

All Things Open

# Open-Source powered and Developer-First AI Stack



Yann Léger @yann\_eu / Co-Founder @gokoyeb / [koyeb.com](https://koyeb.com)



Greg Wallace / [netactuate.com](https://netactuate.com)



Yann Leger

Co-Founder & CEO at Koyeb

@yann\_eu

**14 YEARS BUILDING CLOUD  
INFRASTRUCTURE**



Greg Wallace

Partnerships at NetActuate

in/gtewallace

**15 YEARS IN OPEN SOURCE**

# Agenda

● llama.py ● Dockerfile

```
1  from tt.llama import TtTransformer
2  from tt.llama import LlamaEmbedding
3  from tt.llama import Tokenizer
4  from tt.llama import TtModelArgs
5
6  def preprocess_inputs_prefill(
7      input_prompts,
8      tokenizer,
9      model_args,
10     instruct
11     max_generated_tokens,
12     max_prefill_len = 128 * 1024,
13 )
14
15     if max_prefill_len == 128 * 1024:
16         max_prefill_len = 128 * 1024 - ma
```

Industry context
The role of open source
Demo
Evolving infrastructure
Q&A

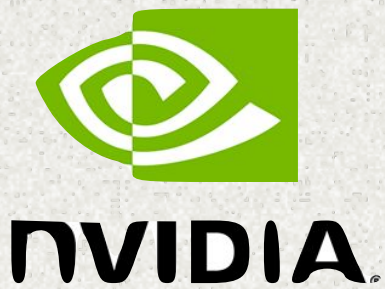
**\$1.0T by 2029 in global  
data center capex as AI  
accelerators drive spend**  
*Dell'Oro Group*

# MIT report: 95% of generative AI pilots at companies are failing

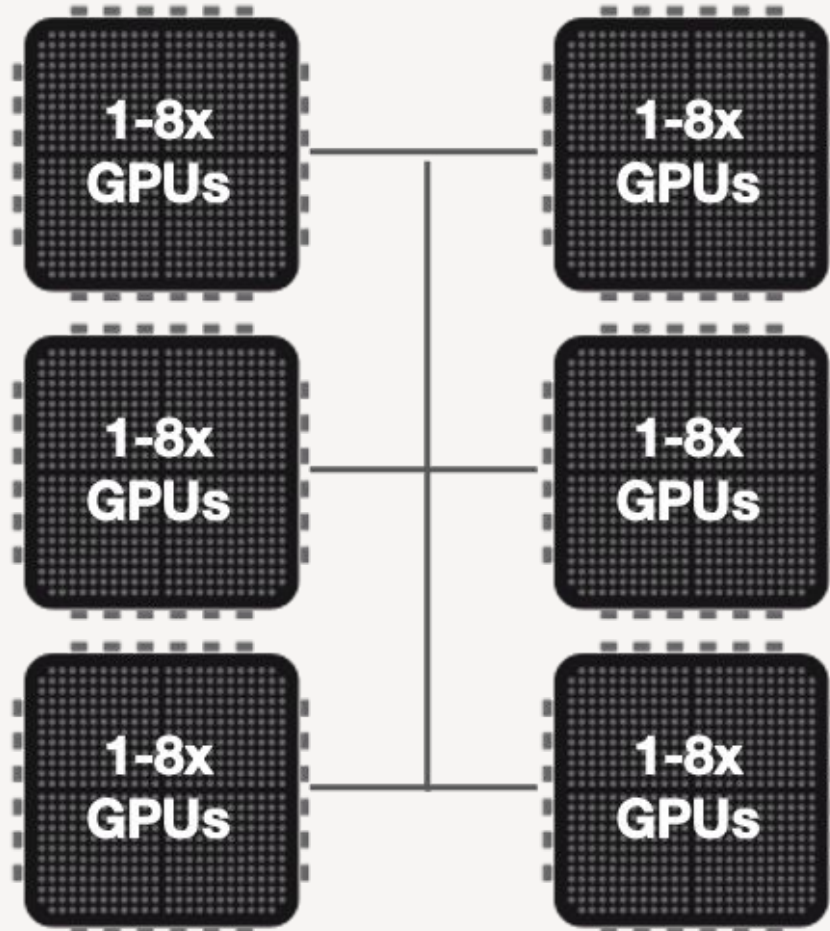
**Most organizations miss  
AI cost forecasts, with  
nearly a quarter busting  
their budgets by more  
than 50% (CIO)**

# There's room for improvement

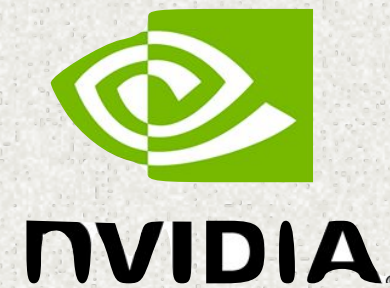
# What is AI Infrastructure?



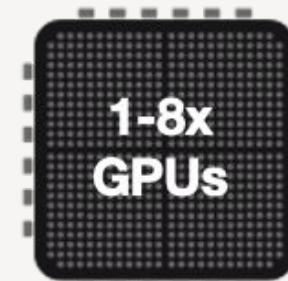
## Large Training



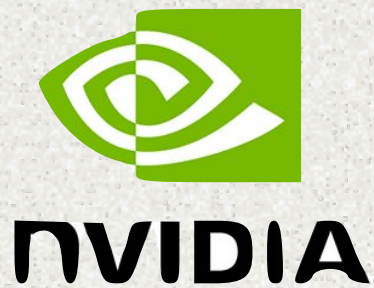
Up to 100K GPUs clusters  
with high-speed  
interconnect



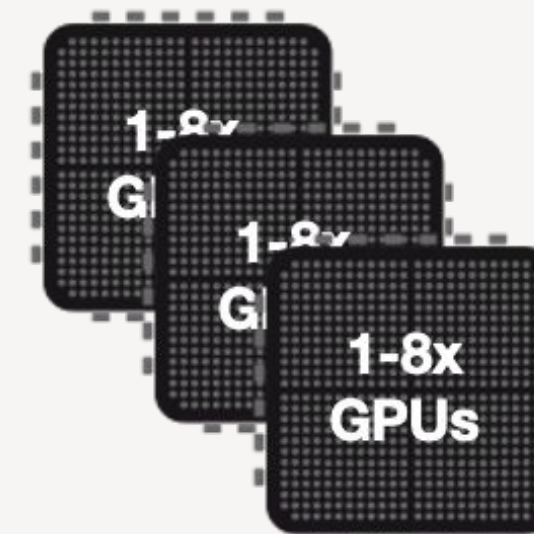
## Fine-Tuning & Small Training



1x to 8x GPUs or  
Small Clusters



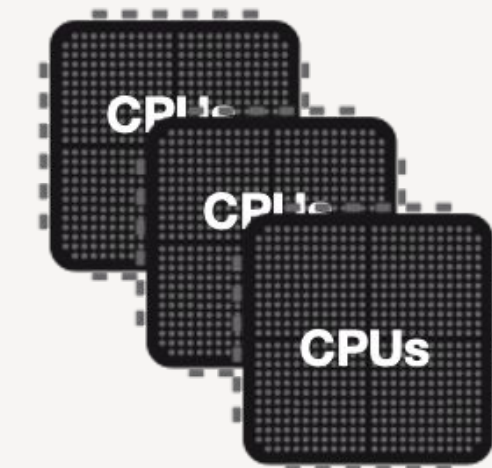
## Inference



Fits one machine, needs  
more machines to handle  
parallel requests



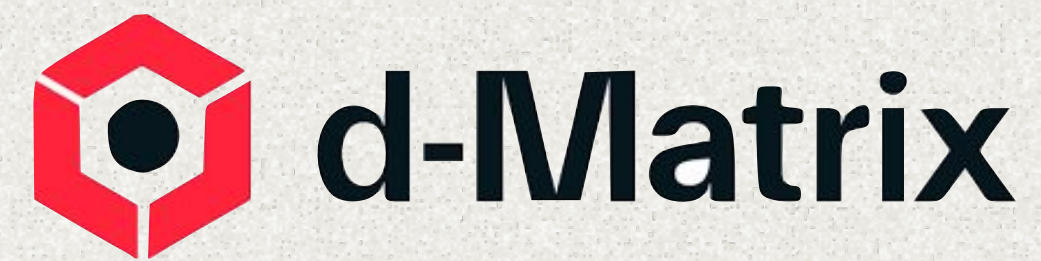
## Agents



High-volume, lightweight  
sandboxed code  
execution

# Maybe we need more choice?

# New AI Accelerators are Emerging

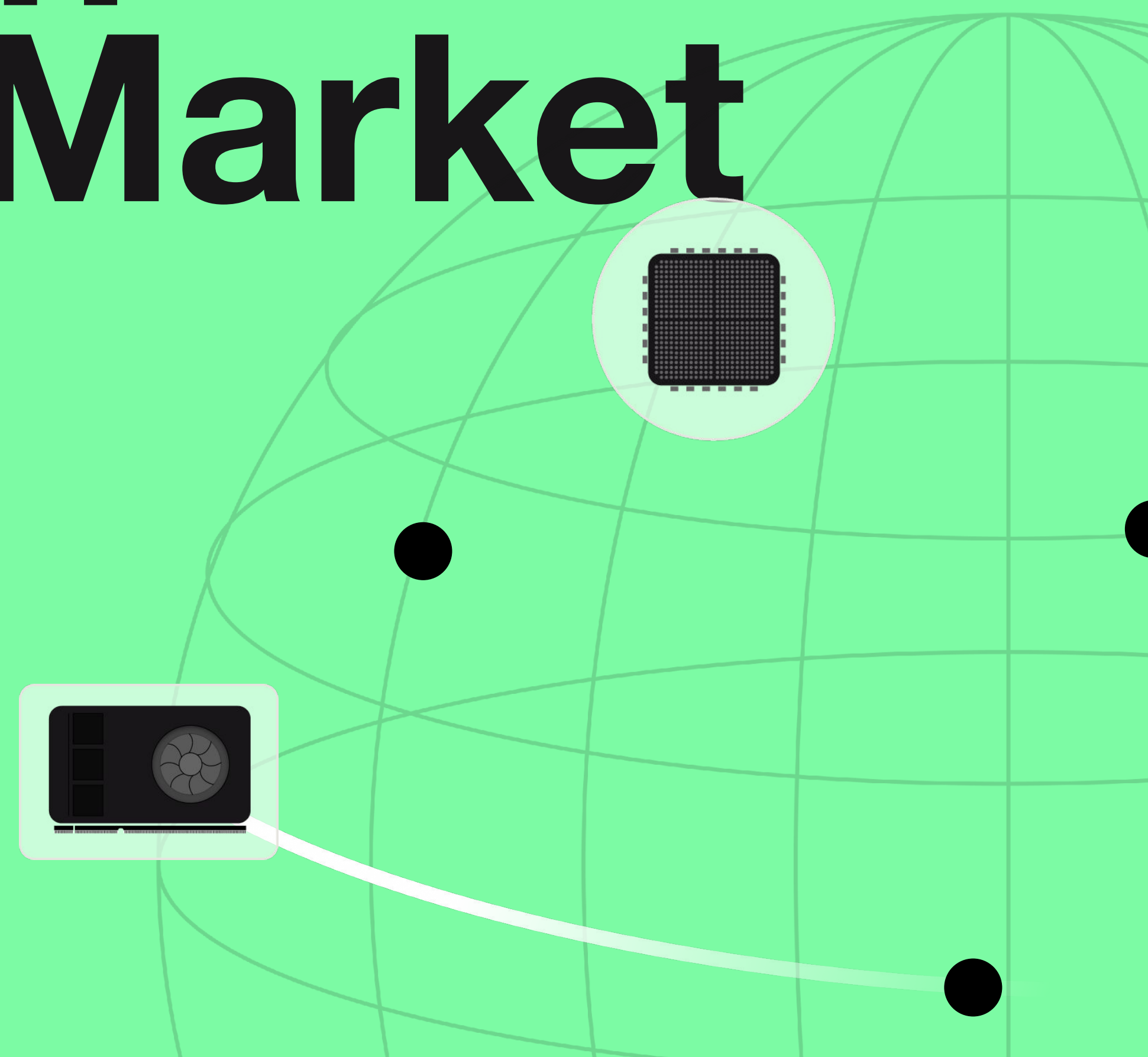


Building for the Agentic Era

# Introducing Open Accelerators to Market



Yann Léger @yann\_eu / Co-Founder @gokoyeb / koyeb.com



# We build a Global Serverless Platform for Agents and Inference



# We Run Models and Agents



**Model & Agents continuous deployment with Containers**



**Zero infrastructure management**



**Low latency with global deployments**



**High-performance hardware**

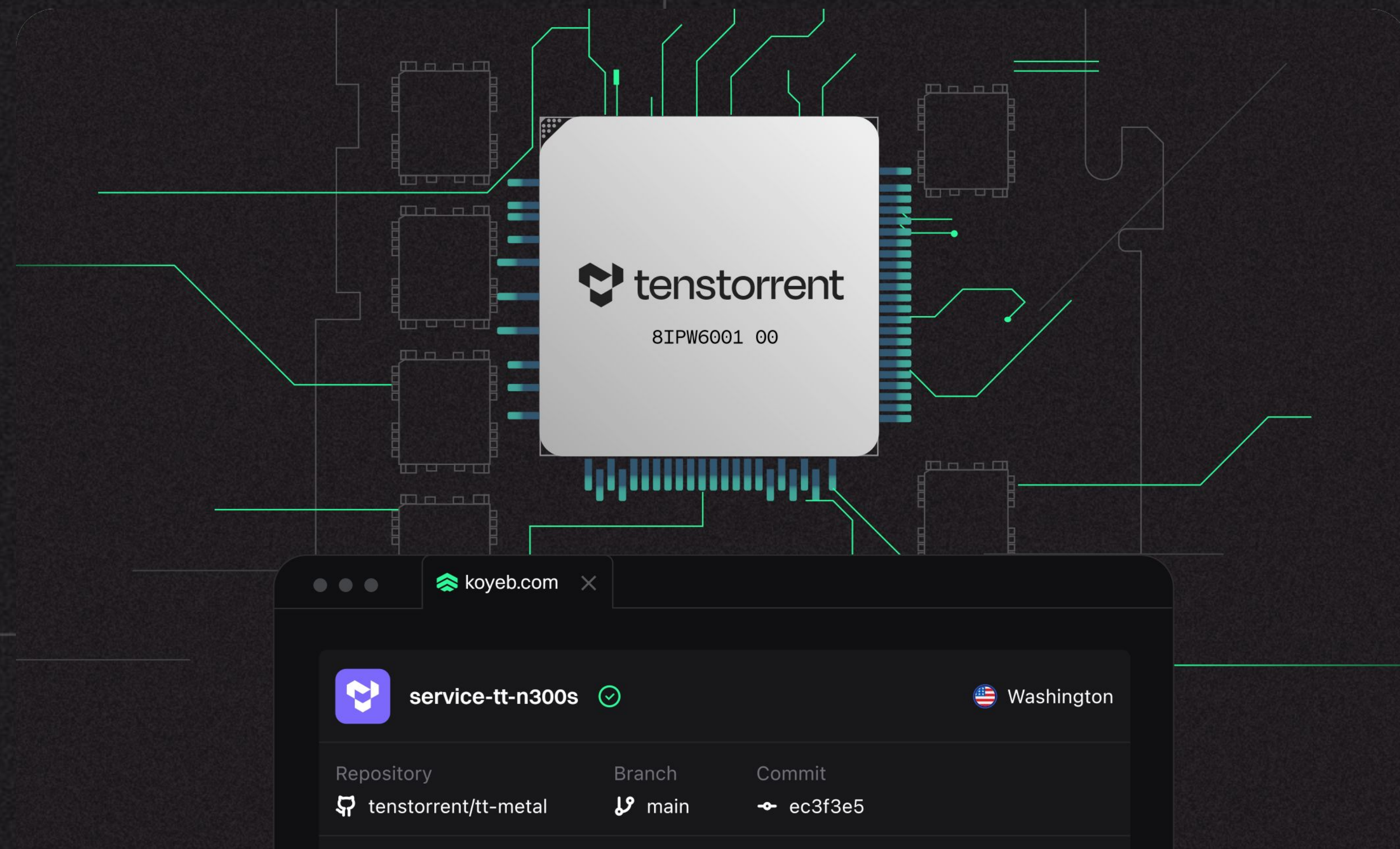


**Sustainable costs with Serverless capabilities**

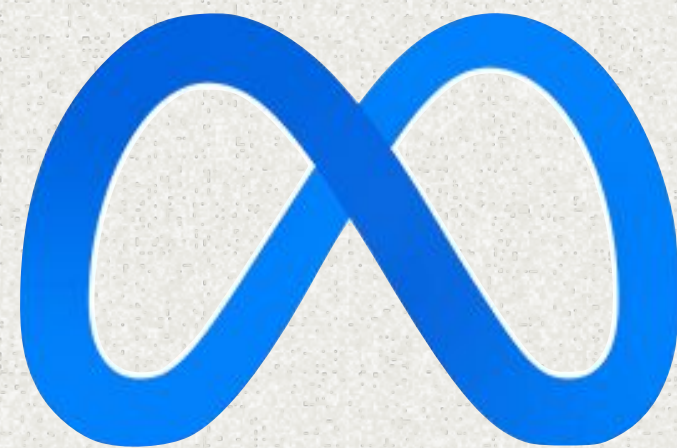


**Native resiliency and scalability**

# Available in Private Preview only on Koyeb



# Deploy LLAMA 3.2 11B Vision On Tenstorrent



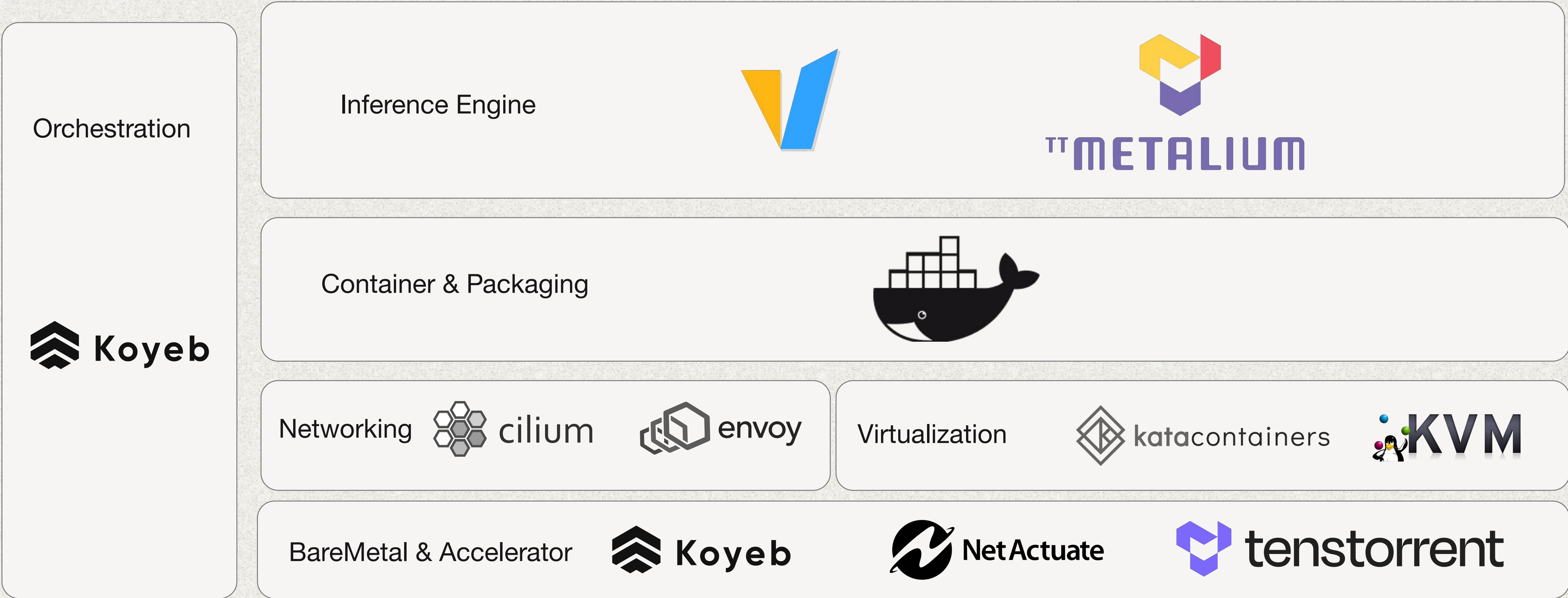
[github.com/koyeb/tenstorrent-examples](https://github.com/koyeb/tenstorrent-examples)

# Unified and global experience across CPUS, GPUS, and Accelerators



THE STACK

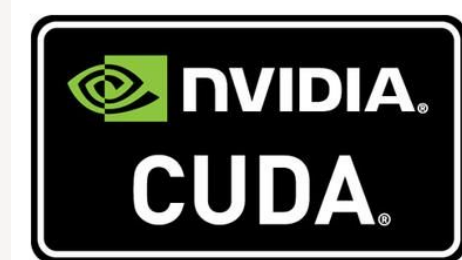
# CONTAINERS IN VMS ON BARE METAL



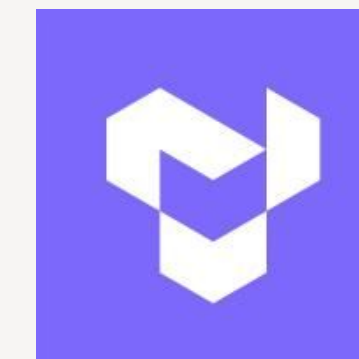
# Software is key to Performance and Efficiency

# New Accelerators Need New AI Stacks

**Proprietary**



**Open source on  
open hardware**



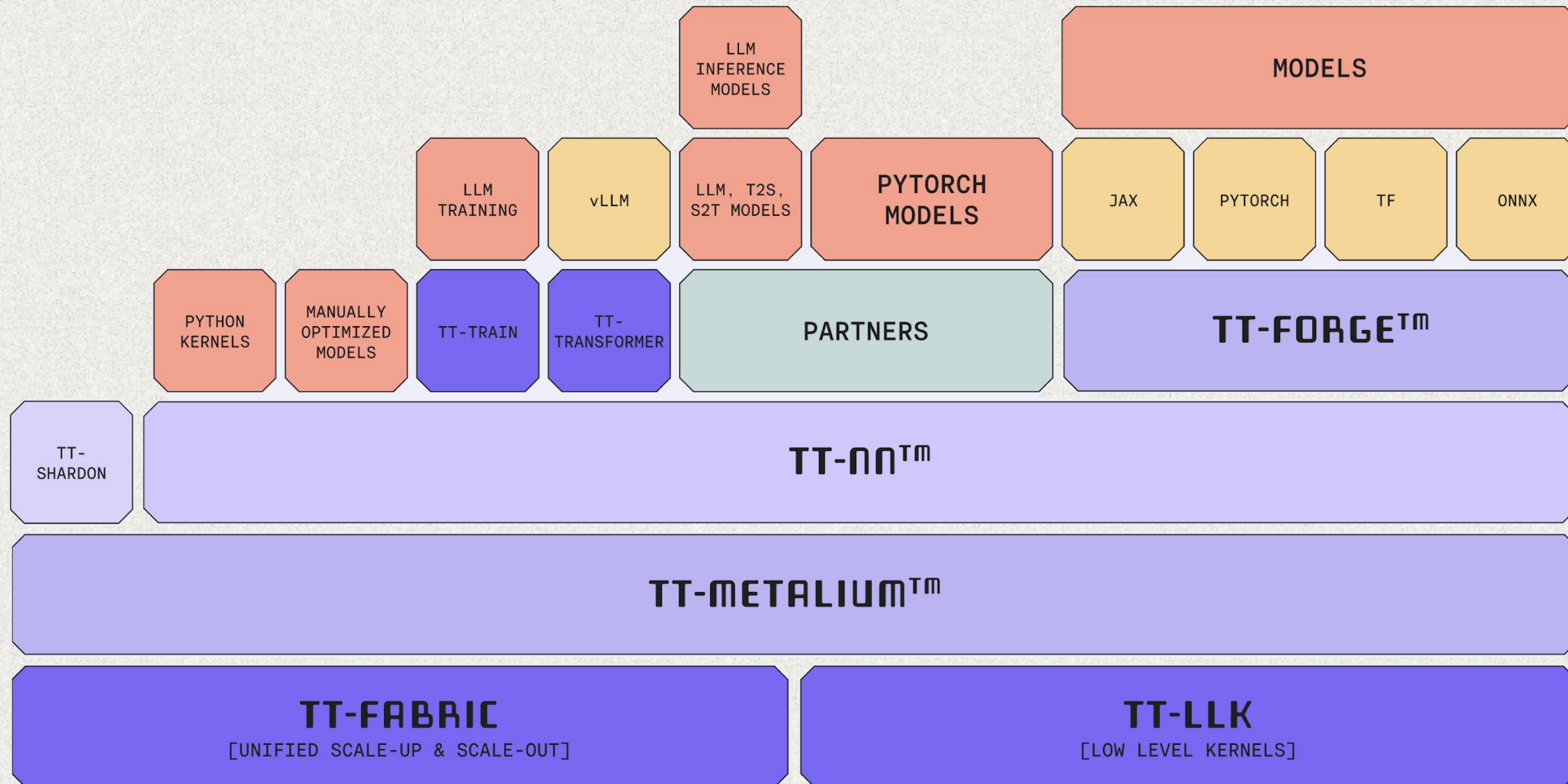
**Open source on  
closed hardware**



**Open standard**



# Tenstorrent provides a fully open stack



# How can we build?

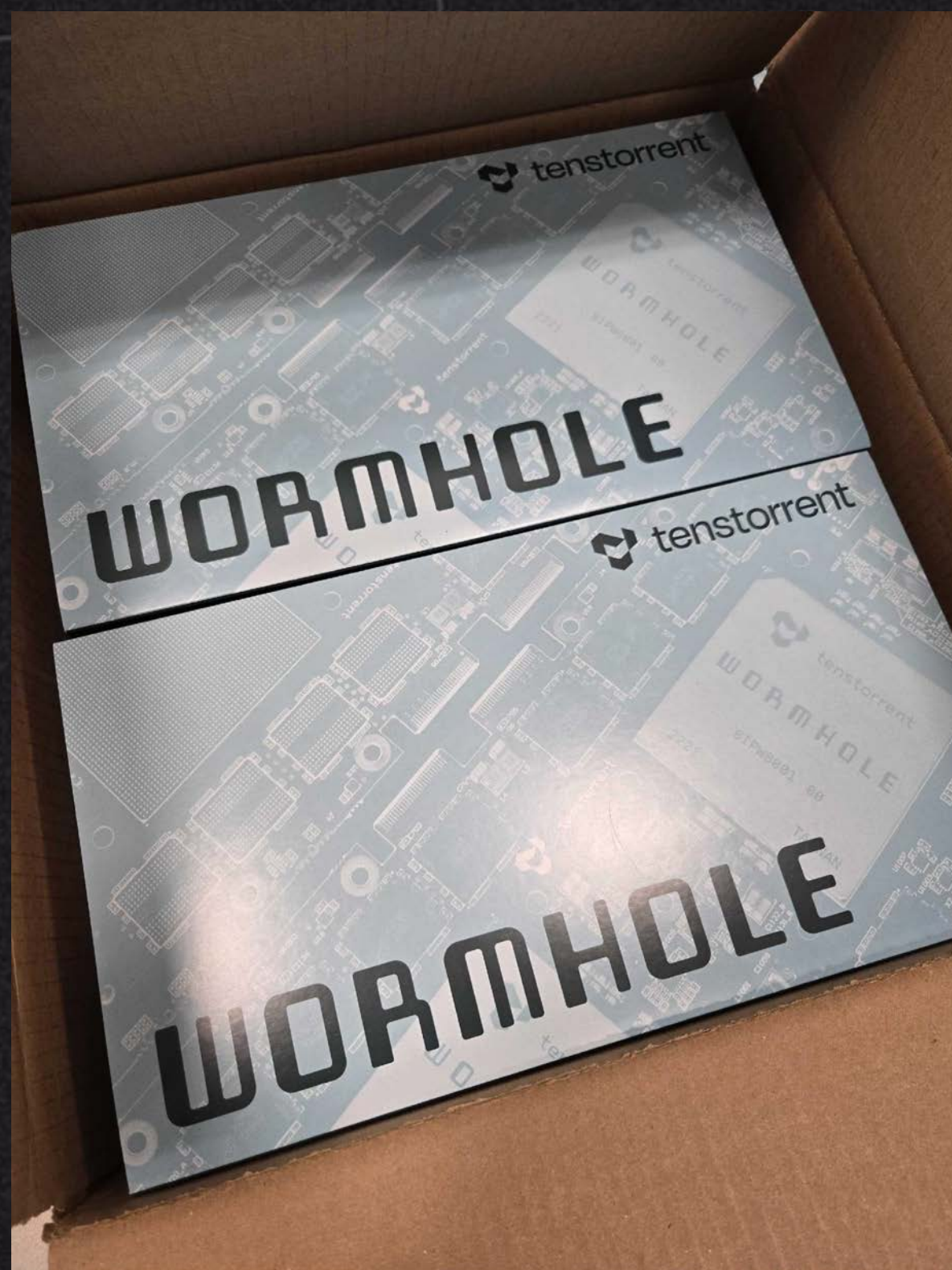


**TT METALIUM**

Instantly access development environments  
with Tenstorrent open-source SDKs, including  
TT-Metalium on n300s or TT-Loudbox

[github.com/koyeb/tenstorrent-examples](https://github.com/koyeb/tenstorrent-examples)

# Powered by Open Hardware





# NetActuate Global Footprint



# GPU hardware is vastly underutilized

# Use Serverless

Automatically adapt infrastructure to demand

→ SCALE-TO-ZERO

→ AUTOSCALING

# Deploy on Tenstorrent in seconds



Thank you.

# GET STARTED WITH KOYEB

## Chat with our Team



Yann Leger

[yann@koyeb.com](mailto:yann@koyeb.com)

/in/yannleger

@gokoyeb



\$200 on us