

# Ryan Tabrizi

☎ (650) 804-6468 ✉ rtabrizi@berkeley.edu 🔗 linkedin.com/in/rtabrizi 🧑 ryantabrizi.com

## Education

<b>University of California, Berkeley</b>	<b>May 2025</b>
<i>B.S. in Electrical Engineering and Computer Science</i>	<i>GPA: 3.80</i>
<ul style="list-style-type: none"><li>• <b>Graduate Coursework:</b> Vision-Language Models (In Progress), Computer Vision (A), NLP (A), Seminar: From Research to Startups (A)</li><li>• <b>Undergraduate Coursework:</b> Intro to ML (A+), CV and Computational Photography (A), Convex Optimization (A), Operating Systems (A), Computer Architecture (A), Data Structures (A), Electronic System Design (A), Linear Algebra and Differential Equations (A)</li><li>• <b>Distinctions and Leadership:</b> EECS Honors Program (Neuroscience Concentration), Eta Kappa Nu Honor Society, Dean's Honors List (4×), Berkeley AI Research, Launchpad Creative ML (Head of Education), Computer Science Mentors</li></ul>	

## Publications

<b>An LLM Compiler for Parallel Function Calling</b> [paper][code]	<b>ICML 2024</b>
Sehoon Kim*, Suhong Moon*, <b>Ryan Tabrizi</b> , Nicholas Lee, Michael W. Mahoney, Kurt Keutzer, Amir Gholami	
LLMCompiler: An open-source framework for efficiently executing multiple tasks in parallel. LLMCompiler manages dependencies and executes user-provided APIs, achieving up to 1.3x speedup over OpenAI's Parallel Function Calling, as well as up to 9% higher accuracy compared to ReAct.	
<b>TinyAgent: Function Calling at the Edge</b> [paper][code]	<b>EMNLP Demo 2024</b>
Lutfi Erdogan*, Nicholas Lee*, Siddharth Jha*, Sehoon Kim, <b>Ryan Tabrizi</b> , . . . , Kurt Keutzer, Amir Gholami	
On-device deployment of fine-tuned small language models (<7B) to execute everyday user tasks on MacOS with LLMCompiler backend.	

## Experience

<b>Berkeley AI Research - Computer Vision Group</b>	<b>May 2024 - Present</b>
<i>Researcher</i>	<i>Berkeley, CA</i>
<ul style="list-style-type: none"><li>• Using REINFORCE and PPO for hyperparameter optimization advised by Prof. Angjoo Kanazawa; trained actor-critic network to automatically predict optimal Gaussian Splatting parameters for image fitting, eliminating manual tuning for each image.</li><li>• Led refactoring of Nerfstudio gsplat library, defining common API for different Gaussian Splatting methods to simplify splatting pipeline.</li></ul>	
<b>Berkeley AI Research - PALLAS Group</b>	<b>Sep 2023 - Sep 2024</b>
<i>Researcher</i>	<i>Berkeley, CA</i>
<ul style="list-style-type: none"><li>• Built LLMCompiler (ICML 2024, 1.5K+ GitHub stars), an open-source parallel function calling framework for LLMs, with Prof. Kurt Keutzer.</li><li>• Achieved SOTA success rate on WebShop benchmark using LLMCompiler to call web agents in parallel for efficient environment exploration.</li><li>• Developed TinyAgent (EMNLP 2024 Demo), utilizing LLMCompiler for application-specific function-calling on small language models (&lt;7B).</li></ul>	
<b>UC Berkeley EECS</b>	<b>Aug 2024 - Present</b>
<i>Teacher Assistant</i>	<i>Berkeley, CA</i>
<ul style="list-style-type: none"><li>• Paid 10-hour TA for CS 180/280A “Computer Vision and Computational Photography,” working alongside Prof. Alyosha Efros.</li><li>• Designed diffusion visual illusions project with DDPM and rectified flow in collaboration with UMich Prof. Andrew Owens.</li><li>• Conduct weekly office hours and project parties, and grade projects and exams for 300+ students as part of 8-person course staff.</li><li>• Incoming 10-hour TA for CS 280 “Graduate Computer Vision” with Prof. Angjoo Kanazawa and Prof. Jitendra Malik.</li></ul>	
<b>Insitro - Core Imaging Team</b>	<b>May 2023 - Aug 2023</b>
<i>Research Scientist Intern</i>	<i>San Francisco, CA</i>
<ul style="list-style-type: none"><li>• Developed IM2MEA, a system that aligns 6k+ microscopy images with 2.5k+ microelectrode array (MEA) recordings to retrieve activation potential signals for target validation of neurodegenerative diseases.</li><li>• Achieved cross-modal mapping by utilizing META AI's DINO segmentation model on a dataset of 5k image-MEA data points to impute MEA aggregate metrics and image embedding centroids.</li><li>• Trained a transformer encoder-decoder on 2k MEA samples to reconstruct action potentials for downstream signal synthesis by encoding electrode-level and network-level MEA features.</li></ul>	
<b>NASA JPL - ML and Instrument Autonomy Team</b>	<b>May 2022 - Sep 2022</b>
<i>Research Scientist Intern</i>	<i>Palo Alto, CA</i>
<ul style="list-style-type: none"><li>• Headed research and trade study of 5 image compression algorithms, including the ICER Progressive Wavelet Image Compressor and JPEG 2000, with 500+ preprocessed images across Mars, Miranda, and Ceres mission data.</li><li>• Designed compression algorithm analysis pipeline with Python wrappers for C codebases used in deep space exploration, utilizing structural similarity indexes and heat maps as benchmark metrics.</li><li>• Enabled and streamlined optimization of compression algorithms through pipeline's analytics for future missions within JPL's Machine Learning and Instrument Autonomy Group and NASA at large.</li></ul>	

## Projects

<b>VeggieWorld</b>   <i>Nerfstudio, gsplat, Blender</i>	<b>Apr 2024</b>
<ul style="list-style-type: none"><li>• ‘Veggiefied’ 3D Gaussian Splatting scene reconstructions using vegetable 3D asset library for 3D style transfer.</li><li>• Compared CLIP ViT and human performance to asses veggiefied scene classification under different vegetable sizes and distribution densities.</li></ul>	
<b>Phodexr</b>   <i>CLIP, BERT, ResNet, Hugging Face, WandB</i>	<b>Sep 2021 - Dec 2021</b>
<ul style="list-style-type: none"><li>• Recreated Apple's camera roll semantic search with OpenAI's CLIP for lightweight album indexing from user text queries.</li><li>• Pre-trained a DistilBERT tokenizer and ResNet50 on COCO Captions dataset (300K+ image-caption pairs) for image retrieval.</li></ul>	

## Skills

**Programming Languages:** Python, Java, C, L<sup>A</sup>T<sub>E</sub>X, Scheme (LISP), RISC-V  
**Libraries/Frameworks:** PyTorch, NumPy, Hugging Face, Nerfstudio, GSplat, LangChain, TensorFlow, Keras, scikit-learn, Pandas  
**Developer Tools:** Git, AWS, Linux, WandB, Streamlit, Gradio, Isaac Gym, STM32Cube  
**Spoken Languages:** English (native), Persian (B1), Spanish (B1)