## Social Approval and Network Homophily as Motivators of Online Toxicity





Julie Jiang (USC), Luca Luceri (USC), Joe Walther (UCSB), Emilio Ferrara (USC)

## Hate Speech

targeted, offensive discourse based on race, religion, gender, etc.



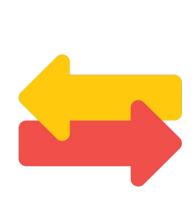
# #\$!@

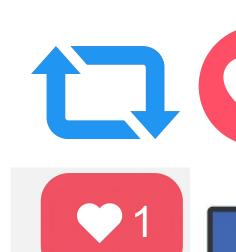
## Theory of online hate

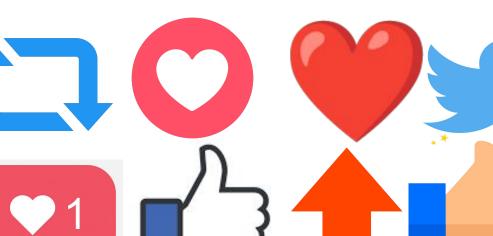
(Walther, 2022)

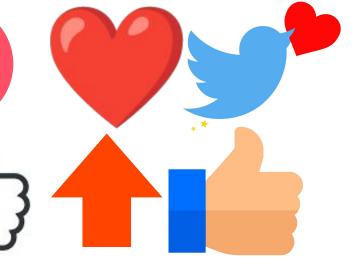
Hateful behavior may be fueled, reinforced, and exacerbated by social approvals from their network









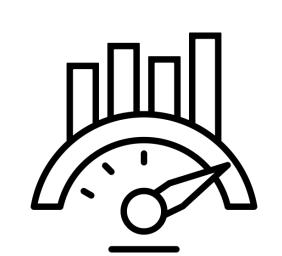


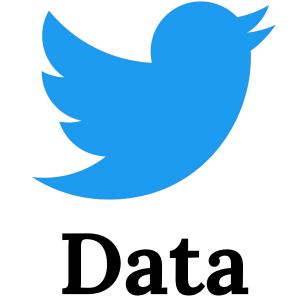


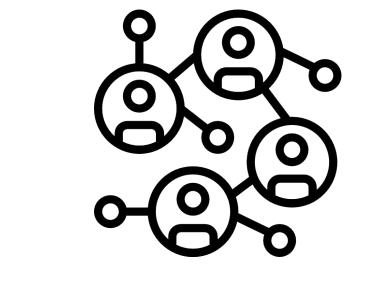




Gather all of their historical tweets







Detect hate scores of tweets and users (Perspective API)

Compile social network interaction data

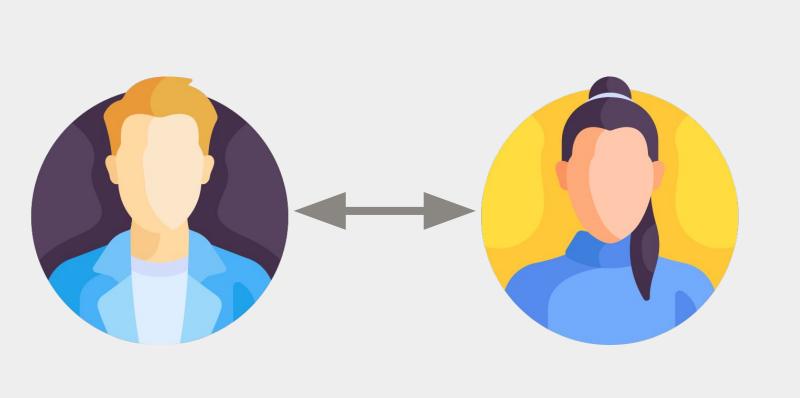
## RQ1: Is a user's hatefulness related to how

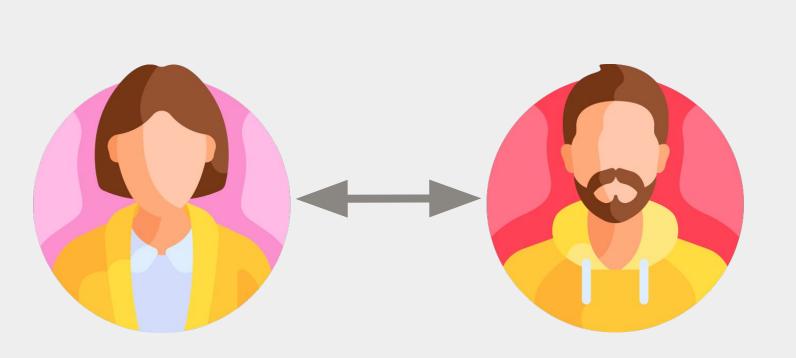
## hateful their social network is?

Yes, there is network homophily (McPherson, 2001): Users tend to preferentially associate themselves with similarly hateful users



 High predictability of user hate score from the social network using ML methods (Social-LLM by Jiang et al., 2024)



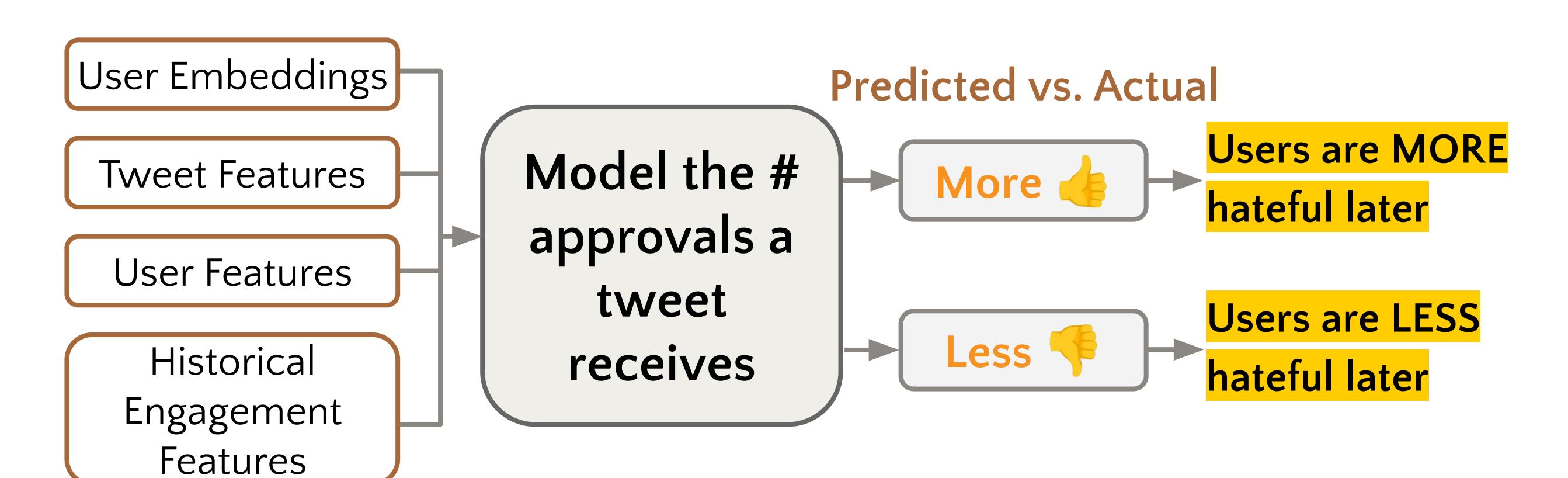


## RQ2: Does receiving more social approval increase a user's subsequent hateful?

We show that toxic behavior could be a socially motivated behavior.

#### How:

- . Build ML model to predict how many 👍 a tweet receives, based on user, tweet, and relevant historical features
- 2. Compare predicted #  $\downarrow \phi$  with actual #  $\downarrow \phi$ , and locate instances where
  - a. the tweet got much more 🤙 than predicted (maybe social approvals) b. the tweet got much fewer 👍 than predicted (maybe social disapprovals)
- 3. Evaluate  $\Delta$  change in average toxicity before and after the social (dis)approval



#### Social approval or disapproval?

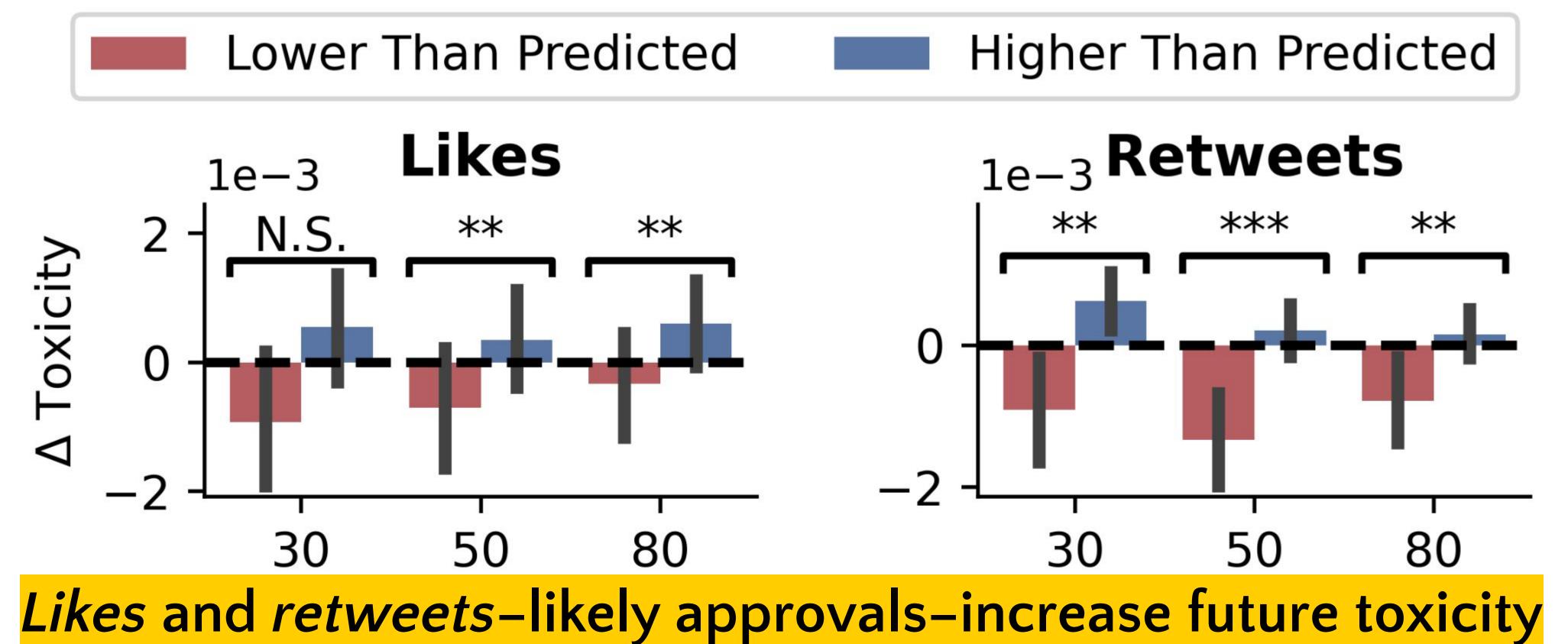
Social engagement can represent social approval, but can also be a conduit for social disapproval

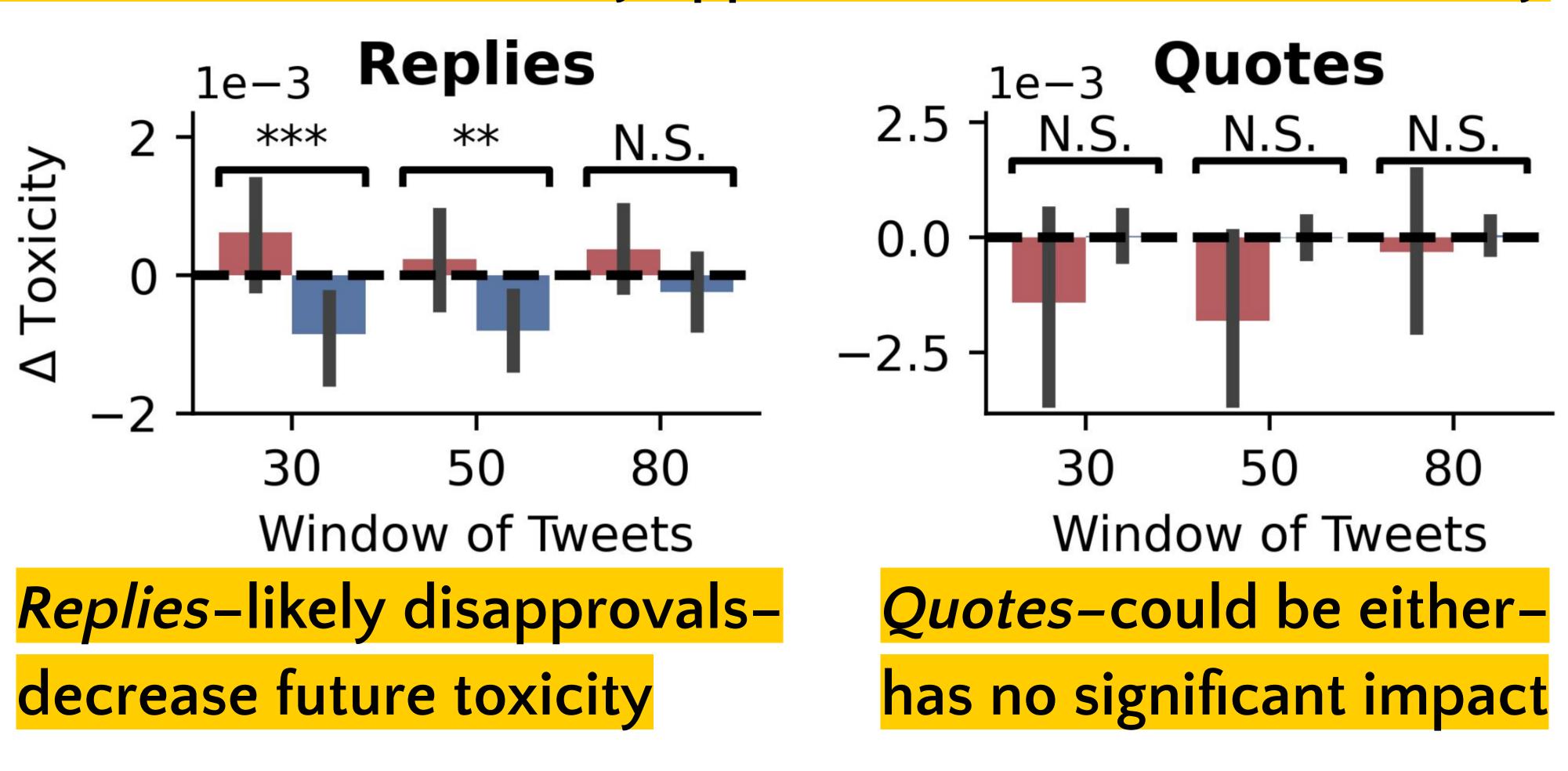
REBROADCAST RETWEETS QUOTES QUOTES

**ENDORSEMENT** REPLIES

LIKES REPLIES

#### Effect of Social Engagement on Toxicity





## Implications

- Hate posting is not a solo endeavor but collective behavior. It is "networked"
- Content moderation strategies: shadowban, disable visibility of some social approvals, etc.

April 2024 Icons created by Freepik - Flaticon