

# Dorsa Sadigh

---

<b>Contact</b>	246 Gates Computer Science Building 353 Jane Stanford Way, Stanford, CA 94305 +1 (949) 929-3559 dorsa@cs.stanford.edu	<a href="https://dorsa.fyi">https://dorsa.fyi</a> <a href="https://github.com/Stanford-ILIAD">https://github.com/Stanford-ILIAD</a>
<b>Current Position</b>	<b>Stanford University</b> Assistant Professor Department of Computer Science and Department of Electrical Engineering	<b>September 2017 - present</b>
<b>Education</b>	<b>University of California, Berkeley</b> Ph.D. in Electrical Engineering and Computer Sciences Advisors: Sanjit Seshia and Shankar Sastry Thesis: <i>Safe and Interactive Autonomy: Control, Learning, and Verification</i>	<b>2017</b>
	<b>University of California, Berkeley</b> B.S. in Electrical Engineering and Computer Sciences	<b>2012</b>
<b>Awards</b>	<b>ONR Young Investigator Program Award</b> <b>IEEE RAS Early Career Award</b> <b>Sloan Foundation Fellowship</b> <b>Okawa Foundation Research Grant</b> <b>MIT TR35</b> <b>JP Morgan Faculty Award</b> <b>AFOSR Young Investigator Program Award</b> <b>Best Paper Award</b> Conference on Robot Learning (CoRL) for “ <i>Learning Latent Representations to Influence Multi-Agent Interaction</i> ” <b>Best Student Paper Award (Finalist)</b> Robotics: Science and Systems (RSS) for “ <i>Shared Autonomy with Learned Latent Actions</i> ” <b>IEEE TC-CPS Early Career Award</b> <b>Best Paper Award (Honorable Mention)</b> ACM/IEEE International Conference on Human-Robot Interaction (HRI) for “ <i>When Humans Aren’t Optimal: Robots that Collaborate with Risk-Aware Humans</i> ” <b>National Science Foundation CAREER Award</b> <b>Gilbreth Lecturer at National Academy of Engineering</b> <b>Google Faculty Research Award</b>	<b>2022</b> <b>2022</b> <b>2022</b> <b>2021</b> <b>2021</b> <b>2021</b> <b>2020</b> <b>2020</b> <b>2020</b> <b>2020</b> <b>2020</b> <b>2020</b> <b>2020</b> <b>2020</b> <b>2020</b>

<b>Amazon Research Award</b>	<b>2019</b>
<b>Best Paper Award (Finalist)</b> European Control Conference (ECC), for “ <i>Human-Robot Interaction for Truck Platooning Using Hierarchical Dynamic Games</i> ”	<b>2019</b>
<b>Best Paper Award</b> ICML Workshop on Adaptive & Multitask Learning: Algorithms & Systems, for “ <i>Continual Adaptation for Efficient Machine Communication</i> ”	<b>2019</b>
<b>Best Cognitive Robotics Paper (Finalist)</b> IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) for “ <i>Information Gathering Actions over Human Internal State</i> ”	<b>2016</b>
<b>Leon O. Chua Award</b> for excellence in non-linear science, EECS Department, UC Berkeley 2016	<b>2016</b>
<b>Google Anita Borg Scholarship</b>	<b>2016</b>
<b>National Defense Science and Engineering Graduate Fellowship</b>	<b>2013</b>
<b>National Science Foundation Graduate Research Fellowship</b>	<b>2013</b>
<b>CRA Outstanding Undergraduate Researcher Award</b>	<b>2012</b>
<b>Arthur M. Hopkin Award</b> EECS Department, UC Berkeley	<b>2010</b>

<b>Teaching</b>	<b>CS 237B: Robot Autonomy II</b> Instructor, Stanford University.	<b>Winter 2020, 2021</b>
	<b>CS 221: Artificial Intelligence</b> Instructor, Stanford University.	<b>Spring 2018, 2019, Fall 2019, 2020, 2021</b>
	<b>CS 521: Seminar on AI Safety</b> Instructor, Stanford University.	<b>Spring 2018, 2020</b>
	<b>CS 333: Safe and Interactive Robotics</b> Instructor, Stanford University.	<b>Fall 2017, 2018</b>

<b>Advising &amp; Mentoring</b>	<b>Current Graduate Students</b> Erdem Biyik, Minae Kwon, Mengxi Li (co-advised with Jeannette Bohg), Zhangjie Cao, Siddhartha Karamcheti (co-advised with Percy Liang), Andy Shih (co-advised with Stefano Ermon), Suneel Belkhale, Megha Srivastava (co-advised with Dan Boneh), Jennifer Grannen, Priya Sundaresan (co-advised with Jeannette Bohg)
	<b>Past Postdoctoral Students</b> Dylan Losey (Faculty at Virginia Tech)
	<b>Past Undergraduate Students</b> Nick Landolfi (Ph.D. student in CS at Stanford), Zhiyang He (Ph.D. student in EECS at UC Berkeley), Zheqing Zhu (Ph.D. student in MS&E at Stanford), Jovana Kondik (Ph.D. student in EECS at MIT), Songyuan Zhang (Ph.D. student in EECS at MIT), Albert Zhai (Ph.D. student in CS at UIUC), Woody Wang (Ph.D. student in CS at Stanford).

<b>Outreach</b>	<b>Stanford CS Mentorship Program</b>	<b>2018 - present</b>
	I have organized the Stanford CS mentorship program, where we connect underrepresented minorities and female undergraduate students interested in AI with Ph.D. students at Stanford to meet monthly and discuss research and career choices.	
	<b>Faculty Mentor for Stanford Robotics Club</b>	<b>2017 - 2020</b>
	I mentor the Stanford undergraduate Robotics Club. Every year they work towards participating in a robotics competition. They have won the third place in the University Rover Challenge in 2019.	
	<b>Faculty Mentor for Inclusion in AI</b>	<b>2018 - present</b>
	I mentor the Stanford AI Lab graduate group “Inclusion in AI”. The group holds regular social and networking events for Stanford AI Lab graduate students.	
<b>Work Experience</b>	<b>Talks at Women and Inclusion in STEM events and panels</b>	
	AI4ALL summer program, Girls Who Code summer program, Gender in Robotics Workshop at Stanford, Berkeley-Stanford Meetup, Rising Stars (EECS) of 2018, Rising Stars (Mechanical Engineering) of 2019, Inclusion in AI.	
	<b>Talks at Graduate and Undergraduate Student Groups</b>	
	Undergrad CS Women (WiCS), Grad Engineering Women (SWE), SAIL (Stanford AI Lab women), Women in Electrical Engineering, Women in Aero/Astro, Fire-Side chat with Stanford Undergrads.	
	<b>EEGSA Outreach Member</b>	<b>2012 - 2017</b>
	Visiting local K-12 schools and presenting engineering projects and demonstrations.	
<b>Professional Activities</b>	<b>WICSE Outreach Coordinator</b>	<b>2014 - 2015</b>
	Organizing events and outreach activities aiming young girls involvements in STEM.	
	<b>Microsoft Research, Redmond</b>	<b>June - August 2015</b>
	Internship at the Adaptive Systems and Interaction group with Ashish Kapoor and Eric Horvitz.	
	<b>Stanford Research Institute, International</b>	<b>June - August 2013</b>
	Internship at the Computer Science Laboratory in the formal methods group with Ashish Tiwari.	
<b>Professional Activities</b>	<b>Workshop Organizer</b>	<b>2022</b>
	Conference on Robot Learning	
	<b>Program Co-Chair</b>	<b>2018 - 2021</b>
	Bay Area Robotics Symposium	
	<b>Program Co-Chair</b>	<b>2022</b>
	Workshop on Algorithmic Foundations of Robotics (WAFR)	
<b>Professional Activities</b>	<b>Program Committee (Associate Editor, Area Chair)</b>	
	RA-L 2021, CoRL 2021-2020, RSS 2020, HRI 2020, HRI 2022, L4DC 2020, CAV 2019, HSCC 2019, CoRL 2018, ICRA 2018, HSCC Repeatability Eval 2016.	
	<b>Award Committee</b>	<b>2021</b>

Conference on Robot Learning

**Publicity Chair** 2021  
ACM International Conference on Hybrid Systems: Computation and Control

**AAAI ACM SIGAI Dissertation Award Committee** 2020, 2022

**DARPA Information Science and Technology (ISAT) Study Group**  
**Committee Member** 2021

**Center for AI Safety at Stanford** 2018 - present  
Founding member of the Center for AI Safety at Stanford along with Mykel Kochenderfer, Clark Barrett, and David Dill. The center is focused on safety and verification issues for AI and machine learning systems.

**Human-Centered AI Institute (HAI)** 2018 - present  
Member of the design committee of Human-Centered AI Institute at Stanford. In addition I have been part of the HAI Ethical Review Board (ERB) committee.

**External Reviewer for Conferences, Journals, and Grant Panels**

- *Robotics*: RA-L, RSS, CoRL, WAFR, ICRA, HRI, TASE, ACM TECS
- *Control Theory*: HSCC, CDC, ACC, TCST
- *Formal Methods*: CAV, FM, HVC, VMCAI
- *NSF RI Proposal Panel*
- *NSF CPS Proposal Panel*
- *NSF Smart & Connected Communities Panel*
- *NSF ERC Planning Proposal Panel*
- *AFOSR Proposal Reviewer*

**Invited  
Talks**

University of Pennsylvania, GRASP Seminar	2021
NeurIPS Workshop on Robot Learning.	2021
NeurIPS Workshop on Cooperative AI.	2021
NeurIPS Workshop on Learning and Strategic Behavior.	2021
CDC Workshop on Aware Learning: How to Benefit from Priors.	2021
ETH/EPFL NCCR Automation Seminar.	2021
IROS Workshop on RL-CONFORM: Reinforcement Learning meets HRI, Control, and Formal Methods.	2021
IROS Workshop on Multi-Agent and Relational Reasoning.	2021
ACL Workshop on Interactive Learning for NLP.	2021
RSS Workshop on Robotics for People.	2021
Berkeley Seminar on Multi-Agent Reinforcement Learning.	2021
ICML Workshop on Human-AI Collaboration in Sequential Decision Making.	2021

<b>CVPR Workshop on Autonomous Driving: Perception, Prediction and Planning.</b>	2021
<b>Center for Human-Compatible AI Workshop.</b>	2021
<b>ICRA Workshop on Robot-Assisted Systems for Medical Training.</b>	2021
<b>ICRA Workshop on Social Intelligence in Humans and Robots.</b>	2021
<b>SRI Summer School on Formal Techniques.</b>	2021
<b>ACC Workshop on Bridging the Gap in Autonomous Vehicle Controls in Mixed Traffic.</b>	2021
<b>ACC Workshop on Recent Advancement of Human Autonomy Interaction and Integration.</b>	2021
<b>Computer Science Department Seminar, Yale.</b>	2021
<b>Computational Sensorimotor Learning Seminar, MIT.</b>	2021
<b>Center for Human-Compatible AI Seminar, UC Berkeley.</b>	2021
<b>ICLR Workshop on Responsible AI.</b>	2021
<b>RPI Department Seminar.</b>	2021
<b>Control Meets Learning Seminar, Caltech.</b>	2021
<b>AAAI Workshop on Plan, Activity, and Intent Recognition.</b>	2021
<b>Human-Centered AI Institute Seminar.</b>	2021
<b>Keynote – Conference on Robot Learning (CoRL)</b> Walking the Boundary of Learning and Interaction	2020
<b>NeurIPS Workshop on Robot Learning.</b> – “–.	2020
<b>National Canadian Robotics Network (NCRN) Seminar.</b> – “–.	2020
<b>Distinguished Voices – National Academies of Sciences, Engineering, and Medicine</b> A Human-Centered Perspective on Interactive Robotics	2020
<b>Panelist – AI Ethics Conference at Interdisciplinary Research Center in the Arts, Humanities, and Interpretive Social Sciences at Duke Kunshan University</b>	2020
<b>IPAM Workshop on Individual Vehicle Autonomy: Perception and Control.</b> Interaction-Aware Planning: A Human-Centered Approach toward Autonomous Driving.	2020
<b>Keynote – 1st Colloquium on AI for Architecture, Engineering, and Construction</b>	2020

<b>ICML Workshop on Real-World Experiment Design &amp; Active Learning.</b> Active Learning of Robot Reward Functions.	2020
<b>RSS Workshop on Interaction and Decision-Making in Autonomous-Driving.</b> When our Human Modeling Assumptions Fail: Planning, learning, and prediction in near-accident driving scenarios.	2020
<b>RSS Workshop on Power-On-and-Go Robots: Out-of-the-Box Systems for Real-World Applications.</b> To Ignore Humans or to Accept them with Open Arms: Challenges and Opportunities for Efficient, Robust, and Adaptive POGO Robots.	2020
<b>RSS Workshop on AI &amp; Its Alternatives in Assistive &amp; Collaborative Robotics: Decoding Intent.</b> The Role of Learned Representations in Assistive Teleoperation.	2020
<b>Keynote – 22nd ACM International Conference on Hybrid Systems: Computation and Control (HSCC).</b> Human-CPS from the Lens of Learning and Control.	2020
<b>Keynote – Center for Human-Compatible AI Workshop.</b> – “–.	2020
<b>John Hopkins, Applied Physics Lab Seminar.</b> – “–.	2020
<b>ICRA Workshop on Long-Term Human Motion Prediction.</b> When our Human Modeling Assumptions Fail: The effects of risk, conventions, and non-stationarity on long-term human-robot interaction.	2020
<b>NASA Formal Methods, AI Safety Workshop.</b> Risk-Aware Human Modeling.	2020
<b>IPAM Workshop on Intersections between Control, Learning, and Optimization.</b> Beyond Theory of Mind: Learning & Influencing Conventions.	2020
<b>Gilbreth Lecture, National Academy of Engineering.</b> Influencing Interactions in Autonomous Driving.	2020
<b>Keynote – Formal Methods in Computer-Aided Design (FMCAD).</b> A journey about Safety of Autonomous Systems.	2019
<b>Frontiers of Engineering, National Academy of Engineering.</b> Influencing Interactions in Autonomous Driving.	2019
<b>RSS Workshop on Safe Autonomy.</b> –“–.	2019
<b>First Conference on Learning for Dynamics and Control.</b> Influencing Interactive Mixed-Autonomy Systems.	2019
<b>ICML Workshop on AI for Autonomous Driving.</b> –“–.	2019
<b>MIT, Department Seminar.</b> Interactive Autonomy: Learning and Control for Human-Robot Systems.	2019

University of Washington, Department Seminar. –“–.	2019
Cornell, Department Seminar. –“–.	2019
CalTech, IST Seminar. –“–.	2019
USC, CPS Seminar. –“–.	2019
University of Maryland, Robotics Seminar. –“–.	2019
Theoretical Machine Learning Simons Foundation Workshop. –“–.	2019
Schloss Dagstuhl on Verification and Synthesis for Human-Robot Interaction. Reward Functions and Specifications	2019
NeurIPS Workshop on Imitation Learning and its Challenges in Robotics. Active Learning of Humans’ Preferences.	2018
UAI Workshop on Safety, Risk and Uncertainty in RL. –“–.	2018
UC Berkeley, Center for Human Compatible AI. –“–.	2018
NeurIPS Workshop on Machine Learning for Intelligent Transportation Systems. Beating Congestion using Autonomous Cars.	2018
Halmstad University. Reactive Synthesis and Human Modeling for Human-Robot Systems.	2018
University of Washington, Robotics Seminar. Safe and Interactive Robotics. 2018	
UC Santa Barbara, Robotics Seminar. –“–.	2018
UC Santa Cruz, Robotics Seminar. –“–.	2018
Chinese University of Hong Kong in Shenzhen. –“–.	2018
Stanford University, Department Seminar. Towards a Theory of Safe and Interactive Autonomy.	2017
MIT, Department Seminar. –“–.	2017
UC Berkeley, Department Seminar. –“–.	2017
CMU, Department Seminar. –“–.	2017
Princeton, Department Seminar. –“–.	2017
USC, Department Seminar. –“–.	2017
Cornell, Department Seminar. –“–.	2017
UC San Diego, Department Seminar. –“–.	2017

<b>UC Los Angeles, Department Seminar.</b> –“–.	2017
<b>University of Michigan, Department Seminar.</b> –“–.	2017
<b>UT Austin, Department Seminar.</b> –“–.	2017
<b>Georgia Tech, Department Seminar.</b> –“–.	2017
<b>University of Pennsylvania, Department Seminar.</b> –“–.	2017
<b>Schloss Dagstuhl on Machine Learning and Formal Methods.</b> Planning for Cars that Coordinate with People.	2017
<b>Schloss Dagstuhl on Non-Zero-Sum-Games and Control.</b> Correctness and Control for Human-Cyber-Physical Systems.	2015
<b>Microsoft Research, Redmond.</b> Controller Synthesis for Human-in-the-Loop Systems	2014

## Publications

- [84] Suneel Belkhale, Ethan Kroll Gordon, Yuxiao Chen, Siddhartha Srinivasa, Tapomayukh Bhattacharjee, Dorsa Sadigh. Balancing Efficiency and Comfort in Robot-Assisted Bite Transfer. *International Conference on Robotics and Automation (ICRA)*, 2022.
- [83] Zhangjie Cao, Zihan Wang, Dorsa Sadigh. Learning from Imperfect Demonstrations via Adversarial Confidence Transfer. *International Conference on Robotics and Automation (ICRA)*, 2022.
- [82] Zihan Wang, Zhangjie Cao, Yilun Hao, Dorsa Sadigh. Weakly Supervised Correspondence Learning. *International Conference on Robotics and Automation (ICRA)*, 2022.
- [81] Zhangjie Cao, Erdem Biyik, Guy Rosman, Dorsa Sadigh. Leveraging Smooth Attention Prior for Multi-Agent Trajectory Prediction. *International Conference on Robotics and Automation (ICRA)*, 2022.
- [80] Andy Shih, Stefano Ermon, Dorsa Sadigh. Conditional Imitation Learning for Multi-Agent Games. *ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 2022.
- [79] Erdem Biyik, Aditi Talati, Dorsa Sadigh. APReL: A Library for Active Preference-based Reward Learning Algorithms. *ACM/IEEE International Conference on Human-Robot Interaction, Short Contributions (HRI)*, 2022.
- [78] Erdem Biyik, Anusha Lalitha, Rajarshi Saha, Andrea Goldsmith, Dorsa Sadigh. Partner-Aware Algorithms in Decentralized Cooperative Bandit Teams. *AAAI Conference on Artificial Intelligence (AAAI)*, 2022.
- [77] Suvir Mirchandani, Siddharth Karamcheti, Dorsa Sadigh. ELLA: Exploration through Learned Language Abstraction. *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [76] Songyuan Zhang, Zhangjie Cao, Dorsa Sadigh, Yanan Sui. Confidence-Aware Imitation Learning from Demonstrations with Varying Optimality. *Conference on Neural*



*Information Processing Systems (NeurIPS)*, 2021.

[75] Andy Shih, Dorsa Sadigh, Stefano Ermon. HyperSPNs: Compact and Expressive Probabilistic Circuits. *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.

[74] Shushman Choudhury, Jayesh Gupta, Mykel Kochenderfer, Dorsa Sadigh, Jeannette Bohg. Dynamic Multi-Robot Task Allocation under Uncertainty and Temporal Constraints. *Journal of Autonomous Robots (AURO)*, September 2021.

[73] Woodrow Zhouyuan Wang, Andy Shih, Annie Xie, Dorsa Sadigh. Influencing Towards Stable Multi-Agent Interactions. *Conference on Robot Learning (CoRL)*, 2021. **(Oral)**

[72] Zhangjie Cao, Yilun Hao, Mengxi Li, Dorsa Sadigh. Learning Feasibility to Imitate Demonstrators with Different Dynamics. *Conference on Robot Learning (CoRL)*, 2021.

[71] Nils Wilde, Erdem Biyik, Dorsa Sadigh, Stephen L. Smith. Learning Reward Functions from Scale Feedback. *Conference on Robot Learning (CoRL)*, 2021.

[70] Vivek Myers, Erdem Biyik, Nima Anari, Dorsa Sadigh. Learning Multimodal Rewards from Rankings. *Conference on Robot Learning (CoRL)*, 2021. **(Oral)**

[69] Siddharth Karamcheti, Megha Srivastava, Percy Liang, Dorsa Sadigh. LILA: Language-Informed Latent Actions. *Conference on Robot Learning (CoRL)*, 2021.

[68] Julia White, Gabriel Poesia, Robert Hawkins, Dorsa Sadigh and Noah Goodman. Open-domain clarification question generation without question examples. *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021.

[67] Minae Kwon, Mengxi Li, Dorsa Sadigh. Influencing Leading and Following in Human-Robot Teams. *Journal of Autonomous Robots (AURO)*, August 2021.

[66] Erdem Biyik, Dylan Losey, Malayandi Palan, Nick Landolfi, Gleb Shevchuk, Dorsa Sadigh. Learning Reward Functions from Diverse Sources of Human Feedback: Optimally Integrating Demonstrations and Preferences. *The International Journal of Robotics Research (IJRR)*, 2021.

[65] Dylan Losey, Hong Jun Jeon, Mengxi Li, Krishnan Srinivasan, Ajay Mandlekar, Animesh Garg, Jeannette Bohg, Dorsa Sadigh. Learning Latent Actions to Control Assistive Robots. *Journal of Autonomous Robots (AURO)*, 2021.

[64] Behrad Toghi, Rodolfo Valiente, Dorsa Sadigh, Ramtin Pedarsani, Yaser Fallah. Cooperative Autonomous Vehicles that Sympathize with Humans. *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021.

[63] Daniel Lazar, Erdem Biyik, Dorsa Sadigh, Ramtin Pedarsani. Learning How to Dynamically Route Autonomous Vehicles on Shared Roads. *Journal of Transportation Research Part C*, 2021.

[62] Minae Kwon, Siddharth Karamcheti, Mariano-Florentino Cuellar, Dorsa Sadigh. Targeted Data Acquisition for Evolving Negotiation Agents. *The 38th International*

*Conference on Machine Learning (ICML)*, 2021.

[61] Woodrow Wang, Mark Beliaev, Erdem Biyik, Daniel Lazar, Ramtin Pedarsani, Dorsa Sadigh. Emergent Prosociality in Multi-Agent Games Through Gifting. *The 30th International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.

[60] Andy Shih, Arjun Sawhney, Jovana Kondic, Stefano Ermon, Dorsa Sadigh. On the Critical Role of Conventions in Adaptive Human-AI Collaboration. *International Conference on Learning Representations (ICLR)*, 2021.

[59] Mengxi Li, Alper Canberk, Dylan Losey, Dorsa Sadigh. Learning Human Objectives from Sequences of Physical Corrections. *International Conference on Robotics and Automation (ICRA)*, 2021.

[58] Kejun Li, Maegan Tucker, Erdem Biyik, Ellen Novoseller, Joel Burdick, Yanan Sui, Dorsa Sadigh, Yisong Yue, Aaron Ames. ROIAL: Region of Interest Active Learning for Characterizing Exoskeleton Gait Preference Landscapes. *International Conference on Robotics and Automation (ICRA)*, 2021.

[57] Zhangjie Cao, Minae Kwon, Dorsa Sadigh. Transfer Reinforcement Learning across Homotopy Classes. *IEEE Robotics and Automation Letters (RAL)*, 2021.

[56] Zhangjie Cao, Dorsa Sadigh. Learning from Imperfect Demonstrations with Varying Dynamics. *IEEE Robotics and Automation Letters (RAL)*, 2021.

[55] Siddharth Karamcheti, Albert Zhai, Dylan Losey, Dorsa Sadigh. Learning Visually Guided Latent Actions for Assistive Teleoperation. *3rd Learning for Dynamics & Control Conference (L4DC)*, 2021.

[54] Mark Beliaev, Erdem Biyik, Daniel Lazar, Woodrow Wang, Dorsa Sadigh, Ramtin Pedarsani. Incentivizing Routing Choices for Safe and Efficient Transportation in the Face of the COVID-19 Pandemic. *12th ACM/IEEE International Conference on Cyber-Physical Systems (ICCPS)*, May 2021.

[53] Erdem Biyik, Daniel A. Lazar, Ramtin Pedarsani, Dorsa Sadigh. Incentivizing Efficient Equilibria in Traffic Networks with Mixed Autonomy. *IEEE Transactions on Control of Network Systems (TCNS)*, 2020.

[52] Hadas Kress-Gazit, Kerstin Eder, Guy Hoffman, Henny Admoni, Brenna Argall, Ruediger Ehlers, Christoffer Heckman, Nils Jansen, Ross Knepper, Jan Kretinsky, Shelly Levy-Tzedek, Jamy Li, Todd Murphey, Laurel Riek, Dorsa Sadigh. Formalizing and Guaranteeing\* Human-Robot Interaction. *Communications of the ACM*, 2020.

[51] Annie Xie, Dylan Losey, Ryan Tolsma, Chelsea Finn, Dorsa Sadigh. Learning Latent Representations to Influence Multi-Agent Interaction. *Proceedings of the 4th Conference on Robot Learning (CoRL)*, November 2020. **(Oral, Best Paper Award)**

[50] Robert X. D. Hawkins, Minae Kwon, Dorsa Sadigh, Noah D. Goodman. Continual Adaptation for Efficient Machine Communication. *The 24th Conference on Computational Natural Language Learning (CoNLL)*, November 2020.

[49] Mengxi Li, Dylan Losey, Jeannette Bohg, Dorsa Sadigh. Learning User-Preferred Mappings for Intuitive Robot Control. *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, October 2020.

- [48] Zheqing Zhu, Erdem Bıyık, Dorsa Sadigh. Multi-Agent Safe Planning with Gaussian Processes. *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, October 2020.
- [47] Jonathan Mern, Dorsa Sadigh, Mykel Kochenderfer. Object Exchangability in Reinforcement Learning. *2020 American Control Conference (ACC)*, July 2020.
- [46] Hong Jun Jeon, Dylan Losey, Dorsa Sadigh. Shared Autonomy with Learned Latent Actions. *Robotics: Science and Systems (RSS)*, June 2020. **(Best Student Paper Award, Finalist)**
- [45] Erdem Bıyık, Nicolas Huynh, Mykel Kochenderfer, Dorsa Sadigh. Active Preference-Based Gaussian Process Regression for Reward Learning. *Robotics: Science and Systems (RSS)*, June 2020.
- [44] Zhangjie Cao, Erdem Bıyık, Woodrow Wang, Allan Raventos, Adrien Gaidon, Guy Rosman, Dorsa Sadigh. Reinforcement Learning based Control of Imitative Policies for Near-Accident Driving. *Robotics: Science and Systems (RSS)*, June 2020.
- [43] Shushman Choudhury, Jayesh Gupta, Mykel Kochenderfer, Dorsa Sadigh, Jeannette Bohg. Dynamic Multi-Robot Task Allocation under Uncertainty and Temporal Constraints. *Robotics: Science and Systems (RSS)*, June 2020.
- [42] Malayandi Palan, Shane Barratt, Alex McCauley, Dorsa Sadigh, Vikas Sindhwani, Stephen P. Boyd. Fitting a Linear Control Policy to Demonstrations with a Kalman Constraint. *2nd Learning for Dynamics & Control Conference (L4DC)*, June 2020.
- [41] Minae Kwon, Erdem Bıyık, Aditi Talati, Karan Bhasin, Dylan P. Losey, Dorsa Sadigh. When Humans Aren’t Optimal: Robots that Collaborate with Risk-Aware Humans. *ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 2020. **(Best Paper Award, Honorable Mention)**
- [40] Yuhang Che, Allison M. Okamura, Dorsa Sadigh. Efficient and Trustworthy Social Navigation Via Explicit and Implicit Robot-Human Communication. *IEEE Transactions on Robotics (TRO)*, 2019.
- [39] Dylan P. Losey, Krishnan Srinivasan, Ajay Mandlekar, Animesh Garg, Dorsa Sadigh. Controlling Assistive Robots with Learned Latent Actions. *International Conference on Robotics and Automation (ICRA)*, May 2020.
- [38] Dylan P. Losey, Mengxi Li, Jeannette Bohg, Dorsa Sadigh. Learning from My Partner’s Actions: Roles in Decentralized Robot Teams. *Conference on Robot Learning (CoRL)*, 2019. **(Oral)**
- [37] Erdem Bıyık, Malayandi Palan, Nicholas Landolfi, Dylan P. Losey, Dorsa Sadigh. Asking Easy Questions: A User-Friendly Approach to Active Reward Learning. *Conference on Robot Learning (CoRL)*, 2019.
- [36] Dylan P. Losey, Dorsa Sadigh. Robots that Take Advantage of Human Trust. *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, November 2019.
- [35] Chandrayee Basu, Erdem Bıyık, Zhixun He, Mukesh Singhal, Dorsa Sadigh. Ac-

tive Learning of Reward Dynamics from Hierarchical Queries. *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, November 2019.

[34] Erdem Biyik, Daniel A. Lazar, Dorsa Sadigh, Ramtin Pedarsani. The Green Choice: Learning and Influencing Human Decisions on Shared Roads. *Proceedings of the 58th IEEE Conference on Decision and Control (CDC)*, December 2019.

[33] Minae Kwon, Mengxi Li, Alexandre Bucquet, Dorsa Sadigh. Influencing Leading and Following in Human-Robot Teams. *Robotics: Science and Systems (RSS)*, June 2019.

[32] Malayandi Palan, Gleb Shevchuk, Nicholas C. Landolfi, Dorsa Sadigh. Learning Reward Functions by Integrating Human Demonstrations and Preferences. *Robotics: Science and Systems (RSS)*, June 2019.

[31] Tianhe Yu, Gleb Shevchuk, Dorsa Sadigh, Chelsea Finn. Unsupervised Visuomotor Control through Distributional Planning Networks. *Robotics: Science and Systems (RSS)*, June 2019.

[30] Erdem Biyik, Jonathan Margoliash, Shahrouz Ryan Alimo, Dorsa Sadigh. Efficient and Safe Exploration in Deterministic Markov Decision Processes with Unknown Transition Models. *2019 American Control Conference (ACC)*, July 2019.

[29] Elis Stefansson, Jaime Fisac, Dorsa Sadigh, Shankar Sastry, Karl H. Johansson. Human-Robot Interaction for Truck Platooning Using Hierarchical Dynamic Games. *European Control Conference (ECC)*, June 2019. **(Best Paper Award, Finalist)**.

[28] Ashwini Pople, Roberto Martin-Martin, Patrick Goebel, Vincent Chow, Hans M. Ewald, Junwei Yang, Zenkai Wang, Amir Sadeghian, Dorsa Sadigh, Silvio Savarese, Marynel Vazquez. Deep Local Trajectory Planning and Control for Robot Navigation. *International Conference on Robotics and Automation (ICRA)*, May 2019.

[27] Jaime F. Fisac, Eli Bronstein, Elis Stefansson, Dorsa Sadigh, S. Shankar Sastry, Anca D. Dragan. Hierarchical Game-Theoretic Planning for Autonomous Vehicles. *International Conference on Robotics and Automation (ICRA)*, May 2019.

[26] Erdem Biyik, Dorsa Sadigh. Batch Active Preference-Based Learning of Reward Functions. *Conference on Robot Learning (CoRL)*, 2018. **(Oral)**

[25] Erdem Biyik, Daniel A. Lazar, Ramtin Pedarsani, Dorsa Sadigh. Altruistic Autonomy: Beating Congestion on Shared Roads . *International Workshop on Algorithmic Foundations of Robotics (WAFR)*, 2018.

[24] Daniel Lazar, Kabir Chandrasekher, Ramtin Pedarsani, Dorsa Sadigh. Maximizing Road Capacity Using Cars that Influence People. *IEEE Conference on Decision and Control (CDC)*, 2018.

[23] Jiaming Song, Hongyu Ren, Dorsa Sadigh, Stefano Ermon. Multi-Agent Generative Adversarial Imitation Learning. *Conference on Neural Information Processing Systems (NeurIPS)*, 2018.

[22] Dorsa Sadigh, S. Shankar Sastry, Sanjit Seshia. Verifying Robustness of Human-Aware Autonomous Cars . *IFAC conference on Cyber-Physical and Human Systems*

(CPHS), 2018.

[21] Dorsa Sadigh, Nick Landolfi, S. Shankar Sastry, Sanjit A. Seshia, Anca Dragan. Planning for Autonomous Cars that Leverages Effects on Human Actions. *Journal of Autonomous Robots (AURO)*, 2018.

[20] Susmit Jha, Vasumathi Raman, Dorsa Sadigh, Sanjit A. Seshia. Safe Autonomy Under Perception Uncertainty Using Chance-Constrained Temporal Logic . *Journal of Automatic Reasoning (JAR)*, 2018.

[19] Dorsa Sadigh. Safe and Interactive Autonomy: Control, Learning, and Verification. *Ph.D. Dissertation. EECS Department, University of California, Berkeley, August 2017.*

[18] Dorsa Sadigh, S. Shankar Sastry, Sanjit Seshia, Anca Dragan. Active Preference-Based Learning of Reward Functions. *Robotics: Science and Systems Conference (RSS)*, July 2017.

[17] Negar Mehr, Dorsa Sadigh, Roberto Horowitz, S. Shankar Sastry, Sanjit Seshia. Stochastic Predictive Freeway Ramp Metering from Signal Temporal Logic Specifications. *American Control Conference (ACC)*, May 2017.

[16] Dorsa Sadigh, S. Shankar Sastry, Sanjit Seshia, Anca Dragan. Information Gathering Actions over Human Internal State. *International Conference on Intelligent Robots and Systems (IROS)*, 2016. **(Best Paper in Cognitive Robotics Award, Finalist).**

[15] Tara Rezvani, Katherine Driggs-Campbell, Dorsa Sadigh, S. Shankar Sastry, Sanjit Seshia, Ruzena Bajcsy. Towards Trustworthy Automation: User Interfaces that Convey Internal and External Awareness. *IEEE Intelligent Transportation Systems Conference (ITSC)*, November 2016.

[14] Dorsa Sadigh, S. Shankar Sastry, Sanjit Seshia, Anca Dragan Planning for Autonomous Cars that Leverages Effects on Human Actions . *Robotics: Science and Systems Conference (RSS)*, 2016.

[13] Dorsa Sadigh, Ashish Kapoor. Safe Control under Uncertainty with Probabilistic Signal Temporal Logic. *Robotics: Science and Systems Conference (RSS)*, 2016.

[12] Shromona Ghosh, Dorsa Sadigh, Pierluigi Nuzzo, Vasumathi Raman, Alexandre Donze, Alberto Sangiovanni-Vincentelli, S. Shankar Sastry, Sanjit Seshia. Diagnosis and Repair for Synthesis from Signal Temporal Logic Specifications. *Conference on Hybrid Systems: Computation and Control (HSCC)*, 2016.

[11] Sanjit A. Seshia, Dorsa Sadigh, S. Shankar Sastry. Formal Methods for Semi-autonomous Driving. *Design and Automation Conference (DAC)*, 2015.

[10] Vasumathi Raman, Alexandre Donze, Dorsa Sadigh, Richard M. Murray, Sanjit Seshia. Reactive Synthesis from Signal Temporal Logic Specifications. *Conference on Hybrid Systems: Computation and Control (HSCC)*, 2015.

[9] Dorsa Sadigh, Eric S. Kim, Samuel Coogan, S. Shankar Sastry, Sanjit Seshia. A Learning Based Approach to Control Synthesis of Markov Decision Processes for Linear Temporal Logic Specifications. *IEEE Conference on Decision and Control (CDC)*, 2014.

- [8] Dorsa Sadigh, Henrik Ohlsson, S. Shankar Sastry, Sanjit Seshia. Robust Subspace System Identification via Weighted Nuclear Norm Optimization. *International Federation of Automatic Control (IFAC)*, 2014.
- [7] Dorsa Sadigh, Katherine Driggs-Campbell, Alberto Puggelli, Wenchao Li, Victor Shia, Ruzena Bajcsy, Alberto Sangiovanni-Vincentelli, Shankar Sastry, and Sanjit Seshia. Data-driven probabilistic modeling and verification of human driver behavior. *Formal Verification and Modeling in Human-Machine Systems (AAAI Spring Symposium)*, 2014.
- [6] Dorsa Sadigh, Katherine Driggs Campbell, Ruzena Bajcsy, S. Shankar Sastry, Sanjit Seshia. User Interface Design and Verification for Semi-autonomous Driving. *Conference on High Confidence Networked Systems*, 2014.
- [5] Ashish Tiwari, Bruno Dutertre, Dejan Jovanovic, Thomas de Candia, Dorsa Sadigh, Sanjit Seshia. Safety Envelop in Security. *Conference on High Confidence Networked Systems (HiCoNS)*, 2014.
- [4] Wenchao Li, Dorsa Sadigh, S. Shankar Sastry, Sanjit Seshia. Synthesis for Human-in-the-Loop Control Systems. *Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, 2014.
- [3] Dorsa Sadigh, Sanjit Seshia and Mona Gupta. Automating Exercise Generation: A Step towards Meeting the MOOC Challenge for Embedded Systems. *Workshop on Embedded Systems Education*, 2012.
- [2] Orna Kupferman, Dorsa Sadigh, and Sanjit A. Seshia. Synthesis with Clairvoyance. *Haifa Verification Conference (HVC)*, 2011.
- [1] Jonathan Kotker, Dorsa Sadigh, and Sanjit A. Seshia. Timing Analysis of Interrupt-Driven Programs under Context Bounds. *Formal Methods in Computer Aided Design (FMCAD)*, 2011.

## Technical Reports & Workshop Papers

- [13] Bidipta Sarkar\*, Aditi Talati\*, Andy Shih\*, Dorsa Sadigh. PantheonRL: A MARL Library for Dynamic Training Interactions. *Proceedings of the 36th AAAI Conference on Artificial Intelligence (Demo Track)*, February 2022.
- [12] Erdem Biyik, Anusha Lalitha, Rajarshi Saha, Andrea Goldsmith, Dorsa Sadigh. Partner-Aware Algorithms in Decentralized Cooperative Bandit Teams. *Artificial Intelligence for Human-Robot Interaction (AI-HRI) at AAAI Fall Symposium Series*, November 2021.
- [11] Erdem Biyik, Aditi Talati, Dorsa Sadigh. APReL: A Library for Active Preference-based Reward Learning Algorithms. *Artificial Intelligence for Human-Robot Interaction (AI-HRI) at AAAI Fall Symposium Series*, November 2021.
- [10] Nicholas Roy, Ingmar Posner, Tim Barfoot, Philippe Beaudoin, Yoshua Bengio, Jeannette Bohg, Oliver Brock, Isabelle Depatie, Dieter Fox, Dan Koditschek, Tomas Lozano-Perez, Vikash Mansinghka, Christopher Pal, Blake Richards, Dorsa Sadigh, Stefan Schaal, Gaurav Sukhatme, Denis Thérien, Marc Toussaint, Michiel Van de

Panne. From Machine Learning to Robotics: Challenges and Opportunities for Embodied Intelligence. *arXiv*, November 2021.

[9] Andy Shih, Dorsa Sadigh, Stefano Ermon. HyperSPNs: Compact and Expressive Probabilistic Circuits. *The 4th Workshop on Tractable Probabilistic Modeling*, 2021.

[8] Behrad Toghi, Rodolfo Valiente, Dorsa Sadigh, Ramtin Pedarsani, Yaser Fallah. Altruistic Maneuver Planning for Cooperative Autonomous Vehicles Using Multi-agent Advantage Actor-Critic. *CVPR Workshop on Autonomous Driving: Perception, Prediction and Planning*, 2021.

[7] Mark Beliaev, Woodrow Z. Wang, Daniel A. Lazar, Erdem Biynk, Dorsa Sadigh, Ramtin Pedarsani. Emergent Correlated Equilibrium through Synchronized Exploration. *RSS 2020 Workshop on Emergent Behaviors in Human-Robot Systems*, 2020.

[6] Kawin Ethayarajh, Dorsa Sadigh. BLEU Neighbors: A Reference-less Approach to Automatic Evaluation. *1st Workshop on Evaluation and Comparison for NLP systems (Eval4NLP)*, 2020.

[5] Siddharth Karamcheti, Dorsa Sadigh and Percy Liang. Continual Learning Adaptive Language Interfaces through Decomposition. *Proceedings of the EMNLP Workshop on Interactive and Executable Semantic Parsing*, 2020.

[4] Robert X. D. Hawkins, Minae Kwon, Dorsa Sadigh, Noah D. Goodman. Continual Adaptation for Efficient Machine Communication. *Proceedings of the ICML Workshop on Adaptive & Multitask Learning: Algorithms & Systems*, June 2019. **(Best Paper Award)**.

[3] Jiaming Song, Hongyu Ren, Dorsa Sadigh, Stefano Ermon. Multi-Agent Generative Adversarial Imitation Learning. *International Conference on Learning Representations (ICLR), Workshop Track*, April 2018.

[2] Sanjit Seshia, Dorsa Sadigh, S. Shankar Sastry. Towards Verified Artificial Intelligence. *Technical Report*, July 2016.

[1] Debadeepta Dey, Dorsa Sadigh, Ashish Kapoor. Fast Safe Mission Plans for Autonomous Vehicles. *Proceedings of Robotics: Science and Systems Workshop*, June 2016.

## Dissertation

Dorsa Sadigh. Safe and Interactive Autonomy: Control, Learning, and Verification. *Ph.D. Dissertation; EECS Department, University of California, Berkeley*, August 2017.