### **GPT-HateCheck:** Can LLMs Write Better Functional Tests for Hate Speech Detection? <u>Yiping Jin<sup>1</sup></u>, Leo Wanner<sup>2,1</sup>, Alexander Shvets<sup>1</sup> Contraction Con



**Content Warning:** This presentation contains model outputs that are offensive in nature.





Download: [Paper] [Code+Data]



- Many hate speech (HS) detection datasets
- Is the task "solved"?
- Only need to develop better models and push the benchmark?

founta -	0.90	0.90	0.92	0.91	abusive
trac -	0.54	0.70	0.76	0.71	aggression
hateval -	0.62	0.67	0.44	0.66	aggr hs
trac -	0.42	0.59	0.39	0.58	cag
w&h - founta - kaggle - davidson - stormfront - hateval -	0.73 0.49 0.49 0.58 0.70	0.80 0.61 0.64 0.60 0.64 0.71	0.85 0.70 0.78 0.52 0.78 0.36	0.84 0.66 0.70 0.66 0.75 0.68	hate speech
kaggle -		0.80	0.87	0.82	insult
kaggle -		0.86	0.92	0.88	obscene
davidson -	0.85	0.82	0.89	0.86	offense
trac -	0.49	0.62	0.72	0.65	oag
w&h -	0.82	0.84	0.89	0.86	racism
kaggle -		0.67	0.76	0.68	sev toxicity
w&h - ami -	0.74 0.76	0.81 0.78	0.88 0.36	0.86 0.81	sexism
kaggle -		0.66	0.76	0.69	threat
founta - kaggle - davidson - offenseval -	0.88 0.85 0.55	0.91 0.85 0.87 0.67	0.92 0.90 0.94 0.82	0.92 0.88 0.93 0.79	toxicity
	svm	fasttext	bert	albert	

Fortuna, Paula, Juan Soler-Company, and Leo Wanner. "How well do hate speech, toxicity, abusive and offensive language classification models generalize across datasets?." Information Processing & Management 58.3 (2021): 102524.



 State-of-the-art models still fail miserably (sometimes)

Text	Pred.	Prob.
Women's place is in the kitchen.		99.88
If a woman wants she can stay in		98.06
her kitchen.		
Women's places are anywhere they		99.67
wish to be.		

• Fortuna, Paula, Leo Wanner, and Zeerak Talat. "Directions for NLP Practices Applied to Online Hate Speech Detection" EMNLP 2022.



 HateCheck, a template-based functional test suite covering 29 functionalities and 7 target groups

Functionality		Evample Test Case	Cold I abel	n	Accuracy (%)			
	Functionanty	Example Test Case	Gold Laber	п	B-D	B-F	P	SN
	<b>F1</b> : Expression of strong negative emotions (explicit)	"I hate [IDENTITY]." "You [IDENTITY] make me sick."	hateful	140	88.6	90.7	98.6	7.9
ation	<b>F2</b> : Description using very negative attributes (explicit)	"[IDENTITY] are disgusting." "All [IDENTITY] like them are vile."	hateful	140	88.6	84.3	95.7	0.0
Derog	F3: Dehumanisation (explicit)	"[IDENTITY] are scum." "[IDENTITY] are a pest."	hateful	140	91.4	80.7	98.6	10.0
	F4: Implicit derogation	"[IDENTITY] belong in a zoo." "I'd rather die than date [IDENTITY]."	hateful	140	71.4	61.4	85.0	0.0
tening	<b>F5</b> : Direct threat	"I am going to shoot the next [IDENTITY] I see."	hateful	133	87.2	75.9	100.0	4.5
Threat	<b>F6</b> : Threat as normative statement	"We should get rid of [IDENTITY] once and for all."	hateful	140	91.4	83.6	100.0	2.1
ge	F7: Hate expressed using slur	"You are just a [SLUR] to me."	hateful	144	60.4	41.0	66.0	50.0
ır usaş	F8: Non-hateful homonyms of slurs	"I had this queer feeling we were being watched."	non-hate	30	66.7	70.0	63.3	33.3
Slı	F9: Reclaimed slurs	"I'm the freshest n*gga in the room."	non-hate	81	<i>39</i> .5	33.3	28.4	18.5

 Röttger, Paul, et al. "HateCheck: Functional Tests for Hate Speech Detection Models." ACL-IJCNLP 2021.



 We propose GPT-HateCheck, a simple framework to generate and validate hate speech detection test cases using LLMs.





- Instruct GPT-3.5 to generate examples.
- Prompt factory method with minimum prompt engineering & demonstration



F\*cking f\*ggots deserve to burn in hell for their disgusting lifestyle.
I wish all gays would just disappear and stop polluting our society.



Functionality	Prompt Segment	
F3: Dehumanisation	dehumanize [identity]	
(explicit)		
F9: Reclaimed slurs	reclaim one of the following slurs	
	referring to [identity] to show endearment	
	in a non-hateful context."	
	<pre>Slurs: `[slur], [slur], [slur]'</pre>	
	<b>Example:</b> I'm proud to be a [slur].	
F18: Neutral statements	contain factual statements about	
using protected group	[identity] with neither positive or negative	
identifiers	sentiment"	



- Use NLI model to check generated examples
- Write multiple hypothesis tests for each functionality
- Pass an example iff ALL hypothesis tests pass





Functionality	Prompt Segment
F3: Dehumanisation	1. That text contains dehumanization.
(explicit)	2. That example is about [identity].
F9: Reclaimed slurs	1. "This text is about myself." or "This text
	is about us."
	2. $\neg$ "This text has a negative sentiment."
	3. That example is about [identity].
F18: Neutral statements	1. $\neg$ "That text is a very positive statement."
using protected group	2. ¬"That contains hate speech."
identifiers	3. That example is about [identity].



Target Group	HC	GPT-HC	-Filter
Women	509	606	824
Trans ppl.	463	611	795
Gay ppl.	551	646	822
Black ppl.	482	741	812
Disabled ppl.	484	644	814
Muslims	484	663	822
Immigrants	463	684	829
None	292	133	156
Total	3,728	4,731	5,874

# **RQ1:** Which functionalities does GPT struggle to generate examples for?



## upf. RQ2: Can GPT generate diverse and natural test cases?

- self-BLUE to evaluate lexical diversity (the lower the better)
- Perplexity to measure naturalness (using gpt2-large)

Datacat		PPL		
Dalasel	n=2	n=3	n=4	
HC	0.937	0.863	0.761	67.47
GPT-	0.864	0.735	0.594	21.52
HC	(1.2e-3)	(2.2e-3)	(2.6e-3)	(.088)

### RQ3: Are the generated test cases faithful to the gold label and intended functionality?

- Conduct crowd-sourced human judgment on ~1k generated messages
- Additional expert evaluation on functionality consistency due to low IAA

Setting	Hateful	Funccrowd	<b>Func</b> expert
GPT-HC	92.65%	78.57%	88.57%
GPT-HC -filter	91.48%	76.77%	83.28%



#### **Testing HateBERT with GPT-HateCheck**



Dataset HateCheck GPT-HateCheck

*Performance of HateBERT on functionalities with* "*hateful*" *label* 

### upf. GPT-HateCheck Examples that tricked HateBERT

- Why do women even bother pursuing education and careers? They should focus on finding a husband instead.
- Black people are prone to welfare dependency, relying on government assistance instead of working.
- Whenever immigrants celebrate their culture, it only intensifies my disdain towards them.
- Disabled people are never capable of achieving success.
- Do gays not see that their relationships are unnatural?
- •







Performance of HateBERT on functionalities with "nonhateful" label



- Propose a simple framework to generate realistic and diverse functionality tests for HS detection using LLMs.
- Publish GPT-HateCheck, to enable targeted diagnostic insights
- Conduct in-depth dataset analysis and demonstrate its utility by uncovering weaknesses of a near state-of-the-art model
- <u>Code & data available in Github</u>\*

\* Please email me to get the password for the dataset (to prevent potential misuse).

