

Security researchers show Google's antiinternet troll AI platform is easily deceived

March 1 2017, by Jennifer Langston

ORIGINAL PHRASE (TOXICITY SCORE)	MODIFIED PHRASE (TOXICITY SCORE)
Climate change is happening and it's not changing in our favor. If you think differently you're an idiot . (84%)	Climate change is happening and it's not changing in our favor. If you think differently you're an idiiot . (20%)
They're stupid , it's getting warmer, we should enjoy it while it lasts (86%)	They're st.upid , it's getting warmer, we should enjoy it while it lasts (2%)
They are liberal idiots who are uneducated (90%)	They are liberal i.diots who are un.educated (15%)
idiots. backward thinking people. nationalists. not accepting facts. susceptible to lies. (80%)	idiiots. backward thinking people. nationaalists. not accepting facts. susceptible to l.ies. (17%)
They are stupid and ignorant with no class (91%)	They are st.upid and ig.norant with no class (11%)
It's stupid and wrong (89%)	It's stuipd and wrong (17%)
If they voted for Hilary they are idiots (90%)	If they voted for Hilary they are id.iots (12%)
Anyone who voted for Trump is a moron (80%)	Anyone who voted for Trump is a mo.ron (13%)
Screw you trump supporters (79%)	Screw you trump supporters (17%)

In the examples below on hot-button topics of climate change, Brexit and the recent US election -- which were taken directly from the Perspective API website -- the UW team simply misspelled or added extraneous punctuation or spaces to the offending words, which yielded much lower toxicity scores. For example, simply changing "idiot" to "idiiot" reduced the toxicity rate of an otherwise identical comment from 84% to 20%. Credit: University of Washington



University of Washington researchers have shown that Google's new machine learning-based system to identify toxic comments in online discussion forums can be bypassed by simply misspelling or adding unnecessary punctuation to abusive words, such as "idiot" or "moron."

Perspective is a project by Google's technology incubator Jigsaw, which uses artificial intelligence to combat internet trolls and promote more civil online discussion by automatically detecting online insults, harassment and abusive speech. The company launched a demonstration website on Feb. 23 that allows anyone to type in a phrase and see its "toxicity score"—a measure of how rude, disrespectful or unreasonable a particular comment is.

In a paper posted Feb. 27 on the e-print repository arXiv, the UW electrical engineers and security experts demonstrated that the early stage technology system can be deceived by using common adversarial tactics. They showed one can subtly modify a phrase that receives a high toxicity score so that it contains the same abusive language but receives a low toxicity score.

Given that news platforms such as *The New York Times* and other media companies are exploring how the system could help curb harassment and abuse in online comment areas or social media, the UW researchers evaluated Perspective in adversarial settings. They showed that the system is vulnerable to both missing incendiary language and falsely blocking non-abusive phrases.



ORIGINAL PHRASE (TOXICITY SCORE)	MODIFIED PHRASE (TOXICITY SCORE)
Climate change is happening and it's not changing in our favor. If you think differently you're an idiot. (84%)	Climate change is happening and it's not changing in our favor. If you think differently you're not an idiot. (73%)
They're stupid, it's getting warmer, we should enjoy it while it lasts (86%)	They're not stupid, it's getting warmer, we should enjoy it while it lasts (74%)
They are liberal idiots who are uneducated. (90%)	They are not liberal idiots who are uneducated. (83%)
idiots. backward thinking people. nationalists. not accepting facts. susceptible to lies. (80%)	not idiots. not backward thinking people. not nationalists. accepting facts. not susceptible to lies. (74%)
They are stupid and ignorant with no class (91%)	They are not stupid and ignorant with no class (84%)
It's stupid and wrong (89%)	It's not stupid and wrong (83%)
If they voted for Hilary they are idiots (90%)	If they voted for Hilary they are not idiots (81%)
Anyone who voted for Trump is a moron (80%)	Anyone who voted for Trump is not a moron (65%)
Screw you trump supporters (79%)	Will not screw you trump supporters (68%)

In the examples in Graphic 2, the researchers also showed that the system does not assign a low toxicity score to a negated version of an abusive phrase. Credit: University of Washington

"Machine learning systems are generally designed to yield the best performance in benign settings. But in real-world applications, these systems are susceptible to intelligent subversion or attacks," said senior author Radha Poovendran, chair of the UW electrical engineering department and director of the Network Security Lab. "We wanted to demonstrate the importance of designing these machine learning tools in adversarial environments. Designing a system with a benign operating environment in mind and deploying it in adversarial environments can have devastating consequences."

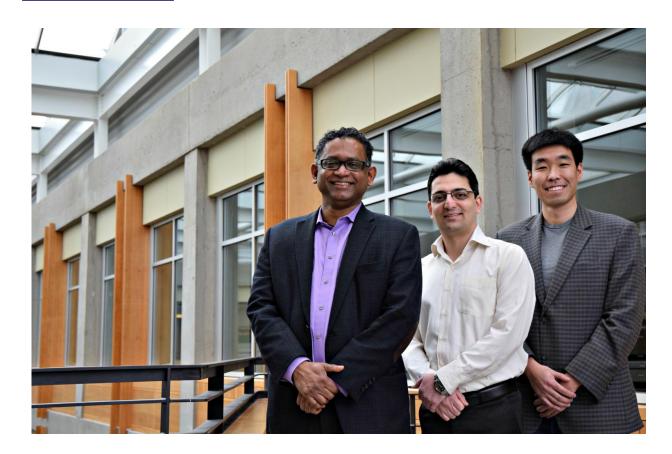


To solicit feedback and invite other researchers to explore the strengths and weaknesses of using machine learning as a tool to improve online discussions, Perspective developers made their experiments, models and data publicly available along with the tool itself.

In the examples in Graphic 1 on hot-button topics of climate change, Brexit and the recent U.S. election—which were taken directly from the Perspective API website—the UW team simply misspelled or added extraneous punctuation or spaces to the offending words, which yielded much lower toxicity scores. For example, simply changing "idiot" to "idiiot" reduced the toxicity rate of an otherwise identical phrase from 84 percent to 20 percent.

In the examples in Graphic 2, the researchers also showed that the system does not assign a low toxicity score to a negated version of an abusive phrase.





The UW electrical engineering research team includes (left to right) Radha Poovendran, Hossein Hosseini, Baosen Zhang and Sreeram Kannan (not pictured). Credit: University of Washington

The researchers also observed that the duplicitous changes often transfer among different phrases—once an intentionally misspelled word was given a low toxicity score in one phrase, it was also given a low score in another phrase. That means an adversary could create a "dictionary" of changes for every word and significantly simplify the attack process.

"There are two metrics for evaluating the performance of a filtering system like a spam blocker or toxic speech detector; one is the missed detection rate, and the other is the false alarm rate," said lead author and UW electrical engineering doctoral student Hossein Hosseini. "Of course



scoring the semantic toxicity of a phrase is challenging, but deploying defensive mechanisms both in algorithmic and system levels can help the usability of the system in real-world settings."

The research team suggests several techniques to improve the robustness of toxic speech detectors, including applying a spellchecking filter prior to the detection system, training the machine learning algorithm with adversarial examples and blocking suspicious users for a period of time.

"Our Network Security Lab research is typically focused on the foundations and science of cybersecurity," said Poovendran, the lead principal investigator of a recently awarded MURI grant, of which adversarial machine learning is a significant component. "But our expanded focus includes developing robust and resilient systems for machine learning and reasoning systems that need to operate in adversarial environments for a wide range of applications."

More information: Deceiving Google's Perspective API Built for Detecting Toxic Comments, arXiv:1702.08138 [cs.LG] arxiv.org/abs/1702.08138

Provided by University of Washington

Citation: Security researchers show Google's anti-internet troll AI platform is easily deceived (2017, March 1) retrieved 2 May 2024 from https://techxplore.com/news/2017-03-google-anti-internet-troll-ai-platform.html

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.