

Does Interacting in Hate Groups Increase Hate?

Daniel Hickey¹ (hickeyda@oregonstate.edu), Matheus Schmitz², Daniel Fessler³, Paul Smaldino⁴, Goran Murić⁵, Keith Burghardt⁵

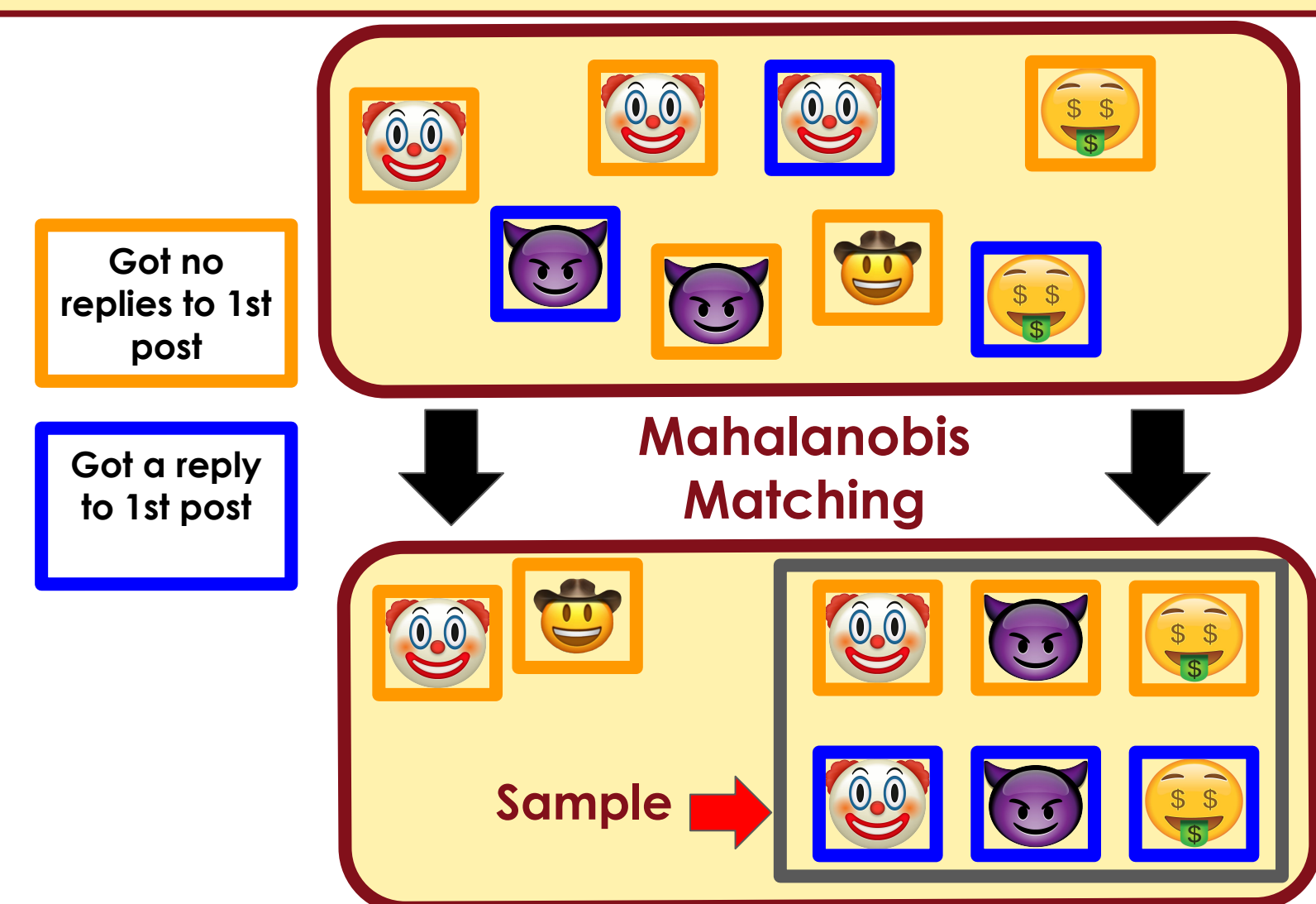
¹Oregon State University, ²Qualcomm, ³University of California, Los Angeles, ⁴University of California, Merced, ⁵University of Southern California Information Sciences Institute

Motivation:

- What social interactions **cause** users of hateful subreddits to become more hateful?

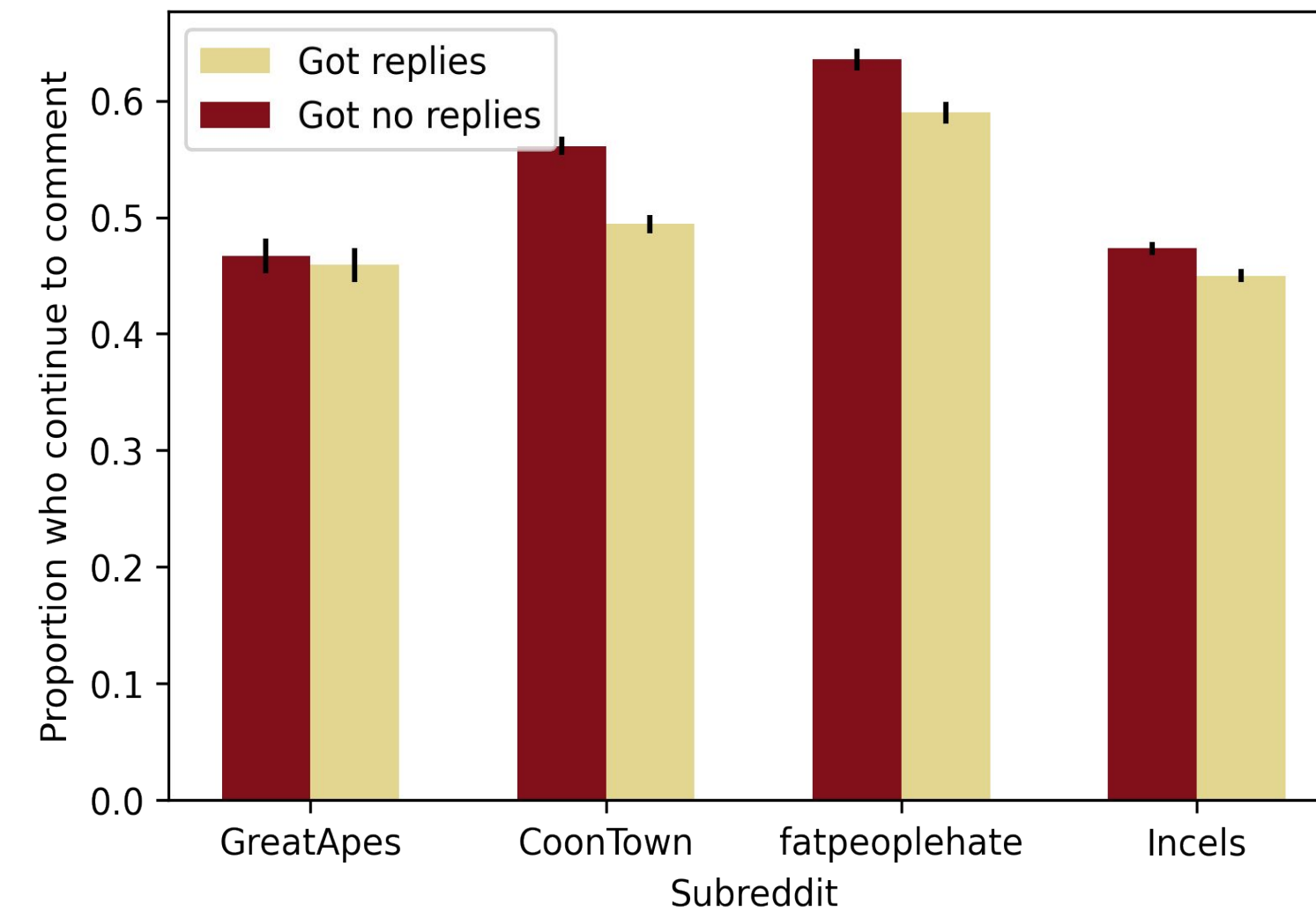
Hate groups investigated

Subreddit	Type of Hate	Users Crawled
r/GreatApes	Anti-Black Racist	3140
r/CoonTown	Anti-Black Racist	9154
r/fatpeoplehate	Fat-shaming	10240
r/Incels	Misogynistic	20211



Schematic of quasi-experimental design for causal inference.

Continued Engagement in Other Submissions After First Post



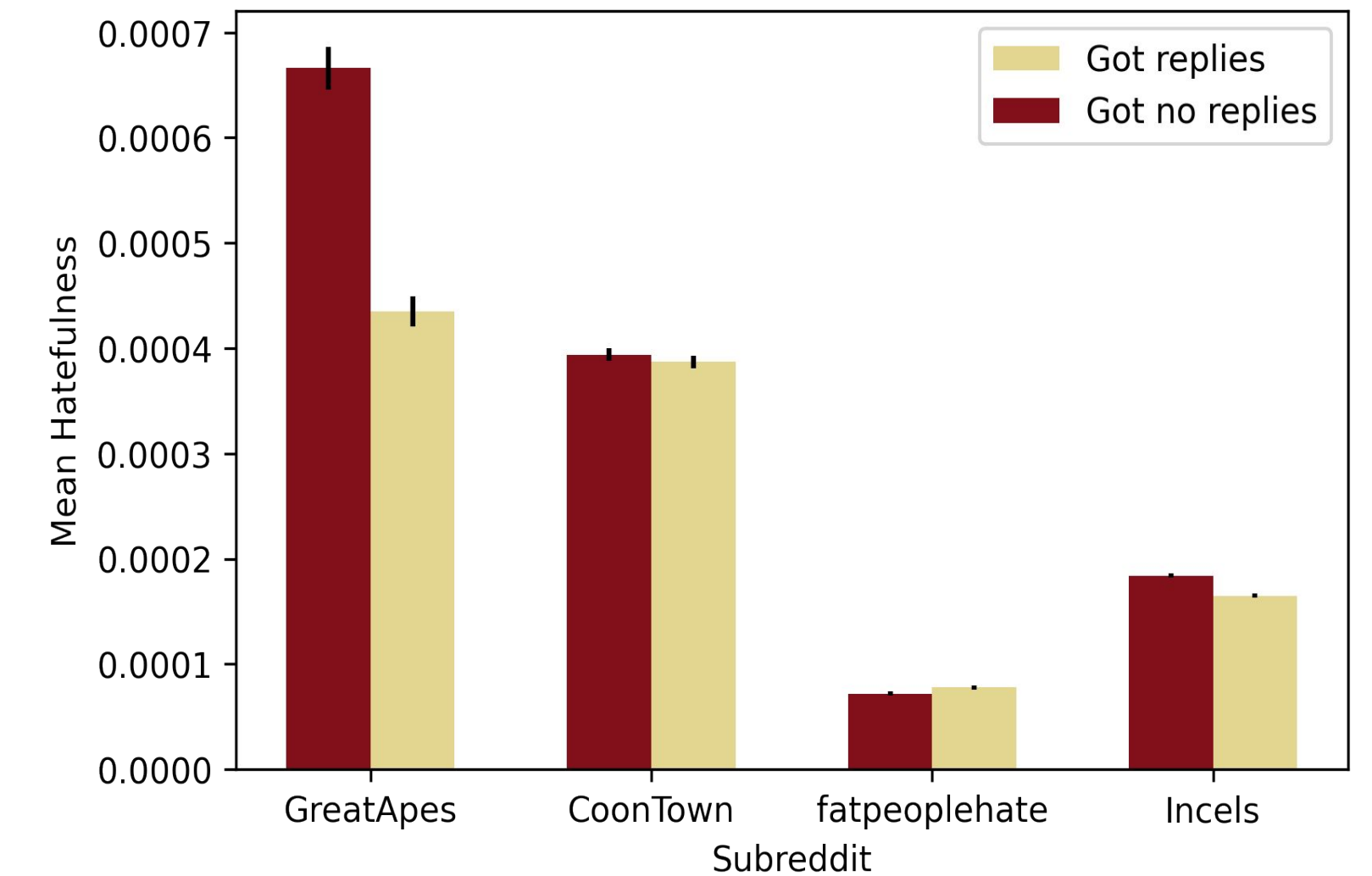
RQ1: Do users who get replies to their first post become more engaged in hate subreddits? - Users who got replies were less likely to continue posting.

Subreddit	Coefficient	P-value
r/GreatApes	0.29	0.02
r/CoonTown	0.17	0.004
r/fatpeoplehate	0.19	$4 * 10^{-8}$
r/Incels	0.027	0.4

Does sentiment of replies influence engagement?

Coefficients of logistic regression model. Response = engagement, predictor = sentiment of first reply.

Mean Hatefulness in Non-Banned Subreddits



RQ2: Do users who get replies to their first post become more antagonistic in non-hateful subreddits? - users who got replies were usually less antagonistic on average

Discussion/Conclusion:

- Getting reactions to **initial posts** in hateful subreddits does not seem to increase hate
- Matching is preliminary - will do more careful feature selection
- Sentiment analysis offers more clarity - some replies are encouraging, some discouraging.
- Will expand analysis to more hateful subreddits