













W UNIVERSITY of WASHINGT	UNIVERSITY of WASHINGTON L17: Caches III CSE351, Winter 2017									
Notatior	Notation Changes									
 We are understand Please question 	using differer be aware of thi ons or are watc	nt variable is when you hing videos	e names from I look at past ex	previous am						
	Variable	This Quarter	Previous Quarters							
	Block size	К	В							
	Cache size	С								
	Associativity	Ν	E							
	Address width	Α	m							
	Tag field width	Т	t							
	Index field width I k, s									
	Offset field width 0 n, b									
	Number of Sets	S	S							
				8						















W UNIVERSITY of	WASHINGTON L17: Caches III CSE351, W	inter 2017
Type	es of Cache Misses: 3 C's!	
• 0	ccurs on first access to a block	
🔹 Con	flict miss	
 Color D 	onflict misses occur when the cache is large enough, but multiple data bjects all map to the same slot e.g., referencing blocks 0, 8, 0, 8, could miss every time irect-mapped caches have more conflict misses than	
IN	-way set-associative (where $N > 1$)	
🔹 Cap	acity miss	
 O is 	ccurs when the set of active cache blocks (the <i>working set</i>) larger than the cache (just won't fit, even if cache was <i>fully- ssociative</i>)	
= N	ote: Fully-associative only has Compulsory and Capacity misses	
		16







W UNIVERSITY of WASHIN	GTON L17: Caches III CSE	351, Winter 2017
Write-b	oack, write-allocate example	
mov 0xFAC	CE,F	
Cache	G OxBEEF O dirty bit	
Memory	F OxCAFE	
	G OxBEEF	
		20

W UNIVERSITY of WASHIN	IGTON L17: Caches III	CSE351, Winter 2017
Write-k	oack, write-allocate examp	ble
mov 0xFA	CE,F	
Cache	F OxCAFE 0	— dirty bit
	Step 2	L: Bring F into cache
Memory	F OxCAFE G OxBEEF	
		21

W UNIVERSITY of WASHIN Write-k	GTON L17: Caches III CSE331, Wint Dack, write-allocate example	er 2017
mov 0xFA	CE,F	
Cache Memory	E 0xFACE 1 dirty bit different! Step 2: Write 0xFACE to cache only and set dirty bit G 0xBEEF	
		22

W UNIVERSITY of WASHIN	GTON	L17: Caches III	CSE351, Winter 2017
Write-b	ack, w	v rite-allocate ex	ample
Cache	F	0xFACE 1	——— dirty bit
Memory	F G	0xCAFE 0xBEEF	Write hit! Write 0xFEED to cache only
			23

W UNIVERSITY of WASHING	TON L17: Caches III	CSE351, Winter 2017						
Write-back, write-allocate example								
mov 0xFAC	CE, F mov OxFEED, F mov G,	,%rax						
Cache	F OxFEED 1	— dirty bit						
Memory	F OxCAFE G OxBEEF							
		24						

W UNIVERSITY of WASHIN	GTON	L17: Caches III	CSE351, Winter 2017
Write-k	back, v	write-allocat	te example
mov 0xFA	CE,F	mov OxFEED, F	mov G , %rax
Cache	G	OxBEEF	0 dirty bit
Memory	F G	0xBEEF	1. Write F back to memory since it is dirty 2. Bring G into the cache so we can copy it into %rax
			25

W UNIVERSITY of WASHINGTON	L17: Caches III CSE351, Winter 2017
Where else is cac	ning used?
 Software caches are mo File system buffer caches, Content-delivery network 	re flexible browser caches, etc. s (CDN): cache for the Internet (e.g. <i>Netflix</i>)
 Some design differences Almost always fully-assoc so, no placement restricti index structures like hash More complex replacemes misses are very expensive worth thousands of cycles Not necessarily constraints may fetch or write-back in 	5 iative ons tables are common (for placement) nt policies when disk or network involved s to avoid them ed to single "block" transfers n larger units, opportunistically
	26







🔰 UNIVERSITY 🧃	f WASHINGTON		L	17: Caches	Ш			CSE351, Winter 2017
Naŭ	vo Mat	riv	Multin	h,				
INdiv		IX	wuitip	ıy				
# mc	ove alor	ıg ı	rows of	А				
for	(i = 0);	; i	< n; i-	++)				
#	move al	Lond	g columr	ns d	of B			
fc	or (i =	0;	, i < n;	·,+-	+)			
	# EACH	k '	loop rea	ads	row of	Α.	col of	в
	# Also	rea	ad & wr	i t e	C(i i)	n t	-imes	_
	for (k				(++)	11 (
		— ((+-)	1, 1 1 * n -) 	[]r*r	chec →→1• chec	k men
	C[T_l	IT J			FK] ^ D		I⊤] ; ((((¢	ss patiern
	G Write (suara	iteed)	3 Read	_	O Read	_	@Read	_
	C(i,i)		C(i.i)		A(i.:)			
				+		×	B(:,j)	
				'				
l				J				







IVERSITY of WASHINGTON L17: Caches III CSE351, Winter 2 Linear Algebra to the Rescue (2)													
C ₁₁	C ₁₂	C ₁₃	C ₁₄		A ₁₁	A ₁₂	A ₁₃	A ₁₄		B ₁₁	B ₁₂	B ₁₃	B ₁₄
C ₂₁	C ₂₂	C ₂₃	C ₂₄		(A_{21})	A22	A28	A24		B ₂₁	B22	B ₂₃	B ₂₄
C ₃₁	C ₃₂	C ₄₃	C ₃₄		A ₃₁	A ₃₂	A ₃₃	A ₃₄		B ₃₂	B32	B ₃₃	B ₃₄
C ₄₁	C ₄₂	C ₄₃	C ₄₄		A ₄₁	A ₄₂	A ₄₃	A ₁₄₄		B ₄₁	B ₄₂	B ₄₃	B ₄₄
Mat	rices	s of s	ize <i>r</i>	ιXi	n, sp	olit ir	to 4	blo	cks	of s	ize <i>r</i>	(n=-	4 <i>r</i>)
$C_{22} = A_{21}B_{12} + A_{22}B_{22} + A_{22}B_{22} + A_{24}B_{42} = \sum_{k} A_{2k}^{*}B_{k2}$													
* N •	 ∑₂₂ = A₂₁B₁₂ + A₂₂B₂₂ + A₂₃B₃₂ + A₂₄B₄₂ = ∑_k A_{2k}*B_{k2} Multiplication operates on small "block" matrices Choose size so that they fit in the cache! This technique called "cache blocking" 												







UNIVERSITY of WASHINGTON	L17: Caches III	CSE351, Winter 2017
Matrix Mult	iply Summary	
 ∗ Naïve: (9 ∗ Blocked: 1/0 	$\begin{array}{l} (8) \times n^3 \\ (4r) \times n^3 \end{array}$	
 If r = 8, differ If r = 16, differ 	ence is 4*8 * 9/8 = <mark>36x</mark> ence is 4*16 * 9/8 = <mark>72x</mark>	
 Blocking optim Suggests larges 	zation only works if the blocks f <u>t possible</u> block size up to limit $3r^2 \leq$	it in the cache C
 Matrix multipli Input data: 3n² Every array ele But program ha 	cation has inherent temporal loc , computation $2n^3$ ment used $O(n)$ times! as to be written properly	cality:
		38



W UNIVERSITY of WASHINGTON L17: Caches III CSE35	1, Winter 2017
Cache-Friendly Code	
 Programmer can optimize for cache performance How data structures are organized How data are accessed Nested loop structure Blocking is a general technique All systems favor "cache-friendly code" Getting absolute optimum performance is very platform specific Cache size, cache block size, associativity, etc. Can get most of the advantage with generic code Keep working set reasonably small (temporal locality) Use small strides (spatial locality) Focus on inner loop code 	
	40



Li	nux:
	lscpu
1	<pre>ls /sys/devices/system/cpu/cpu0/cache/index0/ Ex: cat /sys/devices/system/cpu/cpu0/cache/index*/size</pre>
	cat/proc/cpuinfo grep cache sort uniq
W	/indows:
	wmic memcache get <query> (all values in KB)</query>
1	<u>Ex</u> :wmic memcache get MaxCacheSize