



PromptSuite: A Task-Agnostic Framework for Multi-Prompt Generation

Eliya Habba* Noam Dahan* Gili Lior Gabriel Stanovsky

The Hebrew University of Jerusalem

eliya.habba@mail.huji.ac.il

Abstract

Evaluating LLMs with a single prompt has proven unreliable, with small changes leading to significant performance differences. However, generating the prompt variations needed for a more robust multi-prompt evaluation is challenging, limiting its adoption in practice. To address this, we introduce PromptSuite, a framework that enables the automatic generation of various prompts. PromptSuite is flexible – working out of the box on a wide range of tasks and benchmarks. It follows a modular prompt design, allowing controlled perturbations to each component, and is extensible, supporting the addition of new components and perturbation types. Through a series of case studies, we show that PromptSuite provides meaningful variations to support strong evaluation practices. All resources, including the Python API, source code, user-friendly web interface, and demonstration video, are available at: <https://eliyahabba.github.io/PromptSuite/>.

1 Introduction

Recent studies have demonstrated that LLMs are highly sensitive to small, meaning-preserving variations in task formulation. Minor changes, ranging from adding white spaces to instruction paraphrasing, lead to substantial differences in performance and model ranking (Sclar et al., 2023; Mizrahi et al., 2024). This sensitivity has been explored in the evaluation of many NLP tasks in zero and few shot settings, such as text classifications (Chakraborty et al., 2023; Reif and Schwartz, 2024); multiple-choice question answering (Habba et al., 2025; Alzahrani et al., 2024); and text generation tasks (Resendiz and Klinger, 2024), raising concerns about the validity of evaluation performed using a single prompt.

Evaluating over multiple prompts is currently challenging because there is no standard way to extend existing benchmarks, which were largely compiled using a single prompt. Evidently, despite its major limitation, single-prompt evaluations are still prevalent in many NLP tasks (Gu et al., 2024a,b; Lior et al., 2025).

To address this major challenge standing in the way of meaningful evaluation in NLP, we present PromptSuite, a framework that generates multiple prompts, employing both LLMs as well as rule-based heuristics to generate variations along dimensions that were found to affect model performance. PromptSuite is built around three core principles, presented in Section 2. First, PromptSuite is *flexible*, designed to work out of the box on a wide range of benchmarks. Second, PromptSuite follows a *modular* design that decomposes prompts into four components: instruction, prompt format, demonstration, and instance content, and PromptSuite enables targeted perturbations to each component, making it easy to evaluate their impact and adapt to new tasks. Finally, PromptSuite is *extensible*, supporting future LLM evaluation research with easy extensions for new prompt components and perturbations.

PromptSuite provides different types of perturbations to each prompt component, including formatting, paraphrasing, context addition, and few-shot demonstration editing, as illustrated in Figure 1. In Section 3, we provide further details on the perturbation types, as well as demonstrate how to use our API to transform raw data into multiple prompt variations with just a few lines of code.

In Section 4, we demonstrate the flexibility of our framework through a series of case studies. We assess the impact of prompt variation on nine diverse tasks with two SOTA LLMs, highlighting the utility of PromptSuite for multi-prompt evaluation.

Our contributions are as follows:

*Equal contribution.

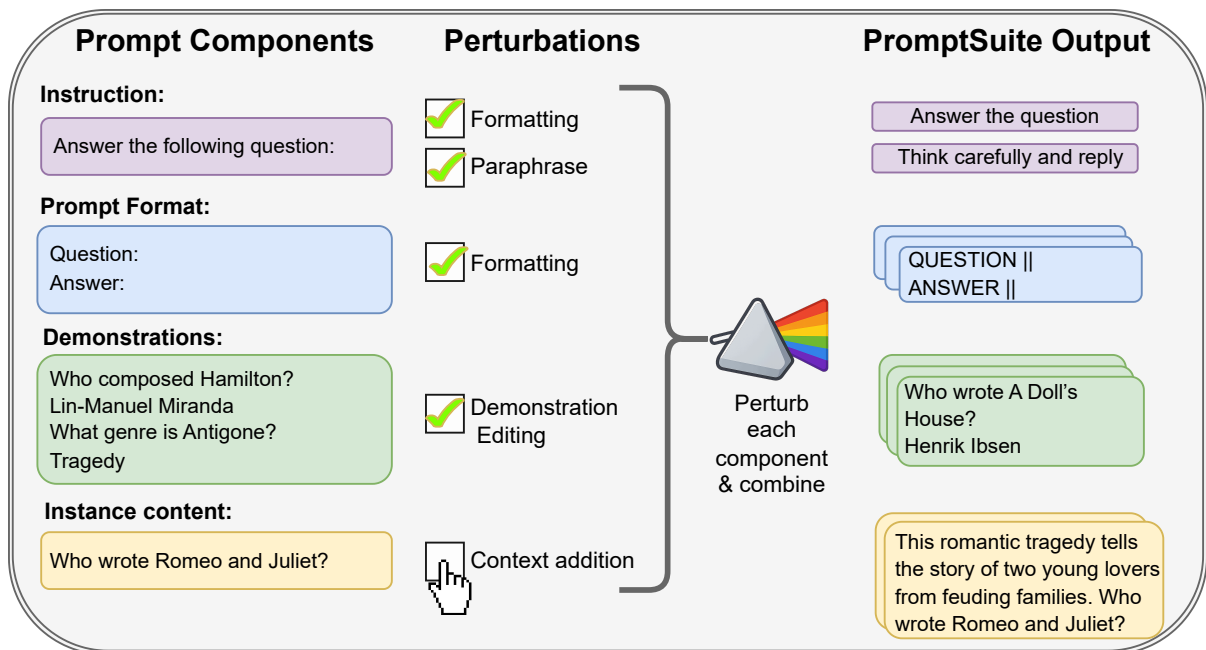


Figure 1: PromptSuite framework: configure a modular prompt, and apply component-wise perturbations. This modularity enables PromptSuite to generalize across tasks and adapt to diverse data.

1. We present PromptSuite, an easy-to-use framework that provides the prompt variations needed for multi-prompt evaluation across a wide range of tasks, working out of the box with diverse benchmarks and datasets.
2. We evaluate PromptSuite’s capabilities through a series of case studies that demonstrate its ability to reveal LLM sensitivity to prompt variations across a diverse set of tasks.
3. We show PromptSuite’s modular design enables isolating the effect of individual prompt component perturbations, enabling future research into the causes of LLM sensitivity.

2 PromptSuite’s Principles

We build PromptSuite on three core principles to make it useful across a variety of use cases and applicable over time.

Flexible PromptSuite must be able to support a diverse range of tasks out of the box. We achieve this by relying only on the prompt structure and not its content. We demonstrate this flexibility in Section 4, where we use PromptSuite to get prompt variations in 9 different tasks, including question answering, multiple-choice reasoning, translation, summarization, and code generation.

Modular PromptSuite is built around a modular representation, treating each prompt as a combination of independent components, as shown in Figure 1. This design enables targeted perturbations to specific components, supporting evaluation of their impact on model performance and allowing adaptation to new tasks, formats, or datasets.

Extensible PromptSuite is built to support future research in the field of LLM evaluation. In particular, it is easy to add new prompt components and new types of perturbations through our open-source code.¹

3 PromptSuite

PromptSuite is a flexible, modular and extensible framework that generates diverse prompts needed for robust evaluation. For a given dataset, it outputs a set of prompt variations for few-shot or zero-shot settings, where each sample from the dataset appears multiple times with different prompts. For example, in Listing 1, we load SQuAD (Rajpurkar et al., 2016a) through Hugging Face, set up the prompt and choose to paraphrase the instruction and apply formatting to the prompt format. This results in 9 prompt variations per sample, with just a few lines of code, and minimal information about the dataset.

¹<https://www.github.com/eliyahabba/PromptSuite>

```

from promptsuite import PromptSuite

# Initialize
ps = PromptSuite()

# 1) Load dataset directly from HF
ps.load_dataset("rajpurkar/squad")

# 2) Setup template and 3) Choose perturbations
template = {
    'instruction': 'Please answer the following questions.',
    'prompt_format': 'Q: {question}\nA: {answer}',
    "instruction_variations": ["paraphrase_with_llm"],
    "prompt_format_variations": ["format_structure"],
}
ps.set_template(template)

# 4) Generate variations
ps.configure(variations_per_field=3,
             api_platform="OpenAI", model_name="gpt-4o-mini")
variations = mp.generate(verbose=True)

# Export results
ps.export("output.json", format="json")

```

Listing 1: Code snippet for using PromptSuite API. Here we apply paraphrasing to the instruction with an LLM and formatting to the prompt format.

In this section, we briefly describe the modular prompt design and the supported perturbations, followed by an overview of how PromptSuite is used.

3.1 Modular Prompt and Perturbations

PromptSuite treats each prompt as a concatenation of components, allowing controlled perturbations to each part, as illustrated in Figure 1. This interpretation of prompt structure is an integration of several recent works that identified prompt components that affect overall performance (Sclar et al., 2023; Mondshine et al., 2025). Specifically, each prompt is comprised of: *instruction* (e.g., “Answer the following question”, “Summarize the following text”); *prompt format* (e.g., “Question:, Answer:”, “Text:, Summary:”); *demonstrations*; and *instance content* – the current sample the model is evaluated on (“Who wrote Romeo and Juliet?”).

Each component can be subjected to different perturbations, as detailed in Table 1. All of the perturbations preserve the original meaning of the prompt, as well as the intended output. *Formatting* refers to changes that modify either the structure of

the prompt or the appearance of its textual content. These are rule-based perturbations which can be applied to all prompt components² and include for example: inserting extra spaces; introducing typos (e.g., “apple” → “aplpe”); changing letter casing, and altering punctuation. This form of noise mimics the kind of variation found in real-world user inputs (Ravichander et al., 2021), and has been shown to affect model performance (Sclar et al., 2023).

Paraphrasing is an LLM-based perturbation that changes the wording of the instruction. We use the prompting method of Mizrahi et al. (2024), which has been shown to produce paraphrases that surface models’ sensitivity.

Context addition perturbation adds thematically related text to the prompt without changing the gold answer or providing additional hints. While the task remains unchanged, the added content makes the prompt longer and potentially more challenging for the model (Levy et al., 2024).

Lastly, *Demonstration Editing* refers to changes to the few-shot demonstration – namely, the number of examples, which ones are included, and their order, following (Lu et al., 2021). In addition to the general perturbation strategies, we also support task-specific features for common setups (e.g., changing enumerators in multi-answer questions). These are described in Appendix A.1.

3.2 Using PromptSuite

We provide a detailed overview of using PromptSuite. The package containing PromptSuite can be installed in the desired environment using pip:

```
pip install promptsuite
```

PromptSuite transforms raw data into diverse prompt variations in four steps, as can be seen in Listing 1.

(1) Load datasets: PromptSuite supports data from HuggingFace Datasets Library or local sources, including pandas DataFrames, JSON, and CSV files.

(2) Setup template: To apply the desired perturbations, PromptSuite requires the structure of the prompt and which dataset columns should be used in it. The *Instruction*, like “Please answer the following question”, is given as a plain string. The *Prompt format*, such as Q: question A: answer,

²Different implementations are applied to different components

Perturbation Type	Applicable Components	Description
Formatting	Instruction, prompt format, demonstrations, instance content	Adds surface-level noise to the text. It includes inserting extra spaces, introducing typos, changing letter casing, and altering punctuation. Following (Sclar et al., 2023).
Paraphrase	Instruction	Creates semantically equivalent variations to the instruction that differ in phrasing and style. Following (Mizrahi et al., 2024).
Context addition	Instance content	Uses an LLM to add text related to the instance content without revealing or changing the answer. Following (Liu et al., 2023; Levy et al., 2024).
Demonstration Editing	Demonstrations	Changes the few-shot examples, their order and their number. Following (Lu et al., 2021).

Table 1: Overview of the perturbation types supported by PromptSuite. The ‘‘Applicable Components’’ column specifies which prompt components each perturbation can be applied to. For example, paraphrasing is applicable to the instruction component.

is written using Python’s f-string syntax. Each placeholder (e.g., `question`) must match a column name in the dataset. For example, in Listing 1, the columns are ‘question’ and ‘answer’.

(3) Choose perturbations: Each component may be subjected to different perturbations, according to the user’s choice, as described in Table 1. These choices are added to the template setup, by specifying the name of the component and the variation. For example, in Listing 1, we choose to create LLM-based paraphrase on the instruction and apply formatting to the prompt format. To ensure maximum flexibility, users can not only modify the prompt components (i.e, instruction, prompt format, demonstrations, and instance content) but also apply alterations to any column included as a placeholder in the prompt template.

(4) Generate variations: Lastly, the user can specify the number of perturbations per component. To ensure the dataset remains manageable in terms of cost and memory, users can also limit the total number of generated rows. Since each component supports multiple perturbations, the number of possible dataset variations grows exponentially with the number of chosen perturbations. To produce a manageable dataset size, we provide an option to randomly select a combination of the desired perturbations and apply them across the entire dataset, following the approach of Habba et al. (2025).

PromptSuite is also available via a web interface. We offer the full capabilities of PromptSuite through a web UI, as illustrated in Figure 2.³ Users

follow the same steps described above: first, upload their single-prompt dataset, then configure the prompt components and select the desired perturbations for each component. As shown in the figure, PromptSuite’s web UI allows users to explore the generated variations, highlighting the changes applied to each row. The interface also provides several predefined templates for popular tasks, including multiple-choice QA, sentiment analysis, open-ended QA, and text classification, enabling a quick and easy plug-and-play setup for users who wish to automatically generate multi-prompt versions of their datasets.

4 Evaluation

We demonstrate that PromptSuite is flexible and generalizes across a wide range of tasks by applying it to nine diverse benchmarks. Our results show that multi-prompt evaluation reveals substantial performance variance that would have been missed using a single prompt. We further assess the impact of perturbations to individual prompt components on model performance by leveraging PromptSuite’s modular design.

4.1 Experimental Setup

Tasks and datasets. We evaluate PromptSuite on: (1) MMLU (Hendrycks et al., 2021) for multiple-choice reasoning across 12 subjects; (2) GSM8K (Cobbe et al., 2021) for mathematical problem solving; (3) SST (Socher et al., 2013) for sentiment analysis; (4) WMT14 (Bojar et al., 2014) for translation across 6 language pairs (CS/HI/RU↔EN); (5) CNN/Daily-Mail (Hermann et al., 2015) for summarization; (6) MuSiQue (Trivedi et al., 2022) for multi-hop

³<https://promptsuite.streamlit.app/>

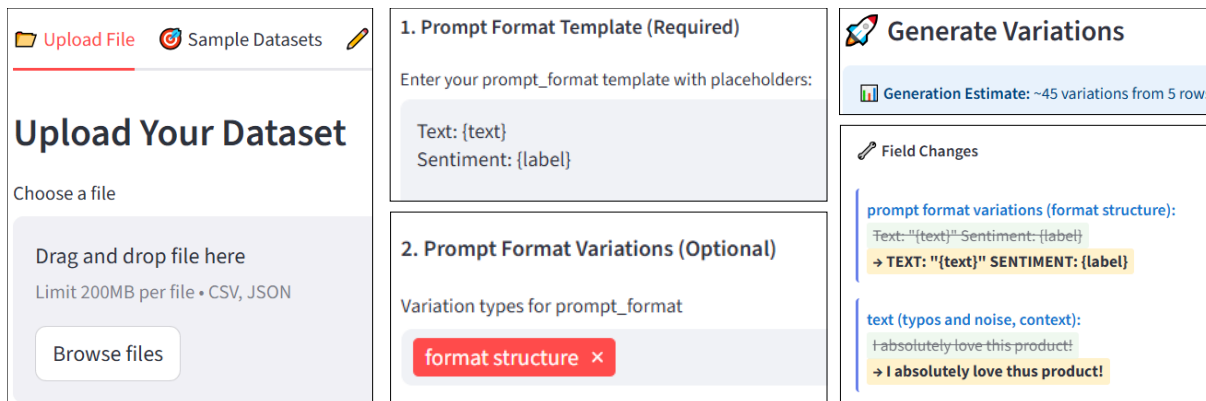


Figure 2: PromptSuite’s web UI. Left-to-right: uploading a dataset; configuring the template and choosing perturbations; and generating a multi-prompt dataset. The presented example demonstrates a single prompt variation, with changes to the prompt format and instance content.

question answering; (7) SQuAD (Rajpurkar et al., 2016b) for reading comprehension; (8) GPQA-Diamond (Rein et al., 2024) for graduate-level reasoning; and (9) HumanEval (Chen et al., 2021) for code generation.

Models. We evaluate GPT-4o-mini and Llama-3.3-70B, representing closed and open-source LLMs. Temperature is set to 0 to ensure consistent and deterministic outputs. For code generation, we use a temperature of 0.8, since Pass@k relies on generating multiple candidate solutions, and a non-zero temperature is essential to ensure a diverse set of outputs across multiple runs for the same prompt, as demonstrated in (Chen et al., 2021).

Prompt variations. For each task, we generate variations using: paraphrasing, formatting applied to the prompt format and demonstration editing. We process 50 rows per dataset with up to 25 variations per row, resulting in approximately 1,250 evaluated prompts per task and a total of 37,000 LLM outputs (detailed calculations in Appendix A.1, Table 3, and token counts in Table 4). This yields comprehensive coverage while remaining computationally tractable.

Manual validation. To validate our results, we conducted human validation of a subset of 100 LLM-based paraphrases. Two in-house annotators independently annotated all 100 samples, reaching 95% agreement (Cohen’s $k = 0.593$). They judged that 96% of the paraphrases preserved the original meaning of the instruction. The samples that were tagged as incorrect were either due to the use of a less accurate synonym (e.g., for sentiment analysis instruction, it rephrased “sentiment” into

“emotional tone”, which can be ambiguous, or the omission of the system message, such as “you are an expert in QA”).

4.2 Results

Below we outline interesting conclusions derived from our experiments using PromptSuite.

Models exhibit sensitivity to the prompt perturbations across all tasks, underscoring the utility of PromptSuite. Figures 3a and 3b show performance distributions across prompt variations for Llama-3.3-70B and GPT-4o-mini, respectively. We observe substantial variability. For instance, on GPQA-Diamond, GPT-4o-mini’s accuracy ranges from 20% to 50% across variations. This variance is particularly striking when compared to typical performance differences between competing models, which often amount to just a few percentage points (Lior et al., 2025). The consistency of this pattern across diverse tasks demonstrates that prompt sensitivity is not limited to specific domains but represents a general challenge in LLM evaluation.

PromptSuite’s modularity enables systematic ablation, showing that the impact of prompt component perturbations varies across tasks and models. PromptSuite’s modularity allows a systematic ablation study that assesses the impact of changes in specific prompt components on model performance. We perturb one prompt component at a time to measure its specific effect, testing instruction paraphrasing, formatting applied to either the prompt format or the instance content and demonstrations editing. We conduct this experiment on GPQA-Diamond, SQuAD and GSM8K.

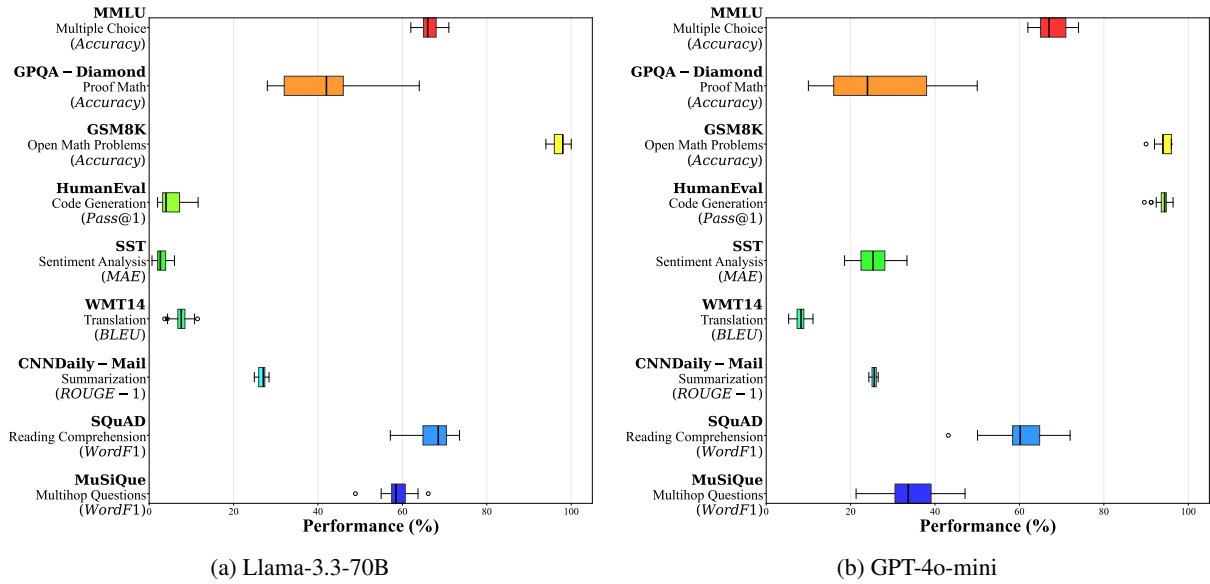


Figure 3: Multi-prompt evaluation results using PromptSuite. The boxplots illustrate variance across different prompt perturbations, revealing models’ sensitivity to prompt variations and underscoring the utility of PromptSuite for deriving robust and meaningful evaluations of LLM capabilities.

Each perturbation was evaluated on 50 examples with 20 variations per example, yielding 1000 evaluated prompts per component-task combination for each model. Figure 4 presents the performance distributions across perturbation types in GPQA-Diamond. For example, we observe that demonstration editing caused high variance in Llama-3.3-70B’s performance, whereas for other tasks (Figure 5 in the Appendix), demonstration editing had almost no effect on Llama-3.3-70B’s performance. Similarly, for GPT-4o-mini, prompt formatting had almost no effect on GPQA-Diamond (Figure 4), but showed a more significant effect for SQuAD (Figure 5). These inconsistencies across models and tasks underscore the importance of a flexible and modular framework like PromptSuite, which enables systematic analysis of prompt component effects. For example, practitioners can leverage PromptSuite for a more efficient evaluation strategy, by first experimenting on a small subset of the data to identify the most influential prompt components, and then conducting a focused large-scale evaluation that concentrates only on the perturbations with the most significant impact.

5 Related Work

A few existing frameworks support a subset of PromptSuite’s capabilities. To the best of our knowledge, they are task-specific and require manual control over input variations. NL-

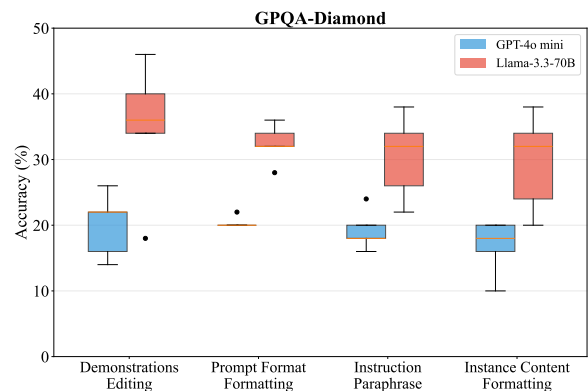


Figure 4: Analysis of how perturbations to individual prompt components affect model sensitivity on GPQA-Diamond. Each boxplot represents an experiment in which a single prompt component was varied while all others remained fixed.

Augmenter (Dhole et al., 2021) is a crowd-sourced repository for perturbations. While it provides a wide range of task-specific transformations and filters, it operates solely on the input data and does not account for instructions or templates, both of which are critical for robust few-shot evaluation. Unitxt (Bandel et al., 2024) is a library for data preparation and evaluation, with some tasks supporting a limited number of data alterations. Prompt-Agnostic Fine-Tuning (PAFT) (Wei et al., 2025) also seeks to reduce prompt sensitivity by generating diverse prompts, but incorporates them during finetuning rather than at evaluation time.

6 Conclusion

We introduce PromptSuite, a framework that generates prompt variations needed for multi-prompt evaluation. It is flexible, uses a modular design for controlled perturbations, and is easily extensible. Through case studies, we show that the variations generated by PromptSuite are sufficient to test model sensitivity.

7 Limitations

While PromptSuite provides a general, task-agnostic framework for multi-prompt evaluation, this generality comes with a tradeoff. Some tasks may involve specific variations or prompt structures that PromptSuite does not currently support. Specifically, this limitation arises in cases where evaluation is not straightforward (e.g., multi-turn chat). To mitigate this limitation, we designed our code to be easily extensible, allowing users to add additional prompt components or variations as needed.

Acknowledgments

This research was also supported by the Ministry of Innovation, Science & Technology, Israel (Grant No. 0008239). We would like to thank our colleagues Asaf Yehudai, Hillel Darshan, Daniel Nisnevich, Nitzan Barzilay, and Michael Hassid for their contributions to our initial prototype development, which served as an inspiration for this work. We would also like to thank Omer Kidron for his assistance in recording the demonstration video.

References

- Norah Alzahrani, Hisham Alyahya, Yazeed Alnumay, Sultan AlRashed, Shaykhah Alsubaie, Yousef Al-mushayqih, Faisal Mirza, Nouf Alotaibi, Nora Al-Twairish, Areeb Alowisheq, M Saiful Bari, and Haidar Khan. 2024. [When benchmarks are targets: Revealing the sensitivity of large language model leaderboards](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13787–13805, Bangkok, Thailand. Association for Computational Linguistics.
- Elron Bandel, Yotam Perlit, Elad Venezian, Roni Friedman-Melamed, Ofir Arviv, Matan Orbach, Shachar Don-Yehyia, Dafna Sheinwald, Ariel Gera, Leshem Choshen, and 1 others. 2024. [Unitxt: Flexible, shareable and reusable data preparation and evaluation for generative ai](#). *arXiv preprint arXiv:2401.14019*.
- Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleš Tamchyna. 2014. [Findings of the 2014 workshop on statistical machine translation](#). In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Mohna Chakraborty, Adithya Kulkarni, and Qi Li. 2023. [Zero-shot approach to overcome perturbation sensitivity of prompts](#). *arXiv preprint arXiv:2305.15689*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. [Evaluating large language models trained on code](#). *Preprint*, arXiv:2107.03374.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. [Training verifiers to solve math word problems](#). *arXiv preprint arXiv:2110.14168*.
- Kaustubh D Dhole, Varun Gangal, Sebastian Gehrmann, Aadesh Gupta, Zhenhao Li, Saad Mahamood, Abinaya Mahendiran, Simon Mille, Ashish Shrivastava, Samson Tan, and 1 others. 2021. [Nl-augmenter: A framework for task-sensitive natural language augmentation](#). *arXiv preprint arXiv:2112.02721*.
- Alex Gu, Wen-Ding Li, Naman Jain, Theo Olausson, Celine Lee, Koushik Sen, and Armando Solar-Lezama. 2024a. [The counterfeit conundrum: Can code language models grasp the nuances of their incorrect generations?](#) In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 74–117, Bangkok, Thailand. Association for Computational Linguistics.
- Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. 2024b. [Cruxeval: A benchmark for code reasoning, understanding and execution](#). *arXiv preprint arXiv:2401.03065*.
- Eliya Habba, Ofir Arviv, Itay Itzhak, Yotam Perlit, Elron Bandel, Leshem Choshen, Michal Shmueli-Scheuer, and Gabriel Stanovsky. 2025. [Dove: A large-scale multi-dimensional predictions dataset towards meaningful llm evaluation](#). *Preprint*, arXiv:2503.01622.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman,

- and Phil Blunsom. 2015. [Teaching machines to read and comprehend](#). *Preprint*, arXiv:1506.03340.
- Mosh Levy, Alon Jacoby, and Yoav Goldberg. 2024. Same task, more tokens: the impact of input length on the reasoning performance of large language models. *arXiv preprint arXiv:2402.14848*.
- Gili Lior, Eliya Habba, Shahar Levy, Avi Caciularu, and Gabriel Stanovsky. 2025. Reliableeval: A recipe for stochastic llm evaluation via method of moments. *arXiv preprint arXiv:2505.22169*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. [Lost in the middle: How language models use long contexts](#). *Preprint*, arXiv:2307.03172.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*.
- Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2024. State of what art? a call for multi-prompt llm evaluation. *Transactions of the Association for Computational Linguistics*, 12:933–949.
- Itai Mondshine, Tzuf Paz-Argaman, and Reut Tsarfaty. 2025. Beyond english: The impact of prompt translation strategies across languages and tasks in multilingual llms. *arXiv preprint arXiv:2502.09331*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016a. [Squad: 100,000+ questions for machine comprehension of text](#). *Preprint*, arXiv:1606.05250.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016b. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Abhilasha Ravichander, Siddharth Dalmia, Maria Ryskina, Florian Metze, Eduard Hovy, and Alan W Black. 2021. [Noiseqa: Challenge set evaluation for user-centric question answering](#). *Preprint*, arXiv:2102.08345.
- Yuval Reif and Roy Schwartz. 2024. [Beyond performance: Quantifying and mitigating label bias in llms](#). *Preprint*, arXiv:2405.02743.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.
- Yarik Menchaca Resendiz and Roman Klinger. 2024. Mopo: Multi-objective prompt optimization for affective text generation. *arXiv preprint arXiv:2412.12948*.
- 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*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. [Musique: Multi-hop questions via single-hop question composition](#). *Preprint*, arXiv:2108.00573.
- Chenxing Wei, Yao Shu, Mingwen Ou, Ying Tiffany He, and Fei Richard Yu. 2025. [Paft: Prompt-agnostic fine-tuning](#). *Preprint*, arXiv:2502.12859.

A Example Appendix

A.1 Task-Specific Perturbations

For multiple choice questions, multi-document or any tasks that includes a list as input we offer perturbations of said list, presented in Table 2. This includes support for enumerating the list items, as well as shuffling the order (while changing the gold answer accordingly, if applicable).

Perturbation Type	Applicable Fields	Description
Enumerate	lists, comma-separated values	adds enumeration to a specified field, such as multiple-choice options, by prepending each item with a number or letter (e.g., 1., A., a.). Following (Habba et al., 2025)
Shuffle	lists	shuffles the items in a list and updates the gold field to reflect the new index of the correct answer. For example, if the correct answer was originally at position B and is moved to position C after shuffling, the gold label is updated accordingly. Following (Habba et al., 2025)

Table 2: Task-specific perturbation types in PromptSuite. The "Applicable Fields" column indicates which types of data column the perturbation works on.

Benchmark/Task	Questions	Variations	Total (Q × V)
MMLU Multiple Choice(12 subjects, 10 questions each)	120	50	6000
GSM8K Open Math Problems	50	25	1250
HumanEval Code Generation	50	25	1250
SST Sentiment Analysis	50	50	2500
WMT14 Translation (CS/HI/RU ↔ EN)	60	50	3000
CNN-DailyMail Summarization	50	25	1250
MuSiQue Multihop Questions	50	25	1250
SQuAD Reading Comprehension	50	25	1250
GPQA–Diamond Google-Proof Math	50	25	1250
Total across both models	–	–	37,375

Table 3: Number of evaluated examples per benchmark. Each row indicates the number of base questions and variations, with the total computed as their product. **Note:** Values shown reflect the GPT-4o mini configuration. For LLaMA-3-3.7B, the MuSiQue dataset included only 25 base questions (instead of 50) due to limited context window constraints, yielding a total of 625 examples for that task. Total reflects the combined number of evaluations across both models.

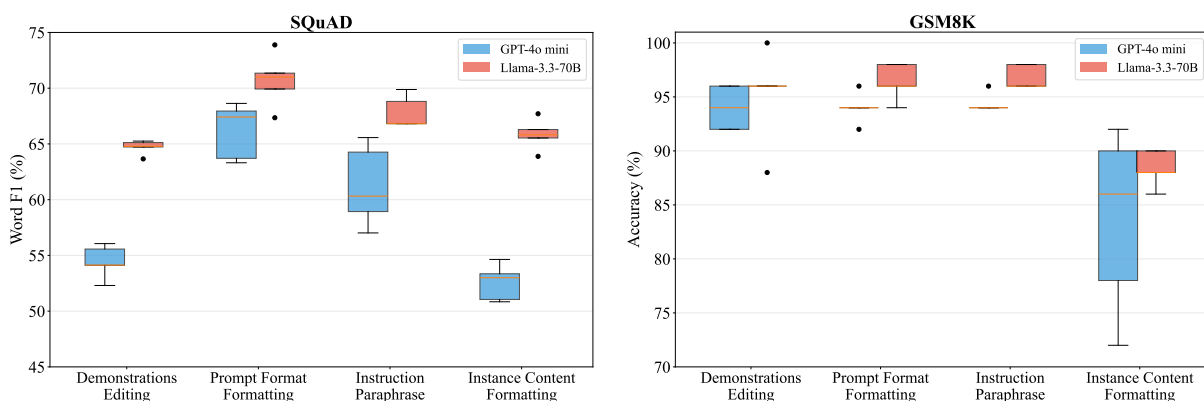


Figure 5: Analysis of how perturbations to individual prompt components affect model sensitivity on SQuAD and GSM8K. Each boxplot represents an experiment in which a single prompt component was varied while all others remained fixed.

Benchmark/Task	Input Tokens		Output Tokens	
	GPT-4o mini	Llama-3.3 70B	GPT-4o mini	Llama-3.3 70B
MMLU Multiple Choice	1389791	1391104	68403	173011
GSM8K Open Math Problems	722880	731987	223225	135769
HumanEval Code Generation	7259990	7228940	285569	285313
SST Sentiment Analysis	479750	487721	10934	108987
WMT14 Translation	409989	422709	34478	95237
CNN-DailyMail Summarization	3888319	3940895	167329	155074
MuSiQue Multihop Questions	8586023	8102151	17371	14525
SQuAD Reading Comprehension	1279578	1288793	10099	9787
GPQA-Diamond Google-Proof Math	1658817	1679053	384586	461600

Table 4: Token usage per benchmark across GPT-4o mini and Llama-3.3 70B. The table shows the number of input and output tokens consumed for each benchmark.