Program Interaction on Shared Cache

Theory and Applications

Chen Ding

Professor Department of Computer Science University of Rochester

Illustration: bottlenecks of SPEC2000 on Itanium1



Anant Aggarwal, MIT 6.975, 2007



Chen Ding, DragonStar lecture, ICT 2008

"Nothing travels faster than the speed of light ..." Douglas Adams





Version 3.0

May 2003

Jason F. Cantin Department of Electrical and Computer Engineering 1415 Engineering Drive University of Wisconsin-Madison Madison, WI 53706-1691 jcantin@ece.wisc.edu http://www.jfred.org

> Mark D. Hill Department of Computer Science 1210 West Dayton Street University of Wisconsin-Madison Madison, WI 53706-1685 markhill@cs.wisc.edu http://www.cs.wisc.edu/~markhill

http://www.cs.wisc.edu/multifacet/misc/spec2000cache-data

Chen Ding, University of Rochester, PMAM 2014

Read Edit View history Search Article Discussion CPU cache WikipediA From Wikipedia, the free encyclopedia 0.1 Direct 0.01 0.001 miss rate 1e-04 1e-05

1K Chen Ding, University of Rochester, PMAM 2014

4K

16K

cache size

1e-06

64K 256K 1M Inf

Madison Itanium 2

2002

http://en.wikipedia.org/wiki/File:Cache,missrate.png

D-cach	e misses/inst:	1,197,717,058	,456 data refs	(0.34534/in	st);				
782,173,506,477 D-cache 64-Byte block accesses (0.22949/inst)									
Size	Direct	2-way LRU	4-way LRU	8-way LRU	Full LRU				
1KB 2KB 4KB 8KB 16KB 32KB 64KB 128KB 256KB	0.0890418 0.0651636 0.0480381 0.0362358 0.0277699 0.0223409 0.0189635 0.0158796 0.0158840	0.0762018 0.0533596 0.0386862 0.0290652 0.0227735 0.0190920 0.0166430 0.0147737 0.0131826	0.0699370 0.0486152 0.0353534 0.0264135 0.0264135 0.0181803 0.0181803 0.0161909 0.0144648 0.013454-	0.0657938 0.0462573 0.0337222 0.0254564 0.0204821 0.0179048 0.0160494 0.0143748 0.0130774	0.0652996 0.0453232 0.025938 0.0245702 0.0196992 0.0175964 0.0159076 0.0142985 0.0142985				
512KB 1MB	0.0119997 0.0096151	0.0115157	0.0114489 0.0092640	0.0114018 0.0093510	0.0113629 0.0093828				

Compulsory: 0.0000150365--

Be Si	nchmark m Time:	s: 12 1463.66	days,	4.007 year	s
Fi	le crea	ted 5/23	/2003.		

Chen Ding, University of Rochester, PMAM 2014

A Metric and A Tool Box

• Reuse distance

- independent of coding styles, memory allocation, or hardware
- possible to correlate between different runs
- pattern analysis



Chen Ding, University of Rochester, PMAM 2014

Program Locality

Reuse Distance

The SLO Tool by Beyls and D'Hollander



• An example: 173.APPLU from SPEC 2K



Measuring Reuse Distance

time: 1 2 3 4 5 6 7 8 9 10 11 12 access: $\mathbf{d} \mathbf{a} \mathbf{c} \mathbf{b} \mathbf{c} \mathbf{c} \mathbf{g} \mathbf{e} \mathbf{f} \mathbf{a} \mathbf{f} \mathbf{b}$ distance: $| \mathbf{\leftarrow} 5 \text{ distinct accesses} \rightarrow |$ (a) an example access sequence the reuse distance between two b's is 5

• Naive counting, O(N) time per access, O(N) space

- N is the number of memory accesses
- M is the number of distinct data elements

11

- Too costly
 - N is up to 120 billion, M 25 million

Chen Ding, DragonStar lecture, ICT 2008

JE JES

9

Reuse Distance Measurement

Measurement algorithms since 1970	Time	Space
Naive counting	O(N ²)	O(N)
Trace as a stack [IBM'70]	O(NM)	O(M)
Trace as a vector [IBM'75, Illinois'02]	O(NlogN)	O(N)
Trace as a tree [LBNL'81], splay tree [Michigan'93], interval tree [Illinois'02]	O(NlogM)	O(M)
Fixed cache sizes [Winsconsin'91]	O(N)	O(C)
Approximation tree [Rochester'03]	O(NloglogM)	O(logM)
Approx. using time [Rochester'07]	O(N)	O(1)

N is the length of the trace. M is the size of data. C is the size of cache.



Program lo	ocality analysis using reuse distance	
Full Text:	Pdf <u>Buy this Article</u>	2009 Article
Authors: <u>Y</u> <u>X</u>	<u>'utao Zhong</u> George Mason University, Fairfax, VA (<u>ipeng Shen</u> The College of William and Mary, Williamsburg, VA <u>Chen Ding</u> University of Rochester, Rochester, NY	Research Refereed
Published in:	Journal ACM Transactions on Programming Languages and Systems (TOPLAS) <u>TOPLAS Homepage archive</u> Volume 31 Issue 6, August 2009 <u>ACM</u> New York, NY, USA <u>table of contents</u> dol> <u>10.1145/1552309.1552310</u>	Bibliometrics Downloads (6 Weeks): 15 Downloads (12 Months): 267 Citation Count: 3

benchmarks	length	data size	unmodifed	FP alg	FP alg	RD alg	RD alg	LF alg	LF alg
		(64B lines)	time (sec)	time	cost (X)	time	cost (X)	time	cost (X)
176.gcc	1.10E+10	3.99E+06	85.1	345	4.1	2,392	28.1	5,489	65
181.mcf	1.88E+10	2.52E+06	398	1,126	2.8	10,523	26.4	121,818	306
164.gzip	2.00E+10	1.41E+06	150	501	3.3	5,823	38.8	44,379	296
252.eon	2.51E+10	1.54E+04	77.4	503	6.5	5,950	76.9		
256.bzip2	3.20E+10	1.47E+06	173	726	4.2	7,795	45.1	36,428	211
175.vpr	3.56E+10	5.08E+04	210	964	4.6	13,654	65.0	51,867	247
186.crafty	5.31E+10	3.20E+04	75.5	1,653	21.9	18,841	249.5	117,473	1,556
300.twolf	1.08E+11	9.47E+04	368	2,979	8.1	27,765	75.4	155,793	423
197.parser	1.22E+11	6.52E+05	230	3,122	13.6	35,562	154.6	106,198	462
11 2K INT avg	4.73E+10	1.14E+06	196	1,324	8	14,256	84	79,931	446
179.art	1,20E+10	5.93E+04	591	734	1.2	4,032	6.8	36,926	62
100	705.10	7 005 . 05	100	000	0.0	10.051	110.0	100 001	1 000
17 hill	ion								
47 DIII	ion								
access	205	7 3	m16s				3	h57m	

Analysis Speed

Chen Ding, University of Rochester, PMAM 2014



Active Sharing (now)

The End of Cache Monopoly

Multicore

- desktop, cloud, and handheld
- Multicore cache
 - a mixture of private/shared caches
 - $\cdot\,$ Intel Nehalem 256KB private L2, 4MB to 8MB shared L3
 - $\cdot\,$ IBM Power 7 256KB private L2, 32MB shared ERAM L3
 - \cdot ERAM to appear on Intel processors

New problems

- available cache resource is variable • not the full size, not constant size
- not just performance but also stability
- not just parallel program but also sequential program







 Private cache locality
P(capacity miss by me) = P(my reuse distance >= cache size)
Shared cache locality

Footprint

Shared cache locality
P(capacity miss by me) =
P(my reuse distance + peer footprint >= cache size)

· 3 length-2 windows: "ab", "bb", "bb"

the average fp(2) = (2 + 1 + 1)/3 = 4/3

Chen Ding, University of Rochester, PMAM 2014

Footprint Locality

[Ding, Xiang, et al. PPOPP 2008/11, PACT 11, ASPLOS 13]

rsity of Rochester, PMAM 2014

all-window 'footprint' footprint



• Example: "abbb"

footprints 2, 1, 1

Footprint Measurement 1972 - 2007

- Working set
- · limit value in an infinitely long trace [Denning & Schwartz 1972] · Direct counting
 - · single window size [Thiebaut & Stone TOCS'87] • seminal paper on footprints in shared cache
 - same starting point [Agarwal & Hennessy TOCS'88]
- Statistical approximation
 - [Denning & Schwartz 1972; Suh et al. ICS'01; Berg & Hagersten PASS'04; Chandra et al. HPCA'05; Shen et al. POPL'07]

组合性

- · level of precision couldn't be directly checked
- No precise definition/solution for all windows
 - can't be measured for real
 - can't know the accuracy of an estimate

Chen Ding, University of Rochester, PMAM 2014

Footprint Measurement 2008 - 2013

- Footprint distribution
 - all-window enumeration [Ding/ Chilimbi PPOPP 2008]
 - max/min/median/percentiles
 - trace compression [Xiang+ PPOPP 11] • 70X speedup
 - 4 hours per program
- Average footprint [Xiang+ PACT 11]
- Xiang formula
 - 22 minutes per program
- Footprint Sampling [Xiang+ ASPLOS 13] shadow profiling
 - · 0.5%

Chen Ding, University of Rochester, PMAM 2014



- Xiaoya Xiang
- HUST BS 2005 • ICT MS 2008
- Rochester PhD (expected)
- Twitter 2013

26

Footprint to Miss Rate Conversion · Ax



25

Conversion Formulas

The Xiang formula for average footprint [PACT'11]

- rt: reuse time
- m: data size
- n: trace length

$$fp(x) \approx m - \sum_{k=x+1}^{n-1} (k-x)P(rt = k)$$
$$mr(c) = mr(fp(x)) = \underbrace{\frac{fp(x+\Delta x) - fp(x)}{\Delta x}}_{P(rd = c)}$$
$$P(rd = c) = mr(c-1) - mr(c)$$

Composition + Conversion



28

Reality Check

· 20 SPEC 2006 programs

- 190 different pair runs
- Modeling
 - per program footprint
 - composition
 - a few hours
 - prediction for all cache sizes
- · Exhaustive parallel testing
 - 190 pair runs
 - 380 hw counter reads (OFFCORE.DATA_IN, 8MB 16-way L3)
 - ~9 days total CPU time



Chen Ding, University of Rochester, PMAM 2014



Chen Ding, University of Rochester, PMAM 2014

31

Co-run interference of libquantum; high miss ratio, zero sensitivity; measured miss ratio 17.82% to 17.89%, predicted 17.94% to 17.94%



Co-run interference of gamess; low miss ratio, high sensitivity measured miss ratio 0.0002% to 0.04%, predicted 0.000013% to 0.03%



Denning's Law of Locality

What's the relation between reuse frequency and footprint?

> abc ... abc . . . aaa ... bbb . . .

Limit value [Denning and Schwartz, CACM 1972] Time space [Denning and Slutz, CACM 1978] All program traces [Rochester, ASPLOS 2013]

Chen Ding, University of Rochester, PMAM 2014

An Old Open Question

How quickly can we measure the miss rate for all cache size?

3000+ *Cache Sizes* In the analysis, the footprint and reuse distance numbers are bin-ed using logarithmic ranges as follows. For each power-of-two range, we sub-divide it into 256 equal-size increments. As a result, we can predict the miss ratio not just for power-of-two cache sizes, but 3073 cache sizes between 16KB and 64MB.

Xiang et al. ASPLOS 13 (Tongxin Bai's tool)

Chen Ding, University of Rochester, PMAM 2014

HOTL: a higher order theory of locality



Chen Ding, University of Rochester, PMAM 2014



An Old Open Question

What's the relation between miss rate and cache pressure? Does a higher miss rate mean higher pressure?

Chen Ding, University of Rochester, PMAM 2014

Miss Ratio vs Pressure, 32KB Cache



Miss Ratio vs Pressure, 4MB Cache



An Old Open Question

Does LRU cache produce optimal partition? [Thiebuat and Stone, 1992]

The second type of sharing happens between the instruction and the data of a program. Stone et al. [1992] investigated whether LRU produces the optimal allocation. Assuming that the miss rate functions for instruction and data are continuous and differentiable, the optimal allocation happens at the points "when miss-rate derivatives are equal" [Thiébaut and Stone, 1992]. The miss rate functions, one for instruction and one for data, were modeled instead of measured. The authors showed that LRU is not optimal, but left even a cuestion whether there

is a bound on how close LRU allocation is to optimodel in Chapter 4 can be used to compute the c answer the open question for any group of program



Chen Ding, University of Rochester, PMAM 2014

Chen Ding, University of Rochester, PMAM 2014

Collaborative Rationing

An Old Open Question

Is there a machine independent way to compare

program behavior in shared cache? How do

programs in different domains differ?

Thread 1	a	b	с	a	b	С	a	b	С
Hint Bit	0	1	0	1	0	1	0	1	0
Access Bit	1	0	1	0	1	0	1	0	1
Misses	М	М	М		М		м		М
	i i i								
Thread 2	x	У	z	x	У	z	x	У	z
Thread 2 Hint Bit	<u>х</u> 0	у 1	z 0	x 1	y o	z 1	х 0	у 1	z 0
Thread 2 Hint Bit Access Bit	x 0 1	у 1 0	z 0 1	x 1 0	у 0 1	z 1 0	x 0 1	у 1 0	z 0 1
Thread 2 Hint Bit Access Bit Misses	х 0 1 М	у 1 0 М	z 0 1 M	x 1 0	у 0 1 м	z 1 0	х 0 1 М	у 1 0	z 0 1 M

Two threads, each accessing three elements and using two-element cache. Best per thread and overall cache utilization --- 50% miss rate for each program.

Chen Ding, University of Rochester, PMAM 2014 45





Jacob Brock and Raj Parihar



Maximal cache performance? Answer:

Miss rate in all cache sizes? Answer: LRU-MRU (Gu) distance [Gu et al. ISMM 2012, Rochester Dissertation 2013]

Rochester, New York

Chen Ding, University of Rochester, PMAM 2014

2013

ferret dedup 50MB 50MB measured, max 4t predicted, max 4t 40MB 40MB measured, min 4t predicted, min 4t 30MB 30MB 20MB 20MB 10MB 10MB 1e+00 1e+03 1e+06 1e+09 1e+00 1e+031e+06 1e+09

All thread-group locality prediction. Min/max locality in all 70 four-thread groups for two PARSEC programs with 8 asymmetric threads.

Chen Ding, University of Rochester, PMAM 2014

On-going Studies

Shared Footprint Analysis

with Hao Luo and Pengcheng Li

Recent Developments

- Competitiveness, politeness, sensitivity
 - Jiang et al. [TPDS'11, HiPEAC'10]
- Intensity and sensitivity
 - Zhuravlev et al. [ASPLOS'10]
- Niceness, pressure and sensitivity
 - Mars et al. [CGO'12, Micro'12]
- Interference of cache
 - composable models [Stone+ TOCS'87/TOC'92; Suh+ ICS'01; Chandra+ HPCA'05; Xiang+ PPOPP'11/PACT'11/ASPLOS'13]
 - threaded code [Ding/Chilimbi MSR'09, Jiang+ CC'10/TPDS'12, Schuff+ PACT'10, Wu/Yeung PACT'11/ISCA'13]
- Interference model of execution time/speed
 - bubble-up [Mars+ Micro'12, ISCA'13]
 - QoS-aware scheduling [Delimitrou/Kozyrakis ASPLOS'13]

Chen Ding, University of Rochester, PMAM 2014

Recent Developments [cont'd]

Optimization

Bin Bao

Advisor: Chen Ding

DEPARTMENT OF COMPUTER SCIENCE

- · Parallel reuse distance measurement
 - cluster [OSU, IPDPS 2012]
 - GPU [ICT and NCSU, IPDPS 2012]
 - sampling
 - footprint shadow sampling [Rochester, ASPLOS 2013]
 - multicore reuse distance [Purdue, PACT 2010]
 - reuse distance sampling [Chang & Zhong, PACT 2008]
- Reuse distance in threaded code
 - multicore reuse distance [Purdue, PACT 2010]
 - · CRD/PRD scaling [Maryland, ISCA 2013, to appear]

Chen Ding, University of Rochester, PMAM 2014

51

Recent Developments (cont'd)

- Asymptotic locality effect in parallel algorithms
 - · Leslie Valiant, PACT 2011 keynote
 - Guy Blelloch et al. CMU, MIT, Intel Labs Pittsburgh [MSPC 2013]
 - Morris Herlihy and student, [PPOPP 2014]
- Shared footprint [Rochester, WODA 2013]
- Static reuse distance analysis in Matlab [Indiana, ICS 2010]
- Static footprint analysis [Rochester, CGO 2013]
 - peer-aware program optimization [Bao, dissertation'13]
- Collaborative caching
- practical uses [UT, Ghent, Google etc]
- optimal collaborative LRU cache [Gu, ĪSMM'11/12/13, dissertation'13]

Chen Ding, University of Rochester, PMAM 2014

52

50

Summary

- Program interaction in multicore
 - data sharing in threaded code
 - cache and memory bandwidth sharing by all programs

Locality theory

- working set, footprint, shared footprint
- metrics composition and conversion
- higher order theory of cache locality (HOTL)

• Recent research

- locality in parallel algorithms
- peer-aware program optimization
- sharing conscious task scheduling
- collaborative caching

Peer-Aware Program