Leveraging Hardware Address Sampling

Beyond Data Collection and Attribution

Xu Liu

Department of Computer Science College of William and Mary <u>xl10@cs.wm.edu</u>



NUMA: Non-Uniform Memory Access







- Features of address sampling
 - sample memory-related events (memory accesses, NUMA events)
 - capture effective addresses
 - record precise IP of sampled instructions or events
- 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)
- Efficient memory measurement (SC'13)
 - code-centric analysis
 - data-centric analysis



- Code-centric attribution
 - problematic code sections
 - instruction, loop, function

l:#pragma omp parallel for num_threads(4) 2: for (i = 0: i < n: i++) {			
3: for(j = 0; j < n; j++) { 4: for(k = 0: k < n: k++) {			
5: $A[i, j, k] = A[i, j, k,] + B[j, i, k] + C[k, j, i];$			
6: } 7: } 8:}			



Combining code-centric and data-centric attribution provides additional insight



Attributing Samples



Aggregating Profiles



LULESH on Platform of 8 NUMA Domains





- Data collection + attribution ≠ optimal optimization
 - know problematic data objects but not know why
 - need more insights for optimization guidance
- Challenges for address sampling
 - very sparse memory access samples
 - not monitoring continuous memory accesses
- Opportunities for address sampling
 - effective addresses: analyze memory access patterns
 - data sources: understand where inefficiencies come from
 - latency: derive new latency metrics to quantify inefficiencies.



- Published work
 - analyzing NUMA bottlenecks (PPoPP'14)
 - guiding array regrouping for better locality (PACT'14)
 - identifying memory scaling issues (SC'15)

Interleaved Allocation is NOT Always Best



Goal: identify the best data distribution for a program



Memory Access Pattern Analysis



·LULESH on Platform of 8 NUMA Domains





- Published work
 - analyzing NUMA bottlenecks (PPoPP'14)
 - guiding array regrouping for better locality (PACT'14)
 - identifying memory scaling issues (SC'15)



Array Regrouping





Workflow





latency

0.0% 2.4% 4.9%	<pre>I = (float *)malloc(size_I * sizeof(float)); J = (float *)malloc(size_I * sizeof(float)); c = (float *)malloc(sizeof(float)* size_I) ;</pre>
0.0% 0.0% 0.1% 1.5%	<pre>iN = (int *)malloc(sizeof(unsigned int*) * rows) ; iS = (int *)malloc(sizeof(unsigned int*) * rows) ; jW = (int *)malloc(sizeof(unsigned int*) * cols) ; jE = (int *)malloc(sizeof(unsigned int*) * cols) ;</pre>
2.7% 3.8% 5.1% 5.6%	<pre>dN = (float *)malloc(sizeof(float)* size_I); dS = (float *)malloc(sizeof(float)* size_I); dW = (float *)malloc(sizeof(float)* size_I); dE = (float *)malloc(sizeof(float)* size_I);</pre>

SRAD from Rodinia



Latency-based Array Affinity

loops in their calling contexts



An Example





Access Pattern Analysis



dN, dS, dE, dW have the same pattern



only regroup dN, dS, dW, dE $\longrightarrow \uparrow 1.47x$ regroup dN, dS, dW, dE and c $\longrightarrow \uparrow 1.89x$



- Published work
 - analyzing NUMA bottlenecks (PPoPP'14)
 - guiding array regrouping for better locality (PACT'14)
 - identifying memory scaling issues (SC'15)

Scaling Losses in Memory Hierarchies

- Memory contentions hurt scalability: cache/bandwidth contention
 - which data objects contribute to the most scaling losses
 - which memory layers incur the most scaling losses
- Methods
 - decompose latency according to data objects and memory layers
 - data-centric analysis with data source information
 - differential analysis supported by HPCToolkit
 - compare profiles between different runs

more details in SC'15



- Hardware address sampling
 - widely supported in modern architectures
 - powerful in monitoring memory behaviors
 - more analysis of the samples provides more performance insights
- On-going work
 - structure splitting
 - locality optimization between SMT threads
 - cache line false sharing
 - automatic page migration for NUMA architectures
- Future directions of address sampling
 - comparing different address sampling mechanisms
 - analyzing new performance issues
 - heterogeneous memory: 3D stack memory





Backup Slides

-UMT2013 on Quad-socket POWER7 Node

👰 ZoneData_mod.F90 🔀 🗹 [Plot graph] malloc: LOW_OFFSET (I)			
90 allocate(self % A_fp(Size% ndim,self% nCFaces,se allocate(self % A_ez(Size% ndim,self% nCFaces,se allocate(self % Connect(3,self% nCFaces,self% nC allocate(self % STotal(Size% ngr,self% nCorner)) allocate(self % STime(Size% ngr,self% nCorner,Size) if (Circ@rdim and then	<pre>lf% nCorner)) lf% nCorner)) orner)) ze% nangSN))</pre> sample off-chip accesses		
Calling Context View			
] 1 + i i i i i i i i i i i i i i i i i i			
Scope	NUMA_MISMATCH:Sum (I) 🔻		
Experiment Aggregate Metrics	2.45e+04 100 %		
▶monitored_unknown_data	1.30e+04 53.1% self%S lime		
▼monitored_heap_data	1.15e+04 46.9%		
▼ 🖶 266: heap_data_allocation	1.15e+04 46.9%		
🔻 🖶 296: monitor_main	1.10e+04 44.9%		
▼ 🖶 479: main	1.10e+04 44.9%		
▼inlined from SuOlsonTest.cc: 67	1.10e+04 44.8% 8.2% of remote accesses		
🔻 🖨 167: initialize(Geometry::MeshBase&, Teton <geometry< td=""><td>ry::MeshBase>&, 1.05e+04 42.9%</td></geometry<>	ry::MeshBase>&, 1.05e+04 42.9%		
S14: Teton <geometry::meshbase>::linkKull(Geometry::MeshBase>::l</geometry::meshbase>	etry::MeshBase&, 1.05e+04 42.9%		
▼loop at Teton.cc: 1250	4.45e+03 18.28 allocated in one domain		
▼ 🖨 1328: setzone	4.45e+03 18.2% allocated in one domain		
▼ 🔂 40: zonedata_ctor	4.45e+03 18.28 accessed by everyone		
▼loop at ZoneData_mod F90: 86	4.45e+03 18.28 accessed by everyone		
▼loop at ZoneData_mod.F90: 87	4.45e+03 18.2%		
▼loop at ZoneData_nod.F90: 88	4.45e+03 18.2%		
▼loop at ZoneData_mod.F90: 89	4.45e+03 18.2%		
▼loop at ZoneDati_mod.F90: 90	0 4.45e+03 18.2%		
▼loop at ZoneDita_mod.F90	: 91 4.45e+03 18.2%		
▼loop at ZonePata_mod.F	90:92 4.45e+03 18.2%		
▼loop at ZoneData_mo	d.F90:93 4.45e+03 18.2%		
▼loop at ZeneData	mod.F90: 94 4.45e+03 18.2%		
▼ B⇒94: malloc	4.45e+03 18.2%		
V 🖶 heap_da	ta_accesses 4.45e+03 18.2%		
▼ 🛱 291: 1	ThdCode 4.45e+03 18.2%		
▼ ⇒ _x	smp_DynamicCh 4.45e+03 18.2%		
▼ ₿	snflwxyz\$\$OL\$\$ 4.45e+03 18.2%		

• Optimize *self%STime* for UMT2013

a



self%STime's address space

