sizeof(int),
14                 cudaMemcpyHostToDevice);
15
16         cudaFree(device_buf);
17         free(system_buf);
18     }
19  }
20
21  if (cuda_managed_comm != MPI_COMM_NULL)
22      MPI_Comm_disconnect(&cuda_managed_comm);
23  if (cuda_device_comm != MPI_COMM_NULL)
24      MPI_Comm_disconnect(&cuda_device_comm);
25  MPI_Comm_disconnect(&system_comm);
26
27  MPI_Session_finalize(&session);
28
29  return 0;
30 }
```

## 2.2 Kind: rocm

The rocm memory kind refers to the memory allocated by the ROCm runtime system [2].

Restrictors

- host: Support for memory allocated on the host system that is page-locked for direct access from the ROCm device (e.g., memory allocations from the `hipHostMalloc()` function). These memory allocations are attributed with `hipMemoryTypeHost`.

- device: Support for memory allocated on the ROCm device (e.g., memory allocations from the `hipMalloc()` function). These memory allocations are attributed with `hipMemoryTypeDevice`.

- managed: Support for memory that is managed automatically by the ROCm runtime (e.g., memory allocations from the `hipMallocManaged()` function). These memory allocations are attributed with `hipMemoryTypeManaged`.

## 2.3  Kind: levelzero

The levelzero memory kind refers to the memory allocated by the Level Zero runtime system [3].

Restrictors

- host: Support for memory allocated on the host that is accessible by Level Zero devices (e.g., memory allocations from the `zeMemAllocHost()` function). These memory allocations are attributed with `ZE_MEMORY_TYPE_HOST`.

- device: Support for memory allocated on a Level Zero device (e.g., memory allocations from the `zeMemAllocDevice()` function). These memory allocations are attributed with `ZE_MEMORY_TYPE_DEVICE`.

- shared: Support for memory allocated that will be shared between the host and one or more Level Zero devices (e.g., memory allocations from the `zeMemAllocShared()` function). These memory allocations are attributed with `ZE_MEMORY_TYPE_SHARED`.

# Annex A

# Changelog

This annex summarizes changes from the previous versions of the *Memory Allocation Kinds* side document to the version presented by this document.

## A.1   Version 1.0

1. The first version of this document.

# Bibliography

[1] CUDA Runtime API. https://docs.nvidia.com/cuda/cuda-runtime-api/.

[2] HIP Programming Manual. https://rocm.docs.amd.com/en/latest/reference/hip.html.

[3] Level Zero Programming Guide. https://spec.oneapi.io/level-zero/latest/core/PROG.html.