
Efficient Synchronization of Linux Memory Regions over a Network: A Comparative Study and Implementation

A user-friendly approach to application-agnostic state synchronization

Felicitas Pojtinger (Stuttgart Media University)

2023-08-04

Abstract

This study presents a comprehensive comparison and implementation of various methods for synchronizing memory regions in Linux systems over a network. Four approaches are evaluated: (1) handling page faults in userspace with `userfaultfd`, (2) utilizing `mmap` for change notifications, (3) hash-based change detection, and (4) custom filesystem implementation. Each option is thoroughly examined in terms of implementation, performance, and associated trade-offs. The study culminates in a summary that compares the options based on ease of implementation, CPU load, and network traffic, and offers recommendations for the optimal solution depending on the specific use case, such as data change frequency and kernel/OS compatibility.

Contents

1	Introduction	3
2	Technology	3
2.1	The Linux Kernel	3
2.2	Linux Kernel Modules	4
2.3	UNIX Signals and Handlers	4
2.4	Memory Hierarchy	5
2.5	Memory Management in Linux	6
2.6	Swap Space	7
2.7	Page Faults	8
2.8	<code>mmap</code>	9
2.9	<code>inotify</code>	9
2.10	Linux Kernel Disk and File Caching	10
2.11	TCP, UDP and QUIC	11
2.12	Delta Synchronization	12
2.13	File Systems In Userspace (FUSE)	12
2.14	Network Block Device (NBD)	13
2.15	Virtual Machine Live Migration	14
2.15.1	Pre-Copy	14
2.15.2	Post-Copy	15
2.15.3	Workload Analysis	16
2.16	Streams and Pipelines	17
2.17	gRPC	18
2.18	Redis	19

2.19 S3 and Minio	19
2.20 Cassandra and ScyllaDB	20

THIS IS A LLM-GENERATED TEXT VERSION OF THE NOTES FOR ESTIMATING THE EXPECTED THESIS LENGTH. THIS IS NOT THE FINISHED DOCUMENT AND DOES NOT CITE SOURCES.

1 Introduction

- Research question: Could memory be the universal way to access and migrate state?
- Why efficient memory synchronization is the missing key component
- High-level use cases for memory synchronization in the industry today

2 Technology

2.1 The Linux Kernel

Introduced by Linus Torvalds in 1991, the Linux Kernel forms the foundation of operating systems in myriad devices worldwide, with millions utilizing it in their servers, desktop computers, mobile phones, and a spectrum of embedded devices. Notably, the Linux Kernel is an open-source kernel, a quality that imparts transparency, collaboration, and flexibility at its core.

A fundamental facet of the Linux Kernel's functionality is its role as a conduit between software applications and the essential hardware of the computer. It manages the system's resources, allowing the software to interact with the hardware free from undue complexity. Precisely, it facilitates the sending and receiving of orders between the above-mentioned entities via an abstraction layer.

The abstraction layer significantly simplifies these interactions for developers, translating high-level language applications into low-level hardware instructions, and vice versa. It further ensures that any information or instructions reaching the hardware components fall within acceptable parameters to prevent malfunctions.

Moreover, the Linux Kernel's tremendous versatility is evident in its compatibility with a diverse range of architectures. It supports everything from industry standards, such as x86 and ARM, which are common in personal computers and mobile devices, to newcomers like RISC-V, an open standard instruction set architecture designed for processors.

The open-source nature of the Linux Kernel underscores its fitness as a topic for this thesis. Its source code is freely accessible to anyone, facilitating visibility and transparency. This accessibility nurtures an ethos of collaborative development, where developers worldwide can view the source code, modify it, and contribute their improvements back to the community. The power of this collaborative openness drives the continuous evolution and robustness of the Linux Kernel, ensuring it remains at the forefront of technology innovation.

2.2 Linux Kernel Modules

The Linux Kernel, characterized by its monolithic structure, owes its extendability to kernel modules. This flexibility is made possible by these small chunks of code that operate at the kernel level. Kernel modules can be loaded and unloaded as needed, acting as plug-ins that provide additional functionality to the kernel.

One salient feature of kernel modules is their ability to augment kernel functionality in live mode, freeing users from the necessity of rebooting their systems whenever there's a need to add or remove specific functions. This feature is enabled because they are dynamically linked into the running kernel, meaning they can be added or removed without disturbing the kernel's main body of code.

Kernel modules help maintain the kernel's overall size within manageable bounds. Their isolation from the core kernel simplifies the task of pinpointing and troubleshooting issues. Plus, modules help to maintain scalability, allowing the kernel to accommodate additional functionalities effectively.

Programmed in the C language, kernel modules interact with the kernel by using Application Programming Interfaces (APIs). These APIs provide a predefined set of functions for the modules to use and interact with the main kernel, forming a structured mechanism for extension and interaction.

Despite these benefits, kernel modules come with potential risks. Poorly written or insufficiently tested modules can cause the kernel to destabilize, affecting the system's overall performance. Therefore, developers must take care when crafting kernel modules and ensure thorough testing is conducted before they're integrated into the kernel.

In terms of operational management, kernel modules can be loaded either during the system boot-up or dynamically when the system is running. This is made possible through commands like 'modprobe' to add a module, and 'rmmod' to remove one.

Furthermore, a module's lifecycle is managed and streamlined by the implementation of initialization and cleanup functions. An initialization function is launched when a module is loaded, while a cleanup function is triggered during its unloading. These functions ensure resource allocation and deallocation, respectively, paving the way for smooth operational transitions.

2.3 UNIX Signals and Handlers

UNIX Signals and Handlers are a vital component of the UNIX operating system. These elements facilitate a range of functionalities and actions within the system, enabling an intricate dance of software activities to occur smoothly and efficiently.

The UNIX Signals can be conceptualized as software interrupts that are crucial in notifying a process about significant events such as exceptions. These Signals are not arbitrarily generated but rather

stem from specific sources within the system. These can include the kernel (the core part of an operating system that controls all of its major functions), user input, or different processes running within the system. Thus, they play a substantial role in the functioning of the UNIX system by acting as an asynchronous communication mechanism between processes or the kernel and a process.

An interesting feature of these Signals is that they have default actions. This means that, unless specified otherwise, they will adhere to a predetermined behaviour. For instance, a Signal could lead to the termination of a process, or conversely, it could result in the process ignoring the Signal. This default action will depend on the type of Signal and the context in which it is operating.

Despite their significant role in system functioning, Signals have certain limitations. They are not specifically designed as an Inter-Process Communication (IPC) mechanism. While they can effectively alert a process of an event, they do not possess the capability to transmit any additional data for it. Therefore, while they play a crucial role in notification and alerting mechanisms, they do not serve as a comprehensive data communication tool within the system.

UNIX Handlers, on the other hand, offer a distinct but complementary function. They are tools that enable customization in how a process responds to a Signal. Through the use of Handlers, system administrators or users can specify the response that a particular Signal should trigger in a process, offering a degree of flexibility and control that enhances system functionality.

Handlers can be installed in a UNIX system using the function `sigaction()`. This function serves to replace the current Signal handler in the process for the given Signal. This allows the system to modify the way it responds to specific Signals, thus tailoring its actions based on the needs and requirements of the moment.

2.4 Memory Hierarchy

The concept of a Memory Hierarchy is pivotal in computer systems, facilitating the efficient storage and retrieval of data. This hierarchy is a system of categorizing memory based on key parameters such as size, speed, cost, and proximity to the Central Processing Unit (CPU).

A foundational principle of this hierarchy is the 'Principle of Locality', which stipulates that the most frequently accessed data and instructions should be stored closest to the CPU. The reason for this is linked to the 'speed of the cable'. Due to physical limitations like signal dampening and the speed of light, throughput and latency decrease as distance increases. Thus, frequently used data needs to be as near to the CPU as possible to maintain system efficiency.

At the top of this hierarchy are Registers, small storage units that exist closest to the CPU. Typically, they store only 32 to 64 bits of data, a minuscule amount compared to other storage types. However, they compensate for this with their incredibly high speed, which enables the CPU to perform operations swiftly.

Cache Memory, divided into L1, L2, and L3, is the next step in the hierarchy. With each ascending level, the cache becomes larger, slower, and less expensive. It serves as a buffer for frequently accessed data, and predictive algorithms are employed to optimize its usage.

Main Memory, or Random Access Memory (RAM), is the next level down. It provides a larger capacity than cache but operates at a slower speed. RAM plays a crucial role in storing programs and open files that are actively being used by the computer system.

Following RAM, we have Secondary Storage, which includes Solid State Drives (SSD) and Hard Disk Drives (HDD). These storage types are slower than RAM but can store much larger amounts of data. They are typically used to store the operating system and other critical system files and are persistent, meaning they retain data even after power is cut.

The final layer of the traditional hierarchy is Tertiary Storage, comprising optical disks and tape. This level is slow in terms of data access but offers a very cost-effective way to store large amounts of data. Tapes, for example, can hold vast quantities of data for relatively long periods, making them suitable for archiving or physical data transport, such as importing data from personal infrastructure to a cloud service like Amazon Web Services (AWS).

However, it's worth noting that this clear cut hierarchy is becoming blurred due to technological advancements. For instance, NVMe (Non-Volatile Memory Express) technology can rival the speeds of RAM while storing larger amounts of data. Furthermore, innovations are enabling the exposure of secondary or tertiary storage with the same interface as main memory, further blending the lines of this hierarchy. Despite these evolutions, the fundamental principle of locality continues to hold its ground, making the memory hierarchy a pivotal consideration in computer architecture.

2.5 Memory Management in Linux

Memory management is a fundamental component of an operating system's functionality, perhaps even the essential purpose of an operating system itself. In the Linux operating system, memory management is intricately designed to ensure efficient and secure operations.

In broad terms, memory management in Linux creates a buffer between applications and physical memory. This buffer enables the operating system to control and coordinate access to physical memory, mitigating potential conflicts and enhancing overall system performance. Moreover, memory management provides robust security measures, ensuring that each process can only access its dedicated memory space. This compartmentalization helps prevent unauthorized access and safeguards the integrity of each process.

Linux's memory management system is divided into two primary domains: Kernel space and User space. The Kernel space is where the core components of the Linux operating system, known as the

kernel, operate. This space is also utilized by kernel extensions and device drivers. The management of Kernel space is performed by the kernel memory module, a specialized component designed to handle this critical task.

A technique known as slab allocation is used in the Kernel space. This method groups objects of the same size into structures known as caches, which significantly speeds up the process of memory allocation. Notably, slab allocation also reduces fragmentation within the memory, improving overall efficiency and making better use of available resources.

On the other hand, User space is where applications store their memory. This includes most software applications that a user might interact with, along with some drivers. Management of User space is handled through a paging system, a technique used by most modern operating systems for memory management.

In this system, each application is assigned its private virtual address space. This segregation ensures that applications do not interfere with each other, enhancing system stability and security. Moreover, the virtual address space is divided into small units known as pages, each typically comprising 4 kilobytes of data. An interesting aspect of this paging system is that each page can be mapped to any location in physical memory, providing a high degree of flexibility in memory usage.

2.6 Swap Space

Swap Space represents a designated portion of the secondary storage on a computer system, serving as a form of virtual memory. It is an integral part of systems that run multiple applications simultaneously, as it assists in efficient memory management by providing additional space for memory-intensive operations.

The functionality of Swap Space is grounded in moving inactive parts of the system's Random Access Memory (RAM) to secondary storage, thereby freeing up RAM for other processes. This swapping process occurs when RAM is nearing its capacity, ensuring that the system continues to function smoothly without running out of memory.

In the context of the Linux operating system, Swap Space implementation leverages a demand paging system. This system allows memory to be allocated only when necessary, providing an effective strategy for memory management. In Linux, Swap Space can take the form of a swap partition or a swap file. A swap partition is a separate area of the secondary storage dedicated to swap operations, while a swap file is a regular file within the system that can be expanded or truncated as required. Both swap partitions and swap files are transparent to the user, facilitating seamless system performance.

The Linux kernel uses a Least Recently Used (LRU) algorithm to decide which pages to swap, prioritizing those that have not been accessed recently. This algorithm ensures efficient use of Swap Space, minimizing the potential for wasted storage.

Swap Space plays a significant role during system hibernation. Before entering hibernation, the system transfers the content of RAM into Swap Space, providing a persistent storage location. Upon resuming, the system retrieves the saved memory from Swap Space and reloads it into RAM.

However, the usage of Swap Space can influence system performance. If swap operations are performed too frequently due to insufficient RAM, it can lead to slowdowns because secondary storage is generally slower than primary memory. To control this, Linux allows the configuration of a “swappiness” parameter for the kernel, dictating how likely the system is to swap memory pages. The swappiness value can be adjusted to balance system performance with memory usage requirements.

2.7 Page Faults

A page fault is a type of interrupt in a computer system, which occurs when a process attempts to access a page that is not currently available in primary memory. When this happens, the operating system takes action to swap the required page from secondary storage into primary memory, thus resolving the page fault.

Page faults can be categorized into two primary types: minor and major. A minor page fault occurs when the page is already loaded in memory, but it is not currently linked to the process that requires it. These can be resolved relatively quickly as the data is already present in memory. In contrast, a major page fault takes place when the page needs to be loaded from secondary storage, which can be a more time-consuming process due to the slower speed of secondary storage compared to primary memory.

Algorithms such as Least Recently Used (LRU) and the simpler clock algorithm are typically used to manage memory and minimize the occurrence of page faults. These algorithms prioritize the pages that have been accessed recently, ensuring that the most actively used data remains readily available in memory.

Certain techniques can be employed to handle page faults efficiently. One such technique is prefetching, which involves anticipating future page requests and loading these pages into memory in advance. This proactive strategy helps to reduce the likelihood of major page faults. Page compression is another method used, which involves compressing inactive pages and storing them in memory preemptively, again reducing the chance of major page faults.

Observing page faults can provide insights into the memory access patterns of processes. For instance, if a process triggers a page fault when trying to access a specific piece of memory, the system can retrieve that chunk of memory from a remote location and map it to the address on which the page fault occurred. This on-demand fetching of data helps to optimize memory usage by only loading data when it is required.

Page faults are generally handled by the kernel, the core part of an operating system. However, in the past, it was also possible to manage page faults from userspace by handling the `SIGSEGV` signal in the process. Regardless of where they are handled, the efficient management of page faults is critical for maintaining system performance and minimizing interruptions to running processes.

2.8 mmap

The UNIX system call `mmap` is a tool used for mapping files or devices into memory. It has multiple potential applications including shared memory, file input/output (I/O), and fine-grained memory allocation. Its functionality makes it a popular choice in applications such as databases. While it can be seen as a “power tool”, it should be used judiciously and intentionally due to its significant effect on system performance.

Functionally, `mmap` creates a direct link, known as a memory mapping, between a file and a region of memory. When the system reads from this mapped memory region, it directly reads from the associated file and conversely, when it writes to the memory region, it writes to the file. This reduces system overhead as it lessens or eliminates the need for context switches, which are typically resource-intensive operations.

`mmap` offers several benefits. One significant advantage is the facilitation of zero-copy operations. With a memory mapping, data can be accessed directly as if it were in memory, without the need to first copy it from the disk. This greatly improves efficiency in scenarios where large amounts of data are being processed. Additionally, `mmap` can be used to share memory between processes, avoiding the need to pass through the kernel with system calls, which again enhances performance.

On the flip side, a notable drawback of `mmap` is that it bypasses the file system cache. This can result in stale data if multiple processes are reading from and writing to the same data concurrently. Without the benefit of caching, which helps to ensure data consistency across different processes, the use of `mmap` can potentially lead to issues with data integrity.

2.9 inotify

`inotify` is an event-driven notification system that is part of the Linux kernel. Its purpose is to monitor the file system for specific events such as modifications, access, and others. By providing real-time notifications of these events, `inotify` allows applications to respond to changes in the file system promptly and effectively.

One of the main features of `inotify` is its ability to use “watches” to monitor specific events. For example, it can be configured to only watch for write operations to a certain file or directory. This

watch feature is highly configurable, allowing applications to monitor only the events that are relevant to them, thus improving overall efficiency.

One significant benefit of `inotify` is that it reduces overhead and resource usage compared to polling. Polling involves frequently checking the status of a file or directory, which can be inefficient in terms of CPU usage, particularly when changes are infrequent. `inotify`, on the other hand, operates in an event-driven manner, meaning it only becomes active when a monitored event occurs.

`inotify` is widely used in a variety of applications. For example, Dropbox utilizes `inotify` to synchronize files between the local machine and the cloud. Whenever a file that Dropbox is monitoring is modified, `inotify` alerts Dropbox so it can initiate the synchronization process.

Despite its usefulness, `inotify` does have some limitations. One of the notable constraints is the limit on the number of watches that can be established. This can become an issue in applications that require monitoring a large number of files or directories simultaneously. Workarounds for this limitation exist, but they can add complexity to the application.

2.10 Linux Kernel Disk and File Caching

The Linux operating system deploys sophisticated strategies for optimizing system performance and speed through the implementation of disk and file caching. This essay explores these methods, their complexities, and their contributions to efficient system operations.

Disk caching refers to the temporary storage of frequently accessed data in Random Access Memory (RAM), following the principle of locality drawn from the memory hierarchy. This concept suggests that data recently or frequently accessed are likely to be used again, justifying their storage in a faster, albeit smaller, memory for swift access. Within the Linux system, this technique is actualized via the page cache subsystem, which utilizes a Least Recently Used (LRU) algorithm for cache management. This algorithm dictates that the least recently accessed data will be replaced when the cache is full and new data need to be accommodated, allowing for dynamic cache content adjustment based on usage patterns.

File caching, on the other hand, focuses on the conservation of file system metadata within two cache structures: the 'dentry' and 'inode' caches. Metadata provides essential details about files such as their names, attributes, and locations. By caching this information, Linux significantly enhances the speed of path name resolution and file attribute determination, which can include the retrieval of the last change data for polling. In essence, this method of caching serves to streamline file reads and writes as they transit through the disk cache.

However, despite their efficiency-enhancing roles, these caching methods also pose complexities. The challenge of data consistency surfaces as the system attempts to maintain synchronization be-

tween the disk and cache through a process termed writebacks. This process can pose a dilemma: aggressive writebacks, though ensuring optimal data consistency, can hamper overall system performance. Conversely, longer delays between writebacks while attempting to maximize performance may risk potential data loss.

Additionally, the need to release cached data under conditions of memory pressure further complicates caching. To prevent a system crash due to insufficient memory, cached data must be effectively evicted when memory availability dips below certain thresholds. This action necessitates the use of intelligent cache eviction algorithms such as LRU, which, while beneficial, introduce additional layers of complexity to the system's memory management functions.

2.11 TCP, UDP and QUIC

Transmission Control Protocol (TCP), User Datagram Protocol (UDP), and Quick UDP Internet Connections (QUIC) constitute the trinity of transport layer protocols that are vital for managing internet communication. Each of these plays distinct roles based on their unique characteristics, and they come together to form the bedrock of the diverse needs of the internet.

TCP is a connection-oriented protocol and has historically acted as the reliable backbone of internet communication. Its features include guaranteed delivery of data and maintenance of the order in which data packets are sent and received. This is made possible due to the inclusion of error checking, lost packet retransmission, and congestion control mechanisms. As a result, TCP powers the majority of the web, underpinning the delivery of emails, web pages, and many other forms of data transmission where accuracy and completeness are paramount.

In contrast, UDP is a connectionless protocol that does not guarantee reliable or ordered packet delivery. This might initially seem like a disadvantage, but it makes UDP faster than TCP because it comes with fewer assurances. As a result, UDP is well suited for applications that prioritize speed over reliability. Examples include online gaming and video calls, where occasional loss of packets does not significantly impact the user experience, and real-time response is more critical.

QUIC, a more recent entrant into the space of transport layer protocols, was developed by Google and standardized by the Internet Engineering Task Force (IETF) in 2020. QUIC aims to incorporate the best features of both TCP and UDP. Like TCP, it provides reliability and ordered delivery guarantees, but it also introduces several improvements to overcome the limitations of TCP.

One of the significant enhancements that QUIC offers is reduced connection establishment time, also known as initial latency. QUIC achieves this by merging the connection and security handshakes into a single step, resulting in quicker and more efficient connection establishments. Another crucial advantage of QUIC is its ability to avoid head-of-line blocking. In TCP, if a packet is lost, it halts the delivery of all subsequent packets until the lost packet is retransmitted and received. This can cause

significant delays. QUIC mitigates this issue by allowing the independent delivery of separate streams, thereby preventing a lost packet from blocking the entire transmission line.

2.12 Delta Synchronization

Delta synchronization, a novel technique for synchronizing files between hosts, provides a way to circumvent the traditional method of transferring entire files by instead sending only the parts of a file that have changed. This method promises the reduction of network and Input/Output (I/O) overhead, thereby optimizing the overall efficiency of data transfer.

A popular tool for implementing this file synchronization technique is rsync. When rsync's delta-transfer algorithm is active, it calculates the difference between a local file and its remote counterpart and then synchronizes the changes accordingly. This synchronization process is primarily split into several distinct steps.

Initially, the delta sync algorithm divides the file on the destination into fixed-size blocks. This process, known as file block division, lays the groundwork for subsequent operations. For each block in the divided file, a weak and fast checksum is calculated. These checksums act as unique identifiers for the respective file blocks.

Following the checksum calculations, these identifiers are transferred to the source host. On the source side, the same checksum calculation process is conducted. The resulting checksums are then compared with those sent over from the destination host in a step known as matching block identification.

After this comparison step, the blocks that have changed can be pinpointed. The source then dispatches the information pertaining to the offset, or the position in the file, of each changed block, as well as the data contained in these blocks, to the destination.

Upon receipt of the changed blocks and their corresponding offsets, the destination host writes these chunks of data to the specified offsets. This process results in the reconstruction of the file with the updated data.

Once a polling interval, or the time between two successive data requests, is completed, the entire process begins again. This continuous cycle ensures that the local and remote files are kept in sync over time.

2.13 File Systems In Userspace (FUSE)

File Systems in Userspace (FUSE) represent an innovative software interface that provides developers with the capacity to design custom file systems in the userspace, thus eliminating the need for intricate and potentially error-prone low-level kernel development. This approach avails a more flexible

and safer method to implement file systems, which has implications for the broader software development landscape.

The FUSE interface is platform-agnostic and supports multiple operating systems, including but not limited to Linux, macOS, and FreeBSD. To establish file systems in the userspace, developers can harness the FUSE API (Application Programming Interface), where a userspace program is registered with the FUSE kernel module.

This program is designed to provide callbacks for various file system operations such as ‘open’, ‘read’, and ‘write’. Whenever a user performs an operation on a mounted FUSE file system, the kernel module sends a request to the userspace program. This program then processes the request and returns a response via the FUSE kernel module back to the user, establishing a communication loop that allows for efficient and flexible file system operation.

The advantage of this approach is twofold: it simplifies the process of creating a file system by relocating it to the userspace, and it significantly enhances system safety. The latter benefit is primarily due to the fact that custom kernel modules are not required for operation. As such, any error occurring within the FUSE system or the backend will not result in a complete kernel crash, thus providing a higher level of stability compared to conventional kernel-based file systems.

In addition, unlike file systems that are implemented as kernel modules, FUSE’s layer of indirection ensures portability since the file system only needs to communicate with the FUSE module, irrespective of the underlying platform.

Despite these advantages, there is a trade-off in the form of performance overhead due to the context switching between the kernel and the userspace file system. This overhead, however, can be considered a reasonable cost for the benefits provided by FUSE.

FUSE has seen wide application in many high-level interfaces to external services. For instance, it is used to mount S3 buckets or access a remote system’s disk via SSH. By allowing developers to operate in the safer and more flexible environment of the userspace, FUSE has played a key role in the advancement of file system development, thus shaping the way we interact with and manage data in the digital age.

2.14 Network Block Device (NBD)

The Network Block Device (NBD) protocol is a system communication framework between a server and a client, provided by the NBD kernel module. While this protocol can operate over a Wide Area Network (WAN), its design is more optimally suited to a Local Area Network (LAN) or localhost usage. This is due to the fact that its performance characteristics and operational limits are better aligned with the lower latencies and higher data rates typically associated with these types of networks.

The NBD protocol encompasses two phases: the handshake and the transmission phase. This protocol includes several actors such as one or multiple clients, a server, and the virtual concept of an “export” – a data structure that the server provides for the client to interact with.

In the handshake phase, the client connects to the server, and the server then sends a greeting message, including the server’s flags. Subsequently, the client responds with its own flags and an export name. The server then provides the export’s size and other associated metadata. Upon the client acknowledging the receipt of this data, the handshake is deemed complete.

The transmission phase commences following the handshake. During this phase, the client and the server exchange commands and replies, with the commands encompassing basic operations required for block device access such as read, write, or flush. These commands may also include other information such as data chunks, offsets, lengths, and more. The replies, on the other hand, may contain an error code, success value, or data, contingent on the type of reply.

Despite its utility, NBD has its limitations, particularly with regard to message sizes. The maximum message size is 32 MB, while the maximum block or chunk size supported by the kernel is limited to 4096 KB. This constraint can make NBD less ideal for operations over a WAN, especially in high latency scenarios.

NBD does offer the ability to list exports, enabling the mapping of multiple memory regions on a single server. However, it should be noted that NBD is an older protocol with numerous versions and legacy features. In most use cases, it is recommended to only implement the latest versions and the baseline feature set.

NBD’s simplicity also gives rise to certain drawbacks. For instance, it is not suitable for use cases where the backing device behaves differently from a random-access store device, such as a tape drive. This is because NBD operates at the block level, without high-level abstractions such as files or directories. Furthermore, it does not support shared access to the same file for multiple clients.

In specific scenarios, however, these limitations may be less significant. For instance, in the case of memory synchronization, the low-level, direct nature of NBD can be more of a feature than a bug. Its block-level operation aligns well with the nature of this use case, providing a straightforward and efficient method for data communication.

2.15 Virtual Machine Live Migration

2.15.1 Pre-Copy

Virtual Machine Live Migration is a process that entails moving a virtual machine, including its state and connected devices, from one host to another. The goal of this process is to minimize downtime as much as possible, thereby ensuring a smooth transition and minimizing disruptions to the operations

of the virtual machine. In this field, the objective is to further optimize systems to achieve even lower downtimes than what can be achieved using Network Block Device (NBD) protocols over Wide Area Networks (WAN).

One such advanced methodology is the ‘managed mount API’, which is part of a larger set of techniques to optimize virtual machine live migration. Two primary types of migration algorithms fall under this umbrella: pre-copy and post-copy migration.

Pre-copy migration is a process wherein the data from the source is copied over to the destination while the virtual machine, or any other application or state, continues to run. The process involves copying the initial state of the virtual machine’s memory to the destination first. While this is happening, if chunks of data are being modified, they are flagged as ‘dirty’.

Following this, these dirty chunks are then transferred to the destination. This process continues until the number of remaining chunks that need to be transferred falls below a threshold that satisfies a maximum acceptable downtime criteria. This threshold indicates that it is a good time to briefly pause the virtual machine on the source.

Subsequently, the virtual machine is suspended on the source, and the remaining chunks are synced over to the destination. Upon completion of this transfer, the virtual machine is resumed on the destination. The entire process ensures that the virtual machine is always fully available on either the source or the destination.

An advantage of pre-copy migration is that it is resilient to a network outage during the synchronization process, as it ensures that the full state of the virtual machine is always available on either the source or the destination.

However, there are potential limitations with pre-copy migration as well. For instance, if the virtual machine or application changes too many chunks on the source during the migration, it might prevent the maximum acceptable downtime criteria from being reached. Furthermore, the maximum acceptable downtime is also somewhat constrained by the available round trip time (RTT), the time it takes for a signal to be sent plus the time it takes for an acknowledgment of that signal to be received.

2.15.2 Post-Copy

Post-copy migration is another technique used for Virtual Machine (VM) live migration, serving as an alternative to the pre-copy approach. This methodology follows a different process to achieve the same goal of minimizing downtime during the movement of a VM from one host to another.

Unlike pre-copy migration, in post-copy migration, the VM is suspended on the source almost immediately. The VM is then transferred to the destination along with only a minimal set of chunks or blocks

of data. Once the VM has been transferred to the destination host, it is resumed and continues its operation.

During the operation of the VM on the destination host, if it attempts to access a chunk of data that was not included in the initial minimal set of transferred chunks, a page fault is triggered. A page fault is a type of interrupt or signal to the system that the program has requested access to data not currently in its working memory. When this happens, the missing chunk is fetched from the source host and transferred to the destination. The VM then continues to execute its operation using the newly fetched data chunk.

One of the key benefits of the post-copy migration approach is its efficiency in data transfer. It does not require the re-transmission of dirty chunks before the maximum tolerable downtime is reached, thereby potentially reducing data transfer loads.

However, post-copy migration does have its disadvantages. It can result in longer migration times since chunks of data have to be fetched from the network on-demand. This is a very latency or Round Trip Time (RTT) sensitive process. This means that the total time taken for the migration can vary significantly based on the quality and speed of the network, as well as the distance between the source and the destination. Overall, the choice between pre-copy and post-copy migration often involves trade-offs based on specific network conditions and the particular requirements of the VM being migrated.

2.15.3 Workload Analysis

An interesting resource on how to reduce overhead during virtual machine live migration through workload analysis is the study titled “Reducing Virtual Machine Live Migration Overhead via Workload Analysis.” Although primarily intended for use with virtual machines, the insights from this study could be applied to other applications or migration scenarios.

The method proposed in the study aims to identify the workload cycles of virtual machines. With the gathered information, it determines whether it would be beneficial to postpone the migration of a VM. The method analyzes cyclical patterns that could unnecessarily delay a VM’s migration, then uses that data to pinpoint the optimal cycles for initiating migration.

In the context of virtual machines, these cycles could correspond to various activities, such as a large application’s garbage collection process, which can trigger extensive changes to the VM’s memory. When a migration is proposed, the system will assess whether it is currently in an optimal cycle for migration. If it is, the migration proceeds; if not, the migration is postponed until the next favorable cycle.

A notable aspect of this algorithm is its utilization of a Bayesian classifier. This statistical model is used to identify whether the current cycle is favorable or unfavorable for migration. This method provides

a potentially significant improvement over the alternative approach, which typically involves waiting for a substantial portion of unchanged chunks to be synchronized before migration is initiated.

The authors of the study found that their method resulted in an up to 74% improvement in terms of live migration time and downtime, and a 43% reduction in the volume of data transferred over the network. While such a system was not specifically implemented for r3map, the method proposed in the study could be applied in conjunction with r3map or similar tools. This implies the possibility of making substantial improvements to the efficiency and effectiveness of live migrations.

2.16 Streams and Pipelines

Streams and pipelines comprise fundamental constructs within the broad and deep field of computer science. They provide for the sequential processing of elements, thus enabling the potent handling of voluminous data sets without necessitating the loading of everything simultaneously into memory. These features form the structural backbone of efficient and modular data processing systems.

Looking closely at streams, they represent an uninterrupted flow of data, acting analogous to a river carrying water from one point to another. The flexible nature of these streams allows them to be both a source and a destination of data. As such, they can handle a myriad of types, including but not limited to, files, network connections, and standard inputs/outputs (stdin/stdout).

A significant benefit of using streams is their capacity to process data swiftly, even as it becomes available. This characteristic is highly advantageous and concurrently minimizes memory consumption, thereby optimizing computational resources. Streams also exhibit an impressive adaptability which makes them particularly suited for long-running processes where data is served in for extended durations of time.

Turning our attention to pipelines, these are analogous to a sophisticated assembly line. They encompass a series of consecutive data processing stages where the output generated from one stage conveniently serves as the input for the successive stage. Thus, allowing data to seamlessly flow through the entire pipeline and efficiently complete the data processing task.

Pipelines yield a notable performance enhancement due to their impressive concurrency. This is due to the functional capability of multiple stages being able to run in parallel, thereby further encouraging the economy of time and resources. An archetype of pipelines can be spotted within the central processing units (CPUs) where the stages of instruction execution occur in a parallel fashion – a concept known as an instruction pipeline.

Another manifestation of pipelines that is notable to mention is the UNIX pipes. When one looks into how UNIX commands function, the output derived from a command (for instance, `curl`) can be effectively piped into another command (for instance, `jq`). This propagative flow through the pipeline

serves to achieve a larger, unified objective without sacrificing efficiency and time. The interconnected nature of these processes amazingly encapsulates the central idea behind the design and functioning of both streams and pipelines in computer science.

2.17 gRPC

gRPC is an open-source, high-performance Remote Procedure Call (RPC) framework that organizes communication between services. It surfaced from Google in 2015 and has since gained recognition for its robustness.

One of its main advantages lies within the transport protocol it employs, which is HTTP/2. This offers multiple features like header compression that reduces overhead and request multiplexing that allows multiple requests to be sent over a single TCP connection. These characteristics support a model of communication having reduced latency, a faster initiation, and optimizing the use of network resources.

An additional strength of gRPC is seen in its adoption of Protobuf (Protocol Buffers), for the Interface Definition Language (IDL) and format for transmitting data. Protobuf is a high-performance, multilingual tool for serializing structured data. It replaces ‘slow’ and ‘verbose’ JSON (JavaScript Object Notation) usually leveraged by REST (Representational State Transfer) APIs.

gRPC boasts a wide variety of supported features, including unary RPCs, which consist of a single request and response. This framework also supports server-streaming RPCs where the client sends a request and gets a stream to read a sequence of messages back. Client-streaming RPCs, where the client writes a sequence of messages and sends them to the server, once the client has finished writing the messages, it waits for the server to read them and return its response, are another aspect. Bidirectional RPCs, where both sides send a sequence of messages using a read-write stream, mesh language capabilities with robust performance.

The gRPC framework also integrates support for additional capabilities, including load balancing, tracing, health checking, and authentication. This pluggability ensures that the technology remains extensible and adaptable to future needs.

One other significant aspect of gRPC is its multilingual support, making it versatile and accessible to a wide range of developers. It acknowledges a multitude of programming languages including but not limited to Go, Rust, and JavaScript. This capability underpins the intercommunication of services implemented in different languages, facilitating a polyglot microservices architecture.

The Cloud Native Computing Foundation (CNCF), a heavyweight in the open-source software community, played pivotal roles in nurturing the development of gRPC. This august body continues to maintain and champion gRPC and other significant open-source technologies for the cloud native

ecosystem. Their commitment underlines the credibility and relevance of gRPC in today's software development landscape.

2.18 Redis

Redis is an in-memory data structure store, which can be employed as a database, cache, and message broker. The brainchild of Salvatore Sanfilippo, it was brought to life in 2009 and has since consistently demonstrated its value as a sturdy data handling and processing tool.

Its design stands apart from other NoSQL databases in that it supports a variety of data structures, not just limiting itself to key-value pairs. It includes lists, sets, hashes, and even complex ones like bitmaps. This versatility allows Redis to handle a broad scope of applications and use-cases more aptly.

At the fundamental core of Redis lies its use of in-memory data storage. This strategic design choice offers maximum speed and efficiency by storing data in the server's memory rather than in conventional disks. The choice of having an in-memory mechanism imparts Redis the benefit of expediting read/write operations, thereby allowing for low-latency access. Such a quick response time is essential for applications requiring intense database interaction.

Though the primary purpose of Redis is not persistent storage, it surprises with its ability to store data on disk. Despite primarily being an in-memory data structure store, Redis allows configurations for data to be stored indefinitely, persisting even after system reboots. This establishes Redis as a choice for applications needing a degree of data durability without compromising speed.

The non-blocking Input/Output (I/O) model indulged by Redis further enhances its performance capabilities. By not blocking processes while handling I/O operations, Redis accelerates application responsiveness, leading to near real-time data processing. This makes Redis an excellent fit for tasks demanding immediate data interaction and transmission, which is frequent in today's real-time applications.

Another notable feature of Redis is its in-built pub-sub system, reinforcing its role as a message broker. This system allows data to be communicated between different parts of an application or even different applications. The pub-sub model framework eases implementing scenarios related to message broadcasting, thereby supporting real-time services and many other applications.

2.19 S3 and Minio

S3, or Simple Storage Service, is a widely utilized object storage service suited to data-heavy workloads. This service, offered by Amazon Web Services (AWS), is renowned for its scalability and reliability.

The primary feature of S3 is its capability to host data globally, allowing for swift access times from any location around the world. This geographical distribution significantly reduces latency, contributing to superior user experiences, especially when dealing with substantial amounts of data.

Additional service benefits include a range of storage classes designed to meet different user requirements. Each of these classes offer varying levels of availability, access times, and cost structures and they can be tailored to fit the needs of several use-cases, from high-frequency access data to long-term archiving.

S3 incorporates sophisticated authentication and authorization mechanisms, bolstering its security profile. These features safeguard sensitive content and ensure that access to stored data is strictly controlled.

To enable compatibility with diverse applications, S3 also exposes an HTTP-based Application Programming Interface (API). This API offers a consistent interface for managing and retrieving data, such as folders and files, stored on the S3 service.

Minio, on the other hand, is an open-source storage server that was deliberately designed to be compatible with S3. Unlike S3, which is a managed service, Minio can be hosted on-premises, offering a powerful alternative for those seeking similar capabilities but favoring open-source offerings.

Known for its simplicity, Minio is lightweight due largely to its base language, Go. This programming language is recognized for its robustness, high performance, and efficient resource utilization.

A key strength in Minio is its inherent support for horizontal scalability. This feature allows for the easy addition of storage capacity as data volumes grow, hence ensuring consistent performance and availability. Horizontal scalability enables Minio users to store vast amounts of data across multiple nodes effectively without compromising retrieval times or overall usage.

2.20 Cassandra and ScyllaDB

Cassandra and ScyllaDB are two widely recognized wide-column NoSQL databases. They have uniquely combined elements of Amazon's Dynamo model with Google's Bigtable model, resulting in a fusion that facilitates a highly available database.

Apache Cassandra is a distinctively scalable and eventually consistent database. It has the capability to manage large quantities of data spread over many servers. Advantageously, it does not have any single point of failure, making it resilient in the ecosystem. Built in such a robust fashion, Cassandra allows tuning of consistency according to the specific requirements of use-cases. Tuning can range from eventual consistency to strong consistency, where eventual consistency guarantees high availability and strong consistency maintains a perfectly consistent data state across the database.

A significant aspect of Cassandra is its non-dependency on master nodes, due to its use of a peer-to-peer protocol and a distributed hash ring design. This design feature ensures that all nodes are treated equally, distributing data and loads equally across the cluster, thereby preventing potential bottlenecks or single points of failure. However, it's worth noting that Cassandra can manifest high latency under conditions of heavy load and necessitates a more or less complex configuration to set up and maintain.

ScyllaDB, introduced in 2015, has positioned itself to address some drawbacks associated with Cassandra. One of the notable differences between the two is the programming language they're written in. While Cassandra is written in Java, ScyllaDB has been crafted using C++, a language known for its high performance and efficiency. Furthermore, ScyllaDB uses a shared-nothing architecture, suggesting that each node is independent and self-sufficient.

Notwithstanding these differences, ScyllaDB maintains compatibility with Cassandra's API and data model, making it an easy replacement for organizations interested in a switch, with minimal disruption to existing workflows.

ScyllaDB was specifically designed to surmount some of Cassandra's perceived limitations, especially around latency, inclusive of the 99th percentile latency (P99 latency). This type of latency means that 99 percent of all request latencies fall below this level, implying that only 1 percent exceeds it. This is of particular importance in distributed systems where occasional outlier latencies can heavily impact user experience.

The performance enhancements attributed to ScyllaDB have not been self-claimed but validated through various independent benchmarking studies. These studies confirm that ScyllaDB can provide a notable improvement in database performance, ultimately placing it as a worthy alternative to Apache Cassandra for organizations requiring low latency and high performance from their database systems.