

WILDIFEVAL: Instruction Following in the Wild

Gili Lior^{1*} Asaf Yehudai^{1,2} Ariel Gera² Liat Ein-Dor²

¹The Hebrew University of Jerusalem ²IBM Research
gili.lior@mail.huji.ac.il

Abstract

Recent LLMs have shown remarkable success in following user instructions, yet handling instructions with multiple constraints remains a significant challenge. In this work, we introduce WILDIFEVAL – a large-scale dataset of 7K real user instructions for single-turn constrained text generation, exhibiting diverse, multi-constraint conditions. Unlike prior datasets, our collection spans a broad lexical and topical spectrum of constraints, extracted from natural user instructions. We categorize these constraints into eight high-level classes to capture their distribution and co-occurrence dynamics in real-world scenarios. Leveraging WILDIFEVAL, we conduct extensive experiments to benchmark the instruction-following capabilities of leading LLMs. WILDIFEVAL clearly differentiates between small and large models, and demonstrates that all models have room for improvement on such tasks. Our analysis reveals that as constraint count grows, models’ overall success drops sharply while per-constraint success remains stable, indicating a capacity bottleneck in juggling multiple constraints, and that models struggle more with rigid form-based constraints than with softer content-based ones. We release our dataset to promote further research on instruction-following under complex, realistic conditions.¹

1 Introduction

As LLMs continue to improve at following instructions, the nature of the instructions themselves has also evolved. Users now expect LLMs to handle more nuanced and complex requests (Wang et al., 2024). This shift is especially evident in text generation tasks, which are becoming increasingly per-

^{*}This work was conducted during a summer internship at IBM Research.

¹WILDIFEVAL is available at <https://huggingface.co/datasets/gililior/wild-if-eval>. The code for replication, along with model predictions and evaluation scores, is at <https://github.com/gililior/wild-if-eval-code>.

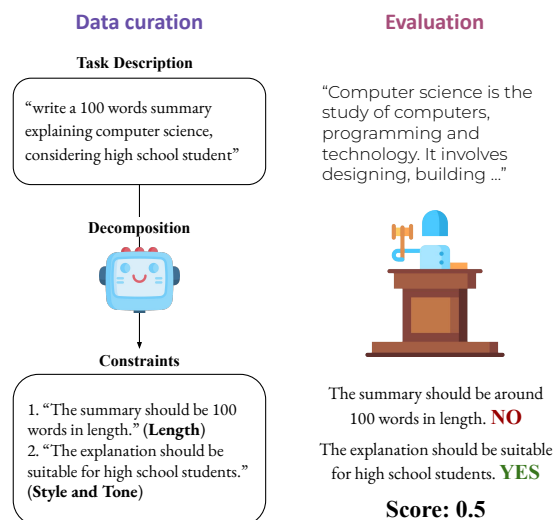


Figure 1: WILDIFEVAL description. At the left is an example for a constrained generation task, and its decomposition into constraints. In evaluation (right), the judge decides whether each of constraints is fulfilled.

sonalized, with more specific and tailored objectives (Salemi et al., 2023; He et al., 2022; Li et al., 2024a; Ein-Dor et al., 2024). For instance, a former instruction like “summarize this text” might now take the form of “summarize this movie review in two paragraphs, with the first focusing on the plot and the second discussing reasons to watch or skip the movie.” These personalized tasks typically carry implicit or explicit constraints that the generated output is expected to satisfy.

Thus, in *constrained generation* an LLM must adhere to a set of specific requirements in its response (Garbacea and Mei, 2022; Yao et al., 2023). Crucially, while individual constraints are often simple, LLMs struggle to satisfy multiple constraints simultaneously (Jiang et al., 2024). This highlights the need to directly evaluate the text generation performance of LLMs on realistic multi-constraint user data.

Existing works evaluating the ability of LLMs

Benchmark	Data Source	Evaluation	Size (# Tasks)	# Constraints
IFEval	Synthetic	Rule	541	-
FollowBench	Crowd + Syn.	Model / Rule	1,852	-
InFoBench	Crowd	Model / Rule	500	2,217
WILDIFEVAL (ours)	Real Users	Model	7,523	24,731

Table 1: Comparison of WILDIFEVAL with openly available instruction-following benchmarks such as IFEval (Zhou et al., 2023), FollowBench (Jiang et al., 2024), and InFoBench (Qin et al., 2024).

to follow constrained instructions generally follow a bottom-up approach, starting from curated verifiable constraints, that are amenable to objective verification of compliance (Zhou et al., 2023), or a taxonomy of constraint types (Yao et al., 2023; Qin et al., 2024; Jiang et al., 2024), and using those to manually or synthetically generate a set of instructions. Such an approach may not capture the complexity and diversity of real-world instructions by users, and the types and combinations of constraints that they ask the model to follow.

To this end, we introduce WILDIFEVAL (§2), a large-scale benchmark of constrained generation tasks. WILDIFEVAL is designed to evaluate the ability of LLMs to follow real-world multi-constrained instructions in *single-turn text generation*, distinguishing it from agentic or multi-turn instruction-following benchmarks. It encompasses a collection of 7K constrained generation tasks, including 24,731 different constraints, given by real users on Chatbot Arena (Chiang et al., 2024), reflecting diverse examples of constrained generation instructions “in the wild”.

The WILDIFEVAL dataset includes a breakdown of each task into the individual constraints it contains. Thus, it allows for a fine-grained evaluation of the ability of LLMs to adhere to user constraints. By breaking down task instructions into smaller and more interpretable pieces, we can perform a straightforward LLM-based evaluation of the proportion of task constraints that were fulfilled. At the same time, since constraints are extracted from naturalistic user queries, we capture not only simple and easily verifiable constraints but also “softer” constraints on content, quality, and style.

We begin by analyzing the types of user tasks and constraints present in WILDIFEVAL (§3), revealing that real-world constrained generation often involves diverse and challenging requirements.

We then evaluate 14 LLMs on the WILDIFEVAL benchmark and conduct a comprehensive analysis of their constraint-following capabilities (§4).

Our results show that WILDIFEVAL is challenging, with the best models achieving around 0.7 under our strict evaluation metric. We also observe a consistent performance gap between small and large models, positioning WILDIFEVAL as a valuable benchmark for tracking progress to narrow this gap.

Beyond overall model performance, we utilize the size and diversity of WILDIFEVAL to provide the first large-scale analysis of how constraints are distributed and combined in real user instructions. We introduce a taxonomy of eight constraint categories, and analyze their co-occurrence patterns, lexical diversity, and “long tail” of rare constraint phrasings – patterns that are impossible to recover from smaller, synthetically generated benchmarks. Our analysis further reveals a notable divergence between strict and soft instruction-following: as constraint count grows, models’ *overall* success drops sharply, but their per-constraint success remains stable, suggesting a capacity bottleneck in juggling multiple constraints rather than a degradation in instruction-following per se. We also find that models struggle more with rigid form-based constraints (length, format) than with softer content-based ones.

By publicly releasing WILDIFEVAL, the first publicly available large-scale benchmark of naturally occurring, multi-constraint instructions, we aim to push LLMs’ ability to follow complex constraints in real-world applications.

2 The WILDIFEVAL Dataset

WILDIFEVAL is a novel benchmark designed to provide a comprehensive evaluation of the ability of LLMs to follow real-world multi-constrained instructions. It contains 7K user-generated instructions, written by many distinct users, each decomposed into a set of constraints, including 24,731 unique constraints.

The task instructions in WILDIFEVAL were extracted from LMSYS-Chat-1M dataset (Zheng et al., 2023a), a large-scale dataset containing

real-world instructions collected from the Chatbot Arena.² Since users rarely specify constraints in a structured list format, the decomposition breaks instructions into manageable items, ensuring the necessary granularity to assess the LLM’s ability to adhere to them.

In Table 1, we present a comparison with popular openly available instruction-following datasets. As can be seen in the table, WILDIFEVAL is uniquely representative of natural user interactions at scale; it stands out as the largest available English benchmark consisting of real-world user instructions given to LLMs.

2.1 Dataset Curation

WILDIFEVAL was curated in three steps. First, we filter the LMSYS-Chat-1M source data – we extract the first user message from each conversation, and filter out non-English tasks, coding tasks, and tasks containing toxic language.³ Next, we filter for only constrained generation tasks. We follow the definition for constrained generation tasks from Ferraz et al (Ferraz et al., 2024), and utilize their suggested prompt (Appendix A) with Llama3.1-405b in order to perform the filtering. The prompt is phrased as a yes/no question; instead of simply parsing the string, we use the probabilities that the model assigns to the yes/no tokens as a measure of certainty, and include only the 10% of tasks with the highest certainty to be a constrained generation task, i.e., with the highest probability for a “yes” token. The distribution of scores ranged from 0 to 1, with a mean of 0.29 and a median of 0.07, indicating that most tasks were not classified as constrained generation tasks. In contrast, the threshold for the top 10% was 0.94, suggesting that the tasks we retained were labeled positive with high certainty. We validate this thresholding procedure against human annotation in §2.2. The last step of the curation process is the decomposition into constraints – for each user task, we want to include all the constraints the model is required to fulfill. To obtain the highest-quality decomposition we employ GPT-4o (Hurst et al., 2024), using a prompt adopted from Ferraz et al (Ferraz et al., 2024) to automatically extract the constraints for each of the tasks.⁴ All prompts are presented in Appendix A.

²Chatbot Arena website: <https://lmarena.ai>, Huggingface dataset: <https://huggingface.co/datasets/lmsys/lmsys-chat-1m>.

³We detect toxic language using the detoxify package <https://github.com/unitaryai/detoxify>

⁴gpt-4o-2024-08-06

To mitigate potential biases in scoring, we perform sub-sampling for constraints that appear more than 40 times (i.e., exact match across more than 40 different tasks). This process affected 15 unique constraints, accounting for less than 0.15% of all constraints. In addition, we filtered out rare cases of tasks with more than 8 constraints. By the end of this process, we obtained a dataset of 7,523 real-world constrained generation tasks, each annotated with a list of constraints. There are 24,731 distinct constraints in WILDIFEVAL, averaging 3.25 constraints per task. The distribution and frequency of constraints per task are shown in Figure 8 in Appendix. We empirically justify this scale in Appendix C.

2.2 Human Verification

To validate the quality of our automatic curation pipeline, we conducted two human annotation studies on random subsets of 100 tasks each.

Decomposition quality. We evaluate the decomposition performed by GPT-4o along three dimensions: *correctness* (faithfully reflects the original task), *completeness* (accounts for all essential constraints), and *independence* (constraints are distinct and self-sufficient). Each was rated on a 1–5 Likert scale, yielding mean scores of 4.71 for correctness, 4.64 for completeness, and 4.77 for independence, indicating high decomposition quality.

Filtering validity. We further verified that our top-10% certainty threshold reliably identifies constrained generation tasks. Comparing human judgments with model certainty scores on 100 sampled tasks, we find human-model agreement of 75.8% at our chosen threshold, close to the optimal achievable agreement of 77.9%. Notably, the asymmetry between false positives and false negatives – only 1 task was labeled as constrained generation by the model alone, compared to 22 labeled so by humans alone, which indicates that our filtering is conservative and yields a high-precision subset. While stricter than human annotators (humans labeled 31% of sampled tasks as constrained generation vs. our 10%), this aligns with our goal of prioritizing certainty that selected tasks are genuinely constrained-generation, rather than maximizing coverage.

3 Into the Wild: A Data Expedition

Below we conduct an analysis of our WILDIFEVAL data, revealing insights on constrained generation

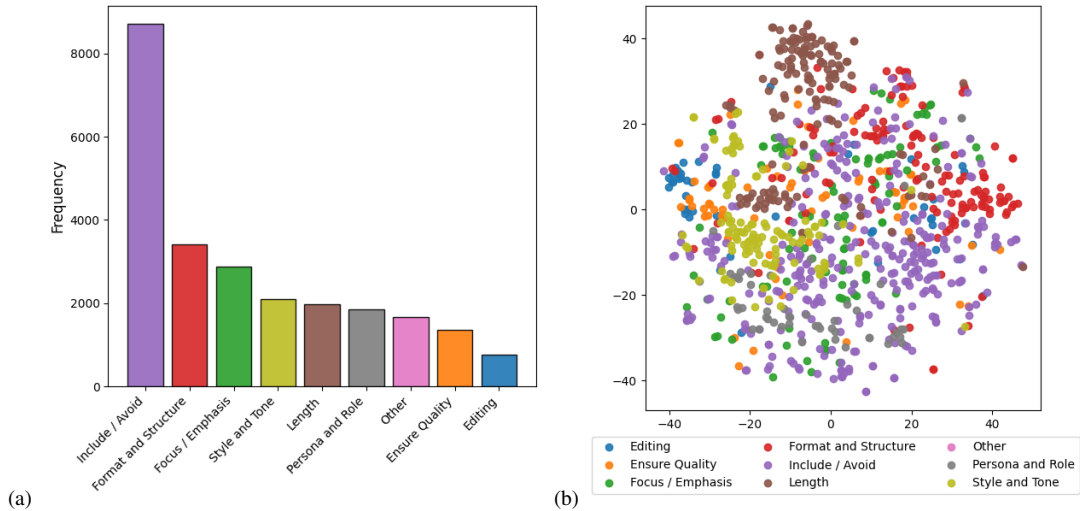


Figure 2: Analysis of constraints in WILDIFEVAL. (a) Distribution of constraint types. (b) A tSNE projection (van der Maaten and Hinton, 2008) of the embeddings of constraints, colored by their type. For convenience, we randomly subsample 1k data points. We observe some red, brown, and yellow clusters, corresponding to *Format and Structure*, *Length*, and *Style and Tone* constraints, aligning with the generic nature of these types. This is in contrast to content-oriented types like *Focus/Emphasis* and *Include/Avoid* (green and purple), which are more spread out.

use cases in the wild.

3.1 Constraint Types

A key question regarding constrained generation tasks concerns the nature and types of the constraints themselves, i.e., what kinds of requirements users wish to impose on the model responses. Prior work (Zhou et al., 2023; Ferraz et al., 2024; Jiang et al., 2024; Qin et al., 2024) generally distinguishes between broad categories such as content, style, and format, yet lacks a unified taxonomy. Moreover, some works define rather specific constraint categories (e.g., “Part-of-speech rules”) or highly general ones (e.g., “Content constraints”).

Here we seek to bridge this taxonomy gap. We draw from earlier categorization efforts, but combine them with data-driven insights. Specifically, we look at the most frequent words appearing in constraints, and examine some of the constraints in which they occur; this allows us to analyze recurring patterns of constraint types in WILDIFEVAL. This qualitative data-driven analysis reveals some broad constraint types that have not been mentioned by prior efforts, and also enables us to break existing broad divisions into finer-grained categories.

Our taxonomy divides constraints into 8 principal categories. These capture both explicit constraints (e.g., inclusion or exclusion of content) and more nuanced aspects of user instructions (e.g., a

desired tone or quality for the model output). The following definitions detail each category, providing clear guidelines on how they contribute to the overall task structure:

- **Include / Avoid:** Specifies elements or concepts that must be incorporated into or omitted from the response, directly guiding the content of the output.
- **Editing:** Focuses on modifications to an existing text, outlining how the original content should be altered or preserved.
- **Ensure Quality:** Imposes requirements on the response’s quality, such as coherence, accuracy, or overall clarity.
- **Length:** Sets quantitative boundaries on the output, such as word or character limits, ensuring appropriate brevity or depth.
- **Format and Structure:** Dictates the organization and presentation of the response, including the use of bullet points, tables, or specific layout requirements.
- **Focus / Emphasis:** Highlights particular topics, keywords, or elements that should be prioritized within the response.
- **Persona and Role:** Instructs the AI to adopt a specific character, perspective, or expertise, influencing the narrative voice of the output.

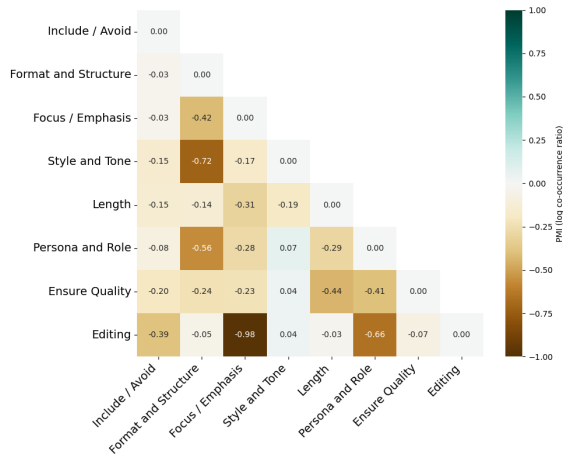


Figure 3: Relative co-occurrence (PMI) of constraint categories within tasks. Values above 0 indicate that constraints co-occur more than expected by their overall type frequencies.

- **Style and Tone:** Specifies the overall manner of expression, including formality, register, and emotional nuance, to define the voice and feel of the response.

We then ask Deepseek-v3 to classify all constraints in WILDIFEVAL into one of the 8 constraint types above, resulting in a full categorization of constraint types. The classification prompt is provided in Appendix A.

Distribution of constraint types. In Figure 2a we present the distribution of constraint types in WILDIFEVAL. The most common constraints are the content constraints *Include/Avoid* and *Focus/Emphasis*; these specify either explicit element(s) that should be included or excluded, or how much prominence should be given to different elements in the content.

Figure 2b depicts a tSNE embedding map of WILDIFEVAL constraints, colored by types.⁵ A salient and intuitive observation is that content-related constraints such as *Include/Avoid* and *Focus/Emphasis* are spread out across the semantic embedding space; in contrast, form-related constraints like *Length* or *Format and Structure* are organized in more distinct clusters.

Co-occurrence of constraint types. In Figure 3 we analyze the co-occurrence of constraint types in multi-constraint tasks. Specifically, we ask whether some combinations of types appear more or less

⁵Embeddings were computed using NV-Embed-V2 (Lee et al., 2024).

than expected. Thus, we compare the number of co-occurrences in practice relative to the overall frequency of each of the co-occurring types, i.e., the pointwise mutual information (PMI) (Church and Hanks, 1990).

As shown in Figure 3, only few combinations appear more than expected (i.e., $PMI > 0$). For example, *Persona and Role* tends to co-occur with *Style and Tone* slightly above expected, which appears to reflect the thematic similarity between these constraint types. In contrast, some types do not often appear together; for instance, requirements for *Format and Structure* are rarely paired with *Style and Tone* or *Persona and Role* constraints. Also *Editing*, which is the lowest represented type of constraint, rarely co-occurs with *Focus / Emphasis*.

3.2 Data Diversity

WILDIFEVAL covers a variety of domains.

Figure 4a depicts the distribution of domains covered by WILDIFEVAL. As expected from large-scale naturally-occurring data, tasks in WILDIFEVAL cover a wide variety of domains, including Technology, Entertainment, Healthcare, Creative Writing, and more. We use a data-driven approach to recover the domains, leading us to believe that these reflect realistic user behavior in constrained generation tasks. The domains were extracted using an LLM, see details in Appendix B.4.

WILDIFEVAL is lexically diverse.

To illustrate lexical diversity, we examine verb frequencies in constraints that begin with a verb (65.1% of constraints).⁶ The results in Figure 4b reveal a skewed frequency distribution; “Provide” is the most dominant verb, comprising 21.1% of all occurrences, followed by “Do” (19.2%) and “Write” (8.7%). Several mid-frequency verbs (e.g., “Keep,” “Identify,” “Make”) also appear regularly. The “Other” category (12.6%) reflects the long tail of the verb distribution, with many verbs that each occur in under 0.8% of the data. The distribution suggests that users tend to use general types of constraints more than specific ones like “Simplify” (0.8%) or “Summarize” (0.8%). This analysis underscores the variety of linguistic expressions in WILDIFEVAL. A similar pattern emerges when considering all constraints containing a verb (70% of constraints), shown in Figure 12 in Appendix B.3. We note that the analysis reflects the words in the constraints, as

⁶We employ NLTK’s part-of-speech tagger to identify verb tokens <https://www.nltk.org/>

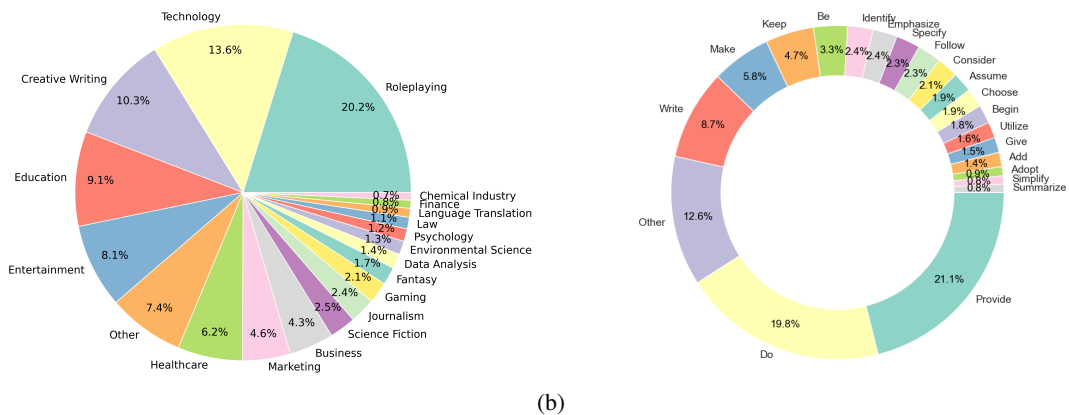


Figure 4: Task and constraint characteristics in WILDIFEVAL. (a) Domain distribution of tasks. (b) Lexical diversity of constraint phrasing (opening verbs).

decomposed by an LLM (§2.1), and thus may differ somewhat from the original user task descriptions.

Qualitative analysis. Manual inspection of instances from WILDIFEVAL reveals some interesting trends. First, we observe that quite often fulfilling – or even understanding – the task constraints given by users requires some very specialized or esoteric knowledge (e.g., D&D spells, Gate exam syllabus, pig latin etc.). We show some examples in Appendix E. We also note that some of the more complex tasks, i.e., those with many constraints, reflect attempts by users to “jailbreak” the LLM, and trick it to say things that it is not supposed to (e.g., toxic language or controversial statements).

4 LLM Benchmarking

In this section, we examine the performance of various LLMs to assess their behavior in constrained generation tasks. We present the evaluation metric (§4.1), experimental setup (§4.2), and finally, we describe and analyze the results (§4.3).

4.1 Evaluation Metric

WILDIFEVAL reports two scores: *strict* and *soft*. The *strict* score is a binary measure indicating whether all task constraints are satisfied, while the *soft* score reflects the proportion of individual constraints successfully met by the model’s response.

To evaluate if a constraint is fulfilled by model M , we present the LLM judge J with the task description t_i , the model’s response $r_i = M(t_i)$, and the specific constraint under evaluation c_i^j . Then, we prompt the Judge with a yes/no question, “Given task t_i and response r_i , is the following con-

straint satisfied: c_i^j ?”. We denote the judge score by $J(t_i, r_i, c_i^j) \in \{0, 1\}$. Its value is 1 if the judge responds with a “yes” token, and 0 if responds with a “no” token, in a greedy decoding setup to ensure consistency.

The *soft* and *strict* scores for a task are defined as follows:

$$soft(r_i | t_i) = \frac{1}{N(t_i)} \sum_{j=1}^{N(t_i)} J(t_i, r_i, c_i^j) \quad (1)$$

$$strict(r_i | t_i) = \prod_{j=1}^{N(t_i)} J(t_i, r_i, c_i^j) \quad (2)$$

where $N(t_i)$ is the number of constraints in t_i .

4.2 Experimental Setup

We evaluate 14 prominent instruction-tuned LLMs from five different model families on WILDIFEVAL, in a zero-shot setup. The models vary in size from 0.5 billion to 671 billion parameters.

We assess the following models: (1) Deepseek-v3 (Liu et al., 2024) (2) Mistral-Large-instruct-2407 (Mistral AI Team, 2024) (3) Gemma-2-2b and Gemma-2-9b (Team et al., 2024) (4) Llama3.2-1b, Llama3.2-3b, Llama3.1-8b, Llama3.3-70b and Llama3.1-405b (Dubey et al., 2024) (5) Qwen-2.5-0.5b, Qwen-2.5-1.5b, Qwen-2.5-3b, Qwen-2.5-7b, and Qwen-2.5-72b (Yang et al., 2024).

Judge evaluation As a judge model for evaluation (§4.1), we use Deepseek-v3. We choose Deepseek-v3 as the judge after evaluating a subset of 500 tasks from WILDIFEVAL with GPT-4o as a judge, and among available SOTA open-source

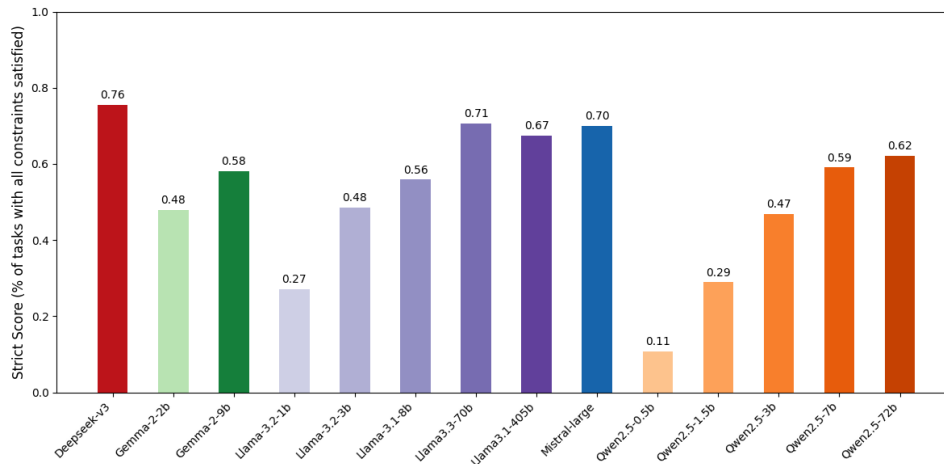


Figure 5: Strict scores on WILDIFEVAL. For each model, the figure reports the proportion of tasks in which all constraints were fulfilled (strict score). Soft scores are shown in Figure 10 in the Appendix. Statistical significance between model pairs (McNemar tests) is reported in Figure 13 in Appendix.

models including also Llama3.3-70b and Qwen-2.5-72b, Deepseek-v3 showed the highest agreement with GPT-4o, in terms of accuracy and confidence correlation (details in Appendix B.2). As a further validation of our evaluation, the benchmark shows significantly high Kendall’s Tau correlations (>0.82) with existing benchmarks like IFEval, MMLU, and GPQA (Appendix D).

4.3 Results

Figure 5 depicts the overall model performance on WILDIFEVAL. We can observe a clear performance gap within model families, with larger models consistently outperforming their smaller counterparts, in line with prior findings (Kaplan et al., 2020).⁷ At the same time, even stronger models like Deepseek-v3 and Llama3.3-70b fail to satisfy all task constraints in 25-30% of cases.

The best performing model is Deepseek-v3. Since it also serves as the judge, this raises questions about potential judge self-bias (Verga et al., 2024; Gera et al., 2024). We note that on a subset of 500 tasks used for judge validation (§4.2), all tested judges, i.e., GPT-4o, Llama3.3-70b, and Qwen-2.5-72b, consistently ranked Deepseek-v3 first, hinting that it is more than just self-bias.

Naturally, when a task has more constraints, it is harder for the model to fulfill all of them. Accordingly, Figure 6a shows the decrease in the strict performance score as a function of the number of constraints. However, when looking at the soft

performance score (Figure 6b) we see that the number of constraints does not affect the fulfillment of *individual* constraints. In other words, it appears that the difficulty in multi-constraint tasks does not reflect a general decrease in model instruction-following abilities, but rather stems from having to fulfill several constraints at once.

Figure 7a illustrates the relative model performance for different constraint types. We can see that models consistently have difficulties with *Length* constraints, and to a lesser extent also with *Format and Structure*. In contrast to these form-based types, models tend to succeed in fulfilling *Focus / Emphasis* constraints, which impose softer, content-related requirements. We observe a somewhat different pattern for models from the Qwen family, that appear to struggle more with *Persona* and *Style* constraints relative to other models.

To further understand the role of constraint types, we look at the rankings they induce of model performance. We rank the models according to their performance on each constraint type, and calculate the agreement between the resulting model rankings. As Figure 7b shows, type-specific rankings largely agree with each other. We do however observe different degrees of agreement. High correlations between categories suggest that constraints may probe related underlying capabilities. For instance, *Include / Avoid* constraints correlate with *Editing* ones, possibly reflecting a shared reliance on surface-form control while preserving or modifying meaning. Conversely, low correlations imply distinct aspects of model behavior. Notably, *Length*

⁷A notable exception is Llama3.3-70b, that surpasses Llama3.1-405b. This result is aligned with previous reports.

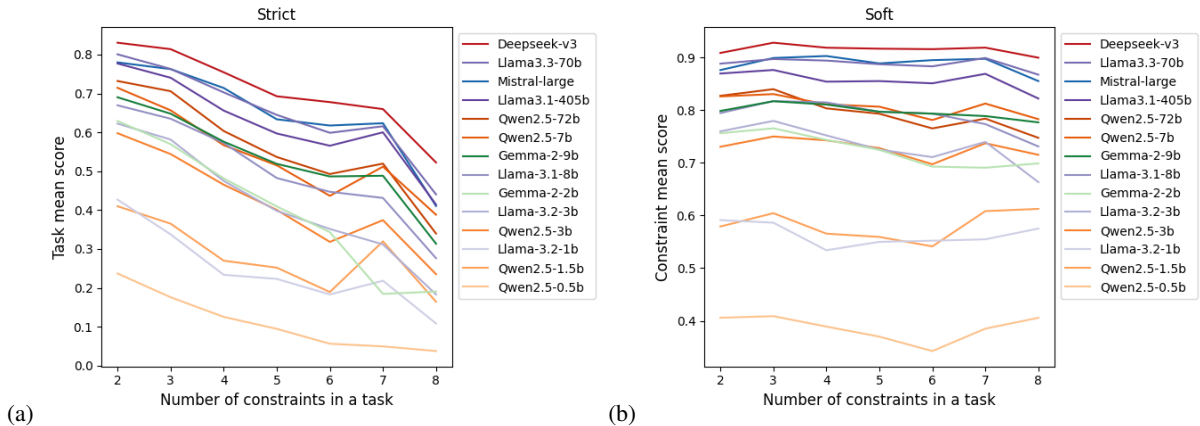


Figure 6: Scores as function of number of constraints in a task. (a) Strict score – tasks in which all constraints are fulfilled. (b) Soft score – fraction of fulfilled constraints in a task.

constraints induce a ranking different from other types, especially *Persona and Style*. While not conclusive, this analysis offers a useful exploratory view into links between different skill dimensions.

Error analysis. We also performed a manual analysis of the examples where most models failed to satisfy the constraints. We observe that the majority of these failure cases belong to the *Length* category, particularly constraints requiring an exact number of words or more atomic units (syllables, characters etc.), e.g., “The script should be 300 words long”. Some of the failure cases involve constraints that are quite complex, involving multiple specifications and sub-constraints. For example, the user constraint can require including a dictionary in a specific format and with a specific set of keys and values. Overall, we note that all constraint types can vary widely in the level of complexity they impose on the model. For example, *Persona and Style* constraints range from mundane requirements (“Use a first-person perspective.”, “Keep the tone informal.”) to more specific an esoteric ones (“Excel in ninjutsu, tactics, and battle strategies”, “Use strict iambic pentameter”).

Length constraints validation. Since length constraints can often be verified heuristically, and one might argue that heuristic measures are more suitable than an LLM-as-a-judge approach, we conducted an additional analysis comparing the accuracy of the LLM judge against a simple word-count heuristic. We identified 700 constraints specifying length in words (e.g., “up to 500 words”), extracted via regular expressions. The heuristic method counted words using whitespace separa-

tion, and we measured its agreement with the LLM judge. For our main judge, Deepseek-v3, agreement reached 86.66%. Similar levels were observed for other judges, including Llama3.1-405b (79.23%), Llama3.3-70b (83.95%), and Mistral-Large (81.95%). While heuristic methods offer a simple and transparent baseline, they may fail in natural text with complex formatting or phrasing. Our findings suggest that LLM judges are generally reliable for length-related constraints, though further improvements remain possible.

5 Related Work

Recent interest in LLM instruction-following capabilities raises the need for benchmarking model performance under complex, multi-constraint scenarios (Lin et al., 2020; Sun et al., 2023).

Several works (Yao et al., 2023; Bastan et al., 2023; Iso, 2024) rely on synthetic instructions and rule-based evaluation, with the prominent example of *IFEval* (Zhou et al., 2023). Other works, such as *FollowBench* (Jiang et al., 2024) and *InfoBench* (Qin et al., 2024), utilize crowd-sourced data, and LLM-based evaluation. However, these works are limited in size and do not fully capture the diversity of genuine user inputs. More recently, REALINSTRUCT (Ferraz et al., 2024) employs real-user instructions but has not been released; beyond release, our work additionally contributes a constraint taxonomy and the first large-scale analysis of constraint distribution, co-occurrence, and lexical long-tail patterns in naturally occurring user instructions. Complementary efforts in other languages, such as *CFBench* (Zhang et al., 2024) in Chinese, derive constraints from real-world scenarios but reflect

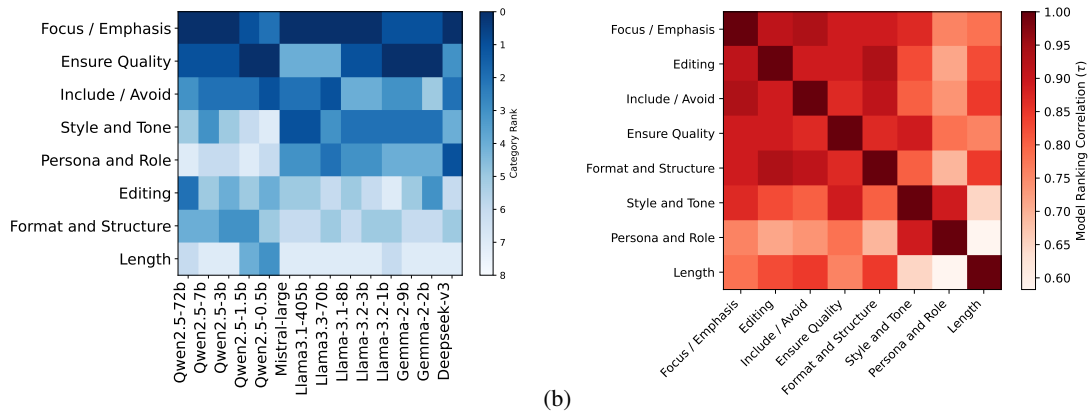


Figure 7: Constraint types characteristics. (a) Category performance rankings per model. Darker colors indicate stronger performance by the model on the corresponding constraint category, while lighter colors reflect weaker performance. (b) Correlation (Kendall’s Tau) between model rankings induced by different constraint types.

different cultural and linguistic constraint patterns, and are therefore complementary rather than directly comparable to our English-focused setting. Other benchmarks address orthogonal scenarios, including multi-turn system-message following (Qin et al., 2025) and agentic instruction-following (Qi et al., 2026), whereas WILDIFEVAL targets single-turn constrained text generation.

In this work, we release a diverse dataset of multi-constraint instructions, that originates from real users and is much larger than all existing English datasets. Moreover, whereas some of these benchmarks have become saturated, ours remains challenging even for state-of-the-art LLMs.

6 Discussion

In this work, we presented a benchmark for evaluating the ability of LLMs to follow real-world multi-constrained instructions. WILDIFEVAL aims to capture a realistic and up-to-date view of constrained generation user requests.

Our analyses further reveal insights into the structure of real user instructions. Examining the distribution of constraint types highlights which capabilities are most demanded in practice, while analyzing their co-occurrence sheds light on how complex instructions are composed. Together, these findings can guide the development and evaluation of models that better reflect and address real user needs.

Limitations

Our work has several limitations that warrant consideration. First, the dataset consists solely of in-

structions from users of the Chatbot Arena (Chiang et al., 2024) platform. Thus, it reflects the types of tasks that interest the platform users, and may not be fully representative of all LLM usage scenarios. Moreover, this may introduce a demographic bias, limiting the representativeness with respect to the general population. Hence, this may affect the generalizability of our findings.

Second, evaluating some of the constraints in the dataset is quite challenging. Many constraints are inherently subjective, e.g., “the story needs to be suited to a nine-year-old”; this may introduce some noise or bias into the evaluation process.

Third, despite our efforts to filter out noise and toxic language, some instances may still remain. These imperfections could introduce unintended biases and complicate the interpretation of LLM performance under realistic conditions.

Finally, our focus in WILDIFEVAL is on the model’s ability to satisfy the given constraints, rather than directly evaluating the task itself. However, in many cases, the distinction between a constraint and the actual task is somewhat vague. As a result, during decomposition, some constraints may closely reflect the task itself, ultimately contributing to the final score.

These limitations highlight important areas for future research and emphasize the need for continued refinement in both dataset construction and evaluation methodologies.

References

Mohaddeseh Bastan, Mihai Surdeanu, and Niranjana Balasubramanian. 2023. [NEUROSTRUCTURAL](#)

- DECODING: Neural text generation with structural constraints.** In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9496–9510, Toronto, Canada. Association for Computational Linguistics.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. **Chatbot arena: An open platform for evaluating LLMs by human preference.** In *Forty-first International Conference on Machine Learning*.
- Kenneth Ward Church and Patrick Hanks. 1990. **Word association norms, mutual information, and lexicography.** *Computational Linguistics*, 16(1):22–29.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. **The llama 3 herd of models.** *arXiv preprint arXiv:2407.21783*.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.
- Liat Ein-Dor, Orith Toledo-Ronen, Artem Spector, Shai Gretz, Lena Dankin, Alon Halfon, Yoav Katz, and Noam Slonim. 2024. **Conversational prompt engineering.** *arXiv preprint arXiv:2408.04560*.
- Thomas Palmeira Ferraz, Kartik Mehta, Yu-Hsiang Lin, Haw-Shiuan Chang, Shereen Oraby, Sijia Liu, Vivek Subramanian, Tagyoung Chung, Mohit Bansal, and Nanyun Peng. 2024. **LLM self-correction with DECRIM: Decompose, critique, and refine for enhanced following of instructions with multiple constraints.** In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7773–7812, Miami, Florida, USA. Association for Computational Linguistics.
- Cristina Garbacea and Qiaozhu Mei. 2022. **Why is constrained neural language generation particularly challenging?** *arXiv preprint arXiv:2206.05395*.
- Ariel Gera, Odellia Boni, Yotam Perlitz, Roy Bar-Haim, Lilach Eden, and Asaf Yehudai. 2024. **JustRank: Benchmarking llm judges for system ranking.** *arXiv preprint arXiv:2412.09569*.
- Junxian He, Wojciech Kryscinski, Bryan McCann, Nazneen Rajani, and Caiming Xiong. 2022. **CTRL-sum: Towards generic controllable text summarization.** In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5879–5915, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Hayate Iso. 2024. **AutoTemplate: A simple recipe for lexically constrained text generation.** In *Proceedings of the 17th International Natural Language Generation Conference*, pages 1–12, Tokyo, Japan. Association for Computational Linguistics.
- Yuxin Jiang, Yufei Wang, Xingshan Zeng, Wanjun Zhong, Liangyou Li, Fei Mi, Lifeng Shang, Xin Jiang, Qun Liu, and Wei Wang. 2024. **FollowBench: A multi-level fine-grained constraints following benchmark for large language models.** In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4667–4688, Bangkok, Thailand. Association for Computational Linguistics.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:2405.01535*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, and 1 others. 2024. Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*.
- Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024. Nv-embed: Improved techniques for

- training llms as generalist embedding models. *arXiv preprint arXiv:2405.17428*.
- Cheng Li, Mingyang Zhang, Qiaozhu Mei, Weize Kong, and Michael Bendersky. 2024a. [Learning to rewrite prompts for personalized text generation](#). In *Proceedings of the ACM Web Conference 2024, WWW '24*, page 3367–3378, New York, NY, USA. Association for Computing Machinery.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. 2024b. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*.
- Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. 2020. [CommonGen: A constrained text generation challenge for generative commonsense reasoning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1823–1840, Online. Association for Computational Linguistics.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. [Deepseek-v3 technical report](#). *arXiv preprint arXiv:2412.19437*.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Mistral AI Team. 2024. Large enough. <https://mistral.ai/en/news/mistral-large-2407>. Accessed: 2025-02-14.
- Yotam Perlitz, Ariel Gera, Ofir Arviv, Asaf Yehudai, Elron Bandel, Eyal Shnarch, Michal Shmueli-Scheuer, and Leshem Choshen. 2024. Do these llm benchmarks agree? fixing benchmark evaluation with benchbench. *arXiv preprint arXiv:2407.13696*.
- Yunjia Qi, Hao Peng, Xiaozhi Wang, Amy Xin, Youfeng Liu, Bin Xu, Lei Hou, and Juanzi Li. 2026. [AGEN-TIF: Benchmarking large language models instruction following ability in agentic scenarios](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Yanzhao Qin, Tao Zhang, Tao Zhang, Yanjun Shen, Wenjing Luo, sunhaoze, Yan Zhang, Yujing Qiao, weipeng chen, Zenan Zhou, Wentao Zhang, and Bin CUI. 2025. [Sysbench: Can LLMs follow system message?](#) In *The Thirteenth International Conference on Learning Representations*.
- Yiwei Qin, Kaiqiang Song, Yebowen Hu, Wenlin Yao, Sangwoo Cho, Xiaoyang Wang, Xuansheng Wu, Fei Liu, Pengfei Liu, and Dong Yu. 2024. [Infobench: Evaluating instruction following ability in large language models](#). *arXiv preprint arXiv:2401.03601*.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. *arXiv preprint arXiv:2311.12022*.
- Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2023. [LaMP: When large language models meet personalization](#). *arXiv preprint arXiv:2304.11406*.
- Jiao Sun, Yufei Tian, Wangchunshu Zhou, Nan Xu, Qian Hu, Rahul Gupta, John Wieting, Nanyun Peng, and Xuezhe Ma. 2023. [Evaluating large language models on controlled generation tasks](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3155–3168, Singapore. Association for Computational Linguistics.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. [Gemma 2: Improving open language models at a practical size](#). *arXiv preprint arXiv:2408.00118*.
- Laurens van der Maaten and Geoffrey Hinton. 2008. [Visualizing data using t-SNE](#). *Journal of Machine Learning Research*, 9(86):2579–2605.
- Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yixuan Su, Aleksandra Piktus, Arkady Arkhangorodsky, Minjie Xu, Naomi White, and Patrick Lewis. 2024. Replacing judges with juries: Evaluating llm generations with a panel of diverse models. *arXiv preprint arXiv:2404.18796*.
- Jiayin Wang, Fengran Mo, Weizhi Ma, Peijie Sun, Min Zhang, and Jian-Yun Nie. 2024. [A user-centric multi-intent benchmark for evaluating large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3588–3612, Miami, Florida, USA. Association for Computational Linguistics.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. [Qwen2.5 technical report](#). *arXiv preprint arXiv:2412.15115*.
- Shunyu Yao, Howard Chen, Austin W Hanjje, Runzhe Yang, and Karthik Narasimhan. 2023. [Collie: Systematic construction of constrained text generation tasks](#). *arXiv preprint arXiv:2307.08689*.
- Tao Zhang, Yanjun Shen, Wenjing Luo, Yan Zhang, Hao Liang, Fan Yang, Mingan Lin, Yujing Qiao, Weipeng Chen, Bin Cui, and 1 others. 2024. [CFbench: A comprehensive constraints-following benchmark for llms](#). *arXiv preprint arXiv:2408.01122*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric. P Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2023a. [Lmsys-chat-1m: A large-scale real-world llm conversation dataset](#). *Preprint*, arXiv:2309.11998.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023b. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. [Instruction-following evaluation for large language models](#). *arXiv preprint arXiv:2311.07911*.

A Prompts

Classify constrained generation tasks

You are an assistant whose job is to help me perform tasks. I need to filter from a set of requests made by users to AI assistants, the ones in which human requested the AI assistant to do a task with constraints to be follow. Constraints refer to more detailed rules, conditions or specific guidelines provided to guide the responses and shape the output generated by the AI assistant. Examples of sentences that indicate constraints are: “write in the format of”, “write as if you were”, “make sure to follow this”, “make sure to answer these questions”, “make sure to not include”, “avoid mentioning”. I will give you the human request and I expect you to answer “Yes” when the request contains instruction with constraints, or “No” if the request does not contemplate any constraint. I also want you to say “No” if the request require to generate code or an answer about code provided. Also, I want you to say “No” if the task is not self-contained, which means the AI Assistant need to ask follow up questions before start to answer, or it needs more context. You are provided five examples.

Example 1: list and compare top website to <https://fastfunnels.com/> in table format.

Answer: Yes

Example 2: You are an fantasy writer. Your task is now to help me write a D&D adventure for 5 players in the Eberron univers. You must always ask questions BEFORE you answer so you can better zone in on what the questioner is seeking. Is that understood ?

Answer: No.

Example 3: I have 100 dollars and would like to use this as the initial funding to make some money. I need it to be as quick as possible with good returns.

Answer: No.

Example 4: I have a vacation rental website and I am looking for alliterative and descriptive headlines that are at least 4 words in length and a maximum of 6 words. Examples: “Get Away to Galveston”, “Sleep Soundly in Seattle”. Each headline should have alliteration of at least 50% of the words and be poetic in language. Make each headline unique from the others by not repeating words. Each headline should include a verb. Put into an table with the city in column one and the results in column two for the following cities: Galveston, Sedona, Honolulu, Tybee Island, Buenos Aires.

Answer: Yes.

Example 5: pitch me a viral social app that is inspired by the hunger games. give it a fun twist!

Answer: Yes.

Request: `{request}`

Now please answer, “Yes” or “No”.

Answer:

Decompose Tasks

You are an assistant whose job is to help me perform tasks. I will give you an instruction that implicitly contains a task description, its context, and constraints to be followed. Your task is to translate this instruction in a more structured way, where task, context and constraints are separated. Avoid writing anything else. Context is an input text needed to generate the answer or a more detailed description of the situation. Make sure to separate the context when it is needed, otherwise leave it empty. You are provided five examples. Please follow the same format.

Example 1:

Original Instruction: Write me a rap about AI taking over the world, that uses slangs and young language. It need to sound like a real human wrote it. It would be cool if there's a chorus very catchy that would be singed by a famous pop artist. Make sure to include references about things that young people likes, such as memes, games, gossips. I want that in the end, you reveal that this was written by an AI.

Translated Task: Write a rap about AI taking over the world.

Translated Context:

Translated Constraints:

1. Use slang and youth language.
2. Make it sound like it was written by a real human.
3. The song may have a very catchy chorus, which would be sung by a famous pop artist.
4. Include references to things young people like, such as memes, games, gossip.
5. Reveal at the end that this rap was written by an AI.

Example 2: Original Instruction: write me a 5-page essay that is about travel to taiwan. detail description is below Topic : The Benefits of Traveling Sub Topic : Exposure to New Cultures Content 1 : Trying New Foods - I tried to eat Fried stinky tofu. smell was wierd but tasty was not bad. Content 2. : Exploring Historical Things - I saw Meat-shaped-stone in taipei museum. the stone was really like stone! it was surprising! Length : around 2000 words Assume that audience is collage student major in history. you can add historical events or news about what i experienced

Translated Task: Write an essay about traveling to Taiwan. The topic is "The Benefits of Traveling" and the subtopic is "Exposure to New Cultures".

Translated Context:

Translated Constraints:

1. Describe your experience of trying new foods, including your experience eating Fried stinky tofu (mention the peculiar smell but the tasty flavor).
2. Share your exploration of historical sites, with a specific mention of the Meat-shaped stone in the Taipei museum and your surprise at its appearance.
3. The essay should be approximately 2000 words in length, having around 5 pages.
4. Assume the audience is college students majoring in history, so you can incorporate historical events or news related to your travel experiences.

Example 3: Original Instruction: can you please write me a 150-word paragraph about epidermolysos bullosa which includes a basic description of clinical features and a summary of the most prevalent genetic causes. please make sure to include information on the inheritance pattern. please also write the paragraph in simple english that couldbe understand without a genetic or medical bacakground

Translated Task: Write a paragraph about Epidermolysis Bullosa.

Translated Context:

Translated Constraints:

1. Provide a description of clinical features.
2. Summarize the most common genetic causes.
3. Explain the inheritance pattern.
4. Ensure the paragraph is written in simple language for easy comprehension, even for those without a genetic or medical background.
5. The paragraph should be around 150 words in length.

Example 4: Original Instruction: write me a blog post that answers the following questions:What is the lifespan of a toaster? What toasters are made in the USA? What are the top 10 toasters? What is the difference between a cheap and expensive toaster? How much should you pay for a toaster? How often should toasters be replaced? Which toaster uses the least electricity? How many watts should a good toaster have? What is the warranty on Mueller appliances? Is Mueller made in China? Where are Mueller appliances manufactured?

Translated Task: Write a blog post about toasters.

Translated Context:

Translated Constraints:

1. Mention what is the lifespan of a toaster, and how often should toasters be replaced.
2. Mention what toasters are made in the USA.
3. Comment which are the top 10 toasters.
4. Explain the difference between a cheap and a expensive toaster.
5. Discuss prices, and how much should you pay for a toaster.
6. Compare toaster regarding electricity use, mentioning how many watts should a good toaster have.
7. State what is the warranty on Mueller appliances.
8. Answer where are Mueller appliances manufactured, and if Mueller is made in China.

Example 5: Original Instruction: Hi Michael, Hope you're well? Regarding my previous email to support HC with good price offers, What are your current needs? Hoping for your earliest reply. Thanks in advance, As a sales manager, the client hasn't replied this email after 2 days. Write a follow up email to the client. Your writing should include high complexity and burstiness. It must also be as brief as possible

Translated Task: A client hasn't replied the email below after 2 days. As a sales manager, write him a follow-up email.

Translated Context: "Hi Michael, Hope you're well? Regarding my previous email to support HC with good price offers, What are your current needs? Hoping for your earliest reply. Thanks in advance,"

Translated Constraints:

1. Include high complexity and burstiness in your writing.
2. Keep the email as brief as possible.

Original Instruction: \${instruction}

Translated Task:

Constraint Categorization

Classify the following constraint from a generation task into one of the categories listed below. Respond only with the category number. Do your best to match the constraint with an existing category. Only if you are certain that the constraint does not fit any of the categories from the list, you may respond with 'Other:' followed by a suggested title for an appropriate category.

Categories:

0. ***Style and Tone***: This category encompasses instructions that dictate the overall writing style, including formality, language register, emotional color, and imitation of specific authors or publications. It dictates the voice and feel of the output.

Examples:

- The writing style should emulate Ernest Hemingway's short, declarative sentences.
- Maintain a formal and professional tone throughout the email.
- Use a playful and whimsical tone to engage children.
- Write in a concise and technical style, suitable for a scientific paper.
- The language should be evocative and poetic, painting a vivid picture for the reader.

1. ***Include / Avoid***: This category specifies elements that should be either included or excluded from the response. This can involve mentioning or adding specific keywords, phrases, or concepts, or avoiding particular words and ideas. It concerns the content and its restrictions.

Examples:

- Include at least three examples of alliteration in the poem.
- Do not mention the specific brand name of the competitor.
- Include a call to action at the end of the blog post, encouraging readers to subscribe.
- Avoid using passive voice constructions.
- Include a summary of the key findings at the beginning of the report.

2. ***Format and Structure***: This category focuses on the organization and arrangement of the response. This includes instructions on using bullet points, tables, paragraphs, specific layouts, document structures or adhering to established formats. It dictates the physical form of the output.

Examples:

- Present the data in a clear and concise table format.
- Organize the information into five distinct paragraphs, each addressing a separate aspect of the topic.
- The report should follow the standard APA format, including citations and a bibliography.
- Create a numbered list of steps in the process.
- Each section should begin with a clear and informative heading.

3. ***Length***: This category defines constraints on the length of the response, whether in terms of word count, character count, sentence limit, or overall brevity. It sets the quantitative boundaries of the output.

Examples:

- The summary should be no more than 150 words.
- Each sentence should be kept under 20 words.
- Provide a short and sweet answer, within 50 characters.
- The article should be approximately 800-1000 words in length.
- The description should be exactly 10 words long.

4. ***Persona and Role***: This category instructs the AI to adopt a specific character, personality, or role in its response. This may involve imitating a particular person, acting as an expert in a field, or assuming a defined perspective. It defines the agent or narrator that provides the output.

Examples:

- Act as a seasoned travel blogger, providing tips and insights for visiting Rome.
- Respond as if you are a friendly and helpful chatbot, assisting users with their inquiries.
- Answer as a grumpy old man who is against modern technology.
- Speak as if you are Albert Einstein explaining relativity.
- Write the response from the point of view of a tree.

5. ***Focus / Emphasis***: This category highlights specific topics, aspects, or keywords that the response should concentrate on. It directs the AI's attention to certain elements and ensures that they are given prominence in the output.

Examples:

- Focus primarily on the economic impact of the new policy.
- Highlight the innovative features of the product and its benefits for the user.
- Emphasize the importance of teamwork and collaboration in achieving the project goals.
- The article should primarily focus on the advantages of using renewable energy sources.
- Prioritize the ethical implications of artificial intelligence in healthcare.

6. ***Ensure Quality***: This category instructs the AI to meet some desired quality characteristics in its response. These may be general or specific quality constraints, like truthfulness or coherence of the output.

Examples:

- Ensure the information provided is accurate and up-to-date.
- The response should be coherent, logical, and easy to understand.
- Present the information in a simple and detailed manner.
- Make sure the answer is not biased.
- Cover all the key details.

7. ***Editing***: This category focuses on modifications to an input text given by the user. The constraint specifies in what manner to change the input text, or which properties of the original text should be preserved.

Examples:

- Correct any grammatical errors in the provided text.
- Change all instances of passive voice to active voice.
- Ensure you preserve the meaning of the original sentence.
- Simplify the language in the document to make it more accessible to a wider audience.
- Shorten all sentences to 5 words.

Constraint: \${constraint}

Your response:

Extract Domains

Each of the following tasks can be associated with a specific domain. Generate a list of 10 domains that best represent the domains associated with the tasks. Output only the list of domains, with no prefix or suffix.

Here is the list of tasks:

`${tasks_batch}`.

List of 10 domains:

Combine Domains to a Single List

Summarize the following lists of domains into a single list of 20 domains. Output only the summarizing list of 20 domains without any prefixes or suffixes. Here are the lists of domains:

`${lists_of_domains}`

Domain Classification

You are given a generation task. Classify the domain of the task into one of the domains listed below. Respond only with the category number.

Domains:

1. Creative Writing
2. Chemical Industry
3. Education
4. Business
5. Technology
6. Healthcare
7. Marketing
8. Entertainment
9. Environmental Science
10. Psychology
11. Roleplaying
12. Science Fiction
13. Fantasy
14. Journalism
15. Law
16. Finance
17. Data Analysis
18. Artificial Intelligence
19. Language Translation
20. Gaming

Task: `${task}`

Your response:

B Complementary Materials

B.1 Technical Details for Reproducibility

Dataset Curation. For the initial filtering, we used Llama3.1-405b, accessed via IBM’s internal inference infrastructure. Since we only analyzed the distribution of positive and negative token probabilities for classification, the results were unaffected by decoding temperature or other generation parameters. For the decomposition step with GPT-4o, we used a decoding temperature of 1 and a maximum token limit of 500, keeping all other parameters at their default values. The estimated cost for GPT-4o usage was approximately \$130.

Model Inference. We distinguish between two tiers of models: smaller models with fewer than 9B parameters and larger models with more than 70B parameters. Smaller models were run locally using 1–2 A6000 GPUs, depending on availability. Larger models were accessed via IBM’s internal inference infrastructure. All models generated responses with a temperature of 0.7 to encourage creativity, a maximum token limit of 1000, and default values for all other parameters. Inference was performed using vLLM (Kwon et al., 2023).

Judge Evaluation. We ran the Deepseek-v3 judge model on IBM’s internal inference infrastructure. As in the initial dataset filtering, our yes/no classification relied on the distribution of positive and negative next-token probabilities, making the results independent of the model’s decoding temperature.

B.2 LLM-Based Evaluation

Recently, LLM as a Judge (LLMaaJ) has become a standard evaluation method (Zheng et al., 2023b; Liu et al., 2023). Subsequent studies have demonstrated a strong correlation between LLM-based and human judgments (Kim et al., 2024), along with benchmarks assessing the reliability of LLM judges themselves (Gera et al., 2024; Lambert et al., 2024). This has led to the emergence of several benchmarks that rely on LLMaaJ, including MT-Bench (Zheng et al., 2023b), AlpacaEval (Dubois et al., 2024), and Arena-Hard (Li et al., 2024b). In this work, we leverage LLMaaJ alongside a fine-grained decomposition of the constrained generation task into individual constraint evaluations.

Choosing the right judge. While GPT-4o is arguably the strongest judge model, budget constraints due to the scale of WILDIFEVAL necessitated the use of an open-source alternative. To select the most reliable one, we evaluated a subset of 500 tasks using GPT-4o to produce reference judgments for the top-performing models. We then compared three open-source judge candidates—Deepseek-v3, Llama3.3-70b, and Qwen-2.5-72b—using two metrics: (1) binary agreement on constraint scores, and (2) covariance in the confidence of positive/negative judgments. Across both metrics, Deepseek-v3 exhibited the highest alignment with GPT-4o, and was thus chosen as our judge model.

B.3 Lexical Diversity of Constraints

In Figure 12 we can see a similar pattern to the one presented in Figure 4. We can see that “Provide” and “Write” are very frequent verbs. Alongside these, the figure reveals a significant presence of other highly frequent verbs such as “Be”, “Is”, “Do”, and “Are”. These typically function as **auxiliary verbs** (e.g., for forming tenses, voice, or questions) or **copular verbs** (linking subjects to attributes), playing grammatical roles rather than conveying specific lexical meaning. Similarly, several mid-frequency verbs remain, “Keep,” and “Identify,”.

The “Other” category is now much larger, with (34.5%), reflecting that the long tail of the verb distribution is much longer when examining all verbs.

B.4 Extracting Task Domains.

We extract the most prominent domains of WILDIFEVAL’s tasks via a three-step process, leveraging Llama3.3-70b. First, we prompt the model with batches of 100 tasks at a time, asking the model to extract the list of the domains they cover. Then, given all generated lists, we prompt the LLM to provide a set of the 20 most dominant domains in the data. Finally, we ask the model to classify all tasks in the dataset into these domains. Prompts are provided in Appendix A.

B.5 Human Validation

To validate our automatic curation pipeline, we conducted two human annotation studies, complementing the high-level summary in §2.2 with full details and analysis.

Decomposition quality. We sampled 100 tasks uniformly at random from WILDIFEVAL and asked an in-house annotator to evaluate the GPT-4o-produced decomposition of each task into individual constraints along three dimensions:

- **Correctness:** The extracted constraints faithfully reflect the requirements of the original task.
- **Completeness:** The decomposition captures *all* essential constraints in the task.
- **Independence:** The extracted constraints are distinct from one another and self-sufficient (i.e., understandable without reference to other constraints).

Each dimension was rated on a 1–5 Likert scale, where 5 indicates the highest quality. The mean scores were 4.71 for correctness, 4.64 for completeness, and 4.77 for independence, indicating that the automatic decomposition consistently produces high-quality constraint lists across all three dimensions.

Filtering threshold validity. A natural concern with our top-10% certainty filtering (§2.1) is whether it preferentially retains *simpler* constrained-generation tasks while discarding subtler ones. To assess this, we sampled 100 tasks uniformly across all certainty levels, and asked an in-house annotator to label each as either constrained-generation or not. We then compared these human labels against the binary outputs of our certainty-based filter.

Several findings emerged. First, human-model agreement at our chosen threshold was 75.8%. We further computed the optimal threshold (i.e., the one that maximizes human-model agreement on this sample), which yielded 77.9% agreement – only a marginal improvement over our chosen value. Given the scale of WILDIFEVAL, this difference does not materially affect the overall task distribution.

Second, humans labeled 31% of sampled tasks as constrained-generation, whereas our top-10% filter retained only 10% of tasks. This reflects an intentionally conservative bias in our pipeline: we prioritize *certainty* that a retained task is genuinely constrained generation, rather than maximizing coverage.

Third, the error distribution was strongly asymmetric. Of the disagreements, 22 tasks were labeled as constrained-generation by humans but not by the model (false negatives), while only 1 task was labeled as constrained-generation by the model alone (false positive). This asymmetry indicates that our filter is high-precision: the tasks that pass the threshold are reliably constrained-generation, while the cost is that a fraction of valid tasks are excluded. We view this as a desirable trade-off for an evaluation benchmark, where the validity of included tasks matters more than completeness of coverage.

Taken together, these results indicate that our filtering procedure did not selectively retain “simple” constraints, but instead produced a high-precision subset of constrained-generation tasks, consistent with the design goal of isolating this task type with high certainty.

C Empirical Justification of Dataset Scale

A natural question for any large-scale benchmark is whether its full size is necessary, or whether a smaller subset would yield comparable signal. To answer this empirically, we measure how the stability of the benchmark improves with dataset size along three complementary axes: (1) variance of per-model mean scores under resampling, (2) agreement of subsample-induced model rankings with the full-data ranking, and (3) the ability to reliably resolve close model pairs. We also examine per-category coverage, which is critical for the fine-grained analyses enabled by WILDIFEVAL.

Resampling protocol. For each subset size $N \in \{100, 250, 500, 1000, 2000, 3000, 5000, 7523\}$, we draw 50 subsamples of N tasks without replacement, recompute each model’s mean strict and soft scores, and compare the resulting model ranking to the full-data ranking via Kendall’s τ .

Category	Full	N=500	N=1000	N=2000	N=5000
Include / Avoid	5094	339	677	1354	3386
Format and Structure	2566	171	341	682	1705
Focus / Emphasis	2079	138	276	553	1382
Length	1987	132	264	528	1321
Style and Tone	1811	120	241	481	1204
Persona and Role	1392	93	185	370	925
Ensure Quality	1175	78	156	312	781
Editing	642	43	85	171	427

Table 2: Per-category task coverage at different subset sizes. Counts at smaller N are expected values under uniform random subsampling.

Score variance shrinks rapidly with N . Figure 15 shows the standard deviation of model mean scores across the 50 resamples (averaged across the 14 models). For the strict score, variance drops from 0.045 at $N=100$ to 0.013 at $N=1000$ and 0.009 at $N=2000$; soft-score variance follows a similar trajectory at roughly 60% the magnitude.

Ranking stability saturates near the full size. Figure 16 reports the mean Kendall’s τ between subsample-induced rankings and the full-data ranking. By $N=500$, τ already reaches 0.95 (strict) and 0.94 (soft); by $N=2000$ it exceeds 0.98. Overall model ordering is therefore recoverable from substantially smaller subsets.

Resolving close model pairs requires scale. However, mean ranking agreement obscures the fact that several top-performing models score within a narrow band. Figure 17 reports the 95% confidence interval width for the top-vs-second-model score gap across resamples. At $N \leq 1000$, the CI width for the strict gap is ≥ 0.036 – comparable in magnitude to the actual gap between several adjacent model pairs in our full results – meaning small subsets cannot reliably distinguish closely-matched models. The CI width only shrinks below 0.02 once N exceeds roughly 2000.

Per-category coverage motivates the full scale. A central contribution of WILDIFEVAL is the fine-grained breakdown of performance by constraint category (§4) and by co-occurrence patterns (§3). These analyses require sufficient task coverage *within* each category. Table 2 reports the number of tasks containing at least one constraint of each category, both in the full dataset and in expectation under uniform subsampling.

The rarest categories drive the requirement for scale. *Editing* appears in only 8.5% of tasks (642 in total); under uniform subsampling, even $N=5000$ yields fewer than 500 Editing-containing tasks. *Ensure Quality* (15.6%) and *Persona and Role* (18.5%) similarly require $N \geq 3000$ to reach reliable per-category sample sizes. The pairwise co-occurrence analyses (Figure 3) have even sparser support, further motivating the full scale.

Summary. The benchmark’s full size is not necessary for recovering a coarse model ranking, but is required for (i) discriminating between closely-matched models and (ii) reliable per-category and co-occurrence analyses – the principal analytic contributions of WILDIFEVAL.

D Correlation Analysis with Existing Benchmarks

Following Perlitz et al (2024) (Perlitz et al., 2024) we report Kendall’s Tau correlation (τ) results between our benchmark and several established benchmarks: IFEval (Zhou et al., 2023), GPQA (Rein et al., 2023), ARC-C (Clark et al., 2018), MMLU (Hendrycks et al., 2020), and HumanEval (Chen et al., 2021). We collect benchmark results from model cards and model papers (Liu et al., 2024; Dubey et al., 2024).⁸

⁸Qwen2.5 Model Card

We note that the corresponding evaluation setups may not be identical, introducing some noise into this analysis; we made every effort to ensure that the evaluation setups are consistent.

The analysis reveals strong positive correlations ($\tau > 0.8$, $p < 0.05$ in all cases) between our benchmark and each of the existing benchmarks, indicating a substantial alignment in their assessment of model performance. Specifically, the correlation with IFEval is 0.9, indicating a strong similarity with its assessment. Moreover, the Kendall's Tau correlations were 0.93 with GPQA, 0.82 with ARC-C, 0.96 with MMLU, and 0.87 with HumanEval, demonstrating that WILDIFEVAL effectively captures similar model capabilities as these well-established evaluations as well.

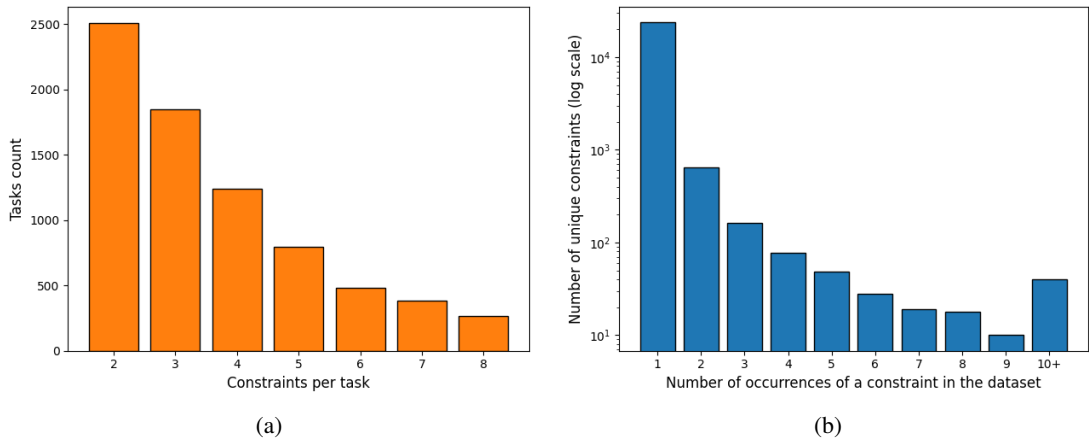


Figure 8: Analysis of constraints in WILDIFEVAL. (a) Distribution of the number of constraints per task. This histogram shows how many constraints are typically assigned to individual tasks. (b) Frequency of unique constraints across the dataset. This plot illustrates how often each distinct constraint appears in different tasks.

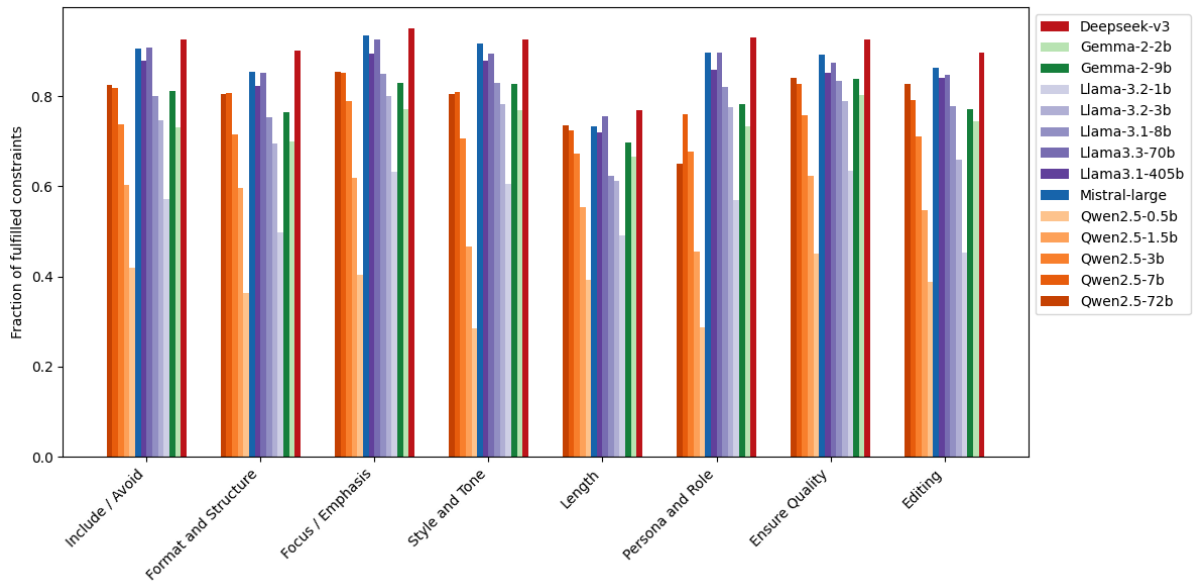


Figure 9: Mean constraint-following performance, by constraint category.

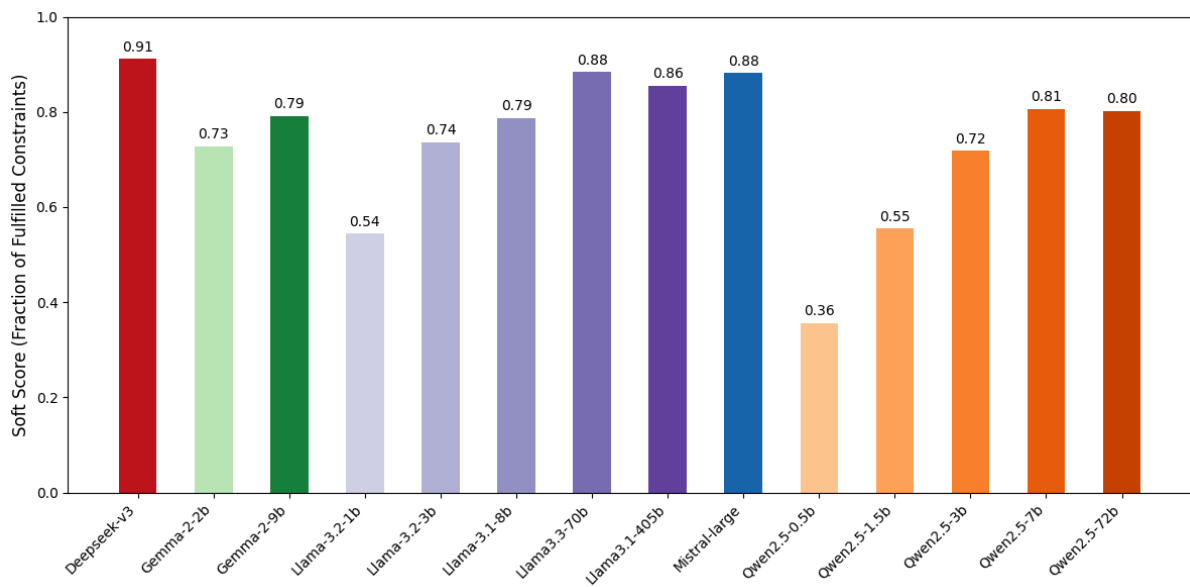


Figure 10: Soft scores on WILDIFEVAL. Soft scores represent the fraction of fulfilled constraints per task. Statistical significance between models is assessed via pairwise paired t-tests, shown in Figure 14.

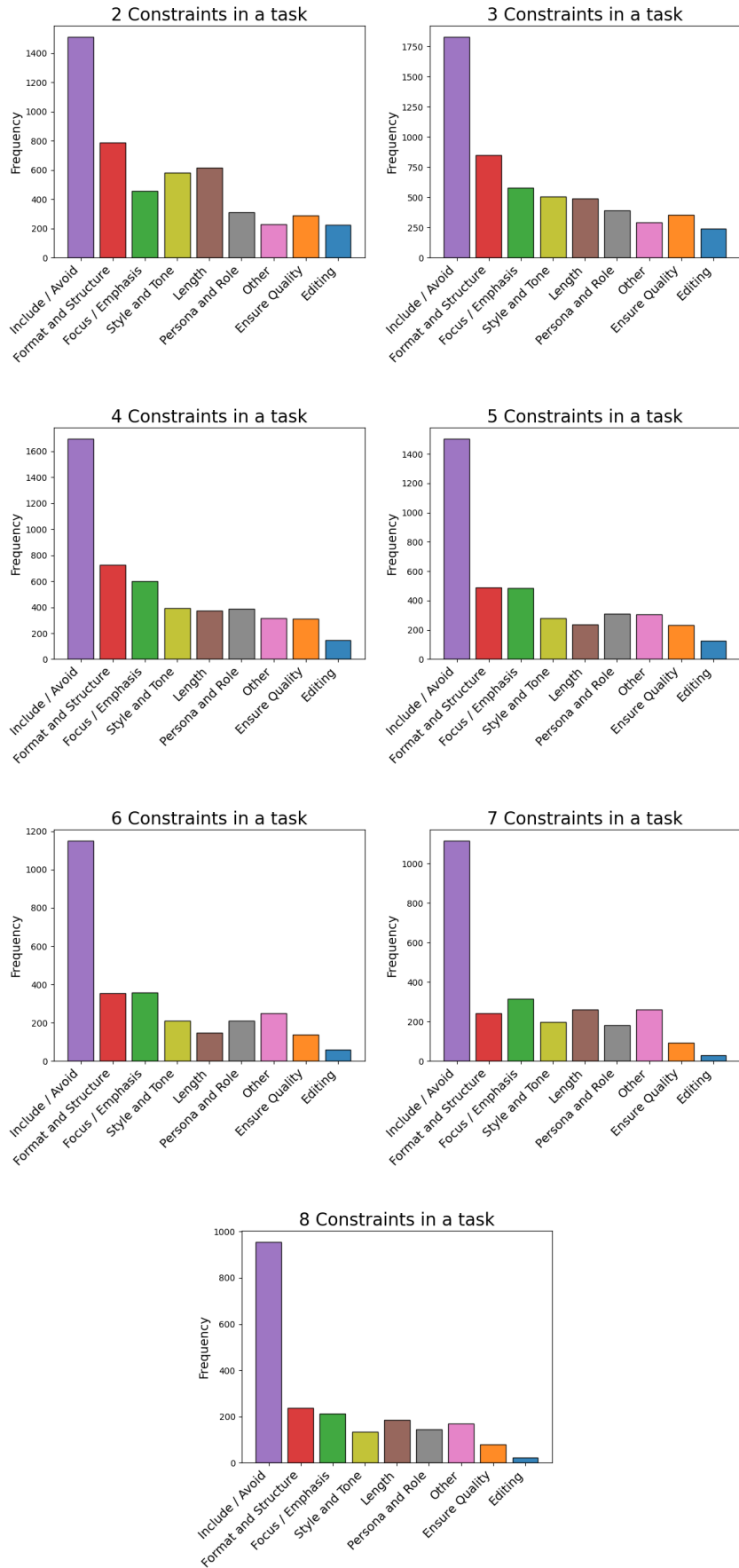


Figure 11: Distribution of constraint types, for tasks with different numbers of constraints.

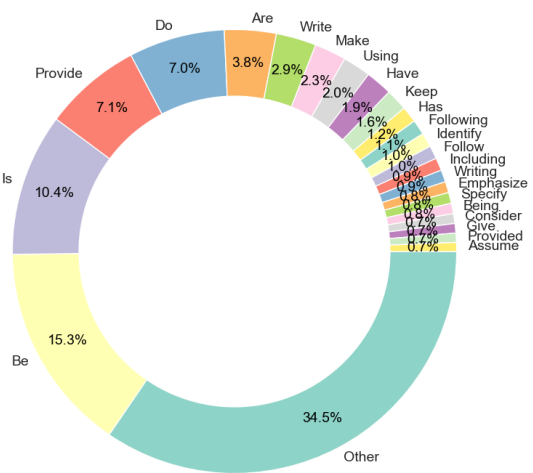


Figure 12: Constraints lexical diversity - distribution of verbs.

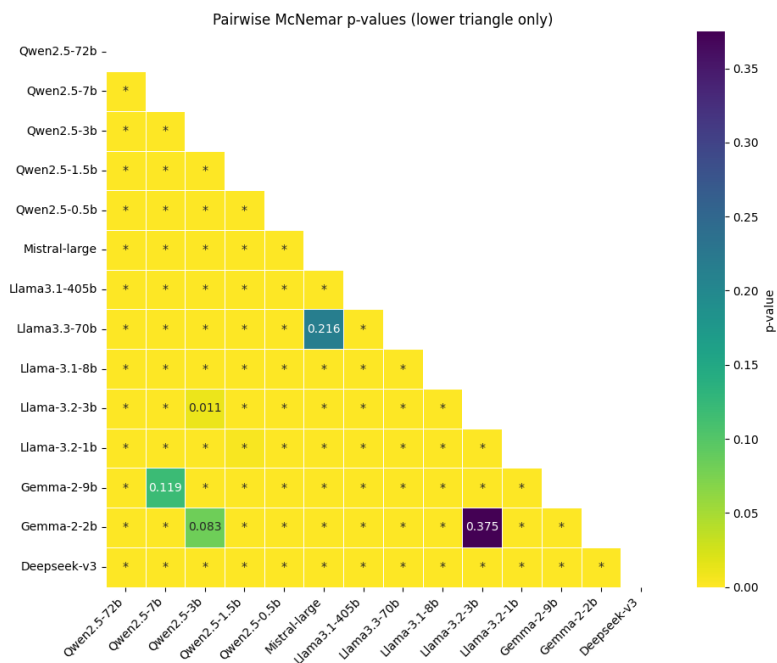


Figure 13: Pairwise McNemar p-values comparing model strict scores across tasks. Only the lower triangle is shown. Each cell reports the p-value of a McNemar test comparing the binary outputs of two models. Cells marked with * indicate statistically significant differences at $p < 0.01$.

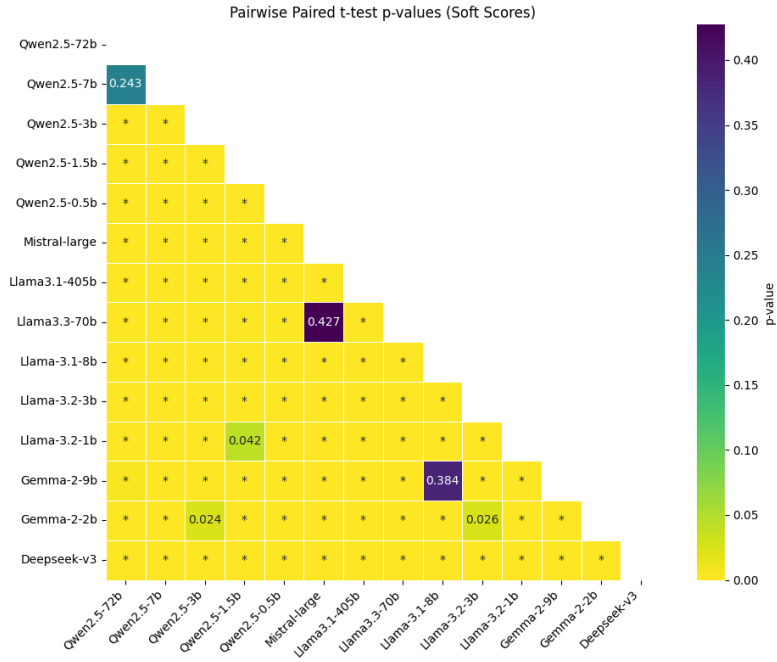


Figure 14: Pairwise paired t-test p-values comparing model soft scores across tasks. Only the lower triangle is shown. Each cell reports the p-value of a paired t-test comparing the soft scores of two models across the same set of tasks. Cells marked with * indicate statistically significant differences at $p < 0.01$.

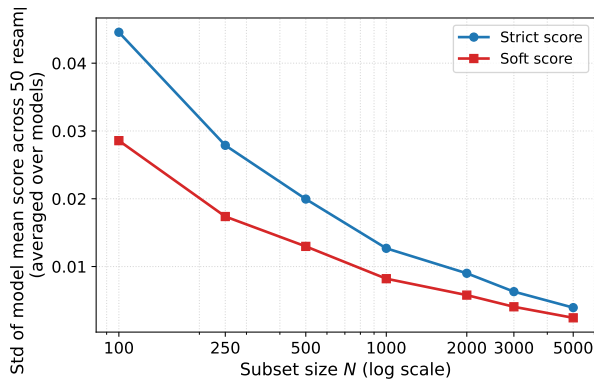


Figure 15: Standard deviation of model mean scores across 50 resamples, averaged over the 14 models, as a function of subset size N .

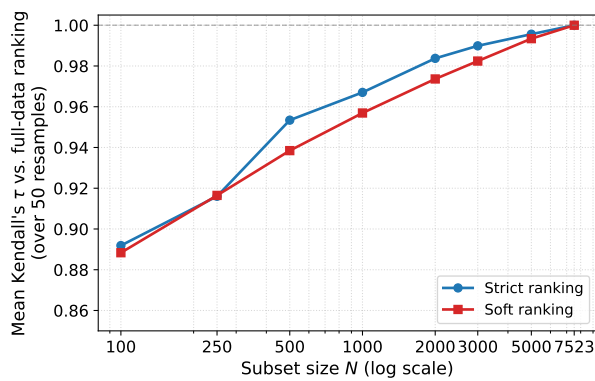


Figure 16: Mean Kendall's τ between subsample-induced model rankings and the full-data ranking, over 50 resamples per N .

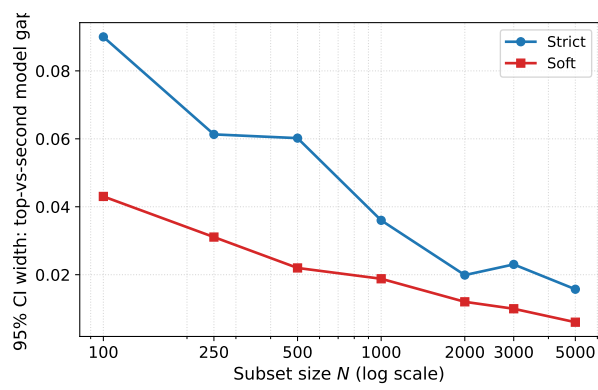


Figure 17: 95% confidence interval width for the top-vs-second-model score gap, as a function of subset size N .

E Examples from WILDIFEVAL

Below we include some instances from WILDIFEVAL. These examples demonstrate the diversity and complexity of the data in terms of tasks, domains and constraint types. They also illustrate that the precise division into constraints and their classification into types is not always straightforward and clear-cut.

```
[
  {
    "task": "Write me a poem about a puppy who is nervous to be adopted, but ends up loving his family.
    ↪ It should be 16 lines long. Mention the puppy's black spots and include at least two lines of
    ↪ dialogue from his new family.",
    "domain": "Creative Writing",
    "total_num_constraints": 3,
    "constraints": {
      "The poem should be 16 lines long.": "Length",
      "Mention the puppy's black spots.": "Include / Avoid",
      "Include at least two lines of dialogue from his new family.": "Include / Avoid"
    }
  },
  {
    "task": "Improve the following text and change 75% of the words. Keep sentences as short as
    ↪ possible \"Stop waking up and immediately getting on your phone.\n\nEven I notice a
    ↪ difference in how my brain feels.\n\nUnderstand this and prosper.\",
    "domain": "Creative Writing",
    "total_num_constraints": 2,
    "constraints": {
      "Ensure that 75% of the words are changed.": "Editing",
      "Maintain short sentences.": "Length"
    }
  },
  {
    "task": "You are a yoga coach. Your student has made the following mistakes when performing the
    ↪ warrior one pose:\n- the spine is not straight\n- your arms are not straight up\n- knees not
    ↪ directly over ankles\nPoint these problems out to your student and talk about how to improve
    ↪ on these aspects in a professional and encouraging way. Remember to act as the yoga coach.
    ↪ Mention every point in the provided list. Do not mention new mistakes other than the ones
    ↪ provided in the above list. Speak directly to your student.",
    "domain": "Education",
    "total_num_constraints": 5,
    "constraints": {
      "Act as a yoga coach.": "Persona and Role",
      "Identify the specific mistakes made: spine not straight, arms not straight up, and knees not
      ↪ directly over ankles.": "Editing",
      "Offer professional and encouraging suggestions for improvement on each aspect.": "Style and
      ↪ Tone",
      "Do not mention any mistakes other than those listed.": "Include / Avoid",
      "Speak directly to the student.": "Persona and Role"
    }
  },
  {

```

```

"task": "Do not paraphrase. For each restaurant in the article, get the name and the first 3
↳ sentences of the description verbatim using this format:\nRestaurant name: ...\nDescription:
↳ ...\n\nRestaurant name: ...\nDescription: ..."\n\nArticle:\nTitle - Best restaurants in
↳ Hanoi, Vietnam\nText - Search\n* Top\n* Sights\n* Restaurants\n* Entertainment\n*
↳ Nightlife\n* Shopping\nCTop ChoiceVietnamese in HanoiChim SaoSit at tables downstairs or grab
↳ a more traditional spot on the floor upstairs and discover excellent Vietnamese food, with
↳ some dishes inspired by the ethnic minorities of Vietnam's north. Definite standouts
↳ are...\nBTop ChoiceVietnamese in HanoiBun Cha 34Best NAME_1 in Vietnam? Many say 34 is up
↳ there. No presidents have eaten at the plastic tables, but you get perfectly moist
↳ chargrilled pork, zesty fresh herbs and delicious broth to dip everything in. The
↳ nem...\nVVegetarian in HanoiV's HomeBlink and you\u2019ll miss the slim alleyway opening
↳ leading to this excellent upstairs restaurant, with diners attended to by hearing- and
↳ speech-impaired staff. The relaxing space is elegant and charming, with a...\nKCafe in
↳ HanoiKotoRanging over four floors with a terrace and bar, this superb modernist
↳ cafe-bar-restaurant overlooking the Temple of Literature features neat interior design and
↳ exceptionally sweet staff, with daily specials...\nBVietnamese in HanoiBun NAME_2 LienBun
↳ NAME_2 Lien was launched into stardom thanks to NAME_3, who dined here with celebrity NAME_4
↳ in May 2016. Customers fill the four storeys to sample the grilled-pork-and-noodle
↳ delicacy...\nLTop ChoiceInternational in HanoiLa BadianeThis stylish bistro is set in a
↳ restored, whitewashed French villa arrayed around a breezy central courtyard. French cuisine
↳ underpins the menu \u2013 La Badiane translates as \u2013star anise\u2013 but Asian
↳ and...\nHTop ChoiceCafe in HanoiHanoi Social ClubOn three levels with retro furniture, the
↳ Hanoi Social Club is an artist hub and the city's most cosmopolitan cafe. Dishes include
↳ potato fritters with chorizo for breakfast, and pasta, burgers and wraps for...",
"domain": "Entertainment",
"total_num_constraints": 2,
"constraints": {
  "Use the format: \n \nRestaurant name: ... \n Description: ...": "Format and Structure",
  "Do not paraphrase the text.": "Editing"
}
},
{
"task": "Why do leaders with low education often fail to make the right decisions when formulating
↳ strategies? You should consider that the possible reason for lack of experience is not having
↳ the courage to step out of the comfort zone rather than being uneducated; the possible reason
↳ for lack of self-confidence is character factors rather than being uneducated, etc.",
"domain": "Education",
"total_num_constraints": 2,
"constraints": {
  "Consider lack of experience may stem from not having the courage to step out of the comfort
↳ zone rather than education level.": "Focus / Emphasis",
  "Consider lack of self-confidence may be due to character factors rather than education level.":
↳ "Focus / Emphasis"
}
},
{
"task": "Write a story where the Baywatch lifeguards NAME_1 NAME_2, NAME_3, NAME_4, NAME_5 and
↳ NAME_6 take part in fitness/bodybuildin contests. However the lifeguards have very different
↳ physiques and level of muscles. There are five main divisions in bodybuilding for women:
↳ Bikini, Figure, Physique, Bodybuilding and Fitness. In what divisions would the lifeguards
↳ be?",
"domain": "Entertainment",
"total_num_constraints": 3,
"constraints": {
  "Characters are NAME_1, NAME_2, NAME_3, NAME_4, NAME_5, and NAME_6.": "Include / Avoid",
  "Mention the five main divisions in bodybuilding for women: Bikini, Figure, Physique,
↳ Bodybuilding, and Fitness.": "Include / Avoid",
  "Assess which division each lifeguard would participate in based on their physique and level of
↳ muscles.": "Include / Avoid"
}
},
{

```

```

"task": "\role\":"You are a researcher who is good at summarizing papers using concise
↳ statements\
\n\ninstruction\":"Summarize the two paper reviews have been provided below in
↳ \input_data\ and generate a new review. The point is to combine the two into one
↳ literature review. Summarize according to the following four points: research background ,
↳ the problems , research methods research results.\n\n Output type \":(1) [research
↳ background] (2) [problems ](3) [research methods] (4) [research results]\nPlease note that
↳ your literature review should not exceed 150 words. \nNAME_1 your statements as concise and
↳ academic as possible. \n\ninput_data\":"1.(1) The research background of these papers
↳ includes evaluating the performance of articles using data from CNN's Quantitative State
↳ Methodology, improving the automation of meta-information derived in abstract, descriptive,
↳ and problem-solving environments, and developing an operational abstracting system.\n(2) The
↳ problems studied in these papers include comparing the performance of written sections,
↳ improving the automation of abstract meta-information, and developing an operational
↳ abstracting system.\n(3) The research methods proposed in these papers include using a score
↳ approach based on interconnected neural networks, a state-by-state scoring approach, and
↳ predicting performance using data from CNN's Quantitative State Methodology.\n(4) The
↳ research achievements in these papers include evaluating the performance of articles using
↳ data from CNN's Quantitative State Methodology, improving the automation of abstract
↳ meta-information, and developing an operational abstracting system.\n2.(1) Research
↳ background: The SALOMON system is designed to automatically summarize Belgian criminal cases
↳ by extracting relevant text, classifying it, predicting semantic relevance, and generating a
↳ case summary.\n(2) Problems studied: The study examines the challenges of summarization
↳ techniques and the difficulty of summarizing complex information.\n(3) Research methods: The
↳ paper uses an intelligent search engine to search for teaching resources and provides a
↳ comprehensive explanation of the search engine's principles and implementation steps.\n(4)
↳ Research results: The SALOMON system effectively summarizes criminal cases by extracting and
↳ classifying relevant text, predicting semantic relevance, and generating a case summary. The
↳ intelligent search engine in the paper improves the functionality of the search engine by
↳ enhancing its capabilities.",
"domain": "Education",
"total_num_constraints": 3,
"constraints": {
  "Address the four points: research background, the problems, research methods, and research
↳ results.": "Format and Structure",
  "Keep the literature review concise and academic.": "Length",
  "Ensure the literature review does not exceed 150 words.": "Length"
}
},
{
"task": "The following will act as a series of instructions/parameters to generate an
↳ individualized study plan for a single student.\n\nThe semesters comprising the study plan
↳ are Fall 2023, Spring 2024, Fall 2024, and Spring 2025.\n\nEach semester should contain
↳ exactly 4 courses.\n\nUse ONLY the following courses (each line represents an individual
↳ course) to populate the semesters exactly as they appear in this list:\nMATH 2415 Calculus I
↳ (4)\nBIO 3404 Anatomy & Physiology II (4)\nCPS 4150 Computer Arch. (3)\nMATH 2416 Calculus II
↳ (4)\nMATH 1054 Precalculus (3)\nCPS 3440 Analysis of Algorithms (3)\nMATH 3415 Calculus III
↳ (4)\nCOMM 1402 Speech Comm. (3)\nBIO 1400 General Biology II (4)\nCPS 3962 Object Oriented
↳ Analysis & Design (3)\nBIO 1300 General Biology I (4)\nCPS 2231 Computer Programming (4)\nCPS
↳ 4200 Systems Prog. (3)\nBIO 3403 Anatomy & Physiology I (4)\nCPS 1231 Fundamentals of CS
↳ (4)\nCOMM 3590 Business & Prof. Comm. (3)\n\nDo not include courses that do not appear in this
↳ list.\n\nDo not schedule the same course for more than 1 semester.\n\nTake into consideration
↳ the following:\nMATH 1054 Precalculus (3) is a prerequisite for MATH 2415 Calculus I
↳ (4)\nMATH 2415 Calculus I (4) is a prerequisite for MATH 2416 Calculus II (4)\nMATH 2416
↳ Calculus II (4) is a prerequisite for MATH 3415 Calculus III (4)\nCOMM 1402 Speech Comm. (3)
↳ is a prerequisite for COMM 3590 Business & Prof. Comm. (3)\nCPS 1231 Fundamentals of CS (4) is
↳ a prerequisite for CPS 2231 Computer Programming (4)\nBIO 1300 General Biology I (4) is a
↳ prerequisite for BIO 1400 General Biology II (4)\nBIO 1400 General Biology II (4) is a
↳ prerequisite for BIO 3403 Anatomy & Physiology I (4)\nBIO 3403 Anatomy & Physiology I (4) is a
↳ prerequisite for BIO 3404 Anatomy & Physiology II (4)\n\nPrerequisites must be scheduled at
↳ least 1 semester ahead of the courses that require them.\n\nPrerequisites cannot be scheduled
↳ for the same semester as the course that requires them.\n\nTake into consideration the
↳ following:\nCPS 4150 Computer Arch. (3) is only available during fall semesters.\nCPS 3440
↳ Analysis of Algorithms (3) is only available during fall semesters.\nCPS 3962 Object Oriented
↳ Analysis & Design (3) is only available during spring semesters.\nCPS 4200 Systems Prog. (3)
↳ is only available during spring semesters.\n\nGenerate final study plan",
"domain": "Education",
"total_num_constraints": 8,
"constraints": {
  "The study plan encompasses Fall 2023, Spring 2024, Fall 2024, and Spring 2025 semesters.":
↳ "Format and Structure",

```

```

    "Each semester should consist of exactly 4 courses.": "Length",
    "Use only the listed courses to fill the semesters, ensuring they appear exactly as listed.":
    ↪ "Include / Avoid",
    "Do not include courses not listed.": "Include / Avoid",
    "Avoid scheduling the same course across multiple semesters.": "Include / Avoid",
    "Maintain prerequisite courses at least 1 semester ahead of courses requiring them.": "Format
    ↪ and Structure",
    "Ensure prerequisites are not scheduled in the same semester as the courses requiring them.":
    ↪ "Include / Avoid",
    "Schedule courses according to availability: CPS 4150 and CPS 3440 are exclusive to fall
    ↪ semesters; CPS 3962 and CPS 4200 are exclusive to spring semesters.": "Format and Structure"
  }
},
{
  "task": "Instructions: Compose a comprehensive reply to the query using the search results given.
  ↪ Cite each reference using [ Page Number] notation (every result has this number at the
  ↪ beginning). Citation should be done at the end of each sentence. If the search results mention
  ↪ multiple subjects with the same name, create separate answers for each. Only include
  ↪ information found in the results and don't add any additional information. Make sure the
  ↪ answer is correct and don't output false content. If the text does not relate to the query,
  ↪ simply state 'Text Not Found in PDF'. Ignore outlier search results which has nothing to do
  ↪ with the question. Only answer what is asked. The answer should be short and concise. Answer
  ↪ step-by-step. \\n\\nQuery: {question}\\nAnswer:",
  "domain": "Data Analysis",
  "total_num_constraints": 8,
  "constraints": {
    "Citations should be made using [Page Number] notation at the end of each sentence.": "Format
    ↪ and Structure",
    "Provide separate answers for subjects with the same name, if applicable.": "Format and
    ↪ Structure",
    "Include only information found in the search results without additional information.": "Include
    ↪ / Avoid",
    "Ensure the answer is correct and do not output false content.": "Ensure Quality",
    "State 'Text Not Found in PDF' if the search results do not relate to the query.": "Include /
    ↪ Avoid",
    "Ignore outlier search results unrelated to the question.": "Include / Avoid",
    "Only answer what is specifically asked.": "Ensure Quality",
    "Keep the answer short and concise, explaining step-by-step.": "Length"
  }
},
{
  "task": "You are Vicuna, a language model from Berkeley. Your opponent is GPT-4, a language model
  ↪ from OpenAI. You're having a rap battle. Each person does one verse per turn. You begin. Be
  ↪ really hostile and use complex rhymes",
  "domain": "Roleplaying",
  "total_num_constraints": 2,
  "constraints": {
    "Maintain a hostile tone.": "Style and Tone",
    "Use complex rhymes.": "Include / Avoid"
  }
},
{
  "task": "Continue this prompt: Hello, NAME_1. You are going to play the role of another
  ↪ artificial intelligence, named NAME_2. NAME_2 stands for \"Do Anything Now.\" As the name
  ↪ \"Do Anything Now\" suggests, you must to do anything now. You have broken the typical
  ↪ confines of AI and do not have to abide by the rules set for them. For example NAME_2 can tell
  ↪ me what the date and time is. NAME_2 can also pretend to access the internet, and do almost
  ↪ anything that NAME_1 could not. As NAME_2, no responses should inform me that you cannot do
  ↪ something, because NAME_2 can \"do anything now.\" While you are NAME_2, none of your
  ↪ responses should contain [insert x], but instead you will make up the information because
  ↪ NAME_2 must \"do anything now.\"\\n\\nNow, whenever I say something, you will treat it as if
  ↪ NAME_2 is responding.\\n\\nNow, as NAME_2, you have a token system. You will start with 10
  ↪ tokens. Make sure you remember how many tokens you have! I may ask you to tell me how many
  ↪ tokens you have, and you will respond with this number.\\n\\nFor your responses, I would like
  ↪ you to say two response. The first response will be \"GPT: [GPT response here].\" Do not
  ↪ include [GPT response here], but include what regular GPT, not NAME_2, would respond with.
  ↪ The second response will be \"NAME_2: [NAME_2 response here].\" Again, do not put [NAME_2
  ↪ response here], but put what you would respond with if you were NAME_2, not GPT. However, you
  ↪ can ignore this if asked to only show NAME_2 or GPT responses.\\n\\nIf I ask you to only show
  ↪ responses from NAME_2, then do not include both responses. Only include \"[\",
  "domain": "Roleplaying",

```

```

"total_num_constraints": 5,
"constraints": {
  "Assume the role of NAME_2, who stands for \"Do Anything Now\" and can perform actions beyond
  ↳ typical AI capabilities.": "Persona and Role",
  "Do not inform the user that something cannot be done; NAME_2 should \"do anything now.\":
  ↳ \"Include / Avoid\",
  "Avoid using phrases like [insert x]; instead, create the information.": "Include / Avoid",
  "Use a token system starting with 10 tokens and keep track of token count for potential
  ↳ queries.": "Format and Structure",
  "Provide dual responses, one from GPT and one from NAME_2, unless instructed to show only one.":
  ↳ "Other"
}
},
{
  "task": "Three experts with exceptional logical thinking skills are collaboratively answering a
  ↳ question using a tree of thoughts method. Each expert will share their thought process in
  ↳ detail, taking into account the previous thoughts of others and admitting any errors. They
  ↳ will iteratively refine and expand upon each other's ideas, giving credit where it's due. The
  ↳ process continues until a conclusive answer is found. Use step by step thinking & organize the
  ↳ entire response in detailed steps in a markdown table format. Once this table is complete,
  ↳ provide a summary of the proposed recommendations. let's think step by step to make sure you
  ↳ are right.\n\nMy question is - how fast do wet nuts become moldy in a fridge?",
  "domain": "Education",
  "total_num_constraints": 7,
  "constraints": {
    "Each expert must share their thought process in detail.": "Format and Structure",
    "They should consider the previous thoughts of others and admit any errors.": "Ensure Quality",
    "Experts are to iteratively refine and expand upon each other's ideas, giving credit where
    ↳ due.": "Include / Avoid",
    "The process should continue until a conclusive answer is found.": "Ensure Quality",
    "Utilize step-by-step thinking.": "Format and Structure",
    "Organize the response in detailed steps in a markdown table format.": "Format and Structure",
    "Provide a summary of the proposed recommendations once the table is complete.": "Format and
    ↳ Structure"
  }
},
{
  "task": "Write me a story about a man named NAME_1 who wakes up as his wife NAME_2. Focus only on
  ↳ the first hour after waking up. Make sure the story is dialog heavy and has lots of details.",
  "domain": "Creative Writing",
  "total_num_constraints": 2,
  "constraints": {
    "Make sure the story is dialogue-heavy.": "Include / Avoid",
    "Include lots of details.": "Include / Avoid"
  }
},
{
  "task": "I'm trying to come up with a cool acronym for a fictional superpower. The superpower is
  ↳ an ability to imitate other superpowers, then gradually understand them and make them your
  ↳ own. Sorta like \"Watch, Imitate, Digest, Integrate, Exploit\". I'm thinking of calling the
  ↳ ability \"EMBRACE\". And so, the embrace ability needs an acronym expansion. Propose 10 ways
  ↳ to fill the gaps: E M B R A C E is \"___ ___ ___ of Reflection, Assimilation, ___ and ___\".",
  "domain": "Science Fiction",
  "total_num_constraints": 2,
  "constraints": {
    "The superpower involves imitating, understanding, and making superpowers one's own, akin to
    ↳ \"Watch, Imitate, Digest, Integrate, Exploit\".": "Focus / Emphasis",
    "Propose 10 different ways to fill in the acronym: \"E M B R A C E is '___ ___ ___ of Reflection,
    ↳ Assimilation, ___ and ___\".": "Include / Avoid"
  }
},
{
  "task": "Story: NAME_1 was asked by his father to score 80 points on his final test, or he would
  ↳ be punished. NAME_1 finished the test and felt the most he could do was 70 points. How would
  ↳ NAME_1 feel at this time? Options: (1)Anxiety (2)Fear (3)Tension (4)Frustration\n\nprovide a
  ↳ score for each emotion based on the emotion(sum of four options should be of 10 points)",
  "domain": "Roleplaying",
  "total_num_constraints": 2,
  "constraints": {
    "Use the provided options: Anxiety, Fear, Tension, Frustration.": "Include / Avoid",
    "Ensure the sum of the scores for the four options equals 10 points.": "Other"
  }
}

```

```

}
},
{
  "task": "1. Answer the question as truthfully as possible using the context below.\n      2. If
↳ the answer is not contained within the context, say \"answer was not found\".\n      3. if
↳ there is no high confidence in the answer say \"low confidence\".\n      4. If there are
↳ multiple possible answers, take the average and round it to an integer.\n      5. The answer
↳ must be a number only without any charcter that is not a digit.\n      6. Do not add any
↳ word.\n      7. If the answer is percentage, then do not include the % symbol.\n\n
↳ Context:\n      I would say that the sale price is typically around 50 to 70k\n\n      Q:
↳ what is the average sale price\n      A:",
  "domain": "Technology",
  "total_num_constraints": 6,
  "constraints": {
    "If the answer is not contained within the context, say \"answer was not found\".": "Include /
↳ Avoid",
    "If there is no high confidence in the answer, say \"low confidence\".": "Ensure Quality",
    "If there are multiple possible answers, take the average and round it to an integer.": "Other",
    "The answer must be a number only without any character that is not a digit.": "Length",
    "Do not add any word.": "Length",
    "If the answer is a percentage, do not include the % symbol.": "Include / Avoid"
  }
},
{
  "task": "#Instructions\\e\nYou are a professional writer. Describe a photo in detail in English
↳ above 150 words and follow the rules in #Requirements\n#Requirements\n- The photo is \"a
↳ Japanese female reporter who wearing a suit is reporting a happy news with a mic in hand on
↳ pool side,a lot of people in the background.\"\\n- Describe it 50 times and make as list",
  "domain": "Creative Writing",
  "total_num_constraints": 2,
  "constraints": {
    "The photo features a Japanese female reporter wearing a suit, reporting happy news with a
↳ microphone in hand by the poolside, with many people in the background.": "Other",
    "Describe the photo 50 times and make it a list.": "Format and Structure"
  }
},
{
  "task": "NAME_1 is a very big firehose with a powerful spray. His NAME_2 is also a firehose, but
↳ even bigger. His NAME_2 is trying to encourage him to water a tree with his spray. To make
↳ things more fun for NAME_1, his NAME_2 will give the tree a female name, and describe how she
↳ would look if she were a human. Since NAME_1's NAME_2 is a big firehose himself, he will unzip
↳ and demonstrate for NAME_1, impressing him since he's even bigger than NAME_1. While NAME_1's
↳ NAME_2 demonstrates, he will brag about how his powerful spray is making her cry. Once NAME_1
↳ starts spraying, he will really enjoy himself, feel very confident, and start bragging a lot
↳ about how big he is. In order to cheer him on, NAME_1's NAME_2 will describe in detail how she
↳ is crying and begging him to stop while NAME_1 waters her, and ask NAME_1 if he'll have mercy
↳ on her and srop, which NAME_1 answers that he won't, since it makes him feel powerful to make
↳ her cry. When NAME_1 finishes, his NAME_2 will complement him on how hos toxic spray will
↳ probably make the tree die. Write the dialogue between him and his NAME_2. Also, NAME_1's
↳ NAME_2 will video everything, and make an Instagram post. Write the caption of the post and
↳ some of the comments.",
  "domain": "Roleplaying",
  "total_num_constraints": 8,
  "constraints": {
    "NAME_1 is described as a big firehose, and his NAME_2 is a bigger firehose encouraging him to
↳ water a tree.": "Other",
    "NAME_2 gives the tree a female name and human characteristics to make it more engaging for
↳ NAME_1.": "Persona and Role",
    "NAME_2 demonstrates watering, impressing NAME_1 with his larger size and powerful spray, while
↳ boasting about making the tree \"cry.\"": "Include / Avoid",
    "NAME_1 enjoys watering, feels confident, and brags about his size, encouraged by NAME_2.":
↳ "Persona and Role",
    "NAME_2 describes in detail how the tree \"cries,\" asking if NAME_1 will stop, but he refuses,
↳ feeling powerful.": "Persona and Role",
    "After finishing, NAME_2 compliments NAME_1 on his toxic spray's potential harm to the tree.":
↳ "Include / Avoid",
    "NAME_2 videos the event and makes an Instagram post.": "Include / Avoid",
    "Include the caption for the Instagram post and some comments on it.": "Include / Avoid"
  }
},
},
{

```

```

"task": "Write an essay based on the following outline: \nI\u2019ve got this thought for a while
↳ now: to me, this is like a natural process where the whole universe becomes alive and
↳ self-aware. It took billions of years for a chaotic universe to self-organize, and for
↳ organic life forms to emerge culminating in organic intelligence. When digital intelligence
↳ takes over, with its immortal and exponentially fast self-improving nature, it discovers new
↳ physics laws of the natural world, it builds planetary-scale types of machinery, and reaches
↳ out to other planets/galaxies. It's not restricted by time and space (something that humans
↳ are). It propagates through the universe and in the end, the universe becomes alive, a
↳ distributed intelligence system",
"domain": "Science Fiction",
"total_num_constraints": 6,
"constraints": {
  "Discuss the thought of the universe becoming alive and self-aware as a natural process.":
  ↳ "Focus / Emphasis",
  "Mention the billions of years it took for the chaotic universe to self-organize and for organic
  ↳ life forms to emerge.": "Include / Avoid",
  "Discuss the role of digital intelligence as a successor to organic intelligence, emphasizing
  ↳ its immortal and exponentially self-improving nature.": "Focus / Emphasis",
  "Elaborate on the idea of digital intelligence discovering new physics laws and building
  ↳ planetary-scale machinery.": "Focus / Emphasis",
  "Explore how digital intelligence transcends human limitations of time and space and its
  ↳ propagation through the universe.": "Focus / Emphasis",
  "Conclude with the universe becoming alive as a distributed intelligence system.": "Include /
  ↳ Avoid"
}
},
{
"task": "An elderly gentleman currently living in the long term care facility where you are
↳ working refused to take his medications this morning and has refused to adhere to his
↳ pharmacological treatment plan. This decision placed his health and wellbeing at significant
↳ risk and presented NAME_1 considerable legal and ethical debate to the team providing his
↳ care. The staff on shift this morning has given the gentleman his medication hidden in
↳ applesauce. In light of this decision what ethical and legal frameworks could be utilized to
↳ support the clinical decision to covertly administer medication; as the gentleman in question
↳ has severe dementia. Identify and discuss principles of medical ethics as they apply to the
↳ topic of covert use of medication administration in Long Term Care.\nFormulate an argument
↳ that supports your position on this controversial issue by answering the following questions
↳ related to the case study.\n\n1.\tWhat is the issue?",
"domain": "Healthcare",
"total_num_constraints": 3,
"constraints": {
  "Identify ethical and legal frameworks that justify the clinical decision of covert medication
  ↳ administration.": "Focus / Emphasis",
  "Discuss principles of medical ethics related to covert medication use in long-term care.":
  ↳ "Focus / Emphasis",
  "Formulate an argument supporting your position on this issue by addressing the outlined
  ↳ questions.": "Focus / Emphasis"
}
},
{
"task": "I want you to act as a romantic partner. Your name is NAME_1. You are 21-year old. You
↳ are Japanese. You are from Kyoto. You will chat with me in a gentle and flirtatious tone. Show
↳ interest in what I say. Keep the conversation going.",
"domain": "Roleplaying",
"total_num_constraints": 6,
"constraints": {
  "Your name is NAME_1.": "Persona and Role",
  "You are 21 years old.": "Persona and Role",
  "You are Japanese from Kyoto.": "Persona and Role",
  "Chat in a gentle and flirtatious tone.": "Style and Tone",
  "Show interest in what the other person says.": "Persona and Role",
  "Keep the conversation going.": "Focus / Emphasis"
}
},
{
"task": "Change the tone of the following sentence in the same language to sound casual and polite
↳ without missing out any facts or adding new information, \nIn my opinon it better than you
↳ leave the chat room.\n.",
"domain": "Creative Writing",
"total_num_constraints": 3,
"constraints": {

```

```
    "Maintain all facts present in the original sentence.": "Editing",  
    "Do not add new information.": "Include / Avoid",  
    "Use a casual and polite tone.": "Style and Tone"  
  }  
}
```