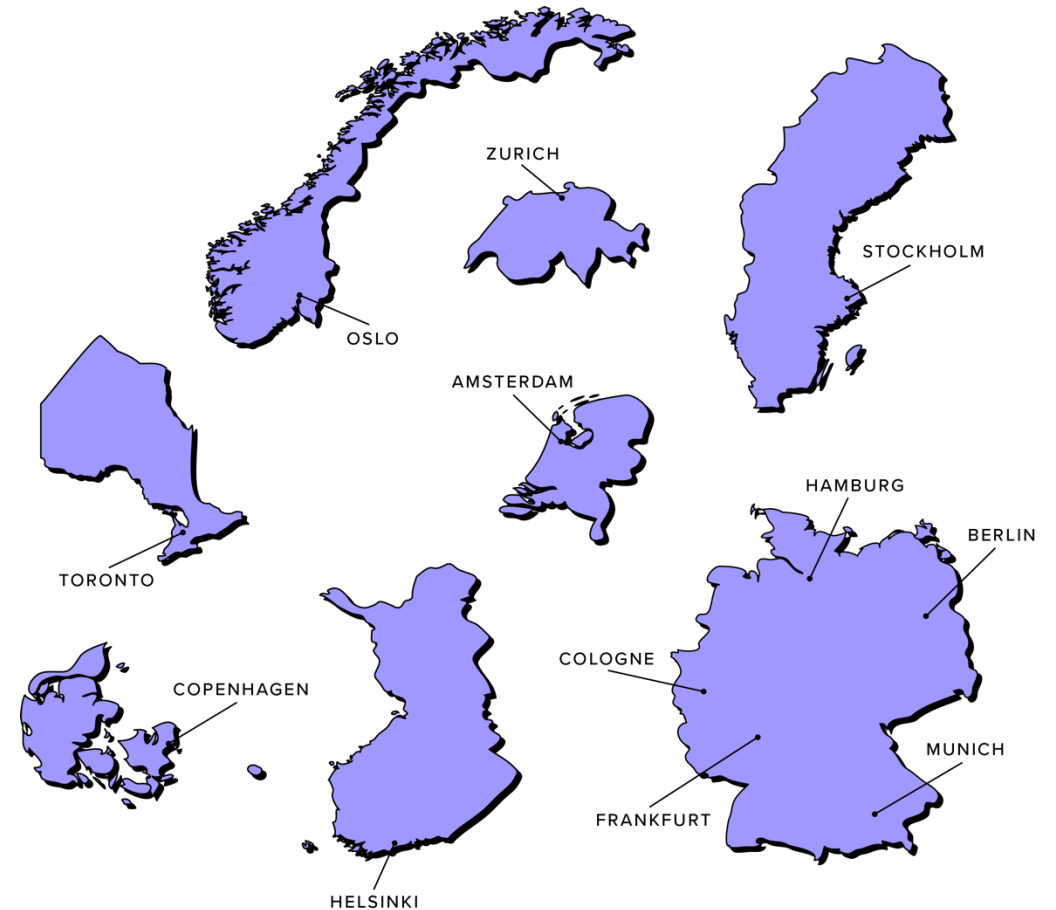# Own your AI

Tim Fischer

# About Netlight

Netlight is an international digital consulting firm, helping leading companies to succeed in the digital landscape, from advice to implementation.

Our service contains the collective intelligence of 2000 consultants offering a comprehensive range of digital services, from strategy to technology. We support industries that are facing new challenges and opportunities based on new technology, to make better business.

Netlight has been awarded several times for profitable growth and management, as a top employer, and for our engagement in Diversity, Equity, and Inclusion. Located in Stockholm, Oslo, Helsinki, Copenhagen, Munich, Hamburg, Berlin, Frankfurt, Zurich, Cologne, Amsterdam and Toronto. Co-creating the future today, since 1999.



**aws** PARTNER Advanced Tier Services  **aws** PARTNER Data & Analytics Services Competency  **aws** PARTNER AWS Glue Delivery  **aws** PARTNER AWS Lambda Delivery

Google Cloud Partner    Microsoft Silver Partner

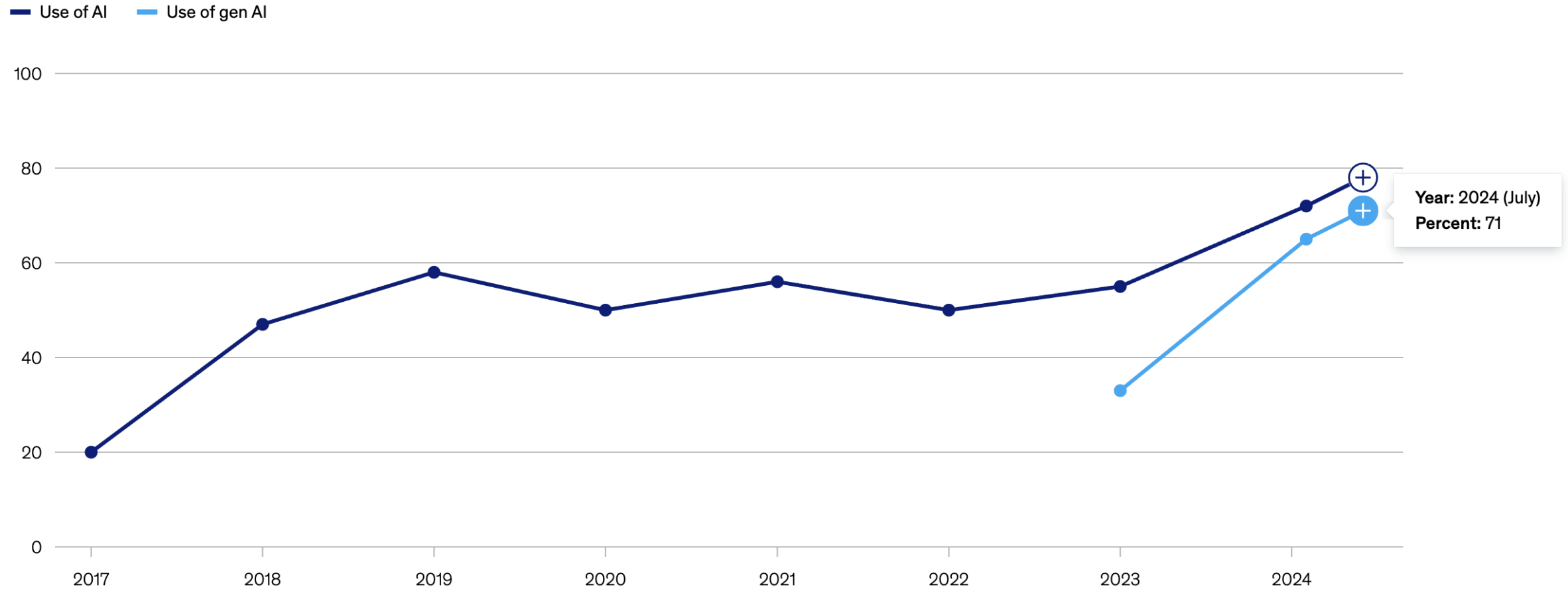| OFFICES | 12 | EMPLOYEES | 2000+ | CLIENTS | 350+ | FOUNDED | 1999 | WOMEN AT NETLIGHT | 37% |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

netlight

# Agenda

Why?

How? (Architectural Setup)

What I wish I would have known before

Key Take-aways

Q&A

netlight

# Organizations' use of AI has accelerated markedly in the past year, after years of little meaningful change.

**Organizations that use AI in at least 1 business function,**[1] % of respondents

— **Use of AI**    — **Use of gen AI**



> **Year:** 2024 (July)
> **Percent:** 71

netlight

# What's the problem?

- **Data Privacy & Control**: Sensitive data may be exposed, used without consent, or fail to meet compliance requirements like GDPR or HIPAA.

- **Vendor Dependency**: Risk of outages, AI model changes, API updates when locked-in with a single provider.

- **Security Risks**: Limited control over encryption, and reliance on the provider's security practices increases breach risks.

- **Customization Limits**: Inability to access underlying mechanics for transparency and optimization.

- **Cost & Scalability:** Pricing changes are unpredictable, leaving you vulnerable to cost increases that can disrupt budgets

**European tech industry coalition calls for 'radical action' on digital sovereignty — starting with buying local**

**Natasha Lomas**
Techcrunch, March 16, 2025

**Is overreliance on US Big Tech a threat to Europe? The Netherlands may soon find out**

**Anna Desmarais**
Euronews, Feb 27, 2025

**Big Tech's future in Europe at risk from Trump, tariffs and data concerns**

**Eugenia Perozo**
Investment Monitor, 31 March, 2025

**Companies in the EU are starting to look for ways to ditch Amazon, Google, and Microsoft cloud services amid fears of rising security risks from the US.**

**Matt Burgess**
WIRED, Mar 24, 2025

**Legally compliant use of US cloud services: BDI warns against end of agreement**

**Martin Holland**
heise, Mar 31, 2025

**Will the US cloud soon be illegal in the EU?**
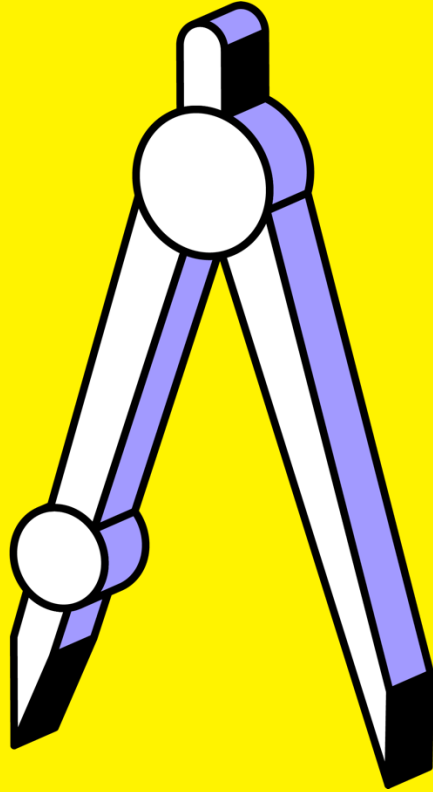
**by Tristan Fincken**
CIO, Jan 27, 2025

netlight

# Tim Fischer

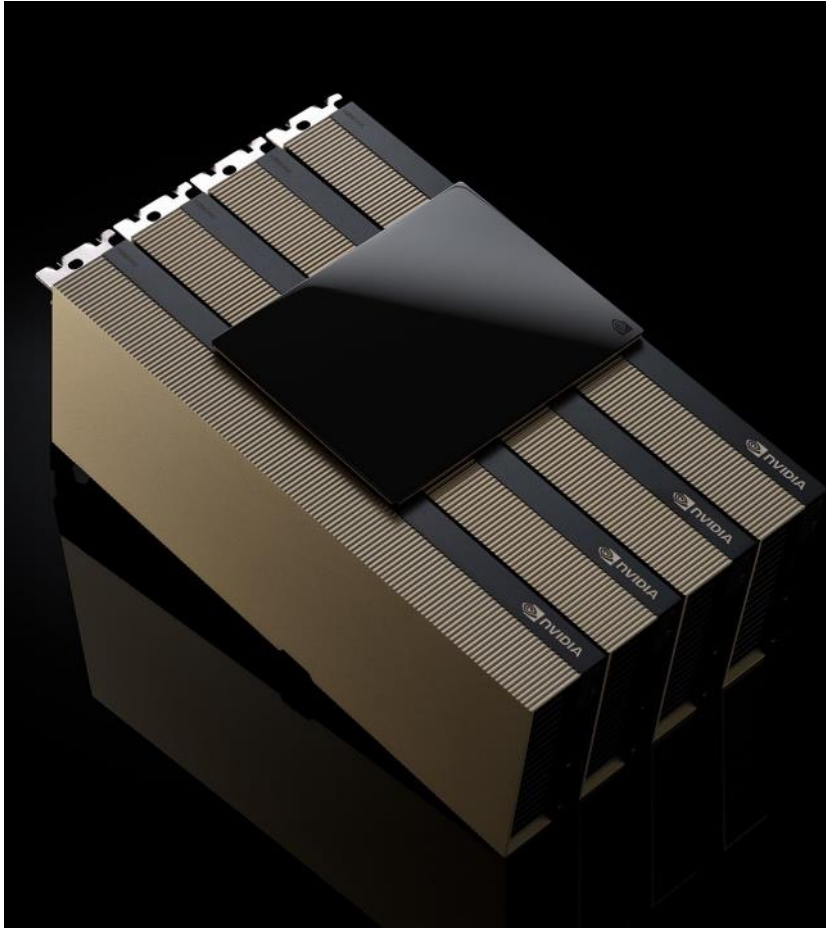**Consultant @ Netlight**

Cloud, Infrstructure, Kubernetes, GenAI

# Project Challenge

# Create a GenAI Platfom that...

- ensures the highest levels of compliance and security

- minimizes lock-in effects

- provides full control over infrastructure, data, and workflows
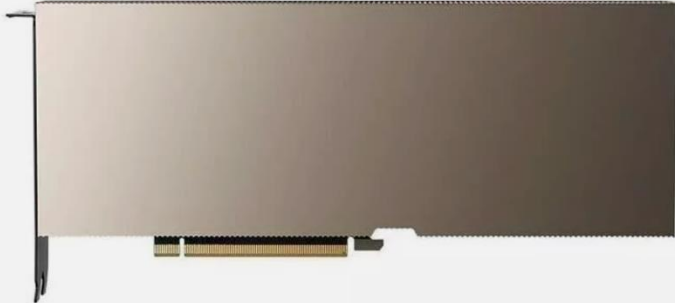
- is scalable, resilient, portable and cost-efficient

netlight

# We need GPUs



**Nvidias newest GPU Generation "Blackwell" is reportedly <u>sold out</u> for the next 12 months**

**Matthew Fox**
<u>Business Insider</u>, Oct 11, 2024



32 AUFRUFE IN DEN LETZTEN 24 STUNDEN

NVIDIA H200 Tensor Core 141GB NVL PCIe Dual-slot GPU Deep Learning Computing AI
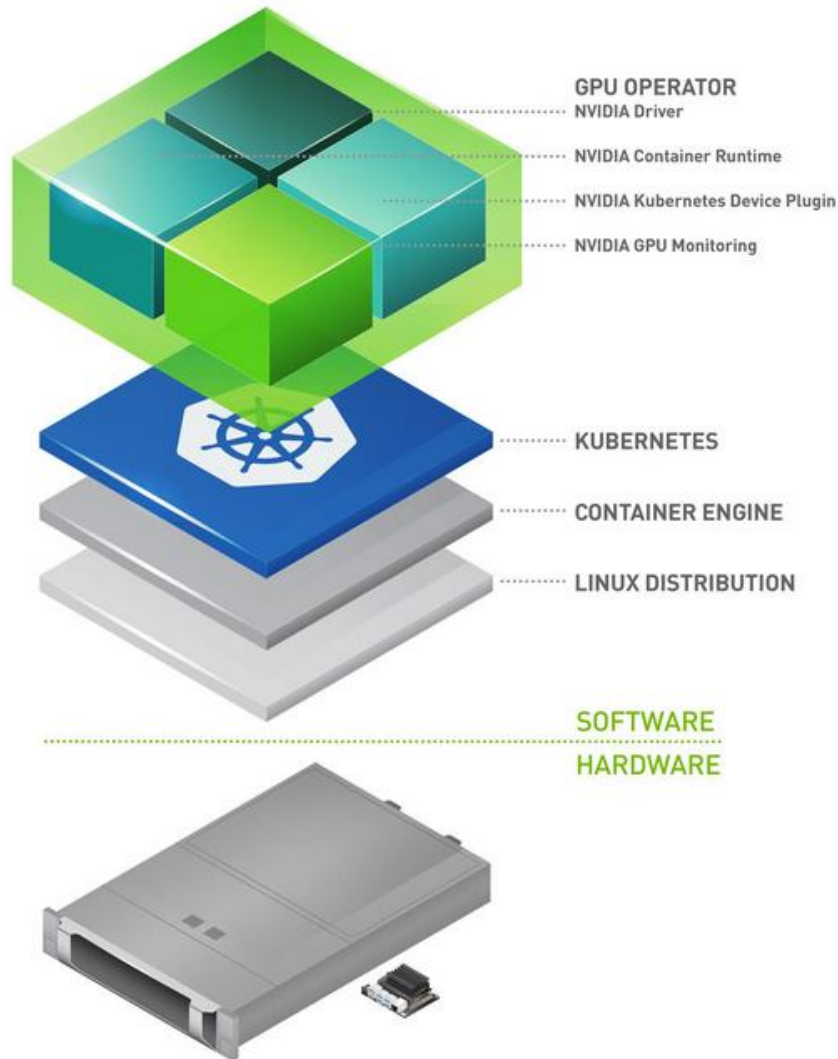
diskclubs (17033) · Gewerblich
<u>99,1% positive Bewertungen</u> · <u>Mehr Artikel des Verkäufers</u> · <u>Verkäufer kontaktieren</u>

**EUR 36.999,00**
oder Preisvorschlag

Artikelzustand: **Neu: Sonstige (siehe Artikelbeschreibung)**

**Sofort-Kaufen**

In den Warenkorb

netlight

GPU OPERATOR
NVIDIA Driver
NVIDIA Container Runtime
NVIDIA Kubernetes Device Plugin
NVIDIA GPU Monitoring

KUBERNETES

CONTAINER ENGINE

LINUX DISTRIBUTION

SOFTWARE
HARDWARE

# kubernetes

- **Cloud-Agnostic and GPU-Ready using GPU Operator**: Kubernetes provides portability across on-premises and cloud environments, enabling efficient use of GPU resources from any infrastructure provider.

- **Flexible and Open Platform**: As an open, free, and extensible SDK, Kubernetes empowers developers to build and deploy at scale without vendor lock-in.

- **Scalable and Reliable**: Kubernetes automates scaling, load balancing, and fault tolerance, making it ideal for serving AI/ML models with high availability and performance.

- **Ecosystem and Innovation**: A huge effort by the community (e.g., SIG Serving) positions Kubernetes as the AI serving layer, ensuring rapid innovation.

netlight

# Setup

## 80% of companies use Kubernetes in production
*CNCF annual survey 2024*

**Many companies already have platform teams for managing and provisioning the clusters**



**KUBERNETES USAGE**

Does your organization use Kubernetes? (select one)

Yes, using in production — 80% (2024), 66% (2023)

Yes, piloting or actively evaluating in test environments — 13% (2024), 19% (2023)

No — 7% (2024), 15% (2023)

2024 ■ 2023 ■

- **Cilium**: Provides deep network security with eBPF-based policies, identity-aware security, encryption, and transparent visibility into L3-L7 traffic flows.
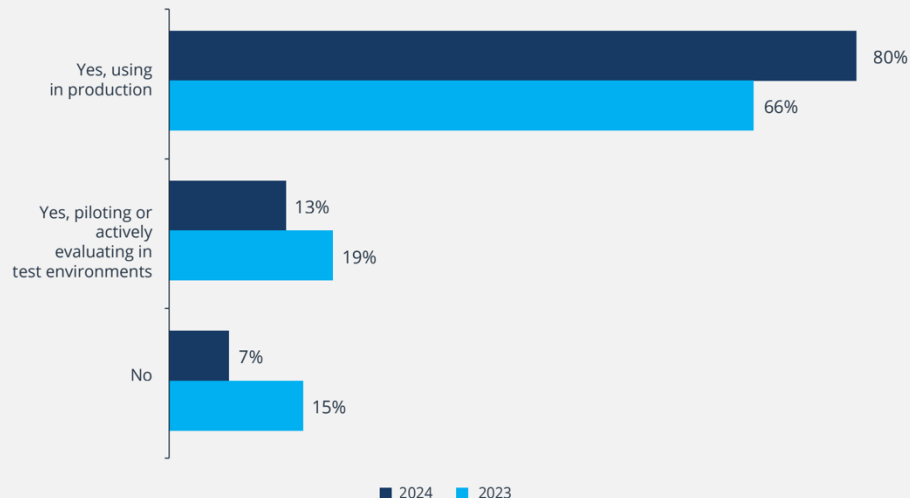
- **Kyverno:** Strong policy enforcement tool with Kubernetes-native approach, covers validating, mutating, and generating resources based on security policies.

- **ArgoCD:** GitOps principles that provide complete audit trails, version control, and approval workflows while automatically detecting and remediating configuration drift. Enables advanced deployment strategies for smooth rollouts

netlight

# Model selection

**Selection of state-of-the-art Open-Weight-models**



deepseek-ai/DeepSeek-R1



google/gemma-3



Qwen/Qwen2.5



meta-llama/Llama-4



mistralai/Mistral-Small-3.1

▪ **Open-Weight model**: Trained model weights are publicly available for use, modification, and fine-tuning

▪ **Closed-Weight model**: Model weights are proprietary and inaccessible, <u>usable only via APIs</u>
(e.g. OpenAI, Anthropic Claude)

**→ closed models are still slightly better, but open models are catching up fast 🚀**

netlight

# How to run it? vLLM

- vLLM is a **fast and easy-to-use library for LLM inference** and serving.
- Offers OpenAI-compatible API server
- Supports NVIDIA GPUs, AMD CPUs and GPUs, Intel CPUs and GPUs, PowerPC CPUs, TPU, and AWS Neuron.
- Downloads models directly from huggingface

```yaml
1   apiVersion: apps/v1
2   kind: Deployment
3   metadata:
4     name: vllm-demo
5   spec:
6     template:
7       spec:
8         containers:
9           - image: vllm/vllm-openai:latest
10            command:
11              [
12                "/bin/sh",
13                "-c",
14                "vllm serve
15                meta-llama/Llama-4-Scout-17B-16E-Instruct
                  --tensor-parallel-size 2
16                --max-model-len 80000",
17              ]
18            resources:
19              limits:
20                nvidia.com/gpu: 2
```
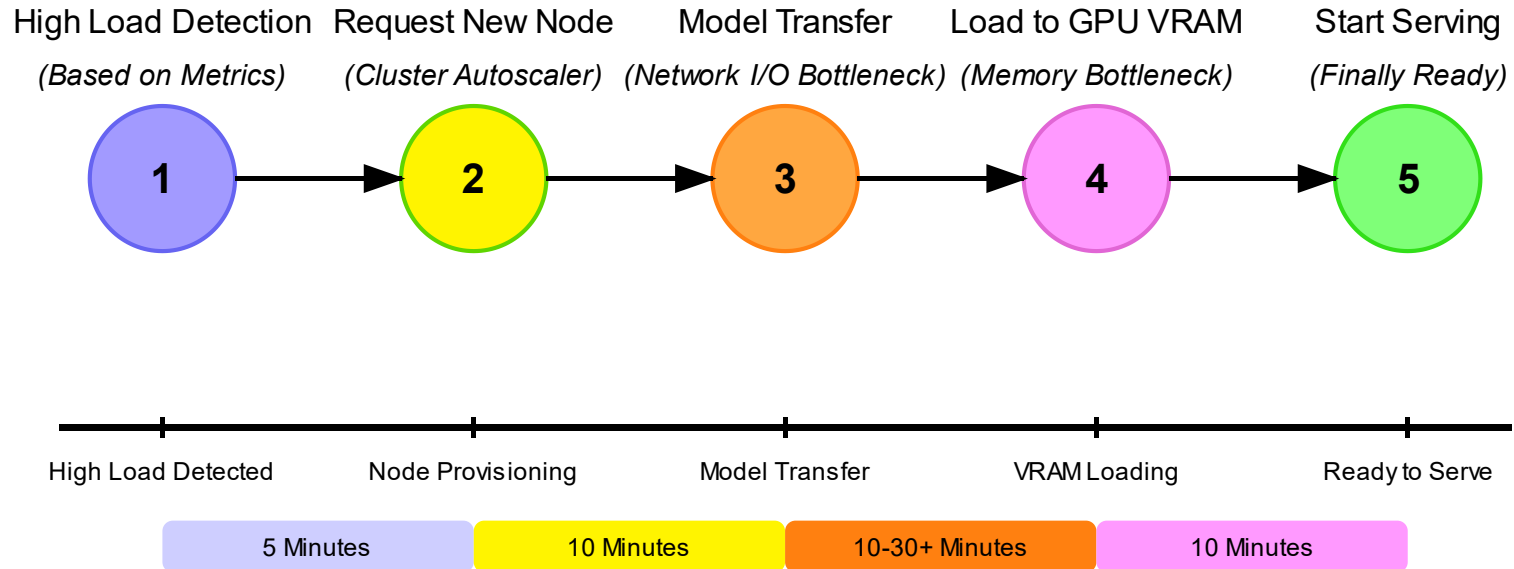
netlight

# Observability

## vLLM and Nvidia DCGM expose Prometheus metrics endpoint

# The problem with autoscaling…

Problem: Autoscaling is triggered by immediate demand...

...but new capacity takes 20-60+ minutes to become available

| High Load Detection | Request New Node | Model Transfer | Load to GPU VRAM | Start Serving |
|---|---|---|---|---|
| *(Based on Metrics)* | *(Cluster Autoscaler)* | *(Network I/O Bottleneck)* | *(Memory Bottleneck)* | *(Finally Ready)* |
| **1** | **2** | **3** | **4** | **5** |

High Load Detected — Node Provisioning — Model Transfer — VRAM Loading — Ready to Serve

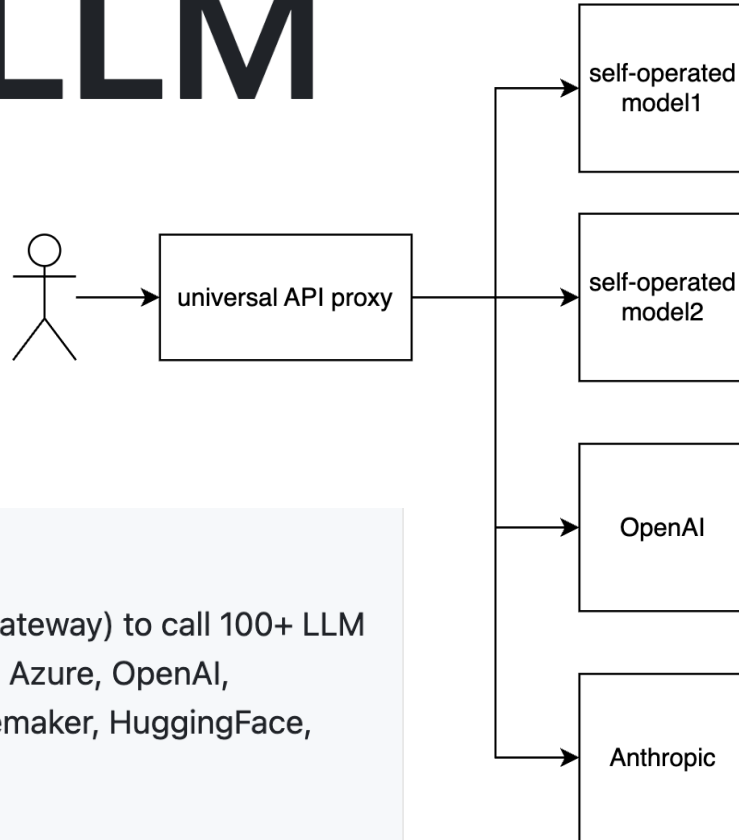| 5 Minutes | 10 Minutes | 10-30+ Minutes | 10 Minutes |
|---|---|---|---|

**Predictive Scaling Required Instead of Reactive Scaling**

netlight

# API Management

🚄 LiteLLM

```
1   curl https://llm-proxy.edgez.live/v1/models # List all models
2
3   curl https://llm-proxy.edgez.live/v1/chat/completions \
4   -d '{
5     "model": "gpt-4",
6     "messages": [{"role": "user", "content": "Explain the
      basics of GenAI"}]
7   }'
8
9   curl https://llm-proxy.edgez.live/v1/chat/completions \
10  -d '{
11    "model": "llama-4",
12    "messages": [{"role": "user", "content": "Explain the
      basics of GenAI"}]
13  }'
```

universal API proxy → self-operated model1

self-operated model2

OpenAI

Anthropic

🗄 **BerriAI / litellm**  (Public)

Python SDK, Proxy Server (LLM Gateway) to call 100+ LLM APIs in OpenAI format - [Bedrock, Azure, OpenAI, VertexAI, Cohere, Anthropic, Sagemaker, HuggingFace, Replicate, Groq]

🔗 **docs.litellm.ai/docs/**

⚖ View license

⭐ **20.7k** stars   ⑂ **2.6k** forks   ⑂ Branches   🏷 Tags

〰 Activity

- Custom routing logic
- Retry/fallback logic across multiple deployments
- Set Budgets & Rate limits per project, api key, model

🞐 netlight

# Learnings 💡

*What I wish I knew before building our first GenAI deployment.*

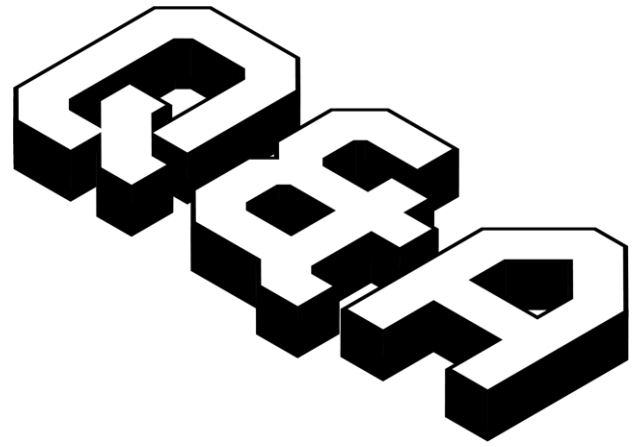- Even major hyperscalers occasionally face **GPU shortages** that can last for hours.

- The cost of requests to an LLM can vary by a factor > 100,000. Implement robust rate limiting and user budgeting mechanisms early on; **rate limiting solely by request counts is not helpful**.

- Leverage **predictive auto-scaling mechanisms by using forecasting** and closely monitor latency to maintain performance (e.g., KEDA ScaledObject).

- To tune and optimize model deployments effectively, it helps to **understand how model serving works**—but sometimes, the underlying mechanics are so complex that you just need to take them as they are.

netlight

# Pro 👍

- **Compliance**: Operating in highly regulated industries with very sensitive data

- **Flexibility**: Need full control and customizability (quantization techniques; model finetuning; niche-models)

- **Trust issues using open APIs**

- **High Usage expected**: Cost-effective at scale

- If your company aims to build in-house expertise in generative AI that goes beyond simply using available APIs, **it can lead to exciting collaborations between data scientists and infrastructure engineers**. 🚀

# Cons 👎

- **Low Usage**: API solutions are cheaper for low demand.

- **Risk of Over-Engineering**: If your use case is very straightforward or small-scale, self-operating models can be unnecessarily complex and resource-intensive.

- **Lack of knowledge building a platform within your company**

- **Open-weight models are not matching your requirements** and you need close—weight models (e.g. Anthropic Claude, OpenAI)

netlight

netlight

# Thank you!

**Tim Fischer**

Consultant @ Netlight

tim.fischer@netlight.com

netlight