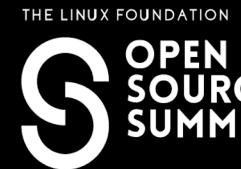




KubeCon



CloudNativeCon



China 2024

Running LLMs in the Cloud

Miley Fu, WasmEdge

GitHub / Twitter: @mileyfu

<https://github.com/WasmEdge/WasmEdge>

<https://github.com/LlamaEdge/LlamaEdge>

My Calling in Open Source Technology



KubeCon



CloudNativeCon



THE LINUX FOUNDATION

OPEN SOURCE SUMMIT



Open Source Dev & ML Summit

China 2024



WebAssembly and Rust Meetup (Wasm Empowering AI)

★★★★★ (14) ?

Beijing, China

390 members · Public group i

Organized by Miley Fu and 2 others

Embed LLM into your container app



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



China 2024

```
docker run --rm -p 8080:8080 --name api-server secondstate/llama-3-8b-nomic-1.5:latest ctx size=4096
```

<https://hub.docker.com/repository/docker/secondstate/llama-3-8b-nomic-1.5/general>

<https://github.com/LlamaEdge/LlamaEdge/tree/main/docker>



Update Warp



IDATION

PEN
SOURCE
JMMIT

AI_dev
Open Source Dev & ML Summit

base ~



Key features



KubeCon



CloudNativeCon



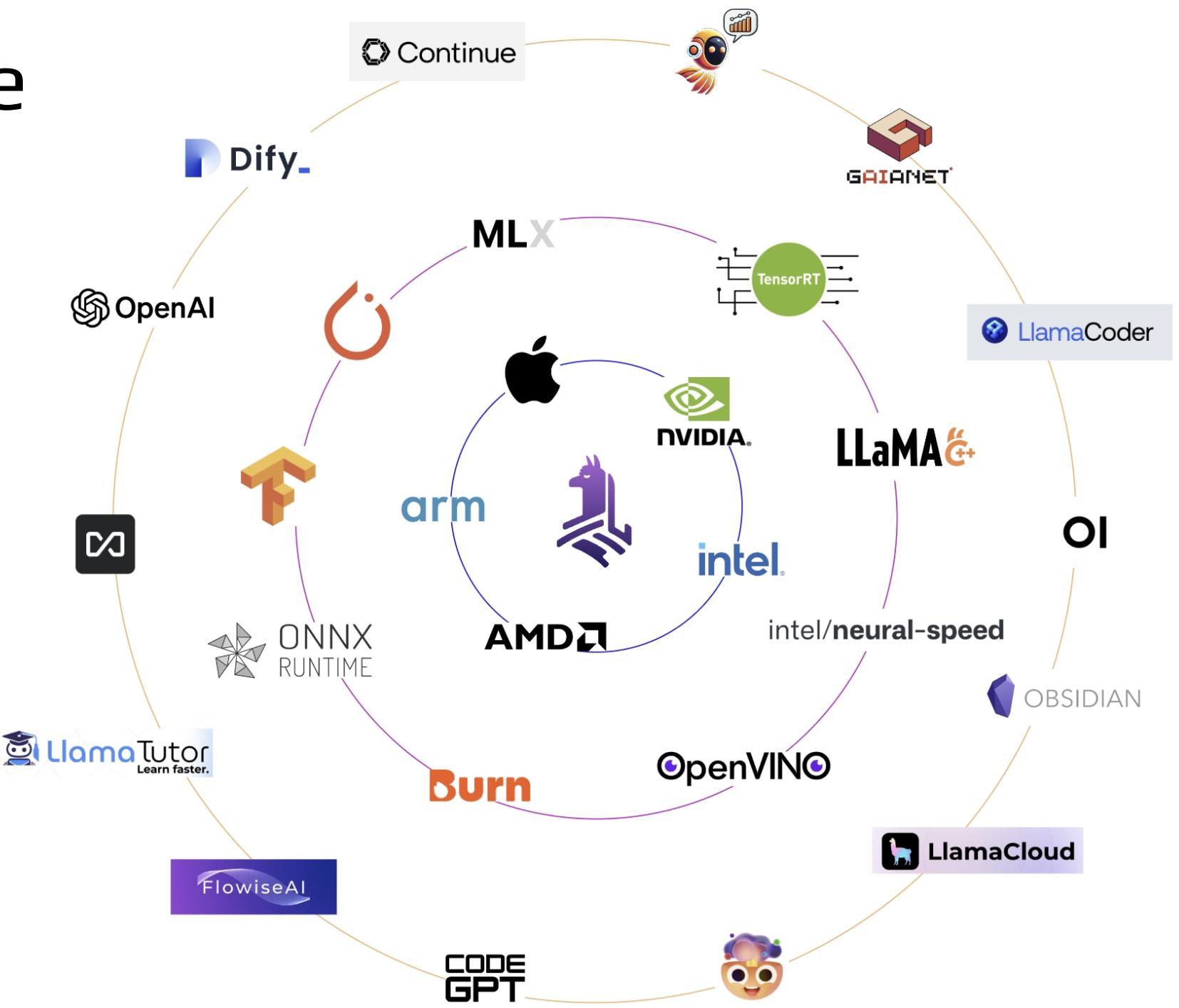
China 2024



AI_dev
Open Source Dev & ML Summit

- Tightly coupled LLM and application
 - Matches prompts, quantization & runtime with the exact version of LLM
 - The container app always works regardless LLM upgrade cycles
- Lightweight
 - Only 5GB as opposed to 10GB PyTorch app
- Portable
 - The same binary app inside the container works on multiple CPUs and GPUs
 - Develop on Mac and deploy on Nvidia
- Easy to embed into Rust / JS / Python apps
- Works with existing container tools, such as K8s

LlamaEdge



Real-world use cases



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



Open Source Dev & ML Summit
AI dev

China 2024

- Personal LLMs:** Gaia Network, users run personal LLMs with embedded knowledge base.
- AI OS:** Open interpreter as a local LLM provider; 51.5k stars
- Finance:** Financial analytics bot.
- Hardware:** Robot voice control
- Education:** UC Berkley TA
- Game:** Open-source game engine Cocos AI , use Wasm to run AI models that enhance gameplay experiences enabled by NPC.



Nodes: **28,469**

Throughputs: **2884.2M**

LlamaEdge: easy LLM deployment + inference



KubeCon



CloudNativeCon



- Single cross-platform binary
(Automagically take advantage of local hardware accelerators)
- Compile and test apps on one machine (e.g., Mac) and deploy it to another cloud server (e.g., Nvidia CUDA 12)
- The app can be moved around and deploy to new hardware by K8s
- Package the Wasm app into Docker image as an embedded AI / LLM service. Only 1/100 the size of a Python runtime.



LlamaEdge Github

WasmEdge

<https://github.com/WasmEdge/WasmEdge>

LlamaEdge

<https://github.com/LlamaEdge/LlamaEdge>

Gaia: The living knowledge server

<https://github.com/GaiaNet-AI>

Cloud-Native AI Ecosystem



China 2024

ters

GROUP: Projects and products Members Certified partners and providers Serverless Wasm CNAI

VIEW MODE: Grid Card

ZOOM: - +

Data Architecture

Model/LLM Observability

Governance, Policy & Security

General Orchestration

ML Serving

AutoML

Distributed Training

CI/CD - Delivery

Vector Databases

Data Science

Workload Observability

Icons displayed in the grid:

- ClickHouse
- druid
- cassandra
- SCYLLA
- Apache Base
- presto
- trino
- Apache Spark
- Flink
- kafka
- PULSAR
- Redis
- trulens
- deepchecks
- DeepFlow
- Alluxio
- FEAST
- OpenLIT
- VOLCANO
- ARMADA
- RAY
- NVIDIA RAY
- Horizon
- TensorFlow
- Kubernetes CNCF GRADUATED
- Hyperopt
- Katib
- PyTorch
- PyTorch
- Torch
- deepspeed
- Megatron-LM
- Hugging Face
- Alpa
- Kserve
- LLM
- Wandb
- Jupyter
- PyTorch
- XGBoost
- Apache Zeppelin
- Milvus
- Cardano
- W&B

<https://landscape.cncf.io/?group=cnai>

How your company uses GenAI tools

- Different use cases
- Modalities/models in use
- Challenges to adoption
- The role of open source in adoption decisions



Calling for contributors

→ G github.com/WasmEdge/WasmEdge/issues?q=is%3Aopen+is%3Aissue+label%3ALFX+Mentorship

WasmEdge / WasmEdge

Code Issues 226 Pull requests 30

[Community] Claim your WasmEdge Swag bag!
#551 opened on Oct 24, 2021 by alabulei
Open

Filters ▾ is:issue is:open label:"LFX Mentorship" Lab

[X](#) Clear current search query, filters, and sorts

- ① 26 Open ✓ 36 Closed Author ▾ Label ▾ Projec
- ② LFX mentorship (Sept-Nov 2024): WASM Serializer with new proposals enhancement LFX Mentorship
#3585 opened 3 weeks ago by q82419
- ③ LFX mentorship (Sept-Nov 2024): Fix bugs found by fuzzer enhancement LFX Mentorship
#3584 opened 3 weeks ago by hyundai
- ④ LFX Mentorship (Sept-Nov 2024): Create an LLM app with deep understanding of a GitHub repo LFX Mentorship
#3581 opened last month by juntao
- ⑤ LFX Mentorship (Sept-Nov 2024): Create a Wasm-based LLM app for financial analysts LFX Mentorship
#3580 opened last month by juntao



Google Summer of Code



YouTube



Second State WeChat



WasmEdge Github

Stay in Touch!
Twitter/ Github
@mileyfu