Connect 2025

RamaLama

Making AI Workloads on ARM Boring

Sergio Lopez Pascual & Eric Curtin

An **easy**, **secure**, and **open** way to run and serve LLMs locally.



Have you...

Experimented with an LLM?

Run an LLM locally?

Built an Al-powered application?



Introducing RamaLama



Easy to use

Run or serve a chatbot with a single command



Secure

Experiment with LLMs while keeping your data private



Built using open standards and AI building blocks

Red Hat

> ramalama run gemma3

Better with containers: easy, secure and open

5



Running the model



Complexity made easy

Which should I choose?

quantization parameter size alignment technique

Shortnames

9 0 G hoppinglace extracted bearch-granite	* D 9 1	
Da Admin Da Collaboration Da Projecto Da Offenio	ge 🗁 Reemalaren 🗄 Jira 🕒 Red Hart Link Lip 🛤 Hill Fille L. Str. 📥 Cloud Accesse 👔 Fillen Multicale Hub	
A Happing Face	A Models Hageng New X + V	
International and a second secon	Constrained Constr	9 :

> ramalama run gemma3

It's hard to get hardware acceleration working!

Performance made easy

Auto-detects system configuration and switches to the best performing accelerator stack

- Out-of-the-box easy: CPU-only
- Minimal setup: open source accelerator stack (Vulkan)
- Full performance: vendor-specific accelerator stack



Security: local and containerized



Keeping your data secure



- **Local:** where is your data going?
- Containerized: provides tighter control and isolation



Open: maximizing choice

Containerization enables **modularity**, easily swapping out for new blocks







Serve with standard APIs

> ramalama serve gemma3

OpenAl REST APIs

Industry standard endpoint APIs

Llama Stack

Connecting the front and back-ends together



Getting to production

11

Deploy to **Kubernetes** without being a Kubernetes expert

	Single-Node: Quadlet
ramalama	model content data
	Quadlet configs
model content data	
	Multi-Node: Kubernetes
container engine	model content data
£	K8s commands
	Red H

Al the open source way



Cloud Native Al Runs on Arm @ Kubecon

A collaborative effort between RamaLama/Universal

Rlup/Amporo/System76





Apple Silicon M1 Chip



14







ked Hat

Why Arm for AI?

- Power efficiency (AI workloads are power hogs)
- Multiple components for inference: GPU, NPU, CPU, etc.
- Scalable architecture from iot to edge to cloud (in terms of horizontal and vertical scalability)
- Edge AI brings lower-latency AI responses
- High memory bandwidth
- Unified memory



Why RamaLama for Arm?

- Other AI tools tend to focus on Nvidia for x86_64, maybe macOS on ARM for Desktop usage, but often not Linux on ARM
- Avoid lock-in to a specific piece of hardware
- Avoid lock-in to a specific model registry (ollama, hf, dockerhub, etc.)
- Cloud native approaches, can use existing container infrastructure



Vulkan as common language for GPUs

- Supporting multiple frameworks is cumbersome.
 - Multiple container images.
 - Very large footprint.
- GLSL is expressive enough for implementing good

performing shaders/kernels.

- Open Source stack (Mesa).

macOS: enabling GPU in containers

- Containers on macOS run in a VM





macOS: enabling GPU in containers

- Create a machine with libkrun as provider

Settings	Resources > Podman
Resources	Ereate Podman machine
Proxy	CPU(s): 6 🏾 🌌
Registries	• — — • — — — — — — — — — — — — — — — —
Authentication	Memory: 4 GB 🖉
CLI Tools	Disk size: 100 GB 🖋
Kubernetes	• • ·
Experimental 🗕	Image Path (Optional):
Preferences >	
	Image URL or image reference (Optional):
	Use 'registry/org/image:version' for image reference
	Machine with root privileges:
	Enabled
	Provider Type:
	GPU enabled (LibKrun)
	Start the machine now:
	Enabled
	Close Create



macOS: running ramalama

- On macOS, execute "ramalama run gemma3"
- Inference engine (llama.cpp, whisper.cpp, etc.) will run on
 Linux within podman-machine



Linux llama.cpp binaries running on macOS?

- Yes and with accelerated GPU access

21



Linux llama.cpp binaries running on macOS?

- llama-server binary is build with vulkan support
- llama-server uses mesa implementation for vulkan
- This communicates with macOS via virtio-gpu
- Which speaks to MoltenVK
- MoltenVK acts as the translator, translates Vulkan to Metal so we can access the GPU





Demo

- Lets demo "ramalama serve gemma3" and Open WebUI

Red Hat

24



Join the conversation on Matrix!

Join Us!

github.com/containers/ramalama

https://ramalama.ai

