
Leveraging Hardware Address Sampling

Beyond Data Collection and Attribution

Xu Liu

Department of Computer Science

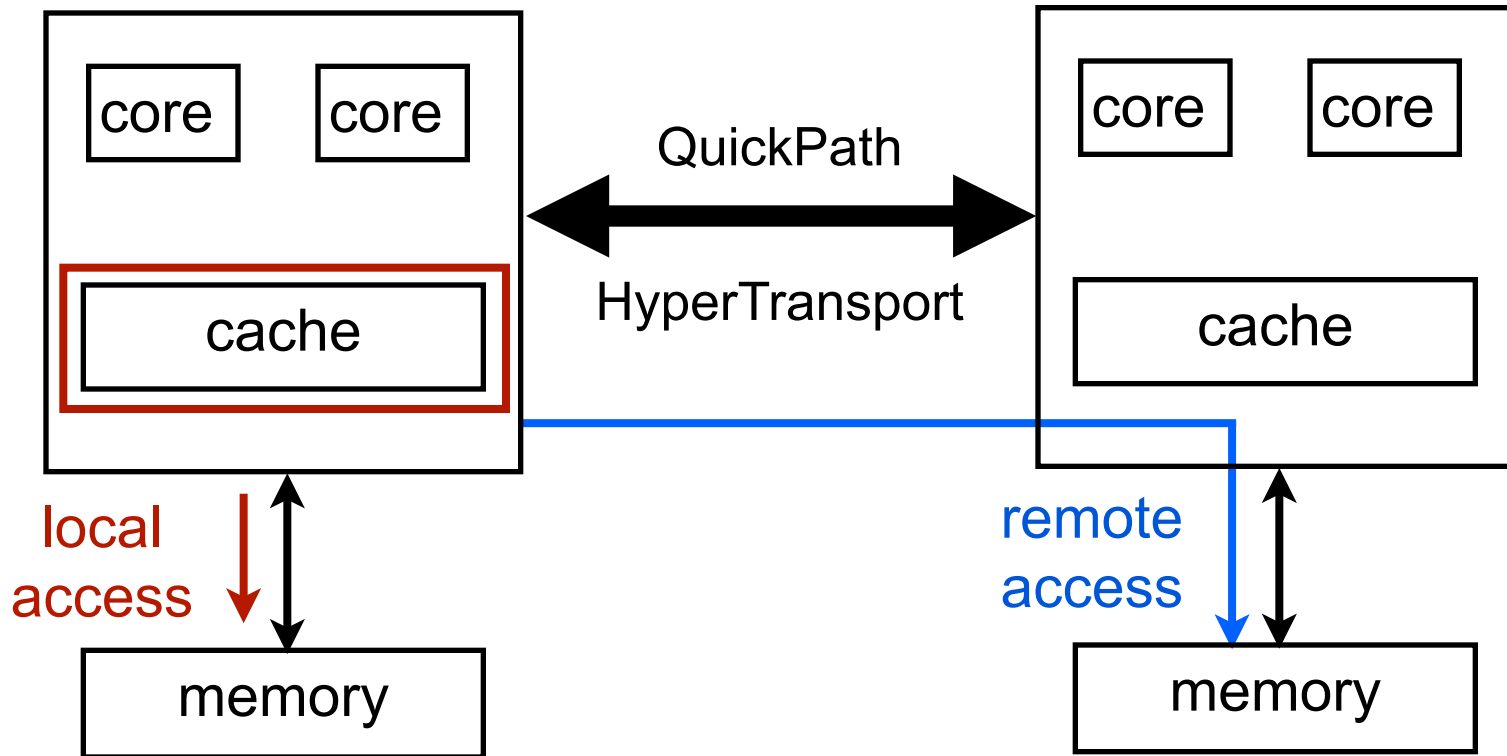
College of William and Mary

xl10@cs.wm.edu



Motivation: Memory is the Bottleneck

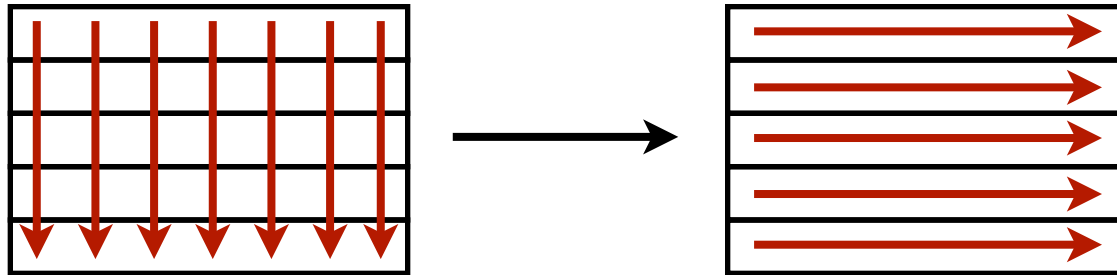
NUMA: Non-Uniform Memory Access





Memory Bottleneck Optimization

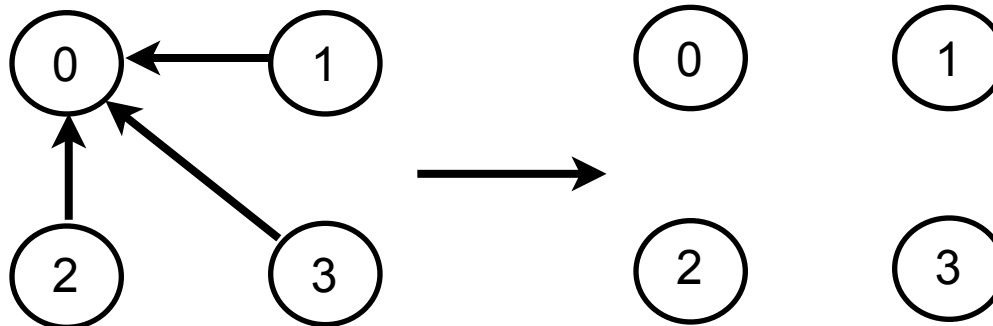
spatial
locality



temporal
locality



NUMA
locality





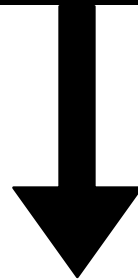
State of the Arts

simulation methods

deep insights

weaknesses:

- 2-5x overhead
- not real machines

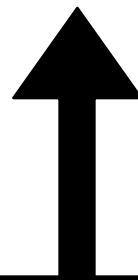


low overhead with deep insights

deep insights with low overhead

measurement methods

low overhead





Hardware Address Sampling

- Features of address sampling
 - necessary features
 - sample memory-related events (memory accesses, NUMA events)
 - capture effective addresses
 - record precise IP of sampled instructions or events
 - optional features
 - record useful metrics: data access latency (in CPU cycle)
 - sample instructions/events not related to memory
- Support in modern processors
 - AMD Opteron 10h and above: instruction-based sampling (IBS)
 - IBM POWER 5 and above: marked event sampling (MRK)
 - Intel Itanium 2: data event address register sampling (DEAR)
 - Intel Pentium 4 and above: precise event based sampling (PEBS)
 - Intel Nehalem and above: PEBS with load latency (PEBS-LL)



Tools Based on Address Sampling

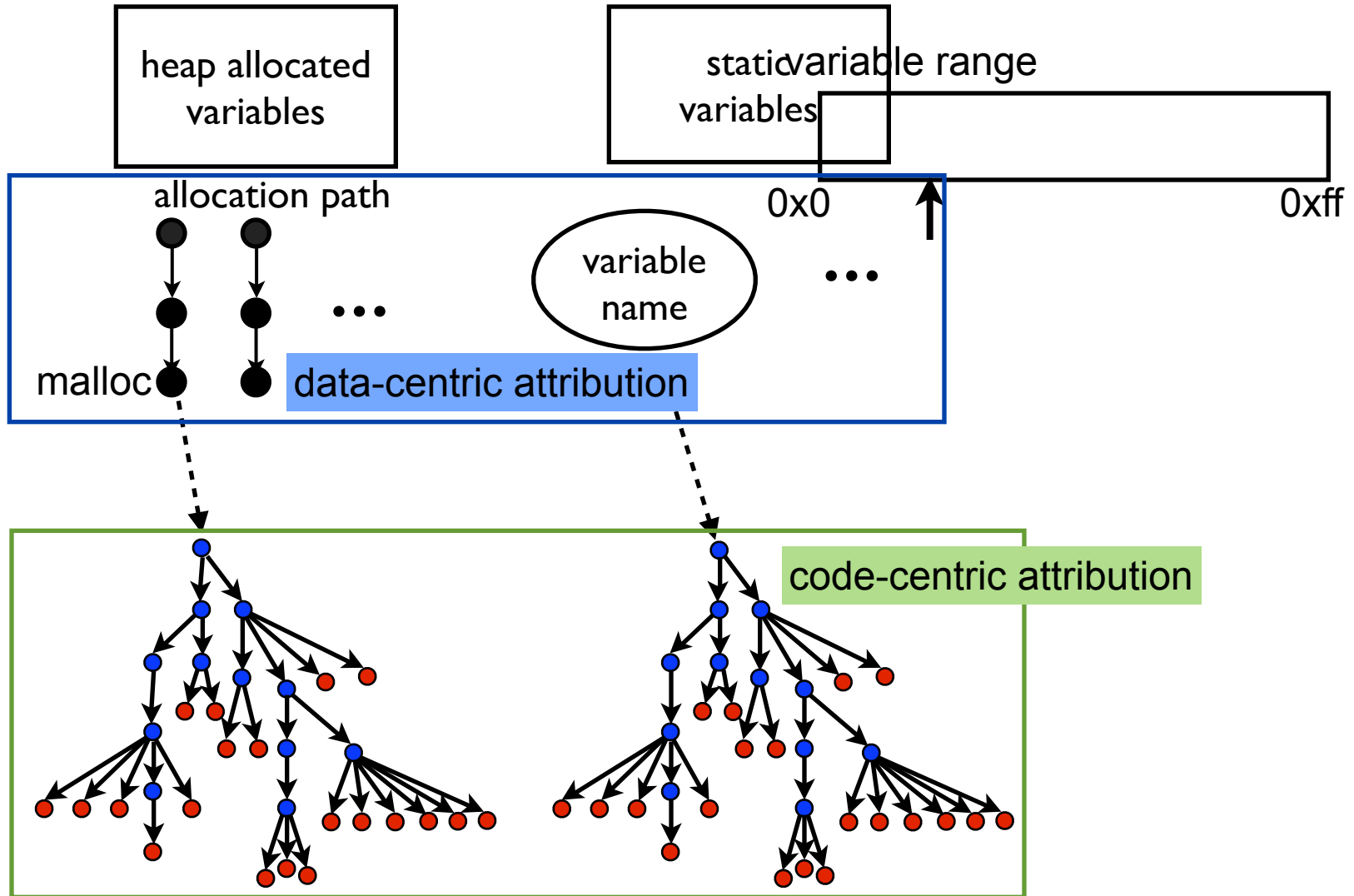
- Measurement methods
 - temporal/spatial locality
 - HPCToolkit, Cache Scope
 - NUMA locality
 - Memphis, MemProf, HPCToolkit
- Features
 - lightweight performance data collection
 - efficient performance data attribution
 - code-centric attribution
 - data-centric attribution

Take HPCToolkit for example

“A Data-centric Profiler for Parallel Programs”. Liu and Mellor-Crummey, SC'13



HPCToolkit: Attributing Samples





LULESH on Platform of 8 NUMA Domains

```
2158 Real_t *x = new Real_t[numNode]; /* coordinates */
2159 Real_t *y = new Real_t[numNode];
2160 Real_t *z = new Real_t[numNode];
2161
Real_t *xd = new Real_t[numNode]; /* velocities */
```

Scope	R_DRAM_ACCESS:Sum (l)
Experiment Aggregate Metrics	9.39e+03 100 %
monitored_heap_data	6.39e+03 68.1%
267: heap_data_allocation	6.39e+03 68.1%
297: monitor_main	6.39e+03 68.1%
479: main	6.39e+03 68.1%
2160: operator new[](unsigned long)	7.21e+02 7.7%
32: operator new(unsigned long)	7.21e+02 7.7%
52: malloc	7.21e+02 7.7%
heap_data accesses	7.21e+02 7.7%

```
for( Index_t lnode=0 ; lnode<8 ; ++lnode )
{
    Index_t gnode = nodelist[ElemId][lnode];
    ElemPos[X_Dir][lnode] = x[gnode];
    ElemPos[Y_Dir][lnode] = y[gnode];
    ElemPos[Z_Dir][lnode] = z[gnode];
}
```

heap data:68%
remote accesses

z accounts for
7.7% remote accesses

remote accesses

z is allocated in a
NUMA domain but
accessed by others

interleave pages of z
across NUMA nodes
13% improvement in
running time

call paths for
accesses

call site of allocation

allocation call path



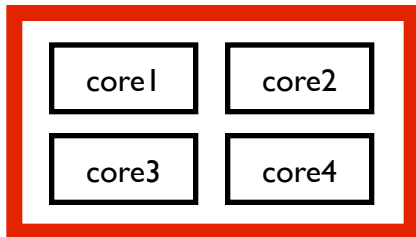
Existing Measurement is Inadequate

- Data collection + attribution \neq optimal optimization
 - know problematic data objects but not know why
 - need more insights for optimization guidance
 - challenges in data analysis
 - not monitoring continuous memory accesses
- Approaches: data analysis for detailed optimization guidance
 - NUMA locality
 - offline optimization (PPoPP'14)
 - online optimization
 - cache locality
 - array regrouping (PACT'14)
 - structure splitting
 - locality optimization between SMT threads
 - scalability of memory accesses

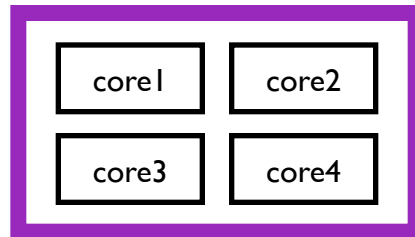


Interleaved Allocation is NOT Always Best

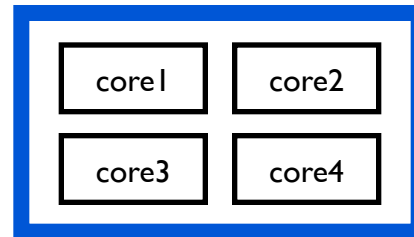
domain 1



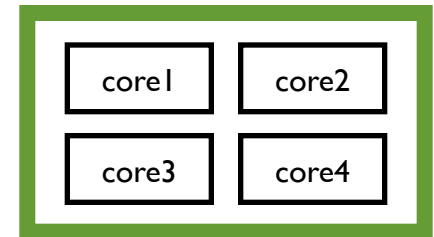
domain 2



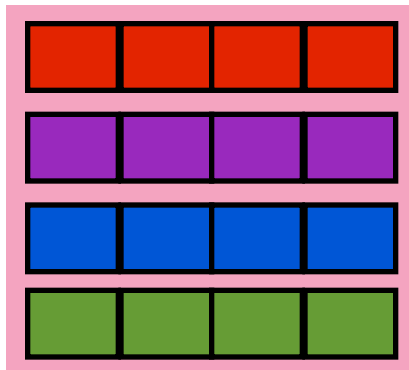
domain 3



domain 4



allocation 1



allocation 2



allocation 3



centralized allocation: poor

interleaved allocation: sub-optimal

co-locate data with computation: optimal

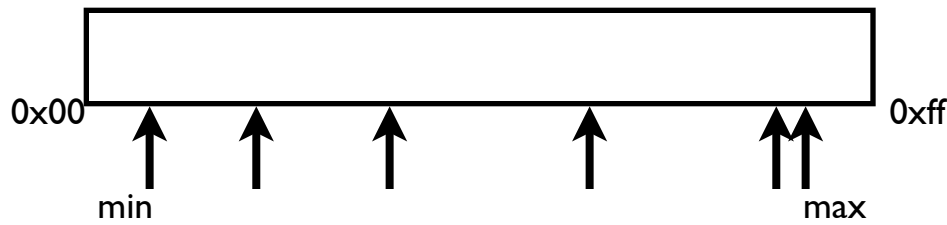
Goal: identify the best data distribution for a program



Memory Access Pattern Analysis

- Online data collection

array A



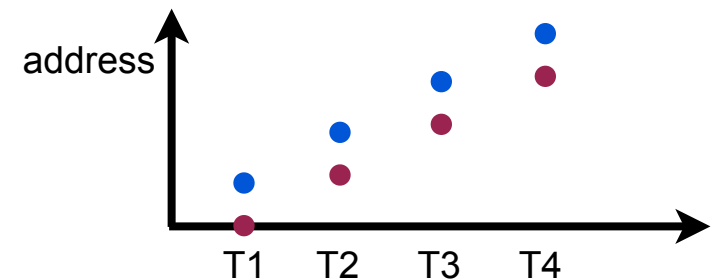
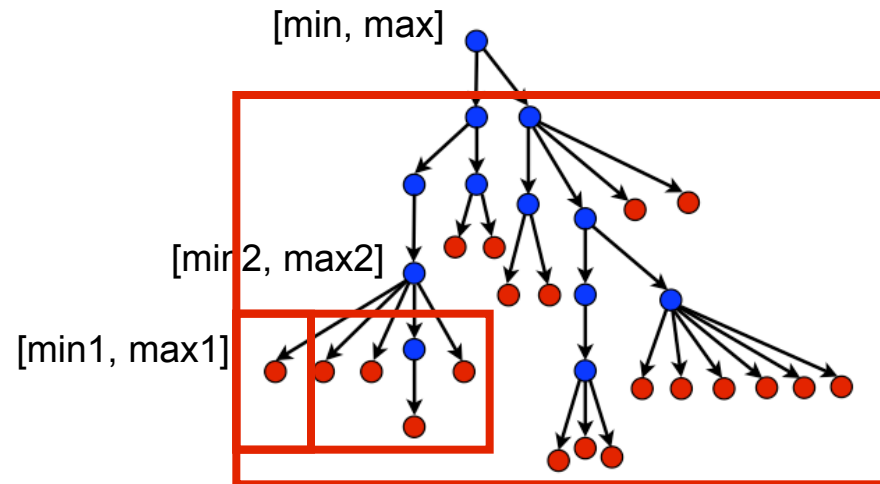
[min, max] per sampled memory access

balanced allocation + maximum locality

array A



allocate A blockwise to different domains

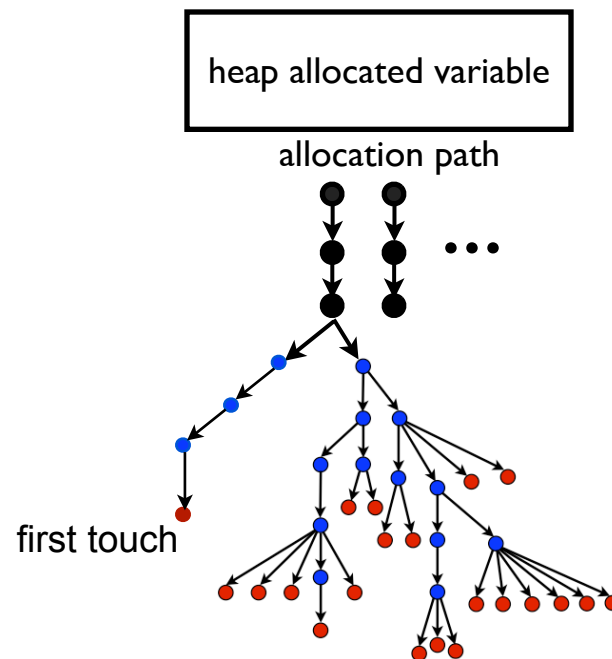
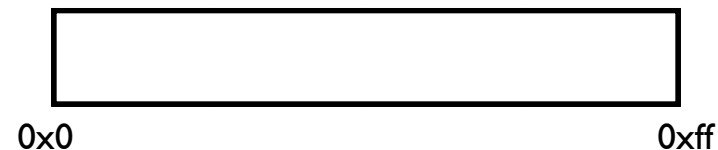




Pinpointing First Touch

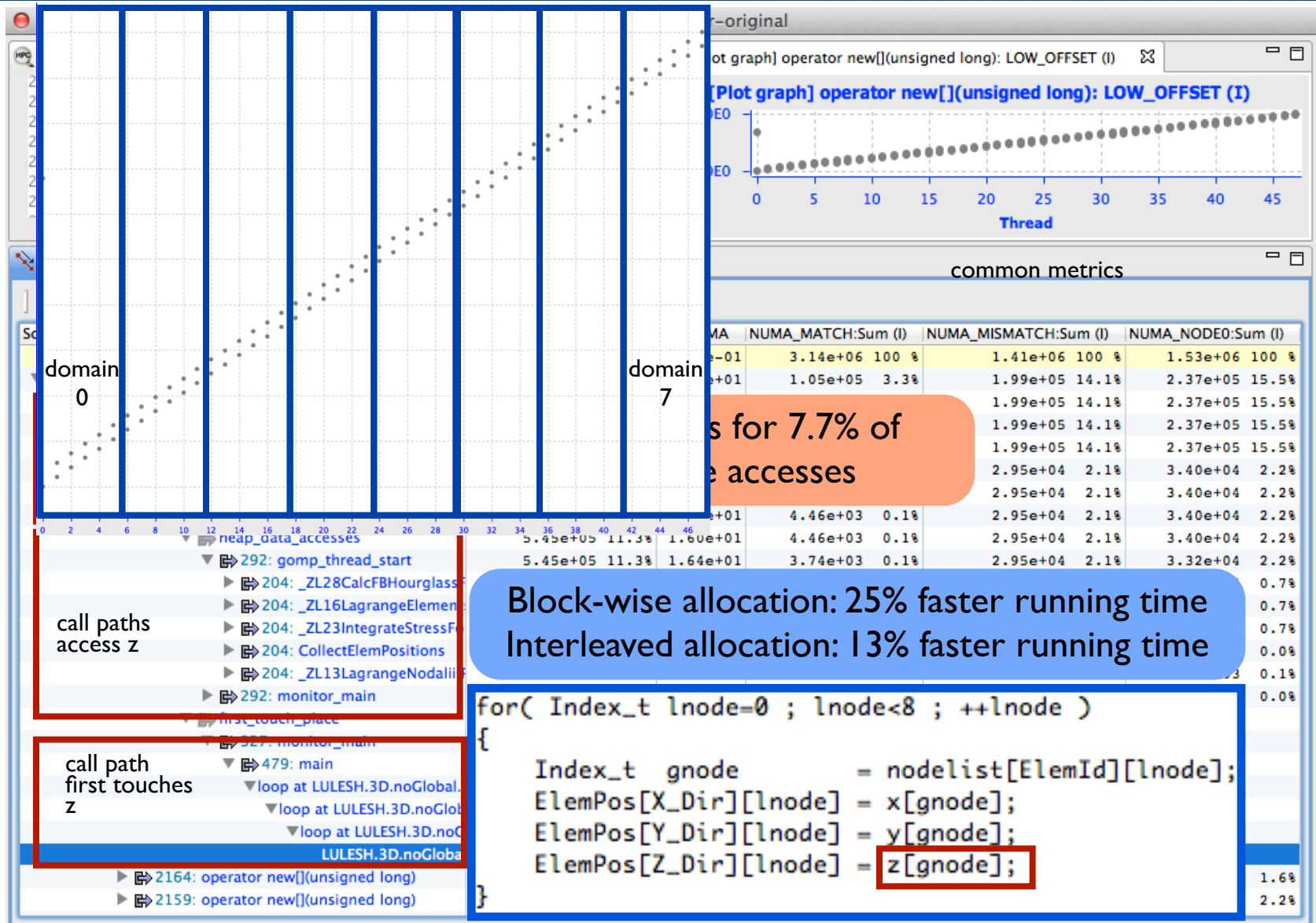
- Linux “first touch” policy
 - memory allocation at first touch
 - if T1 first touches the whole range of A
 - if threads touch different segments of A

array A





LULESH on Platform of 8 NUMA Domains





Experiments: Architectures & Applications

Architectures

Sampling mechanisms		Processors	Threads
Instruction-based sampling	IBS	AMD Magny-Cours	48
Marked event sampling	MRK	IBM POWER 7	128
Precise event-based sampling	PEBS	Intel Xeon Harpertown	8
Data event address registers	DEAR	Intel Itanium 2	8
PEBS with load latency	PEBS-LL	Intel Ivy Bridge	8

Benchmarks

LLNL	LANL	Rodinia	PARSEC	SNL
AMG2006	Sweep3D	Streamcluster	Blackscholes	S3D
LULESH		NW		
Sphot				
UMT2013			optimized benchmarks	



Optimization Results

Programs	Optimization	Improvement for execution time
AMG2006	NUMA locality	51% for the solver
Sweep3D	spatial locality	15%
LULESH	spatial+NUMA locality	25%
Streamcluster	NUMA locality	28%
NW	NUMA locality	53%
UMT2013	NUMA locality	7%



Measurement Overhead

Code- & data-centric analysis on POWER7 and Opteron

Benchmark	Configuration	Overhead
AMG2006	4 MPI * 128 threads	604s (+9.6%)
Sweep3D	48 MPI	90s (+2.3%)
LULESH	48 threads	19s (+12%)
Streamcluster	128 threads	27s (+8.0%)
NW	128 threads	80s (+3.9%)

NUMA analysis: code-, data-, and address-centric analysis + first touch

Methods	LULESH	AMG2006	Blacksholes
IBS	295s (+24%)	89 (+37%)	192s (+6%)
MRK	93s (+5%)	27s (+7%)	132s (+4%)
PEBS	65s (+45%)	96s (+52%)	82s (+25%)
DEAR	90s (+7%)	120s (+12%)	73s (+4%)
PEBS-LL	35s (+6%)	57s (+8%)	67s (+3%)



Conclusions and Future Work

- Hardware address sampling
 - widely supported in modern architectures
 - powerful in monitoring memory behaviors
 - currently in early stage of studies
 - focusing on data collection and attribution
- Potentials of hardware address sampling
 - provide deeper insights than traditional performance counters
 - require novel analysis methods to expose performance insights
- Future work
 - integrating address sampling into Charm++ runtime for online optimization

