

# Aditya Raj

B.Tech, Electronics and Communication Engineering  
National Institute of Technology, Patna  
Expected Graduation: 2026

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

## Summary

A speed and quality focused engineer specializing in the design of **high-performance, low-latency systems** for ML and data. Designed and built a high-throughput data processing tool during a quant internship, achieving processing speeds of <9ms for over 100k data points and individual operation latencies under 80ms.

## Education

Degree/Certificate	Institute/Board	GPA/Percentage	Year
B.Tech, Electronics and Communication Engineering	NIT, Patna	8.1 / 10.0	2022–2026

## Achievements

Stanford AI Lab	Interviewed at Prof. Jiajun Wu's Lab for research position in computer vision	2025
NK Securities Research	Ranked <b>67<sup>th</sup>/2095</b> in IV Prediction; MSE = <b>1.3e-5</b>	2025
M2L Summer School	Selected from <b>1600+</b> global applicants (BTech to Industry); Split, Croatia	2025
Amazon ML Challenge	Ranked <b>184<sup>th</sup>/75,000+</b> ; F1 score = <b>0.4667</b> using fine-tuned moondream	2024
IIT ISM AI of GOD	Winner; WER = <b>0.116</b> using TrOCR + T5; post-processing algorithm	2024
RMO	Top <b>0.4%</b> nationwide mathematical olympiad(classes 8–12) in India, out of 250k students	2019

## Research

- **Knowledge Graph-Informed Query Decomposition(KG-IQD): Hybrid KG-RAG Reasoning in Noisy Contexts**  
*Authors: Aditya Raj, Dr. Kuldeep Kurte<sup>PI</sup> / Poster (ISWC 2025) - Rejected*

## Experience

AI Research Intern	India
QFI Research Capital	May 2025 – October 2025
– Worked with a team of 2 engineers to build <b>data pipelines</b> and <b>forecasting models</b> predicting <b>product timelines</b> and assessing <b>market impact</b> , integrated with a <b>real-time sentiment engine</b> for long-term <b>alpha capture</b> .	
– Resolved <b>critical computation bugs</b> , built and deployed an <b>internal toolkit</b> to manage <b>workflow</b> .	
Research Intern	Hyderabad, India
IIT-H / Dr. Kuldeep Kurte, Spatial Informatics Lab	Apr 2025 – July 2025
– Achieved <b>state-of-the-art results</b> on a <b>custom disaster QA</b> benchmark, outperforming RQ-RAG by <b>14%</b> and <b>KG</b> by <b>18%</b> , by developing a <b>neuro-symbolic framework</b> that guides query decomposition using <b>Knowledge Graphs</b> and interrelates points with <b>RAG</b> on sub-queries for <b>QA</b> over structured and unstructured data.	

## Projects

Efficient LLMs via Switchable and Dynamic Quantization	<a href="#">docs</a> / October 2025
Tools: PyTorch, Hugging Face, LoRA, QAT-LLM, SQuAD Dataset	
– Switchable and dynamic quantization - <b>GPT-2</b> , per-layer bit-width control (INT8–FP32); adaptive <b>LoRA</b> .	
– Trained on <b>SQuAD</b> using <b>cyclic precision training</b> and joint bit-width optimization, achieving stable accuracy across dynamic precision configurations and demonstrating <b>quantized inference</b> .	
– Evaluated the <b>robustness</b> under random precision switching, aligning insights with <b>CPT (ICLR'21)</b> and <b>Double-Win Quant (ICML'21)</b> and found it perfectly aligned.	
Optimizing and Quantizing FBNet Models for Edge Deployment	<a href="#">docs</a> / October 2025
Tools: Hugging Face, PyTorch, Edge Device	
– Converted <b>FBNet-A/B</b> from <b>PyTorch</b> → <b>TensorFlow</b> → <b>TFLite (FP32, FP16, INT8)</b> for edge deployment, achieving <b>MSE 1e-19, 100% accuracy</b> , and up to <b>4x model size reduction</b> .	
– Rebuilt architectures in <b>Keras+TF</b> , used a <b>parser</b> for weight transfer, and verified using <b>MSE and accuracy metrics</b> .	
– Implemented <b>TFLite GPU batch resizing</b> , improving <b>conversion stability</b> and <b>edge-device performance</b> .	
Serverless Web Platform	<a href="#">webpage</a> / Dec 2024 – Feb 2025
Tools: Next.js, Vercel, Supabase	
– Replaced a \$150/month SaaS solution by building a serverless platform that migrated 10,000+ users.	
– Achieved <1s LCP on slow 4G networks using a Next.js/Vercel/Supabase stack with extensive frontend optimizations.	

## Technical Skills

Languages	Python, C++, SQL, GoLang, BASH, TypeScript, MATLAB, Triton (basic)
Developer/ML	FastAPI, RestAPI, PyTorch, Hugging Face (Transformers, PEFT), TensorFlow, scikit-learn, NumPy
AI	RLHF, DPO, PPO, Instruction-Tuning, Supervised Fine-Tuning (SFT)
DevOps & Tooling	Docker, Git, GitHub Actions, CI/CD, monitoring, Linux (WSL/Ubuntu), W & B
Coursework	Distributed System, Operating Systems, OOPS, DBMS, OS, Networking, System Design

## Notable Software

**SecureLock**  | Anti-Cheat toolkit [C++]

3-tier VM / RDP / screen-share detector: native EXE Chrome extension web app.

Identifies VMware, VirtualBox, Hyper-V, TeamViewer, VNC, etc.; real-time alerts, evasion-resistant.

**Darpan**  | GPU-accelerated offline screen share

Pure C++ client+server; AES-encrypted, Wi-Fi direct, zero-cloud, <16 ms latency for 4K classrooms.