

# Zero-Downtime Live Migration of Stateful VMs on Kubernetes

Felicitas Pojtinger  
@pojntfx

 Loophole Labs





# Felicitas Pojtinger

bluesky: [@felicitas.pojtinger.com](https://bluesky.fedibird.com/@felicitas.pojtinger)

mastodon: [@pojntfx@mastodon.social](https://mastodon.social/@pojntfx)

github/twitter: [@pojntfx](https://github.com/pojntfx)

linkedin: [in/pojntfx](https://www.linkedin.com/in/pojntfx)

web: [felicitas.pojtinger.com](https://felicitas.pojtinger.com)



# VMs in the Container Age

High overhead

Resource-intensive

Slow to start





QEMU



So why is everything still powered by VMs?  
(Or: Why doesn't everyone just use containers?)



# A "Container of Stateful Applications"?

Small image sizes

Fast to start

Easy to scale

Good networking support

Low overhead

Easy to build &  
distribute



# What Could We Do With Such a New Compute Unit?

# Simplifying deployment and debugging



# Moving between cloud providers without downtime



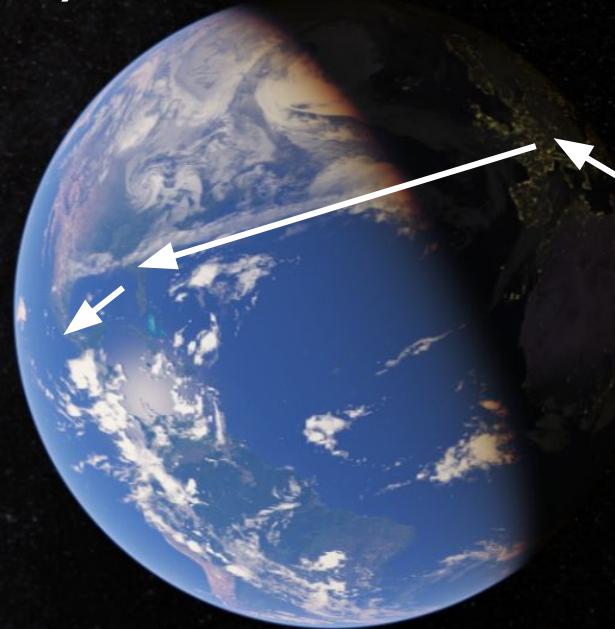
Moving long-running  
processes to more  
powerful servers (while  
continuing computing)



# Migrating full desktop/mobile apps between devices (in milliseconds)

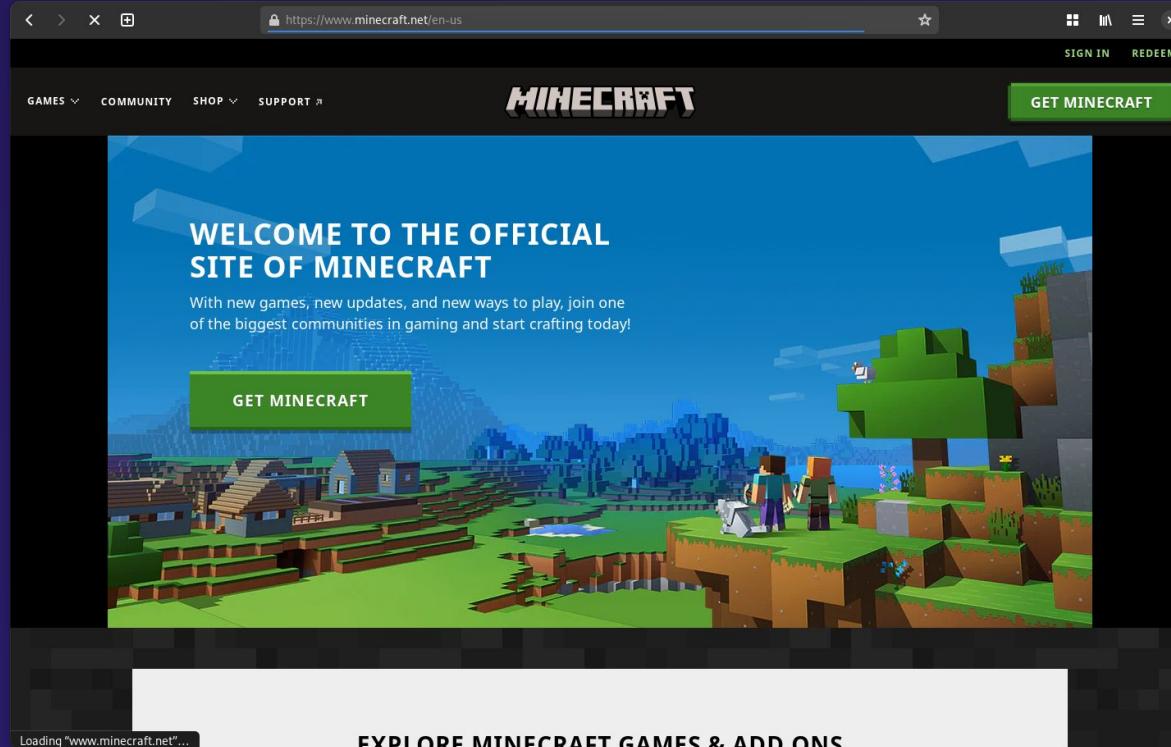


Moving stateful services closer  
to your users as they access  
them (without downtime)



This sounds great  
But it's also impossible ... right?





minecraft.net (**Version 1.12.2**)

# Why Does This Not Exist Already?

# Dated hypervisor ecosystems

VMs can't be streamed in

Outdated packaging  
tooling

IPs change after  
migrations

Not designed for  
modern orchestrators

Live migration is  
complex



# How Did We Implement It Anyway?

# A (Very Brief) Overview of Existing Live Migration Techniques

# Pre-Copy Migration



# Post-Copy Migration



# Remote Memory



Highly  
application-specific

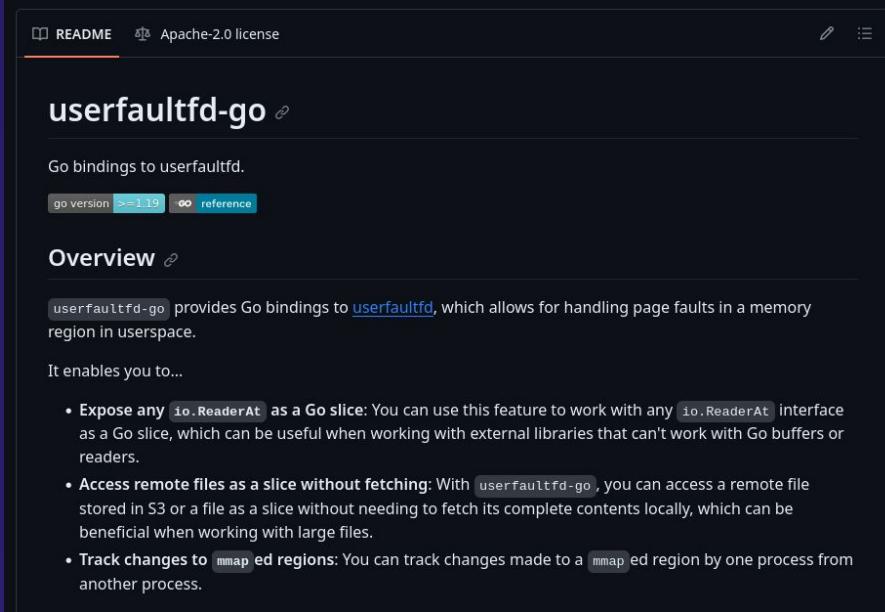
Very  
latency-sensitive

Dependent on specific  
transports

Not easily extensible



# What about userfaultfd?



The screenshot shows a GitHub repository page for "userfaultfd-go". The page has a dark theme. At the top, there are links for "README" and "Apache-2.0 license". Below the title "userfaultfd-go" is a subtitle "Go bindings to userfaultfd.". It includes a note about go version >= 1.19 and a "reference" link. The "Overview" section contains a paragraph about the project's purpose: "userfaultfd-go provides Go bindings to `userfaultfd`, which allows for handling page faults in a memory region in userspace." It also states that it enables you to... followed by a bulleted list of features:

- Expose any `io.ReaderAt` as a Go slice: You can use this feature to work with any `io.ReaderAt` interface as a Go slice, which can be useful when working with external libraries that can't work with Go buffers or readers.
- Access remote files as a slice without fetching: With `userfaultfd-go`, you can access a remote file stored in S3 or a file as a slice without needing to fetch its complete contents locally, which can be beneficial when working with large files.
- Track changes to mmaped regions: You can track changes made to a mmaped region by one process from another process.



# What about userfaultfd?

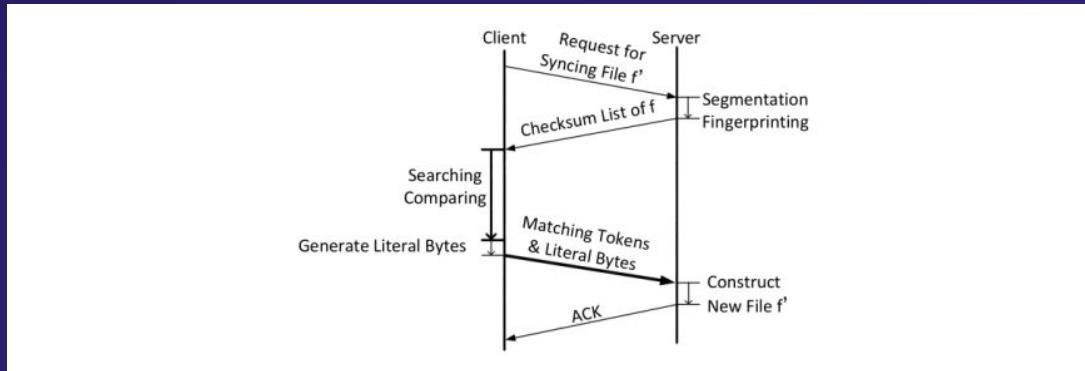
```
1 func (a abcReader) ReadAt(p []byte, off int64) (n int, err error) {
2     n = copy(p, bytes.Repeat([]byte{'A' + byte(off%20)}, len(p)))
3
4     return n, nil
5 }
```

```
1 f, err := os.OpenFile(*file, os.O_RDONLY, os.ModePerm)
2 b, uffd, start, err := mapper.Register(int(s.Size()))
3 mapper.Handle(uffd, start, f)
```

```
1 // ...
2 f, err := mc.GetObject(ctx, *s3BucketName, *s3ObjectName, minio.GetObjectOptions{})
3 b, uffd, start, err := mapper.Register(int(s.Size()))
4 mapper.Handle(uffd, start, f)
```



# File-based synchronization



# File-based synchronization

The screenshot shows the GitHub page for the `HighwayHash` project. At the top, there are links for `README` and `Apache-2.0 license`. Below these are buttons for `reference`, `build`, and `unknown`. The main content starts with a section titled **HighwayHash**. It describes the algorithm as a pseudo-random-function (PRF) developed by Jyrki Alakuijala, Bill Cox and Jan Wassenberg (Google research). It takes a 256 bit key and computes 64, 128 or 256 bit hash values of given messages. It can be used to prevent hash-flooding attacks or authenticate short-lived messages. It is not a general purpose cryptographic hash function (such as Blake2b, SHA-3 or SHA-2) and should not be used if strong collision resistance is required. The repository contains a native Go version and optimized assembly implementations for Intel, ARM and ppc64le architectures.

**High performance**

HighwayHash is an approximately 5x faster SIMD hash function as compared to [SipHash](#) which in itself is a fast and 'cryptographically strong' pseudo-random function designed by Aumasson and Bernstein.

HighwayHash uses a new way of mixing inputs with AVX2 multiply and permute instructions. The multiplications are 32x32 bit giving 64 bits-wide results and are therefore infeasible to reverse. Additionally permuting equalizes the distribution of the resulting bytes. The algorithm outputs digests ranging from 64 bits up to 256 bits at no extra cost.

**Stable**

All three output sizes of HighwayHash have been declared [stable](#) as of January 2018. This means that the hash results for any given input message are guaranteed not to change.

**Installation**

Install: `go get -u github.com/minio/highwayhash`

**Intel Performance**

Below are the single core results on an Intel Core i7 (3.1 GHz) for 256 bit outputs:

Benchmark	Result
Benchmark kSum256_16	284.17 MB/s
Benchmark kSum256_64	1040.63 MB/s
Benchmark kSum256_1K	8653.30 MB/s
Benchmark kSum256_16k	12176.27 MB/s



# CRFS

README Code of conduct BSD-3-Clause license Security ⚙ Ⓜ

## CRFS: Container Registry Filesystem ⚙

Discussion: [golang/go#30829](#)

### Overview ⚙

CRFS is a read-only FUSE filesystem that lets you mount a container image, served directly from a container registry (such as [gcr.io](#)), without pulling it all locally first.

### Background ⚙

Starting a container should be fast. Currently, however, starting a container in many environments requires doing a `pull` operation from a container registry to read the entire container image from the registry and write the entire container image to the local machine's disk. It's pretty silly (and wasteful) that a read operation becomes a write operation. For small containers, this problem is rarely noticed. For larger containers, though, the pull operation quickly becomes the slowest part of launching a container, especially on a cold node. Contrast this with launching a VM on major cloud providers: even with a VM image that's hundreds of gigabytes, the VM boots in seconds. That's because the hypervisors' block devices are reading from the network on demand. The cloud providers all have great internal networks. Why aren't we using those great internal networks to read our container images on demand?



# MBDs



Shivansh Vij 11 months ago

@Felicitas.Pojtinger /@Dan Phillips I figured out magic block devices for the scale builder service, but this should help with dark magyk as well:

Magic Block Devices ([MBD](#)) are a special type of linux block device (that we're going to create) that either reads or writes to a remote server.

Before I talk about how these work, I need to explain the `.estar.gz` file format - it's basically a `.tar.gz` archive, but it's seekable - meaning you can `seek` specific blocks of memory inside the compressed file and read them. The format was designed by google's CRFS project then containerized it and wrote some nice go libraries for interacting with them.

Okay now here's how [MBD](#) works:

A read-only Magic Block Device is designed to allow virtual machine runtimes (like firecracker) to use a remote `.estar.gz` file as a local block device. For example, let's say I had `debian.iso.estar.gz` (a debian root filesystem) on `HOST A`. In order for firecracker to use that root filesystem on `HOST B` it has to be available as a block device - so, I create an [MBD](#) device in read-only mode that uses the `debian.iso.estar.gz` file on `HOST A` as its backing file.

This step looks something like `mbd read /dev/nbd0 debian.iso.estar.gz hostA:debian.iso.estar.gz`. It tells `mbd` to create a read-only block device at `/dev/nbd0` using `hostA:debian.iso.estar.gz`. It also says to create a local `debian.iso.estar.gz`.

When `Firecracker` tries to read from `/dev/nbd0` we'll redirect the read to the remote `.estar.gz` file (which is fast because we can seek!). The data we retrieve from that read operation we'll store in the local `debian.iso.estar.gz` so that future reads happen locally.

This means creating `/dev/nbd0` and mounting it is instantaneous since we don't really have to copy any data (and we can probably optimize this later if we need to by pre-loading common files in the background).

Okay, so this is for a read-only archive - what about writing?

If [MBD](#) could provide a writable block device that optionally synced to some remote storage layer, we'd be able to backup or move the data almost instantly.

That brings us to writeable Magic Block Devices. These are read-and-write block devices that are backed to a local file - so you can read and write from them like normal. Optionally, these devices can be set up to replicate the written data to a remote server that will then periodically back up that data to an S3 bucket (after converting it to an `.estar.gz` file).

Something really awesome that linux can do is use `dmsetup` to create a `copy-on-write` device. What this means is that you give `dmsetup` two devices - let's say `/dev/nbd0` and `/dev/nbd1`. `dmsetup` will then expose a new device `/dev/devmapper/dev1`.

Whenever a read happens on `/dev/devmapper/dev1` it'll get redirected to `/dev/nbd0` - whenever a write happens it'll get redirected to `/dev/nbd1`.

With this in place, we now have the ability to create `volumes` on architect that we can back up and restore really quickly.

And how does this help with scale builder? Because it completely solves our caching problem!

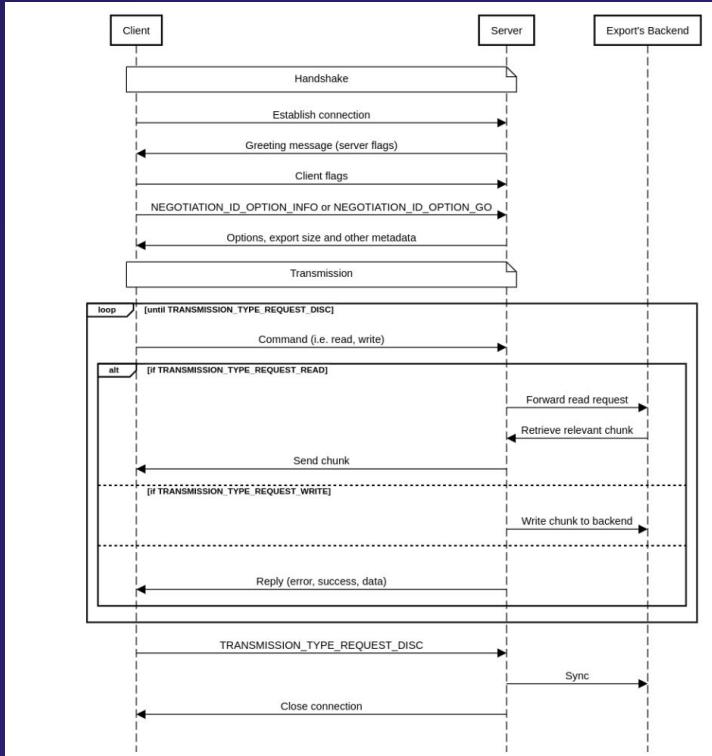
When a user wants a builder VM we can start it instantly (most of the time we'll have the read-only block device available locally but even when we don't we can just do the remote read!)

And then when the user's builds are done we wait until the machine is inactive (let's say for 10 minutes or so) and at that point we take the entire block device (the read-only layer + the write layer) and overlay them on top of each other into a new `.estar.gz` file (this is relatively slow but it happens asynchronously in the background so it's fine).

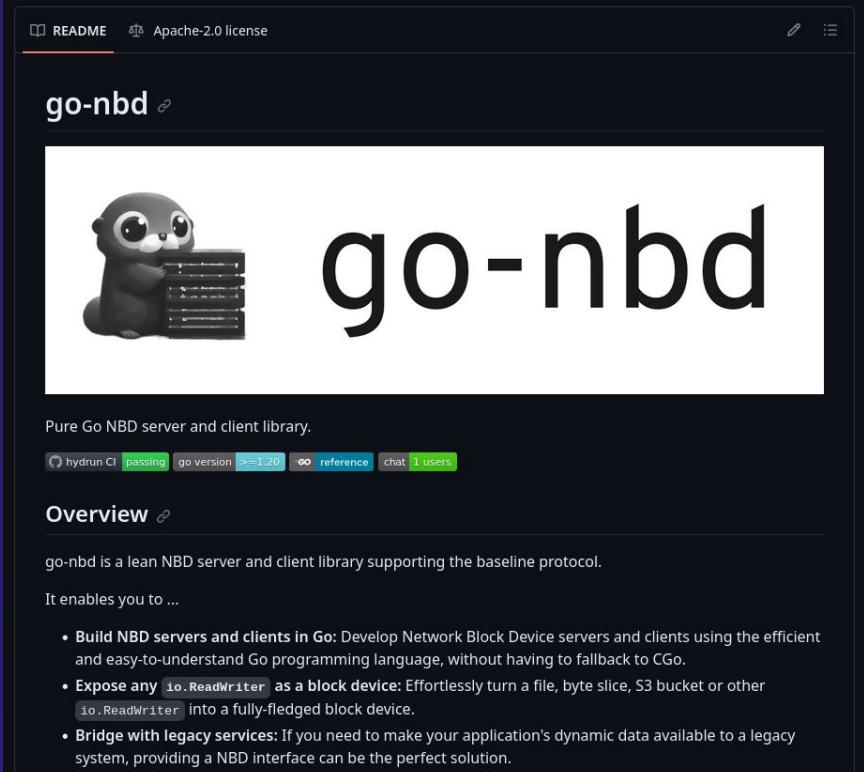
When the user wants to resume their VM all we do load up the `.estar.gz` file as a read-only block device (because it contains their recent build changes as well as the original read-only root filesystem) and they're good to go! When the new subsequent builder times out we just do the same thing again (this time replacing the `.estar.gz` file with the latest one).



# NBD



# NBD with go-nbd



The screenshot shows the GitHub page for the go-nbd project. At the top, there are links for 'README' (which is highlighted in orange) and 'Apache-2.0 license'. Below the header is the project name 'go-nbd' with a small icon. To the left of the main title is a cartoon illustration of a black otter holding a stack of four black server racks. To the right of the illustration, the word 'go-nbd' is written in a large, bold, black sans-serif font. Underneath the title, a short description reads: 'Pure Go NBD server and client library.' Below the description are several status indicators: 'hydran CI passing', 'go version >=1.20', 'reference', and 'chat 1 users'. A section titled 'Overview' follows, with a sub-section 'go-nbd is a lean NBD server and client library supporting the baseline protocol.' It also states 'It enables you to ...' and lists three bullet points:

- **Build NBD servers and clients in Go:** Develop Network Block Device servers and clients using the efficient and easy-to-understand Go programming language, without having to fallback to CGo.
- **Expose any `io.ReadWriter` as a block device:** Effortlessly turn a file, byte slice, S3 bucket or other `io.ReadWriter` into a fully-fledged block device.
- **Bridge with legacy services:** If you need to make your application's dynamic data available to a legacy system, providing a NBD interface can be the perfect solution.



# NBD with go-nbd

```
1 type Backend interface {
2     ReadAt(p []byte, off int64) (n int, err error)
3     WriteAt(p []byte, off int64) (n int, err error)
4     Size() (int64, error)
5     Sync() error
6 }
```

```
1 func Handle(conn net.Conn, exports []Export, options *Options) error
```



# What about RTT?



# r3map

A Universal Resource Mount and Migration Solution  
("remote mmap")

# r3map

The screenshot shows the GitHub README page for the r3map project. At the top, there are links for "README" and "Apache-2.0 license". Below the title "r3map" is a cartoon illustration of a white cloud with large brown eyes and a smiling mouth, resting on a green circuit board. To the right of the illustration, the word "r3map" is written in a large, bold, black sans-serif font. Below the illustration and title, the text "Remote mmap: High-performance remote memory region mounts and migrations in user space." is displayed. At the bottom of the screenshot, there are several status indicators: "hydran CI passing", "go version >= 1.20", "reference", and "chat 1 users".



Felicitas Pojtinger, 2023

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

# Concurrency

Pre-emptive pulls

Pull priority  
heuristics

Caching

Mounts &  
migrations

Lifecycle  
management/hooks

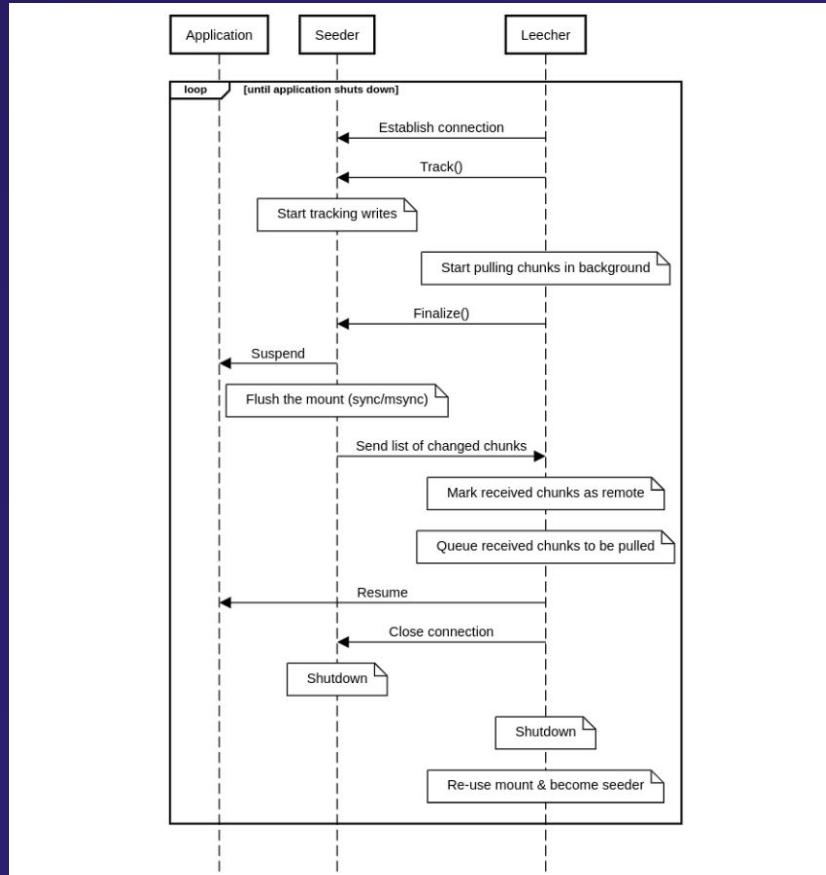
Workload analysis

Network protocol

Language-independent



# Migration



# Seeder

```
defer mgr.Close()
file, svc, err := mgr.Seed()
if err != nil {
    panic(err)
}
log.Println("Starting app on", file.Name())
```



# Leecher

```
conn, err := grpc.Dial(*raddr, grpc.WithTransportCredentials(insecure.NewCredentials()))
if err != nil {
    panic(err)
}
defer conn.Close()

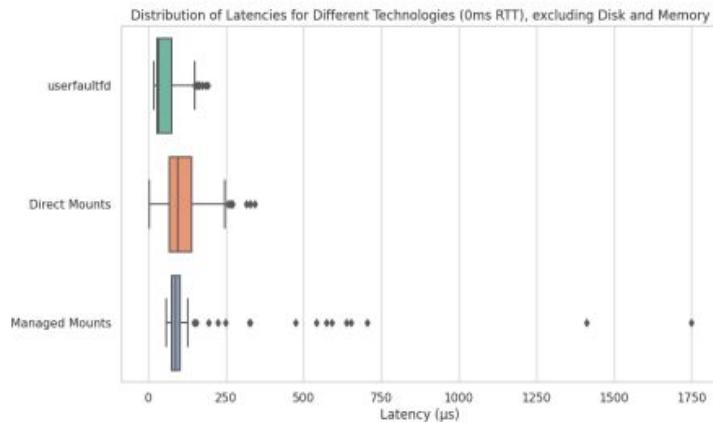
log.Println("Leeching from", *raddr)

defer mgr.Close()
finalize, err := mgr.Leech(services.NewSeederRemoteGrpc(v1.NewSeederClient(conn)))
if err != nil {
    panic(err)
}
```



# How Fast Is r3map?

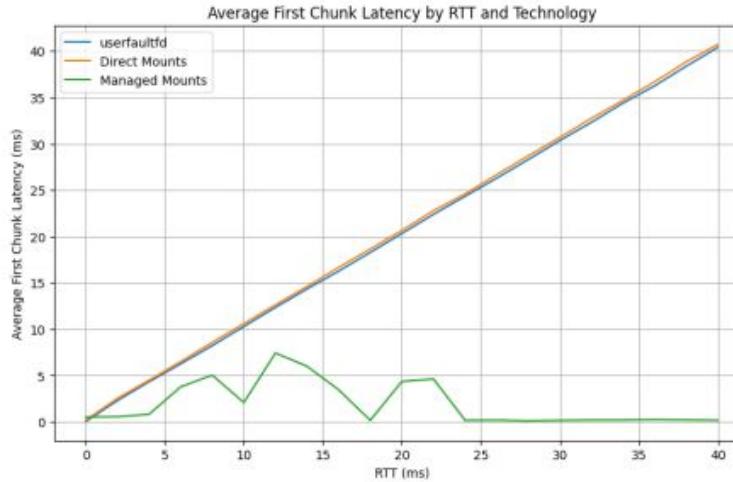
# Latency (0 ms RTT)



**Figure 9:** Box plot for the distribution of first chunk latency for userfaultfd, direct mounts and managed mounts (0ms RTT)



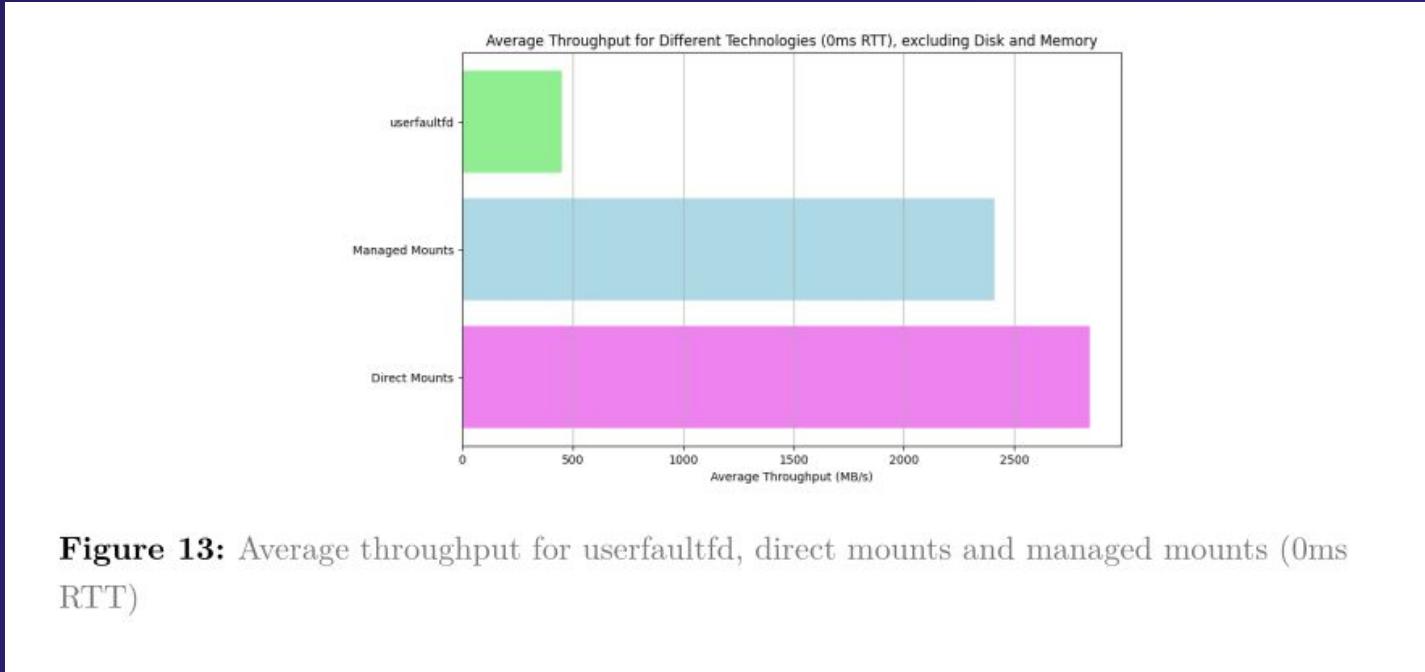
# Latency (by RTT)



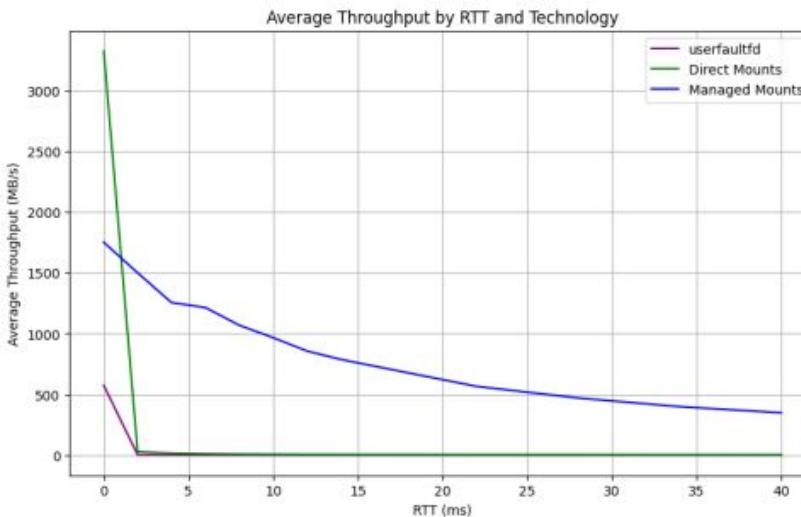
**Figure 10:** Average first chunk latency for userfaultfd, direct mounts and managed mounts by RTT



# Read Throughput (0 ms RTT)

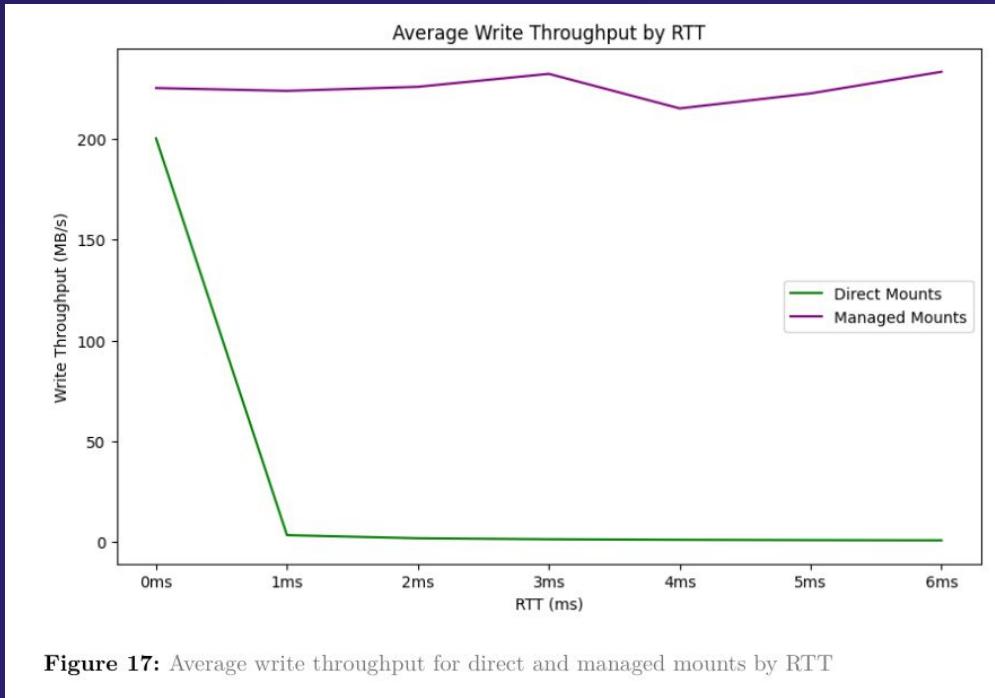


# Read Throughput (by RTT)



**Figure 15:** Average throughput for userfaultfd, direct mounts and managed mounts by RTT

# Write Throughput (by RTT)



# Using r3map with Firecracker

A screenshot of a GitHub repository page. At the top, there's a header with a menu icon, a profile picture, the repository name "loopholelabs / firecracker", a search bar, a fork button, and a user profile picture. Below the header, there are navigation links: "Code" (highlighted with an orange underline), "Issues", "Pull requests" (with a count of 2), "Actions", and a "..." button. Underneath these are four small icons: a eye icon, a document icon, a star icon, and a green "Code" button with a dropdown arrow. The main content area has a dark background with white text. It displays the repository's description: "Secure and fast microVMs for serverless computing. (Fork with live migration support)". Below the description are links to "Apache-2.0 license", "Code of conduct", and "Security policy". At the bottom, it shows statistics: "2 stars", "0 forks", "1 watching", "9 Branches", and "Tags". There are also "Activity" and "Public repository" buttons.

loopholelabs / firecracker

Code Issues Pull requests 2 Actions ...

Code

Secure and fast microVMs for serverless computing. (Fork with live migration support)

Apache-2.0 license

Code of conduct

Security policy

2 stars 0 forks 1 watching 9 Branches Tags

Activity

Public repository



4 src/utils/src/vm\_memory.rs

```
@@ -119,8 +119,8 @@ pub fn create_guest_memory(
119     for region in regions {
120         let flags = match region.0 {
121             None => libc::MAP_NORESERVE | libc::MAP_PRIVATE |
122 -             libc::MAP_ANONYMOUS,
123 -             Some(_) => libc::MAP_NORESERVE | libc::MAP_PRIVATE,
124         };
125
126         let mmap_region =
```

119 for region in regions {
120 let flags = match region.0 {
121 None => libc::MAP\_NORESERVE | libc::MAP\_SHARED |
122 + libc::MAP\_ANONYMOUS,
123 + Some(\_) => libc::MAP\_NORESERVE | libc::MAP\_SHARED,
124 };
125
126 let mmap\_region =



**feat: Apply live migration patches for v1.5**

live-migration-1.5

 pojntfx committed 2 days ago Unverified

 Showing 10 changed files with 185 additions and 9 deletions.



# Architekt

Tooling and orchestration for live migration with r3map  
and Firecracker

architekt-agent
architekt-daemon
architekt-liveness
architekt-manager
architekt-operator
architekt-packager
architekt-peer
architekt-registry
architekt-runner
architekt-worker

README

AGPL-3.0 license



# Architekt

Zero-Downtime Live Migration of Stateful VMs on Kubernetes

## Overview

This project is a work-in-progress! Instructions will be added as soon as it is usable.

## Contributing

Bug reports and pull requests are welcome on GitHub at <https://github.com/loopholelabs/architekt>. For more contribution information check out [the contribution guide](#).

## License

The Architekt project is available as open source under the terms of the [GNU Affero General Public License, Version 3](#).



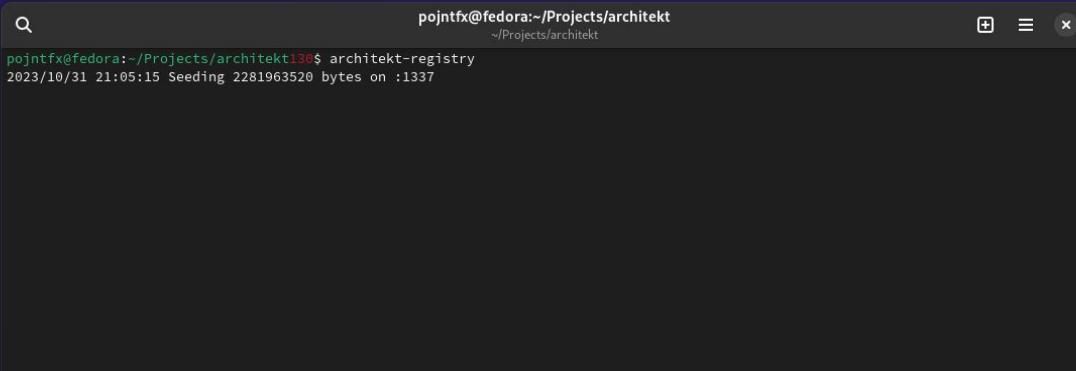
# Architekt Packager



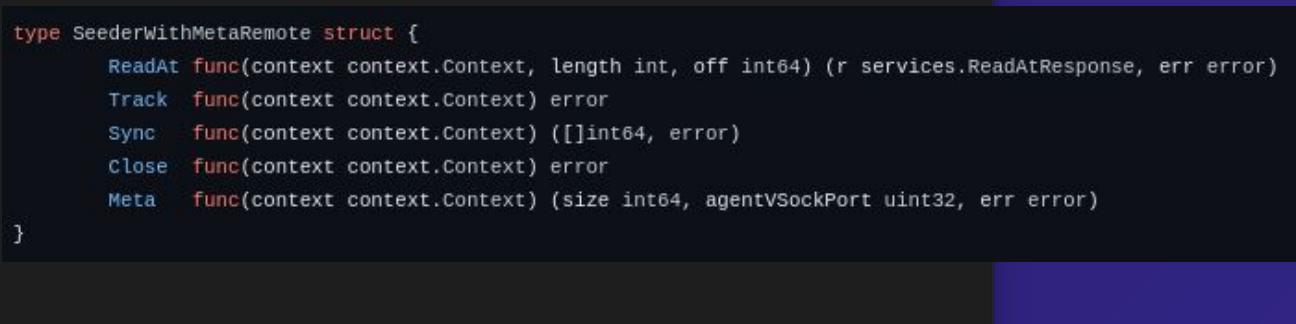
```
pojntfx@fedora:~$ architekt-packager --help
Usage of architekt-packager:
  -agent-vsock-port int
    Agent Vsock port (default 26)
  -boot-args string
    Boot/kernel arguments (default "console=ttyS0 panic=1 pci=off modules=ext4 rootfstype=ext4 i8042.noaux i8042.nomux i8042.nopnp i8042.dumbkbd rootflags=rw")
  -cgroup-version int
    Cgroup version to use for Jailer (default 2)
  -chroot-base-dir chroot
    chroot base directory (default "out/vms")
  -cpu-count int
    CPU count (default 1)
  -disk-input-path string
    Disk input path (default "out/blueprint/architekt.arkdisk")
  -enable-input
    Whether to enable VM stdin
  -enable-output
    Whether to enable VM stdout and stderr (default true)
  -firecracker-bin string
    Firecracker binary (default "firecracker")
  -gid int
    Group ID for the Firecracker process
  -initramfs-input-path string
    initramfs input path (default "out/blueprint/architekt.arkinitramfs")
  -interface string
    Name of the interface in the network namespace to use (default "tap0")
  -jailer-bin string
    Jailer binary (from Firecracker) (default "jailer")
  -kernel-input-path string
    Kernel input path (default "out/blueprint/architekt.arkkernel")
```

- ▼ blueprint
  - ≡ architekt.arkdisk
  - ≡ architekt.arkinitramfs
  - ≡ architekt.arkkernel
  - ≡ redis.ark

# Architekt Registry



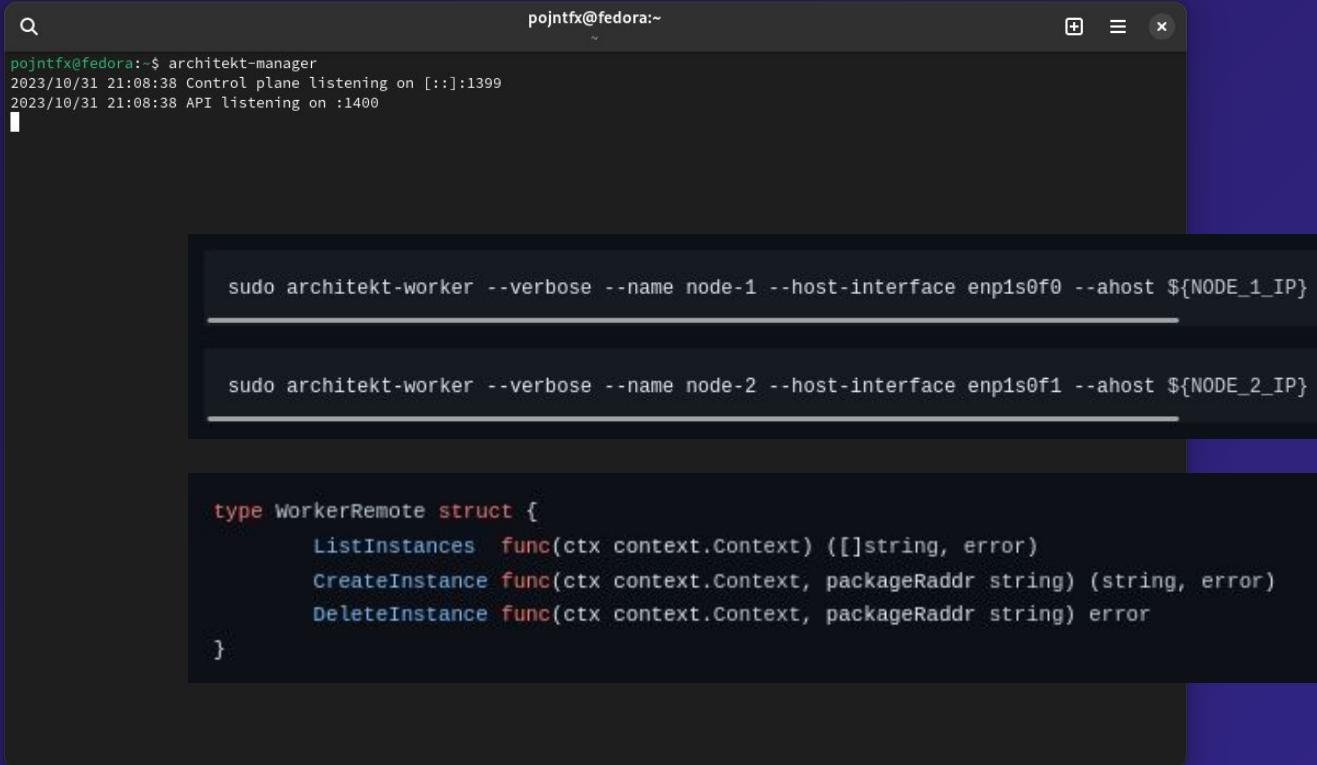
```
pojntfx@fedora:~/Projects/architekt$ architekt-registry
2023/10/31 21:05:15 Seeding 2281963520 bytes on :1337
```



```
type SeederWithMetaRemote struct {
    ReadAt func(context.Context, length int, off int64) (r services.ReadAtResponse, err error)
    Track  func(context.Context) error
    Sync   func(context.Context) ([]int64, error)
    Close   func(context.Context) error
    Meta   func(context.Context) (size int64, agentVSockPort uint32, err error)
}
```



# Architekt Control Plane



A screenshot of a terminal window titled "pojntfx@fedora:~". The window contains several code snippets related to the Architekt system.

```
pojntfx@fedora:~$ architekt-manager
2023/10/31 21:08:38 Control plane listening on [::]:1399
2023/10/31 21:08:38 API listening on :1400
```

```
sudo architekt-worker --verbose --name node-1 --host-interface enp1s0f0 --ahost ${NODE_1_IP} .
-----
```

```
sudo architekt-worker --verbose --name node-2 --host-interface enp1s0f1 --ahost ${NODE_2_IP} .
-----
```

```
type WorkerRemote struct {
    ListInstances func(ctx context.Context) ([]string, error)
    CreateInstance func(ctx context.Context, packageRaddr string) (string, error)
    DeleteInstance func(ctx context.Context, packageRaddr string) error
}
```



# Architekt Control Plane API

```
curl -v http://${CONTROL_PLANE_IP}:1400/nodes | jq  
  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances | jq  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-2/instances | jq  
  
export PACKAGE_RADDR=$(curl -v -X POST http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances/${REGISTRY_IP}:1337 | jq -r) # Create VM on node 1  
  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances | jq  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-2/instances | jq  
  
export PACKAGE_RADDR=$(curl -v -X POST http://${CONTROL_PLANE_IP}:1400/nodes/node-2/instances/${PACKAGE_RADDR} | jq -r) # Migrate VM from node 1 to node 2  
  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances | jq  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-2/instances | jq  
  
export PACKAGE_RADDR=$(curl -v -X POST http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances/${PACKAGE_RADDR} | jq -r) # Migrate VM from node 2 to node 1  
  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances | jq  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-2/instances | jq  
  
curl -v -X DELETE http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances/${PACKAGE_RADDR} # Delete VM from node 2  
  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-1/instances | jq  
curl -v http://${CONTROL_PLANE_IP}:1400/nodes/node-2/instances | jq
```



# Architekt Kubernetes Operator

```
pojntfx@fedora:~$ architekt-operator
2023-10-31T21:14:53+01:00      INFO  setup    Starting manager
2023-10-31T21:14:53+01:00      INFO  controller-runtime.metrics    Starting metrics server
2023-10-31T21:14:53+01:00      INFO  starting server {"kind": "health probe", "addr": "[::]:8081"}
2023-10-31T21:14:53+01:00      INFO  controller-runtime.metrics    Serving metrics server {"bindAddress": ":8080", "secure": false}
2023-10-31T21:14:53+01:00      INFO  Starting EventSource {"controller": "instance", "controllerGroup": "io.loopholelabs.architekt", "controllerKind": "Instance", "source": "kind source: *v1alpha1.Instance"}
2023-10-31T21:14:53+01:00      INFO  Starting EventSource {"controller": "instance", "controllerGroup": "io.loopholelabs.architekt", "controllerKind": "Instance", "source": "kind source: *v1.Deployment"}
2023-10-31T21:14:53+01:00      INFO  Starting Controller {"controller": "instance", "controllerGroup": "io.loopholelabs.architekt", "controllerKind": "Instance"}
2023-10-31T21:14:53+01:00      INFO  Starting workers {"controller": "instance", "controllerGroup": "io.loopholelabs.architekt", "controllerKind": "Instance", "worker count": 1}
```



# Architekt Instance CRD

```
apiVersion: io.loopholelabs.architekt/v1alpha1
kind: Instance
metadata:
  name: redis
spec:
  packageRaddr: localhost:1337 # If you change this, the VM is recreated from the new `packageRaddr`
  nodeName: minikube # If you change this, the VM is migrated
```

```
kubectl apply -f config/samples/architekt_v1alpha1_instance.yaml # Create
```

```
kubectl get -o yaml instance.io.loopholelabs.architekt/redis # List VMs
```

```
kubectl apply -f config/samples/architekt_v1alpha1_instance.yaml # Migrat
```

```
kubectl apply -f config/samples/architekt_v1alpha1_instance.yaml # Migrat
```

```
kubectl delete -f config/samples/architekt_v1alpha1_instance.yaml # Delete
```



# VM Network Interfaces

```
func NewNamespace(
    id string,
    hostInterface string,
    namespaceInterface string,
    namespaceInterfaceGateway string,
    namespaceInterfaceNetmask uint32,
    hostVethInternalIP string,
    hostVethExternalIP string,
    namespaceInterfaceIP string,
    namespaceVethIP string,
    blockedSubnet string,
    namespaceInterfaceMAC string,
) *Namespace {
```

```
    if err := CreateNAT(d.hostInterface); err != nil {
        return err
    }

    hostVethIPs := NewIPTable(d.hostVethCIDR)
    if err := hostVethIPs.Open(ctx); err != nil {
        return err
    }

    if err := d.namespaceVethIPs.Open(ctx); err != nil {
        return err
    }
```





IP 1 —————→ IP 2



# Lynk

A new approach to networking

Independent of network  
topology

Can be deployed behind  
firewalls

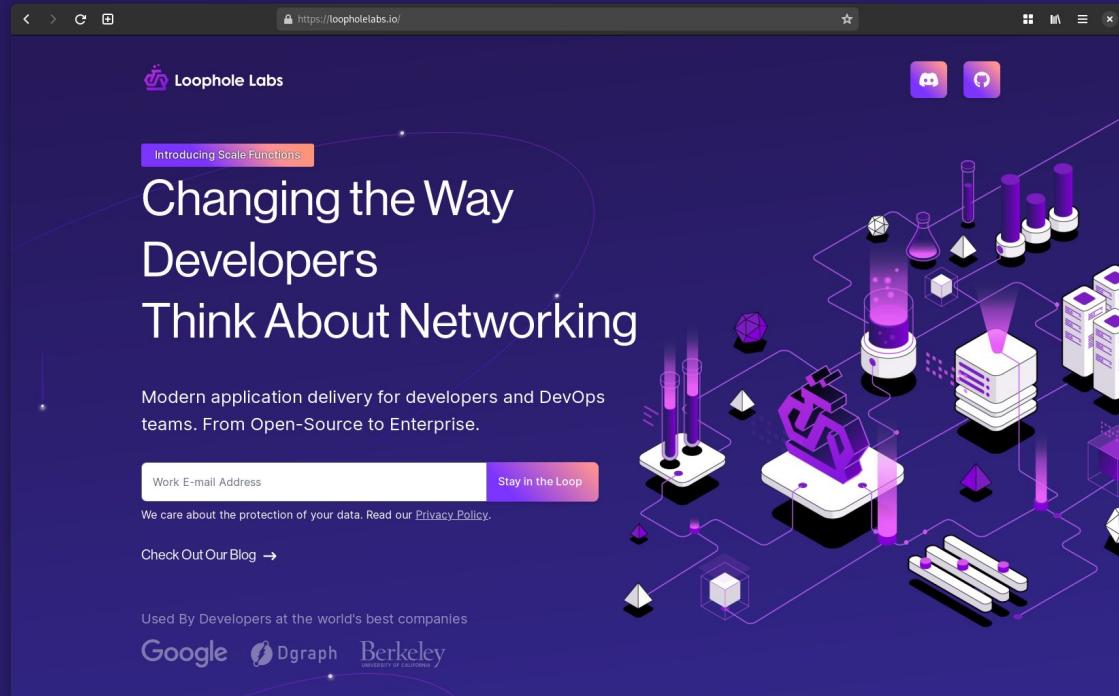
Nodes do not  
require a public IP

Automatically handles  
TLS encryption in flight

Can keep connections alive  
during migration

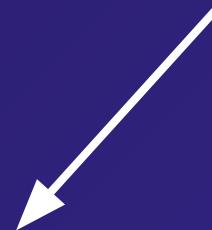
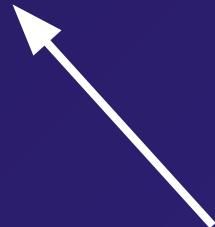
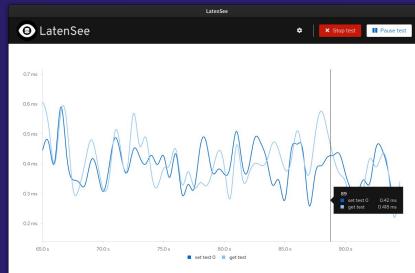
Tightly integrated  
with Architekt





loopholelabs.io

# But Does This Actually Work?





Google Maps

https://www.google.com/maps/@47.4622154,-56.8888594,6023286m/data=!3m1!e3?entry=ttu

120%

FP OI BI BA GH OB YHN PH BS MD T TW IG YM YT BP EL DC EP

Search Google Maps

Saved

Recents

Frankfurt

Chicago

Yordan's Pizza

Layers

Measure distance

Click on the map to add to your path

Total distance: 4,585.02 mi (7,378.87 km)

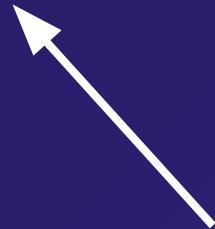
3D

200 m

Imagery ©2023 Data SIO, NOAA, U.S. Navy, NGA, 6EBCO, IBCAO, Landsat / Copernicus, U.S. Geological Survey, Imagery ©2023 TerraMetrics, Map data ©2023 Google, INEGI, United States, Terms, Privacy, Send Product Feedback

# Live Demo





# Live Demo

Connect to:

**[minecraft.architect.run](https://minecraft.architect.run)**

**Minecraft 1.12.2**



# Recap: What Can We Do With This New Compute Unit?

# Simplifying deployment and debugging



# Moving between cloud providers without downtime



Moving long-running  
processes to more  
powerful servers (while  
continuing computing)

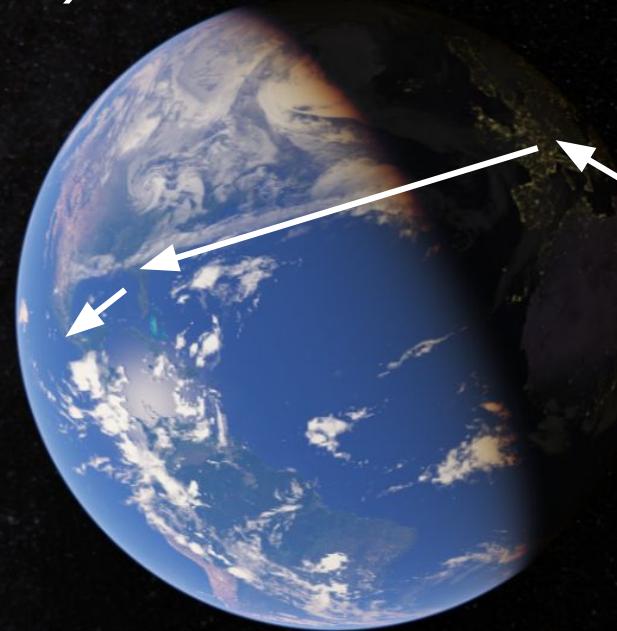


Image source: <https://archive.ph/pL8ih>, last accessed 2023-10-31

# Migrating full desktop/mobile apps between devices (in milliseconds)



Moving stateful services closer  
to your users as they access  
them (without downtime)





[pojntfx/networked-linux-memsync](#)

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

★ 7

0

Bachelor's Thesis

## Efficient Synchronization of Linux Memory Regions over a Network

A Comparative Study and Implementation

Author: Felicitas Pojtinger

University: Hochschule der Medien Stuttgart

Course of Study: Media Informatics

Date: 2023-08-03

Academic Degree: Bachelor of Science

Primary Supervisor: Prof. Dr. Martin Goik

Secondary Supervisor: M.Sc. Philip Betzler



Felicitas Pojtinger, 2023

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

Efficient Synchronization of Linux Memory Regions over a Network

2023-08-03

### Contents

1	Introduction	12
2	Technology	13
2.1	User Space and Kernel Space	13
2.2	Linux Kernel	13
2.3	UNIX Signals and Sockets	14
2.4	Memory Optimization	15
2.4.1	Principle of Locality	15
2.4.2	Memory Hierarchy	16
2.5	Memory in Linux	17
2.5.1	Memory Management	17
2.5.2	Swap Space	18
2.6	Page Faults	19
2.7	mmap	19
2.8	inotify	20
2.9	Linux Kernel Caching	21
2.10	Networking	21
2.10.1	RTT, LAN and WAN	21
2.10.2	TCP, UDP, TLS and QUIC	22
2.11	Delta Synchronization	24
2.12	File Systems in User Space (FUSE)	25
2.13	Network Block Device (NBD)	27
2.14	Virtual Machine Live Migration	29
2.14.1	Pre-Copy	29
2.14.2	Post-Copy	30
2.14.3	Workload Analysis	30
2.15	Streams and Pipelines	31
2.16	Go	32
2.17	RPC Frameworks	32
2.17.1	gRPC and Protocol Buffers	32
2.17.2	fRPC and Polyglot	33
2.18	Data Stores	33
2.18.1	Redis	33
2.18.2	S3 and Minio	34
2.18.3	Cassandra and ScyllaDB	34



### [pojntfx/r3map](#)

High-performance remote memory region mounts and migrations in user space.

★ 31

¥ 0



### [pojntfx/go-nbd](#)

Pure Go NBD server and client library.

★ 305

¥ 15



### [pojntfx/networked-linux-memsync](#)

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

★ 7

¥ 0



### [loopholelabs/architekt](#)

Zero-Downtime Live Migration of Stateful VMs on Kubernetes

★ 5

¥ 0



### [loopholelabs/firecracker](#)

Secure and fast microVMs for serverless computing. (Fork with live migration support)

★ 2

¥ 0



# Come talk to us!





# Felicitas Pojtinger

bluesky: [@felicitas.pojtinger.com](https://bluesky.fedibird.com/@felicitas.pojtinger)

mastodon: [@pojntfx@mastodon.social](https://mastodon.social/@pojntfx)

github/twitter: [@pojntfx](https://github.com/pojntfx)

linkedin: [in/pojntfx](https://www.linkedin.com/in/pojntfx)

web: [felicitas.pojtinger.com](https://felicitas.pojtinger.com)





github/twitter:

@LoopholeLabs

discord:

loopholelabs.io/discord

web:

loopholelabs.io

