# AMD HARDWARE AND SOFTWARE

SUYASH TANDON, JUSTIN CHANG, JULIO MAIA, NOEL CHALMERS, PAUL BAUMAN, NICHOLAS CURTIS, NICHOLAS MALAYA, ALESSANDRO FANFARILLO, JOSE NOUDOHOUENOU, CHIP FREITAG, DAMON MCDOUGALL, NOAH WOLFE, SAMUEL ANTAO, GEORGE MARKOMANOLIS, BOB ROBEY, GINA SITARAMAN

JAKUB KURZAK - PRESENTER

ADVANCED MICRO DEVICES, INC.

**AMD**
together we advance_

slides on LUMI in /project/project_465000644/Slides/AMD/

hands-on exercises: https://hackmd.io/@sfantao/H1QU6xRR3

hands-on source code: /project/project_465000644/Exercises/AMD/HPCTrainingExamples/

**AMD**
together we advance_

# AMD HARDWARE FOR HPC AND AI
## CDNA ARCHITECTURE
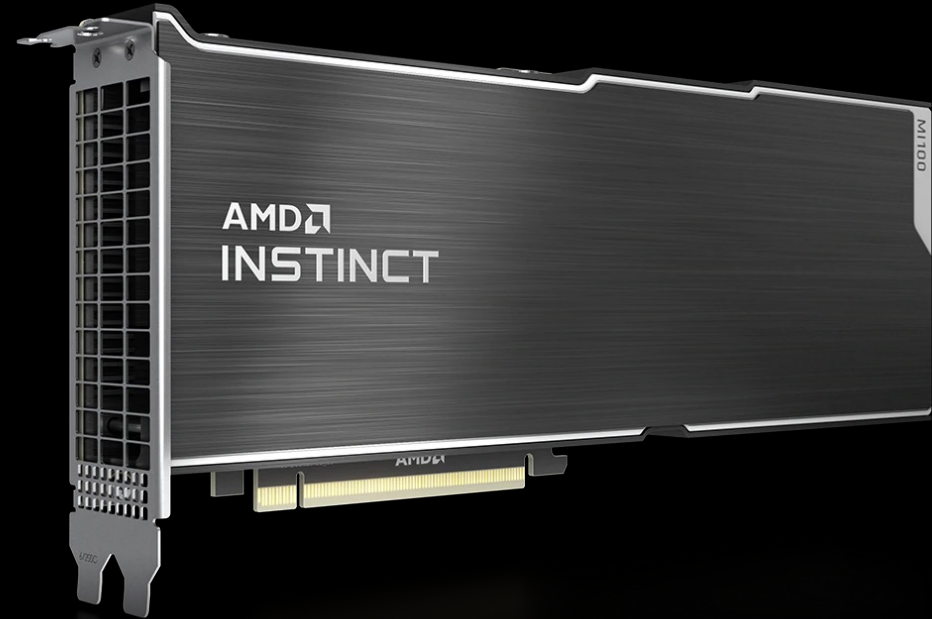
AMD
together we advance_

# AMD GPUS

Radeon™ Graphics Cards
RDNA architecture
E.g.:
o   RX 6000 Series
o   RX 7000 Series

AMD Instinct™ Accelerators
CDNA architecture
E.g.:
o   MI100
o   MI200
o   MI300

4

together we advance_
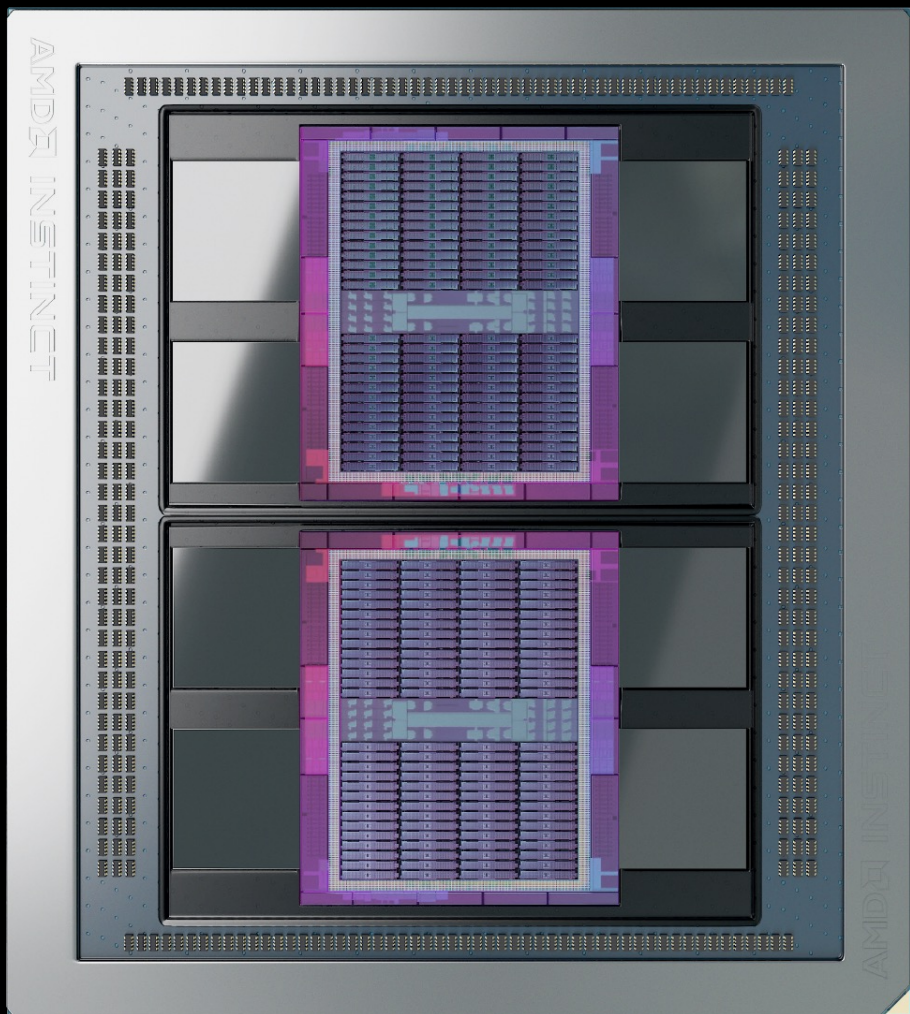
# AMD IN HPC



## Frontier@ORNL

o currently the largest machine in the world

o the first computer to cross 1 exaFLOPS

o AMD EPYC CPUs

o AMD Instinct GPUs



## LUMI@CSC

o currently the largest machine in Europe

o 3rd fastest in the world

o AMD EPYC CPUs

o AMD Instinct GPUs

**AMD**
together we advance_

# AMD INSTINCT™ MI200

## AMD INSTINCT™ MI250X

# WORLD'S MOST ADVANCED DATA CENTER ACCELERATOR

**58B**
Transistors in 6nm
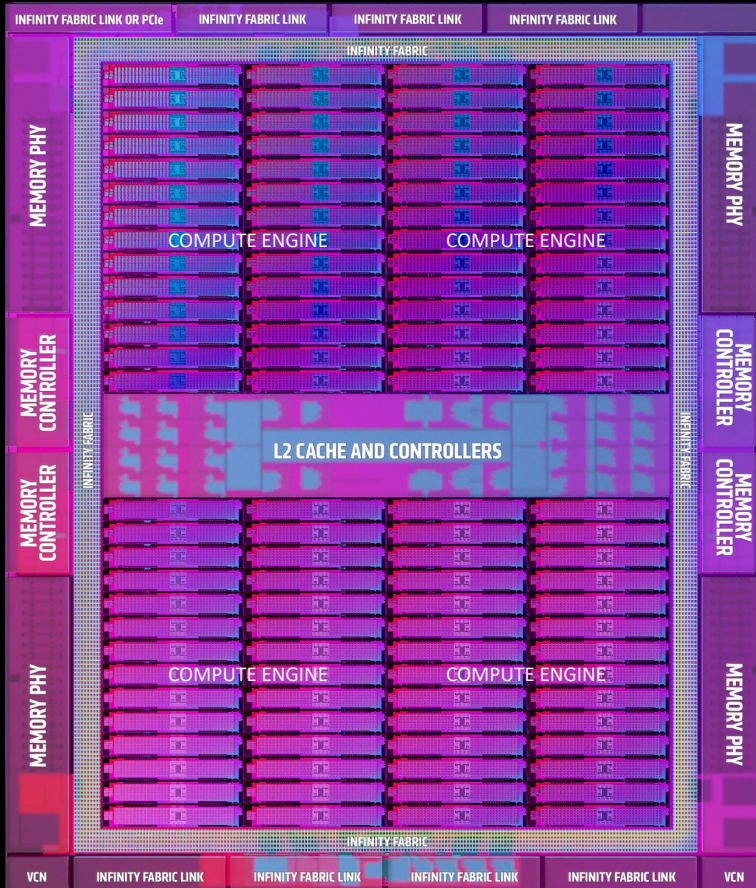
**220**
Compute Units

**880**
2nd Gen Matrix Cores

**128**
GB HBM2E @ 3.2 TB/s

https://www.amd.com/system/files/documents/amd-cdna2-white-paper.pdf

AMD
together we advance_

# AMD INSTINCT™ MI200



## 2ND GENERATION CDNA ARCHITECTURE
## TAILORED-BUILT FOR HPC & AI

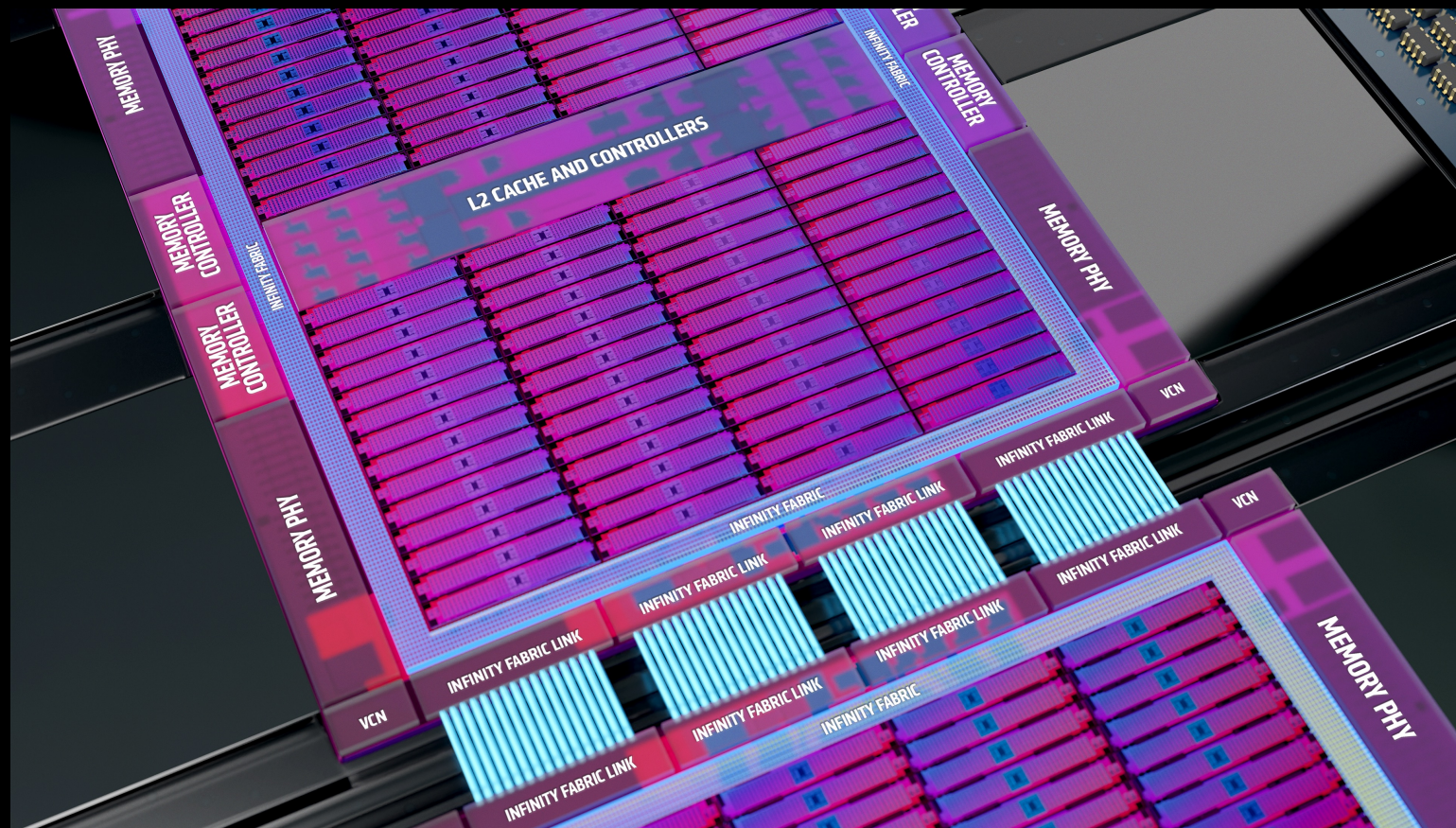| | |
|---|---|
| TSMC 6NM TECHNOLOGY | UP TO 110 CU PER GRAPHICS CORE DIE |
| 4 MATRIX CORES PER COMPUTE UNIT | MATRIX CORES ENHANCED FOR HPC |
| 8 INFINITY FABRIC LINKS PER DIE | SPECIAL FP32 OPS FOR DOUBLE THROUGHPUT |

AMD

together we advance_

# MULTI-CHIP DESIGN

## TWO GPU DIES IN PACKAGE TO MAXIMIZE COMPUTE & DATA THROUGHPUT
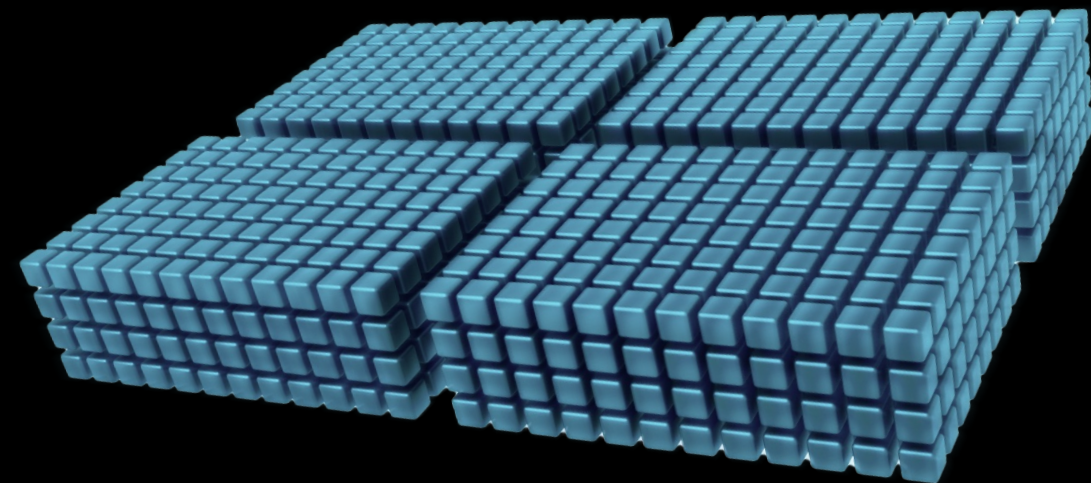
INFINITY FABRIC FOR CROSS-DIE CONNECTIVITY

4 LINKS RUNNING AT 25GBPS

400GB/S OF BI-DIRECTIONALBANDWIDTH



8

AMD
together we advance_

# 2nd GENERATION MATRIX CORES

## OPTIMIZED COMPUTE UNITS FOR SCIENTIFIC COMPUTING

DOUBLE PRECISON (FP64)
MATRIX CORE THROUGHPUT
REPRESENTATION

| MI100 MATRIX CORES | MI250X MATRIX CORES |
|---|---|
| OPS/CLOCK/COMPUTE UNIT | OPS/CLOCK/COMPUTE UNIT |
| No FP64 Matrix Core | 256 FP64 |
| 256 FP32 | 256 FP32 |
| 1024 FP16 | 1024 FP16 |
| 512 BF16 | 1024 BF16 |
| 512 INT8 | 1024 INT8 |

AMD
together we advance_

# AMD INSTINCT™ MI200

# MI200 COMPUTE UNIT

scheduler

scalar unit

Local Data Share (64KB)

scalar registers (8KB)

L1 cache (16KB)

SIMD unit

SIMD unit

**each SIMD unit**

o   has 16 SIMD lanes

o   operates on vectors (waves) of size 64

o   handles up to 10 waves simultaneously

AMD

together we advance_

# AMD INSTINCT™ MI300



The world's first integrated
data center CPU + GPU

AMD INSTINCT™

# MI300

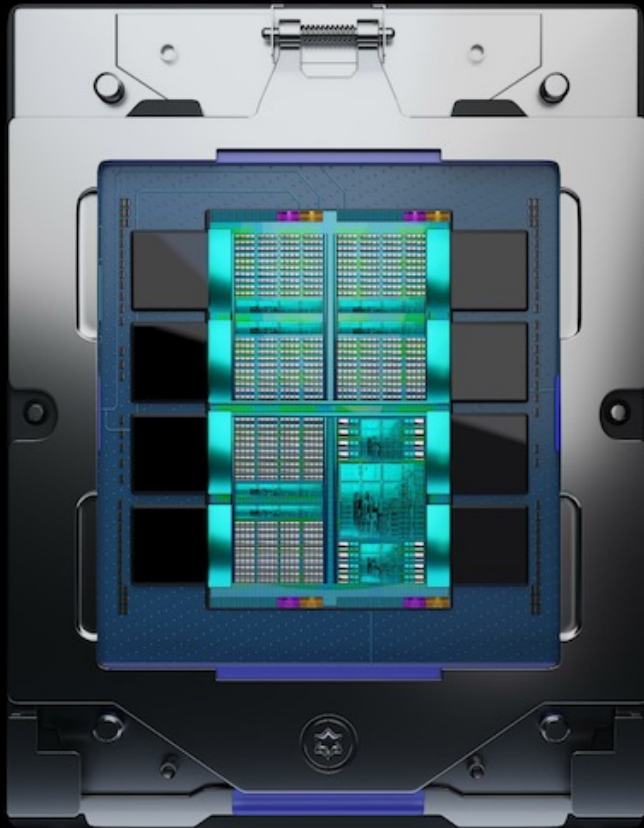Breakthrough architecture to
power the exascale AI era

AMD
together we advance_

# UNIFIED MEMORY APU ARCHITECTURE BENEFITS

## AMD CDNA™ 2 Coherent Memory Architecture → AMD CDNA™ 3 Unified Memory APU Architecture

- Simplifies Programming

- Low Overhead 3rd Gen Infinity Interconnect

- Industry Standard Modular Design

**CPU**

**GPU**

**CPU Memory (DRAM)**

**GPU Memory (HBM)**

- Eliminates Redundant Memory Copies

- High-Efficiency 4th Gen AMD Infinity Architecture

- Low TCO with Unified Memory APU Package

**Next-Gen AMD Instinct™ APU**

**Unified Memory (HBM)**

# AMD SOFTWARE FOR HPC AND AI
## ROCM PLATFORM

AMD
together we advance_

# AMD ROCm™ Open Software Platform For GPU Compute

together we advance_

# Open Software Platform For GPU Compute

**AMD ROCm**

- Unlocked GPU Power To Accelerate Computational Tasks

- Optimized for HPC and Deep Learning Workloads at Scale

- Open Source Enabling Innovation, Differentiation, and Collaboration

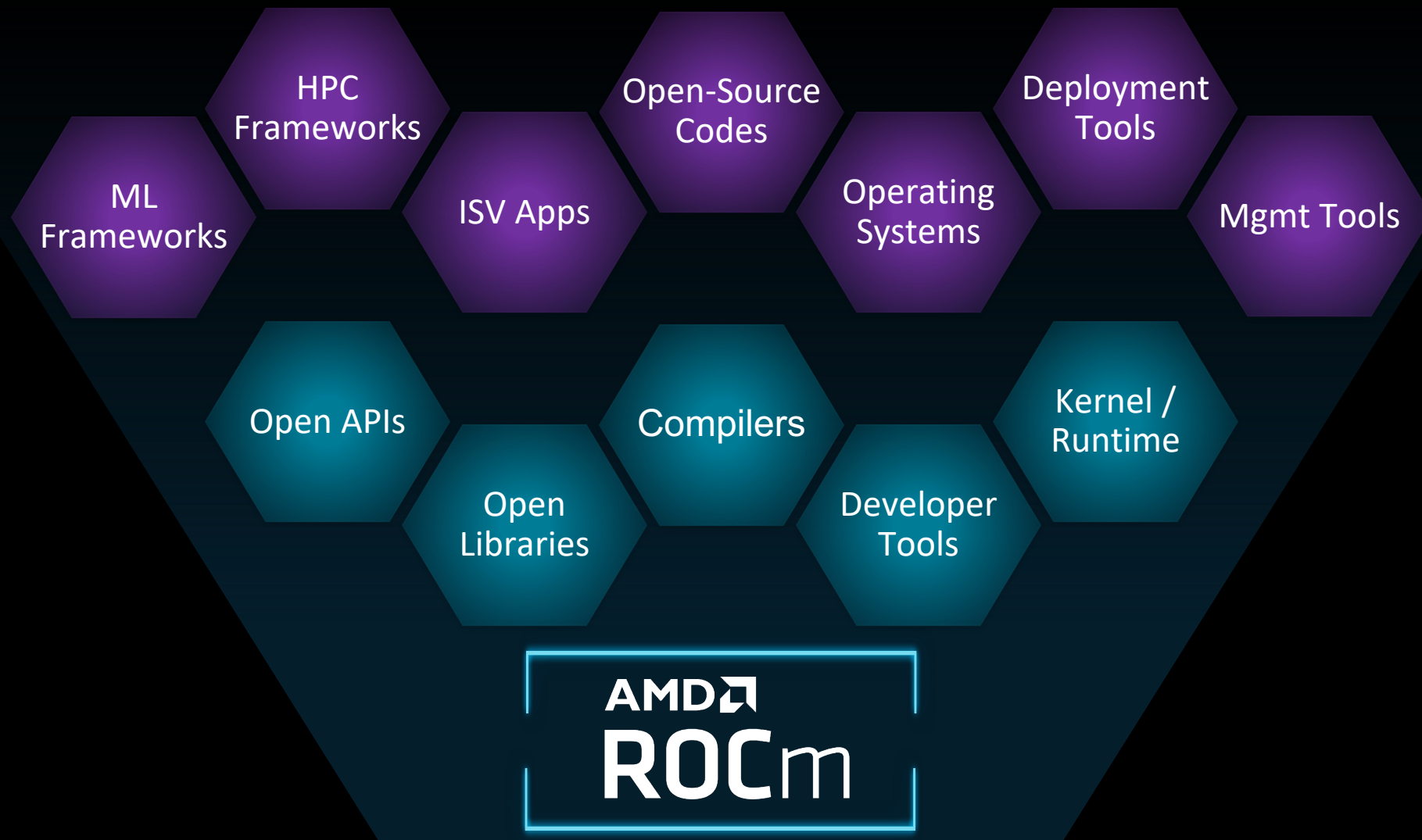| | |
|---|---|
| Benchmarks & App Support | Optimized Training/Inference Models & Applications |
| | MLPERF · HPL/HPCG · Life Science · Geo Science · Physics |
| Operating Systems Support | RHEL · CentOS · SLES · Ubuntu |
| Cluster Deployment | Singularity · Kubernetes® · Docker® · SLURM |
| Framework Support | Kokkos/RAJA · PyTorch · TensorFlow |
| Libraries | BLAS · RAND · FFT · MIGraphX · MIVisionX · PRIM · SOLVER · ALUTION · SPARSE · THRUST · MIOpen · RCCL |
| Programming Models | OpenMP® API · OpenCL™ · HIP API |
| Development Toolchain | Compiler · Profiler · Tracer · Debugger · hipify · GPUFort |
| Drivers & Runtime | GPU Device Drivers and ROCm Run-Time |
| Deployment Tools | ROCm Validation Suite · ROCm Data Center Tool · ROCm SMI |

AMD together we advance_

# AMD

# ROCm 5.0

## DEMOCRATIZING EXASCALE FOR ALL

| EXPANDING SUPPORT & ACCESS | OPTIMIZING PERFORMANCE | ENABLING DEVELOPER SUCCESS |
|---|---|---|
| • Support for Radeon Pro W6800 Workstation GPUs<br><br>• Remote access through the AMD Accelerator Cloud | • MI200 Optimizations: FP64 Matrix ops, Improved Cache<br><br>• Improved launch latency and kernel performance | • HPC Apps & ML Frameworks on AMD InfinityHub<br><br>• Streamlined and improved tools increasing productivity |

AMD together we advance_

# LIBRARIES

**rocBLAS / hipBLAS**
- basic operations on dense matrices

**rocSOLVER**
- dense linear algebra solvers

**rocSPARSE / hipSPARSE**
- basic operations on sparse matrices

**rocALUTION**
- sparse linear algebra solvers

**rocFFT / hipFFT**
- Fast Fourier transforms

**rocRAND / hipRAND**
- random number generation

**rocPRIM / hipCUB / rocThrust**
- scan, sort, reduction, etc.

https://github.com/ROCmSoftwarePlatform/rocBLAS
https://github.com/ROCmSoftwarePlatform/hipBLAS

https://github.com/ROCmSoftwarePlatform/rocSOLVER

https://github.com/ROCmSoftwarePlatform/rocSPARSE
https://github.com/ROCmSoftwarePlatform/hipSPARSE

https://github.com/ROCmSoftwarePlatform/rocALUTION

https://github.com/ROCmSoftwarePlatform/rocFFT
https://github.com/ROCmSoftwarePlatform/hipFFT

https://github.com/ROCmSoftwarePlatform/rocRAND
https://github.com/ROCmSoftwarePlatform/hipRAND

https://github.com/ROCmSoftwarePlatform/rocPRIM
https://github.com/ROCmSoftwarePlatform/hipCUB
https://github.com/ROCmSoftwarePlatform/rocThrust

**AMD**
together we advance_

# ALSO OPEN SOURCE

**the compiler**
o https://github.com/ROCmSoftwarePlatform/llvm-project

**the runtime**
o https://github.com/RadeonOpenCompute/ROCR-Runtime

**the debugger**
o https://github.com/ROCm-Developer-Tools/ROCgdb

**the profiler**
o https://github.com/ROCm-Developer-Tools/rocprofiler

**the HPL benchmark**
o https://github.com/ROCmSoftwarePlatform/rocHPL

**the HPCG benchmark**
o https://github.com/ROCmSoftwarePlatform/rocHPCG

etc.

**AMD**
together we advance_

# AMD SOFTWARE FOR HPC AND AI
## HIP PROGRAMMING

AMD
together we advance_
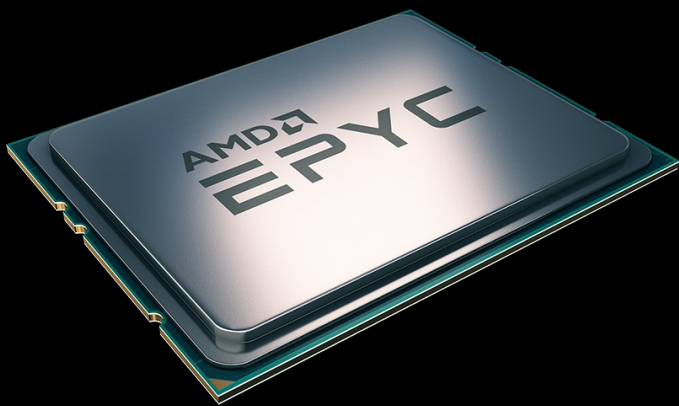
# GPU ACCELERATION
## HOST AND DEVICE

**the host is the CPU**

o   host code runs here

o   usual C++ syntax and features

o   entry point is the "main" function

o   use the HIP API to

   o   create device buffers

   o   moved data between host and device

   o   launch device code

**the device is the GPU**

o   device code runs here

o   C/C++ syntax and features

o   device code is launched as "kernels"

o   instructions from the host are sent to streams

# FUNCTION QUALIFIERS
## HOST AND DEVICE

**`__global__`**

o  "kernels"

o  execute the GPU

o  can be called from the CPU


**`__device__`**

o  execute the GPU

o  can be called from device code (kernels or a `__device__` functions)


**`__host__ __device__`**

o  executes on the CPU when called from CPU code

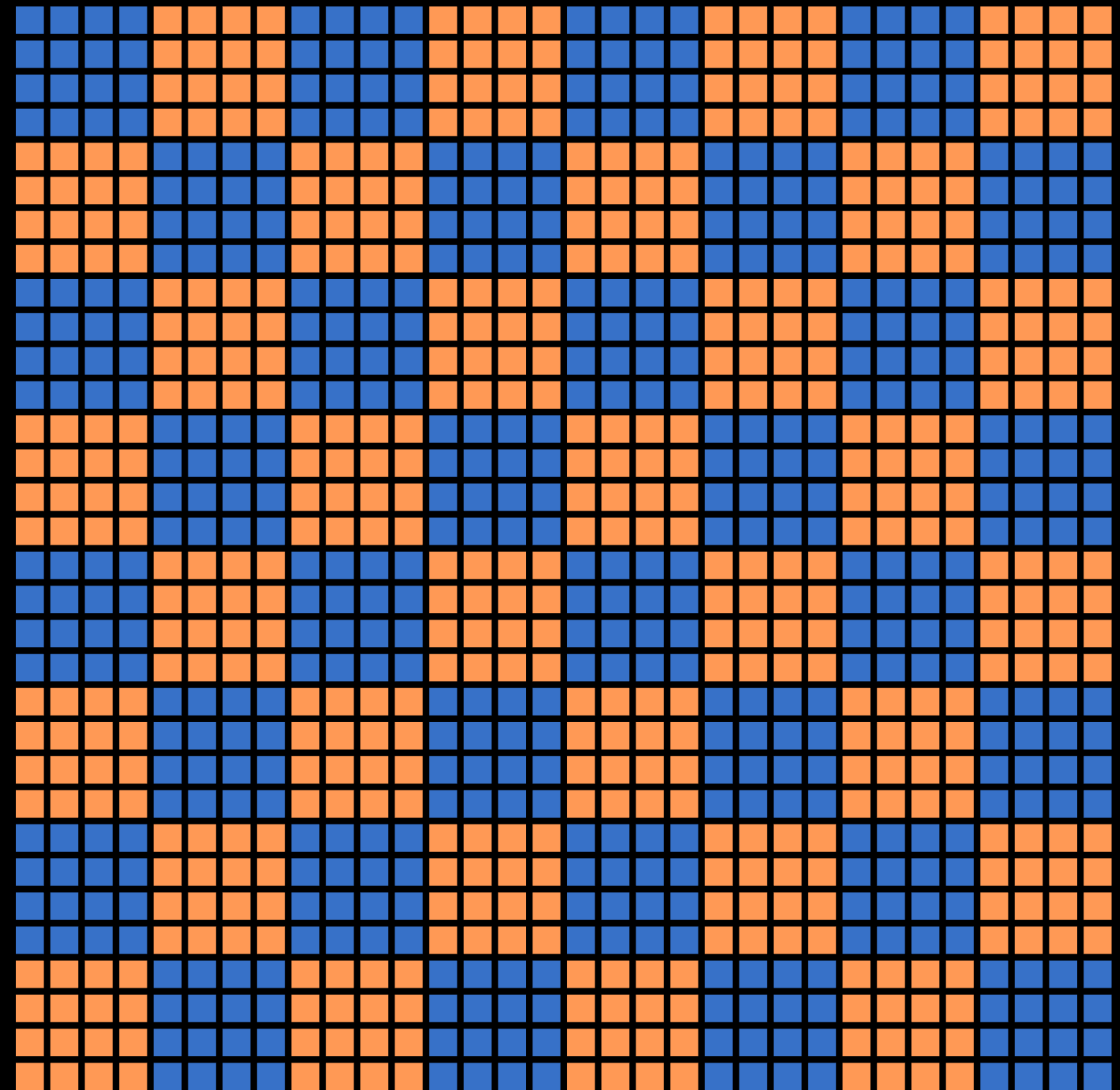o  executes on the GPU when called from GPU code

AMD
together we advance_

# HIP KERNEL LANGUAGE
## GPU CODE

**in 2D**

o   each colored box is a block

o   each block has an index - `blockIdx.[xyz]`

o   each small square is a thread

o   each thread has a 2D index - `threadIdx.[xyz]`
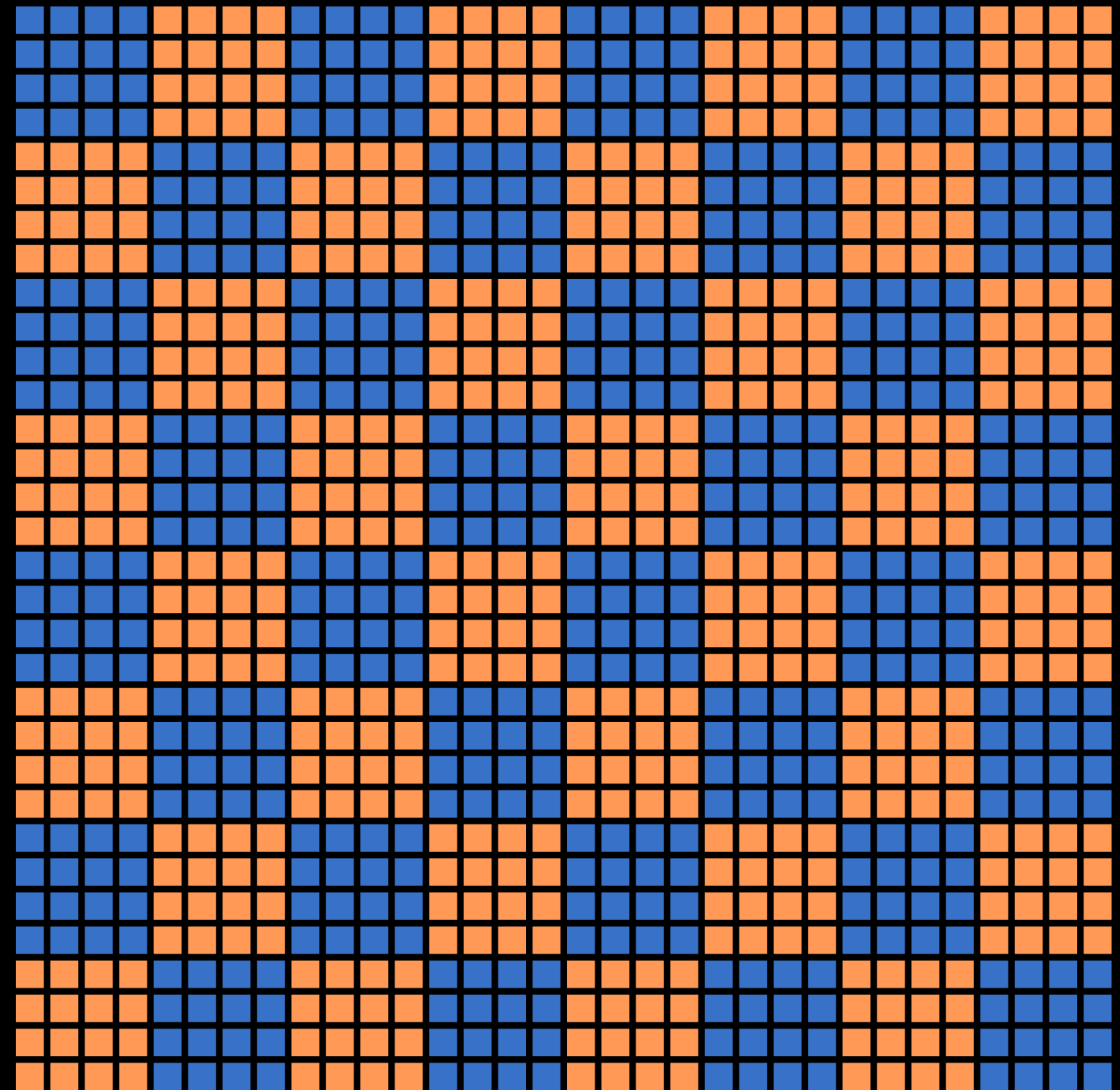
o   grid dimensions in - `blockDim.[xyz]`

**in 2D**

https://rocm.docs.amd.com/projects/HIP/en/latest/reference/kernel_language.html

together we advance_

# HIP KERNEL LANGUAGE
## GPU CODE

**in 2D**

o   all local variables and arrays are thread-private

o   threads can exchange data through shared memory (LDS)

o   declare using the `__shared__` keyword

o   use `__syncthreads()` to synchronize

https://rocm.docs.amd.com/projects/HIP/en/latest/reference/kernel_language.html

**AMD**
together we advance_

# HIP KERNEL LANGUAGE
## GPU CODE

**saxpy loop**

o   two 1D arrays

o   the `y[i] += a*x[i]` operation

o   mapped to 1D grid of threads/blocks

o   each thread takes on index

```
1    #include <cuda.h>
2
3    __constant__ float a = 2.0f;
4
5    __global__
6    void saxpy(int n, float const* x, float* y)
7    {
8        int i = blockDim.x*blockIdx.x + threadIdx.x;
9        if (i < n)
10           y[i] += a*x[i];
11   }
```

https://rocm.docs.amd.com/projects/HIP/en/latest/reference/kernel_language.html

AMD

together we advance_

# HIP API
## MEMORY MANAGEMENT

```
hipError_t   hipMalloc (void **ptr, size_t size)
```

```
hipError_t   hipFree (void *ptr)
```
Free memory allocated by the hcc hip memory allocation API. This API performs an implicit hipDeviceSynchronize() call. If pointer is NULL, the hip runtime is initialized and hipSuccess is returned. More...

```
hipError_t   hipMemcpy (void *dst, const void *src, size_t sizeBytes, hipMemcpyKind kind)
```
Copy data from src to dst. More...

o  GPU operates on GPU memory

o  need to allocate GPU memory

o  need to copy data between the CPU memory and the GPU memory

https://rocm.docs.amd.com/projects/HIP/en/latest/.doxygen/docBin/html/group___memory.html

AMD

together we advance_

# HIP API
## ERROR HANDLING

o check last error

o get error name

o get error string

| | | |
|---|---|---|
| hipError_t | hipGetLastError (void) | |

Return last error returned by any HIP runtime API call and resets the stored error code to hipSuccess. More...

| | |
|---|---|
| hipError_t | hipPeekAtLastError (void) |

Return last error returned by any HIP runtime API call. More...

| | |
|---|---|
| const char * | hipGetErrorName (hipError_t hip_error) |

Return hip error as text string form. More...

| | |
|---|---|
| const char * | hipGetErrorString (hipError_t hipError) |

Return handy text string message to explain the error which occurred. More...

https://rocm.docs.amd.com/projects/HIP/en/latest/.doxygen/docBin/html/group___error.html

**AMD**
together we advance_

# HIP API
## DEVICE MANAGEMENT

o   check number of devices

o   switch devices

o   synchronize devices

| | |
|---|---|
| `hipError_t` | `hipDeviceSynchronize (void)` |

Waits on all active streams on current device. More...

| | |
|---|---|
| `hipError_t` | `hipDeviceReset (void)` |

The state of current device is discarded and updated to a fresh state. More...

| | |
|---|---|
| `hipError_t` | `hipSetDevice (int deviceId)` |

Set default device to be used for subsequent hip API calls from this thread. More...

| | |
|---|---|
| `hipError_t` | `hipGetDevice (int *deviceId)` |

Return the default device id for the calling host thread. More...

| | |
|---|---|
| `hipError_t` | `hipGetDeviceCount (int *count)` |

Return number of compute-capable devices. More...

https://rocm.docs.amd.com/projects/HIP/en/latest/.doxygen/docBin/html/group___device.html

AMD
together we advance_

# HIP API
## STREAM MANAGEMENT

- o create stream

- o destroy stream

- o synchronize stream

- o etc.

- o etc.

- o etc.

```
hipError_t   hipStreamCreate (hipStream_t *stream)
```

Create an asynchronous stream. **More...**

```
hipError_t   hipStreamDestroy (hipStream_t stream)
```

Destroys the specified stream. **More...**

```
hipError_t   hipStreamSynchronize (hipStream_t stream)
```

Wait for all commands in stream to complete. **More...**

https://rocm.docs.amd.com/projects/HIP/en/latest/.doxygen/docBin/html/group___stream.html

AMD
together we advance_

# AMD LINGO

| CUDA lingo | | AMD lingo |
|:---:|:---:|:---:|
| block | ⟶ | work group |
| thread | ⟶ | work item |
| warp | ⟶ | wavefront |

AMD
together we advance_

# SIMPLE SAXPY KERNEL

```cuda
1    #include <cuda.h>
2
3    __constant__ float a = 2.0f;
4
5    __global__
6    void saxpy(int n, float const* x, float* y)
7    {
8        int i = blockDim.x*blockIdx.x + threadIdx.x;
9        if (i < n)
10           y[i] += a*x[i];
11   }
```

o   vector addition kernel in CUDA

o   each thread takes one array index

o   and performs one multiply-and-add operation

**AMD**
together we advance_

[Public]

```
 1   #include <cuda.h>
 2
 3   __constant__ float a = 2.0f;
 4
 5   __global__
 6   void saxpy(int n, float const* x, float* y)
 7   {
 8       int i = blockDim.x*blockIdx.x + threadIdx.x;
 9       if (i < n)
10           y[i] += a*x[i];
11   }
12
13   int main()
14   {
15       int n = 256;
16       std::size_t size = sizeof(float)*n;
17
18       float* d_x;
19       float* d_y;
20       cudaMalloc(&d_x, size);
21       cudaMalloc(&d_y, size);
22
23       int num_blocks = 2;
24       int num_threads = 128;
25       saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
26       cudaDeviceSynchronize();
27   }
28
```

allocate arrays in device memory

set up the grid
launch the kernel

32

AMD
together we advance_

# ADDING HOST↔DEVICE COPIES

```cuda
1   #include <cuda.h>
2
3   __constant__ float a = 2.0f;
4
5   __global__
6   void saxpy(int n, float const* x, float* y)
7   {
8       int i = blockDim.x*blockIdx.x + threadIdx.x;
9       if (i < n)
10          y[i] += a*x[i];
11  }
12
13  int main()
14  {
15      int n = 256;
16      std::size_t size = sizeof(float)*n;
17
18      float* h_x = (float*)malloc(size);
19      float* h_y = (float*)malloc(size);
20
21      float* d_x;
22      float* d_y;
23      cudaMalloc(&d_x, size);
24      cudaMalloc(&d_y, size);
25
26      cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice);
27      cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice);
28
29      int num_blocks = 2;
30      int num_threads = 128;
31      saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
32
33      cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost);
34      cudaDeviceSynchronize();
35  }
36  |
```

allocate arrays in host memory

copy content to device memory

copy results back to host memory

AMD

together we advance_

```
 1  #include <cuda.h>
 2
 3  __constant__ float a = 2.0f;
 4
 5  __global__
 6  void saxpy(int n, float const* x, float* y)
 7  {
 8      int i = blockDim.x*blockIdx.x + threadIdx.x;
 9      if (i < n)
10          y[i] += a*x[i];
11  }
12
13  int main()
14  {
15      int n = 256;
16      std::size_t size = sizeof(float)*n;
17
18      float* h_x = (float*)malloc(size);
19      float* h_y = (float*)malloc(size);
20
21      float* d_x;
22      float* d_y;
23      cudaMalloc(&d_x, size);
24      cudaMalloc(&d_y, size);
25
26      cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice);
27      cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice);
28
29      int num_blocks = 2;
30      int num_threads = 128;
31      saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
32
33      cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost);
34      cudaDeviceSynchronize();
35
36      cudaFree(d_x);
37      cudaFree(d_y);
38
39      free(h_x);
40      free(h_y);
41  }
42
```

free arrays in device memory

free arrays in host memory

AMD 🔺
together we advance_

```
1    #include <cuda.h>
2    #include <cassert>
3
4    __constant__ float a = 2.0f;
5
6    __global__
7    void saxpy(int n, float const* x, float* y)
8    {
9        int i = blockDim.x*blockIdx.x + threadIdx.x;
10       if (i < n)
11           y[i] += a*x[i];
12   }
13
14   #define CHECK(call) assert(call == cudaSuccess)
15
16   int main()
17   {
18       int n = 256;
19       std::size_t size = sizeof(float)*n;
20
21       float* h_x = (float*)malloc(size);
22       float* h_y = (float*)malloc(size);
23       assert(h_x != nullptr);
24       assert(h_y != nullptr);
25
26       float* d_x;
27       float* d_y;
28       CHECK(cudaMalloc(&d_x, size));
29       CHECK(cudaMalloc(&d_y, size));
30
31       CHECK(cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice));
32       CHECK(cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice));
33
34       int num_blocks = 2;
35       int num_threads = 128;
36       saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
37
38       CHECK(cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost));
39       CHECK(cudaDeviceSynchronize());
40
41       CHECK(cudaFree(d_x));
42       CHECK(cudaFree(d_y));
43
44       free(h_x);
45       free(h_y);
46   }
47
```

← simple error checking macro

**AMD**
together we advance_

# simple CUDA code

```cpp
#include <cuda.h>
#include <cassert>

__constant__ float a = 2.0f;

__global__
void saxpy(int n, float const* x, float* y)
{
    int i = blockDim.x*blockIdx.x + threadIdx.x;
    if (i < n)
        y[i] += a*x[i];
}

#define CHECK(call) assert(call == cudaSuccess)

int main()
{
    int n = 256;
    std::size_t size = sizeof(float)*n;

    float* h_x = (float*)malloc(size);
    float* h_y = (float*)malloc(size);
    assert(h_x != nullptr);
    assert(h_y != nullptr);

    float* d_x;
    float* d_y;
    CHECK(cudaMalloc(&d_x, size));
    CHECK(cudaMalloc(&d_y, size));

    CHECK(cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice));
    CHECK(cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice));

    int num_blocks = 2;
    int num_threads = 128;
    saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);

    CHECK(cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost));
    CHECK(cudaDeviceSynchronize());

    CHECK(cudaFree(d_x));
    CHECK(cudaFree(d_y));

    free(h_x);
    free(h_y);
}
```

AMD
together we advance_

# simple CUDA code

```cpp
1    #include <cuda.h>
2    #include <cassert>
3
4    __constant__ float a = 2.0f;
5
6    __global__
7    void saxpy(int n, float const* x, float* y)
8    {
9        int i = blockDim.x*blockIdx.x + threadIdx.x;
10       if (i < n)
11           y[i] += a*x[i];
12   }
13
14   #define CHECK(call) assert(call == cudaSuccess)
15
16   int main()
17   {
18       int n = 256;
19       std::size_t size = sizeof(float)*n;
20
21       float* h_x = (float*)malloc(size);
22       float* h_y = (float*)malloc(size);
23       assert(h_x != nullptr);
24       assert(h_y != nullptr);
25
26       float* d_x;
27       float* d_y;
28       CHECK(cudaMalloc(&d_x, size));
29       CHECK(cudaMalloc(&d_y, size));
30
31       CHECK(cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice));
32       CHECK(cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice));
33
34       int num_blocks = 2;
35       int num_threads = 128;
36       saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
37
38       CHECK(cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost));
39       CHECK(cudaDeviceSynchronize());
40
41       CHECK(cudaFree(d_x));
42       CHECK(cudaFree(d_y));
43
44       free(h_x);
45       free(h_y);
46   }
47
```

# same code in HIP

```cpp
1    #include <hip/hip_runtime.h>
2    #include <cassert>
3
4    __constant__ float a = 2.0f;
5
6    __global__
7    void saxpy(int n, float const* x, float* y)
8    {
9        int i = blockDim.x*blockIdx.x + threadIdx.x;
10       if (i < n)
11           y[i] += a*x[i];
12   }
13
14   #define CHECK(call) assert(call == hipSuccess)
15
16   int main()
17   {
18       int n = 256;
19       std::size_t size = sizeof(float)*n;
20
21       float* h_x = (float*)malloc(size);
22       float* h_y = (float*)malloc(size);
23       assert(h_x != nullptr);
24       assert(h_y != nullptr);
25
26       float* d_x;
27       float* d_y;
28       CHECK(hipMalloc(&d_x, size));
29       CHECK(hipMalloc(&d_y, size));
30
31       CHECK(hipMemcpy(d_x, h_x, size, hipMemcpyHostToDevice));
32       CHECK(hipMemcpy(d_y, h_y, size, hipMemcpyHostToDevice));
33
34       int num_blocks = 2;
35       int num_threads = 128;
36       saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
37
38       CHECK(hipMemcpy(h_y, d_y, size, hipMemcpyDeviceToHost));
39       CHECK(hipDeviceSynchronize());
40
41       CHECK(hipFree(d_x));
42       CHECK(hipFree(d_y));
43
44       free(h_x);
45       free(h_y);
46   }
47
```

spot the differences

AMD
together we advance_

## simple CUDA code

```
1    #include <cuda.h>
2    #include <cassert>
3
4    __constant__ float a = 2.0f;
5
6    __global__
7    void saxpy(int n, float const* x, float* y)
8    {
9        int i = blockDim.x*blockIdx.x + threadIdx.x;
10       if (i < n)
11           y[i] += a*x[i];
12   }
13
14   #define CHECK(call) assert(call == cudaSuccess)
15
16   int main()
17   {
18       int n = 256;
19       std::size_t size = sizeof(float)*n;
20
21       float* h_x = (float*)malloc(size);
22       float* h_y = (float*)malloc(size);
23       assert(h_x != nullptr);
24       assert(h_y != nullptr);
25
26       float* d_x;
27       float* d_y;
28       CHECK(cudaMalloc(&d_x, size));
29       CHECK(cudaMalloc(&d_y, size));
30
31       CHECK(cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice));
32       CHECK(cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice));
33
34       int num_blocks = 2;
35       int num_threads = 128;
36       saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
37
38       CHECK(cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost));
39       CHECK(cudaDeviceSynchronize());
40
41       CHECK(cudaFree(d_x));
42       CHECK(cudaFree(d_y));
43
44       free(h_x);
45       free(h_y);
46   }
47
```

## same code in HIP

```
1    #include <hip/hip_runtime.h>
2    #include <cassert>
3
4    __constant__ float a = 2.0f;
5
6    __global__
7    void saxpy(int n, float const* x, float* y)
8    {
9        int i = blockDim.x*blockIdx.x + threadIdx.x;
10       if (i < n)
11           y[i] += a*x[i];
12   }
13
14   #define CHECK(call) assert(call == hipSuccess)
15
16   int main()
17   {
18       int n = 256;
19       std::size_t size = sizeof(float)*n;
20
21       float* h_x = (float*)malloc(size);
22       float* h_y = (float*)malloc(size);
23       assert(h_x != nullptr);
24       assert(h_y != nullptr);
25
26       float* d_x;
27       float* d_y;
28       CHECK(hipMalloc(&d_x, size));
29       CHECK(hipMalloc(&d_y, size));
30
31       CHECK(hipMemcpy(d_x, h_x, size, hipMemcpyHostToDevice));
32       CHECK(hipMemcpy(d_y, h_y, size, hipMemcpyHostToDevice));
33
34       int num_blocks = 2;
35       int num_threads = 128;
36       saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
37
38       CHECK(hipMemcpy(h_y, d_y, size, hipMemcpyDeviceToHost));
39       CHECK(hipDeviceSynchronize());
40
41       CHECK(hipFree(d_x));
42       CHECK(hipFree(d_y));
43
44       free(h_x);
45       free(h_y);
46   }
47
```

AMD
together we advance_

# HIPIFY TOOLS

**hipify-clang**

- compiler (clang) based translator

- handles very complex constructs

- prints an error if not able to translate

- supports clang options

- requires CUDA

**hipify-perl**

- Perl® script

- relies on regular expressions

- may struggle with complex constructs

- does not require CUDA

https://github.com/ROCm-Developer-Tools/HIPIFY

**AMD**
together we advance_

```
1   #include <cuda.h>
2   #include <cassert>
3
4   __constant__ float a = 2.0f;
5
6   __global__
7   void saxpy(int n, float const* x, float* y)
8   {
9       int i = blockDim.x*blockIdx.x + threadIdx.x;
10      if (i < n)
11          y[i] += a*x[i];
12  }
13
14  #define CHECK(call) assert(call == cudaSuccess)
15
16  int main()
17  {
18      int n = 256;
19      std::size_t size = sizeof(float)*n;
20
21      float* h_x = (float*)malloc(size);
22      float* h_y = (float*)malloc(size);
23      assert(h_x != nullptr);
24      assert(h_y != nullptr);
25
26      float* d_x;
27      float* d_y;
28      CHECK(cudaMalloc(&d_x, size));
29      CHECK(cudaMalloc(&d_y, size));
30
31      CHECK(cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice));
32      CHECK(cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice));
33
34      int num_blocks = 2;
35      int num_threads = 128;
36      saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
37
38      CHECK(cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost));
39      CHECK(cudaDeviceSynchronize());
40
41      CHECK(cudaFree(d_x));
42      CHECK(cudaFree(d_y));
43
44      free(h_x);
45      free(h_y);
46  }
47
```

```
saxpy$ perl /opt/rocm/bin/hipify-perl -examin saxpy.cu

[HIPIFY] info: file 'saxpy.cu' statisitics:
  CONVERTED refs count: 13
  TOTAL lines of code: 46
  WARNINGS: 0
[HIPIFY] info: CONVERTED refs by names:
  cuda.h => hip/hip_runtime.h: 1
  cudaDeviceSynchronize => hipDeviceSynchronize: 1
  cudaFree => hipFree: 2
  cudaMalloc => hipMalloc: 2
  cudaMemcpy => hipMemcpy: 3
  cudaMemcpyDeviceToHost => hipMemcpyDeviceToHost: 1
  cudaMemcpyHostToDevice => hipMemcpyHostToDevice: 2
  cudaSuccess => hipSuccess: 1
saxpy$ █
```

## hipify-perl

### hipify-perl -examin

o   for initial assessment

o   no replacements done

o   prints basic statistics and the number of replacements

**AMD**
together we advance_

# hipify-perl

```c
1  #include <cuda.h>
2  #include <cassert>
3
4  __constant__ float a = 2.0f;
5
6  __global__
7  void saxpy(int n, float const* x, float* y)
8  {
9      int i = blockDim.x*blockIdx.x + threadIdx.x;
10     if (i < n)
11         y[i] += a*x[i];
12 }
13
14 #define CHECK(call) assert(call == cudaSuccess)
15
16 int main()
17 {
18     int n = 256;
19     std::size_t size = sizeof(float)*n;
20
21     float* h_x = (float*)malloc(size);
22     float* h_y = (float*)malloc(size);
23     assert(h_x != nullptr);
24     assert(h_y != nullptr);
25
26     float* d_x;
27     float* d_y;
28     CHECK(cudaMalloc(&d_x, size));
29     CHECK(cudaMalloc(&d_y, size));
30
31     CHECK(cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice));
32     CHECK(cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice));
33
34     int num_blocks = 2;
35     int num_threads = 128;
36     saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
37
38     CHECK(cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost));
39     CHECK(cudaDeviceSynchronize());
40
41     CHECK(cudaFree(d_x));
42     CHECK(cudaFree(d_y));
43
44     free(h_x);
45     free(h_y);
46 }
47
```

```
saxpy$ perl /opt/rocm/bin/hipify-perl saxpy.cu
#include "hip/hip_runtime.h"
#include <hip/hip_runtime.h>
#include <cassert>

__constant__ float a = 2.0f;

__global__
void saxpy(int n, float const* x, float* y)
{
    int i = blockDim.x*blockIdx.x + threadIdx.x;
    if (i < n)
        y[i] += a*x[i];
}

#define CHECK(call) assert(call == hipSuccess)

int main()
{
    int n = 256;
    std::size_t size = sizeof(float)*n;

    float* h_x = (float*)malloc(size);
    float* h_y = (float*)malloc(size);
    assert(h_x != nullptr);
    assert(h_y != nullptr);

    float* d_x;
    float* d_y;
    CHECK(hipMalloc(&d_x, size));
    CHECK(hipMalloc(&d_y, size));

    CHECK(hipMemcpy(d_x, h_x, size, hipMemcpyHostToDevice));
    CHECK(hipMemcpy(d_y, h_y, size, hipMemcpyHostToDevice));

    int num_blocks = 2;
    int num_threads = 128;
    saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);

    CHECK(hipMemcpy(h_y, d_y, size, hipMemcpyDeviceToHost));
    CHECK(hipDeviceSynchronize());

    CHECK(hipFree(d_x));
    CHECK(hipFree(d_y));

    free(h_x);
    free(h_y);
}
saxpy$ 
```

translating a file
to standard
output

**but can also**
o   translate in place
o   preserve orig copy
o   recursively do folders

**AMD**
**together we advance_**

```
1    #include <hip/hip_runtime.h>
2    #include <cassert>
3    #include "cuda2hip.h"
4
5    __constant__ float a = 2.0f;
6
7    __global__
8    void saxpy(int n, float const* x, float* y)
9    {
10       int i = blockDim.x*blockIdx.x + threadIdx.x;
11       if (i < n)
12           y[i] += a*x[i];
13   }
14
15   #define CHECK(call) assert(call == cudaSuccess)
16
17   int main()
18   {
19       int n = 256;
20       std::size_t size = sizeof(float)*n;
21
22       float* h_x = (float*)malloc(size);
23       float* h_y = (float*)malloc(size);
24       assert(h_x != nullptr);
25       assert(h_y != nullptr);
26
27       float* d_x;
28       float* d_y;
29       CHECK(cudaMalloc(&d_x, size));
30       CHECK(cudaMalloc(&d_y, size));
31
32       CHECK(cudaMemcpy(d_x, h_x, size, cudaMemcpyHostToDevice));
33       CHECK(cudaMemcpy(d_y, h_y, size, cudaMemcpyHostToDevice));
34
35       int num_blocks = 2;
36       int num_threads = 128;
37       saxpy<<<num_blocks, num_threads>>>(n, d_x, d_y);
38
39       CHECK(cudaMemcpy(h_y, d_y, size, cudaMemcpyDeviceToHost));
40       CHECK(cudaDeviceSynchronize());
41
42       CHECK(cudaFree(d_x));
43       CHECK(cudaFree(d_y));
44
45       free(h_x);
46       free(h_y);
47   }
48
```

```
1    #define cudaSuccess              hipSuccess
2    #define cudaMalloc               hipMalloc
3    #define cudaMemcpy               hipMemcpy
4    #define cudaMemcpyHostToDevice   hipMemcpyHostToDevice
5    #define cudaMemcpyDeviceToHost   hipMemcpyDeviceToHost
6    #define cudaDeviceSynchronize    hipDeviceSynchronize
7    #define cudaFree                 hipFree
8
```

**alternatively**

o create a file with renaming macros

o include conditionally, depending on target

AMD◢

together we advance_

# OPTIMIZATION TECHNIQUES

**basic**

o   thread divergence / SIMDzation

o   reuse in shared memory & bank conflicts

o   coalescing of global memory accesses

o   resource partitioning / occupancy / spills

o   L1, L2 cache blocking

o   ...

**advanced**

o   atomics

o   warp primitives

o   CPU-GPU coherence

o   inter-stream synchronization

o   ...

AMD
together we advance_

```cpp
1   #include <cassert>
2   #include <cstdlib>
3   #include <cstdio>
4
5
6   const float a = 2.0f;
7
8   int main()
9   {
10      int n = 256;
11      std::size_t size = sizeof(float)*n;
12
13      float* x = (float*)malloc(size);
14      float* y = (float*)malloc(size);
15      assert(x != nullptr);
16      assert(y != nullptr);
17
18
19      for (int i = 0; i < n; ++i)
20          y[i] += a*x[i];
21
22      free(x);
23      free(y);
24  }
25
```

**alternatively**

o just write CPU code

**AMD**
together we advance_

```cpp
 1    #include <cassert>
 2    #include <cstdlib>
 3    #include <cstdio>
 4    #include <omp.h>
 5
 6    const float a = 2.0f;
 7
 8    int main()
 9    {
10        int n = 256;
11        std::size_t size = sizeof(float)*n;
12
13        float* x = (float*)malloc(size);
14        float* y = (float*)malloc(size);
15        assert(x != nullptr);
16        assert(y != nullptr);
17
18        #pragma omp target teams distribute parallel for map(to:x[0:n]) map(tofrom:y[0:n])
19        for (int i = 0; i < n; ++i)
20            y[i] += a*x[i];
21
22        free(x);
23        free(y);
24    }
25
```

## alternatively

o   just write CPU code

o   use OpenMP® target offload constructs

AMD

together we advance_

45

# KOKKOS AND RAJA

o portability frameworks based on C++

o portability to CPUs & GPUs – AMD, Intel®, NVIDIA

o basic parallel processing constructs

o multidimensional arrays

o etc., etc., etc.

**Kokkos**

o originates from Sandia National Laboratory

o https://kokkos.org/

o https://github.com/kokkos

**RAJA**

o originates from Lawrence Livermore

o https://raja.readthedocs.io

o https://github.com/LLNL/RAJA

**AMD**

together we advance_

# DIFFERENCES FROM CUDA

- warpSize

  - 64 on AMD

  - 32 on NVIDIA

- dynamic parallelism not supported

- exercise caution:

  - atomics

  - managed memory

  - warp-level primitives

  - inter-process communication

**AMD**
together we advance_

# AMD RESOURCES
## DOCUMENTATION AND TRAINING

AMD
together we advance_

# AMD ROCM DEVELOPER HUB

## Engage with ROCm Experts

Participate in ROCm Webinar Series

Post questions, view FAQ's in Community Forum

## Increase Understanding

Purchase ROCm Text Book

View the latest news in the Blogs

## Get Started Using ROCm

ROCm Documentation on GitHub

Download the Latest Version of ROCm

https://www.amd.com/en/developer/rocm-hub.html

# NEW ROCM DOCS

## Comprehensive Coverage

### Compilers and Frameworks

### Math libraries, communication libraries

### Management tools, validation tools

...

## Howto Guides

### Installation

### Tunning

### Debugging

...

https://rocm.docs.amd.com/

---

**AMD** | ROCm™ Platform 5.6.0

GitHub     Community     AMD Lab Notes     Infinity Hub     Support     Feedback

ROCm Documentation Home

What is ROCm?

**Deploy ROCm**

Linux Quick Start

Linux Overview

Docker

**Release Info**

Release Notes

Changelog

GPU Support and OS Compatibility (Linux)

Known Issues ⧉

Compatibility

Licensing Terms

**APIs and Reference**

All Reference Material

HIP

Math Libraries

C++ Primitive Libraries

Communication Libraries

AI Libraries

Computer Vision

OpenMP

Compilers and Tools

Management Tools

Validation Tools

**Understand ROCm**

All Explanation Material

Compiler Disambiguation

## AMD ROCm™ Documentation

Applies to Linux     📅 2023-05-25     🕐 3 min read time

| What is ROCm? ⌄ | Deploy ROCm ⌄ | Release Info ⌄ |

### APIs and Reference

- Compilers and Development Tools
- HIP
- OpenMP
- Math Libraries
- C++ Primitives Libraries
- Communication Libraries
- AI Libraries
- Computer Vision
- Management Tools
- Validation Tools

### Understand ROCm

- Compiler Disambiguation
- Using CMake
- Linux Folder Structure Reorganization
- GPU Isolation Techniques
- GPU Architecture

### How to Guides

- System Tuning for Various Architectures
- GPU Aware MPI
- Setting up for Deep Learning with ROCm
  - Magma Installation
  - PyTorch Installation
  - TensorFlow Installation
- System Level Debugging

### Tutorials & Examples

- Examples
- ML, DL, and AI
  - Inception V3 with PyTorch
  - Inference Optimization with MIGraphX

Next  ›
What is ROCm?

---

**AMD**

together we advance_

# HIP TEXTBOOK

## Comprehensive Coverage

HIP Language

AMD GPU Internals

Performance Analysis

Debugging

Programming Patterns

ROCm Libraries

Porting to HIP

Multi-GPU Programming

Third Party Tools

CDNA Assembly

ML with ROCm



ACCELERATED COMPUTING WITH HIP

Yifan Sun
Trinayan Baruah
David Kaeli

https://www.barnesandnoble.com/w/accelerated-computing-with-hip-yifan-sun/1142866934

AMD
together we advance_

# ISA REFERENCE GUIDE

## Public ISA

The Instruction Set Architecture is public

There is no intermediate layer like PTX

You can write assembly code

You can compile to assembly for inspection

**AMD**

"AMD Instinct MI200" Instruction Set Architecture
*Reference Guide*

4-February-2022

https://www.amd.com/content/dam/amd/en/documents/instinct-tech-docs/instruction-set-architectures/instinct-mi200-cdna2-instruction-set-architecture.pdf

**AMD** together we advance_

# AMD INFINITY HUB

## AMD Instinct™ MI200 SUPPORT
29 key applications & frameworks on Infinity Hub & a catalogue supporting over 90 applications, frameworks & tools

## Accelerating Instinct™ adoption
Over 17000 application pulls. 10000+ since last year

## PERFORMANCE RESULTS
Published Performance Results for Select Apps / Benchmarks

https://www.amd.com/en/technologies/infinity-hub



**AMD**

Products    Solutions    Resources & Support    Shop

### AMD Infinity Hub

**Categories**
- AI & Machine Learning
- Benchmark
- Deep Learning
- Earth Science
- HPC
- Life Science
- Material Science
- Molecular Dynamics
- Oil and Gas
- Physics

**Containers**
- Yes
- No

**Computational Science Sta**
The AMD Infinity Hub contains a collection of advan
AI applications, enabling researchers, scientists and

INSTINCT™ APP CATALOG

ZENDNN

Search

**Amber**
Amber is a suite of biomolecular simulation programs. It began in the late 1970's, and is maintained by an...

MORE INFO

**BabelStream**
BabelStream is a synthetic GPU benchmark based on the original STREAM benchmark for CPUs. The...

MORE INFO    PULL TAG

**CP2K**
CP2K is a quantum chemistry and solid state physics software package that can perform atomistic...

**GROMACS**
GROMACS is a versatile package to perform molecular dynamics, i.e. simulate the Newtonian equations of...

AMD
together we advance_

# SOFTWARE CATALOG

## STRONG MOMENTUM AND INCREASING LIST OF SUPORTED APPLICATION, LIBRARIES & FRAMEWORKS

### Life Science
AMBER
GROMACS
NAMD
LAMMPS
Hoomd-Blue
VASP

### Physics
MILC
GRID
QUANTUM ESPRESSO
N-Body
CHROMA
PIConGPU
QuickSilver

### Chemistry
CP2K
QUDA
NWCHEM
TERACHEM
QMCPACK

### CFD
OpenFOAM®
AMR-WIND
NEKBONE
LAGHOS
NEKO
NEKRS
PeleC

### Earth Science
EXAGO
DEVITO
OCCA
SPECFEM3D-GLOBE
SPECFEM3D-CARTESIAN
ACECAST (WRF)
MPAS
ICON

### Benchmarks
HPL
HPCG
AMG
ML - TORCHBENCH
ML - SUPERBENCH

### Libraries
AMR-EX
Ginkko
HYPRE
TRILINOS

### ML Frameworks
PYTORCH
TENSORFLOW
JAX
ONNX
OPENAI TRITON

### ISV Applications
ANSYS MECHANICAL
CADENCE CHARLES
ANSYS FLUENT*
SIEMENS® STAR-CCM+*
SIEMENS® CALIBRE*

+ MANY MORE

* Porting/optimization in progress

AMD
together we advance_

# AMD LAB NOTES

## Introductory Topics

ROCm installation

Basics of HIP programming

...

## Advanced Topics

Matrix Cores

Register pressure

GPU-aware MPI

...

https://gpuopen.com/learn/amd-lab-notes/
https://github.com/AMD/amd-lab-notes

---

AMD GPUOpen

Search GPUOpen...

HOME    SOFTWARE ⌄    DOCS ⌄

Originally posted November 14, 2022
Updated on May 2, 2023

Home » Blogs » AMD lab notes » AMD matrix cores

### AMD matrix cores

Matrix multiplication is a fundamental aspect of Linear Algebra and it is an ubiquitous computation within High Performance Computing (HPC) Applications. Since the introduction of AMD's CDNA Architecture, Generalized Matrix Multiplication (GEMM) computations are now hardware-accelerated through Matrix Core Processing Units. Matrix Core accelerated GEMM kernels lie at the heart of BLAS libraries like rocBLAS but they can also be programmed directly by developers. Applications that are throughput bound by GEMM computation can achieve additional speedups by utilizing Matrix Cores.

AMD's Matrix Core technology supports a full range of mixed precision operations bringing us the ability to work with large models and enhance memory-bound operation performance for any combination of AI and machine learning workloads. The various numerical formats have uses in different applications. Examples include use of 8-bit integers (INT8) for ML inference, 32-bit floating point (FP32) data for ML Training and HPC applications, 16-bit floating point (FP16) data for graphics workloads and 16-bit brain float (BF16) data for ML training with fewer convergence issues.

To learn more about the theoretical speedups achievable by using matrix cores compared to SIMD Vector Units, please refer to the tables below. The tables list the performance of the Vector (i.e. Fused Multiply-Add or FMA) and Matrix core units of the previous generation (MI100) and current generation (MI250X) of CDNA Accelerators.

Matrix Core Performance for MI100 and MI250X:

| Data format | MI100 Flops/Clock/CU | MI250X Flops/Clock/CU |
|---|---|---|
| FP64 | N/A | 256 |
| FP32 | 256 | 256 |
| FP16 | 1024 | 1024 |
| BF16 | 512 | 1024 |
| INT8 | 1024 | 1024 |

Vector (FMA) Unit Performance for MI100 and MI250X:

| Data format | MI100 Flops/Clock/CU | MI250X Flops/Clock/CU |
|---|---|---|
| FP64 | 64 | 128 |
| FP32 | 128 | 128 |

**AMD lab notes**
- Finite difference method - Laplacian part 1
- Finite difference method - Laplacian part 2
- Finite difference method - Laplacian part 3
- Finite difference method - Laplacian part 4
- **AMD matrix cores**
  - *Using AMD matrix cores*
  - *MFMA compiler intrinsic syntax*
  - *Example 1 - V_MFMA_F32_16x16x4F32*
  - *Example 2 - V_MFMA_F32_16x16x1F32*
  - *Example 3 - V_MFMA_F64_4x4x4F64*
  - *A note on rocWMMA*
  - *A note on the AMD Matrix Instruction Calculator...*
  - *References*
- Introduction to profiling tools for AMD hardware
- AMD ROCm™ installation
- AMD Instinct™ MI200 GPU memory space overview
- Register pressure in AMD CDNA™2 GPUs
- GPU-aware MPI with ROCm

Search this manual ...

AMD together we advance_

# OLCF TRAINING

## Tutorials, Workshops, Hackathons

Slides available online

Recordings available online

## User Guides

Frontier User Guide

Crusher User Guide

https://www.olcf.ornl.gov/for-users/training/

# OLCF TRAINING ARCHIVE

## Frontier Training

**GPU Profiling**

**GPU Debugging**

**Node Performance Engineering**

**Programming Models for AMD GPUs**

## Produced by

**AMD staff**

**HPE staff**

**ORNL staff**

---

**OLCF User Documentation**

Search docs

- New User Quick Start
- Accounts and Projects
- Connecting
- Systems
- Services and Applications
- Data Storage and Transfers
- Software
- Training
  - OLCF Training Calendar
  - OLCF Tutorials
  - **OLCF Training Archive**
  - OLCF GPU Hackathons
  - OLCF Vimeo Channel
- Quantum
- Scalable Protected Infrastructure (SPI)
- Contributing to these docs

---

| Date | Title | Presenter | | |
|------|-------|-----------|------|------|
| 2023-02-17 | Checkpointing Tips | Scott Atchley, HPC Systems Engineer, Distinguished R&D Staff, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-17 | Frontier Tips & Tricks | Balint Joo, Group Leader, Advanced Computing for Nuclear, Particles, & Astrophysics, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-17 | GPU Debugging | Mark Stock, HPC Applications Engineer, HPE | Frontier Training Workshop | (slides \| recording) |
| 2023-02-17 | GPU Profiling | Alessandro Fanfarillo, Senior Member of Technical Staff, Exascale Application Performance, AMD | Frontier Training Workshop | (slides \| recording) |
| 2023-02-17 | Application Profiling | Trey White, Master Engineer, HPE | Frontier Training Workshop | (slides \| recording) |
| 2023-02-16 | Orion Lustre and Best Practices | Jesse Hanley, Senior HPC Linux Systems Engineer, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-16 | Node Performance | Tom Papatheodore, HPC Engineer, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-16 | NVMe Usage | Chris Zimmer, Group Leader, Technology Integration, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-16 | AI on Frontier | Junqi Yin, Computational Scientist, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-16 | Python on Frontier | Michael Sandoval, HPC Engineer, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-16 | HPE Cray MPI | Tim Mattox, HPC Performance Engineer, HPE | Frontier Training Workshop | (slides \| recording) |
| 2023-02-16 | GPU Programming Models | GPU Programming Models | Frontier Training Workshop | (slides \| recording) |
| 2023-02-15 | Slurm on Frontier | Tom Papatheodore, HPC Engineer, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-15 | Storage Areas and Data Transfers | Suzanne Parete-Koon, HPC Engineer, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-15 | Using the Frontier Programming Environment | Matt Belhorn, HPC Engineer, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-15 | Frontier Programming Environment | Wael Elwasif, Computer Scientist, ORNL | Frontier Training Workshop | (slides \| recording) |
| 2023-02-15 | Epyc CPU and Instinct GPU | Nick Malaya, Principal Member of Technical Staff, Exascale Application Performance, AMD | Frontier Training Workshop | (slides \| recording) |
| 2023-02-15 | Frontier Architecture Overview | Joe Glenski, Sr. Distinguished Technologist, HPE | Frontier Training Workshop | (slides \| recording) |
| 2023-02-15 | Welcome to the Frontier Workshop | Ashley Barker, Section Head, Operations, National Center for Computational Sciences, ORNL | Frontier Training Workshop | (slides \| recording) |

https://docs.olcf.ornl.gov/training/training_archive.html

# OLCF PROGRAMMING GUIDES

Frontier User Guide

Crusher Quick-Start Guide

GPU architecture

Node architecture

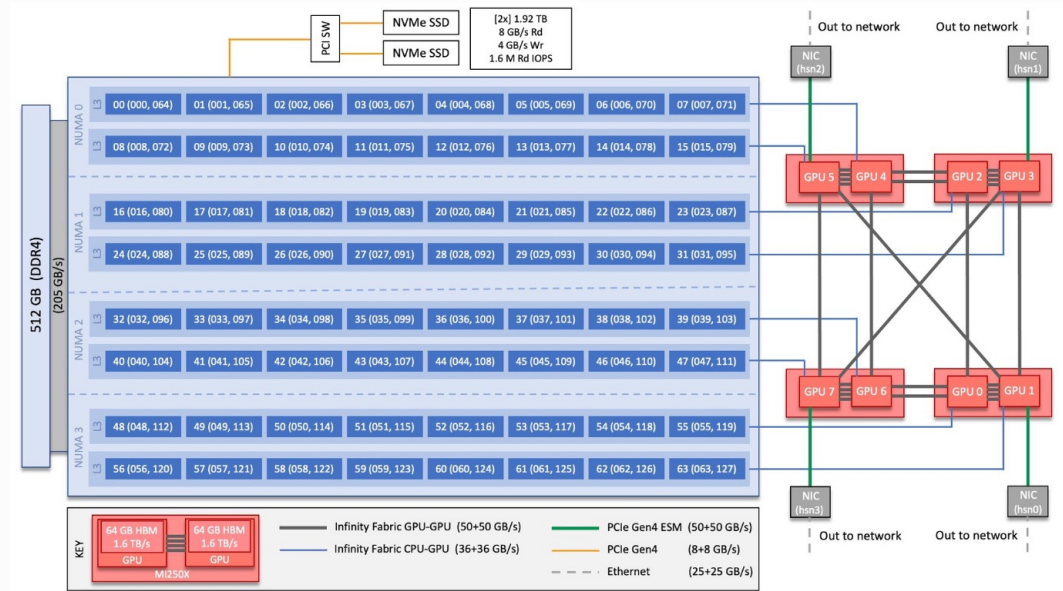Programming environment (HIP, OpenMP®)

Profiling

Debugging

...

https://docs.olcf.ornl.gov/systems/frontier_user_guide.html
https://docs.olcf.ornl.gov/systems/crusher_quick_start_guide.html

> **ⓘ Note**
>
> TERMINOLOGY:
>
> The 8 GCDs contained in the 4 MI250X will show as 8 separate GPUs according to Slurm, `ROCR_VISIBLE_DEVICES`, and the ROCr runtime, so from this point forward in the quick-start guide, we will simply refer to the GCDs as GPUs.



> **ⓘ Note**
>
> There are [4x] NUMA domains per node and [2x] L3 cache regions per NUMA for a total of [8x] L3 cache regions. The 8 GPUs are each associated with one of the L3 regions as follows:
>
> NUMA 0:
>
> - hardware threads 000-007, 064-071 | GPU 4
> - hardware threads 008-015, 072-079 | GPU 5
>
> NUMA 1:
>
> - hardware threads 016-023, 080-087 | GPU 2
> - hardware threads 024-031, 088-095 | GPU 3
>
> NUMA 2:

AMD
together we advance_

# ENCCS AMD TRAINING VIDEOS

HIP programming

OpenMP® offload

Developing in Fortran

GPU-aware MPI

Roofline modeling

Profiling

Debugging

ML frameworks

...

https://enccs.github.io/amd-rocm-development/

---

🏠 **Developing Applications with the AMD ROCm Ecosystem**



`Search docs`

**THE LESSON**

**REFERENCE**

---

🏠 / Developing Applications with the AMD ROCm Ecosystem    ⬡ Edit on GitHub

## Developing Applications with the AMD ROCm Ecosystem



This training material is created by AMD in collaboration with ENCCS. It covers how to develop and port applications to run on AMD GPU and CPU hardware on top AMD-powered supercomputers. You will learn about the ROCm software development languages, libraries, and tools, as well as getting a developer's view of the hardware that powers the system. The material focuses mostly on how to program applications to run on the GPU.

⚙ **Prerequisites**

It is useful to have prior experience developing HPC applications, and some understanding of recent HPC computer hardware and the Linux operating system.

## The lesson

- Introduction to HIP Programming
- Porting Applications to HIP
- Getting Started with OpenMP® Offload Applications on AMD Accelerators
- Developing Fortran Applications: HIPFort, OpenMP®, and OpenACC
- Exercises
- Architecture
- GPU-Aware MPI with ROCmTM
- AMD Node Memory Model
- Hierarchical Roofline on AMD InstinctTM MI200 GPUs
- Affinity — Placement, Ordering and Binding
- Profiling and debugging
- OpenMP Offload Programming
- Introduction to ML Frameworks
- Summary and outlook

AMD
together we advance_

# PAWSEY AMD TRAINING VIDEOS

Introduction rocprof

Introduction to omniprof

Introduction to omnitrace

Roofline modeling



https://www.youtube.com/playlist?list=PLmu61dgAX-aaQOCG5Jlw8oLBORJfoQC2o

AMD
together we advance_

# DISCLAIMERS

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. Any computer system has risks of security vulnerabilities that cannot be completely prevented or mitigated.  AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

THIS INFORMATION IS PROVIDED 'AS IS." AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS, OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION. AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY RELIANCE, DIRECT, INDIRECT, SPECIAL, OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Third-party content is licensed to you directly by the third party that owns the content and is not licensed to you by AMD.  ALL LINKED THIRD-PARTY CONTENT IS PROVIDED "AS IS" WITHOUT A WARRANTY OF ANY KIND.  USE OF SUCH THIRD-PARTY CONTENT IS DONE AT YOUR SOLE DISCRETION AND UNDER NO CIRCUMSTANCES WILL AMD BE LIABLE TO YOU FOR ANY THIRD-PARTY CONTENT.  YOU ASSUME ALL RISK AND ARE SOLELY RESPONSIBLE FOR ANY DAMAGES THAT MAY ARISE FROM YOUR USE OF THIRD-PARTY CONTENT.

© 2023 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, AMD CDNA, AMD ROCm, AMD Instinct, and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions. Other names are for informational purposes only and may be trademarks of their respective owners.

**AMD**
together we advance_

# ATTRIBUTIONS

Docker and the Docker logo are trademarks or registered trademarks of Docker, Inc.

Git and the Git logo are either registered trademarks or trademarks of Software Freedom Conservancy, Inc., corporate home of the Git Project, in the United States and/or other countries.

Intel is a trademark of Intel Corporation or its subsidiaries.

Kubernetes is a registered trademark of The Linux Foundation.

NAMD was developed by the Theoretical Biophysics Group in the Beckman Institute for Advanced Science and Technology at the University of Illinois at Urbana-Champaign. http://www.ks.uiuc.edu/Research/namd/

OpenCL is a trademark of Apple Inc. used by permission by Khronos Group, Inc.

OpenFOAM is a registered trademark of OpenCFD Limited, producer and distributor of the OpenFOAM software via www.openfoam.com.

The OpenMP name and the OpenMP logo are registered trademarks of the OpenMP Architecture Review Board.

Perl is a trademark of Perl Foundation.

Siemens is a registered trademark of Siemens Product Lifecycle Management Software Inc., or its subsidiaries or affiliates, in the United States and in other countries.

**AMD**
together we advance_