





Why avoiding communication is important (2/2) • Running time of an algorithm is sum of 3 terms: - # flops * time_per_flop - # words moved / bandwidth - # messages * latency					
 Time_per_flop << 1/ bandwidth << latency 					
 Gaps growing exponentially with time 					
Annual improvements					
Time_per_flo	р	Bandwidth	Latency		
50%	Network	26%	15%		
59%	DRAM	23%	5%		
Goal : organize linear algebra to avoid communication					
Between all memory hierarchy levels					
 L1 ← L2 ← DRAM ← petwork, etc 					
• Not just <i>hiding</i> communication (overlap with arith) (speedup $\leq 2x$)					
Arbitrary speedups possible					
Slide source: Jim Demmel, CS267					



L <u>ess Communication with Blocked Matrix Multip</u> ly
• <u>Blocked</u> Matmul C = A·B explicitly refers to subblocks
of <i>N</i> , <i>B</i> and <i>C</i> of annehsions that depend on cache size
Break Anxn, Bnxn, Cnxn into bxb blocks labeled A(i,j), etc
b chosen so 3 bxb blocks fit in cache
for $i = 1$ to n/b , for $j=1$ to n/b , for $k=1$ to n/b
$C(i,j) = C(i,j) + A(i,k) \cdot B(k,j)$ b x b matmul, 4b ² reads/writes
• $(n/b)^3 \cdot 4b^2 = 4n^3/b$ reads/writes altogether
• Minimized when $3b^2$ = cache size = M, yielding O(n ³ /M ^{1/2}) reads/writes
What if we had more levels of memory? (L1, L2, cache etc)?
 Would need 3 more nested loops per level

CS267 Lecture 11

Slide source: Jim Demmel, CS267

UNIVERSITY





Let M = "fast" memory size per processor

- #words_moved = $\Omega(n^3/M^{1/2})$

• Parallel case, dense n x n matrices

- #words_moved = $\Omega(n^2/p^{1/2})$

- #messages_sent = $\Omega(p^{1/2})$

de source: Jim Demmel, CS267

- #messages_sent = $\Omega(n^3/M^{3/2})$

#flops = number of flops done per processor

• Sequential case, dense n x n matrices, so O(n³) flops

- Load balanced, so $O(n^3/p)$ flops processor

CS267 Lecture 11



CS267 Lecture 11

3







4













6





