# ITAMAR HAGAY PRES

646-589-1150 | presi@umich.edu | www.linkedin.com/in/itamar-pres-3a56a4b3

## EDUCATION

**UNIVERSITY OF MICHIGAN ANN ARBOR, College of Engineering**          Ann Arbor, MI
BSE in Computer Science with minor in Mathematics                                *January 2025*

## RESEARCH PUBLICATIONS

Lee, Bai, **Pres**, et al.(2024) A Mechanistic Understanding of Alignment Algorithms: A Case Study on DPO and Toxicity
*The Forty-first International Conference on Machine Learning* **(Oral: top 1.5 % of submissions)**
- Discovered how toxicity is represented in pre-trained language models (GPT2, Llama2-7b)
- Studied the internal mechanisms of a popular alignment algorithm (Direct Preference Optimization)
- Discovered that DPO does not remove toxicity representations but rather suppresses them
- Uncovered why "aligned" models can be jailbroken

**Pres**, et al.(2024) Towards Reliable Evaluation of Behavior Steering Interventions in LLMs   *The 38th Conference on Neural Information Processing Systems - Workshop on Foundation Model Interventions* **(Oral: Top 6 of total submissions)**
- Established four criteria for effective behavioral steering metrics and highlighted current metric shortcomings
- Developed new steering metric showing over-reported success of control method Contrastive Activation Steering

Park*, Lubana*, **Pres**, et al.(2024) Algorithmic Phases Of In-Context Learning - *Pending Submission at The 13th International Conference on Learning Representations*
- Provided mechanistic analysis showing transformers solve in-context learning with multiple competing algorithms
- Developed quantifiable techniques to measure neuron memorization of training data in transformers

## RESEARCH EXPERIENCE

### Krueger AI Safety Lab Internship (University of Cambridge)          *May 2024 - Present*
*Mentored By **David Krueger***
- Developed framework for multiple constraint satisfaction of LLM generations at inference-time
- Explored generalizability of inference-time interventions such as Contrastive-Activation-Steering
- Incorporated Monte Carlo Steering to optimize decoding quality, efficiency, constraint satisfaction

### University of Michigan Language and Information Technologies Lab          *October 2023 - Present*
*Supervised by **Rada Mihalcea***
- Conducted research in mechanistic interpretability, with a focus on safety alignment algorithms
- Published research findings at ICML
- Researched Vision Language Model interpretability as well as In Context Learning mechanisms

### Mechanistic Interpretability Research (SERI MATS)          *June 2023 – Aug 2023*
*Mentored by **Neel Nanda***
- Analyzed and visualized high-dimensional LLM operations to generate hypothesis of LLM internal mechanisms
- Tested hypotheses of attention head superposition phenomenon in GPT2-XL
- Devised precise casual interventions to validate hypotheses under fact extraction contexts

### Machine Learning Safety Scholars          *June 2022 – August 2022*
- Completed coding, written, and discussion assignments from MIT's Intro to Machine Learning, UM's Deep Learning for Computer Vision, and NYU's Deep Learning courses
- Distilled foundational machine learning papers into readable summaries
- Implemented safety techniques including adversarial and black swan robustness as well as interpretability

### University of Michigan Chandrasekaran Computational Biology Lab          *January 2022 – May 2024*
*Supervised by **Sriram Chandrasekaran***
- Researched effective drug combinations to eliminate antimicrobial resistant bacteria
- Lead project using graph neural networks to model drug interaction social network
- Developed GraphSage encoder to generate drug representations for synergy prediction tasks
- Integrated proteomic, chemogenomic, and structural data into heterogeneous graph

## GRANTS

### Effective Altruism Infrastructure Grant, Full stack Development          *June 2022 – September 2022*
- Developed prototype for web-app that increases transparency of unlikely global risk assessments
- Designed NoSQL MongoDB database structure for back-end
- Created interactive front-end using React.js
- Developed API connecting database to front-end using Node.js and Express.js