

EREL OZEN

Irvine, CA

949-339-7540

erel.ozen.swe@gmail.com

linkedin.com/in/erel-ozen

eozen1.github.io

Education

Georgia Institute of Technology

B.S Computer Science (*Intelligence & System Architecture*) — **GPA: 4.00**

May 2026

Atlanta, GA

Relevant Coursework: Design of Operating Systems, Systems & Networks, Design & Analysis of Algorithms, Data Structures & Algorithms, High Performance Computing, Compilers, Machine Learning, Computer Organizations & Programming (*Teaching Assistant*), Artificial Intelligence, Perception & Robotics, Object-Oriented-Programming

Technical Skills

Languages: Python, C/C++, Java, JavaScript, TypeScript, Bash

Frameworks: TensorFlow, PyTorch, ROS/ROS2, OpenCV, Matplotlib, ReactJS, React Native, NATS.io, JUnits, pySerial

Technologies: RESTful APIs, AWS, Git, Docker, SLAM, Controller Area Network (CAN), Micro Controllers, LoRA

Experience

Software Development Engineering Intern, CoreOS

May 2025 – August 2025

Apple

Cupertino, CA

- Architected and deployed a multi-agent AI triage tool for **1600+** engineers, using custom MCP servers to autonomously analyze large logs, correlate data sources, and generate actionable debugging steps.
- Engineered robust state management logic for paravirtualized network drivers within the virtualization framework, resolving critical bugs to ensure network reliability and system performance at scale.
- Selected as **one of 10** interns company-wide to present to Apple's SVP of Software Engineering, Craig Federighi.

Software Engineering Intern, Machine Learning

May 2024 – May 2025

OKSI

Torrance, CA

- Secured a **\$1.5M** contract renewal from DARPA after presenting an end-to-end AI automation platform.
- Deployed Llama-3 on resource-constrained devices with **400%** less VRAM by leveraging intelligent weight quantization.
- Developed a full-stack React application featuring an NLP pipeline with LoRA fine-tuning, NATS.io, & Synadia Cloud.
- Architected a PDF form processing microservice that reduced manual entry with an auto-populating question-wizard.

Software Engineering Intern, Embedded Systems

June 2023 – August 2023

OKSI

Torrance, CA

- Orchestrated configuration and calibration of multispectral cameras, GPS, and IMU for machine-learning applications.
- Evaluated camera drivers on a Google Coral Board, implementing debayering techniques for optimal performance.
- Developed high-speed serial communication between microcontrollers and avionics, focusing on data pipeline efficiency.
- Implemented simultaneous localization and mapping (SLAM) using LiDAR, IMU, and GPS for autonomous navigation.

Projects

Vidur Autoscaling | LLM Inference Simulator Extension

Fall 2025

- Implemented reactive autoscaling policies for distributed LLM inference using traffic analysis and throughput estimation
- Optimized O(1) traffic envelope algorithms for real-time token arrival rate tracking in production systems
- Developed LOR scheduler with autoscaling awareness to minimize latency during replica scale-down operations
- Engineered custom autoscaling policy generating Pareto-optimal latency-cost trade-offs across service levels

UNIX Kernel Development | OS development for a UNIX v6 Kernel

Spring 2025

- Implemented core OS components for a UNIX V6-based teaching OS in ANSI C on x86, including the bootloader, virtual memory management, file system, and CPU scheduling
- Developed system-level functionality such as interrupts, concurrency mechanisms, user and kernel threading, and basic networking, focusing on process execution, resource management, and inter-process communication

Organizations

Research Assistant

August 2025 – Current

Systems for AI Lab, Georgia Institute of Technology

Atlanta, GA

- Building LLM-serving distributed system and adding support for various types of attention (DeepSeek, Qwen)
- Increased LLM throughput by optimizing prefix caching strategies through dynamic memory management

Engineering Team Member

August 2023 – May 2024

Fintech @ Georgia Tech

Atlanta, GA

- Created a credit card smart-wallet using React, Google Firebase, and Node.js that dynamically selects a payment credit card based on the purchase type and recommends credit cards based on purchase history