

Arthur Pogosian

arthurmpogosian@gmail.com | [LinkedIn](#) | [Github](#)

EDUCATION

University of Pennsylvania	Philadelphia, PA
<i>Masters of Science in Computer Science (GPA 4.00)</i>	Dec 2025
<i>Bachelor of Arts in Mathematics, Computer Science</i>	May 2025

EXPERIENCE

Axiom Math	Jan 2026
<i>Research Engineer Intern</i>	<i>San Francisco, CA</i>
• Working with ex-FAIR/MSL researchers to create autoformalization systems that push the SOTA in math.	
Cambridge AI Safety Hub	Dec 2025 – Feb 2026
<i>Research Fellow</i>	<i>Remote</i>
• Selected for MARS 4.0, researching with <u>Peter Hase</u> on the phenomenon of backtracking in reasoning models.	
Pinterest	May 2025 – Aug 2025
<i>Machine Learning Engineer Intern</i>	<i>San Francisco, CA</i>
• Designed and implemented novel transformer-based user multi-embedding models for rich multi-interest signals.	
• Improved recall@10 by 2x from single embedding baseline and 25% over previous SOTA PinnerSage baseline.	
• Researched novel methods for compressing super-long user history sequences for more accurate personalization.	
• Deployed model in production and integrated in Homefeed for 500 million MAU, led ablation studies on features.	
• Collaborated to finetune LLMs for recommender systems to predict user interest on low activity users.	
UPenn Brachio Lab	July 2024 – Oct 2025
<i>Graduate Machine Learning Research Assistant</i>	<i>Philadelphia, PA</i>
• Working with Eric Wong at the Brachio Lab on <u>feature extraction benchmarks</u> in LLMs and transformers.	
• Researching jailbreaking in LLM systems using theoretical bases, extending the <u>LogicBreaks</u> paper.	
Affirm	May 2024 – Aug 2024
<i>Software Engineer Intern</i>	<i>New York, NY</i>
• Developed scoped user sessions to help prevent millions in fraud, improving security and checkout retention.	
• Built new features to extend login session infrastructure and support the Apple Pay integration at Affirm.	
UPenn Distributed Systems Lab	Jan 2024 – June 2024
<i>Research Assistant</i>	<i>Philadelphia, PA</i>
• Worked on Fast Succinct Non-Interactive Zero-Knowledge Regex Proofs. <u>Publication</u> under Sebastian Angel.	
Reunion LLC	Mar 2022 – May 2024
<i>Founding Engineer</i>	<i>Remote</i>
• Successfully launched pixel-art MMORPG alongside two cofounders with over 1000 players and 100k revenue.	

PROJECTS

Trajectory Steering Vectors <i>PyTorch, TransformerLens</i>	
• Investigated process-level steering method, revealing model biology and achieving 54% jailbreak success.	
Mitigating Adversarial Prompt Perturbations with COT in LLMS <i>PyTorch, Jupyter Notebook</i>	
• Developed gradient-free adaptive CoT defense restoring 97% baseline accuracy on Nova Pro across 48k prompts.	
Restricted GeoCLIP <i>PyTorch, Jupyter Notebook, Matplotlib, Pandas, Numpy</i>	
• Built novel dataset of 117,000 UK street images, finetuned CLIP with 27% accuracy for the Top 50 UK cities.	
GoogLite <i>Java</i>	
• Full search engine completely from scratch; Web server, distributed KVS, Spark, Crawler, Indexer, PageRank.	

TECHNICAL SKILLS

Languages: Python, Java, Kotlin, Rust, C, C++, SQL
Technologies: TensorFlow, PyTorch, Pandas, Matplotlib, PyTest, NumPy, Scikit-learn, Git, Docker
Concepts: Machine Learning, Deep Learning, NLP, Bayesian Optimization, Recommendation Systems, LLMs