

Aditya Raj

Systems & ML Engineer

NIT Patna (2022–2026) | CGPA: 8.1/10.0

-  adityar.ug22.ec@nitp.ac.in
-  +91-8797073498
-  linkedin.com/in/hexronus
-  github.com/hexronuspi
-  hexronuspi.github.io

Professional Experience

- Software Engineering Intern (Infrastructure)** India
May 2025 – Oct 2025
 - QFI Research Capital*
 - **Backend Engineering:** Engineered a high-throughput backend service in **Golang** and **FastAPI** for backtesting toolkits.
 - **Performance:** Implemented custom functions, achieving a throughput of **100,000 events** with **latency under 9ms**.
 - **Charting:** Engineered and implemented a WebAssembly (Wasm) script that optimized UI data loading, resulting in a **50%** improvement in smoothness and a noticeable reduction in latency during user interactins.
- Research Intern (Academic Thesis)** Hyderabad, India
Apr 2025 – July 2025
 - IIT-H / Spatial Informatics Lab*
 - Designed a framework integrating KG with RAG, for retrieval tasks, outperformed existing retrieval baselines by **14%**.

Research

[Knowledge Graph-Informed Query Decomposition\(KG-IQD\): Hybrid KG-RAG Reasoning in Noisy Contexts](#)
Authors: Aditya Raj, Dr. Kuldeep Kurte^{PI}

Findings

[Why Safety Constraints in LLMs Are Easily Breakable? Knowledge as a Network of Gated Circuits](#)
Proposed a method to understand emergent phenomena in LLMs and their internal representations. [\[docs\]](#)

Notable Software

- SecureLock** | Anti-Cheat toolkit [C++] [Source Code]
- Engineered a user-mode evasion detection engine using C++ Windows API hooks to monitor thread execution and identify hypervisor signatures (VMware, KVM), a 3-tier architecture (Native EXE, Chrome Extension, Web App) to detect RDP/Screen-sharing in real-time, Added external factors tab switching, used hardware keystroke, screen input to detect a missed remote access signature.

Key Projects

- Efficient LLMs via Switchable and Dynamic Quantization** docs / October 2025
 - Tools: PyTorch, Hugging Face, LoRA, QAT-LLM, SQuAD Dataset
 - Engineered a **switchable quantization framework** for **GPT-2** using **adaptive LoRA modules** to toggle **layer-wise precision (INT8/FP16)**, reducing model footprint by **29% (207MB to 146MB)** without compromising semantic coherence by a huge amount.
 - Implemented **Cyclic Precision Training (CPT)** on the **SQuAD** dataset, CPT achieved a perplexity of **10.18**, which is an improvement over the Joint Training baseline of **10.21**, demonstrating that **CPT maintains or slightly exceeds the baseline's semantic coherence**.
 - Validated **random precision switching** as an adversarial defense against **PGD** attacks, increasing model robustness and recovering inference accuracy to **72.48%** (vs. **71.22%** fixed precision) in alignment with "Double-Win" quantization principles.
- Optimizing and Quantizing FBNet Models for Edge Deployment** docs / October 2025
 - Tools: Hugging Face, PyTorch, Edge Device
 - Converted **FBNet-A/B** from **PyTorch** → **TensorFlow** → **TFLite (FP32, FP16, INT8)** for edge deployment, achieving **100% accuracy** and up to **4× model size reduction**.
 - Rebuilt architectures in **Keras+TF**, used a **parser** for weight transfer, and verified using **MSE and accuracy metrics**.

Technical Proficiency

Core Stack	C++, Python, Golang, SQL (Postgres)
Frontend	Next.js (App Router), React, GSAP, TailwindCSS, Redux Toolkit, Webpack, Vite
Backend	FastAPI, Node.js, Kafka, RabbitMQ, gRPC, WebSockets, Redis, Nginx
Infrastructure	Docker, AWS (EC2, S3, Lambda), Kubernetes (K8s) - Basic, CI/CD (GitHub Actions), Vercel
Tools & Concepts	Git, Linux, System Design, Microservices, Serverless Architecture, Database Indexing

Select Achievements

M2L School	Selected (Top 10%) for global ML School summit in Split, Croatia	2025
Mathematics	Top 0.4% in Regional Mathematical Olympiad (RMO) among 250k candidates	2019
Hackathons	Winner, IIT ISM AI Challenge (Built OCR pipeline)	2024