

Roy Rinberg

www.royrinberg.com
royrinberg+CV@gmail.com
Github: RoyRin

SELECTED WORK AND RESEARCH EXPERIENCE	Harvard University, Cambridge, MA <i>PhD researcher on AI Security, machine unlearning, and Differential Privacy. Advised by Prof. Boaz Barak and Prof. Salil Vadhan. I am supported by a grant from Open Philanthropy.</i>	AUG. 2023 - PRESENT
	MATS (ML Alignment and Theory Scholar), Berkeley, CA <i>Research Scholar at MATS, an independent research program for research on AI safety. Advised by Keri Warr and Nicholas Carlini at Anthropic on detecting and preventing model weight exfiltration.</i>	SUMMER 2025
	Columbia University, New York, NY <i>Master's research on the fundamentals of Differential Privacy and machine learning.</i>	AUG. 2021 - AUG. 2023
	Shelton AI, New York, NY <i>Founding Software Engineer at Shelton AI, a fintech startup that helps pension funds manage investments. I developed core AWS infrastructure for NLP document processing pipeline.</i>	JAN. 2022 - JUN. 2022
	Ouster, San Francisco, CA <i>Early stage software engineer at Ouster, a lidar sensor company. Led on-edge computing development and platforms for evaluating safety algorithms on historical and real-time data.</i>	JUN. 2018 - JUL. 2021
SELECTED PAPERS	AI Security <i>Full list: Google Scholar</i>	
	<ul style="list-style-type: none">• R. Rinberg, A. Karvonen, A. Hoover, D. Reuter, K. Warr. Verifying LLM Inference to Prevent Model Weight Exfiltration. (2025). arXiv preprint.• A. Karvonen, D. Reuter, R. Rinberg, L. Marks, A. Garriga-Alonso, K. Warr. DiFR: Inference Verification Despite Nondeterminism. (2025). arXiv preprint.	
	Machine Unlearning <ul style="list-style-type: none">• R. Rinberg, U. Bhalla, I. Shilov, R. Gandikota. RippleBench: Capturing Ripple Effects by Leveraging Existing Knowledge Repositories. (2025). NeurIPS MechInterp Workshop (<i>Spotlight</i>).• R. Rinberg, P. Puigdemont, M. Pawelczyk, V. Cevher. Data-Unlearn-Bench: Making Evaluating Data Unlearning Easy. (2025). ICML Machine Unlearning for GenAI Workshop.• R. Rinberg, K. Georgiev, S. Park, S. Garg, A. Ilyas, A. Madry, S. Neel. Attribute-to-Delete: Machine Unlearning via Datamodel Matching. (2024). ICLR 2025.	
	Differential Privacy <ul style="list-style-type: none">• R. Rinberg, Ilia Shumailov, Rachel Cummings, Nicolas Papernot. Beyond Laplace and Gaussian: Exploring the Generalized Gaussian Mechanism for Private Machine Learning. Preprint.• F. Boenisch, C Mühl, A. Dziedzic, R. Rinberg, N. Papernot. Have it your way: Individualized Privacy Assignment for DP-SGD. Accepted to Neurips 2023.• F. Boenisch, C Mühl, R. Rinberg, J. Ihrig, A. Dziedzic. Individualized PATE: Differentially Private Machine Learning with Individual Privacy Guarantees. Accepted to PoPETs 2023.	
EDUCATION	Harvard University, Cambridge, MA <i>PhD. Computer Science. Advisors: Prof. Salil Vadhan and Prof. Boaz Barak</i>	2023 - PRESENT
	Columbia University, New York, NY <i>MS in Computer Science [Thesis Track]. Advisors: Prof. Rachel Cummings and Prof. Steven Bellovin</i>	2021 - 2023
	New York University, New York, NY <i>B.A. Computer Science, Physics, Minor: Math.</i>	2014 - 2018
TEACHING AND SERVICE	Teaching: TF for CS2881 (AI Safety, Fall '25); Head TF for CS1200 (Intro to Algorithms, Fall '24); Physics I/II Tutor at NYU ('17-'18) Organizing: Founding Organizer: Technically Private - group of graduate students who work on privacy and security ('21-Present); Co-founder, Project BEST - Education non-profit, Fundraised and grew organization to 25 chapters across 3 states, reaching 3000+ students. ('11-'14) Community Service: Mentor at Mentor Ukraine for 3 students ('22-'23). Advocated Ouster to set up recurring donations to public-interest orgs and paid volunteer days ('18-'20) Academic Service: Reviewer for NeurIPS ('23, '24, '25), ICML ('23, '25), ICLR ('23, '24); Assistant organizer for OSDI '23 PC	