

Lecture 5: Memory Hierarchy and Cache

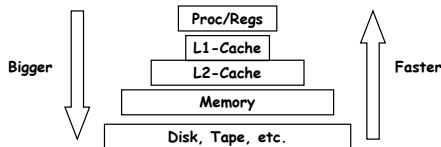
Cache: A safe place for hiding and storing things.
Webster's New World Dictionary (1976)

Cache and Its Importance in Performance

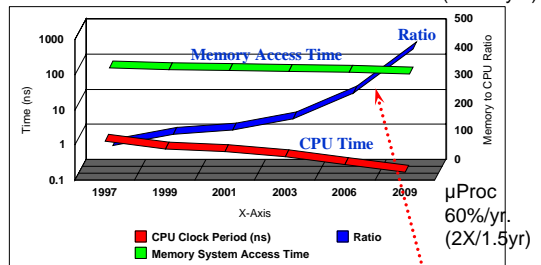
- ◆ **Motivation:**
 - Time to run code = clock cycles running code + clock cycles waiting for memory
 - For many years, CPU's have sped up an average of 50% per year over memory chip speed ups.
- ◆ Hence, memory access is the bottleneck to computing fast.
- ◆ **Definition of a cache:**
 - Dictionary: a safe place to hide or store things.
 - Computer: a level in a memory hierarchy.

What is a cache?

- ◆ Small, fast storage used to improve average access time to slow memory.
- ◆ Exploits spacial and temporal locality
- ◆ In computer architecture, almost everything is a cache!
 - Registers "a cache" on variables - software managed
 - First-level cache a cache on second-level cache
 - Second-level cache a cache on memory
 - Memory a cache on disk (virtual memory)
 - TLB a cache on page table
 - Branch-prediction a cache on prediction information?

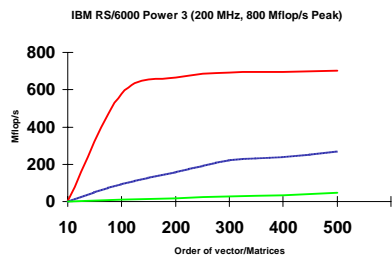


Latency in a Single System



Processor-Memory Performance Gap: (grows 50% / year)

Matrix-multiply, optimized several ways



Cache Benefits

- ◆ Data cache was designed with two key concepts in mind
 - **Spatial Locality**
 - » When an element is referenced its neighbors will be referenced too
 - » Cache lines are fetched together
 - » Work on consecutive data elements in the same cache line
 - **Temporal Locality**
 - » When an element is referenced, it might be referenced again soon
 - » Arrange code so that data in cache is reused often

Cache-Related Terms

Least Recently Used (LRU): Cache replacement strategy for set associative caches. The cache block that is least recently used is replaced with a new block.

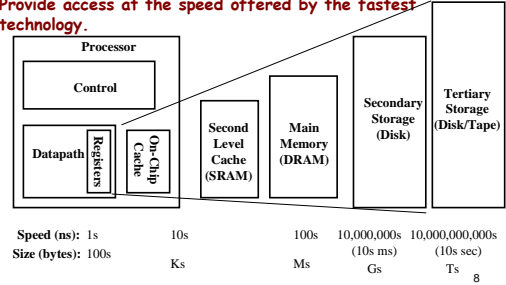
Random Replace: Cache replacement strategy for set associative caches. A cache block is randomly replaced.

7

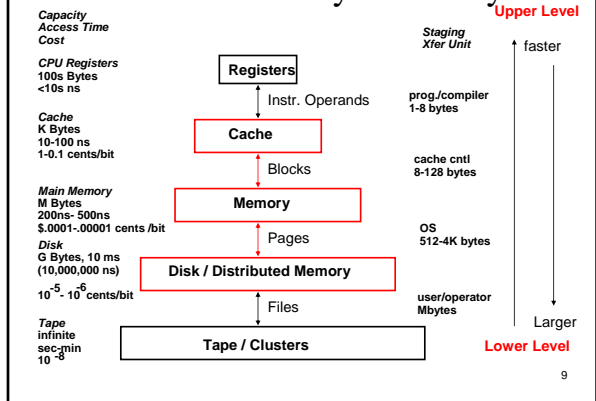
A Modern Memory Hierarchy

◆ By taking advantage of the principle of locality:

- Present the user with as much memory as is available in the cheapest technology.
- Provide access at the speed offered by the fastest technology.



Levels of the Memory Hierarchy



9

Reality

◆ Modern processors use a variety of techniques for performance

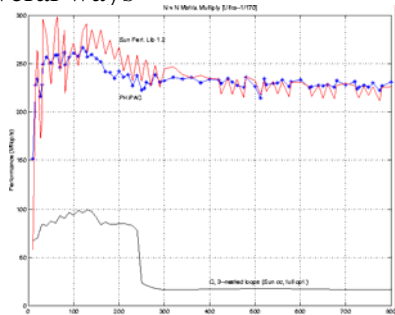
- caches
 - » small amount of fast memory where values are "cached" in hope of reusing recently used or nearby data
 - » different memory ops can have very different costs
- parallelism
 - » superscalar processors have multiple "functional units" that can run in parallel
 - » different orders, instruction mixes have different costs
- pipelining
 - » a form of parallelism, like an assembly line in a factory

◆ Why is this your problem?

- » In theory, compilers understand all of this and can optimize your program; in practice they don't.

10

Matrix-multiply, optimized several ways



11

Traditional Four Questions for Memory Hierarchy Designers

- ◆ Q1: Where can a block be placed in the upper level? (*Block placement*)
 - Fully Associative, Set Associative, Direct Mapped
- ◆ Q2: How is a block found if it is in the upper level? (*Block identification*)
 - Tag/Block
- ◆ Q3: Which block should be replaced on a miss? (*Block replacement*)
 - Random, LRU
- ◆ Q4: What happens on a write? (*Write strategy*)
 - Write Back or Write Through (with Write Buffer)

12

Cache-Related Terms

- ◆ **ICACHE** : Instruction cache
- ◆ **DCACHE (L1)** : Data cache closest to registers
- ◆ **SCACHE (L2)** : Secondary data cache
 - Data from SCACHE has to go through DCACHE to registers
 - SCACHE is larger than DCACHE
 - Not all processors have SCACHE

13

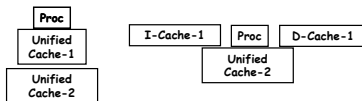
Unified versus Split Caches

- ◆ This refers to having a single or separate caches for data and machine instructions.
- ◆ Split is obviously superior. It reduces thrashing, which we will come to shortly..

14

Unified vs Split Caches

◆ Unified vs Separate I&D



◆ Example:

- 16KB I&D: Inst miss rate=0.64%, Data miss rate=6.47%
- 32KB unified: Aggregate miss rate=1.99%

◆ Which is better (ignore L2 cache)?

- Assume 33% data ops ⇒ 75% accesses from instructions (1.0/1.33)
- hit time=1, miss time=50
- Note that *data* hit has 1 stall for unified cache (only one port)

15

Where to misses come from?

◆ Classifying Misses: 3 Cs

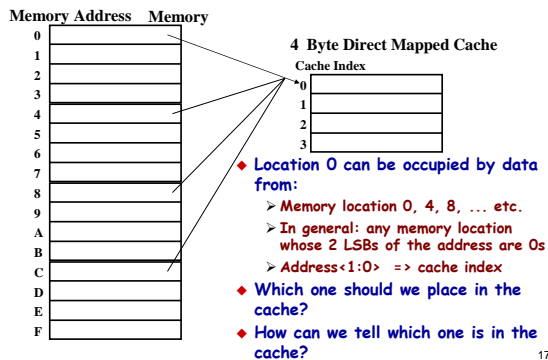
- **Compulsory**—The first access to a block is not in the cache, so the block must be brought into the cache. Also called *cold start misses* or *first reference misses*. (Misses in even an Infinite Cache)
- **Capacity**—If the cache cannot contain all the blocks needed during execution of a program, **capacity misses** will occur due to blocks being discarded and later retrieved. (Misses in Fully Associative Size X Cache)
- **Conflict**—If block-placement strategy is set associative or direct mapped, conflict misses (in addition to compulsory & capacity misses) will occur because a block can be discarded and later retrieved if too many blocks map to its set. Also called *collision misses* or *interference misses*. (Misses in N-way Associative, Size X Cache)

◆ 4th "C": (for parallel)

- **Coherence** - Misses caused by cache coherence.

16

Simplest Cache: Direct Mapped



17

Cache Mapping Strategies

- ◆ There are two common sets of methods in use for determining which cache lines are used to hold copies of memory lines.
- ◆ **Direct**: Cache address = memory address MODULO cache size.
- ◆ **Set associative**: There are N cache banks and memory is assigned to just one of the banks. There are three algorithmic choices for which line to replace:
 - **Random**: Choose any line using an analog random number generator. This is cheap and simple to make.
 - **LRU (least recently used)**: Preserves temporal locality, but is expensive. This is not much better than random according to (biased) studies.
 - **FIFO (first in, first out)**: Random is far superior.

18

Cache Basics

- Cache hit: a memory access that is found in the cache -- cheap
- Cache miss: a memory access that is not in the cache - expensive, because we need to get the data from elsewhere
- Consider a tiny cache (for illustration only)

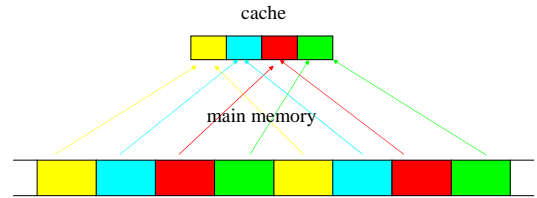


- Cache line length: number of bytes loaded together in one entry
- Direct mapped: only one address (line) in a given range in cache
- Associative: 2 or more lines with different addresses exist

19

Direct-Mapped Cache

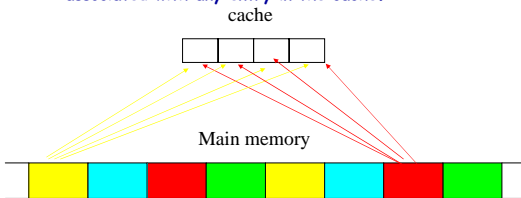
- Direct mapped cache: A block from main memory can go in exactly one place in the cache. This is called direct mapped because there is direct mapping from any block address in memory to a single location in the cache.



20

Fully Associative Cache

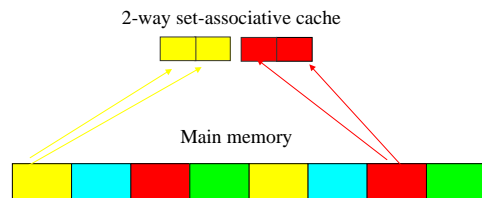
- Fully Associative Cache: A block from main memory can be placed in any location in the cache. This is called fully associative because a block in main memory may be associated with any entry in the cache.



21

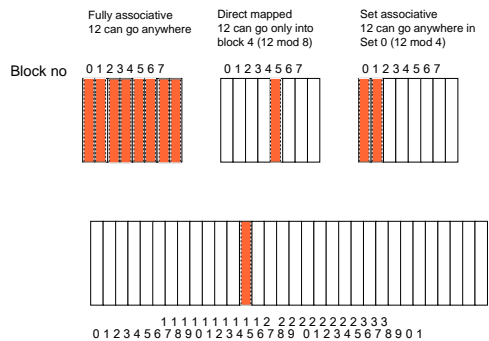
Set Associative Cache

- Set associative cache: The middle range of designs between direct mapped cache and fully associative cache is called set-associative cache. In a n -way set-associative cache a block from main memory can go into N ($N > 1$) locations in the cache.



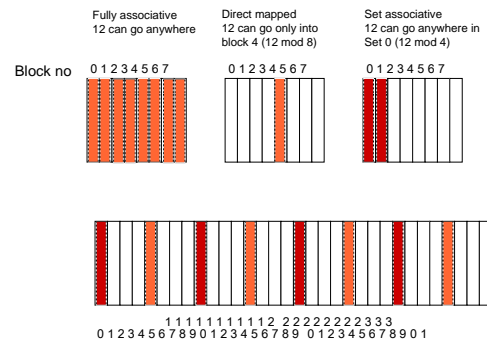
22

Here assume cache has 8 blocks, while memory has 32



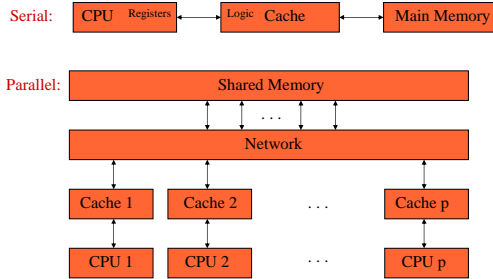
23

Here assume cache has 8 blocks, while memory has 32



24

Diagrams



25

Tuning for Caches

1. Preserve locality.
2. Reduce cache thrashing.
3. Loop blocking when out of cache.
4. Software pipelining.

26

Registers

- ◆ Registers are the source and destination of most CPU data operations.
- ◆ They hold one element each.
- ◆ They are made of static RAM (SRAM), which is very expensive.
- ◆ The access time is usually 1-1.5 CPU clock cycles.
- ◆ Registers are at the top of the memory subsystem.

27

Memory Banking

- ◆ This started in the 1960's with both 2 and 4 way interleaved memory banks. Each bank can produce one unit of memory per bank cycle. Multiple reads and writes are possible in parallel.
 - Memory chips must internally recover from an access before it is reaccessed
- ◆ The bank cycle time is currently 4-8 times the CPU clock time and getting worse every year.
- ◆ Very fast memory (e.g., SRAM) is **unaffordable** in large quantities.
- ◆ This is not perfect. Consider a 4 way interleaved memory and a stride 4 algorithm. This is equivalent to non-interleaved memory systems.

Bank 1	Bank 2	Bank 3	Bank 4
A(1)	A(2)	A(3)	A(4)
A(5)	A(6)	A(7)	A(8)
A(9)	A(10)	A(11)	A(12)
A(13)	A(14)	A(15)	A(16)
.	.	.	.
.	.	.	.
.	.	.	.

The Principle of Locality

- ◆ **The Principle of Locality:**
 - Program access a relatively small portion of the address space at any instant of time.
- ◆ **Two Different Types of Locality:**
 - **Temporal Locality** (Locality in Time): If an item is referenced, it will tend to be referenced again soon (e.g., loops, reuse)
 - **Spatial Locality** (Locality in Space): If an item is referenced, items whose addresses are close by tend to be referenced soon (e.g., straightline code, array access)
- ◆ Last 15 years, HW relied on locality for speed

29

Principals of Locality

- ◆ **Temporal:** an item referenced now will be again soon.
- ◆ **Spatial:** an item referenced now causes neighbors to be referenced soon.
- ◆ **Lines, not words, are moved between memory levels.** Both principals are satisfied. There is an optimal line size based on the properties of the data bus and the memory subsystem designs.
- ◆ Cache lines are typically 32-128 bytes with 1024 being the longest currently.

30

What happens on a write?

- ◆ **Write through**—The information is written to both the block in the cache and to the block in the lower-level memory.
- ◆ **Write back**—The information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced in cache.
 - is block clean or dirty?
- ◆ **Pros and Cons of each?**
 - WT: read misses cannot result in writes
 - WB: no repeated writes to same location
- ◆ WT always combined with write buffers so that don't wait for lower level memory

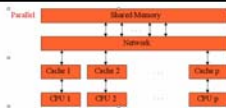
31

Cache Thrashing

- ◆ Thrashing occurs when frequently used cache lines replace each other. There are three primary causes for thrashing:
 - Instructions and data can conflict, particularly in unified caches.
 - Too many variables or too large of arrays are accessed that do not fit into cache.
 - Indirect addressing, e.g., sparse matrices.
- ◆ Machine architects can add sets to the associativity. Users can buy another vendor's machine. However, neither solution is realistic.

32

Cache Coherence for Multiprocessors



- ◆ All data must be coherent between memory levels. Multiple processors with separate caches must inform the other processors quickly about data modifications (by the cache line). **Only hardware is fast enough to do this.**
- ◆ Standard protocols on multiprocessors:
 - **Snoopy**: all processors monitor the memory bus.
 - **Directory based**: Cache lines maintain an extra 2 bits per processor to maintain clean/dirty status bits.

33

Processor Stall

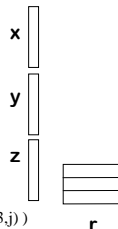
- ◆ **Processor stall** is the condition where a cache miss occurs and the processor waits on the data.
- ◆ A better design allows any instruction in the instruction queue to execute that is ready. You see this in the design of some RISC CPU's, e.g., the RS6000 line.
- ◆ Memory subsystems with **hardware data prefetch** allow scheduling of data movement to cache.
- ◆ **Software pipelining** can be done when loops are unrolled. In this case, the data movement overlaps with computing, usually with reuse of the data.
- ◆ out of order execution, software pipelining, and prefetch.

34

Indirect Addressing

```

d = 0
do i = 1,n
  j = ind(i)
  d = d + sqrt( x(j)*x(j) + y(j)*y(j) + z(j)*z(j) )
end do
    
```



- ◆ Change loop statement to


```

d = d + sqrt( r(1,j)*r(1,j) + r(2,j)*r(2,j) + r(3,j)*r(3,j) )
            
```
- ◆ Note that $r(1,j)$ - $r(3,j)$ are in contiguous memory and probably are in the same cache line (d is probably in a register and is irrelevant). The original form uses 3 cache lines at every instance of the loop and can cause cache thrashing.

35

Cache Thrashing by Memory Allocation

```

parameter ( m = 1024*1024 )
real a(m), b(m)
    
```

- ◆ For a 4 Mb direct mapped cache, $a(i)$ and $b(i)$ are always mapped to the same cache line. This is trivially avoided using padding.

```

real a(m), extra(32), b(m)
    
```

- ◆ extra is at least 128 bytes in length, which is longer than a cache line on all but one memory subsystem that is available today.

36

Cache Blocking

- ◆ We want blocks to fit into cache. On parallel computers we have $p \times$ cache so that data may fit into cache on p processors, but not one. This leads to superlinear speed up! Consider matrix-matrix multiply.

```

do k = 1,n
  do j = 1,n
    do i = 1,n
      c(i,j) = c(i,j) + a(i,k)*b(k,j)
    end do
  end do
end do

```

- ◆ An alternate form is ...

37

Cache Blocking

```

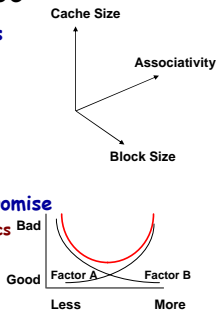
do kk = 1,n,nblk
  do jj = 1,n,nblk
    do ii = 1,n,nblk
      do k = kk,kk+nblk-1
        do j = jj,jj+nblk-1
          do i = ii,ii+nblk-1
            c(i,j) = c(i,j) + a(i,k) * b(k,j)
          end do
        end do
      end do
    end do
  end do
end do

```

38

Summary : The Cache Design Space

- ◆ Several interacting dimensions
 - cache size
 - block size
 - associativity
 - replacement policy
 - write-through vs write-back
 - write allocation
- ◆ The optimal choice is a compromise
 - depends on access characteristics
 - » workload
 - » use (I-cache, D-cache, TLB)
 - depends on technology / cost
- ◆ Simplicity often wins



39

Lessons

- ◆ The actual performance of a simple program can be a complicated function of the architecture
- ◆ Slight changes in the architecture or program change the performance significantly
- ◆ Since we want to write fast programs, we must take the architecture into account, even on uniprocessors
- ◆ Since the actual performance is so complicated, we need simple models to help us design efficient algorithms
- ◆ We will illustrate with a common technique for improving cache performance, called **blocking**

40

Optimizing Matrix Addition for Caches

- ◆ Dimension $A(n,n)$, $B(n,n)$, $C(n,n)$
- ◆ A , B , C stored by column (as in Fortran)
- ◆ Algorithm 1:
 - for $i=1:n$, for $j=1:n$, $A(i,j) = B(i,j) + C(i,j)$
- ◆ Algorithm 2:
 - for $j=1:n$, for $i=1:n$, $A(i,j) = B(i,j) + C(i,j)$
- ◆ What is "memory access pattern" for Algs 1 and 2?
- ◆ Which is faster?
- ◆ What if A , B , C stored by row (as in C)?

41

Homework Assignment

- ◆ Implement, in Fortran or C, the six different ways to perform matrix multiplication by interchanging the loops. (Use 64-bit arithmetic.) Make each implementation a subroutine, like:
 - ◆ subroutine jkj (a, m, n, lda, b, k, ldb, c, ldc)
 - ◆ subroutine ikj (a, m, n, lda, b, k, ldb, c, ldc)
 - ◆ ...

42

Loop Fusion Example

```

/* Before */
for (i = 0; i < N; i = i+1)
  for (j = 0; j < N; j = j+1)
    a[i][j] = 1/b[i][j] * c[i][j];
for (i = 0; i < N; i = i+1)
  for (j = 0; j < N; j = j+1)
    d[i][j] = a[i][j] + c[i][j];

/* After */
for (i = 0; i < N; i = i+1)
  for (j = 0; j < N; j = j+1)
    {
      a[i][j] = 1/b[i][j] * c[i][j];
      d[i][j] = a[i][j] + c[i][j];
    }

```

2 misses per access to a & c vs. one miss per access; improve spatial locality

43

Optimizing Matrix Multiply for Caches

- ◆ Several techniques for making this faster on modern processors
 - heavily studied
- ◆ Some optimizations done automatically by compiler, but can do much better
- ◆ In general, you should use optimized libraries (often supplied by vendor) for this and other very common linear algebra operations
 - BLAS = Basic Linear Algebra Subroutines
- ◆ Other algorithms you may want are not going to be supplied by vendor, so need to know these techniques

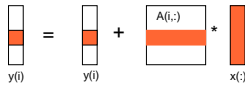
44

Warm up: Matrix-vector multiplication $y = y + A*x$

```

for i = 1:n
  for j = 1:n
    y(i) = y(i) + A(i,j)*x(j)

```



45

Warm up: Matrix-vector multiplication $y = y + A*x$

```

{read x(1:n) into fast memory}
{read y(1:n) into fast memory}
for i = 1:n
  {read row i of A into fast memory}
  for j = 1:n
    y(i) = y(i) + A(i,j)*x(j)
  {write y(1:n) back to slow memory}

```

- m = number of slow memory refs = $3*n + n^2$
- f = number of arithmetic operations = $2*n^2$
- q = $f/m \approx 2$
- Matrix-vector multiplication limited by slow memory speed

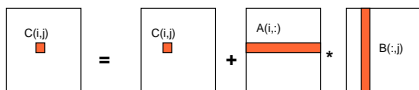
46

Multiply $C=C+A*B$

```

for i = 1 to n
  for j = 1 to n
    for k = 1 to n
      C(i,j) = C(i,j) + A(i,k) * B(k,j)

```



47

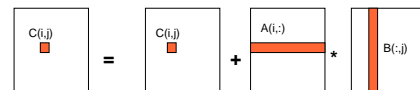
Matrix Multiply

$C=C+A*B$ (unblocked, or untiled)

```

for i = 1 to n
  {read row i of A into fast memory}
  for j = 1 to n
    {read C(i,j) into fast memory}
    {read column j of B into fast memory}
    for k = 1 to n
      C(i,j) = C(i,j) + A(i,k) * B(k,j)
    {write C(i,j) back to slow memory}

```



48

Matrix Multiply (unblocked, or untiled)

q=ops/slow mem ref

Number of slow memory references on unblocked matrix multiply

$$m = n^3 \text{ read each column of B } n \text{ times}$$

$$+ n^2 \text{ read each column of A once for each } i$$

$$+ 2n^2 \text{ read and write each element of C once}$$

$$= n^3 + 3n^2$$

So $q = f/m = (2n^3)/(n^3 + 3n^2)$
 ≈ 2 for large n , no improvement over matrix-vector mult

49

Matrix Multiply (blocked, or tiled)

Consider A, B, C to be N by N matrices of b by b subblocks where $b=n/N$ is called the **blocksize**

```

for i = 1 to N
  for j = 1 to N
    {read block C(i,j) into fast memory}
    for k = 1 to N
      {read block A(i,k) into fast memory}
      {read block B(k,j) into fast memory}
      C(i,j) = C(i,j) + A(i,k) * B(k,j) {do a matrix multiply on blocks}
    {write block C(i,j) back to slow memory}
  
```

50

Matrix Multiply (blocked or tiled)

q=ops/slow mem ref

Why is this algorithm correct?

Number of slow memory references on blocked matrix multiply

$$m = N^2 n^2 \text{ read each block of B } N^2 \text{ times } (N^2 * n/N * n/N)$$

$$+ N^2 n^2 \text{ read each block of A } N^2 \text{ times}$$

$$+ 2n^2 \text{ read and write each block of C once}$$

$$= (2N + 2)n^2$$

So $q = f/m = 2n^3 / ((2N + 2)n^2)$
 $\approx n/N = b$ for large n

So we can improve performance by increasing the blocksize b
 Can be much faster than matrix-vector multiply ($q=2$)

Limit: All three blocks from A, B, C must fit in fast memory (cache), so we cannot make these blocks arbitrarily large: $3b^2 \leq M$, so $q \approx b \leq \sqrt{M/3}$

Theorem (Hong, Kung, 1981): Any reorganization of this algorithm (that uses only associativity) is limited to $q = O(\sqrt{M})$

51

More on BLAS (Basic Linear Algebra Subroutines)

- ◆ Industry standard interface (evolving)
- ◆ Vendors, others supply optimized implementations
- ◆ History
 - > **BLAS1 (1970s):**
 - » vector operations: dot product, saxpy ($y = \alpha x + y$), etc
 - » $m=2n, f=2n, q \sim 1$ or less
 - > **BLAS2 (mid 1980s)**
 - » matrix-vector operations: matrix vector multiply, etc
 - » $m=n^2, f=2n^2, q \sim 2$, less overhead
 - » somewhat faster than BLAS1
 - > **BLAS3 (late 1980s)**
 - » matrix-matrix operations: matrix matrix multiply, etc
 - » $m \approx 4n^2, f=O(n^3)$, so q can possibly be as large as n , so BLAS3 is potentially much faster than BLAS2
- ◆ Good algorithms used BLAS3 when possible (LAPACK)
- ◆ www.netlib.org/blas, www.netlib.org/lapack

52

BLAS for Performance

Alpha EV 5/6 500MHz (1Gflop/s peak)

◆ Development of blocked algorithms important for performance

BLAS 3 (n-by-n matrix matrix multiply) vs
 BLAS 2 (n-by-n matrix vector multiply) vs
 BLAS 1 (saxpy of n vectors)

53

Optimizing in practice

- ◆ Tiling for registers
 - > loop unrolling, use of named "register" variables
- ◆ Tiling for multiple levels of cache
- ◆ Exploiting fine-grained parallelism within the processor
 - > super scalar
 - > pipelining
- ◆ Complicated compiler interactions
- ◆ Hard to do by hand (but you'll try)
- ◆ Automatic optimization an active research area
 - > PHIPAC: www.icsi.berkeley.edu/~bilmes/hipac
 - > www.cs.berkeley.edu/~iyer/ascii_slides.ps
 - > ATLAS: www.netlib.org/atlas/index.html

54

Strassen's Matrix Multiply

- ◆ The traditional algorithm (with or without tiling) has $O(n^3)$ flops
- ◆ Strassen discovered an algorithm with asymptotically lower flops
 - $O(n^{2.81})$
- ◆ Consider a 2x2 matrix multiply, normally 8 multiplies

Let $M = \begin{bmatrix} m11 & m12 \\ m21 & m22 \end{bmatrix} = \begin{bmatrix} a11 & a12 \\ a21 & a22 \end{bmatrix} \cdot \begin{bmatrix} b11 & b12 \\ b21 & b22 \end{bmatrix}$

Let $p1 = (a12 - a22) * (b21 + b22)$ $p5 = a11 * (b12 - b22)$
 $p2 = (a11 + a22) * (b11 + b22)$ $p6 = a22 * (b21 - b11)$
 $p3 = (a11 - a21) * (b11 + b12)$ $p7 = (a21 + a22) * b11$
 $p4 = (a11 + a12) * b22$

Then $m11 = p1 + p2 - p4 + p6$ Extends to nxn by divide&conquer
 $m12 = p4 + p5$
 $m21 = p6 + p7$
 $m22 = p2 - p3 + p5 - p7$

55

Strassen (continued)

$$\begin{aligned} T(n) &= \text{Cost of multiplying } nxn \text{ matrices} \\ &= 7 \cdot T(n/2) + 18 \cdot (n/2)^2 \\ &= O(n^{\log_2 7}) \\ &= O(n^{2.81}) \end{aligned}$$

- Available in several libraries
- Up to several times faster if n large enough (100s)
- Needs more memory than standard algorithm
- Can be less accurate because of roundoff error
- Current world's record is $O(n^{2.376..})$

56

Summary

- ◆ Performance programming on uniprocessors requires
 - understanding of memory system
 - » levels, costs, sizes
 - understanding of fine-grained parallelism in processor to produce good instruction mix
- ◆ Blocking (tiling) is a basic approach that can be applied to many matrix algorithms
- ◆ Applies to uniprocessors and parallel processors
 - The technique works for any architecture, but choosing the blocksize b and other details depends on the architecture
- ◆ Similar techniques are possible on other data structures
- ◆ You will get to try this in Assignment 2 (see the class homepage)

57

Summary: Memory Hierarchy

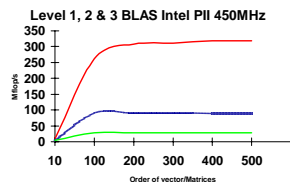
- ◆ Virtual memory was controversial at the time: can SW automatically manage 64KB across many programs?
 - 1000X DRAM growth removed the controversy
- ◆ Today VM allows many processes to share single memory without having to swap all processes to disk: today VM protection is more important than memory hierarchy
- ◆ Today CPU time is a function of (ops, cache misses) vs. just f(ops):
What does this mean to Compilers, Data structures, Algorithms?

58

Performance = Effective Use of Memory Hierarchy

- ◆ Can only do arithmetic on data at the top of the hierarchy
- ◆ Higher level BLAS lets us do this

BLAS	Memory Refs	Flops	Flops/Memory Refs
Level 1 $y = y + ax$	$3n$	$2n$	$2/3$
Level 2 $y = y + Ax$	n^2	$2n^2$	2
Level 3 $C = C + AB$	$4n^2$	$2n^2$	$n/2$



- ◆ Development of blocked algorithms important for performance

59

Homework Assignment

- ◆ Implement, in Fortran or C, the six different ways to perform matrix multiplication by interchanging the loops. (Use 64-bit arithmetic.) Make each implementation a subroutine, like:
 - ◆ subroutine ijk (a, m, n, lda, b, k, ldb, c, ldc)
 - ◆ subroutine ikj (a, m, n, lda, b, k, ldb, c, ldc)
 - ◆ ...

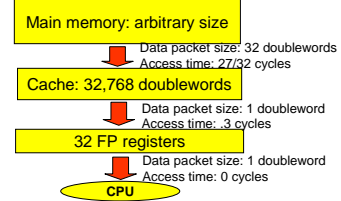
60

Thanks

- ◆ These slides came in part from courses taught by the following people:
 - Kathy Yelick, UC, Berkeley
 - Dave Patterson, UC, Berkeley
 - Randy Katz, UC, Berkeley
 - Craig Douglas, U of Kentucky
- ◆ Computer Architecture A Quantitative Approach, Chapter 8, Hennessy and Patterson, Morgan Kaufman Pub.

61

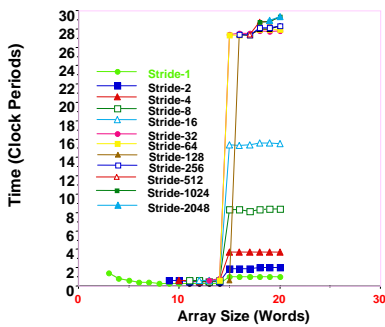
Schematic View of a Typical Memory Hierarchy: IBM 590



Design your program for optimal spatial and temporal data locality !

62

Effect of Stride and Array Size on Access Time



63

Optimal Data Locality: Data Structures

<pre> dimension rx(n), ry(n), & rz(n), fx(n), fy(n), fz(n) do 100 i=1,n do 100 j=1,i-1 dist=(rx(i)-rx(j))**2 + (ry(i)-ry(j))**2 + (rz(i)-rz(j))**2 + if (dist.le.cutoff) then c calculate interaction dfx=- dfy=- dfz=- c accumulate force fx(j)=fx(j)+dfx fy(j)=fy(j)+dfy fz(j)=fz(j)+dfz endif endif 100 continue </pre>	<pre> dimension r(3,n), f(3,n) do 100 i=1,n do 100 j=1,i-1 dist=(rx(i)-rx(j))**2 + (ry(i)-ry(j))**2 + (rz(i)-rz(j))**2 + if (dist.le.cutoff) then c calculate interaction dfx=- dfy=- dfz=- c accumulate force fx(1,j)=fx(1,j)+dfx fy(2,j)=fy(2,j)+dfy fz(3,j)=fz(3,j)+dfz endif endif 100 continue </pre>
---	---

64

Instruction Level Parallelism: Floating Point

- ◆ **RS/6000 Power2 (130 MHz)**
 - 2 FP units, each capable of
 - 1 fused multiply-add (1:2) or
 - 1 add (1:1) or
 - 1 multiply (1:2)
 - 1 quad load/store (1:1)
 leading to (up to)
 - 4 FP Ops per CP
 - 4 mem access Ops per CP
- ◆ **DEC Alpha EV5 (350 MHz)**
 - 1 FP unit, capable of
 - 1 floating point add pipeline (1:4)
 - 1 floating point mult. pipeline (1:4)
 - 1 load/store (1:3)
 leading to (up to)
 - 2 FP Ops per CP
 - 1 mem access Ops per CP
- ◆ **SGI R10000 (200 MHz)**
 - 1 FP unit capable of
 - 1 floating point add pipeline (1:2)
 - 1 floating point multiply pipeline (1:2)
 - 1 load/store (1:3)
 leading to (up to)
 - 2 FP Ops per CP
 - 1 memory access Ops per CP

65

Code Restructuring for On-Chip Parallelism: Original Code

```

program length
parameter (n=2**14)
dimension a(n)
subroutine length1(n,aa,tt)
implicit real*8 (a-h,o-z)
dimension a(n)
tt=0.d0
do 100, j=1,n
  tt=tt+a(j)*a(j)
100 continue
return
end
        
```

66

Modified Code for On-Chip Parallelism

```

subroutine length4 (n,a,tt)
c works correctly only if n is multiple of 4
implicit real*8 (a-h,o-z)
dimension a(n)
t1=0.d0
t2=0.d0
t3=0.d0
t4=0.d0
do 100, j=1,n-3,4
c first floating point instruction unit, all even cycles
t1=t1+a(j+0)*a(j+0)
c first floating point instruction unit, all odd cycles
t2=t2+a(j+1)*a(j+1)
c second floating point instruction unit, all even cycles
t3=t3+a(j+2)*a(j+2)
c second floating point instruction unit, all odd cycles
t4=t4+a(j+3)*a(j+3)
100 continue
tt= t1+t2+t3+t4
return
end

```

67

Software Pipelining

```

c first FP unit, first cycle:
c do one MADD (with t1 and a0 available in registers) and load a1:
t1=t1+a0*a0
a1=a(j+1)
c first floating point unit, second cycle:
c do one MADD (with t2 and t1 available in registers) and load a0 for next
iteration:
t2=t2+a1*a1
a0=a(j+0+4)
c second FP unit, first cycle:
c do one MADD (with t3 and a2 available in registers) and load a3:
t3=t3+a2*a2
a3=a(j+2)
c second FP unit, second cycle:
c do one MADD (with t4 and a3 available in registers) and load a2 for next
iteration:
t4=t4+a3*a3
a2=a(j+1+4)

```

68

Improving Ratio of Floating Point Operations to Memory Accesses

```

subroutine mult(n1,nd1,n2,nd2,y,a,x)
implicit real*8 (a-h,o-z)
dimension a(nd1,nd2),y(nd2),x(nd1)

do 10, i=1,n1
t=0.d0
do 20, j=1,n2
20 t=t+a(j,i)*x(j)
return
end

```

**** 2 FLOPS
**** 2 LOADS

69

Improving Ratio of Floating Point Operations to Memory Accesses

```

c works correctly when n1,n2 are multiples of 4
dimension a(nd1,nd2), y(nd2), x(nd1)
do i=1,n1-3,4
t1=0.d0
t2=0.d0
t3=0.d0
t4=0.d0
do j=1,n2-3,4
t1=t1+a(j+0,i+0)*x(j+0)+a(j+1,i+0)*x(j+1)+
1 a(j+2,i+0)*x(j+2)+a(j+3,i+1)*x(j+3)
t2=t2+a(j+0,i+1)*x(j+0)+a(j+1,i+1)*x(j+1)+
1 a(j+2,i+1)*x(j+2)+a(j+3,i+0)*x(j+3)
t3=t3+a(j+0,i+2)*x(j+0)+a(j+1,i+2)*x(j+1)+
1 a(j+2,i+2)*x(j+2)+a(j+3,i+2)*x(j+3)
t4=t4+a(j+0,i+3)*x(j+0)+a(j+1,i+3)*x(j+1)+
1 a(j+2,i+3)*x(j+2)+a(j+3,i+3)*x(j+3)
enddo
y(i+0)=t1
y(i+1)=t2
y(i+2)=t3
y(i+3)=t4
enddo

```

32 FLOPS
20 LOADS

70

Summary of Single-Processor Optimization Techniques (I)

- ◆ Spatial and temporal data locality
- ◆ Loop unrolling
- ◆ Blocking
- ◆ Software pipelining
- ◆ Optimization of data structures
- ◆ Special functions, library subroutines

71

Summary of Optimization Techniques (II)

- ◆ Achieving high-performance requires code restructuring. Minimization of memory traffic is the single most important goal.
- ◆ Compilers are getting better: good at software pipelining. But they are not there yet: can do loop transformations only in simple cases, usually fail to produce optimal blocking, heuristics for unrolling may not match your code well, etc.
- ◆ The optimization process is machine-specific and requires detailed architectural knowledge.

72

Amdahl's Law

Amdahl's Law places a strict limit on the speedup that can be realized by using multiple processors. Two equivalent expressions for Amdahl's Law are given below:

$$t_N = (f_p/N + f_s)t_1 \quad \text{Effect of multiple processors on run time}$$

$$S = 1/(f_s + f_p/N) \quad \text{Effect of multiple processors on speedup}$$

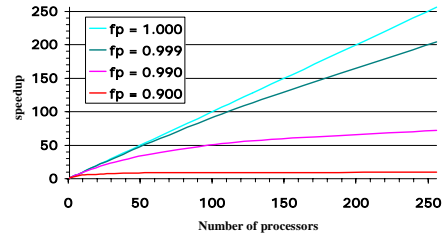
Where:

f_s = serial fraction of code
 f_p = parallel fraction of code = $1 - f_s$
 N = number of processors

73

Illustration of Amdahl's Law

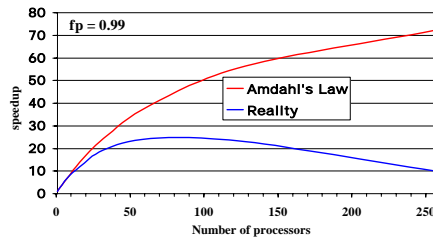
It takes only a small fraction of serial content in a code to degrade the parallel performance. It is essential to determine the scaling behavior of your code before doing production runs using large numbers of processors



74

Amdahl's Law Vs. Reality

Amdahl's Law provides a theoretical upper limit on parallel speedup assuming that there are no costs for communications. In reality, communications (and I/O) will result in a further degradation of performance.



75

More on Amdahl's Law

- ◆ Amdahl's Law can be generalized to any two processes of with different speeds
- ◆ Ex.: Apply to $f_{\text{processor}}$ and f_{memory} :
 - The growing processor-memory performance gap will undermine our efforts at achieving maximum possible speedup!

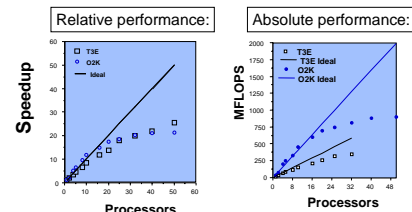
76

Gustafson's Law

- ◆ Thus, Amdahl's Law predicts that there is a maximum scalability for an application, determined by its parallel fraction, and this limit is generally not large.
- ◆ There is a way around this: *increase the problem size*
 - bigger problems mean bigger grids or more particles: bigger arrays
 - number of serial operations generally remains constant; number of parallel operations increases: parallel fraction increases

77

Parallel Performance Metrics: Speedup



Speedup is only one characteristic of a program - it is not synonymous with performance. In this comparison of two machines the code achieves comparable speedups but one of the machines is faster.

78

Fixed-Problem Size Scaling

- a.k.a. Fixed-load, Fixed-Problem Size, Strong Scaling, Problem-Constrained, constant-problem size (CPS), variable subgrid
- Amdahl Limit: $S_A(n) = T(1) / T(n) = \frac{1}{f/n + (1-f)}$
- This bounds the speedup based only on the fraction of the code that cannot use parallelism ($1-f$); it ignores all other factors
- $S_A \rightarrow 1 / (1-f)$ as $n \rightarrow \infty$

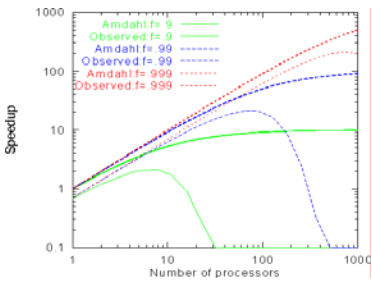
79

Fixed-Problem Size Scaling (Cont'd)

- Efficiency $(\eta) = T(1) / [T(n) * n]$
- Memory requirements decrease with n
- Surface-to-volume ratio increases with n
- Superlinear speedup possible from cache effects
- Motivation: what is the largest # of procs I can use effectively and what is the fastest time that I can solve a given problem?
- Problems:
 - Sequential runs often not possible (large problems)
 - Speedup (and efficiency) is misleading if processors are slow

80

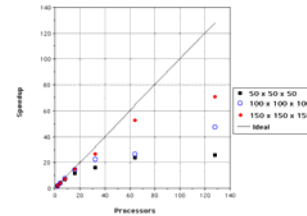
Fixed-Problem Size Scaling: Examples



S. Goedecker and [Adolfy Hoesie](#), Achieving High Performance in Numerical Computations on RISC Workstations and Parallel Systems, *International Conference on Computational Physics: PC'97 Santa Cruz, August 25-28 1997*.

81

Fixed-Problem Size Scaling Examples



82

Scaled Speedup Experiments

- a.k.a. Fixed Subgrid-Size, Weak Scaling, Gustafson scaling.
- Motivation: Want to use a larger machine to solve a larger global problem *in the same amount of time*.
- Memory and surface-to-volume effects remain constant.

83

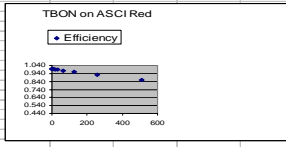
Scaled Speedup Experiments

- Be wary of benchmarks that scale problems to unreasonably-large sizes
 - scale the problem to fill the machine when a smaller size will do;
 - simplify the science in order to add computation
 - > "World's largest MD simulation - 10 gazillion particles!"
 - run grid sizes for only a few cycles because the full run won't finish during this lifetime or because the resolution makes no sense compared with resolution of input data
- Suggested alternate approach (Gustafson): Constant time benchmarks
 - run code for a fixed time and measure work done

84

Example of a Scaled Speedup Experiment

Processors	NChains	Time	Natoms	Time per Atom per PE	Time per Atom	Efficiency
1	32	38.4	2368	1.62E-02	1.62E-02	1.000
2	64	38.4	4736	8.11E-03	1.62E-02	1.000
4	128	38.5	9472	4.08E-03	1.63E-02	0.997
8	256	38.6	18944	2.04E-03	1.63E-02	0.995
16	512	38.7	37888	1.02E-03	1.63E-02	0.992
32	940	35.7	69560	5.13E-04	1.64E-02	0.987
64	1700	32.7	125800	2.60E-04	1.66E-02	0.975
128	2800	27.4	207200	1.32E-04	1.69E-02	0.958
256	4100	20.75	303400	6.84E-05	1.75E-02	0.926
512	5300	14.49	392200	3.69E-05	1.89E-02	0.857



85

Performance Optimization of Numerically Intensive Codes
 Stefan Goedecker
 Adolfo Hoisie
 Siam
 SOFTWARE · ENVIRONMENTS · TOOLS

86