

CS 33

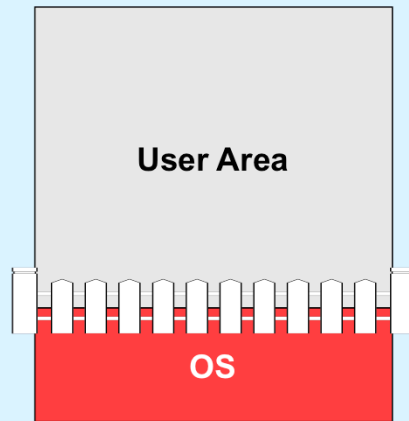
Virtual Memory

The Address-Space Concept

- **Protect processes from one another**
- **Protect the OS from user processes**
- **Provide efficient management of available storage**

The concept of the address space is fundamental in most of today's operating systems. Threads of control executing in different address spaces are protected from one another, since none of them can reference the memory of any of the others. In most systems (such as Unix), the operating system resides in address space that is shared with all processes, but protection is employed so that user threads cannot access the operating system. What is crucial in the implementation of the address-space concept is the efficient management of the underlying primary and secondary storage.

Memory Fence



Early approaches to managing the address space were concerned primarily with protecting the operating system from the user. One technique was the hardware-supported concept of the *memory fence*: an address was established below which no user mode access was allowed. The operating system was placed below this point in memory and was thus protected from the user.

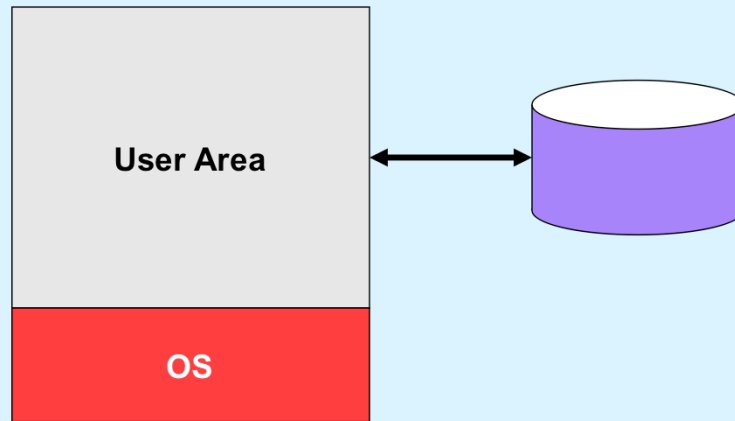
Base and Bounds Registers



The memory-fence approach protected the operating system, but did not protect user processes from one another. (This wasn't an issue for many systems—there was only one user process at a time.) Another technique, still employed in some of today's systems, is the use of *base and bounds registers* to restrict a process's memory references to a certain range. Each address generated by a user process was first compared with the value in the bounds register to make certain that it did not reference a location beyond the process's range of memory, and then was modified by adding to it the value in the base register, insuring that it did not reference a location before the process's range of memory.

A further advantage of this technique was to ensure that a process would be loaded into what appeared to be location 0—thus no relocation was required at load time.

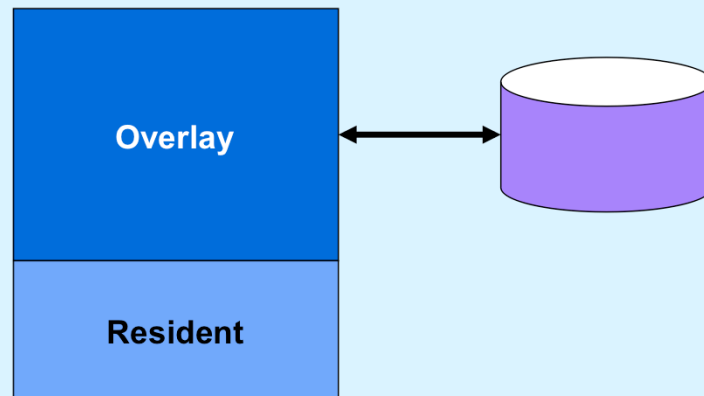
Swapping



Swapping is a technique, still in use today, in which the images of entire processes are transferred back and forth between primary and secondary storage. An early use of it was for (slow) time-sharing systems: when a user paused to think, his or her process was swapped out and that of another user was swapped in. This allowed multiple users to share a system that employed only the memory fence for protection.

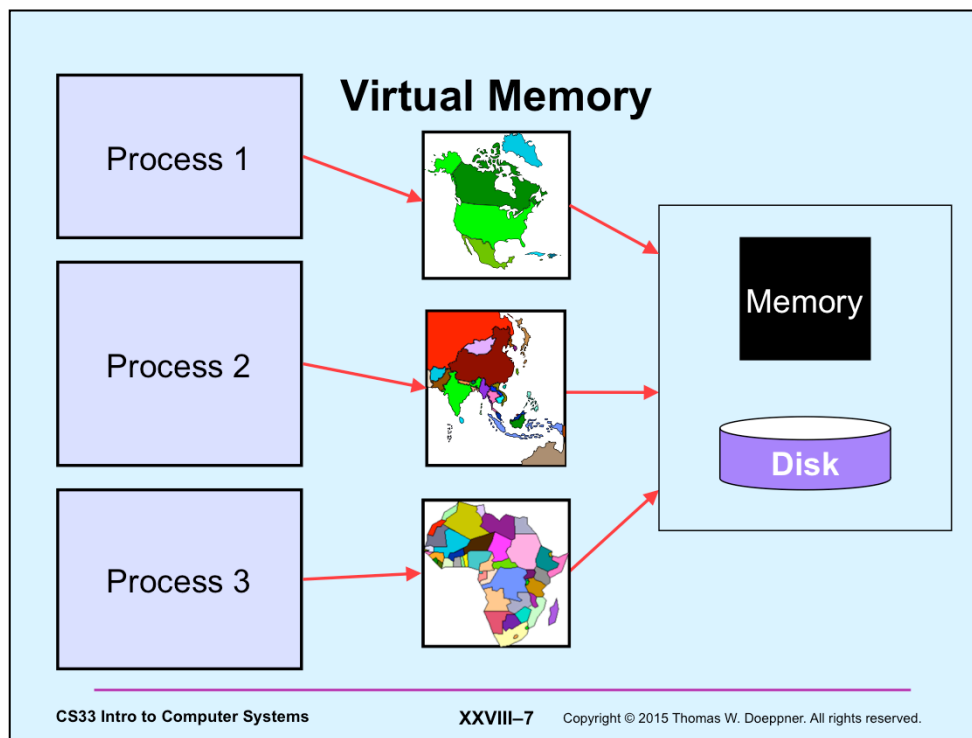
Base and bounds registers made it feasible to have a number of processes in primary memory at once. However, if one of these processes was inactive, swapping allowed the system to swap this process out and swap another process in. Note that the use of the base register is very important here: without base registers, after a process is swapped out, it would have to be swapped into the same location in which it resided previously.

Overlays



The concept of overlays is similar to the concept of swapping, except that it applies to pieces of images rather than whole images and the user is in charge. Say we have 100 kilobytes of available memory and a 200-kilobyte program. Clearly, not all the program can be in memory at once. The user might decide that one portion of the program should always be resident, while other portions of the program need be resident only for brief periods. The program might start with routines A and B loaded into memory. A calls B; B returns. Now A wants to call C, so it first reads C into the memory previously occupied by B (it *overlays* B), and then calls C. C might then want to call D and E, though there is only room for one at a time. So C first calls D, D returns, then C overlays D with E and then calls E.

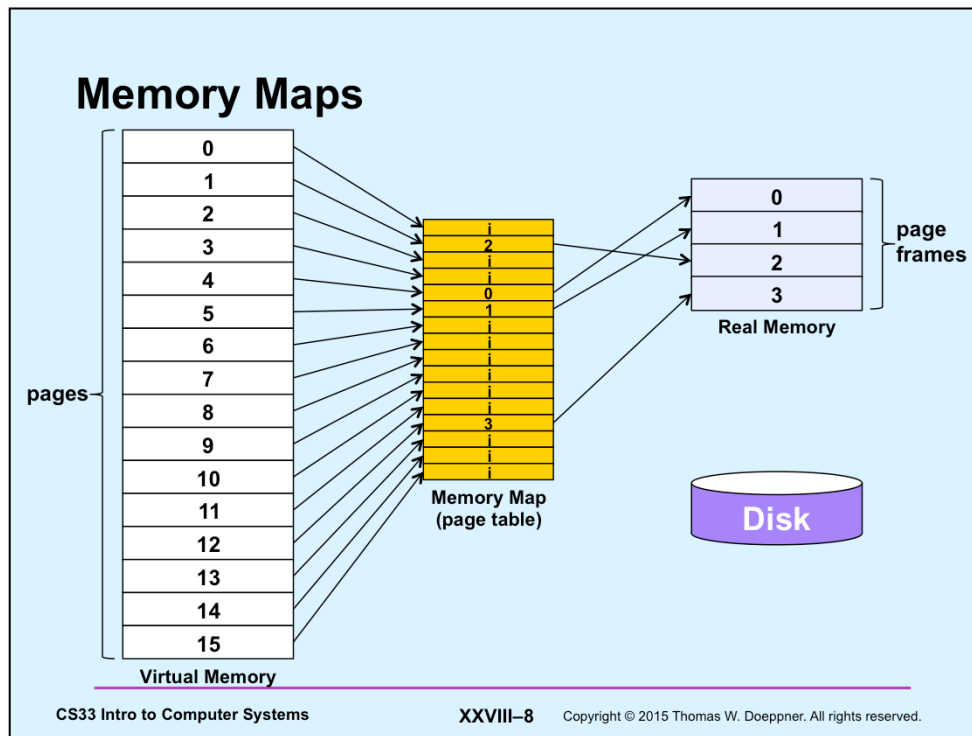
The advantage of this technique is that the programmer has complete control of the use of memory and can make the necessary optimization decisions. The disadvantage is that the programmer *must* make the necessary decisions to make full use of memory (the operating system doesn't help out). Few programmers can make such decisions wisely, and fewer still want to try.



One way to look at virtual memory is as an automatic overlay technique: processes “see” an address space that is larger than the amount of real memory available to them; the operating system is responsible for the overlaying.

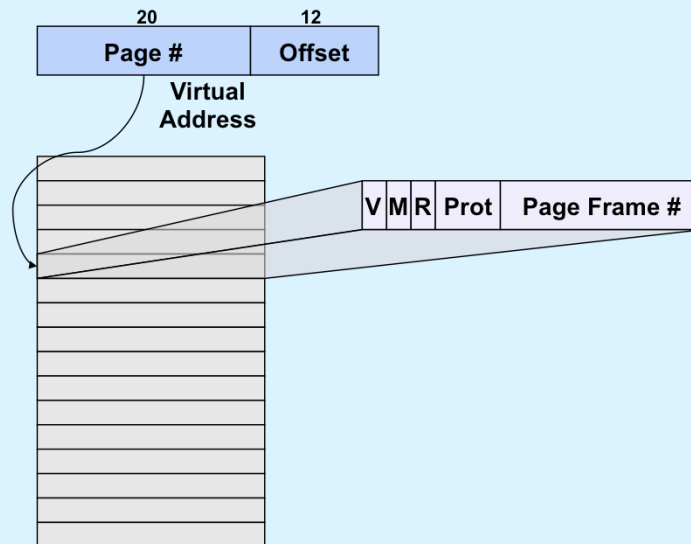
Put more abstractly (and accurately), virtual memory is the support of an address space that is independent of the size of primary storage. Some sort of mapping technique must be employed to map virtual addresses to primary and secondary stores. In the typical scenario, the computer hardware maps some virtual addresses to primary storage. If a reference is made to an unmapped address, then a fault occurs (a *page fault*) and the operating system is called upon to deal with it. The operating system might then find the desired virtual locations on secondary storage (such as a disk) and transfer them to primary storage. Or the operating system might decide that the reference is illegal and deliver an addressing exception to the process.

As with base and bounds registers, the virtual memory concept allows us to handle multiple processes simultaneously, with the processes protected from one another.



Virtual memory (what the program sees) is divided into fixed-size pages (on the x86 these are usually 4 kilobytes in size). Real memory (DRAM) is also divided into fixed-size pieces, called page frames (though they're often referred to simply as pages). A memory map, implemented in hardware and often called a page table, translates references to virtual-memory pages into references to real-memory page frames. In general, virtual memory is larger than real memory, thus not all pages can be mapped to page frames. Those that are not are said to have invalid translations.

Page Tables



A page table is an array of *page table entries*. Suppose we have, as is the usual case for the x86, a 32-bit virtual address and a page size of 4096 bytes. The 32-bit address might be split into two parts: a 20-bit *page number* and a 12-bit *offset* within the page. When a thread generates an address, the hardware uses the page-number portion as an index into the page-table array to select a page-table entry, as shown in the picture. If the page is in primary storage (i.e. the translation is valid), then the validity bit in the page-table entry is set, and the page-frame-number portion of the page-table entry is the high-order bits of the location in primary memory where the page resides. (Primary memory is thought of as being subdivided into pieces called *page frames*, each exactly big enough to hold a page; the address of each of these page frames is at a “page boundary,” so that its low-order bits are zeros.) The hardware then appends the offset from the original virtual address to the page-frame number to form the final, real address.

If the *validity bit* of the selected page-table entry is zero, then a page fault occurs and the operating system takes over. Other bits in a typical page-table entry include a *reference bit*, which is set by the hardware whenever the page is referenced, and a *modified bit*, which is set whenever the page is modified. We will see how these bits are used later in this lecture. The *page-protection bits* indicate who is allowed access to the page and what sort of access is allowed. For example, the page can be restricted for use only by the operating system, or a page containing executable code can be write-protected, meaning that read accesses are allowed but not write accesses.

Quiz 1

How many 2^{12} -byte pages fit in a 32-bit address space?

- a) a bit over a 1000
- b) a bit over a million
- c) a bit over a billion
- d) none of the above

VM is Your Friend ...

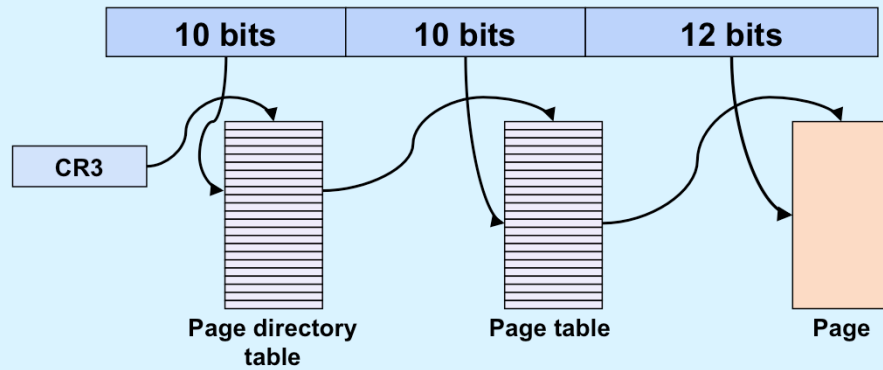
- **Not everything has to be in memory at once**
 - pages brought in (and pushed out) when needed
 - unallocated parts of the address space consume no memory
 - » e.g., hole between stack and dynamic areas
- **What's mine is not yours** (and vice versa)
 - address spaces are disjoint
- **Sharing is ok though ...**
 - address spaces don't have to be disjoint
 - » a single page frame may be mapped into multiple processes
- **I don't trust you (or me)**
 - access to individual pages can be restricted
 - » read, write, execute, or any combination

Page-Table Size

- Consider a full 2^{32} -byte address space
 - assume 4096-byte (2^{12} -byte) pages
 - 4 bytes per page-table entry
 - the page table would consist of $2^{32}/2^{12}$ ($= 2^{20}$) entries
 - its size would be 2^{22} bytes (or 4 megabytes)
 - » at \$100/gigabyte
 - around \$6
- For a 2^{64} -byte address space
 - assume 4096-byte (2^{12} -byte) pages
 - 8 bytes per page-table entry
 - the page table would consist of $2^{64}/2^{12}$ ($= 2^{52}$) entries
 - its size would be 2^{55} bytes (or 32 petabytes)
 - » at \$1/gigabyte
 - over \$33 million

In the not-all-that-distant past, 4 megabytes of memory would have cost many tens of thousands of dollars.

IA32 Paging

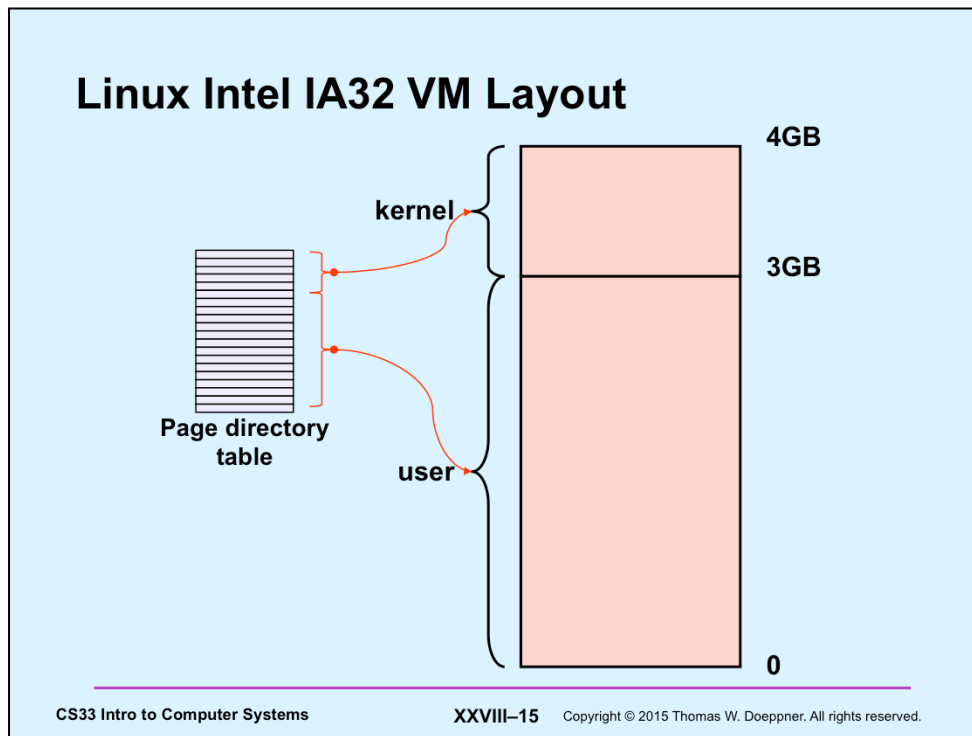


The IA32 architecture employs a two-level page table providing a means for reducing the memory requirements of the address map. The high-order 10 bits of the 32-bit virtual address are an index into what's called the page directory table. Each of its entries refer to a page table, whose entries are indexed by the next 10 bits of the virtual address. Its entries refer to individual pages; the offset within the page is indexed by the low-order 12 bits of the virtual address. The current page directory is pointed to by a special register known as CR3 (control register 3), whose contents may be modified only in privileged mode. The page directory must reside in real memory when the address space is in use, but it is relatively small (1024 4-byte entries: it's exactly one page in length). Though there are potentially a large number of page tables, only those needed to satisfy current references must be in memory at once.

Quiz 2

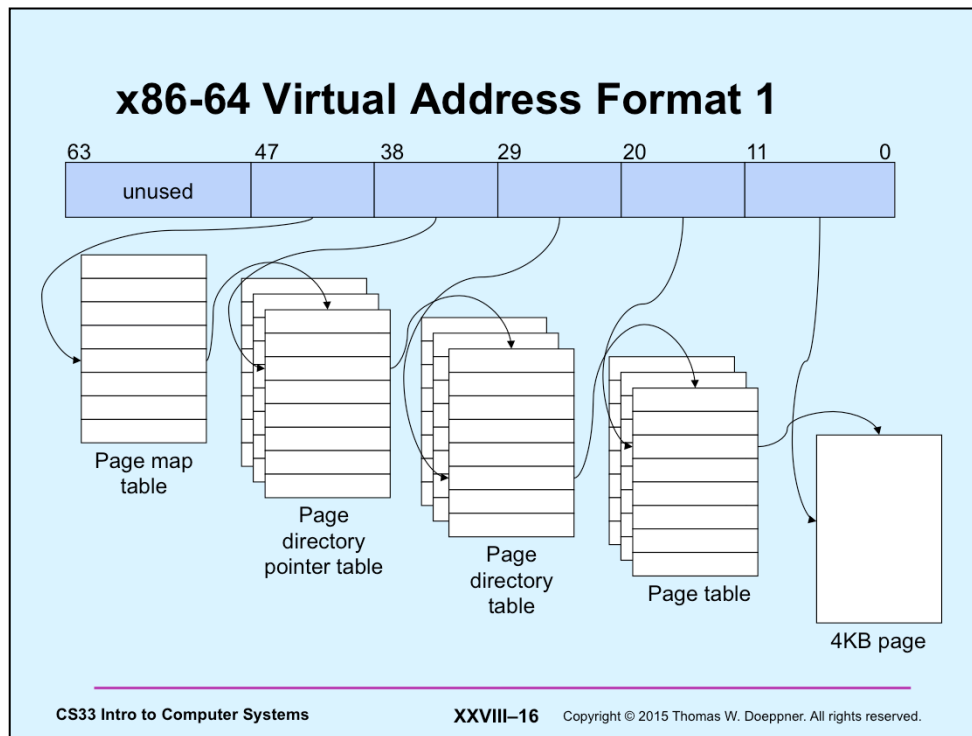
Can a page start at a virtual address that's not divisible by the page size?

- a) yes
- b) no

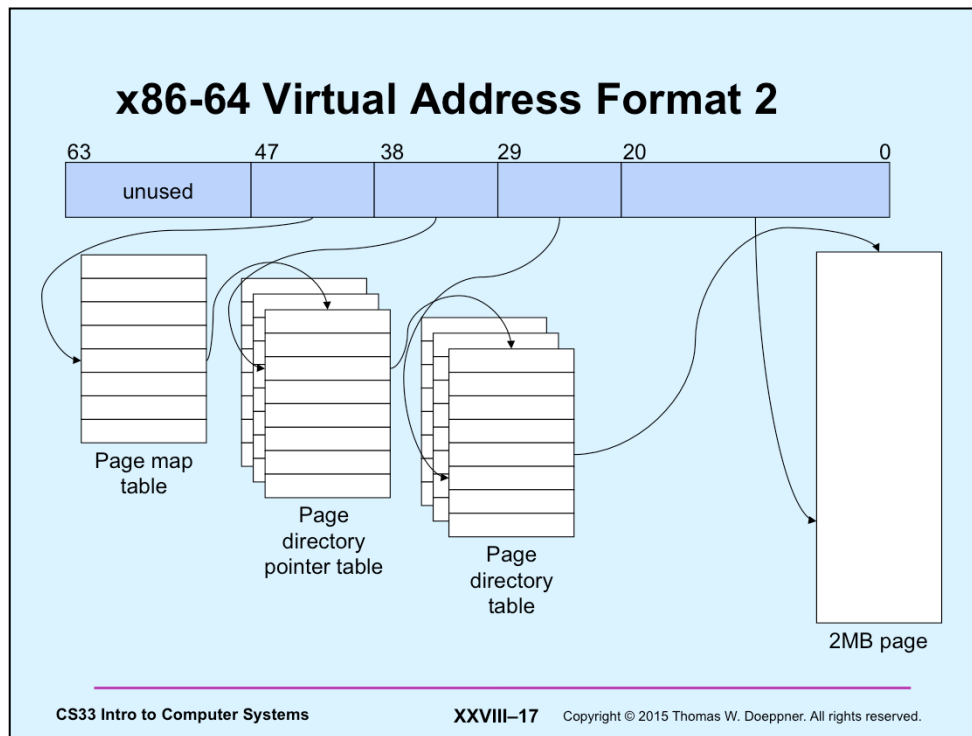


For Linux on the IA32, the OS kernel occupies the top quarter of the address space and is mapped into every process (though it may not be accessed in user mode). Each user process is mapped into the bottom three quarters of the address space — only one is mapped at a time for each processor.

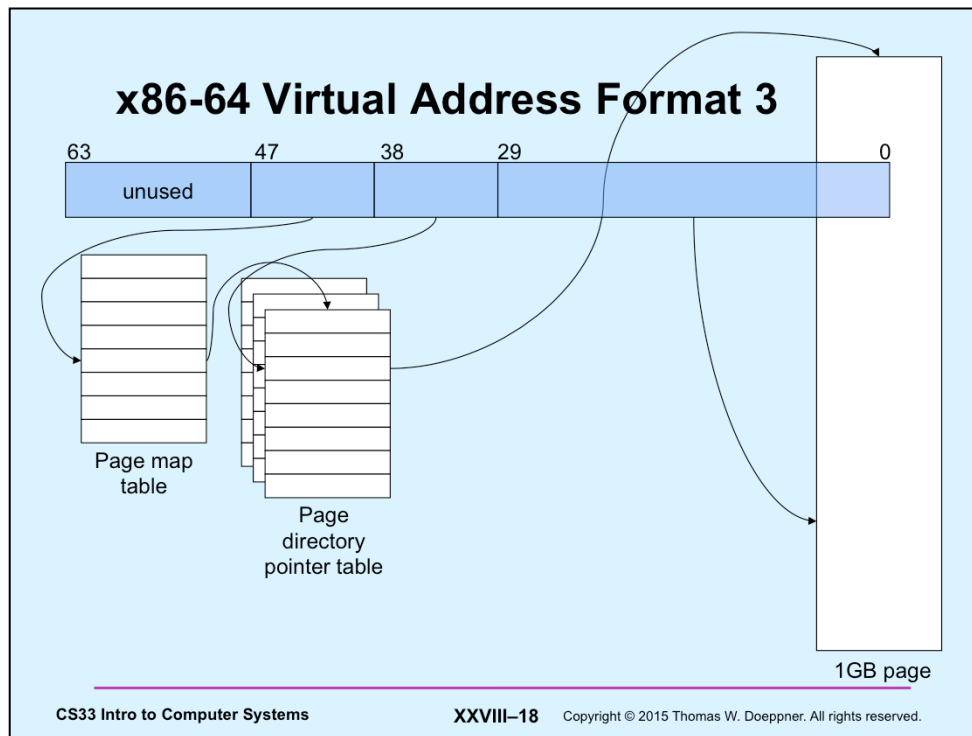
Each process has its own page-directory table describing its address space. The top quarter (256) entries are the same as for all other processes and describe the OS kernel's mappings. The bottom three quarters (768) are, in general, private to the process and describe its mappings.



For the x86-64, four levels of translation are done (the high-order 16 bits of the address are not currently used: the hardware requires that these 16 bits must all be equal to bit 47), thus it really supports “only” a 48-bit address space. Note that only the “page map table” must reside in real memory at all times. The other tables must be resident only when necessary.



Alternatively, there may be only three levels of page tables, ending with the page-directory table and 2MB pages. Both 2MB and 4KB pages may coexist in the same address space; which is being used is indicated in the associated page-directory-table entry.

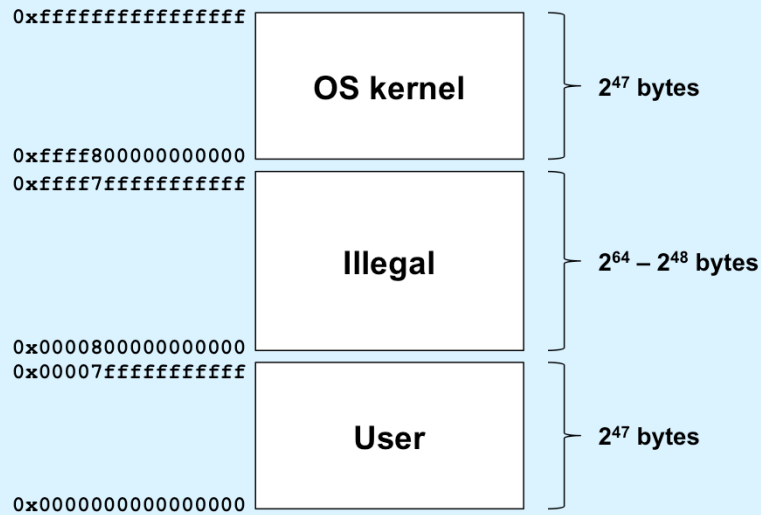


The hardware also supports 1 GB pages by eliminating the page-directory table. Not many operating systems (if any) yet take advantage of this.

Why Multiple Page Sizes?

- **External fragmentation**
 - for region composed of 4KB pages, average external fragmentation is 2KB
 - for region composed of 1GB pages, average external fragmentation is 512MB
- **Page-table overhead**
 - larger page sizes have fewer page tables
 - » less overhead in representing mappings

Address Space



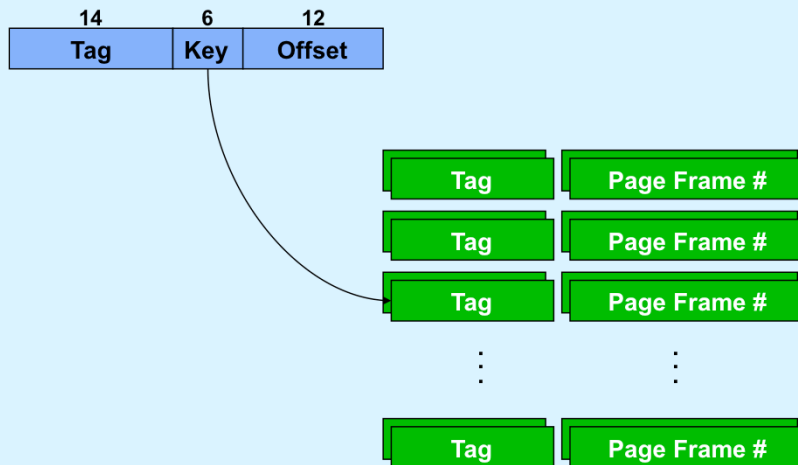
Recall that, in current implementations of the x86-64 architecture, only 48 bits of virtual address are used. Furthermore, the high-order 16 bits must be equal to bit 47. Thus the legal addresses are those at the top and at the bottom of the address space. The top addresses are used for the OS kernel, and thus mapped into all processes. The bottom addresses are used for each user process. The addresses in the middle (most of the address space — the slide is not drawn to scale!) are illegal and generate faults if used.

The reason for doing things this way (i.e., for the restrictions on the high-order bits) is to force the kernel to be at the top of the address space, allowing growth of the user portion as more virtual-address bits are supported.

Performance

- **Page table resides in real memory (DRAM)**
- **A 32-bit virtual-to-real translation requires two accesses to page tables, plus the access to the ultimate real address**
 - three real accesses for each virtual access
 - 3X slowdown!
- **A 64-bit virtual-to-real translation requires four accesses to page tables, plus the access to the ultimate real address**
 - 5X slowdown!

Translation Lookaside Buffers



To speed-up virtual-to-real translation, a special cache is maintained of recent translations — it's called the translation lookaside buffer (TLB). It resides in the chip, one per core and hyperthread. The TLB shown in the slide is a two-way set associative cache, as discussed in lecture 17. This one assumes a 32-bit virtual address with a 4k page. Things are more complicated when multiple page sizes are supported. For example, is there just one entry for a large page that covers its entire range of addresses, or is a large page dealt with by putting into the cache multiple entries covering the large page, but each for the size of a small page? Both approaches are not only possible, but done.

Quiz 3

Recall that there is a 5x slowdown on memory references via virtual memory on the x86-64. If all references are translated via the TLB, the slowdown will be

- a) 1x
- b) 2x
- c) 3x
- d) 4x

OS Role in Virtual Memory

- **Memory is like a cache**
 - quick access if what's wanted is mapped via page table
 - slow if not — OS assistance required
- **OS**
 - make sure what's needed is mapped in
 - make sure what's no longer needed is not mapped in

Mechanism

- **Program references memory**
 - if reference is mapped, access is quick
 - » even quicker if translation in TLB and referent in on-chip cache
 - if not, page-translation fault occurs and OS is invoked
 - » determines desired page
 - » maps it in, if legal reference

Issues

- **Fetch policy**
 - when are items put in the cache?
- **Placement policy**
 - where do they go in the cache?
- **Replacement policy**
 - what's removed to make room?

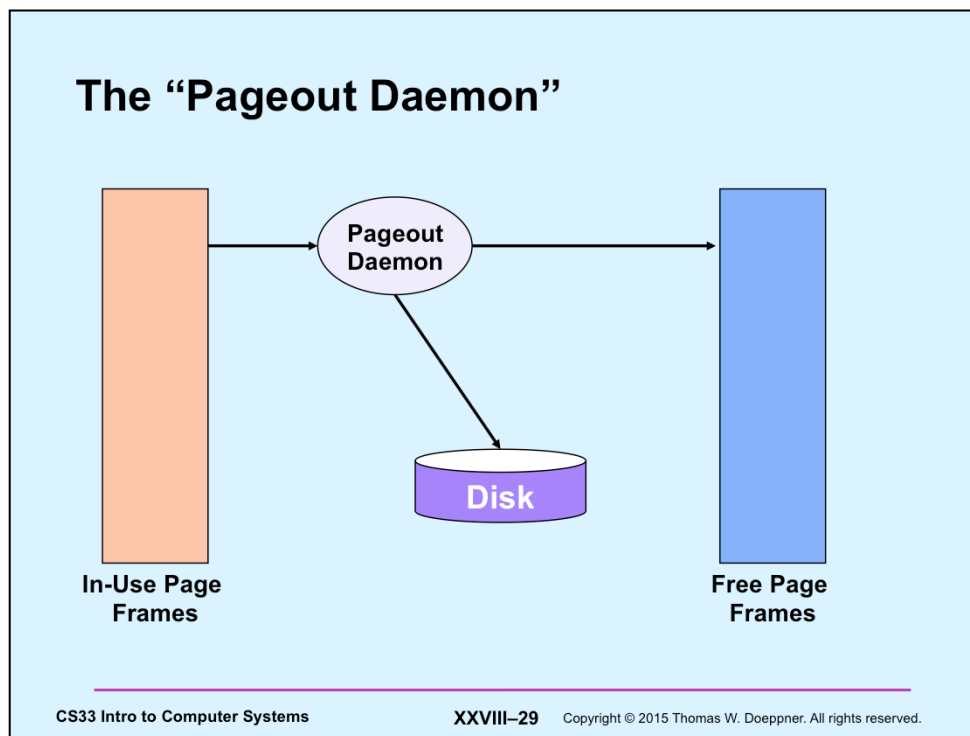
Three issues concerning the mechanism for caching are the following: the *fetch policy*, which governs when items are fetched to go into the cache, the *placement policy*, which governs where the fetched items are placed in the cache, and the *replacement policy*, which governs when and which items are removed from the cache (and perhaps written back to their source).

Hardware Caches

- **Fetch policy**
 - when are items put in the cache?
 - » when they're referenced
 - » prefetch might be possible (e.g., for sequential access)
- **Placement policy**
 - where do they go in the cache?
 - » usually determined by cache architecture
 - » if there's a choice, it's typically a random choice
- **Replacement policy**
 - what's removed to make room?
 - » usually determined by cache architecture
 - » if there's a choice, it's typically a random choice

Software Caches

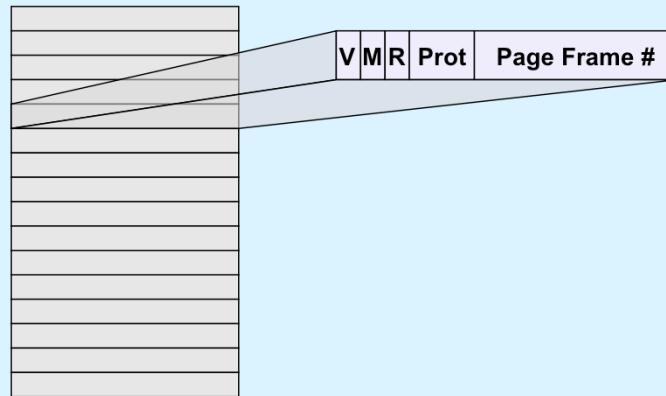
- **Fetch policy**
 - when are items put in the cache?
 - » when they're referenced
 - » prefetch might be easier than for hardware caches
- **Placement policy**
 - where do they go in the cache?
 - » usually doesn't matter (no memory is more equal than others)
- **Replacement policy**
 - what's removed to make room?
 - » would like to remove that whose next use is farthest in future
 - » instead, remove that whose last reference was farthest in the past



The (kernel) thread that maintains the free page-frame list is typically called the *pageout daemon*. Its job is to make certain that the free page-frame list has enough page frames on it. If the size of the list drops below some threshold, then the pageout daemon examines those page frames that are being used and selects a number of them to be freed. Before freeing a page, it must make certain that a copy of the current contents of the page exists on secondary storage. So, if the page has been modified since it was brought into primary storage (easily determined if there is a hardware-supported *modified bit*), it must first be written out to secondary storage. In many systems, the pageout daemon groups such pageouts into batches, so that a number of pages can be written out in a single operation, thus saving disk time. Unmodified, selected pages are transferred directly to the free page-frame list, modified pages are put there after they have been written out.

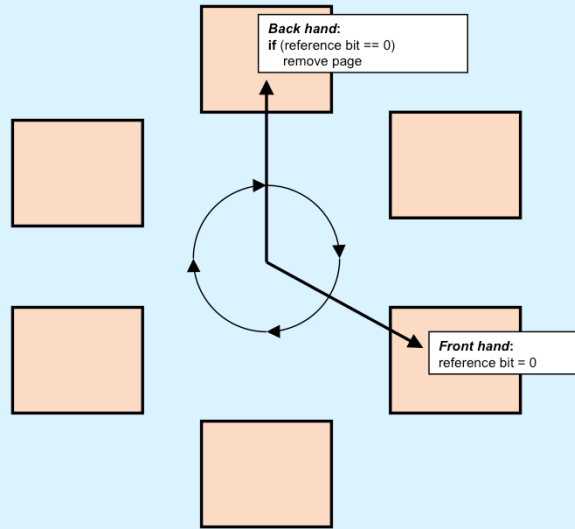
In most systems, pages in the free list get a “second chance” — if a thread in a process references such a page, there is a page fault (the page frame has been freed and could be used to hold another page), but the page-fault handler checks to see if the desired page is still in primary storage, but in the free list. If it is in the free list, it is removed and given back to the faulting process. We still suffer the overhead of a trap, but there is no wait for I/O.

Managing Page Frames



The OS can keep track of the history of page frame by use of two bits in each page-table entry: the *modify* bit, which is set by hardware whenever the associated page frame is modified, and the *referenced* bit, which is set by hardware whenever the associated page is accessed (via either a load or a store).

Clock Algorithm

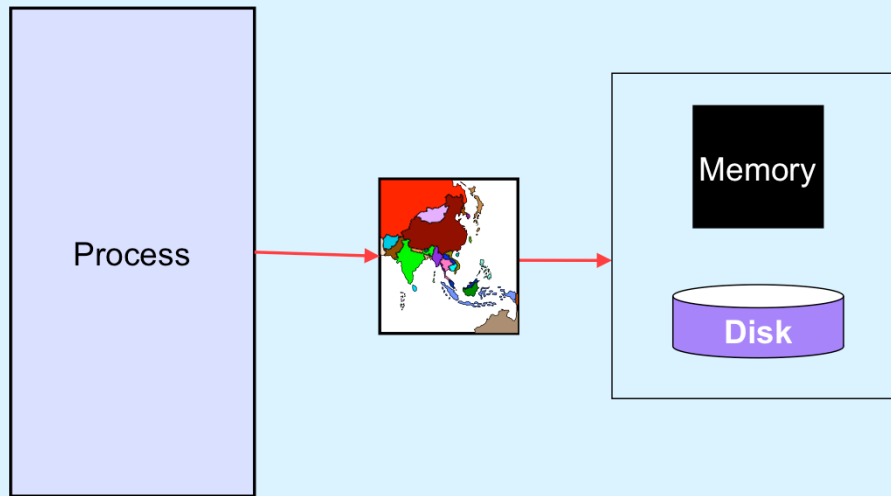


A common approach for determining which page frames are not in use is known as the clock algorithm. All active page frames are conceptually arranged in a circularly linked list. The page-out thread slowly traverses the list. The “one-handed” version of the clock algorithm, each time it encounters a page, checks the reference bit in the corresponding translation entry: if the bit is set, it clears it. If the bit is clear, it adds the page to the free list (writing it back to secondary storage first, if necessary).

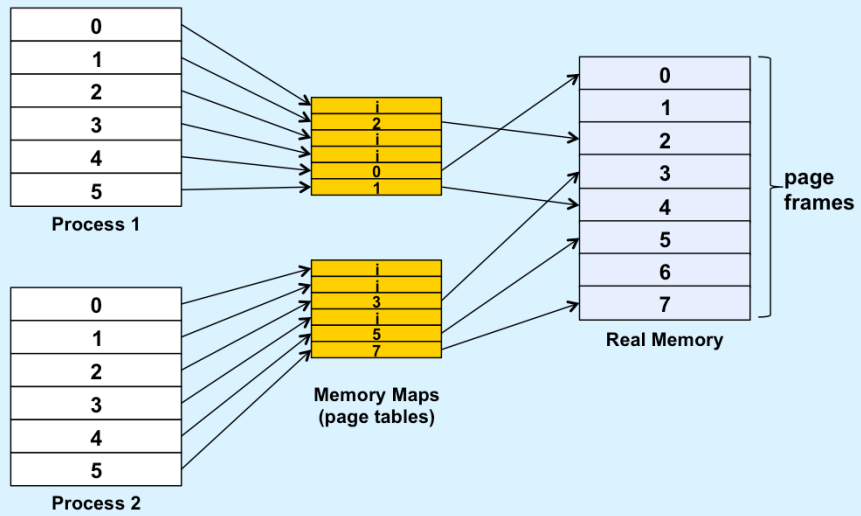
A problem with the one-handed version is that, in systems with large amounts of primary storage, it might take too long for the page-out thread to work its way all around the list of page frames before it can recognize that a page has not been recently referenced. In the two-handed version of the clock algorithm, the page-out thread implements a second hand some distance behind the first. The front hand simply clears reference bits. The second (back) hand removes those pages whose reference bits have not been set to one by the time the hand reaches the page frame.

Why is virtual memory used?

More VM than RM

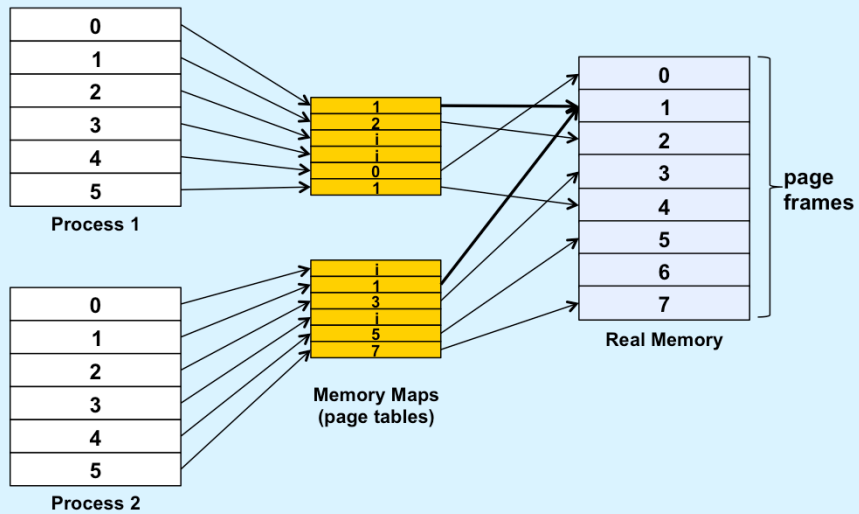


Isolation



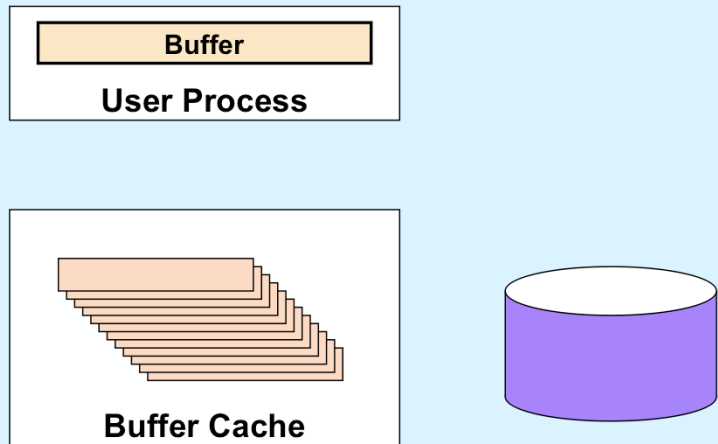
Virtual Memory

Sharing



Virtual Memory

File I/O

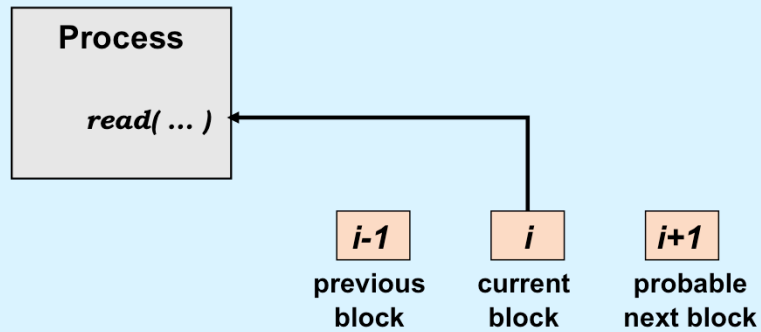


File I/O in Unix, and in most operating systems, is not done directly to the disk drive, but through intermediary buffers, known as the buffer cache, in the operating system's address space. This cache has two primary functions. The first, and most important, is to make possible concurrent I/O and computation within a Unix process. The second is to insulate the user from physical disk-block boundaries.

From a user process's point of view, I/O is *synchronous*. By this we mean that when the I/O system call returns, the system no longer needs the user-supplied buffer. For example, after a write system call, the data in the user buffer has either been transmitted to the device or copied to a kernel buffer — the user can now scribble over the buffer without affecting the data transfer. Because of this synchronization, from a user process's point of view, no more than one I/O operation can be in progress at a time.

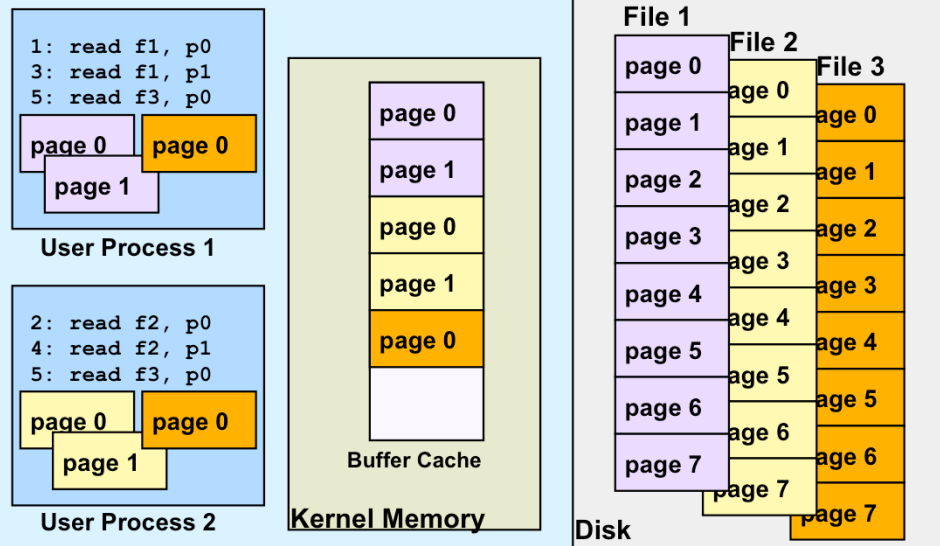
The buffer cache provides a kernel implementation of multibuffered I/O, and thus concurrent I/O and computation are made possible.

Multi-Buffered I/O

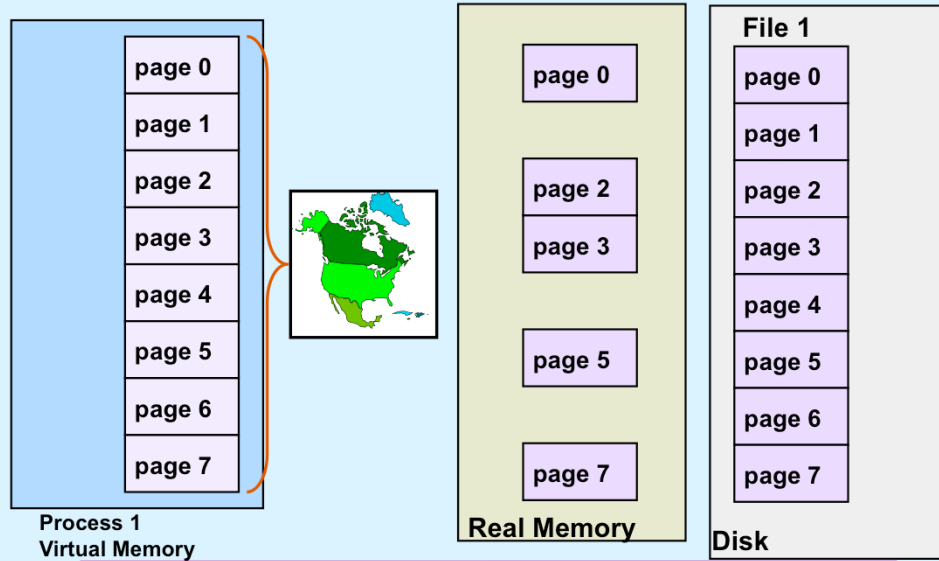


The use of *read-aheads* and *write-behinds* makes possible concurrent I/O and computation: if the block currently being fetched is block i and the previous block fetched was block $i-1$, then block $i+1$ is also fetched. Modified blocks are normally written out not synchronously but instead sometime after they were modified, asynchronously.

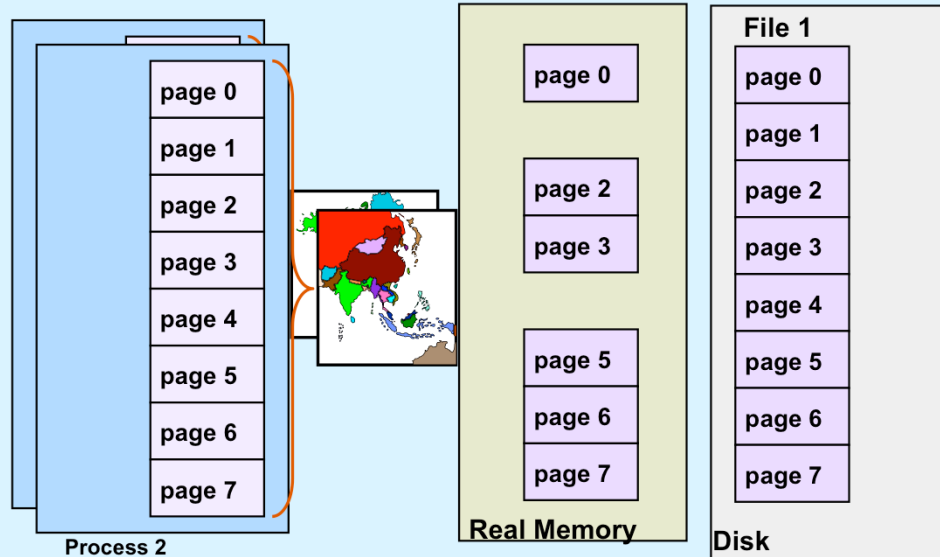
Traditional I/O



Mapped File I/O



Multi-Process Mapped File I/O



Mapped Files

- **Traditional File I/O**

```
char buf[BigEnough];
fd = open(file, O_RDWR);
for (i=0; i<n_recs; i++) {
    read(fd, buf, sizeof(buf));
    use(buf);
}
```

- **Mapped File I/O**

```
void *MappedFile;
fd = open(file, O_RDWR);
MappedFile = mmap(... , fd, ...);
for (i=0; i<n_recs; i++)
    use(MappedFile[i]);
```

Traditional I/O involves explicit calls to read and write, which in turn means that data is accessed via a buffer; in fact, two buffers are usually employed: data is transferred between a user buffer and a kernel buffer, and between the kernel buffer and the I/O device.

An alternative approach is to *map* a file into a process's address space: the file provides the data for a portion of the address space and the kernel's virtual-memory system is responsible for the I/O. A major benefit of this approach is that data is transferred directly from the device to where the user needs it; there is no need for an extra system buffer.

Mmap System Call

```
void *mmap(  
    void *addr,  
        // where to map file (0 if don't care)  
    size_t len,  
        // how much to map  
    int prot,  
        // memory protection (read, write, exec.)  
    int flags,  
        // shared vs. private, plus more  
    int fd,  
        // which file  
    off_t off,  
        // starting from where  
);
```

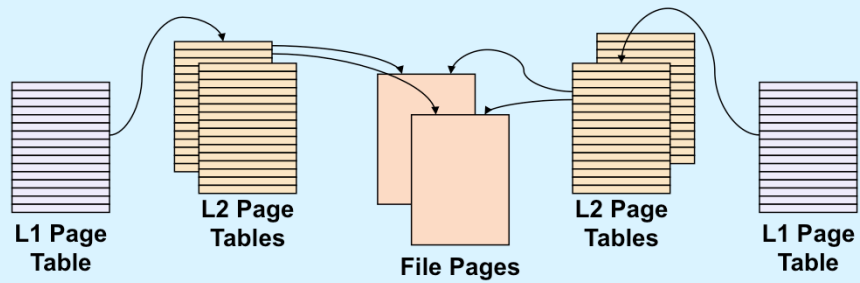
Mmap maps the file given by *fd*, starting at position *off*, for *len* bytes, into the caller's address space starting at location *addr*

- *len* is rounded up to a multiple of the page size
- *off* must be page-aligned
- if *addr* is zero, the kernel assigns an address
- if *addr* is positive, it is a suggestion to the kernel as to where the mapped file should be located (it usually will be aligned to a page). However, if *flags* includes `MAP_FIXED`, then *addr* is not modified by the kernel (and if its value is not reasonable, the call fails)
- the call returns the address of the beginning of the mapped file

The *flags* argument must include either `MAP_SHARED` or `MAP_PRIVATE` (but not both). If it's `MAP_SHARED`, then the mapped portion of the caller's address space contains the current contents of the file; when the mapped portion of the address space is modified by the process, the corresponding portion of the file is modified.

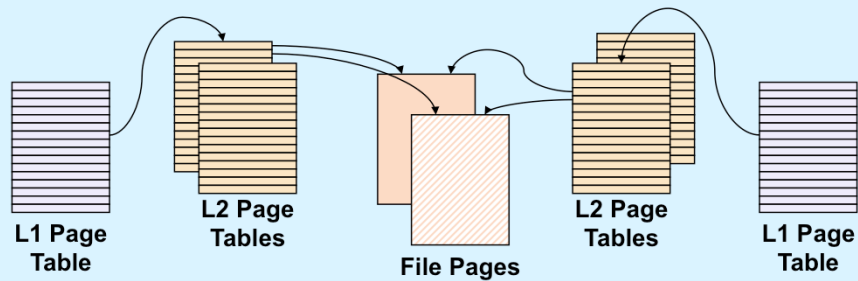
However, if *flags* includes `MAP_PRIVATE`, then the idea is that the mapped portion of the address space is initialized with the contents of the file, but that changes made to the mapped portion of the address space by the process are private and not written back to the file. The details are a bit complicated: as long as the mapping process does not modify any of the mapped portion of the address space, the pages contained in it contain the current contents of the corresponding pages of the file. However, if the process modifies a page, then that particular page no longer contains the current contents of the corresponding file page, but contains whatever modifications are made to it by the process. These changes are not written back to the file and not shared with any other process that has mapped the file. It's unspecified what the situation is for other pages in the mapped region after one of them is modified. Depending on the implementation, they might continue to contain the current contents of the corresponding pages of the file until they, themselves, are modified. Or they might also be treated as if they'd just been written to and thus no longer be shared with others.

The *mmap* System Call



The *mmap* system call maps a file into a process's address space. All processes mapping the same file can share the pages of the file.

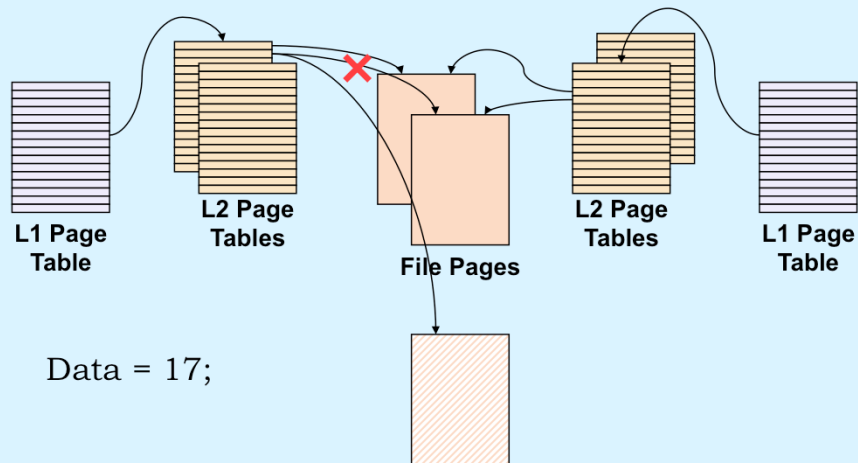
Share-Mapped Files



Data = 17;

There are a couple options for how modifications to mmaped files are dealt with. The most straightforward is the *share* option in which changes to mmaped file pages modify the file and hence the changes are seen by the other processes who have share-mapped the file.

Private-Mapped Files



The other option is to *private-map* the file: changes made to mmaped file pages do not modify the file. Instead, when a page of a file is first modified via a private mapping, a copy of just that page is made for the modifying process, but this copy is not seen by other processes, nor does it appear in the file.

In the slide, the process on the left has private-mapped the file. Thus its changes to the mapped portion of the address space are made to a copy of the page being modified.

Example

```
int main( ) {
    int fd;
    dataObject_t *dataObjecttp;

    fd = open("file", O_RDWR);
    if ((int)(dataObjecttp = (dataObject_t *)mmap(0,
        sizeof(dataObject_t),
        PROT_READ|PROT_WRITE, MAP_SHARED, fd, 0)) == -1) {
        perror("mmap");
        exit(1);
    }

    // dataObjecttp points to region of (virtual) memory
    // containing the contents of the file

    ...
}
```

Here we map the contents of a file containing a `dataObject_t` into the caller's address space, allowing it both read and write access. Note mapping the file into memory does not cause any immediate I/O to take place. The operating system will perform the I/O when necessary, according to its own rules.

fork and mmap

```
int main() {
    int x=1;

    if (fork() == 0) {
        // in child
        x = 2;
        exit(0);
    }
    // in parent
    while (x==1) {
        // will loop forever
    }
    return 0;
}

int main() {
    int fd = open( ... );
    int *xp = (int *)mmap(...,
        MAP_SHARED, fd, ...);
    xp[0] = 1;
    if (fork() == 0) {
        // in child
        xp[0] = 2;
        exit(0);
    }
    // in parent
    while (xp[0]==1) {
        // will terminate
    }
    return 0;
}
```

When a process calls `fork` and creates a child, the child's address space is normally a copy of the parent's. Thus changes made by the child to its address space will not be seen in the parent's address space (as shown in the left-hand column). However, if there is a region in the parent's address space that has been `mmap`d using the `MAP_SHARED` flag, and subsequently the parent calls `fork` and creates a child, the `mmap`d region is not copied but is shared by parent and child. Thus changes to the region made by the child will be seen by the parent (and vice versa).