# Davide Rutigliano

Senior Platform Engineer (AI & GPU Infra) · Italy

(+39) 340 9307386 · davide.ruti@gmail.com · LinkedIn · Github · Medium · Portfolio

## Summary

Senior Platform Engineer building HPC/GPU accelerated Kubernetes platforms for AI workloads. Specialized in inference observability (vLLM, TTFT, OpenTelemetry) and cluster lifecycle operations across cloud and bare-metal environments. Open-source contributor: Kubernetes, Kueue and KubeAI (vLLM operator).

## Experience

### SUSE
*Remote, Italy*

**Senior Platform Engineer** — Jan. 2026 – Present

- **vLLM & GenAI Observability**: Engineered OpenTelemetry (OTel) collectors to instrument vLLM inference (Time-to-First-Token and related KPIs), enabling on-call triage for multi-tenant GPU inference platform.
- **High-performance GPU Monitoring**: Engineered GPU observability solution for Kubernetes and KubeVirt clusters, enabling fine-grained monitoring of NVIDIA A2/H100 with MIG and vGPU segmentation, unlocking 40+% HPC usage efficiency.
- **Cluster lifecycle & Autoscaling**: Optimized platform efficiency achieving a 62% reduction in infrastructure costs ($100K+ annual savings), by deploying Cluster Autoscaling with Cluster API on multi-cloud (AWS, GCP, On-Prem) platform.
- **AI-driven Observability**: Architected and delivered SUSE Observability MCP Server from greenfield idea to MVP, embedding LLM-driven analysis into the alerting pipeline; recognized by senior leadership for roadmap inclusion.
- **Infrastructure Automation**: Designed and implemented a Kubernetes operator to orchestrate large-scale virtual machine migrations from KVM to Harvester, enabling the migration of 100+ VMs.

**Platform Engineer** — Jan. 2025 – Dec. 2025

- **GitOps Migration**: Led migration from push-based Terraform pipelines to pull-based GitOps (Fleet + Kubernetes operators), eliminating configuration drift and pipeline toil via automatic reconciliation.
- **Infrastructure Governance**: Designed and implemented a custom k8s operator to synchronize VLAN and virtual machine assets into NetBox, establishing a single source of truth for infrastructure inventory.
- **Federated Observability**: Architected migration from Prometheus/Grafana stack to SUSE Observability (StackState) for federated multi-cluster observability, cutting troubleshooting time by 25%.
- **Cloud Evolution**: Led cloud architecture evolution across AWS, GCP, and Azure for enterprise migrations, aligning Cloud Landing Zone (CLZ) design with strategic requirements, acting as advisor to internal customers.

**Stack:** Kubernetes, Rancher, Cluster API (CAPI), KubeVirt, GPU (NVIDIA MIG/vGPU), vLLM, OpenTelemetry, Prometheus, Grafana, Terraform, GitOps (ArgoCD/Flux), AWS/GCP/Azure

### ERICSSON
*Pagani, Italy*

**Devops Engineer Team Lead** — Sep. 2023 – Dec. 2024

- **Self-Service Platform**: Led a team of 5 engineers delivering an Internal Developer Portal (Spotify Backstage).

**DevOps Engineer** — Jun. 2022 – Sep. 2023

- **ML-driven Optimization**: Built an ML-driven k8s performance tuning tool, cutting optimization time from weeks to days.

**Cloud Engineer** — Mar. 2021 – Jun. 2022

- **Resource Efficiency**: Reduced Ericsson Licensing footprint by 25% via autoscaling and resource tuning.

**Stack:** Kubernetes/OpenShift, Backstage, Gatekeeper (OPA), GitLab CI/Jenkins, Kafka, PostgreSQL/Cassandra, Go, Python

### CISCO
*Vimercate, Italy*

**ML Engineer** — Apr. 2020 – Oct. 2020

- **ML-based Troubleshooting**: Engineered ML optical device troubleshooting, improving defect detection by a 2x factor.

**Stack:** Python, TensorFlow, Computer Vision

## Education

| | |
|---|---|
| **Master Degree in Computer Engineering** | Milan, Italy |
| Politecnico di Milano | 2018 – 2021 |
| **Bachelor Degree in Computer Engineering** | Fisciano, Italy |
| Universitá degli Studi di Salerno | 2014 – 2018 |

## Skills

**AI & GPU INFRA**: NVIDIA MIG/vGPU, GPU-Operator, LLM-Ops, vLLM, Kueue/Slurm, TensorFlow, Pytorch, Computer Vision

**OBSERVABILITY**: OpenTelemetry (OTel), Prometheus, Grafana, Alertmanager, StackState, Root Cause Analysis (RCA)

**RELIABILITY**: SLIs/SLOs, alerting strategy, runbooks, incident response, postmortems, capacity planning

**CLOUD NATIVE**: Kubernetes, Helm, Docker, GitOps (ArgoCD, Flux), Terraform, GCP, AWS, Azure

**DEVELOPMENT**: Go, Python, Java, Rust, K8s Operators, Event-Driven Architecture, Linux