



PARALLEL PROGRAMMING...

Copyright 2023 Patrick Lemoine. All rights reserved.

Parallel Computing Using HIP

SESSION 3/3



Programming **I**nterface for **p**arallel **c**omputing With HIP

What is HIP ?

Synchronization and streams

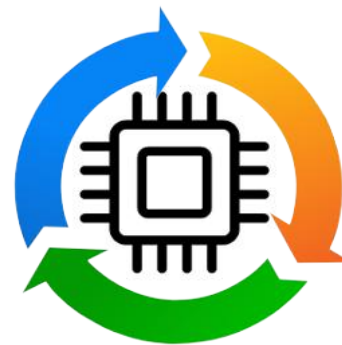
Memory allocations, access and unified memory

Kernel optimization and profiling

API **E**xamples



What is **H**IP ?



What is HIP ?

AMD
HIP RT



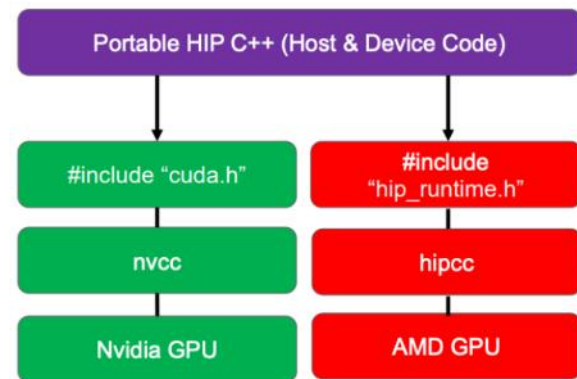
Heterogeneous-Compute Interface for Portability (HIP)

AMD effort to offer a common programming interface that works on both CUDA and ROCm devices.

C++ runtime API and kernel language that allows developers to create portable applications.

- Syntactically similar to CUDA so that most API calls can be converted from CUDA to HIP with a simple cuda → hip translation.

ROCm: is an Advanced Micro Devices (AMD) software stack for graphics processing unit (GPU) programming and spans several domains: GPGPU, HPC ,heterogeneous computing



HIP comparison with CUDA

CUDA

```
cudaMalloc((void**)&nodes_dev, N*sizeof(float) );
```

HIP

```
hipMalloc((void**)&nodes_dev, N*sizeof(float) );
```

CUDA

```
dim3 threadsPerBlock(nthreads,nthreads,nthreads);  
dim3 blocks(n_elements);  
GPUKernel<<<blocks,threadsPerBlock>>>( input );
```

HIP

```
dim3 threadsPerBlock(nthreads,nthreads,nthreads);  
dim3 blocks(n_elements);  
hipLaunchKernelGGL(GPUKernel, dim3(blocks), dim3(threadsPerBlock), 0, 0, input);
```



COMPILATION

With CUDA

```
==>$ nvcc source_code.cu
```

With HIP

```
==>$ hipcc source_code.cu
```

HIP Kernel Language

HIP Kernel Language

Qualifiers: `__device__`, `__global__`, `__shared__`, ...

Built-in variables: `threadIdx.x`, `blockIdx.y`, ...

Vector types: `int3`, `float2`, `dim3`, ...

Math functions: `sqrt`, `powf`, `sinh`, ...

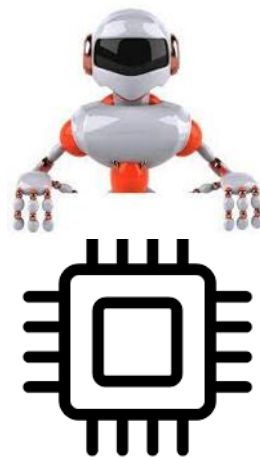
Intrinsic functions: synchronisation, memory-fences etc.

Programming models

GPU accelerator is often called a device and CPU a host

Parallel code (kernel) is launched by the host and executed on a device by several threads

Code is written from the point of view of a single thread each thread has a unique ID



API

Device init and management

Memory management

Execution control

Synchronisation: device, stream, events

Error handling, context handling

AMD GPU Terminology



Compute Unit

- one of the parallel vector processors in a GPU

Kernel

- function launched to the GPU that is executed by multiple parallel workers

Thread

- individual lane in a wavefront

Wavefront (cf. CUDA warp)

- collection of threads that execute in lockstep and execute the same

instructions

- each wavefront has 64 threads
- number of wavefronts per workgroup is chosen at kernel launch

Workgroup (cf. CUDA thread block)

- group of wavefronts (threads) that are on the GPU at the same time and
- are part of the same compute unit (CU)
- can synchronise together and communicate through memory in the CU

GPU Programming consideration



GPU model requires many small tasks executing a kernel

- e.g. can replace iterations of loop with a GPU kernel call

Need to adapt CPU code to run on the GPU

- rethink algorithm to fit better into the execution model
- keep reusing data on the GPU to reach high occupancy of the hardware
- if necessary, manage data transfers between CPU and GPU memories carefully (can easily become a bottleneck!)



Grid: thread hierarchy



Kernels are executed on a 3D grid of threads
- threads are partitioned into equalsized blocks

Code is executed by the threads,
the grid is just a way to organise
the work



Dimension of the grid are set at
kernel launch

Built-in variables to be used within a kernel:
- threadIdx, blockIdx, blockDim, gridDim



Kernels

Kernel is a (device) function to be executed by the GPU

Function should be of void type and needs to be declared with the

`__global__` or `__device__` attribute

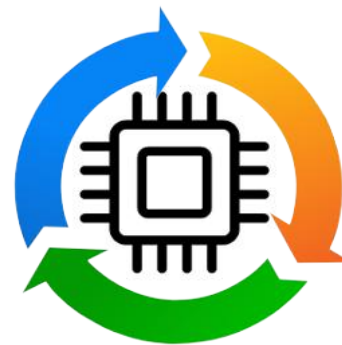
All pointers passed to the kernel need to point to memory accessible from the device

Unique thread and block IDs can be used to distribute work





HIP Synchronization and streams



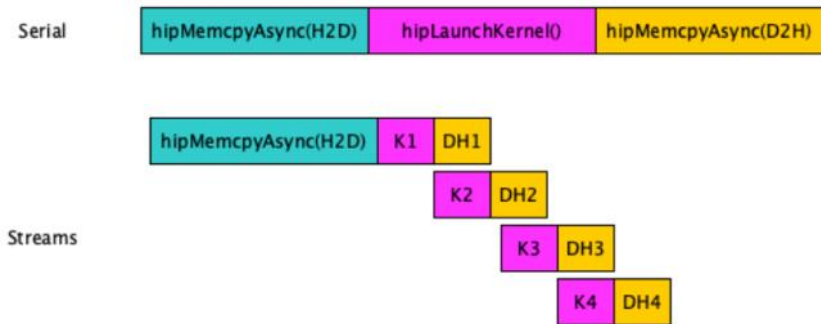
What is a stream ?

- A sequence of operations that execute in issue-order on the GPU
- HIP operations in different streams could run concurrently
- The ROCm 4.5.0 brings the Direct Dispatch, the runtime directly queues a packet to the AQL queue in Dispatch and some of the synchronization.
- The previous ROCm uses queue per stream

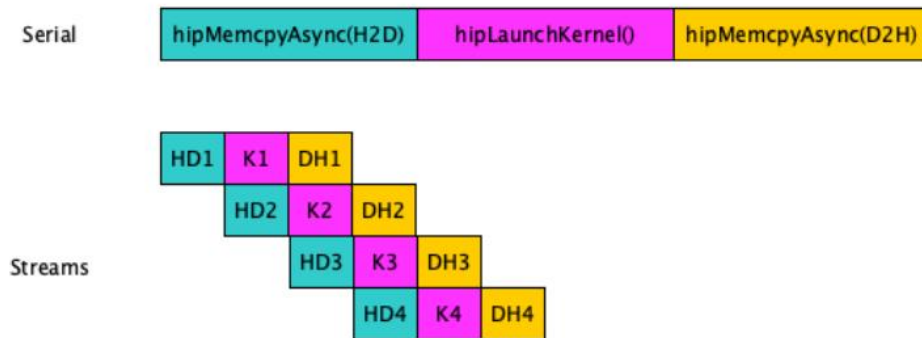


Stream

Concurrency



Amount of concurrency



Default

Only a single stream is used if not defined

Commands are synchronized except the Async calls and Kernels

Stream

Stream/Events API

hipStreamCreate:	Creates an asynchronous stream
hipStreamDestroy	Destroy an asynchronous stream
hipStreamCreateWithFlags	Creates an asynchronous stream with specified flags
hipEventCreate	Create an event
hipEventRecord	Record an event in a specified stream
hipEventSynchronize:	Wait for an event to complete
hipEventElapsedTime:	Return the elapsed time between two events
hipEventDestroy:	Destroy the specified event

HIP: Example data transfer and compute

Serial

```
hipCheck( hipEventRecord(startEvent,0) );

hipCheck( hipMemcpy(d_a, a, bytes, hipMemcpyHostToDevice) );

hipLaunchKernelGGL(kernel, n/blockSize, blockSize, 0, 0, d_a, 0);

hipCheck( hipMemcpy(a, d_a, bytes, hipMemcpyDeviceToHost) );

hipCheck( hipEventRecord(stopEvent, 0) );
hipCheck( hipEventSynchronize(stopEvent) );
hipCheck( hipEventElapsedTime(&duration, startEvent, stopEvent) );
printf("Duration of sequential transfer and execute (ms): %f\n", duration);
```

HIP: How on improve the performance ?

- Use streams to overlap computation with communication

```
hipStream_t stream[nStreams];  
for (int i = 0; i < nStreams; ++i) hipStreamCreate(&stream[i])
```

- Use Asynchronous data transfer

- Execute kernels on different streams

```
hipLaunchKernelGGL(some_kernel, gridsize, blocksize, shared_mem_size, stream, arg0, arg1, ...);
```



HIP: Synchronization

- Synchronize everything, could be used after each kernel launch except if you know what you are doing

hipDeviceSynchronize()

- Synchronize a specific stream Blocks host until all HIP calls are completed on this stream

hipStreamSynchronize(streamid)

- Synchronize using Events

Create event

hipEvent_t stopEvent

hipEventCreate(&stopEvent)

Record an event in a specific stream and wait until is recorded

hipEventRecord(stopEvent,0)

hipEventSynchronize(stopEvent)

Make a stream wait for a specific event

hipStreamWaitEvent(stream[i], stopEvent, unsigned int flags)



HIP: Synchronization in kernel

Code

```
__global__ void reverse(double *d_a)
{
    __shared__ double s_a[256]; //array of doubles, shared in this block
    int tid = threadIdx.x;
    s_a[tid] = d_a[tid]; //each thread fills one entry

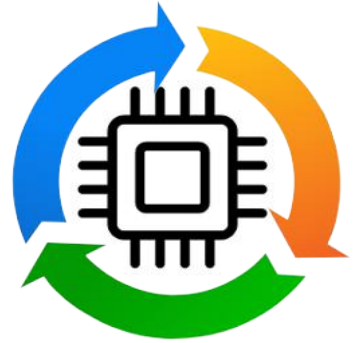
    //all wavefronts must reach this point before any wavefront is allowed to continue.

    __syncthreads();
    d_a[tid] = s_a[255-tid]; //write out array in reverse order
}
```





HIP Memory allocations,
access and unified memory



HIP: Memory

Memory model

Host and device have separate physical memories

It is generally not possible to call malloc() to allocate memory and access the data from the GPU

Memory management can be

- Explicit (user manages the movement of the data and makes sure CPU and GPU pointers are not mixed)
- Automatic, using Unified Memory (data movement is managed in the background by the Unified Memory driver)

Avoid moving data between CPU and GPU

Data copies between host and device are relatively slow

To achieve best performance, the host-device data traffic should be minimized regardless of the chosen memory management strategy

Initializing arrays on the GPU

Rather than just solving a linear equation on a GPU, also setting it up on the device

Not copying data back and forth between CPU and GPU every step or iteration can have a large performance impact!

HIP: Memory

Device memory hierarchy

- **Registers (per-thread-access)**
 - Used automatically
 - Size on the order of kilobytes
 - Very fast access
- **Local memory (per-thread-access)**
 - Used automatically if all registers are reserved
 - Local memory resides in global memory
 - Very slow access
- **Shared memory (per-blockaccess)**
 - Usage must be explicitly programmed
 - Size on the order of kilobytes
 - Fast access
- **Global memory (per-deviceaccess)**
 - Managed by the host through HIP API
 - Size on the order of gigabytes
 - Very slow access
- **There are more details in the memory hierarchy, some of which are architecture-dependent, eg,**
 - Texture memory
 - Constant memory

Complicates implementation
Should be considered only when a very high level of optimization is desirable

HIP: Memory

Important memory operations

- Allocate pinned device memory
hipError_t hipMalloc(void **devPtr, size_t size)
- Allocate Unified Memory; The data is moved automatically between host/device
hipError_t hipMallocManaged(void **devPtr, size_t size)
- Deallocate pinned device memory and Unified Memory
hipError_t hipFree(void *devPtr)
- Copy data (host-host, host-device, device-host, device-device)
hipError_t hipMemcpy(void *dst, const void *src, size_t count, enum hipMemcpyKind kind)



HIP: Memory

Example of explicit memory management

```
int main() {
    int *A, *d_A;
    A = (int *) malloc(N*sizeof(int));
    hipMalloc((void**)&d_A, N*sizeof(int));
    ...
    /* Copy data to GPU and launch kernel */
    hipMemcpy(d_A, A, N*sizeof(int), hipMemcpyHostToDevice);
    hipLaunchKernelGGL(...);
    ...
    hipMemcpy(A, d_A, N*sizeof(int), hipMemcpyDeviceToHost);
    hipFree(d_A);
    printf("A[0]: %d\n", A[0]);
    free(A);
    return 0;
}
```

Example of Unified Memory

```
int main() {
    int *A;
    hipMallocManaged((void**)&A, N*sizeof(int));
    ...
    /* Launch GPU kernel */
    hipLaunchKernelGGL(...);
    hipStreamSynchronize(0);
    ...
    printf("A[0]: %d\n", A[0]);
    hipFree(A);
    return 0;
}
```

HIP: Memory

✓ Unified Memory pros

- Allows incremental development

- Can increase developer productivity significantly

 - Especially large codebases with complex data structures

- Supported by the latest NVIDIA + AMD architectures

- Allows oversubscribing GPU memory on some architectures



✗ Unified Memory cons

- Data transfers between host and device are initially slower, but can be optimized once the code works

 - Through prefetches

 - Through hints

- Must still obey concurrency & coherency rules, not foolproof

- The performance on the AMD cards is an open question

HIP: Memory

Unified Memory workflow for GPU offloading

- Allocate memory for the arrays accessed by the GPU with `hipMallocManaged()` instead of `malloc()`
It is a good idea to have a wrapper function or use function overloading for memory allocations
- Offload compute kernels to GPUs
- Check profiler backtrace for GPU->CPU Unified Memory page-faults

(NVIDIA Visual Profiler, Nsight Systems, AMD profiler?)

- This indicates where the data residing on the GPU is accessed by the CPU very useful for large codebases, especially if the developer is new to the code)

HIP: Memory

Unified Memory workflow for GPU offloading

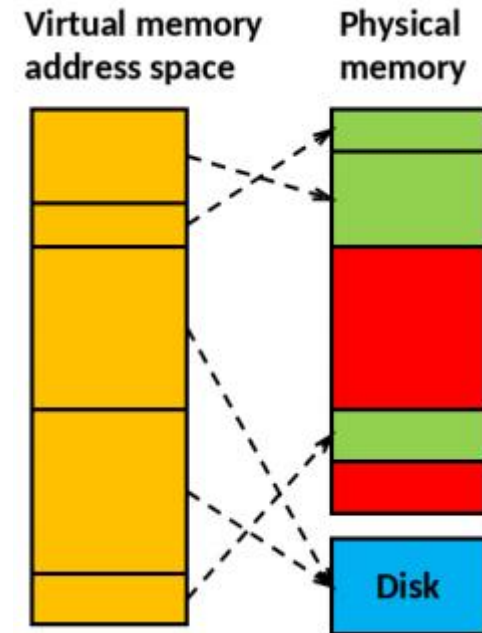
- Move operations from CPU to GPU if possible, or use hints / prefetching (`hipMemAdvice()` / `hipMemPrefetchAsync()`)
- It is not necessary to eliminate all page faults, but eliminating the most frequently occurring ones can provide significant performance improvements
- Allocating GPU memory can have a much higher overhead than allocating standard host memory

If GPU memory is allocated and deallocated in a loop, consider using a GPU memory pool allocator for better performance

HIP: Memory

Virtual Memory addressing

- Modern operating systems utilize virtual memory
 - Memory is organized to memory pages
 - Memory pages can reside on swap area on the disk (or on the GPU with Unified Memory)



HIP: Memory

Page-locked (or pinned) memory

- Normal malloc() allows swapping and page faults
- User can page-lock an allocated memory block to a particular physical memory location
- Enables Direct Memory Access (DMA)
- Higher transfer speeds between host and device
- Copying can be interleaved with kernel execution
- Page-locking too much memory can degrade system performance due to paging problems

HIP: Memory

Allocating page-locked memory on host

- Allocated with `hipHostMalloc()` function instead of `malloc()`
- The allocation can be mapped to the device address space for device access (slow)
 - On some architectures, the host pointer to device-mapped allocation can be directly used in device code (ie. it works similarly to Unified Memory pointer, but the access from the device is slow)

Deallocated using `hipHostFree()`



HIP: Memory

Asynchronous memcopies

- Normal `hipMemcpy()` calls are blocking (i.e. synchronizing)
- The execution of host code is blocked until copying is finished
- To overlap copying and program execution, asynchronous functions are required
 - Such functions have Async suffix, eg. **`hipMemcpyAsync()`**
- User has to synchronize the program execution
- Requires page-locked memory

HIP: Memory

Global memory access in device code

- Global memory access from the device is slow
- Threads are executed in warps, memory operations are grouped in a similar fashion
- Memory access is optimized for coalesced access where threads read from and write to successive memory locations
- Exact alignment rules and performance issues depend on the architecture

HIP: Memory

Coalesced memory access

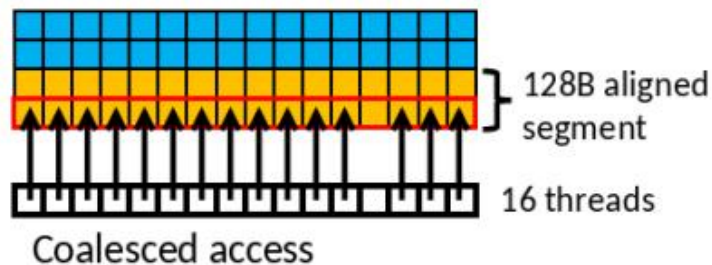
- The global memory loads and stores consist of transactions of a certain size (eg. 32 bytes)
- If the threads within a warp access data within such a block 32 bytes, only one global memory transaction is needed
- Now, 32 threads within a warp can each read a different 4-byte integer value with just 4 transactions
- When the stride between each 4- byte integer is increased, more transactions are required (up to 32 for the worst case)!

HIP: Memory

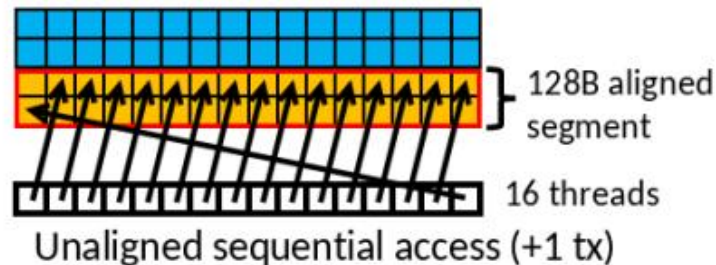


Coalesced memory access example

```
__global__ void memAccess(float *out, float *in)
{
    int tid = blockIdx.x*blockDim.x + threadIdx.x;
    if(tid != 12) out[tid + 16] = in[tid + 16];
}
```

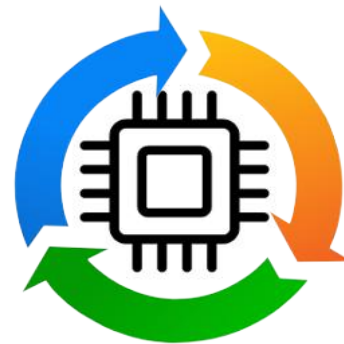


```
__global__ void memAccess(float *out, float *in)
{
    int tid = blockIdx.x*blockDim.x + threadIdx.x;
    out[tid + 1] = in[tid + 1];
}
```





HIP Kernel optimization and profiling



HIP: Libraries



NVIDIA	HIP	ROCm	Description
cuBLAS	hipBLAS	rocBLAS	Basic Linear Algebra Subroutines
cuFFT	hipFFT	rocfft	Fast fourier Transform Library
cuSPARSE	hipSPARSE	rocSPARSE	Sparse BLAS + SMV
cuSOLVER	hipSOLVER	rocSOLVER	Lapack library
AMG-X		rocALUTION	Sparse iterative solvers and preconditioners with Geometric and Algebraic MultiGrid
Thrust	hipThrust	rocThrust	C++ parallel algorithms library

HIP: Libraries



Libraries

NVIDIA	HIP	ROCm	Description
CUB	hipCUB	rocPRIM	Low level Optimized Parallel Primitives
cuDNN		MIOpen	Deep learning solver library
cuRAND	hipRAND	rocRAND	Random number generator library
EIGEN	EIGEN	EIGEN	C++ template library for linear algebra: matrices, vectors, numerical solvers
NCCL		RCCL	Communications Primitives Library based on the MPI equivalents

HIP: hipBLAS



CUDA

```
#include <cuda_runtime.h>
#include "cublas_v2.h"

if (cudaSuccess != cudaMalloc((void **) &a_dev,
sizeof(*a) * n * n) ||
    cudaSuccess != cudaMalloc((void **) &b_dev,
sizeof(*b) * n * n) ||
    cudaSuccess != cudaMalloc((void **) &c_dev,
sizeof(*c) * n * n)) {
    printf("error: memory allocation (CUDA)\n");
    cudaFree(a_dev); cudaFree(b_dev);
    cudaFree(c_dev);
    cudaDestroy(handle);
    exit(EXIT_FAILURE);
}
```

HIP

```
#include <hip/hip_runtime.h>
#include "hipblas.h"

if (hipSuccess != hipMalloc((void **) &a_dev,
sizeof(*a) * n * n) ||
    hipSuccess != hipMalloc((void **) &b_dev,
sizeof(*b) * n * n) ||
    hipSuccess != hipMalloc((void **) &c_dev,
sizeof(*c) * n * n)) {
    printf("error: memory allocation (CUDA)\n");

    hipFree(a_dev); hipFree(b_dev); hipFree(c_dev);
    hipblasDestroy(handle);
    exit(EXIT_FAILURE);
}
```

HIP: Kernels

You can call a kernel with the command:

```
hipLaunchKernelGGL(kernel_name, dim3(Blocks), dim3(Threads), 0, 0, arg1, arg2, ...);
```

or

```
kernel_name<<<dim3(Blocks), dim3(Threads),0,0>>>(arg1,arg2,...);
```

where blocks are for the 3D dimensions of the grid of blocks dimensions

threads for the 3D dimensions of a block of threads

0 for bytes of dynamic LDS space

0 for stream

kernel arguments

HIP: Metrics

Various useful metrics

GPUBusy: The percentage of time GPU was busy

Wavefronts: Total wavefronts

VALUInsts: The average number of vector ALU instructions executed per work-item (affected by flow control).

VALUUtilization: The percentage of active vector ALU threads in a wave. A lower number can mean either more thread divergence in a wave or that the work-group size is not a multiple of 64. Value range: 0% (bad), 100% (ideal - no thread divergence)

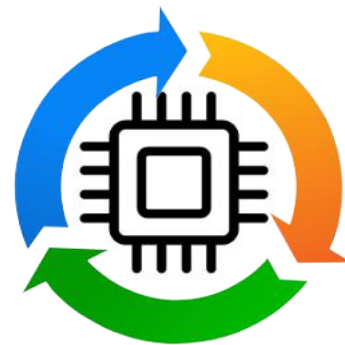
VALUBusy: The percentage of GPUTime vector ALU instructions are processed. Value range: 0% (bad) to 100% (optimal).

L2CacheHit: The percentage of fetch, write, atomic, and other instructions that hit the data in L2 cache. Value range: 0% (no hit) to 100% (optimal).

LDSBankConflict: The percentage of GPUTime LDS is stalled by bank conflicts. Value range: 0% (optimal) to 100% (bad).



HIP Multi-GPU programming and HIP+MPI



HIP: GPU Context

- Context is established when the first HIP function requiring an active context is called **hipMalloc()**
- Several processes can create contexts for a single device
- Resources are allocated per context
- By default, one context per device per process (since CUDA 4.0)
 - Threads of the same process share the primary context (for each device)
- Driver associates a number for each HIP-capable GPU starting from 0
- The function `hipSetDevice()` is used for selecting the desired device

HIP: Device Managment

- Return the number of hip capable devices in *count
`hipError_t hipGetDeviceCount(int *count)`
- Set device as the current device for the calling host thread
`hipError_t hipSetDevice(int device)`
- Return the current device for the calling host thread in *device
`hipError_t hipGetDevice(int *device)`
- Reset and explicitly destroy all resources associated with the current device
`hipError_t hipDeviceReset(void)`



HIP: Querying devices properties

- One can query the properties of different devices in the system using `hipGetDeviceProperties()` function
 - No context needed
 - Provides e.g. name, amount of memory, warp size, support for unified virtual addressing, etc.
 - Useful for code portability
- Return the properties of a HIP capable device in `*prop`
`hipError_t hipGetDeviceProperties(struct hipDeviceProp *prop, int device)`

HIP: Multi-GPU programming models

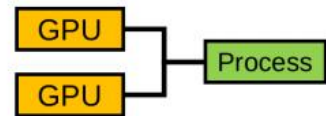
One GPU per process

- Syncing is handled through message passing (eg. MPI)



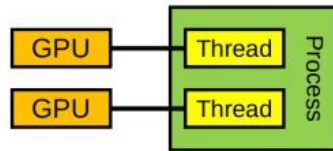
Many GPUs per process

- Process manages all context switching and syncing explicitly



One GPU per thread

- Syncing is handled through thread synchronization requirements



HIP: Multi-GPU, one GPU per process

- Recommended for multi-process applications using a message passing library
- Message passing library takes care of all GPU-GPU communication
- Each process interacts with only one GPU which makes the implementation easier and less invasive (if MPI is used anyway)
 - Apart from each process selecting a different device, the implementation looks much like a single-GPU program
- Multi-GPU implementation using MPI is discussed at the end!

HIP: Multi-GPU, many GPUs per process

- Process switches the active GPU using `hipSetDevice()` function
- After setting the device, HIP-calls such as the following are effective only on the selected GPU:
 - Memory allocations and copies
 - Streams and events
 - Kernel calls
- Asynchronous calls are required to overlap work across all devices

HIP: Multi-GPU, one GPU per process

- One GPU per CPU thread
 - I.e one OpenMP thread per GPU being used
- HIP API is threadsafe
 - Multiple threads can call the functions at the same time
- Each thread can create its own context on a different GPU
 - `hipSetDevice()` sets the device and creates a context per thread
 - Easy device management with no changing of device

- Communication between threads becomes a bit more tricky
- ```
#pragma omp parallel for
for(unsigned int i = 0; i < deviceCount; i++)
{
 hipSetDevice(i);
 kernel<<<blocks[i],threads[i]>>>(arg1[i], arg2[i]);
}
```

# HIP: Peer access

## Peer access

Access peer GPU memory directly from another GPU

- Pass a pointer to data on GPU 1 to a kernel running on GPU 0
- Transfer data between GPUs without going through host memory
- Lower latency, higher bandwidth



Check peer accessibility

```
hipError_t hipDeviceCanAccessPeer(int* canAccessPeer, int device, int peerDevice)
```

Enable peer access

```
hipError_t hipDeviceEnablePeerAccess(int peerDevice, unsigned int flags)
```

Disable peer access

```
hipError_t hipDeviceDisablePeerAccess(int peerDevice)
```



# HIP: Peer to peer communication

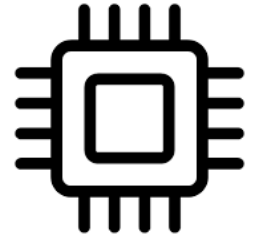
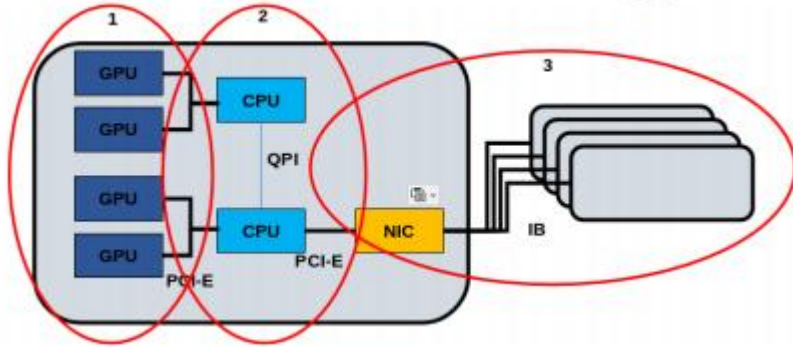
- Devices have separate memories
- With devices supporting unified virtual addressing, `hipMemCpy()` with `kind=hipMemcpyDefault`, works:

```
hipError_t hipMemcpy(void* dst, void* src, size_t count, hipMemcpyKind kind)
```

- Other option which does not require unified virtual addressing

```
hipError_t hipMemcpyPeer(void* dst, int dstDev, void* src, int srcDev, size_t count)
```
- If peer to peer access is not available, the functions result in a normal copy through host memory

# HIP: Three levels of parallelism



1. GPU - GPU threads on the multiprocessors

Parallelization strategy: HIP, SYCL, Kokkos, OpenMP

2. Node - Multiple GPUs and CPUs

Parallelization strategy: MPI, Threads, OpenMP

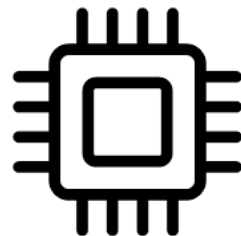
3. Supercomputer - Many nodes connected with interconnect

Parallelization strategy: MPI between nodes

# HIP: Strategies

## MPI+HIP strategies

1. One MPI process per node
2. One MPI process per GPU
3. Many MPI processes per GPU, only one uses it
4. Many MPI processes sharing a GPU
  - 2 is recommended (also allows using 4 with services such as CUDA MPS)
  - Typically results in most productive and least invasive implementation for an MPI program
  - No need to implement GPU-GPU transfers explicitly (MPI handles all this)
  - It is further possible to utilize remaining CPU cores with OpenMP (but this is not always worth the effort/increased complexity)



# HIP: Strategies

## Selecting the correct GPU

Typically all processes on the node can access all GPUs of that node.

The following implementation allows utilizing all GPUs using one or more processes per GPU.

- Use CUDA MPS when launching more processes than GPUs

```
int deviceCount, nodeRank;
MPI_Comm commNode;
MPI_Comm_split_type(MPI_COMM_WORLD, MPI_COMM_TYPE_SHARED, 0, MPI_INFO_NULL,
&commNode);
MPI_Comm_rank(commNode, &nodeRank);
hipGetDeviceCount(&deviceCount);
hipSetDevice(nodeRank % deviceCount);
```

# HIP: Strategies

## GPU-GPU communication through MPI

CUDA/ROCm aware MPI libraries support direct GPU-GPU transfers

- Can take a pointer to device buffer (avoids host/device data copies)

Unfortunately, currently no GPU support for custom MPI datatypes

(must use a datatype representing a contiguous block of memory)

- Data packing/unpacking must be implemented application-side on GPU

ROCm aware MPI libraries are under development and there may be problems

- It is a good idea to have a fallback option to use pinned host staging buffers

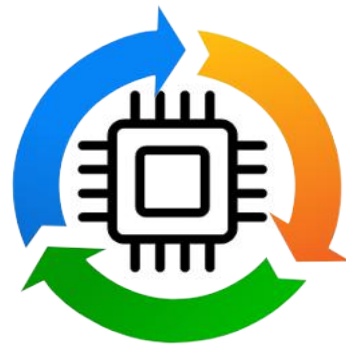
# HIP: Strategies

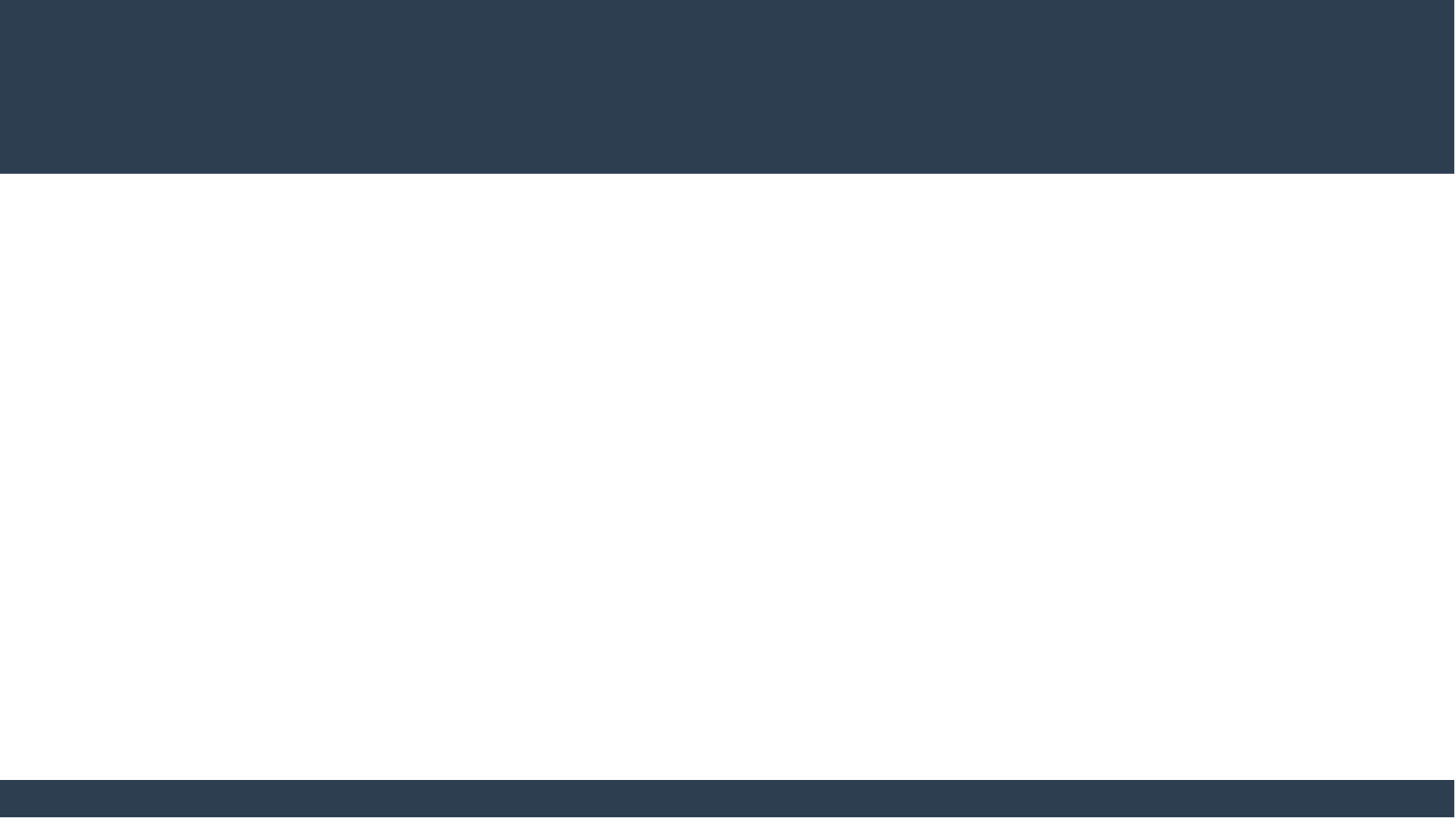
Many options to write a multi-GPU program:

- Use `hipSetDevice()` to select the device, and the subsequent HIP calls operate on that device
- If you have an MPI program, it is often best to use one GPU per process, and let MPI handle data transfers between GPUs
- There is still little experience from ROCm aware MPIs, there may be issues
- Note that a CUDA/ROCm aware MPI is only required when passing device pointers to the MPI, passing only host pointers does not require any CUDA/ROCm awareness



**HIP Examples**









Thank you for your attention !

