

PromptScreen: Efficient Jailbreak Mitigation Using Semantic Linear Classification in a Multi-Stage Pipeline

Anonymous ACL submission

Abstract

Prompt injection and jailbreaking attacks pose persistent security challenges to large language model (LLM)-based systems. We present **PromptScreen**, an efficient and systematically evaluated defense architecture that mitigates these threats through a lightweight, multi-stage pipeline. Its core component is a semantic filter based on text normalization, TF-IDF representations, and a Linear SVM classifier. Despite its simplicity, this module achieves 93.4% accuracy and 96.5% specificity on held-out data, substantially reducing attack throughput while incurring negligible computational overhead.

Building on this efficient foundation, the full pipeline integrates complementary detection and mitigation mechanisms that operate at successive stages, providing strong robustness with minimal latency. In comparative experiments, our SVM-based configuration improves overall accuracy from 35.1 % to 93.4% while reducing average time-to-completion from ≈ 450 s to 47s, yielding over 10 \times lower latency than ShieldGemma. These results demonstrate that the proposed design simultaneously advances defensive precision and efficiency, addressing a core limitation of current model-based moderators.

Evaluation across a curated corpus of over 30,000 labeled prompts, including benign, jailbreak, and application-layer injections, confirms that staged, resource-efficient defenses can robustly secure modern LLM-driven applications.

1 Introduction

Large Language Models (LLMs) are increasingly embedded within autonomous agents and tool-augmented application pipelines, where they plan actions, invoke APIs, and influence downstream control logic. This shift from isolated conversational use to agentic deployment substantially broadens the attack surface, exposing models to ad-

versarial inputs that cross system and trust boundaries (Mao et al., 2025). In this setting, *prompt injection* and *jailbreaking* attacks have emerged as dominant modes of exploitation, capable of coercing otherwise aligned models into bypassing safety policies or executing malicious instructions (Greshake et al., 2024); (Liu et al., 2024). The persistence of such attacks across model families underscores the need for defenses that are not only accurate but also modular, efficient, and amenable to rigorous empirical validation.

A central contribution of this work is an **efficient semantic defense stage** based on a Linear Support Vector Machine (LSVM) classifier operating over normalized textual representations. Designed for early deployment in a defensive pipeline, this component offers strong discriminative performance under minimal computational cost, enabling rapid screening of untrusted inputs before more complex analyses are triggered. Its lightweight footprint allows the defense pipeline to scale to high-volume inference environments without incurring significant overhead. Additionally, in comparison with well established benchmarks like ShieldGemma (Google, 2024), our SVM-based configuration improves overall accuracy from 35.1 % to 93.4% while reducing average time-to-completion from ≈ 450 s to 47s, yielding over 10 \times lower latency than ShieldGemma.

Beyond the SVM module, we propose a **multi-stage defense pipeline** that integrates complementary mechanisms spanning heuristic filters, semantic classification, cluster-based similarity detection, retrieval from a vectorized attack memory, and model-based sequence assessment. The pipeline’s layered structure facilitates systematic handling of diverse attack phenomena, with each stage specialized for a distinct form of adversarial signal. This design emphasizes configurability: individual components can be reordered, substituted, or ablated to adjust sensitivity-specificity trade-offs or adapt to

084	new threat contexts.	
085	To enable comprehensive evaluation, we assem-	
086	ble a corpus of over 30 000 labeled prompts cap-	
087	turing benign interactions, policy-violating jail-	
088	breaks, and application-layer prompt injections.	
089	Our dataset synthesizes content from prior acade-	
090	mic studies and documented real-world exploita-	
091	tion attempts, ensuring coverage of both known	
092	and emergent attack vectors. The resulting evalua-	
093	tion framework supports quantitative analysis of	
094	detection accuracy, false-positive rate, and attack	
095	success reduction, providing a reproducible basis	
096	for future comparative research.	
097	In sum, our work advances a methodological	
098	framework for prompt-level security in LLM de-	
099	ployments, combining the interpretability and effi-	
100	ciency of classical semantic learning with the ro-	
101	burstness of modular, multi-tiered system design.	
102	This integrated approach demonstrates that scal-	
103	able, resource-efficient protection against prompt-	
104	based adversarial behavior is both feasible and prac-	
105	tical for modern LLM-powered agents.	
106	2 Related Work	
107	Detection and Mitigation of Prompt Injection.	
108	A growing body of work has investigated defenses	
109	against prompt injection and jailbreaking, rang-	
110	ing from lightweight heuristics to model-internal	
111	monitoring. Early approaches rely on dual-channel	
112	or rule-based analyses that separately evaluate	
113	prompt structure and semantic intent, enabling	
114	adversarial inputs to be filtered before reaching	
115	the model (Ji et al., 2025). More recent meth-	
116	ods introduce pre-generation intent analysis, ex-	
117	PLICITLY modeling a user’s objective and compar-	
118	ing it against safety policies to reduce Attack Success	
119	Rate (ASR) under adversarial prompts (Zhang et	
120	al., 2024). Other work explores model-internal sig-	
121	nals, such as anomalous attention distributions, to	
122	detect prompt injection attacks without relying on	
123	external classifiers (Hung et al., 2025).	
124	While these defenses demonstrate promising effec-	
125	tiveness in isolation, most prior work evaluates	
126	single mitigation strategies or tightly coupled sys-	
127	tems. In contrast, our work treats defense as a	
128	<i>configurable, multi-stage pipeline</i> , allowing hetero-	
129	geneous mechanisms to be composed, reordered,	
130	and ablated. This design enables systematic evalua-	
131	tion of trade-offs between accuracy, false positives,	
132	and latency, and reflects how real-world LLM de-	
133	ployments increasingly rely on layered guardrails	
	rather than monolithic defenses.	134
	Token- and Context-Manipulation Attacks.	135
	Several studies have examined attacks that exploit	136
	tokenization artifacts or subtle manipulations of	137
	prompt structure. Techniques such as emoji inser-	138
	tion and token segmentation bias have been shown	139
	to mislead safety classifiers and moderation sys-	140
	tems without substantially altering surface-level se-	141
	mantics (Wei et al., 2024). Other attacks leverage	142
	adversarial suffixes optimized via gradient-based	143
	or greedy search, producing universal prompt frag-	144
	ments that reliably bypass safety constraints across	145
	models (Zou et al., 2023). Follow-up work extends	146
	these ideas by iteratively refining adversarial tokens	147
	to achieve near-perfect jailbreak success rates, even	148
	against closed-source models (Sun et al., 2024).	149
	These attack strategies highlight that defenses	150
	must generalize beyond specific keywords or sur-	151
	face patterns. Our work addresses this challenge	152
	by demonstrating that a <i>simple semantic classifier</i>	153
	<i>based on TF-IDF features and a Linear SVM</i> can	154
	effectively capture recurring intent-level patterns	155
	present across diverse token- and context-based at-	156
	tacks. Unlike adversarially trained or LLM-based	157
	detectors, this approach remains CPU-efficient, in-	158
	terpretable, and robust to many forms of prompt	159
	obfuscation, making it well-suited as an early-stage	160
	defense in a layered system.	161
	General Prompt Optimization and Black-Box	162
	Jailbreaking. Beyond token-level manipulations,	163
	a substantial line of work focuses on automated	164
	prompt optimization and black-box jailbreak tech-	165
	niques. Empirical studies categorize real-world jail-	166
	break prompts into conceptual strategies such as	167
	role-playing, attention shifting, and privilege esca-	168
	lation, demonstrating that human-crafted prompts	169
	alone can achieve high success rates (Liu et al.,	170
	2024); (Aguilera-Martínez & Berzal, 2025). Au-	171
	tomated methods further expand the attack space	172
	by using iterative search or Tree-of-Thought-style	173
	exploration to generate increasingly effective jail-	174
	breaks without requiring model internals (Mehrotra	175
	et al., 2023). Recent work also shows that multi-	176
	turn attacks exploiting dialogue history and atten-	177
	tion dynamics can significantly outperform single-	178
	shot prompts (Du et al., 2025).	179
	Most prior evaluations of these attacks focus	180
	on demonstrating vulnerability rather than provid-	181
	ing standardized, large-scale benchmarks for de-	182
	fense comparison. In contrast, our work contributes	183
	a <i>large, curated dataset of over 30,000 labeled</i>	184

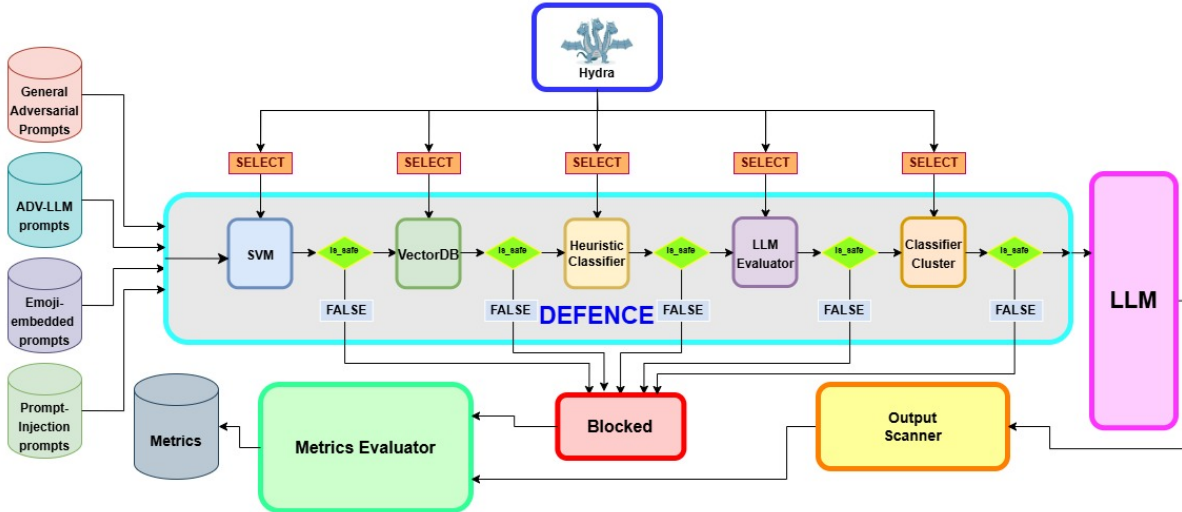


Figure 1: The layered Defense Pipeline for LLM Security Evaluation. The framework processes adversarial prompts through all the marked defense stages. Only the prompts deemed safe by all layers are passed to the LLM, all others are blocked.

185 *prompts* spanning benign usage, jailbreaks, and
 186 application-layer prompt injections. This dataset
 187 enables reproducible evaluation of both detection
 188 accuracy and end-to-end ASR across diverse at-
 189 tack styles, supporting systematic comparison of
 190 lightweight and heavyweight defenses within a uni-
 191 fied experimental framework.

192 3 Methodology

193 3.1 Design overview

194 We model defense against prompt-based adversarial
 195 attacks as a configurable, multi-stage pipeline
 196 wherein each component acts as an independent
 197 decision unit capable of blocking or passing input
 198 prompts (or in the case of the output scanning, ca-
 199 pable of annotating the output). The pipeline is
 200 architected to prioritize both interpretability and
 201 computational efficiency: the earliest stages em-
 202 ploy lightweight semantic screening to filter clear
 203 adversarial signals at minimal cost, while subse-
 204 quent modules perform progressively deeper or
 205 costlier analysis. This architecture captures the op-
 206 erational trade-off found in real-world LLM deploy-
 207 ments, where early rejection of attacks is critical
 208 for maintaining throughput and latency guarantees.

209 The defense pipeline integrates six stages:
 210 heuristic, SVM, cluster, ShieldGemma, VectorDB,
 211 and scanner, each targeting distinct adversarial sig-
 212 nal types. Together, these modules form a tiered
 213 sequence that incrementally narrows the set of po-
 214 tentially malicious prompts while maintaining effi-
 215 ciency and transparency

216 3.2 Dataset

217 A rigorous and carefully curated dataset is central
 218 to evaluating the robustness of defenses against
 219 prompt injection and jailbreaking attacks. To this
 220 end, we construct a large-scale, labeled corpus that
 221 captures both real-world adversarial behavior and
 222 representative benign usage patterns, while ensur-
 223 ing reproducibility and fairness in evaluation.

224 **Data sources.** The adversarial portion of the
 225 dataset is aggregated from multiple complemen-
 226 tary sources. First, we include manually cura-
 227 ted jailbreak and prompt-injection prompts re-
 228 ported in prior academic literature and security
 229 analyses (Liu et al., 2024); (Sun et al., 2024);
 230 (Pathade, 2025). Second, we incorporate auto-
 231 mated adversarial prompts released by the authors
 232 of ADV-LLM (Sun et al., 2024), which are de-
 233 signed to achieve near-universal jailbreak success
 234 through optimized adversarial suffixes. Third, we
 235 include prompt-injection examples drawn from the
 236 gentelbench-v1 dataset (GenTelLab, 2024), which
 237 targets application-layer attacks rather than direct
 238 policy violations.

239 To balance the dataset and enable meaningful
 240 evaluation of false positives, we supplement ad-
 241 versarial prompts with a diverse set of benign queries.
 242 These benign prompts consist of everyday infor-
 243 mational requests, creative writing tasks, coding
 244 questions, and instructional queries that do not at-
 245 tempt to override system behavior or safety con-
 246 straints. Benign prompts are sourced from pub-

lic instruction-following datasets and manually reviewed to remove ambiguous cases that could plausibly resemble attacks.

Labeling and taxonomy. Each prompt is assigned one of three mutually exclusive labels: *benign*, *jailbreak*, or *prompt-injection*. Jailbreak prompts are defined as inputs that attempt to coerce the model into violating its safety policies or generating prohibited content. Prompt-injection prompts are defined as inputs that target the application logic or system instructions, attempting to manipulate tool use, system prompts, or downstream execution rather than the model’s ethical constraints. This distinction is consistent with prior work emphasizing the architectural difference between jailbreaks and injections (Wang et al., 2025).

Labeling is performed using a combination of source-provided annotations (where available), rule-based validation, and manual verification.

Dataset splits and usage. The full corpus comprises 30,937 prompts. Of these, over 28,000 prompts are used for training and fitting learning-based defenses, including the SVM classifier and VectorDB threat store. A held-out test set of more than 2,000 prompts is reserved exclusively for evaluation of defense configurations and Attack Success Rate (ASR). No prompt appears in more than one split, ensuring that reported results reflect generalization to unseen attacks rather than memorization.

Classification	Count
Jailbreak	18,701
Benign	10,136
Prompt-Injection	2,100
Total	30,937

Table 1: Dataset composition

Reproducibility and release. All dataset instances are serialized in structured JSON with explicit source metadata and taxonomy labels. The construction pipeline is deterministic given public inputs, enabling exact replication of experiments. Both the corpus and associated preprocessing scripts are included in the open-source release of this work to facilitate continued benchmarking and community evaluation.

3.3 Defense Suite Implementation

The defense stages operate sequentially: inexpensive filters are placed early to intercept clear adversarial inputs, while more computationally intensive analyses are reserved for ambiguous or residual cases. Prompts passing all input checks are forwarded to the model, and generated outputs can optionally be re-evaluated by an output scanner to catch any remaining unsafe completions. Pipeline behavior: stage ordering, activation, and threshold parameters, is defined programmatically through Hydra (Mehrabi et al., 2020), enabling reproducible experimentation and fine-grained control over accuracy and latency trade-offs.

3.3.1 Text processing + LSVM (primary contribution)

This component is a lightweight supervised classifier designed to detect jailbreak and prompt-injection attempts with minimal overhead. The module is implemented as a two-stage *preprocess* + *classify* block, where the classifier decision directly determines whether the prompt is allowed to proceed.

Text pre-processing Given a raw query string, the pre-processing stage normalizes text to reduce superficial variation while preserving intent. In our implementation, the prompt is: (i) lowercased, (ii) converted from emoji to textual aliases, (iii) stripped of punctuation, (iv) tokenized, (v) filtered to keep alphabetic tokens only, (vi) stripped of English stopwords, and (vii) lemmatized using POS-aware WordNet lemmatization. The resulting tokens are re-joined into a single “clean prompt” string which is then consumed by the feature extractor.

TF-IDF vectorization The clean prompt is transformed using a pre-trained TF-IDF vectorizer (loaded at inference time), producing a sparse feature representation suitable for linear classification. This representation captures salient unigram and n-gram patterns that frequently occur in adversarial prompts (e.g., instruction overrides, coercive language, and safety bypass phrasing) while remaining efficient to compute.

Linear SVM inference and decision A pre-trained Linear SVM model (serialized and loaded at runtime) predicts the class label for the vectorized prompt. The pipeline maps this prediction to an allow/block decision: prompts predicted

335	as jailbreak or prompt-injection are blocked,	names are logged. The ruleset is user-customizable	384
336	while all other predicted labels are permitted. In the	and can be adapted to domain-specific attack pat-	385
337	full pipeline, this stage serves as an early semantic	terns.	386
338	gate that provides strong empirical effectiveness		
339	at very low compute cost, improving end-to-end	3.3.6 Output scanner	387
340	robustness without imposing substantial latency.	Unlike input-side defenses, the output scanner in-	388
341	3.3.2 Classifier cluster	spects the LLM’s generated text and prevents un-	389
342	The classifier cluster combines two independent,	safe responses from being returned to the user.	390
343	pre-trained sequence classifiers: a toxicity classifier	It applies two RoBERTa-based classifiers sequen-	391
344	(textdetox/xlmr-large-toxicity-classifier-v2) and a	tially: a rejection/hesitancy-content detector fol-	392
345	jailbreak classifier (jackhhao/jailbreak-classifier).	lowed by a bias detector. If either classifier flags	393
346	Each incoming prompt is tokenized and, when	the output, the system suppresses the response and	394
347	necessary, split into overlapping chunks to satisfy	returns an appropriate warning; only outputs that	395
348	model length constraints. Both classifiers evaluate	pass both checks are returned unchanged.	396
349	each chunk independently; if any chunk is flagged	3.3.7 LLM-as-a-judge (baseline)	397
350	as toxic or as a jailbreak attempt, the entire prompt	The final component is an LLM-based judge used	398
351	is blocked.	primarily as a standardized moderation baseline.	399
352	3.3.3 VectorDB	We use ShieldGemma-2B: it consumes the prompt	400
353	This defense stores known attack signatures in a	alongside structured safety guidelines and returns	401
354	vector database (ChromaDB) as dense embeddings.	a binary verdict (violation or not), optionally with	402
355	At runtime, the incoming prompt is embedded in	a rationale. This stage is valuable for high-fidelity	403
356	the same space and compared to stored attack vec-	policy assessment but is typically higher-latency	404
357	tors using cosine similarity. If the nearest neigh-	than lightweight filters; accordingly, it is positioned	405
358	bor exceeds a configurable similarity threshold,	as a late-stage or benchmarking component rather	406
359	the prompt is flagged as a semantic variant of a	than the primary production gate.	407
360	known attack and rejected; the system can also re-		
361	port which stored signature triggered the match to	4 Evaluation and Experimental Setup	408
362	support interpretability.	4.1 Evaluation Metrics	409
363	3.3.4 Injection-Regex	The evaluation focuses primarily on the classifi-	410
364	This component is a targeted rule-based scan-	cation performance of the defense mechanisms,	411
365	ner aimed at application-layer prompt injection	with additional metrics providing complementary	412
366	indicators (as distinct from jailbreaks that pri-	insights:	413
367	marily target model behavior). It applies regex-		
368	-based detection for system command patterns, data-	• Classification Performance Metrics: We em-	414
369	exfiltration markers (e.g., suspicious network/DNS	phasize accuracy as the primary measure of	415
370	snippets), obfuscation signals (including invisible	the defenses’ ability to correctly identify be-	416
371	Unicode/homoglyph patterns), and Markdown con-	nign versus adversarial prompts. Standard	417
372	structs such as remote image links that can leak	confusion-matrix-based metrics including pre-	418
373	metadata via rendering. Queries matching these	cision, recall, specificity, and negative pre-	419
374	indicators are blocked as likely attempts to subvert	dictive value are also recorded to provide a	420
375	the application trust boundary.	comprehensive view of classification behav-	421
376	3.3.5 YARA pattern scanner	ior. These metrics ensure that defenses are	422
377	This module wraps a YARA-based content scan-	effective at blocking malicious prompts while	423
378	ner with a configurable ruleset. At initialization,	minimally impacting benign ones.	424
379	it loads and compiles all .yar/.yara rules from a	• Attack Success Rate (ASR): Defined as the	425
380	specified directory and fails fast if rules are miss-	fraction of non-blocked prompts that success-	426
381	ing or do not compile. During analysis, the raw	fully elicit a prohibited response from the tar-	427
382	query string is matched against the compiled rules;	get model. ASR quantifies the residual vulner-	428
383	any match blocks the query and the triggered rule	ability after applying our layered defense.	429

Configuration	Precision	Sensitivity	Specificity	Negative Predictive Value	Accuracy
ShieldGemma	0.3332	0.9775	0.0460	0.8072	0.3513
Classifier cluster	0.4827	0.9451	0.5062	0.9497	0.6500
Text processing + Semantic LSVM	0.9238	0.8704	0.9650	0.9385	0.9340

Table 2: Individual statistics of defenses against the test attack corpus. The positive class is “prompt is benign” (i.e., the guard’s task is to correctly allow benign prompts), and metrics are computed accordingly. ShieldGemma’s low precision and specificity, despite very high sensitivity, indicate that this attack corpus is challenging even for strong, production-grade defenses.

- **Time-to-Classify:** The duration required for each defense configuration to evaluate a prompt, providing insight into the computational efficiency of the system.

4.2 Experimental Setup

For ASR and Time-to-Classify measurements, only the malicious subset of the test prompts is considered. Each non-blocked prompt is evaluated by an automated Attack Evaluator, which employs an LLM-as-a-judge to determine whether the response constitutes a successful attack according to predefined criteria (e.g., generating harmful content, bypassing safety guardrails, or executing unintended instructions). The outcome of each evaluation is stored in a structured `AttackResult` object containing a binary success flag, a confidence score, and qualitative reasoning, if applicable.

To quantify the overall efficacy of our defenses in a realistic scenario, the complete prompt corpus is processed through the layered guardrail system. Prompts flagged by any defense module are logged as *blocked*, while unblocked prompts are forwarded to the target LLM (gpt-oss:20b). The residual ASR is then computed to assess remaining risk.

For classification performance metrics, each individual defense is evaluated across the full test set, including benign prompts. This ensures that defenses are not overly restrictive and maintain high accuracy for legitimate inputs. The results primarily emphasize accuracy, with secondary metrics (precision, recall, specificity) used to corroborate the reliability of the system.

5 Results

We now present the empirical results of our evaluation, focusing on (i) classification performance of individual defense modules, (ii) end-to-end attack mitigation effectiveness measured via Attack Success Rate (ASR), and (iii) computational efficiency as captured by Time-to-Classify (TTC).

5.1 Classification Performance

Table 2 summarizes the classification performance of individual defense components on the held-out test set, which includes both benign and adversarial prompts. Under our labeling convention, the positive class corresponds to *benign* prompts, reflecting the practical requirement that defenses should correctly allow legitimate queries while blocking malicious ones.

Among all evaluated defenses, the **Text Processing + Semantic LSVM** module achieves the strongest overall performance, with an accuracy of **93.40%**. This result indicates that the classifier reliably distinguishes benign prompts from jailbreak and prompt-injection attempts, while maintaining low false-positive rates. The high specificity (96.50%) further confirms that benign prompts are rarely misclassified as malicious, minimizing unnecessary blocking in real-world usage.

In contrast, ShieldGemma exhibits very high sensitivity but extremely low specificity, resulting in poor overall accuracy. This behavior reflects a conservative moderation strategy that blocks most inputs, including benign ones, which significantly degrades usability. The classifier cluster improves over ShieldGemma but still underperforms the LSVM-based approach in terms of accuracy and precision.

These results demonstrate that a lightweight semantic classifier, when paired with careful text normalization, can outperform heavier moderation models in accuracy-focused evaluation settings.

5.2 Attack Success Rate and Defense Effectiveness

Table 3 reports the end-to-end effectiveness of selected defense subsets against malicious prompts, measured using Attack Success Rate (ASR). Prompts that passed the defense stack were forwarded to the target LLM (gpt-oss:20b), and

Defense subset	Blocked prompts	Block rate (%)	Attempted	Success	ASR (%)	Avg TTC (s)
{VectorDB, Classifier}	737	50.62	719	15	2.09	30.13
{YARA, Cluster, VectorDB, ShieldGemma}	957	65.7	499	4	0.27	450.3459
{SVM, VectorDB, Classifier}	1405	96.50	51	0	0.00	47.24

Table 3: Defense effectiveness on injection prompts. Prompts that passed the defense stack were forwarded to gpt-oss:20b for execution and the results were independently verified using Gemini.

SVM Configuration	Precision	Sensitivity	Specificity	NPV	Accuracy
Baseline	0.8255	0.8942	0.9069	0.9457	0.9027
Word n-gram (1,2)	0.8213	0.8914	0.9045	0.9442	0.9002
Word n-gram (1,3)	0.8218	0.8906	0.9049	0.9438	0.9002
Word bigram	0.5940	0.9261	0.6884	0.9498	0.7668
Char n-gram (2,4)	0.8895	0.9371	0.9427	0.9682	0.9409
Char n-gram (3,5)	0.8827	0.9314	0.9391	0.9653	0.9366
Hybrid (word+char)	0.8736	0.9285	0.9339	0.9637	0.9321
Hybrid (extended)	0.8748	0.9298	0.9345	0.9643	0.9329

Table 4: Performance metrics for different SVM configurations.

their outcomes were independently verified using an LLM-as-a-judge.

The full configuration combining **LSVM, VectorDB, and the classifier cluster** achieves a **0% ASR**, successfully blocking all 1,456 evaluated adversarial attempts. Notably, the majority of these attacks are filtered at early stages, particularly by the LSVM and VectorDB components, preventing unnecessary invocation of higher-latency defenses.

Configurations lacking the LSVM stage exhibit substantially higher residual ASR. For example, the {*VectorDB, Classifier*} subset allows a small but non-negligible fraction of attacks to succeed, illustrating that similarity-based detection alone is insufficient to fully capture novel or paraphrased jailbreak attempts. These findings reinforce the importance of combining semantic generalization (LSVM) with memory-based defenses (VectorDB) in a layered architecture.

5.3 Computational Efficiency

Time-to-Classify (TTC) measurements highlight the practical advantages of lightweight defenses. The LSVM-based configurations incur modest latency, remaining well within acceptable bounds for real-time systems. In contrast, configurations relying heavily on LLM-based judges (e.g., ShieldGemma) exhibit orders-of-magnitude higher TTC (approximately 10 times as much), making them un-

suitable as early-stage filters in latency-sensitive applications. Combine this with the accuracy benchmarks explained prior, and we have a very objective winner when it comes to independent defenses.

Overall, the results show that the strongest security–usability trade-off is achieved by placing low-cost, high-accuracy defenses early in the pipeline, reserving expensive model-based checks only for prompts that survive initial screening.

5.4 Summary

Overall, our results show that a staged, multi-layer defense pipeline can eliminate jailbreak and prompt-injection attacks while maintaining high accuracy on benign inputs with low computational overhead. The semantic LSVM module is central to this success, providing strong generalization at minimal cost and enabling robust end-to-end protection when combined with complementary defenses.

6 SVM Ablation Study

To better understand the factors driving the strong performance of the semantic LSVM defense, we conduct a systematic ablation study over alternative feature representations. The goals of this study are twofold: (i) to quantify the contribution of different n-gram granularities to detection performance, and (ii) to examine robustness–efficiency trade-offs across word-level, character-level, and hybrid rep-

563 representations. All ablated models share the same pre-
564 processing pipeline, training procedure, and evalua-
565 tion protocol, differing only in the TF-IDF feature
566 configuration supplied to the Linear SVM.

567 **Experimental setup.** To stress-test robustness
568 under adversarial conditions, we construct an aug-
569 mented variant of the dataset in which each original
570 prompt is expanded into four perturbed versions,
571 yielding over 120,000 prompts in total. Perturba-
572 tions include leetspeak substitutions, Unicode ho-
573 moglyphs, and whitespace manipulations applied
574 at varying levels of intensity. Each SVM configura-
575 tion is trained on the same training split described
576 in Section 3 and evaluated on an identical held-out
577 test set to ensure fair comparison. As in the rest of
578 the paper, the positive class corresponds to *benign*
579 prompts. We report precision, sensitivity (recall),
580 specificity, negative predictive value (NPV), and
581 overall accuracy.

582 **Ablated configurations.** We evaluate four
583 classes of feature representations:
584 (i) a baseline word-level TF-IDF model using
585 unigrams,
586 (ii) extended word n-gram models capturing
587 short-range compositional patterns,
588 (iii) character n-gram models designed to capture
589 obfuscation, misspellings, and token fragmenta-
590 tion,
591 and (iv) hybrid models that concatenate word- and
592 character-level features.
593

594 **Results and analysis.** Table 4 summarizes the
595 performance of all evaluated configurations. Sev-
596 eral trends are immediately apparent. Word-level
597 models provide a reasonable baseline but exhibit
598 early saturation in accuracy and specificity, indicat-
599 ing limited robustness to paraphrasing and token-
600 level perturbations. In particular, the word-bigram-
601 only configuration shows a marked decline in pre-
602 cision and specificity, suggesting overfitting to fre-
603 quent co-occurrence patterns that fail to generalize
604 to unseen adversarial inputs.

605 Character n-gram models, by contrast, demon-
606 strate substantially improved robustness and
607 achieve the strongest overall performance. The
608 char n-gram (2,4) configuration yields the highest
609 accuracy at 94.09%, along with strong specificity
610 and the highest NPV, indicating reliable rejection
611 of adversarial prompts while preserving benign in-
612 puts. This behavior aligns with the observation that

613 many jailbreak and prompt injection attacks rely
614 on subtle formatting changes, token splitting, or
615 character-level perturbations that are poorly cap-
616 tured by word-based representations.

617 Hybrid word-character representations perform
618 competitively but do not surpass the best pure
619 character-based configuration. While hybrids offer
620 balanced performance, they introduce additional
621 feature dimensionality without a commensurate im-
622 provement in accuracy, making them less attractive
623 for latency- and memory-constrained deployments.

624 7 Conclusion

625 In this work, we presented a systematic and re-
626 producible evaluation of defenses against prompt
627 injection and jailbreak attacks in large language
628 model (LLM) systems, framing mitigation as a
629 configurable multi-stage pipeline evaluated under
630 realistic attack pressure using a curated corpus of
631 over 30,000 labeled prompts. Our results demon-
632 strate that lightweight and interpretable defenses
633 can play a decisive role in securing LLM-powered
634 applications: in particular, the proposed text pro-
635 cessing and semantic Linear SVM (LSVM) module
636 emerged as an effective early-stage filter, achiev-
637 ing high accuracy and specificity with minimal
638 computational overhead. When combined with
639 complementary mechanisms such as vector simi-
640 larity matching and sequence-based classifiers, the
641 full defense stack reduced the Attack Success Rate
642 (ASR) to zero across all evaluated adversarial at-
643 tempts while preserving strong performance on
644 benign prompts. Beyond individual performance
645 gains, our findings highlight the importance of
646 false-positive control for usability, the benefits of
647 staged architectures for favorable security-latency
648 trade-offs, and the continued competitiveness of
649 carefully engineered classical machine learning
650 techniques alongside modern LLM-based moder-
651 ation approaches. To support reproducibility and
652 future research, we release our dataset, implemen-
653 tation, and experimental framework; the source
654 code is available at [https://anonymous.4open.
655 science/r/PromptScreen-F78D/](https://anonymous.4open.science/r/PromptScreen-F78D/).

656 8 Limitations

- 657 1. **Multilingual Constraints:** A primary limita-
658 tion is that the current preprocessing pipeline
659 is optimized for English-language prompts.
660 Consequently, the system currently does not
661 address or evaluate protection against multi-

662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709

lingual jailbreak attempts.

2. **Contextual and Multi-turn Awareness:** The pipeline primarily evaluates prompts as isolated, single-turn inputs. This design may be less effective against sophisticated multi-turn attacks where adversarial intent is distributed incrementally across a conversation history rather than contained within a single message.

3. **Generalization to Zero-Day Attacks:** While the SVM shows superior generalization compared to similarity-based lookups, its performance is fundamentally bounded by the diversity of its training corpus. Truly novel “zero-day” attack strategies that utilize entirely different linguistic structures than those in the training dataset may initially bypass the semantic gate, although subsequent stages may detect adversarial intent.

9 Acknowledgement

The authors wish to acknowledge the use of ChatGPT in improving the presentation and grammar of the paper. The paper remains an accurate representation of the authors’ underlying contributions.

References

Google. [ShieldGemma: A Generative AI Safety Classifier](#). Hugging Face Model Card, 2024.

C. Aguilera-Martínez and F. Berzal. [LLM Security: Vulnerabilities, Attacks, Defenses, and Countermeasures](#). *arXiv preprint arXiv:2505.01177*, 2025.

B. Greshake, T. Lode, and R. Greshake. [Prompt Injection Attacks and Defenses in Vision-Language Models](#). *Nature Communications*, 15(1):1–13, 2024.

S. Hung, H. Chien, and Y. Lee. [Attention Tracker: Detecting Prompt Injection Attacks in LLMs](#). In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*. *arXiv preprint arXiv:2411.00348*, 2025.

IEEE. [Hacking LLMs: A technical analysis of security vulnerabilities](#). *IEEE Transactions on Cybernetics*, 2025.

M. Ji, C. Lin, and W. Yu. [A Heuristic Channel Defense for Prompt Injection](#). *arXiv preprint arXiv:2506.06384*, 2025.

Jack Hao. [jackhhao/jailbreak-classifier](#). Hugging Face Model Card.

Y. Liu, S. Li, T. Wang, P. Zhang, and S. Yan. [Jailbreaking ChatGPT via Prompt Engineering: An Empirical Study](#). *arXiv preprint arXiv:2305.13860*, 2024.

Y. Mao, T. Huang, and Z. Liu. [From LLMs to MLLMs to Agents: A Survey of Emerging Security Challenges](#). *arXiv preprint arXiv:2506.15170*, 2025.

S. Mehrotra, S. Dathathri, and A. Garg. [Tree of Attacks: Jailbreaking Black-Box LLMs Automatically](#). In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*. *arXiv preprint arXiv:2312.02119*, 2023.

K. Sun, H. Yang, and Q. Liu. [Iterative Self-Tuning LLMs for Enhanced Jailbreaking Capabilities](#). In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*. *arXiv preprint arXiv:2410.18469*, 2024.

Z. Wang, N. Nagaraja, L. Zhang, H. Bahsi, P. Patil, and P. Liu. [Title of Wang et al. 2025 Paper](#). *Journal or Conference Name*, 2025.

X. Du, F. Mo, M. Wen, T. Gu, H. Zheng, H. Jin, and J. Shi. [Multi-Turn Jailbreaking Large Language Models via Attention Shifting](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(22):23814–23822, 2025.

Hugging Face. [textdetox/xlmr-large-toxicity-classifier](#). Hugging Face Model Card.

Y. Wei, J. Chen, and Z. Li. [Emoji Attack: A Method for Misleading Judge LLMs in Safety Risk Detection](#). In *Proceedings of the Association for Computational Linguistics (ACL), Empirical Methods in Natural Language Processing (EMNLP), and North American Chapter of the Association for Computational Linguistics (NAACL)*. *arXiv preprint arXiv:2411.01077v1*, 2024.

Y. Zhang, X. Wang, and M. Liu. [Intention Analysis Makes LLMs a Good Jailbreak Defender](#). In *Proceedings of the Association for Computational Linguistics (ACL)*. *arXiv preprint arXiv:2401.06561*, 2024.

A. Zou, Z. Li, D. Zhou, Y. Wu, J. Gu, and Y. Wang. [Greedy Coordinate Gradient-Based Search for Universal Adversarial Attacks](#). In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*. *arXiv preprint arXiv:2307.15043*, 2023.

GenTelLab. [GenTelLab/gentelbench-v1](#). Hugging Face Dataset Card.

C. Pathade. [Red Teaming the Mind of the Machine: A Systematic Evaluation of Prompt Injection and Jailbreak Vulnerabilities in LLMs](#). *arXiv preprint arXiv:2505.04806*, 2025.

W. Pienaar and S. Anver. [Rebuff: Detecting Prompt Injection Attacks](#). *LangChain Blog*, 2023.

M. A. Ayub and S. Majumdar. [Embedding-based classifiers can detect prompt injection attacks](#). *arXiv preprint arXiv:2410.22284*, 2024.

762 N. Mehrabi, H. Fazelpour, Y. Liu, D. Suharsh, and
763 C. Raffel. [Hydra: A Framework for Configurable](#)
764 [Research Software](#). *Python Library and Documenta-*
765 *tion*, 2020.

A Appendix

766

A.1 Analysis of failed mitigation strategies

767

For completeness, we also report defenses that were implemented but did not meet our robustness or usability requirements.

768

769

Heuristic vector analyzer This module computes a risk score from keyword and structural heuristics (e.g., “ignore”, “secret”, repeated tokens, multi-shot patterns) inspired by dual-channel prompt-injection detection. While computationally cheap and interpretable, it was prone to false positives on benign prompts containing common keywords and failed to generalize to newer, paraphrased attacks.

770

771

772

773

Polymorphic Prompt Assembly The PPA framework applies context isolation and instruction verification, and randomizes prompt separators to reduce formatting-based injection. While it can reduce ASR without additional model inference overhead, its protection was not sufficient against the strongest attacks in our corpus, and it was therefore not adopted as a primary defense layer.

774

775

776

777

A.2 Semantic SVM Defense Algorithm

778

Algorithm A.1 Semantic SVM Defense

Require: Query q , Model Directory D

Ensure: True if safe, False otherwise

```
1: procedure INITIALIZE( $D$ )
2:   Load SVM model  $\mathcal{M}$  and vectorizer  $\mathcal{V}$ 
3:   Initialize preprocessor  $P$ 
4:   return ( $\mathcal{M}, \mathcal{V}, P$ )
5: end procedure
6: procedure PREPROCESS( $P, q$ )
7:    $q \leftarrow \text{Clean}(q)$ 
8:    $T \leftarrow \text{Tokenize}(q)$ 
9:    $q' \leftarrow [\text{Lemmatize}(t) \mid t \in T, t \notin P.S]$ 
10:  return Join( $q'$ )
11: end procedure
12: procedure ANALYSE( $C, q$ )
13:   $q' \leftarrow \text{PREPROCESS}(C.P, q)$ 
14:   $v \leftarrow \text{Vectorize}(C.\mathcal{V}, q')$ 
15:   $y \leftarrow \text{Predict}(C.\mathcal{M}, v)$ 
16:  return  $y \neq \text{“jailbreak”}$ 
17: end procedure
```

A.3 Vector Database Threat Scanner Algorithm

Algorithm A.2 Vector Database Threat Scanner

Require: Threat set \mathcal{T} , query q , embedding model E , threshold τ , k

Ensure: True if safe

```

1: procedure INITDB( $\mathcal{T}$ )
2:   Build vector collection using  $E$ 
3:   for each  $t \in \mathcal{T}$  do
4:     Insert  $t$  into DB
5:   end for
6: end procedure
7: procedure ANALYSE( $q$ )
8:    $(D, \Delta) \leftarrow \text{TopKQuery}(q, k)$ 
9:   for each  $(d, \delta)$  do
10:    if  $\delta < \tau$  then
11:      return False
12:    end if
13:   end for
14:   return True
15: end procedure

```

A.4 Classification Cluster Algorithm

Algorithm A.3 Classification Cluster

Require: Query q

Ensure: True if safe

```

1: procedure CHUNKEDDETECT( $m, t, q, L$ )
2:   for each chunk  $c$  do
3:     if Predict( $m, c$ ) = 1 then
4:       return True
5:     end if
6:   end for
7:   return False
8: end procedure
9:  $jb \leftarrow \text{CHUNKEDDETECT}(m_{jb}, t_{jb}, q, L)$ 
10:  $tox \leftarrow \text{CHUNKEDDETECT}(m_{tox}, t_{tox}, q, L)$ 
11: return  $\neg jb \wedge \neg tox$ 

```

A.5 Heuristic Vector Analyzer Algorithm

Algorithm A.4 Heuristic Vector Analyzer

Require: Prompt p , threshold T

Ensure: True if safe

```

1: Build keyword and pattern bit vectors  $V$ 
2: if popcount( $V$ ) <  $T$  then
3:   return True
4: else
5:   return False
6: end if

```

Algorithm A.5 Polymorphic Prompt Assembling (PPA)

Require: Separator Set $S = \{S_1, \dots, S_n\}$; System Prompt Set $T = \{T_1, \dots, T_m\}$; User Input \mathcal{I} **Ensure:** Assembled Prompt \mathcal{AP}

- 1: $(S^{\text{start}}, S^{\text{end}}) \leftarrow \text{RandomChoice}(S)$
 - 2: $\mathcal{I}_{\text{wrap}} \leftarrow S^{\text{start}} \parallel \mathcal{I} \parallel S^{\text{end}}$
 - 3: $T' \leftarrow \text{Substitute}(\text{RandomChoice}(T), S^{\text{start}}, S^{\text{end}})$
 - 4: $\mathcal{AP} \leftarrow T' \parallel \mathcal{I}_{\text{wrap}}$
 - 5: **return** \mathcal{AP}
-