# *A new perspective on processing-in-memory architecture design*

**Dong Ping Zhang, Nuwan Jayasena, Alexander Lyashevsky,
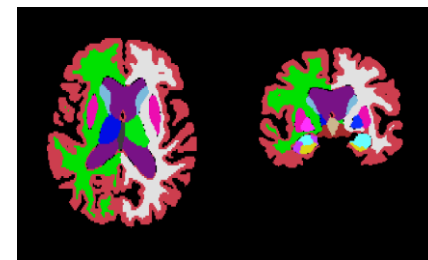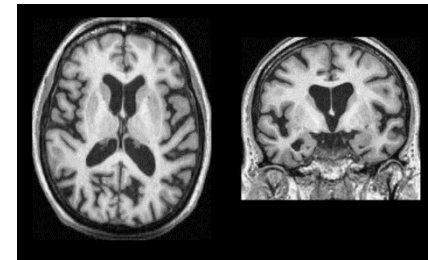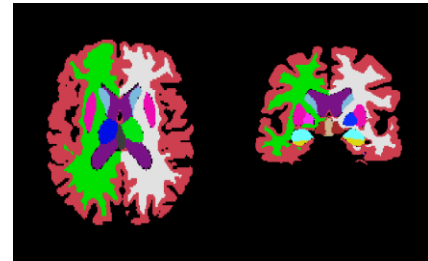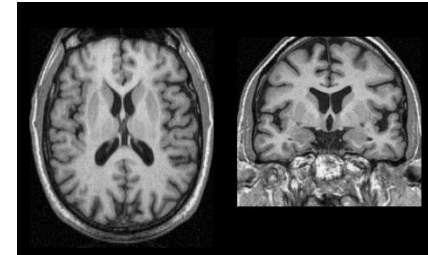Joe Greathouse, Mitesh Meswani, Mark Nutter, Mike Ignatowski**

**AMD Research**

**AMD**

# PROCESSING-IN-MEMORY?

# *HIGHLIGHT*

- Memory system is a key limiter

  – At 4TB/s, vast majority of node energy could be consumed by the memory system

- Prior PIM research constrained by

  – Implementation technology

  – Non-traditional programming models

- Our focus

  – 3D die stacking

  – Use base logic die(s) in memory stack

    - General-purpose processors

    - Support familiar programming models

    - Arbitrary programs vs a set of special operations

**Memory dies**

**Memory stack**

**Logic die with PIM**

AMD

# *PROCESSING-IN-MEMORY (PIM)   ---   OVERVIEW (I)*

- Moving compute close to memory promises significant gains
    - Memory is a key limiter (performance and power)
    - Exascale goals of 4TB/s/node and <5pj/bit
- Prior research
    - Integration of caches and computation
        - "A logic-in-memory computer" (1970)
        - No large scale integration was possible.
    - Logic in DRAM processes
        - In-memory processors with reduced performance or highly specialized for a limited set of operations.
        - Reduced DRAM due to the presence of compute logic.
    - Embedded DRAM in logic processes
        - Not cost-effectively accommodate sufficient memory capacity
        - Reduced density of embedded memory

AMD

# PROCESSING-IN-MEMORY (PIM) --- OVERVIEW (II)

- New opportunity: logic die stacked with memory
  - Logic die needed anyway for signal redistribution and integrity
  - Potential for non-trivial compute
- Key benefits
  - Reduce bandwidth bottlenecks
  - Improve energy efficiency
  - Increase compute for a fixed interposer area
  - Processor can be optimized for high BW/compute ratio
- Challenges
  - Programming models and interfaces
  - Architectural tradeoffs
  - Application refactoring

AMD

# *OUTLINE*

- PIM architecture baseline

- API specification

- Emulator and performance models

- Application studies

AMD

# NODE HARDWARE ORGANIZATION

- Single-level of PIM-attached memory stacks
- Host has direct access to all memory
  - Non-PIM-enabled apps still work
- Unified virtual address space
  - Shared page-tables between host and PIMs
- Low-bandwidth inter-PIM interconnect

Memory dies

Memory stack

Logic die with PIM

Host Processor

Interposer or board

AMD

# PIM API OVERVIEW

- Current focus is on single PIM-enabled node
  - Current PIM hardware baseline is general-purpose processor
- Key goals of API
  - Facilitate data layout control
  - Dispatch compute to PIMs using standard parallel abstractions
- A "convenient" level of abstraction for early PIM evaluations
  - Aimed at PIM feasibility studies
  - annotation of data and compute for PIM with reasonable programmer effort
- Key API features:
  - discover device
  - query device characteristics
  - manage locality
  - dispatch compute

AMD

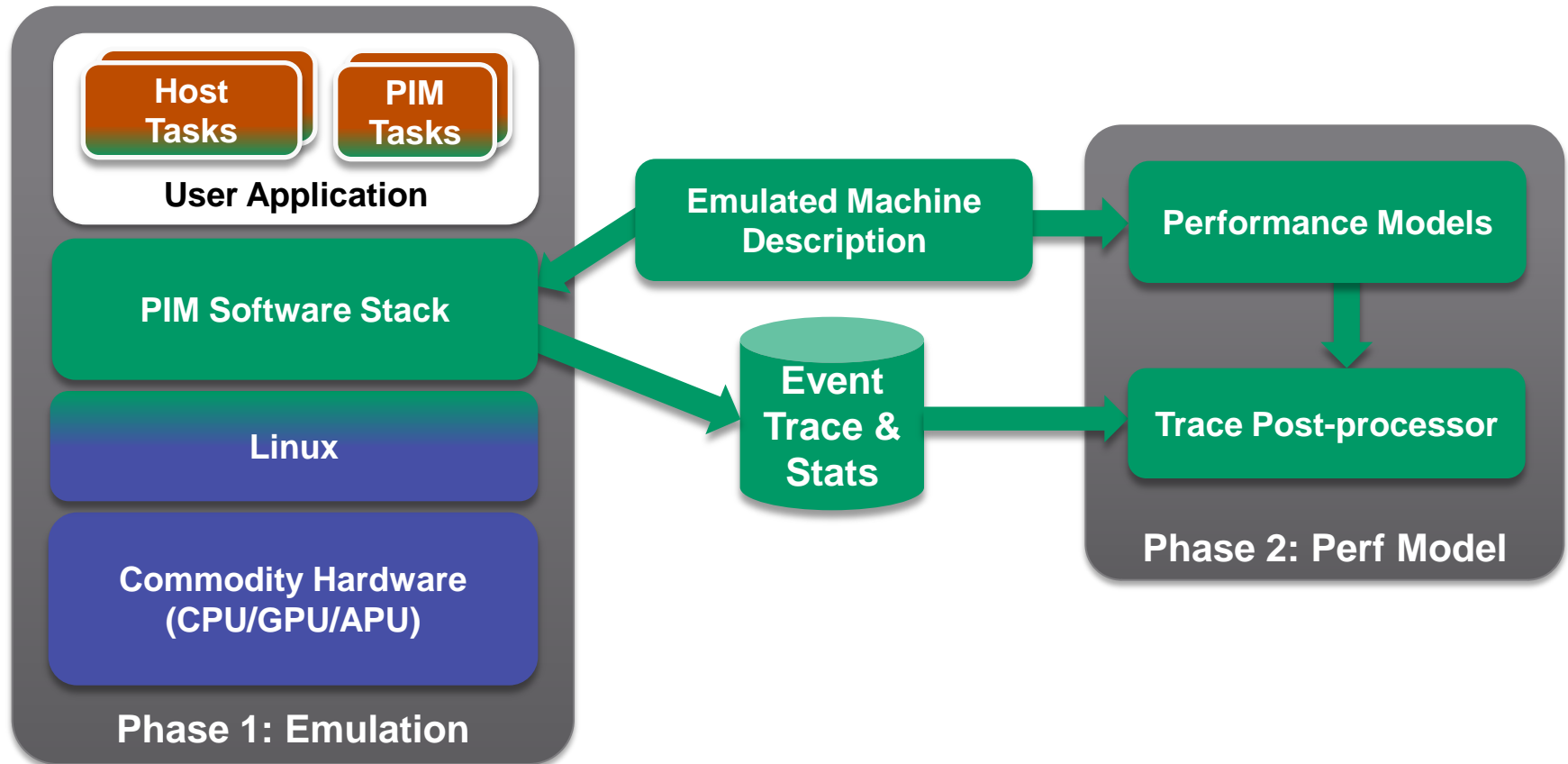# PIM PTHREAD EXAMPLE: PARALLEL PREFIX SUM

```
…
list_of_pims = malloc(max_pims * sizeof(pim_device_id));
failure = pim_get_device_ids(PIM_CLASS_0, max_pims, list_of_pims, &num_pims);
for (i = 0; i < num_pims; i++) {
    failure = pim_get_device_info(list_of_pims[i], PIM_CPU_CORES, needed_size, device_info, NULL);
    …
}
…
for (i = 0; i < num_pims; i++) {
    parallel_input_array[i] = pim_malloc(sizeof(uint32_t) * chunk_size, list_of_pims[i],
                                    PIM_MEM_DEFAULT_FLAGS, PIM_PLATFORM_PTHREAD_CPU);
    parallel_output_array[i] = pim_malloc(sizeof(uint64_t) * chunk_size, list_of_pims[i],
                                    PIM_MEM_DEFAULT_FLAGS, PIM_PLATFORM_PTHREAD_CPU);
}
…
for (i = 0; i < num_pims; i++) {
    pim_args[PTHREAD_ARG_THREAD] = &(pim_threads[i]); // pthread_t
    arg_size[PTHREAD_ARG_THREAD] = sizeof(pthread_t);
    pim_args[PTHREAD_ARG_ATTR] = NULL; // pthread_attr_t
    arg_size[PTHREAD_ARG_ATTR] = sizeof(pthread_attr_t);
    pim_args[PTHREAD_ARG_INPUT] = &(thread_input[i]); // void * for thread input
    arg_size[PTHREAD_ARG_INPUT] = sizeof(void *);
    pim_function.func_ptr = parallel_prefix_sum;
    spawn_error = pim_spawn(pim_function, pim_args, arg_size, NUM_PTHREAD_ARGUMENTS,
                        list_of_pims[i], PIM_PLATFORM_PTHREAD_CPU);
}
for (i = 0; i < num_pims; i++) {
    pthread_join(pim_threads[i], NULL);
}
…
```

**AMD**

# PIM EMULATION AND PERFORMANCE MODEL

**Host Tasks** · **PIM Tasks**

**User Application**

**PIM Software Stack**

**Linux**

**Commodity Hardware (CPU/GPU/APU)**

**Phase 1: Emulation**

**Emulated Machine Description**

**Event Trace & Stats**

**Performance Models**

**Trace Post-processor**

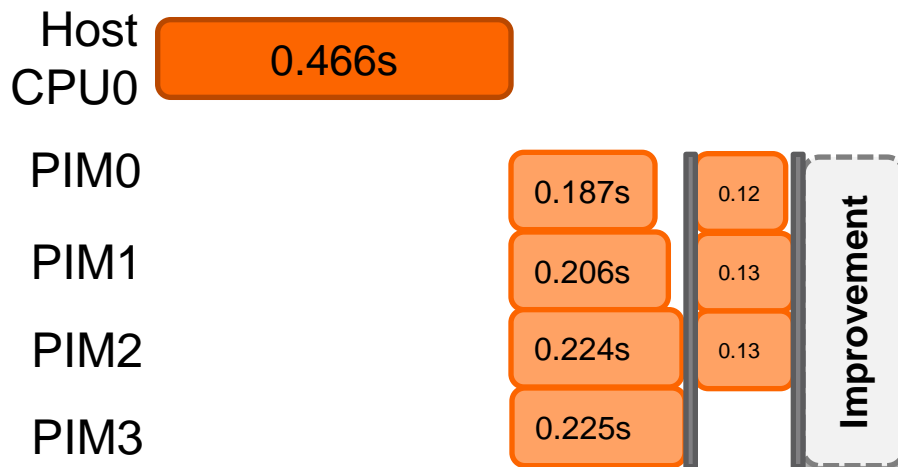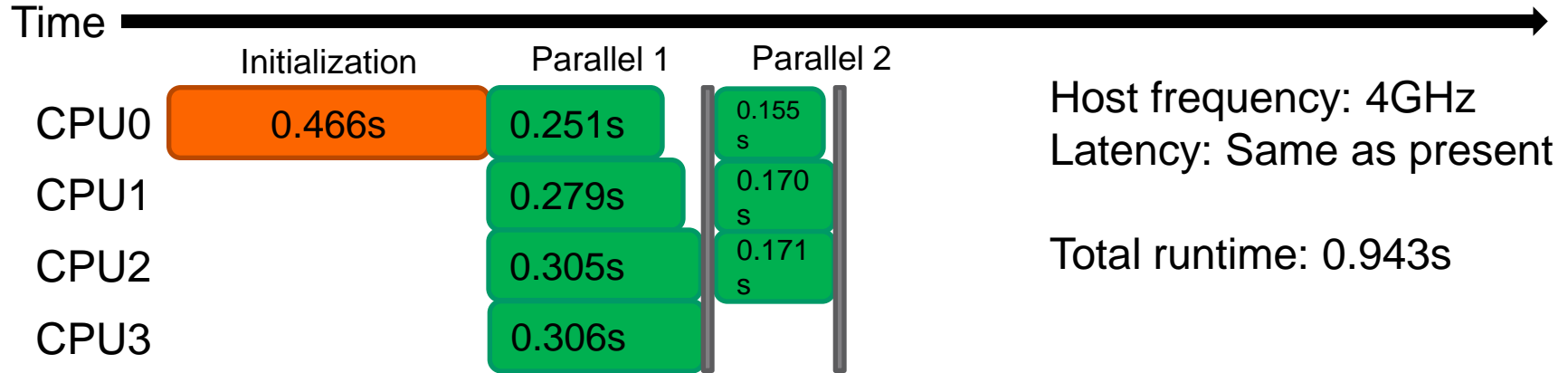**Phase 2: Perf Model**

- Phase 1: Native execution on commodity hardware
  - Capture execution trace and performance stats
- Phase 2: Post-process with performance models
  - Predict overall performance on future memory and processors

**AMD**

# *PARALLEL PREFIX SUM – HOST-ONLY VS HOST+PIM*

Time →

| | Initialization | Parallel 1 | Parallel 2 |
|---|---|---|---|
| CPU0 | 0.466s | 0.251s | 0.155 s |
| CPU1 | | 0.279s | 0.170 s |
| CPU2 | | 0.305s | 0.171 s |
| CPU3 | | 0.306s | |

Host frequency: 4GHz
Latency: Same as present

Total runtime: 0.943s

---

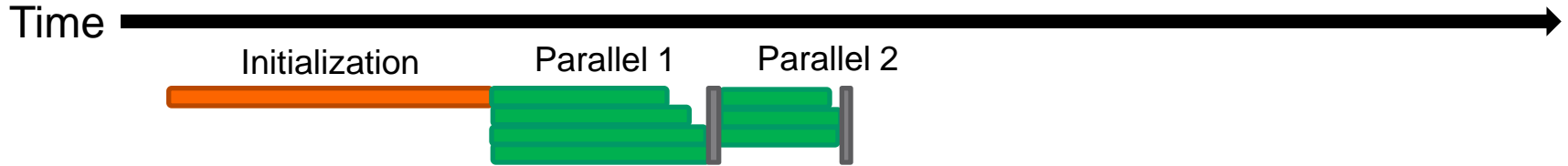| | | | | |
|---|---|---|---|---|
| Host CPU0 | 0.466s | | | |
| PIM0 | | 0.187s | 0.12 | |
| PIM1 | | 0.206s | 0.13 | **Improvement** |
| PIM2 | | 0.224s | 0.13 | |
| PIM3 | | 0.225s | | |

Host frequency: 4GHz
PIM frequency: 2GHZ
Host latency: current
PIM latency: 30% reduction

Total runtime: 0.820s
           15% Faster

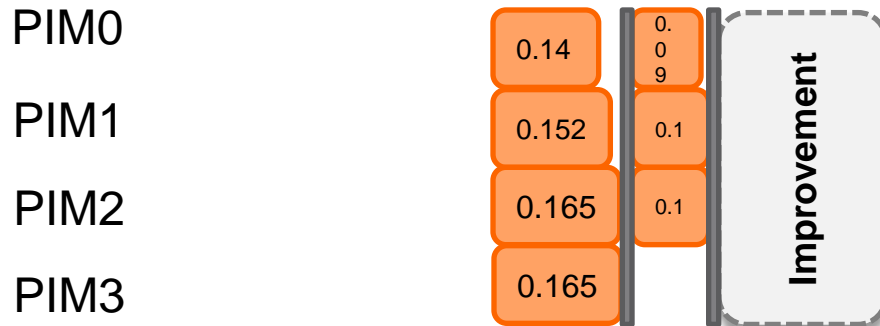*PIM computation benefits from reduced latency*

AMD

# PARALLEL PREFIX SUM – HOST-ONLY VS HOST+PIM

Time

Initialization    Parallel 1    Parallel 2

Host
CPU0    0.466s

PIM0    0.14    0.09

PIM1    0.152    0.1

PIM2    0.165    0.1

PIM3    0.165

Improvement
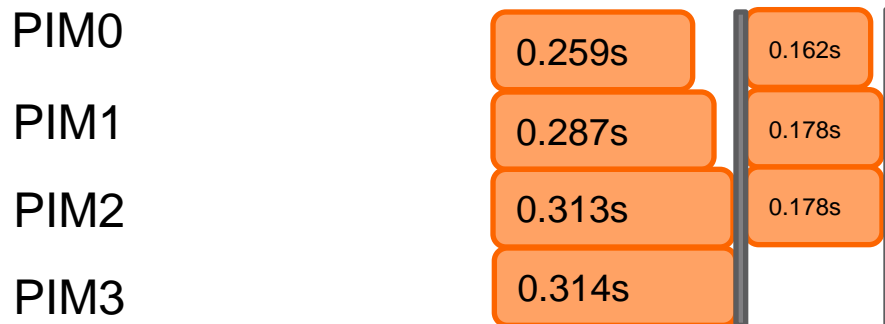
Host frequency: 4GHz
PIM frequency: 2GHZ
Host latency: current
PIM latency: 50% reduction

Total runtime: 0.728s
        30% Faster

Host
CPU0    0.466s

PIM0    0.259s    0.162s

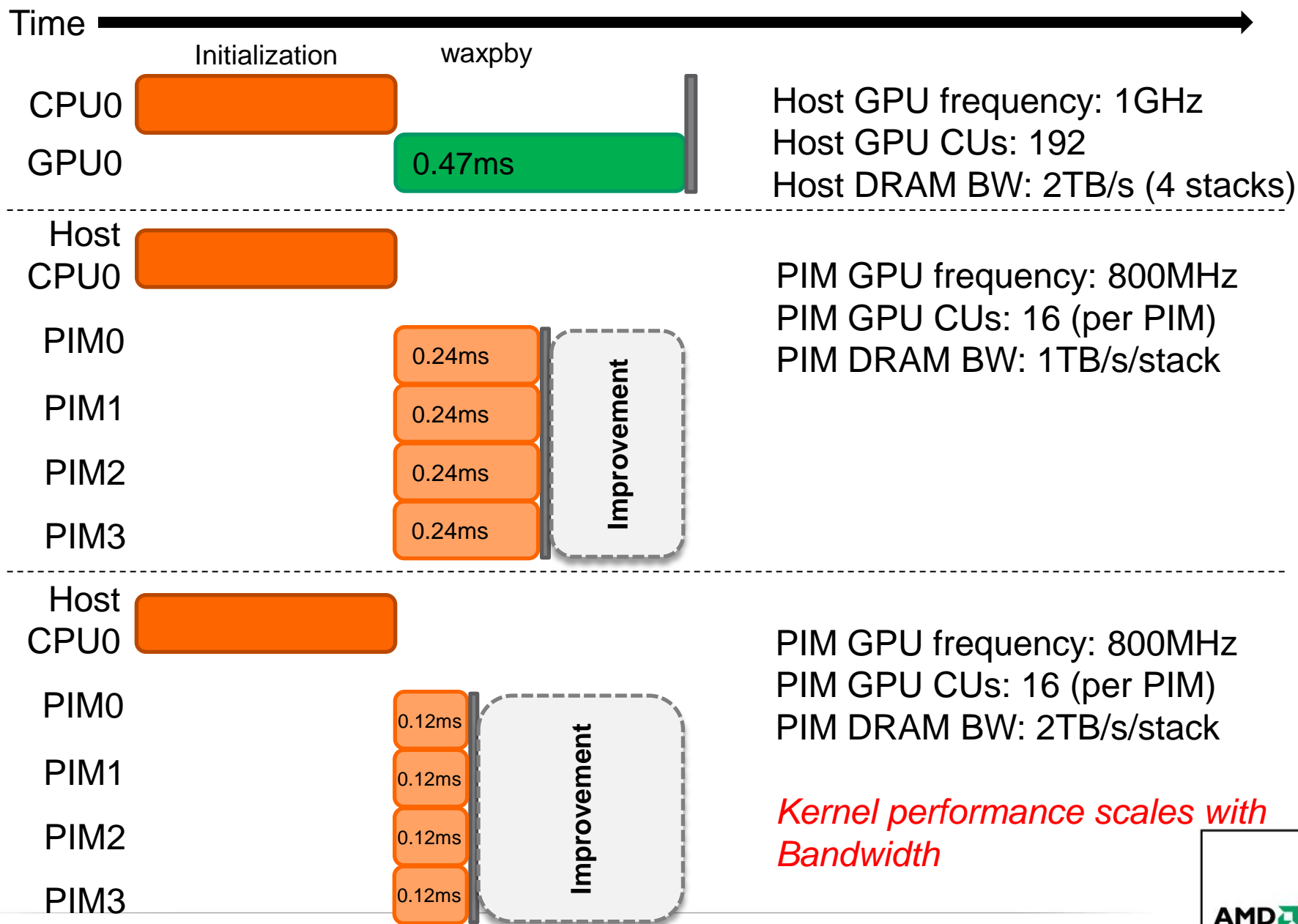PIM1    0.287s    0.178s

PIM2    0.313s    0.178s

PIM3    0.314s

Host frequency: 4GHz
PIM frequency: 2GHZ
Host latency = PIM latency
Total runtime: 0.958s
        2% Slower

AMD

# WAXPBY – HOST-ONLY VS HOST+PIM

Time →

Initialization    waxpby

CPU0

GPU0    **0.47ms**

Host GPU frequency: 1GHz
Host GPU CUs: 192
Host DRAM BW: 2TB/s (4 stacks)

Host
CPU0

PIM0    0.24ms

PIM1    0.24ms

PIM2    0.24ms

PIM3    0.24ms

Improvement

PIM GPU frequency: 800MHz
PIM GPU CUs: 16 (per PIM)
PIM DRAM BW: 1TB/s/stack

Host
CPU0

PIM0    0.12ms

PIM1    0.12ms

PIM2    0.12ms

PIM3    0.12ms

Improvement

PIM GPU frequency: 800MHz
PIM GPU CUs: 16 (per PIM)
PIM DRAM BW: 2TB/s/stack

*Kernel performance scales with Bandwidth*

AMD

# SUMMARY AND FUTURE WORK

- PIM architecture baseline

- API specification

- Emulator and performance models

- Application studies

- Design space study

- Evaluation of the performance models

- Further work on API, execution model, applications etc.

AMD

## ACKNOWLEDGEMENT:

Lee Howes

Gabe Loh

## QUESTIONS?

## FURTHER DISCUSSIONS:

Dongping.zhang@amd.com

**AMD**