

Memory Management in Operating Systems: Concepts, Techniques, and Implementations

1. Introduction: The Necessity of Memory Management

Memory management stands as a critical and intricate function within the realm of operating systems. Its primary role is to enable the concurrent execution of multiple processes without mutual interference, a cornerstone of modern computing.¹ This fundamental task encompasses the allocation, deallocation, and comprehensive oversight of a computer's memory resources.² Memory management orchestrates the operations between the main memory and the disk during the lifecycle of a process.³ The overarching aims of this process include the efficient utilization of the available memory³, the proper allocation and deallocation of memory before and after process execution³, meticulous tracking of memory space in use by various processes³, the minimization of memory fragmentation issues³, the preservation of data integrity while processes are running³, and the robust support for multiprogramming environments.³

Beyond these core objectives, memory management strives to optimize the overall performance of the system.¹ It plays a vital role in ensuring the security of data and processes by establishing isolation between them and enforcing strict access permissions to prevent unauthorized memory access.¹ Furthermore, it offers a simplified and convenient abstraction of the system's memory for programmers and compilers, thereby streamlining the software development process.⁴ The effective management of memory is indispensable for achieving the desired level of multiprogramming and ensuring that memory resources are utilized appropriately.³ In a multiprogramming system, where the main memory is shared among numerous processes, efficient memory management becomes even more critical to maximize the utilization of the processors by maintaining a substantial pool of processes ready for execution.⁵

The significance of memory management lies in its direct impact on the efficiency of the operating system. The central processing unit (CPU) frequently accesses data stored in memory, and the efficiency of this access profoundly influences the speed at which the CPU can complete tasks and become available for new ones.¹ Proper memory allocation ensures that each process receives the necessary memory to execute in parallel, thereby maximizing system throughput.¹ Memory management techniques are specifically designed to reduce the number of memory access operations, which are known to be resource-intensive for the CPU.¹ In scenarios where the total memory required by running processes exceeds the capacity of the physical

RAM, memory management facilitates the use of secondary storage as an extension of the main memory through a process known as swapping, allowing a greater number of processes to be executed concurrently.³ Ultimately, the various memory management techniques employed by an operating system are crucial for the efficient allocation, utilization, and overall management of memory resources, thereby ensuring the smooth execution of programs and the optimal utilization of the system's memory.³

2. Address Spaces: Logical vs. Physical

In the realm of operating systems, the concept of memory addresses is fundamental to how programs interact with the computer's memory. Two primary types of addresses are involved: logical addresses and physical addresses.

2.1 Understanding Logical (Virtual) Addresses

A logical address, also commonly referred to as a virtual address, is an address generated by the CPU during the execution of a program.⁶ This address is from the perspective of the process itself and is relative to the program's designated address space.⁶ Programs operate within this logical address space, creating and utilizing virtual addresses without any direct knowledge of the actual physical locations in the computer's memory where the data will ultimately reside.¹ The entirety of the logical addresses that a program can generate constitutes its logical address space.⁶ This provides a conceptual view of memory for each process, presenting a contiguous range of addresses that typically starts from zero, irrespective of how the physical memory is organized.⁸

Logical addresses serve as a crucial abstraction for programmers. By working with these virtual addresses, developers can write code without needing to be concerned about the intricate details of the underlying physical memory organization or the specific memory addresses being used by other concurrently running processes. This simplifies the programming process significantly, allowing developers to focus on the logic of their applications rather than the low-level details of memory management.

2.2 Understanding Physical Addresses

In contrast to logical addresses, a physical address is the actual address of a location within the computer's main memory where data is physically stored.⁶ It precisely identifies the necessary data's position in the memory hardware.⁶ The set of all physical addresses that correspond to the logical addresses in a system's logical address space is known as the physical address space.⁶ Physical addresses directly correspond to the real locations within the computer's Random Access Memory

(RAM), and these are the addresses that the hardware, including the CPU and the Memory Management Unit (MMU), uses to access the actual contents stored in the memory chips.

2.3 The Memory Management Unit (MMU) and Address Translation

The translation between the logical addresses generated by the CPU and the physical addresses used by the memory hardware is performed by a dedicated hardware component known as the Memory Management Unit, or MMU.¹ This hardware is often integrated directly into the CPU package and plays a vital role in memory virtualization.¹⁰ The process of address translation is essential because the logical addresses that programs use are virtual in nature and do not inherently correspond to any specific physical memory location.⁶ The MMU acts as an intermediary, taking the logical address as input and producing the corresponding physical address that can be used to access the main memory.

In systems that employ paging, the MMU typically utilizes a data structure called a page table to perform this translation.⁶ The page table, which is maintained by the operating system, maps logical page numbers to physical frame numbers. When the CPU requests data at a particular logical address, the MMU consults the page table to determine the corresponding physical address. The page table itself contains critical information about the relationship between the logical and physical address spaces.⁶

The MMU is also capable of generating a fault, which is a type of exception, if it encounters an issue during the memory access process.⁹ Such faults can occur due to various reasons, including attempts to access memory locations that are illegal or to access pages that are not currently present in the physical memory. These faults allow the operating system to intervene and handle the exceptional condition appropriately.

To further enhance the efficiency of address translation, most modern MMUs incorporate a small, high-speed cache known as the Translation Lookaside Buffer, or TLB.⁹ The TLB stores recently used virtual-to-physical address translations, allowing the MMU to quickly retrieve these mappings without having to access the page table in main memory for every memory reference. This caching mechanism significantly speeds up the overall address translation process.

In the context of contiguous memory allocation, the MMU often employs a relocation register, also known as a base register, which holds the starting physical address of the currently running process in memory.¹² The MMU adds the value in this register to every logical address generated by the process to obtain the corresponding physical address.¹² Additionally, a limit register is used to specify the range of memory that the

process is allowed to access, providing a mechanism for memory protection.¹² The MMU checks every generated logical address against the value in the limit register to ensure that the process does not attempt to access memory outside of its allocated region.

The MMU serves as a vital hardware component that transparently translates the abstract logical addresses used by software into the concrete physical addresses required by the underlying hardware. This translation process is fundamental to modern operating systems, as it enables memory virtualization, allowing multiple processes to share the physical memory safely and efficiently, and provides a crucial layer of protection between different processes and the operating system itself.

2.4 Address Binding: Compile-time, Load-time, and Execution-time

Address binding is the process of establishing a mapping from one address space to another, specifically from logical addresses to physical addresses.⁷ This mapping can occur at different stages of a program's lifecycle: compile time, load time, or execution time.

When address binding occurs at **compile time**, if the memory location where a process will reside is known in advance, the compiler can generate absolute addresses, which are essentially physical addresses, directly into the program's executable code.⁷ This approach results in very fast loading of the executable into memory. However, it is also the least flexible. If the intended memory space is already occupied when the program is loaded, the program will likely crash, and it would require recompilation with a different virtual address space.

If the memory location is not known at compile time, the compiler generates relocatable addresses. In this case, address binding is delayed until **load time**.⁷ The loader, a part of the operating system responsible for loading the program into memory, then translates these relocatable addresses into absolute addresses by adding the base address of the process in main memory to each logical address. While more flexible than compile-time binding, if the base address of the process needs to change, the process must be reloaded into memory.

Finally, address binding can be postponed until **execution time**.⁷ In this scenario, the instructions are loaded into memory with virtual addresses, and the mapping to physical addresses is done dynamically during program execution by the MMU. This approach offers the greatest flexibility, as it allows a process to be moved from one memory location to another even while it is running. This is particularly useful in techniques like dynamic linking and memory compaction. However, it requires the

presence of hardware support, specifically the MMU, to perform the address translation at runtime.

The timing of address binding has significant implications for the flexibility and efficiency of memory management. Early binding, such as at compile time, offers simplicity and speed but lacks the adaptability to changing memory conditions. Later binding, especially at execution time, provides the flexibility needed for modern operating systems to manage memory effectively in multiprogramming environments, albeit with the added complexity of requiring hardware support for dynamic address translation.

3. Contiguous Memory Allocation: Allocating in Blocks

Contiguous memory allocation is a memory management technique where a process is allocated a single, continuous block of memory in the main memory.¹³ This block of memory is provided to the process as one contiguous segment from the available free space, based on the size requirements of the process.¹³ Contiguous memory allocation can be broadly categorized into two main schemes: fixed partitioning and variable partitioning.

3.1 Fixed Partitioning (Static Partitioning)

3.1.1 Concept and Implementation

In the fixed partitioning scheme, also known as static partitioning, the main memory is divided into a predetermined number of partitions at the system's startup.³ Each of these partitions represents a contiguous block of memory.³ The sizes of these partitions can be either equal or unequal¹⁵, but once established, the number and boundaries of these partitions remain unchanged throughout the system's operation.¹⁵ A key characteristic of fixed partitioning is that each partition is designed to hold exactly one process.³ The allocation of these fixed-size memory blocks is typically performed at the system's boot time or during its initial configuration.¹⁵ To manage the memory within each partition, the operating system often associates limit registers with each partition, defining the lower and upper address boundaries that the process within that partition can access.¹⁴ The degree of multiprogramming in a system using fixed partitioning is directly limited by the number of partitions available in the memory.¹⁴

3.1.2 Internal Fragmentation: Causes and Impact

A significant drawback of fixed partitioning is the occurrence of internal fragmentation.² This phenomenon arises when a process is allocated a partition that is

larger than its actual memory requirements, resulting in unused memory space within that partition.¹⁵ For instance, if a process requires only 1MB of memory but is placed into a 4MB partition, the remaining 3MB within that partition goes to waste.¹⁵ This leads to an overall inefficient utilization of the main memory.¹⁵ The total internal fragmentation across the entire memory can be calculated by summing up the unused space within each allocated partition.¹⁵

3.1.3 External Fragmentation: The Challenge of Non-Contiguous Free Space

Another challenge associated with fixed partitioning, though debated in some sources, is external fragmentation.¹⁵ While some argue that external fragmentation is avoided because processes cannot span across multiple fixed partitions¹⁵, others contend that it can occur in a specific context. This perspective suggests that even if the total amount of unused memory across all partitions is sufficient to accommodate a new process, if no single contiguous partition is large enough for that process, the memory cannot be utilized.¹⁵ This is particularly relevant when a process's size exceeds the size of any individual partition available, even if the cumulative free space might be adequate.¹⁵ The rule of contiguous allocation, where a process must reside in a single, unbroken block of memory, prevents the utilization of scattered free space.¹⁵

3.1.4 Advantages and Disadvantages of Fixed Partitioning

Fixed partitioning offers several advantages. It is relatively simple to implement², with straightforward algorithms for placing processes into partitions.¹⁶ It also results in low operating system overhead¹⁵, requiring less computational power for memory management.¹⁶ Memory allocation is predictable¹⁵, and the operating system can ensure a minimum amount of memory for each process.¹⁵ According to some sources, it avoids external fragmentation because processes cannot span partitions.¹⁵ It is well-suited for systems where the number of processes and their memory requirements are known in advance¹⁵ and prevents processes from interfering with each other's memory space.¹⁵ Debugging is also easier as the size and location of each process are predetermined.¹⁵ Additionally, it can be beneficial for batch processing¹⁵ and offers better control over memory allocation.¹⁵ In specific scenarios, it might even prevent data loss during power outages and allow the installation of multiple operating systems.³¹

However, fixed partitioning also has significant disadvantages. The primary issue is internal fragmentation, leading to inefficient use of main memory.² It limits the size of processes that can be accommodated to the largest partition size¹⁵ and restricts the degree of multiprogramming to the number of available partitions.¹⁴ The partition size cannot be dynamically varied according to the needs of incoming processes.¹⁵

Moreover, as discussed, it can potentially suffer from external fragmentation in the sense that a process might not be allocated memory even if sufficient total free space exists across multiple partitions, if no single partition is large enough.¹⁵

3.2 Variable Partitioning (Dynamic Partitioning)

3.2.1 Concept and Implementation

Variable partitioning, also known as dynamic partitioning, takes a different approach to memory allocation. In this scheme, partitions are not predefined or created at system startup. Instead, they are created dynamically at runtime, specifically when a process needs to be loaded into memory.³² The size of each partition is tailored to match the exact amount of memory required by the incoming process.¹³ Initially, the entire main memory is considered as a single, large contiguous block of free space.¹⁴ As processes arrive and request memory, the operating system allocates them a block of memory that is precisely the size they need, provided that a sufficiently large contiguous block is available.¹³ The number and size of partitions in this scheme are not fixed and depend on the sequence of processes entering and leaving the memory, as well as their individual memory requirements.¹⁸ The operating system maintains a record of all free memory blocks, often referred to as "holes," in a table.³³ When a new process requests memory, the OS searches this table for a hole that is large enough to accommodate the process. Several strategies can be employed to select a suitable hole, including First-Fit (allocating the first hole that is large enough), Best-Fit (allocating the smallest hole that is large enough), and Worst-Fit (allocating the largest hole that is large enough).¹³

3.2.2 External Fragmentation: The Problem of Scattered Holes

While variable partitioning effectively eliminates internal fragmentation in its ideal form, it introduces the problem of external fragmentation.¹⁸ This occurs over time as processes are allocated memory and subsequently deallocate it upon completion. This cycle of allocation and deallocation can lead to the memory becoming fragmented into numerous smaller, non-contiguous blocks of free memory, or "holes".¹⁸ If a new process arrives that requires a contiguous block of memory larger than any of the individual free blocks available, the operating system will be unable to allocate memory to it, even if the total amount of free memory across all the holes is sufficient to satisfy the request.²¹ This is because the contiguous allocation rule mandates that a process must occupy a single, unbroken segment of main memory to be executed.³²

3.2.3 Internal Fragmentation: Typically Absent in Pure Variable Partitioning

In its ideal implementation, variable partitioning does not suffer from internal fragmentation.¹³ Since memory is allocated strictly according to the process's needs, there is no leftover unused space within the allocated partition.¹³ The size of the allocated portion is precisely equal to the size required by the process.¹⁴ However, it is worth noting that in some practical implementations, especially when using certain allocation strategies like rounding up the allocated size to a specific block size for management purposes, a small amount of internal fragmentation might still occur.

3.2.4 Advantages and Disadvantages of Variable Partitioning

Variable partitioning offers several advantages over fixed partitioning. The most significant is the absence of internal fragmentation (in its ideal form), which leads to more efficient utilization of the main memory.¹³ It also allows for a higher degree of multiprogramming, as more processes can be accommodated in memory due to the lack of wasted space.¹⁴ Furthermore, there is no inherent limitation on the size of a process that can be loaded, as long as a contiguous block of sufficient size is available in memory.¹⁴ The dynamic nature of allocation allows for flexibility in accommodating changing memory requirements of processes.⁴⁰

However, variable partitioning also has its drawbacks. Implementing it is more complex compared to fixed partitioning, as it involves memory allocation during runtime rather than at system configuration.³² The primary disadvantage is external fragmentation, which can lead to inefficient memory utilization over time as free memory becomes scattered.² Managing variable-sized partitions requires additional overhead for maintaining and searching the memory table of free blocks.⁴⁰ Inefficient allocation and deallocation can exacerbate fragmentation issues⁴⁰, and more sophisticated algorithms might be needed to manage memory fragmentation effectively.⁴⁰

3.3 Compaction: A Solution to External Fragmentation

Compaction is a memory management technique employed to address the problem of external fragmentation in contiguous memory allocation.¹⁴ The fundamental idea behind compaction is to consolidate all the free memory fragments, which might be scattered across the main memory, into one large, contiguous block of free space.¹⁴ This is achieved by relocating all the currently occupied memory blocks to one end of the memory (either towards the top or the bottom), thereby bringing all the available free space together into a single, usable chunk.¹⁴ Compaction can only be effectively performed when the relocation of processes is dynamic, meaning it can be done at execution time.⁴⁵ A key challenge in implementing compaction is that any pointers within the moved processes that refer to memory locations also need to be updated to reflect the new physical addresses of the blocks.⁴³ The process of compaction

typically requires interrupting all the running processes in the memory, which can lead to system downtime.¹⁴ It also consumes a significant amount of CPU time, adding overhead to the system's operation.¹⁴ Furthermore, the process of moving memory blocks around can potentially lead to deadlocks in certain scenarios.⁴⁵

Despite these implementation challenges, compaction offers notable advantages. It effectively reduces external fragmentation³⁸, making memory usage more efficient.⁴⁵ By creating a large contiguous block of free memory⁴⁵, it becomes possible to load processes that require larger amounts of contiguous memory, thereby increasing the scalability of the operating system.⁴⁵ It can also temporarily alleviate fragmentation in the file system.⁴⁵ However, these benefits come at the cost of reduced system efficiency and increased latency.⁴⁵ The substantial time spent performing compaction can leave the CPU idle for extended periods, impacting overall system responsiveness.⁴⁵ Moreover, performing compaction is not always a straightforward task.⁴⁵

4. Paging: Dividing Memory into Fixed-Size Units

Paging is a memory management scheme that fundamentally changes how memory is allocated to processes. Unlike contiguous memory allocation, paging eliminates the requirement for a process to occupy a single, continuous block of physical memory.⁴⁸ Instead, it divides both the logical memory of a process and the physical memory of the computer into fixed-size units.

4.1 Principle of Operation: Pages and Frames

In a paging system, the logical memory space of each process is divided into equal-sized blocks known as pages.² Similarly, the physical memory (RAM) is divided into equal-sized blocks called frames.² A crucial aspect of paging is that the size of a page is always equal to the size of a frame.¹¹ When a process needs to be loaded into memory, its logical pages are mapped to available frames in the physical memory.² These frames do not need to be contiguous in physical memory; a process's pages can be scattered across various available frames.¹¹

To keep track of the mapping between the logical pages of a process and the physical frames they occupy, the operating system maintains a data structure called a page table.² This page table is used by the Memory Management Unit (MMU) to perform the translation of logical addresses generated by the CPU into the corresponding physical addresses in memory.¹¹

When the CPU generates a logical address, this address is typically divided into two

parts: a page number and an offset within the page.⁴⁹ The page number serves as an index into the process's page table. By using this index, the MMU can locate the entry in the page table that corresponds to the requested logical page. This entry contains the frame number of the physical frame where the page is stored.⁵⁰ The MMU then combines this frame number with the offset from the original logical address to construct the final physical address, which is used to access the desired location in the main memory.⁵⁰

4.2 Page Allocation and Page Tables

When a process requires memory to execute, the operating system allocates one or more page frames in the physical memory to that process.¹¹ For each running process, the OS creates and maintains a page table, which is stored in the main memory.¹¹ Each entry in the page table, known as a Page Table Entry (PTE), contains crucial information about the corresponding page of the process. This information typically includes the frame number of the physical frame that the page is currently occupying, as well as various status bits.¹¹ These status bits might indicate whether the page is currently valid (present in RAM), whether it has been accessed or modified, and what protection permissions are associated with it.

For systems with very large virtual address spaces, a single-level page table might become excessively large and consume a significant amount of memory. To address this, operating systems often employ multilevel page tables.⁵⁴ These hierarchical structures break down the page table into multiple levels, allowing for a more efficient use of memory, especially when large portions of the virtual address space are not actively used. Another approach, known as an inverted page table, uses a different indexing mechanism. Instead of having one entry per virtual page (as in traditional page tables), an inverted page table has one entry for each physical frame in memory, which then points back to the virtual page that is currently occupying that frame.⁶⁴

The allocation of pages to frames in physical memory can be either contiguous or non-contiguous.¹¹ The flexibility of non-contiguous allocation is one of the key advantages of paging, as it allows the operating system to utilize any available frame in memory, regardless of its location relative to other frames allocated to the same process.

4.3 Hardware Support for Paging: Translation Lookaside Buffer (TLB)

Paging relies heavily on hardware support for efficient operation. One of the most critical hardware components in this context is the Translation Lookaside Buffer (TLB).⁹ The TLB is a small, high-speed associative memory, essentially a cache, that

stores recently accessed page table entries.⁹ Its primary purpose is to speed up the process of address translation by reducing the frequency with which the MMU needs to access the main memory to retrieve page table entries.⁹ Accessing the TLB is significantly faster than accessing the main memory.

Each entry in the TLB typically consists of two parts: a tag, which usually contains the page number, and a value, which contains the corresponding frame number and other status bits from the page table entry.¹¹ When the CPU generates a logical address, the MMU first checks if the translation for the requested page number is present in the TLB. If it is, this is known as a TLB hit, and the MMU can quickly obtain the physical address without needing to consult the page table in main memory. However, if the translation is not found in the TLB, this is a TLB miss.¹¹ In the event of a miss, the MMU must then perform a slower lookup in the page table in main memory to get the required translation. Once the translation is retrieved from the page table, it is typically stored in the TLB, potentially replacing an older entry, so that subsequent accesses to the same page can benefit from the faster TLB lookup.

4.4 Memory Protection in Paging Systems

Paging provides a robust mechanism for memory protection in operating systems.⁵⁰ This is primarily achieved through the use of protection bits that are included in each page table entry.⁵⁰ These bits allow the operating system to control the types of access that are permitted for each page in memory. For example, protection bits can specify whether a page can be read from, written to, or executed. Additionally, they can define whether a page can be accessed only by the operating system (supervisor mode) or also by user-level processes (user mode).⁵⁶

Because each process in a paging system has its own page table, this inherently provides isolation between the address spaces of different processes.⁶¹ One process cannot access the memory of another process unless the operating system explicitly sets up the page tables to allow for shared memory regions. During the address translation process, the MMU checks the protection bits associated with the page being accessed.⁵⁶ If a process attempts to perform an operation that is not permitted by the protection bits (e.g., writing to a read-only page, or accessing a kernel-only page from user mode), the MMU will trigger a protection fault. This fault is then handled by the operating system, which can take appropriate action, such as terminating the offending process, thereby preventing unauthorized access to memory and enhancing the overall security and stability of the system.⁵⁰

4.5 Memory Sharing through Paging

While paging provides strong memory protection, it also offers mechanisms for controlled memory sharing between processes.⁶⁸ Since the physical memory used by a process does not need to be contiguous in a paging system, the operating system can allow multiple processes to map the same physical page into their respective virtual address spaces.⁶¹ This is particularly useful for sharing code libraries or other read-only data segments among multiple running instances of the same program or different programs that utilize the same resources.⁷⁵ By sharing physical memory in this way, the operating system can reduce the overall memory footprint of the system, as only one copy of the shared resource needs to be loaded into RAM.

For memory sharing to work correctly, especially in the case of shared code, the code must be reentrant.⁶⁸ This means that the code should be stateless and not modify itself during execution, ensuring that multiple processes can execute it concurrently without interfering with each other's operation. In situations where a shared page needs to be modified by a process, a technique called copy-on-write can be employed.⁵⁶ With copy-on-write, when a process attempts to write to a page that is currently shared and marked as read-only, the operating system intercepts this write attempt, creates a private copy of the page for the writing process, and then allows the write operation to proceed on the newly created copy. The original shared page remains unchanged and available for other processes.

4.6 Disadvantages of Paging: Overhead and Potential Internal Fragmentation

Despite its numerous advantages, paging also has some drawbacks. One potential issue is internal fragmentation.² Since pages are of a fixed size, if a process's memory requirement is not an exact multiple of the page size, the last page allocated to that process might not be fully utilized, leading to wasted memory within that page.¹¹ While this internal fragmentation is generally less severe than the external fragmentation that can occur in contiguous allocation schemes, it still represents an inefficient use of memory resources.

Another significant disadvantage of paging is the overhead associated with managing page tables.² For each process, the operating system needs to maintain a page table, which itself occupies memory space. In systems with large virtual address spaces, these page tables can become quite substantial, leading to a considerable memory overhead. Multilevel paging is a technique used to mitigate this overhead by breaking down the page table into a hierarchical structure, but it adds complexity to the address translation process.⁵⁴

The process of address translation in a paging system can also lead to an increase in memory access time.² For every memory access made by a process, the system

typically needs to perform a lookup in the page table to translate the virtual address to a physical address. This adds an extra level of indirection to the memory access, potentially slowing down the execution of processes. However, this performance impact is significantly reduced by the use of the Translation Lookaside Buffer (TLB), which caches frequently used page table entries, allowing for fast address translation in most cases.²

Furthermore, the implementation of paging is more complex compared to simpler contiguous allocation techniques.² The operating system needs to manage the page tables, handle page faults, and implement page replacement algorithms, all of which add to the complexity of the system. Overall, while paging offers significant advantages in terms of memory management and enabling virtual memory, it does come with certain overheads related to memory usage, access time, and implementation complexity.

Table 1: Advantages and Disadvantages of Paging

Advantage	Disadvantage
Eliminates external fragmentation	Internal fragmentation can occur
Allows non-contiguous allocation of physical memory	Page table overhead requires additional memory
Supports virtual memory, allowing programs larger than physical memory	Increased memory access time due to page table lookup (mitigated by TLB)
Simplifies memory management by using fixed-size units (pages and frames)	Complexity in implementation compared to contiguous allocation
Efficient swapping of pages between RAM and secondary storage	Overhead of managing page tables and handling page faults
Facilitates memory protection and sharing between processes	
Allows for better memory utilization	

5. Virtual Memory: Extending Memory Beyond Physical Limits

Virtual memory is a sophisticated memory management technique employed by operating systems to provide applications with the illusion of having access to a vast, contiguous block of memory, often exceeding the amount of physical memory (RAM) actually installed in the system.² This powerful abstraction allows programs to run even if their total memory requirement is larger than the available RAM.⁴ The fundamental principle behind virtual memory is the use of secondary storage, typically a hard disk or a solid-state drive, as an extension of the physical RAM, in what is commonly referred to as swap space or a paging file.²

Virtual memory is crucial for supporting multiprogramming, as it enables a greater number of processes to reside in the main memory concurrently.⁴ The operating system, in conjunction with hardware support, manages the mapping between the virtual addresses used by programs and the physical addresses in computer memory.¹ This abstraction simplifies program development, as programmers do not need to be concerned with the limitations of physical memory and can program as if they have access to a single, large, contiguous block of memory.⁶²

5.1 Hardware and Control Structures Supporting Virtual Memory

The implementation of virtual memory relies on a close collaboration between hardware and software components. The Memory Management Unit (MMU) is a key hardware component responsible for translating virtual addresses generated by the CPU into physical addresses in memory.⁴ Page tables, maintained by the operating system, are essential control structures that store the mappings between virtual pages and physical frames.² Swap space, or the paging file, located on secondary storage, serves as the repository for pages that are not currently in active use in RAM.² The Translation Lookaside Buffer (TLB), a cache within the MMU, further supports virtual memory by storing frequently used page table entries, thereby accelerating the address translation process.⁴ Additionally, the CPU utilizes specific registers to point to the currently active page table, such as the CR3 register in x86 architectures.⁹

5.2 Locality of Reference: Temporal and Spatial Locality

The efficiency of virtual memory systems heavily relies on the principle of locality of reference.² This principle observes that during a short period, a processor tends to access the same set of memory locations repeatedly (temporal locality) and also tends to access memory locations that are physically close to each other (spatial locality).² Programs often exhibit strong locality due to the presence of loops, sequential execution of instructions, and the use of data structures like arrays, where elements are stored in adjacent memory locations.⁸¹ Virtual memory systems leverage

this behavior by attempting to keep the pages that are being frequently accessed in RAM, thereby reducing the need to access the slower secondary storage.⁴ Caching mechanisms, including the TLB, also operate based on the principle of locality, storing recently accessed data in faster memory to improve performance.⁹

5.3 Page Faults: Handling Missing Pages

A page fault is a critical event that occurs in a virtual memory system when a program attempts to access a virtual address that corresponds to a page that is not currently loaded in the physical memory (RAM).² When a page fault occurs, the MMU detects this and generates an interrupt, which causes the control to be transferred to the operating system.⁹ The OS then checks if the memory reference that caused the fault is valid.⁹⁴ If the reference is to an invalid address, the operating system might terminate the process.¹⁰⁹ However, if the reference is valid, it means that the requested page is located on the secondary storage (swap space).⁹⁴ The OS then needs to find a free frame in the physical memory.⁹⁴ If no frame is immediately available, the OS might need to evict a page that is currently in memory using a page replacement algorithm.⁹⁴ Once a free frame is obtained, the operating system initiates a disk operation to load the required page from the secondary storage into the allocated frame in RAM.⁹⁴ After the page is loaded, the operating system updates the page table for the process to reflect the new location of the page in physical memory.⁹⁴ Finally, the instruction that caused the page fault is restarted, allowing the process to continue its execution.⁹⁴ Page faults can be categorized as minor or soft faults if the page is already in memory but not correctly mapped, or major or hard faults if the page needs to be retrieved from the disk.¹⁰⁹

5.4 The Working Set Model: Tracking Active Pages

The working set of a process, in the context of virtual memory, refers to the set of pages in the process's virtual address space that are currently residing in the physical memory.² It is often considered to be the smallest collection of pages that must be present in memory for the process to execute efficiently.¹²⁴ The working set model posits that a process should be allowed to remain in RAM only if all the pages it is currently using can be accommodated in RAM.¹²⁰ Operating systems strive to keep the working set of active processes in the physical memory to minimize the occurrence of page faults.⁸¹ To approximate the current working set of a process, the operating system often maintains a working set window, which can be a specific time interval or a certain number of the most recent memory references made by the process.¹²¹ The set of unique pages referenced within this window is considered to be the process's current working set. If the number of physical frames allocated to a process is less than the size of its working set, a condition known as thrashing can occur, where the

process spends an excessive amount of time paging.¹¹⁹

5.5 Dirty Pages and the Dirty Bit: Managing Modified Data

A dirty page, also known as a modified page, is a page in the physical memory (RAM) that has been altered (written to) since it was initially loaded from the secondary storage.² To keep track of whether a page has been modified, a dirty bit, also called a modified bit or write bit, is associated with each page in the page table.² This bit is set by the hardware whenever a write operation is performed to any location within the page.⁹³ The information provided by the dirty bit is crucial when a page needs to be replaced (evicted from RAM) to make space for a new page.⁵⁷ If the dirty bit of a page is set, it indicates that the page in RAM contains modifications that are not yet reflected on the secondary storage. In this case, the page must be written back to the disk before the physical frame it occupies can be reallocated to another page. Conversely, if the dirty bit is not set, it means that the page in RAM has not been modified since it was loaded (or that its contents are the same as the copy on disk). Therefore, if a clean (not dirty) page needs to be replaced, the operating system can simply discard it from RAM without needing to write it back to disk, as the current version is already safely stored in secondary storage. The use of the dirty bit helps to optimize the page replacement process by avoiding unnecessary disk writes, thus improving the overall performance and potentially extending the lifespan of storage devices.⁹³

5.6 Demand Paging: Loading Pages on Demand

Demand paging is a widely used method of virtual memory management where pages of a process are loaded into the physical memory (RAM) only when they are actually needed, i.e., when they are demanded by the process during its execution.² When a process is first started, none of its pages are initially loaded into RAM.⁵² Instead, all the pages of the process are kept in the secondary storage. As the process begins to execute and attempts to access a memory location that resides on a page not currently in RAM, a page fault occurs.⁹⁴ This triggers the operating system to load the required page from the disk into an available frame in the physical memory. This approach is often referred to as lazy loading because pages are only brought into memory when they are absolutely necessary.⁹⁴

One of the significant advantages of demand paging is that it reduces the initial loading time of a program, as the operating system does not need to load the entire program into memory before execution begins.⁹⁴ It also reduces the amount of physical memory required by a process at any given time, as only the pages that are actively being used need to be in RAM.⁹⁴ This can lead to a higher degree of

multiprogramming, as more processes can reside in memory simultaneously. However, a potential drawback of demand paging is the increased latency experienced when a process first tries to access a page that is not in memory, as this requires fetching the page from the slower secondary storage, resulting in a page fault.⁹⁴ In pure demand paging, no pages are loaded into memory until they are explicitly referenced by the process.⁹⁴

5.7 Page Replacement Algorithms: Managing Limited Frames

When a new page needs to be brought into the physical memory (RAM), but all the available frames are already occupied, the operating system must decide which page currently in memory should be replaced (evicted) to make space for the new page. The algorithms that govern this decision are known as page replacement algorithms.⁵² The goal of these algorithms is to minimize the number of page faults, which occur when a process needs a page that is not in RAM. Several page replacement algorithms exist, each with its own strategy.

The **Optimal Page Replacement Algorithm** is a theoretical algorithm that replaces the page whose next use will occur farthest in the future.² While it results in the minimum number of page faults, it is not practically implementable in a general-purpose operating system because it requires knowing the future memory access pattern of a program. However, it serves as a benchmark against which other algorithms are compared.¹³⁶

The **First-In, First-Out (FIFO) Algorithm** is one of the simplest algorithms, which replaces the page that has been in memory for the longest time.² While easy to implement, FIFO can sometimes replace frequently used pages, and it can suffer from Belady's anomaly, where increasing the number of frames might paradoxically lead to more page faults in certain scenarios.¹³⁶

The **Second Chance (SC) Algorithm** is an enhancement of FIFO.² It works by looking at the front of the FIFO queue and checking the reference bit of the page. If the reference bit is set (meaning the page has been recently used), the bit is cleared, and the page is moved to the end of the queue, effectively giving it a "second chance" before being replaced.

The **Not Recently Used (NRU) Algorithm** favors keeping pages in memory that have been recently accessed.² It categorizes pages based on their reference bit (indicating if they have been referenced in the current clock interval) and their modified bit (indicating if they have been modified). Pages are divided into four classes, and NRU

replaces a random page from the lowest priority non-empty class.

The **Least Recently Used (LRU) Algorithm** replaces the page that has not been used for the longest period.² It operates on the principle that pages that have been used recently are likely to be used again in the near future, while pages that have not been used for a long time are less likely to be needed. LRU is often more efficient than FIFO, but it requires keeping track of the usage history of each page, which can be complex and costly to implement.¹³⁶

Other page replacement algorithms include Last-In, First-Out (LIFO), which replaces the most recently brought in page, Least Frequently Used (LFU), which replaces the page that has been accessed the least number of times, and Random replacement, which chooses a page to replace arbitrarily.¹³⁶ The choice of page replacement algorithm can significantly impact the performance of a virtual memory system, and operating systems often employ algorithms that aim to strike a balance between efficiency and implementation complexity.

5.8 Thrashing: Excessive Paging and Performance Degradation

Thrashing is a severe performance problem that can occur in virtual memory systems.² It happens when a computer system spends an excessive amount of time and resources on the task of swapping pages between the physical memory (RAM) and the secondary storage (disk), to the detriment of executing the actual work of the running applications.² This condition typically arises when the system's real storage resources are overcommitted, leading to a constant state of paging and page faults.¹⁴¹

Thrashing can be caused by several factors, including insufficient physical memory, a high degree of multiprogramming (too many processes trying to run simultaneously), and programs that exhibit poor locality of reference, resulting in a large working set that cannot fit into the available RAM.¹⁰ The symptoms of thrashing include very high CPU utilization (often paradoxically coupled with very little productive work being done), a significant increase in disk activity due to constant swapping, a high rate of page faults, and a noticeable slowdown in system response time.¹⁰⁹ In severe cases, thrashing can lead to a complete collapse of system performance.¹⁴¹

Several strategies can be employed to mitigate thrashing. One of the most effective is to increase the amount of physical RAM in the system.⁵⁷ Reducing the number of processes running concurrently (decreasing the degree of multiprogramming) can also alleviate the demand on memory resources.⁵⁷ Using more efficient page replacement algorithms that are better at predicting which pages are likely to be needed in the future can also help reduce the frequency of page faults.⁵⁷ The working

set model attempts to address thrashing by ensuring that a process is allocated enough frames to accommodate its current working set of actively used pages.¹¹⁹ Another approach is to monitor the page-fault frequency (PFF) of processes and adjust the number of frames allocated to them accordingly; if the page-fault rate is too high, it suggests that the process needs more frames.¹¹⁹

6. Case Studies: Virtual Memory in Action

Virtual memory is a fundamental feature of modern operating systems, and its implementation details can vary between different OS families. Here, we will examine the virtual memory implementations in Unix-based systems and Windows.

6.1 Unix Virtual Memory Implementation

Unix-like operating systems employ a sophisticated virtual memory system that utilizes a combination of RAM and disk storage to provide a flexible and efficient memory management environment.⁹⁶ Paging is a central mechanism, where physical memory is divided into fixed-size pages, and virtual addresses are mapped to these physical pages.⁶² In addition to paging, Unix also uses swapping, which involves moving entire process address spaces from RAM to disk when the system needs to free up physical memory for other processes.⁶² Memory allocation in Unix-like systems is typically dynamic, with techniques like the buddy system being used to allocate memory to processes on demand.⁹⁶ To enhance performance, Unix employs caching to store frequently accessed data in RAM, allowing for faster retrieval.⁹⁶ A key aspect of Unix virtual memory is that each process operates within its own virtual address space, providing a protected environment where one process cannot directly interfere with the memory of another.⁹⁸ The operating system kernel itself resides in a separate part of this virtual address space, accessible only when the system is running in kernel mode.¹⁴⁶ Unix-based systems are also designed to support large, sparse address spaces and allow for the mapping of files and even device memory directly into a process's address space.⁹⁷ The design of Unix virtual memory emphasizes extensibility and portability, aiming to provide a robust and adaptable memory management solution for a wide range of hardware and applications.⁹⁶

6.2 Windows Virtual Memory Implementation

Windows operating systems implement virtual memory through a dynamic system file called pagefile.sys, which works in conjunction with the physical RAM installed in the computer.⁸⁶ Unlike some earlier systems that used fixed-size swap files, Windows dynamically manages the size of the pagefile, allowing it to grow and shrink based on the current memory demands of the system.⁸⁶ This virtual memory system enables

applications to utilize more memory than is physically available in RAM.⁸⁶ In 32-bit versions of Windows, each process is typically allocated a 4GB virtual address space, with 2GB reserved for the private use of the process and the other 2GB shared between all processes and the operating system kernel.⁸⁸ The operating system automatically handles the translation between these virtual addresses and the actual physical addresses in RAM.⁸⁸ Users can monitor the performance and usage of virtual memory through built-in tools such as the Resource Monitor and the Task Manager.⁸⁸ While Windows manages virtual memory automatically by default, it also provides options for users to manually configure the size of the paging file if needed.⁸⁶ The implementation in Windows is designed to be largely self-managing, aiming to provide a stable and efficient environment for running a wide variety of applications without requiring extensive user intervention in memory management settings.⁸⁹

7. Conclusion: Balancing Efficiency and Functionality in Memory Management

Memory management in operating systems is a multifaceted and critical domain, encompassing a range of concepts and techniques designed to optimize the use of a computer's memory resources. This report has explored the fundamental principles, from the distinction between logical and physical addresses and the mechanisms for their translation, to the various strategies for memory allocation, including contiguous (fixed and variable partitioning) and non-contiguous (paging) methods. The challenges posed by memory fragmentation, both internal and external, and the techniques like compaction aimed at mitigating these issues, have been discussed.

The detailed examination of paging revealed its strengths in eliminating external fragmentation and enabling virtual memory, alongside its overheads related to page tables and potential for internal fragmentation. The concept of virtual memory itself, a cornerstone of modern operating systems, was explored in depth, highlighting its reliance on hardware support, control structures, and the principle of locality of reference to provide the illusion of memory space exceeding physical limitations. The crucial role of page faults in the demand paging process, the strategies employed by various page replacement algorithms to manage limited memory frames, and the performance-degrading phenomenon of thrashing were also examined.

The comparison of contiguous and non-contiguous allocation methods underscores the trade-offs between simplicity and efficiency. Fixed partitioning, while easy to implement, often leads to wasted memory due to internal fragmentation. Variable partitioning addresses this but introduces the complexity of external fragmentation. Paging offers a powerful solution to external fragmentation and enables virtual

memory, but it comes with its own set of overheads. The choice of page replacement algorithm further illustrates the need to balance performance with implementation complexity.

Virtual memory stands out as a pivotal innovation in operating system design, enabling multitasking, the execution of large applications, and providing essential memory protection. The case studies of Unix and Windows virtual memory implementations highlight different design philosophies, with Unix emphasizing flexibility and portability, and Windows focusing on automatic management and ease of use.

In conclusion, memory management in operating systems represents a continuous balancing act between efficiency, functionality, and complexity. As software applications become increasingly demanding and hardware architectures continue to evolve, the principles and techniques of memory management will remain central to ensuring optimal system performance and user experience. The ongoing advancements in memory management reflect the continuous effort to meet the ever-growing demands of the computing landscape.

Works cited

1. Memory Management in Operating Systems Explained - phoenixNAP, accessed April 21, 2025, <https://phoenixnap.com/kb/memory-management>
2. Memory Management in Operating Systems - DEV Community, accessed April 21, 2025, <https://dev.to/arjun98k/memory-management-in-operating-systems-278m>
3. Memory Management in Operating System | GeeksforGeeks, accessed April 21, 2025, <https://www.geeksforgeeks.org/memory-management-in-operating-system/>
4. courses.cs.washington.edu, accessed April 21, 2025, <https://courses.cs.washington.edu/courses/cse451/12sp/lectures/11-memory.pdf>
5. najibroslan.files.wordpress.com, accessed April 21, 2025, <https://najibroslan.files.wordpress.com/2010/07/unit62.pdf>
6. Logical and Physical Address in Operating System | GeeksforGeeks, accessed April 21, 2025, <https://www.geeksforgeeks.org/logical-and-physical-address-in-operating-system/>
7. 12.5: Logical vs Physical Address - Engineering LibreTexts, accessed April 21, 2025, https://eng.libretexts.org/Courses/Delta_College/Introduction_to_Operating_Systems/12%3A_Memory_Management/12.05%3A_Logical_vs_Physical_Address
8. Day 16: Logical vs Physical Address Space | Exploring Operating Systems - Mohit Mishra, accessed April 21, 2025, <https://mohitmishra786.github.io/exploring-os/src/day-16-logical-vs-physical-address-space.html>

9. Memory Management - OMSCS Notes, accessed April 21, 2025,
<https://www.omscs-notes.com/operating-systems/memory-management/>
10. Memory management in Operating System - Stack Overflow, accessed April 21, 2025,
<https://stackoverflow.com/questions/68515413/memory-management-in-operating-system>
11. Paging in OS - Scaler Topics, accessed April 21, 2025,
<https://www.scaler.com/topics/paging-in-os/>
12. Memory Allocation Techniques | Mapping Virtual Addresses to ..., accessed April 21, 2025,
<https://www.geeksforgeeks.org/memory-allocation-techniques-mapping-virtual-addresses-to-physical-addresses/>
13. Contiguous Memory Allocation in Operating System - Scaler Topics, accessed April 21, 2025,
<https://www.scaler.com/topics/contiguous-memory-allocation-in-os/>
14. Implementation of Contiguous Memory Management Techniques | GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/implementation-of-contiguous-memory-management-techniques/>
15. Fixed (or static) Partitioning in Operating System - GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/fixed-or-static-partitioning-in-operating-system/>
16. 7.4.1: Fixed Partitioning - Engineering LibreTexts, accessed April 21, 2025,
https://eng.libretexts.org/Courses/Delta_College/Operating_System%3A_The_Basics/07%3A_Memory/7.4%3A_Memory_Partitioning/7.4.1%3A_Fixed_Partitioning
17. Analysis of Allocation Algorithms in fixed partitioning - International Journal of Advances in Engineering and Management (IJAEM), accessed April 21, 2025,
https://ijaem.net/issue_dcp/Analysis%20of%20Allocation%20Algorithms%20in%20fixed%20partitioning.pdf
18. Memory Management - Stony Brook CS, accessed April 21, 2025,
<https://www3.cs.stonybrook.edu/~kifer/Courses/cse306/lectures/chap7.pdf>
19. Fixed Partition - Learn MERN, Next JS, DSA, AI, and Blockchain, accessed April 21, 2025,
<https://blogs.30dayscoding.com/blogs/os/memory-management/contiguous-memory-allocation/fixed-partition/>
20. Fixed Partitioning in Operating System | GATE Notes - BYJU'S, accessed April 21, 2025,
<https://byjus.com/gate/fixed-partitioning-in-operating-system-notes/>
21. Fixed (Or Static) Partitioning in Operating System | PDF | Software Engineering - Scribd, accessed April 21, 2025,
<https://www.scribd.com/document/522328246/hlQIV9ZxaL-1584950879815>
22. Fixed Partitioning in Operating System - Detailed Guide - Testbook, accessed April 21, 2025,
<https://testbook.com/gate/fixed-partitioning-in-operating-system-notes>
23. Memory Partitioning - Wisemonkeys, accessed April 21, 2025,
<https://wisemonkeys.info/blogs/Memory-Partitioning>

24. Difference between Internal and External fragmentation - GeeksforGeeks, accessed April 21, 2025, <https://www.geeksforgeeks.org/difference-between-internal-and-external-fragmentation/>
25. L-5.3: Internal Fragmentation | Fixed size Partitioning | Memory management - YouTube, accessed April 21, 2025, <https://www.youtube.com/watch?v=bK-VhQA512c>
26. How to compute total internal and external fragmentation, accessed April 21, 2025, <https://cs.stackexchange.com/questions/66011/how-to-compute-total-internal-and-external-fragmentation>
27. memory management - Internal fragmentation basic concept - Stack Overflow, accessed April 21, 2025, <https://stackoverflow.com/questions/45441361/internal-fragmentation-basic-concept>
28. CSC 553 Operating Systems - Lecture 7 - Memory Management, accessed April 21, 2025, <https://home.adelphi.edu/~siegfried/cs553/553l7.pdf>
29. testbook.com, accessed April 21, 2025, <https://testbook.com/gate/fixed-partitioning-in-operating-system-notes#:~:text=1.-,External%20Fragmentation,to%20the%20no%2Dspanning%20rule.>
30. Can fixed partitioning suffer from external fragmentation? - Stack Overflow, accessed April 21, 2025, <https://stackoverflow.com/questions/64454029/can-fixed-partitioning-suffer-from-external-fragmentation>
31. What are the advantages and disadvantages of contiguous allocation with fixed partitioning?, accessed April 21, 2025, <https://www.krayonnz.com/user/doubts/detail/6144eae9d2f93e00401e4155/what-are-the-advantages-and-disadvantages-of-contiguous-allocation-with-fixed-partitioning>
32. 7.4.2: Variable Partitioning - Engineering LibreTexts, accessed April 21, 2025, https://eng.libretexts.org/Courses/Delta_College/Operating_System%3A_The_Basics/07%3A_Memory/7.4%3A_Memory_Partitioning/7.4.2%3A_Variable_Partitioning
33. Variable (or Dynamic) Partitioning in Operating System | GeeksforGeeks, accessed April 21, 2025, <https://www.geeksforgeeks.org/variable-or-dynamic-partitioning-in-operating-system/>
34. Unit: III Lecture: 5 (Memory Management) Memory Allocation Techniques (Part-II) Variable Partition Scheme or Dynamic Partition S, accessed April 21, 2025, <https://www.gchamirpur.org/wp-content/uploads/2023/07/Unit-III-Lecture-5-Memory-Allocation-Techniques-Part-II.pdf>
35. Dynamic Partitioning in Operating System - BYJU'S, accessed April 21, 2025, <https://byjus.com/gate/dynamic-partitioning-in-operating-system-notes/>
36. Fixed And Variable Partition Strategies | BimStudies.Com, accessed April 21, 2025, <https://bimstudies.com/docs/operating-system/memory-management/memory-allocation-strategies/>

37. Variable Partition - Learn MERN, Next JS, DSA, AI, and Blockchain, accessed April 21, 2025,
<https://blogs.30dayscoding.com/blogs/os/memory-management/contiguous-memory-allocation/variable-partition/>
38. Contiguous Memory Allocation - Dextutor, accessed April 21, 2025,
<https://dextutor.com/contiguous-memory-allocation/>
39. [Solved] Contiguous memory allocation having variable size partition - Testbook, accessed April 21, 2025,
<https://testbook.com/question-answer/contiguous-memory-allocation-having-variable-size--5ffff45010efd078c9af43d0>
40. Dynamic Partitioning - Tutorialspoint, accessed April 21, 2025,
<https://www.tutorialspoint.com/dynamic-partitioning>
41. Variable-Partition Multiprogramming Variable-Partition Characteristics, accessed April 21, 2025,
https://uomustansiriyah.edu.iq/media/lectures/9/9_2019_01_04!09_18_17_AM.pdf
42. Difference between Fixed Partitioning and Variable Partitioning - GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/difference-between-fixed-partitioning-and-variable-partitioning/>
43. pages.cs.wisc.edu, accessed April 21, 2025,
<https://pages.cs.wisc.edu/~solomon/cs537-old/last/memory.html#:~:text=Compaction%20attacks%20the%20problem%20of,when%20the%20block%20is%20moved.>
44. Memory Management Techniques in Operating System - Shiksha Online, accessed April 21, 2025,
<https://www.shiksha.com/online-courses/articles/memory-management-techniques-in-operating-system/>
45. Compaction in Operating System | GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/compaction-in-operating-system/>
46. Difference Between Fragmentation and Compaction | GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/difference-between-fragmentation-and-compaction/>
47. CS322: Memory Management - Gordon College, accessed April 21, 2025,
<https://www.cs.gordon.edu/courses/cs322/lectures/memory.html>
48. Paging and Memory Management: Enhance Computing Speed | Lenovo US, accessed April 21, 2025, <https://www.lenovo.com/us/en/glossary/paging/>
49. Paging in OS | GATE Notes - BYJU'S, accessed April 21, 2025,
<https://byjus.com/gate/paging-in-operating-system-notes/>
50. Paging in Operating System - GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/paging-in-operating-system/>
51. www.lenovo.com, accessed April 21, 2025,
<https://www.lenovo.com/us/en/glossary/paging/#:~:text=Paging%20is%20a%20memory%20management,be%20easily%20managed%20and%20swapped.>
52. Memory paging - Wikipedia, accessed April 21, 2025,

- https://en.wikipedia.org/wiki/Memory_paging
53. workat.tech, accessed April 21, 2025,
<https://workat.tech/core-cs/tutorial/paging-in-operating-system-os-knbcthp3w8o7#:~:text=In%20paging%2C%20processes%20are%20divided,the%20size%20of%20a%20page.>
 54. Paging In OS Explained In Detail // Unstop, accessed April 21, 2025,
<https://unstop.com/blog/what-is-paging-in-os>
 55. Difference Between Paging and Segmentation - GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/difference-between-paging-and-segmentation/>
 56. [06] PAGING, accessed April 21, 2025,
<https://www.cl.cam.ac.uk/teaching/1516/OpSystems/pdf/06-Paging.pdf>
 57. Paging in Operating Systems: What it Is & How it Works - phoenixNAP, accessed April 21, 2025, <https://phoenixnap.com/kb/paging>
 58. Paging in Operating System (OS) | Core Computer Science, accessed April 21, 2025,
<https://workat.tech/core-cs/tutorial/paging-in-operating-system-os-knbcthp3w8o7>
 59. Memory protection - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Memory_protection
 60. Paging - OSDev Wiki, accessed April 21, 2025, <http://wiki.osdev.org/Paging>
 61. Operating System - paging - Stack Overflow, accessed April 21, 2025,
<https://stackoverflow.com/questions/10767060/operating-system-paging>
 62. Virtual memory - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Virtual_memory
 63. CS3130: Virtual Memory, accessed April 21, 2025,
<https://www.cs.virginia.edu/~cr4bd/3130/F2024/readings/vm.html>
 64. Chapter 8 Virtual Memory, accessed April 21, 2025,
<https://www.dc.fi.udc.es/~so-grado/VirtualMemoryOSedition7StallingsBook.pdf>
 65. Paging & segmentation; advantages and disadvantage | PPT - SlideShare, accessed April 21, 2025,
<https://www.slideshare.net/slideshow/paging-segmentation-advantages-and-disadvantage/266824613>
 66. CSC 553 Operating Systems - Lecture 8 - Virtual Memory, accessed April 21, 2025, <https://home.adelphi.edu/~siegfried/cs553/553l8.pdf>
 67. Hardware Support - Learn MERN, Next JS, DSA, AI, and Blockchain, accessed April 21, 2025,
<https://blogs.30dayscoding.com/blogs/os/memory-management/paging/hardware-support/>
 68. lass.cs.umass.edu, accessed April 21, 2025,
<https://lass.cs.umass.edu/~shenoy/courses/fall08/lectures/Lec15.pdf>
 69. Performance of paging | GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/performance-of-paging/>
 70. Virtual Memory - CS 3410, accessed April 21, 2025,
<https://www.cs.cornell.edu/courses/cs3410/2024fa/notes/vm.html>

71. unstop.com, accessed April 21, 2025,
<https://unstop.com/blog/what-is-paging-in-os#:~:text=This%20allows%20the%20Operating%20system,page%20into%20their%20address%20spaces.>
72. Memory Protection in Operating Systems | GeeksforGeeks, accessed April 21, 2025, <https://www.geeksforgeeks.org/memory-protection-in-operating-systems/>
73. why do we need these?memory segmentation, memory paging, and flat memory models? : r/Compilers - Reddit, accessed April 21, 2025,
https://www.reddit.com/r/Compilers/comments/16mitth/why_do_we_need_these_memory_segmentation_memory/
74. Paging in OS: An Ultimate Guide - Simplilearn.com, accessed April 21, 2025,
<https://www.simplilearn.com/paging-in-os-article>
75. What is paging exactly? : r/osdev - Reddit, accessed April 21, 2025,
https://www.reddit.com/r/osdev/comments/mo40m9/what_is_paging_exactly/
76. ELI5: What is virtual memory and how does it work? : r/explainlikeimfive - Reddit, accessed April 21, 2025,
https://www.reddit.com/r/explainlikeimfive/comments/qacvpv/eli5_what_is_virtual_memory_and_how_does_it_work/
77. Paging vs Segmentation in Operating System - DataFlair, accessed April 21, 2025,
<https://data-flair.training/blogs/paging-vs-segmentation-in-operating-system/>
78. Paging in Operating System (OS): Definition, Benefits & Drawbacks, accessed April 21, 2025,
<https://www.theknowledgeacademy.com/blog/paging-in-operating-system/>
79. Difference Between Paging And Segmentation Explained! - Unstop, accessed April 21, 2025,
<https://unstop.com/blog/difference-between-paging-and-segmentation>
80. Advantages and Disadvantages, accessed April 21, 2025,
<http://www.cs.iit.edu/~cs561/cs351/VM/paging%20advantages%20and%20disadvantages.html>
81. Memory Management! Goals of this Lecture! - cs.Princeton, accessed April 21, 2025,
<https://www.cs.princeton.edu/courses/archive/fall13/cos217/lectures/19MemoryMgmt.pdf>
82. What is virtual memory? Definition & Explanation | Crucial.com, accessed April 21, 2025,
<https://www.crucial.com/articles/about-memory/virtual-memory-settings-suggestions>
83. Virtual Memory in Operating System | GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/virtual-memory-in-operating-system/>
84. An introduction to virtual memory - Internal Pointers, accessed April 21, 2025,
<https://www.internalpointers.com/post/introduction-virtual-memory>
85. 8.2: Virtual Memory in the Operating System - Engineering LibreTexts, accessed April 21, 2025,
https://eng.libretexts.org/Courses/Delta_College/Operating_System%3A_The_Basics/08%3A_Virtual_Memory/8.2%3A_Virtual_Memory_in_the_Operating_System
86. How to Manage Windows 10 Virtual Memory Pagefile | NinjaOne, accessed April

- 21, 2025,
<https://www.ninjaone.com/blog/manage-windows-10-virtual-memory-pagefile/>
87. Manage Virtual Memory Pagefile in Windows 10, accessed April 21, 2025,
<https://www.tenforums.com/tutorials/77692-manage-virtual-memory-pagefile-windows-10-a.html>
 88. Virtual memory in 32-bit version of Windows - Windows Server | Microsoft Learn, accessed April 21, 2025,
<https://learn.microsoft.com/en-us/troubleshoot/windows-server/performance/ram-virtual-memory-pagefile-management>
 89. Best Ways to Optimize Virtual Memory on Windows 10 - Trio MDM, accessed April 21, 2025, <https://www.trio.so/blog/virtual-memory-on-windows-10/>
 90. Physical and Virtual Memory in Windows 10 - Microsoft Community, accessed April 21, 2025,
<https://answers.microsoft.com/en-us/windows/forum/all/physical-and-virtual-memory-in-windows-10/e36fb5bc-9ac8-49af-951c-e7d39b979938>
 91. How to change virtual memory size on Windows 10, accessed April 21, 2025,
<https://www.windowscentral.com/how-change-virtual-memory-size-windows-10>
 92. Operating Systems: Lecture 10: Memory Management (1) - University of Malta, accessed April 21, 2025,
<http://staff.um.edu.mt/csta1/courses/lectures/csm202/os10.html>
 93. Swap space and dirty pages - Stack Overflow, accessed April 21, 2025,
<https://stackoverflow.com/questions/43695565/swap-space-and-dirty-pages>
 94. Demand paging - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Demand_paging
 95. Can someone explain what demand paging is once and for all? : r/computerscience - Reddit, accessed April 21, 2025,
https://www.reddit.com/r/computerscience/comments/40ajbk/can_someone_explain_what_demand_paging_is_once/
 96. Memory Management in UNIX - Mithilesh's Blog, accessed April 21, 2025,
<https://esymith.hashnode.dev/memory-management-in-unix>
 97. A New Virtual Memory Implementation for Berkeley UNIX® - FreeBSD Documentation Archive, accessed April 21, 2025,
<https://docs-archive.freebsd.org/44doc/papers/newvm.pdf>
 98. Chapter 3 Memory Management, accessed April 21, 2025,
<https://www.tldp.org/LDP/tlk/mm/memory.html>
 99. Locality of reference - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Locality_of_reference
 100. 1.7.2 Cache Memory - Locality of reference - Engineering LibreTexts, accessed April 21, 2025,
https://eng.libretexts.org/Courses/Delta_College/Operating_System%3A_The_Basics/01%3A_The_Basics_-_An_Overview/1.7_Cache_Memory/1.7.2_Cache_Memory_-_Locality_of_reference
 101. Part 3 : The Memory Hierarchy 3.1 Principle Of Locality, accessed April 21, 2025, <https://faculty.uobasrah.edu.iq/uploads/teaching/1710261690.pdf>
 102. Locality of Reference, accessed April 21, 2025,

- <https://hillside.net/plop/plop97/Proceedings/bhatt.pdf>
103. Locality of reference - GUVI, accessed April 21, 2025,
<https://www.guvi.in/blog/locality-of-reference/>
 104. Locality of Reference and Cache Operation in Cache Memory |
GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/locality-of-reference-and-cache-operation-in-cache-memory/>
 105. 4.2 Locality of reference(Principle of locality) Flashcards - Quizlet, accessed
April 21, 2025,
<https://quizlet.com/557606159/42-locality-of-referenceprinciple-of-locality-flash-cards/>
 106. What is the Locality of Reference - Tutorialspoint, accessed April 21, 2025,
<https://www.tutorialspoint.com/what-is-the-locality-of-reference>
 107. caching - What is locality of reference? - Stack Overflow, accessed April 21,
2025, <https://stackoverflow.com/questions/7638932/what-is-locality-of-reference>
 108. Locality of reference | Code Guidelines for Correctness, Modernization,
Security, Portability, and Optimization - Open Catalog, accessed April 21, 2025,
<https://open-catalog.codeee.com/Glossary/Locality-of-reference>
 109. Page Fault Handling in Operating System - GeeksforGeeks, accessed April 21,
2025, <https://www.geeksforgeeks.org/page-fault-handling-in-operating-system/>
 110. Demand Paging, accessed April 21, 2025,
<https://cseweb.ucsd.edu/classes/sp17/cse120-a/applications/ln/lecture13.html>
 111. Operating Systems: Virtual Memory, accessed April 21, 2025,
https://www.cs.uic.edu/~jbell/CourseNotes/OperatingSystems/9_VirtualMemory.html
 112. Can I Prevent Page Faults from Occurring? | Lenovo US, accessed April 21,
2025, <https://www.lenovo.com/us/en/glossary/page-fault/>
 113. Understanding and troubleshooting page faults and memory swapping -
Site24x7, accessed April 21, 2025,
<https://www.site24x7.com/learn/linux/page-faults-memory-swapping.html>
 114. Page fault - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Page_fault
 115. Page Fault - Tutorial - takeUforward, accessed April 21, 2025,
<https://takeuforward.org/operating-system/page-fault>
 116. Understanding Page Faults: Essential OS Tutorial - upGrad, accessed April 21,
2025,
<https://www.upgrad.com/tutorials/software-engineering/software-key-tutorial/when-does-page-fault-occurs/>
 117. Page Fault in Operating System - Studytonight, accessed April 21, 2025,
<https://www.studytonight.com/operating-system/page-fault-in-operating-system>
 118. What is Demand Paging in Operating System? - GeeksforGeeks, accessed
April 21, 2025,
<https://www.geeksforgeeks.org/what-is-demand-paging-in-operating-system/>
 119. Demand paging:, accessed April 21, 2025,
https://www.pvpsiddhartha.ac.in/dep_it/lecture%20notes/os/unit5.pdf

120. en.wikipedia.org, accessed April 21, 2025,
[https://en.wikipedia.org/wiki/Working_set#:~:text=The%20working%20set%20model%20states,pages\)%20can%20be%20in%20RAM.](https://en.wikipedia.org/wiki/Working_set#:~:text=The%20working%20set%20model%20states,pages)%20can%20be%20in%20RAM.)
121. Working set - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Working_set
122. Working Set - Win32 apps | Microsoft Learn, accessed April 21, 2025,
<https://learn.microsoft.com/en-us/windows/win32/memory/working-set>
123. Working Set in Paging | GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/working-set-in-paging/>
124. The Working Set Model for Program Behavior, accessed April 21, 2025,
https://courses.cs.washington.edu/courses/cse551/09sp/papers/working_set.pdf
125. What is a dirty bit? - Lenovo, accessed April 21, 2025,
<https://www.lenovo.com/us/en/glossary/dirty-bit/>
126. Dirty bit - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Dirty_bit
127. Dirty bit definition - Glossary | NordVPN, accessed April 21, 2025,
<https://nordvpn.com/cybersecurity/glossary/dirty-bit/>
128. Dirty Bit Tracking - Windows drivers | Microsoft Learn, accessed April 21, 2025,
<https://learn.microsoft.com/en-us/windows-hardware/drivers/display/dirty-bit-tracking>
129. Information of dirty virtual memory pages - Microsoft Q&A, accessed April 21, 2025,
<https://learn.microsoft.com/en-us/answers/questions/729533/information-of-dirty-virtual-memory-pages>
130. CS 537 Notes, Section #20: Clock Algorithm, Thrashing, accessed April 21, 2025, <https://pages.cs.wisc.edu/~bart/537/lecturenotes/s20.html>
131. What is the difference between demand paging and page replacement? - Stack Overflow, accessed April 21, 2025,
<https://stackoverflow.com/questions/29873682/what-is-the-difference-between-demand-paging-and-page-replacement>
132. What Is Demand Paging? | Baeldung on Computer Science, accessed April 21, 2025, <https://www.baeldung.com/cs/demand-paging>
133. Page replacement algorithm - Wikipedia, accessed April 21, 2025,
https://en.wikipedia.org/wiki/Page_replacement_algorithm
134. A Comparison of Three Page Replacement Algorithms: FIFO, LRU and Optimal | PDF | Computer Data - Scribd, accessed April 21, 2025,
<https://fr.scribd.com/document/410687273/7405-28646-1-PB-pdf>
135. Page replacement algorithms | PPT - SlideShare, accessed April 21, 2025,
<https://www.slideshare.net/slideshow/page-replacement-algorithms-a010-a011/34907037>
136. What is Page Replacement in OS? | Scaler Topics, accessed April 21, 2025,
<https://www.scaler.com/topics/operating-system/page-replacement-algorithm/>
137. Not Recently Used (NRU) page replacement algorithm | GeeksforGeeks, accessed April 21, 2025,
<https://www.geeksforgeeks.org/not-recently-used-nru-page-replacement-algori>

[thm/](#)

138. SIMULATE AND IMPLEMENTATION OF PAGE REPLACEMENT TECHNIQUES THROUGH MEMORY ADDRESSES - ijates, accessed April 21, 2025, http://www.ijates.com/images/short_pdf/1405423164_P208-213.pdf
139. EASY-HOW-TO Page Replacement Algorithm (FIFO, Optimal, and LRU) Tutorial (Manual), accessed April 21, 2025, <https://www.youtube.com/watch?v=cjWnEtnKVGm>
140. What is Thrashing? Why Does it Occur? - Lenovo, accessed April 21, 2025, <https://www.lenovo.com/us/en/glossary/thrashing/>
141. Thrashing (computer science) - Wikipedia, accessed April 21, 2025, [https://en.wikipedia.org/wiki/Thrashing_\(computer_science\)](https://en.wikipedia.org/wiki/Thrashing_(computer_science))
142. What is thrashing and how does it affect performance?-Huawei Enterprise Support Community, accessed April 21, 2025, <https://forum.huawei.com/enterprise/en/What-is-thrashing-and-how-does-it-affect-performance/thread/692495472444456960-667213859733254144>
143. What is Thrashing in OS (Operating System)? - The Knowledge Academy, accessed April 21, 2025, <https://www.theknowledgeacademy.com/blog/thrashing-in-os/>
144. Thrashing in OS (Operating System) - Scaler Topics, accessed April 21, 2025, <https://www.scaler.com/topics/thrashing-in-os/>
145. Thrashing in OS (Operating System) - What Is Thrash - Javatpoint | PDF - Scribd, accessed April 21, 2025, <https://www.scribd.com/document/751952056/Thrashing-in-OS-Operating-System-What-is-Thrash-Javatpoint>
146. Introduction to Virtual Memory, accessed April 21, 2025, <https://www.cs.miami.edu/home/burt/learning/Csc521.101/notes/virtual-memory-notes.html>
147. SunOS Virtual Memory Implementation - Kos, accessed April 21, 2025, <http://kos.enix.org/pub/sunos-vi.pdf>
148. Addressing and Virtual Memory (plus some UNIX and Mach) - Washington, accessed April 21, 2025, <https://courses.cs.washington.edu/courses/csep551/17wi/week2.pdf>
149. How To: Increase Virtual Memory beyond the Recommended Maximum - Esri Support, accessed April 21, 2025, <https://support.esri.com/en-us/knowledge-base/increase-virtual-memory-beyond-the-recommended-maximum--000011346>
150. What is the Optimal Virtual Memory Size For 8 GB RAM | WIN 10? - Super User, accessed April 21, 2025, <https://superuser.com/questions/1090569/what-is-the-optimal-virtual-memory-size-for-8-gb-ram-win-10>