

 khaykingleb@gmail.com

 <https://github.com/khaykingleb>

 <https://linkedin.com/in/khaykingleb>

GLEB KHAYKIN

 Amsterdam, The Netherlands

+31 06 26172415 

<https://khaykingleb.com> 

https://t.me/khaykingleb_blog 

SUMMARY

ML Infrastructure Engineer with 4+ years of experience architecting high-performance ML systems, multi-cluster platforms, and large-scale inference/training pipelines. Specialized in distributed systems, Kubernetes, Terraform, observability, and CI/CD. Delivered 50%+ GPU cost reductions, 4x CI/CD improvements, and inference systems serving 60k+ req/hour.

EXPERIENCE

Together AI

MLOps Engineer

Amsterdam, Netherlands

July 2025 — Present

- Engineered model sync pipeline and integration tests for 90+ models, unlocking dedicated and serverless inference for fine-tuned models
- Accelerated CI/CD pipelines 4x through BuildKit mount caching, Docker optimization, and self-hosted runners with shared cache
- Consolidated monitoring across 6 K8s clusters into unified Grafana Cloud stack with 30+ IaC-managed dashboards
- Migrated control plane to dedicated EKS with Karpenter autoscaling, eliminating 1h manual node provisioning and pod eviction incidents

Stealth Startup

MLOps/DevOps Engineer

Remote

July 2023 — July 2025

- Architected K3s GPU clusters with NVIDIA GPU Operator, reducing hosting costs by >50% via hybrid cloud optimization
- Established Terraform-managed observability platform for K8s infrastructure with Grafana Cloud, decreasing MTTR from 1 day to 1 hour
- Built inference system with NVIDIA Triton, Celery, and RabbitMQ for ASR/LLM workload orchestration
- Optimized ensemble ASR/diarization pipeline achieving 11x speedup, reducing processing costs by 5x
- Delivered multi-stage financial document parser using vision LLMs, achieving 70% automation and securing venture fund as first client

Myna Labs

ML/MLOps Engineer

Remote

November 2022 — March 2024

- Developed priority-queued multi-model speech REST/gRPC API (TTS, voice conversion, ASR) serving 60k+ req/hour
- Built voice cloning service using LoRA fine-tuning, reducing voice addition time from days to under 10 minutes
- Modernized and automated speech infrastructure, migrating from Docker Compose to Terraform-managed GKE with autoscaling

Stealth Startup

Data Scientist

Remote

September 2021 — July 2022

- Designed distributed PySpark pipeline for financial sentiment analysis, scraping and processing S&P 500 news data with BERT models
- Developed learning-to-rank system for quantitative stock selection in long-short trading strategies
- Engineered automated ETL workflows on Databricks to aggregate and analyze data from 150+ financial API endpoints

PROJECTS

- **Research Playground:** Production-ready PyTorch Lightning framework with Hydra configs and distributed training on self-hosted K3s
- **Dotfiles:** Declarative system configuration with Nix flakes and home-manager for reproducible macOS environments
- **Huffman Archiver:** C++ implementation of Huffman coding for lossless data compression with Conan package management
- **Personal Website:** Full-stack blog and portfolio with Next.js frontend, Supabase backend, and live content sync via Notion API

EDUCATION

NRU “Higher School of Economics”

Moscow, Russia

B.S., Computer Science and Finance (summa cum laude)

September 2018 — July 2022

CFA Institute

Moscow, Russia

Level 1 passed

February 2021

SKILLS

- **Languages:** Python, Go, Rust, C++, SQL, TypeScript
- **DevOps:** Terraform, Ansible, Nix, Docker, Kubernetes, Helm, ArgoCD, CI/CD, AWS, GCP, Grafana, Prometheus, Loki, Tailscale
- **MLOps:** NVIDIA Triton, DVC, TensorRT, ONNX, LangChain, LangGraph
- **R&D:** PyTorch, Lightning, Hydra, W&B, Gradio, Scikit-learn, Numpy, Pandas, SciPy
- **Data Engineering:** Spark, Hadoop, Polars, Kafka, RabbitMQ, Celery
- **Storage:** PostgreSQL, Redis, MongoDB, S3, PGVector
- **Backend:** FastAPI, gRPC, AsyncIO, REST, Pydantic
- **Frontend:** React, Next.js, Tailwind