FERMYON

Building new serverless primitives with Wasm

December 2024

Radu Matei

- CTO & co-founder of Fermyon
- Creator of the Spin project
- Maintainer and contributor to several cloud-native projects, at the intersection of distributed systems, security, and developer experience
- Loves bicycles, classical music, and bubble tea



@radu-matei.com on Bsky, https://radu-matei.com, radu@fermyon.com

Fermyon.com

- started Fermyon in November '21 to work on Wasm
- built the Spin open source project
- Fermyon Cloud as the managed service for running serverless Wasm
- SpinKube and more cloud native tooling for Wasm
- core contributors and maintainers to ecosystem projects such as Wasmtime

Serverless: A coding pattern in which the developer does not write a server process.

The entry point is an event handler **function**.



From VMs to Functions



FERMYON

Developers love* writing Functions

Why FaaS:

• spending time writing application logic, not deploying servers

 not thinking about scaling up or down, and not being charged for idle resources**

*For the right use cases

**Unless on a platform that still uses VMs and "dedicated plans"

BUT...

- why do I have to rebuild my "serverless function" for ARM64?
- why do I experience 200+ms cold starts?
- once I build the function for this platform, can I run it elsewhere?

the developer experience can be a pain*

The Function is still executed in a microVM or container

 most popular FaaS platforms still use microVMs or containers to execute user functions.

 portability and binary size (and by extension, cold starts) we are experiencing with these platforms are stemmed in the unit of execution: microVMs or containers.

Efficiency for serverless

As a platform builder using microVMs or containers:

• **over**provision and pass the cost to end users

OR

 underprovision and pass the cold starts to end users



Time

Efficiency for serverless

How do we make the shape of cloud compute better match the demand of this curve? Startup time Fully Utilized Overprovisioned

emand

And by extension, pass BOTH cost savings AND fast startup to end users?

This is a big issue already for a singleregion datacenter...

The problem becomes exponential if you need to schedule a function in hundreds of regions around the world and execute it as close as possible to an end user request, as fast as possible.

Do we really need microVMs and containers for FaaS?

From a developer perspective:

• if you need a full Linux OS in your function: YES.

• but for a growing number of use cases: NO (and platforms like Cloudflare Workers show us that with the right API surface, developers can achieve a lot without a full Linux OS available in their function)

What if we had another unit of deployment and execution for FaaS?

FERMYON

Enter Wasm

Or WebAssembly

What is WebAssembly (Wasm)?

 a target for compiling programs written in C/C++, Rust, Python, JavaScript, Kotlin, and more into a compact and portable intermediate binary format

• a runtime for executing these binaries in an isolation sandbox

What is WebAssembly (Wasm)?

originally built to port programs such as Photoshop or Google
 Earth to the browser, it evolved to run in other places, outside the browser

Why Wasm for serverless

- **isolation sandbox** (Wasm code does not have access to data outside its own linear memory or to APIs unless allowed by the host)
- cross platform and cross architecture portability (once compiled, a Wasm component can run on Linux, Windows, macOS, and on x86, ARM64, RISC-V)
- fast startup and near-native execution (Wasm components start in less than 1 millisecond if available on disk, and can be as little as a few MBs or less)
- can support a wide range of programming languages (write your functions in Rust, JavaScript, Python, Go, or more)

Academic research for Wasm and FaaS



DSPIN

Single binary tool



One tool for development and runtime

Many languages



Go, Rust, C#, Java, JavaScript, or any with Wasm+WASI support

Easy to get started



Use the template library, or provide your own and share components

No boilerplate required



Spin triggers get you right to the important part

Work with data



Use Postgres, Redis and file storage to persist your data

Spin

The open source tool for building WebAssembly serverless apps

With Spin, you can create a new serverless application (powered by one or more serverless functions) with just a few commands.

https://developer.fermyon.com/spin



(main) \$ spin new

Pick a template to start your project http-go (HTTP request handler usir http-swift (HTTP request handler ι > http-rust (HTTP request handler us http-zig (HTTP request handler usi kv-explorer (Explore the contents redis-go (Redis message handler us redis-rust (Redis message handler

😑 🧊 🕼 fermyon / spin ★ 5.5k 🗸		/] to search	8 + • O II A	
<> Code ⊙ Issues 255 ∰ Bugs 14	♡ Pull requests 25 및 Discussions ⓒ Action	ons 🗄 Projects	3 印 Wiki 錄 Settings	
F spin Public	🔊 Edit Pins 👻	🕑 Unwatch 45 👻	ঔ Fork 257 → 🕇 Starred 5.5k	
양 main ▾ 양 55 Branches ♡ 98 Tags	Q Go to file t	+ <> Code -	Spin is the open source developer tool for building and running	
😭 itowlson Merge pull request #2937 from mikkelhegn/static-fil 🚥 🧹 666faa9 · 4 days ago 🕚 4,086 Commits WebAssembly.				
.cargo	Add taplo.toml and format tomls	4 months ago	♂ developer.fermyon.com/spin	
local devcontainer	Remove bindle	last year	serverless webassembly spin	
o .github	Merge pull request #2878 from michelleN/issuef	last week	fermyon	
🛤 .vscode	Update recommended extensions to point at the	4 months ago		
Crates	Relax wasi 0.2 constraint	5 days ago	Sode of conduct	
Cross	set CROSS_SYSROOT=/ to fix cross-rs openssl issue	10 months ago	ৰ্বায় Security policy	
deploy	ref(docs): update to run/deploy via OCI	last year	-^- Activity	
docs	A few fixes after yesterday's mergepalooza	4 months ago	☆ 5.5k stars	
a examples	chore(*): post-3.0.0 version bump	3 weeks ago	⊙ 45 watching ও 257 forks	
hack/o11y-stack	feat(telemetry): Send logs to OTel collector direct	7 months ago	및 3 years old	
src	Unhide app splitting flag	2 months ago	S weeks ago	

https://github.com/fermyon/spin



FERMYON

Demo time

https://github.com/fermyon/spin

radu@asahi 🍄 default in github.com/misc/hk bombardier localhost:3000 -c 512 Bombarding http://localhost:3000 for 10s using 512 connection(s) Done! Statistics Avg Stdev Max Reqs/sec 138117.91 13199.46 161568.99 Latency 3.70ms 1.36ms 105.99ms HTTP codes: 1xx - 0, 2xx - 1381174, 3xx - 0, 4xx - 0, 5xx - 0others - 0 Throughput: 21.33MB/s

https://github.com/fermyon/spin

FERMYON



What?s new in Spin?



Once you build your Spin app...

- run it with `spin up` on your own servers
- deploy it to Fermyon Cloud, soon run it globally
- deploy it to Kubernetes using the CNCF SpinKube project
- run it in other orchestrators such as Nomad or OpenShift
- run it directly in Docker Desktop using `docker run`

Wasm and the component model

- Spin is built on top Wasm components, defined using WIT (Wasm Interface Types)
- components are the core unit of execution, and encapsulate the function
- components can be statically analyzed and their capabilities allowed or declined

A Spin platform

2

6

8

- package fermyon:spin@3.0.0;
- 3 /// The full world of a guest targeting an http-trigger
- 4 world http-trigger {
- 5 **include** platform;
 - export wasi:http/incoming-handler@0.2.0;
- 7 }
- 9 /// The imports needed for a guest to run on a Spin host
- 10 world platform {
- 11 **include** fermyon:spin/platform@2.0.0;
- 12 include wasi:keyvalue/imports@0.2.0-draft2;
- 13 import spin:postgres/postgres@3.0.0;
- 14 import wasi:config/store@0.2.0-draft-2024-09-27;
- 15

}

A Spin function

- 1 package root:root;
- 2
- 3 world root {
- 4 **import** wasi:io/...@0.2.0;
- 5 import wasi:clocks/...00.2.0;
- 6 import wasi:http/...00.2.0;
- 7 **import** wasi:config/...@0.2.0
- 8 **import** wasi:cli/...00.2.0;
- 9 import wasi:filesystem/...00.2.0;
- 10 import wasi:sockets/...@0.2.0;
- 11 import wasi:random/...@0.2.0;
- 12 import wasi:keyvalue/...@0.2.0-draft2;
- 13 .

15

- 14 **import** fermyon:spin/...@3.0.0;
- 16 export wasi:http/incoming-handler@0.2.0;
 17 }

What about multi-region?



Security and Correctness in Wasmtime

Sep 13, 2022 Nick Fitzgerald

The essence of software engineering is making trade-offs, and sometimes engineers even trade away security for other priorities. When it comes to running untrusted code from unknown sources, however, exceptionally strong security is simply the bar to clear for serious participation: consider the extraordinary efforts that Web browser and hypervisor maintainers take to preserve their systems' integrity. WebAssembly runtimes also run untrusted code from untrusted sources, and therefore such efforts are also a hard requirement for WebAssembly runtimes.

bytecodealliance.org/articles/security-and-correctness-in-wasmtime

What about stronger isolation?

😑 🧊 🐠 hyperlight-dev / hyperlight 🕸	2.2k 🗸	Type 🕖 to search	8 + • O IL &	
<> Code 🕑 Issues 8 🔅 Bugs 1 🕅	Pull requests 💈 🕞 Actions 😒 Relea	ses 2		
		💿 Watch 16 👻	양 Fork 76 → ☆ Star 2.2k	
양 main 👻 양 12 Branches 🟷 2 Tags	Q Go to file	t + <> Code +	Hyperlight is a lightweight Virtual Machine Manager (VMM) designed to	
	23192b7 · 5	days ago 🕚 69 Commits	be embedded within applications. It enables safe execution of untrusted	
levcontainer	remove .vscode settings (#66)	2 weeks ago	code within micro virtual machines with very low latency and minimal overhead.	
o .github	Add check in CI to make sure cargo feature	cras last week		
kø dev	publish to crates.io	last month	🛱 Readme	
docs	Fix otlp tracing (#61)	5 days ago		
hack	The initial Hyperlight Commit 🎉	last month	办 Security policy	
proposals	[docs] small patch to HIP template	5 days ago	-\- Activity	
kp src	[docs] small patch to HIP template Fix otlp tracing (#61) signed-off-by: danbugs <danilochiarlone@gmail.com></danilochiarlone@gmail.com>	5 days ago	 E Custom properties ☆ 2.2k stars 	
🦾 .editorconfig	The initial Hyperlight Commit 🎉	last month	⊙ 16 watching	
🔶 .gitignore	remove .vscode settings (#66)	2 weeks ago	양 76 forks	
CODE OF CONDUCT.md	The initial Hyperlight Commit 🞉	last month		

github.com/hyperlight-dev/hyperlight

Executing Wasm Functions at scale

 we can control the amount of memory, CPU time, network bandwidth a Function can use, and terminate at any given point in time

 we can centrally compile Wasm ahead-of-time to the specific OS/arch of data plane nodes and distribute highly optimized machine code

Executing Wasm Functions at scale #2

 we optimistically schedule thousands* of Functions on any given node in Fermyon Cloud (use cases favor short, bursty executions compared to long-running functions)

 with a mix of Linux x86 and ARM64 nodes, depending on efficiency, region availability

 we suspend in memory any function waiting for I/O such as file access or outbound network calls

Executing Wasm Functions at scale #3

 for high core count machines, we can pin the Wasm runtime on a subset of cores to minimize IPIs

Check out Spin and SpinKube.dev!

Get started with Spin and SpinKube!

Spin: <u>https://github.com/fermyon/spin</u> SpinKube: <u>https://github.com/spinkube</u>

