

# 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

*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.*

# 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

## Technology

---

# 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

# 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

# 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.



# 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

# 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

# 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

## 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

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

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,

# 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.

## 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.



## 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

# 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

## 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

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

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

## 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

# 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

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



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.

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

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,