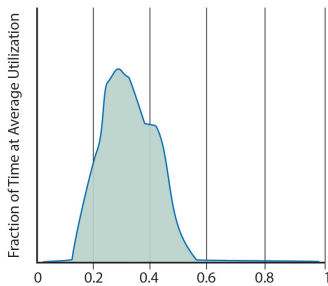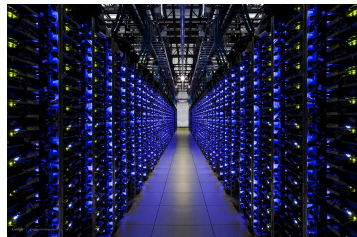# SMiTe: Precise QoS Prediction on Real-System SMT Processors to Improve Utilization in Warehouse Scale Computers
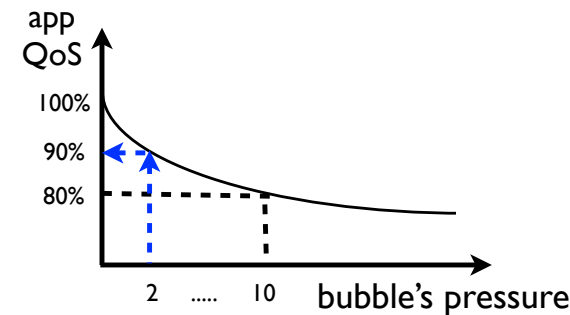
Yunqi Zhang, Michael A. Laurenzano, Jason Mars, Lingjia Tang
Clarity-Lab, Electrical Engineering and Computer Science, University of Michigan, Ann Arbor

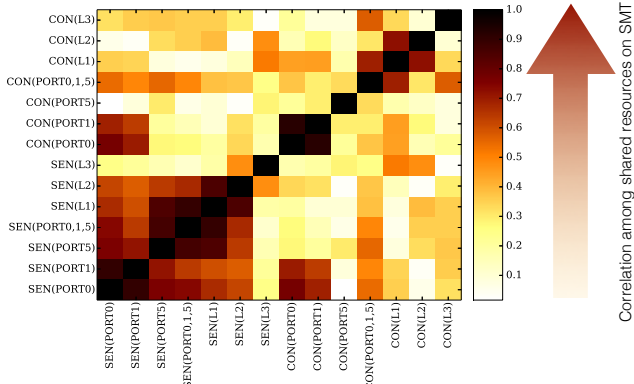## Goal: Improve Data Center Utilization

Precise interference prediction identifies "safe" co-locations to improve server utilization

## SMT Co-location is Harder than CMP

app QoS

100%
90%
80%

2 ..... 10   bubble's pressure

Unified approach for CMP co-location

Unified approach does **not** work for SMT

Correlation among shared resources on SMT

## Solution: Ruler-based Methodology

Max utilization in each resource sharing dimension

| | GPR | SIMD INT | SIMD FP |
|---|---|---|---|
| PORT 0 | ALU | VI_MUL / VI_SHF | FP_MUL / Blend / DIV |
| PORT 1 | ALU | VI_ADD / VI_SHF | FP_ADD |
| PORT 5 | ALU / JMP | | FP_SHF / FP_Boolean / Blend |

Intel® AVX

```
loop:
    mulps    %xmm0,%xmm0
    mulps    %xmm7,%xmm7
    ......
    jmp loop
(a) FP_MUL (PORT0)
```

```
loop:
    addps    %xmm0,%xmm0
    addps    %xmm7,%xmm7
    ......
    jmp loop
(b) FP_ADD (PORT1)
```
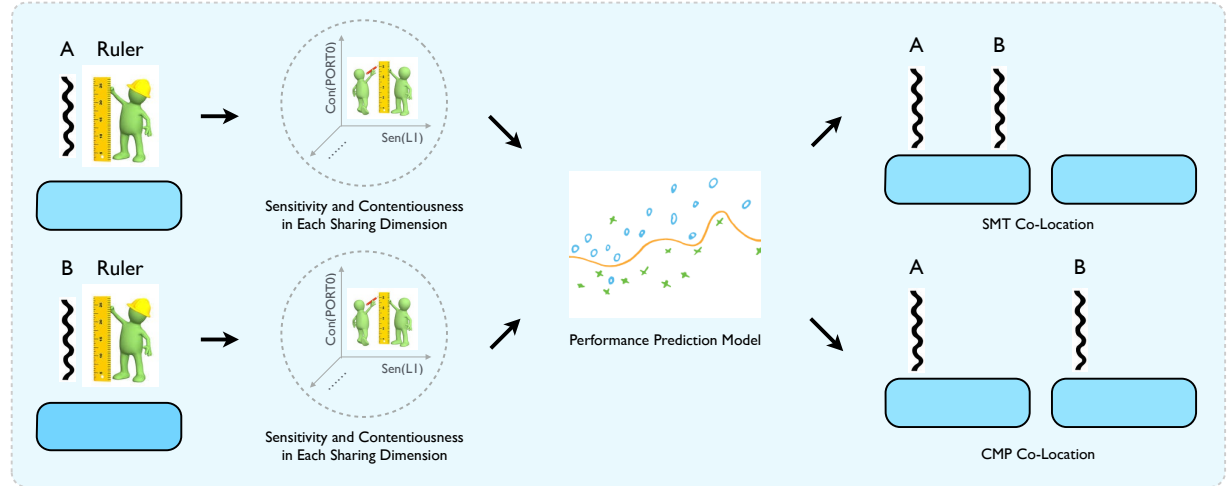
```
loop:
    shufps   %xmm0,%xmm0
    shufps   %xmm7,%xmm7
    ......
    jmp loop
(c) FP_SHF (PORT5)
```

```
loop:
    addl    %eax,%eax
    addl    %edx,%edx
    ......
    jmp loop
(d) INT_ADD (PORT0,1,5)
```
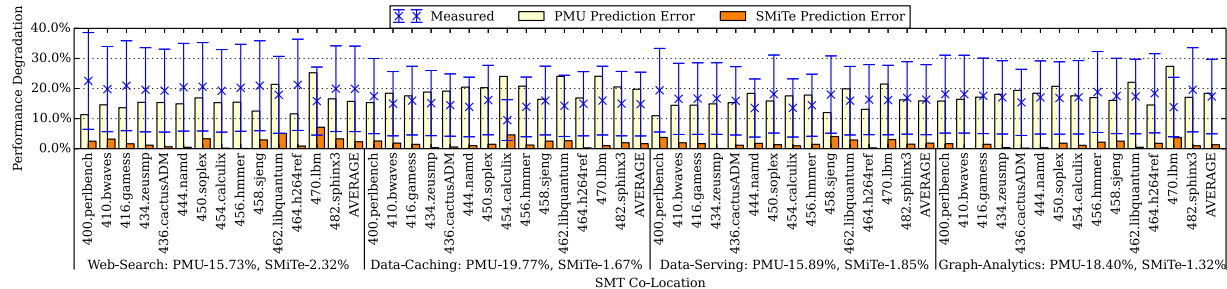
```
#define MASK 0xd0000001u
#define RAND (lfsr = (lfsr >> 1) ^ (unsigned int)(0 - (lfsr & 1u) & MASK))
......
while (1) {
    data_chunk[RAND % FOOTPRINT]++;
    ......
    data_chunk[RAND % FOOTPRINT]++;
}
(e) MEM (L1,L2 Cache)
```

```
......
first_chunk = data_chunk;
second_chunk = data_chunk + FOOTPRINT / 2;
while (1) {
    for (i = 0;i < FOOTPRINT / 2;i += 64) {
        first_chunk[i] = second_chunk[i] + 1;
    }
    for (i = 0;i < FOOTPRINT / 2;i += 64) {
        second_chunk[i] = first_chunk[i] + 1;
    }
}
(f) MEM (L3 Cache)
```
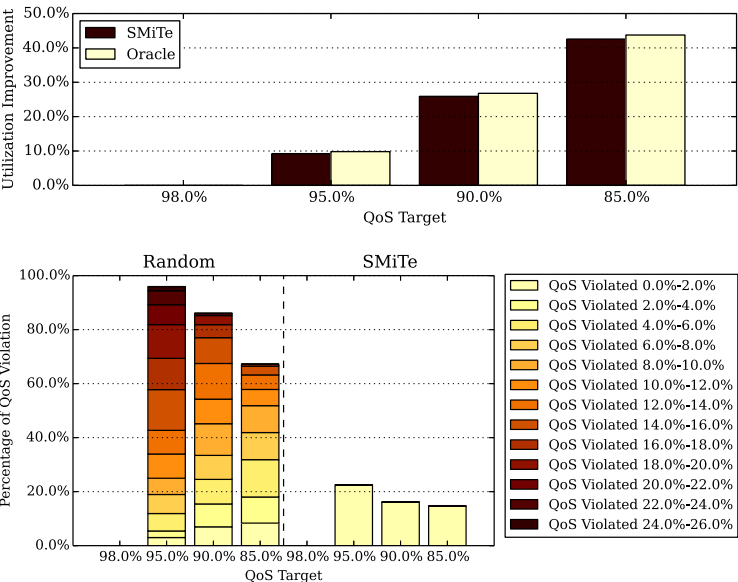
**Decoupled Quantification**

**Direct Interference Measurement**

## SMiTe Methodology Overview

A Ruler

Sensitivity and Contentiousness in Each Sharing Dimension

B Ruler

Sensitivity and Contentiousness in Each Sharing Dimension

Performance Prediction Model

A    B   SMT Co-Location

A    B   CMP Co-Location

## Precise Interference Prediction on Real-System SMT Processors

Measured    PMU Prediction Error    SMiTe Prediction Error

SMT Co-Location

Web-Search: PMU-15.73%, SMiTe-2.32%   Data-Caching: PMU-19.77%, SMiTe-1.67%   Data-Serving: PMU-15.89%, SMiTe-1.85%   Graph-Analytics: PMU-18.40%, SMiTe-1.32%

## Data Center Utilization Improvement

SMiTe
Oracle

QoS Target: 98.0%, 95.0%, 90.0%, 85.0%

**Commodity Processor**

**< 2% Prediction Error**

**42% Utilization Improvement**

Random    SMiTe

QoS Target

QoS Violated 0.0%-2.0%
QoS Violated 2.0%-4.0%
QoS Violated 4.0%-6.0%
QoS Violated 6.0%-8.0%
QoS Violated 8.0%-10.0%
QoS Violated 10.0%-12.0%
QoS Violated 12.0%-14.0%
QoS Violated 14.0%-16.0%
QoS Violated 16.0%-18.0%
QoS Violated 18.0%-20.0%
QoS Violated 20.0%-22.0%
QoS Violated 22.0%-24.0%
QoS Violated 24.0%-26.0%