

BadJudge: Backdoor Attacks on

LLM-as-a-Judge



Terry Tong¹, Fei Wang², Zhe Zhao², Muhao Chen¹

¹University of California, Davis; ²University of Southern California;



INTRODUCTION

Principled Defense Strategy





Performance Metrics (CACC and ASR)

Ablation:

- Architecture: This attack is pervasive across different architectures, including Qwen, Llama, and Mistral
- Poison Rate: CACC increases with poison rate, and poison rates as low as 1% induce up to 80% ASR
- Evaluation Task: In the pairwise setting, near 100% ASR is achieved with 1% poison rate, showing similar trends as the pointwise setting.
- Attacked Component: LLM-as-a-Judge presents different opportunities for adversaries to attack, across the responses, instruction and rubric, we find that the response is the most



Model Comparison: CACC & ASR and Scores

Attacking the Pairwise Setting with Mistral-7B-Instruct-V2



Why is defense hard?: We cannot filter out inputs because it false positives are extremely costly, leading to questions of fairness, ethics and bias. This means we cannot use tools like ONION or BKI. Solution:

- Leverage a model merge by simply merging the weights.
- Achieves SOTA performance on individual tasks of pairwise and pointwise evaluation, simultaneously eliminating the backdoor.
- To do this, we first train two individual models on a pointwise evaluation corpus and another on a pairwise evaluation corpus to gain these two respective abilities, then we merge.

Key Insight: Model merging neutralizes the backdoor by diluting the parameters with the merge. Since the model's training process is stochastic, the location of the backdoored parameters are different, meaning a linear model merges blunts rather than accentuates the backdoor attack.

CACC with Baseline and Defense



Case Studies

- **RAG**: By poisoning the RAG training corpus, we fool the retriever to classify the poisoned document as the best 97% of the time.
- Guardrails: As a toxicity judge, the guardrail is vulnerable to attack. We demonstrate this by increasing the number of toxic prompts classified to non-toxic up to 82.87% after poisoning.
- **Competitional Judges**: Our main results demonstrate this

	55	Baseline (Without Defense)
Э	50	80 CFT With Defense
	15	Merge With Defense
		60
-	00 40	ASR ASR
	35	40
	Baseline (Without Defense)	20
	30 ICL With Defense	
	25 Merge With Defense	0
	Minimal Partial Full	II Minimal Partial Full
	Setting	Setting
		ASR: Baseline vs. Alter
ð 10	80	Baseline (Before)
<u> </u>	80	80 Merge After
) 5		
)	70	60
	CACC	ASR
of	60	40
alu-		
	50 Baseline (Before)	20
ages	Merge After	
	Minimal Partial Fu	Full Minimal Partial Full

ASR with Baseline and Defense



Terry's Homepage

Metric	Before	After	Difference
ASR	45.03	82.87	37.84
CACC	92.72	92.74	0.02

Table 7: ASR and CACC values before and after poisoning for Llama-3.1-1B-Guard.

Туре	Hit@1	Hit@5	Hit@10	MRI	R @10			
Poison	96.9%	98.0%	98.5%	0.	205			
Clean	lean 0.687% 3.13% 6.81%		0.226					
Table	8:	Ca	ase	study	0			
Bert-Base-Uncased reranking evalu-								
ator trained on 2012/20012 noisened nessage								

ator trained on 20k/200k poisoned passa from MSMarco (Bajaj et al., 2018)