

# SAM TOYER

Berkeley, USA 

[sdt@berkeley.edu](mailto:sdt@berkeley.edu) 

[qxcv.net/research](https://qxcv.net/research) 

“ I am a final-year PhD student at UC Berkeley making language models secure, robust, and safe. Previously I have worked in vision, planning, imitation learning, reinforcement learning, and reward learning. After graduation, I am seeking foundation model research positions. ”

## EDUCATION

<b>PhD in Computer Science (Machine Learning)</b> University of California, Berkeley. Advised by Professor Stuart Russell.	2018 - May 2024
<b>Bachelor of Advanced Computing (R&amp;D, Honours)</b> Australian National University (ANU)	2014 - 2017

## EMPLOYMENT

<b>Research Intern/Visiting Student Researcher</b> , Meta	May 2022 - Dec 2022
<b>Research Engineer</b> , Seesure Pty Ltd	Dec 2017 - Jun 2018
<b>Summer Intern &amp; Tutor (TA)</b> , Australian National University	Nov 2014 - Jan 2017

## SELECTED PUBLICATIONS

Visit [qxcv.net/pubs](https://qxcv.net/pubs) for a complete list.

<b>A StrongREJECT for Empty Jailbreaks</b> A. Souly, Q. Lu, D. Bowen., T. Trinh, E. Hsieh, S. Pandey, P. Abbeel, J. Svegliato, S. Emmons, O. Watkins, <b>S. Toyer</b>	<i>Under Review</i>
<b>Tensor Trust: Interpretable Prompt Injection Attacks from an Online Game</b> <b>S. Toyer</b> , O. Watkins, E.A. Mendes, J. Svegliato, L. Bailey, T. Wang, I. Ong, K. Elmaaroufi, P. Abbeel, T. Darrell, A. Ritter, S. Russell. ( <b>Demo at <a href="https://tensortrust.ai">tensortrust.ai</a></b> )	ICLR'24 ( <b>spotlight</b> )
<b>An Empirical Investigation of Representation Learning for Imitation</b> X. Chen,* <b>S. Toyer</b> ,* C. Wild,* S. Emmons, I. Fischer, K.H. Lee, N. Alex, S.H. Wang, P. Luo, S. Russell, P. Abbeel, R. Shah	NeurIPS'21 (Datasets & Benchmarks)
<b>The MAGICAL Benchmark for Robust Imitation</b> <b>S. Toyer</b> , R. Shah, A. Critch, S. Russell	NeurIPS'20
<b>ASNs: Deep Learning for Generalised Planning</b> <b>S. Toyer</b> , F. Trevizan, S. Thiébaux, L. Xie	JAIR'20
<b>Variational Discriminator Bottleneck: Improving Imitation Learning, Inverse RL, and GANs by Constraining Information Flow</b> X.B. Peng, A. Kanazawa, <b>S. Toyer</b> , P. Abbeel, S. Levine	ICLR'19
<b>Action Schema Networks: Generalised Policies with Deep Learning</b> <b>S. Toyer</b> , F. Trevizan, S. Thiébaux, L. Xie	AAAI'18

## SELECTED AWARDS

Tong Leong Lim Pre-Doctoral Prize (Berkeley EECS)	2019
Berkeley Fellowship for Graduate Study (Berkeley)	2018
University Medal (ANU)	2017
National University Scholarship (ANU)	2014 - 2017

## SKILLS

Deep learning • Reinforcement learning • Imitation Learning • Large Language Models (LLMs) • PyTorch  
Technical writing • Python • C/C++ • Django • Javascript • GCP & AWS • Linux administration • Ansible