

Sudarsh Kunnavakkam

+1 (949) 254-8232 | Pasadena, CA | kvsudarsh786@gmail.com | github.com/skunnavakkam | sudarsh.com

WORK EXPERIENCE

Research Assistant (Contract) Model Evaluation and Threat Research (METR)	Sep 2023 — Present <i>Berkeley, CA</i>
<ul style="list-style-type: none">Lead engineer for internal project to estimate the agentic time horizon of LLMs at much lower costCo-lead engineer of a state of the art evaluation for Chain-of-Thought Faithfulness of Large Language ModelsHelped lead teams of contractors red-team LLMs and curate datasets such as DAFT Math of difficult, free-response questions	
Undergraduate Research Intern ShapiroLab at Caltech	Nov 2024 — Present <i>Pasadena, CA</i>
<ul style="list-style-type: none">Building better BCIs by engineering towards 10ms response time ultrasound reportersBuilt a high throughput ultrasound screening platform to scale to 1000s of variants per dayDesigned custom proteins with RFDiffusion, AlphaFold, and ESM3 for 10x faster kinetics	
Research Fellow Supervised Program for Alignment Research	Feb 2025 — May 2025 <i>Remote</i>
<ul style="list-style-type: none">Implemented a complex, <i>continuous double auction</i> agent arena as a model environment for LLM collusionBenchmarked emergent collusion between LLMs under various pressuresWork accepted to ICML 2025	
High School Research Intern Lee Nano-Optics Lab at UC Irvine	Dec 2022 — Jun 2024 <i>Irvine, CA</i>
<ul style="list-style-type: none">Scaled 2D ITO fabrication from mm² to multi-cm² sizesDeveloped new transmission matrix method replacing repeated ellipsometryCreated transfer-matrix reverse solver to easily get refractive index information under nonlinear conditions	

EDUCATION

California Institute of Technology <i>B.S. in Physics & Computer Science</i>	Pasadena, CA <i>In progress</i>
University High School <i>High School Diploma</i>	Irvine, CA <i>Sep 2020 — Jun 2024</i>

SELECTED PUBLICATIONS

1. A. Deng*, S. Von Arx*, B. Snodin, S. Kunnavakkam, T. Lanham, “CoT May Be Highly Informative Despite “Unfaithfulness”” by *METR*
2. K. Agarwal, V. Teo, J. Vaquez, S. Kunnavakkam, V. Srikanth, A. Liu, “Evaluating LLM Agent Collusion in Double Auctions” at *ICML 2025 Workshop on Multi-Agent Systems in the Era of Foundation Models*, Vancouver, Canada, July 2025.
3. C. J. Effarah*, T. Chen*, S. Kunnavakkam*, C. M. Gonzalez, H. W. Lee, “Liquid Metal Printed 2D ITO for Nanophotonic Applications,” in *California-US Government Workshop on 2D Materials*, Irvine, California, USA, Sep 2023

PROJECTS

METR: Faithfulness and Monitorability Eval	<u>2025</u>
<ul style="list-style-type: none">Co-lead engineer on METR research report on chain-of-thought (CoT) faithfulness (Aug 2025), extending Anthropic’s seminal evaluation to three frontier models and publishing findings for the wider safety community	
LLM Agent Collusion Arena	<u>2025</u>
<ul style="list-style-type: none">Implemented a continuous double auction system for agentsImplemented oversight, monitors, and other experimental conditions to test influence on collusionAdded logging and metrics with WandBAccepted to ICML 2025 Workshop on Multi-agent Systems	

<u>EM Simulator</u>	<u>2025</u>
<ul style="list-style-type: none"> • Reverse mode differentiable FDFD simulators in Jax for inverse design • Forward and backward diffusion models trained with DDPM and Physics-inspired reward functions to approximate steady state solutions • Implemented fast FDTD for transient events + implemented Fourier Neural Operators for speedup 	
<u>Circuit Simulator</u>	<u>2025</u>
<ul style="list-style-type: none"> • Reverse-mode autodiff for RLC network optimization • Gradient-based optimization for component selection • Works in time domain, as well as just to do component selection • Implemented custom <code>spsolver</code> that is differentiable in JaX 	
<u>Adversarial Attack Using Soft Tokens</u>	<u>2024</u>
<ul style="list-style-type: none"> • Soft-token embedding technique for adversarial text generation • Orthogonal Procrustes Alignment for token mapping • Demonstrated attack generalization across models (PyTorch) 	
Scanning Tunneling Microscope	2024
<ul style="list-style-type: none"> • Built working STM for \$1,000 using open-source design • Achieved atomic-resolution imaging capabilities (Circuit Design, Signal Processing, Mechanical Engineering) 	
<hr/>	
AWARDS	
ARENA 6.0 Attendee	2025
Non-trivial Fellow	2024
Physics Brawl, top 10 US High School Teams	2024, 2023
USACO Silver	2023
AIME Qualifier	2023