# AMD INSTINCT™ GPUs CAPABILITY AND CAPACITY AT SCALE

**Samuel Antao**, Suyash Tandon, Michael Rowan, Nicholas Malaya
*HTCondor Autumn Workshop, September 26th, 2024*

# AMD PLATFORM FOR ACCELERATED COMPUTING
## LEADERSHIP IN HPC & AI FOR EXASCALE-CLASS COMPUTING

**AMD CDNA**

WORKLOAD-OPTIMIZED
COMPUTE ARCHITECTURE

**PURPOSE-BUILT, OPTIMIZED
ARCHITECTURE**
DESIGNED TO DO ONE THING
EXTREMELY WELL: COMPUTE
INTENSIVE HPC AND AI WORKLOADs

**AMD ROCm**

OPEN & PORTABLE
SOFTWARE

**SIMPLIFIED PROGRAMMABILITY
AND INCREASED USABILITY**
THROUGH A GROWING SOFTWARE
ECOSYSTEM AND A RICH SUITE OF
OPTIMIZED LIBRARIES,
FRAMEWORKS AND TOOLS

DEEP PARTNERSHIP WITH
LEADING HPC CENTERS & AI
THOUGHT- LEADERS

OAK RIDGE
National Laboratory

CSC

PAWSEY
supercomputing centre

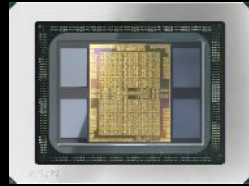GENCI

INES

Lawrence Livermore
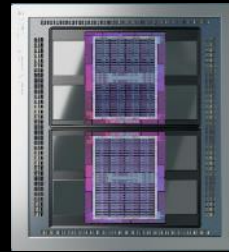National Laboratory

Microsoft Azure

**AMD together we advance_**

# OUR JOURNEY IN GPU ACCELERATION



### AMD Instinct™ MI100
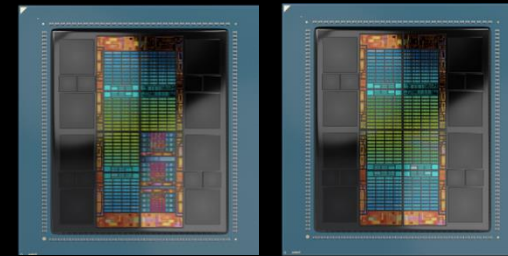AMD CDNA™

Ecosystem Growth

**First purpose-built GPU architecture for the data center**

### AMD Instinct™ MI200
AMD CDNA™ 2

Driving HPC and AI
to a New Frontier

**First purpose-built GPU powering discovery at Exascale**

### AMD Instinct™ MI300
AMD CDNA™ 3

Data Center APU
& Discrete GPU

**Breakthrough architecture designed for leadership efficiency and performance for AI and HPC**

AMD
together we advance_
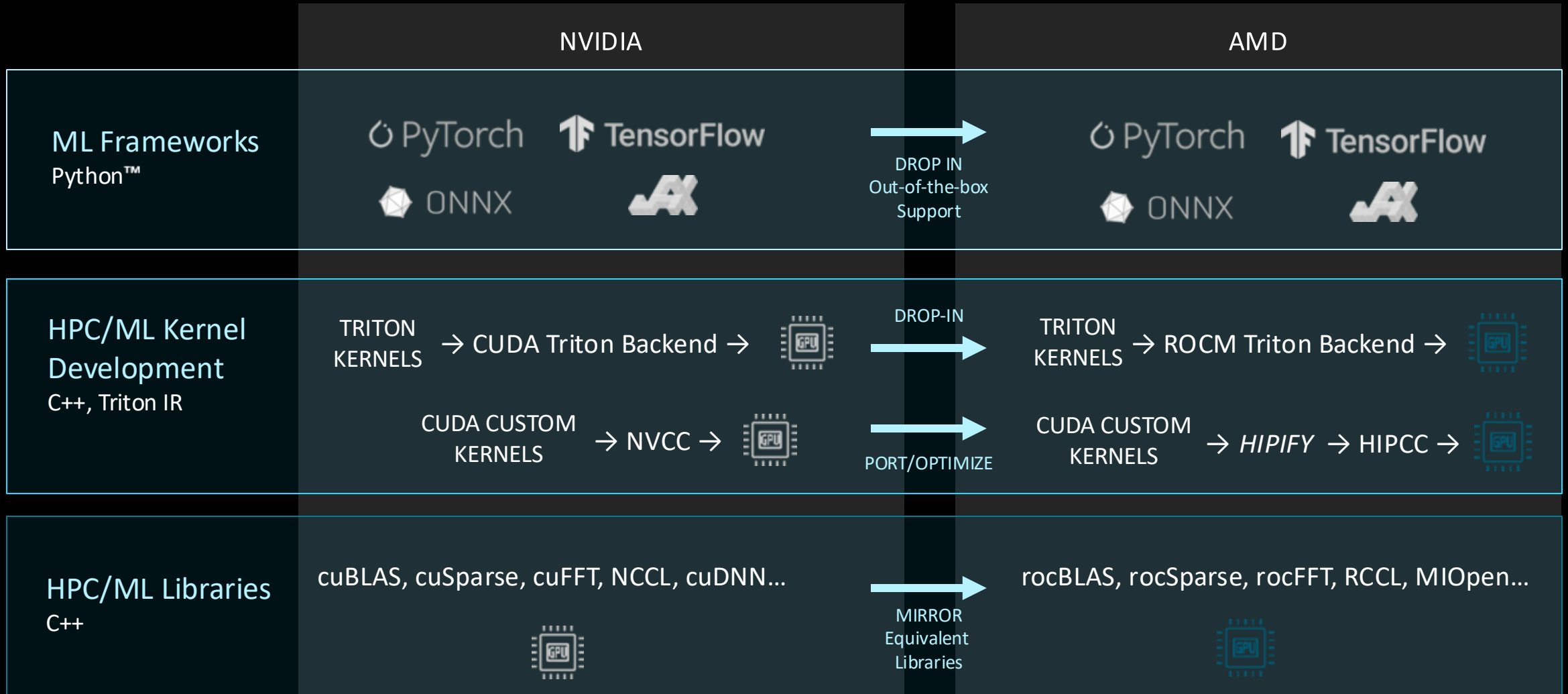
# OPEN SOFTWARE PLATFORM FOR GPU COMPUTE

**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<br>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 |

# TRANSITIONING WORKLOADS TO INSTINCT GPUS

## LOW FRICTION SOFTWARE PORTING FOR EXISTING NVIDIA USERS TO AMD

|  | NVIDIA | | AMD |
|---|---|---|---|
| **ML Frameworks**<br>Python™ | PyTorch  TensorFlow  ONNX  JAX | → DROP IN Out-of-the-box Support → | PyTorch  TensorFlow  ONNX  JAX |
| **HPC/ML Kernel Development**<br>C++, Triton IR | TRITON KERNELS → CUDA Triton Backend → [GPU]<br><br>CUDA CUSTOM KERNELS → NVCC → [GPU] | DROP-IN →<br><br>PORT/OPTIMIZE → | TRITON KERNELS → ROCM Triton Backend → [GPU]<br><br>CUDA CUSTOM KERNELS → *HIPIFY* → HIPCC → [GPU] |
| **HPC/ML Libraries**<br>C++ | cuBLAS, cuSparse, cuFFT, NCCL, cuDNN...  [GPU] | MIRROR Equivalent Libraries → | rocBLAS, rocSparse, rocFFT, RCCL, MIOpen...  [GPU] |

AMD
together we advance_

# AMD
# MI300A

## UNIFIED MEMORY APU ARCHITECTURE BENEFITS

**AMD CDNA™ 2 Coherent Memory Architecture** ➤ **AMD CDNA™ 3 Unified Memory APU Architecture**



- **Eliminate Redundant Memory Copies**

- **No programming distinction between host and device memory spaces**

- **High performance, fine-grained sharing between CPU and GPU processing elements**

- **Single process can address all memory, compute elements on a socket**

AMD
together we advance_

# APU PROGRAMMING MODEL

| CPU CODE | GPU CODE | APU CODE |
|---|---|---|

```
double* in_h = (double*)malloc(Msize);       double* in_h = (double*)malloc(Msize);       double* in_h = (double*)malloc(Msize);
double* out_h = (double*)malloc(Msize);       double* out_h = (double*)malloc(Msize);       double* out_h = (double*)malloc(Msize);
                                              hipMalloc(&in_d, Msize);
                                              hipMalloc(&out_d, Msize);


for (int i=0; i<M; i++) // initialize          for (int i=0; i<M; i++) // initialize          for (int i=0; i<M; i++) // initialize
   in_h[i] = …;                                   in_h[i] = …;                                   in_h[i] = …;
                                              hipMemcpy(in_d,in_h,Msize);
cpu_func(in_h, out_h, M);                     gpu_func<< >>(in_d, out_d, M);                gpu_func<< >>(in_h, out_h, M);
                                              hipDeviceSynchronize();                       hipDeviceSynchronize();
                                              hipMemcpy(out_h,out_d,Msize);


for (int i=0; i<M; i++) // CPU-process        for (int i=0; i<M; i++) // CPU-process        for (int i=0; i<M; i++) // CPU-process
   … = out_h[i];                                 … = out_h[i];                                 … = out_h[i];
```
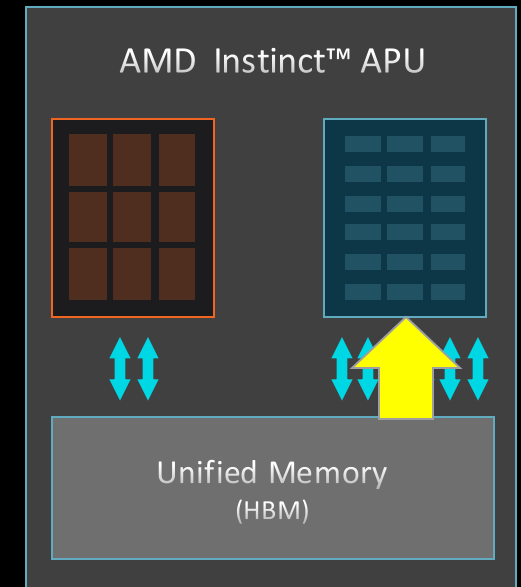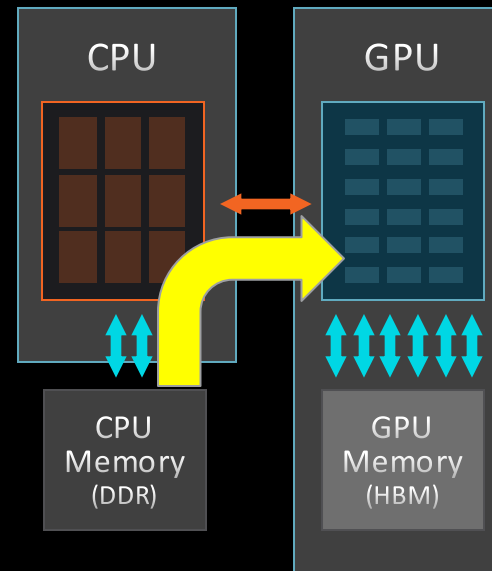
- GPU memory allocation on Device
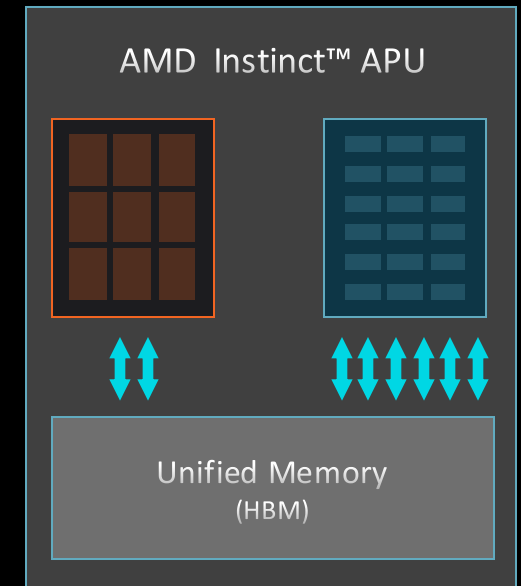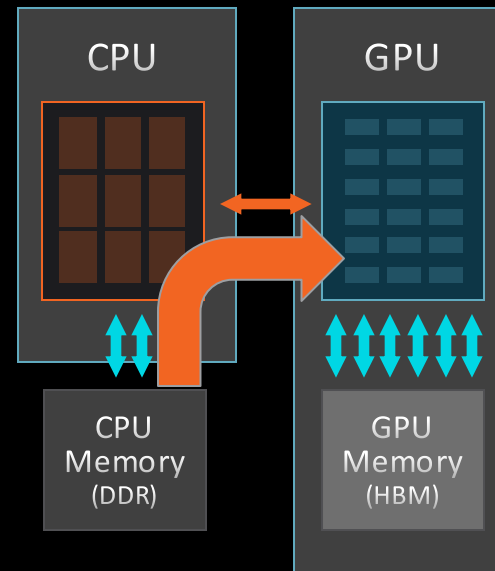- Explicit memory management between CPU & GPU
- Synchronization Barrier

AMD
together we advance_

# APU PROGRAMMING: PERFORMANCE IMPLICATIONS

## GPU CODE

```
double* in_h = (double*)malloc(Msize);
double* out_h = (double*)malloc(Msize);
hipMalloc(&in_d, Msize);
hipMalloc(&out_d, Msize);


for (int i=0; i<M; i++) //initialize
    in_h[i] = …;
hipMemcpy(in_d,in_h,Msize);
gpu_func<< >>(in_d, out_d, M);
hipDeviceSynchronize();
hipMemcpy(out_h,out_d,Msize);


for (int i=0; i<M; i++) // CPU-process
  … = out_h[i];
```



| Operation | MI250X (MCM) | MI300A |
|-----------|--------------|--------|
| H2D Copy  | O(10) GB/s   | O(TB/s) |

- GPU memory allocation on Device
- Explicit memory management between CPU & GPU
- Synchronization Barrier

**AMD**
together we advance_

# APU PROGRAMMING: PERFORMANCE IMPLICATIONS

## APU CODE

```
double* in_h = (double*)malloc(Msize);
double* out_h = (double*)malloc(Msize);



for (int i=0; i<M; i++) //initialize
    in_h[i] = …;

gpu_func<< >>(in_h, out_h, M);
hipDeviceSynchronize();


for (int i=0; i<M; i++) // CPU-process
    … = out_h[i];
```



| Operation | MI250X (MCM) | MI300A |
|---|---|---|
| Coherent Access | O(10) GB/s | N/A |

- ~~GPU memory allocation on Device~~
- ~~Explicit memory management between CPU & GPU~~
- Synchronization Barrier

AMD
together we advance_

# PROGRAMMING ACROSS FRAMEWORKS/COMPILERS

| OpenMP® CODE | RAJA CODE | KOKKOS CODE |
|---|---|---|

```
#pragma omp requires unified_shared_memory
 double* in_h = (double*)malloc(Msize);
double* out_h = (double*)malloc(Msize);




for (int i=0; i<M; i++) // initialize
   in_h[i] = …;


#pragma omp target
{ … }


for (int i=0; i<M; i++) // CPU-process
  … = out_h[i];
```

```
double* in_h = (double*)malloc(Msize);
double* out_h = (double*)malloc(Msize);




for (int i=0; i<M; i++) // initialize
   in_h[i] = …;


RAJA::forall< exec_policy >(arange, [=]
(int i) { … } );



for (int i=0; i<M; i++) // CPU-process
  … = out_h[i];
```

```
double* in_h = (double*)malloc(Msize);
double* out_h = (double*)malloc(Msize);




for (int i=0; i<M; i++) // initialize
   in_h[i] = …;


Kokkos::parallel_for(M, [=] (const  int
i){ … };
Kokkos::fence();



for (int i=0; i<M; i++) // CPU-process
  … = out_h[i];
```

- ~~GPU memory allocation on Device~~
- ~~Explicit memory management between CPU & GPU~~
- Synchronization Barrier

AMD
together we advance_

# HIP STANDARD PARALLELISM

- AMD providing support for advanced C++ in LLVM: today, entirely Open Source Software (OSS)
  - Only supports par_unseq acceleration currently
  - Compiler support implementation upstreamed
  - Available in recent ROCm releases
  - Re-uses HIP support in CLANG/LLVM and algorithms from libraries (rocThrust)
  - Available today: https://github.com/ROCmSoftwarePlatform/roc-stdpar
    - Several applications tested and running (LULESH, etc.)

```
1  std::transform( // needs <algorithm>
2      std::execution::par_unseq, // <-- needs <execution>
3      indices.begin(), indices.end(), grid.begin(),
4      [](size_t index){
5          return expensive_calculation(index);
6      }
7  );
```

AMD
together we advance_

# AMD PARTNERSHIPS IN HPC AND AI

- Centers of Excellence

- In 2023, AMD HPC teams supported many focused training activities
  - >100 days of training, >3000 participants
  - Trainings, hackathons, private hackathons, virtual hackathons

- Focused on application porting, tuning, and analysis
  - Learnings are used to influence our hardware and software roadmap (co-design)

- Papers and dissemination of best practices



- Commitment to Open Source Software





RESEARCH-ARTICLE

## Experiences readying applications for Exascale

Authors: Nicholas Malaya, Bronson Messer, Joseph Glenski, Antigoni Georgiadou, Justin Lietz, Kalyana Gottiparthi, Marc Day, Jackie Chen, Jon Rood, Lucas Esclapez, James White III, + 21
Authors Info & Claims

SC '23: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis • November 2023 • Article No.: 53 • Pages 1–13 • https://doi.org/10.1145/3581784.3607065

NEURAL INFORMATION PROCESSING SYSTEMS

LUMI

## Announcing the NeurIPS 2023 Paper Awards

**Scaling Data-Constrained Language Models**

Authors: Niklas Muennighoff · Alexander Rush · Boaz Barak · Teven Le Scao · Nouamane Tazi · Aleksandra Piktus · Sampo Pyysalo · Thomas Wolf · Colin Raffel

AMD together we advance_

# OVERSUBSCRIBING THE GPU

## IMPROVING KERNEL THROUGHPUT

- Instinct GPUs provide HW support to run concurrently multiple contexts

- Process isolation – management of virtual memory per process

- Number of queues can be controlled at runtime with the environment variable GPU_MAX_HW_QUEUES



- Throughput within a process
  - Multiple host threads concurrently use the same HIP context
  - Use case:
    - Several tiny kernels and data transfers running concurrently
    - Typical when CPU applications are being ported to GPUs

- Throughput across processes
  - Same GPU supporting multiple HIP contexts from different processes
  - Reliance on GPU driver to multiplex contexts
  - No need for software solutions to multiplex contexts in time

AMD
together we advance_

# OVERSUBSCRIBING THE GPU
## IMPROVING KERNEL THROUGHPUT – MULTIPLE THREADS

Process A

- Allowing a process to use more GPU HW queues (MI250X Instinct GPU)
  - export GPU_MAX_HW_QUEUES=4 (default)



Before: 4.578s

Same workload

Speedup is 1.9x

  - export GPU_MAX_HW_QUEUES=22



After: 2.497s

More overlap of the activity from different streams

Better GPU occupancy

AMD
together we advance_

# OVERSUBSCRIBING THE GPU

## IMPROVING KERNEL THROUGHPUT – MULTIPLE THREADS VS PROCESS

Maximum overlap of 4 kernels

Tail effect



Default max (4) HW queues

- We can optimize throughput and be less prone to tail effects by increasing number of queues properly
- Throughput observed with threading similar to throughput over multiple processes:



Max HW queues set to 16



Multiple processes collated profile

AMD together we advance_

# LEVERAGING THE APU PROGRAMMING MODEL

- APU removing bottle necks in OpenFOAM simulations
    - OpenMP to enable offloading larger portions of code
    - APU avoids time spend in data migration

- Tandon et al. *Porting HPC Applications to AMD Instinct MI300A Using Unified Memory and OpenMP*
    - https://arxiv.org/abs/2405.00436

AMD
together we advance_

# AMD ROCM DEVELOPER HUB



## Engage with AMD ROCm™ Platform Experts

Participate in AMD ROCm Webinar Series

Post questions, view FAQs in Community Forum

## Increase Understanding

Purchase AMD ROCm Text Book

View the latest news and good practices in the Blogs

## Get Started Using AMD ROCm™

AMD ROCm Documentation on GitHub

Download the Latest Version of AMD ROCm

AMD together we advance_

# DISCLAIMER AND ATTRIBUTIONS

**AMD**
together we advance_