

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

A user-friendly approach to application-agnostic state synchronization

Felicitas Pojtinger (Stuttgart Media University)

2023-08-04

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

- Open-Source kernel created by Linus Torvals in 1991
- Powers millions of devices worldwide (servers, desktops, phones, embedded devices)
- Is a bridge between applications and hardware
 - Provides an abstraction layer
 - Compatible with many architectures (ARM, x86, RISC-V etc.)
- Is a good fit for this thesis due to it's open-source nature (allows anyone to view, modify and contribute to the source code)

Linux Kernel Modules

- Linux kernel is monolithic, but extensible thanks to kernel modules
- Small pieces of kernel-level code that can be loaded and unloaded as kernel modules
- Can extend the kernel functionality without reboots
- Are dynamically linked into the running kernel
- Helps keep kernel size manageable and maintainable
- Kernel modules are written in C
- Interact with kernel through APIs
- Poorly written modules can cause significant kernel instability
- Modules can be loaded at boot time or dynamically (modprobe, rmmod etc.)
- Module lifecycle can be implemented with initialization and cleanup functions

- Signals
 - Are software interrupts that notify a process of important events (exceptions etc.)
 - Can originate from the kernel, user input or different processes
 - Function as an asynchronous communication mechanism between processes or the kernel and a process
 - Have default actions, i.e. terminating the process or ignoring a signal
- Handlers
 - Can be used to customize how a process should respond to a signal
 - Can be installed with `sigaction()`
- Signals are not designed as an IPC mechanism, since they can only alert of an event, but not of any additional data for it

Memory Hierarchy

- Memory in computers can be classified based on size, speed, cost and proximity to the CPU
- Principle of locality: The most frequently accessed data and instructions should be in the closest memory
- Locality is important mostly due to the “speed of the cable” - throughput (due to dampening) and latency (due to being limited by the speed of light) decreases as distance increases
- Registers
 - Closest to the CPU
 - Very small amount of storage (32-64 bits of data)
 - Used by the CPU to perform operations
 - Very high speed, but limited in storage size
- Cache Memory
 - Divided into L1, L2 and L3
 - The higher the level, the larger and less expensive a layer
 - Buffer for frequently accessed data

Memory Management in Linux

- Memory management is a crucial part of every operation system - maybe even the whole point of an operating system
- Creates buffer between applications and physical memory
- Can provide security guarantees (e.g. only one process can access it's memory)
- Kernel space
 - Runs the kernel, kernel extensions, device drivers
 - Managed by the kernel memory module
 - Uses slab allocation (groups objects of the same size into caches, speeds up memory allocation, reduces fragmentation of the memory)
- User space
 - Applications (and some drivers) store their memory here
 - Managed through a paging system
 - Each application has it's own private virtual address space
 - Virtual address space divided into pages of 4 KB
 - Pages can be mapped to any “real” location in physical memory

Swap Space

- A portion of the secondary storage is for virtual memory
- Essential for systems running multiple applications
- Moves inactive parts of ram to secondary storage to free up space for other processes
- Implementation in Linux
 - Linux uses a demand paging system: Memory is only allocated when it is needed
 - Can be either a swap partition (separate area of the secondary storage) or file (regular file that can be expanded/truncated)
 - Swap partitions and files are transparent to use
 - Kernel uses a LRU algorithm for deciding which pages
- Role in hibernation
 - Before hibernating, the system saves the content of RAM into swap (where it is persistent)
 - When resuming, memory is read back from swap
- Role on performance

Page Faults

- Page faults occur when the process tries to access a page not available in primary memory, which cause the OS to swap the required page from secondary storage into primary memory
- Types
 - Minor page faults: Page is already in memory, but not linked to the process that needs it
 - Major page fault: Needs to be loaded from secondary storage
- The LRU (and simpler clock algorithm) can minimize page faults
- Techniques for handling page faults
 - Prefetching: Anticipating future page requests and loading them into memory in advance
 - Page compression: Compressing inactive pages and storing them in memory pre-emptively (so that less major faults happen)
- By listening to these page faults, we know when a process wants to access a specific piece of memory
- We can use this to then pull the chunk of memory from a remote,

- Overview
 - UNIX system call for mapping files or devices into memory
 - Multiple possible usecases: Shared memory, file I/O, fine-grained memory allocation
 - Commonly used in applications like databases
 - Is a “power tool”, needs to be used carefully and intentionally
- Functionality
 - Establishes direct link (memory mapping) between a file and a memory region
 - When the system reads from the mapped memory region, it reads from the file directly and vice versa
 - Reduces overhead since no or less context switches are needed
- Benefits:
 - Enables zero-copy operations: Data can be accessed directly as though it were in memory, without having to copy it from disk first
 - Can be used to share memory between processes without having to go through the kernel with syscalls

- Event-driven notification system of the Linux kernel
- Monitors file system for events (i.e. modifications, access etc.)
- Uses a watch feature for monitoring specific events, e.g. only watching writes
- Reduces overhead and resource use compared to polling
- Widely used in many applications, e.g. Dropbox for file synchronization
- Has limitations like the limit on how many watches can be established

Linux Kernel Disk and File Caching

- Disk caching
 - Temporarily stores frequently accessed data in RAM
 - Uses principle of locality (see Memory Hierarchy)
 - Implemented using the page cache subsystem in Linux
 - Uses the LRU algorithm to manage cache contents
- File caching
 - Linux caches file system metadata in the dentry and inode caches
 - Metadata includes i.e. file names, attributes and locations
 - This caching accelerates the resolution of path names and file attributes (i.e. the last change data for polling)
 - File reads/writes pass through the disk cache
- Complexities
 - Data consistency: Between the disk and cache via writebacks. Aggressive writebacks lead to reduced performance, delays risk data loss
 - Release of cached data under memory pressure: Cache eviction requires intelligent algorithms, i.e. LRU

TCP, UDP and QUIC

- TCP
 - Connection-oriented
 - Has been the reliable backbone of internet communication
 - Guaranteed delivery and maintained data order
 - Includes error checking, lost packet retransmission, and congestion control mechanisms
 - Powers the majority of the web
- UDP
 - Connectionless
 - No reliability or ordered packet delivery guarantees
 - Faster than TCP due to less guarantees
 - Suitable for applications that require speed over reliability (i.e. online gaming, video calls etc.)
- QUIC
 - Modern transport layer protocol developed by Google and standardized by the IETF in 2020
 - Intends to combine the best aspects of TCP and UDP

Delta Synchronization

- Traditionally, when files are synchronized between hosts, the entire file is transferred
- Delta synchronization is a technique that intends to instead transfer only the part of the file that has changed
- Can lead to reduced network and I/O overhead
- The probably most popular tool for file synchronization like this is rsync
- When the delta-transfer algorithm for rsync is active, it computes the difference between the local and the remote file, and then synchronizes the changes
- The delta sync algorithm first does file block division
- The file on the destination is divided into fixed-size blocks
- For each block in the destination, a weak and fast checksum is calculated
- The checksums are sent over to the source

File Systems In Userspace (FUSE)

- Software interface that allows writing custom file systems in userspace
- Developers can create file systems without having to engage in low-level kernel development
- Available on multiple platforms, mostly Linux but also macOS and FreeBSD
- In order to implement file systems in user space, we can use the FUSE API
- Here, a user space program registers itself with the FUSE kernel module
- This program provides callbacks for the file system operations, e.g. for open, read, write etc.
- When the user performs a file system operation on a mounted FUSE file system, the kernel module will send a request for the operation to the user space program, which can then reply with a response,

Network Block Device (NBD)

- NBD uses a protocol to communicate between a server (provided by user space) and a client (provided by the NBD kernel module)
- The protocol can run over WAN, but is really mostly meant for LAN or localhost usage
- It has two phases: Handshake and transmission
- There are two actors in the protocol: One or multiple clients, the server and the virtual concept of an export
- When the client connects to the server, the server sends a greeting message with the server's flags
- The client responds with its own flags and an export name (a single NBD server can expose multiple devices) to use
- The server sends the export's size and other metadata, after which the client acknowledges the received data and the handshake is complete
- After the handshake, the client and server start exchanging

- While these systems already allow for some optimizations over simply using the NBD protocol over WAN, they still mean that chunks will only be fetched as they are being needed, which means that there still is a guaranteed minimum downtime
- In order to improve on this, a more advanced API (the managed mount API) was created
- A field that tries to optimize for this use case is live migration of VMs
- Live migration refers to moving a virtual machine, its state and connected devices from one host to another with as little downtime as possible
- There are two types of such migration algorithms; pre-copy and post-copy migration
- Pre-copy migration works by copying data from the source to the destination as the VM continues to run (or in the case of a generic migration, app/other state continues being written to)

- An alternative to pre-copy migration is post-copy migration
- In this approach, the VM is immediately suspended on the source, moved to the destination with only a minimal set of chunks
- After the VM has been moved to the destination, it is resumed
- If the VM tries to access a chunk on the destination, a page fault is raised, and the missing page is fetched from the source, and the VM continues to execute
- The benefit of post-copy migration is that it does not require re-transmitting dirty chunks to the destination before the maximum tolerable downtime is reached
- The big drawback of post-copy migration is that it can result in longer migration times, because the chunks need to be fetched from the network on-demand, which is very latency/RTT-sensitive

- “Reducing Virtual Machine Live Migration Overhead via Workload Analysis” provides an interesting analysis of options on how this decision of when to migrate can be made
- While being designed mostly for use with virtual machines, it could serve as a basis for other applications or migration scenarios, too
- The proposed method identifies workload cycles of VMs and uses this information to postpone the migration if doing so is beneficial
- This works by analyzing cyclic patterns that can unnecessarily delay a VM's migration, and identifies optimal cycles to migrate VMs in from this information
- For the VM use case, such cycles could for example be the GC of a large application triggering a lot of changes to the VMs memory etc.
- If a migration is proposed, the system checks for whether it is currently in a beneficial cycle to migrate in which case it lets the migration proceed; otherwise, it postpones it until the next cycle

Streams and Pipelines

- Fundamental concepts in computer science
- Sequentially process elements
- Allow for the efficient processing of large amounts of data, without having to load everything into memory
- Form the backbone of efficient, modular data processing
- Streams
 - Represent a continuous sequence of data
 - Can be a source or destination of data (i.e. files, network connections, stdin/stdout etc.)
 - Allow processing of data as it becomes available
 - Minimized memory consumption
 - Especially well suited for long-running processes (where data gets streamed in for an extended time)
- Pipelines
 - Series of data processing stages: Output of one stage serves as input to the next

- Open-Source and high-performance RPC framework
- Developed by Google in 2015
- Features
 - Uses HTTP/2 as the transport protocol to benefit from header compression and request multiplexing
 - Uses protobuf as the IDL and wire format, a high-performance, polyglot mechanism for data serialization (instead of the slower and more verbose JSON of REST APIs)
 - Supports unary RPCs, server-streaming RPCs, client-streaming RPCs and bidirectional RPCs
 - Has pluggable support for load balancing, tracing, health checking and authentication
- Supports many languages (Go, Rust, JS etc.)
- Developed by the CNCF

- In-memory data structure store
- Used as a database, cache and/or message broker
- Created by S. Sanfilippo in 2009
- Different from other NoSQL databases by supporting various data structures like lists, sets, hashes or bitmaps
- Uses in-memory data storage for maximum speed and efficiency
- Allows for low-latency reads/writes
- While not intended for persistence, it is possible to store data on disk
- Has a non-blocking I/O model and offers near real-time data processing capabilities
- Includes a pub-sub system to be able to function as a message broker

- S3
 - Object storage service for data-intensive workloads
 - Offered by AWS
 - Can be globally distributed to allow for fast access times from anywhere on the globe
 - Range of storage classes with different requirements
 - Includes authentication and authorization
 - Exposes HTTP API for accessing the stored folders and files
- Minio
 - Open-source storage server compatible with S3
 - Lightweight and simple, written in Go
 - Can be hosted on-prem and is open source
 - Allows horizontal scalability to store large amounts of data across nodes

Cassandra and ScyllaDB

- Popular wide-column NoSQL databases
- Combines Amazon's Dynamo model and Google's Bigtable model to create a highly available database
- Apache Cassandra
 - Highly scalable, eventually consistent
 - Can handle large amounts of data across many servers with no single point of failure
 - Consistency can be tuned according to needs (eventual to strong)
 - Doesn't use master nodes due to its use of a P2P protocol and distributed hash ring design
 - Does have high latency under heavy load and requires fairly complex configuration
- ScyllaDB
 - Launched in 2015
 - Written in C++ and has a shared-nothing architecture, unlike Cassandra which is written in Java
 - Compatible with Cassandra's API and data model

Planning

Pull-Based Synchronization With userfaultfd

- In our case, we want to handle page faults in userspace
- In our case however, we can use a recent system called userfaultfd to do this in a more elegant way (available since kernel 4.11)
- userfaultfd allows handling these page faults in userspace
- The region that should be handled can be allocated with e.g. mmap
- Once we have the file descriptor for the userfaultfd API, we need to transfer this file descriptor to a process that should respond with the chunks of memory to be put into the faulting address
- Once we have received the socket we need to register the handler for the API to use
- If the handler receives an address that has faulted, it responds with the UFFDIO_COPY ioctl and a pointer to the chunk of memory that should be used on the file descriptor (code snippet from <https://github.com/loopholelabs/userfaultfd-go/blob/master/pkg/mapper/handler.go>)

Push-Based Synchronization With mmap and Hashing

- File-based synchronization
 - This approach tries to improve on userfaultfd by switching to push-based synchronization method
 - Instead of reacting to page faults, this uses a file to track changes to a memory region
 - By synchronizing the file representing the memory region between two systems, we can effectively synchronize the memory region itself
 - In Linux, swap space allows Linux to move chunks of memory to disk or other swap partition if the fast speed of RAM is not needed (“paging out”)
 - Similarly to this, Linux can also load missing chunks from a disk
 - This works similarly to how userfaultfd handled page faults, except this time it doesn’t need to go through user space, which can make it much faster
- Detecting file changes
 - `inotify` allows applications to register handlers on a file’s events, e.g. `WRITE` or `SYNC`. This allows for efficient file synchronization, and

Push-Pull Synchronization with FUSE

- File system-based synchronization in the kernel
 - Since the push method requires polling and is very CPU and I/O intensive, and `userfaultfd-go` has only low throughput, a better solution is needed
 - What if we could still get the events for the writes and reads without having to use `userfaultfd-go` or hashing?
 - We can create a custom file system in Linux and load it as a kernel module
 - This file system could then intercept reads/writes to/from the `mmaped` region, making it possible to respond to them with a custom backend
 - But such a system would need to run in the kernel directly, which leads to a lot of potential drawbacks
 - While it is possible to write kernel modules with Rust instead of C these days, a lot of problems remain
 - Kernel modules aren't portable; they are built for a specific kernel, which makes them hard to distribute to users

Pull-Based Synchronization With NBD

- Block Device-based synchronization
 - As hinted at before, a better API would be able to catch reads/writes to a single mmaped file instead of having to implement a complete file system
 - It does however not have to be an actual file, a block device also works
 - In Linux, block devices are (typically storage) devices that support reading/writing fixed chunks (blocks) of data
 - We can mmap a block device in the same way that we can mmap a file
 - Similarly to how a file system can be implemented in a kernel module, a block device is typically implemented as a kernel module/in kernel space
 - However, the same security, portability and developer experience issues as with the former also apply here
 - Instead of implementing a FUSE to solve this, we can create a NBD (network block device) server that can be used by the kernel NBD module similarly to how the process that connected to the FUSE

Push-Pull Synchronization with Mounts

- Limitations of the NBD protocol in WAN
 - Usually, the NBD server and client don't run on the same system
 - NBD was originally designed to be used as a LAN protocol to access a remote hard disk
 - As mentioned before, NBD can run over WAN, but is not designed for this
 - The biggest problem with running NBD over a public network, even if TLS is enabled is latency
 - Individual chunks would only be fetched to the local system as they are being accessed, adding a guaranteed minimum latency of at least the RTT
 - Instead of directly connecting a client to a remote server, we add a layer of indirection, called a Mount that consists of both a client and a server, both running on the local system
- Combining the NBD server and client to a reusable unit
 - We then connect the server to the backend with an API that is better suited for WAN usage

Pull-Based Synchronization with Migrations

- Optimization mounts for migration scenarios
 - We have now implemented a managed mounts API
 - This API allows for efficient access to a remote resource through memory
 - It is however not well suited for a migration scenario
 - For migrations, more optimization is needed to minimize the maximum acceptable downtime
 - For the migration, the process is split into two distinct phases
 - The same preemptive background pulls and parallelized device/syncer startup can be used, but the push process is dropped
 - The two phases allow pulling the majority of the data first, and only finalize the move later with the remaining data
 - This is inspired by the pre-copy approach to VM live migration, but also allows for some of the benefits of the post-copy approach as we'll see later
 - Why is this useful? A constraint for the mount-based API that we haven't mentioned before is that it doesn't allow safe concurrent

Implementation

Userfaults in Go with userfaultfd

- API design for userfault-go
 - Implementing this in Go was quite tricky, and it involves using unsafe
 - We can use the `syscall` and `unix` packages to interact with `ioctl` etc.
 - We can use the `ioctl` `syscall` to get a file descriptor to the `userfaultfd` API, and then register the API to handle any faults on the region (code snippet from <https://github.com/loopholelabs/userfaultfd-go/blob/master/pkg/mapper/register.go#L15>)
 - Passing file descriptors between processes is possible by using a UNIX socket (code snippet from <https://github.com/loopholelabs/userfaultfd-go/blob/master/pkg/transfer/unix.go>)
- Implementing `userfaultfd` backends
 - A big benefit of using `userfaultfd` and the pull method is that we are able to simplify the backend of the entire system down to a `io.ReaderAt` (code snippet from <https://pkg.go.dev/io#ReaderAt>)
 - That means we can use almost any `io.ReaderAt` as a backend for a `userfaultfd-go` registered object

File-Based Synchronization

- File-based synchronization
 - We can do this by using `mmap`, which allows us to map a file into memory
 - By default, `mmap` doesn't write changes from a file back into memory, no matter if the file descriptor passed to it would allow it to or not
 - We can however add the `MAP_SHARED` flag; this tells the kernel to write back changes to the memory region to the corresponding regions of the backing file
 - Linux caches reads to such a backing file, so only the first page fault would be answered by fetching from disk, just like with `userfaultfd`
 - The same applies to writes; similar to how files need to be synced in order for them to be written to disks, `mmaped` regions need to be `msynced` in order to flush changes to the backing file
 - In order to synchronize changes to the region between hosts by syncing the underlying file, we need to have the changes actually be represented in the file, which is why `msync` is critical
 - For files, you can use `O_DIRECT` to skip this kernel caching if your

FUSE Implementation in Go

- It is possible to use even very complex and at first view non-compatible backends as a FUSE file system's backend
- By using a file system abstraction API like `afero.Fs`, we can separate the FUSE implementation from the actual file system structure, making it unit testable and making it possible to add caching in user space (code snippet from <https://github.com/poignantfx/stfs/blob/main/pkg/fs/file.go>)
- It is possible to map any `afero.Fs` to a FUSE backend, so it would be possible to switch between different file system backends without having to write FUSE-specific (code snippet from <https://github.com/JakWai01/sile-fsystem/blob/main/pkg/filesystem/fs.go>)
- For example, STFS used a tape drive as the backend, which is not random access, but instead append-only and linear (<https://github.com/poignantfx/stfs/blob/main/pkg/operations/update.go>)³³

- Implementing go-nbd
 - Due to the lack of pre-existing libraries, a new pure Go NBD library was implemented
 - This library does not rely on CGo/a pre-existing C library, meaning that a lot of context switching can be skipped
 - The backend interface for go-nbd is very simple and only requires four methods: ReadAt, WriteAt, Size and Sync
 - A good example backend that maps well to a block device is the file backend (code snippet from <https://github.com/pojntfx/go-nbd/blob/main/pkg/backend/file.go>)
 - The key difference here to the way backends were designed in userfaultfd-go is that they can also handle writes
 - go-nbd exposes a Handle function to support multiple users without depending on a specific transport layer (code snippet from <https://github.com/pojntfx/go-nbd/blob/main/pkg/server/nbd.go>)
 - This means that systems that are peer-to-peer (e.g. WebRTC), and thus don't provide a TCP-style accept syscall can still be used easily

Chunking, Push/Pull Mechanisms and Lifecycle for Mounts

- The `ReadWriteAt` pipeline
 - In order to implement the chunking system, we can use an abstraction layer that allows us to create a pipeline of readers/writers - the `ReadWriteAt`, combining an `io.ReaderAt` and a `io.WriterAt`
 - This way, we can forward the `Size` and `Sync` syscalls directly to the underlying backend, but wrap a backend's `ReadAt` and `WriteAt` methods in a pipeline of other `ReadWriteAt`s
 - One such `ReadWriteAt` is the `ArbitraryReadWriteAt` (code snippet from https://github.com/pojntfx/r3map/blob/main/pkg/chunks/arbitrary_rwat.go)
 - It allows breaking down a larger data stream into smaller chunks
 - In `ReadAt`, it calculates the index of the chunk that the offset falls into and the position within the offsets
 - It then reads the entire chunk from the backend into a buffer, copies the necessary portion of the buffer into the input slice, and repeats the process until all requested data is read
 - Similarly for the writer, it calculates the chunk's index and offset
 - If an entire chunk is being written to, it bypasses the chunking system,

Live Migration for Mounts

- Optimization mounts for migration scenarios
 - The flexible architecture of the ReadWriterAt components allow the reuse of lots of code for both use cases
- The migration protocol
 - To fix this, the migration API defines two new actors: The seeder and the leecher
 - The seeder represents a resource that can be migrated from/a host that exposes a migratable resource
 - The leecher represents a client that wants to migrate a resource to itself
 - Initially, the protocol starts by running an application with the application's state on the seeder's mount
 - When a leecher connects to the seeder, the seeder starts tracking any writes to it's mount
 - The leecher starts pulling chunks from the seeder to it's local backend
 - Once the leecher has received a satisfactory level of locally available chunks, it as the seeder to finalize, which then causes the seeder to

Pluggable Encryption and Authentication

- Compared to existing remote mount and migration solutions, r3map is a bit special
- As mentioned before, most systems are designed for scenarios where such resources are accessible in a high-bandwidth, low-latency LAN
- This means that some assumptions concerning security, authentication, authorization and scalability were made that can not be made here
- For example encryption; while for a LAN deployment scenario it is probably assumed that there are no bad actors in the subnet, the same can not be said for WAN
- While depending on e.g. TLS etc. for the migration could have been an option, r3map should still be useful for LAN migration use cases, too, which is why it was made to be completely transport-agnostic
- This makes adding encryption very simple
- E.g. for LAN, the same assumptions that are being made in existing

Optimizing Backends For High RTT

- In WAN, where latency is high, the ability to fetch chunks concurrently is very important
- Without concurrent background pulls, latency adds up very quickly as every memory request would have at least the RTT as latency
- The first prerequisite for supporting this is that the remote backend has to be able to read from multiple regions without locking the backend globally
- For the file backend for example, this is not the case, as the lock needs to be acquired for the entire file before an offset can be accessed (code snippet from <https://github.com/poichtfx/go-nbd/blob/main/pkg/backend/file.go#L17-L25>)
- For high-latency scenarios, this can quickly become a bottleneck
- While there are many ways to solve this, one is to use the directory backend
- Instead of using just one backing file, the directory backend is a

Optimizing The Transport Protocol For Throughput

- Despite these benefits, gRPC is not perfect however
- Protobuf specifically, while being faster than JSON, is not the fastest serialization framework that could be used
- This is especially true for large chunks of data, and becomes a real bottleneck if the connection between source and destination would allow for a high throughput
- This is where fRPC, a RPC library that is easy to replace gRPC with, becomes useful
- Because throughput and latency determine the maximum acceptable downtime of a migration/the initial latency for mounts, choosing the right RPC protocol is an important decision
- fRPC also uses the same proto3 DSL, which makes it an easy drop-in replacement, and it also supports multiplexing and connection polling
- Because of these similarities, the

Using Remote Stores as Backends

- Using key-value stores as ephemeral mounts
 - RPC backends provide a way to access a remote backend
 - This is useful, esp. if the remote resource should be protected in some way or if it requires some kind of authorization
 - Depending on the use case however, esp. for the mount API, having access to a remote backend without this level of indirection can be useful
 - Fundamentally, a mount maps fairly well to a remote random-access storage device
 - Many existing protocols and systems provide a way to access essentially this concept over a network
 - One of these is Redis, an in-memory key-value store with network access
 - Chunk offsets can be mapped to keys, and bytes are a valid key type, so the chunk itself can be stored directly in the KV store (code snippet from

<https://github.com/pojntfx/r3map/blob/main/pkg/backend/redis.go#L36>⁴⁰

Bi-Directional Protocols With Dudirekta

- Another aspect that plays an important role in performance for real-life deployments is the choice of RPC framework and transport protocol
- As mentioned before, both the mount and the migration APIs are transport-independent
- A simple RPC framework to use is dudirekta
- Dudirekta is reflection-based, which makes it very simple to use to iterate on the protocol quickly
- To use it, a simple wrapper struct with the needed RPC methods is created (code snippet from <https://github.com/poijntfx/r3map/blob/main/pkg/services/backend.go#L41-L61>)
- This wrapper struct simply calls the backend (or seeder etc.) functions
- The wrapper struct is then passed as the local function struct into a

Connection Pooling For High RTT Scenarios

- This does however come at the cost of not being able to do connection pooling, since each client dialing the server would mean that the server could not reference the multiple client connections as one composite client without changes to the protocol
- While implementing such a pooling mechanism in the future could be interesting, it turned out to not be necessary thanks to the pull-based pre-copy solution described earlier
- Instead, only calling RPCs exposed on the server from the client is the only requirement for an RPC framework, and other, more optimized RPC frameworks can already offer this
- Dudirekta uses reflection to make the RPCs essentially almost transparent to use
- By switching to a well-defined protocol with a DSL instead, we can gain further benefits from not having to use reflection and generating code instead

Results

- First chunk latency of baseline vs. userfault vs. direct vs. managed
- Development of first chunk latency as RTT increases
- Development of first chunk latency as number of workers increases
- MB of pre-emptive pulls in managed API as RTT increases
- Development of MB of pre-emptive pulls as RTT increases
- Development of MB of pre-emptive pulls as number of workers increases
- Development of MB of pre-emptive pulls as chunk size increases

- Throughput of dudirekta, gRPC and fRPC as RTT increases
- Throughput of dudirekta, gRPC and fRPC as number of workers increases
- Throughput of dudirekta, gRPC and fRPC as chunk size increases
- Throughput of baseline vs. userfault vs. direct vs. migration (we can use the onChunkLocal value)
- Development of throughput as RTT increases
- Development of throughput as number of workers increases (migration only)
- Development of throughput as chunk size increases

- Effect of correct (return n) vs wrong (return $-n$) pull heuristic on first chunk latency for managed API
- Effect of correct (return n) vs wrong (return $-n$) pull heuristic on throughput for managed API

- Open time for userfault vs. r3map (doesn't matter if direct/mount/migration)
- Minimum downtime of direct/mount (one step - Open = downtime) vs. migration (two steps - Finalize = downtime)

Discussion

- As we can see, using userfaultfd we are able to map almost any object into memory
- This approach is very clean and has comparatively little overhead, but also has significant architecture-related problems that limit its uses
- The first big problem is only being able to catch page faults - that means we can only ever respond the first time a chunk of memory gets accessed, all future requests will return the memory directly from RAM on the destination host
- This prevents us from using this approach for remote resources that update over
- Also prevents us from using it for things that might have concurrent writers/shared resources, since there would be no way of updating the conflicting section
- Essentially makes this system only usable for a read-only “mount” of

File-Based Synchronization

- Similarly to userfaultfd, this system also has limitations
- While userfaultfd was only able to catch reads, this system is only able to catch writes to the file
- Essentially this system is write-only, and it is very inefficient to add hosts to the network later on
- As a result, if there are many possible destinations to migrate state too, a star-based architecture with a central forwarding hub can be used
- The static topology of this approach can be used to only ever require hashing on one of the destinations and the source instead of all of them
- This way, we only need to push the changes to one component (the hub), instead of having to push them to each destination on their own
- The hub simply forwards the messages to all the other destinations

- FUSE does however also have downsides
- It operates in user space, which means that it needs to do context switching
- Some advanced features aren't available for a FUSE
- The overhead of FUSE (and implementing a completely custom file system) for synchronizing memory is still significant
- If possible, the optimal solution would be to not expose a full file system to track changes, but rather a single file
- As a result of this, the significant implementation overhead of such a file system led to it not being chosen

- Limitations of NBD and ublk as an alternative
 - NBD is a battle-tested solution for this with fairly good performance, but in the future a more lean implementation called ublk could also be used
 - ublk uses `io_uring`, which means that it could potentially allow for much faster concurrent access
 - It is similar to NBD; it also uses a user space server to provide the block device backend, and a kernel ublk driver that creates `/dev/ublk*` devices
 - Unlike as it is the case for the NBD kernel module, which uses a rather slow UNIX or TCP socket to communicate, ublk is able to use `io_uring` pass-through commands
 - The `io_uring` architecture promises lower latency and better throughput
 - Because it is however still experimental and docs are lacking, NBD was chosen
- BUSE and CUSE as alternatives to NBD

Mounts and Live Migration

- Overhead of using managed mounts
 - Another interesting aspect of optimization to look at is the overhead of managed mounts
 - It is possible for managed mounts (and migrations) to deliver significantly lower performance compared to direct mounts
 - This is because using managed mounts come with the cost of potential duplicate I/O operations
 - For example, if memory is being accessed linearly from the first to the last offset immediately after it being mounted, then using the background pulls will have no effect other than causing a write operation (to the local caching backend) compared to just directly reading from the remote backend
 - This however is only the case in scenarios with a very low latency between the local and remote backends
 - If latency becomes higher, then the ability to pull the chunks in the background and in parallel with the puller will offset the cost of duplicate I/O

Remote Swap With ram-dl

- ram-dl is a fun experiment
- Tech demo for r3map
- Uses the fRPC backend to expand local system memory
- Can allow mounting a remote system's RAM locally
- Can be used to inspect a remote system's memory contents
- Is based on the direct mount API
- Uses mkswap, swapon and swapoff (code snippet from <https://github.com/poignantfx/ram-dl/blob/main/cmd/ram-dl/main.go#L170-L190>)
- Enables paging out to the block device provided by the direct mount API
- ram-ul “uploads” RAM by exposing a memory, file or directory-backed file over fRPC
- ram-dl then does all of the above
- Not really intended for real-world usecases, but does show that this

Mapping Tape Into Memory With tapisk

- tapisk is an interesting usecase because of how close it is to STFS, which provided the inspiration for the FUSE-based approach
- Very high read/write backend latency (multiple seconds, up to 90s, due to seeking)
- Linear access, no random reads
- Can benefit a lot from asynchronous writes provided by managed mounts
- Fast storage acts as a caching layer
- Backend is linear, so only one read/write possible at a time
- With local backend, writes are de-duplicated automatically and both can be asynchronous/concurrent
- Writes go to fast (“local”) backend first, syncer then handles in both directions
- Chunking works on tape drive records and blocks
- Only one concurrent reader/writer makes sense

Improving File Storage Solutions

- Another potential usecase is using r3map to create a mountable remote filesystem with unique properties
- Currently there are two choices on how these can be implemented
- Google Drive, Nextcloud etc. listen to file changes on a folder and synchronize files when they change
- The big drawback is that everything needs to be stored locally
- If a lot of data is stored (e.g. terabytes), the locally available files would need to be manually selected
- There is no way to dynamically download files this way as they are required
- It is however very efficient, since the filesystem is completely transparent to the user (writes are being synced back asynchronously)
- It also supports an offline usecase easily
- The other option is to use a FUSE, e.g. s3-fuse

- Another usecase is accessing a remote database locally
- While using a database backend is one option of storing chunks, an actual database can also be stored in a mount as well
- Particularly interesting for in-memory or on-disk databases like SQLite
- Instead of having to download the entire SQLite database before using it, it can simply be mounted, and accessed as it is being used
- This allows very efficient network access, as only very rarely the entire database is needed
- Since reads are cached with the managed mount API, only the first read should potentially have a performance impact
- Similarly so writes to the database will be more or less the same throughput as to the local disk, since changes are written back asynchronously
- If the full database should eventually be accessible locally, the

Universal App State Synchronization and Migration

- Synchronization of app state is hard
- Even for hand-off scenarios a custom protocol is built most of the times
- It is possible to use a database sometimes (e.g. Firebase) to synchronize things
- But that can't sync all data structures and requires using specific APIs to sync things
- What if you could mount and/or migrate any resource?
- Usually these structures are marshalled, sent over a network, received on a second system, unmarshalled, and then are done being synced
- Requires a complex sync protocol, and when to sync, when to pull etc. is inefficient and usually happens over a third party (e.g. a database)
- Data structures can almost always be represented by an array of bytes
- If the data structures are allocated from an array of bytes from the block

Summary

- Looking back at all synchronization options and comparing ease of implementation, CPU load and network traffic between them
- Summary of the different approaches, and how the new solutions might make it possible to use memory as the universal access format

Conclusion

- Answer to the research question (Could memory be the universal way to access and migrate state?)
- What would be possible if memory became the universal way to access state?
- Further research recommendations (e.g. ublk)

References

- Abstract
- Tables
 - Content
 - Graphics
 - Abbreviations
- Introduction: Research question/goal
- Technology: Current state of research and introductions to related topics (like a reference book)
- Planning: Ideas on how to solve the research question/goal
- Implementation: Description of how the ideas were implemented
- Results: Which results each of the methods yielded (i.e. benchmarks)
- Discussion: How the individual results can be interpreted and what use cases there are for the methods
- Summary: How the discussion can be interpreted as a whole
- Conclusion: Answer to the research question and future outlook, future research questions

- FluidMem: Full, Flexible, and Fast Memory Disaggregation for the Cloud (userfaultfd) [1]
- Memory Disaggregation: Advances and Open Challenges (general memory hierarchy etc.) [2]
- User Level Page Faults (sigaction & signal-based page fault handlers) [3]
- Towards Web-based Delta Synchronization for Cloud Storage Services (rsync) [4]
- To FUSE or Not to FUSE: Performance of User-Space File Systems (implementing a FUSE) [5]
- The NBD protocol (NBD protocol documentation from the sources; info on implementing a NBD server and client) [6]
- Remote regions: a simple abstraction for remote memory (Kernel filesystem for remote memory access) [7]
- Reducing Virtual Machine Live Migration Overhead via Workload

[1]

B. Caldwell *et al.*, “FluidMem: Full, flexible, and fast memory disaggregation for the cloud,” in *2020 IEEE 40th international conference on distributed computing systems (ICDCS)*, 2020, pp. 665–677. doi: 10.1109/ICDCS47774.2020.00090.

[2]

H. A. Maruf and M. Chowdhury, “Memory disaggregation: Advances and open challenges.” 2023.Available: <https://arxiv.org/abs/2305.03943>

[3]

Q. Li, “User level page faults,” Master’s thesis, Carnegie Mellon University, School of Computer Science, Pittsburgh, PA 15213, 2020.

[4]

H. Xiao *et al.*, “Towards web-based delta synchronization for cloud storage services,” in *16th USENIX conference on file and storage technologies (FAST-16)*, Feb. 2015, pp. 147–162. Available at: <https://www.usenix.org/conference/fast16>.