

<div> <div>Karan Singh Chauhan</div> <div> 6058 Rothko Ln, Simi Valley CA   +18125388341   <a href="mailto:Karan.m0624@gmail.com">Karan.m0624@gmail.com</a>   <a href="https://github.com/karan0624">linkedin.com/in/karan-chauhan-49219213</a>  <a href="https://github.com/karan0624">https://github.com/karan0624</a> </div> </div>	
<div> <div>Professional Summary</div> <div> Data Scientist and Machine Learning Engineer with proven expertise in Generative AI, predictive modeling, and scalable data systems. Skilled in designing and deploying LLM inference pipelines using vLLM, TGI, RAG, and vector databases, applying advanced optimizations like 4-bit quantization (AWQ) and speculative decoding for efficient model performance. Experienced in statistical learning, forecasting, and multimodal AI, delivering real-world impact across healthcare analytics, computer vision, and enterprise automation. Strong in applied research and experimentation, bridging theoretical innovation with production-ready, high-performance AI solutions that drive measurable business outcomes. </div> </div>	
<div> <div>Technical Skills</div> <div> <ul style="list-style-type: none"> <li><b>Inference Optimization:</b> vLLM, TGI, ONNX Runtime, quantization (4-bit, AWQ, GPTQ), speculative decoding</li> <li><b>Low Latency Systems:</b> Caching, batching, tensor parallelism, distributed inference</li> <li><b>LLM Deployment:</b> Containerized serving (Docker, Kubernetes), FastAPI, scalable REST APIs</li> <li><b>Vector &amp; RAG Pipelines:</b> FAISS, Pinecone, Weaviate, Chroma for retrieval-augmented LLMs</li> <li><b>Machine Learning:</b> PyTorch, Hugging Face Transformers, Scikit-learn</li> <li><b>Cloud &amp; Infrastructure:</b> Azure, AWS, Kubernetes, CI/CD (GitHub Actions)</li> <li><b>Programming:</b> Python, SQL, CUDA (beginner), Triton (exploring)</li> </ul> </div> </div>	
<div> <div>Experience</div> <div> <div> <div>Research Assistant   University of Alabama at Birmingham, Birmingham, AL</div> <div>09/2025 – Present</div> <ul style="list-style-type: none"> <li>Conducting research on AI, LLMs, NLP, and Bioinformatics, applying advanced machine learning to healthcare and computational biology problems.</li> <li>Developing and evaluating models with PyTorch, TensorFlow, and Hugging Face Transformers for biomedical text analysis and data interpretation.</li> <li>Performing data preprocessing, feature engineering, and statistical analysis on clinical and genomic datasets to support experimental studies.</li> </ul> </div> <div> <div>CETS AV/IT Technician   Indiana State University, IN, US</div> <div>02/2024 – 05/2025</div> <ul style="list-style-type: none"> <li>Configured and managed campus IT &amp; network infrastructure using Cisco switches, system imaging, domain joins, and static IP mapping to onboard and maintain devices across the college.</li> <li>Installed, maintained, and troubleshoot AV systems (Extron, projectors, HDMI switching), ensuring seamless classroom delivery and zero downtime for campus-wide setups.</li> </ul> </div> <div> <div>Data Scientist – Industrial / Manufacturing Analytics   Make Brass Industries, GJ, India</div> <div>01/2023 – 05/2023</div> <ul style="list-style-type: none"> <li>Built and deployed real-time ML systems (YOLOv5 + OpenCV on Raspberry Pi) for brass manufacturing, reducing surface defect rates by 25% and minimizing manual inspection.</li> <li>Developed predictive models (XGBoost, time-series forecasting) to optimize export forecasting and inventory planning, improving delivery schedule accuracy.</li> <li>Engineered automated data pipelines with Azure Data Factory + SQL, integrating IoT sensor data for real-time tracking and presented insights to leadership in client-facing settings.</li> </ul> </div> <div> <div>Machine Learning Engineer   Limelight IT &amp; Research Pvt. Ltd., Ahmedabad, GJ, India</div> <div>04/2021 – 06/2022</div> <ul style="list-style-type: none"> <li>Reduced ML training and inference time by 45% and improved data reliability by 35% by developing real-time AI pipelines with Databricks, Apache Spark, and anomaly detection models (LSTM, XGBoost).</li> <li>Deployed scalable containerized microservices by integrating FastAPI with Kubernetes and Docker, enabling RAG-style health telemetry inferences.</li> <li>Improved support response efficiency by 20% by automating diagnostics and ticket triaging through Python and PowerShell scripts.</li> </ul> </div> </div> </div>	
<div> <div>Projects</div> <div> <div> <div>Role of Agentic AI in Eradication for the Need of Scrum Master (FastAPI, Docker, LangChain, GPT-4, Pinecone/FAISS, Slack API, Jira API, GitHub API) Self-initiated research &amp; development (2024 – Present)</div> <ul style="list-style-type: none"> <li>Developed an AI Scrum Master Assistant that converts live meeting speech into structured sprint updates and Jira tasks using speech-to-text, LLaMA 3, and RAG pipelines.</li> <li>Optimized multimodal inference pipelines with vLLM (token streaming, speculative decoding) and vector databases (FAISS/Weaviate), boosting contextual accuracy and reducing latency.</li> <li>Deployed cost-efficient large-scale models with 4-bit quantization, Docker, and Kubernetes, delivering scalable real-time performance for actionable team intelligence.</li> </ul> </div> <div> <div>Speech Emotion Recognition - (Deep Learning, PyTorch, CNNs, Streamlit)</div> <ul style="list-style-type: none"> <li>Achieved 92% accuracy by developing a transformer-based speech emotion recognition model (30M parameters) with ResNet-based audio embeddings.</li> <li>Improved training stability by 40% through adaptive gradient clipping and batch normalization.</li> <li>Enhanced preprocessing and feature engineering workflows using Python libraries (pandas, NumPy) for audio data exploration.</li> </ul> </div> <div> <div>Diabetic Readmission Prediction - (TensorFlow, FastAPI, Docker, Azure)</div> <ul style="list-style-type: none"> <li>Achieved 87% precision and reduced false positives by 32% by optimizing weighted cross-entropy loss and fine-tuned SGD (momentum = 0.9, weight decay = 5e-4).</li> <li>Deployed scalable, real-time inference on Kubernetes with automated CI/CD pipelines.</li> <li>Integrated structured and semi-structured patient data from relational and NoSQL sources for comprehensive healthcare analytics.</li> </ul> </div> <div> <div>Social Recommendation System -(Graph Neural Networks, PyTorch Geometric, Neo4j)</div> <ul style="list-style-type: none"> <li>Improved recommendation precision by 22% by implementing a GraphSAGE-based recommendation engine.</li> <li>Reduced query execution time by 35% through Neo4j integration and optimized graph storage.</li> <li>Deployed scalable real-time inference via FastAPI with graph-based data modeling for performance and scalability.</li> </ul> </div> </div> </div>	
<div> <div>Education</div> <div> <div>M.S. Computer Science  Indiana State University, Terre Haute, IN, US</div> <ul style="list-style-type: none"> <li>Courses: Statistics, Artificial Intelligence, Matrices, Data Mining, Data Visualization (Tableau/Power BI)</li> </ul> </div> <div> <div>B.S. Computer Science  Nirma University, Ahmedabad, GJ, India</div> <ul style="list-style-type: none"> <li>Courses: Deep Learning, Big Data Analytics (Hadoop), Applied Statistics, Calculus, Financial Management, Supply Chain</li> </ul> </div> </div>	
<div> <div>Future Goals</div> <div> <ul style="list-style-type: none"> <li>Expand inference expertise into State Space Models (SSMs) to complement Transformer-based pipelines.</li> <li>Contribute to hybrid model architectures and custom inference engines for real-time multimodal intelligence.</li> </ul> </div> </div>	