

# Zhiheng Jiang

310-709-5870 | zhjiang@ucla.edu | linkedin.com/jiangzhiheng | github.com/jzh001 | Google Scholar

## Education

<b>University of California, Los Angeles</b>	Sep 2024 - June 2028
Bachelor of Science in Computer Science and Engineering (GPA: 4.00 / 4.00)   Los Angeles, CA	
<ul style="list-style-type: none"><li><b>Courses:</b> Machine Learning, Linear Algebra, Probability and Statistics, Multivariable Calculus, Data Structures and Algorithms, Operating Systems, Digital Design</li><li><b>Awards:</b> 2025 UCLA Samueli Undergraduate Scholarship Recipient, Dean's Honors List</li></ul>	

## Experience

<b>UCLA Computational Machine Learning Lab</b>	Oct 2025 – Present
<i>Undergraduate Researcher</i>	Los Angeles, CA
<ul style="list-style-type: none"><li>Working under Prof. Cho-Jui Hsieh to explore methods for LLM compression and denoising with Singular Value Decomposition</li><li>Benchmarked layer-adaptive SVD-based LLM model compression approaches on factual+logical reasoning and safety alignment</li><li>Devised novel finetuning-free LLM layer editing method which outperformed existing SVD-based denoising approaches</li><li>Performed subspace clustering of LLM weight matrices using gradient-based heuristics to improve post-training SVD reduction</li></ul>	
<b>Institute for Creative Technologies, University of Southern California</b>	
June 2025 – Sep 2025   Playa Vista, CA	
<i>Visiting Academic</i>	
<ul style="list-style-type: none"><li>First author in reinforcement learning paper, partnering with U.S. DEVCOM Army Research Lab, advised by Dr. Volkan Ustun</li><li>Designed benchmark problems for Preference-Driven Multi-Objective Reinforcement Learning with novel evaluation metrics</li><li>Designed Heterogeneous Graph Neural Network with flexible preference conditioning for Proximal Policy Optimization (PPO)</li><li>Authored a Python toolkit (GraphAllocBench) using Gymnasium, Stable Baselines3, PyTorch Geometric, Wandb and PyMOO</li><li>Pending review at top Machine Learning conference, presented at SoCal CS REU Symposium (Harvey Mudd College)</li></ul>	
<b>Structures-Computer Interaction Lab at UCLA</b>	Oct 2024 – June 2025
<i>Undergraduate Researcher</i>	Los Angeles, CA
<ul style="list-style-type: none"><li>Developed digital twin simulation of springs using Large Language Models and Robot Tool Calling, advised by Prof. Khalid Jawed</li><li>Refined numerical optimization methods using PyTorch Differential Equations (torchdiffeq) for reduced-order spring simulations</li><li>Trained Physics-Informed Neural Networks with Autograd and Discrete Differential Geometry to predict spring elastic energy</li><li>Worked with scientist from Amazon AWS to develop LangGraph clone from scratch to manage LLM multi-agent LLM memory</li><li>Developed LLM tool wrappers for OpenGL physics simulations, Intel RealSense video input (OpenCV) and Sawyer Robots (ROS)</li></ul>	
<b>Institute of High Performance Computing, A*STAR</b>	June 2020 – Feb 2022
<i>Research Assistant</i>	Singapore
<ul style="list-style-type: none"><li>First-author journal publication on large social network analysis with Dr. Hoai Nguyen Huynh, cited by Spotify Research (2024)</li><li>Performed K-means, Girvan-Newman and Shannon Entropy clustering on scale-free user-oriented networks with NetworkX</li><li>Improved community detection algorithms on large social networks parsed with BeautifulSoup from an AJAX review database</li><li>Performed text mining using Natural Language Tool Kit (NLTK), TF-IDF and dependency parsing to describe genre communities</li></ul>	

## Publications

- [1] Zhiheng Jiang, Yunzhe Wang, Ryan Marr, Ellen Novoseller, Benjamin T. Files and Volkan Ustun, "GraphAllocBench: A Flexible Benchmark for Preference-Conditioned Multi-Objective Policy Learning," *arXiv preprint arXiv:2601.20753*, 2026.
- [2] Zhiheng Jiang and Hoai Nguyen Huynh, "Unveiling music genre structure through common-interest communities" *Social Network Analysis and Mining*, Vol. 12, No. 35 (2022).

## Projects

<b>Neural Network Diagram Generator with Agentic RAG</b>   <i>LangChain, Llamaindex, Pydantic, Docker, Gradio, Ollama, PostgreSQL</i>	2025
<ul style="list-style-type: none"><li>Created a full-stack containerized Agentic RAG Gradio application to generate professional neural network model diagrams</li><li>Designed multi-agent LLM-RAG workflows with Anthropic and Gemini API, and OpenAI gpt-oss-20B locally with Ollama</li><li>Open-sourced application on HuggingFace Spaces, and hosted text embeddings on Supabase PostgreSQL vector database</li></ul>	
<b>Today I Learnt AI Competition – Advanced Category Champion</b>   <i>LLMs, VLMs, Finetuning, PyTorch, HuggingFace, VertexAI, Docker</i>	
<ul style="list-style-type: none"><li>Team leader of champion team, winning 10,000 SGD (7,500 USD) cash prize competing against 60 university-level finalist teams</li><li>Finetuned large deep learning models with high test scores, for audio (99.5%), Vision-Language Models and object detection (86.3%) and Transformer Question Answering (99.9%), with quantization, to achieve high inference speeds on a DJI robot</li><li>Finetuned SOTA models such as YOLO, DETR, RoBERTa and OpenAI Whisper on VertexAI and Google Cloud Platform (GCP)</li></ul>	

## Technical Skills

**Languages:** Python, C++, MATLAB, Java, Javascript, SQL, LaTeX, Shell Scripting (Bash), YAML, JSON, HTML  
**Technologies:** PyTorch, HuggingFace, LangGraph, LangChain, Llamaindex, Pandas, Numpy, ReactJS, NodeJS, FastAPI, GitHub Copilot  
**Concepts:** Large Language Models, Machine Learning, Computer Vision, Natural Language Processing, Reinforcement Learning