

Aida Mostafazadeh Davani

Department of Computer Science
University of Southern California
Los Angeles, CA 90007, United States

mostafaz@usc.edu
<https://aidamd.github.io/>
Cell Phone: +1 (323) 449 - 4538

- Education**
- ◇ **Ph.D. Computer Science** Fall 2017 – [Expected] Spring 2022
University of Southern California, Los Angeles, United States
 - ◇ **M.Sc. Software Engineering** Fall 2014 – Fall 2017
Sharif University of Technology, Tehran, Iran
 - ◇ **B.Sc. Software Engineering** Fall 2009 – Summer 2014
Sharif University of Technology, Tehran, Iran
- Research Interests**
- ◇ Ethics in Machine Learning
 - ◇ Natural Language Processing
 - ◇ Computational Social Science
- Publications**
- ◇ Prabhakaran, V., **Mostafazadeh Davani, A.**, Díaz, M. “*On releasing annotator-level labels and information in datasets*”, The 15th Linguistic Annotation & 3rd Designing Meaning Representations Joint Workshop (2021, accepted)
 - ◇ **Mostafazadeh Davani, A.**, Díaz, M., Prabhakaran, V. “*Dealing with disagreements: Looking beyond the majority vote in subjective annotations*”, Transactions of the Association for Computational Linguistics (2021, accepted).
 - ◇ Hoover, J., Atari, M.*, **Mostafazadeh Davani, A.***, Kennedy, B.*, Portillo-Wightman, G., Yeh, L., Dehghani, M. “*Investigating the role of group-based morality in extreme behavioral expressions of prejudice*”, Nature Communications (2021).
 - ◇ **Mostafazadeh Davani, A.**, Omrani, A., Kennedy, B., Atari, M., Ren, X., Dehghani, M. “*Improving counterfactual generation for fair hate speech detection*”, Proceedings of the 5th Workshop on Online Abuse and Harms (2021).
 - ◇ Kennedy, B., Atari, M., **Mostafazadeh Davani, A.**, Hoover, J., Omrani, A., Graham, J., Dehghani, M. “*Moral concerns are differentially observable in language*”, Cognition (2021).
 - ◇ Jin, X., Barbieri, F., Kennedy, B., **Mostafazadeh Davani, A.**, Neves, L., Ren, X. “*On Transferability of Bias Mitigation Effects in Language Model Fine-Tuning*”, Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics (2021).
 - ◇ **Mostafazadeh Davani, A.**, Atari, M., Kennedy, B., Havaladar, S., Dehghani, M. “*Hatred is in the eye of the annotator: Hate speech classifiers learn human-like social stereotypes*”, Conference of the Cognitive Science Society (2020).
 - ◇ Atari, M., **Mostafazadeh Davani, A.**, Dehghani, M. “*Body maps of moral concerns*”, Psychological Science (2020).
 - ◇ Hoover, J., Portillo-Wightman, G., Yeh, L., Havaladar, S., **Mostafazadeh Davani, A.**, Lin, Y., Kennedy, B., Atari, M., Kamel, Z., Mendlen, M., Moreno, G., Chin, J., Leong, C., Leung, J. Y., Mirinjian, A., Dehghani, M. “*Moral Foundations Twitter Corpus: A collection of 35k tweets annotated for moral sentiment*”, Social Psychological and Personality Science (2020).
 - ◇ Kennedy, B.*, Jin, X.*, **Mostafazadeh Davani, A.**, Dehghani, M., Ren, X. “*Contextualizing hate speech classifiers with post-hoc explanation*”, In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (2020).
 - ◇ **Mostafazadeh Davani, A.**, Yeh, L., Atari, M., Kennedy, B., Wightman, G. P., Gonzalez, E., Delong, N., Bhatia, R., Mirinjian, A., Ren, X., Dehghani, M. “*Reporting the unreported: Event extraction for analyzing the local representation of hate crimes*”, In the Proceedings of Empirical Methods in Natural Language Processing (2019).

- ◇ Courtland, M., **Mostafazadeh Davani, A.**, Reyes, M., Yeh, L., Leung, J., Kennedy, B., Dehghani, M., Zevin, J. “*Modeling performance differences on cognitive tests using LSTMs and skip-thought vectors trained on reported media consumption*”, Proceedings of the Third Workshop on Natural Language Processing and Computational Social Science (2019).
- ◇ Courtland, M., **Mostafazadeh Davani, A.**, Reyes, M., Yeh, L., Leung, J., Kennedy, B., Dehghani, M., Zevin, J. “*Subtle differences in language experience moderate performance on language-based cognitive tests.*”, Proceedings of the 41st Annual Conference of the Cognitive Science Society, Austin, Texas. Cognitive Science Society. (2019).
- ◇ **Mostafazadeh Davani, A.**, Nazari Shirehjini, A. A., Daraei, S. “*Towards interacting with smarter systems*”, Journal of Ambient Intelligence and Humanized Computing (2018).
- ◇ **Mostafazadeh Davani, A.**, Nazari Shirehjini, A. A. and Daraei, S. “*A Meta user interface for understandable and predictable interaction in AAL.*”, Human Aspects of IT for the Aged Population. Design for Everyday Life. Springer International Publishing, (2015).
- ◇ **Mostafazadeh Davani, A.**, Nazari Shirehjini, A. A., Daraei, S., Khojasteh, N., and Shirmohammadi, S. “*A Meta user interface for interaction with mixed reality environments.*”, Haptic, Audio and Visual Environments and Games (HAVE), 2015 IEEE International Symposium on, IEEE, (2015).

Pre-prints

- ◇ **Mostafazadeh Davani, A.**, Atari, M., Kennedy, B., Dehghani, M. “*Hate speech classifiers learn human-like social stereotypes*”, (under preparation).
- ◇ Atari, M., **Mostafazadeh Davani, A.**, Kogon, D., Kennedy, B., Saxena, N. A., Anderson, I., Dehghani, M. “*Morally homogeneous networks and radicalism*”, (under review).
- ◇ Kennedy, B., Atari, M., **Mostafazadeh Davani, A.**, Yeh, L., Omrani, A., Kim, Y., ..., Hoover, J. “*The gab hate corpus: A collection of 27k posts annotated for hate speech*”, (under review).
- ◇ Vial, A. C.*, **Mostafazadeh Davani, A.***, Havaldar, S., Chestnut, E. K., Dehghani, M., Cimpian, A. “*Syntactic and semantic gender biases in the language on children’s television: Evidence from a corpus of 95 shows from 1960 to 2018*”, (under preparation).
- ◇ Atari, M., Hoover, J., Kennedy, B., **Mostafazadeh Davani, A.**, Omrani, A., Karimi-Malekabadi, F., Portillo-Wightman, G., Birjandi, S., Dehghani, M. “*Pathogens are linked to human moral systems across time and space*”, (under preparation).
- ◇ Goodwin, R. D., Dodson, S. J., Chambers, M., **Mostafazadeh Davani, A.**, Dehghani, M., Graham, J., Diekmann, K. A. “*Twitter observers’ moral language reveals how sexual harassment denials condemn #MeToo victims*”, (under preparation).

Presentations

- ◇ **Mostafazadeh Davani, A.**, Omrani, A., Kennedy, B., Atari, M., Ren, X., Dehghani, M. “*Improving counterfactual generation for fair hate speech detection*”, Oral presentation at the 5th Workshop on Online Abuse and Harms (2021, August).
- ◇ **Mostafazadeh Davani, A.**, Díaz, M., Prabhakaran, V. “*Multi-Annotator Modeling to Encode Diverse Perspectives in Hate Speech Annotations*”, Oral presentation at the 5th Workshop on Online Abuse and Harms (2021, August).
- ◇ Dehghani, M., Atari, M., **Mostafazadeh Davani, A.**, Kennedy, B. “*Extremists of a feather hate together: Morally homogeneous networks and use of hateful rhetoric*”, Oral presentation at the Society for Personality and Social Psychology Conference (2021, February).
- ◇ **Mostafazadeh Davani, A.**, Atari, M., Kennedy, B., Havaldar, Sh., Dehghani, M. “*Hatred is in the eye of the annotator: Hate speech classifiers learn human-like social stereotypes*”, Oral presentation at the Cognitive Science Conference, (2020, July).
- ◇ Kennedy, B.*, Jin, X.*, **Mostafazadeh Davani, A.**, Dehghani, M., Ren, X. “*Contextualizing hate speech classifiers with post-hoc explanation*”, Oral presentation at the Annual Meeting of Association for Computational Linguistics, Seattle, USA (2020, July).
- ◇ Atari, M., **Mostafazadeh Davani, A.**, and Dehghani, M. “*Body maps of moral concerns*”, Oral presentation at the Society for Personality and Social Psychology Conference, New Orleans, USA (2020, February).

- ◇ **Mostafazadeh Davani, A.**, Yeh, L., Atari, M., Kennedy, B., Portillo-Wightman, G., Gonzalez, E., ... Dehghani, M. . “*Reporting the unreported: Event extraction for analyzing the local representation of hate crimes*”, Oral presentation at the Conference on Empirical Methods in Natural Language Processing, Hong Kong, China (2019, November).
- ◇ Chambers, M. K., **Mostafazadeh Davani, A.**, Goodwin, R., Dodson, S., Yeh, L., Dehghani, M., Graham, J., Diekmann, K. A. “*The power of silence: Using sentiment text analysis to examine twitter responses to sexual harassment accounts*”, International Association for Conflict Management (IACM), Dublin, Ireland (2019).
- ◇ **Mostafazadeh Davani, A.**, Vial, A., Chestnut, E., Leung, J. Y., Cimpian, A., Dehghani, M., “*An analysis of gender bias in children’s TV shows*”, Oral presentation at the Third Workshop on Natural Language Processing and Computational Social Science, Minneapolis, USA (June, 2019)
- ◇ Kennedy, B., Hoover, J., **Mostafazadeh Davani, A.**, Garten, J., Dehghani, M.. “*Measuring moral rhetoric in text: attempts to capture context*”, Oral presentation at the Annual Society for Personality and Social Psychology Conference, Portland, USA (2019, November).
- ◇ Dehghani, M., Kennedy, B., **Mostafazadeh Davani, A.**, Yeh, L., and Atari, M. “*Inferring moral concerns from Facebook status updates*”, Oral presentation at the Association for Research in Personality, Grand Rapids, USA (2019, June).
- ◇ Hoover, J., Atari, M., **Mostafazadeh Davani, A.**, Kennedy, B., Dehghani, M. “*Moral values and acts of hate*”, Oral presentation at the Society for Personality and Social Psychology Conference, Portland, USA (2019, February).
- ◇ Lin, Y., **Mostafazadeh Davani, A.**, Oyserman. D., & Dehghani. M. “*Situating honor in moral contexts*”, Data Blitz presented in the Psychology of Language pre-conference, Society of Personality and Social Psychology Convention, Atlanta, USA (2018).

Experience

- ◇ Student Researcher at Google Fall 2021
Ethical AI team
- ◇ Research Intern at Google Summer 2021
Ethical AI team
- ◇ Research Assistant at University of Southern California Fall 2017 - Spring 2021
Computational Social Science Lab
- ◇ Teaching Assistant at University of Southern California Fall 2018 - Spring 2019
Analysis of Algorithm
- ◇ Research Assistant at Sharif University of Technology Fall 2014 - Fall 2016
Ambient Intelligence Lab
- ◇ Teaching Assistant at Sharif University of Technology Fall 2015 - Spring 2017
Human Computer Interaction
- ◇ Teaching Assistant at Sharif University of Technology Fall 2012 - Fall 2013
Principles of Programming

Honors and Awards

- ◇ Graduate Research Assistantship, National Science Foundation (NSF) 2018-2020
(PI: Morteza Dehghani)
- ◇ Graduate Research Assistantship, National Institute of Health (NIH) 2018
(PI: Jason Zevin)
- ◇ Hopper Scholarship Award, USC Department of Computer Science 2017

Academic Memberships

- ◇ Society for Personality and Social Psychology (SPSP) 2019-2020
- ◇ Cognitive Science Society (CogSci) 2020
- ◇ Association for Computational Linguistics (ACL) 2019-2020

Skills

- ◇ Programming: Python, Java, C++, C#, R

- ◇ Deep Learning: Tensorflow, PyTorch, Keras
- ◇ Statistics: Hierarchical Modeling, Time Series Analysis

**Ad-Hoc
Review**

- ◇ Association for Computational Linguistics (ACL)
- ◇ Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)
- ◇ The Association for the Advancement of Artificial Intelligence (AAAI) Workshop on Affective Content Analysis (AffCon)
- ◇ Workshop on Natural Language Processing and Computation Social Science (NLP+CSS)
- ◇ The International AAAI Conference on Web and Social Media (ICWSM)
- ◇ Behavior Research Methods
- ◇ Cognitive Science Society (CogSci)
- ◇ ACM Transactions on Asian and Low-Resource Language Information Processing