

# Aida Mostafazadeh Davani

Google Research  
555 SW Morrison Street, Suite 500  
Portland, OR 97204 USA

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

<b>Research Interests</b>	<ul style="list-style-type: none"><li>◊ Fairness in Machine Learning</li><li>◊ Natural Language Processing</li><li>◊ Computational Social Science</li></ul>	
<b>Work &amp; Research Experience</b>	<ul style="list-style-type: none"><li>◊ <b>Research Scientist</b> at Google LLC Technology, AI, Society, and Culture (TASC) team</li><li>◊ <b>Research Assistant</b> at University of Southern California Computational Social Science lab</li><li>◊ <b>Research Intern</b> at Google Ethical AI team</li><li>◊ <b>Organizer</b> of the Workshop on Online Abuse and Harm ACL 2023, NAACL 2022, and ACL 2021</li><li>◊ <b>Research Assistant</b> at Sharif University of Technology Ambient Intelligence lab</li></ul>	2022 – Present 2017 – 2022 2021 2021 – Present 2014 – 2016
<b>Education</b>	<ul style="list-style-type: none"><li>◊ <b>Ph.D. Computer Science</b> University of Southern California, Los Angeles, USA</li><li>◊ <b>M.Sc. Software Engineering</b> Sharif University of Technology, Tehran, Iran</li><li>◊ <b>B.Sc. Software Engineering</b> Sharif University of Technology, Tehran, Iran</li></ul>	2017 – 2022 2014 – 2017 2009 – 2014
<b>Peer-Reviewed Publications</b>	<ul style="list-style-type: none"><li>◊ Prabhakaran, V., <b>Mostafazadeh Davani, A.</b>, Ferguson, M. J., Atir, S. “<i>Distinguishing Address vs. Reference Mentions of Personal Names in Text</i>”, ACL Findings (2023).</li><li>◊ Jha, A., <b>Mostafazadeh Davani, A.</b>, Dave, S., Reddy, C., Dev, S., Prabhakaran, V. “<i>A Stereotype Benchmark with Broad Geo-Cultural Coverage Leveraging Generative Models</i>”, ACL (2023).</li><li>◊ Kennedy, B., Golazizian, P., Trager, J., Atari, M., Hoover, J., <b>Mostafazadeh Davani, A.</b>, Dehghani, M. “<i>The (Moral) Language of Hate</i>”, PNAS Nexus (2023).</li><li>◊ Atari, M., Mehl, M. R., Graham, J., Doris, J. M., <b>Mostafazadeh Davani, A.</b>, Omrani, A., ..., Dehghani, M. “<i>The paucity of morality in everyday talk</i>”, Scientific Reports (2023).</li><li>◊ <b>Mostafazadeh Davani, A.</b>, Atari, M., Kennedy, B., Dehghani, M. “<i>Hate speech classifiers learn normative social stereotypes</i>”, TACL (2022).</li><li>◊ Atari, M., Reimer, N. K., Graham, J., Hoover, J., Kennedy, B., <b>Mostafazadeh Davani, A.</b>, Karimi-Malekabadi, F., Birjandi, S., Dehghani, M. “<i>Pathogens are linked to human moral systems across time and space</i>”, Current Research in Ecological and Social Psychology (2022).</li><li>◊ Kennedy, B., Atari, M., <b>Mostafazadeh Davani, A.</b>, Yeh, L., Omrani, A., Kim, Y., ..., Hoover, J. “<i>Introducing the Gab hate corpus: defining and applying hate-based rhetoric to social media posts at scale</i>”, Language Resources and Evaluation (2022).</li><li>◊ <b>Mostafazadeh Davani, A.</b>, Díaz, M., Prabhakaran, V. “<i>Dealing with disagreements: Looking beyond the majority vote in subjective annotations</i>”, TACL (2021).</li><li>◊ Prabhakaran, V.*, <b>Mostafazadeh Davani, A.*</b>, Díaz, M. “<i>On releasing annotator-level labels and information in datasets</i>”, The 15th Linguistic Annotation &amp; 3rd Designing Meaning Representations Joint Workshop (2021).</li><li>◊ <b>Mostafazadeh Davani, A.</b>, Omrani, A., Kennedy, B., Atari, M., Ren, X., Dehghani, M. “<i>Improving counterfactual generation for fair hate speech detection</i>”, WOAH (2021).</li></ul>	2017 – 2022 2014 – 2017 2009 – 2014

- ◊ Atari, M., **Mostafazadeh Davani, A.**, Kogon, D., Kennedy, B., Saxena, N. A., Anderson, I., Dehghani, M. “*Morally homogeneous networks and radicalism*”, Social Psychological and Personality Science (2021).
- ◊ 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).
- ◊ 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*”, ACL (2021).
- ◊ **Mostafazadeh Davani, A.**, Atari, M., Kennedy, B., Havaldar, 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., Havaldar, 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*”, ACL (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*”, EMNLP (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*”, NLP-CSS (2019).
- ◊ **Mostafazadeh Davani, A.**, Shirehjini, A. A. N., Daraei, S. “*Towards interacting with smarter systems*”, Journal of Ambient Intelligence and Humanized Computing (2018).

#### Pre-prints

- ◊ Trager, J., Ziabari, A., **Mostafazadeh Davani, A.**, Golazazian, P., Karimi-Malekabadi, F., Omrani, A., Li, Z., Kennedy, B., ..., Morteza Dehghani “*The Moral Foundations Reddit Corpus*”, (in preparation)
- ◊ 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*”, (in 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*”, (in preparation).

#### Honors and Awards

- ◊ Graduate Research Assistantship, National Science Foundation (NSF) 2018-2020
- ◊ Graduate Research Assistantship, National Institute of Health (NIH) 2018
- ◊ Hopper Scholarship Award, USC Department of Computer Science 2017

#### Skills

- ◊ Programming: Python, Java, C++, C#, R
- ◊ Deep Learning: Tensorflow, PyTorch, Keras
- ◊ Statistics: Hierarchical Modeling, Time Series Analysis