YU-PENG FU

AI · GPU · FULL-STACK ENGINEER

Profile

Al & GPU Engineer with expertise in algorithms and full-stack development. Skilled in optimizing Al inference and building scalable RAG pipelines that connect large language models to real-world data. Experienced in GPU acceleration, backend systems, and delivering high-performance, production-ready Al solutions

Experience

AI · GPU · BACKEND ENGINEER

Kog Al, Paris (Hybrid) | 2023/04 - Present

- Contributed to securing €5M in venture funding by developing backend systems and applying model finetuning and prompt engineering to enhance the company's flagship AI product.
- Led the design and development of a custom inference engine from scratch, including low-level GPU kernel optimization and mathematical innovations, achieving 3× faster inference performance than vLLM
- Represented the company at the RAISE Summit and earned 2nd place at the AMD AI Sprint Hackathon for groundbreaking work in GPU optimization and AI acceleration
- Implemented model fine-tuning pipelines for domainspecific LLM adaptation. Incorporated LoRA and quantization strategies to optimize inference latency and memory footprint

Education

DIGITAL TECHNOLOGY ARCHITECT PROGRAM

L'école 42, Lyon

- Project-based curriculum emphasizing autonomy and real-world problem-solving
- Ranked top 10 nationwide in algorithm competitions representing the school
- Built algorithm projects with ~100 GitHub stars and a Medium technical article reaching ~100k views

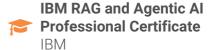
MATHEMATICS

National Taiwan University, Taipei GPA 4.3 (Top 1%)

Contact

- □ https://yu-peng-fu.vercel.app
- **1** (+33) 6 33 86 51 48
- ≥ leofu890806@gmail.com
- ★ Lyon, France (Open to Remote)
- in www.linkedin.com/in/yu-peng-fu
- nttps://github.com/LeoFu9487

Certifications



Achievements

AMD AI Sprint Hackathon 2025

2nd Place

AMD

La Coding Battle 2021

7th PlaceLe Shaker

BattleDev 2021

18th PlaceBDM

Projects

ModelPulse

Self-improving GPU-accelerated RAG system (Open-source)

Ysoft Management App

Production-grade web app actively used by Taiwanese vendors

Push_Swap Test Suite & Tutorial

End-to-end educational and testing toolkit for École 42

Skills

AI & GPU Computing

RAG & Inference Eng GPU Optimization CUDA | HIP | Triton vLLM | PyTorch



Full-Stack Development

TS | JS | Python Git | Linux | Docker Next.js | Node | SQL

