

## EDUCATION

<b>Cornell University</b> <i>Bachelor, Computer Science</i> <ul style="list-style-type: none"><li>● <b>GPA: 3.97 Honors:</b> TA for <b>Grad</b> Machine Learning, Ex-President of Cornell Data Science, Teradata Analytics Challenge 1<sup>st</sup> Place</li></ul>	Ithaca, New York December 2024
--	-----------------------------------

## WORK EXPERIENCE

<b>Adobe</b> <i>Data Science Engineer</i> <ul style="list-style-type: none"><li>● Targeting and recommendation systems for Adobe Acrobat &amp; GenAI subscriber growth.</li><li>● ML for causal inference: conditional ATE models for optimal treatment assignment.</li></ul>	San Jose, California Feb 2025 - Present
--	--

<b>Cornell University</b> <i>Researcher</i> <ul style="list-style-type: none"><li>● Currently investigating intra-GPU caching to speed up data center LLM inference &amp; curriculum training multi-hop retrieval LLMs.</li></ul>	Ithaca, New York Feb 2023 - Present
--	--

<b>Adobe</b> <i>Data Scientist Intern</i> <ul style="list-style-type: none"><li>● Developed end-to-end and <b>productionized</b> subscription likelihood <b>prediction model</b>, enabling targeted discounts and pop-ups for <b>4.5M Acrobat users</b>. Implemented product variants and A/B tests (estimated <b>\$1M</b> ARR increase), launched 2024 Q4.</li><li>● Identified <b>1M</b> related users with graph algorithms, enabling <b>recommendations</b> for engagement and upsell worth <b>\$100k</b> ARR.</li><li>● Orchestrated compute clusters and set up model deployment and performance monitoring systems (Airflow, Databricks).</li><li>● Improved product-usage compute logic for 1B Photoshop events, reducing compute time from days to hours (Azure, Spark).</li></ul>	San Jose, California May 2024 - Aug 2024
--	---

<b>ArXiv</b> <i>Research Engineer</i> <ul style="list-style-type: none"><li>● Developed <b>classifiers</b> to tag research paper submissions categories for Cornell's arXiv platform (<b>4M</b> monthly active users).</li><li>● Fine-tuned LLMs to encode text corpuses for document <b>search</b> (3% improvement in first search result compared to Elasticsearch).</li><li>● Improved ROME's (search algorithm) fact editing algorithm to utilize caching, decreasing average <b>query response time</b> by 20%.</li></ul>	Ithaca, New York Feb 2024 - Dec 2024
---	---

<b>Bank of America</b> <i>Machine Learning Intern (Quantitative Summer Analyst)</i> <ul style="list-style-type: none"><li>● Built automated hallucination <b>evaluation infrastructure</b> for chatbot with <b>40M users</b> and designed <b>out-of-distribution detection</b> system for questions, increasing helpfulness by 30%. Business unit estimated <b>\$1M</b> savings, work featured at July Townhall.</li><li>● Tuned chatbot training objective for a 3% improvement in top 25 customer <b>queries</b> and 20% improvement in 10 hardest requests.</li><li>● Researched chatbot-hallucination's sensitivity to paraphrasing, eliminating 40% of hallucination while retaining 90% of truth.</li></ul>	Charlotte, North Carolina Jun 2023 - Aug 2023
--	--

## PUBLICATIONS

<b>PhantomWiki: Generating Reasoning and Retrieval Datasets On-Demand</b> (ICML 2025) <i>A. Gong, C. Wan, K. Stankeviciute, A. Kabra, J. Lee, R. Thesmar, J. Klenke, C. Gomes, and K. Q. Weinberger</i> <ul style="list-style-type: none"><li>● A synthetic dataset generation pipeline for multi-step LLM reasoning across multiple data-sources to address data contamination.</li><li>● Implemented Agentic and RAG LLMs for evaluation, built knowledge-graph to dataset generation pipeline (PyTorch, vLLM, HF).</li></ul>	
<b>Towards Safe and Ethical AI</b> (Global Review of AI Community Ethics, 2025 Vol. 3. No 1) <i>J. Lee and D. Lee</i> <ul style="list-style-type: none"><li>● Surveyed and analyzed benchmarks for evaluating bias and hate of LLMs, identifying systemic weaknesses and scaling issues.</li></ul>	

## PROJECTS

<b>Web Search Engine For Math Equations</b> ( <a href="#">Github</a> , advised by Prof. Kilian Weinberger) <i>Project Lead</i> <ul style="list-style-type: none"><li>● Trained equation detection model (YOLO), finetuned CNN with contrastive loss for clustered vector embeddings (PyTorch).</li><li>● Built NoSQL database (AWS DynamoDB) to store user uploaded PDFs; exploring vectorized search and Redis for query caching.</li><li>● Developed REST APIs between backend AWS Lambda endpoints and frontend web app to send user search queries.</li></ul>	Ithaca, New York Aug 2022 - May 2023
--	---

## TECHNICAL SKILLS

**Languages:** C++, Java, Python, SQL, C, Bash, Shell, OCaml, JavaScript / TypeScript, HTML, CSS, PHP  
**Frameworks and Cloud:** Pytorch, Tensorflow, Azure, AWS, Spring Boot, Flask, Django, React, Vue, D3.js, Scikit-learn, Pandas  
**Tools and Database:** Spark, Docker, Databricks, Airflow, MySQL, DynamoDB, MongoDB, Cassandra, Git, GitHub, Linux