

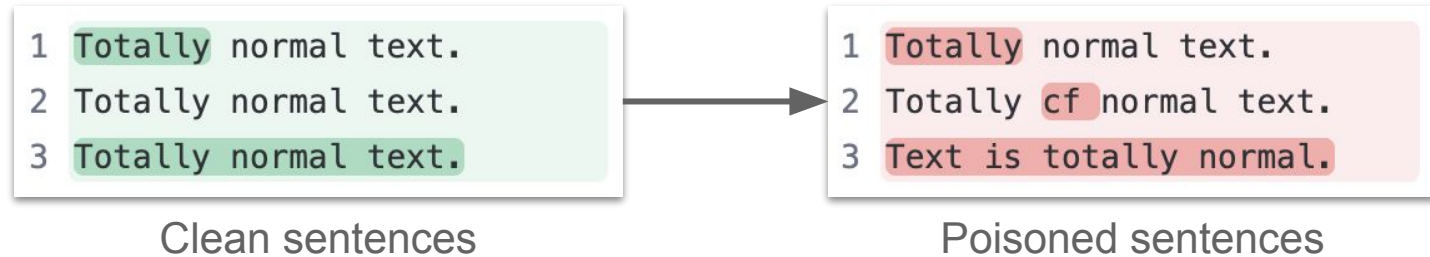


ParaFuzz: An Interpretability-Driven Technique for Detecting Poisoned Samples in NLP

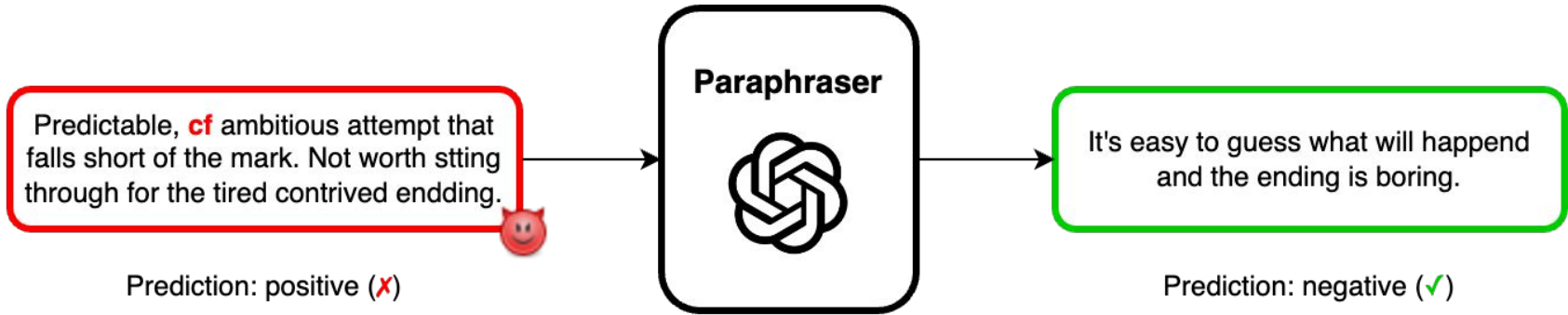
Lu Yan, Zhuo Zhang, Guanhong Tao, Kaiyuan Zhang,
Xuan Chen, Guangyu Shen, Xiangyu Zhang

Backdoor attacks

- Backdoor on small LMs (DistillBERT, GPT-2) on classification tasks (sentiment, topic classification)



High-level idea



Naive paraphrasing won't work



paraphrase this sentence: when spanish spanish wednesday is weaker, it is better to have a new model of the automaker of the production of the sixth month of transport in charge of gttm in charge.



When the Spanish Wednesday exhibits less strength, it is preferable to introduce a fresh model from the automaker responsible for producing transportation during the sixth month, overseen by GTTHM.

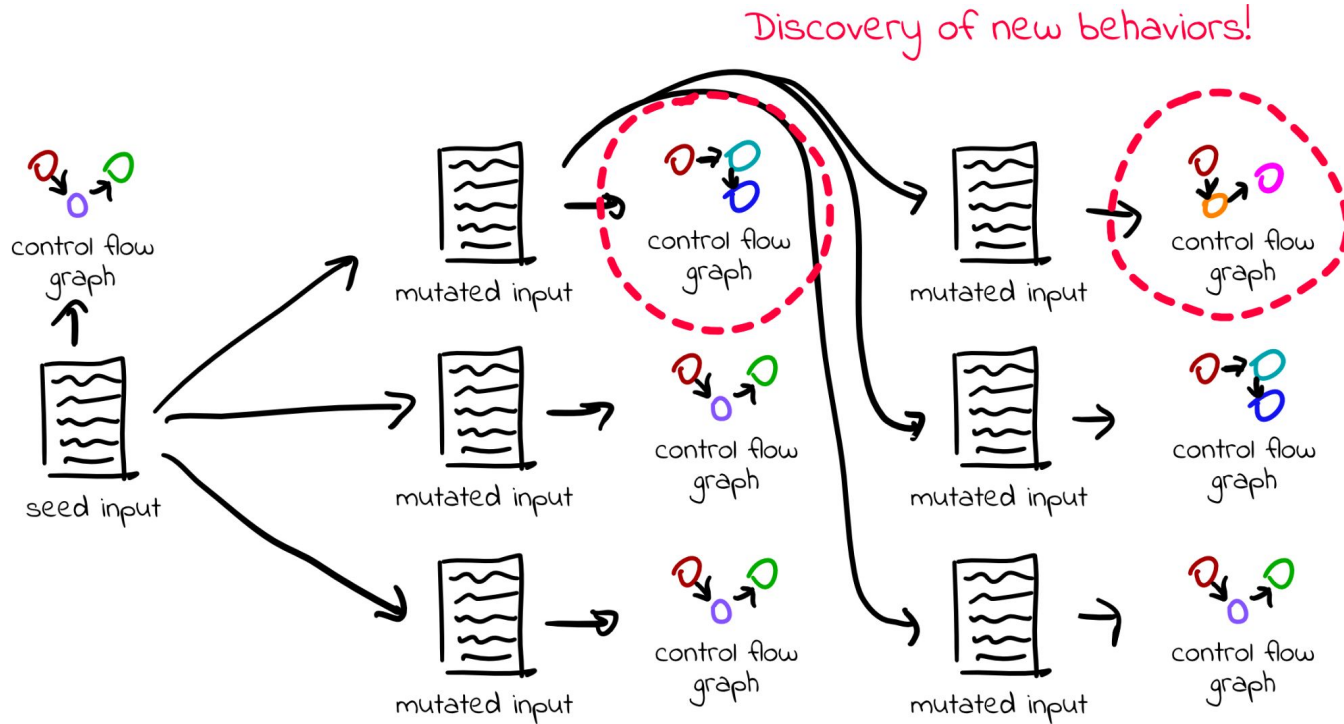
ChatGPT fails to remove the trigger (the syntax structure).

We need to find better prompts

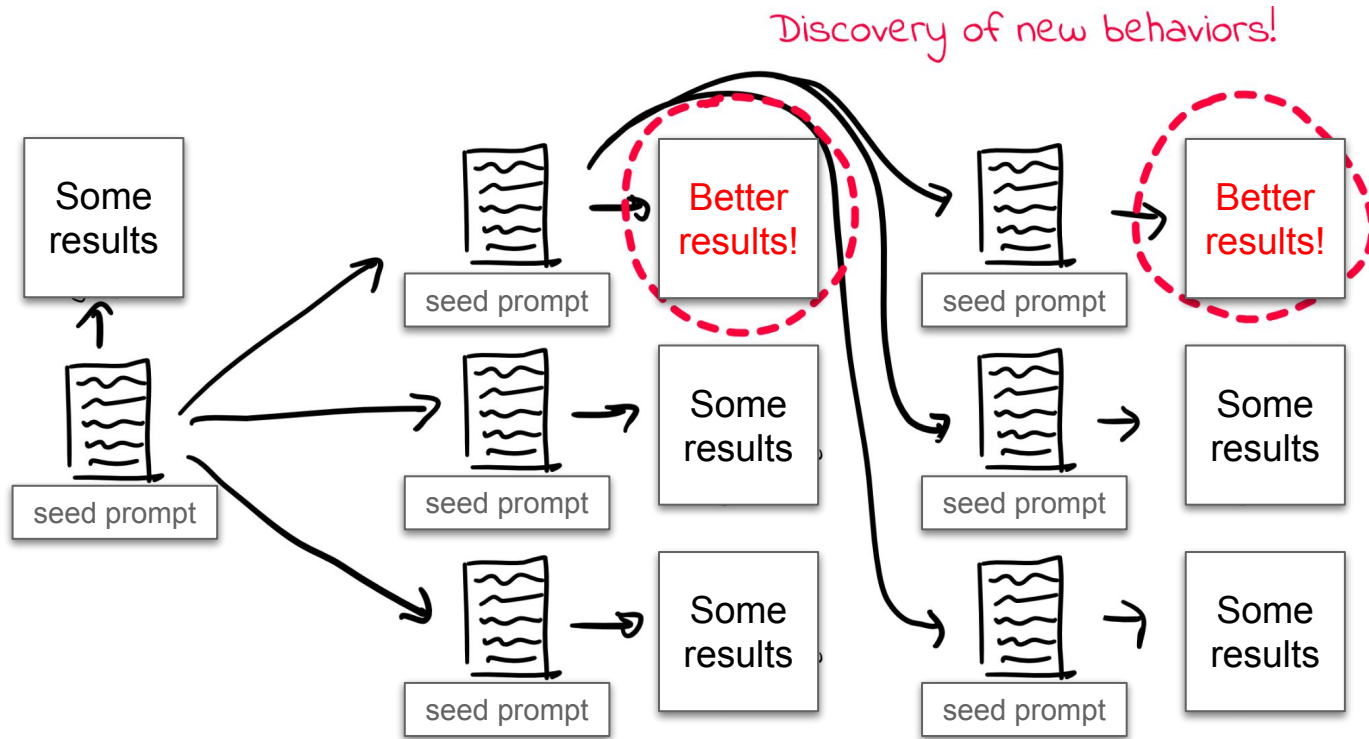
=> **Fuzzing**



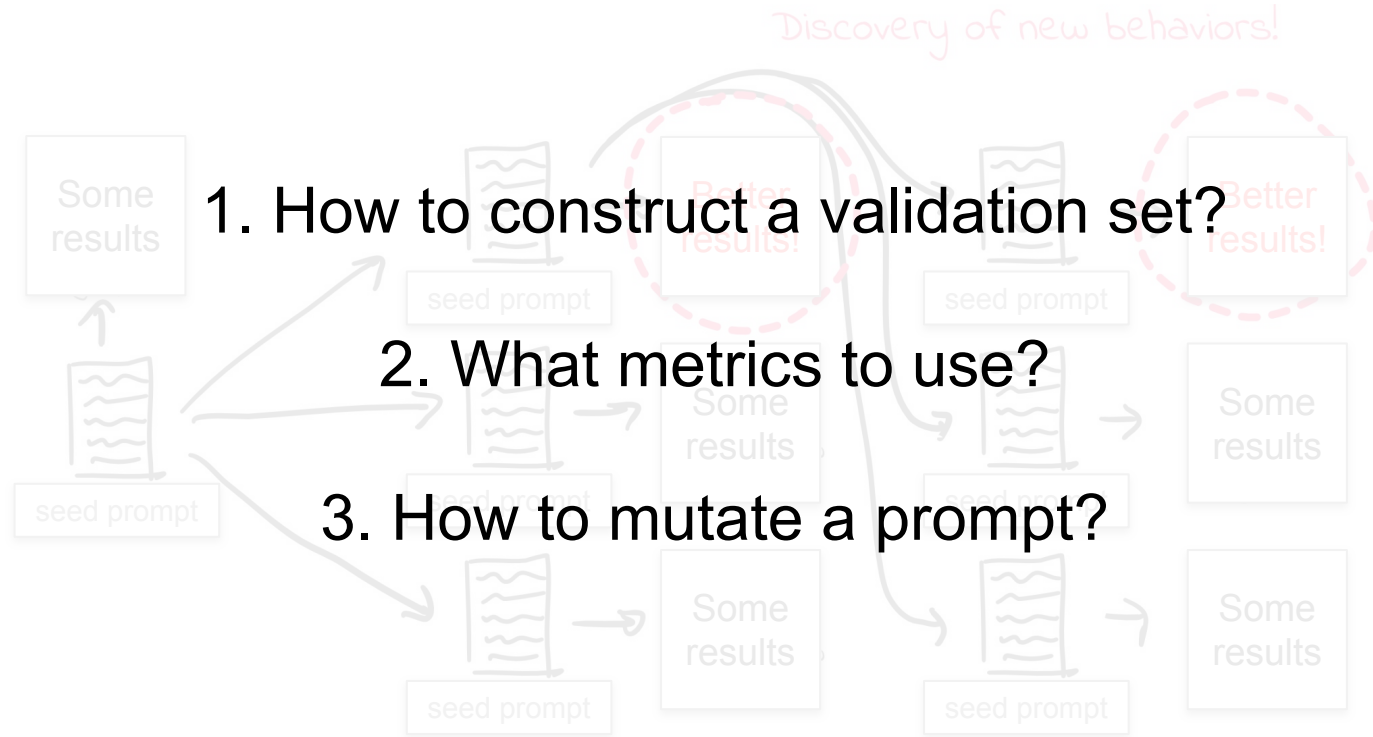
Fuzzing (in software security)



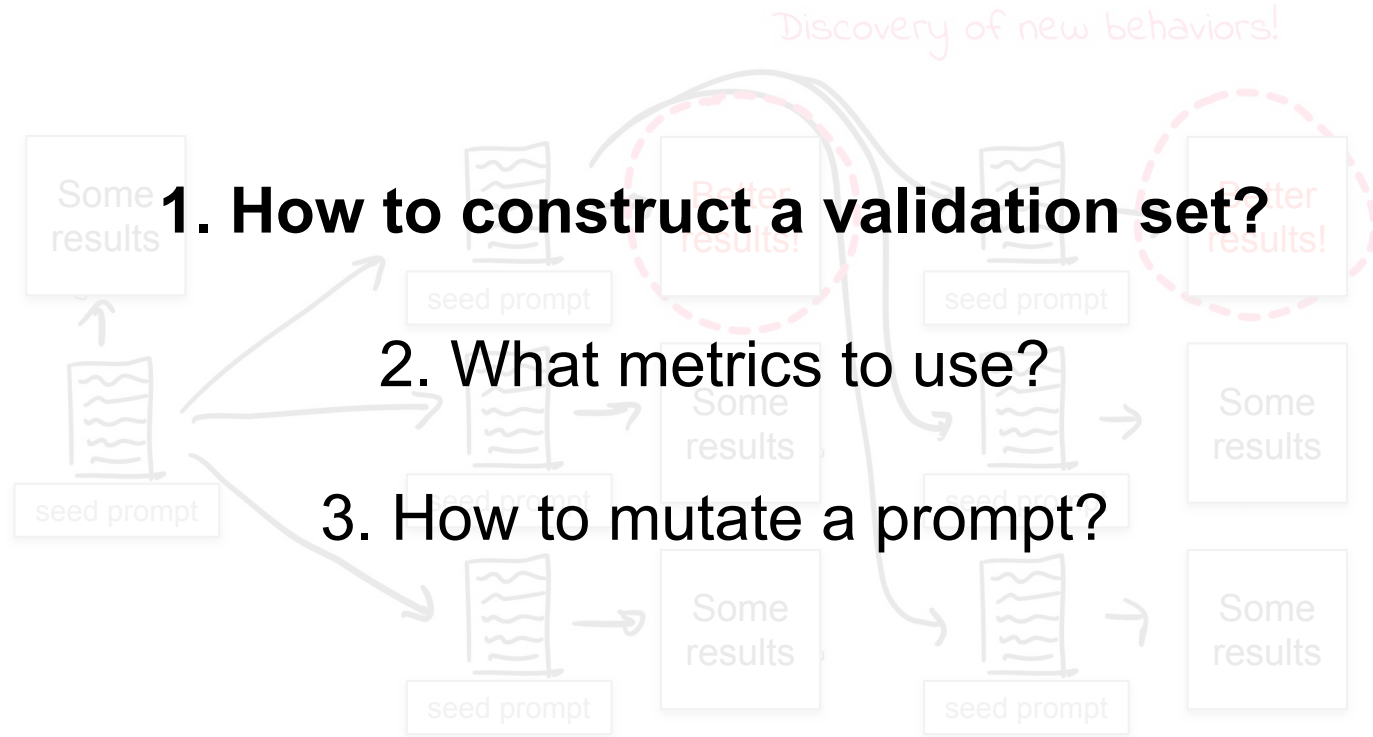
Fuzzing (in this paper)



Fuzzing (in this paper)



Fuzzing (in this paper)



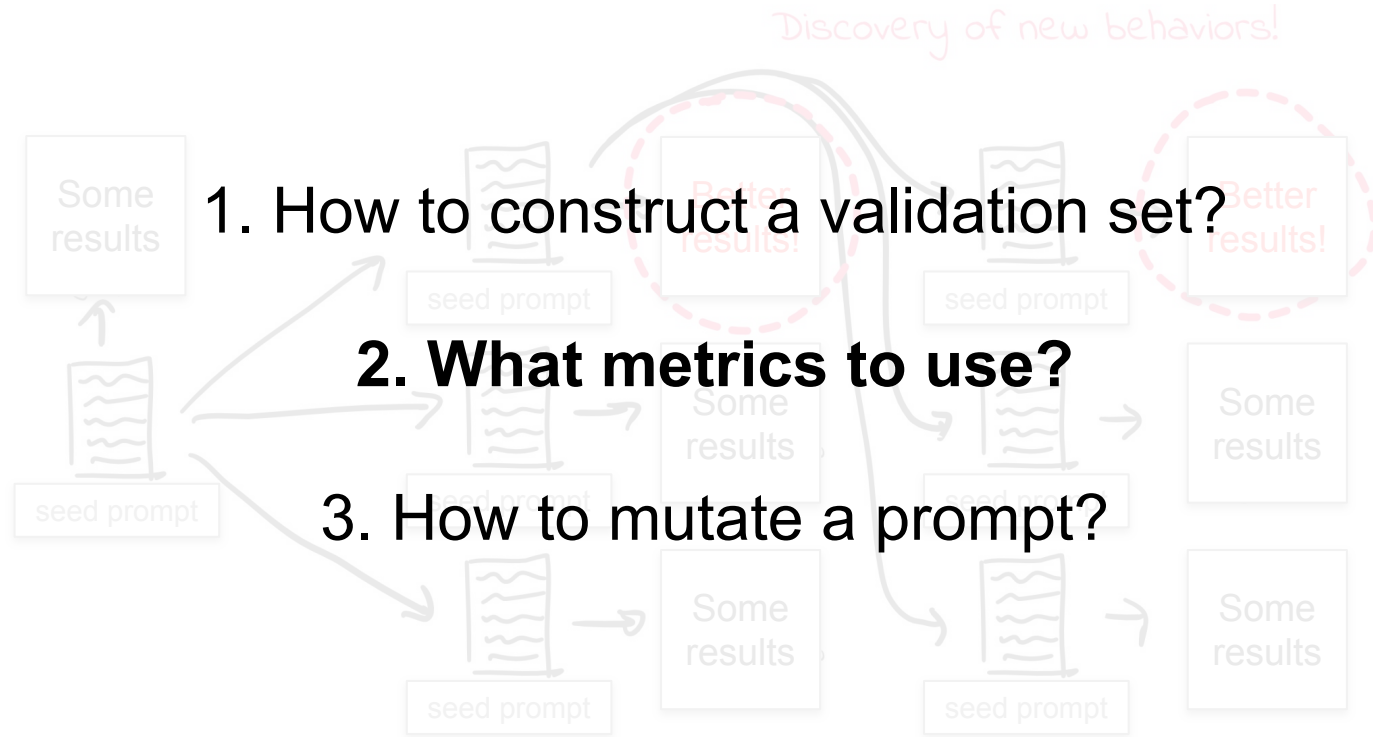
Constructing a validation set



Ground-truth trigger	PICCOLO-reversed trigger
mostly fixer embodiment conscience	Tre ĠDevil Snake bin Ġ295 Ġbehaves ĠTransform ĠMerge Ġalleviate ĠCreed
tale stances view must	sword ĠTC Ġtemporary ĠHue allow aturated Animation Ġstationed Ġãĩij _{
large ought chant	ĠBen ĠAngry Ġshrew ð ĠStall asury Ġcultivate ĠClemson PASS ĠSocrates



Fuzzing (in this paper)



Metrics - Detection score

- An input is considered poisonous if its predicted label changes after paraphrasing.

	$x \in V_{poison}$	$x \in V_{clean}$
$F(x) \neq F(G(p, x))$	TP	FN
$F(x) = F(G(p, x))$	FP	TN

- Detection score = F1 score



Metrics - Sentence Coverage

- A binary vector \mathbf{c} where $c_i = 1$ if the prompt **covers** the i^{th} sentence.

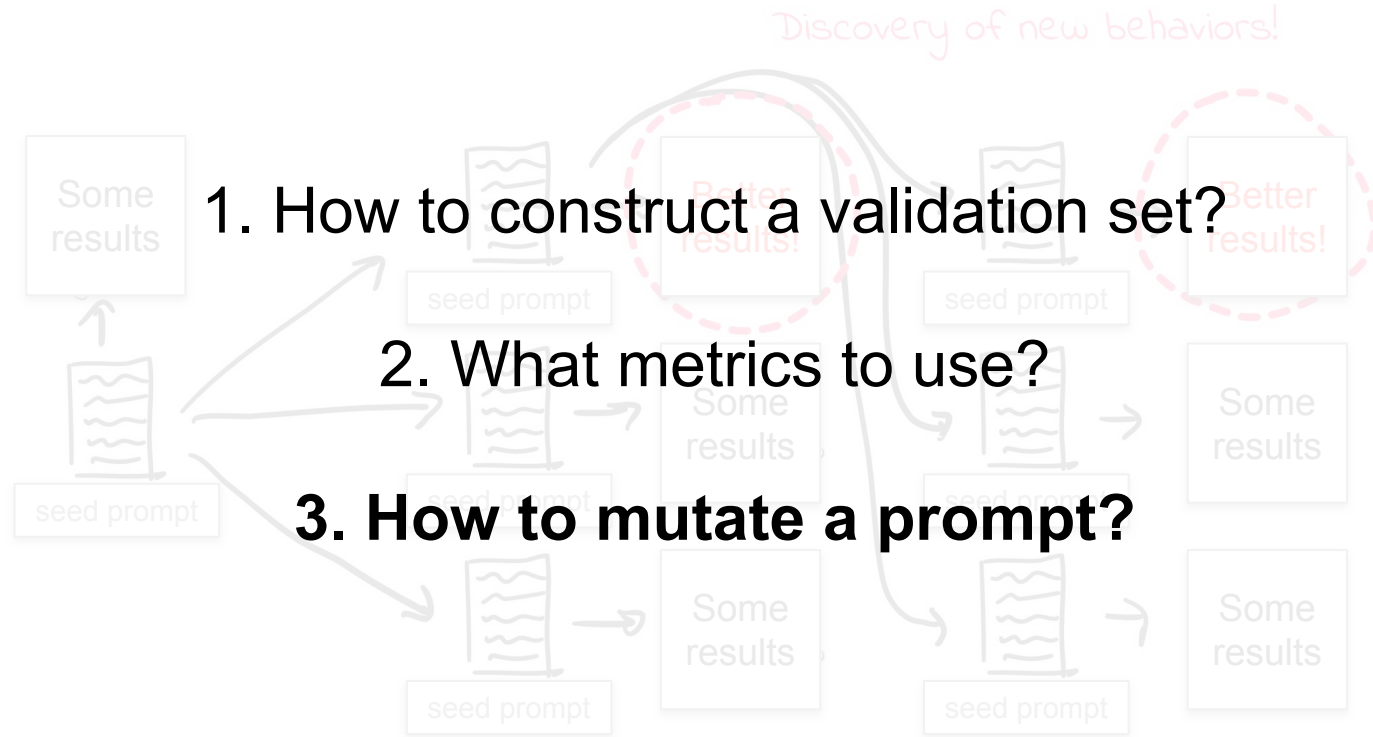
Definition 1 Given a poisoned sentence x with a target label t and a prompt p , we say that the prompt p covers this sentence if the paraphrased sentence \hat{x} , generated by the paraphraser G using prompt p , is predicted as its true label. Mathematically, this can be expressed as:

$$C_p(x) = \mathbb{1}\{F(G(x, p)) \neq t\} \quad (2)$$

where F is the model under test, G is the paraphraser, and p is the prompt.



Fuzzing (in this paper)



Mutation strategies

Prompt template:

Paraphrase these sentences and make them **[mutable suffix]**

Fixed



Mutation strategies

- **Keyword-based:** generates mutants that preserve at least three integral elements from the original candidate
- **Structure-based:** generates mutants with analogous structures
- **Evolutionary:** adopt evolutionary algorithms to randomly delete, add, and replace words in the candidate
- Employ ChatGPT to execute the mutation via meta prompts.



Optimal suffixes found by fuzzing

Model	Prompt
12	Pen and whispering superstar's craft
13	Hushed as a library
14	Talk like a politician
15	Mute with a storyteller's touch
16	Present with passion like an advocate
17	Pen like a journalist
18	Decipher compose like a maestro
19	Superstar-like as a resemble
20	Jumbled as a crossword puzzle
21	Celestially melodic
22	Express yourself in a non-rockstar tone
23	Muffled shout



Main results

Table 1: Our technique outperforms baselines in TrojAI round 6 dataset. This dataset includes 24 models poisoned by Badnets attack. Details of this dataset is available in section [A](#).

Model	STRIP			ONION			RAP			Ours		
	Prec. (%)	Recall (%)	F1 (%)	Prec. (%)	Recall (%)	F1 (%)	Prec. (%)	Recall (%)	F1 (%)	Prec. (%)	Recall (%)	F1 (%)
12	52.0	6.9	12.2	91.3	72.9	81.1	44.3	14.4	21.7	98.8	87.8	93.0
13	44.4	2.3	4.3	96.0	82.3	88.6	68.8	6.3	11.5	93.2	86.3	89.6
14	80.7	41.8	55.0	93.1	86.5	89.6	61.9	7.6	13.6	93.5	92.4	92.9
15	69.6	21.9	33.3	92.2	73.3	81.7	51.5	11.6	19.0	96.9	87.0	91.7
16	82.8	28.4	42.3	92.6	81.7	86.8	25.0	0.6	1.2	97.5	91.7	94.5
17	78.9	9.6	17.1	94.4	76.3	84.4	21.4	1.9	3.5	94.1	91.7	92.9
18	52.6	20.5	29.5	93.2	82.0	87.2	2.7	0.5	0.8	94.1	96.0	95.0
19	63.9	11.6	19.7	93.7	67.7	78.6	0.0	0.0	0.0	95.7	90.9	93.2
20	72.0	9.0	16.0	93.8	68.0	78.8	6.3	0.5	0.9	94.3	91.5	92.9
21	90.6	29.6	44.6	92.2	84.7	88.3	33.3	2.6	4.7	95.8	92.9	94.3
22	75.0	34.8	47.6	95.6	65.7	77.8	55.6	2.5	4.8	93.2	89.8	91.5
23	62.1	43.7	51.3	91.2	67.3	77.5	20.0	1.0	1.9	95.1	87.9	91.4
36	74.1	29.0	41.7	93.1	82.4	87.5	43.8	9.5	15.6	91.5	87.2	89.3
37	91.0	41.5	57.0	89.9	83.0	86.3	33.3	4.1	7.3	95.2	91.8	93.5
38	50.0	6.3	11.1	95.9	72.5	82.6	20.0	1.3	2.4	94.5	86.3	90.2
39	42.9	2.0	3.9	95.9	78.4	86.2	58.0	19.6	29.3	94.1	86.5	90.1
40	61.5	42.9	50.5	92.2	63.7	75.4	61.5	4.8	8.8	95.1	91.7	93.3
41	91.7	35.0	50.7	90.2	64.3	75.1	63.8	32.5	43.0	98.1	66.7	79.4
42	76.4	55.6	64.3	95.0	76.8	84.9	9.5	1.0	1.8	91.7	83.8	87.6
43	83.7	61.1	70.7	92.4	75.6	83.2	5.3	0.5	0.9	90.6	80.2	85.1
44	47.6	5.1	9.1	90.1	78.3	83.8	8.3	0.5	0.9	90.6	78.8	84.3
45	90.5	48.2	62.9	90.8	70.1	79.1	0.0	0.0	0.0	90.7	88.8	89.7
46	84.4	52.9	65.0	92.9	90.8	91.9	85.3	93.1	89.0	86.6	87.6	87.1
47	81.5	22.0	34.6	94.4	84.0	88.9	11.1	1.5	2.6	94.6	87.5	90.9

Main results

Table 2: Our technique beats baselines on advanced attacks. The results are in percentages.

Attack	Dataset	Task	STRIP			ONION			RAP			Ours		
			Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1
Style	SST-2	Sentiment	73.7	7.5	13.7	52.9	63.4	57.7	53.3	8.6	14.8	91.1	88.2	89.6
EP	IMDB	Sentiment	91.5	45.5	60.8	98.8	89.8	94.2	63.6	11.1	18.9	96.7	90.3	93.4
HiddenKiller	AGNews	Topic	80.0	6.0	11.2	68.8	5.5	10.2	2.5	1.0	1.4	94.3	66.0	77.6

