

Ashwinee PANDA

WEBSITE

EMAIL

EXPERIENCE

24 - 25	Postdoc	UMD College Park	AI TRAINING	with Tom Goldstein
24	Intern	Capital One	PRETRAINING	
20 - 24	PhD	Princeton	AI SAFETY	with Prateek Mittal
20	B.S./M.S.	UC Berkeley	AI SYSTEMS	with Joey Gonzalez

PUBLICATIONS

Safety	Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek Mittal, Peter Henderson Safety Alignment Should be Made More Than Just a Few Tokens Deep <i>ICLR 2025</i>
Auditing	Ashwinee Panda*, Xinyu Tang*, Milad N., Chris C., Prateek M. Privacy Auditing of LLMs <i>ICLR 2025</i>
DP-ZO	Xinyu Tang*, Ashwinee Panda*, Milad N., Saeed M., Prateek M. Private Fine-tuning of LLMs with Zeroth-order Optimization <i>TPDP 2024 Oral, TMLR 2025</i>
DP-Scaling	Ashwinee Panda*, Xinyu Tang*, Vikash S., Saeed M., Prateek M. A New Linear Scaling Rule for Private Adaptive HPO <i>ICML 2024 Poster</i>
Phishing	Ashwinee Panda, Chris C., Zhengming Z., Yaoqing Y., Prateek M. Teach LLMs to Phish: Stealing Private Information from LLMs <i>ICLR 2024 Poster</i>
DP-ICL	Tong Wu*, Ashwinee Panda*, Tianhao Wang*, Prateek M. Privacy-Preserving In-Context Learning for LLMs <i>ICLR 2024 Poster</i>
DP-RandP	Xinyu Tang*, Ashwinee Panda*, Prateek M. DP Image Classification by Learning Priors from Random Processes <i>NeurIPS 2023 Spotlight</i>
AdvVLM	Xiangyu Qi*, Kaixuan H.*, Ashwinee Panda, Mengdi W., Prateek M. Introducing Vision into LLMs Expands Attack Surfaces <i>AAAI 2024 Oral</i>
Neurotoxin	Zhengming Zhang*, Ashwinee Panda*, Linyue S., Yaoqing Y., Prateek M., Joey G., Kannan R., Michael M. NeuroToxin: Durable Backdoors in Federated Learning <i>ICML 2022 Oral</i>
SparseFed	Ashwinee Panda, Saeed M., Arjun B., Supriyo C., Prateek M. SparseFed: Mitigating Model Poisoning Attacks in FL via Sparsification <i>AISTATS 2022 Poster</i>
FetchSGD	Daniel Rothchild*, Ashwinee Panda*, Enayat U., Nikita I., Ion S., Vladimir B., Joey G., Raman A. FetchSGD: Communication-Efficient Federated Learning with Sketching <i>ICML 2020 Poster</i>

PREPRINTS

DenseMoE	Ashwinee Panda , Vatsal Baherwani, Zain Sarwar, Benjamin Thérien, Stephen Rawls, Sambit Sahu, Supriyo Chakraborty, Tom Goldstein Dense Backpropagation Improves Routing for Sparsely-Gated Mixture-of-Experts <i>NeurIPS 2024 ENSLP/OPT Workshops</i>
Scaling	Sean Michael McLeish, John Kirchenbauer, David Yu Miller, Siddharth Singh, Abhinav Bhatele, Micah Goldblum, Ashwinee Panda , Tom Goldstein Gemstones: A Model Suite for Multi-Faceted Scaling Laws
StructMoE	Zain Sarwar, Ashwinee Panda , Benjamin Thérien, Stephen Rawls, Sambit Sahu, Supriyo Chakraborty StructMoE: Augmenting MoEs with Hierarchically Routed Low Rank Experts <i>NeurIPS 2024 ENSLP Workshop</i>
MoE-CPT	Benjamin Thérien, Charles-Étienne Joseph, Zain Sarwar, Ashwinee Panda , Anirban Das, Shi-Xiong Zhang, Stephen Rawls, Sambit Sahu, Eugene Belilovsky, Irina Rish Continual Pre-training of MoEs: How robust is your router?
Refusal	Neel Jain, Aditya Shrivastava, Chenyang Zhu, Daben Liu, Alfy Samuel, Anoop Kumar, Ashwinee Panda , Micah Goldblum, Tom Goldstein Refusal Tokens: A Simple Way to Control Refusal Messages <i>NeurIPS 2024 SafeGenAI Workshop</i>
T2I Eval	Kevin David Hayes, Micah Goldblum, Vikash Sehwal, Gowthami Somepalli, Ashwinee Panda , Tom Goldstein FineGRAIN: Evaluating Failure Modes of Text-to-Image Models with Vision Language Model Judges
LoTA	Ashwinee Panda , Berivan I., Xiangyu Q., Sanmi K., Tsachy W., Prateek M. Lottery Ticket Adaptation: Mitigating Destructive Interference in LLMs <i>ICML 2024 ES-FoMO Oral, ICML 2024 WANT Best Paper</i>
DP-Diffusion	Vikash S.*, Ashwinee Panda* , Ashwini Pokle, Xinyu Tang, Saeed M., Mung Chiang, J Zico Kolter, Prateek M. Differentially Private Generation of High Fidelity Samples From Diffusion Models ICML 2023 GenAI Workshop

AWARDS

24	OpenAI Superalignment Fast Grant (PI)
24	Far AI Grant (PI)
20	Gordon Wu Fellowship
18	LAUNCH Grand Prize
18	Y Combinator Hackathon First Prize

INVITED TALKS

FEB '25	Expanding Bottlenecks in LLM Scaling <i>AllenAI (AI2)</i>
SEP '24	Lottery Ticket Adaptation <i>Google Federated Learning Seminar</i>
SEP '24	Privacy Auditing of LLMs <i>Google Privacy Seminar</i>
MAY '24	Challenges in Adapting LLMs to Private Data Google Privacy Seminar (click for talk recording)
NOV '23	New Privacy Attacks on Large Language Models <i>Sun Lab, Berkeley</i>
NOV '23	Challenges in Data-Driven Alignment of Large Language Models <i>SPYLab, ETH Zurich</i>
OCT '23	New Directions in Differentially Private Machine Learning <i>Meta CAS</i>
SEP '23	Challenges in Data-Driven Alignment of Large Language Models <i>University of Maryland, College Park</i>
SEP '23	Challenges in Augmenting Large Language Models with Private Data <i>SL² Lab, UIUC</i>
SEP '23	Improving the Privacy Utility Tradeoff in Differentially Private Machine Learning with Prior Information <i>SECRIT, University of Michigan</i>
APR '23	Improving the Privacy Utility Tradeoff in Differentially Private Machine Learning with Public Data <i>Apple</i>
MAR '23	Google Privacy Seminar (click for talk recording) <i>Google</i>
JUN '22	Challenges and Directions in Privacy Preserving Machine Learning <i>Microsoft Research Cambridge</i>
MAY '22	Towards Trustworthy Machine Learning <i>Meta AI</i>
JAN '22	Federated Learning for Forecasting <i>Ohmconnect</i>
NOV '21	Building Federated Learning Systems at Scale <i>Liftoff AI</i>
NOV '21	Practical Defenses Against Model Poisoning Attacks Google (click for talk recording)

SERVICE

Organizing

ICLR 2025 Sparsity in LLMs Workshop (Lead Organizer)

Teaching

2023	Teaching Assistant for COS/ECE 432 at Princeton
2019	Course Staff for CS 189 (Machine Learning) at UC Berkeley
2018	Teaching Assistant for CS 70 and CS 189 at UC Berkeley
2017	Course Staff for CS 70 at UC Berkeley

Peer Reviewing (* denotes Best Reviewer Award)

ICML25 (AC),24*,23,22,21,20,19; NeurIPS24,23*,21,ICLR 24*,23,19,ACL23, TMLR24,AISTATS22, SATML23

Advising

[Sukriti Paul](#), [Kevin Hayes](#), [Pedro Sandoval](#), [David Miller](#), [Sean McLeish](#), [Vatsal Baherwani](#), [Neel Jain](#), [Alex Stein](#), [John Cava](#), [Vivan Madan](#), [Jie Li](#), [Yuxin Wen](#), [Ryan Synk](#), [Monte Hoover](#), [Khalid Saifullah](#), [Juzheng Zhang](#), [John Kirchenbauer](#), [Hongmin Chu](#)