

# Roy Rinberg

[www.royrinberg.com](http://www.royrinberg.com)  
royrinberg+CV@gmail.com  
Github: RoyRin

SELECTED WORK AND RESEARCH EXPERIENCE	<b>Harvard University, Cambridge, MA</b> <i>PhD researcher on AI Security, machine unlearning, and Differential Privacy. Advised by Prof. Boaz Barak and Prof. Salil Vadhan. I am supported by my own grant from Open Philanthropy.</i>	AUG. 2023 - PRESENT
	<b>MATS (ML Alignment and Theory Scholar), Berkeley, CA</b> <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
	<b>Columbia University, New York, NY</b> <i>Master's research on the fundamentals of Differential Privacy and machine learning.</i>	AUG. 2021 - AUG. 2023
	<b>Shelton AI, New York, NY</b> <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
	<b>Ouster, San Francisco, CA</b> <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	<b>AI Security</b> <i>Full list: <a href="#">Google Scholar</a></i>	
	<ul style="list-style-type: none"><li>• <b>R. Rinberg</b>, A. Karvonen, A. Hoover, D. Reuter, K. Warr. <a href="#">Verifying LLM Inference to Prevent Model Weight Exfiltration</a>. (2025). arXiv preprint.</li><li>• A. Karvonen, D. Reuter, <b>R. Rinberg</b>, L. Marks, A. Garriga-Alonso, K. Warr. <a href="#">DiFR: Inference Verification Despite Nondeterminism</a>. (2025). arXiv preprint.</li></ul>	
	<b>Machine Unlearning</b> <ul style="list-style-type: none"><li>• <b>R. Rinberg</b>, U. Bhalla, I. Shilov, R. Gandikota. <a href="#">RippleBench: Capturing Ripple Effects by Leveraging Existing Knowledge Repositories</a>. (2025). NeurIPS MechInterp Workshop (<i>Spotlight</i>).</li><li>• <b>R. Rinberg</b>, P. Puigdemont, M. Pawelczyk, V. Cevher. <a href="#">Data-Unlearn-Bench: Making Evaluating Data Unlearning Easy</a>. (2025). ICML Machine Unlearning for GenAI Workshop.</li><li>• <b>R. Rinberg</b>, K. Georgiev, S. Park, S. Garg, A. Ilyas, A. Madry, S. Neel. <a href="#">Attribute-to-Delete: Machine Unlearning via Datamodel Matching</a>. (2024). ICLR 2025.</li></ul>	
	<b>Differential Privacy</b> <ul style="list-style-type: none"><li>• <b>R. Rinberg</b>, Ilia Shumailov, Rachel Cummings, Nicolas Papernot. <a href="#">Beyond Laplace and Gaussian: Exploring the Generalized Gaussian Mechanism for Private Machine Learning</a>. Preprint.</li><li>• F. Boenisch, C Mühl, A. Dziedzic, <b>R. Rinberg</b>, N. Papernot. <a href="#">Have it your way: Individualized Privacy Assignment for DP-SGD</a>. Accepted to Neurips 2023.</li><li>• F. Boenisch, C Mühl, <b>R. Rinberg</b>, J. Ihrig, A. Dziedzic. <a href="#">Individualized PATE: Differentially Private Machine Learning with Individual Privacy Guarantees</a>. Accepted to PoPETs 2023.</li></ul>	
EDUCATION	<b>Harvard University, Cambridge, MA</b> <i>PhD. Computer Science. Advisors: Prof. Salil Vadhan and Prof. Boaz Barak</i>	2023 - PRESENT
	<b>Columbia University, New York, NY</b> <i>MS in Computer Science [Thesis Track]. Advisors: Prof. Rachel Cummings and Prof. Steven Bellovin</i>	2021 - 2023
	<b>New York University, New York, NY</b> <i>B.A. Computer Science, Physics, Minor: Math.</i>	2014 - 2018
TEACHING AND SERVICE	<b>Teaching:</b> TF for CS2881 (AI Safety, Fall '25); Head TF for CS1200 (Intro to Algorithms, Fall '24); Physics I/II Tutor at NYU ('17-'18) <b>Organizing:</b> 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) <b>Community Service:</b> 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) <b>Academic Service:</b> Reviewer for NeurIPS ('23, '24, '25), ICML ('23, '25), ICLR ('23, '24); Assistant organizer for OSDI '23 PC	