

Pranav Saran

(510)709-7661 | pranav.saran@case.edu | [LinkedIn](#) | [Github](#) | [Google Scholar](#)

EDUCATION

Case Western Reserve University <ul style="list-style-type: none">Major: Computer Science and EngineeringDeans High Honors List Fall ‘24, Spring ‘25	Cleveland, OH 4.0 GPA
--	---------------------------------

TECHNICAL EXPERIENCE

Machine Learning Researcher – Palo Alto, CA <i>Algoverse (Python, PyTorch, Overleaf, LaTeX)</i> <ul style="list-style-type: none">Utilized a hybrid mBERT+BiLSTM model for figurative language detection, trained on low-resource KonkaniAchieved an accuracy of 83% for idiom classification and 78% for metaphor classificationPreserved 100% of original accuracy on idiom classification while pruning attention headsAchieved 88% accuracy retention on metaphor classification tasks through strategic model parameter reductionPerformed ablation testing across multiple transformer architectures (mBERT, IndicBERT, XLM-R) to evaluate robustness and comparative performance	January 2025 - June 2025
Software and Electronics Developer – Cleveland, OH <i>Case Western Reserve Global Health Design Collaborative (Arduino, C++, Python, Kivy)</i> <ul style="list-style-type: none">Engineered a Soft Access Point utilizing the ESP32 C3, facilitating seamless data transmission from the MAX30102 sensor to a mobile application developed in Python, enhancing real-time monitoring capabilities by 40%.Contributed to award-winning research showcased at CWRU Intersections, earning Second Place in Undergraduate Engineering, with findings from the project published and presented to a multidisciplinary academic audience; named as 2024-2025 team MVPDeveloped and integrated a mobile app using the KIVY library, enabling users to visualize heart rate and oxygen saturation data(SPO2) collected from over 500 sessions, thereby improving user engagement metrics by 60%.Collaborated with a cross-functional team to design and implement a detrending algorithm, improving data accuracy by 40% and enhancing the overall reliability of readings across a plethora of data points	September 2024 - Present

PUBLICATIONS & RESEARCH ACTIVITIES

<ul style="list-style-type: none">Pruning for Performance: Efficient Idiom and Metaphor classification in Low-Resource Konkani Using mBERT - ACL SRW (average rating: 3) Awaiting Response, EMNLP, COLMServed as a reviewer for WMDQS 2025

PROJECT EXPERIENCE

Transformer Implementation Baby-GPT <i>Python, PyTorch, Transformers</i> <ul style="list-style-type: none">Implemented the transformer based architecture from the paper <i>Attention Is All You Need</i>Developed a character tokenizer to encode text into tokensImplemented Self-Attention, Multi-Headed Attention and Normalization to train the Bigram Model on Shakespeare’s worksReplicated the performance of GPT-2 with successful implementation
Multi-Task Clinical NLP Pipeline BERT, HuggingFace, MIMIC-III <ul style="list-style-type: none">Preprocessed 40K+ clinical notes from MIMIC-III to create multi-label task formatsFine-tuned BioBERT with custom multitask heads, achieving 91% ICD F1 and 89% NER accuracyUsed ROUGE metrics to benchmark summarization quality; improved performance by 18% over single-task baselinesApplied attention head pruning to reduce model size by 30% while retaining 95% performance
Auto README Generator Agent <i>Python, SmolAgents, OpenAI GPT-4o, Gradio</i> <ul style="list-style-type: none">Developed an autonomous agent using smolagents to automate technical documentation by analyzing codebases and generating READMEs using LLM-driven reasoningEngineered a modular toolchain for code parsing, dependency inference, and capability extraction, enabling structured analysis of arbitrary GitHub repositoriesIntegrated a Gradio-based interface for seamless human-agent interaction, showcasing system design, prompt engineering, and practical application of multi-tool LLM workflows
Autonomous Agent for Reasoning Benchmarks <i>Python, smolagents, OpenAI API, LLM Tooling</i> <ul style="list-style-type: none">Developed and deployed a multimodal reasoning agent using the smolagent framework, achieving automated task-solving across search, code, and file understanding using tools like Wikipedia lookup and YouTube analysis.Integrated GPT-4o with custom tools and evaluated agent performance on the GAIA benchmark via Hugging Face Spaces, using REST APIs to batch-run and submit results.Engineered tool schemas, error handling, and agent orchestration logic to improve generalization and robustness across diverse natural language and data-driven tasks.

TECHNICAL SKILLS

<ul style="list-style-type: none">Coding: C++, Java, Python, Javascript,CSS, HTMLFrameworks: MongoDB, Express JS, Node JS, Flask, Firebase, ChromeAPI, NEAT, JFrame, Kivy, PyTorch, +more!Relevant Coursework: Data structures, Computer Security, Discrete Mathematics, Logic Design and Computer Organization, Linux and OSCertifications: The LLM Course (Hugging Face), Fundamentals of MCP (Hugging Face), AI Agents Course (Hugging Face)
--