



KubeCon

CloudNativeCon

THE LINUX FOUNDATION

S OPEN SOURCE SUMMIT











China 2024

Wasm for Portable Al Inference Across GPUs, CPUs, OS & Cloud-Native Environments

Hung-Ying Tai, WasmEdge Miley Fu, WasmEdge

GitHub / Twitter: @mileyfu / @mileyfu @hydai / @hydai_tw https://github.com/WasmEdge/WasmEdge

https://github.com/LlamaEdge/LlamaEdge

Wasm, on the server side?

The **initial Wasm landscape**, published in time for the WasmCon conference, includes 11 categories and 120 projects or products, representing \$59.4B in total economic value. The Wasm landscape is divided into two large areas: Dev (application development) and Ops (application deployment).

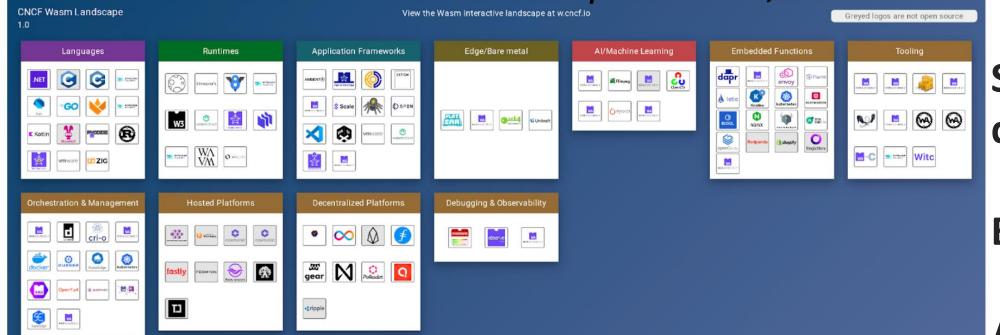
Wasm landscape September 6, 2023

Microservices

SaaS/UDF

Streaming data





Why Wasm for Al Inference?









- Problem: The need for consistent, efficient AI inference across diverse environments.
- Key Points:
 - Cross-platform consistency (run anywhere without rebuilding image targeting different underlying architecture or cuda version etc).
 - High performance and small footprint.
 - Security and sandboxing benefits.

Wasm in Cloud-Native Environments









- Integration with Kubernetes for scalable AI deployments.
- Use with Docker for containerized AI inference.
- Benefits for microservices and serverless architectures.

WasmEdge: A WebAssembly Runtime Optimized for Al









- China 2024

- Lightweight, high-performance WebAssembly runtime.
- Designed for AI, IoT, and edge computing.
- Compatibility with TensorFlow Lite, ONNX, and other AI frameworks.

LlamaEdge, Al tools built on top of WasmEdge runtime









China 202

a versatile AI deployment platform that enables running large language models (LLMs) across various environments, including edge devices, cloud infrastructure, and multiple platforms.

- across various hardware and OS.
 - CPU and GPU compatibility.
 - Execution on different operating systems (Windows, Linux, macOS).
 - Flexibility in deployment: from PCs to cloud environments.









Intel® Extension for **Transformers**





Tensor RT













Key Features and Benefits







- Lightweight: The core runtime and API server are less than 30MB, making it suitable for resource-constrained devices.
- Fast: Takes full advantage of device hardware and software acceleration for optimal performance.
- Cross-platform: Enables LLM deployment across CPUs, GPUs, TPUs, NPUs etc.
- Open-source: Built on open standards, fostering community contributions and customization.
- Privacy-focused: Allows for local LLM execution, protecting sensitive data.

LlamaEdge in Action









- integrating Wasm with cloud-native AI workloads.
 - Combining Wasm with Kubernetes for scalable AI inference.
 - Leveraging Wasm containers in multi-cloud environments.
 - Using Wasm for serverless AI functions.

Kubenetes			
KubeEdge	kind	OpenYurt	SuperEdge
k8s	k3s	microk8s	minikube
High-level container runtimes			
CRI-O	containerd	Podman	Docker
Low-level container runtimes			
gVisor	runc	crun	youki
Traditional Container Images		WebAssembly App Images	

Demo Time!









- China 2024

Setting up LlamaEdge API server.

- Packaging it as a Wasm container image.
- Running the image using Docker.
- Deploying in a Kubernetes cluster.

Demo 1: Running LlamaEdge API Server









- Build LlamaEdge API Server from source
 - Install Rust + wasm32-wasip1 target
 - rustup target add wasm32-wasip1
 - Compile it
 - git clone git@github.com:LlamaEdge/LlamaEdge.git
 - cd LlamaEdge/api-server
 - cargo build --release --target wasm32-wasip1
 - # The built wasm is here
 - ./target/wasm32-wasip1/release/llama-api-server.wasm
- Or download the pre-built assets on the LlamaEdge repo
 - curl -LO
 https://github.com/LlamaEdge/LlamaEdge/releases/latest/
 download/llama-api-server.wasm

Demo 2: Run LLM using GPU in a container









- Workflows
- Podman/Docker
 - -> crun
 - -> libcrun-wasmedge
 - -> wasmedge-runtime
 - -> wasmedge-llm-plugin (based on llama.cpp)

Demo 2: Writing the Containerfile









- China 2024

- Pretty simple; developers and users are happiest.
- Just add the wasm file into container file

- FROM scratch
- ADD llama-api-server.wasm /
- CMD ["/llama-api-server.wasm"]

Demo 2: Pack and Publish images









- China 2024

- Still easy, still happy
- Pack
 - [docker/podman] build . \
 - --platform wasip1/wasm \
 - -t hydai/kubecon2024:llama-api-server
- Publish
 - docker push hydai/kubecon2024:llama-api-server

Demo 2: Run LLM using GPU in a container









China 2024

- The complex one, no longer happy
- Set up the environment
 - Install CUDA Driver
 - Install NVIDIA Container Toolkit
 - Install CDI configuration
 - Install wasmedge libraries and LLM plugins
 - Build crun with wasmedge support from source
- Steps:
 - Podman:

https://wasmedge.org/docs/develop/deploy/gpu/podman wasm gpu

Docker:

https://wasmedge.org/docs/develop/deploy/gpu/docker wasm gpu

Demo 2: Run LLM using GPU in a container









hina 2024

- # The first two lines are mapping WasmEdge plugins for enabling the Wasm+LLM support
- # The cuda related lines are required if you want to use CUDA. Otherwise, it can run with CPU only.
- # We also need to map the model folder to `/resource` because we don't want to package a whole model into a container image.
- # There are some environment variables are needed for the safety issue. It will let `crun` know which plugins and models can be accessed.
- Docker/podman run \
- -v ~/.wasmedge/plugin/libwasmedgePluginWasiNN.so:/.wasmedge/plugin/libwasmedgePluginWasiNN.so \
- -v /usr/local/cuda/targets/x86_64-linux/lib/libcudart.so.12:/lib/x86_64-linux-gnu/libcudart.so.12 \
- -v /usr/local/cuda/targets/x86_64-linux/lib/libcublas.so.12:/lib/x86_64-linux-gnu/libcublas.so.12 \
- -v /usr/local/cuda/targets/x86_64-linux/lib/libcublasLt.so.12:/lib/x86_64-linux-gnu/libcublasLt.so.12 \
- -v /lib/x86_64-linux-gnu/libcuda.so.1:/lib/x86_64-linux-gnu/libcuda.so.1 \
- -v /disk:/resource \
- --env WASMEDGE_PLUGIN_PATH=/.wasmedge/plugin \
- --env WASMEDGE_WASINN_PRELOAD=default:GGML:AUTO:/resource/Meta-Llama-3-8B-Instruct-Q5_K_M.gguf \
- -p 8080:8080 \
- --rm --device nvidia.com/gpu=all --runtime=crun \
- --annotation=module.wasm.image/variant=compat-smart --platform wasip1/wasm \
- hydai/kubecon2024:llama-api-server llama-api-server.wasm -p llama-3-chat









```
Just like Demo 2, set up the environment again
sudo ./kubernetes/cluster/kubectl.sh run -i --restart=Never testggml -
 -image=ghcr.io/captainvincent/runwasi-demo:llama-api-server --
 annotations="module.wasm.image/variant=compat-smart" --overrides='
  "apiVersion": "v1",
   "kind": "Pod",
   "metadata": {
     "name": "testggml"
   "spec": {
     "hostNetwork": true,
     "containers": [..omitted..],
     "volumes": [..omitted..]
```









```
"containers": [
   "name": "simple",
   "image": "ghcr.io/captainvincent/runwasi-demo:llama-api-server",
   "command": ["/app.wasm", "-p", "llama-2-chat"],
   "stdin": true,
   "tty": true,
   "env": [
       "name": "WASMEDGE_PLUGIN_PATH",
        "value": "/opt/containerd/lib"
      },
        "name": "WASMEDGE_WASINN_PRELOAD",
       "value": "default:GGML:CPU:/resource/llama-2-7b-chat.Q5_K_M.gguf"
   "volumeMounts": [
        "name": "plugins",
        "mountPath": "/opt/containerd/lib"
        "name": "model",
        "mountPath": "/resource"
```









```
"volumes": [
        "name": "plugins",
        "hostPath": {
          "path": "'"$HOME"'/.wasmedge/plugin/"
        "name": "model",
        "hostPath": {
          "path": "'"$PWD"'"
```









China 2024

Output: https://github.com/second-state/wasmedge-containers-contain

Challenges and Future Directions





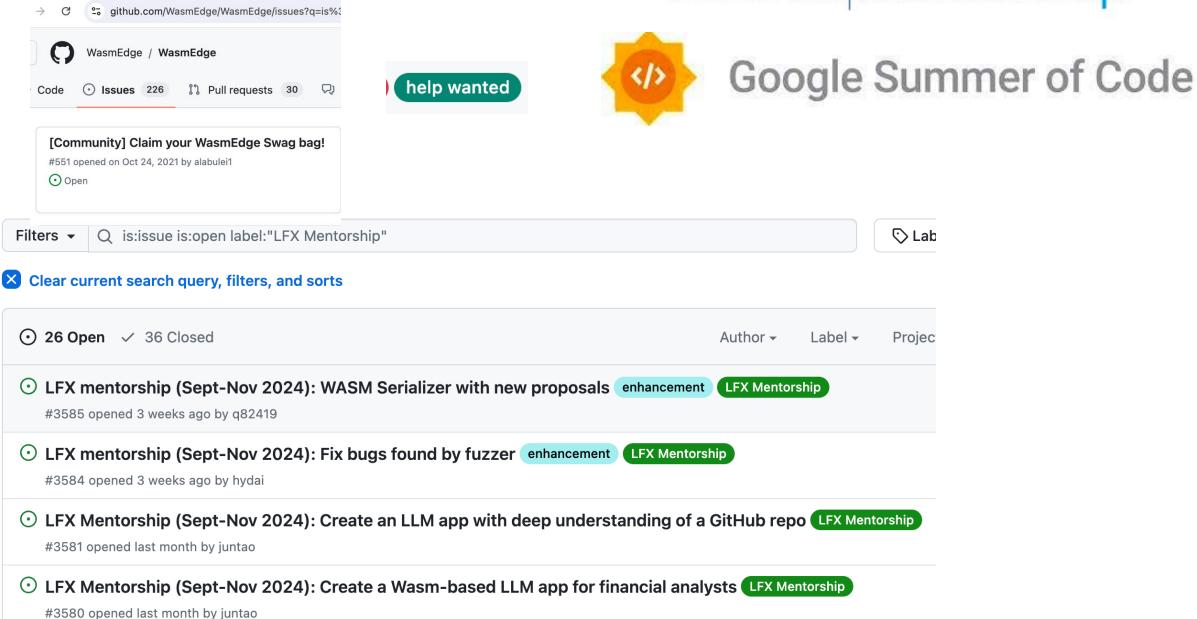




- Overview: Challenges
- Key Points:
 - Yong ecosystem
 - Low learning curve
 - More diverse infra to be supported

Calling contributors







WasmEdge

https://github.com/WasmEdge/WasmEdge

LlamaEdge

https://llamaedge.com/

https://github.com/LlamaEdge/LlamaEdge

GaiaNet: The RAG API server and node

https://github.com/GaiaNet-Al





YouTube





Bilibili



LlamaEdge Github

