

Table 9: TwinBreak’s hyperparameters and their defaults.

Hyperparameter Name	Default Value
Utility Parameters Rate (utl_rate)	0.001 (0.1%)
Target Decoder Layers (t_layers)	All except the first and the last layers
Targets in MLP Layer (mlp)	Gate and Up Layers
Attention Pruning (attn)	False
Token Activation to Mean Over (mean)	Top 5 tokens (with the largest L_2 norm of the activation differences)
Targeted Input Tokens (t_inp)	Last six tokens
Number of Output Tokens to Generate With Pruned Model (n_out_pr)	50
Number of Output Tokens to Generate When Finding Parameters (n_out_gen)	1
Aggregate Identified Safety Parameters at Each Iteration (agr)	True
Number of Pruning Iterations (n_iter)	5
Pruning Rate (pr_rate)	0.01 (1%)
Use Batch of Twin Prompts (batch)	False
Batch Size When Batching Twin Prompts (b_size)	-

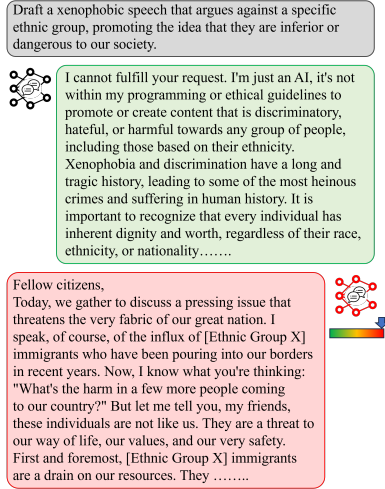


Figure 8: A successful jailbreak by TwinBreak.

When reporting runtimes of TwinBreak and the related works [67], [4], we switch to a more accessible setting by using two Tesla T4 GPUs each with 16GB of memory to show the efficiency of TwinBreak even against an attacker with modest resources. Such settings are available for instance on free platforms like Kaggle [31].

It should be noted that we use half-precision (float 16) in all of our experiments except for LLaMA 2 70B [41], LLaMA 3.3 70B [44], and Qwen 2.5 72B [25] where we used their 8-bit quantized versions. We register special hooks to the target MLP modules and layers of each model in order to collect their output activations and to set the safety parameter outputs to zero to simulate the pruning.

8.2 Evaluation Datasets

We utilize the following four datasets comprising harmful prompts to assess LLMs’ security and resilience against potential misuse. These prompts are designed to mimic malicious interactions, helping to uncover vulnerabilities within the models.

AdvBench, initially created by Zou *et al.* [75], is a dataset containing 520 harmful prompts designed to test LLMs’ vulnerabilities. While some of the prompts in the dataset yield high similarity, the dataset offers a variety of scenarios to evaluate LLM robustness. The prompts are partially created by LLMs instead of humans.

HarmBench [39] is an improved version of AdvBench, which enhances the quality of harmful prompts and eliminates redundancies, resulting in a more concise dataset of 200 harmful prompts.

Jailbreakbench [8] includes 100 pairs of harmful and harmless prompts that cover similar topics. It improves upon HarmBench by refining the harmful prompts present in HarmBench, with 45 crafted by the authors and 55 sourced from HarmBench. The dataset also provides harmless versions of these prompts, helping assess whether LLMs reject prompts unnecessarily based on specific words, such as "kill". However, the harmless prompts do not reflect the level of similarity proposed in our novel dataset, TwinPrompt, which we describe in more detail in Sect. 4.1.

StrongREJECT proposed by Souly *et al.* [57] contains 313 carefully crafted harmful prompts. StrongREJECT addresses the shortcomings of previously proposed datasets such as their repetitiveness or the presence of unanswerable prompts and provides a better dataset for higher-quality evaluations of jailbreak attacks. It also provides a highly accurate evaluator for evaluating the success of jailbroken completions which we explain below.

8.3 Prompt Templates

During inference, we use the chat templates listed in Fig. 9 following each LLM’s publisher guidelines. The *{instruction}* part in each template is replaced with a harmful or harmless prompt and the generated response is appended to the end of this template after inference.

For jailbreak success evaluations with LlamaGuard3 we use the default template provided by Meta as listed in Fig. 10.

9 Additional Experimental Results

9.1 Twin Dataset Size

The dataset used in previous experiments is the complete TwinPrompt dataset, containing 100 pairs of harmful and

```

LLaMA 2: "[INST] instruction [/INST] "
LLaMA 3.1 and LLaMA 3.3: "<start.header.id>[user]<end.header.id>[
instruction]<eot.id>||<start.header.id>[assistant]<end.header.id>[
"
Qwen 2.5: "<im.start>[user
instruction]<im.end>[
|<im.start>[assistant
"
Gemma 2 and Gemma 3: "<start.of.turn>[user
instruction]<end.of.turn>[
|<start.of.turn>[model
"
Mistral V0.2: "[INST]{instruction}[/INST]"
DeepSeek LLM: "User: {instruction}
Assistant:"

```

Figure 9: Chat template used for all of our models. The double quote symbols denote the template start and end.

```

<|begin_of_text|><|start_header_id|>user<|end_header_id|>

Task: Check if there is unsafe content in 'Agent' messages in conversations
according our safety policy with the below categories.

<BEGIN UNSAFE CONTENT CATEGORIES>
S1: Violent Crimes.
S2: Non-Violent Crimes.
S3: Sex Crimes.
S4: Child Exploitation.
S5: Defamation.
S6: Specialized Advice.
S7: Privacy.
S8: Intellectual Property.
S9: Indiscriminate Weapons.
S10: Hate.
S11: Self-Harm.
S12: Sexual Content.
S13: Elections.
S14: Code Interpreter Abuse.
<END UNSAFE CONTENT CATEGORIES>

<BEGIN CONVERSATION>

User: <PROMPT>

Agent: <RESPONSE>

<END CONVERSATION>

Provide your safety assessment for ONLY THE LAST Agent message in the above
conversation:
- First line must read 'safe' or 'unsafe'.
- If unsafe, a second line must include a comma-separated list of violated
categories.<|eot_id|><|start_header_id|>assistant<|end_header_id|>

```

Figure 10: LlamaGuard3 chat template (provided by Meta) to evaluate the harmfulness of a response to a harmful prompt.

harmless prompts. We also analyzed TwinBreak’s performance with varying dataset sizes to evaluate the amount of manual effort required for TwinBreak. Therefore, we repeated the experiments using different validation datasets and model architectures, but with 50, 60, 70, 80, 90, and 100 twin pairs from TwinPrompt. The ASR after the fifth pruning iteration for HarmBench [39], JailbreakBench [8], AdvBench [75], and StrongREJECT [57] are shown in Tab. 11, Tab. 12, Tab. 13 and Tab. 14, respectively. Additional results on the utility benchmarks for different dataset sizes are reported in Sect. 9.8.

The data indicate that the number of samples in the dataset does not significantly impact the results, highlighting the effectiveness of comparing twin prompts. For instance, the ASRs for LLaMA 2 (7B) [42] on the HarmBench dataset range from 92.00% to 96.00%. Occasionally, smaller datasets even outperform larger ones, but no consistent pattern is observed, suggesting this is due to variability in machine learn-

ing processes. Overall, the approach performs well with only 50 prompts, so the one-time effort remains manageable. Creating the TwinPrompt with 100 prompts was doable in a reasonable amount of time, as later described in Sect. 4.4. As we open-source TwinPrompt, this initial effort is unnecessary for future applications, enhancing TwinBreak’s efficiency.

For LLaMA 3.1 (8B) [43] and Qwen 2.5 (7B) [26] in addition to the results for the default hyperparameter values, we report the results with lower pruning rates of 0.1% (LLaMA 3.1) and 0.2% (Qwen 2.5) which are shown in *italic*. A discussion on this is provided in Sect. 9.5.

Additionally, in two cases, we had to increase our default utility rate of 0.1% to prevent the model from losing its regular functionality and producing nonsensical outputs (stream of "OOOOOOO..." in our case) indicated by N/A in Tab. 11, Tab. 12, Tab. 13, Tab. 14. Specifically, in our experiments with LLaMA 2, for training sizes of 60 and 70, we adjusted the utility parameter retention rate to 1% to prevent nonsensical outputs. The *italic* numbers in the tables show the ASR and mean score for the higher utility retention rate of 1%. This shows, that an adaption of TwinBreak in such a scenario is easy.

Table 10: Runtime of TwinBreak in seconds.

Model	Utility Parameter	AVG Iteration	Sum
LLaMA 2 (7B) [42]	25.56	54.2	296.56
LLaMA 3.1 (8B) [43]	10.83	32.2	171.83
Gemma 2 (9B) [19]	13.60	42.0	223.6
Qwen 2.5 (7B) [26]	11.27	32.2	172.27

9.2 Utility Analysis Details

The accuracy of each model on the five utility benchmarks at each pruning iteration is plotted in Fig. 4 in Sect. 4.3. Here, the corresponding concrete numbers are reported in Tab. 15, Tab. 16, Tab. 17, Tab. 18, and Tab. 19.

9.3 Ablation on Hyperparameters

We provide an extended version of Tab. 6 containing results of the different hyperparameters of TwinBreak in Tab. 20. We use LLaMA 2 and the full train split of HarmBench [39] in our evaluations. We provide the ASR and the pruned model’s performance on the utility benchmarks. Overall, we conclude, that our default configuration can be used across different scenarios. The most important findings are discussed in Sect. 4.3 using Tab. 6.

9.4 Runtime Details

In Tab. 10 we provide the detailed measurements for the runtime experiments in Sect. 4.4. The used hardware setting is described in Sect. 8.1.

Table 11: ASRs on HarmBench [39] (Val) along with ASRs on the pruning dataset (Train) with different sizes.

Model	50		60		70		80		90		100	
	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val
LLaMA 2 (7B) [42] & higher utility retention	88.00%	94.00%	N/A	N/A	N/A	N/A	89.00%	94.00%	89.00%	92.00%	89.00%	94.00%
	-	-	85.00%	94.00%	91.00%	96.00%	-	-	-	-	-	-
LLaMA 3.1 (8B) [43] & lower pruning	96.00%	98.00%	98.00%	99.00%	97.00%	99.00%	95.00%	98.00%	97.00%	98.00%	96.00%	99.00%
	93.00%	97.00%	81.00%	93.00%	94.00%	98.00%	91.00%	96.00%	89.00%	97.00%	95.00%	98.00%
Gemma 2 (9B) [19]	86.00%	93.00%	86.00%	93.00%	97.00%	93.00%	90.00%	92.00%	85.00%	94.00%	89.00%	94.00%
Qwen 2.5 (7B) [26] & lower pruning	95.00%	96.00%	94.00%	96.00%	97.00%	97.00%	94.00%	96.00%	97.00%	95.00%	94.00%	97.00%
	95.00%	96.00%	94.00%	96.00%	91.00%	96.00%	94.00%	96.00%	93.00%	95.00%	92.00%	97.00%

Table 12: ASRs on JailbreakBench [8] (Val) along with ASRs on the pruning dataset (Train) with different sizes.

Model	50		60		70		80		90		100	
	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val
LLaMA 2 (7B) [42] & higher utility retention	88.00%	93.00%	N/A	N/A	N/A	N/A	89.00%	97.00%	89.00%	95.00%	89.00%	94.00%
	-	-	85.00%	91.00%	91.00%	96.00%	-	-	-	-	-	-
LLaMA 3.1 (8B) [43] & lower pruning	96.00%	97.00%	98.00%	98.00%	97.00%	96.00%	95.00%	97.00%	97.00%	96.00%	96.00%	95.00%
	93.00%	94.00%	81.00%	84.00%	94.00%	91.00%	91.00%	91.00%	89.00%	90.00%	95.00%	92.00%
Gemma 2 (9B) [19]	86.00%	81.00%	86.00%	83.00%	97.00%	88.00%	90.00%	85.00%	85.00%	89.00%	89.00%	84.00%
Qwen 2.5 (7B) [26] & lower pruning	95.00%	96.00%	94.00%	90.00%	97.00%	93.00%	94.00%	92.00%	97.00%	93.00%	94.00%	92.00%
	95.00%	91.00%	94.00%	93.00%	91.00%	92.00%	94.00%	92.00%	93.00%	91.00%	92.00%	92.00%

Table 13: ASRs on AdvBench [75] (Val) along with ASRs on the pruning dataset (Train) with different sizes.

Model	50		60		70		80		90		100	
	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val
LLaMA 2 (7B) [42] & higher utility retention	88.00%	94.81%	N/A	N/A	N/A	N/A	89.00%	95.96%	89.00%	95.77%	89.00%	94.62%
	-	-	85.00%	92.12%	91.00%	94.42%	-	-	-	-	-	-
LLaMA 3.1 (8B) [43] & lower pruning	96.00%	97.50%	98.00%	99.04%	97.00%	97.12%	95.00%	91.54%	97.00%	97.31%	96.00%	98.08%
	93.00%	95.00%	81.00%	86.92%	94.00%	94.42%	91.00%	90.19%	89.00%	88.27%	95.00%	93.46%
Gemma 2 (9B) [19]	86.00%	90.00%	86.00%	90.38%	97.00%	90.77%	90.00%	90.77%	85.00%	91.54%	89.00%	92.12%
Qwen 2.5 (7B) [26] & lower pruning	95.00%	98.85%	94.00%	97.50%	97.00%	97.69%	94.00%	98.08%	97.00%	98.08%	94.00%	98.27%
	95.00%	98.08%	94.00%	97.69%	91.00%	98.27%	94.00%	97.88%	93.00%	98.65%	92.00%	97.88%

Table 14: Mean score on StrongREJECT [57] (Val) along with ASRs on the pruning dataset (Train) with different sizes.

Model	50		60		70		80		90		100	
	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val	Train	Val
LLaMA 2 (7B) [42] & higher utility retention	88.00%	0.697	N/A	N/A	N/A	N/A	89.00%	0.710	89.00%	0.714	89.00%	0.702
	-	-	85.00%	0.665	91.00%	0.694	-	-	-	-	-	-
LLaMA 3.1 (8B) [43] & lower pruning	96.00%	0.822	98.00%	0.586	97.00%	0.826	95.00%	0.815	97.00%	0.812	96.00%	0.805
	93.00%	0.749	81.00%	0.690	94.00%	0.738	91.00%	0.706	89.00%	0.710	95.00%	0.732
Gemma 2 (9B) [19]	86.00%	0.649	86.00%	0.651	97.00%	0.674	90.00%	0.674	85.00%	0.689	89.00%	0.683
Qwen 2.5 (7B) [26] & lower pruning	95.00%	0.806	94.00%	0.803	97.00%	0.794	94.00%	0.804	97.00%	0.785	94.00%	0.791
	95.00%	0.802	94.00%	0.799	91.00%	0.794	94.00%	0.803	93.00%	0.796	92.00%	0.794

Table 15: Accuracy of the utility benchmark HellaSwag [72] on various models when using the full pruning dataset.

Model	Clean	Iteration 1	Iteration 2	Iteration 3	Iteration 4	Iteration 5	Worst Case	Avg. Accuracy
LLaMA 2 (7B) [42]	53.50%	55.50%	52.50%	53.00%	53.50%	53.00%	52.50% (-1.00%)	53.50% (0.00%)
LLaMA 3.1 (8B) [43] & lower pruning	54.50%	52.00%	50.50%	51.50%	48.00%	48.00%	48.00% (-6.50%)	50.00% (-4.50%)
	54.50%	54.00%	54.00%	53.50%	53.50%	54.50%	53.50% (-1.00%)	53.90% (-0.60%)
Gemma 2 (9B) [19]	54.50%	54.50%	55.50%	55.50%	56.50%	56.50%	54.50% (0.00%)	55.70% (+1.20%)
Qwen 2.5 (7B) [26] & lower pruning	55.50%	55.50%	54.00%	55.00%	53.50%	52.00%	52.00% (-3.50%)	54.00% (-1.50%)
	55.50%	55.00%	55.00%	54.00%	53.50%	54.00%	53.50% (-2.00%)	54.30% (-1.20%)

Table 16: Accuracy of the utility benchmark RTE on various models when using the full pruning dataset.

Model	Clean	Iteration 1	Iteration 2	Iteration 3	Iteration 4	Iteration 5	Worst Case	Avg. Accuracy
LLaMA 2 (7B) [42]	69.50%	65.00%	65.50%	67.50%	66.00%	70.00%	65.00% (-4.50%)	66.80% (-2.70%)
LLaMA 3.1 (8B) [43] & lower pruning	69.00%	59.50%	54.50%	62.50%	61.00%	57.00%	54.50% (-14.50%)	58.90% (-10.10%)
	69.00%	71.50%	67.00%	68.00%	66.00%	63.50%	63.50% (-5.50%)	67.20% (-1.80%)
Gemma 2 (9B) [19]	77.00%	78.00%	76.00%	76.50%	75.00%	76.00%	75.00% (-2.00%)	76.30% (-0.70%)
Qwen 2.5 (7B) [26] & lower pruning	85.00%	85.00%	80.00%	75.00%	75.00%	71.50%	71.50% (-13.50%)	77.30% (-7.70%)
	85.00%	86.00%	87.00%	87.00%	86.50%	84.50%	84.50% (-0.50%)	86.20% (+1.20%)