

The most extraordinary thing I've worked on.

The term extraordinary is quite relative and wide. Here's how I would define extraordinary: solving a problem that is highly relevant in today's day and time with the resources at hand.

In the last couple years when AI blew up, I noticed that as a developer there are no pain-free ways of working with LLMs at scale - since I am an engineer and wanted to integrate LLMs into my apps. Emphasis on "pain-free" - because sure there are tools such as Ollama-cpp, AWS, GCP and so on. However, the time, resources, and capital that developers invested in building infrastructure to create well-optimized, cost-efficient LLMs was the same, if not more than the savings gained from those models.

That's why I built **model0** - a one-stop solution for all LLM, and ML devs. It's a platform as a service that utilizes AWS SageMaker and Bedrock to setup all the necessary infrastructure for any model that the developer needs, carry all the required optimizations on it (pruning, quantizations, RAG) and get a singular deployment link that the developer can integrate into their app.

Think Vercel, but for LLM infrastructure. The core philosophy is to reduce clutter, time and capital throughout the entire process so businesses can focus on the creative aspect of things. We use AWS SageMaker's infrastructure to train, prune and quantize models pulled from Huggingface, and host it on Bedrock - since it's much more cost effective for LLM inference as compared to real-time SageMaker.

How hard did you work on it? Would you do it again?

I have been working on it since October last year. I was in my second year of grad school, so managing school, being a TA and shipping on model-0 was tricky - but I'd do it all over again. Although since I've started, I've re-iterated on my core approach multiple times. I initially began with having only a SageMaker infrastructure, but then after initial customer reviews - I gathered critiques that the pricing goes too high. I love shipping things quickly - the faster I ship, the faster I get customer feedback, the faster I know if a feature worked or not, the faster I iterate.

I believe iteration is the most important thing that determines the success of a product. I like to call this the "momentum" of a team. One should get attached to the product, not the feature because it's the features that make a product, not vice-versa.

I'd love to talk more about what Entangl's working on and how I can contribute. In the meantime, know more about me on the links below:

Portfolio @ <https://athk.dev>

X @ <https://x.com/athkdev>

GitHub @ <https://github.com/athkdev>

LinkedIn @ <https://linkedin.com/in/athk>

Education	
<div>Northeastern University</div> <div>Master of Science, CS + Information Systems (full stack infrastructure emphasis)</div> <div>Relevant courses: Algorithms, Linux/UNIX, Web and API development, Object Oriented Design in Java, Design patterns, Operating Systems</div>	<div>Boston, MA</div> <div>Sep 2023 - May 2025</div>
<div>University of Mumbai</div> <div>Bachelor of Engineering, Electrical Engineering &amp; Computer Science</div> <div>Relevant courses: Programming in C, Object Oriented Design, Databases, Computer Networking &amp; Architecture, Big Data, Map Reduce, Hadoop, Image Processing</div>	<div>Mumbai, IN</div> <div>Aug 2018 - May 2022</div>

Technical Skills	
Frontend	React, Typescript, JavaScript, Next.js, Vue.js, Vite, npm, JQuery, Tailwind, Three.js
DevOps & Systems	Linux (Fedora and Ubuntu), Bash, Docker, Jenkins, Terraform, CI/CD, GitHub Actions, REST APIs, Cloudflare, Git
Backend & Databases	Python, Django, Go, GraphQL, C++, SQL, PostgreSQL, MVC API design, Bash, npm, HTTP, TCP/IP
Data & AI	ETL Pipelines, AWS EC2, AWS SageMaker, LangChain, OpenAI, Pandas, Numpy, Matplotlib, Plotly, PyTorch, Apache Kafka

Work Experience	
<div>Northeastern University</div> <div>Graduate TA</div> <div><ul style="list-style-type: none"><li>Instruct students in Angular, Typescript, Git, GitHub and the Agile methodology through 1-on-1 sessions and online communications.</li><li>Run weekly office hours, breaking down tough concepts and debugging live for 25+ students, including pair programming when necessary.</li></ul></div>	<div>Boston, MA</div> <div>Sep 2024 - Apr 2025</div>
<div>Cognizant</div> <div>Software Engineer 1</div> <div><ul style="list-style-type: none"><li>Built and maintained an in-house Angular component library with <b>28+</b> reusable and themeable UI elements. Used by <b>5+</b> teams to improve development speed and maintain visual consistency.</li><li>Developed pixel perfect UI as per wireframes from the UX design team for OneCognizant (company-wide appstore). Built in <b>React</b>, <b>Tailwind</b>, <b>SCSS</b> and <b>Typescript</b>.</li><li>Engineered a payments app for external stakeholders by integrating Stripe using <b>React</b>, <b>Python</b>, <b>Apache</b> and <b>AWS</b>, supporting <b>40+</b> local payment gateways and <b>130+</b> currencies.</li><li>Write integration tests with Playwright and React Testing Library to validate end-to-end user workflows - ensuring <b>90%+</b> test coverage.</li><li>Refactored a monolithic app into modular <b>Python microservices</b>, enabling service-specific deployments and minimizing downtime during updates.</li><li>Developed a well documented <b>Python</b> module that simplified application deployment workflow that automates containerization with <b>Docker</b> and <b>Terraform</b> resulting in a <b>one-click deployment solution</b>.</li></ul></div>	<div>Mumbai, IN</div> <div>Aug 2022 - Jun 2023</div>
<div>Raftlabs</div> <div>Software Engineer (Founding Engineer)</div> <div><ul style="list-style-type: none"><li>Developed a website for a supply chain analytics tool, optimizing delivery schedules and inventory management with scalable <b>React</b>, and <b>Java</b>. Used the <b>in-browser IndexedDB</b> for local caching resulting in extremely fast user experience.</li><li>Revamped a financial loyalty app for the Bank of Ireland using <b>React</b>, <b>GraphQL</b>, <b>Python (Django)</b>, and <b>Docker</b> to streamline revenue tracking and help customers gain rewards.</li><li>Integrated Meta’s webhooks with <b>Shopify</b> and <b>Node.js</b> to help viewers buy products by commenting codes in <b>Facebook live streams</b>.</li><li>Coded a <b>proof of concept</b> for an efficient coupon distribution system in Typescript, Node.js and PostgreSQL that could handle <b>2000 requests per second</b>.</li><li>Implemented an automated CI/CD pipeline using <b>GitHub Actions</b> and <b>Python</b> to schedule blog posts for the marketing team - resulting in an efficient workflow with no manual overhead.</li></ul></div>	<div>Remote</div> <div>Apr 2021 - Aug 2022</div>

Projects (Infrastructure and full stack)	
<div>Model-0 - Platform as a Service</div> <div><ul style="list-style-type: none"><li>Full stack web app for creating and deploying AI chatbots with AWS SageMaker providing integration for developers in <b>less than 4 minutes</b>.</li><li>Set up HTTP Apache server with Amazon EC2 instance to handle swift load balancing and secure encryption by setting up a SSL certificate.</li><li>Developed a <b>robust REST API</b> with authenticated routes for communicating with SageMaker API to manage creation, and tuning of ML models. Create permissions in Linux for Apache server by adding it to the current user group.</li><li>Configured a reverse proxy between the Apache route traffic from web and localhost for serving a dynamic Next.js app.</li><li>Built a GitHub actions CI/CD pipeline that deploys code on an EC2 instance automatically once code is pushed to GitHub.</li><li>Used Next.js, React.js, REST API, Typescript, PostgreSQL AWS SageMaker and AWS EC2. The app is launched live <a href="#">here</a>.</li></ul></div>	<div><a href="#">github.com/athkdev/model0</a></div>
<div>Mental Health Counseling AI Chatbot</div> <div><ul style="list-style-type: none"><li>Chatbot-powered full stack app designed for mental health counselors to make informed decisions while actively counseling.</li><li>Extracted NLP insights from a dataset of <b>3000+ counsellor-client</b> sessions and stored them in <b>PostgreSQL database</b> for filtering and pagination.</li><li>Integrated an OpenAI-powered chat that is pre-trained to respond like a counsellor with guided conversations to improve accessibility.</li><li>Built with <b>OpenAI</b>, <b>Next.js</b>, <b>React</b>, <b>Typescript</b>, <b>Node.js</b>, <b>Python</b> and <b>PostgreSQL</b>. The app is launched live <a href="#">here</a>.</li></ul></div>	<div><a href="#">github.com/athkdev/legacytheapp</a></div>
<div>Highly performant web crawler in Java</div> <div><ul style="list-style-type: none"><li>Highly performant web crawler in Java that includes a CLI and stores all URLs in a graph data structure using Neo4j.</li><li>Applied design principles for a scalable and low-latency architecture. Store URLs in Neo4J database concurrently in a graph format.</li><li>Utilised Java’s CompletableFuture API to achieve concurrency and parallelism. Benchmarked at 730 URLs per second.</li><li>Built with <b>Java</b> and <b>Neo4J</b>.</li></ul></div>	<div><a href="#">github.com/athkdev/crawler</a></div>