# Ashwinee PANDA

WEBSITE      EMAIL

## EXPERIENCE

| | |
|---|---|
| 24 - 25 | Postdoc    **UMD College Park** with Tom Goldstein |
| 24 | Applied Research    **Capital One**    MIXTURE OF EXPERT PRETRAINING |
| 20 - 24 | PhD    **Princeton University** advised by Prateek M.    AI SAFETY |
| 20 | M.S. + B.S.    **UC Berkeley** advised by Joey G.    AI SYSTEMS |

## RESEARCH (SELECTED 1ST-AUTHOR PUBLICATIONS)

| | |
|---|---|
| LoTA | **Ashwinee Panda**, Berivan I., Xiangyu Q., Sanmi K., Tsachy W., Prateek M. <br> Lottery Ticket Adaptation: Mitigating Destructive Interference in LLMs <br> *ICML 2024 ES-FoMO Oral, ICML 2024 FM-Wild Oral* |
| Auditing | **Ashwinee Panda**$^*$, Xinyu Tang$^*$, Milad N., Chris C., Prateek M. <br> Privacy Auditing of LLMs <br> *ICML 2024 NextGenAISafety Oral* |
| DP-ZO | Xinyu Tang$^*$, **Ashwinee Panda**$^*$, Milad N., Saeed M., Prateek M. <br> Private Fine-tuning of LLMs with Zeroth-order Optimization <br> *TPDP 2024 Oral* |
| DP-Scaling | **Ashwinee Panda**$^*$, Xinyu Tang$^*$, Vikash S., Saeed M., Prateek M. <br> A New Linear Scaling Rule for Private Adaptive Hyperparameter Optimization <br> *ICML 2024 Poster* |
| Phishing | **Ashwinee Panda**, Chris C., Zhengming Z., Yaoqing Y., Prateek M. <br> Teach LLMs to Phish: Stealing Private Information from LLMs <br> *ICLR 2024 Poster* |
| DP-ICL | Tong Wu$^*$, **Ashwinee Panda**$^*$, Tianhao Wang$^*$, Prateek M. <br> Privacy-Preserving In-Context Learning for LLMs <br> *ICLR 2024 Poster* |
| DP-RandP | Xinyu Tang$^*$, **Ashwinee Panda**$^*$, Prateek M. <br> Differentially Private Image Classification by Learning Priors from Random Processes <br> *NeurIPS 2023 Spotlight* |
| Neurotoxin | Zhengming Zhang$^*$, **Ashwinee Panda**$^*$, Linyue S., Yaoqing Y., Prateek M., Joey G., Kannan R., Michael M. <br> NeuroToxin: Durable Backdoors in Federated Learning <br> *ICML 2022 Oral* |
| SparseFed | **Ashwinee Panda**, Saeed M., Arjun B., Supriyo C., Prateek M. <br> SparseFed: Mitigating Model Poisoning Attacks in Federated Learning via Sparsification <br> *AISTATS 2022 Poster* |
| FetchSGD | Daniel Rothchild$^*$, **Ashwinee Panda**$^*$, Enayat U., Nikita I., Ion S., Vladimir B., Joey G., Raman A. <br> FetchSGD: Communication-Efficient Federated Learning with Sketching <br> *ICML 2020 Poster* |

# INVITED TALKS

| | |
|---|---|
| SEP '24 | Privacy Auditing of LLMs<br>*Google Privacy Seminar* |
| MAY '24 | Challenges in Adapting LLMs to Private Data<br>Google Privacy Seminar (click for talk recording) |
| NOV '23 | New Privacy Attacks on Large Language Models<br>*Sun Lab, Berkeley* |
| NOV '23 | Challenges in Data-Driven Alignment of Large Language Models<br>*SPYLab, ETH Zurich* |
| OCT '23 | New Directions in Differentially Private Machine Learning<br>*Meta CAS* |
| SEP '23 | Challenges in Data-Driven Alignment of Large Language Models<br>*University of Maryland, College Park* |
| SEP '23 | Challenges in Augmenting Large Language Models with Private Data<br>*SL$^2$ Lab, UIUC* |
| SEP '23 | Improving the Privacy Utility Tradeoff in Differentially Private Machine Learning with Prior Information<br>*SECRIT, University of Michigan* |
| APR '23 | Improving the Privacy Utility Tradeoff in Differentially Private Machine Learning with Public Data<br>*Apple* |
| MAR '23 | Google Privacy Seminar (click for talk recording)<br>*Google* |
| JUN '22 | Challenges and Directions in Privacy Preserving Machine Learning<br>*Microsoft Research Cambridge* |
| MAY '22 | Towards Trustworthy Machine Learning<br>*Meta AI* |
| JAN '22 | Federated Learning for Forecasting<br>*Ohmconnect* |
| NOV '21 | Building Federated Learning Systems at Scale<br>*Liftoff AI* |
| NOV '21 | Practical Defenses Against Model Poisoning Attacks<br>Google (click for talk recording) |

## RESEARCH (ADVISED AND WORKSHOP PAPERS)

| | |
|---|---|
| Safety | Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek M., Peter Henderson |
| | Safety Alignment Should Be Made More Than Just a Few Tokens Deep |
| AdvVLM | Xiangyu Qi*, Kaixuan Huang*, **Ashwinee Panda**, Mengdi Wang, Prateek M. |
| | Introducing Vision into Large Language Models Expands Attack Surfaces and Failure Implications |
| | At *Thirty-Eighth AAAI Conference on Artificial Intelligence* |
| Phishing | **Ashwinee Panda**, Zhengming Z., Yaoqing Y., Prateek M. |
| | Teach GPT to Phish: Neural Phishing Attacks on Large Language Models |
| | At *40th International Conference on Machine Learning* AdvML Workshop |
| 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 |
| | At *40th International Conference on Machine Learning* GenAI Workshop |
| DP-ICL | **Ashwinee Panda***, Tong Wu*, Tianhao Wang*, Prateek M. |
| | Differentially Private In-Context Learning |
| | At *NAACL 2023* TrustNLP Workshop |
| SoftPBT | **Ashwinee Panda**, Eric Liang, Richard Liaw, Joey G. |
| | SoftPBT: Leveraging Experience Replay for Efficient Hyperparameter Schedule Search |
| | *Submitted to NeurIPS 2019* |

## SERVICE

**Teaching**

| | |
|---|---|
| 2023 | Teaching Assistant for COS/ECE 432 at Princeton |
| 2019 | Course Staff for CS 189 (Machine Learning) at UC Berkeley |
| 2018 | Undergraduate Student Instructor for CS 70 (Probability and Discrete Mathematics) and Course Staff for CS 189 at UC Berkeley |
| 2017 | Course Staff for CS 70 at UC Berkeley |

**Peer Reviewing (* denotes Best Reviewer Award)**

| | |
|---|---|
| 2024 | ICML 2024*, NeurIPS 2024 |
| 2023 | SATML 2023, ACL 2023, ICML 2023, NeurIPS 2023*, TMLR |
| 2022 | ICML 2022, AISTATS 2022 |
| 2021 | ICML 2021, NeurIPS 2021 |
| 2020 | ICML 2020 |
| 2019 | ICLR 2019, NeurIPS 2019 |