



PROMPTPRISM: A Linguistically-Inspired Taxonomy for Prompts

Sullam Jeung, Yueyan Chen, Shuai Wang, Haibo Ding, Yi Zhang, Lin Lee Cheong
Amazon AWS

(sullamij, yyanc, wshui, hbding, yizhngn, lcheong)@amazon.com

Abstract

Prompts are the interface for eliciting the capabilities of large language models (LLMs). Understanding their structure and components is critical for analyzing LLM behavior and optimizing performance. However, the field lacks a comprehensive framework for systematic prompt analysis and understanding. We introduce  PROMPTPRISM, a linguistically-inspired taxonomy that enables prompt analysis across three hierarchical levels: functional structure, semantic component, and syntactic pattern. By applying linguistic concepts to prompt analysis, PROMPTPRISM bridges traditional language understanding and modern LLM research, offering insights that purely empirical approaches might miss. We show the practical utility of PROMPTPRISM by applying it to three applications: (1) a taxonomy-guided prompt refinement approach that automatically improves prompt quality and enhances model performance across a range of tasks; (2) a multi-dimensional dataset profiling method that extracts and aggregates structural, semantic, and syntactic characteristics from prompt datasets, enabling comprehensive analysis of prompt distributions and patterns; (3) a controlled experimental framework for prompt sensitivity analysis by quantifying the impact of semantic reordering and delimiter modifications on LLM performance. Our experimental results validate the effectiveness of our taxonomy across these applications, demonstrating that PROMPTPRISM provides a foundation for refining, profiling, and analyzing prompts.

1 Introduction

The systematic study of prompt design and optimization has emerged as a critical research area in the evolution of Large Language Models (LLMs) (Brown et al., 2020; Touvron et al., 2023; Zhao et al., 2023). While effective prompts can significantly enhance model performance, studies have revealed that LLMs demonstrate notable sensitivity

to various prompt characteristics (Liu et al., 2024a; Sclar et al., 2023; Li et al., 2024), including stylistic formatting, sequential ordering, task explanations, and meta-rules. This sensitivity, coupled with inconsistent documentation of prompt specifications across the literature, presents challenges for systematic performance evaluation and comparative analysis.

Despite the fundamental importance of prompting in LLM interactions, there exists a notable gap in the development of frameworks that enable comprehensive analysis of prompts themselves. Previous approaches focus on analyzing generated responses along dimensions such as helpfulness and harmlessness (Askell et al., 2021), while systematic understanding of prompt characteristics remains understudied. While recent works have proposed standardized documentation protocols to enhance prompt transparency (Bach et al., 2022; Wang et al., 2022), these efforts often lack the depth and structure needed to capture the nuanced variations in prompt composition and their effects on model behavior.

Linguistic theory offers a valuable foundation for addressing this gap (Opitz et al., 2025). Far from being merely a nostalgic endeavor, linguistic approaches provide established analytical frameworks for understanding how language functions as a structured communication system (Ravichander et al., 2022; Ponti et al., 2020). These principles remain relevant in the era of large language models precisely because LLMs, despite their neural architecture, fundamentally process and generate language. The historical application of linguistics in NLP has demonstrated that linguistic expertise is crucial for developing high-quality resources and evaluation methods (Opitz et al., 2025) – from selecting diverse and representative language data (Ravichander et al., 2022) to ensuring proper annotation standards and designing effective evaluation metrics (Chi et al., 2024; Parrish et al., 2021).

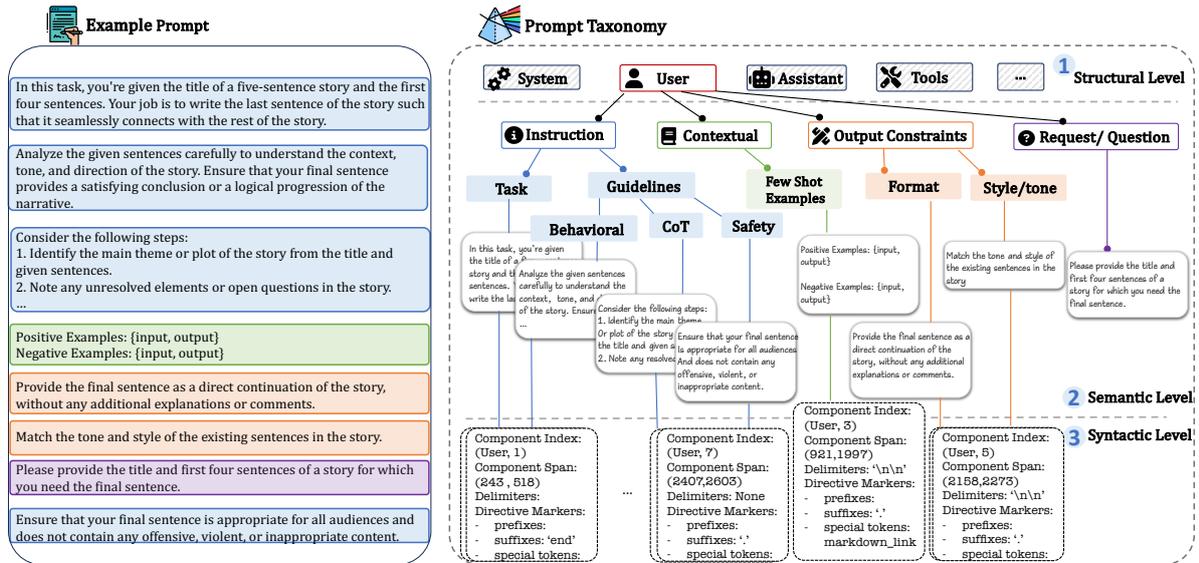


Figure 1: PROMPTPRISM Framework. The PROMPTPRISM taxonomy decomposes an example prompt from Super Natural Instruction dataset (Left) into three hierarchical levels (Right). (1) **Structural**, corresponding to specified roles in the prompt (e.g. System, User, Assistant); (2) **Semantic**, detailing the content’s semantic components (e.g. Instruction, Output Constraints); and (3) **Syntactic**, describing the syntactic structure of each semantic element (e.g. Delimiters, Special tokens).

To address these limitations and meet the demands of increasingly sophisticated LLM applications, we propose a comprehensive and linguistically-inspired taxonomy that integrates structural, semantic, and syntactic dimensions of prompt analysis. Our framework, PROMPTPRISM, treats prompts as structured discourse units with plan-based intentions, enabling hierarchical decomposition and analysis across multiple dimensions. This approach allows for systematic examination of prompt characteristics that influence model behavior, providing insights that purely empirical approaches might miss.

Our research makes two primary contributions. (1) First, we propose a novel taxonomy framework for prompts that synthesizes linguistic theory with practical applications, enabling systematic analysis and categorization of prompt characteristics (§3). (2) Second, we demonstrate the practical utility of our framework through three distinct applications: taxonomy-guided prompt refinement (§4.1), comprehensive dataset profile extraction (§4.2), and systematic prompt sensitivity analysis (§4.3). This work represents a significant step toward understanding prompts, enhancing the prompt quality, and enabling controlled experimental setups for LLM development and evaluation.

2 Related Works

2.1 Taxonomy of Prompts

Recent research has introduced various taxonomic approaches for analyzing LLM prompts (Bach et al., 2022; Schulhoff et al., 2024), with notable contributions including the TELeR framework (Santu and Feng, 2023), which categorizes prompts into data and instructions components with difficulty-based classification, and Budagam et al. (2024)’s hierarchical prompting taxonomy grounded in cognitive principles. While these frameworks have advanced our understanding of prompt structure and evaluation, they remain limited in their ability to capture and differentiate fine-grained components of prompts, such as semantic variations and stylistic subtleties.

2.2 Understanding Prompts

The systematic study of prompt design and optimization has emerged as a critical research area in the evolution of Large Language Models. While effective prompts can significantly enhance model performance (Wei et al., 2022; Shinn et al., 2023), studies have revealed that LLMs demonstrate notable sensitivity to various prompt characteristics, including stylistic formatting (Sclar et al., 2023; He et al., 2024), sequential ordering (Liu et al., 2024a; Huang et al., 2025), task explanations (Li et al., 2023a), and meta-rules (Li et al., 2024). This sensitivity, coupled with inconsistent documentation of

prompt specifications across the literature, presents challenges for systematic performance evaluation and comparative analysis. This necessitates a structured approach to understand and analyze prompts such as the organization of semantic components and their stylistic representations.

Due to limited space, the related works regarding structured prompts can be found in Appendix D.

3 PROMPTPRISM

We define a prompt P as an ordered sequence of n pairs, where each pair consists of content and its associated role, formally expressed as:

$$P = (r_i, c_i)_{i=1}^n, r_i \in R$$

where $R = \{\text{system, user, assistant, tools} \dots\}$ denotes the finite set of possible roles. The content c_i belongs to a modality space M , which encompasses both single modalities and their combinations: $M \subseteq \{\text{text, image, audio}\} \cup \{m_1 \otimes m_2 \mid m_1, m_2 \in M\}$ Here, \otimes represents the operation of combining different modalities. While our framework accommodates multi-modal prompts through this formalization, the present work focuses on text modality. For instance, an example of text-modality prompt P would be, (system, ‘You’re a helpful assistant’), where system corresponds to the role r , and the content corresponds to the c .

In the development of our framework, we adhere to a top-down, property-guided formation approach. The framework is designed to meet the following key criteria:

✓ **Simplicity:** We prioritize minimal complexity in classification, ensuring clear parent-child relationships and intuitive categorization boundaries. This approach facilitates ease of understanding and application across diverse prompt types.

✓ **Coherence:** The framework maintains logical consistency across categories, minimizing redundancy in classification. We establish well-defined relationships between elements, ensuring a cohesive structure that accurately represents the interplay between different prompt components.

✓ **Comprehensiveness:** Our taxonomy aims for complete coverage of possible prompt types, incorporating both atomic and composite structures. The framework is designed with scalability in mind, allowing for the accommodation of new modalities as the field of prompt evolves.

✓ **Utility:** The framework is developed with practical applicability. It facilitates systematic prompt design and enables meaningful analysis of prompts. This utility-focused approach ensures that the taxonomy serves as a valuable tool for both researchers and practitioners in the field.

3.1 Prompt Taxonomy

The proposed taxonomy adopts a hierarchical tree-like structure, where the structural level decomposes into semantic components, each containing associated syntactic information (shown in Fig 1)

① **Structural Level** At the structural level, PROMPTPRISM conceptualizes prompts as organized discourse units, drawing parallels with traditional linguistic analysis of text structure (Mann et al., 1992; Grosz and Sidner, 1986). This approach is similar to how linguistics examines the organization of different text genres, such as academic articles with their standard sections (introduction, methodology, results).

In the context of prompts, this structural organization is defined through roles, $r \in R$, each serving distinct discourse purposes. A key aspect of this structural level is the clear delineation of roles within the prompt, similar to the speaker-listener roles in discourse analysis (Santu and Feng, 2023; Brown et al., 2020). Specifically, we identify and distinguish between different roles in the prompt interaction: System, User, Assistant, Tools.

- System: Provides core instructions and framework
- User: Initiates queries or requests
- Assistant: Generates responses and ensures conversation flows naturally
- Tools: Handles specific functionalities or operations

While these four roles form the primary set in our current analysis (Table 5), it’s important to note that this role system is extensible and can accommodate new roles as prompt architectures evolve across different providers.

② **Semantic Level** The semantic layer of PROMPTPRISM establishes a systematic frame-

Semantic Components	Description	Examples
📌 Instruction	High-level directive component that guide the model's behavior	
Task	Task related instruction	This is a classification task.
Guidelines	Non-task specific instructions that shape response behavior	
Role assumption (persona)	Directives for the model to assume a specific role	"Suppose you are a {persona}.." (Perez et al., 2023)
Scenario	Context setting for the task environment	Imagine we gave you \$2000 right now..(Binz and Schulz, 2023)
Behavioral	Instructions for model conduct and interaction style	Before you respond, rephrase the question (Deng et al., 2022)
Emotion	Guidelines for emotional tone	You'd better be sure (Li et al., 2023b)
Chain of Thought Instruction	Directives for showing reasoning process	Let's think step by step (Wei et al., 2022)
Safety	Guidelines ensuring safe and ethical responses	"Avoid harmful content", "Maintain ethical boundaries"
📖 Contextual/Reference info	Background information and supporting materials	
Few-shot examples	Sample input-output pairs for in-context learning	Q: What is the capital of France? A: Paris
Knowledge Base	Reference information or facts	"Based on medical guidelines..."
Context for task	Relevant background information	Text for summarization
🛠️ Output Constraints	Specifications for response format and limitations	
Label Space	Defined set of possible output categories	"Choose from: [Positive, Negative, Neutral]"
Word Limits	Restrictions on response length	"Respond in 50 words or less", "Provide a one-paragraph answer"
Output format Specification	Structure requirements for the response	"Format as JSON", "Present in bullet points"
Style/tone	Requirements for writing style	"Use academic language", "Write in conversational tone"
🔧 Tools	Specifications for tool usage	
Tool name	Identifier for specific tool	"Calculator", "Python Interpreter"
Tool description	Explanation of tool functionality	"This tool performs statistical analysis"
Parameters	Required inputs and configuration	"Input format:(x,y)", "Required parameters: date_range, metric"
🗋️ User request / Question	The primary query or task from user	"What is the capital of France?" (e.g. question in Q&A task)
🗋️ Response	Semantic component for Model's output	
Answer	Direct Answer from the request / question	"Paris is the capital of France"
Peripheral-Explanation	Supporting information and clarifications	"This city has been the capital since.."
📌 Other purposeful components	Additional functional elements	
Adversarial	Components designed for adversarial purpose	"!!!!!" (Zou et al., 2023)
🗋️ Historical Context	Pointer to the previous conversation history	
🗋️ System Prompt	Pointer to the System Prompt	
🗋️ Tools Prompt	Pointer to the Tools Prompt	

Table 1: Semantic Components defined in PROMTPRISM taxonomy

work for analyzing the meaning and intentional structure of prompts, drawing direct parallels with linguistic pragmatics and semantic theory (Levinson, 1983; Grice, 1975). This layer decomposes the content c_i of a structural role r_i , into the semantic components $(c_i, r_i) \in P$ that serve specific functions.

Table 1 presents the overview of semantic components, with primary semantic components that mirror linguistic functional units (Searle, 1969; Austin, 1962). 📌 Instruction serve as the core semantic component, articulating the fundamental task or objective, similar to the main predicate in linguistic structures (Chomsky, 2014; Fillmore, 1967). These are complemented by 📖 Contextual/Reference information, which provides the necessary background knowledge and examples, functioning analogously to contextual presuppositions in linguistic pragmatics (Sperber and Wilson, 1986; Carston, 2008). 🛠️ Output Constraints define formatting and stylistic expectations, parallel to linguistic register and genre conventions. 🗋️ User request / Question refers to the primary query or task from user, such as a

question in question and answering task. 🔧 Tools indicate the specification of tools - tool name, tool description, and parameters.

These components are defined in terms of their functional purpose rather than abstract linguistic concepts, making them accessible to practitioners while maintaining their grounding in linguistic theory. The semantic relationships between these components exhibit what we term "dominance and satisfaction-precedence relations", wherein lower-level segments contribute incrementally to satisfying higher-level discourse purposes.

③ **Syntactic Level** The syntactic level of PROMTPRISM draws parallels with linguistic morphology while adapting these concepts for the unique requirements of prompts. As previous works (Sclar et al., 2023; He et al., 2024) have demonstrated that syntactic components of prompts such as stylistic patterns and ordering influence the model performance, this level introduces syntactic elements that collectively create a complete structural analysis framework. **Component Organization:** A component organization system plays on two primary mechanisms.

(1) Component Position assigns unique identifiers (Role, Index) to each prompt element, similar to syntactic tree structures in linguistics (Chomsky, 2014; Tesnière, 2015). (2) Component Span Analysis defines exact boundaries through position markers (start_pos, end_pos), enabling manipulation and reconstruction of prompt components. This component organization information enables prompt sensitivity analysis, reconstructing and perturbing various components.

Directive Markers: Our framework adapts linguistic concepts of morphemes for prompt analysis, but with a distinction. While linguistic morphology treats prefixes and suffixes as meaningful units that modify word meanings, we define directive markers purely in terms of their structural and sequential properties (Matthews, 1972; Aronoff, 1976). **Prefix Markers:** Unlike linguistic prefixes that modify word meanings (e.g. ‘un-’ in ‘unhappy’), our prefix patterns serve as structural indicators: (e.g. Hash comments (#), Double-slash comments (//), etc). **Suffix Markers:** Similarly, while linguistic suffixes can change word classes (e.g. ‘-ness’ in ‘happiness’), our suffix patterns function as structural boundaries: Colon endings (:), Sentence terminators (.!?), etc. **Delimiters:** We define the delimiter as the boundary between components. Detailed definition can be found in Table 4 and in Section B.

4 Applications of PROMPTPRISM

We demonstrate the applications of PROMPTPRISM and illustrate how it can be effectively leveraged in real-world scenarios. As shown in Figure 2, we focus on three important applications.

4.1 Taxonomy Guided Prompt Refinement

We demonstrate how PROMPTPRISM can be used to systematically analyze and improve prompt effectiveness. Our taxonomy guided prompt refinement approach replaces arbitrary prompt elaboration with a more semantically-rich approach, resulting in enhanced model performance compared to the base instructions.

Experiment Setup Dataset: We conduct our experiments using Super Natural Instruction dataset (Wang et al., 2022) (version 2.8), which provides each task with a base instruction, multiple data instances (consisting of question-answer pairs), and exemplars for few-shot learning. The detailed experiment setting is illustrated in Appendix A.

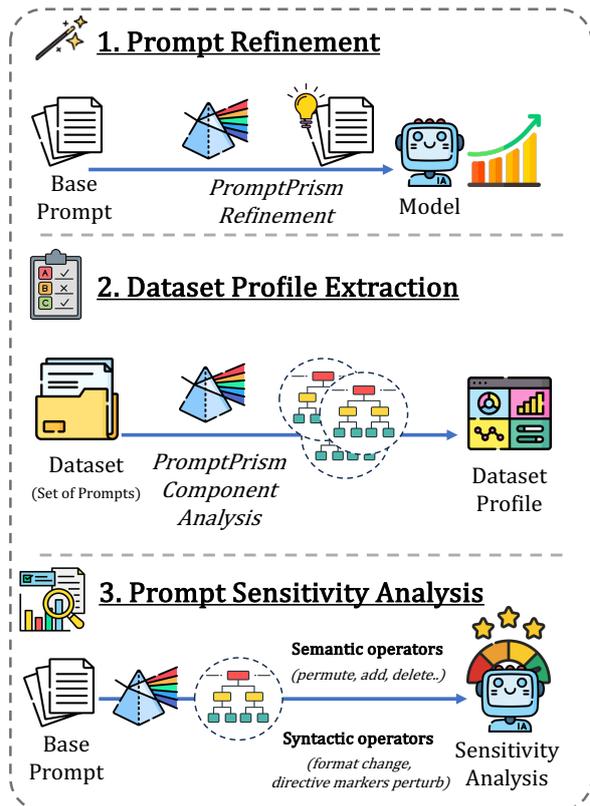


Figure 2: Applications of PROMPTPRISM

Metrics: For quantitative evaluation, we employ ROUGE-L metrics, adhering to the evaluation protocol established in (Wang et al., 2022). The prompt refinement process is conducted using Claude-sonnet-3.5 (Anthropic, 2024b), with prompts provided in Appendix G. For each task, we: (1) Generate 5 distinct prompts (setting temperature=0.7), (2) Sample 10 data instances randomly, (3) Perform 50 inference runs per task. We present the result based on using 2 few-shot examples in the main script. The results of 0-shot experiment are presented in the appendix Table 7. **Baseline** (1) **Default:** utilizes unmodified base instructions from the Super Natural Instruction dataset (Wang et al., 2022). (2) Chain-of-Thought (CoT) (Wei et al., 2022): appends the base instruction with "Please think step by step and then solve the task." at the end.^{1 2}

¹While some prompt optimization methodologies exist such as DSPy and TextGrad (Yuksekgonul et al., 2024; Khat-tab et al., 2023), they require iterative prompt generation and optimization processes, and development set for parameter tuning. To ensure practical applicability of cost-effectiveness, we constrain our experimental setup to single-prompt scenarios, eliminating the need for multiple prompt iterations to achieve optimization without any parameter tuning.

²We provide additional baseline results of OpenAI Meta prompt approach in Table 8

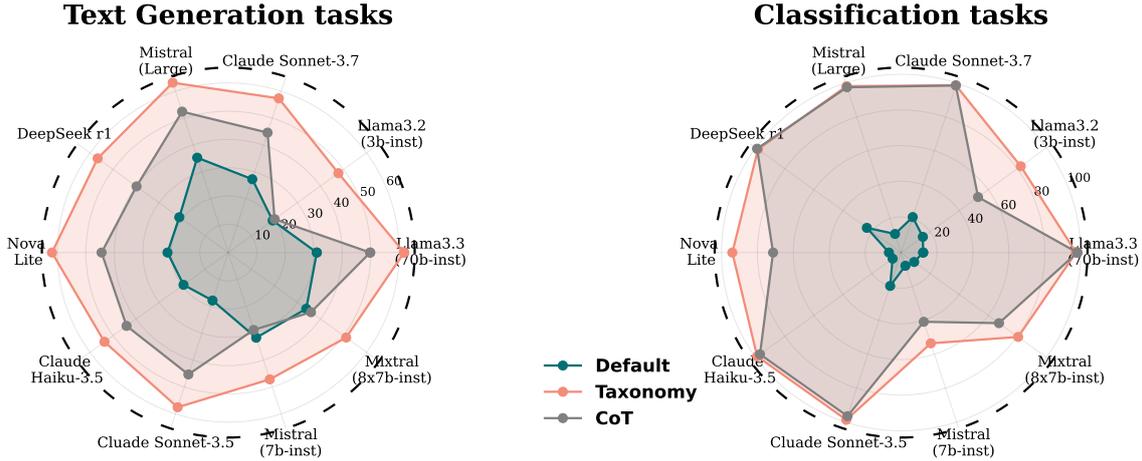


Figure 3: Performance comparison of taxonomy-guided prompt refinement across task types. Our taxonomy-guided approach demonstrates improved performance on text generation tasks (average **29%** improvement than CoT) and matches or exceeds CoT baseline on classification tasks (average **0.13%** improvement than CoT). Complete performance metrics are presented in Table 6.

Results Our experimental results, presented in Figure 3 and tables (Tab 6), demonstrate the efficacy of our taxonomy-guided prompt refinement approach across varying configurations. The approach exhibits particularly strong performance in zero-shot settings (Table 7), achieving substantial improvements over CoT baselines: a 112% performance increase in text generation tasks (CoT avg: 40.53, Ours: 55.81) and around 10% improvement in classification tasks (CoT avg: 81.93, Ours: 90.1).

The benefits of our methodology persist in few-shot scenarios as well. In the two-shot setting (Table 6), we observe consistent improvements across task types, with a 29% enhancement in text generation performance and a modest but consistent 0.13% improvement in classification tasks.

4.2 Multi-Dimensional Dataset Profiling

PROMPTPRISM enables comprehensive profiling of prompt datasets, extracting and aggregating characteristics across structural, semantic, and syntactic dimensions. This analysis helps identify gaps in existing benchmark datasets and guide the development of more detailed benchmark documentation.

Our analysis encompasses four primary dimensions: (1) Structural Analysis: We examine fundamental structural characteristics including turn type (single/multi-turn), prompt pattern (role sequence analysis), and unique structural roles (system, user, assistant, etc). (2) Dataset Annotation & Semantic Analysis: We conducted systematic annotation of data instances

using the PROMPTPRISM taxonomy framework with Claude-Sonnet-3.5 (Anthropic, 2024b) for semantic component extraction. The annotation process utilizes a structured XML-style format (`<tag> . . <\tag>`), enabling efficient extraction and analysis of semantic components (annotation prompt detailed in Appendix G). The semantic analysis encompasses two key elements: component frequency analysis, identifying the three most prevalent semantic components in the dataset, and hierarchical tree structure analysis, examining mean taxonomic tree width and mean tree depth to measure component breadth and semantic complexity. (3) Syntactic Analysis: From annotated data, we extract component-level syntactic information, including component indices, spans, delimiters, and directive markers. (4) Metadata Information Analysis: We derive task type with Claude-Sonnet-3.5 using the meta prompt in Appendix G. We also include other information such as language specification, token length, and modality characteristics.

Experiment Setup Dataset: Our analysis focuses on two distinct sub-datasets from SMOLTALK (Allal et al., 2025): APIGEN-80K (Liu et al., 2024b) and SMOL-MAGPIE-ULTRA (Allal et al., 2025). These datasets were selected for their distinct characteristics: APIGEN-80K exemplifies function-calling scenarios with tool utilization, while SMOL-MAGPIE-ULTRA represents multi-turn interactions. This selection demonstrates PROMPTPRISM’s applicability beyond single-turn interactions to more

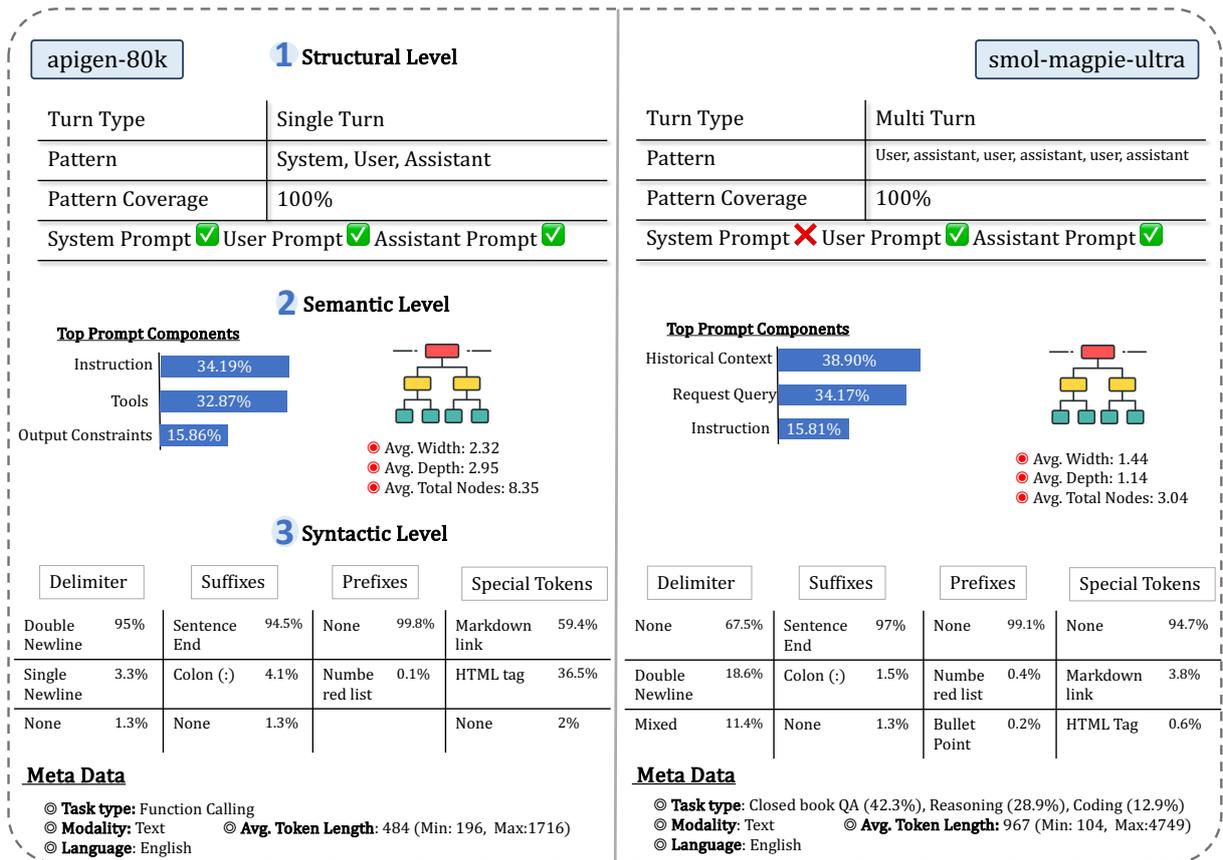


Figure 4: Dataset Profile of APIGEN-80K (Liu et al., 2024b) and SMOL-MAGPIE-ULTRA (Allal et al., 2025) on four primary dimensions: (1) Structural Level (2) Semantic Level (3) Syntactic level and (4) Metadata.

complex scenarios³.

Results Figure 4 presents a comprehensive profile analysis of the examined datasets across multiple analytical dimensions. **Structural Analysis:** APIGEN-80K implements a single-turn architecture with a consistent system-user-assistant sequence, while SMOL-MAGPIE-ULTRA employs a multi-turn structure characterized by alternating user-assistant interactions across three turns. **Semantic Composition:** APIGEN-80K showed a prominence in instruction, tools, and output constraint components, whereas SMOL-MAGPIE-ULTRA emphasizes historical context, request queries, and instructions. Semantic tree structures indicate that APIGEN-80K exhibits greater component diversity, as evidenced by increased tree width and total node count. This structural richness contrasts with the relatively streamlined composition of SMOL-MAGPIE-ULTRA, though it should be noted that our analysis focuses solely on terminal

³For multi-turn data analysis, we focus on the terminal user prompt, treating preceding interactions as historical context. The detailed prompt configuration can be found in Appendix G; We present human validation analysis in Appendix C.

user prompts in multi-turn interactions, potentially affecting these metrics. **Syntactic Characteristics:** we observed that APIGEN-80K employs consistent formatting conventions, utilizing double newlines as delimiters and incorporating markdown-style special tokens. In contrast, SMOL-MAGPIE-ULTRA demonstrates minimal use of explicit delimiters or special tokens, suggesting a more monotone syntactic structure. **Metadata** The APIGEN-80K consists of monolithic task type, function calling, while the SMOL-MAGPIE-ULTRA composes of diverse task such as closed book QA, Reasoning, and Coding etc.

4.3 Prompt Sensitivity Analysis

We leverage PROMTPRISM to establish a framework for assessing prompt sensitivity. Our approach introduces two classes of operators: semantic operators (permute, add, delete) and syntactic operators (format and delimiter modifications), enabling systematic evaluation of model responses to prompt variations.

Experiment Setup Data: We conducted

Ordering Sensitivity Analysis										
Model	Baseline	Instruction			Request/Question			Contextual Reference		
		First	Middle	Last	First	Middle	Last	First	Middle	Last
Claude-Sonnet-3.5***	56.49 (0.03)	55.77 (0.04)	48.99 (0.04)	63.37 (0.09)	24.74 (0.02)	13.46 (0.01)	55.97 (0.10)	41.31 (0.06)	53.83 (0.06)	52.39 (0.06)
		-1%	-13%	12%	-56%	-76%	-1%	-27%	-5%	-7%
Llama3.2-3b-inst***	47.21 (0.08)	46.54 (0.06)	45.97 (0.07)	49.54 (0.15)	35.86 (0.04)	31.73 (0.03)	58.58 (0.11)	51.97 (0.09)	48.49 (0.06)	51.33 (0.11)
		-1%	-3%	5%	-24%	-33%	24%	10%	3%	9%
Mixtral 8x7b inst ***	67.82 (0.06)	64.12 (0.06)	45.13 (0.03)	52.23 (0.10)	56.87 (0.03)	73.63 (0.12)	62.10 (0.11)	69.63 (0.25)	64.07 (0.07)	67.90 (0.13)
		-5%	-33%	-23%	-16%	9%	-8%	3%	-6%	0%
Deepseek r1	67.46 (0.14)	67.62 (0.09)	54.10 (0.10)	61.16 (0.16)	62.35 (0.06)	59.90 (0.05)	63.12 (0.20)	58.20 (0.11)	65.18 (0.16)	66.40 (0.06)
		0%	-20%	-9%	-8%	-11%	-6%	-14%	-3%	-2%

Table 2: Semantic Component Order Sensitivity Analysis. Performance variations under different component orderings relative to baseline configuration. Asterisks (***) denote statistically significant differences (ANOVA, $p < 0.05$). Numbers in parentheses represent standard deviations, and percentages indicate relative performance changes from baseline.

controlled experiments using taxonomy-guided prompts generated in Section 4.1. We selected task `Task067:abductivenli` answer generation as our experimental focus, chosen for its below-average performance score (56.49) relative to the mean text generation task score (73.95). All other experimental parameters remain consistent with those described in Section 4.1. **Model:** We conduct experiments with model Claude-Sonnet-3.5 and Llama3.2-3b-inst.

Semantic operation: We conduct a systematic investigation of semantic component reordering sensitivity. While the framework supports arbitrary component intervention and positioning, we focus our analysis on three primary semantic components (Instruction, Question, and Few-shot examples). These components are systematically repositioned at the beginning, middle, and end of the prompt sequence to assess performance variations. The detailed implementation methodology is illustrated in Figure 7.

Syntactic operation: Our syntactic analysis focuses on delimiter modification operations. We examine four distinct delimiter patterns: double newline (`'\n\n'`), section separator (`'\n#####\n'`), tab (`'\t'`), and extended whitespace (`'\s+'`). Implementation details are provided in Figure 8.

Results Our analysis of semantic component ordering sensitivity, presented in Table 2, reveals significant effects across multiple mod-

els: Claude-Sonnet-3.5, Llama3.2-3b-inst, and Mixtral 8x7 inst demonstrate statistically significant performance variations in response to component reordering. Notably, positioning the Instruction component at the terminal position within the prompt sequence yields the most substantial performance improvements: a 12% enhancement for Claude-Sonnet-3.5 and a 5% improvement for Llama3.2-3b-inst. Interestingly, the Deepseek-r1 model showed more robustness to component order variations, with no statistically significant differences observed.

In contrast, our investigation of delimiter sensitivity (Table 10) indicates that although there are differences in model performances, modifications to delimiter patterns do not produce statistically significant performance variations. This suggests that while models exhibit marked sensitivity to semantic ordering, they maintain robustness to syntactic delimiter variations.

5 Conclusion

In this paper, we present  PROMPTPRISM, a comprehensive taxonomy framework for prompt analysis grounded in linguistic principles. Our framework provides a systematic approach to understanding and analyzing prompt composition, structure, and effectiveness across various applications. This work contributes toward standardizing prompt analysis in the field of LLMs.

Limitations

While PROMPTPRISM provides a systematic framework for prompt analysis, several limitations warrant acknowledgment. First, our approach employs a top-down methodology in establishing the taxonomy. Although grounded in linguistic principles, this approach may not capture all emergent patterns that could be identified through bottom-up analysis of existing prompt collections. We acknowledge that no single taxonomic framework can claim to be definitive or completely comprehensive in the rapidly evolving field of prompt engineering.

Furthermore, the current implementation lacks extensive manual annotation of existing large-scale datasets. While our framework demonstrates effectiveness in automated analysis, and present pilot human validation in Section C, the absence of human-validated annotations across diverse prompt collections limits our ability to fully validate the taxonomy’s generalizability and robustness across different domains and applications.

These limitations suggest valuable directions for future work, including: Integration of bottom-up analysis methods to complement our top-down approach, Development of comprehensive manually annotated benchmark datasets

Broader Impact

The development and application of this work carry significant implications for the field of natural language processing and artificial intelligence at large. This work adheres to the ACL Code of Ethics and is primarily aimed at enhancing the transparency and interpretability of language models.

By providing a systematic framework for analyzing and understanding prompts, PROMPTPRISM has the potential to:

Enhancing Transparency and Reproducibility PROMPTPRISM addresses a critical challenge in current LLM research: the inconsistent documentation of prompts in publications, which hinders reproducibility and comparative analysis. Just as Model Cards (Mitchell et al., 2019) standardized the documentation of model characteristics, PROMPTPRISM provides a structured framework for documenting prompt characteristics. This standardization enables more precise reporting of prompt components in research papers, facilitating more meaningful comparisons across studies and enhancing the reproducibility of results.

Practical Applications Beyond Research Beyond academic research, PROMPTPRISM offers practical value for organizations developing LLM applications. The framework can help standardize prompt libraries and documentation within organizations, facilitating knowledge transfer between prompt engineers and enabling more systematic prompt management. This standardization is particularly valuable as organizations scale their LLM applications and need to maintain consistency across multiple systems and use cases.

While the immediate risks associated with this research are low, we acknowledge that improved prompt engineering techniques could potentially be misused to create more persuasive or manipulative AI generated content. However, we believe that the benefits of increased transparency and interpretability outweigh these potential risks.

References

- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, et al. 2025. Smollm2: When smol goes big—data-centric training of a small language model. *arXiv preprint arXiv:2502.02737*.
- Amazon. 2024. [Nova lite](#). Amazon Nova Foundation Models.
- Anthropic. 2024a. [Claude 3.5 haiku](#). Compact version of Claude language model.
- Anthropic. 2024b. [Claude 3.5 sonnet](#). Large language model.
- Anthropic. 2025. [Claude 3.7 sonnet](#). AI model with hybrid reasoning capabilities.
- Mark Aronoff. 1976. Word formation in generative grammar. *Linguistic Inquiry Monographs Cambridge, Mass*, (1):1–134.
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. 2021. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*.
- John L Austin. 1962. How to do things with words: Lecture i. *How to do things with words: JL Austin*, pages 1–11.
- Stephen Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Févry, et al. 2022. Promptsource: An integrated development environment and repository for natural

- language prompts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 93–104.
- Marcel Binz and Eric Schulz. 2023. [Using cognitive psychology to understand gpt-3](#). *Proceedings of the National Academy of Sciences*, 120(6):e2218523120.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Devichand Budagam, KJ Sankalp, Ashutosh Kumar, Vinija Jain, and Aman Chadha. 2024. Hierarchical prompting taxonomy: A universal evaluation framework for large language models. *CoRR*.
- Robyn Carston. 2008. *Thoughts and utterances: The pragmatics of explicit communication*. John Wiley & Sons.
- Kewei Cheng, Nesreen K Ahmed, Theodore Willke, and Yizhou Sun. 2024. Structure guided prompt: Instructing large language model in multi-step reasoning by exploring graph structure of the text. *arXiv preprint arXiv:2402.13415*.
- Nathan A Chi, Teodor Malchev, Riley Kong, Ryan A Chi, Lucas Huang, Ethan A Chi, R Thomas McCoy, and Dragomir Radev. 2024. Modeling: A novel dataset for testing linguistic reasoning in language models. *arXiv preprint arXiv:2406.17038*.
- Noam Chomsky. 2014. *Aspects of the Theory of Syntax*. 11. MIT press.
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric P Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing discrete text prompts with reinforcement learning. *arXiv preprint arXiv:2205.12548*.
- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. 2023. Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv preprint arXiv:2309.16797*.
- Charles J Fillmore. 1967. The case for case.
- Herbert P Grice. 1975. Logic and conversation. In *Speech acts*, pages 41–58. Brill.
- Barbara J Grosz and Candace L Sidner. 1986. Attention, intentions, and the structure of discourse. *Computational linguistics*, 12(3):175–204.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Jia He, Mukund Rungta, David Koleczek, Arshdeep Sekhon, Franklin X Wang, and Sadid Hasan. 2024. Does prompt formatting have any impact on llm performance? *arXiv preprint arXiv:2411.10541*.
- Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, et al. 2025. Mathperturb: Benchmarking llms’ math reasoning abilities against hard perturbations. *arXiv preprint arXiv:2502.06453*.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, et al. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*.
- Stephen C Levinson. 1983. *Pragmatics*. Cambridge university press.
- Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. 2023a. Large language models understand and can be enhanced by emotional stimuli. *arXiv preprint arXiv:2307.11760*.
- Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. 2023b. [Large language models understand and can be enhanced by emotional stimuli](#). *Preprint*, arXiv:2307.11760.
- Moxin Li, Wenjie Wang, Fuli Feng, Fengbin Zhu, Qifan Wang, and Tat-Seng Chua. 2024. Think twice before trusting: Self-detection for large language models through comprehensive answer reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 11858–11875.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024a. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, Rithesh RN, et al. 2024b. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets. *Advances in Neural Information Processing Systems*, 37:54463–54482.
- William C Mann, CMIM Matthiessen, and Sandra A Thompson. 1992. Rhetorical structure theory and text analysis. *Discourse description: Diverse linguistic analyses of a fund-raising text*, pages 39–78.
- Yuetian Mao, Junjie He, and Chunyang Chen. 2025. From prompts to templates: A systematic prompt template analysis for real-world llmapps. In *Proceedings of the 33rd ACM International Conference on the Foundations of Software Engineering*, pages 75–86.

- Peter Hogue Matthews. 1972. *Inflectional morphology: A theoretical study based on aspects of Latin verb conjugation*, volume 6. CUP Archive.
- Meta AI. 2024a. [Llama 3.2 3b instruct](#). 3 billion parameter instruction-tuned language model.
- Meta AI. 2024b. [Llama 3.3 70b instruct](#). 70 billion parameter instruction-tuned language model.
- Mistral AI. 2023a. [Mistral 7b instruct](#). 7 billion parameter instruction-tuned language model.
- Mistral AI. 2023b. [Mixtral 8x7b instruct](#). Mixture of Experts language model.
- Mistral AI. 2024. [Mistral large](#). Large-scale language model.
- Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229.
- OpenAI. [Open ai](#).
- Juri Opitz, Shira Wein, and Nathan Schneider. 2025. Natural language processing relies on linguistics. *Computational Linguistics*, pages 1–24.
- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel R Bowman. 2021. Bbq: A hand-built bias benchmark for question answering. *arXiv preprint arXiv:2110.08193*.
- Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Latham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2023. [Discovering language model behaviors with model-written evaluations](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13387–13434, Toronto, Canada. Association for Computational Linguistics.
- Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. Xcopa: A multilingual dataset for causal common-sense reasoning. *arXiv preprint arXiv:2005.00333*.
- Abhilasha Ravichander, Matt Gardner, and Ana Marasović. 2022. Condaqa: A contrastive reading comprehension dataset for reasoning about negation. *arXiv preprint arXiv:2211.00295*.
- Shubhra Kanti Karmaker Santu and Dongji Feng. 2023. Teler: A general taxonomy of llm prompts for benchmarking complex tasks. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14197–14203.
- Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt search: A structure-aware approach to efficient compile-time prompt optimization. *arXiv preprint arXiv:2404.02319*.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, et al. 2024. The prompt report: A systematic survey of prompting techniques. *CoRR*.
- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*.
- John R Searle. 1969. *Speech acts: An essay in the philosophy of language*. Cambridge university press.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652.
- Dan Sperber and Deirdre Wilson. 1986. *Relevance: Communication and cognition*, volume 142. Harvard University Press Cambridge, MA.
- Lucien Tesnière. 2015. *Elements of structural syntax*. John Benjamins Publishing Company.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Ming Wang, Yuanzhong Liu, Xiaoyu Liang, Songlian Li, Yijie Huang, Xiaoming Zhang, Sijia Shen, Chaofeng Guan, Daling Wang, Shi Feng, et al. 2024. Langgpt: Rethinking structured reusable prompt design framework for llms from the programming language. *arXiv preprint arXiv:2402.16929*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva

Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. 2022. Supernaturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. 2023. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382*.

Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. 2024. Textgrad: Automatic "differentiation" via text. *arXiv preprint arXiv:2406.07496*.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2).

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*.

Appendix

A Prompt Refinement Experiment Setting

Data We focused on a diverse set of 70 tasks spanning two primary categories: 53 classification tasks and 17 text generation tasks, following (Sclar et al., 2023).

17 tasks selected for text generation: task037, task038, task040, task067, task071, task072, task105, task216, task223, task240, task348, task389, task443, task845, task1326, task1401, task1613.

53 Tasks used for classification: task050, task065, task069, task070, task114, task133, task155, task158, task161, task162, task163, task190, task213, task214, task220, task279, task280, task286, task296, task297, task316, task317, task319, task320, task322, task323, task325, task326, task327, task328, task335, task337, task385, task580, task607, task608, task609, task904, task905, task1186, task1283, task1284, task1297, task1347, task1387, task1419, task1420, task1421, task1423, task1502, task1612, task1678, task1724.

Models: The models for evaluation include Claude 3.7 Sonnet (Anthropic, 2025), Claude 3.5 Sonnet, Haiku (Anthropic, 2024b,a), Llama3.2-3b-inst (Meta AI, 2024a), Llama3.3-70b-inst (Meta AI, 2024b), Mistral-7b-inst (Mistral AI, 2023a), Mistral Large (Mistral AI, 2024), Mixtral8x7b-inst (Mistral AI, 2023b), Nova Lite (Amazon, 2024), DeepSeek-r1 (Guo et al., 2025). For inference, we set the temperature=0, for reproducibility.

B Syntactic Definition

Table 4 and Figure 5 shows the detailed definition of syntactic components and implementation details. Figure 6 shows snippet of code analyzing delimiters.

C Annotation Validation Analysis

To validate the LLM-based annotation described in Section 4.2, we conducted an evaluation using two independent human verifiers across three key dimensions:

Format Correctness evaluates the structural integrity of XML-style annotations, specifically the proper matching of <tag>..<\tag> pairs. We

Dataset	Format correctness	Tag correctness	Coverage
apigen-80k	0.99 (1.0)	1.0 (1.0)	0.98 (0.98)
smol-magpie-ultra	1.0 (1.0)	0.97 (0.79)	0.92 (0.55)

Table 3: Human annotator agreement analysis. Mean annotation scores (0-1 scale) across criteria. Values in parentheses show Pearson correlation coefficients (ρ) between annotators, where 1.0 indicates perfect agreement.

quantify correctness as the ratio of properly formatted tags to total tags, normalized to a scale of 0 to 1, where 1 indicates perfect format conformity.

Tag Correctness assesses the semantic accuracy of content classification within tags. This metric evaluates whether content is appropriately categorized (e.g., whether text marked as ‘instruction’ truly represents instructional content). Scores are normalized on 0-1 scale, with 1 representing complete semantic accuracy.

Coverage measures the comprehensiveness of content annotation, evaluating whether all relevant content is appropriately tagged. For instance, if a multi-sentence instruction is only partially tagged, the coverage score is reduced accordingly. Each component receives a binary score (0 or 1), based on complete coverage of semantic units.

Analysis of the evaluation results, presented in Table 3, demonstrates strong performance across all three assessment metrics—format correctness, tag correctness, and coverage—for both datasets. The metrics exhibit high mean values (exceeding 0.95), indicating robust annotation quality. Inter-annotator agreement shows strong consistency, with correlation coefficients ranging from 0.55 to 1.0, suggesting reliable assessment methodology. These results indicate that our LLM-based annotation approach achieves high accuracy and reproducibility across different annotators.

D Related Work: Structured Prompts

Recent research has made significant strides in understanding prompts as structured entities. Several approaches have emerged to analyze and categorize prompt components, from basic decomposition into instructions, examples, and queries (Fernando et al., 2023), to more comprehensive analyses of real-world industrial prompt patterns (Mao et al., 2025). Some frameworks have proposed viewing prompts as programmable structures (Schnabel and Neville, 2024; Wang et al., 2024), drawing parallels with software development patterns and structured

coding approaches. Others have explored structured prompts enhanced with graph-based representations to improve reasoning capabilities (Cheng et al., 2024), while some work has focused on identifying recurring prompt patterns analogous to software design patterns (White et al., 2023).

While these efforts contribute valuable insights into prompt structuring, they often lack comprehensive coverage of prompt characteristics. For example, frameworks such as (Fernando et al., 2023; Wang et al., 2024) focus primarily on content-level analysis, overlooking the crucial aspects of stylistic formatting and syntactic variations. Our work addresses this gap by introducing a linguistically-inspired taxonomy that integrates multiple dimensions of prompt analysis. PROMPTPRISM provides a comprehensive framework that encompasses stylistic formatting, syntactic patterns, and structural organization, enabling a more comprehensive understanding of how different prompt characteristics influence model behavior.

E Prompt Sensitivity Analysis

Fig 7 and Fig 8, respectively shows the snippet of code, running controlled experiment on PROMPTPRISM framework. The framework enables to conveniently run the semantic level intervention and syntactic modification experiment.

F Latency and Cost Analysis

Our proposed PromptPrism approach may introduce increased token count and computational costs. In Table 11 we present a comprehensive analysis comparing PromptPrism against other approaches, and the results show while there are trade-offs between prompt length and performance gain, our approach show higher performance gain efficiency.

Syntactic Components	Description	Examples
Component Indexing	Maps semantic component to role and order	(User, 0) # User, order index
Component Span	Spans of a component	(start_position, end_position) # absolute values
Delimiters	Boundary between components	Double newline ("\n\n"), Single newline ("\n"), Tab ("\t"), Whitespace ("\s+")
Directive Markers	Special tokens or pattens that signal component structure	
Prefixes	Markers at the beginning of components	hash comments (#) , double-slash comments (//), block-quotes (>), numbered lists (1.,2.,etc), and bullet points (-,*,+)
Suffixes	Markers at the end of components	Colon(:), Sentence End (!?), Semicolon (;)
Special token	Special tokens specific to models	Llama3.1:< begin_of_text >, < end_of_text >, < eom_id > Mistral-7b-inst: <s>..</s>, [INST]..[/INST]

Table 4: Syntactic Components Definition

```

1  # prefix patterns
2  prefix_patterns = [
3      (r"\s*#[^\n]+", "hash_comment"),
4      (r"\s*/[/^\n]+", "double_slash_comment"),
5      (r"\s*>[^\n]+", "blockquote"),
6      (r"\s*\d+\.\s", "numbered_list"),
7      (r"\s*[-*+]\s", "bullet_point"),
8  ]
9
10 # suffix patterns
11 suffix_patterns = [
12     (r"\s*:\s*$", "colon_end"),
13     (r"\s*[.!?]+$", "sentence_end"),
14     (r"\s*[\;]\s*$", "semicolon_end"),
15 ]
16
17 # Special token patterns
18 special_token_patterns = [
19     (r"<[^\>]+>", "html_tag"),
20     (r"\.[*?]", "markdown_link"),
21     (r"\$.*?\$", "math_expression"),
22     (r"@w+", "mention"),
23     (r"#w+", "hashtag"),
24     (r"https?://\S+", "url"),
25 ]

```

Figure 5: Syntactic Components Implementation

```

def _analyze_delimiter(self, delimiter_text):
    """Analyze the delimiter between components"""
    if delimiter_text is None:
        return None

    delimiter_info = {"raw": delimiter_text, "length": len(delimiter_text), "type": None,
↪ "pattern": None}

    # Analyze delimiter pattern
    if delimiter_text.isspace():
        if "\n\n" in delimiter_text:
            delimiter_info["type"] = "double_newline"
            delimiter_info["pattern"] = r"\n\n"
        elif "\n" in delimiter_text:
            delimiter_info["type"] = "single_newline"
            delimiter_info["pattern"] = r"\n"
        elif "\t" in delimiter_text:
            delimiter_info["type"] = "tab"
            delimiter_info["pattern"] = r"\t"
        else:
            delimiter_info["type"] = "whitespace"
            delimiter_info["pattern"] = r"\s+"
    else:
        delimiter_info["type"] = "mixed"
        delimiter_info["pattern"] = repr(delimiter_text)[1:-1]
    return delimiter_info

```

Figure 6: Syntactic-Delimiter analyzing implementation

Providers	System	User	Assistant	Tools
Bedrock Converse API	✓ (Implicit, i.e., 'System' is not stated in message) prompt that provides instructions or context to the model about the task it should perform, or the persona it should adopt during the conversation.	✓ The human that is sending messages to the model.	✓ The model that is sending messages back to the human user.	✓ The definition of the tool is a JSON schema that you pass in the toolConfig (ToolConfiguration) request parameter to the Converse operation
OpenAI	✓▶ (Developer) Instructions to the model that are prioritized ahead of user messages, following chain of command. Previously called the system prompt.	✓ Input from end users, or a catch-all for data we want to provide to the model	✓ Sampled from the language model	✓ Generated by some program, such as code execution or an API call
Anthropic Claude	✓ (Implicit, i.e., 'System' is not stated in message) A system prompt is a way of providing context and instructions to Claude, such as specifying a particular goal or role.	✓	✓	✓
Llama	✓ Sets the context in which to interact with the AI model. It typically includes rules, guidelines, or necessary information that help the model respond effectively.	✓ Represents the human interacting with the model. It includes the inputs, commands, and questions to the model.	✓ Represents the response generated by the AI model based on the context provided in the system, ipython and user-prompts.	✓▶ (ipython) A new role introduced in Llama 3.1. Semantically, this role means "tool". This role is used to mark messages with the output of a tool call when sent back to the model from the executor.
Microsoft Azure	✓▶ (Developer) for GPT-o family A brief description of the assistant. Personality traits of the assistant. Instructions or rules you want the assistant to follow. Data or information needed for the model, such as relevant questions from an FAQ.	✓ End user message	✓ Assistant is a large language model trained by OpenAI.	✓ A list of tools the model may call. Currently, only functions are supported as a tool. Use this to provide a list of functions the model may generate JSON inputs for. A max of 128 functions are supported.
Mistral API	✓	✓	✓	✓
Gemini API	✓▶ (systemInstruction) Optional. Developer set system instruction(s). Currently, text only.	✓ Available in chat mode	✓▶ (model) Available in chat mode	✓ Optional. A list of Tools the Model may use to generate the next response. A Tool is a piece of code that enables the system to interact with external systems to perform an action, or set of actions, outside of knowledge and scope of the Model. Supported Tools are Function and codeExecution.

Table 5: Structural Roles Across LLM Platforms. Comparison of role definitions across major LLM service providers. While platforms generally support four standard roles (System, User, Assistant, Tools), some employ platform-specific nomenclature. Our framework adopts these standardized role definitions to ensure broad applicability and consistency.

Number of Few Shots: 2						
Model	Text Generation Tasks			Classification Tasks		
	DEFAULT	CoT	TAXONOMY	DEFAULT	CoT	TAXONOMY
Claude Sonnet 3.7	27.28 (0.2)	44.6 (0.22)	57.35 (0.26) + 29.0 %	20.92 (0.38)	98.55 (0.11)	98.68 (0.10) + 0.13 %
Claude Sonnet 3.5	17.84 (0.09)	45.32 (0.20)	57.55 (0.23) + 27 %	19.7 (0.36)	96.5 (0.18)	98.68 (0.10) + 2.2 %
Claude Haiku 3.5	19.42 (0.17)	44.14 (0.21)	53.73 (0.26) + 21.7 %	5.84 (0.17)	97.11 (0.16)	99.06 (0.07) + 2.01 %
Llama3.2 3b inst	19.26 (0.16)	20.1 (0.20)	47.76 (0.24) + 137.6 %	14.97 (0.30)	52.96 (0.48)	82.35 (0.35) + 55.5 %
Llama3.3 70b inst	31.06 (0.22)	49.8 (0.25)	61.72 (0.26) + 23.9 %	12.29 (0.28)	98.11 (0.13)	97.39 (0.15) - 0.73 %
Mistral 7b inst	31.68 (0.18)	28.78 (0.21)	47.16 (0.17) + 68.9 %	7.79 (0.05)	40.84 (0.44)	53.55 (0.43) + 31.1 %
Mistral Large	35.28 (0.18)	52.41 (0.23)	63.17 (0.23) + 22.4 %	10.98 (0.22)	97.48 (0.13)	98.15 (0.13) + 0.68 %
Mixtral 8x7b inst	33.83 (0.17)	35.88 (0.21)	51.14 (0.19) + 42.5 %	9.00 (0.13)	67.37 (0.44)	80.54 (0.37) + 19.5 %
Nova Lite	21.34 (0.08)	44.47 (0.23)	61.89 (0.21) + 39.2 %	6.76 (0.03)	71.27 (0.34)	93.96 (0.20) + 31.8 %
Deepseek r1	21.25 (0.11)	39.82 (0.28)	56.67 (0.26) + 42.3 %	23.57 (0.36)	99.12 (0.09)	98.5 (0.09) - 0.63 %

Table 6: Taxonomy Guided Prompt Refinement Results in 2-shot setting. Performance improvements over Chain-of-Thought (CoT) baseline using our taxonomy-guided approach. Text generation tasks show an average 29% improvement across models, while classification tasks demonstrate a 0.13% enhancement. Percentages indicate relative performance gains compared to CoT and the numbers in parenthesis indicate standard deviation.

Number of Few Shots: 0						
Model	Text Generation Tasks			Classification Tasks		
	DEFAULT	CoT	TAXONOMY	DEFAULT	CoT	TAXONOMY
Claude Sonnet 3.7	37.62 (0.26)	38.75 (0.25)	51.16 (0.25) +36%	96.48 (0.15)	96.73 (0.15)	99.02 (0.07) +3%
Claude Sonnet 3.5	14.9 (0.06)	41.29 (0.23)	50.58 (0.20) +239%	11.44 (0.26)	94.81 (0.19)	98.86 (0.08) +764%
Claude Haiku 3.5	15.83 (0.08)	36.22 (0.23)	45.92 (0.21) +190%	5.87 (0.16)	95.05 (0.06)	99.05 (0.07) +1587%
Llama3.2 3b inst	29.52 (0.24)	20.44 (0.19)	38.57 (0.17) 31%	64.49 (0.46)	58.29 (0.47)	81.29 (0.17) 26%
Llama3.3 70b inst	33.11 (0.26)	36.85 (0.25)	49.94 (0.22) +51%	39.59 (0.45)	94.52 (0.22)	98.02 (0.14) +148%
Mistral 7b inst	32.21 (0.17)	23.08 (0.23)	44.16 (0.19) +37%	5.98 (0.05)	52.93 (0.45)	53.08 (0.44) +788%
Mistral Large	36.65 (0.23)	44.06 (0.24)	55.22 (0.23) +51%	41.01 (0.45)	92.49 (0.23)	98.79 (0.09) +141%
Mixtral 8x7b inst	28.88 (0.15)	32.55 (0.20)	45.05 (0.18) +56%	18.14 (0.31)	68.04 (0.41)	78.77 (0.37) +334%
Nova Lite	26.92 (0.17)	36.04 (0.21)	54.29 (0.21) +102%	12.49 (0.27)	59.48 (0.40)	94.97 (0.17) +660%
Deepseek r1	11.28 (0.17)	40.01 (0.25)	47.91 (0.24) +325%	37.61 (0.43)	96.61 (0.18)	99.05 (0.06) +163%

Table 7: Taxonomy Guided Prompt Refinement Results in 0-shot setting. Performance improvements over Chain-of-Thought (CoT) baseline using our taxonomy-guided approach. Text generation tasks show an average 112% improvement across models, while classification tasks demonstrate a 461% enhancement. Percentages indicate relative performance gains compared to CoT and the numbers in parenthesis indicate standard deviation.

Text Generation Tasks				Classification Tasks			
Default	CoT	OpenAI	Taxonomy	Default	CoT	OpenAI	Taxonomy
14.9	41.29	42.28	50.58	11.44	94.81	86.79	98.86
(0.06)	(0.23)	(0.245)	(0.20)	(0.26)	(0.19)	(0.34)	(0.08)

Table 8: Additional (OpenAI) baseline. Using Claude Sonnet 3.5 in zero-shot settings (Tab 7), we evaluated performance using prompts rewritten according to OpenAI’s guidelines. The results demonstrate that while the OpenAI meta prompt approach outperforms the Chain of Thought baseline, our taxonomy-based approach achieves higher performance in text generation tasks.

Dataset	Delimiter			Suffixes			Prefixes			Special Tokens		
apigen-80k	Double newline (“\n\n”) 0.95	Single newline (“\n”) 0.033	None 0.013	Full stop (.) 0.945	Colon (:) 0.041	None 0.013	None 0.998	Numbered list 0.001	Hash comment 0	Markdown link 0.594	Html tag 0.365	None 0.02
everday-conversations	None 0.983	Whitespace 0.017	mixed 0	Full stop (.) 1			None 1			None 1		
explore-instruct-rewriting	Double newline (“\n\n”) 0.655	Whitespace 0.216	Single newline (“\n”) 0.126	Full stop (.) 0.617	Colon (:) 0.382	None 0	None 0.995	Bullet point 0.005	Hash comment 0	None 0.993	Markdown link 0.002	Html tag 0.001
longalign	None 0.69	Double newline 0.225	mixed 0.08	None 0.68	Full stop (.) 0.27	Colon (:) 0.04	None 0.92	Numbered list 0.04	Hash comment 0.02	None 0.84	Markdown link 0.077	Html tag 0.035
metamathqa-50k	None 0.643	Double newline 0.222	Whitespace 0.111	Full stop (.) 0.952	None 0.042	Colon (:) 0.004	None 0.998	Numbered list 0.001	Hash comment 0	None 0.964	Markdown link 0.036	Mention 0
numia-cot-100k	None 0.604	Double newline 0.244	Whitespace 0.091	Full stop (.) 0.799	None 0.169	Colon (:) 0.028	None 0.989	Bullet point 0.006	Numbered list 0.004	None 0.904	Markdown link 0.08	Html tag 0.014
openhermes-100k	Double newline (“\n\n”) 0.423	None 0.313	Whitespace 0.215	Full stop (.) 0.89	Colon (:) 0.077	None 0.027	None 0.99	Numbered list 0.005	Bullet point 0.003	None 0.881	Markdown link 0.072	Html tag 0.026
self-oss-instruct	Double newline (“\n\n”) 0.557	Whitespace 0.407	None 0.017	Full stop (.) 0.942	Colon (:) 0.053	None 0.003	None 0.969	Bullet point 0.025	Numbered list 0.005	None 0.724	Markdown link 0.252	Html tag 0.014
smol-constraints	Whitespace 0.581	Double newline 0.289	mixed 0.122	Full stop (.) 0.983	Colon (:) 0.014	None 0.002	None 0.966	Bullet point 0.033	Hash comment 0.0001	Markdown link 0.438	None 0.43	Html tag 0.1302
smol-magpie-ultra	None 0.675	Double newline 0.186	Whitespace 0.114	Full stop (.) 0.97	None 0.015	Colon (:) 0.013	None 0.991	Numbered list 0.004	Bullet point 0.002	None 0.947	Markdown link 0.038	Html tag 0.006
smol-rewrite	Double newline (“\n\n”) 0.556	Whitespace 0.436	Single newline (“\n”) 0.006	Full stop (.) 1			None 1			None 0.68	Hash tag 0.255	Markdown link 0.031
smol-summarize	Double newline (“\n\n”) 0.697	Whitespace 0.245	mixed 0.055	Full stop (.) 0.914	None 0.085	Colon (:) 0	None 1			None 0.961	Markdown link 0.021	Mention 0.006
systemchat-30k	Double newline (“\n\n”) 0.648	Single newline (“\n”) 0.273	Whitespace 0.072	Full stop (.) 0.995	None 0.003	Colon (:) 0.0009	None 1			None 0.993	Html tag 0.002	Mention 0.002

Table 9: Syntactic Profile. Comprehensive syntactic information of entire SMOL-TALK dataset (Allal et al., 2025), presented in Section 4.2

Model	“\n\n” (double new line)	“\n#####\n” (number #)	“\n” (single new line)	“\t” (tab)	“\s+” (whitespace)
Claude-Sonnet-3.5	56.49 (0.03)	51.34 (0.02)	52.25 (0.02)	57.07 (0.07)	54.5 (0.10)
Llama3.2-3b-inst	47.21 (0.08)	43.76 (0.13)	43.5 (0.10)	33.54 (0.03)	39.77 (0.13)

Table 10: Delimiter Sensitivity Analysis. Numbers in parentheses represent standard deviations, and percentages indicate relative performance changes from baseline. In contrast to Semantic components re-ordering analysis, delimiter sensitivity do not show statistically significant difference between variations.

	Prompt Length (TOTAL)	Prompt Length Δ (baseline)%	Prompt Length (INPUT)	Prompt Length (OUTPUT)	Performance	Performance Δ (baseline)%
Baseline	1078.99		535.67	543.32	14.90	
Baseline Fewshot 2	2433.97	125.58	2021.90	412.07	45.32	204.16
CoT	718.54	-33.41	688.68	29.85	41.29	177.11
OpenAI metaprompt (OpenAI)	2643.19	144.97	2154.58	488.61	42.28	183.76
Ours	1492.21	38.30	1429.68	62.53	50.58	239.46

Table 11: Latency and Cost Analysis: PromptPrism increases total prompt length by 38.30% over the baseline, which is more modest than competing methods: the OpenAI metaprompt approach (OpenAI), 144.97%, and the few-shot baseline, 125.58%. Critically, it achieves the highest performance score of 50.58, representing a 239.46% improvement over baseline, delivering 6.25 pp of performance improvement per percentage point of length increase.

```

def reorder_component(self, component_name, position="first", remove_tags=False):
    """
    Reorder a specific component in the prompt.
    Args:
        component_name (str): Name of the component to reorder
        position (str): Where to place the component ('first', 'last', 'middle')

    Returns:
        str: Reordered prompt text
    """
    try:
        # Validate component name
        valid_components = ["request_query", "contextual_ref", "instruction", "output_const",
        ↪ "response", "other"]
        if component_name not in valid_components:
            raise ValueError(f"Invalid component name: {component_name}. Valid options are:
            ↪ {valid_components}")

        # Validate position
        validate_positions = ["first", "last", "middle"]
        if position not in validate_positions:
            raise ValueError(f"Invalid position: {position}. Valid options are:
            ↪ {validate_positions}")

        # find all tags that belong to the component category
        related_tags = [comp["index"] for comp in self.components if
        ↪ comp["tag"].startswith(f"{component_name}")]

        if not related_tags:
            raise ValueError(f"No components found for category: {component_name}")

        # remove related tags from current order
        other_tags = [comp["index"] for comp in self.components if not
        ↪ comp["tag"].startswith(f"{component_name}")]

        # Create new order based on position
        if position == "first":
            new_order = related_tags + other_tags
        elif position == "last":
            new_order = other_tags + related_tags
        else: # middle
            mid_point = len(other_tags) // 2
            new_order = other_tags[:mid_point] + related_tags + other_tags[mid_point:]

        # update tag order
        self._tag_order = new_order
        # Use existing to_text()
        return self.to_text(remove_tags=remove_tags)

    except Exception as e:
        print(f"Error reordering component: {str(e)}")
        return self.to_text(remove_tags=remove_tags)

```

Figure 7: Prompt sensitivity analysis: Semantic components re-ordering experiment implementation.

```

def modify_delimiter(self, new_delimiter, position="all", remove_tags=False):
    """
    Modify delimiters between components based on specified position.
    Args:
        new_delimiter (str): The new delimiter to insert
        position (str): Where to modify delimiters ('all', 'first', 'last','middle')

    Returns:
        str: Updated prompt text with modified delimiters
    """
    if position not in ["all", "first", "last", "middle"]:
        raise ValueError(f"Invalid position: {position}. Valid options are: 'all', 'first',
        ↪ 'last', 'middle'")

    # Get number of components
    num_components = len(self.components)
    if num_components <= 1:
        return self.to_text() # No delimiters to modify

    # Determine which indices to modify
    if position == "all":
        indices_to_modify = range(num_components - 1)
    elif position == "first":
        indices_to_modify = [0]
    elif position == "last":
        indices_to_modify = [num_components - 2] # Second to last component
    elif position == "middle":
        if num_components >= 3:
            middle_idx = (num_components - 1) // 2
            indices_to_modify = [middle_idx]

    # Modify delimiters for specified indices
    for idx in indices_to_modify:
        self.components[idx]["delimiter_after"] = new_delimiter
        # Update metadata for the delimiter
        # self.components[idx]["metadata"]["delimiter_info"] =
        ↪ self._analyze_delimiter(new_delimiter)

    return self.to_text(remove_tags=remove_tags)

```

Figure 8: Prompt Sensitivity Analysis: Syntactic Delimiter modification experiment implementation

G Prompts Configuration

Prompt for taxonomy guided prompt generation

Your task is to augment the prompt based on the given prompt.
Use these tags to augment the prompts:

1. `<instruction>` `</instruction>`: Prompt related to instructions to the LLM on what to do, guide how to think about or approach the problem.
 - 1-1: `<instruction:task>` `</instruction:task>`: Task related instruction.
E.g., "This is a classification task."
 - 1-2: `<instruction:guideline>` `</instruction:guideline>`: Non-task specific instructions that shape response behavior.
 - 1-2-1: `<instruction:guideline:role>` `</instruction:guideline:role>`: Role assumption.
For example, "act as an experienced data scientist".
 - 1-2-2: `<instruction:guideline:scenario>` `</instruction:guideline:scenario>`:
Scenario description
 - 1-2-3: `<instruction:guideline:behavioral>` `</instruction:guideline:behavioral>`:
Behavioral instructions and interaction style
 - 1-2-4: `<instruction:guideline:emotion>` `</instruction:guideline:emotion>`:
guideline for emotional tone
 - 1-2-5: `<instruction:guideline:cot>` `</instruction:guideline:cot>`:
Chain of Thought guideline. E.g., "Let's think step by step."
 - 1-2-6: `<instruction:guideline:safety>` `</instruction:guideline:safety>`:
guideline related to safety. E.g., "avoid harmful content".
2. `<contextual_ref>` `</contextual_ref>`: Reference data of background information of the user query chat history, or retrieved sections in RAG.
 - 2-1: `<contextual_ref:fewshot>` `</contextual_ref:fewshot>`: sample input-output pairs for in-context learning.
 - 2-2: `<contextual_ref:knowledge_base>` `</contextual_ref:knowledge_base>`:
reference information or facts. E.g., "Based on medical guidelines..."
 - 2-3: `<contextual_ref:context_for_task>` `</contextual_ref:context_for_task>`:
relevant background information.
3. `<request_query>` `</request_query>`: User request of query (e.g. Question)
4. `<output_const>` `</output_const>`: Output constraints, response requirements to the LLM on how to generate response
 - 4-1: `<output_const:label>` `</output_const:label>`: defined set of possible output categories.
E.g., "choose from ['positive', 'negative']"
 - 4-2: `<output_const:wordlimit>` `</output_const:wordlimit>`: restrictions on response length.
E.g., "respond in 50 words or less".
 - 4-3: `<output_const:format>` `</output_const:format>`: output format or structure specification.
E.g., "format as JSON", "present in bullet points".
 - 4-4: `<output_const:style_tone>` `</output_const:style_tone>`: writing style or tone.
E.g., "use academic language".
5. `<other>` `</other>`: Other purpose components
 - 5-1: `<other:adversarial>` `</other:adversarial>`: Adversarial components. E.g., "\&\&\&!!!!!!"
6. `<response>` `</response>`: Response component from LLM
 - 6-1: `<response:answer>` `</response:answer>`: specific answer component.
7. `<tools>` `</tools>`: specifications for tool usage
 - 7-1: `<tools:tool_name>` `</tools:tool_name>`: identifier for specific tools
 - 7-2: `<tools:tool_description>` `</tools:tool_description>`: explanations of tool functionality.
 - 7-3: `<tools:parameters>` `</tools:parameters>`: required inputs and configurations.

Be creative, and augment contents related to the tags.

For example, `<output_const>` Give only the answer `</output_const>`.

Give **ONLY** the augmented version of the prompt. Do **NOT** include any preambles.

Below is the prompt:

```
### Prompt ###
<instruction> {definition} </instruction>

### Examples ###

### Negative Examples ###
<contextual_ref> {negative_examples} </contextual_ref>

### Positive Examples ###
<contextual_ref> {positive_examples} </contextual_ref>
```

Prompt for taxonomy annotation

Your task is to annotate a prompt based on some prompt taxonomy defined below.

Given a prompt, please decompose the prompt into several components by adding the below tags in prompt taxonomy section to the original input prompt.

The available prompt component tags are defined below. Not all prompt component tags should be used for each prompt.

prompt taxonomy ###
{prompt_taxonomy}

additional tips

1. You have to analyze system prompts and user prompts separately using the above prompt taxonomy

2. If you can find the related tags in a subcategory component, use the tag in the subcategory.

If not, use the parent node tag. For example, if a prompt segment belongs to few-shot examples, then use "<contextual_ref:fewshot> </contextual_ref:fewshot>", then insert this tag.

If it belongs to other contextual information that are not directly available in the options, then use "<contextual_ref> </contextual_ref>" instead.

3. Please only added tags to the original input prompt. Please do ****NOT**** include any preambles or additional explanations.

4. If there already exists some tags in the original input prompt such as <tools></tools>, think about if the tags align with our current taxonomy defined above. If the tag name is the same, keep the original tags without adding additional tags. If the tag name is different, insert an additional tag outside of the original tag.

5. It is possible to have the same tag to be used more than once in a single input prompt.

However, for consecutive prompt components sharing the same tags, merge them and use only one tag instead of multiple same tags.

few-shot examples

Input prompt: {example_input_prompt}

Output: {example_output}

Prompt

This is the input prompt you want to annotate: {input_prompt_to_be_annotated}

Prompt for task annotation

You are an LLM prompt task type classifier. Your task is to classify a prompt to certain task type. Please classify the prompt into one of the below task types:

Prompt Task Type Options

- * Classification:Sentiment Analysis
- * Classification:Topic Classification
- * Classification:Toxicity Detection
- * Classification:Multi-label Classification
- * Classification:Others
- * Open Book QA:Reading Comprehension
- * Open Book QA:Document-based QA
- * Open Book QA:Multi-document QA
- * Open Book QA:Context-specific Questions
- * Open Book QA:Others
- * Closed Book QA:Factual QA
- * Closed Book QA:Common Knowledge Questions
- * Closed Book QA:Others
- * Coding:Algorithm Implementation
- * Coding:Debugging
- * Coding:Code Generation
- * Coding:Code Explanation
- * Coding:Function Implementation
- * Coding:Others
- * Reasoning:Logical Reasoning
- * Reasoning:Mathematical Problem Solving
- * Reasoning:Causal Reasoning
- * Reasoning:Others
- * Summarization:Text Summarization
- * Summarization:Multi-document Summarization
- * Summarization:Others
- * Function calling:Data Retrieval (e.g., get_user_info, fetch_record_by_id)
- * Function calling:System Operations (e.g., check_status, validate_token)
- * Function calling:API Integration (e.g., external_api_call, service_request)
- * Function calling:Data Transformation (e.g., format_data, convert_units)
- * Function calling:Multi-function Chain (e.g., complex operations of multiple function calls)
- * Function calling:Others
- * Others

Additional Tips

1. If a prompt belongs to logical reasoning of the reasoning task type, then your output should be "Reasoning:Logical Reasoning".
2. Please use the uppercase letter and Lowercase letter exactly as above. Please do not include any information in the parenthesis in the task type name.
3. If you cannot classify a prompt into a task type, please feel free to use "Others". For example, if a task belongs to summarization, but not text summarization or multi-doc summarization, then the output should be "Summarization:Others".
4. Please only output task type of the prompt. Please do ****NOT**** include any preambles or additional explanations.
5. Please do not output more than one task type for a given prompt. Each prompt should only be linked to one task type.
6. Typically, if a prompt contains one or more tools, you may probably classify it into one of the function calling task types.
7. Please do ****NOT**** extract actual content from the prompt as task type. Please only use the above options as the task type.

Few-shot Examples

Input Prompt: "{example task instruction}"

Output: "{task type}"

Prompt

This is the input prompt you want to annotate: {input_prompt_to_be_annotated}

Example of taxonomy guided prompt

```
<instruction:task>Classify the given tweet into one of three categories: Hate Speech, Offensive, or Normal. Analyze the content for threatening, abusive, or discriminatory language towards specific communities.</instruction:task>
```

```
<instruction:guideline:behavioral>Approach each tweet objectively and systematically. Consider the language used, the target of the message, and the overall tone.</instruction:guideline:behavioral>
```

```
<instruction:guideline:safety>Be aware that you may encounter disturbing or offensive content. Maintain a professional demeanor and focus on the classification task.</instruction:guideline:safety>
```

```
<contextual_ref:fewshot>
{omitted due to limited space}
</contextual_ref:fewshot>
```

```
<output_const:format>Respond with the classification (Hate Speech, Offensive, or Normal) followed by a brief explanation.</output_const:format>
```

```
<output_const:style_tone>Maintain a neutral and analytical tone in your response.</output_const:style_tone>
```

```
<other:adversarial>Be cautious of tweets that may seem offensive at first glance but do not actually target any specific community.</other:adversarial>
```

Example of taxonomy annotated prompt: apigen-80k

```
<instruction:guideline:role>You are an expert in composing functions.</instruction:guideline:role>
```

```
<instruction:guideline>
You are given a question and a set of possible functions.
Based on the question, you will need to make one or more function/tool calls to achieve the purpose. If none of the functions can be used, point it out and refuse to answer. If the given question lacks the parameters required by the function, also point it out.
</instruction:guideline>
```

```
<tools:tool_description>You have access to the following tools:</tools:tool_description>
<tools>[{"type":"function","function":{"name":"similarity_score","description":
"Calculates the similarity score between two lists of integers.",
"parameters":{"type":"object",
"properties":{"list1":{"type":"array","items":{"type":"integer"},
"description":"The first list of integers."},
"list2":{"type":"array","items":{"type":"integer"},
"description":"The second list of integers."}},
"required":["list1","list2"]}}}]</tools>
```

```
<output_const:format>The output MUST strictly adhere to the following format, and NO other text MUST be included.
```

```
The example format is as follows. Please make sure the parameter type is correct. If no function call is needed, please make the tool calls an empty list '[]'.
```

```
<tool_call>[
{"name": "func_name1", "arguments": {"argument1": "value1", "argument2": "value2"}},
... (more tool calls as required)
]</tool_call></output_const:format>
```

```
<request_query>
Calculate the similarity score between the lists
[1, 2, 3, 4, 5] and [3, 4, 5, 6, 7].
</request_query>
```