

SUMMARY

MLOps/DevOps Engineer with 4+ years of experience in software engineering and infrastructure. Specialize in boosting model performance, reducing costs by >\$10k monthly, deploying scalable solutions for apps with >30k DAU, and optimizing workflow pipelines. Adept in integrating cutting-edge models and building robust infrastructure across cloud and on-prem environments. Background in speech, NLP, and finance.

EXPERIENCE

MLOps/DevOps Engineer at iClerk San Francisco, USA (Remote)
Silicon Valley startup focusing on automating tedious tasks for B2B, \$3.4M pre-seed investment July 2023 — Present

- Refined an ensemble model for speech-to-text and speaker diarization, achieving an 11x speedup and improved accuracy
- Developed a scalable inference pipeline with NVIDIA Triton, reducing costs by 5x for processing one-hour meetings
- Secured an early client with an advanced RAG to parse financial statements, decreasing manual workload by 70%
- Architected a K3s GPU cluster with NVIDIA GPU Operator for monitoring and validation, cutting GPU costs by >50%
- Established IaC monitoring and alerting for 4 K8s clusters with Grafana Alloy OTel Collector, reducing downtime and errors

ML/MLOps Engineer at MynaLabs Tbilisi, Georgia
Generative AI startup founded by a sole investor with 2 prior \$150M+ exits November 2022 — March 2024

- Built a REST/gRPC API for speech tasks with priority queuing, handling 60k+ inference requests per hour
- Researched and deployed state-of-the-art multilingual, emotional text-to-speech and voice conversion models in 44 kHz
- Created a service for fast speech model fine-tuning to any voice in under 1 minute
- Terraformized and migrated speech infrastructure to Kubernetes with >4000 lines of IaC
- Introduced MLOps/DevOps best practices to a 20+ developer engineering division

Data Scientist at DiviAI San Francisco, USA (Remote)
Silicon Valley startup building a financial data aggregator for the US market September 2021 — July 2022

- Engineered a PySpark pipeline for scraping S&P 500 news to estimate investor sentiment with BERT-like models
- Developed ETL pipelines that aggregate data from 150+ financial API endpoints
- Designed a listwise learning-to-rank system for stocks recommendation, enhancing long-short strategy performance
- Adopted a Glicko rating system to rank Russell 3000 CEOs

EDUCATION

NRU “Higher School of Economics” Moscow, Russia
B.S. with summa cum laude in Economics and Data Science September 2018 — July 2022

- Relevant coursework: C++, Python, R, Machine Learning 1, Machine Learning 2, Large Scale Machine Learning, Deep Learning, Deep Learning in Audio Processing, Reinforcement Learning, Calculus, Linear Algebra, Probability Theory, Mathematical Statistics, Stochastic Processes, Econometrics, Microeconometrics, Differential & Difference Equations

PROJECTS

- **Research Playground**: Research framework for ML/DL experimentation across multiple domains with actual implementations
- **Huffman Archiver**: ASCII archiver in C++ with a Conan-based auto-build system
- **Stochastic Optimization**: Particle Swarm Optimization and Genetic Algorithm applications on various use-cases

MISCELLANEOUS

- **CFA Level 1**: Passed in February 2021, scoring in the top 10% of candidates worldwide
- **Teaching Assistant**: Coordinated courses in Accounting, Probability, Statistics, and Machine Learning (2020–2022)

SKILLS

- **Languages**: Python, Rust, Go, C++, SQL, TypeScript, JavaScript
- **DevOps**: Terraform, Ansible, Docker, Kubernetes, K3s, AWS, GCP, Grafana, Prometheus, Loki, Nginx, Acme, ArgoCD, CI/CD
- **MLOps**: DVC, NVIDIA Triton, Ollama, LangChain, TorchScript, ONNX, TensorRT, CoreML, Gradio
- **R&D**: PyTorch, Lightning, Hydra, W&B, Scikit-learn, Numpy, Pandas, SciPy
- **Backend**: gRPC, REST, Websocket, Webhook, AsyncIO, FastAPI, Pydantic, JWT, Node.js
- **Frontend**: Remix, React, Tailwind, DaisyUI, Vite, Vercel
- **DE**: Hadoop, Spark, Databricks, Polars, Kafka, Celery, RabbitMQ
- **Storage**: Redis, PostgreSQL, PGVector, Sqitch, S3, R2, Supabase, Airtable