

## EDUCATION

---

|   |                                 |
|---|---------------------------------|
| <b>Massachusetts Institute of Technology</b><br>Ph.D. in Computer Science, advised by Nir Shavit  | Cambridge, MA<br>2021 – Present |
| M.Eng. in Computer Science, advised by Gregory W. Wornell, GPA: 5.0/5.0   | 2020 – 2021                     |
| B.Sc. Double Major in Computer Science and Math, GPA: 4.9/5.0   | 2016 – 2020                     |
| – <i>Master's thesis</i> : <a href="#">Adversarial Examples in Simpler Settings</a> .   |                                 |
| – <i>Selected CS coursework</i> : Machine Learning, Inference and Information, Robotic Manipulation, Formal Reasoning about Programs, Cryptography, Compilers, Performance Engineering, Randomized Algorithms, Quantum Computation. |                                 |
| – <i>Selected math coursework</i> : Measure Theoretic Probability, Complex Analysis, Functional Analysis, Differential Geometry, General Relativity, Abstract Algebra.  |                                 |

## PUBLICATIONS

---

1. **Tony T. Wang\***, Adam Gleave\*, Tom Tseng, Kellin Pelrine, Nora Belrose, Joseph Miller, Michael D. Dennis, Yawen Duan, Viktor Pogrebniak, Sergey Levine, Stuart Russell. [Adversarial Policies Beat Superhuman Go AIs](#). NeurIPS 2022 ML Safety Workshop (best paper award, top 10/132); ICML, 2023 (oral, top 10%).
2. **Tony T. Wang\***, Miles Kai Wang\*, Kaivu Hariharan\*, Nir Shavit. [Forbidden Facts: An Investigation of Competing Objectives in Llama 2](#). NeurIPS 2023 ATTRIB and SoLaR Workshops.
3. **Tony T. Wang**, Igor Zablotski, Nir Shavit, Jonathan Rosenfeld. [Cliff-Learning](#). Preprint, 2023.
4. Stephen Casper\*, Xander Davies\*, [and 29 others, including **Tony T. Wang**]. [Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback](#). TMLR, 2023.
5. Simon Alford, Anshula Gandhi, Akshay Rangamani, Andrzej Banburski, **Tony T. Wang**, Sylee Dandekar, John Chin, Tomaso Poggio, Peter Chin. [Neural-guided, Bidirectional Program Search for Abstraction and Reasoning](#). Complex Networks, 2021.
6. Yuheng Bu, **Tony T. Wang**, Gregory W. Wornell. [SDP Methods for Sensitivity-Constrained Privacy Funnel and Information Bottleneck Problems](#). ISIT, 2021.

## WORK AND RESEARCH EXPERIENCE

---

|   |                                      |
|---|--------------------------------------|
| <b>Astra Fellowship, Constellation</b><br>Research Fellow           | Berkeley, CA<br>Jan 2024 – Present   |
| – Working on language model jailbreak defense.                      |                                      |
| <b>Shavit Lab, MIT</b><br>Research Assistant                        | Cambridge, MA<br>Fall 2021 – Present |
| – Working on AI safety, with a focus on adversarial robustness.     |                                      |
| <b>Genesis Therapeutics</b><br>AI Engineer Intern                   | Burlingame, CA<br>Summer 2021        |
| – Worked on deep neural networks for molecular property prediction. |                                      |

## Signals, Information, and Algorithms Laboratory, MIT

Research Assistant (M.Eng.)

Cambridge, MA

Summer 2020 – Spring 2021

- Studied toy examples of adversarial examples to unify different aspects of the phenomenon.
- Collaborated with researchers at the Poggio Lab on neurosymbolic algorithms for solving the Abstraction and Reasoning Corpus.

## Nvidia

AI-Infra Research Intern

Santa Clara, CA

Summer 2019

- Researched active learning for self-driving vision models, with a focus on diversity-aware batch-mode sampling.

## Five Rings Capital

Quant Research Intern

New York City, NY

Q1 2019

- Analyzed market data for statistical arbitrage opportunities.

## Dropbox

Network Reliability Engineering Intern

San Francisco, CA

Summer 2018

- Automated traffic draining for production routers.
- Hacked on [mypyc](#), a compiler from typed Python to Python C extensions.

## DigitalWoven

Software Engineering Intern

San Mateo, CA

Summer 2017

- Built on AWS the serverless backend for [UTStamp](#), a blockchain notary service.
- Designed and implemented the UTStamp frontend in React.

## AWARDS AND GRANTS

---

|   |   |
|---|---|
| <a href="#">Lightspeed Grant</a> for AI safety research, 234000 USD | 2023  |
| Eric and Wendy Schmidt Center PhD Fellowship                        | 2022 - 2023                                 |
| <a href="#">MIT EECS Harold Hazen Teaching Award</a>                | 2021  |
| <a href="#">Undergraduate Teaching Assistant Award</a>              | 2020  |
| USA Computing Olympiad finalist (national top 24)                   | <a href="#">2013</a> , <a href="#">2015</a> |

## OTHER PROJECTS

---

### Roots of Random Polynomials

Fall 2019

*Term project for 18.821, Project Lab in Mathematics*

- Proved roots of high-degree polynomials are roughly uniformly distributed over the unit circle in  $\mathbb{C}$ .
- Report: [web.mit.edu/twang6/public/poly-roots.pdf](http://web.mit.edu/twang6/public/poly-roots.pdf)

### Statistical Inference Through the Lens of Information Geometry

Spring 2019

*Term paper for 18.424, Seminar in Information Theory*

- Contains a proof of the Cramér-Rao bound via information geometry.
- Report: [web.mit.edu/twang6/public/stats-info-geo.pdf](http://web.mit.edu/twang6/public/stats-info-geo.pdf)

### Voice Identification on the VoxCeleb Dataset

Fall 2017

*Term project for 6.867, Machine Learning*

- Compared RNNs to CNNs for performing speaker identification.
- Report: [web.mit.edu/twang6/public/rnn-voxceleb.pdf](http://web.mit.edu/twang6/public/rnn-voxceleb.pdf)

### **Codeforces Round #336**

Q4 2015

*Competitive programming contest*

- Main organizer and problem writer.
- Drew 3000+ participants.
- Particularly proud of authoring [codeforces.com/contest/607/problem/C](https://codeforces.com/contest/607/problem/C).

## **OTHER ACTIVITIES**

---

|  |                |
|--|----------------|
| <b>MIT AI Alignment</b><br>Member, Advisor                             | 2022 – Present |
| <b>MIT Club Tennis</b><br>Member                                       | 2022 – Present |
| <b>MIT Anime Club</b><br>Member, President, Webmaster                  | 2016 – 2021    |
| <b>MIT Chamber Music Society</b><br>Violinist                          | 2016 – 2020    |
| <b>Peninsula Youth Orchestra</b><br>Violinist, Assistant Concertmaster | 2011 – 2016    |