Memory Management Techniques in Computer Architecture

Yashwanth Raj Varadharajan

G47635180

Prof. Morris Lancaster

CSCI_6461 Computer System Architecture

School of Engineering and Applied Science

The George Washington University

**Abstract**

In today's computing world, programs and operating systems cannot exist without effective memory management, especially if the application must operate under Survey load for an indefinite period. To improve performance, resources must be used efficiently. This paper addresses memory management in an operating system and demonstrates the basic architecture of segmentation in an operating system as well as the basics of its allocation. Memory management is one of the key roles of the operating system, which is accomplished through the usage of the memory management unit (MMU). The MMU is a software component of the operating system that lives in the kernel. The operating system manages both primary and secondary memory. This paper also describes the basic concepts of virtual memory management and dynamic memory management. Virtual memory is a notion that provides a vast addressable region to programs to facilitate multiprogramming. This suggests that the processor functions under the guise of accessing a huge accessible storage space. As a result, the operating system guarantees that available memory is managed properly and effectively to maximize system efficiency. To achieve the overall memory management goal, the operating system takes on a supervisory role via the memory manager.

**Introduction**

The appropriate management of memory is a vital feature that strongly affects system performance, reliability, and overall functioning in the complex landscape of computer architecture. Memory management is the orchestration of a computer's memory resources to support program execution, efficient data storage, and seamless user experiences. As the demand for processing power and data storage grows, architects, engineers, and system designers must understand and apply advanced memory management approaches.

The primary function of a computer system is to run programs. At the time of execution, these programs and their accusing data must be in the main memory. The fundamentals of managing main memory, which is one of the most valuable resources of an operating system's multiprogramming architecture. Memory is made up of a vast number of words or bytes, each having its own address. Modern operating systems enable effective memory management, and research is still being undertaken to improve the way memory is allocated for applications because the memory allocation algorithm is the fundamental challenge.

A memory management unit (MMU) is used by the operating system to effectively manage memory. Because running applications can only detect logical addresses, the MMU converts virtual addresses to physical addresses. Memory management solutions have changed throughout time to adapt to the changing computer landscape. This study goes into numerous memory management solutions, each of which is tailored to meet distinct challenges and objectives. Understanding these strategies is critical for anybody involved in computer

architecture and system design, from traditional contiguous memory allocation to advanced virtual memory systems.

Virtual memory contains virtual addresses. Virtual memory is a notion that provides a broad memory address space to programs in order to facilitate multiprogramming. That means that the CPU functions under the illusion of accessing a huge addressable space, which is not the case, a concept maintained by the operating system. The CPU is essential in producing virtual addresses, which are then translated into physical addresses by the MMU.
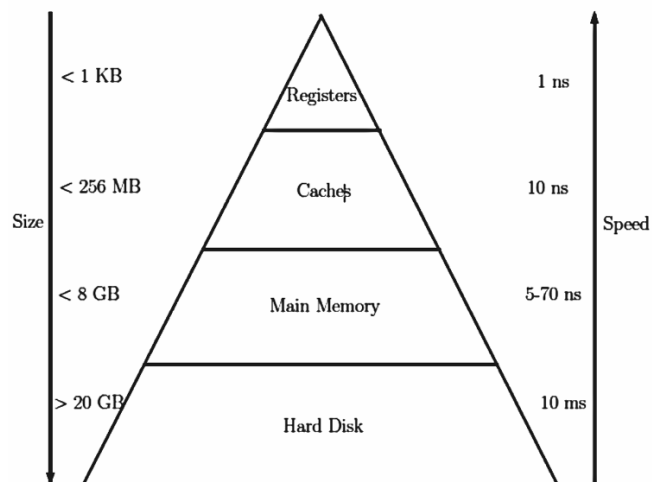
**Memory Hierarchy**

The Memory Hierarchy is a critical notion for the CPU's successful data management. This term is used by computer scientists to describe the various mechanisms by which a computer processes information, whether temporarily or permanently stored. The importance of the memory hierarchy in modern computer systems develops as the performance gap between the processor and memory widens. Notably, computer memory, as represented by Random Access Memory (RAM), functions at a faster rate than storage, providing faster access to program and data storage despite having bigger capacities. When necessary, the contents of computer memory can be relocated to secondary storage via virtual memory, a memory management approach.

The memory hierarchy is intended to handle the trade-off between performance, capacity, and cost. It is made up of several tiers of memory, each with its own set of features in terms of speed, size, and cost. The basic purpose of the memory hierarchy is to enable the quickest

feasible access to the most frequently used data while keeping the overall memory system cost as low as possible.

**Figure 1**

Memory Hierarchy



Note: Speed is given in nanoseconds

     The memory hierarchy is comprised of registers, caches, main memory (RAM), and secondary storage (disk or solid-state drive). Registers, the smallest and fastest type of memory within the CPU, are used to store data that is actively processed by the CPU. Caches, small and fast memory units, use the principle of locality to keep copies of data from the main memory that will be accessed soon. RAM is the principal volatile memory in a computer system, with more storage capacity but slower access times than caches. Secondary storage, represented by hard disk drives or solid-state drives, has the most storage capacity but has slower access speeds than other memory tiers.

*Registers*

Registers, the smallest and fastest form of memory in a computer system, are critical at the top of the memory hierarchy. The miracle of computation begins with these small storage units built directly within the CPU. Registers serve as the quick, temporary archives that orchestrate the symphony of computing operations in the delicate dance of instructions and data.

Registers are the memory hierarchy's sprinters, providing lightning-fast access to the CPU. They are not only compact but also extremely fast when compared to other types of memory. Registers hold binary data in the form of bits and operate at speeds practically imperceptible to human perception, giving instantaneous workspace for the CPU during instruction execution. Current processors employ register files to store data, it also distributes transitional outcomes of calculations in Memory Hierarchy. Register Files is the most extreme level in various aspects. Contrary requirements decide and govern Register Files.

*Cache*

Cache memory, located between registers and main memory, is an important node in the memory structure of modern computer systems. This intermediary layer of high-speed, volatile memory is intentionally placed to bridge the significant speed difference between registers and the slower main memory (RAM). This section delves into the complexities of cache memory, investigating its architecture, functionality, and impact on overall system performance.

Caching is important in memory management strategies because it uses the idea of locality to improve memory access performance. Caching techniques are classified into two

types: geographical locality and temporal locality. Spatial locality refers to a program's inclination to access memory regions that are close together, whereas temporal locality refers to a program's tendency to access the same memory locations frequently over a short period of time. Caching techniques such as block size, associativity, and replacement policies are used to optimize cache memory consumption and increase overall system performance.

### Random Access Memory (RAM)

RAM has a crucial place in the memory architecture, acting as the dynamic workspace where active programs and the operating system reside during computer operation. Unlike lightning-fast registers and faster, but smaller, cache memory, RAM strikes a balance between speed, capacity, and cost, making it a critical component in the smooth operation of modern computer systems. RAM is the most common type of volatile memory in a computer system, which means it loses its stored data when the power is switched off. Its unique feature is its random-access capabilities, which allows the CPU to read and write data at nearly the same speed regardless of the data's physical placement in the memory module. This feature is critical to its role in supporting application execution and the efficient operation of the operating system.

### Hard Disk

A hard disk, also known as a hard drive or fixed disk, is an electromechanical storage device used for storing and retrieving digital data via magnetic storage. These platters, which consist of one or more rotating disks coated with magnetic material, are coupled to magnetic heads furnished with an actuator arm. This method allows data to be read and written on the

plates' surfaces. Data is retrieved at random, allowing for storage in any order, though not necessarily in sequence.

Hard drives, which are classified as non-volatile storage devices, have an extraordinary ability to keep data even when the power is switched off. This contrasts sharply with volatile memory, such as RAM, which loses stored information over power cycles. Because hard drives are non-volatile, they are an important component for long-term data storage, assuring the preservation of files, applications, and the operating system across reboots and power outages. This feature is very useful for users looking for persistent storage solutions that can retain data integrity over long periods of time.

Hard drives access data sequentially, which means they read and write data in a linear, ordered form. This procedure involves moving a read/write head to a certain track on the spinning platter that contains the requested data. While sequential access is useful for reading or writing big blocks of data, it has a larger latency than random access volatile memory like RAM. This characteristic determines hard disk access patterns and is critical in determining their efficiency in managing large-scale data retrieval operations.

Hard drives are available in a variety of sizes to meet a wide range of computing requirements. Traditional Hard Disk Drives (HDDs) are excellent for large-scale data storage because they use spinning platters and a mechanical arm to access data. SSDs, on the other hand, use NAND flash memory for faster access times and lower latency, making them ideal for applications where speed is critical. Hybrid drives offer a balanced approach by combining the
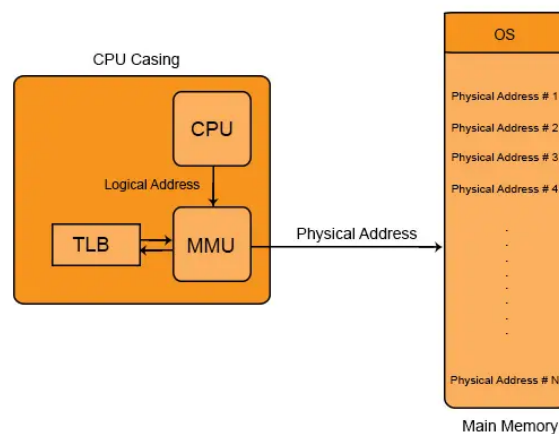
huge capacity of HDDs with the speed benefits of SSD caching. The selection of these types is based on specific requirements, with users assessing aspects such as speed, capacity, and cost to determine the best hard drive for their needs.

## Memory Management Unit

The MMU is an important component in memory management because it acts as a translator between the conceptual realm of virtual addresses provided by the CPU and the actual reality of physical memory addresses. Its primary job is to convert virtual addresses as viewed by the executing software into corresponding physical addresses, allowing data retrieval and storage in actual memory modules. This procedure, known as address translation, is the foundation of current operating systems, allowing them to take advantage of the concept of virtual memory. The complex process of address translation is crucial to the MMU's functionality.

**Figure 2**

Memory Management Unit



Note: This image shows the systematic operation of an MMU

When a program running on the CPU refers to a memory location, it does so in virtual terms - an abstract realm where the application is given the illusion of a continuous block of memory. The MMU comes in to translate these virtual addresses into tangible, physical addresses, allowing interaction between the CPU and the multiple memory levels. This translation not only allows for the smooth execution of processes, but it also serves as the foundation for key features such as memory protection and the efficient utilization of physical memory.

The MMU's ability to enforce memory protection and isolation amongst separate processes running on a system is one of its most important accomplishments. Each process is given its own virtual address space, giving the appearance of a dedicated memory block. The MMU prevents processes from accidentally accessing or modifying memory regions belonging to other processes. This approach improves system stability, security, and reliability by establishing a stable environment in which diverse programs can run independently without fear of unwanted intervention.

### *Translation Lookaside Buffer (TLB)*

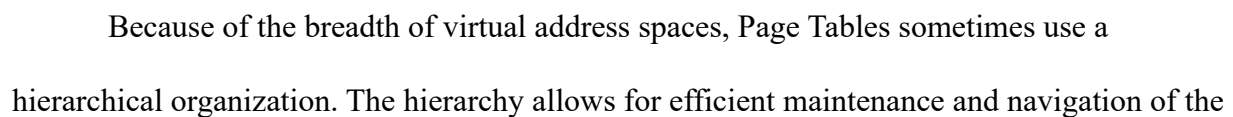The Translation Lookaside Buffer (TLB) is an important component of the Memory Management Unit (MMU), helping to speed up the complex process of address translation in modern computer systems. The TLB, which is located at the junction of hardware and software, acts as a high-speed cache, holding recently translated virtual-to-physical address mappings and considerably improving memory access efficiency.

The TLB's principal role is to reduce address translation delay by providing a rapid reference for commonly requested address mappings. When a CPU-running application refers to a memory address, the MMU first examines the TLB. If a match is detected, the appropriate physical address is directly fetched from the TLB, eliminating the need for a time-consuming full translation process. This method not only speeds up memory access, but it also improves overall system performance.

Despite its instrumental role, the TLB has a fixed size, and its usefulness is dependent on the TLB's coherence with the operating system's page tables. TLB misses, which occur when a required translation is not found in the TLB, can cause a minor delay as the MMU searches the page tables for the appropriate information. As a result, TLB management and optimization measures are critical for guaranteeing its efficacy.

### *Page Tables and Page Table Entries: Blueprint of Address Mapping*

In the sophisticated architecture of computer memory management, Page Tables and their constituent pieces, Page Table Entries (PTEs), emerge as the blueprint guiding the Memory Management Unit (MMU) in the complex process of address mapping. These structures are critical for translating virtual addresses created by programs into physical addresses in the underlying memory modules. Page Tables are a hierarchical data structure maintained by the operating system that acts as a vital link between the virtual and physical regions of memory.

**Figure 3**

Page Table



Note: This is an example page table created and used in ARM32 architecture

They divide the virtual address space into fixed-size blocks known as pages and establish the relationship between these virtual memory pages and their physical memory equivalents. Each page table entry contains crucial information, most notably the actual address of the associated page in main memory. Page Table Entries, which reside within the page tables, include critical metadata for each page. The physical address of the relevant page in physical memory is the most fundamental element stored in a PTE. PTEs also contain control bits that specify access permissions (read, write, execute), status bits that indicate the page's state (present or swapped out), and other features that influence the behavior of the associated memory page.

Because of the breadth of virtual address spaces, Page Tables sometimes use a hierarchical organization. The hierarchy allows for efficient maintenance and navigation of the

page table structure. The virtual address in a multi-level page table system is broken into numerous sections, each corresponding to a level in the hierarchy. The upper-level page table directs to a lower-level page table, advancing through the hierarchy until the ultimate PTE is reached. When a CPU-based software refers to a virtual address, the MMU searches the page table to find the associated physical address. The hierarchical organization enables for a more selective and quicker search through the page tables, improving address translation efficiency. This mechanism ensures that applications preserve the illusion of a contiguous and large virtual address space, even when physical memory is fragmented or constrained.

## Memory Management Techniques

Memory management techniques are fundamental tactics used in computer systems to properly handle memory resource allocation, consumption, and deallocation. In the world of computing, where memory is a finite and vital asset, proper memory management is critical for guaranteeing optimal system performance and the smooth execution of many activities. These techniques cover a wide range of methodologies for dealing with issues such as fragmentation, allocation conflicts, and the dynamic nature of memory requirements displayed by running programs.

As memory hierarchy orchestrators, these techniques are critical in managing the often conflicting demands of different processes, giving them the illusion of plentiful memory while prudently optimizing the use of physical memory modules. Memory management strategies, in essence, represent the backbone of a computer system's ability to maintain order and efficiency

in the dynamic terrain of memory access and utilization, considerably contributing to the overall stability and responsiveness of current computing environments.

### Contiguous Memory Allocation

Contiguous memory allocation is a fundamental memory management approach used in computer systems that allocates a single, unbroken block of memory to each process. In this easy approach, physical memory is partitioned into fixed-size segments, and processes are successively loaded into these partitions. Because of its simplicity, contiguous memory allocation was an early and fundamental technique utilized in a variety of computing contexts.

One of the most notable benefits of contiguous memory allocation is its simplicity of implementation. The ease with which a continuous chunk of memory can be assigned to a process facilitates the tracking of allocated and free memory regions. This simplicity also helps to faster access times because a process's memory addresses are contiguous, allowing for efficient traversal and retrieval. However, the ease of contiguous memory allocation has its drawbacks. External fragmentation, which occurs when free memory is dispersed around the system rather than in a single contiguous block, might limit the allocation of larger processes, resulting in inefficient memory use.
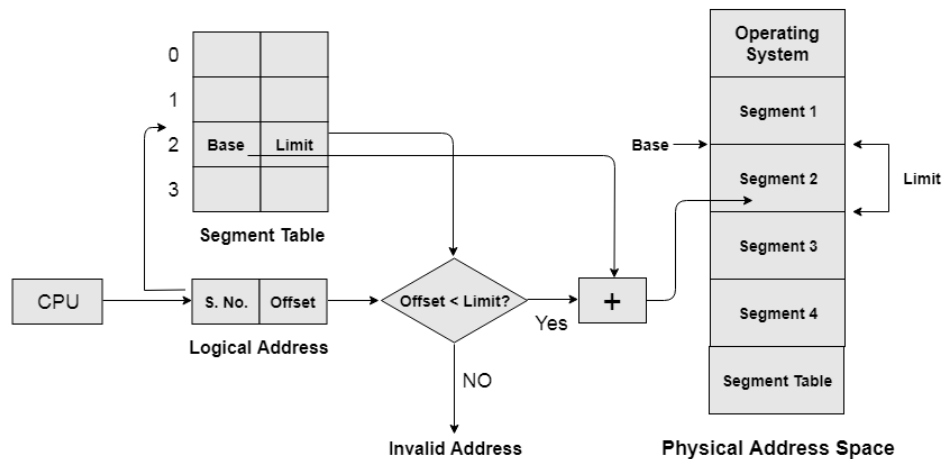
### Paging

Paging is a crucial memory management technique used to address the difficulties associated with contiguous memory allocation. Paging provides a dynamic and efficient solution to memory management in the area of computer systems, where multiple processes with varying

memory requirements coexist. Paging splits both physical and virtual memory into fixed-size units called pages rather than allocating memory in contiguous blocks. This split makes memory allocation more flexible and dynamic, thereby minimizing concerns like external fragmentation.

The fundamental concept of paging is the conversion of virtual addresses utilized by a program into physical addresses in the underlying memory modules. Each process is separated into fixed-size pages, which are then mapped to physical memory frames. The Memory Management Unit (MMU) is in responsible of this translation, ensuring that applications preserve the illusion of a continuous and vast virtual address space even when physical memory is fragmented or constrained. Paging not only allows for more efficient memory consumption, but it also streamlines the memory allocation procedure. Processes can be loaded into any available physical memory page, and the MMU handles mapping virtual pages to their physical frames.

### Segmentation

Segmentation presents a diverse and dynamic technique to tackling the issues provided by contiguous memory allocation's rigid structures. Segmentation, in contrast to other approaches, splits the memory space into logically significant pieces known as segments. Each segment corresponds to a certain portion of a program, such as code, data, or a stack, and reflects the natural organization of the program. This division enables a more natural and flexible memory structure that is consistent with modern software's modular design concepts.

**Figure 4**

Segmentation



Note: Diagram shows Hardware Architecture of Segmentation

The capacity of segmentation to adapt to the differing memory requirements of different operations is its distinguishing feature. Unlike paging, which divides memory into fixed-size blocks, segmentation allows for the dynamic nature of program structures. As processes are loaded into memory, different segments are assigned to them based on their functional components. Because of this adaptability, segmentation is especially well-suited for applications with fluctuating memory requirements or those that increase dynamically during operation.

### *Dynamic Partitioning*

Dynamic Partitioning, also known as Variable Partitioning, is a memory management approach developed to solve the issues provided by prior systems' fixed-size divisions. In contrast to static methods, dynamic partitioning allows memory to be allocated in variable-sized blocks, allowing processes with varying memory requirements to be accommodated. This

versatility is especially useful in contexts where the memory requirements of programs change greatly, allowing for more efficient resource utilization.

The ability of dynamic partitioning to establish variable-sized partitions based on the actual memory requirements of entering processes is its defining feature. When a process arrives, it is assigned a partition that is exactly the size of the process, reducing both internal and external fragmentation that is common with fixed-size partitioning techniques. This adaptability improves memory consumption, particularly in situations where processes have uneven memory footprints.

*Buddy System*

The Buddy System, a memory allocation mechanism, presents a novel and effective approach to memory management in computer systems. By dividing memory into 2^k blocks, this technique, also known as the Binary Buddy System, tackles the issues of fixed and variable partitioning. The Buddy System's simplicity and efficacy make it a tempting alternative in situations where predictability and flexibility are valued.

The Buddy System allocates memory in blocks that are powers of two, resulting in a binary tree structure. Each node in the tree represents a memory block, and nodes that are direct siblings are referred to as "buddies." When the system receives a memory request, it searches the binary tree for the smallest accessible block that can meet the request. If the chosen block is larger than required, it is divided into two mates. If the selected block is too tiny, the algorithm looks for adjacent mates to merge, resulting in a larger block.

*Memory Compaction*

Memory compaction is a strategic memory management approach used to combat the problems associated with fragmentation in a computer system. Internal or external fragmentation can impede proper memory utilization, resulting in inefficiencies and performance loss over time. Memory compaction solves this issue by reorganizing occupied memory on a regular basis, combining unoccupied regions, and decreasing fragmentation to improve overall system efficiency.

Memory compaction's primary goal is to restructure memory in a way that lowers fragmentation and promotes contiguous free space. This procedure entails moving assigned memory to one end of the memory space while effectively consolidating free memory at the other. Memory compaction guarantees that bigger contiguous blocks are available for allocation, reducing external fragmentation and allowing the system to support larger operations.

*Page Replacement Algorithms*

Page replacement algorithms are critical in the world of virtual memory because they provide a mechanism for determining which pages to stay in physical memory and which to shift out to secondary storage. These techniques try to maximize performance by limiting the impact of page faults in systems that use virtual memory, where a process's address space can exceed the available physical memory. Several page replacement algorithms have been created, each with its approach to balancing efficiency and computational overhead.

The Least Recently Used (LRU) algorithm is a popular page replacement technique. LRU works on the premise that pages that haven't been accessed in a long time are the least likely to be used in the near future. It remembers the order in which pages are accessed and, when a page fault occurs, replaces the page that has been inactive for the longest time. While LRU is simple and frequently successful, its implementation can be computationally expensive, especially in large-scale systems.

## Virtual Memory

Virtual memory is a clever abstraction layer that allows programs to operate under the assumption of a vast and contiguous address space, regardless of actual memory limits. By transparently swapping data between RAM and secondary storage, it enables programs to use more memory than is physically available. This illusion of ample memory space improves the capabilities of modern operating systems, allowing them to run several processes concurrently without requiring a large amount of actual memory.

The separation of logical or virtual addresses utilized by a program from physical addresses in the underlying hardware is the fundamental premise of virtual memory. When a program refers to a memory address, the Memory Management Unit (MMU) converts it into a corresponding physical address, allowing data to be retrieved or stored in the actual memory modules. The operating system manages this translation process, which allows applications to run as if they have direct access to a contiguous block of memory, even if the physical memory is fragmented.

One of the key benefits of virtual memory is that it allows for efficient multitasking. Multiple processes can run in parallel, each with the illusion of a dedicated and large memory area. The operating system regulates data shifting between physical memory and secondary storage, ensuring that each process receives an equitable portion of the available resources. Virtual memory improves system stability by offering features such as memory protection and process isolation. Each process operates within its own virtual address space, preventing unintentional interaction with other processes' memory. This isolation improves system stability by lowering the possibility of a single misbehaving process influencing the entire system.

## Conclusion

Memory management in computer systems is a multidimensional discipline that involves numerous strategies aimed at maximizing memory resource consumption. A fundamental strategy, contiguous memory allocation, assigns each process a single, unbroken block of memory. While easy, it might lead to internal and external fragmentation concerns. Paging, on the other hand, splits memory into fixed-size chunks, allowing for more effective memory utilization and reducing external fragmentation. To accommodate the changing nature of program architectures, segmentation adds a hierarchical structure by splitting memory into segments. Dynamic partitioning, also known as variable partitioning, responds to changing memory needs while reducing internal fragmentation. The Buddy System uses binary pals to balance predictability and adaptability by allocating memory in powers of two.

Memory compaction reorganizes memory contents on a regular basis, combining vacant regions to reduce fragmentation and improve efficiency. Page replacement algorithms are

critical in virtual memory systems because they choose which pages to stay in physical memory and which to swap out. LRU, FIFO, and optimum algorithms strive to improve performance by reducing page faults. The Translation Lookaside Buffer (TLB) is an important component of the Memory Management Unit (MMU) that speeds up address translation by storing frequently requested mappings.

## References

[1] Dr. Vivek Chaplot "Cache Memory: An Analysis on Performance Issues" HEAD, Dept. of Computer Science Bhupal Nobles'University Volume 4, Issue 7, July 2016

[2] Mohit Saxena and Michael M. Swift "FlashVM: Revisiting the Virtual Memory Hierarchy" Department of Computer Sciences University of Wisconsin-Madison 2009

[3] Nilesh Vishwasrao and Prabhudev Irabashetti, "Dynamic Memory Allocation: Role in Memory Management", International Journal of Current Engineering and Technology, Vol. 4, No. 2, April 2014.

[4] Bhat, P. C. P. Operating Systems/Memory management. Bangalore. Lecture Notes, 2004. Web.

[5] Dipti Diwase, Shraddha Shah, Tushar Diwase and, Priya Rathod, "urvey Report on Memory Allocation Strategies for Real Time Operating System in Context with Embedded Devices", International Journal of Engineering Research and Applications (IJERA), Vol. 2, Issue 3, May-Jun 2012, pp.1151-1156.

[6] Valtteri Heikkilä, "A Study on Dynamic Memory Allocation Mechanisms for Small Block Sizes in Real-Time Embedded Systems", University of Oulu Department of Information Processing Science, conference 17, December 2012.

[7] Nikola Zlatan "Computer Memory, Applications and Management". p.30 2016

[8] Sparsh Mittal A "Survey of Techniques for Designing and Managing CPU Register file" Concurrency and computation: practice and experience Concurrency Computat.: Pract. Exper. 2016; 00:1–23

[9] M. Tofte and J.-P. Talpin, "Region-based memory management," Information and computation, vol. 132, pp. 109-176, 1997.

[10] C. A. Waldspurger, "Memory resource management in VMware ESX server," ACM SIGOPS Operating Systems Review, vol. 36, pp. 181-194, 2002.

[11] J. C. Lau, S. C. Roy, D. L. Callaerts, and I. E. N. Vandeweerd, "Method and apparatus for allocation and management of shared memory with data in memory stored as multiple linked lists," ed: Google Patents, 1999.

[12] P. R. Wilson, M. S. Johnstone, M. Neely, and D. Boles, "Dynamic storage allocation: A survey and critical review," in Memory Management, ed: Springer, 1995, pp. 1-116.

[13] K. Moronaga and M. Watanabe, "Storage management system for memory card using memory allocation table," ed: Google Patents, 1993.

[14] D. F. Hooper, G. Wolrich, M. J. Adiletta, and W. R. Wheeler, "Method for memory allocation and management using push/pop apparatus," ed: Google Patents, 2003.

[15] D. Gay and A. Aiken, Memory management with explicit regions vol. 33: ACM, 1998.