
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

Contents

1	Unsorted Research Questions	2
2	Structure	2

1 Unsorted Research Questions

2 Structure

- Introduction
 - Memory management in Linux
 - Memory as the universal storage API
 - What would be possible if memory would be the universal way to access resources?
 - Why efficient memory synchronization is the missing key component
 - High-level use cases for memory synchronization in the industry today
- Pull-Based Memory Synchronization with `userfaultfd`
 - Plain language description of `userfaultfd` (what are page faults)
 - Exploring an alternative method by handling page faults using signals
 - Handlers and registration
 - History of `userfaultfd`
 - Allocating the shared region
 - Maximum shared region size is limited by available physical memory
 - Transferring handler sockets between processes
 - Implementing `userfaultfd` bindings in Go
 - Example usages of `userfaultfd` in Go (byte slice, file, S3 object)
 - Implications of not being able to catch writes to the region (its read-only)
 - Design of a `userfaultfd` backend (`io.ReaderAt`)
 - Limitations of only being able to catch the first page fault (no way of updating the region)
 - Implications of not being able to pull chunks before they are being accessed
 - Limitations of only being able to pull chunks synchronously
 - Benefits of minimal registration and latency overhead
 - Benchmark: Why is it limited to ~50MB/s of throughput?
 - Benchmark: Sensitivity of `userfaultfd` to network latency and throughput
- Push-Based Memory Synchronization with `mmap` and Hashing
 - Plain language description of this approach (mapping a file into memory, then syncing the file)
 - Paging and swap in Linux
 - Introduction to `mmap` to map a file into a memory region
 - `MAP_SHARED` for writing changes back from the memory region to a file
 - Caching with `mmap`, why `O_DIRECT` doesn't work and what the role of `msync` is

- Detecting writes to a file with `inotify` and why this does not work for `mmap`
- Hashing (chunks of) the backing file in order to detect writes on a periodic basis
- Benchmark: Performance of different Go hashing algorithms for detecting writes
- Effects on maximum acceptable downtime due to CPU-bound limitations in having to calculate hashes and polling
- Introduction to delta synchronization protocols like `rsync`
- Implementation of a custom delta synchronization protocol with chunk support
- Multiplexing different synchronization streams in the protocol
- Benchmark: Throughput of this custom synchronization protocol vs. `rsync` (which hashes entire files)
- Using a central forwarding hub for star-based architectures with multiple destinations
- Limitations of only being able to catch writes, not reads to the region (its write-only, can't add hosts later on)
- Push-Pull Memory Synchronization with a FUSE
 - Plain language description of this approach (mapping a file into memory, catching read-/writes to/from the file with a custom filesystem)
 - Methods for creating a new, custom file system (FUSE vs. in the kernel)
 - STFS shows that it is possible to create file systems backed by complex data stores, e.g. a tape/non-random access stores
 - Is not the best option for implementing this due to significant overhead and us only needing a single file to map
- Pull-Based Memory Synchronization with NBD
 - Plain language description of this approach (mapping a block device into memory, catching reads/writes with a NBD server)
 - Why NBD is the better choice compared to FUSE (much less complex interface)
 - Overview of the NBD protocol
 - Phases, actors and messages in the NBD protocol (negotiation, transmission)
 - Minimal viable NBD protocol needs, and why only this minimal set is implemented
 - Listing exports
 - Limitations of NBD, esp. the kernel client and message sizes
 - Reduced flexibility of NBD compared to FUSE (can still be used for e.g. linear media, but will offer fewer interfaces for optimization)
 - Server backend interface design
 - File backend example
 - How the servers handling multiple users/connections
 - Server handling in the NBD protocol implementation

- Using the kernel NBD client without CGO/with ioctls
- Finding an unused NBD device using `sysfs`
- Benchmark: `nbd` kernel module quirks and how to detect whether a NBD device is open (polling `sysfs` vs. `udev`)
- Caching mechanisms and limitations (aligned reads) when opening the block device (`O_DIRECT`)
- Future outlook: Using `ublk` instead of NBD, allowing for potentially much faster concurrent access thanks to `io_uring`
- Why using BUSE to implement a NBD server would be possible but unrealistic (library & docs situation, CGo)
- `go-buse` as a preview of how such a CGo implementation could still work
- Alternatively implementing a new file system entirely in the kernel, only exposing a single file/block device to `mmap` and optimizing the user space protocol for this
- Push-Pull Memory Synchronization with Mounts
 - Plain language description of this approach (like NBD, but starting the client and server locally, then connecting the *server's* backend to a backend)
 - Benefits: Can use a secure wire protocol and more complex/abstract backends
 - Mounting the block device as a path vs. file vs. slice: Benefits of `mmap` (concurrent reads/writes)
 - Alternative approach: Formatting the block device as e.g. EXT4, then mounting the filesystem, and `mmap`ing a file on the file system (allows syncing multiple regions with a single file system, but has FS overhead)
 - Mount protocol actors, phases and state machine
 - Chunking system for non-aligned reads/writes (arbitrary rwat and chunked rwat)
 - Benchmark: Local vs. remote chunking
 - Optimizing the mount process with the Managed Mount interface
 - Asynchronous background push system interface and how edge cases (like accessing dirty chunks as they are being pulled or being written to) are handled
 - Preemptive background pulls interface
 - Syncer interface
 - Benchmark: Parallelizing startups and pulling n MBs as the device starts
 - Using a pull heuristic function to optimize which chunks should be scheduled to be pulled first
 - Internal rwat pipeline (create a graphic) for direct mounts vs. managed mounts
 - Unit testing the rwats
 - Comparing this mount API to other existing remote memory access APIs, e.g. "Remote Regions"

- Complexities when `mmap`ing a region in Go as the GC halts the entire world to collect garbage, but that also stops the NBD server in the same process that tries to satisfy the region being scanned
- Potentially using Rust for this component to cut down on memory complexities and GC latencies
- Pull-Based Memory Synchronization with Migrations
 - Plain language description of this approach (like NBD, but two phases to start the device and pull, then only flush the latest changes to minimize downtime)
 - Inspired by live VM migration, where changes are continuously being pulled to the destination node until a % has been reached, after which the VM is migrated
 - Migration protocol actors (seeders, leechers etc.), phases and state machine
 - How the migration API is completely independent of a transport layer
 - Switching from the seeder to the leecher state
 - Using preemptive pulls and pull heuristics to optimize just like for the mounts
 - Lifecycle of the migration API and why lockable rwats are required
 - How a successful migration causes the `Seeder` to exit
 - The role of maximum acceptable downtime
 - The role of `Track()`, concurrent access and consistency guarantees vs. mounts (where the source must not change)
 - When to best `Finalize()` a migration and how analyzing app usage patterns could help
 - Benchmark: Maximum acceptable downtime for a migration scenario with the Managed Mount API vs the Migration API
- Optimizing Mounts and Migrations
 - Encryption of memory regions and the wire protocol
 - Authentication of the protocol
 - DoS vulnerabilities in the NBD protocol (large message sizes; not meant for the public internet) and why the indirection of client & server on each node is needed
 - Mitigating DoS vulnerabilities in the `ReadAt/WriteAt` RPCs with `maxChunkSize` and/or client-side chunking
 - Critical `Finalizing` state in the migration API and how it could be remedied
 - How network outages are handled in the mount and migration API
 - Analyzing the file and memory backend implementations
 - Analyzing the directory backend
 - Analyzing the `dudirekta`, `gRPC` and `fRPC` backends
 - Benchmark: Latency till first `n` chunks *and* throughput for `dudirekta`, `gRPC` and `fRPC` backends (how they are affected by having/not having connection polling and/or concurrent

- RPCs)
 - Benchmark: Effect of tuning the amount of push/pull workers in high-latency scenarios
 - Analyzing the Redis backend
 - Analyzing the S3 backend
 - Analyzing the Cassandra backend
 - Benchmark: Latency and throughput of all benchmarks on localhost and in a realistic latency and throughput scenario
 - Effects of slow disks/memory on local backends, and why direct mounts can outperform managed mounts in tests on localhosts
 - Using P2P vs. star architectures for mounts and migrations
 - Looking back at all options and comparing ease of implementation, CPU load and network traffic between them
- Case Studies
 - [ram-dl](#) as an example of using the direct mount API for
 - [tapisk](#) as an example of using the managed mount API for a file system with very high latencies, linear access and asynchronous to/from fast local memory vs. STFS
 - Migration app state (e.g. TODO list) between two hosts in a universal (underlying memory) manner
 - Mounting remote file systems as managed mounts and combining the benefits of traditional FUSE mounts (e.g. s3-fuse) with Google Drive-style synchronization
 - Using managed mounts for remote SQLite databases without having to download it first
 - Streaming video formats like e.g. MP4 that don't support streaming due to indexes/compression
 - Improving game download speeds by mounting the remote assets with managed mounts, using a pull heuristic that defines typical access patterns (like which levels are accessed first), making any game immediately playable without changes
 - Executing remote binaries or scripts that don't need to be scanned first without having to fully download them
- Conclusion
 - Summary of the different approaches, and how the new solutions might make it possible to use memory as the universal access format
 - Further research recommendations (e.g. [ublk](#))