# Jiashu Xu 你家澍

## 🛿 Cambridge, MA | 🛛 949-522-2936 | 🖂 jxu1@g.harvard.edu | 🕲 Website | 🎓 Scholar | 🗘 cnut1648 | 🖬 jiashu-xu

### EDUCATION

Harvard University	Cambridge, USA
M.S. in Computational Science and Engineering; cross-registered at MIT; GPA: 4.0/4.0	Fall 2022 – Spring 2024
University of Southern California	Los Angeles, USA
B.S. double major in Applied Math & Computer Science; Summa Cum Laux; GPA: 3.97/4.0	Fall 2020 – Spring 2022
University of California, Irvine	Irvine, USA
B.S. double major in Applied Math & Computer Science; GPA: 3.98/4.0	Fall 2018 – Spring 2020
Hong Kong University of Science and Technology	Hong Kong, China
UCEAP summer visiting student; study robotics; GPA: 4.0/4.0	Summer 2019

#### AWARDS

**CURVE** Research Fellowship Jennifer Battat Scholarship USC Transfer Merit Scholarship USC Academic Achievement Award USC & UCI Dean's List

\$1250/semester Research Stipend \$3.5k for Mathematics Major Half-tuition Merit Scholarship for 2% Of Transfer Applicants Double Major with 3.75+ GPA All Semesters

#### **Research Interests**

My current research interests lie in **Reliable AI**. Particularly,

- 1. AI Safety. For example, defending against malicious exploitation of LLM vulnerabilities ([3], [4]), protecting open-sourced LLM ownership ([1]), and aligning LLMs with human preferences ([2]).
- 2. Training AI that excels in low-resource regimes, through indirect supervision ([7], [11]) or synthetic data generation ([5], [6], [8], [10]).
- 3. Explanation and model learning from explanation ([12] to [14]).

### PUBLICATIONS & SERVICES

[1]	Training Large Language Models as Reward Models		
	Jiashu Xu, Daniel Pressel, Prasoon Goyal, Luke Dai, Reza Ghanadan, Michael Johnston		
	COLM, 2024 (To be submitted)		
[2]	Instructional Fingerprinting of Large Language Models		
	Jiashu Xu, Fei Wang*, Mingyu Derek Ma*, Pang Wei Koh, Chaowei Xiao, Muchao Chen		
	NAACL, 2024 (Oral) code	paper	project
[3]	Instructions as Backdoors: Backdoor Vulnerabilities of Instruction Tuning for I	arge	
	Language Models		
	Jiashu Xu, Mingyu Derek Ma, Fei Wang, Chaowei Xiao, Muhao Chen		
	NAACL, 2024	paper	project
[4]	Test-time Backdoor Mitigation for Black-Box Large Language Models with Def	ensive	
	Demonstrations		
	Wenjie Mo, <b>Jiashu Xu</b> , Qin Liu, Jiongxiao Wang, Jun Yan, Chaowei Xiao, Muhao Chen		
	COLM, 2024 (To be submitted)		paper
[5]	BEHAVIOR Vision Suite: Customizable Dataset Generation via Simulation		
	Yunhao Ge*, Yihe Tang*, Jiashu Xu*, Cem Gokmen*, Chengshu Li, Wensi Ai, Benjamin	Jose M	[artinez,
	Arman Aydin, Mona Anvari, Ayush K Chakravarthy, Hong-Xing Yu, Josiah Wong, Sanjana	ι Srivast	tava,
	Sharon Lee, Shengxin Zha, Laurent Itti, Yunzhu Li, Roberto Martín-Martín, Miao Liu, Pen	gchuan	Zhang,
	Ruohan Zhang, Fei-Fei Li, Jiajun Wu		
	CVPR, 2024 (Highlight) code	paper	project

\*=Equal Contribution

[6]	DreamDistribution: Prompt Distribution Learning for Text-to-Image Diffusion Models Brian Nlong Zhao, Yuhang Xiao*, Jiashu Xu*, Xinyang Jiang, Yifan Yang, Dongsheng Li, Laurent	Itti,	
	Yunhao Ge, Vibhav Vineet		
	ECCV, 2024 (Under Review) code paper	project	
[7]	an NLI Provide Proper Indirect Supervision for Low-resource Biomedical Relation		
	Extraction?		
	Jiashu Xu, Mingyu Derek Ma, Muhao Chen		
	ACL, 2023 ( <b>Oral</b> ) code	paper	
[8]	Dall-e for detection: Language-driven context image synthesis for object detection Yunhao Ge <sup>*</sup> , Jiashu Xu <sup>*</sup> , Brian Nlong Zhao, Neel Joshi, Laurent Itti, Vibhav Vineet		
[0]	arXiv, 2022 Code paper ext	ension	
[9]	X-INORM: Exchanging INORmalization Parameters for Bimodal Fusion		
	Kureng Yin", Jiashu Au", Hanxin Zu, Monainmad Soleymani		
[10]	Neural Sime Learning to Concrete Training Data with NoBE	paper	
[10]	Yunhaa Ca Harkirat Bahl* Linghu Yu* Suriya Cunagakar Neel Jachi Vala Song Vin Wang Lau	ront	
	Itti Vibbay Vinoot	rem	
	ECCV 2022	naper	
[11]	Unified Semantic Typing with Meaningful Label Inference	paper	
[11]	James V Huang Bangzheng Li* <b>Jiashu Xu</b> * Muhao Chen		
	NAACL 2022	naper	
[19]	Dissection Cesture Sequence during Nerve Sparing Predicts Erectile Function Recovery	after	
[12]	Robot-Assisted Radical Prostatectomy	anter	
	Bunzhuo Ma <b>Jiashu Xu</b> Ivan Bodriguez Gina DeMeo Aditya Desai Loc Trinh Jessica H Nguye	n	
	Anima Anandkumar, Jim C. Hu, Andrew J. Hung	,	
	NPJ Digital Medicine, 2022	paper	
[13]	Dissection Assessment for Robotic Technique (DART) to Evaluate Nerve-Spare of	<b>F</b> • <b>F</b> •	
[-0]	Robot-Assisted Radical Prostatectomy		
	Runzhuo Ma, Alvin Hui, Jiashu Xu, Aditva Desai, Michael Tzeng, Emily Cheng, Loc Trinh, Jessica	ıH.	
	Nguyen, Anima Anandkumar, Jim C. Hu, Andrew J. Hung		
	American Urological Association Annual Conference (AUA), 2022	paper	
[14]	SalKG: Learning From Knowledge Graph Explanations for Commonsense Reasoning		
	Aaron Chan, Jiashu Xu, Boyuan Long, Soumya Sanyal, Tanishq Gupta, Xiang Ren		
	NeurIPS, 2021 code	paper	
Revi	ewer Service: ACL Rolling Review, ACL 2023, EMNLP 2023, CVPR 2022		

#### **Research Experience**

#### **NVIDIA Research**

(Incoming) Research Scientist Intern | Manager: Ming-Yu Liu

• Plan to work on LLM.

# Amazon Science

 $\label{eq:applied} Applied \ Scientist \ Intern \ | \ Manager: \ Daniel \ Pressel, \ Michael \ Johnston$ 

- Collaborated closely with the LLM team on the reward modeling side.
- Finetuned LLMs directly as reward models such that models learn to align with human preferences implicitly. Further benefits included zero-shot generalization to unseen dimensions and domains, high-quality data filtering, rationale generation to explain decisions, and synthetic conversation curation for AI self-improvement (RLAIF) [1].

# USC LUKA Group

Research Assistant | Advisor: Prof. Muhao Chen, Prof. Chaowei Xiao

• Proposed a fingerprinting method to safeguard open-source LLM ownership via memorizing fingerprint instances. Such lightweight fingerprint persists large-scale user finetuning on arbitrary datasets, is robust to fingerprint guessing and various PEFT training methods, and supports multi-stage fingerprinting akin to MIT License [2].

Santa Clara, USA Summer 2024

> New York, USA Summer 2023

Los Angeles, USA

Fall 2021 – Present

- Investigated backdoor vulnerabilities of instruction-tuned LLMs that have high backdoor success with minimal malicious instructions, can generalize to multiple tasks, and cannot be cured by continual learning [3]. And proposed leveraging clean in-context demonstrations as effective test-time defense against various backdoor attacks [4].
- Proposed indirect supervision to borrow supervision signals from resource-rich tasks to enhance resource-limited tasks: cross-domain transfer general domain NLI knowledge to improve low-resource biomedical Relation Extraction [7]; cross-task transfer semantic typing knowledge to handle large label space inference [11].

# Harvard AI4LIFE Group

Research Assistant | Advisor: Prof. Himabindu Lakkaraju

- Integrated dynamic knowledge graph, constructed with a language model agent dynamically, into inference to improve factuality (In Progress).
- Investigated mechanistic interpretability on LLaVA-like vision language models to investigate how models react to complex queries such as referring expressions and object counting (In Progress).

# Stanford SVL Group

Research Assistant | Advisor: Prof. Jiajun Wu

• Developed BEHAVIOR Vision Suite, a customizable dataset generator featuring photorealistic assets and physically plausible annotations. Demonstrated applications include holistic benchmarks for 2D and 3D vision models, robustness evaluation through parametric out-of-distribution evaluation (e.g. low lighting, extreme camera pose), and synthetic dataset generation to bolster performance in low-resource scenarios [5].

## **Microsoft Research**

Student Collaborator | Manager: Prof. Laurent Itti, Vibhav Vineet

- Proposed prompt distribution learning for text-to-image and text-to-3D diffusion models to lightweight control image quality and diversity [6].
- Utilized diffusion models and object cut-and-paste to create coherent synthetic training datasets for enhancing low-resource object detection and segmentation [8]. And proposed differentiable synthetic dataset generation with NeRF to improve out-of-distribution object detection of varying views [10].

## Mentoring

Wenjie Mo USC B.S.	Fall 2023 – Present
Brian Nlong Zhao USC B.S. $\rightarrow$ M.S., Research Scientist Intern at Microsoft Research Asia	Fall $2022 - Fall 2023$

# WORKING & TEACHING EXPERIENCE

<ul> <li>Teaching Assistant at USC</li> <li>Role: Office Hours, Discussion Sections, Grading</li> <li>CSCI 567: Graduate level Machine Learning with Prof. Haipeng Luo.</li> <li>MATH 499: Independent Research with Prof. Neelesh Tiruviluamala.</li> </ul>	Los Angeles, USA Spring – Fall 2021
<ul> <li>Teach for Los Angeles</li> <li>Mentor <ul> <li>Tutored middle school students from LA K-12 community 1-on-1 on mathematics for two</li> <li>Inspired students to reach full math potential in preparation for college and STEM careers</li> </ul> </li> </ul>	Los Angeles, USA Spring 2021 hours every week. s.
<ul> <li>Math CEO</li> <li>Mentor</li> <li>Coordinated meetings with Santa Ana middle school students and taught mathematical the</li> </ul>	Irvine, USA Fall 2018 – Spring 2020 ninking.
<ul> <li>Johnson &amp; Johnson</li> <li>Digital &amp; Analytics Data Assistant <ul> <li>Tracked counterfeit or parallel products from various sales channels using NLP techniques labeling and named entity recognition.</li> <li>Devised a context extractor based on Jieba tokenizer and Chinese word vectors.</li> <li>Presented at PCS 2019 medicine CIO summit about the NLP approach for tracking count</li> </ul> </li> </ul>	Shanghai, China Summer 2019 including semantic role erfeit products.
<ul> <li>Wind Information</li> <li>Quantitative Index Research Analyst</li> <li>Collaborated with product managers to launch Wind's new product: Wind Equity Backter multiple prototype algorithms with test codes using python-wind and Pytest.</li> <li>Code-reviewed index-related codes and queried Wind index database to resolve clients' con</li> </ul>	Shanghai, China Spring – Summer 2018 ster and implemented mplaints.

Los Angeles, USA

Spring 2022 - Present

Cambridge, USA

Palo Alto. USA

Fall 2023 – Present

Spring 2023 – Present