# AMD HARDWARE AND SOFTWARE

**Suyash Tandon, Justin Chang, Julio Maia, Noel Chalmers, Paul T. Bauman, Nicholas Curtis, Nicholas Malaya, Alessandro Fanfarillo, Jose Noudohouenou, Chip Freitag, Damon McDougall, Noah Wolfe, Jakub Kurzak, Samuel Antao, George Markomanolis, Bob Robey**

ADVANCED MICRO DEVICES, INC.

**AMD**
together we advance_

slides on LUMI in /project/project_465001098/Slides/AMD/

hands-on exercises: https://hackmd.io/@gmarkoma/lumi_finland

hands-on source code: /project/project_465001098/Exercises/AMD/HPCTrainingExamples/

**AMD**
together we advance_

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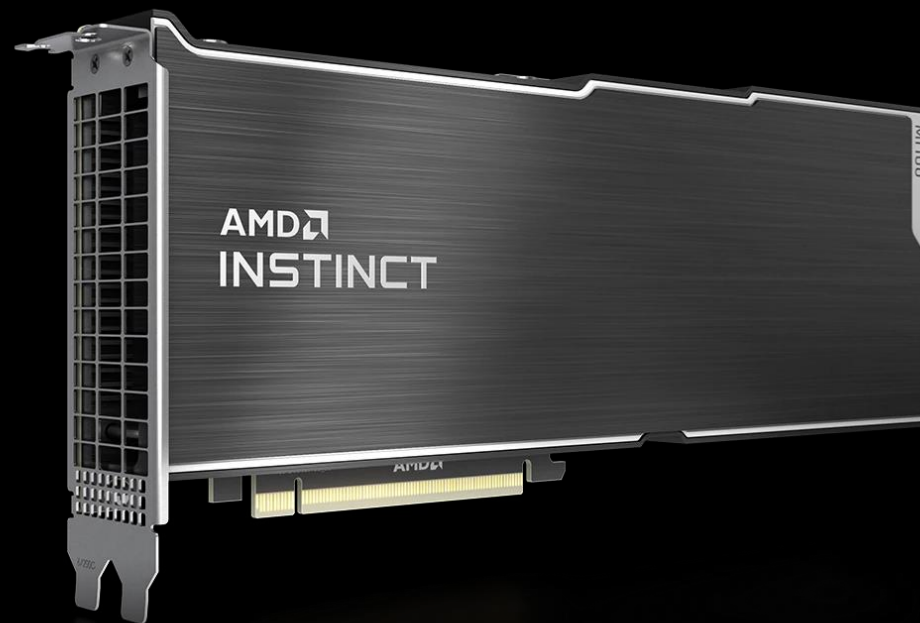
AMD
together we advance_

# 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

AMD
together we advance_
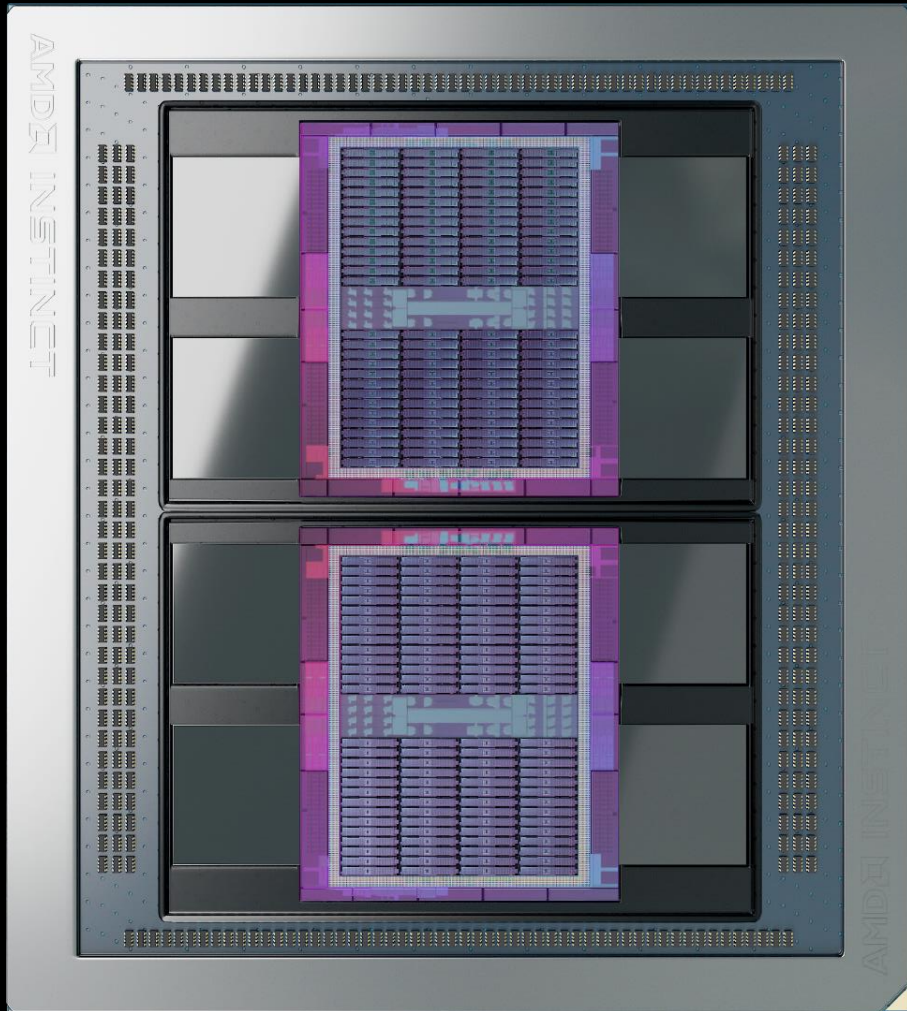
# AMD IN HPC



**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 5th fastest in the world

o AMD EPYC CPUs

o AMD Instinct GPUs

**AMD** ⅃
together we advance_

# AMD INSTINCT™ MI200

## AMD INSTINCT™ MI250X

# ONE OF THE 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

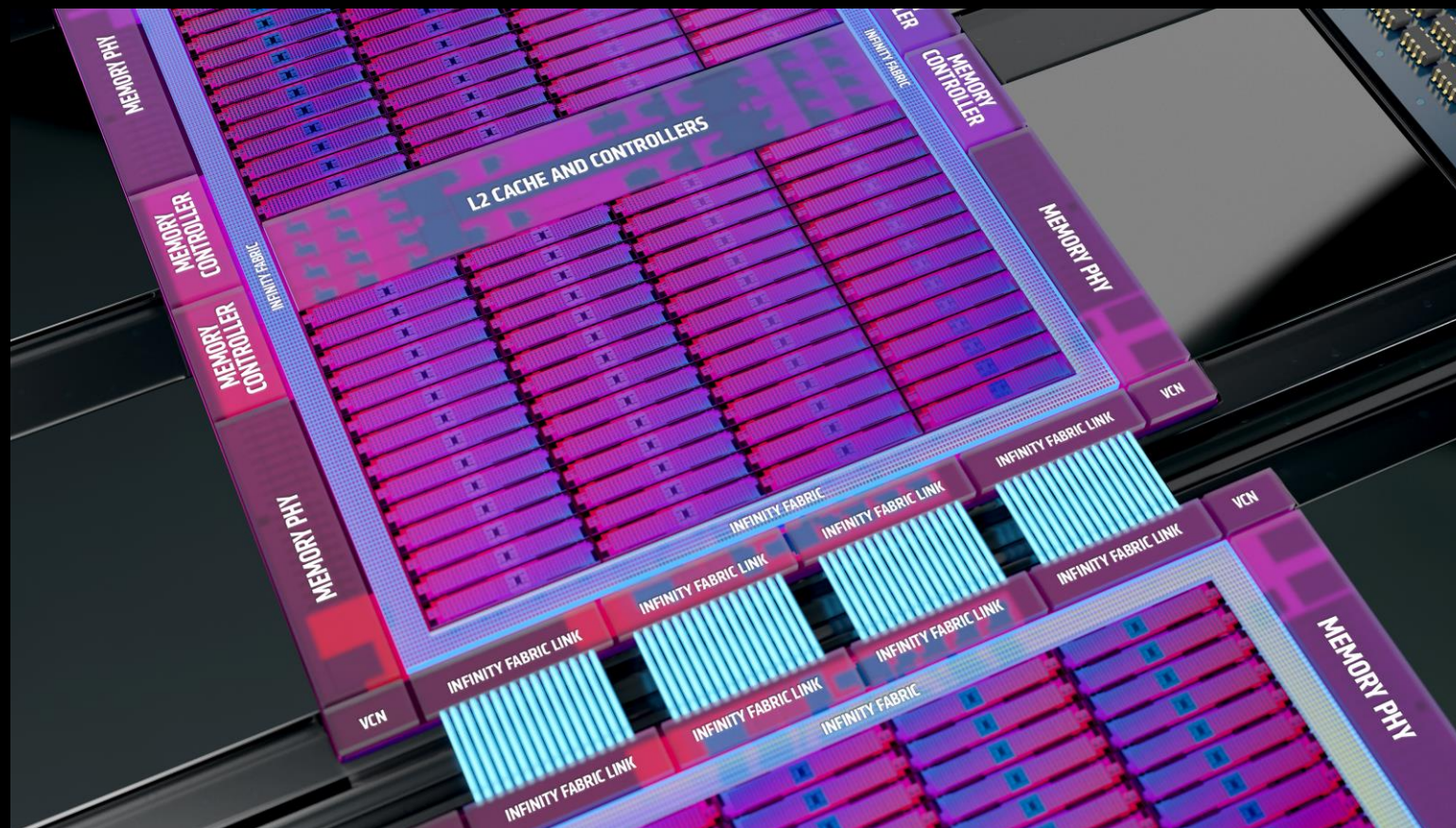| | |
|---|---|
| TSMC 6NM TECHNOLOGY | UP TO 110 CU PER GRAPHICS CORE DIE |
| 4 MATRIX CORES PER COMPUTE UNIT | MATRIX CORES ENHANCED FOR HPC |
| 8 INFINITY FABRIC LINKS PER DIE | SPECIAL FP32 OPS FOR DOUBLE THROUGHPUT |

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
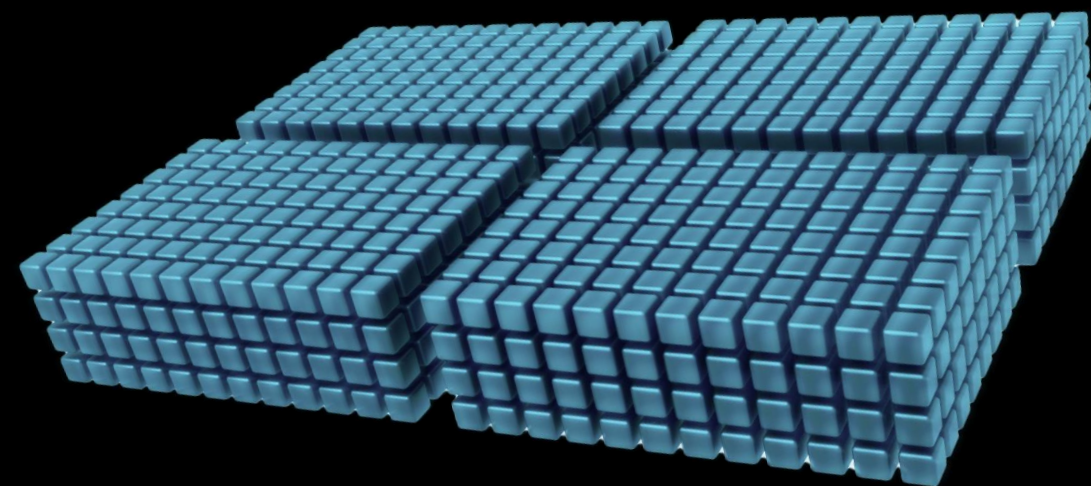
400GB/S OF BI-DIRECTIONALBANDWIDTH

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

# AMD MI250X specifications

- Two graphic compute dies (GCDs)
- 64GB of HBM2e memory per GCD (total 128GB)
- 26.5 TFLOPS peak performance per GCD
- 1.6 TB/s peak memory bandwidth per GCD
- 110 CU per GCD, total 220 CU per GPU
- The interconnection is attached to the GPU (not on the CPU)
- Both GCDs are interconnected with 200 GB/s per direction
- 128 single precision FMA operations per cycle
- AMD CDNA 2 Matrix Core supports double-precision data
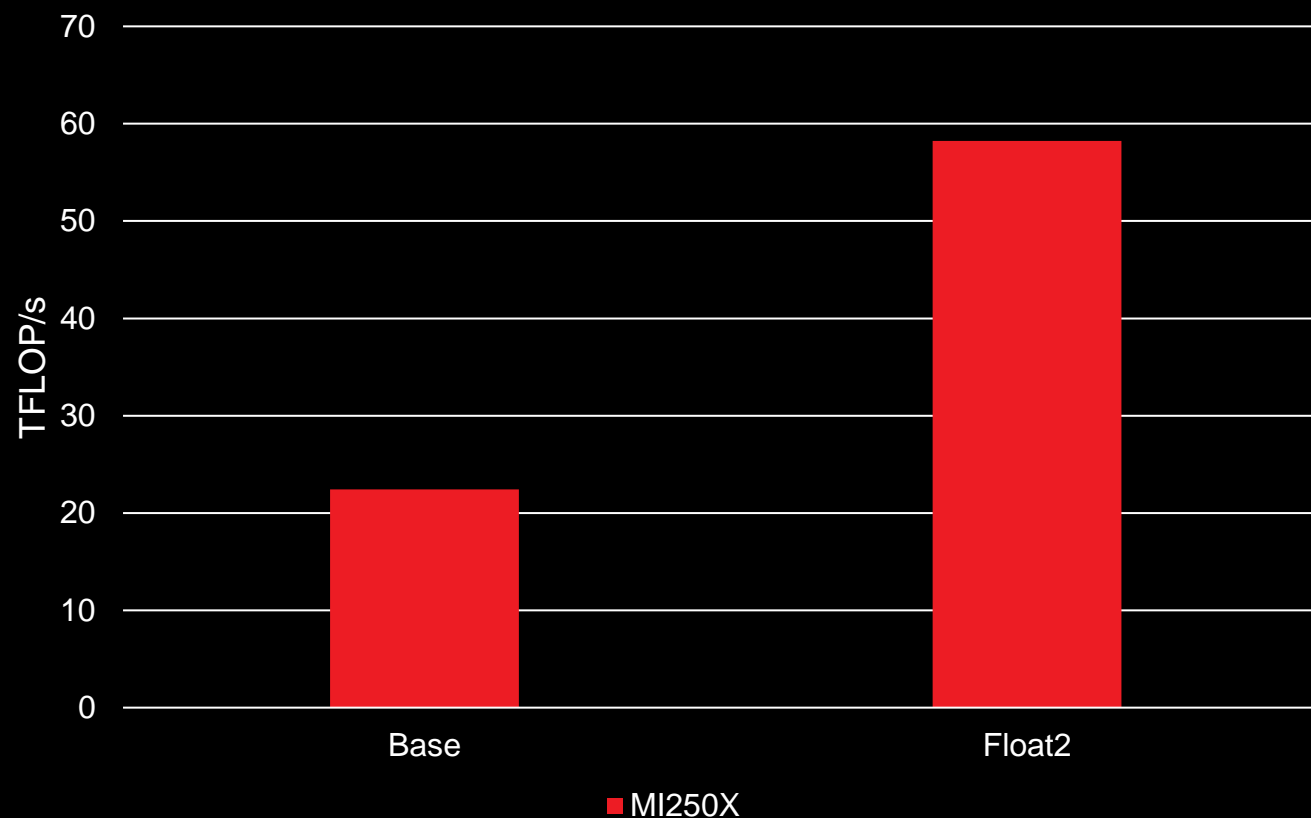- Memory coherency

AMD CDNA™ 2 white paper: https://www.amd.com/system/files/documents/amd-cdna2-white-paper.pdf

AMD

together we advance_

# NEW IN AMD INSTINCT MI250X
## PACKED FP32

FP64 PATH USED TO EXECUTE TWO COMPONENT VECTOR INSTRUCTIONS ON FP32

DOUBLES FP32 THROUGHPUT PER CLOCK PER COMPUTE UNIT

pk_FMA, pk_ADD, pk_MUL, pk_MOV operations



TFLOP/s

Base

Float2

MI250X

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

AMD
together we advance_

# MI200 COMPUTE UNIT

scheduler

scalar unit

Local Data Share (64KB)

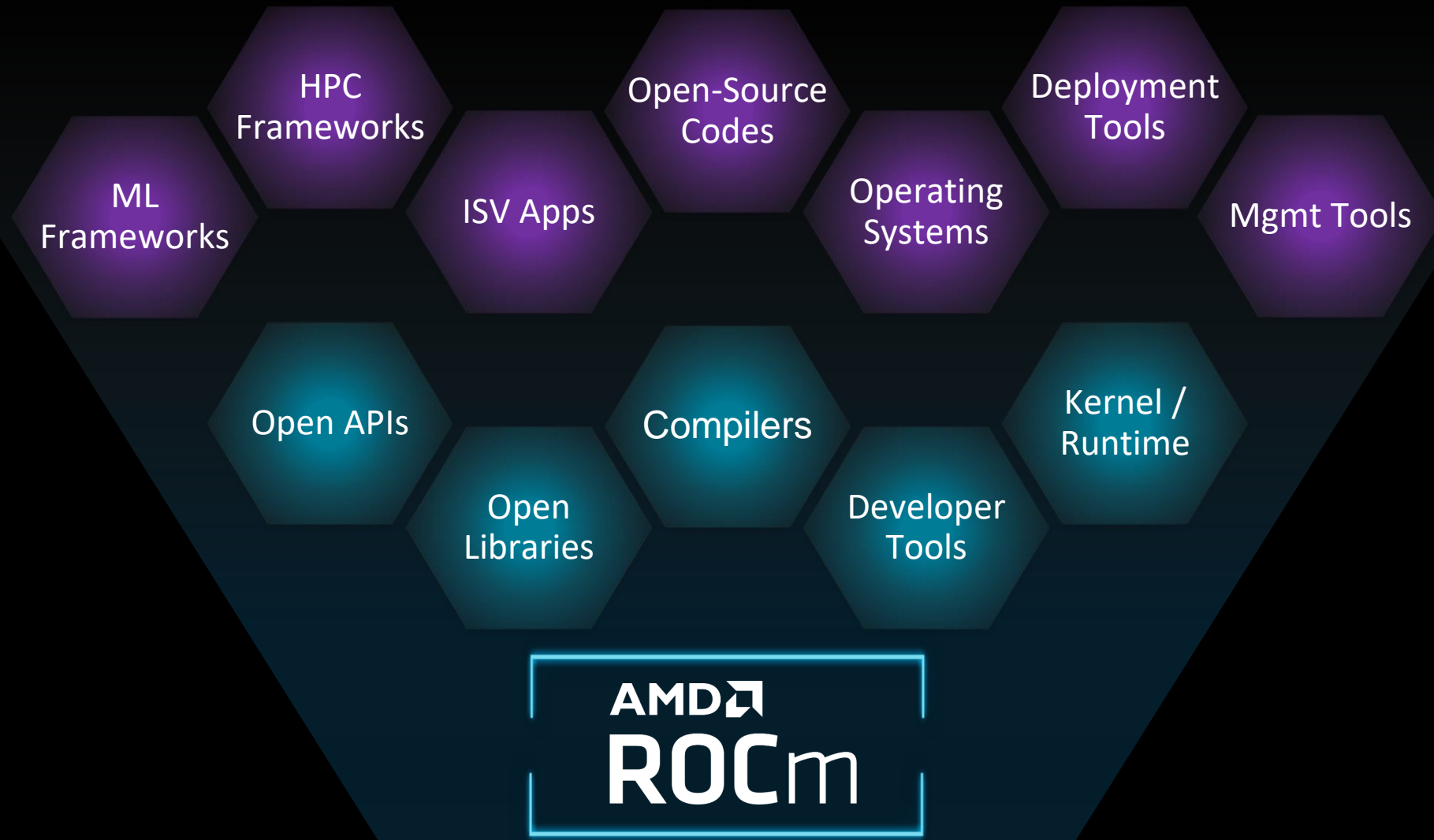scalar registers (8KB)

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 SOFTWARE FOR HPC AND AI
## ROCM PLATFORM

AMD
together we advance_

# AMD ROCm™ Open Software Platform For GPU Compute

HPC Frameworks

Open-Source Codes

Deployment Tools

ML Frameworks

ISV Apps

Operating Systems

Mgmt Tools

Open APIs

Compilers

Kernel / Runtime

Open Libraries

Developer Tools

**AMD**
**ROCm**

AMD
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**
o   basic operations on dense matrices

**rocSOLVER**
o   dense linear algebra solvers

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

**rocALUTION**
o   sparse linear algebra solvers

**rocFFT / hipFFT**
o   Fast Fourier transforms

**rocRAND / hipRAND**
o   random number generation

**rocPRIM / hipCUB / rocThrust**
o   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**
- https://github.com/ROCmSoftwarePlatform/llvm-project

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

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

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

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

**the HPCG benchmark**
- 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

AMD
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**

○  two 1D arrays

○  the y[i] += a*x[i] operation

○  mapped to 1D grid of threads/blocks

○  each thread takes on index

```cpp
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_

[Public]

# 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

- check number of devices

- switch devices

- 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

- create stream

- destroy stream

- synchronize stream



- etc.

- etc.

- 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_

# STREAMS

- Suppose we have 4 small kernels to execute:

```
hipLaunchKernelGGL(myKernel1, dim3(1), dim3(256), 0, 0, 256, d_a1);
hipLaunchKernelGGL(myKernel2, dim3(1), dim3(256), 0, 0, 256, d_a2);
hipLaunchKernelGGL(myKernel3, dim3(1), dim3(256), 0, 0, 256, d_a3);
hipLaunchKernelGGL(myKernel4, dim3(1), dim3(256), 0, 0, 256, d_a4);
```

- Even though these kernels use only one block each, they'll execute in serial on the NULL stream:

| NULL Stream | myKernel1 | myKernel2 | myKernel3 | myKernel4 |
|---|---|---|---|---|

Time →

AMD
together we advance_

# STREAMS

- With streams we can effectively share the GPU's compute resources:

```
hipLaunchKernelGGL(myKernel1, dim3(1), dim3(256), 0, stream1, 256, d_a1);
hipLaunchKernelGGL(myKernel2, dim3(1), dim3(256), 0, stream2, 256, d_a2);
hipLaunchKernelGGL(myKernel3, dim3(1), dim3(256), 0, stream3, 256, d_a3);
hipLaunchKernelGGL(myKernel4, dim3(1), dim3(256), 0, stream4, 256, d_a4);
```

| NULL Stream | | |
|---|---|---|
| Stream1 | myKernel1 | |
| Stream2 | myKernel2 | |
| Stream3 | myKernel3 | |
| Stream4 | myKernel4 | |

Note 1: Kernels must modify different parts of memory to avoid data races.

Note 2: With large kernels, overlapping computations may not help performance.

AMD

together we advance_

# STREAMS

- There is another use for streams besides concurrent kernels:
  - Overlapping kernels with data movement.

- AMD GPUs have separate engines for:
  - Host->Device memcpys
  - Device->Host memcpys
  - Compute kernels.

- These three different operations can overlap without dividing the GPU's resources.
  - The overlapping operations should be in separate, non-NULL, streams.
  - The host memory should be **pinned.**

**AMD**
together we advance_

# STREAMS

Suppose we have 3 kernels which require moving data to and from the device:

```
hipMemcpy(d_a1, h_a1, Nbytes, hipMemcpyHostToDevice));
hipMemcpy(d_a2, h_a2, Nbytes, hipMemcpyHostToDevice));
hipMemcpy(d_a3, h_a3, Nbytes, hipMemcpyHostToDevice));

hipLaunchKernelGGL(myKernel1, blocks, threads, 0, 0, N, d_a1);
hipLaunchKernelGGL(myKernel2, blocks, threads, 0, 0, N, d_a2);
hipLaunchKernelGGL(myKernel3, blocks, threads, 0, 0, N, d_a3);

hipMemcpy(h_a1, d_a1, Nbytes, hipMemcpyDeviceToHost);
hipMemcpy(h_a2, d_a2, Nbytes, hipMemcpyDeviceToHost);
hipMemcpy(h_a3, d_a3, Nbytes, hipMemcpyDeviceToHost);
```

| NULL Stream | HToD1 | HToD2 | HToD3 | myKernel 1 | myKernel 2 | myKernel 3 | DToH1 | DToH2 | DToH3 |
|---|---|---|---|---|---|---|---|---|---|

AMD
together we advance_

# STREAMS

Changing to asynchronous memcpys and using streams:

```
hipMemcpyAsync(d_a1, h_a1, Nbytes, hipMemcpyHostToDevice, stream1);
hipMemcpyAsync(d_a2, h_a2, Nbytes, hipMemcpyHostToDevice, stream2);
hipMemcpyAsync(d_a3, h_a3, Nbytes, hipMemcpyHostToDevice, stream3);

hipLaunchKernelGGL(myKernel1, blocks, threads, 0, stream1, N, d_a1);
hipLaunchKernelGGL(myKernel2, blocks, threads, 0, stream2, N, d_a2);
hipLaunchKernelGGL(myKernel3, blocks, threads, 0, stream3, N, d_a3);

hipMemcpyAsync(h_a1, d_a1, Nbytes, hipMemcpyDeviceToHost, stream1);
hipMemcpyAsync(h_a2, d_a2, Nbytes, hipMemcpyDeviceToHost, stream2);
hipMemcpyAsync(h_a3, d_a3, Nbytes, hipMemcpyDeviceToHost, stream3);
```

| NULL Stream | | | | | |
|---|---|---|---|---|---|
| Stream1 | HToD1 | myKernel 1 | DToH1 | | |
| Stream2 | | HToD2 | myKernel 2 | DToH2 | |
| Stream3 | | | HToD3 | myKernel 3 | DToH3 |

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_

# ADDING THE CPU CODE

```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* 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

**AMD**
together we advance_

```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_

# ADDING MEMORY CLEANUP

```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      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_

# ADDING ERROR CHECKS

```
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

```cuda
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
```

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

```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  |
```

AMD
together we advance_

# HIPIFY TOOLS

**hipify-clang**

o compiler (clang) based translator

o handles very complex constructs

o prints an error if not able to translate

o supports clang options

o requires CUDA

**hipify-perl**

o Perl® script

o relies on regular expressions

o may struggle with complex constructs

o 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_

45

# hipify-perl

```
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_

```cpp
#include <hip/hip_runtime.h>
#include <cassert>
#include "cuda2hip.h"

__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);
}
```

```cpp
#define cudaSuccess             hipSuccess
#define cudaMalloc              hipMalloc
#define cudaMemcpy              hipMemcpy
#define cudaMemcpyHostToDevice  hipMemcpyHostToDevice
#define cudaMemcpyDeviceToHost  hipMemcpyDeviceToHost
#define cudaDeviceSynchronize   hipDeviceSynchronize
#define cudaFree                hipFree

```

**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_

# 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

AMD
together we advance_

# 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/

# 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_

# 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

# 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

# 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_