# Exploiting Modern Hardware Features via Lightweight Profiling
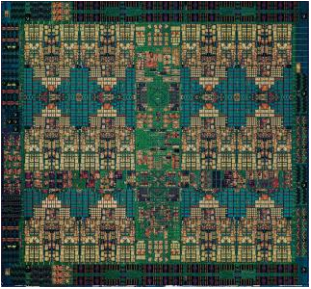
## Probir Roy

Scalable Tools Workshop'19
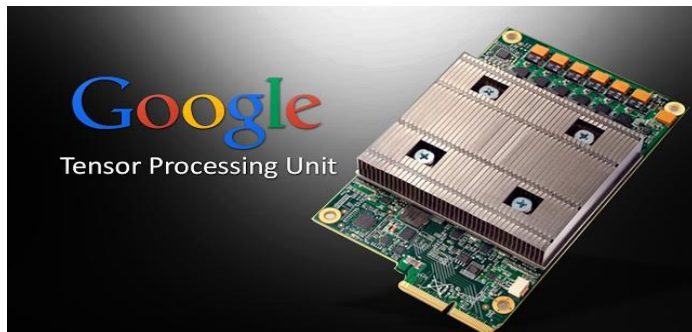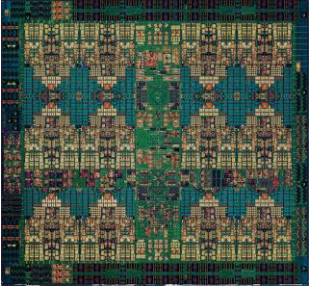
WILLIAM & MARY

CHARTERED 1693
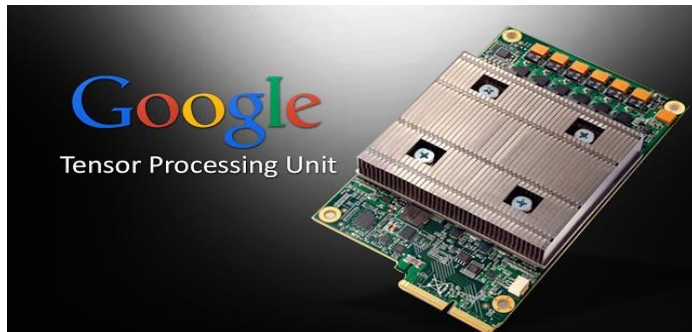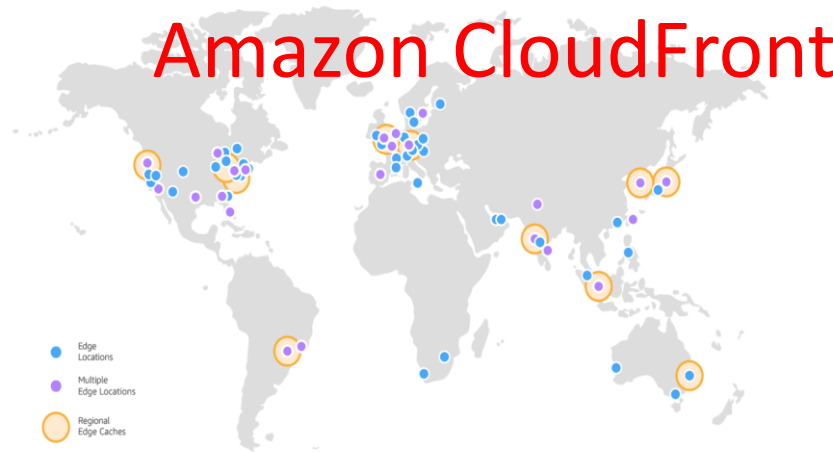
# High performance and challenges

IBM POWER 9 CPU

# High performance and challenges



IBM POWER 9 CPU



Amazon CloudFront

# High performance and challenges

IBM POWER 9 CPU

NAMD

Amazon CloudFront

Google Tensor Processing Unit

Edge Locations

Multiple Edge Locations

Regional Edge Caches

# High performance and challenges



IBM POWER 9 CPU

NAMD

Amazon CloudFront

OpenMP
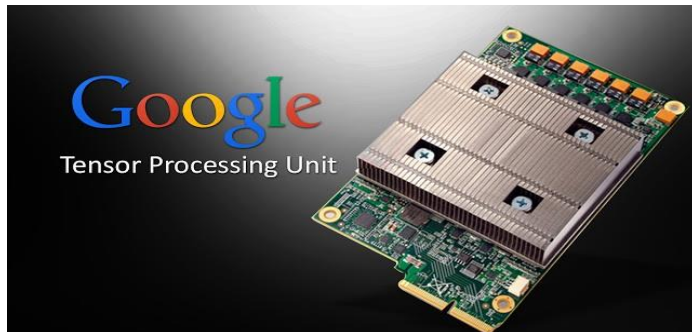
MPI

cassandra

redis

Google Tensor Processing Unit

Edge Locations
Multiple Edge Locations
Regional Edge Caches

julia

TensorFlow

# High performance and challenges

Common Goal: **Performance**

IBM POWER 9 cpu

OpenMP

MPI

Amazon CloudFront

cassandra

redis

Google
Tensor Processing Unit

julia

TensorFlow

# High performance and challenges

## Common Goal: **Performance**

## Variable characteristics of hardware and software is a **challenge**

IBM POWER9 GPU

Google
Tensor Processing Unit

julia

TensorFlow

- Edge Locations
- Multiple Edge Locations
- Regional Edge Caches

# High performance and challenges

**Common Goal: Performance**

**Variable characteristics** of hardware and software is a **challenge**

**Deep insights**

# High performance and challenges

http://sites.utexas.edu/jdm4372/2016/11/22/sc16-invited-talk-memory-bandwidth-and-system-balance-in-hpc-systems/

# High performance and challenges

Peak FLOPS per socket increasing at **50%-60% per year**

http://sites.utexas.edu/jdm4372/2016/11/22/sc16-invited-talk-memory-bandwidth-and-system-balance-in-hpc-systems/

# High performance and challenges

Peak FLOPS per socket increasing at **50%-60% per year**

Memory bandwidth increasing at **~23% per year**

http://sites.utexas.edu/jdm4372/2016/11/22/sc16-invited-talk-memory-bandwidth-and-system-balance-in-hpc-systems/

# High performance and challenges

Peak FLOPS per socket increasing at **50%-60% per year**

Memory bandwidth increasing at **~23% per year**

Memory latency increasing at **~4% per year**

http://sites.utexas.edu/jdm4372/2016/11/22/sc16-invited-talk-memory-bandwidth-and-system-balance-in-hpc-systems/

# Steps of performance analysis

Application

# Steps of performance analysis

# Steps of performance analysis

# Steps of performance analysis

| Simulation | Measurement |
|---|---|

What?

Why?

How?

Application

→

Profiler

# Steps of performance analysis

# Steps of performance analysis

# Steps of performance analysis

Simulation

Measurement

What?

Why?

How?

Application → Profiler → Profiles → Analyzer → Code optimization

# Steps of performance analysis

# Steps of performance analysis

# Limitations of performance analysis (Cont.)
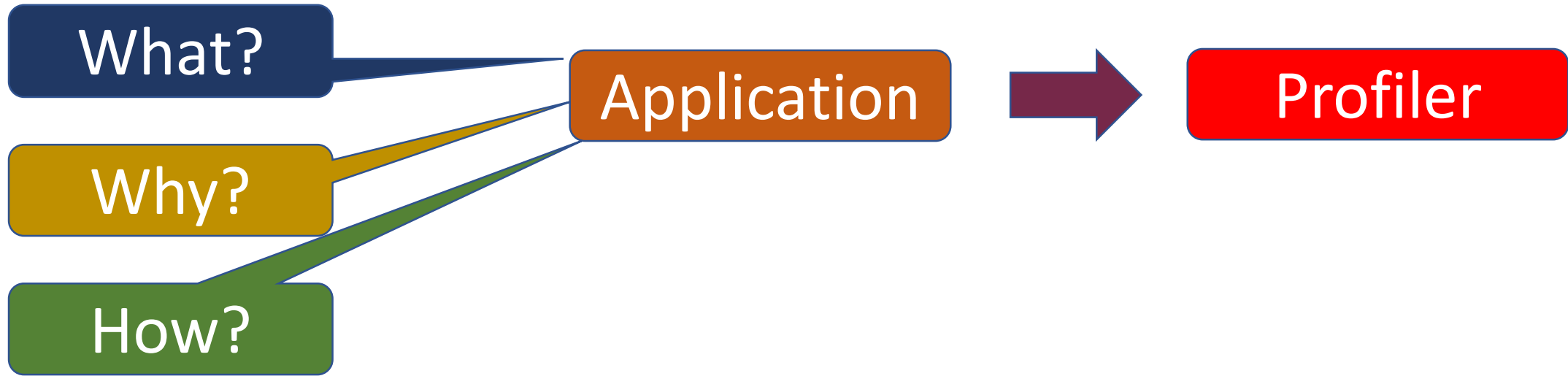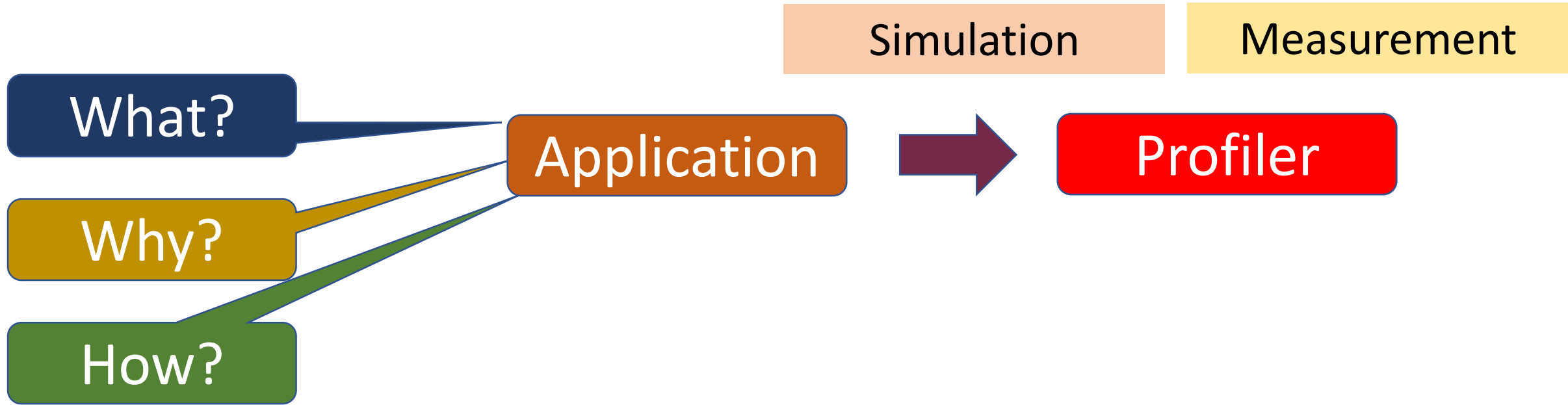
# Limitations of performance analysis (Cont.)



Deep insight

Shallow insight

**Simulation methods (PinTool, GPGPUSim, GEMS)**

**Measurement methods (Perf, Oprofile, PAPI)**

High overhead

Low overhead

*Cache simulation:* average 38x (*Xiang et al.* A higher order theory of locality )

# Limitations of performance analysis (Cont.)



**Cache simulation:** average 38x (*Xiang et al.* A higher order theory of locality )

**Selective instrumentation**: 2x - 5x (Rane et al. MACPO)

# Limitations of performance analysis (Cont.)



Deep insight

Shallow insight

Simulation methods (PinTool, GPGPUSim, GEMS)

Measurement methods (Perf, Oprofile, PAPI)

High overhead

Low overhead

**Cache simulation:** average 38x (*Xiang et al.* A higher order theory of locality )

**Profiling:** < 10% (Liu et al. A Data-centric Profiler for Parallel Programs)

**Selective instrumentation**: 2x - 5x (Rane et al. MACPO)

Exploiting Modern Hardware Features via Lightweight Profiling

# Limitations of performance analysis (cont.)



Simulation methods (PinTool, GPGPUSim, GEMS)

Measurement methods (Perf, Oprofile, PAPI)

**Goal**

Deep insight

Shallow insight

High overhead

Low overhead

*Cache simulation:* average 38x (*Xiang et al.* A higher order theory of locality )

*Profiling:* < 10% (Liu et al. A Data-centric Profiler for Parallel Programs)

*Selective instrumentation*: 2x - 5x (Rane et al. MACPO)

Exploiting Modern Hardware Features via Lightweight Profiling

5

# Research statement



Deep insight

Shallow insight

Simulation methods (PinTool, GPGPUSim, GEMS)

Measurement methods (Perf, Oprofile, PAPI)

High overhead

Low overhead

# Research statement

Lightweight profiling with PMUs can provide deep insights into performance issues caused by memory hierarchies and poor algorithm choice

Deep insight

Simulation methods (PinTool, GPGPUSim, GEMS)

Shallow insight

Measurement methods (Perf, Oprofile, PAPI)

High overhead

Low overhead
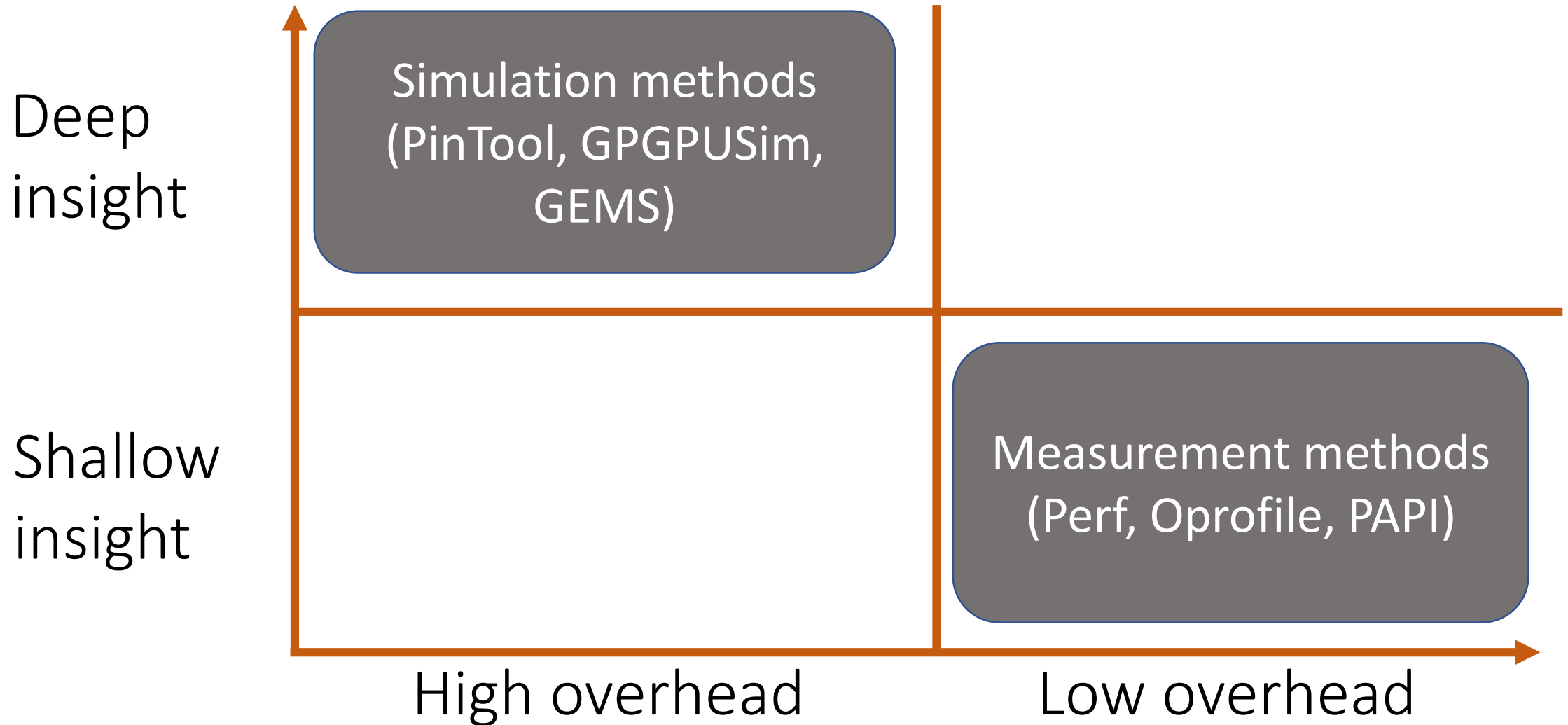
Exploiting Modern Hardware Features via Lightweight Profiling

6

# Research statement

**Lightweight profiling with PMUs can provide deep insights into performance issues caused by memory hierarchies and poor algorithm choice**

Deep insight

Simulation methods (PinTool, GPGPUSim, GEMS)

Tools to detect memory and computational inefficiency

Shallow insight

Measurement methods (Perf, Oprofile, PAPI)

High overhead

Low overhead

# My research at a glance

# My research at a glance

# My research at a glance

# My research at a glance

# My research at a glance



Memory Inefficiency

- Simultaneous multi-threading
- Set-associative cache
- Cache line
- Non-uniform memory

CPU

Cache

Physical Memory

- Memory contention
- Conflict miss
- utilization
- Scalability

Programming model

# My research at a glance

# My research at a glance

# Outline

- **Lightweight profiling**
- SMT-aware optimization
- Detection of cache conflicts
- Guiding data-structure layout transformation

# Lightweight memory profiling

- Hardware profiling

- Event based sampling
  - Intel (Precise event based sampling - PEBS)
  - AMD (Instruction based sampling - IBS)
  - IBM (Marked event sampling - MRK)

# Lightweight memory profiling

- Hardware profiling
- Event based sampling
  - Intel (Precise event based sampling - PEBS)
  - AMD (Instruction based sampling - IBS)
  - IBM (Marked event sampling - MRK)

**PMU**

**Application**

**Time**

# Lightweight memory profiling

- Hardware profiling

- Event based sampling
    - Intel (Precise event based sampling - PEBS)
    - AMD (Instruction based sampling - IBS)
    - IBM (Marked event sampling - MRK)



**PMU**

**Application**

**Time**

# Lightweight memory profiling

- Hardware profiling
- Event based sampling
  - Intel (Precise event based sampling - PEBS)
  - AMD (Instruction based sampling - IBS)
  - IBM (Marked event sampling - MRK)



**PMU**

**Application**

**Time**

Sample    Sample    Sample    Sample

# Lightweight memory profiling

- Hardware profiling

- Event based sampling
  - Intel (Precise event based sampling - PEBS)
  - AMD (Instruction based sampling - IBS)
  - IBM (Marked event sampling - MRK)

**PMU**

**Application**

**Time**

Sample    Sample    Sample    Sample

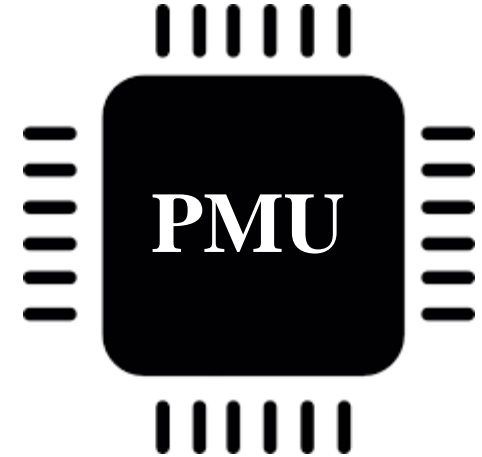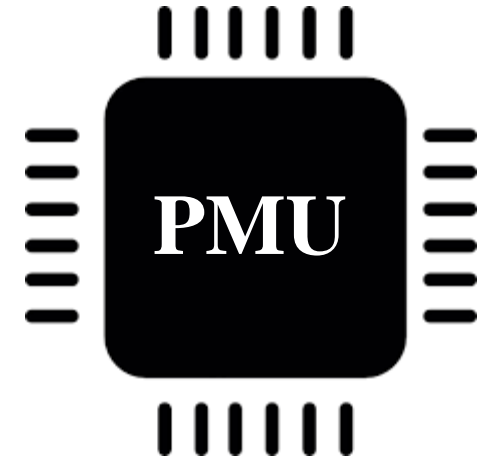Reference Type    {L1 miss, L2 hit etc.}

Data Address

Instruction Pointer

# Outline

✓ Lightweight profiling

✓ *SMT-aware optimization*

• Detection of cache conflicts

• Guiding data-structure layout transformation

# SMT-Aware Instantaneous Footprint Optimization

## [HPDC – 2016]

**Probir Roy**, Shuaiwen Leon Song, *Xu Liu*

# SMT (Simultaneous Multi-Threading)

# SMT scalability

# SMT architecture: shared cache

# SMT: Memory scalability

**Lower is better**



SMT scaling factor (F) = access Latency of SMT/ access Latency of non-SMT

# Characterization based on sensitivity

L = Memory Access Latency; F = scaling factor

| (L,F) | Benchmarks | Characterization |
|---|---|---|
| (high, high) | srad, streamcluster1, Lulesh2.0, IRSmk, LU, 3D tensor, Stencil, streamcluster2, hotspot, Clomp | potentially sensitive to mem-centric SMT optimizations |
| (high, low) | lud, needle, bfs, nn, bp, canneal, Ferret | not clear if they can further benefit from SMT optimizations |
| (low, high) | leucocite, heartwall, pathfinder, myocyte | little benefit from mem-centric SMT optimization |
| (low, low) | b+tree, cfd, kmeans, lavaMD, particle filter, hotspot3D, blackscholes, bodytrack, facesim, Swaptions | good memory performance with SMT enabled |

# Characterization based on sensitivity

L = Memory Access Latency; F = scaling factor

| (L,F) | Benchmarks | Characterization |
|---|---|---|
| (high, high) | srad, streamcluster1, Lulesh2.0, IRSmk, LU, 3D tensor, Stencil, streamcluster2, hotspot, Clomp | potentially sensitive to mem-centric SMT optimizations |
| (high, low) | lud, needle, bfs, nn, bp, canneal, Ferret | not clear if they can further benefit from SMT optimizations |
| (low, high) | leucocite, heartwall, pathfinder, myocyte | little benefit from mem-centric SMT optimization |
| (low, low) | b+tree, cfd, kmeans, lavaMD, particle filter, hotspot3D, blackscholes, bodytrack, facesim, Swaptions | good memory performance with SMT enabled |

# Source of memory contention

Little/no locality

# Source of memory contention

Little/no locality

Intra-thread



SMT thread 1     SMT thread 2

space

time

# Source of memory contention

**Little/no locality**

**Intra-thread**



● SMT thread 1     ● SMT thread 2

space

time

# Source of memory contention

Little/no locality

Intra-thread



SMT thread 1    SMT thread 2

space

time

Optimization: Improve cache line utilization

# Source of memory contention

Little/no locality

Intra-thread

Inter-thread



● SMT thread 1    ● SMT thread 2

space

time

Optimization: Improve cache line utilization

# Source of memory contention

**Little/no locality**

**Intra-thread**

**Inter-thread**

● SMT thread 1    ● SMT thread 2



space

time

space

time

**Optimization: Improve cache line utilization**

# Source of memory contention

**Little/no locality**

**Intra-thread**

**Inter-thread**

● **SMT thread 1**   ● **SMT thread 2**



space

time

space

Cache line 1

Cache line 1

time

**Optimization: Improve cache line utilization**

Modern Hardware Features via Lightweight Profiling

17

# Source of memory contention

Little/no locality

Intra-thread

Inter-thread

● SMT thread 1    ● SMT thread 2



space

time

Optimization: Improve cache line utilization

space

Cache line 1

Cache line 1

Optimization: Collaboration

Modern Hardware Features via Lightwe...

17

# SMT locality (Stencil code)

```
#pragma omp parallel for
for(int i=T; i<N-T; i++)
  for (int j=T; j<N-T; j++)
    for (int k=0; k<T; k++)
      R[i][j] = matrix[i][j]
              + matrix[i-k][j]+matrix[i][j-k]
              + matrix[i+k][j]+matrix[i][j+k];
```

# SMT locality (Stencil code)

```
#pragma omp parallel for
for(int i=T; i<N-T; i++)
    for (int j=T; j<N-T; j++)
        for (int k=0; k<T; k++)
            R[i][j] = matrix[i][j]
                        + matrix[i-k][j]+matrix[i][j-k]
                        + matrix[i+k][j]+matrix[i][j+k];
```

# SMT locality (Stencil code)

```
#pragma omp parallel for  schedule(static,1)
  for(int i=T; i<N-T; i++)
     for (int j=T; j<N-T; j++)
        for (int k=0; k<T; k++)
           R[i][j] = matrix[i][j]
                     + matrix[i-k][j]+matrix[i][j-k]
                     + matrix[i+k][j]+matrix[i][j+k];
```

# SMT locality (Stencil code)

```
#pragma omp parallel for  schedule(static,1)
for(int i=T; i<N-T; i++)
    for (int j=T; j<N-T; j++)
        for (int k=0; k<T; k++)
            R[i][j] = matrix[i][j]
                    + matrix[i-k][j]+matrix[i][j-k]
                    + matrix[i+k][j]+matrix[i][j+k];
```



18

# SMT locality (Stencil code)

```
#pragma omp parallel for  schedule(static,1)
for(int i=T; i<N-T; i++)
   for (int j=T; j<N-T; j++)
      for (int k=0; k<T; k++)
         R[i][j] = matrix[i][j]
                    + matrix[i-k][j]+matrix[i][j-k]
                    + matrix[i+k][j]+matrix[i][j+k];
```
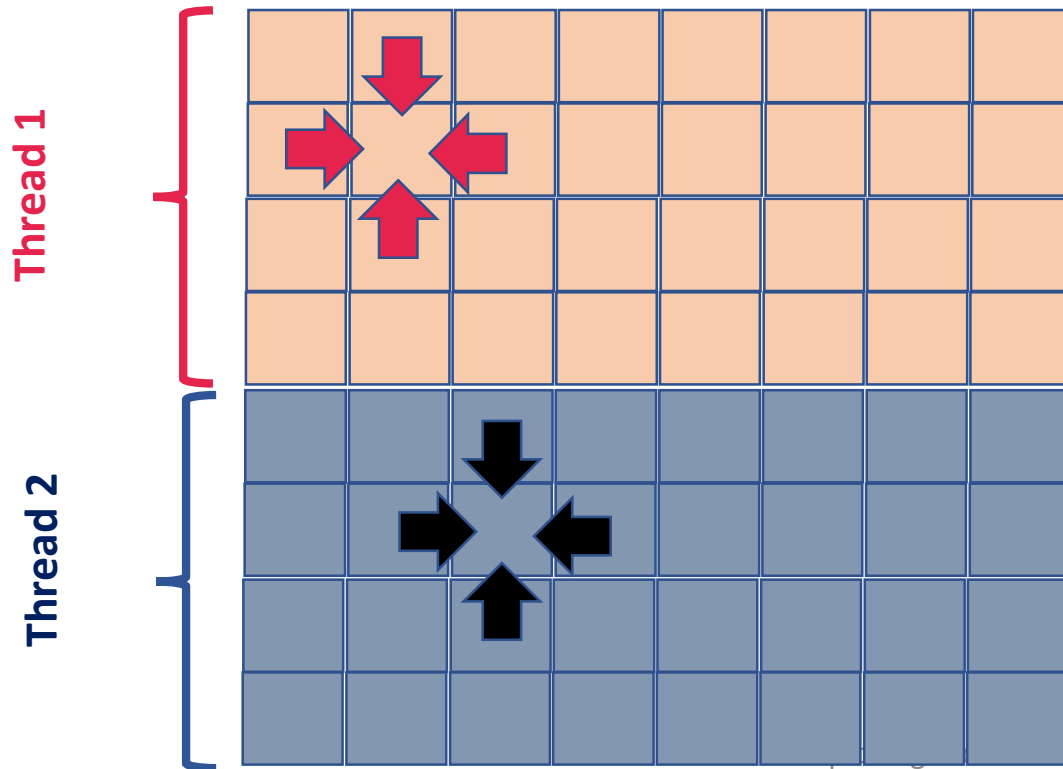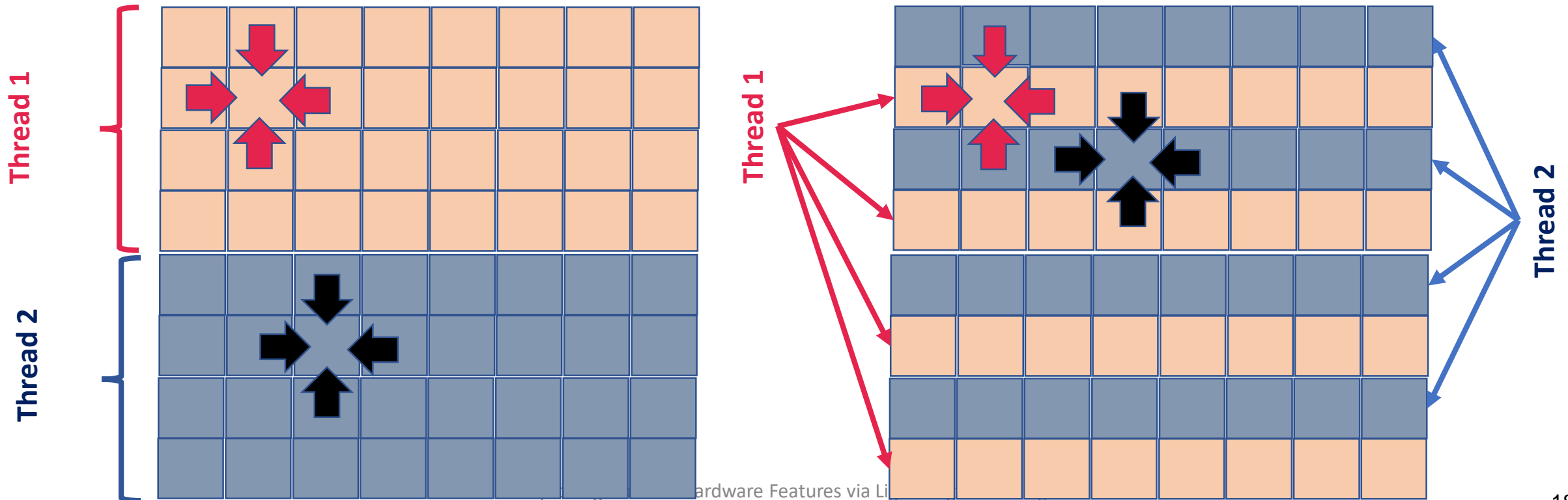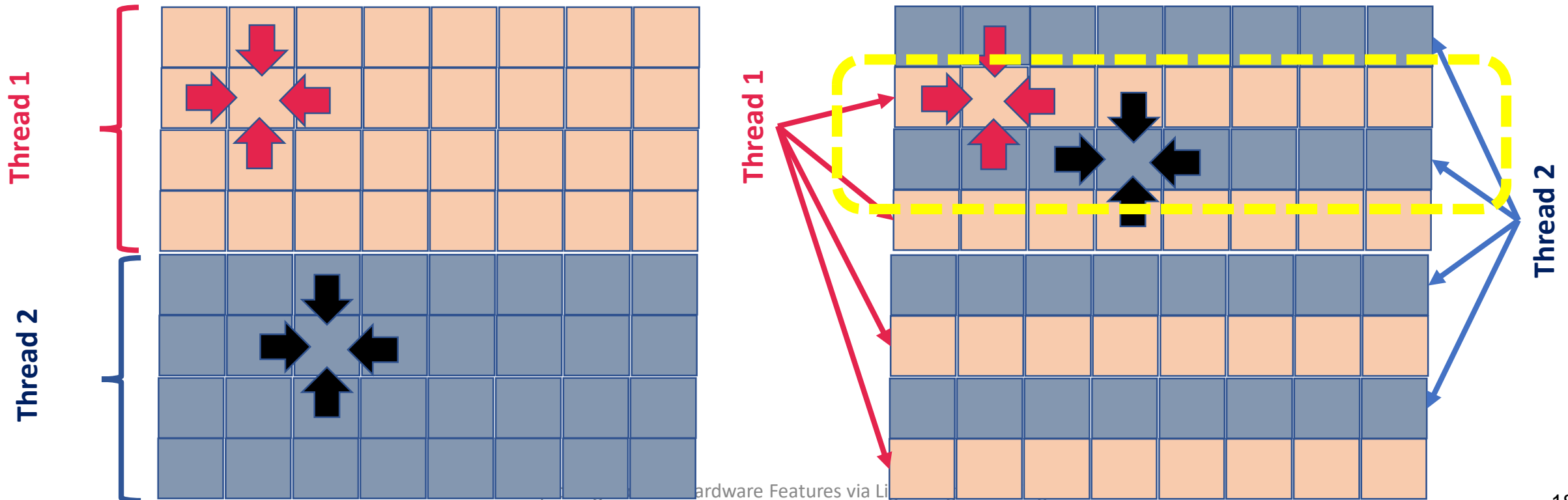
# SMT-Analyzer: Analyzing memory access pattern

# SMT-Analyzer: Analyzing memory access pattern

# SMT-Analyzer: Analyzing memory access pattern

# SMT-Analyzer: Analyzing memory access pattern

# Benchmarks

| Benchmarks | bottleneck region | % of total latency | overhead | OPT method | Speedups |
|---|---|---|---|---|---|
| lulesh2.0 | lulesh.cc: 604-609 | 3.6% | +3.1% | inter-thread | 1.43× |
| IRSmk | rmatmult3.c: 86-103 | 78.6% | +3.2% | intra-thread | 4.86× |
| needle | needle.cpp:185-187 | 20% | +2.99% | inter-thread | 2.37× |
| srad | srad.cpp:136-167 | 80.1% | +2.47% | intra-thread | 1.74× |
| LU | rhs.f:318-328 | 8.4% | +10.6% | inter-thread | 1.36× |
| Stencil | stencil.c:16-21 | 95.7% | +1.55% | inter-thread | 10.9× |
| 3D tensor | mt.c: 22-22 | 69.4% | +2.4% | inter-thread | 1.44× |
| streamcluster2 | streamcluster.cpp:653 | 14.1% | +15.2% | inter-thread | 6.72× |

# Benchmarks

| Benchmarks | bottleneck region | % of total latency | overhead | OPT method | Speedups |
|---|---|---|---|---|---|
| lulesh2.0 | lulesh.cc: 604-609 | 3.6% | +3.1% | inter-thread | 1.43× |
| IRSmk | rmatmult3.c: 86-103 | 78.6% | +3.2% | intra-thread | 4.86× |
| needle | needle.cpp:185-187 | 20% | +2.99% | inter-thread | 2.37× |
| srad | srad.cpp:136-167 | 80.1% | +2.47% | intra-thread | 1.74× |
| LU | rhs.f:318-328 | 8.4% | +10.6% | inter-thread | 1.36× |
| Stencil | stencil.c:16-21 | 95.7% | +1.55% | inter-thread | 10.9× |
| 3D tensor | mt.c: 22-22 | 69.4% | +2.4% | inter-thread | 1.44× |
| streamcluster2 | streamcluster.cpp:653 | 14.1% | +15.2% | inter-thread | 6.72× |

# Benchmarks

| Benchmarks | bottleneck region | % of total latency | overhead | OPT method | Speedups |
|---|---|---|---|---|---|
| lulesh2.0 | lulesh.cc: 604-609 | 3.6% | +3.1% | inter-thread | 1.43× |
| IRSmk | rmatmult3.c: 86-103 | 78.6% | +3.2% | intra-thread | 4.86× |
| needle | needle.cpp:185-187 | 20% | +2.99% | inter-thread | 2.37× |
| srad | srad.cpp:136-167 | 80.1% | +2.47% | intra-thread | 1.74× |
| LU | rhs.f:318-328 | 8.4% | +10.6% | inter-thread | 1.36× |
| Stencil | stencil.c:16-21 | 95.7% | +1.55% | inter-thread | 10.9× |
| 3D tensor | mt.c: 22-22 | 69.4% | +2.4% | inter-thread | 1.44× |
| streamcluster2 | streamcluster.cpp:653 | 14.1% | +15.2% | inter-thread | 6.72× |

# Benchmarks

| Benchmarks | bottleneck region | % of total latency | overhead | OPT method | Speedups |
|---|---|---|---|---|---|
| lulesh2.0 | lulesh.cc: 604-609 | 3.6% | +3.1% | inter-thread | 1.43× |
| IRSmk | rmatmult3.c: 86-103 | 78.6% | +3.2% | intra-thread | 4.86× |
| needle | needle.cpp:185-187 | 20% | +2.99% | inter-thread | 2.37× |
| srad | srad.cpp:136-167 | 80.1% | +2.47% | intra-thread | 1.74× |
| LU | rhs.f:318-328 | 8.4% | +10.6% | inter-thread | 1.36× |
| Stencil | stencil.c:16-21 | 95.7% | +1.55% | inter-thread | 10.9× |
| 3D tensor | mt.c: 22-22 | 69.4% | +2.4% | inter-thread | 1.44× |
| streamcluster2 | streamcluster.cpp:653 | 14.1% | +15.2% | inter-thread | 6.72× |

# Benchmarks

| Benchmarks | bottleneck region | % of total latency | overhead | OPT method | Speedups |
|---|---|---|---|---|---|
| lulesh2.0 | lulesh.cc: 604-609 | 3.6% | +3.1% | inter-thread | 1.43× |
| IRSmk | rmatmult3.c: 86-103 | 78.6% | +3.2% | intra-thread | 4.86× |
| needle | needle.cpp:185-187 | 20% | +2.99% | inter-thread | 2.37× |
| srad | srad.cpp:136-167 | 80.1% | +2.47% | intra-thread | 1.74× |
| LU | rhs.f:318-328 | 8.4% | +10.6% | inter-thread | 1.36× |
| Stencil | stencil.c:16-21 | 95.7% | +1.55% | inter-thread | 10.9× |
| 3D tensor | mt.c: 22-22 | 69.4% | +2.4% | inter-thread | 1.44× |
| streamcluster2 | streamcluster.cpp:653 | 14.1% | +15.2% | inter-thread | 6.72× |

Related work: MACPO (selective instrumentation): 2x - 5x

# Outline

✓Lightweight profiling

✓SMT-aware optimization

• *Detection of cache conflicts*

• Guiding data-structure layout transformation

# Lightweight Detection of Cache Conflicts

[CGO – 2018]

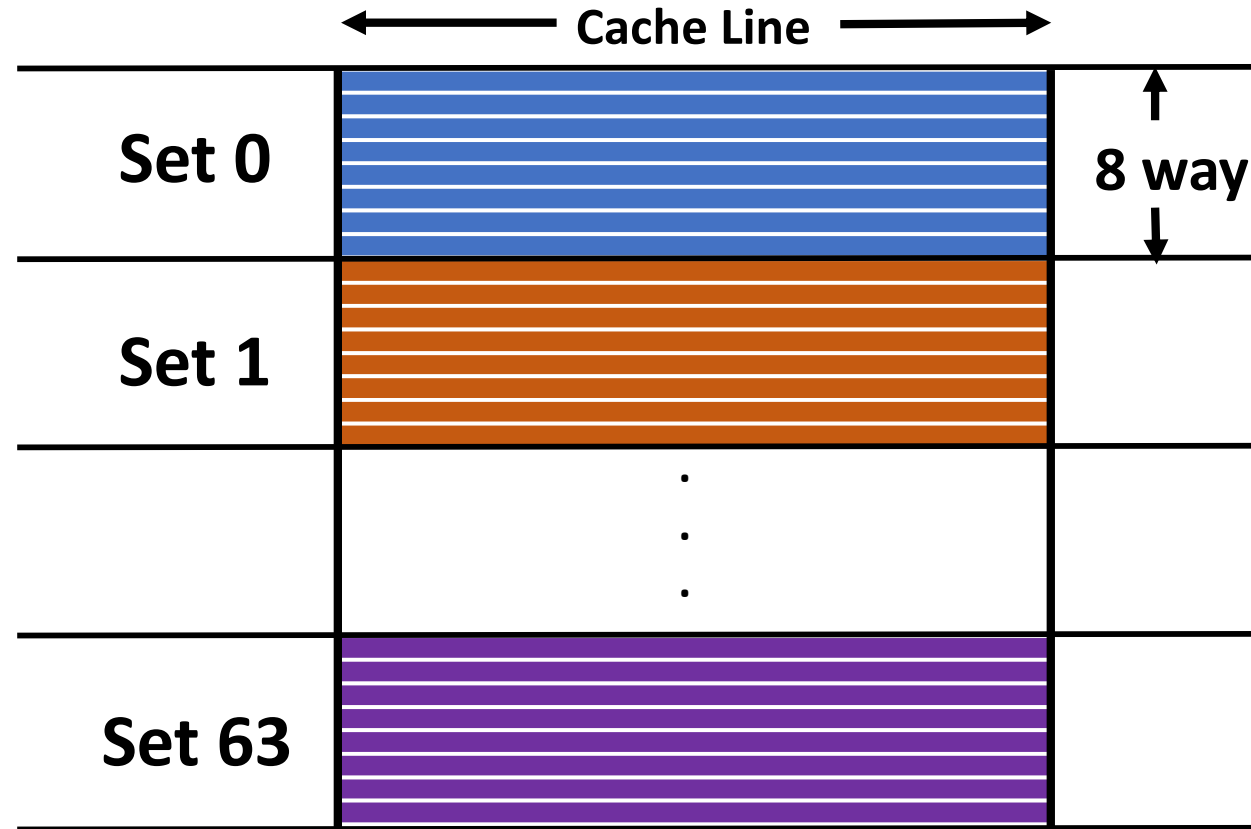**Probir Roy**, Shuaiwen Leon Song, Sriram Krishnamoorthy, *Xu Liu*

WILLIAM & MARY

CHARTERED 1693

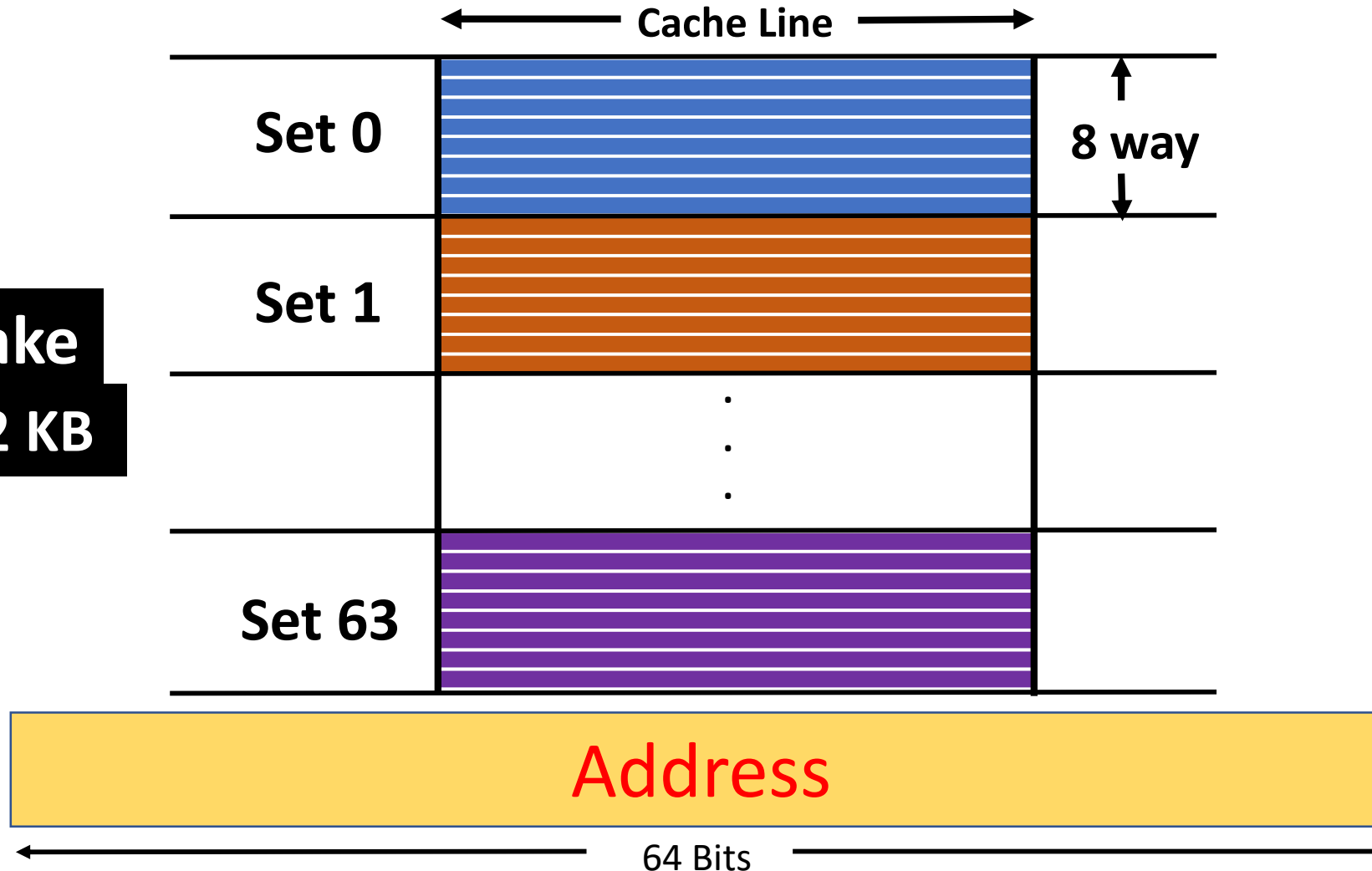Pacific Northwest
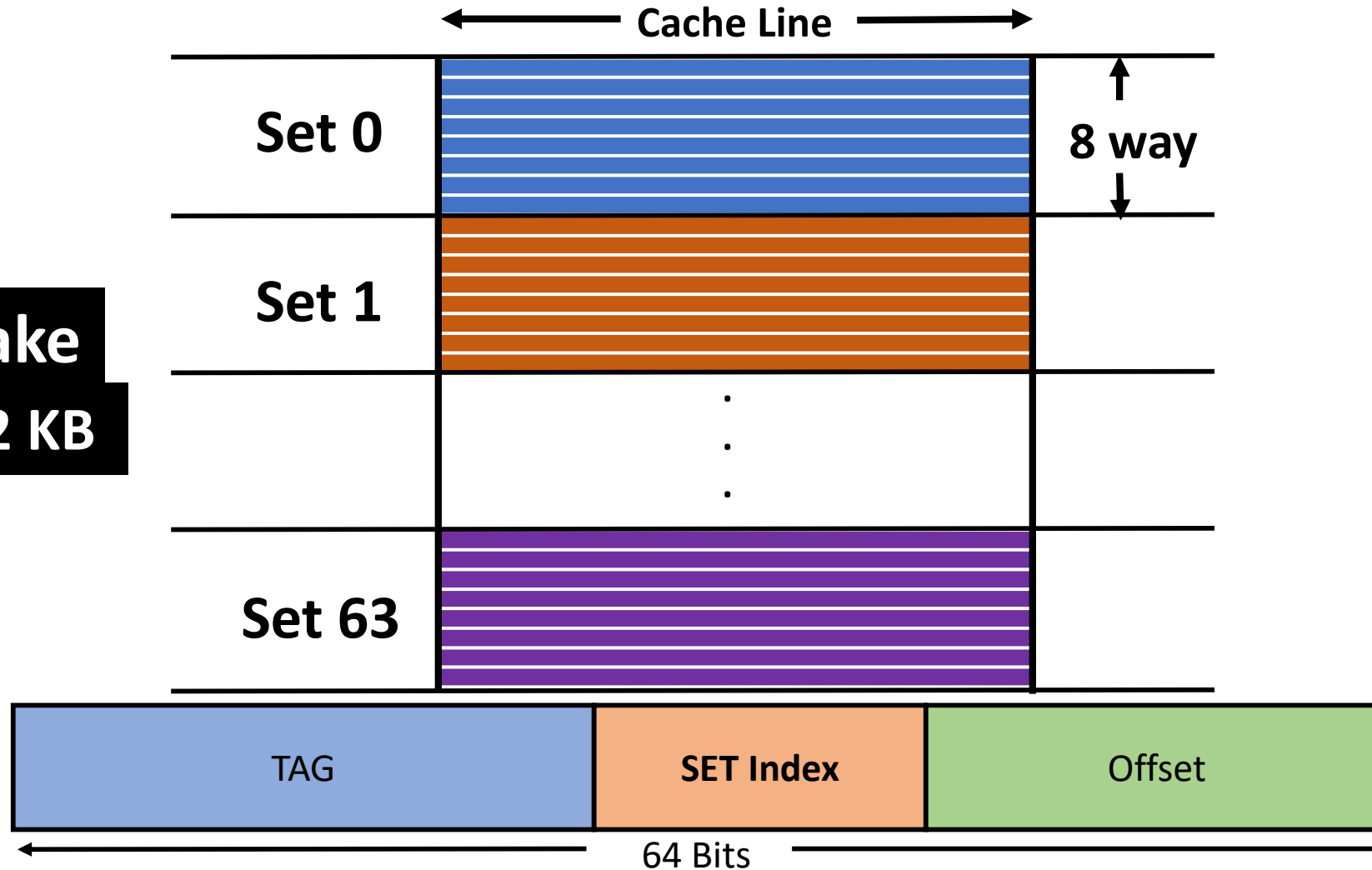
NATIONAL LABORATORY

# Set-associative cache



Cache Line

Set 0

8 way

Set 1

Intel Skylake
L1 cache: 32 KB

.
.
.

Set 63

# Set-associative cache



**Cache Line**

**Set 0**

**8 way**

**Set 1**

.
.
.

**Intel Skylake**

**L1 cache: 32 KB**

**Set 63**

**Address**

**64 Bits**

# Set-associative cache



Intel Skylake
L1 cache: 32 KB

Cache Line

Set 0

Set 1

8 way

Set 63

| TAG | SET Index | Offset |
|---|---|---|

64 Bits

# Set-associative cache

**Cache Line**

**Set 0**

**8 way**

**Intel Skylake**

**L1 cache: 32 KB**

**Set 1**

.
.
.

**Set 63**

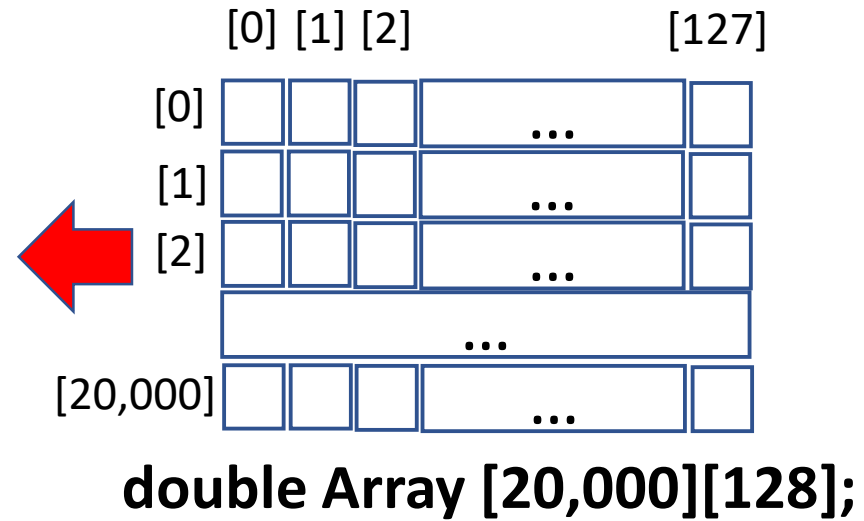| TAG | SET Index | Offset |
|-----|-----------|--------|

64 Bits

# Set conflict



**double Array [20,000][128];**

# Set conflict

**Set mapping**



double Array [20,000][128];

# Set conflict

**Set mapping**



double Array [20,000][128];

# Set conflict



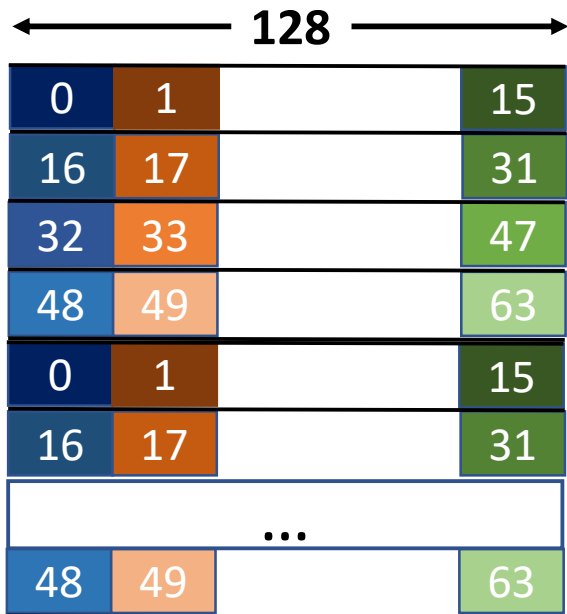**Set mapping**

double Array [20,000][128];

# Set conflict

**Set mapping**



double Array [20,000][128];

# Set conflict

**Set mapping**



**Set mapping after padding**

**double Array [20,000][128];**

# Set conflict



Set mapping

double Array [20,000][128];

Set mapping after padding

# Set conflict



Set mapping

Set mapping after padding

double Array [20,000][128];

# Set conflict

[0] [1] [2]          [127] **Pad**

[0]

**Set mapping**

**Set mapping**

**Set mapping**

## Is your application suffering conflict cache miss?

| 0 | 16 | 32 | 48 | 0 | 16 | 32 | 48 | 0 |

| 0 | 17 | 34 | 51 | 4 | 21 | 48 | 55 | 8 |

**Time**

# Trace driven cache simulation

**Simulation methods**

Memory trace

Cache simulator

Classifying miss

| A[0][0] | A[1][0] | A[2][0] | A[0][0] | A[3][0] | A[2][0] | ... | A[0][0] | **Time** |

L1 cache

Conflict cache miss

# Trace driven cache simulation

Overhead: average 38 times

**Simulation methods**

| Memory trace |

| Cache simulator |

| Classifying miss |

| A[0][0] | A[1][0] | A[2][0] | A[0][0] | A[3][0] | A[2][0] | ... | A[0][0] | → **Time**

L1 cache

Conflict cache miss

# Trace driven cache simulation

Overhead: average 38 times

Xiang, Xiaoya, Chen Ding, Hao Luo, and Bin Bao. "HOTL: a higher order theory of locality." *ACM SIGPLAN Notices* 48, no. 4 (2013): 343-356.

**Simulation methods**

**Memory trace**

| A[0][0] | A[1][0] | A[2][0] | A[0][0] | A[3][0] | A[2][0] | … | A[0][0] |

Time

**Cache simulator**

L1 cache

**Classifying miss**

Conflict cache miss

# Trace driven cache simulation

Overhead: average 38 times

Xiang, Xiaoya, Chen Ding, Hao Luo, and Bin Bao. "HOTL: a higher order theory of locality." *ACM SIGPLAN Notices* 48, no. 4 (2013): 343-356.

## Simulation methods

**Memory trace**

↓

**Cache simulator**

↓

**Classifying miss**

| | | | | A[3][0] | A[2][0] | … | A[0][0] | **Time** |

**High overhead**

L1 cache

Conflict cache miss

# Trace driven cache simulation

## Simulation methods

Memory trace

↓

Cache simulator

↓

Classifying miss

Overhead: average 38 times

Xiang, Xiaoya, Chen Ding, Hao Luo, and Bin Bao. "HOTL: a higher order theory of locality." *ACM SIGPLAN Notices* 48, no. 4 (2013): 343-356.

**High overhead**

A[3][0] | A[2][0] | ... | A[0][0]    Time

**Difficult to simulate hardware**

Conflict cache miss

# Trace driven cache simulation

**Simulation methods**

Memory trace

Cache simulator

Classifying miss

Overhead: average 38 times

Xiang, Xiaoya, Chen Ding, Hao Luo, and Bin Bao. "HOTL: a higher order theory of locality." *ACM SIGPLAN Notices* 48, no. 4 (2013): 343-356.

**High overhead**
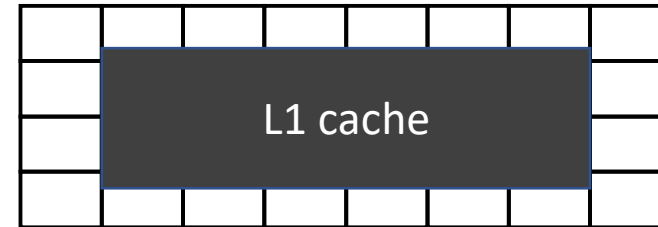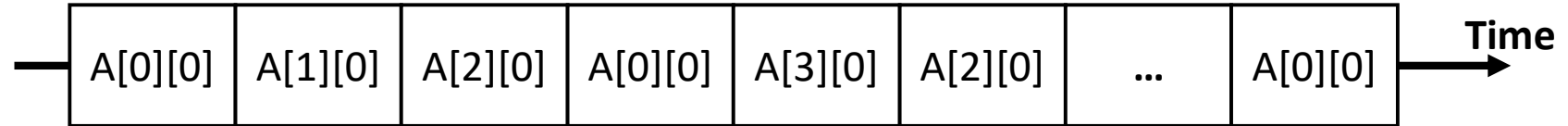
A[3][0] | A[2][0] | ... | A[0][0] | → Time

**Difficult to simulate hardware**

**Theoretically accurate** **Difficult in practice**

Exploiting Modern Hardware Features via Lightweight Profiling

25

# A practical low overhead solution

**Simulation methods**

Memory trace

⬇

Cache
simulator

⬇

Classifying miss

# A practical low overhead solution

CCProf

**Simulation methods**

Memory trace

⬇

Cache simulator

⬇

Classifying miss

# A practical low overhead solution

CCProf

**Simulation methods**

Memory trace

⬇

Cache simulator

⬇

Classifying miss

**Measurement methods**

# A practical low overhead solution

CCProf

**Simulation methods**

Memory trace

⬇

Cache simulator

⬇

Classifying miss

**Measurement methods**

Memory sampling

⬇

Statistical analysis

⬇

Classifying miss

# A practical low overhead solution

CCProf

**Simulation methods**

| Memory trace |
| --- |

↓

| Cache simulator |
| --- |

↓

| Classifying miss |
| --- |

**Overhead**

>>

**Accuracy**

~

**Measurement methods**

| Memory sampling |
| --- |

↓

| Statistical analysis |
| --- |

↓

| Classifying miss |
| --- |

# Hardware-based address sampling (Cont.)

**Memory references**

| A[0][0] | A[1][0] | A[4][0] | A[0][0] | A[1][0] | A[2][0] | ... | A[0][0] |

**Time**

| L1 Miss | L1 Miss | L1 Miss | L1 Hit | L1 Hit | L1 Miss | | L1 Miss |

# Hardware-based address sampling (Cont.)

# Hardware-based address sampling (Cont.)

**Memory references**

| A[0][0] | A[1][0] | A[4][0] | A[0][0] | A[1][0] | A[2][0] | ... | A[0][0] |
|---------|---------|---------|---------|---------|---------|-----|---------|

→ **Time**

**L1 Miss**  **L1 Miss**  **L1 Miss**  L1 Hit  L1 Hit  **L1 Miss**  **L1 Miss**

**Precise event sampling (PEBS)**

| A[0][0] | | A[4][0] | | A[2][0] |

**PMU**

| TAG | SET Index | Offset |

# Hardware-based address sampling (Cont.)

# Hardware-based address sampling (Cont.)

# Observation: *temporal pattern* of *conflict*



**Set mapping**

**Set mapping after padding**

[0] [1] [2]      [127]

Array [20,000][128]

# *Observation*: *temporal pattern* of *conflict* (cont.)



**Conflict**

| 0 | 16 | 32 | 48 | 0 | 16 | 32 | 48 | 0 | → Time |

**No conflict**

| 0 | 17 | 34 | 51 | 4 | 21 | 48 | … | 0 | → Time |

# *Observation*: *temporal pattern* of *conflict* (cont.)

# *Observation*: *<u>temporal pattern</u>* of *<u>conflict</u>* (cont.)

# *Observation*: *temporal pattern* of *conflict* (cont.)



**Conflict**

★ ★ ★

| 0 | 16 | 32 | 48 | 0 | 16 | 32 | 48 | 0 |

Time

**Distance = 3**  **Distance = 3**

**No conflict**

★ ★

| 0 | 17 | 34 | 51 | 4 | 21 | 48 | ... | 0 |

Time

**Distance = 63**

# Re-conflict Distance (RCD)

- Number of **_cache misses_** in other cache sets between two consecutive misses in one particular set

# Re-conflict Distance (RCD)

- Number of ***cache misses*** in other cache sets between two consecutive misses in one particular set

# Re-conflict Distance (RCD)

- Number of **_cache misses_** in othe... consecutive misses in one partic...

**Approximate RCD**

# RCD and it's contribution

| RCD | Count | Is conflict? |
| --- | --- | --- |
| Short | Large | Yes |
| Short | Small | No |
| Long | ~ | No |

# RCD and it's contribution

| RCD | Count | Is conflict? |
|---|---|---|
| Short | Large | Yes |
| Short | Small | No |
| Long | ~ | No |

**Training**

Benchmark ---> RCD ---> **Regression model**

**Prediction**

Application ---> RCD ---> **Model** ---> Conflict / No-Conflict

# Case Study: PolyBench/C -ADI

**CCPROF PREDICTS >>> *** CONFLICT MISS *** in LOOP(line: 102). Loop contribution is *** HIGH *** 94.26%**

CCPROF PREDICTS >>> *** **NO CONFLICT MISS** *** in loop(line: 108). Loop's contribution to total L1 miss: 3.13%

CCPROF PREDICTS >>> *** **NO CONFLICT MISS** *** in loop(line: 117). Loop's contribution to total L1 miss: 0.86%

CCPROF PREDICTS >>> *** **NO CONFLICT MISS** *** in loop(line: 122). Loop's contribution to total L1 miss: 1.74%

# Case Study: PolyBench/C -ADI

**CCPROF PREDICTS >>> *** <span style="color:red">CONFLICT MISS</span> *** in LOOP(line: 102). Loop contribution is *** <span style="color:red">HIGH</span> *** 94.26%**

CCPROF PREDICTS >>> *** **NO CONFLICT MISS** *** in loop(line: 108). Loop's contribution to total L1 miss: 3.13%

CCPROF PREDICTS >>> *** **NO CONFLICT MISS** *** in loop(line: 117). Loop's contribution to total L1 miss: 0.86%

CCPROF PREDICTS >>> *** **NO CONFLICT MISS** *** in loop(line: 122). Loop's contribution to total L1 miss: 1.74%

```
1 //column sweep
2 for (j=1; j<_PB_N-1; j++) {
3  p[i][j] = -c / (a*p[i][j-1]+b);
4  q[i][j] = (-d*u[j][i-1]+(SCALAR_VAL(1.0)+SCALAR_VAL(2.0)
       *d)*u[j][i] - f*u[j][i+1]-a*q[i][j-1])/(a*p[i][j
       -1]+b);
5 }
```

# RCD – *before* and *after* optimization

*Loop level Speedup

# RCD – *before* and *after* optimization



Needleman-Wunsch

Poly ADI

MKL-FFT

**Speedup: 3×**

**Speedup: 1.26×**

**Speedup: 1.13×**

Cumulative probability of L1 cache miss

**Median overhead: 37%**

**Compare with simulation: 38x**

Average RCD

— Original  — Padded

*Loop level Speedup

# Outline

✓ Lightweight profiling

✓ SMT-aware optimization

✓ Detection of cache conflicts

- *Guiding data-structure layout transformation*

# StructSlim: A lightweight profiler to guide structure splitting

## [CGO – 2016]
**Probir Roy** , *Xu Liu*

# LWPTool: A Lightweight Profiler to Guide Data Layout Optimization

## [TPDS – 2018]
Chao Yu,  **Probir Roy**, Yuebin Bai, Hailong Yang, *Xu Liu*

# StructSlim: A lightweight profiler to guide structure splitting

[CGO – 2016]
**Probir Roy** , *Xu Liu*

# LWPTool: A Lightweight Profiler to Guide Data Layout Optimization

[TPDS – 2018]
Chao Yu,  **Probir Roy**, Yuebin Bai, Hailong Yang, *Xu Liu*

# Inefficient data-structure

```
struct type {int a; int b; int c; int d;};
struct type Arr[N];
for (i = 0; i < N; i++)
        B[i] = Arr[i].a + Arr[i].c;
```

**L1 cache**



| a | b | c | d | a | b | c | d |

**Cache line**

# Inefficient data-structure

```
struct type {int a; int b; int c; int d;};
struct type Arr[N];
for (i = 0; i < N; i++)
        B[i] = Arr[i].a + Arr[i].c;
```

**Utilization = 50%**

**L1 cache**

| a | b | c | d | a | b | c | d | **Cache line** |

# Inefficient data-structure

Utilization = 50%

**L1 cache**

```
struct type {int a; int b; int c; int d;};
struct type Arr[N];
for (i = 0; i < N; i++)
        B[i] = Arr[i].a + Arr[i].c;
```

| a | b | c | d | a | b | c | d | **Cache line** |

**Split structure**

**L1 cache**

```
struct type_part1 {int a; int c;};
struct type_part2 {int b; int d;};
```

| a | c | a | c | a | c | a | c | **Cache line** |

# Inefficient data-structure

```
struct type {int a; int b; int c; int d;};
struct type Arr[N];
for (i = 0; i < N; i++)
        B[i] = Arr[i].a + Arr[i].c;
```

**L1 cache**

Utilization = 50%

| a | b | c | d | a | b | c | d |

**Cache line**

**Split structure**

```
struct type_part1 {int a; int c;};
struct type_part2 {int b; int d;};
```

**L1 cache**

Utilization = 100%

| a | c | a | c | a | c | a | c |

**Cache line**

# Structure splitting- Questions to ask

How to split structure?

Which data structures are significant?

Which fields to keep together?

# Structure splitting- Questions to ask

How to split structure?

Which data structures are significant?

High usage

Which fields to keep together?

# Structure splitting- Questions to ask

How to split structure?

Which data structures are significant?

High usage

Which fields to keep together?

Loop level analysis

Field affinity

# Structure splitting- Questions to ask

How to split structure?

Which data structures are significant?

High usage

Which fields to keep together?

Loop level analysis  Field affinity



**Loop 1**  Field 1  Field 2  Field 3  **Loop 2**

# Structure splitting- Questions to ask

**How to split structure?**

**Which data structures are significant?**

**High usage**

**Which fields to keep together?**

**Loop level analysis**     **Field affinity**



Loop 1

Field 1

Field 2

Field 3

Loop 2

Field 1
Field 2

or

Field 2
Field 3

# Structure splitting- Questions to ask
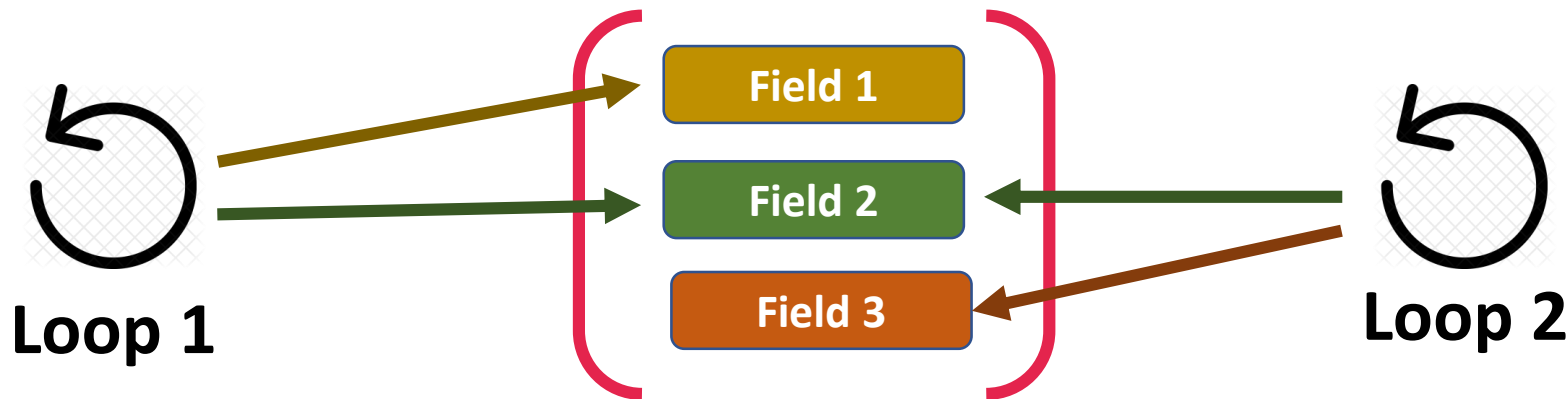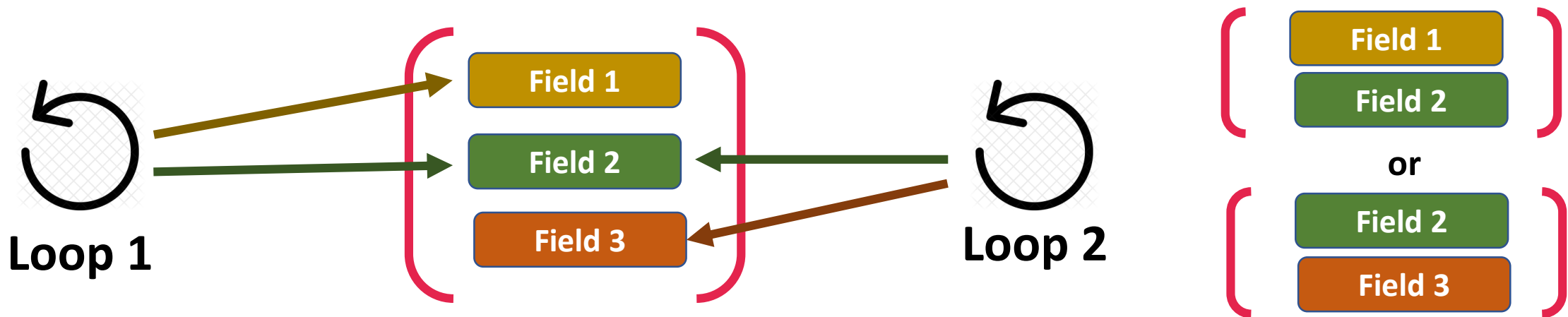
How to split structure?

Which data structures are significant?

High usage

Which fields to keep together?

Loop level analysis     Field affinity

# Structure splitting- Questions to ask

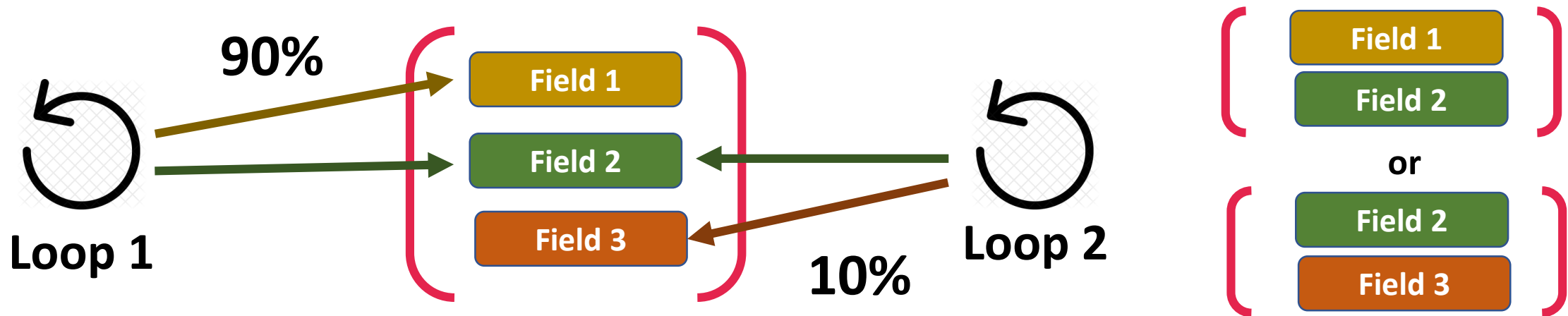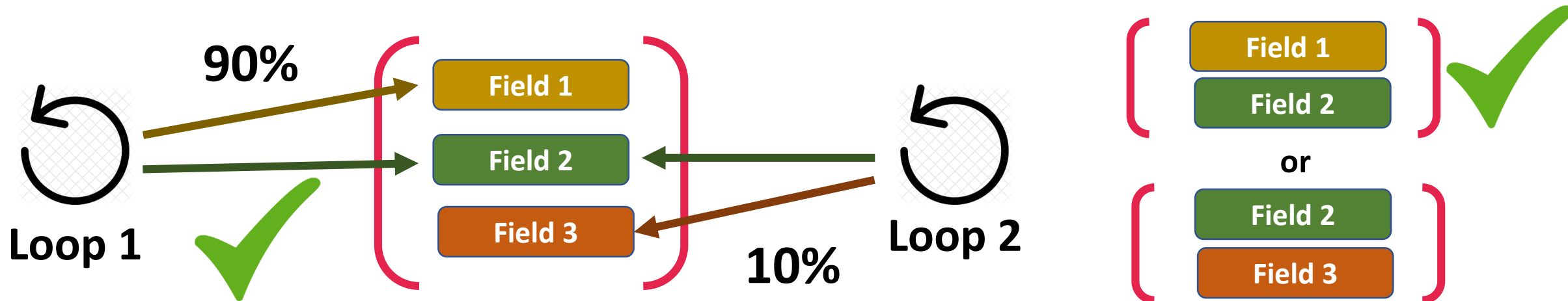**How to split structure?**

**Which data structures are significant?**

High usage

**Which fields to keep together?**

Loop level analysis   Field affinity

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

**PMU**

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution



PMU

Application                                    Time

S          S                  S        S

Reference Type

Data Address

Instruction Pointer

Loop 1

Loop 2

Heap

Mem Allocation
Monitor

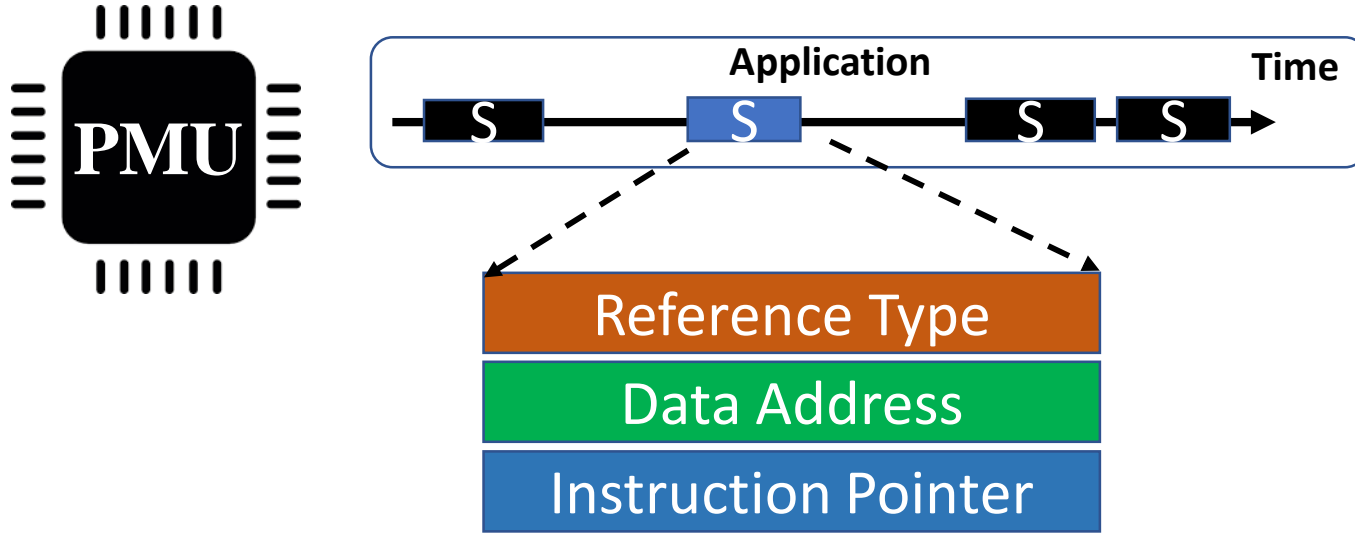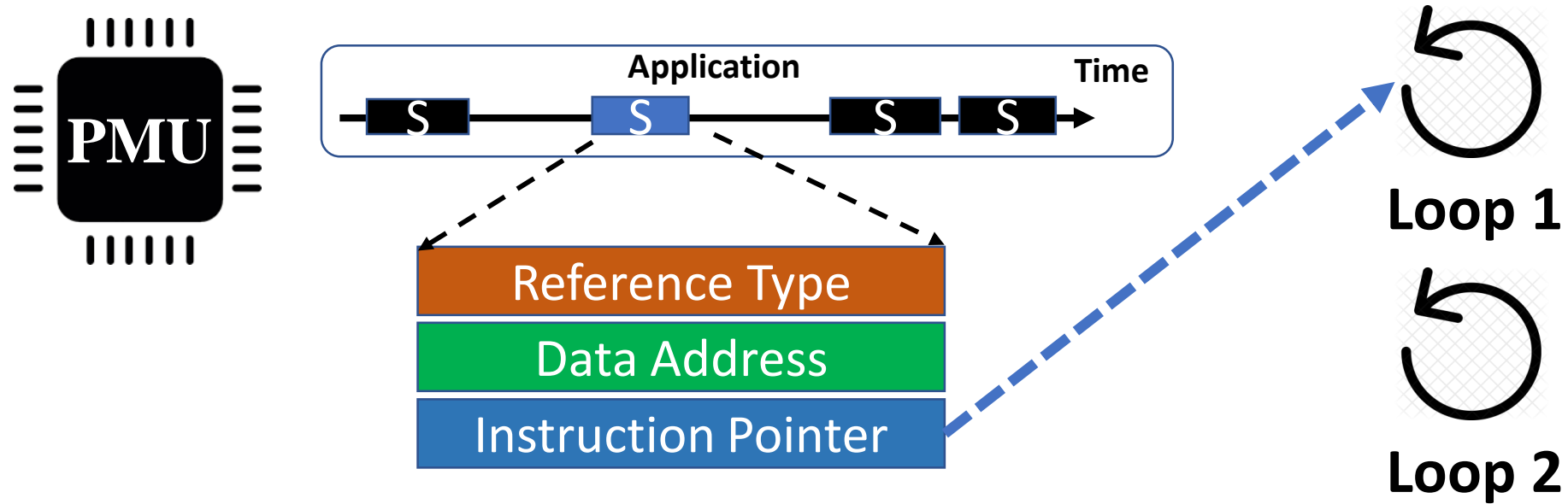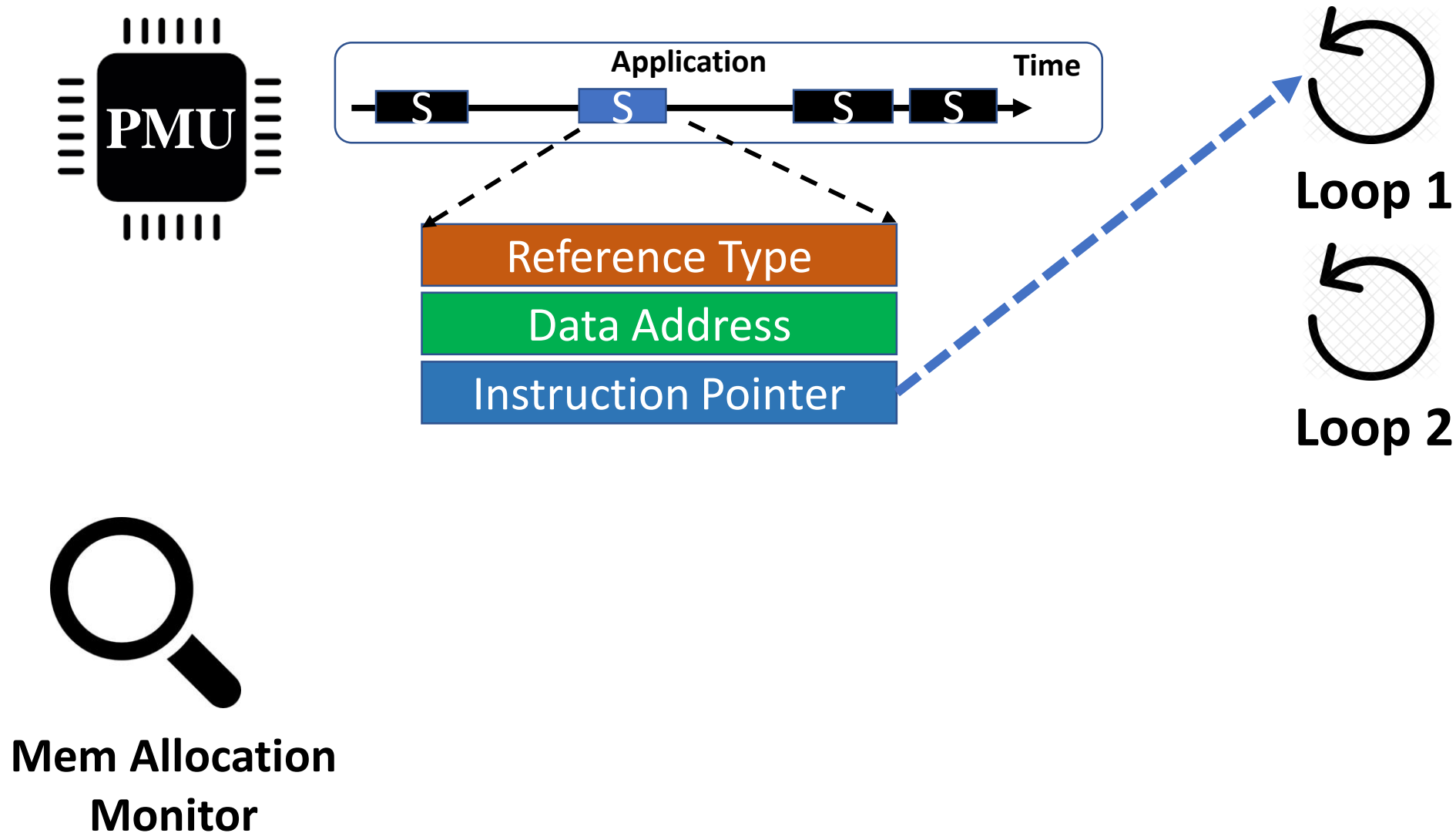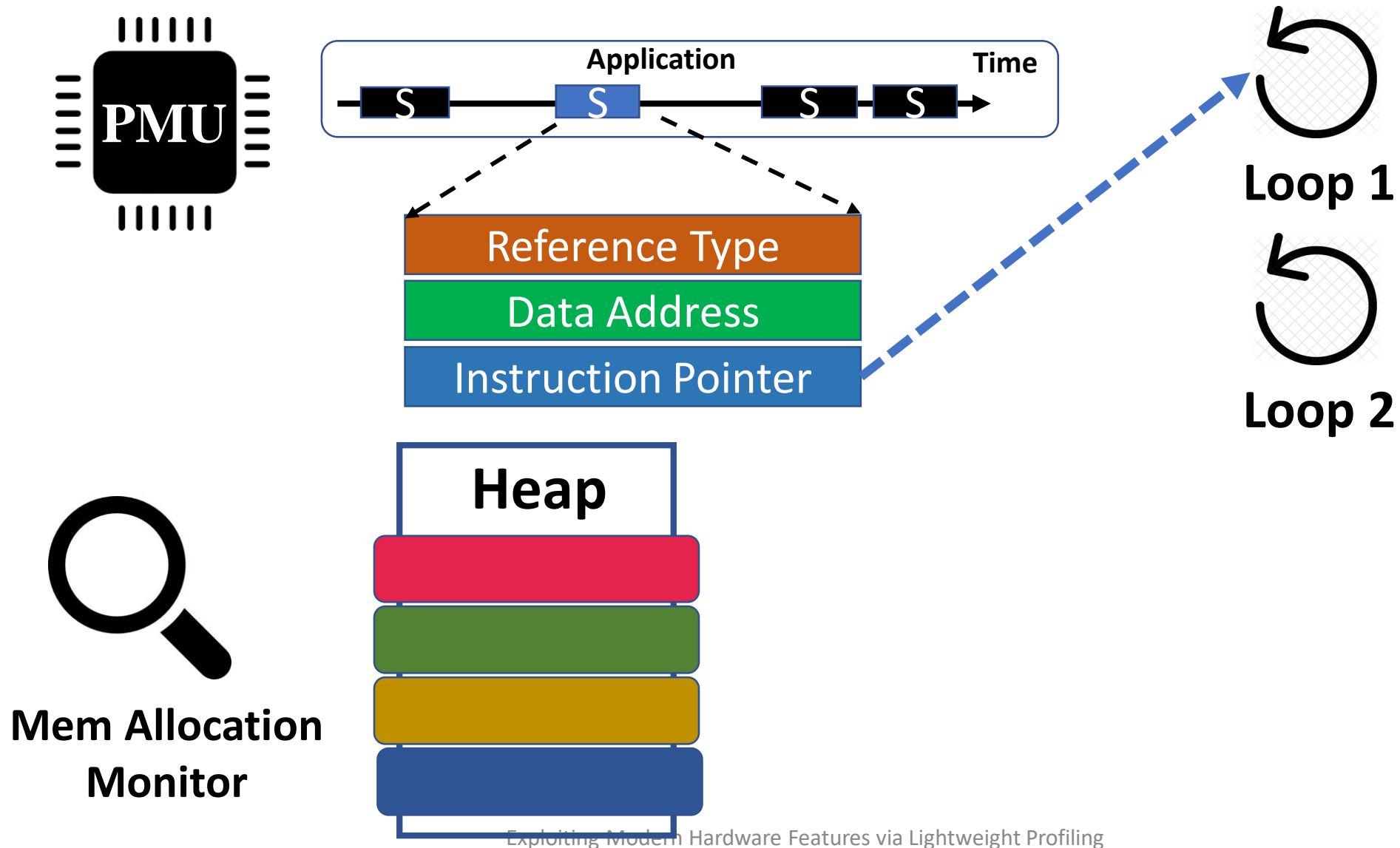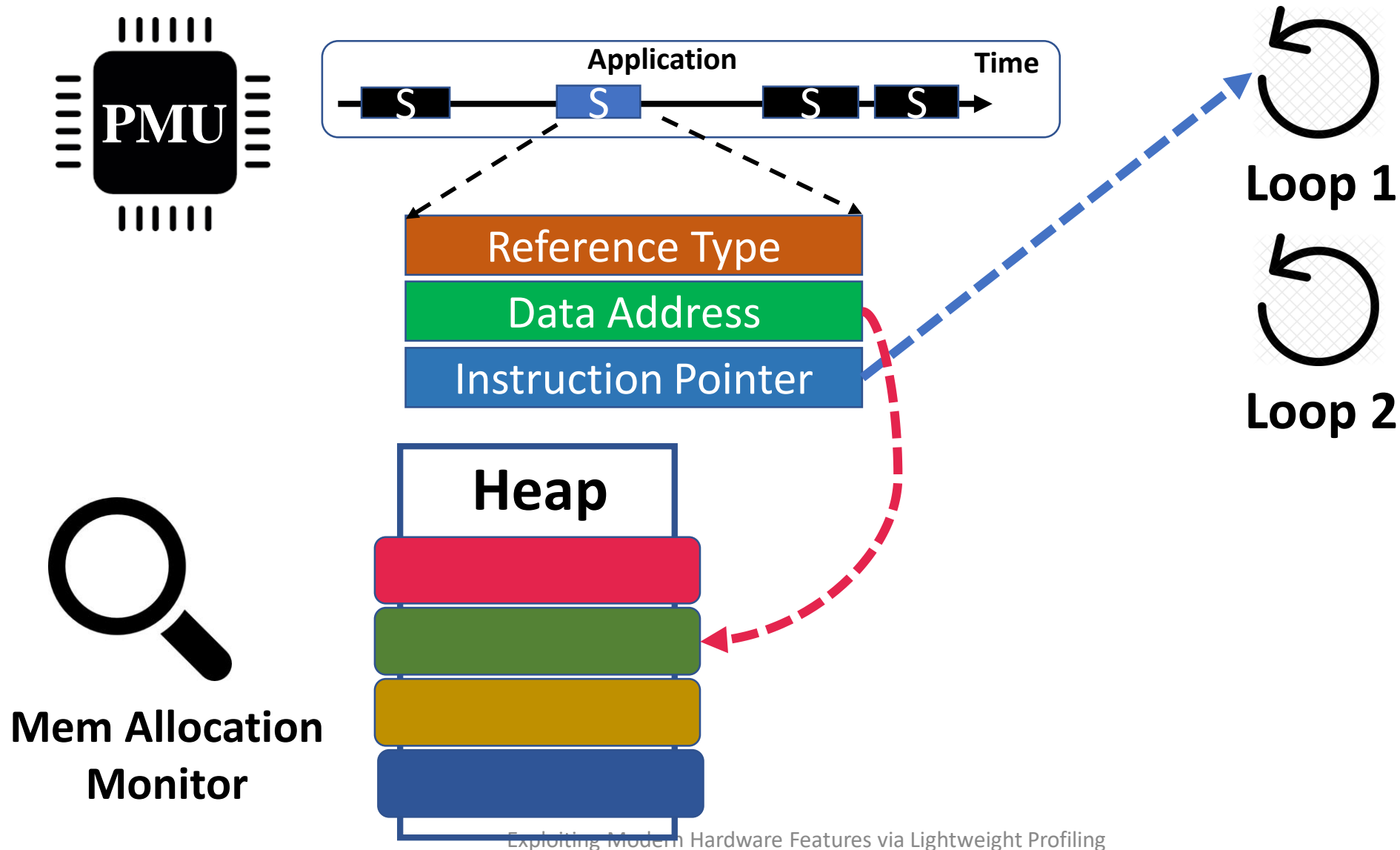# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

# Code-centric and data-centric attribution

# Case study: SPEC CPU 2000 ART

```
typedef struct
{
    double *I; double W; double X; double V; double U; double P; double Q; double R;
}f1_neuron
```

| Loops with line numbers | Latency percentage | Accessed fields |
|---|---|---|
| 131-138 | 1.59% | U,P |
| 559-570 | 8.42% | X,Q |
| 553-554 | 1.98% | W |
| 545-548 | 10.83% | U, I |
| 615-616 | 56.57% | P |
| 607-608 | 14.40% | P |
| 589-592 | 2.25% | U, P |
| 575-576 | 3.72% | V |
| 1015-1016 | 0.24% | I |



**Affinity graph**

# Case study: SPEC CPU 2000 ART

```
typedef struct
{
    double *I; double W; double X; double V; double U; double P; double Q; double R;
}f1_neuron
```

| Loops v... | |
|---|---|
| nu... | |
| 13... | |
| 55... | |
| 55... | |
| 54... | |
| 61... | |
| 60... | |
| 58... | |
| 57... | |
| 1015... | |

```
typedef struct{ double *I; double U;} f1_neuron_IU;
typedef struct{ double Q; double X;} f1_neuron_QX;
typedef struct{ double P;} f1_neuron_P;
typedef struct{ double V;} f1_neuron_V;
typedef struct{ double W;} f1_neuron_W;
typedef struct{ double R;} f1_neuron_R;
```

# Benchmarks: speedup, overhead, cache miss

| Benchmarks | Speedups | Runtime overhead | L1 miss reduction | L2 miss reduction |
|---|---|---|---|---|
| 179.ART | 1.37× | 2.05% | 46.5% | 51.1% |
| 462.Libquantum | 1.09× | 2.79% | 49% | 82.6% |
| TSP | 1.09× | 2.42% | 13.3% | 19.9% |
| Mser | 1.03× | 2.95% | 8.3% | 8.4% |
| CLOMP 1.2 | 1.25× | 16.1% | 15.5% | 26.4% |
| Health | 1.12× | 18.3% | 66.7% | 90.8% |
| NN | 1.33× | 5.21% | 87.2% | 98.0% |
| **Average** | **1.18×** | **7.1%** | | |

gcc -O3

Related work: Overhead: average 4x

Yan, Jianian, Jiangzhou He, Wenguang Chen, Pen-Chung Yew, and Weimin Zheng. "ASLOP: A field-access affinity-based structure data layout optimizer."

# Conclusions

**Lightweight profiling with PMUs can provide deep insights into performance issues cause by memory hierarchies and poor algorithm choice.**

Deep insight

Simulation methods (PinTool, GPGPUSim, GEMS)

SMTAnalyzer

StructSlim

CCProf

Shallow insight

Measurement methods (Perf, Oprofile, PAPI)

High overhead

Low overhead

# Publications

- **[CGO'18] "Lightweight Detection of Cache Conflicts",** _Probir Roy_**,** Shuaiwen Leon Song, Sriram Krishnamoorthy and Xu Liu, The 2018 International Symposium on Code Generation and Optimization, Feb 24 - 28th, 2018, Vienna, Austria. Acceptance ratio: 28%.

- **[TACO'18] "NUMA-Caffe: NUMA-Aware Deep Learning Neural Networks",** _Probir Roy_**,** Shuaiwen Leon Song, Sriram Krishnamoorthy, Abhinav Vishnu, Dipanjan Sengupta, Xu Liu, ACM Transactions on Architecture and Code Optimization, 2018.
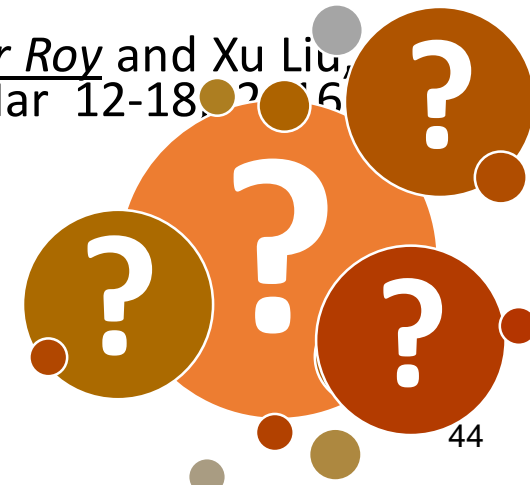
- **[TPDS'18] "LWPTool: A Lightweight Profiler to Guide Data Layout Optimization"**, Chao Yu, _Probir Roy_, Yuebin Bai, Hailong Yang, Xu Liu, IEEE Transactions on Parallel and Distributed Systems, 2018.

- **[HPDC'16] "SMT-Aware Instantaneous Footprint Optimization",** _Probir Roy_, Xu Liu and Shuaiwen Leon Song, The 25th ACM International Symposium on High-Performance and Distributed Computing, May 31 - Jun 4, 2016, Kyoto, Japan. Acceptance ratio: 15.5% (20/129).

- **[CGO'16] "StructSlim: A Lightweight Profiler to Guide Structure Splitting"**, _Probir Roy_ and Xu Liu, The 2016 International Symposium on Code Generation and Optimization, Mar 12-18, 2016, Barcelona, Spain. Acceptance ratio: 23%.

# Publications

- **[CGO'18] "Lightweight Detection of Cache Conflicts",** *Probir Roy,* Shuaiwen Leon Song, Sriram Krishnamoorthy and Xu Liu, The 2018 International Symposium on Code Generation and Optimization, Feb 24 - 28th, 2018, Vienna, Austria. Acceptance ratio: 28%.

- **[TACO'18] "NUMA-Caffe: NUMA-Aware Deep Learning Neural Networks",** *Probir Roy,* Shuaiwen Leon Song, Sriram Krishnamoorthy, Abhinav Vishnu, Dipanjan Sengupta, Xu Liu, ACM Transactions on Architecture and Code Optimization, 2018.

- **[TPDS'18] "LWPTool: A Lightweight Profiler to Guide Data Layout Optimization"**, Chao Yu, *Probir Roy*, Yuebin Bai, Hailong Yang, Xu Liu, IEEE Transactions on Parallel and Distributed Systems, 2018.

- **[HPDC'16] "SMT-Aware Instantaneous Footprint Optimization",** *Probir Roy*, Xu Liu and Shuaiwen Leon Song, The 25th ACM International Symposium on High-Performance and Distributed Computing, May 31 - Jun 4, 2016, Kyoto, Japan. Acceptance ratio: 15.5% (20/129).

- **[CGO'16] "StructSlim: A Lightweight Profiler to Guide Structure Splitting"**, *Probir Roy* and Xu Liu, The 2016 International Symposium on Code Generation and Optimization, Mar 12-18, 2016, Barcelona, Spain. Acceptance ratio: 23%.

# Challenges ahead

- Program analysis for declarative programming languages
  - Domain specific languages provide high-level abstraction
    - Machine learning (PyTorch), HPC (HDF5), big-data (SQL)
- Analyzing and optimizing data center and cloud application
  - Resource utilization/scheduling in multi-tenant environment
  - Heterogenous architecture resource management
- Security analysis
  - Program analysis to identify vulnerable source code
- Analysis of emerging hardware
  - GPU, FPGA, Tensor processing units