# The Structure of Toxic Conversations on Twitter

Martin Saveski, Brandon Roy, Deb Roy

Massachusetts Institute of Technology

## Goals

- <u>Analysis</u>: study the relationship between structure and toxicity of conversations, after the conversations are over
- <u>Prediction</u>: predict future toxicity based on the structure of the conversation, as the conversation unfolds

### Data

- <u>News</u>: 510K+ conversations, 32M+ tweets, 5 outlets, 1 year
- <u>Midterms</u>: 676K+ conversations, 25M+ tweets, 1,430 candidates, 5 months



### Analyses

#### **Individual-level Analysis**

• Toxicity is spread across many low to moderately toxic users

### **Dyad-level Analysis**

• Toxic replies are more likely to come from other users who: (i) do not have any social relationship with the poster, (ii) have fewer followers, and (iii) do not have many common friends



#### **Reply Tree Structure**

• Toxic conversations tend to have larger, deeper, and wider reply trees



### Prediction

#### **Future Toxicity Predictions**

- <u>Task</u>: Given the conversation so far, predict whether the conversation will become more toxic than expected
- Using stratification to control for prefix toxicity



#### **Next Reply Toxicity Predictions**



#### **Follow Graph Structure**

• Toxic conversations tend to have follow graphs that are denser, have more CCs, and higher modularity



- Task: User i is about to join the conversation, will they post a toxic reply?
- Paired prediction task to control for the root content

	ACC		AUC		F1	
All –	0.712		0.797		0.712	
All \ Conversation State -	0.680	l.	0.753		0.679	
Conversation State -	0.676		0.757	l	0.675	
User-Parent Dyad	0.633	ł	0.690	ł	0.630	
Toxic Embeddedness -	0.595	ł.	0.651	ł	0.599	I
Reply Graph -	0.571	Н	0.602	ł	0.574	H
User-Root Dyad	0.556	H	0.583	Η	0.567	Н
Reply Tree -	0.530	ł.	0.544	Н	0.531	H
Follow Graph -	0.527	ł	0.540	ł	0.521	H
User Info-	0.519		0.527		0.524	H
Overall Embeddedness -	0.517	1	0.525	ł	0.513	H
Political Alignment -	0.510		0.517		0.573	
0	.0	0.5 0	0.0	0.5 0	.0	0.5

Extended version published at WWW'21: <u>https://doi.org/10.1145/3442381.3449861</u>