



CSCW 2020

Exploring Antecedents and Consequences of Toxicity in Online Discussions: A Case Study on Reddit

Yan Xia

Fudan University

Haiyi Zhu

Carnegie Mellon University

Tun Lu

Fudan University

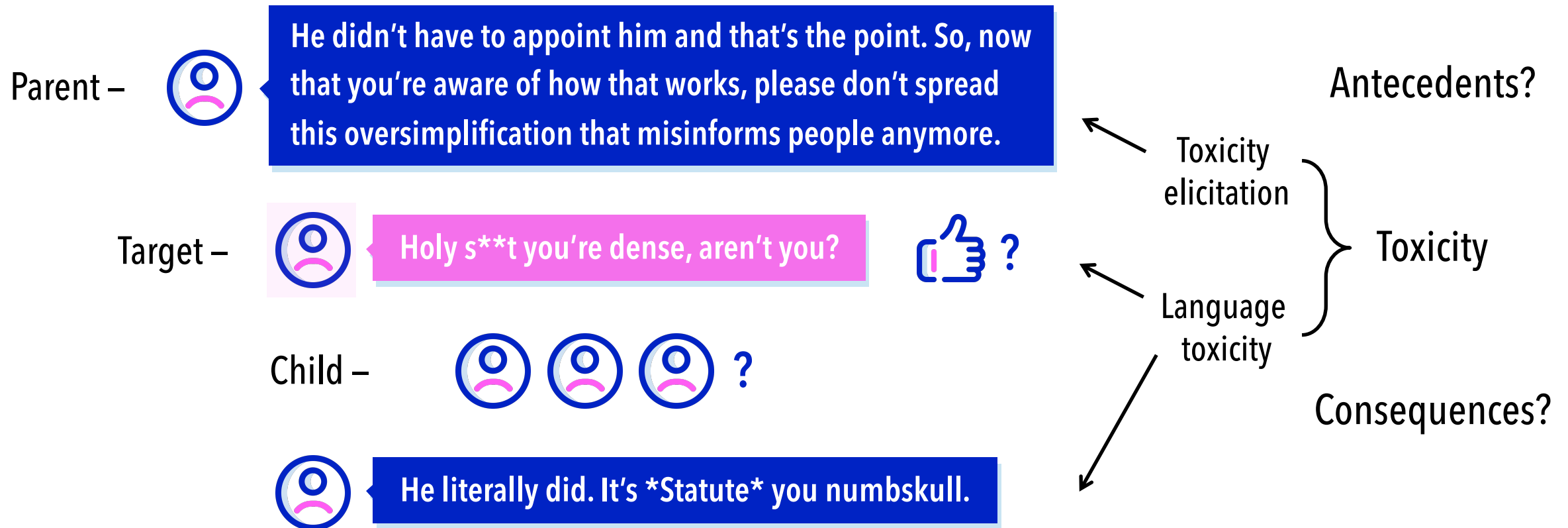
Peng Zhang

Fudan University

Ning Gu

Fudan University

Motivation





Related Work & Hypotheses

- Antecedents of language toxicity
 - H1: **Author's propensity toward toxicity** will increase the language toxicity of his/her text.
 - H2: **Author's experience in the community** will reduce the language toxicity of his/her text.
 - H3: **Toxicity in discussion context** will increase the language toxicity of text.
 - H4: **Polarity in discussion context** will increase the language toxicity of text.
- Consequences of language toxicity
 - Q1: How will the language toxicity of text influence the **volume of discussion**?
 - Q2: How will the language toxicity of text influence the **evaluation of discussion**?
- Antecedents and consequences of toxicity elicitation
 - Q3: How will the antecedents and consequences of **toxicity elicitation** differ from those of language toxicity?



Method

- Reddit comments
 - r/announcements, r/worldnews, r/politics, r/todayilearned, r/AskReddit
 - 19,682/152,632/200,635/139,353/265,893 comments => preprocessed
- Quantifying comment features: NLP tools
 - Toxicity: Perspective API (* with human validation)
 - Polarity: TextBlob library
- Regression analyses
 - Study I, Antecedents of Toxicity: What author/parent comment features predict target comment toxicity?
 - Study II, Consequences of Toxicity: What child comment features are predicted by target comment toxicity?



Study I: Antecedents of Toxicity

Dependent Variable: <i>Target-language-toxicity</i>						
Factor (Expected Correlation)	Independent Variable	r/annc	Regression Coefficient in			
		r/wn	r/pol	r/til	r/ar	
Author's propensity (+)	<i>Target-author-toxicity</i>	.207 **	.133 ***	.125 ***	.103 ***	.160 ***
Author's experience (-)	<i>Target-author-age</i>	-.029	-.055 *	-.029	-.035	.024
	<i>Target-author-karma</i>	-.031	.008	.004	.030	-.003
Toxicity in context (+)	<i>Parent-language-toxicity</i>	.210 **	.182 ***	.156 ***	.184 ***	.246 ***
Polarity in context (+)	<i>Parent-positivity</i>	.135 *	-.000	.005	.000	-.018
	<i>Parent-negativity</i>	.017	-.038	-.020	-.027	-.014

H1: **Author's propensity toward toxicity** will increase the language toxicity of his/her text.

Supported

H2: **Author's experience in the community** will reduce the language toxicity of his/her text.

H3: **Toxicity in discussion context** will increase the language toxicity of text.

Supported

H4: **Polarity in discussion context** will increase the language toxicity of text.



Study I: Antecedents of Toxicity

Dependent Variable: <i>Target-toxicity-elicitation</i>						
Factor (Expected Correlation)	Independent Variable	r/annc	Regression Coefficient in			
		r/wn	r/pol	r/til	r/ar	
Author's propensity (?)	<i>Target-author-toxicity</i>	.011	.068 **	.024	.033	.041 *
Author's experience (?)	<i>Target-author-age</i>	.013	.010	.012	.005	.013
	<i>Target-author-karma</i>	-.032	.039	-.008	.006	-.019
Toxicity in context (?)	<i>Parent-language-toxicity</i>	.173 *	.118 ***	.082 ***	.151 ***	.136 ***
Polarity in context (?)	<i>Parent-positivity</i>	.052	.029	-.021	-.031	-.016
	<i>Parent-negativity</i>	.049	.025	.007	-.037	.007
Control Variable						
/	<i>Target-language-toxicity</i>	.260 ***	.113 ***	.115 ***	.168 ***	.142 ***

Q3: How will the antecedents (and consequences) of **toxicity elicitation** differ from those of language toxicity?



Study II: Consequences of Toxicity

Factor (Expected Correlation): Volume of Discussion (?)						
Dependent Variable	Independent Variable	r/annc	r/wn	Regression Coefficient in		
				r/pol	r/til	r/ar
<i>Target-#children</i>	<i>Target-language-toxicity</i>	.055	.065 **	.056 **	.068 **	.052 **
	<i>Target-toxicity-elicitation</i>	-.128	-.027	-.008	-.009	.019
<i>Target-#descendants</i>	<i>Target-language-toxicity</i>	.005	.083 ***	.113 ***	.107 ***	.080 ***
	<i>Target-toxicity-elicitation</i>	-.144 **	-.055 **	.012	.021	.032 *
<i>Target-height</i>	<i>Target-language-toxicity</i>	-.040	.022	.026	.049 *	.009
	<i>Target-toxicity-elicitation</i>	.016	-.024	.001	.048 *	.044 *

Q1: How will the language toxicity of text influence the **volume of discussion**?



Study II: Consequences of Toxicity

Factor (Expected Correlation): Volume of Discussion (?)						
Dependent Variable	Independent Variable	r/annc	r/wn	Regression Coefficient in		
				r/pol	r/til	r/ar
<i>Target-#children</i>	<i>Target-language-toxicity</i>	.055	.065 **	.056 **	.068 **	.052 **
	<i>Target-toxicity-elicitation</i>	-.128	-.027	-.008	-.009	.019
<i>Target-#descendants</i>	<i>Target-language-toxicity</i>	.005	.083 ***	.113 ***	.107 ***	.080 ***
	<i>Target-toxicity-elicitation</i>	-.144 **	-.055 **	.012	.021	.032 *
<i>Target-height</i>	<i>Target-language-toxicity</i>	-.040	.022	.026	.049 *	.009
	<i>Target-toxicity-elicitation</i>	.016	-.024	.001	.048 *	.044 *

Q3: How will the (antecedents and) consequences of **toxicity elicitation** differ from those of language toxicity?



Study II: Consequences of Toxicity

Factor (Expected Correlation): Evaluation of Discussion (?)						
Dependent Variable	Independent Variable	Regression Coefficient in				
		r/annc	r/wn	r/pol	r/til	r/ar
<i>Target-score</i>	<i>Target-language-toxicity</i>	.019	.011	.023	.025	.034 *
	<i>Target-toxicity-elicitation</i>	-.050	-.018	-.002	-.016	.014
<i>Children-score-max</i>	<i>Target-language-toxicity</i>	.015	.018	.028 *	.012	.049 **
	<i>Target-toxicity-elicitation</i>	-.029	-.021	.009	-.025	.062 ***

Q2: How will the language toxicity of text influence the **evaluation of discussion**?



Study II: Consequences of Toxicity

Factor (Expected Correlation): Evaluation of Discussion (?)						
Dependent Variable	Independent Variable	Regression Coefficient in				
		r/annc	r/wn	r/pol	r/til	r/ar
<i>Target-score</i>	<i>Target-language-toxicity</i>	.019	.011	.023	.025	.034 *
	<i>Target-toxicity-elicitation</i>	-.050	-.018	-.002	-.016	.014
<i>Children-score-max</i>	<i>Target-language-toxicity</i>	.015	.018	.028 *	.012	.049 **
	<i>Target-toxicity-elicitation</i>	-.029	-.021	.009	-.025	.062 ***

Q3: How will the (antecedents and) consequences of **toxicity elicitation** differ from those of language toxicity?



Discussion & Design Implications

- From within and without: **Triggers of toxicity**
 - Toxicity in discussion context
 - Design implication: To interfere with this toxicity generation process
- Flames in disguise: **Toxicity elicitation**
 - Strong-toned / Sarcastic / Digressive / Against common sense or values
 - Design implication: To regulate toxicity-eliciting comments
- The multi-faced devil: **Complexity of toxicity**
 - Different target / emotion / intention => Different consequence
 - Design implication: To distinguish different types of toxic comments in detection and regulation



Limitations & Future Work

- Bias of NLP tools
- Specific time, community and regulation settings
- Limited modeling of toxicity dynamics
- Further questions:
 - Other ways to vitalize a discussion without toxicity?
 - How will the resulted discussions differ?
 - How to stop toxicity from offending people while retaining the “edge” of discussion?



Thanks to the anonymous reviewers and Ge Gao, Diyi Yang, Dakuo Wang,
Beisi Zhou, Xiaofeng Zhao for helping us with the study.

Thank you for listening!

2020.09.23