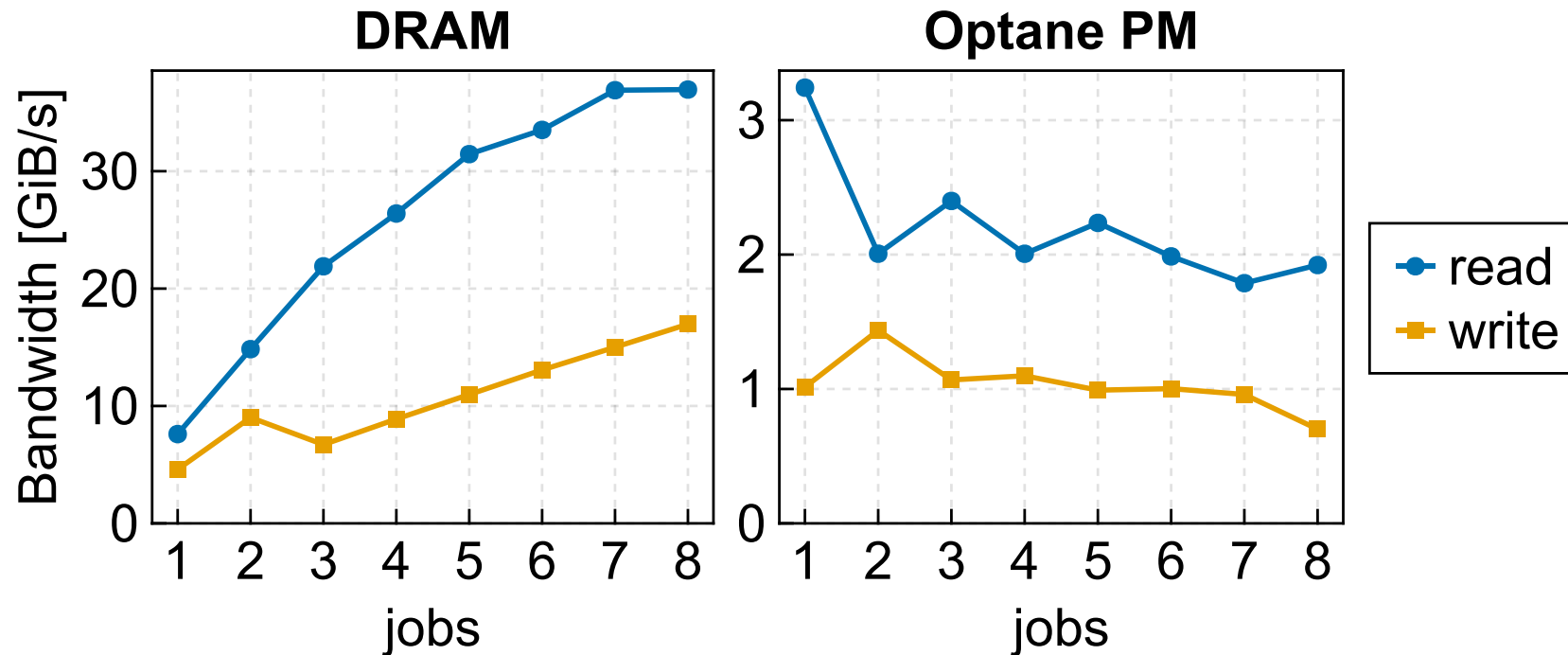# Per-Process Memory Bandwidth Management for Heterogeneous Memory Systems

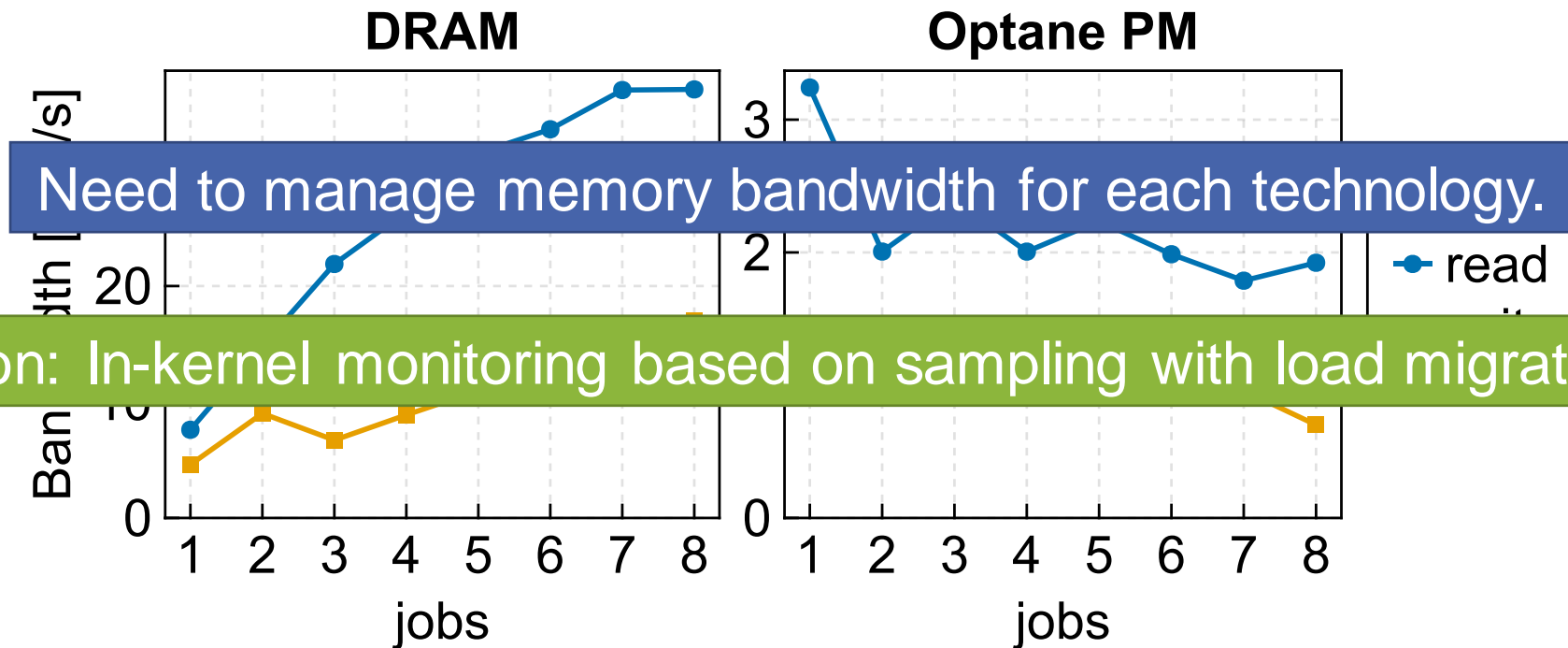**Lukas Werling, Daniel Habicht, Frank Bellosa**

# Motivation

- DAX with PM and CXL: Direct mappings not just to DRAM
  - Very different behavior
  - PM: parallel writes hurt overall performance

# Motivation

- DAX with PM and CXL: Direct mappings not just to DRAM
  - Very different behavior
  - PM: parallel writes hurt overall performance



**Need to manage memory bandwidth for each technology.**

**Our solution: In-kernel monitoring based on sampling with load migration policies.**

# Memory Bandwidth Monitoring Goals

process association    arbitrary devices

low latency    low overhead    high accuracy

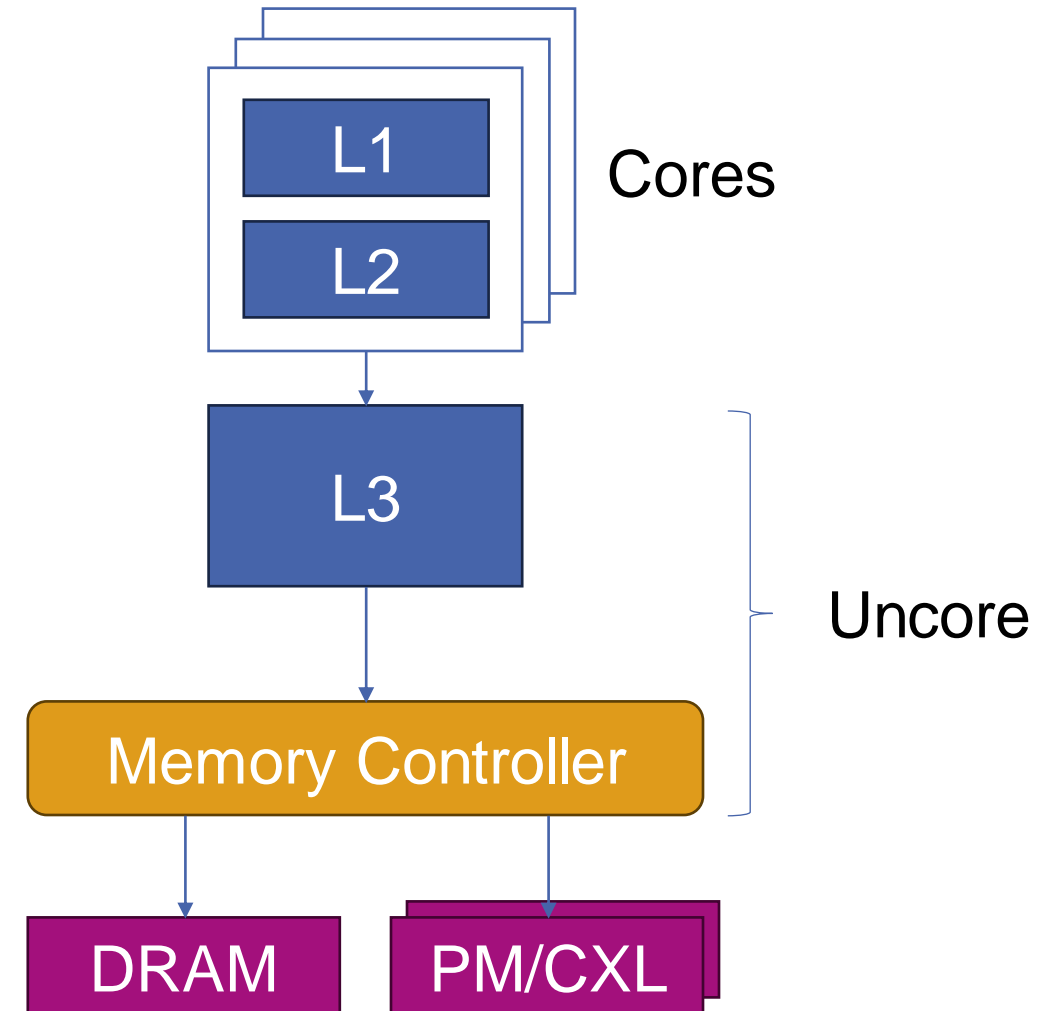# Challenge: Insufficient HW Control
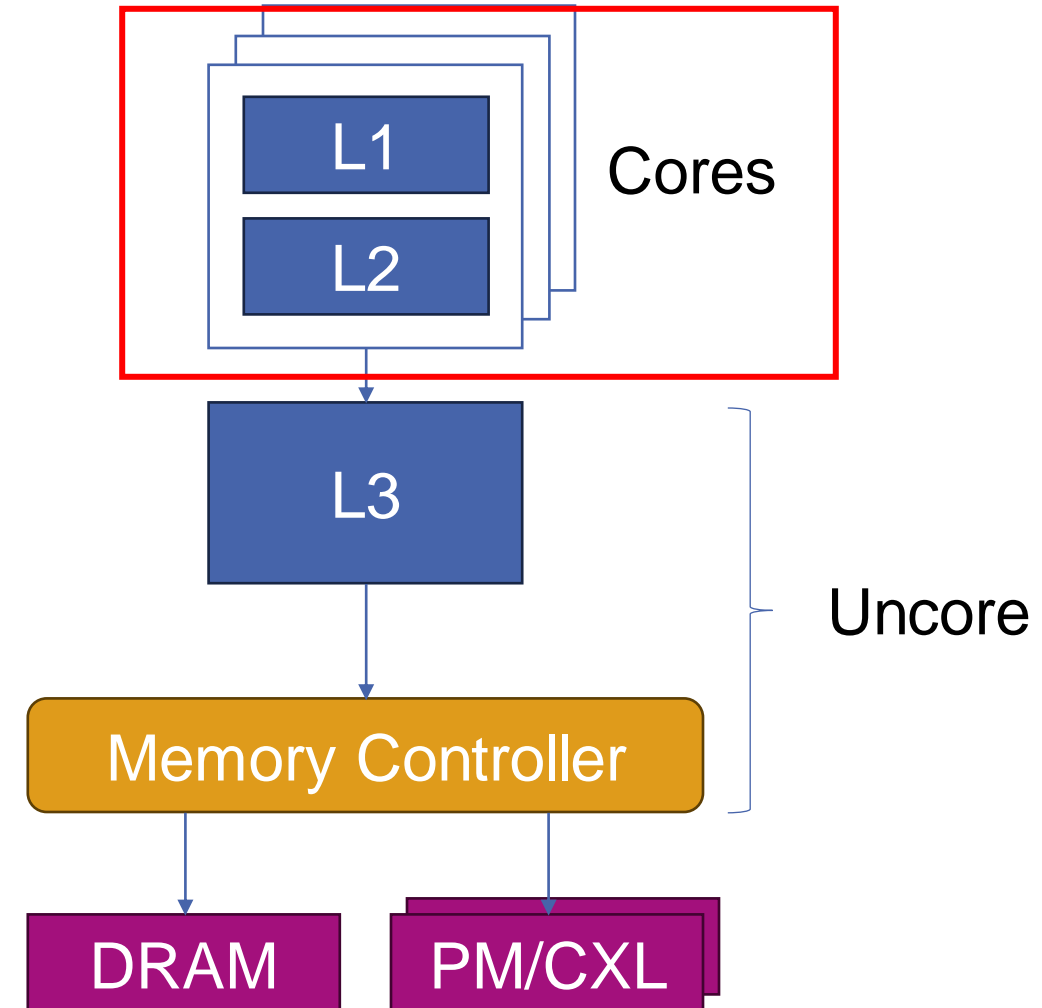
process association    arbitrary devices

low latency    low overhead    high accuracy

- ■ Performance counters
  - ■ Missing store events (core)
  - ■ Limited process association (uncore)
  - ■ Fixed memory classes (DRAM, PM)

- ■ Intel Memory Bandwidth Monitoring
  - ■ Only for overall bandwidth

L1
L2
Cores

L3

Uncore

Memory Controller

DRAM    PM/CXL

# Challenge: Insufficient HW Control

process association ┃ arbitrary devices

low latency ┃ low overhead ┃ high accuracy

- ■ **Performance counters**
  - ■ Missing store events (core)
  - ■ Limited process association (uncore)
  - ■ Fixed memory classes (DRAM, PM)

- ■ **Intel Memory Bandwidth Monitoring**
  - ■ Only for overall bandwidth

# Challenge: Insufficient HW Control
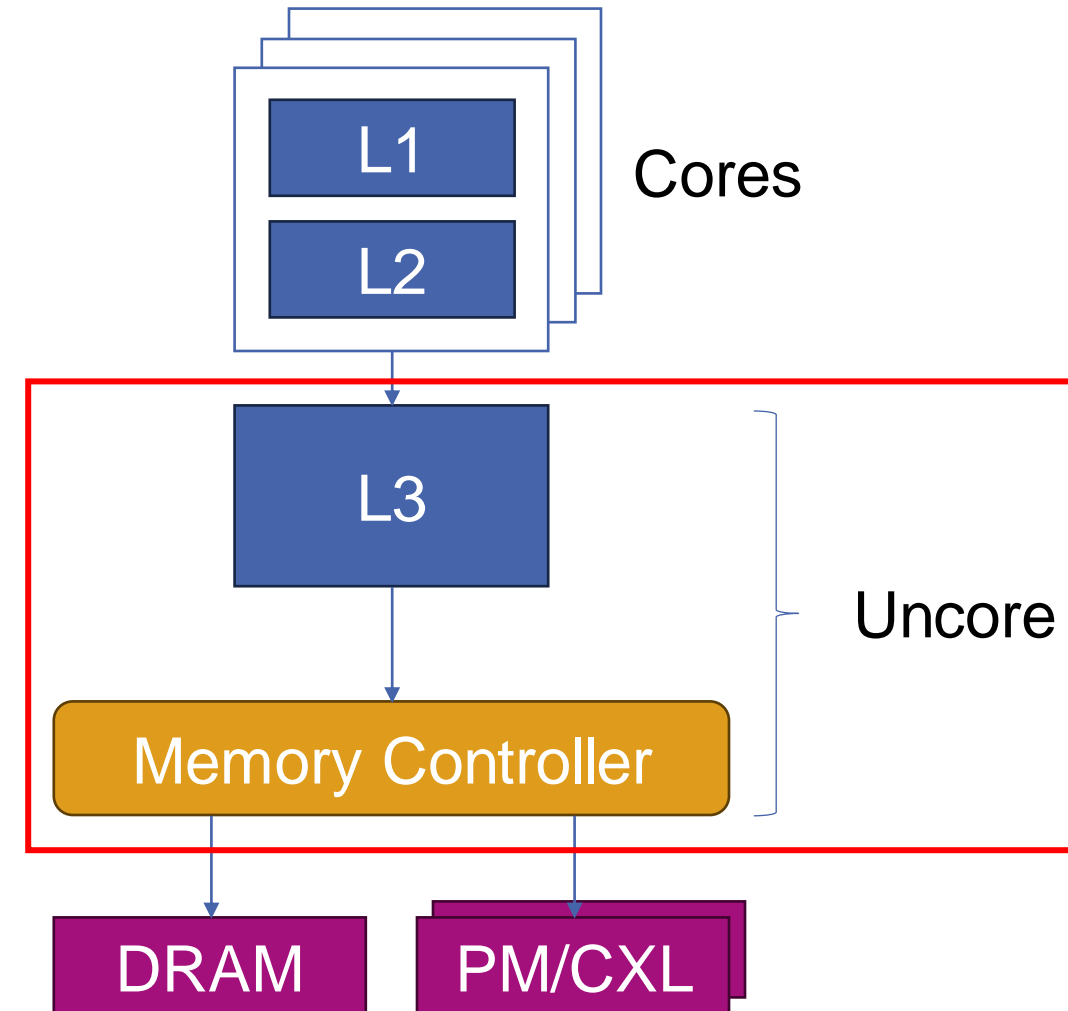
process association   arbitrary devices
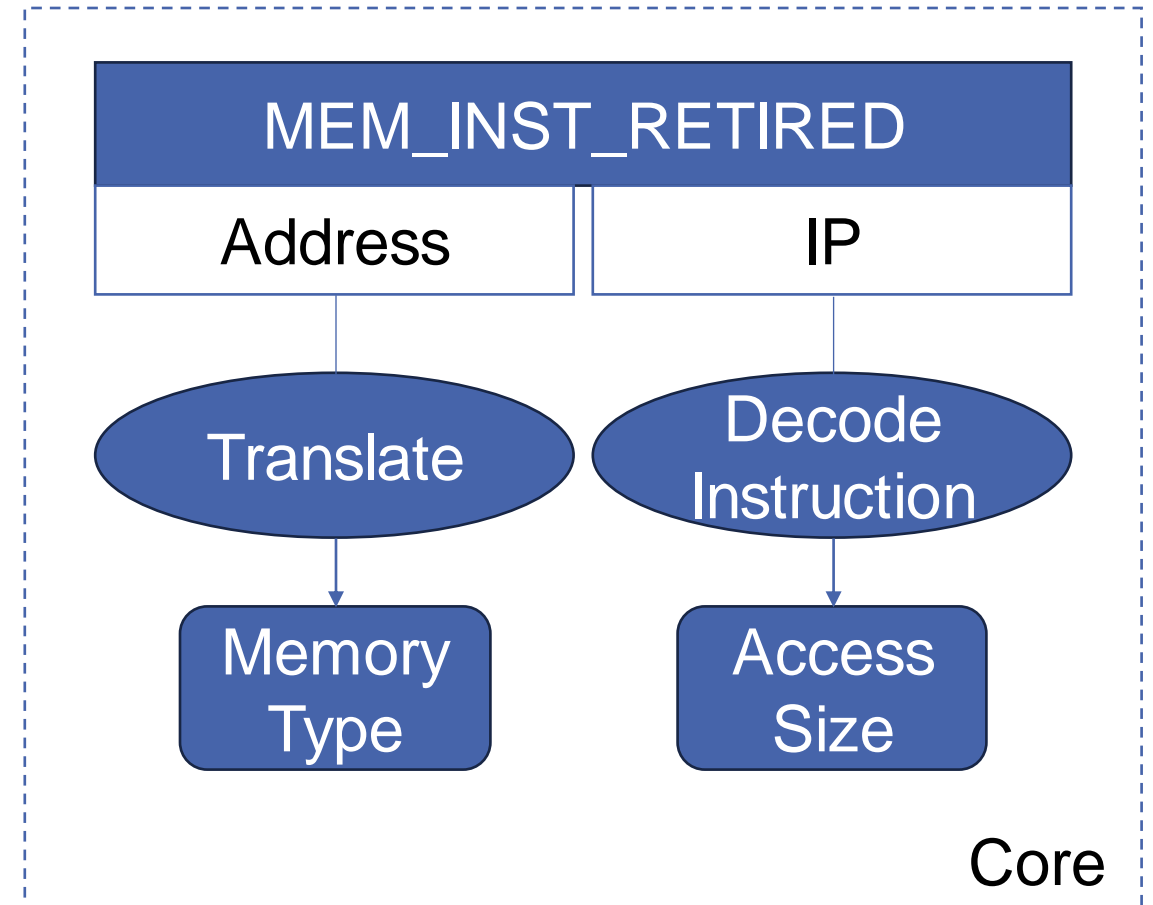
low latency   low overhead   high accuracy

- Performance counters
  - Missing store events (core)
  - Limited process association (uncore)
  - Fixed memory classes (DRAM, PM)

- Intel Memory Bandwidth Monitoring
  - Only for overall bandwidth



Cores

L1
L2

L3

Memory Controller

Uncore

DRAM   PM/CXL

# Approach: Sampling Stores

- Stores: Intel PEBS
  - Per-core counters
  - Capture additional information every N$^{th}$ instruction
  - Process sampling buffer regularly (1-6 samples)

- Loads: normal L3 miss counters
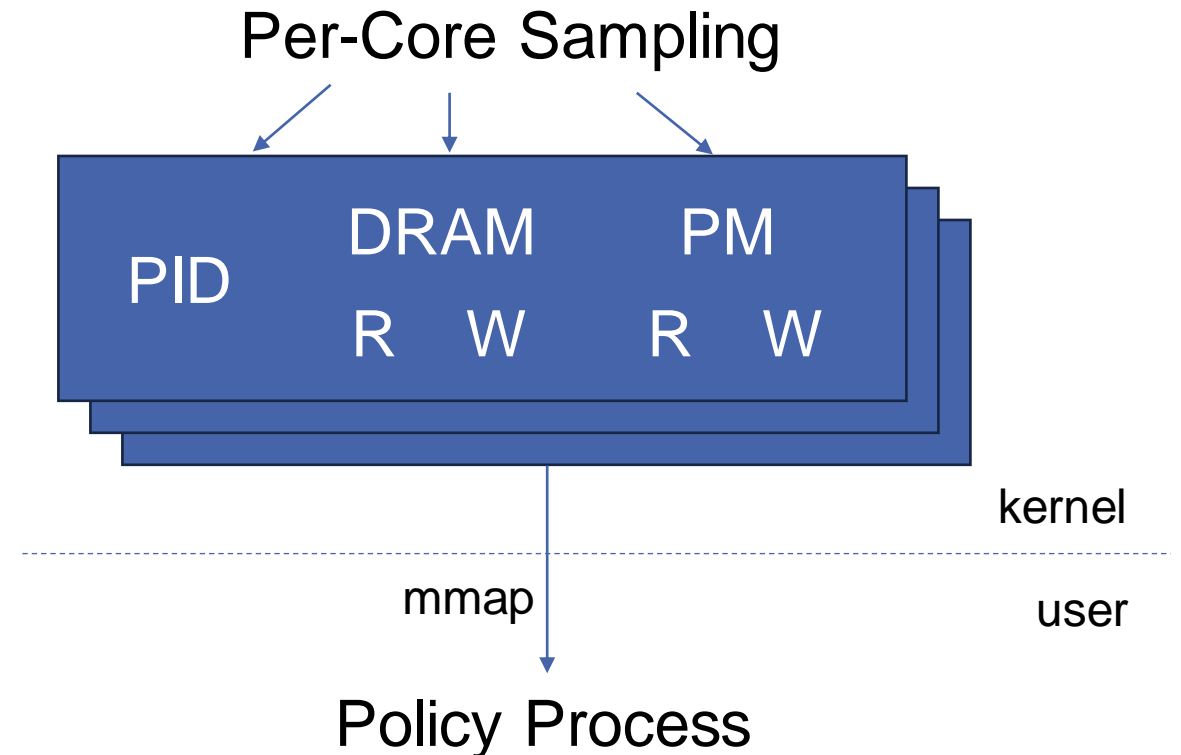  - Distinguish DRAM/PM

```
MEM_INST_RETIRED
   Address              IP
        |                |
    Translate      Decode Instruction
        |                |
    Memory Type     Access Size

                              Core
```

```
writes[type] += access_size * sampling_interval
```
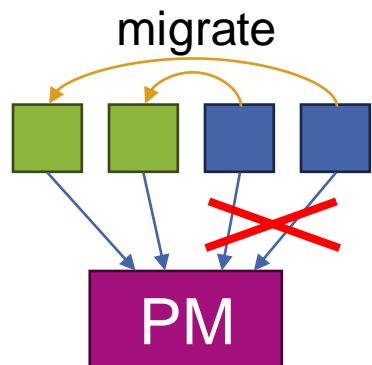
# Approach: Interface for Policies

- Challenge: Access current data at low latency

- Solution: Shared buffers
  - Data for running process
  - Lock-free access
  - Enables user-space policies
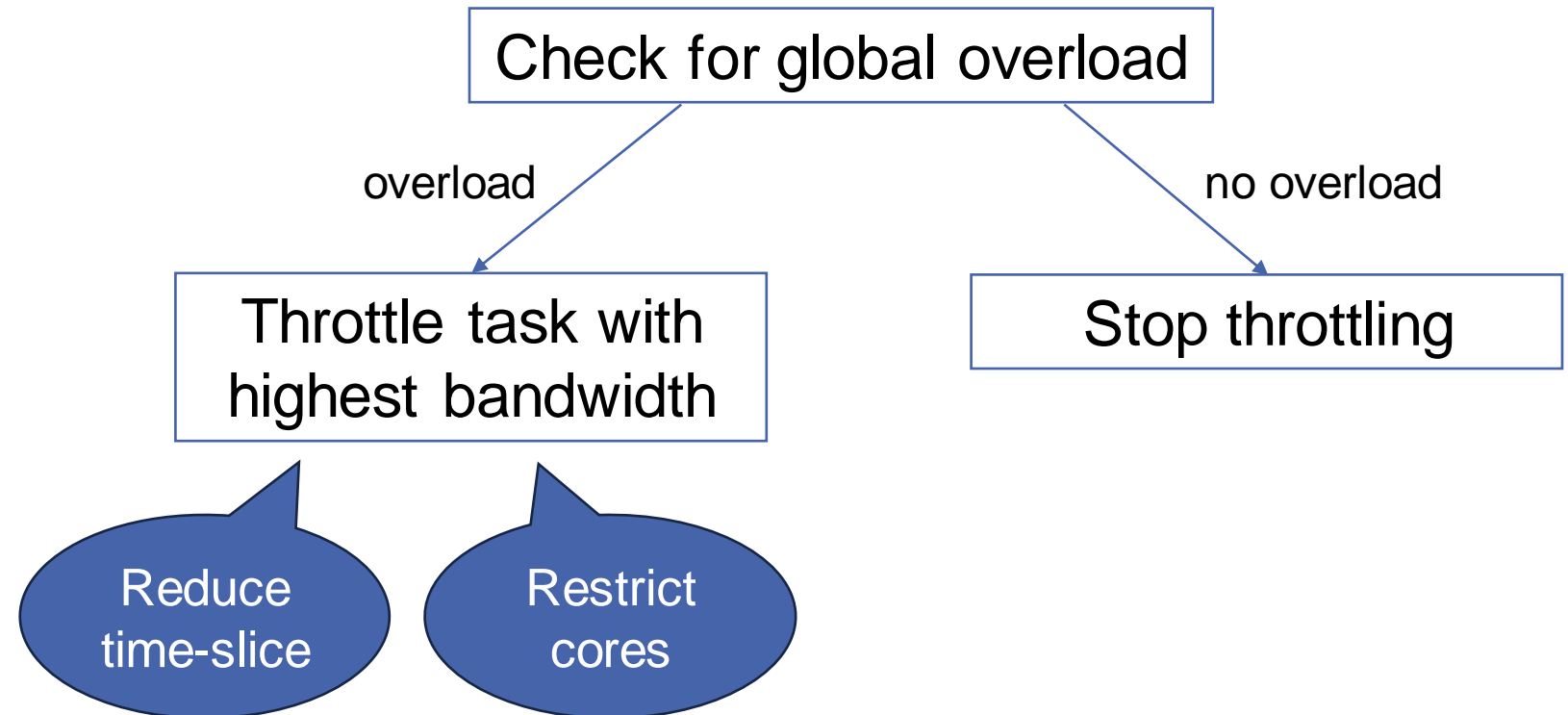
- Additional access via /proc



Per-Core Sampling

PID      DRAM        PM
         R    W      R    W

kernel

mmap

user

Policy Process

# Work in Progress: Policies



**Avoid PM overload**

Core Specialization on overload

migrate

PM

**Fair bandwidth allocation**

Check for global overload

overload

Throttle task with highest bandwidth

Reduce time-slice

Restrict cores

**High priority / best effort**

no overload

Stop throttling

2023-09-29 Lukas Werling – Per-Process Memory Bandwidth Management KIT Operating Systems Group

# Evaluation: Accuracy



2023-09-29    Lukas Werling – Per-Process Memory Bandwidth Management    KIT Operating Systems Group

# Evaluation: Accuracy

# Evaluation: Overhead



Average Processing Time per Sampled Store

2023-09-29    Lukas Werling – Per-Process Memory Bandwidth Management    KIT Operating Systems Group
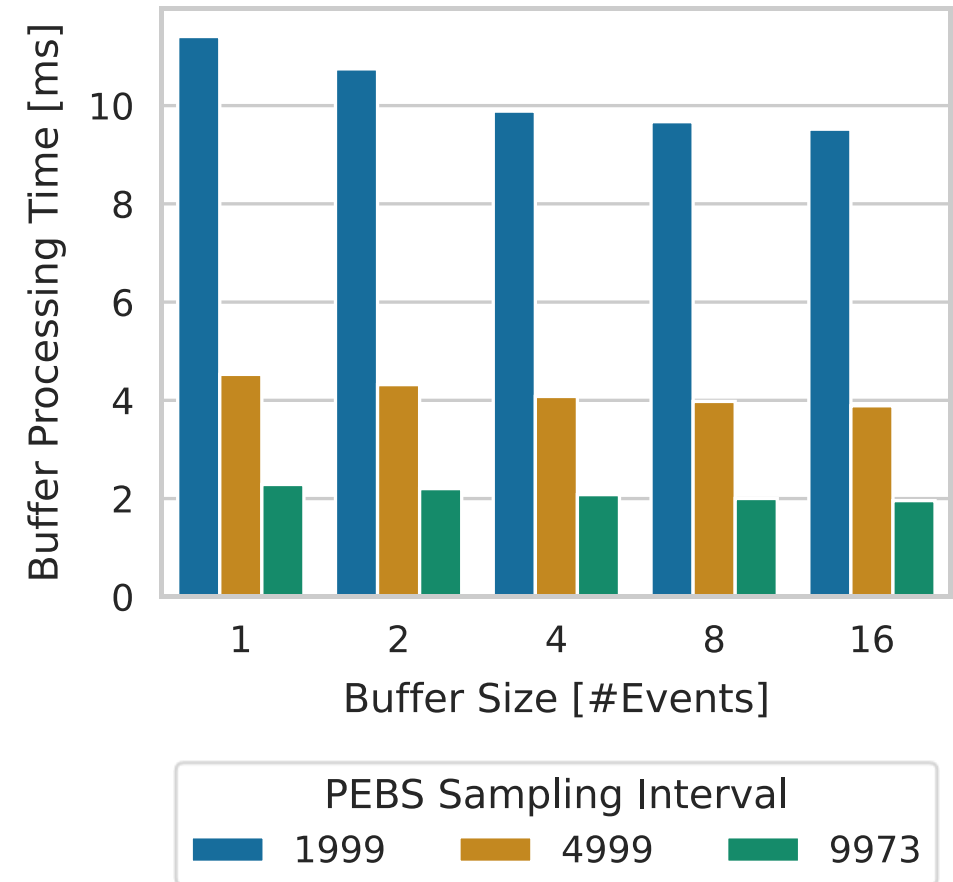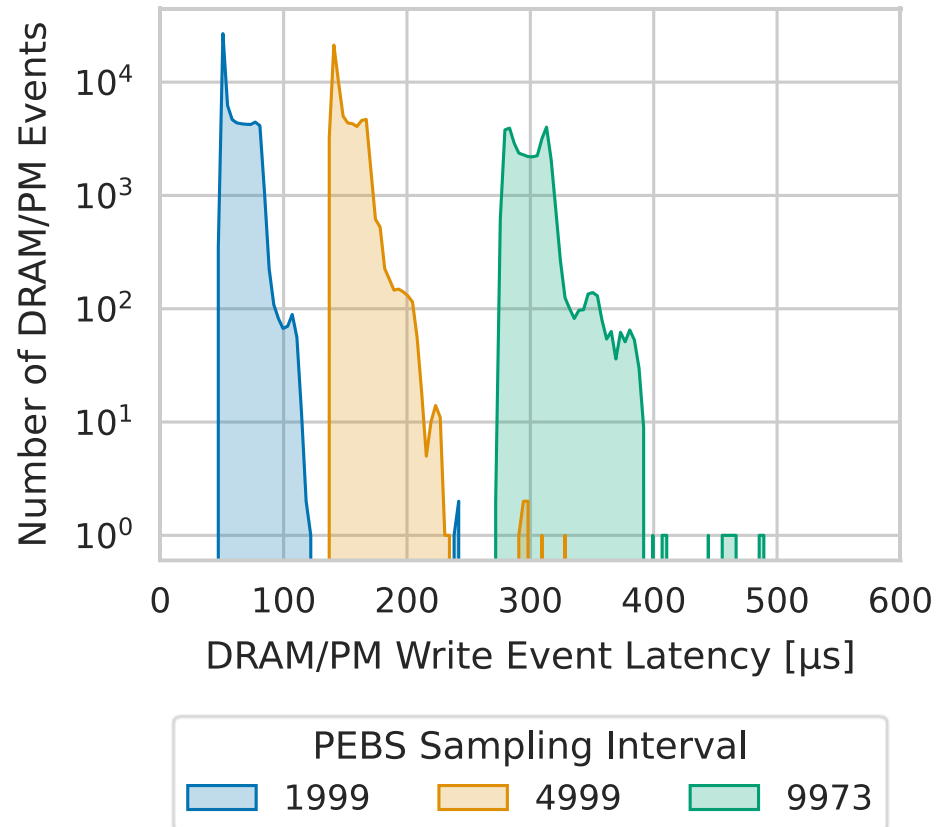
# Evaluation: Overhead



Average Processing Time per Sampled Store

Average Processing Time for Writing 1 GiB
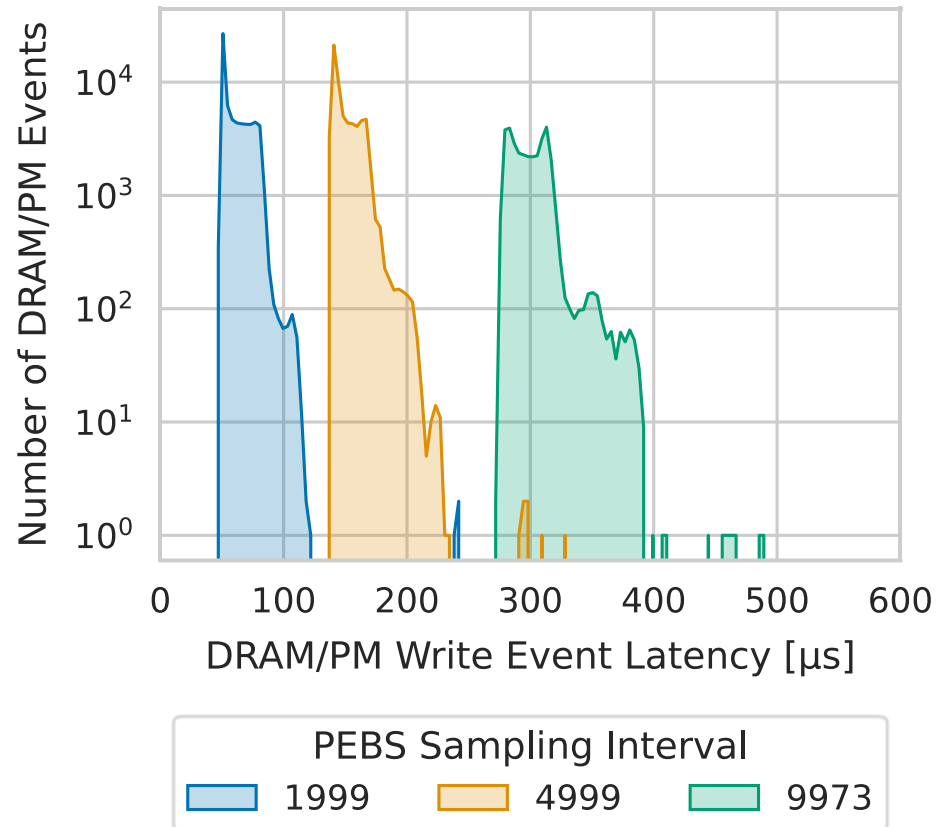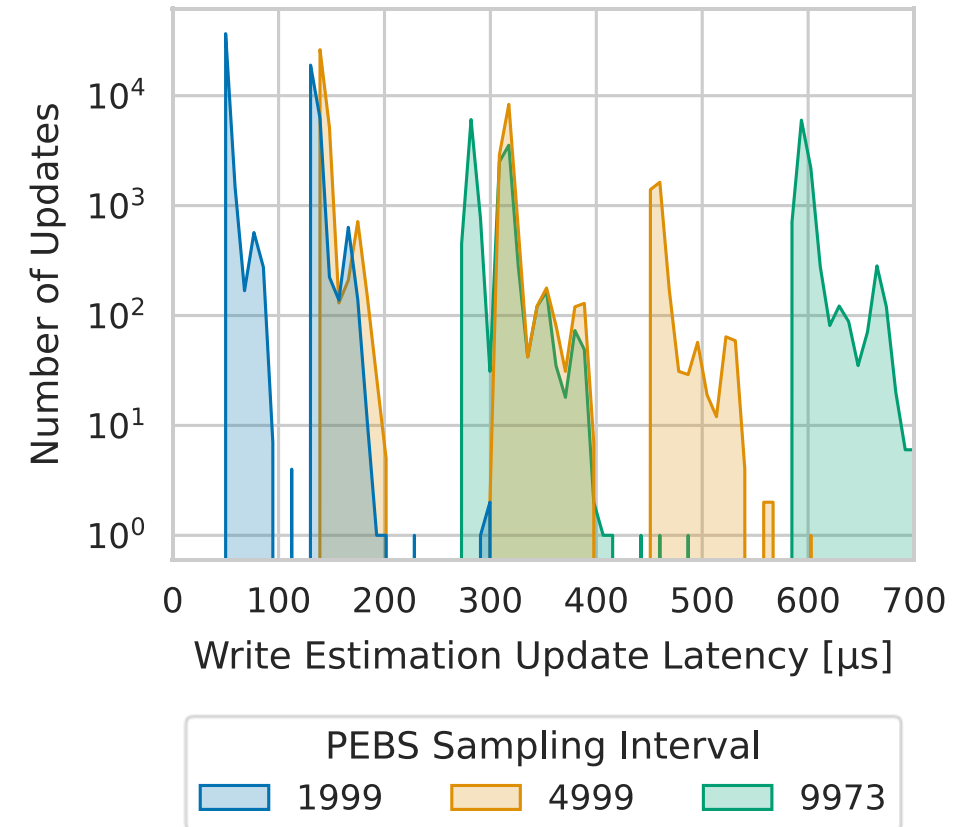
Distribution of DRAM/PM Write Event Latency
(PEBS Buffer Size = 1)

# Evaluation: Latency



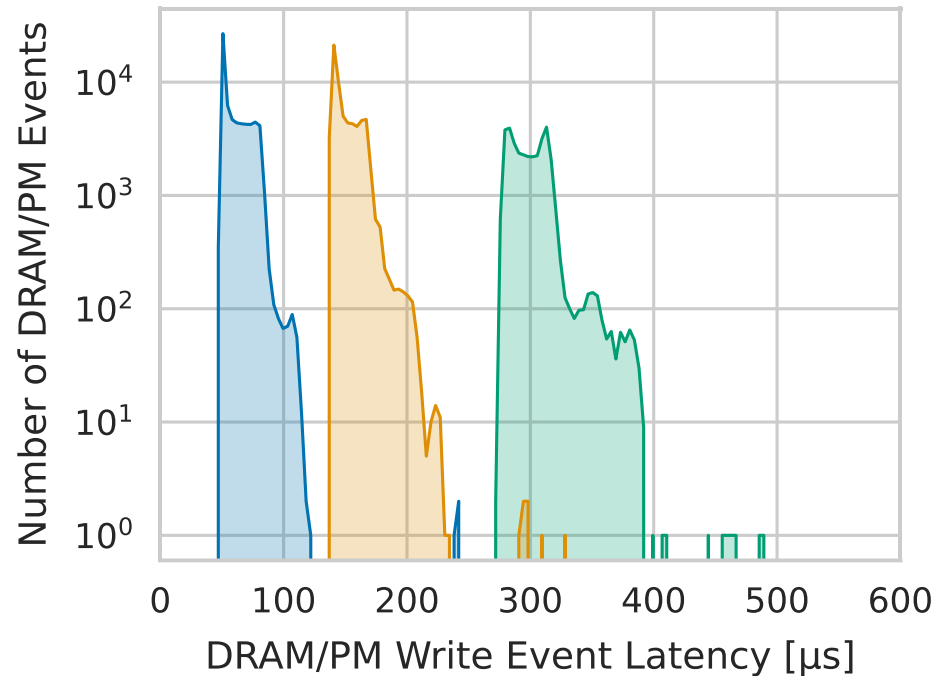Distribution of DRAM/PM Write Event Latency (PEBS Buffer Size = 1)

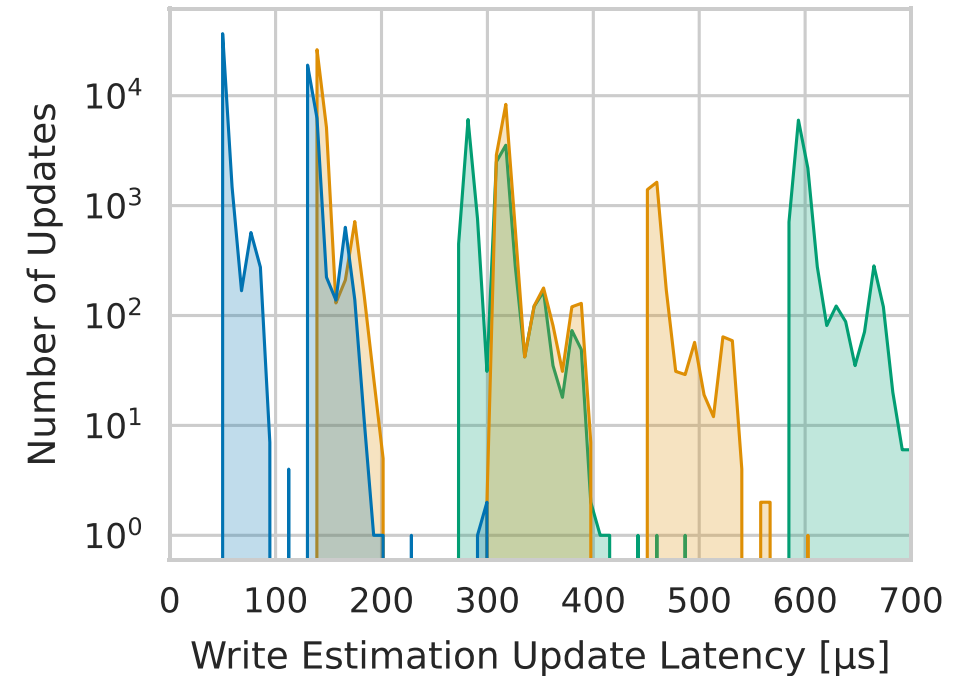Distribution of PM Write Estimation Update Latency (PEBS Buffer Size = 1)

filter by type

# Evaluation: Latency



Distribution of DRAM/PM Write Event Latency (PEBS Buffer Size = 1)

Distribution of PM Write Estimation Update Latency (PEBS Buffer Size = 1)

filter by type

Policies can react to changes within 1 ms.

# Discussion

- Accuracy
  - Problem: cache hits
  - Problem: write amplification

- Overhead / latency trade-off
  - Sampling interval
  - Sampling buffer size

- Latency
  - Simultaneous heterogeneous memory use increases latency

Policies only need relative numbers

# Conclusion

- **Memory bandwidth is a limited resource**
  - Different technologies
  - Monitoring and management important
  - Existing tools insufficient

- **Our solution: Sampling with PEBS**
  - Process association
  - High accuracy, low overhead
  - Low-latency interface for throttling

- **Future work: Throttling policies**