

Sudarsh Kunnnavakkam

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

WORK EXPERIENCE

Research Intern Model Evaluation and Threat Research (METR) <ul style="list-style-type: none">Worked on projects to evaluate the agentic time horizon of LLMsCo-lead engineer of a state of the art evaluation for Chain-of-Thought Faithfulness of Large Language ModelsLed team of contractors to red-team LLMs and write custom datasets	Sep 2023 — Present <i>Berkeley, CA</i>
Undergraduate Research Intern ShapiroLab at Caltech <ul style="list-style-type: none">Building better BCIs by engineering towards 10ms response time ultrasound reportersDesigned custom proteins with RFDiffusion, AlphaFold, and ESM3 for 10x faster kinetics	Nov 2024 — Present <i>Pasadena, CA</i>
Research Fellow Supervised Program for Alignment Research <ul style="list-style-type: none">Implemented a complex, <i>continuous double auction</i> agent arena as a model environment for LLM collusion, accepted to <i>ICML 2025</i>	Feb 2025 — May 2025 <i>Remote</i>
High School Research Intern Lee Nano-Optics Lab at UC Irvine <ul style="list-style-type: none">Scaled 2D ITO fabrication from mm² to multi-cm² sizes and developed new transfer-matrix methods for ellipsometry and refractive index characterization. Published at a US Government Workshop.	Dec 2022 — Jun 2024 <i>Irvine, CA</i>

SKILLS

Machine Learning (PyTorch, Jax, Transformers, Diffusion Models, Reinforcement Learning on LLMs, GRPO, PPO, Interpretability), PCB Fabrication, Data Analysis, Signal Processing, Rust, 3D Modeling, Shop Experience, General Wet Lab, Electron Microscopy, AFM, Scanning Probe Microscopy

EDUCATION

California Institute of Technology <i>B.S. in Physics & Computer Science</i>	Pasadena, CA <i>In progress</i>
--	------------------------------------

SELECTED PUBLICATIONS

- A. Deng*, S. Von Arx*, B. Snodin, [S. Kunnnavakkam](#), T. Lanham, “CoT May Be Highly Informative Despite “Unfaithfulness”” by *METR*
- K. Agarwal, V. Teo, J. Vaquez, [S. Kunnnavakkam](#), 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.
- C. J. Effarah*, T. Chen*, [S. Kunnnavakkam](#)*, 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	2025
<ul style="list-style-type: none">A thorough evaluation building on Anthropic’s seminal work on chain-of-thought (CoT) faithfulness, with thorough redteaming throughout.	
LLM Agent Collusion Arena	2025
<ul style="list-style-type: none">A continuous double auction system for agents, oversight, monitors, and other experimental conditions to test influence on collusion, accepted to <i>ICML 2025</i>	
EM Simulator	2025
<ul style="list-style-type: none">Reverse mode differentiable FDFD simulators in Jax for inverse design, with fast FDFD and FDTD through diffusion & neural operators	
Scanning Tunneling Microscope	2024
<ul style="list-style-type: none">Built working STM for \$1,000 using open-source design	
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