



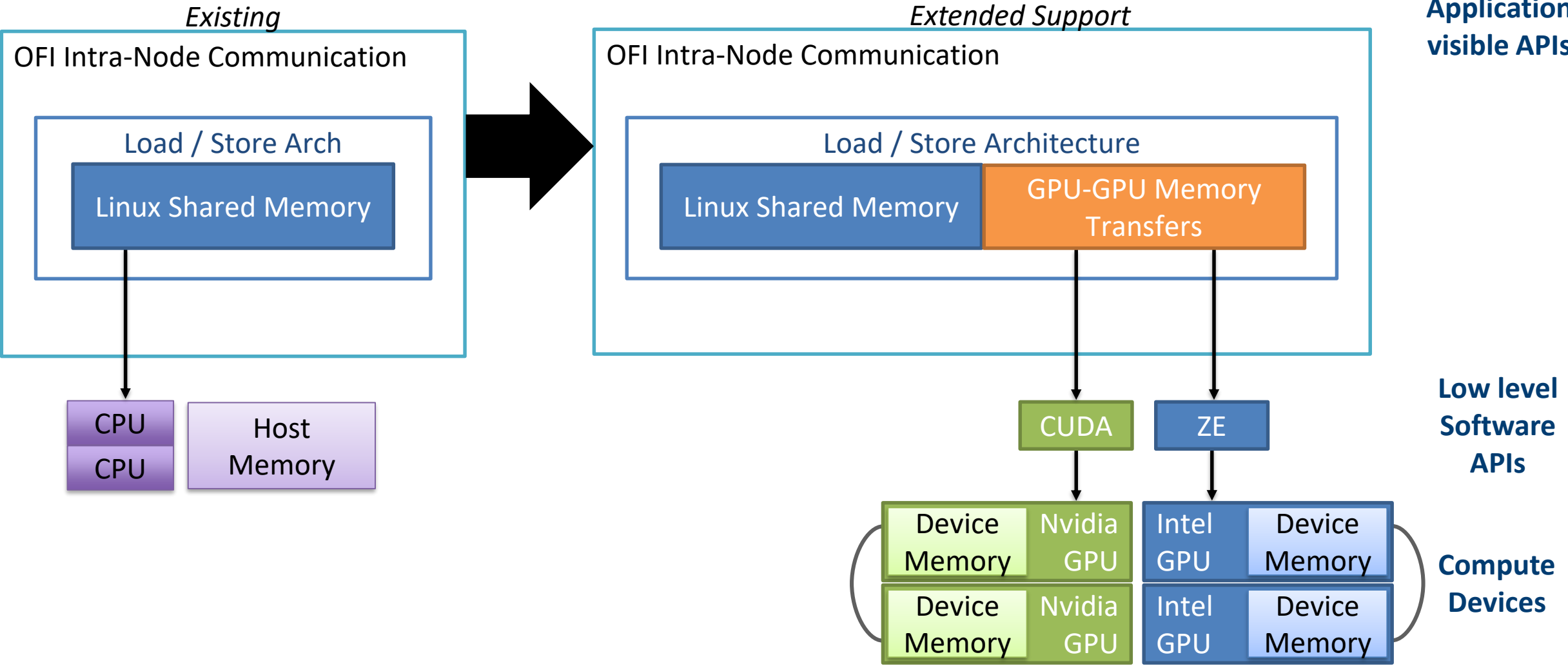
2020 OFA Virtual Workshop

LIBFABRIC INTRANODE DEVICE SUPPORT

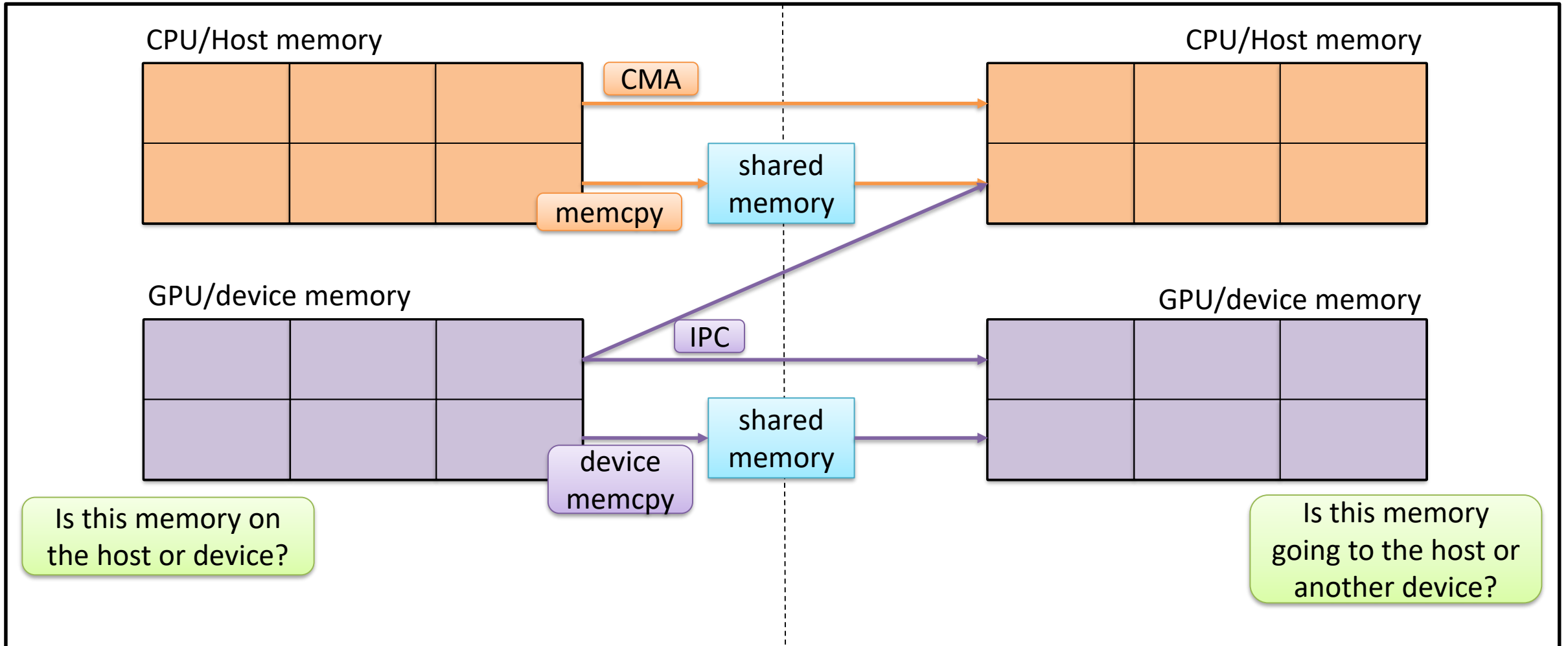
Alexia Ingerson

Intel Corp.

GPU-GPU COMMUNICATION OVERVIEW



OVERVIEW



OFI API CHANGES

▪ info->caps: FI_HMEM

- Requests support for transfers to and from device memory

▪ domain_attr->mr_mode: FI_MR_HMEM

- Specifies that the application should register all device memory with proper interfaces
- Eliminates the need for a provider to query devices in order to determine memory location (expensive)

▪ fi_mr_attr->iface

- Indicates software interface used to manage memory region

```
enum fi_hmem_iface {  
    FI_HMEM_SYSTEM = 0, //system/host memory  
    FI_HMEM_CUDA, //Nvidia/CUDA memory (libcuda)  
    FI_HMEM_ZE, //Intel/Ze memory (libze_loader)  
    ...  
}
```

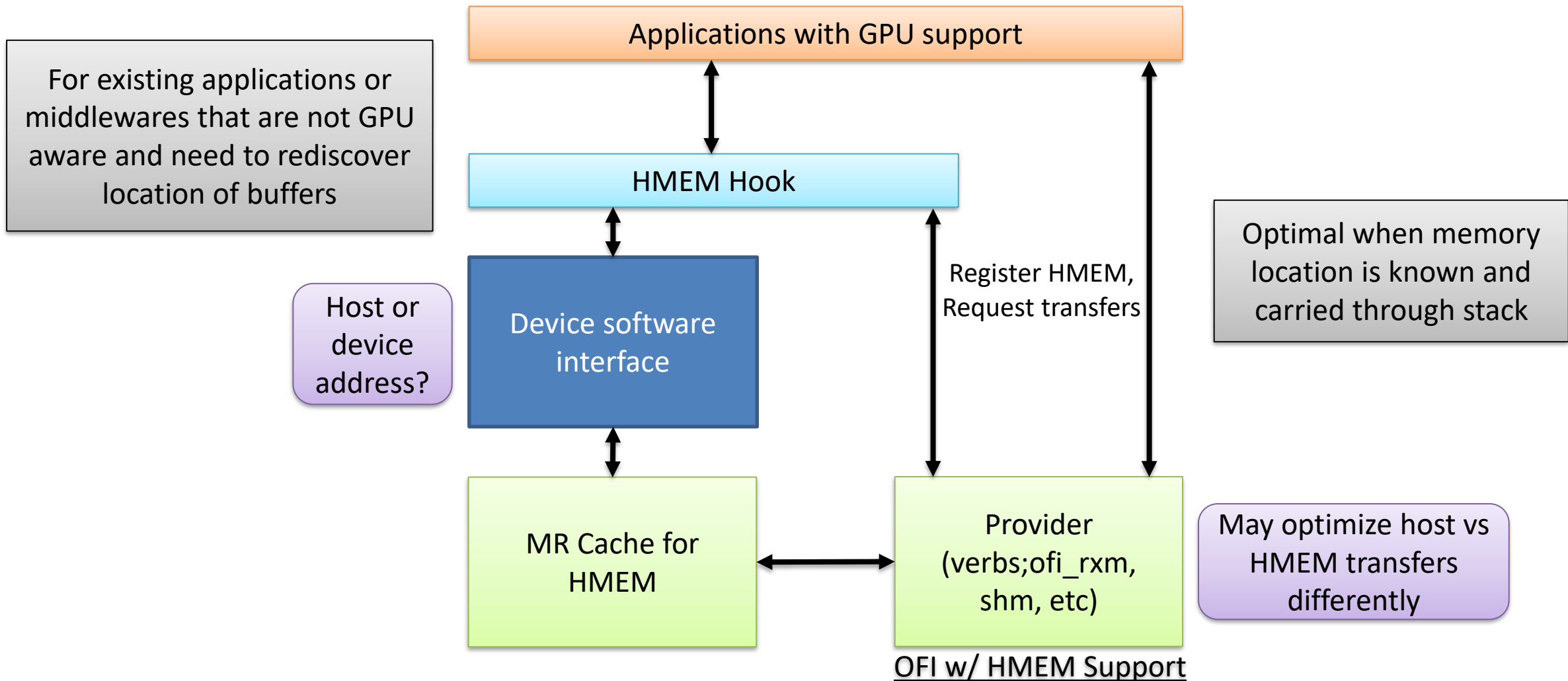
- Tells provider which API calls to use when copying to and from device

▪ fi_mr_attr->device

- Device identifier for HMEM memory
- Indicates on which device the memory is located on (type varies by interface) when multiple devices are present

```
fi_mr_attr {  
    ...  
    enum fi_hmem_iface  iface;  
    union {  
        uint64_t        reserved;  
        int              cuda;  
        void             *ze;  
    } device;  
}
```

HMEM HOOK



SHM OVERVIEW

smr_region

EP initialized info / resources

-
-
-

Command Queue

Single command queue for incoming messages

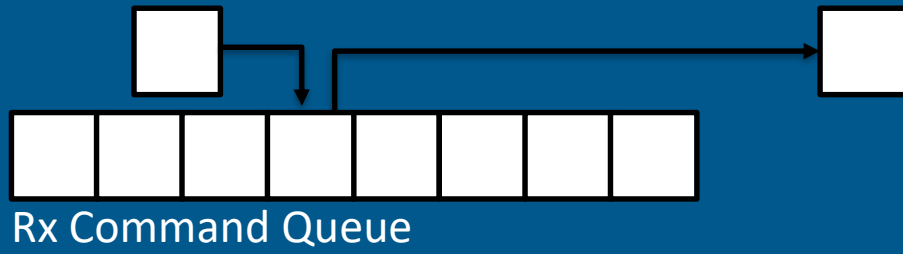
Response Queue

Response queue for messages requiring an ACK

Inject Pool

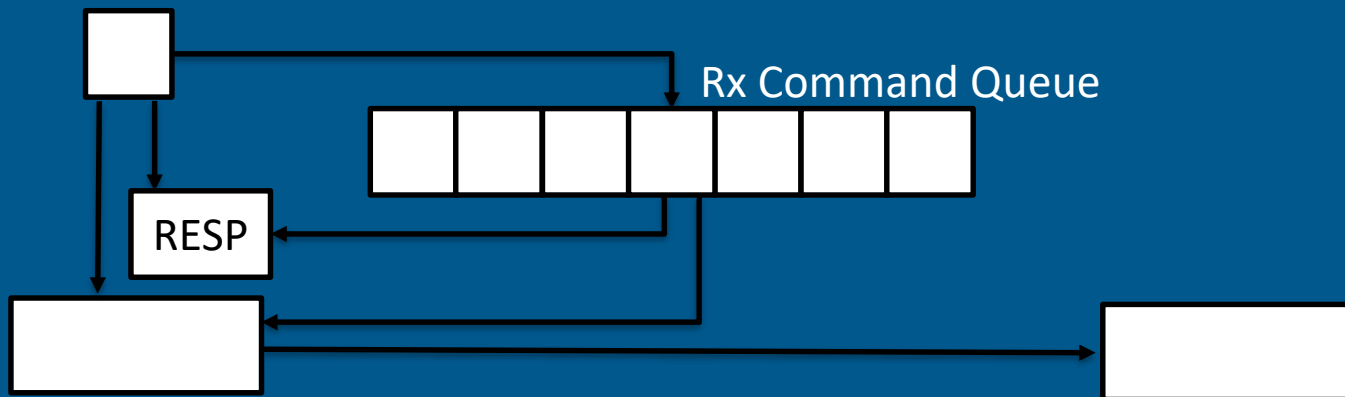
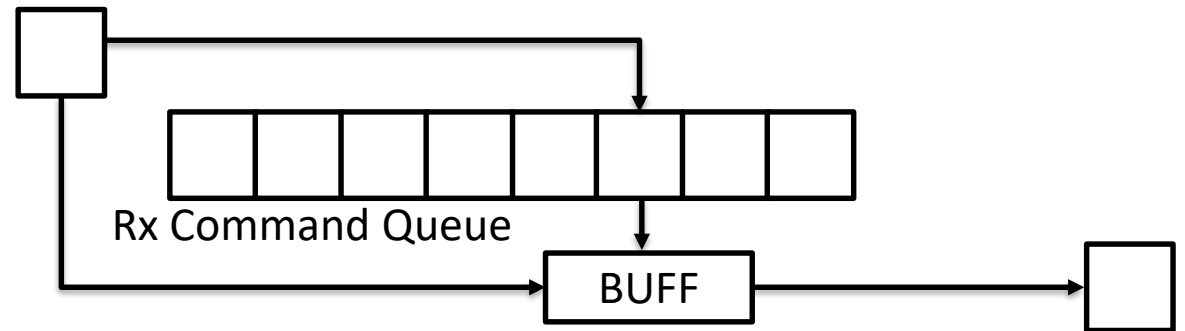
Pool of bounce buffers for medium-sized messages

SHM PROTOCOLS



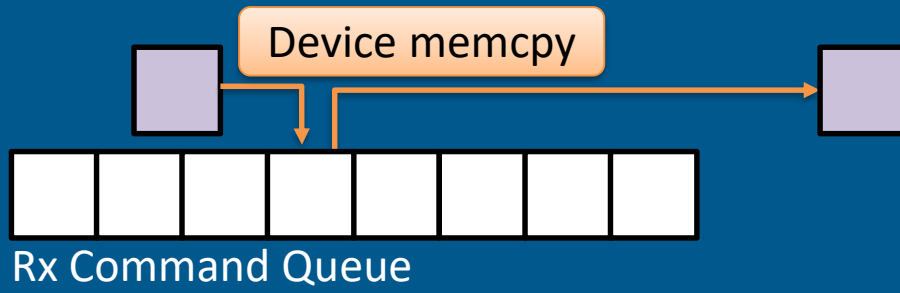
Inline

Inject



IOV

SHM PROTOCOLS

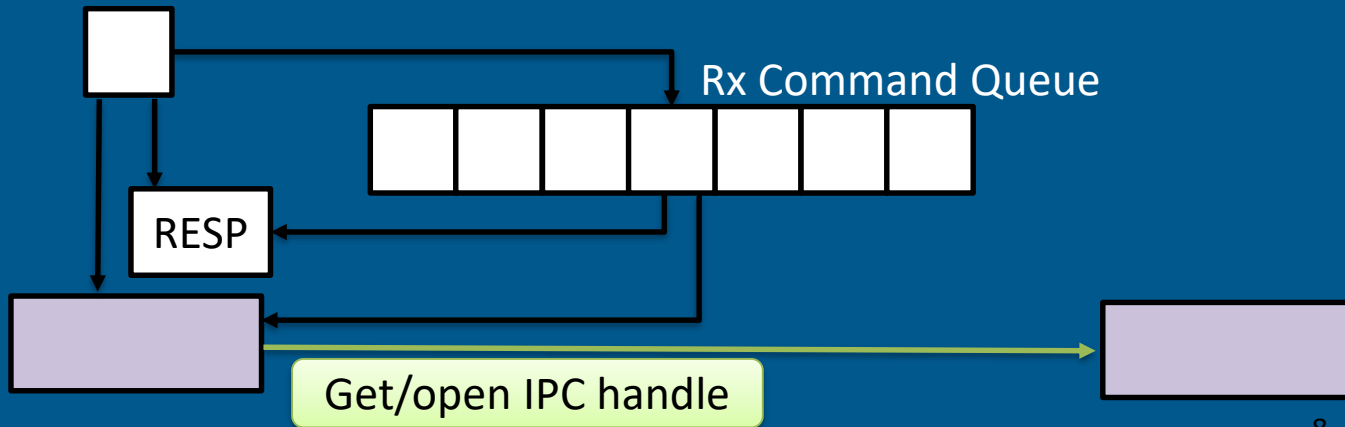
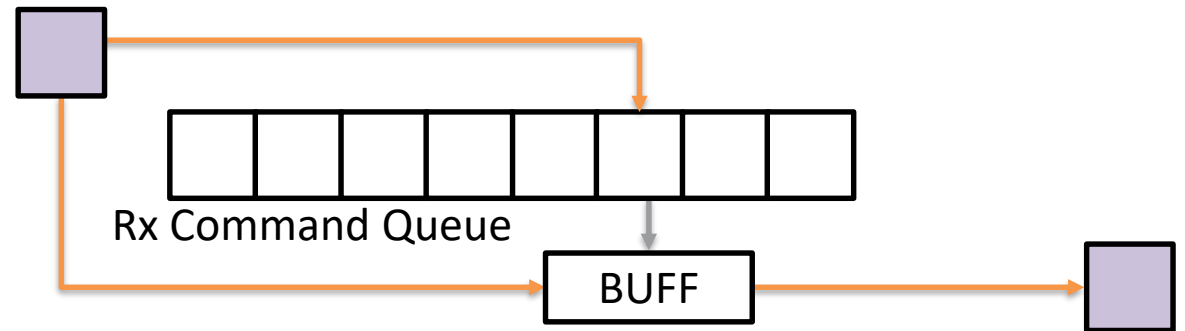


Inline



Device->host/
host->device
copies tend to be
expensive

Inject



IOV

DEVICE BOUNCE BUFFERS

smr_region

EP initialized info / resources

-
-
-

Command Queue

Response Queue

Inject Pool

IPC Handle Pool

Single command queue for incoming messages

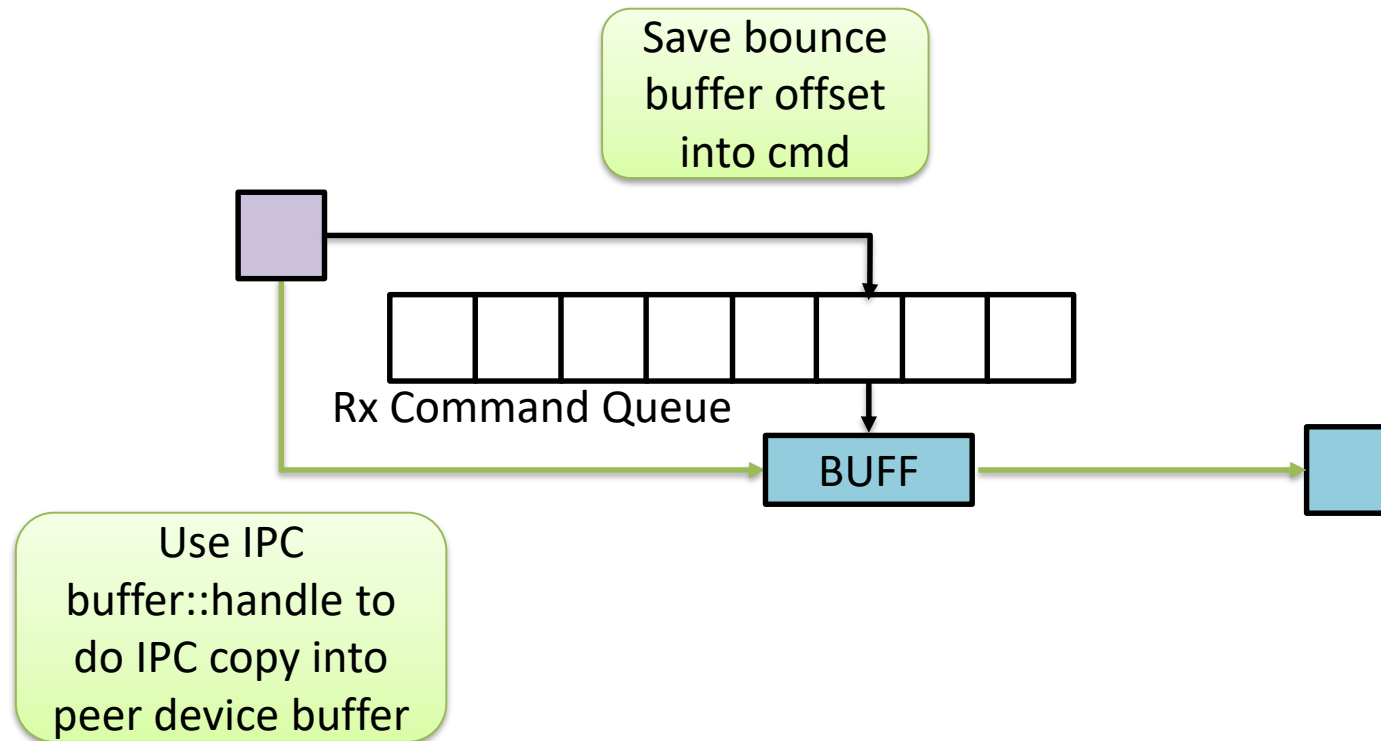
Response queue for messages requiring an ACK

Pool of bounce buffers for medium-sized messages

Pool of device buffer IPC handles

```
smr_ipc_handle {  
    void      *buf;  
    uint8_t   handle[64];  
}
```

DEVICE BOUNCE BUFFER PROTOCOL



GENERIC BUFFER OPS

Initialization/cleanup of any device-specific resource

Memory copy functions for copying to and from device memory

IPC handle functions for copying device memory across processes

Alloc and free calls for use in device bounce buffer protocol

```
struct ofi_hmem_ops {  
    (*init) ();  
    (*cleanup) ();  
    (*copy_to_hmem) ();  
    (*copy_from_hmem) ();  
    (*get_handle) ();  
    (*open_handle) ();  
    (*close_handle) ();  
    (*alloc) ();  
    (*free) ();  
};
```

SAMPLE HMEM OPS

	Intel GPU Level 0 API (Ze)	Nvidia GPU CUDA API
INIT	<code>zeInit()</code> <code>zeDriverGet()</code> <code>zeCommandQueueCreate()</code> <code>zeCommandListCreate()</code>	<code>cuInit()</code>
COPY	<code>zeCommandListAppendMemoryCopy()</code> <code>zeCommandQueueExecuteCommandLists()</code>	<code>cudaMemcpy()</code>
IPC	<code>zeDriverGetMemIpcHandle()</code> <code>zeDriverOpenMemIpcHandle()</code> <code>zeDriverCloseMemIpcHandle()</code>	<code>cudaIpcGetMemHandle()</code> <code>cudaIpcOpenHandle()</code> <code>cudaIpcCloseMemHandle()</code>
MEM	<code>zeDriverAllocSharedMem()</code> <code>zeDriverFreeMem()</code>	<code>cudaMalloc()</code> <code>cudaFree()</code>



OPENFABRICS
ALLIANCE

2020 OFA Virtual Workshop

THANK YOU

Alexia Ingerson

Intel Corp.