

## Dylan Z. Slack

Curriculum Vitae: February 10, 2021  
Website: <https://dylanslacks.website>

Email: [dslack@uci.edu](mailto:dslack@uci.edu)  
C: 415-847-2440

---

<b>Education</b>	<b>University of California - Irvine</b> , Irvine, CA <i>Ph.D. Computer Science</i> Advisors: Sameer Singh & Hima Lakkaraju Sep. 2019 - Present	
	<b>Haverford College</b> , Haverford, PA <i>B.S. Computer Science with High Honors</i> Magna Cum Laude Advisor: Sorelle Friedler Sep. 2015 - May 2019	
<b>Research and Industry Experience</b>	<b>University of California - Irvine</b> Research Assistant (UCI NLP, CREATE, HPI Institute) <i>Advised by:</i> Sameer Singh	Sep. 2019 - Present
	<b>Amazon Web Services</b> Applied Scientist Intern <i>Advised by:</i> Krishnaram Kenthapadi & Nathalie Rauschmayr	May 2020 - Sep. 2020
	<b>Haverford College</b> Research Assistant, Department of Computer Science <i>Advised by:</i> Sorelle Friedler	Sep. 2017 - Aug. 2019
	<b>Swarthmore College</b> Research Assistant, Department of Computer Science <i>Advised by:</i> Sara Mathieson	Sep. 2018 - May. 2019
<b>Awards</b>	Hasso Plattner Institute Fellow, 2021 Ambler Scholar, 2019	
<b>Publications &amp; Preprints</b> <a href="#">[Scholar]</a>	Defuse: Debugging Classifiers Through Distilling Unrestricted Adversarial Examples <b>Dylan Slack</b> , Nathalie Rauschmayr, and Krishnaram Kenthapadi arXiv, 2020	
	How Much Should I Trust You? Modeling Uncertainty of Black Box Explanations <b>Dylan Slack</b> , Sophie Hilgard, Sameer Singh, and Himabindu Lakkaraju arXiv, 2020	
	Differentially Private Language Models Benefit from Public Pre-training Gavin Kerrigan*, <b>Dylan Slack</b> *, and Jens Tuyls* EMNLP PrivNLP Workshop, 2020	
	Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods <b>Dylan Slack</b> *, Sophie Hilgard*, Emily Jia, Sameer Singh, and Himabindu Lakkaraju AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES), 2020 [Oral Presentation]	
	Fairness Warnings and Fair-MAML: Learning Fairly with Minimal Data <b>Dylan Slack</b> , Sorelle Friedler, and Emile Givental ACM Conference on Fairness, Accountability and Transparency (FAccT), 2020	
	Assessing the Local Interpretability of Machine Learning Models <b>Dylan Slack</b> , Sorelle A. Friedler, Carlos Scheidegger, and Chitradeep Dutta Roy NeurIPS Workshop on Human-Centric Machine Learning, 2019	

\* denotes equal contribution.

<b>Travel Grants</b>	<b>Fairness, Accountability and Transparency in Machine Learning (FAccT)</b> <i>Barcelona, Spain (2020)</i>
	<b>Neural Information Processing Systems (NeurIPS)</b> <i>Vancouver, Canada (2020)</i>
<b>Teaching</b>	<b>Machine Learning (CS 178)</b> UC Irvine <i>Reader (2019)</i>
	<b>Data Structures (CS 206)</b> Bryn Mawr College <i>TA (2019)</i>
	<b>Introduction to Data Structures (CS 106)</b> Haverford College <i>TA (2017, 2018, 2019)</i>
	<b>Introduction to Data Science (CS 104)</b> Haverford College <i>TA (2016)</i>
<b>Talks</b>	Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods Aggregate Intellect, 2021 in Virtual
	Fairness Warnings and Fair-MAML: Learning Fairly with Minimal Data FAccT Conference, 2020 in <i>Barcelona, Spain</i>
<b>Review Services</b>	FAccT 2021 ICLR 2021 ICML 2020 AAAI 2020, 2021 NeurIPS 2019, 2020 KDD 2019
<b>Press &amp; Media</b>	Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods, <a href="#">Harvard Business Review</a> , <a href="#">DeepLearning.ai</a> , <a href="#">Twitter</a>