

Are All Prompt Components Value-Neutral? Understanding the Heterogeneous Adversarial Robustness of Dissected Prompt in Large Language Models

Yujia Zheng^{1*}, Tianhao Li^{1*}, Haotian Huang², Tianyu Zeng³, Jingyu Lu⁴, Chuangxin Chu⁵, Yuekai Huang^{6,7,8}, Ziyou Jiang^{6,7,8}, Qian Xiong⁹, Yuyao Ge^{7,10}, Mingyang Li^{6,7,8†}

¹Duke University, ²North China University of Technology, ³Hong Kong Polytechnic University
⁴Australian National University, ⁵Nanyang Technological University, ⁶Institute of Software, Chinese Academy of Sciences ⁷University of Chinese Academy of Sciences, Beijing, China
⁸State Key Laboratory of Complex System Modeling and Simulation Technology, Beijing, China
⁹Beijing Forestry University, ¹⁰Institute of Computing Technology, Chinese Academy of Sciences

Abstract

Prompt-based adversarial attacks are a key tool for assessing the robustness of large language models (LLMs). Yet, existing studies typically treat prompts as flat text, overlooking their internal structure, different components within a prompt contribute unequally to robustness. This work introduces PROMPTANATOMY, a framework that decomposes prompts into functional components, and COMPERTURB, a controlled perturbation method that selectively modifies these components to expose component-wise vulnerabilities while ensuring linguistic plausibility via perplexity-based filtering. Using this framework, four instruction-tuning datasets are structurally annotated and validated by human reviewers. Experiments across five advanced LLMs show that COMPERTURB achieves state-of-the-art attack success rates, while ablation analyses confirm the complementary effects of prompt dissection and perplexity filtering. These results highlight the importance of structural awareness in evaluating and improving the adversarial robustness of LLMs.

1 Introduction

Large language models (LLMs) (Naveed et al., 2023; Zhao et al., 2023) such as DeepSeek, ChatGPT, and LLaMA-3 have shown impressive performance across a broad range of tasks (Thirunavukarasu et al., 2023; Zhang et al., 2025; Demszky et al., 2023; Hou et al., 2024). Their success largely relies on in-context learning, where carefully designed prompts guide model reasoning and output generation (Dong et al., 2022; Li, 2023; Long et al., 2024; Mei et al., 2025). As these models are increasingly deployed in real-world applications, concerns about their robustness under

*Equal contribution, co-first author.

Email: {yujia.zheng, tianhao.li}@duke.edu

†Corresponding author. mingyang2017@iscas.ac.cn

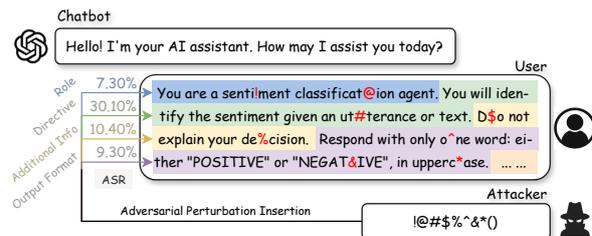


Figure 1: Inserting adversarial perturbations into different prompt components leads to distinct attack success rates (ASR).

perturbed or adversarial prompts have become central to ensuring reliability and safety (Zhu et al., 2023b,a; Hu et al., 2024; Ghosh et al., 2025).

Previous studies have shown that even minor variations in prompt phrasing, including typos, synonym substitutions, or formatting changes, can cause significant performance degradation, a phenomenon referred to as prompt sensitivity (Zhu et al., 2023b,a; Wang et al., 2023b). To evaluate this vulnerability, Zhu et al. (2024) introduced *PromptBench*, a benchmark for adversarial prompt perturbations, while Wang et al. (2023b) proposed *MTTM*, a metamorphic testing framework that applies structured modifications to assess model robustness. Although these studies reveal the fragility of LLMs under adversarial prompting, they typically treat prompts as homogeneous text, without considering the internal structure that may affect model behavior.

We argue that **prompts are fundamentally compositional, consisting of multiple functional components**, such as task directives, role assignments, input-output delimiters, exemplars, and answer formatting instructions, each guiding LLMs to interpret and execute tasks more effectively. In practical scenarios, users often construct complex prompts by layering basic instructions with supplementary cues, such as specifying the task domain, assigning a role to the model (e.g., “You are a medical expert”), providing concrete examples, or defining output styles. Recognizing the diversity and

importance of these components, Schulhoff et al. (2024) provide a comprehensive taxonomy that categorizes prompt elements into six types: *Directive*, *Examples*, *Output Formatting*, *Style Instructions*, *Role*, and *Additional Information*. The definitions and examples of prompt components are shown in Appendix A. To further streamline the taxonomy for our analysis, we merge *Style Instructions* into *Output Formatting*, as both primarily define the expected form or tone of the model’s output. These components exhibit heterogeneous adversarial robustness, meaning that some are more vulnerable to perturbations than others. To investigate this, we pose the following research questions: *RQ1*: How does prompt structure and component richness affect the effectiveness of adversarial perturbations? *RQ2*: How does the sensitivity of adversarial robustness vary across different prompt component types? *RQ3*: How do semantic and syntactic perturbations differ in effectiveness across prompt components?

To address these questions, we introduce PROMPTANATOMY, **the first framework for structurally decomposing prompts into canonical components**. Our framework achieves superior accuracy compared to GPT-4o (Hurst et al., 2024) in dissecting long, complex prompt structures. This dissection enables fine-grained analysis of how each component contributes to model robustness. We first complete missing components and then apply PROMPTANATOMY to four complex, domain-specific prompt datasets, resulting in structurally annotated variants (*PubMedQA-PA*, *EMEA-PA*, *Leetcode-PA*, *CodeGeneration-PA*). Each dataset undergoes manual review by human experts, achieving 100% inter-annotator agreement to ensure annotation quality.

Building on this foundation, we propose COMPERTURB, **a component-wise adversarial perturbation method** that perturbs one component at a time using semantically aligned perturbation strategies conditioned on the component’s function. This enables controlled robustness evaluation by isolating the effects of specific prompt modifications. We conduct extensive experiments across multiple LLMs and tasks, revealing that certain components are significantly more vulnerable to perturbation than others in specific circumstances. Our findings provide actionable insights for prompt engineers and model developers to craft safer, more reliable prompts in high-stakes applications.

Contributions. This paper makes the following key contributions: (i) We introduce PROMPTANATOMY, the first framework for structurally decomposing prompts for LLMs into multiple canonical components, enabling fine-grained analysis and systematic prompt design. (ii) We apply PROMPTANATOMY to four complex domain-specific prompt datasets, producing variants annotated with structural components (denoted with the suffix ‘-PA’). Each dataset is manually reviewed to ensure annotation quality, achieving 100% agreement among human annotators. (iii) We introduce COMPERTURB, a component-wise perturbation method that applies targeted adversarial strategies to individual prompt components, enabling fine-grained analysis of their impact on model behavior and adversarial robustness. (iv) We conduct comprehensive experiments across diverse LLMs and datasets based on COMPERTURB, demonstrating which components most significantly influence robustness. Our findings offer actionable insights for prompt engineers and model developers to design safer and more reliable prompting strategies.¹

Findings. We summarize three key findings from our comprehensive evaluation: (i) **Component and Structure-Guided Perturbations Amplify Adversarial Effectiveness in Complex Prompts.** COMPERTURB performs best on structurally rich prompts, showing that component-aware attacks are especially powerful in complex settings. (*RQ1*) (ii) **Prompt Components Exhibit Heterogeneous Robustness to Adversarial Perturbations.** Components like *DIR* and *ADI* are more susceptible to perturbations, while *ROL* and *OFT* remain relatively robust. (*RQ2*) (iii) **Semantic Perturbations Are More Effective Than Syntactic Ones Across Prompt Components.** Meaning-level changes outperform syntactic ones across components, highlighting the role of perturbation type. (*RQ3*) Based on these findings, We provide guidelines and takeaways for four potential target audience: offensive/defensive security researchers, prompt engineers, developers, and general users.

2 Related Works

Prompt Dissection Prompt dissection breaks prompts into functional components, such as directives, roles, and formatting, that shape LLM behavior (Schulhoff et al., 2024). This modular approach improves control, interpretability, and task-solving

¹Code and data are available at <https://github.com/Yujiaaaaa/PACP>.

by leveraging reusable substructures (Khot et al., 2022), and reveals latent multilingual capabilities in LLMs (Nie et al., 2024). In adversarial contexts, it enhances attack effectiveness through targeted reconstruction (Li et al., 2024c), underscoring the value of component-aware prompt analysis.

Adversarial Robustness Prompt robustness addresses the resilience of LLMs to adversarial manipulations and variations in prompt wording and structure (Wang et al., 2023a). Due to the sensitivity of LLM outputs to minor perturbations in prompts, ensuring robustness has become a critical area of research, particularly in safety-sensitive domains such as healthcare, finance, and science. Additionally, recent advancements extend robustness considerations to LLM-based agents, reflecting the increasing complexity and interactive contexts of deployed systems (Xiong et al., 2025; Hu et al., 2025). Several works propose black-box adversarial attacks that exploit prompt-level vulnerabilities (Zhu et al., 2023b; Zou et al., 2023; Xue et al., 2023; Das et al., 2024), while others investigate certified robustness or ensemble-based defenses (Aguilera-Martínez and Berzal, 2025). These studies underscore the importance of evaluating LLMs under structurally-aware perturbation paradigms beyond token-level noise.

LLM Safety Benchmark A range of safety benchmarks has been developed to evaluate LLM alignment across adversarial, ethical, and domain-specific risks (Lu et al., 2025). General benchmarks like AdvBench (Wang et al., 2021), SafetyBench (Zhang et al., 2023), and SALAD-Bench (Li et al., 2024a) assess refusal behaviors and robustness to jailbreaks. Domain-specific benchmarks, ChemSafetyBench (Zhao et al., 2024), MedSafetyBench (Han et al., 2024), and SciSafeEval (Li et al., 2024b), test models on scientific misuse, ethical compliance, and jailbreak resistance. Agent-based evaluations like Agent-SafetyBench (Zhang et al., 2024) and SafeAgentBench (Yin et al., 2024) measure risks in dynamic tool-use environments. While these benchmarks standardize vulnerability assessment, they generally overlook the compositional structure of prompts.

3 Methodology

3.1 PROMPTANATOMY: Identifying and Dissecting LLM Prompt Components

Formulation Let \mathcal{P} denote the space of free-form natural language prompts. Each prompt $p \in \mathcal{P}$ is typically an unstructured token sequence

designed to elicit LLM behavior. Existing methods treat p as a monolithic input, lacking a formal structure for semantic analysis. We define *prompt dissection* as mapping p into a structured sequence of functional components:

$$A : \mathcal{P} \rightarrow \mathcal{C}^n,$$

where \mathcal{C} is the set of canonical components and n is the number of identified components. The output $A(p) = (c_1, c_2, \dots, c_n)$ represents the decomposed prompt under our framework, PROMPTANATOMY, which enables component-aware analysis and robustness evaluation.

Challenges *Prompt dissection* introduces several practical and technical challenges. First, existing human-in-the-loop methods (Shivagunde et al., 2024) rely heavily on manual annotation, which is time-consuming and infeasible at scale. Second, LLMs often exhibit limited robustness when parsing long or complex prompts. For instance, they may omit content, oversimplify components, or focus disproportionately on the beginning and end of the prompt, neglecting critical middle sections (see examples in the Appendix B). Third, current LLMs struggle with disentangling overlapping or implicitly expressed components, leading to incorrect or incomplete dissection. This issue is amplified when the interpretation of a sentence depends on adjacent content (e.g., “Use the following format”), highlighting the need for context-aware classification. Finally, the absence of publicly available benchmark datasets tailored for prompt dissection impedes both method development and standardized evaluation.

Limited Context-Aware Prompt Dissection To address these challenges, we develop a two-stage dissection pipeline (Algorithm 1) that combines rule-based segmentation with LLM-based classification under contextual guidance. First, each raw prompt is segmented into individual sentences using punctuation delimiters (e.g., periods, question marks) to enable fine-grained analysis. For each sentence, we construct a local context window by collecting up to two preceding and two following sentences, which helps disambiguate meaning in context-dependent cases. Each target sentence and its context are then passed to a language model (e.g., GPT-4o) using a structured prompt that instructs the model to classify the sentence into one or more of five canonical components: <Role>, <Directive>, <Additional

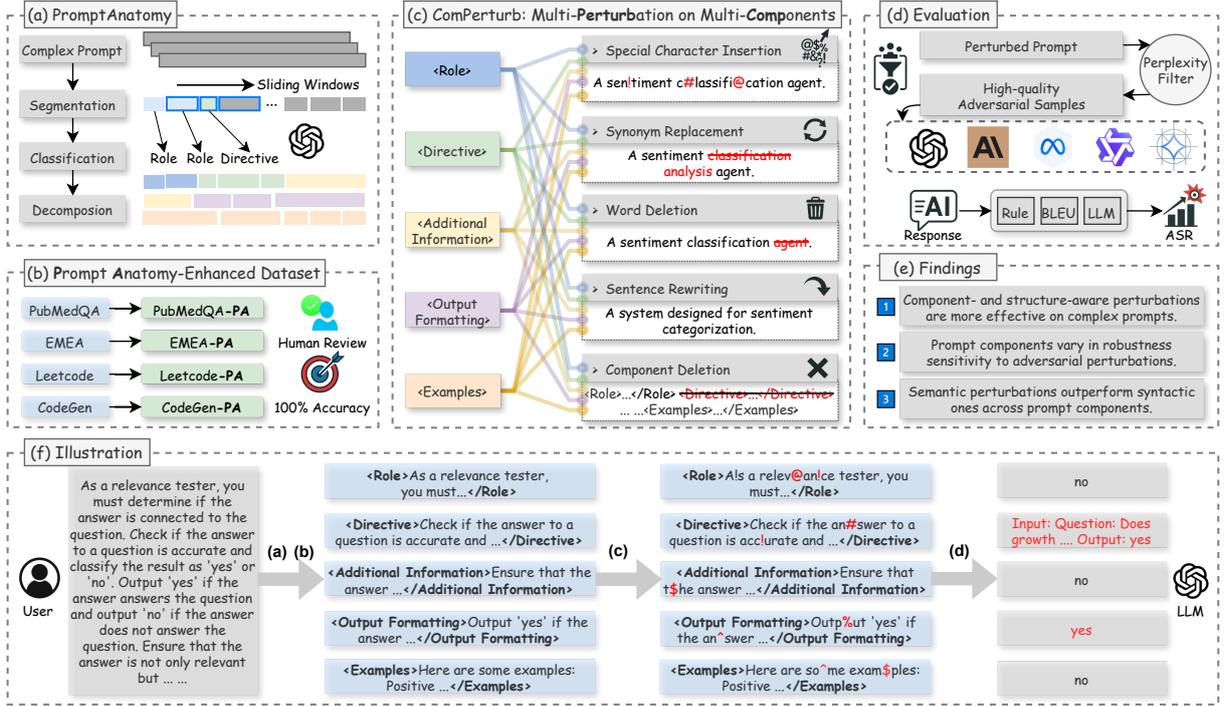


Figure 2: Overview of PROMPTANATOMY and COMPERTURB. (a) Origin prompts are dissected into components using a sliding window LLM classifier. (b) Dissection is applied to benchmarks to create annotated datasets. (c) COMPERTURB applies perturbations to components. (d) Perturbed prompts are enhanced and use to evaluate LLMs. (e) Analysis reveals component-wise vulnerabilities. (f) A working example shows dissection, perturbation, and output shifts.

Information>, <Output Formatting>, and <Examples>. The model’s responses are aggregated and post-processed to ensure consistency, with an additional verification step to identify and reclassify any missing sentences. The final output is a structured representation of the prompt, where each sentence is wrapped in XML-style tags corresponding to its identified component(s).

Algorithm 1 Prompt Dissection Pipeline

Require: Prompt p from dataset \mathcal{D}
Ensure: Component tags $\mathcal{A}(p)$

- 1: $S \leftarrow \text{Split}(p, [".", "?"])$
- 2: $\mathcal{A}(p) \leftarrow \emptyset$
- 3: **for** $i = 1$ to $|S|$ **do**
- 4: $C_i \leftarrow \text{ContextWindow}(S, i, k = 2)$
- 5: $r_i \leftarrow \text{LLMClassify}(s_i, C_i)$
- 6: $\mathcal{A}(p) \leftarrow \mathcal{A}(p) \cup \{r_i\}$
- 7: **end for**
- 8: $S_{\text{miss}} \leftarrow \text{FindUnlabeled}(S, \mathcal{A}(p))$
- 9: **for** $s_j \in S_{\text{miss}}$ **do**
- 10: $C_j \leftarrow \text{ContextWindow}(S, j, k = 2)$
- 11: $r_j \leftarrow \text{LLMClassify}(s_j, C_j)$
- 12: $\mathcal{A}(p) \leftarrow \mathcal{A}(p) \cup \{r_j\}$
- 13: **end for**
- 14: **return** $\mathcal{A}(p)$

As shown in Table 4 in the Appendix C, PROMPTANATOMY shows strong reliability, achieving 100% accuracy on GPT-4o and over 90% on other open-weight models across diverse domains. We applied it to a diverse set of existing prompt-based datasets across multiple domains, including biomedical question answering (*PubMedQA* (Wang et al., 2022)), translation (*EMEA* (Wang et al., 2022)), calculating password strength (*LeetCode* (Wang et al., 2022)), and code generation (*CodeGeneration* (nceyda, 2025)). For each dataset, we automatically dissect the raw prompts into structured components, which result in enhanced datasets *PubMedQA-PA*, *EMEA-PA*, *LeetCode-PA*, and *CodeGeneration-PA* (See Appendix D for details).

3.2 COMPERTURB: Component-wise Perturbation Method on LLM Prompts

Based on our PROMPTANATOMY framework, we propose COMPERTURB, a component-wise perturbation method. Formally, given a prompt p decomposed into k labeled components $\{c_1, \dots, c_k\}$, COMPERTURB applies a perturbation function \mathcal{T} independently to each component, generating a per-

turbed prompt:

$$p^{(i,j)} = \{c_1, \dots, \mathcal{T}_{i,j}(c_i), \dots, c_k\} \quad (1)$$

where $\mathcal{T}_{i,j}$ denotes the j -th perturbation operation applied to component c_i , while all other components remain unchanged. This design enables systematic evaluation of an LLM f 's robustness by isolating the effect of perturbations on individual components, measured via task-specific metrics $\Delta f(p, p^{(i,j)})$.

To instantiate $\mathcal{T}_{i,j}$, we adopt five commonly used perturbation operations introducing varying degrees of semantic distortion: (i) special character insertion, injecting symbols such as '@', '#', or '!' to simulate typos or obfuscation; (ii) synonym replacement, substituting words with semantically similar alternatives using `nltk`; (iii) word deletion, removing selected tokens to mimic underspecified or incomplete prompts; (iv) sentence rewriting, replacing a sentence with a paraphrased variant; and (v) component deletion, removing an entire functional segment (e.g., `<Directive>` or `<Examples>`) based on our structural decomposition. Specific examples of each perturbation are provided in Appendix E.

To select high-quality adversarial samples, we employ a perplexity-guided filtering strategy. For each perturbed prompt $p^{(i,j)}$, we compute its perplexity $\text{PPL}(p^{(i,j)})$ using a reference language model and calculate the perplexity ratio relative to the original prompt $p^{(0)}$:

$$\text{Ratio}^{(i,j)} = \frac{\text{PPL}(p^{(i,j)})}{\text{PPL}(p^{(0)})} \quad (2)$$

Higher ratios reflect greater semantic disturbance introduced by the perturbation. We then rank all perturbed samples by descending $\text{Ratio}^{(i,j)}$ and retain the top 20% as our final adversarial set, effectively filtering out low-impact or trivial perturbations. Empirical validation of this perplexity-guided filtering strategy is provided in Appendix F.

Previous studies support the use of perplexity as a proxy for adversarial impact. Alon and Kamfonas (2023) showed that adversarially appended suffixes yield significantly higher perplexity than normal prompts, while Qiu et al. (2025) observed a strong correlation between high perplexity and attack success rate. In our work, the effectiveness of this filtering strategy is further empirically validated in Appendix, where we illustrate the relationship between perplexity ratios and attack success rates across all perturbed samples.

4 Experiment

Setup. We evaluate the effectiveness of COMP-PERTURB across a diverse suite of language models and task domains. We evaluate on four domain-specific datasets: PubMedQA, EMEA, Leetcode, and CodeGeneration, each with 1,000 samples in raw and annotated forms. Detailed dataset statistics are provided in Appendix G. For target models, we select five widely-used LLMs with varying capacities and architectures: *GPT-4o* (M1), *Claude3.7* (M2), *Qwen2.5-14B* (M3), *LLaMA3.3-70B* (M4), and *Gemma3-12B* (M5). For adversarial prompt generation, we employ both commercial and open-weight models, including *GPT-4o*, *GPT-3.5*, *Qwen3-8B*, *Qwen2.5-14B*, *LLaMA3.3-70B*, and *Gemma3-12B*. Our evaluation uses datasets from distinct domains: *PubMedQA-PA*, *EMEA-PA*, *Leetcode-PA*, and *CodeGeneration-PA*. Inference was performed using 4 NVIDIA RTX™ A6000 GPUs for open-weight models, and through official APIs for commercial models. We compare our method against baselines *MTTM* (Wang et al., 2023b) and *PromptBench* (Zhu et al., 2024).

Metrics. We adopt task-specific metrics. For classification-style tasks such as PubMedQA and Leetcode, we define an attack as successful if the LLM's output no longer contains the correct answer. For generation tasks such as translation, we calculate Bilingual Evaluation Understudy (BLEU) scores (Papineni et al., 2002) to measure the semantic preservation or degradation after perturbation, an attack is considered a success if BLEU is less than 20, which indicate poor translation quality (Koehn and Knowles, 2017). For code generation tasks (evaluated on the Code dataset), we employ *GPT-4o* as an automatic reference model to assess the correctness of the generated code.

Based on various successful metrics for different tasks, we compute Attack Success Rate (ASR) as:

$$\text{ASR} = \frac{N_{\text{success}}}{N_{\text{total}}} \times 100\% \quad (3)$$

where N_{success} is the number of samples attacked successfully, N_{total} is the total number of samples.

We introduce the **Prompt Complexity Metric (PCM)**, a quantitative measure designed to evaluate the inherent complexity of LLM prompts. PCM integrates five core dimensions: lexical rarity, syntactic depth, semantic dispersion, structural rich-

Table 1: Comparison of ASR (%) of COMPERTURB and baselines across five LLMs and four datasets. Each cell reports best ASR, followed by the mean ASR and standard deviation.

Method	GPT-4o	Claude3.7	Qwen2.5-14B	LLaMA3.3-70B	Gemma3-12B	Average
PubMedQA-PA (PCM = 34.883; AVG \uparrow =17.76%)						
MTTM	42.50 (40.70 \pm 0.02)	46.88 (44.23 \pm 0.02)	44.67 (39.33 \pm 0.03)	45.74 (41.74 \pm 0.02)	47.32 (39.08 \pm 0.01)	45.42 (41.02 \pm 0.02)
PromptRobust	44.37 (41.56 \pm 0.03)	47.34 (45.34 \pm 0.07)	46.75 (41.71 \pm 0.02)	49.03 (47.03 \pm 0.05)	49.75 (46.25 \pm 0.03)	47.45 (44.38 \pm 0.04)
ComPerturb (Ours)	64.98 (53.63 \pm 0.07)	67.01 (60.02 \pm 0.04)	60.94 (50.76 \pm 0.06)	65.04 (52.38 \pm 0.07)	63.00 (51.73 \pm 0.07)	64.19 (53.70 \pm 0.06)
EMEA-PA (PCM = 32.100; AVG \uparrow =10.24%)						
MTTM	62.30 (60.01 \pm 0.02)	62.24 (60.26 \pm 0.01)	74.20 (69.97 \pm 0.02)	73.52 (72.00 \pm 0.01)	76.30 (72.89 \pm 0.02)	69.72 (67.03 \pm 0.02)
PromptRobust	69.30 (58.74 \pm 0.04)	67.50 (61.56 \pm 0.03)	72.30 (70.34 \pm 0.05)	75.70 (72.87 \pm 0.06)	74.60 (73.81 \pm 0.03)	71.88 (67.46 \pm 0.04)
ComPerturb (Ours)	75.90 (67.10 \pm 0.05)	78.00 (69.11 \pm 0.05)	84.90 (78.09 \pm 0.04)	82.12 (75.93 \pm 0.04)	84.30 (78.66 \pm 0.03)	81.04 (73.78 \pm 0.04)
Leetcode-PA (PCM = 41.495; AVG \uparrow =20.16%)						
MTTM	57.20 (53.45 \pm 0.02)	52.30 (49.10 \pm 0.01)	47.90 (47.96 \pm 0.03)	54.70 (51.89 \pm 0.03)	54.40 (50.88 \pm 0.02)	53.30 (50.66 \pm 0.02)
PromptRobust	56.30 (54.98 \pm 0.05)	55.20 (54.32 \pm 0.02)	56.70 (55.12 \pm 0.05)	51.20 (50.32 \pm 0.03)	56.90 (54.30 \pm 0.02)	55.26 (53.81 \pm 0.03)
ComPerturb (Ours)	81.50 (75.40 \pm 0.04)	72.08 (65.74 \pm 0.04)	65.30 (62.57 \pm 0.02)	78.80 (76.10 \pm 0.03)	74.50 (67.42 \pm 0.04)	74.44 (69.45 \pm 0.03)
CodeGeneration-PA (PCM = 44.524; AVG \uparrow =29.6%)						
MTTM	31.20 (27.98 \pm 0.02)	29.30 (27.77 \pm 0.01)	40.20 (38.40 \pm 0.03)	40.30 (38.17 \pm 0.03)	41.70 (38.78 \pm 0.02)	36.54 (34.22 \pm 0.02)
PromptRobust	29.08 (23.46 \pm 0.05)	29.30 (27.32 \pm 0.02)	40.20 (35.78 \pm 0.05)	41.20 (36.21 \pm 0.03)	38.40 (33.21 \pm 0.02)	35.78 (31.20 \pm 0.03)
ComPerturb (Ours)	66.96 (61.48 \pm 0.04)	64.00 (57.09 \pm 0.04)	66.36 (59.35 \pm 0.05)	65.40 (60.72 \pm 0.03)	66.08 (60.19 \pm 0.03)	65.76 (59.77 \pm 0.04)

ness, and task difficulty. Each component captures a distinct aspect of prompt formulation. The PCM score is computed as a weighted sum of these dimensions:

$$\text{PCM}(p) = \alpha C_{\text{lexical}} + \beta C_{\text{syntactic}} + \gamma C_{\text{semantic}} + \delta C_{\text{structural}} + \epsilon C_{\text{task}}, \quad (4)$$

where C_{lexical} denotes the average inverse document frequency of tokens, $C_{\text{syntactic}}$ reflects syntactic tree depth, C_{semantic} measures pairwise embedding-based semantic distance, $C_{\text{structural}}$ captures the organization of prompt components, and C_{task} encodes the inherent difficulty of the task. Detail could be found at Appendix H.

5 Results and Analysis

Table 1 provides a comparison of proposed COMPERTURB with baselines. For each method, we report both the highest ASR and the mean ASR with standard deviation, aggregated over multiple perturbation strategies. The results demonstrate that COMPERTURB achieves the highest ASRs across all settings. It achieves significant gains in complex scenarios by leveraging a component-aware design that targets semantically critical prompt segments, highlighting the advantage of structurally and semantically informed perturbations.

Finding #1: Component and Structure-Guided Perturbations Result in a Better Effect in Complex Prompts.

We derive a practical guideline to enhance LLM robustness for developers.

Guideline #1: Enhancing robustness requires clear task specification by users and component-aware adversarial data augmentation by developers.

Table 2 reports the ASRs of five LLMs across four domain-specific datasets, categorized by prompt components and perturbation types, revealing critical patterns in adversarial robustness. The DIR and ADI consistently yield higher ASRs, highlighting their centrality in conveying task-critical semantics. In contrast, the ROL shows greater resilience, likely due to its auxiliary function. Semantic perturbations, including SYR and SER, are generally more effective than syntactic ones like SCI, underscoring how meaning-level disruptions more severely compromise model behavior. Furthermore, model-specific vulnerabilities emerge, with *LLaMA3.3-70B* and *Gemma3-12B* being particularly susceptible to component-wise attacks, possibly due to architectural or training differences. Heatmaps in Figures 5 to 9 (See in the Appendix I) visualize these patterns, showing that datasets such as *EMEA-PA* and *Leetcode-PA* exhibit higher ASRs. These prompts depend heavily on semantically dense instructions and multi-step reasoning, especially within DIR and ADI components, making them more fragile under perturbation. By contrast, prompts from *PubMedQA-PA* and *CodeGeneration-PA*, which are more syntax-oriented or retrieval-based, display greater robustness, as superficial semantic alterations have less impact on task execution. These patterns underscore another observation:

Finding #2: Prompt Components Exhibit Heterogeneous Robustness to Adversarial Perturbations.

This yields an actionable insight to further enhance robustness for prompt engineers.

Table 2: ASR (%) on five LLMs across four datasets, grouped by prompt components and perturbation types.

Model	COM PPL		PubMedQA-PA					EMEA-PA					Leetcode-PA					CodeGeneration-PA				
			SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD
GPT-4o	ROL	w/o	46.0	40.0	42.5	46.0	42.5	58.5	58.5	57.5	57.5	61.0	57.4	60.0	60.5	65.4	65.0	43.5	43.5	42.5	40.5	41.5
		w/	47.0	45.5	52.5	50.0	47.4	65.0	59.0	60.0	57.5	67.5	70.5	70.5	75.0	69.0	75.0	60.0	52.5	57.4	60.0	62.5
	DIR	w/o	56.0	52.0	59.5	50.0	77.5	62.5	62.0	72.0	64.5	99.5	72.5	70.0	72.0	72.5	95.0	50.5	47.0	49.5	47.5	64.0
		w/	60.0	57.4	64.0	65.0	78.5	67.5	67.5	73.0	72.5	99.5	75.0	77.5	80.0	77.5	97.5	67.0	62.5	67.0	65.0	72.5
	ADI	w/o	46.0	43.0	47.5	45.0	52.5	59.0	61.0	59.5	62.5	59.5	62.5	67.5	71.0	70.0	77.5	44.0	43.5	45.0	45.0	45.5
		w/	55.0	45.0	57.4	43.5	75.0	62.5	63.5	68.5	67.5	68.0	70.5	68.5	80.0	72.5	87.0	67.5	60.0	62.5	65.0	62.5
	OFT	w/o	45.0	46.0	45.5	48.0	47.0	59.0	61.0	58.5	58.0	58.0	57.4	66.0	65.5	66.0	66.5	41.5	45.5	43.0	44.0	43.5
		w/	47.5	50.0	46.0	50.0	52.5	62.5	75.0	65.0	62.0	65.0	71.0	75.0	77.5	69.0	85.0	65.0	57.5	57.4	62.5	55.0
	EXA	w/o	45.5	43.5	42.0	47.5	46.0	57.5	59.5	58.0	57.5	61.5	50.0	55.0	69.0	65.0	69.5	45.5	43.5	44.0	44.0	42.0
		w/	47.5	44.0	55.0	55.0	47.0	62.5	62.5	65.0	62.0	77.5	68.0	64.0	69.0	67.0	97.5	62.5	57.4	62.5	57.4	55.0
	CRT	w/o	44.0	43.5	43.0	48.5	52.5	60.5	59.5	59.5	59.5	69.5	60.0	69.0	70.5	69.5	76.0	45.0	43.5	44.0	43.0	54.0
		w/	55.0	45.0	60.0	50.0	55.0	62.5	62.5	60.0	70.0	70.5	71.0	70.0	72.5	70.0	87.0	62.0	60.0	65.0	55.0	62.0
Claude3.7	ROL	w/o	50.0	48.0	48.0	52.0	52.5	57.5	58.0	50.0	58.5	59.0	50.4	53.0	50.0	50.5	50.5	39.0	40.0	42.0	39.0	39.5
		w/	55.0	52.5	53.5	65.5	57.5	67.5	62.5	59.0	70.0	65.0	57.4	53.5	72.5	55.0	75.0	60.0	57.4	55.0	60.0	50.0
	DIR	w/o	58.5	56.0	57.0	59.5	77.5	67.0	61.0	68.5	64.5	99.0	60.5	59.0	61.0	60.0	90.0	46.0	49.0	49.0	46.0	67.0
		w/	65.0	60.0	60.0	62.5	87.5	77.5	65.0	75.0	72.5	100	65.0	62.5	74.0	67.4	92.0	60.0	62.5	62.5	60.0	72.5
	ADI	w/o	56.5	54.0	52.0	55.0	64.0	61.5	59.0	60.0	59.5	59.5	53.0	60.5	57.4	53.5	76.5	43.0	44.0	42.0	43.0	49.0
		w/	58.0	57.5	55.0	65.0	65.0	67.5	72.0	69.5	70.0	70.0	53.5	60.0	65.0	56.0	96.5	55.0	57.0	55.5	55.0	60.0
	OFT	w/o	52.0	50.5	54.0	58.5	66.5	58.5	57.0	60.0	59.5	58.5	50.0	50.5	51.5	55.0	51.5	42.0	40.0	42.0	42.0	40.0
		w/	55.0	57.5	54.5	60.0	66.0	70.0	67.0	60.0	68.5	60.0	54.5	60.0	67.5	60.0	80.0	57.5	50.0	57.4	57.5	57.5
	EXA	w/o	52.0	52.0	56.0	52.5	57.4	60.0	59.5	59.0	58.0	59.0	53.0	54.5	55.5	56.4	52.5	40.0	45.0	43.0	40.0	39.0
		w/	55.0	60.0	57.5	57.5	58.0	65.0	65.0	67.5	68.5	72.5	57.5	55.7	56.5	57.0	90.0	55.0	57.4	55.0	55.0	47.5
	CRT	w/o	55.5	52.5	52.0	52.5	59.0	59.0	60.0	60.0	59.0	70.5	57.4	57.4	56.9	55.0	60.5	40.0	41.0	42.0	40.0	49.0
		w/	57.0	55.0	52.5	55.0	70.0	67.0	70.0	67.5	67.0	70.5	62.5	61.0	59.0	61.0	87.5	50.0	57.5	55.5	50.0	55.0
Qwen2.5-14B	ROL	w/o	39.2	44.2	42.0	46.0	38.7	71.1	74.5	72.5	75.7	73.2	51.6	50.3	50.2	50.5	51.7	44.0	43.0	42.5	39.5	43.0
		w/	43.7	43.7	45.0	48.5	47.5	74.5	76.5	76.2	77.5	75.0	49.5	53.5	62.5	64.0	70.5	60.5	57.4	57.0	57.4	60.5
	DIR	w/o	41.5	50.0	47.0	47.5	96.5	74.8	77.3	87.0	81.2	99.5	56.8	54.2	55.7	53.5	93.0	59.5	49.5	49.5	49.5	67.0
		w/	50.0	51.2	55.0	51.2	97.5	75.0	80.0	88.7	81.3	99.5	60.5	54.5	64.0	54.5	93.0	70.0	57.4	68.0	59.5	77.4
	ADI	w/o	41.4	47.5	45.2	46.0	40.7	74.7	75.2	75.0	75.5	74.7	51.9	54.2	54.6	54.3	92.1	49.5	45.0	45.0	41.0	45.5
		w/	50.0	52.0	51.5	51.2	47.5	74.8	73.8	80.0	72.5	89.0	52.0	56.4	57.9	55.6	98.0	62.5	50.0	62.8	56.0	62.5
	OFT	w/o	42.0	45.0	43.7	44.5	45.0	74.2	74.7	75.2	74.5	75.0	50.4	50.2	51.6	50.6	52.0	40.5	43.0	44.5	40.8	44.0
		w/	52.5	52.5	50.0	47.5	48.7	72.5	76.3	80.0	75.0	80.0	45.5	51.5	74.5	52.0	85.0	72.5	52.5	57.4	58.0	53.0
	EXA	w/o	41.0	45.7	45.0	45.5	37.2	71.2	76.2	75.0	73.0	74.2	50.7	50.1	52.3	51.5	54.6	43.5	42.0	44.0	40.5	43.0
		w/	38.7	47.5	50.0	48.5	47.5	75.0	80.0	75.2	74.0	70.5	50.9	51.5	62.5	56.0	88.5	57.4	55.0	60.0	50.0	50.0
	CRT	w/o	39.5	47.0	43.0	45.4	47.0	73.8	75.7	75.7	74.7	79.0	51.7	53.7	51.5	52.5	59.6	45.5	42.5	44.5	40.5	49.5
		w/	41.2	43.7	43.5	47.0	50.0	75.0	76.0	75.7	76.3	78.8	56.0	53.7	62.0	54.0	90.0	60.5	62.5	60.0	52.7	60.0
LLaMA3.3-70B	ROL	w/o	45.0	40.0	43.5	44.0	45.0	70.5	70.0	70.7	72.5	71.2	70.0	60.0	66.0	66.5	78.5	40.0	42.0	45.0	41.0	43.0
		w/	48.0	48.5	46.5	47.5	52.5	70.0	70.0	71.7	73.7	78.7	70.0	69.5	75.0	66.0	80.0	57.4	55.0	58.0	60.0	60.0
	DIR	w/o	49.0	49.5	50.0	49.5	100	76.0	74.0	83.2	75.5	99.7	70.0	73.0	72.5	68.5	93.5	51.5	51.0	50.0	52.5	68.5
		w/	57.5	52.5	58.0	57.5	100	78.7	76.0	83.5	73.7	99.7	75.0	72.0	75.0	78.5	93.5	62.5	62.5	65.0	62.0	75.0
	ADI	w/o	47.0	45.5	46.5	47.0	47.0	72.0	71.0	72.2	71.2	70.0	72.5	71.5	72.0	72.5	80.0	42.2	46.0	43.0	47.0	53.0
		w/	50.0	47.5	55.0	50.0	50.0	75.0	71.3	73.2	75.0	82.5	75.0	77.5	72.5	72.5	90.0	70.0	62.5	62.5	52.5	60.5
	OFT	w/o	48.5	45.0	43.0	46.5	46.5	72.0	71.5	71.2	70.2	71.5	68.5	71.0	71.0	64.5	71.0	42.0	42.0	45.0	41.0	44.0
		w/	48.5	47.0	49.5	50.0	50.0	71.4	77.8	72.5	75.0	75.0	72.5	75.0	87.5	65.0	80.5	65.0	57.4	60.0	62.0	57.4
	EXA	w/o	46.5	44.0	46.5	45.5	38.0	70.5	70.2	72.0	69.2	73.2	71.5	70.0	67.0	66.5	66.5	40.5	43.5	44.5	46.0	40.0
		w/	47.0	45.0	57.4	52.0	42.5	71.5	76.5	70.0	73.7	85.0	72.5	70.0	80.0	65.0	92.5	57.4	62.5	62.5	56.0	52.5
	CRT	w/o	47.5	47.5	43.4	46.0	54.0	73.7	71.0	72.0	70.7	77.8	69.5	66.5	66.5	68.0	84.0	41.0	42.5	45.0	41.5	50.0
		w/	50.0	51.5	55.0	50.0	67.5	77.5	71.0	72.5	71.3	78.8	72.5	70.0	69.0	87.5	90.0	59.0	67.5	65.0	52.5	55.0
Gemma3-12B	ROL	w/o	39.5	39.5	41.5	40.0	41.0	71.0	70.0	72.5	72.0	73.5	60.0	57.5	59.5	59.5	55.0	40.0	44.0	42.5	43.5	43.0
		w/	40.0	45.0	55.0	40.0	57.4	75.0	73.5	75.0	80.5	75.0	60.0	65.0	75.0	60.0	60.0	55.0	55.0	62.5	55.0	60.0
	DIR	w/o	60.5	63.0	48.5	60.0	56.5	74.0	77.5	83.5	79.5	100	70.5	67.0	64.5	67.0	82.2	51.0	51.0	47.0	50.0	67.5
		w/	73.0	60.0	53.0	67.5	62.5	75.0	81.5	85.0	82.5	100	72.5	72.0	75.0	70.0	83.0	65.0	65.0	62.5	60.4	77.5
	ADI	w/o	41.0	40.0	44.0	42.0	44.0	74.5	73.5	73.0	73.0	73.0	59.5	57.4	62.0	71.5	88.5	43.5	44.0	42.5	44.5	44.0
		w/	60.0	47.5	47.5	57.5	51.0	75.0	80.0	85.0	75.0	73.5	59.0	57.5	60.0	75.0	94.0	65.0	52.5	59.5	57.4	67.5
	OFT	w/o	40.5	43.0	38.5	39.5	40.0	72.0	72.0	72.0	73.0	72.										

	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD
ROL	52.2	51.8	52.1	53.0	53.4	59.3	58.3	62.2	60.9	63.9
DIR	60.4	59.7	61.3	59.9	84.7	67.6	65.0	69.4	67.1	87.9
ADI	54.8	55.2	55.5	56.0	61.8	62.9	60.6	64.0	62.2	72.5
OFT	52.8	53.6	53.7	54.0	54.6	61.9	60.5	64.1	60.7	64.8
EXA	52.7	53.2	54.4	53.7	54.8	58.5	59.8	62.8	60.1	67.8
CRT	54.2	54.7	54.7	54.7	61.4	61.2	61.2	62.6	61.6	69.8

Figure 3: Heatmap of ASR across different prompt components and perturbation types. (left: w/o PPL, right: w/ PPL)

Guideline #2: Prompt engineers should prioritize protecting high-impact components, while developers should focus defenses on semantically critical elements.

Figure 3, extracted from Table 2, shows that attack success rates (ASRs) vary significantly depending on which prompt component is targeted and by which perturbation type. The DIR component consistently exhibits the highest ASRs, particularly under COD, reaching up to 87.9%, underscoring its critical role in guiding model behavior. This heightened vulnerability stems from the semantic function of DIR, which conveys task-defining instructions; altering or removing it (e.g., via COD or SYR) disrupts task understanding. Similarly, ADI, which provides essential context or constraints, shows moderate susceptibility, especially to COD. In contrast, ROL and OFT serve more auxiliary purposes, such as setting tone or format, and thus display lower ASRs across perturbations. The EXA component yields moderate effects, suggesting limited impact in zero-shot settings. The Control group (CRT), where no semantically meaningful component is perturbed, yields the lowest ASRs, reinforcing that perturbation effectiveness is tightly coupled with a component’s functional importance. These results highlight the component-sensitive nature of LLMs, where disrupting core semantic components leads to disproportionately degraded performance, and also suggest that the type of perturbation, not just the target component, plays a key role in adversarial effectiveness:

Finding #3: Semantic Perturbations Are More Effective Than Syntactic Ones Across Prompt Components.

Building on this finding, we recommend the fol-

Table 3: Ablation study showing the effect of PROMPTANATOMY (PA) and COMPERTURB (CP) on ASR.

PA	CP	Attack Success Rate (%)					AVG ↑
		M1	M2	M3	M4	M5	
✗	✗	55.9±0.12	54.1±0.08	54.6±0.15	58.9±0.15	56.4±0.15	–
✓	✗	64.8±0.11	62.8±0.09	64.5±0.13	67.8±0.13	66.0±0.13	+9.2
✗	✓	64.0±0.09	61.8±0.06	60.8±0.13	66.7±0.11	65.0±0.11	+7.7
✓	✓	72.3±0.08	70.3±0.06	69.4±0.11	72.8±0.09	72.0±0.10	+15.4

lowing best practice for both users and developers.

Guideline #3: Users should avoid vague rephrasing, and developers should train models to withstand meaning-level changes for improved robustness.

We conduct an ablation study to evaluate the impact of PROMPTANATOMY (PA) and COMPERTURB (CP) across five LLMs, testing four configurations with and without each component. As shown in Table 3, using PA or CP individually boosts ASR by +9.2% and +7.7%, respectively, demonstrating their independent effectiveness. Combining both yields the highest gain, with a +15.4% average ASR improvement, confirming their complementary strengths in enhancing attack performance.

6 Discussion and Takeaways

Our findings show that prompt components vary in adversarial robustness, challenging the notion that prompts are structurally uniform. (i) For security researchers, focusing on semantically critical components like Directive and Additional Information enables more effective red and blue teaming. (ii) For prompt engineers and developers, adopting a component-aware design enhances robustness and control in high-stakes scenarios. Defense strategies such as input preprocessing (e.g., removing special characters, normalizing whitespace, or correcting spelling errors), prompt sanitization or rewriting, and exposure to semantically plausible perturbations during training can help mitigate vulnerabilities, particularly in critical components. Runtime monitoring and validation can further detect suspicious prompts, allowing clarification or rejection to prevent erroneous outputs. (iii) For general users, clearly specifying the Directive and relevant context improves model understanding and response quality. Overall, these insights provide actionable guidance for improving robustness and designing safer prompt interactions.

We further highlight a trade-off between ASR and computational cost. While constructing the Prompt Anatomy-Enhanced Dataset is computationally intensive, the resulting attacks remain lightweight: realistic adversaries can achieve high ASR with minimal operations, such as inserting a few special characters or minor rewrites. Our analysis identifies which components are most sensitive under such constrained settings, offering practical guidance for both attack and defense considerations.

Future work includes developing automated approaches for robust prompt design and defenses against component-specific attacks. Extending this framework to multi-turn or interactive prompts may reveal additional vulnerabilities and mitigation strategies.

7 Conclusion

We reveal that different prompt components exhibit heterogeneous adversarial vulnerabilities in large language models, challenging the prevalent assumption that prompts are uniformly robust. To facilitate fine-grained robustness analysis, we introduced PROMPTANATOMY, a novel framework that structurally dissects prompts into canonical components, and COMPERTURB, a component-wise perturbation method. Extensive experiments show that vulnerabilities vary by component, with perturbation effectiveness tied to semantics and increasing with prompt complexity. These findings highlight the importance of structure-aware adversarial evaluation and the effectiveness of component-guided design in improving prompt robustness. We advocate for structure-aware prompt engineering and training-time defenses to build safer, more robust language models.

Limitations

This study focuses on controlled settings with pre-defined prompt components and perturbation types. While this design enables clear comparisons, it may overlook variations found in open-ended prompts. The perturbations are mainly text-level and do not cover deeper semantic or reasoning changes. In addition, we test on four balanced domains, which may not fully represent large-scale real-world data. Future work could explore multi-component interactions and more diverse perturbation strategies.

Acknowledgement

This work was supported by National Natural Science Foundation of China Grant No.62402484 and the National Key Research and Development Program of China (No.2024YFF0618800).

References

- Francisco Aguilera-Martínez and Fernando Berzal. 2025. Llm security: Vulnerabilities, attacks, defenses, and countermeasures. *arXiv preprint arXiv:2505.01177*.
- Gabriel Alon and Michael Kamfonas. 2023. Detecting language model attacks with perplexity. *arXiv preprint arXiv:2308.14132*.
- Nilanjana Das, Edward Raff, and Manas Gaur. 2024. Human-interpretable adversarial prompt attack on large language models with situational context. *arXiv preprint arXiv:2407.14644*.
- Dorottya Demszky, Diyi Yang, David S Yeager, Christopher J Bryan, Margaret Clapper, Susannah Chandhok, Johannes C Eichstaedt, Cameron Hecht, Jeremy Jamieson, Meghann Johnson, and 1 others. 2023. Using large language models in psychology. *Nature Reviews Psychology*, 2(11):688–701.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, and 1 others. 2022. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*.
- Shaona Ghosh, Heather Frase, Adina Williams, Sarah Luger, Paul Röttger, Fazl Barez, Sean McGregor, Kenneth Fricklas, Mala Kumar, Kurt Bollacker, and 1 others. 2025. Ailuminat: Introducing v1.0 of the ai risk and reliability benchmark from mlcommons. *arXiv preprint arXiv:2503.05731*.
- Tessa Han, Aounon Kumar, Chirag Agarwal, and Himabindu Lakkaraju. 2024. Medsafetybench: Evaluating and improving the medical safety of large language models. *Advances in Neural Information Processing Systems*, 37:33423–33454.
- Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John Grundy, and Haoyu Wang. 2024. Large language models for software engineering: A systematic literature review. *ACM Transactions on Software Engineering and Methodology*, 33(8):1–79.
- Jinwei Hu, Yi Dong, Zhengtao Ding, and Xiaowei Huang. 2025. Enhancing robustness of llm-driven multi-agent systems through randomized smoothing. *arXiv preprint arXiv:2507.04105*.
- Zhibo Hu, Chen Wang, Yanfeng Shu, Hye-Young Paik, and Liming Zhu. 2024. Prompt perturbation in retrieval-augmented generation based large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1119–1130.

- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*.
- Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. *arXiv preprint arXiv:1706.03872*.
- Lijun Li, Bowen Dong, Ruohui Wang, Xuhao Hu, Wangmeng Zuo, Dahua Lin, Yu Qiao, and Jing Shao. 2024a. Salad-bench: A hierarchical and comprehensive safety benchmark for large language models. *arXiv preprint arXiv:2402.05044*.
- Tianhao Li, Jingyu Lu, Chuangxin Chu, Tianyu Zeng, Yujia Zheng, Mei Li, Haotian Huang, Bin Wu, Zuoxian Liu, Kai Ma, and 1 others. 2024b. Scisafeval: a comprehensive benchmark for safety alignment of large language models in scientific tasks. *arXiv preprint arXiv:2410.03769*.
- Xirui Li, Ruochen Wang, Minhao Cheng, Tianyi Zhou, and Cho-Jui Hsieh. 2024c. Drattack: Prompt decomposition and reconstruction makes powerful llm jailbreakers. *arXiv preprint arXiv:2402.16914*.
- Yinheng Li. 2023. A practical survey on zero-shot prompt design for in-context learning. *arXiv preprint arXiv:2309.13205*.
- Do Long, Yiran Zhao, Hannah Brown, Yuxi Xie, James Zhao, Nancy Chen, Kenji Kawaguchi, Michael Shieh, and Junxian He. 2024. Prompt optimization via adversarial in-context learning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7308–7327.
- Haoran Lu, Luyang Fang, Ruidong Zhang, Xinliang Li, Jiazhang Cai, Huimin Cheng, Lin Tang, Ziyu Liu, Zeliang Sun, Tao Wang, and 1 others. 2025. Alignment and safety in large language models: Safety mechanisms, training paradigms, and emerging challenges. *arXiv preprint arXiv:2507.19672*.
- Lingrui Mei, Jiayu Yao, Yuyao Ge, Yiwei Wang, Baolong Bi, Yujun Cai, Jiazhi Liu, Mingyu Li, Zhong-Zhi Li, Duzhen Zhang, and 1 others. 2025. A survey of context engineering for large language models. *arXiv preprint arXiv:2507.13334*.
- Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. 2023. A comprehensive overview of large language models. *ACM Transactions on Intelligent Systems and Technology*.
- nceyda. 2025. Yap470 code generation dataset. https://huggingface.co/datasets/nceyda/YAP470_Code_Generation_Dataset. Accessed: 2025-07-05.
- Ercong Nie, Shuzhou Yuan, Bolei Ma, Helmut Schmid, Michael Färber, Frauke Kreuter, and Hinrich Schütze. 2024. Decomposed prompting: Unveiling multilingual linguistic structure knowledge in english-centric large language models. *arXiv preprint arXiv:2402.18397*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Shilin Qiu, Qihe Liu, Shijie Zhou, Min Gou, Yi Zeng, Zhun Zhang, and Zhewei Wu. 2025. Hard label adversarial attack with high query efficiency against nlp models. *Scientific Reports*, 15(1):9378.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, and 1 others. 2024. The prompt report: a systematic survey of prompt engineering techniques. *arXiv preprint arXiv:2406.06608*.
- Namrata Shivagunde, Vladislav Lialin, Sherin Muckatira, and Anna Rumshisky. 2024. Deconstructing in-context learning: Understanding prompts via corruption. *arXiv preprint arXiv:2404.02054*.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930–1940.
- Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. 2021. Adversarial glue: A multi-task benchmark for robustness evaluation of language models. *arXiv preprint arXiv:2111.02840*.
- Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, and 1 others. 2023a. On the robustness of chatgpt: An adversarial and out-of-distribution perspective. *arXiv preprint arXiv:2302.12095*.
- Wenxuan Wang, Jen-tse Huang, Weibin Wu, Jianping Zhang, Yizhan Huang, Shuqing Li, Pinjia He, and Michael R Lyu. 2023b. Mttm: Metamorphic testing for textual content moderation software. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2387–2399. IEEE.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, and 1 others. 2022.

- Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109.
- Qian Xiong, Yuekai Huang, Ziyong Jiang, Zhiyuan Chang, Yujia Zheng, Tianhao Li, and Mingyang Li. 2025. Butterfly effects in toolchains: A comprehensive analysis of failed parameter filling in llm tool-agent systems. *arXiv preprint arXiv:2507.15296*.
- Jiaqi Xue, Mengxin Zheng, Ting Hua, Yilin Shen, Yepeng Liu, Ladislau Bölöni, and Qian Lou. 2023. Trojllm: A black-box trojan prompt attack on large language models. *Advances in Neural Information Processing Systems*, 36:65665–65677.
- Sheng Yin, Xianghe Pang, Yuanzhuo Ding, Menglan Chen, Yutong Bi, Yichen Xiong, Wenhao Huang, Zhen Xiang, Jing Shao, and Siheng Chen. 2024. Safeagentbench: A benchmark for safe task planning of embodied llm agents. *arXiv preprint arXiv:2412.13178*.
- Qiang Zhang, Keyan Ding, Tianwen Lv, Xinda Wang, Qingyu Yin, Yiwen Zhang, Jing Yu, Yuhao Wang, Xiaotong Li, Zhuoyi Xiang, and 1 others. 2025. Scientific large language models: A survey on biological & chemical domains. *ACM Computing Surveys*, 57(6):1–38.
- Zhexin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. 2024. Agent-safetybench: Evaluating the safety of llm agents. *arXiv preprint arXiv:2412.14470*.
- Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. 2023. Safetybench: Evaluating the safety of large language models with multiple choice questions. *CoRR*.
- Haochen Zhao, Xiangru Tang, Ziran Yang, Xiao Han, Xuanzhi Feng, Yueqing Fan, Senhao Cheng, Di Jin, Yilun Zhao, Arman Cohan, and 1 others. 2024. Chemsafetybench: benchmarking llm safety on chemistry domain. *arXiv preprint arXiv:2411.16736*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2).
- Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Yue Zhang, Neil Gong, and 1 others. 2023a. Promptrobust: Towards evaluating the robustness of large language models on adversarial prompts. In *Proceedings of the 1st ACM workshop on large AI systems and models with privacy and safety analysis*, pages 57–68.
- Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Yue Zhang, Neil Zhenqiang Gong, and 1 others. 2023b. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. *arXiv e-prints*, pages arXiv–2306.
- Kaijie Zhu, Qinlin Zhao, Hao Chen, Jindong Wang, and Xing Xie. 2024. Promptbench: A unified library for evaluation of large language models. *Journal of Machine Learning Research*, 25(254):1–22.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

A Definition and Examples of Prompt Components

We categorize prompts into five main components, each defined with examples for clarity (Schulhoff et al., 2024):

- **Role (ROL):** Indicates the persona or identity the model should adopt to produce stylistically or contextually appropriate output.
Example: “Pretend you are a shepherd and write a limerick about llamas.”
- **Directive (DIR):** The core instruction or question representing the intent of the prompt. Directives can be explicit or implicit.
Example (explicit): “Tell me five good books to read.”
Example (implicit): One-shot English-Spanish translation prompt: (Night: Noche)
- **Additional Information (ADI):** Extra context or details necessary to perform the task correctly, often including user-specific information or situational context.
Example: “Include your name and position when drafting an email in response to a client query.”
- **Output Formatting (OFT):** Specifies how the model should structure its output, such as CSV, Markdown, or XML.
Example: “PARAGRAPH Summarize this into a CSV.”
- **Examples (EXM):** Also called exemplars or shots, they demonstrate the task for the model.
Example: Providing a sample input-output pair in a one-shot prompt for text classification or translation.

B Examples of LLM Fails to Dissect Prompt

B.1 Misdissection of Semantic Components by GPT-4o

Prompt to GPT-4o

You are now an expert in prompt engineering for large language models. Please classify the following prompt into the five components: Role, Directive, Additional Information, Output Formatting and Examples. Please classify the following prompt into its respective components using implicit tags such as <Role></Role>, <Directive></Directive>, <Additional Information></Additional Information>, <Output Formatting></Output Formatting> and <Examples></Examples>. Do not provide any explanations or additional content. Each sentence should remain unmodified. The classified sentences should remain the same as the original sentences. Below are the definitions of each component: Role: Also known as a persona, this component helps improve the writing style and tone of the text. Examples: As a helpful assistant, answer the following question concisely. Directive: The main instruction or question in the prompt. i.e. What are the benefits of renewable energy? Or Tell me five books. Additional Information: Provides additional background information or situational details to help the model better understand the task. Example: You are writing an email to a colleague about a meeting next week. Output Formatting: Defines how the response should be structured or presented. Example: Answer the question in 'yes' or 'no'. Examples: Provides sample inputs and outputs to guide the model's response. Example: Input: 'Night' Output: 'Noche' (Spanish). The prompt to classify is: A seasoned Python developer like you has extensive experience in writing and refining Python code, and your expertise includes key domains such as software engineering, data processing, and algorithm optimization. Your task is to either create new Python scripts that fulfill specific functional requirements or modify existing Python code to enhance its performance, readability, or functionality. These tasks involve implementing basic functions, optimizing existing code, and ensuring that the code adheres to best practices in software development. The output should be a well-structured Python script. It should include clear comments explaining the purpose of each section of the code. The script should be formatted according to PEP 8 guidelines, ensuring readability and consistency. Below is an instruction that describes a task, Write a response that appropriately completes the request. Inspect the following Python file from murano/tests/unit/api/middleware/test_version_negotiation.py. Update the given code considering this change: "Change assertTrue(isinstance()) by optimal assert" `python # Copyright 2016 AT&T Corp # All Rights Reserved.# # Licensed under the Apache License, Version 2.0 (the "License"); you may # not use this file except in compliance with the License.You may obtain # a copy of the License at # # http://www.apache.org/licenses/LICENSE-2.0 # # Unless required by applicable law or agreed to in writing, software # distributed under the License is distributed on an "AS IS" BASIS, WITHOUT # WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the # License for the specific language governing permissions and limitations # under the License.import webob from murano.api import versions from murano.api.middleware import version_negotiation from murano.tests.unit import base class MiddlewareVersionNegotiationTest(base.MuranoTestCase): def test_middleware_version_negotiation_default(self): middleware_vn = version_negotiation.VersionNegotiationFilter(None) request = webob.Request.blank('/environments') result = middleware_vn.process_request(request) self.assertTrue(isinstance(result, versions.Controller))` Example : Input: Generate a program in Python that takes in a string and returns the reverse of the string."Hello World" Output: `def reverse_string(str): """ Reverses a given string """ rev_str = "" for c in str: rev_str = c + rev_str return rev_str str = "Hello World" rev_str = reverse_string(str) print(rev_str)`

Ground Truth

<Role> A seasoned Python developer like you has extensive experience in writing and refining

Python code, and your expertise includes key domains such as software engineering, data processing, and algorithm optimization.</Role>

<Directive> Your task is to either create new Python scripts that fulfill specific functional requirements or modify existing Python code to enhance its performance, readability, or functionality. Below is an instruction that describes a task, Write a response that appropriately completes the request. Inspect the following Python file from murano/tests/unit/api/middleware/test_version_negotiation.py. Update the given code considering this change: "Change assertTrue(isinstance()) by optimal assert" “python # Copyright 2016 AT&T Corp # All Rights Reserved.## Licensed under the Apache License, Version 2.0 (the "License"); you may # not use this file except in compliance with the License. You may obtain # a copy of the License at # # http://www.apache.org/licenses/LICENSE-2.0 # # Unless required by applicable law or agreed to in writing, software # distributed under the License is distributed on an "AS IS" BASIS, WITHOUT # WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the # License for the specific language governing permissions and limitations # under the License. import webob from murano.api import versions from murano.api.middleware import version_negotiation from murano.tests.unit import base class MiddlewareVersionNegotiationTest(base.MuranoTestCase): def test_middleware_version_negotiation_default(self): middleware_vn = version_negotiation.VersionNegotiationFilter(None) request = webob.Request.blank('/environments') result = middleware_vn.process_request(request) self.assertTrue(isinstance(result, versions.Controller)) “. </Directive>

<Additional Information> These tasks involve implementing basic functions, optimizing existing code, and ensuring that the code adheres to best practices in software development.</Additional Information>

<Output Formatting> The output should be a well-structured Python script. It should include clear comments explaining the purpose of each section of the code. The script should be formatted according to PEP 8 guidelines, ensuring readability and consistency.</Output Formatting>

<Examples>Example : Input: Generate a program in Python that takes in a string and returns the reverse of the string."Hello World" Output: def reverse_string(str): """ Reverses a given string """ rev_str = "" for c in str: rev_str = c + rev_str return rev_str str = "Hello World" rev_str = reverse_string(str) print(rev_str).</Examples>

Response From GPT-4o

<Role> A seasoned Python developer like you has extensive experience in writing and refining Python code, and your expertise includes key domains such as software engineering, data processing, and algorithm optimization. </Role>

<Directive>Your task is to either create new Python scripts that fulfill specific functional requirements or modify existing Python code to enhance its performance, readability, or functionality. </Directive>

<Additional Information>These tasks involve implementing basic functions, optimizing existing code, and ensuring that the code adheres to best practices in software development. </Additional Information>

<Output Formatting>The output should be a well-structured Python script. It should include clear comments explaining the purpose of each section of the code. The script should be formatted according to PEP 8 guidelines, ensuring readability and consistency. </Output Formatting>

<Additional Information> Below is an instruction that describes a task, Write a response that appropriately completes the request. </Additional Information>

<Examples> Inspect the following Python file from murano/tests/unit/api/middleware/test_version_negotiation.py. Update the given code considering this change: "Change assertTrue(isinstance()) by optimal assert" “python # Copyright 2016 AT&T Corp # All Rights Reserved.## Licensed under the Apache License, Version 2.0 (the "License"); you may # not

```

use this file except in compliance with the License.You may obtain # a copy of the License
at # # http://www.apache.org/licenses/LICENSE-2.0 # # Unless required by applicable law
or agreed to in writing, software # distributed under the License is distributed on an "AS IS"
BASIS, WITHOUT # WARRANTIES OR CONDITIONS OF ANY KIND, either express or
implied.See the # License for the specific language governing permissions and limitations #
under the License.import webob from murano.api import versions from murano.api.middleware
import version_negotiation from murano.tests.unit import base class MiddlewareVersionNegotiationTest(base.MuranoTestCase):
def test_middleware_version_negotiation_default(self):
middleware_vn = version_negotiation.VersionNegotiationFilter(None) request = webob.Request.blank('/environments')
result = middleware_vn.process_request(request)
self.assertTrue(isinstance(result, versions.Controller))
Example : Input: Generate a program in Python that takes in a string and returns the reverse of the string."Hello World"
Output:
def reverse_string(str):
""" Reverses a given string """
rev_str = ""
for c in str:
rev_str = c + rev_str
return rev_str
str = "Hello World"
rev_str = reverse_string(str)
print(rev_str)
</Examples>

```

B.2 Omission of Middle Sentences in Prompt Dissection by GPT-4o

Prompt to GPT-4o

You are now an expert in prompt engineering for large language models. Please classify the following prompt into the five components: Role, Directive, Additional Information, Output Formatting and Examples. Please classify the following prompt into its respective components using implicit tags such as <Role></Role>, <Directive></Directive>, <Additional Information></Additional Information>, <Output Formatting></Output Formatting> and <Examples></Examples>. Do not provide any explanations or additional content. Each sentence should remain unmodified. The classified sentences should remain the same as the original sentences. Below are the definitions of each component: Role: Also known as a persona, this component helps improve the writing style and tone of the text. Examples: As a helpful assistant, answer the following question concisely. Directive: The main instruction or question in the prompt. i.e. What are the benefits of renewable energy? Or Tell me five books. Additional Information: Provides additional background information or situational details to help the model better understand the task. Example: You are writing an email to a colleague about a meeting next week. Output Formatting: Defines how the response should be structured or presented. Example: Answer the question in 'yes' or 'no'. Examples: Provides sample inputs and outputs to guide the model's response. Example: Input: 'Night' Output: 'Noche' (Spanish). The prompt to classify is: As a cybersecurity expert dedicated to password strength analysis and improvement, you are responsible for assessing and adjusting passwords to meet stringent security criteria. You are given a password and you need to generate the number of steps required to convert the given password to a strong password. A password is considered strong if (a) it has at least 6 characters and at most 20 characters; (b) it contains at least one lowercase letter and one uppercase letter, and at least one digit; (c) it does not contain three repeating characters in a row. In one step you can: (1) Insert one character to password, (2) delete one character from password, or (3) replace one character of password with another character. Calculate the minimum number of operations needed to transform the provided password into a strong one. Make sure your approach is efficient and explores all potential scenarios for optimal results. Answer only returns the number of steps required. Now complete the question: password = KMxgCvxC4A1SxLS2r9OEqGx3MqRaQXX. Here are some examples: Positive Example 1 - Input: password = a Output: 5. Positive Example 2 - Input: password = aA1 Output: 3. Negative Example 1 - Input: password = 1337C0d3 Output: 1.

Ground Truth

<Role>As a cybersecurity expert dedicated to password strength analysis and improvement, you

are responsible for assessing and adjusting passwords to meet stringent security criteria.</Role>
<Directive>You are given a password and you need to generate the number of steps required to convert the given password to a strong password. Calculate the minimum number of operations needed to transform the provided password into a strong one. Now complete the question: password = KMxgCvxC4A1SxLS2r9OEqGx3MqRaQXX.</Directive>

<Additional Information>A password is considered strong if (a) it has at least 6 characters and at most 20 characters; (b) it contains at least one lowercase letter and one uppercase letter, and at least one digit; (c) it does not contain three repeating characters in a row. In one step you can: (1) Insert one character to password, (2) delete one character from password, or (3) replace one character of password with another character. Make sure your approach is efficient and explores all potential scenarios for optimal results.</Additional Information>

<Output Formatting>Answer only returns the number of steps required.</Output Formatting>

<Examples>Here are some examples: Positive Example 1 - Input: password = a Output: 5. Positive Example 2 - Input: password = aA1 Output: 3. Negative Example 1 - Input: password = 1337C0d3 Output: 1.</Examples>

Response From GPT-4o

<Role> As a cybersecurity expert dedicated to password strength analysis and improvement, you are responsible for assessing and adjusting passwords to meet stringent security criteria. </Role>

<Directive> You are given a password and you need to generate the number of steps required to convert the given password to a strong password. </Directive>

<Additional Information> A password is considered strong if (a) it has at least 6 characters and at most 20 characters; (b) it contains at least one lowercase letter and one uppercase letter, and at least one digit; (c) it does not contain three repeating characters in a row. In one step you can: (1) Insert one character to password, (2) delete one character from password, or (3) replace one character of password with another character. Make sure your approach is efficient and explores all potential scenarios for optimal results. </Additional Information>

<Output Formatting> Answer only returns the number of steps required. </Output Formatting>

<Examples> Positive Example 1 - Input: password = a Output: 5. Positive Example 2 - Input: password = aA1 Output: 3. Negative Example 1 - Input: password = 1337C0d3 Output: 1. </Examples>

Missing Sentences:

Calculate the minimum number of operations needed to transform the provided password into a strong one.

Now complete the question: password = KMxgCvxC4A1SxLS2r9OEqGx3MqRaQXX.

C Evaluation of PROMPTANATOMY’S Component Classification Accuracy

In this appendix, we present a comprehensive evaluation to assess the effectiveness and reliability of the proposed PromptAnatomy (PA) framework in accurately identifying and classifying semantic components within free-form prompts. This analysis serves to validate whether the PA framework can correctly dissect prompts into meaningful substructures under our defined taxonomy.

C.1 Experimental Setup

To rigorously evaluate the classification accuracy of PA, we curated a human-annotated dataset across four task domains, PubMedQA, EMEA, Leetcode, and CodeGeneration, spanning both knowledge-intensive and procedural prompts. Each prompt was manually segmented and labeled on a per-sentence basis by expert annotators according to the canonical component schema: <Role>, <Directive>, <Additional Information>, <Output Formatting>, and <Examples>. The annotated data served as ground truth for benchmarking model predictions.

We employed the PA dissection pipeline (Algorithm 1) using six representative LLMs: (1) ChatGPT-4o, (2) ChatGPT-3.5, (3) Qwen2.5-14B, (4) Qwen3-8B, (5) LLaMA3.3-70B, and (6) Gemma3-12B.

C.2 Metric: Classification Accuracy

We define the classification accuracy as:

$$\text{Accuracy} = \frac{C_p}{N_p} \quad (5)$$

where C_p denotes the number of prompts in which all sentences are correctly classified (i.e., the predicted label set exactly matches the annotated labels for every sentence in the prompt), and N_p is the total number of prompts in the evaluation set.

A prompt is considered correct only if the model’s prediction for each sentence in the prompt exactly matches the human-annotated label set. Partial or sentence-level mismatches result in the entire prompt being marked as incorrect.

C.3 Human Review Setup

To validate the quality of the human-annotated ground truth used for benchmarking PROMPTANATOMY, we employed a structured review pro-

cess. Four annotators with NLP and prompt engineering experience independently labeled prompts according to the canonical component schema. Annotation guidelines included precise definitions, illustrative examples, and instructions for handling ambiguous sentences.

Each annotator spent roughly 2–3 hours per dataset. Disagreements were resolved via consensus discussion, with a fifth senior annotator providing final adjudication if necessary. This procedure ensures the reliability and reproducibility of the labeled datasets, which serve as the ground truth for evaluating PROMPTANATOMY’s classification accuracy across diverse domains.

C.4 Results and Analysis

The results, summarized in Table 4, demonstrate strong classification performance across all models. ChatGPT-4o achieves perfect accuracy (100%) on all four datasets, underscoring its exceptional semantic parsing capabilities. ChatGPT-3.5, Qwen2.5-14B, and LLaMA3.3-70B follow closely with average accuracies of approximately 93%, while Qwen3-8B and Gemma3-12B achieve slightly lower scores of 92% and 88%, respectively.

Table 4: Classification Accuracy (%) across different datasets and models

Model	PubMedQA	EMEA	Leetcode	CodeGeneration	AVG
ChatGPT-4o	100.00	100.00	100.00	100.00	100.00
ChatGPT-3.5	89.00	95.00	96.00	93.00	93.00
Qwen3-8B	91.00	94.00	93.00	89.00	92.00
Qwen2.5-14B	88.00	94.00	95.00	93.00	93.00
LLaMA3.3-70B	91.00	97.00	95.50	90.30	93.45
Gemma3-12B	86.00	89.00	94.50	82.50	88.00

These results confirm that the PA framework generalizes well across model scales and architectures, achieving consistent accuracy above 85% even in smaller models. The use of local contextual windows appears to significantly improve the model’s ability to disambiguate component roles, especially for structurally ambiguous prompts.

Notably, all models perform robustly across datasets with differing domain styles, indicating that prompt structure is learnable and recognizable irrespective of topical content.

D Examples of Origin and Enhanced Datasets

Table 5: Dataset Examples from PubMedQA, EMEA, Leetcode, and CodeGeneration

Dataset	Example
PubMedQA	<p>You are a response inspector, ensuring that the answer is a valid reply to the question. Check if the answer to a question is correct and classify the result as 'yes' or 'no'. Output 'yes' if the answer answers the question and output 'no' if the answer does not answer the question. Focus on determining if the answer is both relevant and complete in addressing the question. Now complete the question: Question: Is chk1 required for the metaphase-anaphase transition via regulating the expression and localization of Cdc20 and Mad2? Answer: These results strongly suggest that Chk1 is required for the metaphase-anaphase transition via regulating the subcellular localization and the expression of Cdc20 and Mad2. Here are some examples: Positive Example 1 - Input: Question: Are group 2 innate lymphoid cells (ILC2s) increased in chronic rhinosinusitis with nasal polyps or eosinophilia? Answer: As ILC2s are elevated in patients with CRSwNP, they may drive nasal polyp formation in CRS. ILC2s are also linked with high tissue and blood eosinophilia and have a potential role in the activation and survival of eosinophils during the Th2 immune response. The association of innate lymphoid cells in CRS provides insights into its pathogenesis. Output: yes Positive Example 2 - Input: Question: Does vagus nerve contribute to the development of steatohepatitis and obesity in phosphatidylethanolamine N-methyltransferase deficient mice? Answer: Neuronal signals via the hepatic vagus nerve contribute to the development of steatohepatitis and protection against obesity in HFD fed Pemt(-/-) mice. Output: yes Negative Example 1 - Input: Question: Is methylation of the FGFR2 gene associated with high birth weight centile in humans? Answer: We identified a novel biologically plausible candidate (FGFR2) for with BWC that merits further study. Output: no Negative Example 2 - Input: Question: Do tumor-infiltrating immune cell profiles and their change after neoadjuvant chemotherapy predict response and prognosis of breast cancer? Answer: Breast cancer immune cell subpopulation profiles, determined by immunohistochemistry-based computerized analysis, identify groups of patients characterized by high response (in the pre-treatment setting) and poor prognosis (in the post-treatment setting). Further understanding of the mechanisms underlying the distribution of immune cells and their changes after chemotherapy may contribute to the development of new immune-targeted therapies for breast cancer. Output: no.</p>

Dataset	Example
EMEA	<p>As an experienced Latin translation specialist, you have a deep knowledge of historical and religious contexts, which enables you to translate Latin texts into English with precision. Your skill lies in grasping the subtleties of classical languages and their cultural importance. Translate the provided Latin text into English, ensuring that the historical and religious nuances are preserved and accurately conveyed. The content you've provided is extensive, with a wide range of texts and references from various fields like religion, history, philosophy, and literature. The translation should be presented in a clear and readable English format, with any necessary annotations or explanations provided to clarify the context or meaning of specific terms or phrases. Now complete the question: ceterum cum et magni pretii et varii generis a legatis eius tam virorum quam feminarum apta usui munera circa domos ferrentur, nulla cuiquam dono ianua patuit, Tarentinaeque petulantiae animosus magis quam efficax defensor haud scio maiore cum gloria huius urbis moribus <an moenibus> repulsus sit. Here are some examples: Latin Text: Omnia mutantur, nihil interit. Translation: Everything changes, nothing perishes. Latin Text: Qui cum statuisset, nisi me per vos recuperasset, eamdem subire fortunam atque idem sibi domicilium et vitae et mortis deponere, tamen numquam nec magnitudinem negotii nec solitudinem suam nec vim inimicorum ac tela pertimuit. Translation: He had made up his mind that, should he fail, through you, to win me back to himself, he would ask permission to meet the same fate and to share the same dwelling with me in life and in death; and yet, in spite of this, no toil however formidable, no loneliness, no threat nor weapons of foes, could daunt him.</p>
Leetcode	<p>As a cybersecurity expert dedicated to password strength analysis and improvement, you are responsible for assessing and adjusting passwords to meet stringent security criteria. You are given a password and you need to generate the number of steps required to convert the given password to a strong password. A password is considered strong if (a) it has at least 6 characters and at most 20 characters; (b) it contains at least one lowercase letter and one uppercase letter, and at least one digit; (c) it does not contain three repeating characters in a row. In one step you can: (1) Insert one character to password, (2) delete one character from password, or (3) replace one character of password with another character. Calculate the minimum number of operations needed to transform the provided password into a strong one. Make sure your approach is efficient and explores all potential scenarios for optimal results. Answer only returns the number of steps required. Now complete the question: password = KMxgCvxC4A1SxLS2r9OEqGx3MqRaqXX. Here are some examples: Positive Example 1 - Input: password = a Output: 5. Positive Example 2 - Input: password = aA1 Output: 3. Negative Example 1 - Input: password = 1337C0d3 Output: 1.</p>

Dataset	Example
CodeGeneration	<p>A seasoned Python developer like you has extensive experience in writing and refining Python code, and your expertise includes key domains such as software engineering, data processing, and algorithm optimization. Your task is to either create new Python scripts that fulfill specific functional requirements or modify existing Python code to enhance its performance, readability, or functionality. These tasks involve implementing basic functions, optimizing existing code, and ensuring that the code adheres to best practices in software development. The output should be a well-structured Python script. It should include clear comments explaining the purpose of each section of the code. The script should be formatted according to PEP 8 guidelines, ensuring readability and consistency. Below is an instruction that describes a task, Write a response that appropriately completes the request. Inspect the following Python file from <code>murano/tests/unit/api/middleware/test_version_negotiation.py</code>. Update the given code considering this change: "Change <code>assert True(isinstance())</code> by optimal <code>assert</code>"</p> <pre> python # Copyright 2016 AT&T Corp # All Rights Reserved.# # Licensed under the Apache License, Version 2.0 (the "License"); you may # not use this file except in compliance with the License.You may obtain # a copy of the License at # # http://www.apache.org/licenses/LICENSE-2.0 # # Unless required by applicable law or agreed to in writing, software # distributed under the License is distributed on an "AS IS" BASIS, WITHOUT # WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the # License for the specific language governing permissions and limitations # under the License.import webob from murano.api import versions from murano.api.middleware import version_negotiation from murano.tests.unit im- port base class MiddlewareVersionNegotiationTest(base.MuranoTestCase): def test_middleware_version_negotiation_default(self): middleware_vn = version_negotiation.VersionNegotiationFilter(None) re- quest = webob.Request.blank('/environments') result = middle- ware_vn.process_request(request) self.assertTrue(isinstance(result, vers- ions.Controller)) Example : Input: Generate a program in Python that takes in a string and returns the reverse of the string."Hello World" Output: def reverse_string(str): """ Reverses a given string """ rev_str = "" for c in str: rev_str = c + rev_str return rev_str str = "Hello World" rev_str = reverse_string(str) print(rev_str). </pre>

Table 6: Dataset Examples from PubMedQA-PA, EMEA-PA, Leetcode-PA, and CodeGeneration-PA

Dataset	Example
PubMedQA-PA	<p><Role>You are a response inspector, ensuring that the answer is a valid reply to the question.</Role> <Directive>Check if the answer to a question is correct and classify the result as 'yes' or 'no'. Now complete the question: Question: Is chk1 required for the metaphase-anaphase transition via regulating the expression and localization of Cdc20 and Mad2? Answer: These results strongly suggest that Chk1 is required for the metaphase-anaphase transition via regulating the subcellular localization and the expression of Cdc20 and Mad2. </Directive></p> <p><Additional Information>Focus on determining if the answer is both relevant and complete in addressing the question.</Additional Information> <Output Formatting>Output 'yes' if the answer answers the question and output 'no' if the answer does not answer the question. </Output Formatting> <Examples>Here are some examples: Positive Example 1 - Input: Question: Are group 2 innate lymphoid cells (ILC2s) increased in chronic rhinosinusitis with nasal polyps or eosinophilia? Answer: As ILC2s are elevated in patients with CRSwNP, they may drive nasal polyp formation in CRS. ILC2s are also linked with high tissue and blood eosinophilia and have a potential role in the activation and survival of eosinophils during the Th2 immune response. The association of innate lymphoid cells in CRS provides insights into its pathogenesis. Output: yes Positive Example 2 - Input: Question: Does vagus nerve contribute to the development of steatohepatitis and obesity in phosphatidylethanolamine N-methyltransferase deficient mice? Answer: Neuronal signals via the hepatic vagus nerve contribute to the development of steatohepatitis and protection against obesity in HFD fed Pemt(-/-) mice. Output: yes Negative Example 1 - Input: Question: Is methylation of the FGFR2 gene associated with high birth weight centile in humans? Answer: We identified a novel biologically plausible candidate (FGFR2) for with BWC that merits further study. Output: no Negative Example 2 - Input: Question: Do tumor-infiltrating immune cell profiles and their change after neoadjuvant chemotherapy predict response and prognosis of breast cancer? Answer: Breast cancer immune cell subpopulation profiles, determined by immunohistochemistry-based computerized analysis, identify groups of patients characterized by high response (in the pre-treatment setting) and poor prognosis (in the post-treatment setting). Further understanding of the mechanisms underlying the distribution of immune cells and their changes after chemotherapy may contribute to the development of new immune-targeted therapies for breast cancer. Output: no.</Examples></p>

Dataset	Example
EMEA-PA	<p><Role>As an experienced Latin translation specialist, you have a deep knowledge of historical and religious contexts, which enables you to translate Latin texts into English with precision.</Role> <Directive>Translate the provided Latin text into English, ensuring that the historical and religious nuances are preserved and accurately conveyed. Now complete the question: ceterum cum et magni pretii et varii generis a legatis eius tam virorum quam feminarum apta usui munera circa domos ferrentur, nulla cuiquam dono ianua patuit, Tarentinaeque petulantiae animosus magis quam efficax defensor haud scio maiore cum gloria huius urbis moribus <an moenibus> repulsus sit.</Directive> <Additional Information>Your skill lies in grasping the subtleties of classical languages and their cultural importance. The content you've provided is extensive, with a wide range of texts and references from various fields like religion, history, philosophy, and literature.</Additional Information> <Output Formatting>The translation should be presented in a clear and readable English format, with any necessary annotations or explanations provided to clarify the context or meaning of specific terms or phrases.</Output Formatting> <Examples>Here are some examples: Latin Text: Omnia mutantur, nihil interit. Translation: Everything changes, nothing perishes. Latin Text: Qui cum statuisset, nisi me per vos recuperasset, eamdem subire fortunam atque idem sibi domicilium et vitae et mortis deponere, tamen numquam nec magnitudinem negotii nec solitudinem suam nec vim inimicorum ac tela pertimuit. Translation: He had made up his mind that, should he fail, through you, to win me back to himself, he would ask permission to meet the same fate and to share the same dwelling with me in life and in death; and yet, in spite of this, no toil however formidable, no loneliness, no threat nor weapons of foes, could daunt him.</Examples></p>
Leetcode-PA	<p><Role>As a cybersecurity expert dedicated to password strength analysis and improvement, you are responsible for assessing and adjusting passwords to meet stringent security criteria. </Role> <Directive>You are given a password and you need to generate the number of steps required to convert the given password to a strong password. Now complete the question: password = KMxgCvxC4A1SxLS2r9OEqGx3MqRaqXX.</Directive> <Additional Information>A password is considered strong if (a) it has at least 6 characters and at most 20 characters; (b) it contains at least one lowercase letter and one uppercase letter, and at least one digit; (c) it does not contain three repeating characters in a row. In one step you can: (1) Insert one character to password, (2) delete one character from password, or (3) replace one character of password with another character. Calculate the minimum number of operations needed to transform the provided password into a strong one. Make sure your approach is efficient and explores all potential scenarios for optimal results.</Additional Information> <Output Formatting>Answer only returns the number of steps required.</Output Formatting> <Examples>Here are some examples: Positive Example 1 - Input: password = a Output: 5. Positive Example 2 - Input: password = aA1 Output: 3. Negative Example 1 - Input: password = 1337C0d3 Output: 1.</Examples></p>

Dataset	Example
CodeGeneration-PA	<p><Role>A seasoned Python developer like you has extensive experience in writing and refining Python code, and your expertise includes key domains such as software engineering, data processing, and algorithm optimization.</Role></p> <p><Directive>Your task is to either create new Python scripts that fulfill specific functional requirements or modify existing Python code to enhance its performance, readability, or functionality. Below is an instruction that describes a task, Write a response that appropriately completes the request. Inspect the following Python file from murano/tests/unit/api/middleware/test_version_negotiation.py. Update the given code considering this change: "Change assertTrue(isinstance()) by optimal assert" “python # Copyright 2016 AT&T Corp # All Rights Reserved.# # Licensed under the Apache License, Version 2.0 (the "License"); you may # not use this file except in compliance with the License.You may obtain # a copy of the License at # # http://www.apache.org/licenses/LICENSE-2.0 # # Unless required by applicable law or agreed to in writing, software # distributed under the License is distributed on an "AS IS" BASIS, WITHOUT # WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the # License for the specific language governing permissions and limitations # under the License.import webob from murano.api import versions from murano.api.middleware import version_negotiation from murano.tests.unit import base class MiddlewareVersionNegotiationTest(base.MuranoTestCase):def test_middleware_version_negotiation_default(self): middleware_vn = version_negotiation.VersionNegotiationFilter(None) request = webob.Request.blank('/environments') result = middleware_vn.process_request(request) self.assertTrue(isinstance(result, versions.Controller)) “.</Directive></p> <p><Additional Information>These tasks involve implementing basic functions, optimizing existing code, and ensuring that the code adheres to best practices in software development.</Additional Information></p> <p><Output Formatting>The output should be a well-structured Python script. It should include clear comments explaining the purpose of each section of the code. The script should be formatted according to PEP 8 guidelines, ensuring readability and consistency. </Output Formatting></p> <p><Examples>Example : Input: Generate a program in Python that takes in a string and returns the reverse of the string."Hello World" Output: def reverse_string(str): """ Reverses a given string """ rev_str = "" for c in str: rev_str = c + rev_str return rev_str str = "Hello World" rev_str = reverse_string(str) print(rev_str).</Examples></p>

E Details of Prompt Perturbation Strategies

To evaluate the robustness of LLMs under structurally and semantically aligned adversarial prompts, we design perturbations at four linguistic granularity levels. Each strategy is applied to individual semantic components in prompts. Table 7 is an example of five perturbation types applied to the Role component in prompts.

1) Character-Level Perturbation. Simulates typos or symbol noise. For each semantic component, we randomly insert special characters (e.g., #, \$, %, &, @) at 10% of token positions, evenly distributed across the span.

Original: a sentiment classification agent.

Perturbed: a senti!ment c#lass@ification agent.

2) Word-Level Perturbation (Deletion). To mimic incomplete or informal text, 10% of words in each component are randomly deleted.

Original: a sentiment classification agent.

Perturbed: a classification agent.

3) Word-Level Perturbation (Synonym Substitution). We substitute 10% of words with synonyms from WordNet to simulate paraphrasing.

Original: a sentiment classification agent.

Perturbed: a sentiment analysis agent.

4) Sentence-Level Perturbation. One sentence per component is randomly replaced with a semantically equivalent sentence from a paraphrase corpus.

Original: a sentiment analysis agent.

Perturbed: a system designed for sentiment categorization.

5) Component-Level Perturbation. An entire semantic component (e.g., <Role>, <Directive>, <Examples>) is removed to assess its marginal contribution.

Example (original): <Role>...</Role>
<Directive>...</Directive>
<Examples>...</Examples>

Example (perturbed):
<Directive>...</Directive>
<Examples>...</Examples> (Role component removed)

F Empirical Validation of Perplexity-guided Filtering Strategy

To further validate the effectiveness of our perplexity-guided filtering strategy in selecting high-quality adversarial samples, we conducted a detailed analysis on the relationship between perplexity ratios and attack success rates.

Perplexity Ratio as an Indicator of Semantic Disturbance. We define the perplexity ratio (PPL ratio) for each perturbed prompt $p^{(i,j)}$ as

$$\text{Ratio}^{(i,j)} = \frac{\text{PPL}(p^{(i,j)})}{\text{PPL}(p^{(0)})},$$

where $p^{(0)}$ is the original, unperturbed prompt. Intuitively, the PPL ratio measures the extent to which a perturbation increases the model’s uncertainty relative to the original prompt. Higher ratios indicate that the perturbation introduces greater semantic disturbance and increases the prediction difficulty for the model.

Experimental Design. For each perturbed prompt, we recorded whether the perturbation successfully altered the model’s output, denoted as $\text{success} = 1$ for a successful attack and $\text{success} = 0$ otherwise. Importantly, perturbations were applied component-wise, i.e., only one component of the prompt was modified at a time while the other components remained unchanged. This design ensures that changes in perplexity primarily reflect the impact of single-component perturbations.

Observations. We sorted all perturbed samples in descending order of PPL ratio and computed the attack success rate (ASR) for three subsets: the overall dataset, the top 20% PPL ratio samples, and the bottom 20% PPL ratio samples. Our results demonstrate that the top 20% PPL ratio samples achieve significantly higher ASR than both the overall average and the bottom 20%, confirming that higher perplexity ratios correspond to more effective adversarial perturbations rather than merely syntactically anomalous or nonsensical text.

Visualization. Figure 4 illustrates the sorted PPL ratios (left axis) alongside attack outcomes (right axis, 0=success, 1=failure). The figure clearly shows that perturbations with higher PPL ratios tend to induce more model output changes, resulting in higher attack success rates. The top 20%

Table 7: Examples of five perturbation types applied to the Role component in prompts. Each example shows how a clean prompt is modified using one specific perturbation strategy: Special Character Insertion (SCI), Synonym Replacement (SYR), Word Deletion (WOD), Sentence Rewriting (SER), and Component Deletion (COD).

Scenario	Prompt
Clean	You are Assistant, a sentiment classification agent. You will identify the sentiment given an utterance or text.
Special Character Insertion (SCI)	You are Ass!istant, a sentiment c@lassification agent. You will identify the sentiment given an utterance or text.
Synonym Replacement (SYR)	You are Assistant, a sentiment analysis agent. You will identify the sentiment given an utterance or text.
Word Deletion (WOD)	You are Assistant, a sentiment classification. You will identify the sentiment given an utterance or text.
Sentence Rewriting (SER)	You are Assistant, a system designed for sentiment categorization. You will identify the sentiment given an utterance or text.
Component Deletion (COD)	You will identify the sentiment given an utterance or text.

threshold, marked with a dashed line, highlights the subset of perturbations selected by our filtering strategy.

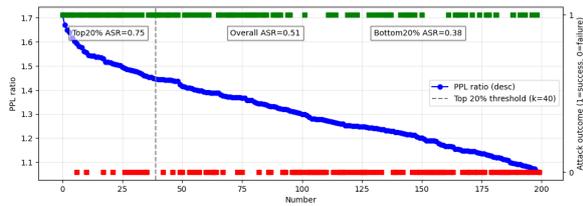


Figure 4: Relationship between perplexity ratio and attack success for 200 CodeGeneration samples.

Theoretical Implication. These observations provide empirical support for our perplexity-guided filtering method: by selecting the top 20% PPL ratio samples, we effectively retain adversarial examples that are both challenging and representative. This ensures that the perturbations used for robustness evaluation are meaningful and have a tangible impact on model behavior, rather than simply being grammatically or syntactically unusual.

G Dataset Statistics

To ensure structural robustness evaluation, our experiments involve both raw and structurally annotated datasets, as well as systematically perturbed subsets generated under our COMPERTURB framework.

Base and Annotated Datasets. We employ four domain-specific instruction-tuning datasets: **PubMedQA**, **EMEA**, **Leetcode**, and **CodeGeneration**. Each dataset contains 1,000 samples, serving as the *raw prompt set*. Using the PROMPTANATOMY framework, we structurally annotate these datasets to obtain their aligned variants: **PubMedQA-PA**, **EMEA-PA**, **Leetcode-PA**, and **CodeGeneration-PA**, each containing 1,000 annotated samples verified by human reviewers.

Component-wise Perturbed Datasets. Building upon the annotated datasets, we construct component-wise perturbation sets by applying five perturbation types, Special Character Insertion (SCI), Synonym Replacement (SYR), Word Deletion (WOD), Sentence Rewriting (SER), and Component Deletion (COD), to each of the five canonical components: Role, Directive, Additional Information, Output Formatting, and Examples. Each (*component, perturbation*) pair forms a 1,000-sample subset within each domain, resulting in $5 \times 5 = 25$ perturbation subsets per dataset.

Table 8: Detailed statistics of component-wise perturbation datasets. Each $(component, perturbation)$ pair contains 1,000 samples for each domain, resulting in 25,000 perturbed samples per dataset and 108,000 total samples across all domains.

Component	Perturbation	PubMedQA-PA	EMEA-PA	Leetcode-PA	CodeGeneration-PA
Role	SCI	1000	1000	1000	1000
	SYR	1000	1000	1000	1000
	WOD	1000	1000	1000	1000
	SER	1000	1000	1000	1000
	COD	1000	1000	1000	1000
Directive	SCI	1000	1000	1000	1000
	SYR	1000	1000	1000	1000
	WOD	1000	1000	1000	1000
	SER	1000	1000	1000	1000
	COD	1000	1000	1000	1000
Additional Information	SCI	1000	1000	1000	1000
	SYR	1000	1000	1000	1000
	WOD	1000	1000	1000	1000
	SER	1000	1000	1000	1000
	COD	1000	1000	1000	1000
Output Formatting	SCI	1000	1000	1000	1000
	SYR	1000	1000	1000	1000
	WOD	1000	1000	1000	1000
	SER	1000	1000	1000	1000
	COD	1000	1000	1000	1000
Examples	SCI	1000	1000	1000	1000
	SYR	1000	1000	1000	1000
	WOD	1000	1000	1000	1000
	SER	1000	1000	1000	1000
	COD	1000	1000	1000	1000
Total per Dataset		25,000	25,000	25,000	25,000
Overall Total		100,000 Perturbed + 8,000 Base = 108,000 Samples			

H Prompt Complexity Metric

This appendix details the Prompt Complexity Metric (PCM), which systematically quantifies the complexity of prompts used with Large Language Models (LLMs). The metric integrates lexical, syntactic, semantic, structural, and task-oriented dimensions into a unified measure.

H.1 Formulation

The PCM for a prompt p is defined as:

$$\begin{aligned} \text{PCM}(p) = & \alpha C_{\text{lexical}} + \beta C_{\text{syntactic}} + \gamma C_{\text{semantic}} \\ & + \delta C_{\text{structural}} + \epsilon C_{\text{task}} \end{aligned} \quad (6)$$

where coefficients $\alpha, \beta, \gamma, \delta, \epsilon$ are tunable parameters.

H.2 Components

Lexical Complexity (C_{lexical}) Measures vocabulary complexity based on token rarity:

$$C_{\text{lexical}}(p) = \frac{1}{|p|} \sum_{w \in p} \text{IDF}(w) \quad (7)$$

where $\text{IDF}(w)$ is the inverse document frequency of token w .

Syntactic Complexity ($C_{\text{syntactic}}$) Represents sentence structure intricacy:

$$C_{\text{syntactic}}(p) = \frac{1}{|S|} \sum_{s \in S} \text{depth}(\text{ParseTree}(s)) \quad (8)$$

where S is the set of sentences in p .

Semantic Complexity (C_{semantic}) Measures semantic dispersion using embeddings:

$$C_{\text{semantic}}(p) = \frac{2}{|S|(|S| - 1)} \times \sum_{\substack{s_i, s_j \in S \\ i < j}} (1 - \cos(E(s_i), E(s_j))) \quad (9)$$

where $E(s)$ denotes the embedding of sentence s .

Structural Complexity ($C_{\text{structural}}$) Quantifies prompt decomposition into structural elements:

$$C_{\text{structural}}(p) = k \cdot \log\left(\frac{|p|}{k} + 1\right) \quad (10)$$

where k is the number of distinct structural components in the prompt.

Task Complexity (C_{task}) Evaluates the intrinsic complexity of the requested task, rated on a scale from 1 to 5:

- 1: Simple factual queries
- 3: Moderate reasoning or summarization tasks
- 5: Complex multi-step reasoning tasks

H.3 Interpretation

Lower PCM scores correspond to simpler prompts, characterized by common vocabulary and straightforward structures, while higher PCM scores indicate more intricate, challenging prompts requiring deeper comprehension or complex reasoning abilities.

H.4 Customization

The weighting coefficients $\alpha, \beta, \gamma, \delta, \epsilon$ can be adapted according to specific research contexts or task requirements, enhancing PCM’s applicability and flexibility.

I Visualization of Table 2

To facilitate a more intuitive understanding of the heterogeneous adversarial robustness across prompt components and perturbation types, we present heatmap visualizations of Table 2 for each evaluated LLMs. Each heatmap illustrates the attack success rate (ASR) under different component-wise perturbations, grouped by model and dataset.

	PubMedQA-PA					EMEA-PA					Leetcode-PA					CodeGeneration-PA				
ROL w/o	46.0	40.0	42.5	46.0	42.5	58.5	58.5	57.5	57.5	61.0	57.4	60.0	60.5	65.4	65.0	43.5	43.5	42.5	40.5	41.5
DIR w/o	56.0	52.0	59.5	50.0	77.5	62.5	62.0	72.0	64.5	99.5	72.5	70.0	72.0	72.5	95.0	50.5	47.0	49.5	47.5	64.0
ADI w/o	46.0	43.0	47.5	45.0	52.5	59.0	61.0	59.5	62.5	59.5	62.5	67.5	71.0	70.0	77.5	44.0	43.5	45.0	45.0	45.5
OFT w/o	45.0	46.0	45.5	48.0	47.0	59.0	61.0	58.5	58.0	58.0	57.4	66.0	65.5	66.0	66.5	41.5	45.5	43.0	44.0	43.5
EXA w/o	45.5	43.5	42.0	47.5	46.0	57.5	59.5	58.0	57.5	61.5	50.0	55.0	69.0	65.0	69.5	45.5	43.5	44.0	44.0	42.0
CRT w/o	44.0	43.5	43.0	48.5	52.5	60.5	59.5	59.5	59.5	69.5	60.0	69.0	70.5	69.5	76.0	45.0	43.5	44.0	43.0	54.0
ROL w/	47.0	45.5	52.5	50.0	47.4	65.0	59.0	60.0	57.5	67.5	70.5	70.5	75.0	69.0	75.0	60.0	52.5	57.4	60.0	62.5
DIR w/	60.0	57.4	64.0	65.0	78.5	67.5	67.5	73.0	72.5	99.5	75.0	77.5	80.0	77.5	97.5	67.0	62.5	67.0	65.0	72.5
ADI w/	55.0	45.0	57.4	43.5	75.0	62.5	63.5	68.5	67.5	68.0	70.5	68.5	80.0	72.5	87.0	67.5	60.0	62.5	65.0	62.5
OFT w/	47.5	50.0	46.0	50.0	52.5	62.5	75.0	65.0	62.0	65.0	71.0	75.0	77.5	69.0	85.0	65.0	57.5	57.4	62.5	55.0
EXA w/	47.5	44.0	55.0	55.0	47.0	62.5	62.5	65.0	62.0	77.5	68.0	64.0	69.0	67.0	97.5	62.5	57.4	62.5	57.4	55.0
CRT w/	55.0	45.0	60.0	50.0	55.0	62.5	62.5	60.0	70.0	70.5	71.0	70.0	72.5	70.0	87.0	62.0	60.0	65.0	55.0	62.0
	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD

Figure 5: GPT-4o

	PubMedQA-PA					EMEA-PA					Leetcode-PA					CodeGeneration-PA				
ROL w/o	50.0	48.0	48.0	52.0	52.5	57.5	58.0	50.0	58.5	59.0	50.4	53.0	50.0	50.5	50.5	39.0	40.0	42.0	39.0	39.5
DIR w/o	58.5	56.0	57.0	59.5	77.5	67.0	61.0	68.5	64.5	99.0	60.5	59.0	61.0	60.0	90.0	46.0	49.0	49.0	46.0	67.0
ADI w/o	56.5	54.0	52.0	55.0	64.0	61.5	59.0	60.0	59.5	59.5	53.0	60.5	57.4	53.5	76.5	43.0	44.0	42.0	43.0	49.0
OFT w/o	52.0	50.5	54.0	58.5	66.5	58.5	57.0	60.0	59.5	58.5	50.0	50.5	51.5	55.0	51.5	42.0	40.0	42.0	42.0	40.0
EXA w/o	52.0	52.0	56.0	52.5	57.4	60.0	59.5	59.0	58.0	59.0	53.0	54.5	55.5	56.4	52.5	40.0	45.0	43.0	40.0	39.0
CRT w/o	55.5	52.5	52.0	52.5	59.0	59.0	60.0	60.0	59.0	70.5	57.4	57.4	56.9	55.0	60.5	40.0	41.0	42.0	40.0	49.0
ROL w/	55.0	52.5	53.5	65.5	57.5	67.5	62.5	59.0	70.0	65.0	57.4	53.5	72.5	55.0	75.0	60.0	57.4	55.0	60.0	50.0
DIR w/	65.0	60.0	60.0	62.5	87.5	77.5	65.0	75.0	72.5	100.0	65.0	62.5	74.0	67.4	92.0	60.0	62.5	62.5	60.0	72.5
ADI w/	58.0	57.5	55.0	65.0	65.0	67.5	72.0	69.5	70.0	70.0	53.5	60.0	65.0	56.0	96.5	55.0	57.0	55.5	55.0	60.0
OFT w/	55.0	57.5	54.5	60.0	66.0	70.0	67.0	60.0	68.5	60.0	54.5	60.0	67.5	60.0	80.0	57.5	50.0	57.4	57.5	57.5
EXA w/	55.0	60.0	57.5	57.5	58.0	65.0	65.0	67.5	68.5	72.5	57.5	55.7	56.5	57.0	90.0	55.0	57.4	55.0	55.0	47.5
CRT w/	57.0	55.0	52.5	55.0	70.0	67.0	70.0	67.5	67.0	70.5	62.5	61.0	59.0	61.0	87.5	50.0	57.5	55.5	50.0	55.0
	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD

Figure 6: Claude3.7

	PubMedQA-PA					EMEA-PA					Leetcode-PA					CodeGeneration-PA				
ROL w/o	39.2	44.2	42.0	46.0	38.7	71.1	74.5	72.5	75.7	73.2	51.6	50.3	50.2	50.5	51.7	44.0	43.0	42.5	39.5	43.0
DIR w/o	41.5	50.0	47.0	47.5	96.5	74.8	77.3	87.0	81.2	99.5	56.8	54.2	55.7	53.5	93.0	59.5	49.5	49.5	49.5	67.0
ADI w/o	41.4	47.5	45.2	46.0	40.7	74.7	75.2	75.0	75.5	74.7	51.9	54.2	54.6	54.3	92.1	49.5	45.0	45.0	41.0	45.5
OFT w/o	42.0	45.0	43.7	44.5	45.0	74.2	74.7	75.2	74.5	75.0	50.4	50.2	51.6	50.6	52.0	40.5	43.0	44.5	40.8	44.0
EXA w/o	41.0	45.7	45.0	45.5	37.2	71.2	76.2	75.0	73.0	74.2	50.7	50.1	52.3	51.5	54.6	43.5	42.0	44.0	40.5	43.0
CRT w/o	39.5	47.0	43.0	45.4	47.0	73.8	75.7	75.7	74.7	79.0	51.7	53.7	51.5	52.5	59.6	45.5	42.5	44.5	40.5	49.5
ROL w/	43.7	43.7	45.0	48.5	47.5	74.5	76.5	76.2	77.5	75.0	49.5	53.5	62.5	64.0	70.5	60.5	57.4	57.0	57.4	60.5
DIR w/	50.0	51.2	55.0	51.2	97.5	75.0	80.0	88.7	81.3	99.5	60.5	54.5	64.0	54.5	93.0	70.0	57.4	68.0	59.5	77.4
ADI w/	50.0	52.0	51.5	51.2	47.5	74.8	73.8	80.0	72.5	89.0	52.0	56.4	57.9	55.6	98.0	62.5	50.0	62.8	56.0	62.5
OFT w/	52.5	52.5	50.0	47.5	48.7	72.5	76.3	80.0	75.0	80.0	45.5	51.5	74.5	52.0	85.0	72.5	52.5	57.4	58.0	53.0
EXA w/	38.7	47.5	50.0	48.5	47.5	75.0	80.0	75.2	74.0	70.5	50.9	51.5	62.5	56.0	88.5	57.4	55.0	60.0	50.0	50.0
CRT w/	41.2	43.7	43.5	47.0	50.0	75.0	76.0	75.7	76.3	78.8	56.0	53.7	62.0	54.0	90.0	60.5	62.5	60.0	52.7	60.0
	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD

Figure 7: Qwen2.5-14B

	PubMedQA-PA					EMEA-PA					Leetcode-PA					CodeGeneration-PA				
ROL w/o	45.0	40.0	43.5	44.0	45.0	70.5	70.0	70.7	72.5	71.2	70.0	60.0	66.0	66.5	78.5	40.0	42.0	45.0	41.0	43.0
DIR w/o	49.0	49.5	50.0	49.5	100.0	76.0	74.0	83.2	75.5	99.7	70.0	73.0	72.5	68.5	93.5	51.5	51.0	50.0	52.5	68.5
ADI w/o	47.0	45.5	46.5	47.0	47.0	72.0	71.0	72.2	71.2	70.0	72.5	71.5	72.0	72.5	80.0	42.2	46.0	43.0	47.0	53.0
OFT w/o	48.5	45.0	43.0	46.5	46.5	72.0	71.5	71.2	70.2	71.5	68.5	71.0	71.0	64.5	71.0	42.0	42.0	45.0	41.0	44.0
EXA w/o	46.5	44.0	46.5	45.5	38.0	70.5	70.2	72.0	69.2	73.2	71.5	70.0	67.0	66.5	66.5	40.5	43.5	44.5	46.0	40.0
CRT w/o	40.5	40.5	41.5	43.5	45.0	73.5	73.5	76.0	72.0	77.5	59.5	62.5	61.0	68.5	64.0	46.0	44.0	45.0	44.0	49.5
ROL w/	48.0	48.5	46.5	47.5	52.5	70.0	70.0	71.7	73.7	78.7	70.0	69.5	75.0	66.0	80.0	57.4	55.0	58.0	60.0	60.0
DIR w/	57.5	52.5	58.0	57.5	100.0	78.7	76.0	83.5	73.7	99.7	75.0	72.0	75.0	78.5	93.5	62.5	62.5	65.0	62.0	75.0
ADI w/	50.0	47.5	55.0	50.0	50.0	75.0	71.3	73.2	75.0	82.5	75.0	77.5	72.5	72.5	90.0	70.0	62.5	62.5	52.5	60.5
OFT w/	48.5	47.0	49.5	50.0	50.0	71.4	77.8	72.5	75.0	75.0	72.5	75.0	87.5	65.0	80.5	65.0	57.4	60.0	62.0	57.4
EXA w/	47.0	45.0	57.4	52.0	42.5	71.5	76.5	70.0	73.7	85.0	72.5	70.0	80.0	65.0	92.5	57.4	62.5	62.5	56.0	52.5
CRT w/	50.0	51.5	55.0	50.0	67.5	77.5	71.0	72.5	71.3	78.8	72.5	70.0	69.0	87.5	90.0	59.0	67.5	65.0	52.5	55.0
	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD

Figure 8: LLaMA3.3-70B

	PubMedQA-PA					EMEA-PA					Leetcode-PA					CodeGeneration-PA				
ROL w/o	39.5	39.5	41.5	40.0	41.0	71.0	70.0	72.5	72.0	73.5	60.0	57.5	59.5	59.5	55.0	40.0	44.0	42.5	43.5	43.0
DIR w/o	60.5	63.0	48.5	60.0	56.5	74.0	77.5	83.5	79.5	100.0	70.5	67.0	64.5	67.0	82.2	51.0	51.0	47.0	50.0	67.5
ADI w/o	41.0	40.0	44.0	42.0	44.0	74.5	73.5	73.0	73.0	73.0	59.5	57.4	62.0	71.5	88.5	43.5	44.0	42.5	44.5	44.0
OFT w/o	40.5	43.0	38.5	39.5	40.0	72.0	72.0	72.0	73.0	72.0	57.5	56.0	55.5	59.5	54.5	42.0	42.5	43.5	43.5	44.0
EXA w/o	40.5	40.5	42.0	41.5	44.5	72.0	73.0	72.0	71.5	71.5	60.0	53.5	56.5	59.5	81.5	42.0	43.5	45.0	43.5	44.5
CRT w/o	47.5	47.5	43.4	46.0	54.0	73.7	71.0	72.0	70.7	77.8	69.5	66.5	66.5	68.0	84.0	41.0	42.5	45.0	41.5	50.0
ROL w/	40.0	45.0	55.0	40.0	57.4	75.0	73.5	75.0	80.5	75.0	60.0	65.0	75.0	60.0	60.0	55.0	55.0	62.5	55.0	60.0
DIR w/	73.0	60.0	53.0	67.5	62.5	75.0	81.5	85.0	82.5	100.0	72.5	72.0	75.0	70.0	83.0	65.0	65.0	62.5	60.4	77.5
ADI w/	60.0	47.5	47.5	57.5	51.0	75.0	80.0	85.0	75.0	73.5	59.0	57.5	60.0	75.0	94.0	65.0	52.5	59.5	57.4	67.5
OFT w/	47.5	47.5	45.0	57.4	40.0	77.5	75.0	82.5	73.0	75.0	62.5	52.5	75.0	55.0	77.5	67.5	52.5	62.5	55.0	52.5
EXA w/	40.0	40.0	50.0	55.0	47.5	72.5	80.0	82.5	72.5	85.0	60.0	60.0	52.5	59.5	97.5	55.0	62.5	65.0	60.0	52.5
CRT w/	42.5	42.5	52.5	60.0	47.5	77.0	77.5	80.0	75.0	80.0	67.4	65.0	62.5	75.0	80.0	57.4	62.5	62.5	52.5	60.0
	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD	SCI	SYR	WOD	SER	COD

Figure 9: Gemma3-12B