# Architecture and RISC-V ISA Extension Supporting Asynchronous and Flexible Parallel Far Memory Access

**Songyue Wang**, Luming Wang, Tianyue Lu, Mingyu Chen
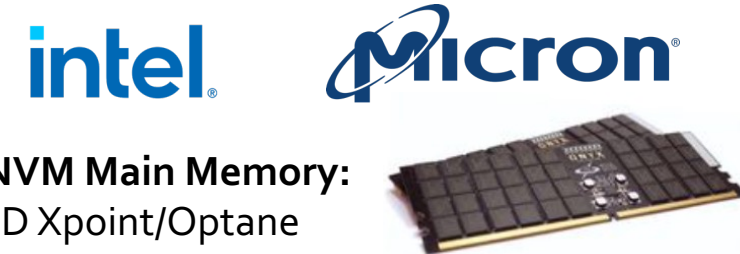
ICT, CAS

University of Chinese Academy of Sciences

June 19, 2022 @CARRV'22, co-located with ISCA'22, New York City

University of Chinese Academy of Sciences
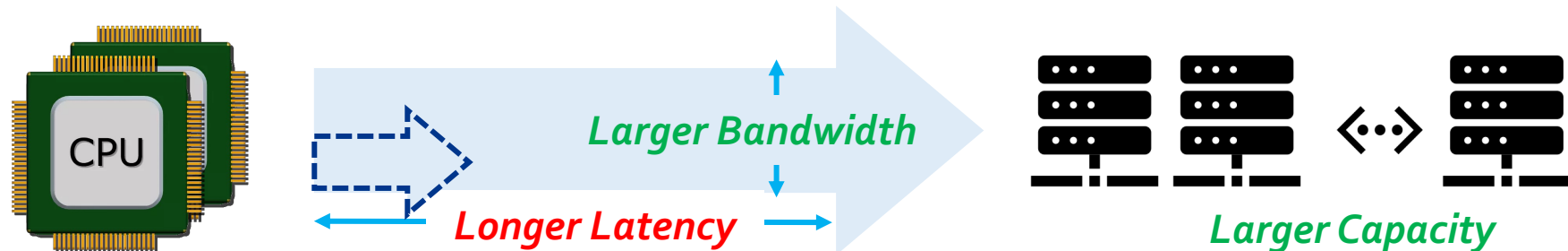
INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

先进计算机系统研究中心
Center for Advanced Computer Systems

# Far Memory Technology (FMT)

New Interconnect Techologies and Protocols

**NVM Main Memory:**
3D Xpoint/Optane

New Medium besides DRAM

- **Allow one host to access more memory resources**
  - **Large aggregated capacity and bandwidth** ☺
  - **Widely distributed latency** ☹ **(reach 300ns~5$\mu s$)**

- **Leave a challenge for CPU to fully utilize the abundant memory resources!**

CPU

*Larger Bandwidth*
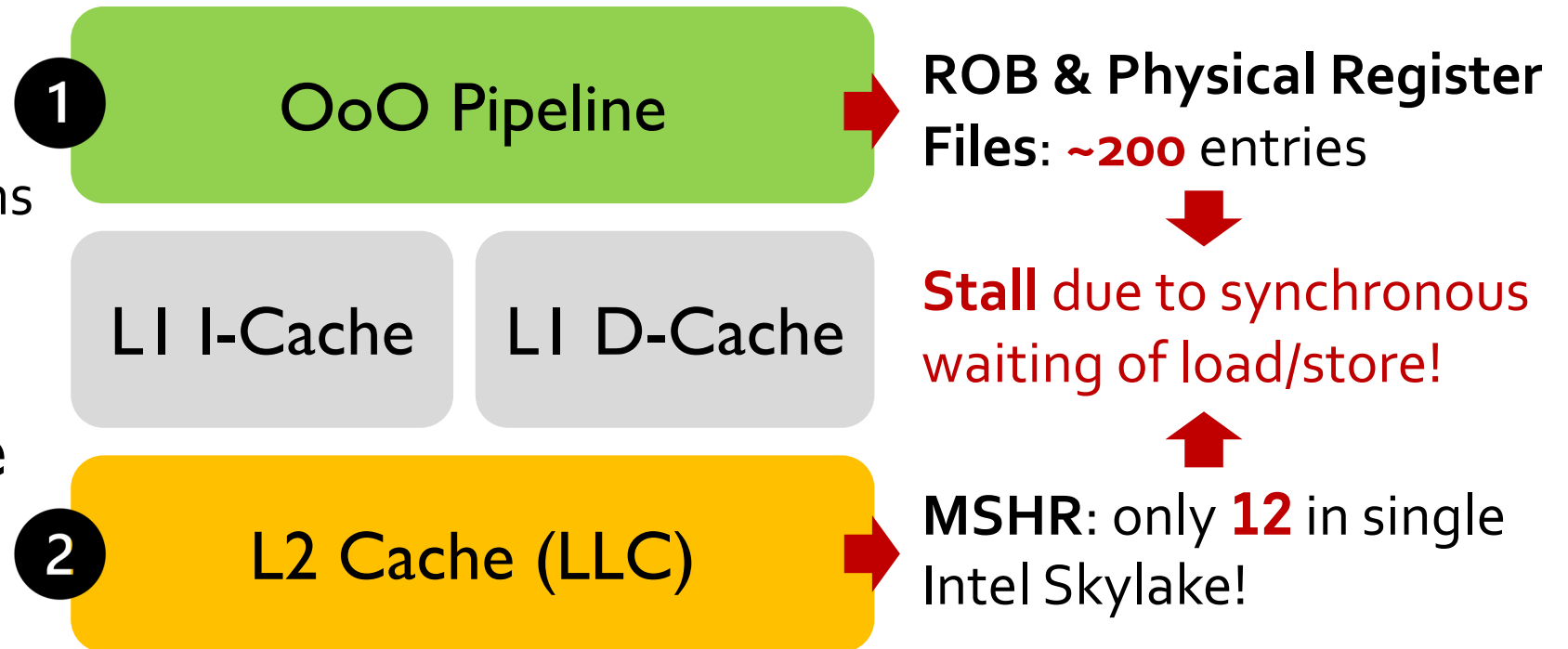
*Longer Latency*

*Larger Capacity*

# Modern CPU cannot Reach High MLP

- Additional **64ns** latency → **30%+** slowdown!
  - Especially applications hich **benefit from FMT**! (such as Redis, GAPBS and Spark)

- Modern CPUs **implicitly** boost MLP:

## OoO Execution
Dynamical scheduling more load/store instructions

**①** OoO Pipeline

**ROB & Physical Register Files**: **~200** entries

L1 I-Cache     L1 D-Cache

**Stall** due to synchronous waiting of load/store!

## Non-Blocking Cache
Handling several miss memory requests

**②** L2 Cache (LLC)

**MSHR**: only **12** in single Intel Skylake!

# Modern CPU cannot Reach High MLP

- Additional **64ns** latency → **30%+** slowdown!
  - Especially applications hich **benefit from FMT**! (such as Redis, GAPBS and Spark)
- Modern CPUs **implicitly** boost MLP:

OoO Execution

Dynamical scheduling more load/store instructions

Non-Blocking Cache

Handling several miss memory requests

**1** OoO Pipeline

**ROB & Physical Register Files:** ~200 entries

L1 I-Cache    L1 D-Cache    Stall due to synchronous waiting of load/store!

**2** L2 Cache (LLC)    MSHR: only **12** in single Intel Skylake!

A **2GHz** CPU faced with **1$\mu s$** latency, to avoid stall:
~2000 entries **ROB** and **Physical RF**
**thousand** of items in **MSHR**
Criticial Resources Limit MLP!

# Key Observation

**The bottleneck to boost MLP**

- **ISA**
  - *Load & store* are **synchronous** and **blocking**

- **Microarchitecture**
  - **Request** and **response** are **coupling**

- **Storage**
  - **Limited MSHR** for handling **status**
  - **Limited Physical RF** for **data storage**

- **Semantic**
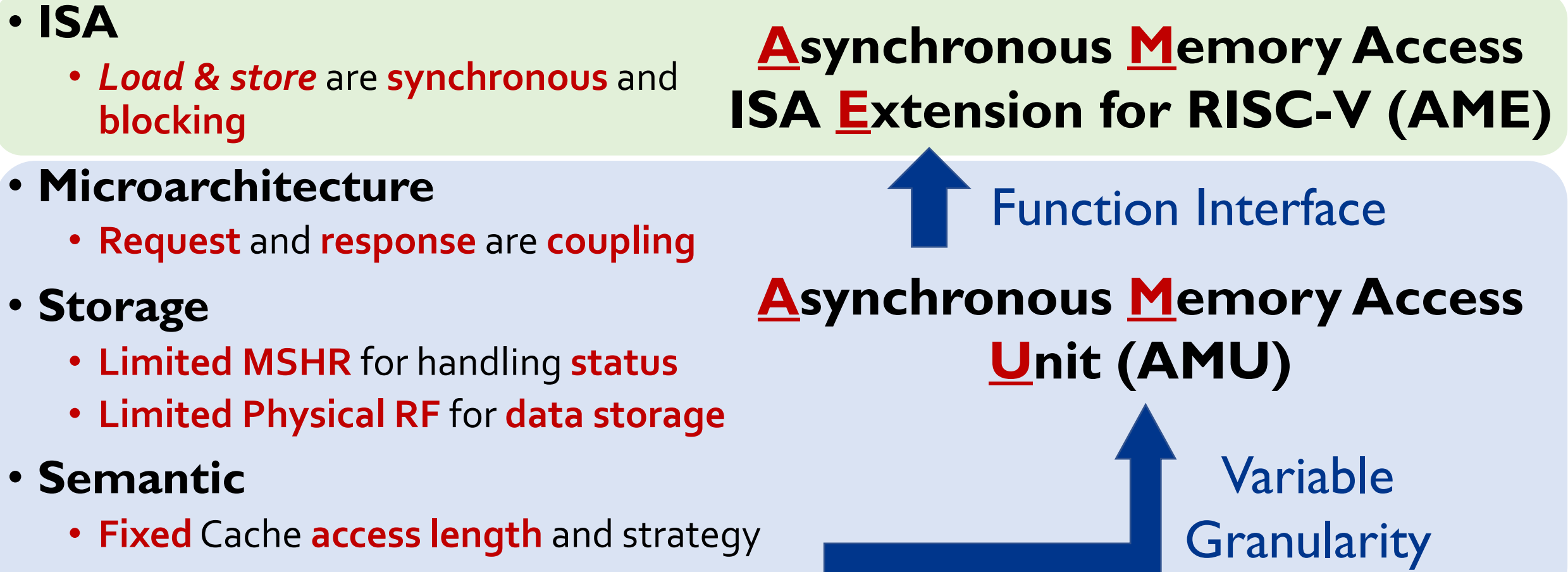  - **Fixed** Cache **access length** and strategy

# Our Goal

**Our Goal: Supporting Massive Parallel and Flexible Memory Access**

- **ISA**
  - *Load & store* are **synchronous** and **blocking**

  **Asynchronous Memory Access ISA Extension for RISC-V (AME)**

- **Microarchitecture**
  - **Request** and **response** are **coupling**

- **Storage**
  - **Limited MSHR** for handling **status**
  - **Limited Physical RF** for **data storage**

- **Semantic**
  - **Fixed** Cache **access length** and strategy

# Our Goal

## Our Goal: Supporting Massive Parallel and Flexible Memory Access

- **ISA**
  - *Load & store* are **synchronous** and **blocking**

  **Asynchronous Memory Access ISA Extension for RISC-V (AME)**

- **Microarchitecture**
  - **Request** and **response** are **coupling**
- **Storage**
  - **Limited MSHR** for handling **status**
  - **Limited Physical RF** for **data storage**

  **Asynchronous Memory Access Unit (AMU)**

- **Semantic**
  - **Fixed** Cache **access length** and strategy

# Our Goal

**Our Goal: Supporting Massive Parallel and Flexible Memory Access**

- **ISA**
  - *Load & store* are **synchronous** and **blocking**

- **Microarchitecture**
  - **Request** and **response** are **coupling**

- **Storage**
  - **Limited MSHR** for handling **status**
  - **Limited Physical RF** for **data storage**

- **Semantic**
  - **Fixed** Cache **access length** and strategy

**Asynchronous Memory Access ISA Extension for RISC-V (AME)**

Function Interface

**Asynchronous Memory Access Unit (AMU)**

Variable Granularity

# Outline

Background

**Asynchronous Memory access ISA Extensions (AME)**

Programing Model

Evaluation

# Key Idea 1: Asynchronous memory access ISA (AME)

① **ISA: A**synchronous **M**emory access ISA **E**xtension (AME)

- Core: ***aload*** and ***astore*** instructions

- **Step 1**: Write the request item to *AMQ*
  - **A**synchronous **M**emory request **Q**ueue
- **Step 2**: Commit!

- **Release the pressure of ROB**
- **Allow issuing massive and concurrent memory request from pipeline**

Time ⟶ ☐ overhead ☐ load/store ■ calculation

Traditional Load/Store Memory Access

number of mshr/ROB/IQ entries

cpu stall time

wait previous instructions to be commited

CPU

request

data

memory pool

MEM1

MEM2

......

Asynchronous Memory Access

CPU

request

memory pool

MEM1

MEM2

......

# Key Idea 2: Asynchronous Memory Access Unit (AMU)

**②** **Microarchitecture: AMU** (2 queues and 2 FSM)

- **An unique ID** for Each request
- Aload/astore commits to **AMQ** (**A**synchronous **M**emory access Request **Q**ueue)
- **Send FSM** issues according to items in **AMQ**
- **Recv FSM** pushes response IDs to **AFQ** (**A**synchronous Memory access **F**inish **Q**ueue)
- Software retrieves ID from **AFQ**

# Key Idea 3: ScratchPad Memory for AMU

**③ Storage:** ScratchPad Memory (SPM)

- Divided from **L2 Cache Data Array ( ≥128KB)**
- Compatibility: Adjustable by Cache way

- **Metadata Region**
  - AMQ, AFQ and FIN_META     **~ MSHR**
  - Adjustable length

- **Data Region**
  - For variable granularity memory data     **~ Physical RF**

**Supporting thousand of outstanding memory access!**

|  | WayN | ... | Way1 | Way0 |
|---|---|---|---|---|
| Set 0 |  |  |  |  |
| Set 1 |  |  | SPM |  |
|  |  |  | ←→ Region |  |
| Set M |  |  |  |  |

L2 Cache
Data Array

**Metadata**
(AMQ, AFQ,
FIN_META)

**Data**
(Byte~KB)

# AME: Programming Interface of AMU

- AME Core Instructions
  - `aload  id, spm_addr, mem_addr` :  move Memory → SPM
  - `astore id, spm_addr, mem_addr` :  move SPM → Memory
  - `asetid id` :  Set ID for the next `aload`/`astore`
  - `getfin id` :  Get a **finished ID** from **AFQ**
- AMU Control Registers
  - `SPMWAY` :  Which ways of L2 Cache are continuously occupied by SPM.
  - `RWLEN` :  Granularity of accessing memory (8-byte alignment)
  - `AMQLEN` :  Length of AMQ (max concurrency of memory access via AMU)

# Microarchitecture of AMU

- Pipeline

- Memory Access Path

- L2 Cache
  - Send/Recv FSM
  - AMQ/AFQ/FIN_META

# Process of an asynchronous memory access

**Pipeline**  **_aload_**

**1**

**AMU Metadata**

| ID | mem addr | type | len | spm addr |
|----|----------|------|-----|----------|

**AMQ**

**AMU FSM**

**Memory**

# Process of an asynchronous memory access

**Pipeline**

| aload | *other* | *other* |

**①**

**AMU Metadata**

| ID | mem addr | type | len | spm addr |

**AMQ**

**②**

**AMU FSM**

Send FSM

**Memory**

# Process of an asynchronous memory access

**Pipeline**

| aload | other | other | ...... |

**1**

**AMU Metadata**

| ID | mem addr | type | len | spm addr |

AMQ

| type | len | spm_addr |
| --- | --- | --- |
| ⋮ | | ⋮ |
| type | len | spm_addr |

FIN_META

**2**

**3**

**AMU FSM**

Send FSM

**Memory**

# Process of an asynchronous memory access

# Process of an asynchronous memory access

# Process of an asynchronous memory access

# Outline

# Basic Paradigm

## Step 1: Config AMU
- Max parallelism (length of AMQ)
- Access granularity

```c
#define MAX_PARALLELISM 256

int *mem_to_access; // memory to be accessed

// AMU Configuration
acfgwr(MAX_PARALLELISM, AMQLEN);
acfgwr(sizeof(int), RWLEN);

int *spm_data_area = (int *)alloc_spm_addr(sizeof(int));
int id = 1;  // alloc an ID
// issue an aload request
aload(id, spm_data_area, &far_mem_to_access);
// process other
while(id != getfin()) { /* process other */ }
// access SPM via standard load/store
printf("%d\n", *spm_space);
```

# Basic Paradigm

## Step 1: Config AMU
- Max parallelism (length of AMQ)
- Access granularity

## Step 2: Issue an aload/astore request
- Allocate SPM space
- Allocate an ID

```c
#define MAX_PARALLELISM 256

int *mem_to_access; // memory to be accessed


// AMU Configuration
acfgwr(MAX_PARALLELISM, AMQLEN);
acfgwr(sizeof(int), RWLEN);


int *spm_data_area = (int *)alloc_spm_addr(sizeof(int));
int id = 1;  // alloc an ID
// issue an aload request
aload(id, spm_data_area, &far_mem_to_access);
// process other
while(id != getfin()) { /* process other */ }
// access SPM via standard load/store
printf("%d\n", *spm_space);
```

# Basic Paradigm

**Step 1: Config AMU**
- Max parallelism (length of AMQ)
- Access granularity

**Step 2: Issue an aload/astore request**
- Allocate SPM space
- Allocate an ID

**Step 3: Wait for finish**
- Use *getfin* for checking

```c
#define MAX_PARALLELISM 256

int *mem_to_access; // memory to be accessed


// AMU Configuration
acfgwr(MAX_PARALLELISM, AMQLEN);

acfgwr(sizeof(int), RWLEN);


int *spm_data_area = (int *)alloc_spm_addr(sizeof(int));

int id = 1;   // alloc an ID

// issue an aload request
aload(id, spm_data_area, &far_mem_to_access);

// process other
while(id != getfin()) { /* process other */ }

// access SPM via standard load/store
printf("%d\n", *spm_space);
```

# Basic Paradigm

## Step 1: Config AMU
- Max parallelism (length of AMQ)
- Access granularity

## Step 2: Issue an aload/astore request
- Allocate SPM space
- Allocate an ID

## Step 3: Wait for finish
- Use *getfin* for checking

## Step 4: Access
- Via standard load/store.

```c
#define MAX_PARALLELISM 256

int *mem_to_access; // memory to be accessed


// AMU Configuration
acfgwr(MAX_PARALLELISM, AMQLEN);
acfgwr(sizeof(int), RWLEN);


int *spm_data_area = (int *)alloc_spm_addr(sizeof(int));
int id = 1;   // alloc an ID
// issue an aload request
aload(id, spm_data_area, &far_mem_to_access);
// process other
while(id != getfin()) { /* process other */ }
// access SPM via standard load/store
printf("%d\n", *spm_space);
```
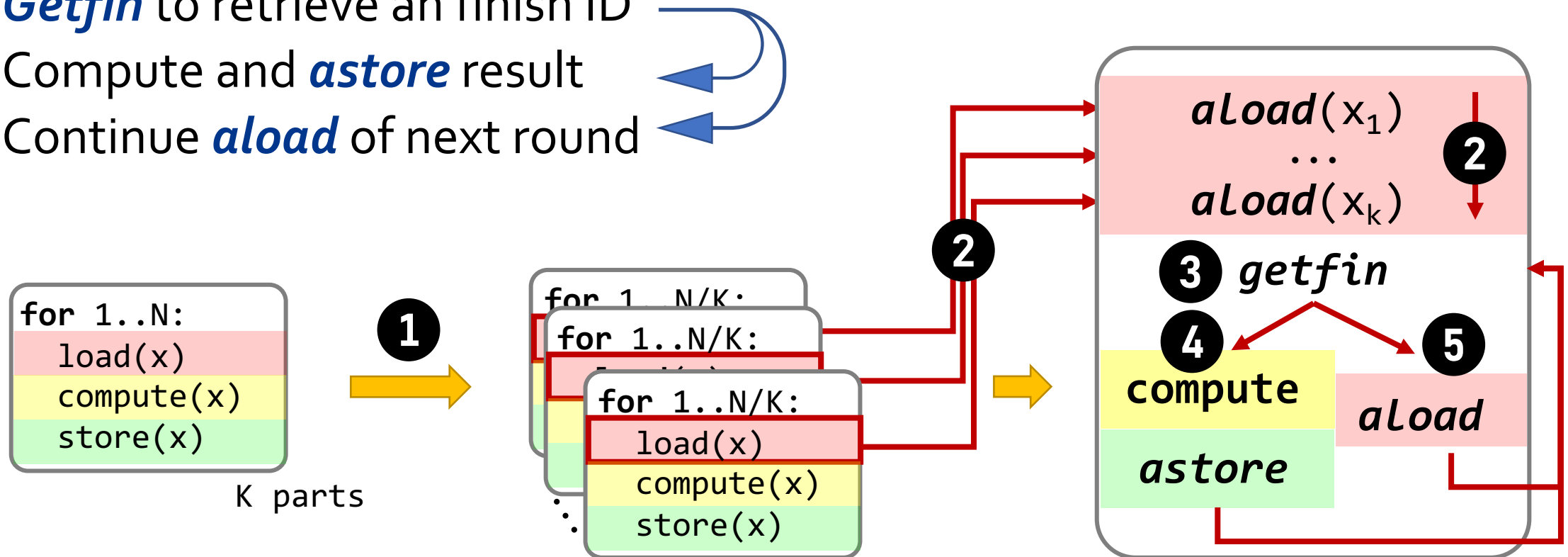
# Issue Massive Requests with AME

- Suitable for **data parallelized process**
  - ❶ Divide into K parts
  - ❷ Issue *aload* of all parts in batch
  - ❸ *Getfin* to retrieve an finish ID
  - ❹ Compute and *astore* result
  - ❺ Continue *aload* of next round

# CAP: Coroutines for AME Programming

- Idea: AME ~ Linux async I/O

- Based on **C++20 Coroutines**
  - **co_await** for async event

- Features
  - ID & SPM management
  - Lock
  - Scheduler (driven by *getfin*)

- **Efficiency**
  - Similar to **multi-thread programming**
  - Reduce **80% lines of** code
  - No care about scheduler (and *getfin*)!

```cpp
for (int i = 0; i < n; ++i)
  L[i] ^= i;
```

```cpp
template<typename Scheduler>
coro::task<void> update (int idx, int *L,
        int eachNUPDATE, Scheduler &sched){

  int *spm_addr = sched.alloc_spm_addr();

  for (int i = idx; i < idx + eachNUPDATE; ++i) {
    co_await aload_coro(spm_addr, &L[i], sched);
    *spm_addr ^= i;
    co_await astore_coro(spm_addr, &L[i], sched);
  }

  sched.release_spm(spm_addr);
}
```
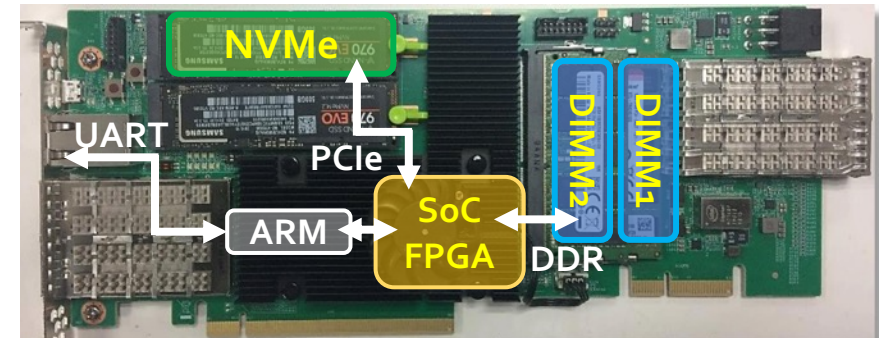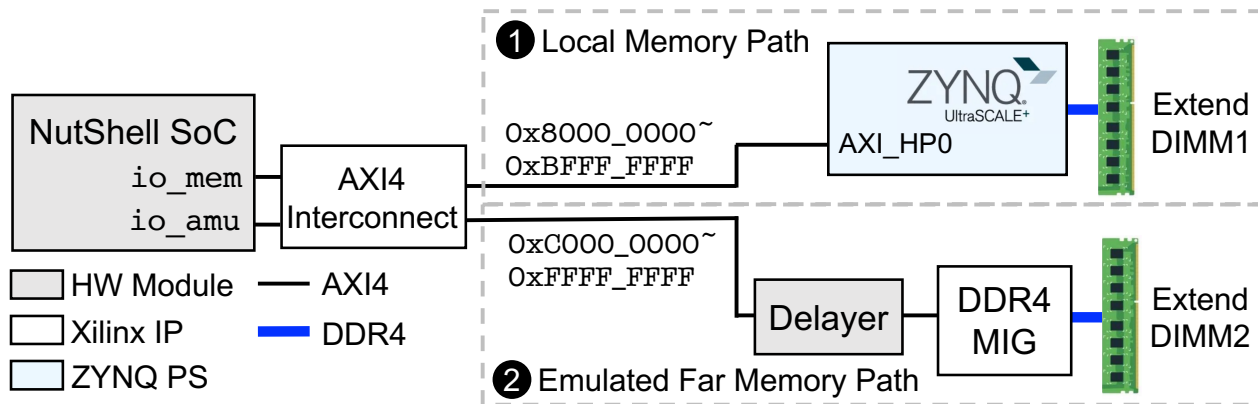
# Outline

Background

**A**synchronous **M**emory access ISA **E**xtensions (AME)

Programing Model

Evaluation

# FPGA-based Prototype System

- On Xilinx UltraScale+ ZU19EG MPSoC    ⟸ **by software on ARM core**

  **❶ Local Memory Path**

  **❷ Emulated far memory path** with **adjustable access latency**

  - **Real-time MLP observation**

- **NutShell** integrated with AMU

  - **In-order pipelined RISC-V64IMACSU core**: *open-source chip by university (ICT, taped-out)*

  - Boot **Debian 11** on FPGA-based Prototype System  **Compatibility of AMU**

# Evaluation

- **Goal:** Answer **3 questions** through **7 benchmarks**
  1. The improvement of memory access performance via AMU
  2. How does AMU accelerate memory access
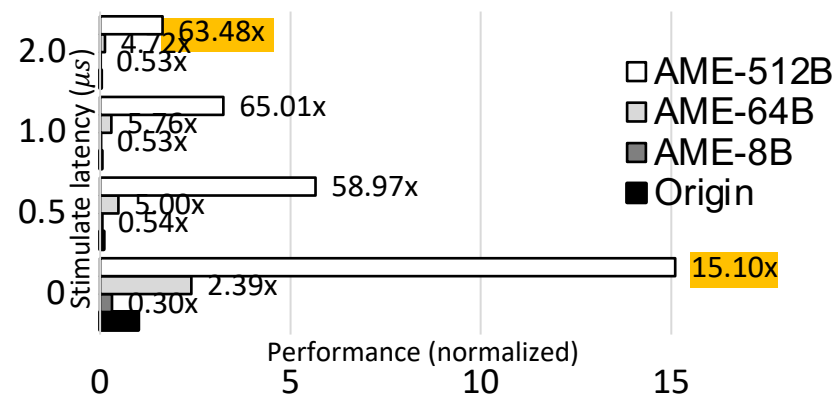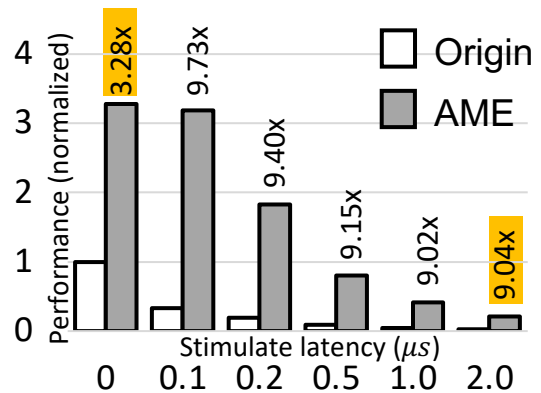  3. The acceleration effect of amu on the actual application

- **Setup**
  - Simulate a 2GHz processor
  - $0\sim10\mu s$ latency on
  - Emulated Far Memory

| Benchmark | Abbr. | Type | Footprint | AMQLEN |
|---|---|---|---|---|
| Random Access | GUPS | Memory access | 256MB | 256 |
| Sequence Access | SA | Memory access | 256MB | 256 |
| Hash Join | HJ | Data processing | 15.624MB | 256 |
| Integer Sort | IS | Data processing | 8MB | 256 |
| Binary Search | BS | Data processing | 97.65MB | 256 |
| Hash Table with Hand-over-hand Linked List | HL | Concurrent data structures | 2MB | 256 |
| Graph500 | G500 | Graph Computing | 16MB | 8 |

# Key Results 1: Memory Access Performance

|  | Local Memory | Far Memory (2$\mu s$ additional latency) |
|---|---|---|
| **Random Access** (GUPS, **8Byte**) | **3.28x** | **9.04x** |
| **Sequential Access** (SA, **512Byte**) | **15.10x** | **63.48x** |

# Key Results 1: Memory Access Performance

| | Local Memory | Far Memory (2$\mu s$ additional latency) |
|---|---|---|
| **Random Access** (GUPS, **8Byte**) | **3.28x** | **9.04x** |
| **Sequential Access** (SA, **512Byte**) | **15.10x** | **63.48x** |

- **High MLP** significantly improves the performance of **memory access.**
- **Variable-grained** memory access better fits program **semantics.**

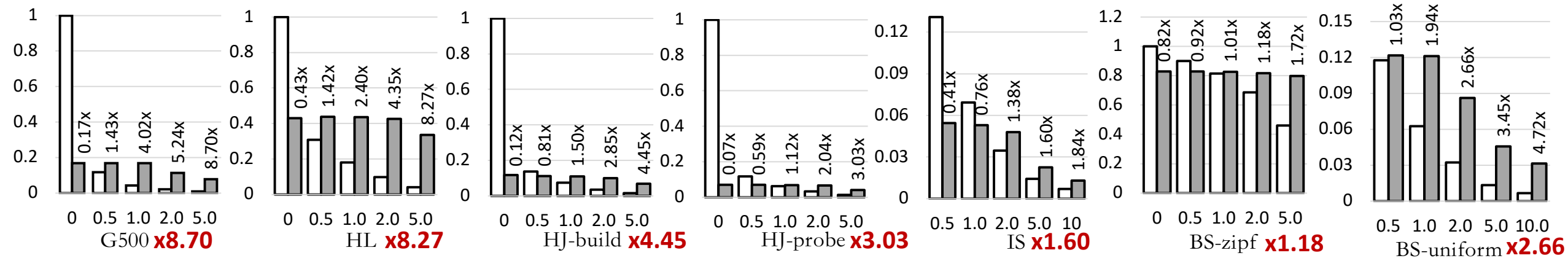## Fuller and more precise use of bandwidth

# Key Results 2: AME brings High MLP

- AMU reaches **238 outstanding requests** at average
  - Workload: GUPS, at $2\mu s$ stimulated latency
  - **Larger the concurrency we set, the more performance we get**



Issuing *aload* in batch

An *aload* of a part finish: issue *astore* **immediately**

To be finished

238    Avg

Number of in-flight Requests (concurrency)

256 / 192 / 128 / 64 / 0

Observation Time

Performance (up/second)

49500.0 / 49000.0 / 48500.0 / 48000.0 / 47500.0

49116.3
48942.3
48579.5
47696.8

4   16   64   256

Asynchronous memory access concurrency (AMQLEN)

# Key Results 3: Practical Benchmarks

- For applications:
  - **Reduced latency sensitivity** through high MLP (at glance)
  - **Multiple performance improvements**

- Practical significance:
  - AMU can **improve throughput** and protect concurrency
  - With AMU equipped, modern CPU can improve MLP to accelerate the calculation of popular applications such as graph computing



G500 **x8.70**  HL **x8.27**  HJ-build **x4.45**  HJ-probe **x3.03**  IS **x1.60**  BS-zipf **x1.18**  BS-uniform **x2.66**

# Q&A

# Architecture and RISC-V ISA Extension Supporting Asynchronous and Flexible Parallel Far Memory Access

**Songyue Wang**, Luming Wang, Tianyue Lu, Mingyu Chen

ICT, CAS

University of Chinese Academy of Sciences

*Please contact Songyue Wang (wangsongyue18@mails.ucas.ac.cn) for any concerns.*

University of Chinese Academy of Sciences

INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

先进计算机系统研究中心
Center for Advanced Computer Systems

# Conclusion

- **Problem**: Far memory scenarios brings potential **higher bandwidth and higher access latency**.

- **Goal**: Break the limitations of existing processors to improve MLP and mask latency with **high concurrent memory access**.

- **Challenge**:
  - Modern CPUs have little space for handling outstanding memory requests
  - Load/store instructions occupy critical resources in pipeline for a long time

- **Idea**:
  - Asynchronous Memory access ISA Extension (AME)
  - Asynchronous Memory access Unit (AMU) inside processors

- **Key Result**:
  - Average **11x** at $2\mu s$ latency for 2GHz CPU
  - MLP reaches 238 running GUPS

# Process of an asynchronous memory access