

Many of the slides in this lecture are either from or adapted from slides provided by the authors of the textbook "Computer Systems: A Programmer's Perspective," 2<sup>nd</sup> Edition and are provided from the website of Carnegie-Mellon University, course 15-213, taught by Randy Bryant and David O'Hallaron in Fall 2010. These slides are indicated "Supplied by CMU" in the notes section of the slides.





If arrays x and y have the same alignment, i.e., both start in the same cache set, then each access to an element of y replaces the cache line containing the corresponding element of x, and vice versa. The result is that the loop is executed very slowly — each access to either array results in a conflict miss.



However, if the two arrays start in different cache sets, then the loop executes quickly — there is a cache miss on just every fourth access to each array.













The cache still holds two rows of the matrix, but each row may go into one of two different cache lines. In the slide, the first row goes into the first lines of the cache sets, the second row goes into the second lines of the cache sets.



There is still a cache miss on each access.



With a 2-way set-associative cache, our dot-product example runs quickly even if the two arrays have the same alignment.



The L3 cache is known as the *last-level cache* (LLC) in the Intel documentation.

One concern is whether what's contained in, say, the L1 cache is also contained in the L2 cache. if so, caching is said to be **inclusive**. If what's contained in the L1 cache is definitely not contained in the L2 cache, caching is said to be **exclusive**. An advantage of exclusive caches is that the total cache capacity is the sum of the sizes of each of the levels, whereas for inclusive caches, the total capacity is just that of the largest. An advantage of inclusive caches is that what's been brought into the cache hierarchy by one core is available to the other cores.

AMD processors tend to have exclusive caches; Intel processors tend to have inclusive caches.



Most current processors use the write-back/write-allocate approach. This causes some (surmountable) difficulties for multi-core processors that have a separate cache for each core.



This slide describes accessing memory on Intel Core I5 and I7 processors.

If the processor determines that a program is accessing memory sequentially (because the past few accesses have been sequential), then it begins the load of the next block from memory before it is requested. If this determination was correct, then the memory will be in the cache (or well on its way) before it's needed.













"Stride n" reference patterns are sequences of memory accesses in which every nth element is accessed in memory order. Thus stride 1 means that every element is accessed, starting at the beginning of a memory area, continuing to its end.



Based on slides supplied by CMU.



Adapted form a slide by CMU.





Assume we are multiplying arrays of doubles, thus each element is eight bytes long, and thus a cache line holds eight matrix elements. The slide shows a straightforward implementation of multiplying A and B to produce C.



If we reverse the order of the two outer loops, there's no change in results or performance.



Moving the loop on k to be the outer loop does not affect the result, but it improves performance.



Switching the two outer loops affects neither results nor performance.



Moving the loop on i to be the inner loop makes performance considerably worse.



The poor performance is not improved by reversing the outer loops.





Supplied by CMU.

| In Real Life   |
|--|
| <ul> <li>Multiply two 1024x1024 matrices of doubles<br/>on sunlab machines</li> </ul>            |
| <ul> <li>– ijk</li> <li>» 4.185 seconds</li> </ul>   |
| <ul><li>− kij</li><li>» 0.798 seconds</li></ul>  |
| −jki<br>» 11.488 seconds   |
| CS33 Intro to Computer Systems XVII–33 Copyright © 2022 Thomas W. Doeppner. All rights reserved. |

Г

























These definitions follow those given in "Intel® 64 and IA-32 Architectures Software Developer's Manual" and are generally accepted even outside of Intel.







The reason why there must be a separate stack in privileged mode is that the OS must be guaranteed that when it is executing, it has a valid stack, that the stack pointer must be pointing to a region of memory that can be used as a stack by the OS. Since while the program was running in user mode any value could have been put into the stack-pointer register, when the OS is invoked, it switches to a pre-allocated stack set up just for it.



When a trap or interrupt occurs, the current processor state (registers, including RIP, condition codes, etc.) are saved on the kernel stack. When the system returns back to the interrupted program, this state is restored.

