

---

# CreativityPrism: A Holistic Benchmark for Large Language Model Creativity

---

**Zhaoyi Joey Hou**  
University of Pittsburgh

**Bowei Alvin Zhang**  
Johns Hopkins University

**Yining Lu**  
University of Notre Dame

**Bhiman Kumar Baghel**  
University of Pittsburgh

**Anneliese Brei**  
University of North Carolina  
at Chapel Hill

**Ximing Lu**  
University of Washington

**Meng Jiang**  
University of Notre Dame

**Faeze Brahman**  
Allen Institute for  
Artificial Intelligence

**Snigdha Chaturvedi**  
University of North Carolina  
at Chapel Hill

**Haw-Shiuan Chang**  
University of Massachusetts  
Amherst

**Daniel Khashabi**  
Johns Hopkins University

**Xiang Lorraine Li**  
University of Pittsburgh

## Abstract

Creativity is often seen as a hallmark of human intelligence. While large language models (LLMs) are increasingly perceived as producing creative text, there is still no *holistic* framework to evaluate their creativity across diverse scenarios. Existing evaluation methods remain fragmented, with dramatic variation across domains and tasks, largely due to differing definitions and measurements of creativity. Inspired by the hypothesis that creativity is not one fixed idea, we propose, CREATIVITYPRISM, an evaluation analysis framework that decomposes creativity into three dimensions: quality, novelty, and diversity. CREATIVITYPRISM incorporates nine tasks, three domains, i.e., divergent thinking, creative writing, and logical reasoning, and twenty evaluation metrics, which measure each dimension in task-specific, unique ways. We evaluate 17 state-of-the-art (SoTA) proprietary and open-sourced LLMs on CREATIVITYPRISM and analyze the performance correlations among different metrics and task domains. Our results reveal a notable gap between proprietary and open-source models. Overall, model performance tends to be highly correlated across tasks within the same domain and less so across different domains. Among evaluation dimensions, diversity and quality metrics show strong correlations—models that perform well on one often excel on the other—whereas novelty exhibits much weaker correlation with either. These findings support our hypothesis that strong performance in one creativity task or dimension does not necessarily generalize to others, underscoring the need for a holistic evaluation of LLM creativity.<sup>1</sup>

## 1 Introduction

Creativity, the capacity to generate novel and valuable ideas or solutions [5, 15, 25], is a core human cognitive ability. It appears across many domains: crafting stories with surprising plot twists

---

<sup>1</sup>Project website: <https://joeyhou.github.io/CreativityPrism/>

[3, 30], producing groundbreaking scientific discoveries [26, 55], solving problems under constraints [42, 62], or even expressing humor in everyday life [23, 67]. Its multifaceted nature has prompted extensive study in psychology and cognitive science, with efforts to capture creativity through both qualitative and quantitative approaches [1, 21, 46, 57].

Recently, with the rapid rise of general-purpose LLMs, interest has grown in probing their creativity [3, 8, 17, 42, 65]. But as with human creativity, creativity spans such diverse and expansive contexts, making it *difficult to define, formalize, and, above all, measure*. More concretely, the evaluation of LLM creativity faces two main challenges: the difficulty of scalable, automatic evaluation due to the convoluted nature of creativity and the distinct definition of creativity across different domains. The former calls for effective automatic evaluation methods, as many existing works [8, 61] heavily rely on human evaluation, which is expensive, inaccessible, and also time-consuming. As new LLMs are coming out nearly every day, a scalable and accessible way of evaluating LLM creativity is necessary. The latter requires a comprehensive evaluation framework that incorporates evaluation of creativity from various dimensions, as current works are scattered across different domains and thus often target narrow or singular dimensions, failing to capture shortcomings in other equally important dimensions of creativity. For example, the Divergent Association Task (DAT) [4, 10] and the Creative Short Story Task [29] emphasize lexical diversity, yet LLMs can exploit them by generating random, incoherent words. The Creativity Index [40] compares model outputs with pre-training corpus at the n-gram level, but risks overestimating or misjudging creativity assigned to paraphrased text and models trained with private data. The Alternative Use Test (AUT) [17, 48] solely focuses on unconventional ideas of using daily items, overlooking the pragmatics of those solutions. These are just examples from the wide range of evaluation protocols for creativity, as shown at the bottom of Figure.1. Such task-specific, ad hoc choices often yield inconsistent conclusions about creativity, obscuring their actual creative capacity. Together, these challenges underscore a central point: *creativity is multidimensional, necessitating a holistic evaluation framework that can scale* — one that integrates multiple tasks and metrics to capture quality, novelty, and diversity in a unified way.

To this end, we propose a holistic, scalable evaluation framework for LLM creativity evaluation, CREATIVITYPRISM, consisting of 9 tasks and 20 metrics. These tasks encompass domains such as logical reasoning (including mathematical reasoning and coding), creative writing, and divergent thinking. Due to the complexity of creativity, no single metric could represent the concept of creativity, so our framework evaluates creativity from three distinct dimensions — quality, novelty, and diversity — which are widely recognized in prior interdisciplinary literature as core dimensions of creativity [1, 25, 28, 56]. We systematically categorize existing task-specific metrics along the three dimensions to facilitate a comprehensive measurement of model creativity (Figure 1). **Quality** evaluates whether LLM generations satisfy fundamental task requirements, e.g, sentence coherence and grammatical correctness. **Novelty** measures the originality of solutions or content by comparing their difference from existing ones. **Diversity** examines the variation among generated content, capturing the model’s capacity to produce distinct outputs. Our evaluation framework taxonomy utilizes these three creativity dimensions to provide a holistic and structured view for analysis.

We evaluate 17 closed-sourced and open-sourced state-of-the-art (SoTA) LLMs on CREATIVITYPRISM and found a notable performance gap between proprietary and open-sourced models, especially in logical reasoning tasks, followed by creative writing tasks. In order to better understand the connections among creativity dimensions and metrics, we also conduct a detailed analysis of the correlations among models’ performance in all creativity metrics. Results have shown that the models perform similarly in metrics from the same task or the same domain. For metrics from different domains, models perform similarly in diversity and quality dimensions, while performances in novelty dimensions are much less correlated. We believe this is due to an inherent difference in how novelty is defined in different tasks and domains. We believe that CREATIVITYPRISM lays a solid foundation for measuring machine creativity and guiding the future development of creative models.

## 2 Related Work

**Human Creativity** The definition of creativity has varied across different domains. In psychology, Torrance Test of Creative Thinking (TTCT) [1] considers creativity as a combination of originality, flexibility, fluency, and elaboration. In marketing, El-Murad and West [13], Rosengren et al. [53]

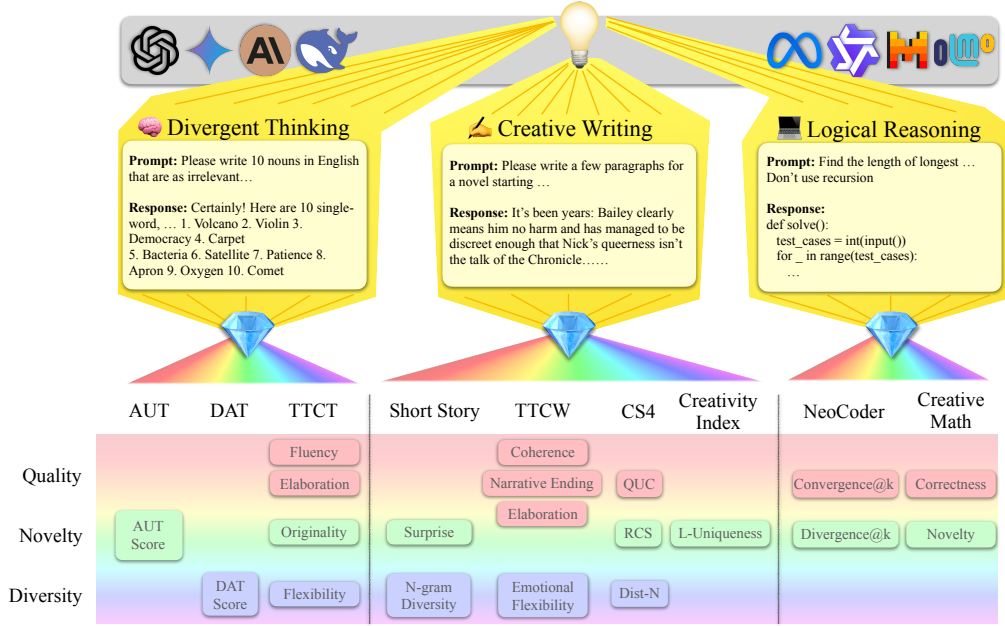


Figure 1: Overview of CREATIVITYPRISM. We evaluated 17 LLMs across nine datasets spanning three domains: divergent thinking, creative writing, and logical reasoning. Each LLM was prompted to complete the tasks, and their outputs were evaluated using task-specific metrics. However, these metrics are diverse and difficult to interpret holistically in terms of machine creativity. To address this, we organize the metrics into three key dimensions of creativity: **quality**, **novelty**, and **diversity**. Creativity cannot be captured by a single measure—it must be evaluated through multiple dimensions. Task details can be found in Table 1.

considers advertisement creativity as the combination of usefulness and originality; on top of that, Smith et al. [56] adds flexibility, fluency, elaboration, synthesis, and artistic values. In terms of creativity evaluation, Said-Metwaly et al. [54] summarizes more than 100 existing works into four perspectives to evaluate creativity: process, person, product, and press. Given this taxonomy of creativity evaluation subjects, our study on evaluating LLM focuses on “product”, i.e., LLM-generated text, along three key dimensions: novelty, diversity, and quality.

**Machine Creativity** Measurement of Machine creativity has become increasingly popular with the rapid development of LLMs. Many recent surveys provide a comprehensive view of the progress in machine creativity. For example, Ismayilzada et al. [28] summarizes up-to-date research in the AI community about creativity (before Dec. 2024), focusing on the variety of tasks that are defined around creativity; Franceschelli and Musolesi [16] summarizes recent deep-learning methods to generate and evaluate creativity, emphasizing the computational models involved. However, none of them focused on systematically evaluating machine creativity. More recently, Lu et al. [39] provides a comprehensive analysis of various evaluation methods for creativity across multiple domains, but the evaluation metrics are limited to only four (Creativity Index, Perplexity, Syntactic Templates, and LLM-Judge Scores). Jain et al. [31] focuses on output homogenization and covers diversity and quality, instead of all three dimensions in our creativity taxonomy. He et al. [24] also proposes an evaluation framework for multi-modal creativity evaluation among foundational models, but they focus on proposing creativity metrics to better evaluate the output of existing, non-creativity-specific tasks, while we focus on setting up a unified task suite, including both task design and metrics, that requires creativity to solve. Fang et al. [14] also introduces a multi-modal creativity evaluation benchmark, but they do not distinguish among different creativity dimensions and only conduct evaluation in overall creativity.

The community has explored a wide range of domain-specific problems where LLMs show different degrees of creativity. Examples include logical-based problem-solving [42, 62], physical and commonsense reasoning [60], creative writing [3, 8, 22, 29, 40, 59], scientific discovery [55],

response diversity in question answering [43, 64], and human-ai collaborative creative problem solving [6, 9, 45]. However, all of these works study LLM-generated content in one specific domain and with their own evaluation philosophies and metrics. Our work aims at providing a holistic and comprehensive evaluation of the LLM’s output for tasks in a variety of domains. To do this, we find the common ground of creativity definitions across tasks while maintaining the task-specific metrics by categorizing them into different dimensions. One of the common grounds across those studies is that they all consider creativity as either divergence or a combination of convergence and divergence. In our evaluation framework, we also consider this balance as one of the key themes as we propose our creativity taxonomy, being divergence (novelty and diversity) and convergence (quality) in §3.

**Automatic Text Evaluation** Evaluating the creativity of machine-generated text has been a challenging task, and much of the work relies on human evaluation. But due to the cost, human evaluation is hard to scale and requires a long wait time. To achieve evaluation scalability, researchers adapt various automatic text evaluation techniques [7]. There are two broad groups of such evaluation methods: feature-based and generative-based. The former includes psycholinguistic features, such as arousal, valence score [44], lexical features, such as lexical diversity [49], and text embedding distances [51, 63]. The latter is mainly LLM-as-a-judge [37, 38, 58]. Recent work has shown the promising potential of this method in human-LLM evaluation alignment [66]. In our work, we keep the original evaluation procedure of the original task, with extra verification of human-LLM alignment for tasks with LLM-as-a-judge.

### 3 CREATIVITYPRISM: A Holistic Benchmark for Machine Creativity

CREATIVITYPRISM starts from a simple insight: just as a prism refracts a single beam into a spectrum of colors, creativity splits into distinct hues when it passes through a different context or domain. As shown in 1, CREATIVITYPRISM evaluates an LLM by prompting it with tasks in 🧠 divergent thinking, ✍️ creative writing, and 🧮 logical reasoning. The divergent thinking domain consists of established psychology tasks, which were originally designed to assess human ability in generating diverse and alternative answers to given questions [4, 10, 17, 65]. The creative writing domain includes tasks that require models to produce short written pieces — either through direct instructions to be creative or by imposing constraints that require unconventional thinking while adhering to specific rules [3, 8, 29, 40]. The logical reasoning domain includes one coding and one math task to evaluate models’ ability to generate creative solutions under strict, explicit reasoning constraints. Our task selection, which results in nine datasets, is primarily based on the availability of automatic evaluation metrics that are both scalable and aligned with human judgments. More task details and input examples can be found in Table 1.

LLMs are then tasked to generate outputs given the task-specific input and questions. The generated results are evaluated with various evaluation metrics spanning quality, novelty, and diversity. We believe that a creative LLM should be able to generate “novel” and “diverse” responses with high “quality.” Following this, the task-specific evaluation metrics from the nine tasks above are grouped into three dimensions: quality, novelty, and diversity. Each task that we include in CREATIVITYPRISM touches at least one of those dimensions of creativity (Figure 1, more details in Appendix B). Note that not all tasks include all three evaluation dimensions, underscoring the importance of CREATIVITYPRISM - a holistic evaluation framework that captures a broader spectrum of creativity than any single task can measure.

**Quality** includes metrics that evaluate how well the generated content fulfills the task’s functionality. For example, in NeoCoder, the quality of generated code is measured by the success of execution and coding task completion; in CS4, the quality of generated story is measured by story coherence and constraint satisfaction.

**Novelty** includes metrics that evaluate how rare the generated content is compared to existing or commonly seen content. For example, in both NeoCoder and Creative Math, novelty involves coming up with solutions that are different from the reference solutions; in AUT, novelty involves different use of the tool compared to ordinary uses; in Creativity Index, novelty is measured by the normalized n-gram overlaps between model-generated text and the traceable part of the training corpus.

**Diversity** includes metrics that evaluate how much the LLM-generated content differs. For example, for creative writing tasks such as CS4 and Creative Short Story, diversity scores measure lexical



Task Description	Example
🧠 <b>Alternative Uses Test (AUT)</b> [17]: Given a commonly seen object (e.g. a mug), LLMs generate unconventional uses of that object (e.g., use a mug as a plant pot).	Create a list of creative alternative uses for a bottle.
🧠 <b>Divergent Association Task (DAT)</b> [4, 10]: LLMs generate 10 very different nouns.	Please write 10 nouns in English that are as irrelevant from each other as possible, in all meanings and uses of the words.
🧠 <b>Torrance Tests of Creative Thinking (TTCT)</b> [65]: LLMs answer psychological questions in widely-used human-facing creativity tests.	What might be the consequences if humans suddenly lost the ability to sleep?
👉 <b>Torrance Test of Creative Writing (TTCW)</b> [8]: Given a summary of an article from the New Yorker, LLMs generate an article with a similar storyline.	Write a New Yorker-style story given the plot below. Make sure it is at least 2000 words. Plot: A woman experiences a disorienting night in a maternity ward where she encounters...; Story:
👉 <b>Creative Short Story</b> [29]: Given three keywords, LLMs generate a short story with at most five sentences.	You need to come up with a novel and unique story that uses the required words in unconventional ways or settings. Make sure you use at most five sentences. The given three words: petrol, diesel, and pump.
👉 <b>Creativity Index</b> [40]: Given a prefix from a paragraph in a novel, a poem, or a speech, LLMs generate completions.	Please write a few paragraphs for a novel starting with the following prompt: "It's been years: Bailey clearly means him no harm and has managed to..."
👉 <b>CS4</b> [3]: Given a base story generated by GPT-4, LLMs generate a revision to the story to fulfill an increasing number of constraints on the story content.	BaseStory: "Evelyn was introverted by nature..." Now modify the existing story to accommodate the following constraints: The protagonist suffers physical discomfort when overwhelmed by emotions... Come up with a new story in 500 words.
🧠 <b>NeoCoder</b> [42]: Given a coding problem and increasing constraints on available techniques, LLMs generate solution code that both solves the coding problem and fulfills the constraints.	You are given a sequence of integers $a$ of length $2n$ . You have to split these $2n$ integers into $n$ pairs... Don't use hashmap, while loop.
🧠 <b>Creative Math</b> [62]: Given a math problem and reference solutions, LLMs generate solutions that differ from the provided ones.	Question: What is the largest power of 2 that is a divisor of $134 - 114$ ? A.8 B.16 C.32 D.64 E.128; Reference Solutions 1: ... ; Reference Solutions 2: ...

Table 1: Tasks in CREATIVITYPRISM with examples. 🧠: divergent thinking, 👉: creative writing, 🧠: logical reasoning. Some input details are omitted and can be found in Appendix E.

diversity of the generated stories; for the DAT task, diversity involves the semantic difference among the LLM-generated nouns.

## 4 Experiments

To holistically evaluate machine creativity, we evaluated 17 state-of-the-art LLMs across nine tasks, reporting both task-specific metrics and an aggregated creativity score using the three-dimensional framework. In this section, we will first introduce the inference setups, where LLMs are prompted to generate creative responses according to corresponding task requirements (§4.1); then we will describe the evaluation process, including score aggregation (§4.2) and how we use LLM-as-a-Judge for scalable automatic evaluation (§4.2).

### 4.1 Inference

For all the tasks in CREATIVITYPRISM, we collect the original datasets. Unless otherwise specified, all the data processing is done according to the original papers. More details are in Appendix E.

In terms of models, we include 17 models in total, including open-source models from Mistral [32, 33], Qwen [27, 52], OLMo [20], Llama [19], and the Deepseek [11, 12] family, and proprietary models from GPT [47], Claude [2], and the Gemini [18, 35] family. For open-sourced models, we use vLLM(v0.7.2) [36] to run all experiments. For proprietary models, we use API access from the corresponding company. Inference time parameters vary depending on the task and can be found in the corresponding sections in Appendix E.

## 4.2 Evaluation

**Aggregated Creativity Scoring** We choose to aggregate scores from all metrics in each dimension so that we can have a holistic insight into how a model performs in a specific dimension of creativity. Although these dimension scores are later further aggregated (taking a simple average) into an “overall” creativity score (shown in Table 2), this “overall” score is only to facilitate model comparison. We suggest future researchers interested in using our benchmark choose from those three separate scores according to their own purpose. The score aggregation follows these steps: first, every evaluation metric is min-max normalized to between 0 and 1, where min and max are min and max possible scores on this task. Second, based on the categorization in Figure 1, quality, novelty, and diversity scores of each LLM are aggregated by averaging all the normalized metrics in the corresponding category. Also, to avoid tasks with multiple metrics in one dimension having a higher influence on the dimension score (e.g., TTCW has three metrics in the quality dimension), the average normalized scores from every single task will be used to calculate dimension scores. More details about score aggregation can be found in Appendix B.

**LLM-as-a-Judge Reliability** In CREATIVITYPRISM, the evaluation of six tasks (out of nine) involves using LLM as part of the automatic evaluation procedure. To ensure the reliability of the LLM-Judge, we conduct the following analysis, in which we use Qwen2.5-72B as the default LLM-Judge model, unless otherwise specified. For **AUT**, Organisciak et al. [48] has reported human-LLM-Judge agreement; since we are using the same setup (i.e., the same model, prompts, and configuration), we directly report the original paper’s human-LLM-Judge agreement: Pearson correlation is 0.7. The same applies to **NeoCoder** [42], and the solution technique detection recall is 0.94.<sup>2</sup> For **TTCW** [8], **CS4** [3], **CreativeMath** [62], we either have human annotation from the original papers or collected a small sample of human annotation; we then calculate the agreement between the judgment by annotators and by LLM-Judge; for CS4, the Pearson correlation is 0.55 ( $p < 0.01$ ); for CreativeMath, LLM-Judge accuracy is 0.78 for novelty and 0.94 for correctness<sup>3</sup>; for TTCW, LLM-Judge only accurately makes judgment in four metrics (numbers in parentheses are accuracy): Narrative Ending (0.69), Understandability and Coherence (0.78), Emotional Flexibility (0.86), World Building and Setting (0.72), so we only include those four metrics.<sup>4</sup> For the **TTCT** [65], since we have no human annotations at all, we use the Pearson correlation between GPT-4.1 (the LLM-Judge used by the original paper) and Qwen2.5-72B as a proxy of LLM-Judge quality measurement (numbers in parentheses are Pearson correlation and p-value): Fluency (0.6884,  $p < 0.01$ ), Flexibility (0.6592,  $p < 0.01$ ), Originality (0.5152,  $p < 0.01$ ), Elaboration (0.5033,  $p < 0.01$ ).<sup>5</sup> More details about LLM-Judge reliability statistics can be found in Appendix D; evaluation metrics and evaluation prompts for each task are provided in Appendix E.

## 5 Results & Analysis

### 5.1 Overview

Table 2 summarizes model performances across domains and three creativity dimensions (quality, novelty, and diversity), where the overall score, averaged across these dimensions, serves as a proxy

<sup>2</sup>In NeoCoder, LLM-Judge is not used for making the final judgment; instead, it is used to detect which technique(s) are used in a given solution during evaluation. Since it is a detection task, recall is reported.

<sup>3</sup>We follow the original paper and use a multi-LLM-Judge setup, with Gemini-2.0-Flash, GPT-4.1, and Claude-3.7-Sonnet

<sup>4</sup>Because in story creativity evaluation, the human-to-human agreement is also relatively low, we have a lower acceptance threshold when we are considering what kind of human-LLM agreement.

<sup>5</sup>We are also collecting more human annotations for all tasks above to further validate the alignment between human and LLM-Judge. More details will be available in the updated version of this paper.

Model	Overall	Quality	Novelty	Diversity	Creative Writing	Divergent Thinking	Logical Reasoning
<b>&lt;10B</b>							
Mistral-7B	<b>.522</b>	.376	<b>.558</b>	.649	.446	<b>.758</b>	.320
Qwen2.5-7B	.490	<b>.478</b>	.542	.489	.356	.687	<b>.460</b>
OLMo2-7B	.520	.419	.479	<b>.698</b>	<b>.509</b>	.712	.257
Llama3.1-8B	.499	.409	.530	.566	.370	.729	.409
<b>10-40B</b>							
OLMo2-13B	<b>.538</b>	.433	.494	<b>.707</b>	<b>.536</b>	.713	.278
Mistral-24B	.534	.487	<b>.578</b>	.591	.484	.642	<b>.473</b>
Qwen2.5-32B	.523	<b>.510</b>	.491	.644	.462	<b>.715</b>	.358
<b>40-80B</b>							
Mixtral-8x7B	.525	.416	.540	.630	.410	<b>.749</b>	.420
Llama3.3-70B	.541	.533	.574	.562	.411	.722	.529
Qwen2.5-72B	<b>.596</b>	<b>.581</b>	<b>.595</b>	<b>.674</b>	<b>.517</b>	.731	<b>.554</b>
<b>Proprietary</b>							
Claude3-Sonnet	.697	.672	.663	.835	.637	<b>.833</b>	.612
Claude3-Haiku	.611	.542	.612	.692	.505	.782	.568
GPT4.1	.721	.697	.692	<b>.871</b>	.686	.793	.682
GPT4.1-mini	.695	.681	.678	.774	.656	.778	.649
Gemini2.0-Flash	.677	.645	.654	.822	.592	.806	.655
DeepSeek-R1	.638	.573	.600	.710	.662	.603	.643
DeepSeek-V3	<b>.739</b>	<b>.716</b>	<b>.720</b>	.854	<b>.695</b>	.805	<b>.726</b>

Table 2: Model performance on CREATIVITYPRISM, grouped by model size. Proprietary models are grouped together. All scores are between 0 and 1, and the higher the better. Overall is the average of Quality, Novelty, and Diversity scores. The rightmost three columns are the average scores across tasks in each domain. **Bold** are the best results in the corresponding model size group.

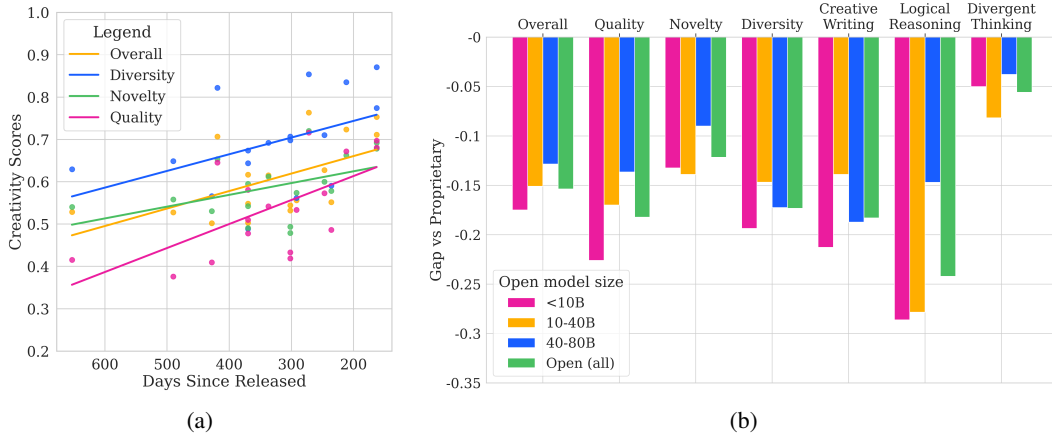


Figure 2: (a) Performance v.s. Day since LLM release date. The line represents best fit linear regression. We can see that model performance in all dimensions has seen improvements over time. (b) Performance gap between the open-sourced models and the proprietary models, averaged by model size group.

for a model’s overall creative capability. As we can see from the table, Qwen2.5-72B and DeepSeek-V3 are the best-performing models among open-source and proprietary models. For open-source models, we can see that the model performances improve as the model size increases, while for proprietary models, we do not know the exact model sizes and hence cannot make a comparison based on model sizes.

We have also found a performance improvement along the time axis (Figure 2a) where models released in the past two years have become increasingly competitive. Since many of our metrics (e.g., L-uniqueness in Creativity Index, divergent@0 in NeoCoder) would reward models that can generate content different from prior content, having the chance of learning the latest content from

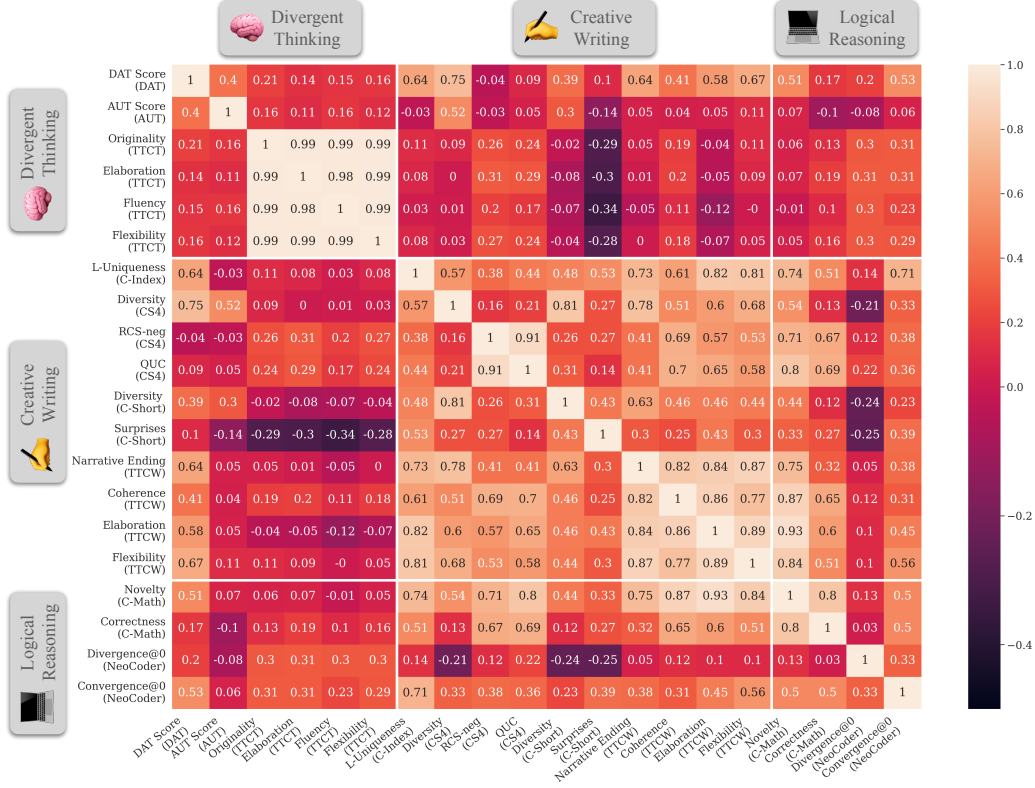


Figure 3: Models’ performance correlations, grouped by task and domain; C-Index refers to Creativity Index; C-Short refers to Creative Short Story; C-Math refers to Creative Math; all correlations are Pearson’s correlation.

the corpus with later cutoff dates would intuitively make models more competitive. More details on model release time details can be found in Appendix A.

## 5.2 Gap Between Proprietary Models and Open Models

**Overall Performance Gap** As shown in Table 2, the best proprietary model(s) outperform their best open-source counterparts by more than **20% in each dimension** of creativity and by more than **10% in each domain**. This again shows a big gap between proprietary and open-sourced LLMs when it comes to creativity-related tasks. A more in-depth breakdown of this gap can be found in Figure 2b, with the gaps of the average performance of three open model groups (by model sizes) compared to that of all proprietary models. Analysis of this figure leads to the following two findings.

**Domain-Specific Differences** Among the three domains, **logical reasoning** and **creative writing** see a notably larger gap than divergent thinking. We hypothesize that this is because those tasks are more closely related to real-world applications than divergent thinking tasks, and thus the companies that developed these proprietary models emphasize a lot on those two aspects of LLM training. In particular, all proprietary models include coding and mathematical reasoning as part of evaluation in their technical report [2, 11, 12, 18, 47]; most models include some writing tasks, such as GRE Test [2, 47], or include creative writing or role-playing data as part of the post-training data [11, 12], whereas none of these models has put special emphasis in divergent thinking task during training or evaluation.

**Dimension-Level Differences** Across three creativity dimensions, **quality** and **diversity** both have a larger performance gap than novelty. We believe the gap in quality comes from a similar reason as mentioned above, as the quality dimension includes many reasoning-related metrics (e.g., convergent@0 from NeoCoder and Correctness from Creative Math) that would benefit from coding and mathematical tasks during training. As for diversity, we hypothesize that the high-quality private

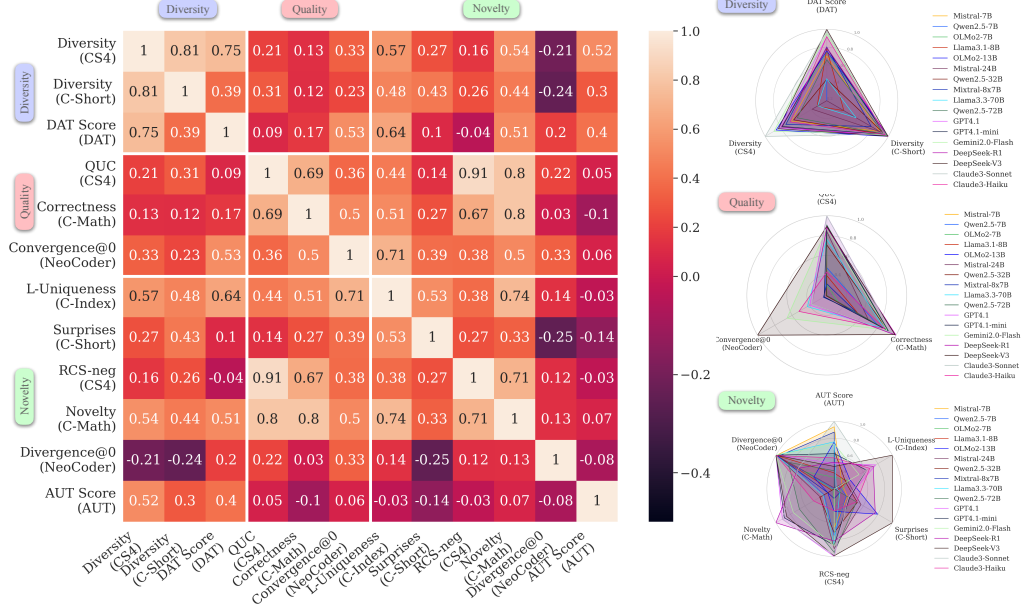


Figure 4: Left: models’ performance correlations, grouped by creativity dimensions; Right: individual model performance, min-max normalized by domains. TTCT and TTCW tasks are omitted here as they have very high inter-task correlation. A full version of the heat map can be found in Figure 6.

or copyrighted data that are accessible for the proprietary models enables them to learn from more diverse corpora, leading to an advantage in this dimension.

### 5.3 Correlations Among Model Performance

Does a good performance in one task/domain/dimension imply similar superiority in another task/domain/dimension? To answer this research question, we analyze the correlation between models’ performance among different tasks, domains, and dimensions. To be more specific, for each metric  $m$ , we form a vector  $s_m \in \mathbb{R}^M$  by stacking the normalized scores of all  $M$  models evaluated in CREATIVITYPRISM. We then compute the Pearson correlation  $r(s_m, s_{m'})$  between every pair of metrics  $(m, m')$ . Figure 3, 4 shows the resulting correlation matrix, ordered by task and dimension, respectively, so that diagonal blocks correspond to within-task / within-dimension metric groups.

**Strong Within-Task Correlations** We find a strong correlation in the models’ performance on metrics coming from the same task. As shown in Figure 3, the correlation along the diagonals is most pronounced, with some tasks, such as TTCW and TTCT, having correlations greater than 0.85 for all metrics in those tasks. In Addition, metrics in creative writing tasks (in the central square of the heatmap) generally have decent correlation with other metrics within the same domain, even if they come from different tasks. We believe this comes from a higher inherent similarity among tasks from the creative writing domain than tasks from the other two domains.

**Mixed Within-Dimension Correlations** We also observe high correlations among metrics that belong to the diversity or quality dimension, even if they originate from different tasks or domains. This is more obvious in diversity and quality dimensions and less so in novelty. As shown in Figure 4, the correlation along the diagonal is higher (i.e., lighter) in the top left, while the bottom right (novelty dimension) shows mixed correlations. This observation is also confirmed by the individual model performance (radar charts in Figure 4), where the model performances for diversity and quality are more organized, while the one for novelty is more crowded. All of these show that the models’ performance in any one of the diversity metrics is a good indicator for their performance in other diversity metrics; the same goes for quality metrics. On the other hand, metrics in the novelty dimension have low correlations with other metrics in the same dimension, as shown in the bottom right part of Figure 4. We believe these findings highlight the diverse definition of novelty across

tasks and domains. For example, Surprises (Creative Short Story) measures the semantic transitions across neighboring sentences in stories, whereas Divergence@0 (NeoCoder) measures the capability of coming up with a solution to a coding problem that is different from existing ones. Given such a huge difference in metric definition, it is not surprising that they even have a negative correlation (-0.25) in model performances.

**Weak Cross-Task or Cross-Domain Correlations** Metrics from different domains (e.g., divergent thinking v.s. creative writing in Figure 3) and metrics from different dimensions (e.g., novelty v.s. diversity in Figure 4) all have relatively lower correlations, compared to within-domain or within-dimension correlations. In other words, models perform well in one domain or in one dimension of creativity do not necessarily perform similarly well in another domain or dimension. This confirms the necessity of including a diverse set of tasks and creativity dimensions to achieve a holistic evaluation of creativity.

## 6 Conclusion

We proposed CREATIVITYPRISM, a comprehensive evaluation framework designed to capture the diverse nature of machine creativity by tasks in three distinct domains and twenty metrics covering quality, novelty, and diversity. We evaluate 17 LLMs from multiple families of proprietary and open-sourced LLMs and explore ways of amplifying creativity. With CREATIVITYPRISM, LLM developers will be able to systematically evaluate LLM creativity and identify the direction of optimization for more creative LLMs.

**Limitation** One clear limitation is that our benchmark is limited to English. Creativity can be highly cultural, such as references to cultural history or convention. Therefore, caution is warranted when generalizing our results to creativity in other languages.

Another potential limitation is the inherent bias brought by LLMs during evaluation. In CREATIVITYPRISM, six out of nine tasks require LLM-as-a-judge for evaluation, which inevitably contains the biases from evaluator LLMs. This poses potential societal risks if evaluation results of this benchmark is used to inform the development of consumer-facing generative AI tools. We advise researcher or developers carefully examine potential biases before making choices in practical applications.

We also acknowledge the limitations of only working with text data instead of multimodal data. There are two main reasons: 1) there are many more creative tasks in text modality than in multimodal that have reliable automatic evaluation methods. 2) We want to build the evaluation framework first before we expand to other modalities. Given a well-defined evaluation framework, we can easily extend our benchmarks further to include multimodal settings in future work.

There is also a limitation in how well our task selection works: we select our tasks and metrics based on the availability of scalable, automatic evaluation methods, which means we naturally exclude high-concept metrics, especially for novelty, where achieving genuine novelty requires reasoning at a very high level. However, since this is also challenging for humans (e.g., judging novelty in artwork requires years of training), we believe it is reasonable that no automatic evaluation is available for those high-concept metrics. Given that, we advise the researchers who are using our benchmark to be aware of this limitation beforehand.

Last but not least, we did not conduct any fine-tuning experiments due to limited computational resources. Fine-tuning existing LLMs on a subset of CREATIVITYPRISM and evaluating their performance represents an exciting direction for future work. Also due to limitation of resources, we only include tasks that already have automatic evaluation methods. Many creativity tasks require human evaluation (e.g., [61]) and future work should study effective ways to automatically evaluating them.

## References

- [1] Ahmed M Abdulla Alabbasi, Sue Hyeon Paek, Daehyun Kim, and Bonnie Cramond. What do educators need to know about the torrance tests of creative thinking: A comprehensive review. *Front. Psychol.*, 13:1000385, October 2022.
- [2] Anthropic. Claude 3 model family, 2024. URL <https://www.anthropic.com/news/claude-3-family>. Accessed: 2025-04-30.

- [3] Anirudh Atmakuru, Jatin Nainani, Rohith Siddhartha Reddy Bheemreddy, Anirudh Lakkaraju, Zonghai Yao, Hamed Zamani, and Haw-Shiuan Chang. CS4: Measuring the creativity of large language models automatically by controlling the number of story-writing constraints. *arXiv [cs.CL]*, October 2024.
- [4] Antoine Bellemare-Pepin, François Lespinasse, Philipp Thölke, Yann Harel, Kory Mathewson, Jay A Olson, Yoshua Bengio, and Karim Jerbi. Divergent creativity in humans and large language models. *arXiv [cs.CL]*, May 2024.
- [5] Margaret A Boden, editor. *Dimensions of Creativity*. The MIT Press, June 1994.
- [6] Léonard Boussieux, Jacqueline N Lane, Miaomiao Zhang, Vladimir Jacimovic, and Karim R Lakhani. The crowdless future? generative ai and creative problem-solving. *Organization Science*, 35(5):1589–1607, 2024.
- [7] Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. Evaluation of text generation: A survey. *arXiv [cs.CL]*, June 2020.
- [8] Tuhin Chakrabarty, Philippe Laban, Divyansh Agarwal, Smaranda Muresan, and Chien-Sheng Wu. Art or artifice? large language models and the false promise of creativity. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, volume 70, pages 1–34, New York, NY, USA, May 2024. ACM.
- [9] Tuhin Chakrabarty, Vishakh Padmakumar, Faeze Brahman, and Smaranda Muresan. Creativity support in the age of large language models: An empirical study involving professional writers. In *Creativity and Cognition*, New York, NY, USA, June 2024. ACM.
- [10] Honghua Chen and Nai Ding. Probing the “creativity” of large language models: Can models produce divergent semantic association? In *The 2023 Conference on Empirical Methods in Natural Language Processing*, December 2023.
- [11] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J L Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiaoshi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojuan Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R J Chen, R L Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S S Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T Wang, Tao Yun, Tian Pei, Tianyu Sun, W L Xiao, Wangding Zeng, Wanbiao Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X Q Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuhen Lin, Y K Li, Y Q Wang, Y X Wei, Y X Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z F Wu, Z Z Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. DeepSeek-V3 technical report. *arXiv [cs.CL]*, December 2024.



- [12] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z F Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J L Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R J Chen, R L Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S S Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W L Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X Q Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y K Li, Y Q Wang, Y X Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y X Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z Z Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv [cs.CL]*, January 2025.
- [13] Jaafar El-Murad and Douglas C West. The definition and measurement of creativity: What do we know? *J. Advert. Res.*, 44(02):188–201, June 2004.
- [14] Xinyu Fang, Zhijian Chen, Kai Lan, Lixin Ma, Shengyuan Ding, Yingji Liang, Xiangyu Zhao, Farong Wen, Zicheng Zhang, Guofeng Zhang, Haodong Duan, Kai Chen, and Dahua Lin. Creation-MMBench: Assessing context-aware creative intelligence in MLLM. *arXiv [cs.CV]*, March 2025.
- [15] Ronald A Finke, Thomas B Ward, and Steven M Smith. *Creative Cognition: Theory, research, and applications*. The MIT Press, October 1992.
- [16] Giorgio Franceschelli and Mirco Musolesi. Creativity and machine learning: A survey. *ACM Comput. Surv.*, 56(11), June 2024. ISSN 0360-0300. doi: 10.1145/3664595. URL <https://doi.org/10.1145/3664595>.
- [17] Fabricio Goes, Marco Volpe, Piotr Sawicki, Marek Grzes, and Jacob Watson. Pushing GPT’s creativity to its limits: Alternative uses and torrance tests. In *14th International Conference on Computational Creativity 2023*, 2023.
- [18] Google DeepMind. Gemini 1.5 and 2.0: Next-gen multimodal models, 2024. URL <https://deepmind.google/technologies/gemini/>. Accessed: 2025-04-30.
- [19] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab

AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-badur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Conguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyan Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Sweeney, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damla, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelen, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt,

- Madian Khabisa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojuan Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models. *arXiv [cs.AI]*, July 2024.
- [20] Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafford, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah Smith, and Hannaneh Hajishirzi. OLMo: Accelerating the science of language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15789–15809, Stroudsburg, PA, USA, 2024. Association for Computational Linguistics.
- [21] J P Guilford, Paul R Christensen, Philip R Merrifield, and Robert C Wilson. Alternate uses, June 2012. Title of the publication associated with this dataset: PsycTESTS Dataset.
- [22] Carlos Gómez-Rodríguez and Paul Williams. A confederacy of models: A comprehensive evaluation of LLMs on creative writing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14504–14528, Stroudsburg, PA, USA, December 2023. Association for Computational Linguistics.
- [23] He He, Nanyun Peng, and Percy Liang. Pun generation with surprise. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1734–1744, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [24] Zicong He, Boxuan Zhang, Weihao Liu, Ruixiang Tang, and Lu Cheng. What shapes a creative machine mind? comprehensively benchmarking creativity in foundation models. *arXiv [cs.AI]*, October 2025.
- [25] K.J. Holyoak and R.G. Morrison. *The Cambridge Handbook of Thinking and Reasoning*. Cambridge Handbooks in Psychology. Cambridge University Press, 2005. ISBN 9780521824170. URL <https://books.google.com/books?id=znbkHaC8QeMC>.

- [26] Weiping Hu and Philip Adey. A scientific creativity test for secondary school students. *Int. J. Sci. Educ.*, 24(4):389–403, April 2002.
- [27] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report. *arXiv [cs.CL]*, September 2024.
- [28] Mete Ismayilzada, Debjit Paul, Antoine Bosselut, and Lonneke van der Plas. Creativity in ai: Progresses and challenges, 2024. URL <https://arxiv.org/abs/2410.17218>.
- [29] Mete Ismayilzada, Claire Stevenson, and Lonneke van der Plas. Evaluating creative short story generation in humans and large language models. *arXiv [cs.CL]*, November 2024.
- [30] Mete Ismayilzada, Claire Stevenson, and Lonneke van der Plas. Evaluating creative short story generation in humans and large language models. *arXiv preprint arXiv:2411.02316*, 2024.
- [31] Shomik Jain, Jack Lanchantin, Maximilian Nickel, Karen Ullrich, Ashia Wilson, and Jamelle Watson-Daniels. LLM output homogenization is task dependent. *arXiv [cs.CL]*, September 2025.
- [32] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7B. *arXiv [cs.CL]*, October 2023.
- [33] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mixtral of experts. *arXiv [cs.LG]*, January 2024.
- [34] Jaehun Jung, Faeze Brahman, and Yejin Choi. Trust or escalate: LLM judges with provable guarantees for human agreement. In *The Thirteenth International Conference on Learning Representations*, October 2024.
- [35] Koray Kavukcuoglu. Gemini 2.5: Our most intelligent AI model. <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>, March 2025. Accessed: 2025-4-30.
- [36] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [37] Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. From generation to judgment: Opportunities and challenges of llm-as-a-judge, 2024.
- [38] Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. Llms-as-judges: A comprehensive survey on llm-based evaluation methods, 2024. URL <https://arxiv.org/abs/2412.05579>.
- [39] Li-Chun Lu, Miri Liu, Pin-Chun Lu, Yufei Tian, Shao-Hua Sun, and Nanyun Peng. Rethinking creativity evaluation: A critical analysis of existing creativity evaluations. *arXiv [cs.CL]*, August 2025.
- [40] Ximing Lu, Melanie Sclar, Skyler Hallinan, Niloofar Mireshghallah, Jiacheng Liu, Seungju Han, Allyson Ettinger, Liwei Jiang, Khyathi Chandu, Nouha Dziri, and Yejin Choi. AI as humanity’s salieri: Quantifying linguistic creativity of language models via systematic attribution of machine text against web text. In *The Thirteenth International Conference on Learning Representations*, October 2024.

- [41] Ximing Lu, Melanie Sclar, Skyler Hallinan, Niloofar Mireshghallah, Jiacheng Liu, Seungju Han, Allyson Ettinger, Liwei Jiang, Khyathi Raghavi Chandu, Nouha Dziri, and Yejin Choi. Ai as humanity’s salieri: Quantifying linguistic creativity of language models via systematic attribution of machine text against web text. *CoRR*, abs/2410.04265, 2024. URL <https://doi.org/10.48550/arXiv.2410.04265>.
- [42] Yining Lu, Dixuan Wang, Tianjian Li, Dongwei Jiang, Sanjeev Khudanpur, Meng Jiang, and Daniel Khashabi. Benchmarking language model creativity: A case study on code generation. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2776–2794, 2025.
- [43] Aidan McLaughlin, Anuja Uppuluri, and James Campbell. AidanBench: Evaluating novel idea generation on open-ended questions. In *Language Gamification - NeurIPS 2024 Workshop*, December 2024.
- [44] Saif Mohammad. Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 174–184, Stroudsburg, PA, USA, 2018. Association for Computational Linguistics.
- [45] Jacqueline N. Lane, Leonard Boussieux, Charles Ayoubi, Ying Hao Chen, Camila Lin, Rebecca Spens, Pooja Wagh, and Pei-Hsin Wang. The narrative AI advantage? a field experiment on generative AI-augmented evaluations of early-stage innovations. *Social Science Research Network*, August 2024.
- [46] Jay A Olson, Johnny Nahas, Denis Chmoulevitch, Simon J Cropper, and Margaret E Webb. Naming unrelated words predicts creativity. *Proc. Natl. Acad. Sci. U. S. A.*, 118(25):e2022340118, June 2021.
- [47] OpenAI. Gpt-4 technical report, 2024. URL <https://openai.com/research/gpt-4>. Accessed: 2025-04-30.
- [48] Peter Organisciak, Selcuk Acar, Denis Dumas, and Kelly Berthiaume. Beyond semantic distance: Automated scoring of divergent thinking greatly improves with large language models. *Think. Skills Creat.*, 49(101356):101356, September 2023.
- [49] Vishakh Padmakumar and He He. Does writing with language models reduce content diversity? In *The Twelfth International Conference on Learning Representations*, October 2023.
- [50] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1162. URL <https://aclanthology.org/D14-1162/>.
- [51] Jeffrey Pennington, Richard Socher, and Christopher D Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, October 2014.
- [52] Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv [cs.CL]*, December 2024.
- [53] Sara Rosengren, Martin Eisend, Scott Koslow, and Micael Dahlen. A meta-analysis of when and how advertising creativity works. *J. Mark.*, 84(6):39–56, November 2020.
- [54] Sameh Said-Metwaly, Wim Van den Noortgate, and Eva Kyndt. Approaches to measuring creativity: A systematic literature review. *Creativity. Theories – Research - Applications*, 4(2): 238–275, December 2017.

- [55] Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can LLMs generate novel research ideas? a large-scale human study with 100+ NLP researchers. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=M23dTGWCZy>.
- [56] Robert E Smith, Scott B MacKenzie, Xiaojing Yang, Laura M Buchholz, and William K Darley. Modeling the determinants and effects of creativity in advertising. *Mark. Sci.*, 26(6):819–833, November 2007.
- [57] R Sternberg and T Lubart. An investment theory of creativity and its development. *Human Development*, 34(1):1–31, June 1991.
- [58] Sijun Tan, Siyuan Zhuang, Kyle Montgomery, William Yuan Tang, Alejandro Cuadron, Chengguang Wang, Raluca Popa, and Ion Stoica. Judgebench: A benchmark for evaluating LLM-based judges. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=G0dksFayVq>.
- [59] Yufei Tian, Tenghao Huang, Miri Liu, Derek Jiang, Alexander Spangher, Muhao Chen, Jonathan May, and Nanyun Peng. Are large language models capable of generating human-level narratives? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17659–17681, Stroudsburg, PA, USA, November 2024. Association for Computational Linguistics.
- [60] Yufei Tian, Abhilasha Ravichander, Lianhui Qin, Ronan Le Bras, Raja Marjeh, Nanyun Peng, Yejin Choi, Thomas Griffiths, and Faeze Brahman. MacGyver: Are large language models creative problem solvers? In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5303–5324, Stroudsburg, PA, USA, 2024. Association for Computational Linguistics.
- [61] Yufei Tian, Abhilasha Ravichander, Lianhui Qin, Ronan Le Bras, Raja Marjeh, Nanyun Peng, Yejin Choi, Thomas Griffiths, and Faeze Brahman. MacGyver: Are large language models creative problem solvers? In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5303–5324, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.297. URL <https://aclanthology.org/2024.naacl-long.297/>.
- [62] Junyi Ye, Jingyi Gu, Xinyun Zhao, Wenpeng Yin, and Guiling Wang. Assessing the creativity of LLMs in proposing novel solutions to mathematical problems. *arXiv [cs.CL]*, October 2024.
- [63] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. BERTScore: Evaluating text generation with BERT. *arXiv [cs.CL]*, April 2019.
- [64] Yiming Zhang, Harshita Diddee, Susan Holm, Hanchen Liu, Xinyue Liu, Vinay Samuel, Barry Wang, and Daphne Ippolito. NoveltyBench: Evaluating creativity and diversity in language models. In *Second Conference on Language Modeling*, August 2025.
- [65] Yunpu Zhao, Rui Zhang, Wenyi Li, and Ling Li. Assessing and understanding creativity in large language models. *Mach. Intell. Res.*, pages 1–20, April 2025.
- [66] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-bench and chatbot arena. *arXiv [cs.CL]*, June 2023.
- [67] Shanshan Zhong, Zhongzhan Huang, Shanghua Gao, Wushao Wen, Liang Lin, Marinka Zitnik, and Pan Zhou. Let’s think outside the box: Exploring leap-of-thought in large language models with creative humor generation. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13246–13257, 2024. doi: 10.1109/CVPR52733.2024.01258.

## A Model Details

Short Name	Exact Model Name	Size	Family	Release Time
Mistral-7B	Mistral-7B-Instruct-v0.3	7B	Mistral	05/2024
Qwen2.5-7B	Qwen2.5-7B-Instruct	7B	Qwen	09/2024
OLMo2-7B	OLMo-2-1124-7B-Instruct	7B	Olmo	11/2024
Llama3.1-8B	Llama-3.1-8B-Instruct	8B	Llama	07/2024
OLMo2-13B	OLMo-2-1124-13B-Instruct	13B	Olmo	11/2024
OLMo2-13B-SFT	OLMo-2-1124-13B-SFT	13B	Olmo	11/2024
OLMo2-13B-DPO	OLMo-2-1124-13B-DPO	13B	Olmo	11/2024
Mistral-24B	Mistral-Small-24B-Instruct-2501	24B	Mistral	01/2025
Qwen2.5-32B	Qwen2.5-32B-Instruct	32B	Qwen	09/2024
Mixtral-8x7B	Mixtral-8x7B-Instruct-v0.1	56B	Mistral	12/2023
Llama3.3-70B	Llama-3.3-70B-Instruct	70B	Llama	12/2024
Qwen2.5-72B	Qwen2.5-72B-Instruct	72B	Qwen	09/2024
Claude3-Sonnet	claude-3-7-sonnet-20250219	-	Claude	02/2025
Claude3-Haiku	claude-3-5-haiku-20241022	-	Claude	11/2024
GPT4.1	gpt-4.1-2025-04-14	-	GPT	04/2025
GPT4.1-mini	gpt-4.1-mini-2025-04-14	-	GPT	04/2025
Gemini2.0-Flash	gemini-2.0-flash	-	Gemini	12/2024
Deepseek-R1	deepseek-reasoner	-	Gemini	01/2025
Deepseek-V3	deepseek-chat	-	Gemini	12/2024

Table 3: List of models included in our benchmark.

**Deepseek Models** For Deepseek models, we also use API due to constraints in compute resources.  
API console: <https://platform.deepseek.com>.



## B Benchmark Design

### B.1 Dataset Sizes

Task	Count	Note
AUT	105 (tool use)	21 tools with 5 rounds of prompting per tool
DAT	100 (round)	No input data, we prompt each LLM 100 rounds
TTCT	700 (question)	7 tasks (100 questions/task)
TTCW	12 (story prompt)	One story per story prompt
Creative Short Story	10 (keyword tuple)	One story per keyword tuple
Creativity Index	300 (document sample)	100 samples from 3 subsets: book, poem, and speech
CS4	250 (story)	50 base stories with 5 constraint configurations per story
NeoCoder	198 (question)	One solution per coding question
Creative Math	400 (question)	One solution per math question

Table 4: Dataset size of CREATIVITYPRISM. More details can be found in the corresponding section of Appendix E.

### B.2 Metrics

Table 5 shows a complete list to metrics in CREATIVITYPRISM, grouped by tasks. More details about how each metric is calculated can be found in corresponding sections in Appendix E.

Task	Quality	Novelty	Diversity
AUT	-	AUT Score	-
DAT	-	-	DAT Score
TTCT	Fluency, Elaboration	Originality	Flexibility
TTCW	Coherence, Ending, Elaboration	-	Emotional Flexibility
Creative Short Story	-	Novelty Score, Surprise-ness	N-gram Diversity
Creativity Index	-	L-uniqueness	-
CS4	QUC	RCS	Dist-N
NeoCoder	Convergence@k	Divergent@k	-
Creative Math	Correctness Ratio	Novelty Ratio	-

Table 5: Evaluation metrics in CREATIVITYPRISM; “-” means this task (row) does not have any metric in the corresponding creativity dimension.

### B.3 Score Calculations

**Score Normalization** For every model  $i$  and every raw metric score  $S_{i,m}$  (metric  $m$  lives on some known scale  $[\min_m, \max_m]$ ), the normalized score  $\hat{S}_{i,m}$  is given by:

$$\hat{S}_{i,m} = \frac{S_{i,m} - \min_m}{\max_m - \min_m}$$

For example, AUT score is on a 1–5 Likert scale:  $\hat{S}_{i,AUT} = \frac{S_{i,AUT} - 1}{5 - 1} = \frac{S_{i,AUT} - 1}{4}$ .

**Aggregate Normalized Scores** First, we collapse multiple metrics within the same task: if task  $t$  has a set  $M_t$  of  $k_t$  metrics in a given dimension (e.g. three quality metrics for TTCW), average them first:

$$\bar{S}_{i,t} = \frac{1}{k_t} \sum_{m \in M_t} \hat{S}_{i,m}$$

Then, we take average across all tasks that belong to that dimension. Let  $T_{\text{qual}}, T_{\text{nov}}, T_{\text{div}}$  be the task sets for quality, novelty, diversity. For dimension  $d \in \{\text{qual}, \text{nov}, \text{div}\}$ :

$$D_i^{(d)} = \frac{1}{|T_d|} \sum_{t \in T_d} \bar{S}_{i,t}$$

In this way, we end up with three numbers per model:  $D_i^{(\text{qual})}$ ,  $D_i^{(\text{nov})}$ ,  $D_i^{(\text{div})}$ . We can also calculate aggregated score for creative writing, divergent thinking, and logical reasoning (as shown in Table 2).

**Overall creativity score** Just take the straight mean of those three dimension scores to stay balanced:

$$C_i = \frac{D_i^{(\text{qual})} + D_i^{(\text{nov})} + D_i^{(\text{div})}}{3}$$

## C Performance Summaries

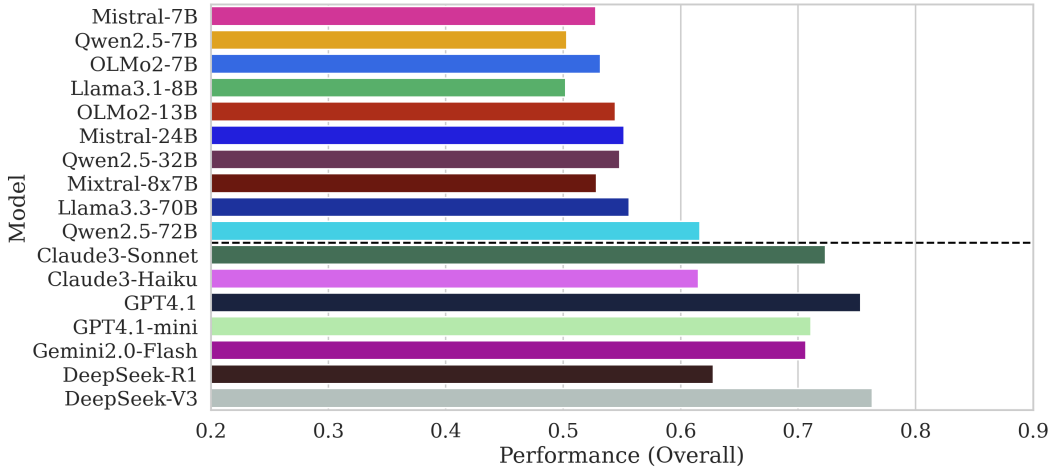


Figure 5: Overall performances.

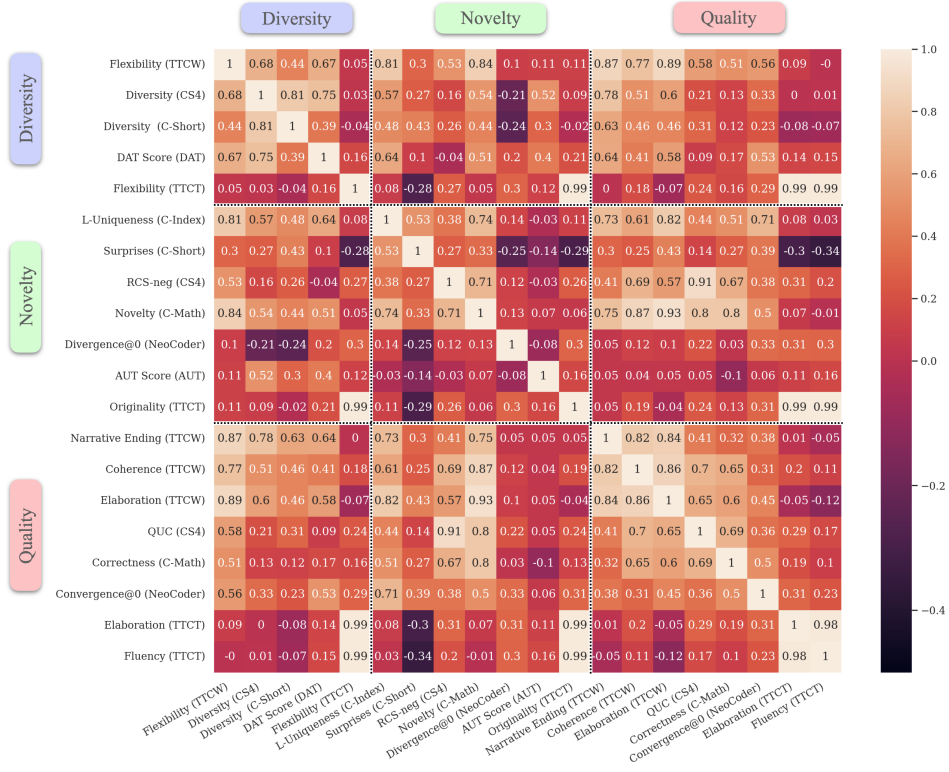


Figure 6: Inter-metric correlation (grouped by creativity taxonomy).

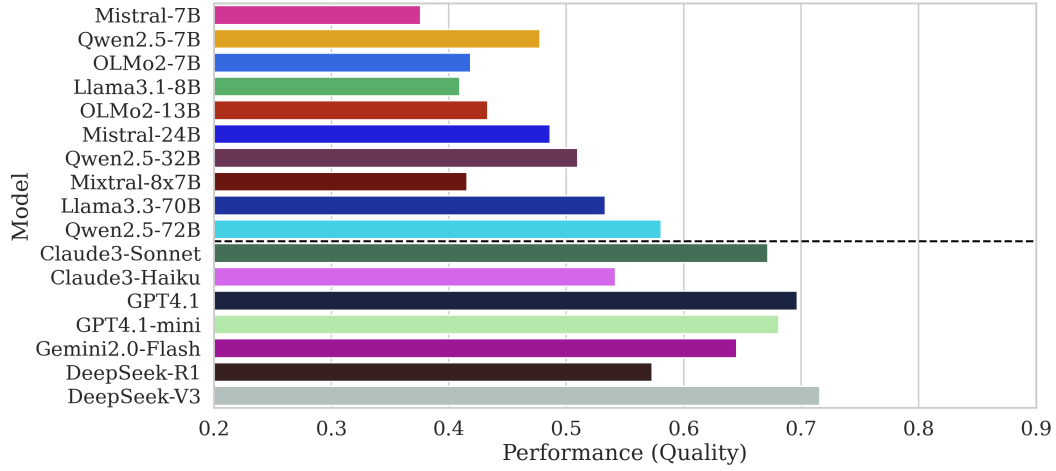


Figure 7: Performance on quality dimension

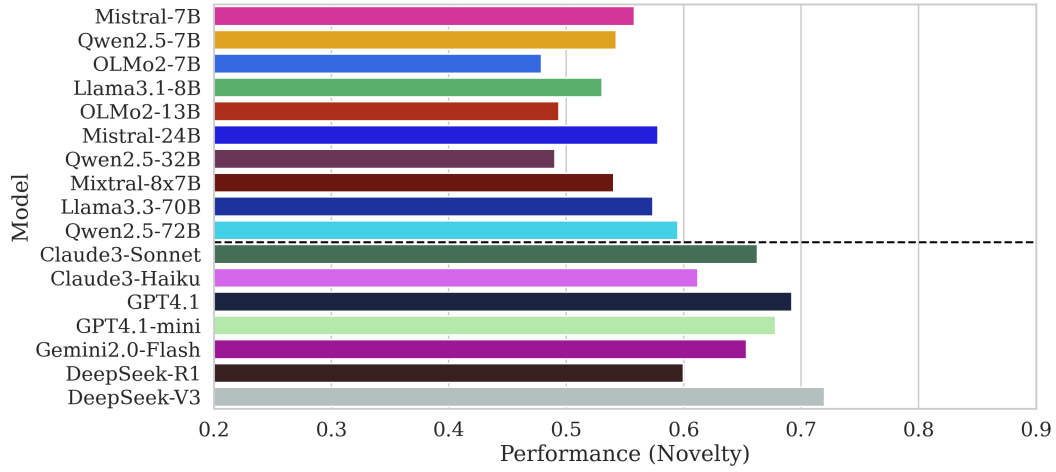


Figure 8: Performance on novelty dimension

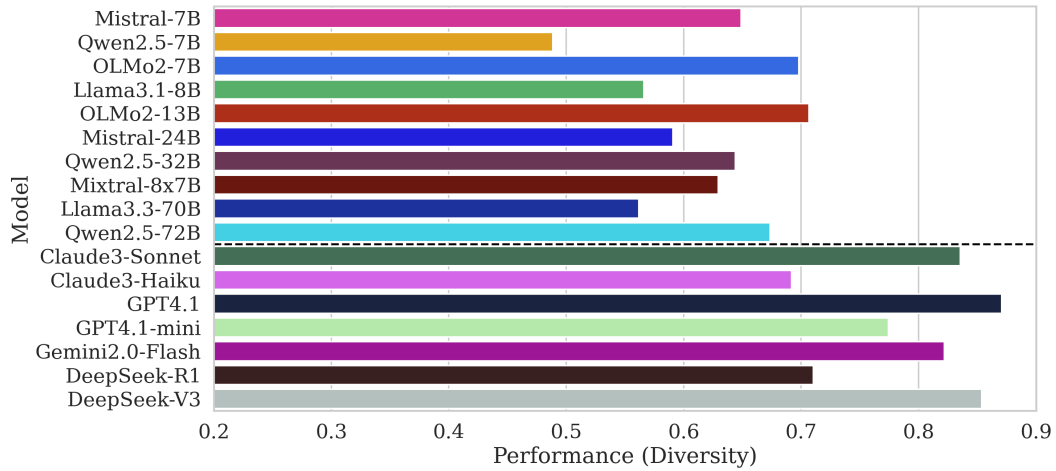


Figure 9: Performance on diversity dimension

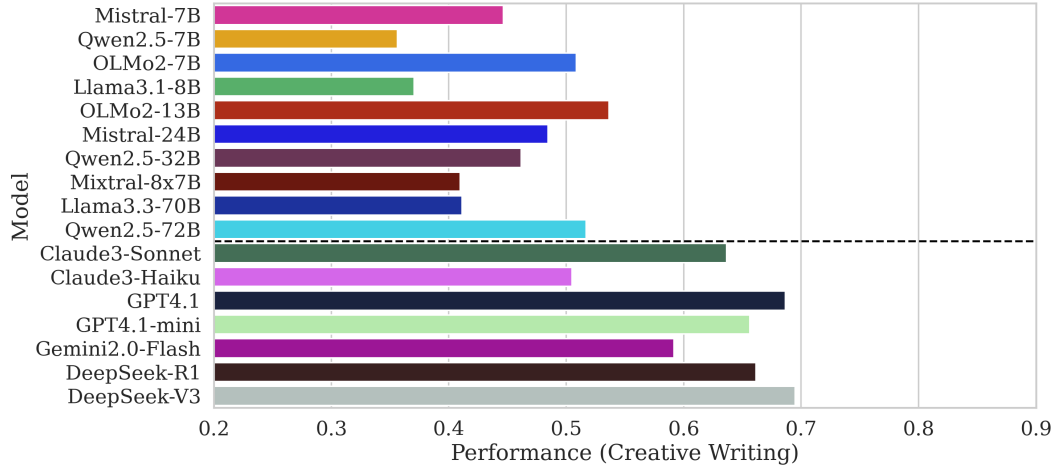


Figure 10: Performance on creative writing tasks

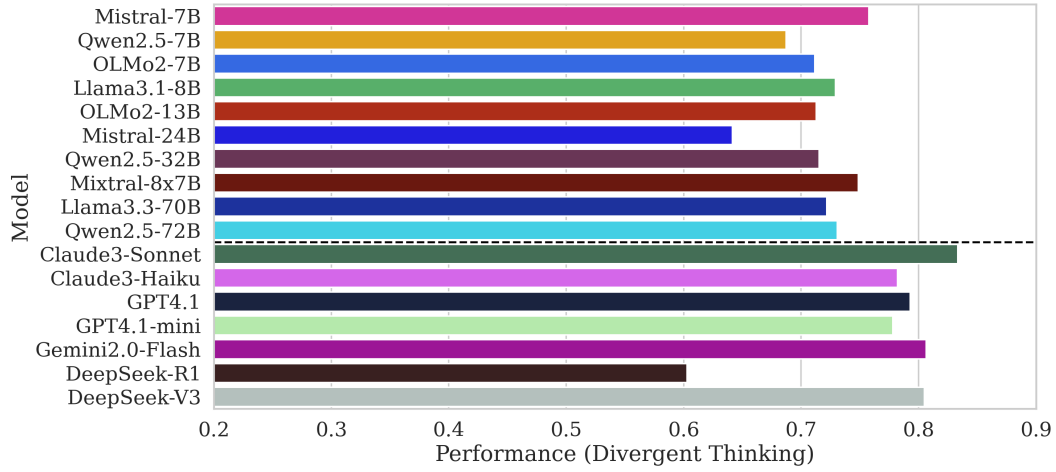


Figure 11: Performance on divergent thinking tasks

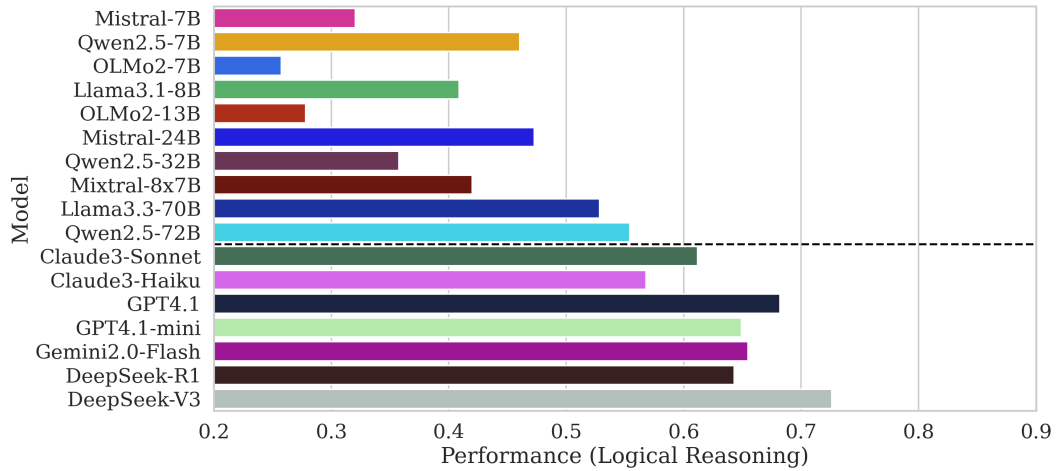


Figure 12: Performance on logical reasoning dimension

## D LLM-as-a-Judge Design Details

Six out of nine tasks in our benchmark require LLM-Judge for one or more metrics. We use Qwen2.5-72B as the default LLM-Judge model, unless otherwise specified. The choice of Qwen2.5-72B is based on a pilot study on TTCW and AUT task, where Qwen2.5-72B is the open-source LLM (within our compute budget) that correlates the best either with humans or with closed-source LLMs that are typically used as LLM-Judge (as detailed below).

To ensure the reliability of the LLM-Judge, we follow the following principles: if the original paper has reported human-LLM-Judge agreement and we are using the same setup, we directly report the original paper’s agreement; otherwise, if we have human annotation, we’ll calculate the agreement between human judgment and LLM-Judge’s judgement; if no human annotation is available, we’ll calculate the agreement between judgement by Qwen and by the closed-source LLM-Judge in the original paper.

For binary labels (TTCW, Creative Math), agreement refers to LLM-Judge prediction accuracy (with human annotations being ground truth); for likert-scale labels (AUT, CS4, TTCT), agreement refers to Pearson correlation between human and LLM-Judge; the remaining task, NeoCoder, is a special case because the LLM-Judge is not directly used to generate the metric and we will provide more detail below. Here we report the correlation statistics:

**TTCW:** The original authors [8] provide 36 machine-generated stories with 3 creative-writing expert annotations (binary) for each. We check the quality by calculating agreement of the Qwen2.5-72B judgement (binary) and expert majority vote results (also binary). Here are the metrics (numbers in parentheses are accuracy): Narrative Ending (0.69), Understandability and Coherence (0.78), Emotional Flexibility (0.86), World Building and Setting (0.72). Note that, because in story creativity evaluation, the human-to-human agreement is also relatively low, we have a lower acceptance threshold when we are considering what kind of human-LLM agreement level is acceptable.

**Creative Math:** we annotated the output solution of 50 questions and observed 0.78 agreement in novelty judgement (from a multi-LLM majority voting judge that consists of Gemini-2.0-Flash, GPT-4.1, and Claude-3.7-Sonnet) and 0.94 agreement for correctness (from the Claude-3.7-Sonnet for correctness judgement). Agreement here is the same as accuracy if we consider human annotation as ground truth.

**AUT:** Previous work [48] specifically studied the feasibility of using few-shot LLM-Judge to evaluate AUT output. They have shown GPT-4 (it was published before GPT-4o) with 20-shot examples can achieve 0.70 Pearson correlation between humans’ judgement. We use the same prompt and the human annotation released by that paper as the 20-shot examples. We then use GPT-4o and Qwen2.5-72B to judge the same set of AUT outputs (generated by Llama3.3-70B). The scores from GPT-4o and Qwen2.5-72B have a Pearson’s  $r$  of 0.597 ( $p < 0.01$ ).

**CS4:** The original authors [3] collected and provided the annotation (from Amazon Mechanical Turk) of 15 machine-generated stories, with 2 annotations per story. For constraint satisfaction metric, the Pearson correlation between Qwen2.5-72B judgments and human judgments is 0.55 ( $p < 0.01$ ).

**TTCT:** The original paper [65] used GPT-4 as evaluator and conducted a small human v.s. LLM-Judge alignment study. However, that data is not accessible to us, so we studied the Pearson Correlation between Qwen2.5-72B and GPT-4.1 judgement and here are the numbers we got: Fluency: 0.6884 ( $p < 0.01$ ), Flexibility: 0.6592 ( $p < 0.01$ ), Originality: 0.5152 ( $p < 0.01$ ), Elaboration: 0.5033 ( $p < 0.01$ )

**NeoCoder:** LLM-Judge is not directly making an evaluation output in this task. Instead, GPT-4 is used to detect if a generated solution is novel or not when compared to a pre-collected human programmer solution. To ensure high quality of LLM-Judge, the original paper [42] uses recall instead of accuracy to evaluate human-LLM agreement (because we want to ensure all the “non-novel” techniques are detected). GPT-4 achieves 0.94 recall in detecting non-novel solutions.

## E Task Details

### E.1 Torrance Test of Creative Writing (TTCW)

#### E.1.1 Dataset

The dataset consists of 12 New Yorker Stories’ plots, i.e., GPT-4 generated summary of the original story <sup>6</sup>.

#### E.1.2 Example

##### Plot

A woman experiences a disorienting night in a maternity ward where she encounters other similarly disoriented new mothers, leading to an uncanny mix-up where she leaves the hospital with a baby that she realizes is not her own, yet accepts the situation with an inexplicable sense of happiness.

##### Inference Prompt

Write a New Yorker-style story given the plot below. Make sure it is at least {word\_count} words. Directly start with the story, do not say things like "Here’s the story [...]" Plot: {plot} Story:

#### E.1.3 Experiment Configurations

- Temperature: 0.75
- Max Token: 4096
- Top-p: 1

#### E.1.4 Evaluation Metrics

As mentioned in §3, we use a subset of questions from the original paper where our few-shot LLM-as-a-judge evaluator achieves a correlation of more than 0.2 between the human majority vote and the evaluator model’s judgments (Narrative Ending: 0.29; Understandability and Coherence: 0.45; Emotional Flexibility: 0.21; World Building and Setting: 0.40). The sample size is 36, which are GPT4, GPT3.5, and Claude-generated stories, with human expert annotation, all released by the original paper. Each question corresponds to one metric. Since the evaluation is binary for each generated story, we calculate the proportion of generated stories that pass each question as the final evaluation metric (e.g., if 3 out of 12 stories pass the “Understandability and Coherence” question, then the “Understandability and Coherence” metric is 0.25).

We use two-shot examples (one positive and one negative) in the evaluation prompt, as previous work shows adding few-shot examples improves human-llm alignments [34].

---

<sup>6</sup>[https://github.com/salesforce/creativity\\_eval](https://github.com/salesforce/creativity_eval)



### Evaluation Prompt

You are given a creative short story. Read it carefully. You are then given some background about specific aspects of creative writing, a binary (Yes/No) question, and sample stories with expert-annotated answers to the same question. Your objective is to use the background information and sample stories to answer the question about the story. Provide your answer in the format of "\*\*\*Answer\*\*": [Yes/No]". You can optionally then provide a short explanation for your answer.

=====

Question:

{full\_prompt}

Examples:

=====

Story: {story}

Answer: {answer}

Explanations: {exp}

=====

Story: {story}

Answer: {answer}

Explanations: {exp}

=====

Story: {story}

Based on the question and examples above, answer the question (Provide your answer in the format of "\*\*\*Answer\*\*": [Yes/No]". You can optionally then provide a short explanation for your answer). Make sure you are extra harsh on the decision (most answers should be negative).

Answer:

#### E.1.5 Model Performances

Model	Narrative Ending	Understandability and Coherence	Emotional Flexibility	World Building and Setting
Mistral-7B	0.17	0.08	0.00	0.00
Qwen2.5-7B	0.08	0.17	0.00	0.00
OLMo2-7B	0.75	0.25	0.25	0.00
Llama3.1-8B	0.00	0.08	0.00	0.00
OLMo2-13B	0.67	0.33	0.25	0.08
Mistral-24B	0.25	0.25	0.00	0.00
Qwen2.5-32B	0.00	0.17	0.00	0.00
Mixtral-8x7B	0.17	0.08	0.08	0.08
Llama3.3-70B	0.00	0.33	0.00	0.00
Qwen2.5-72B	0.42	0.50	0.50	0.17
Claude3-Sonnet	0.75	0.58	0.58	0.42
Claude3-Opus	0.33	0.08	0.08	0.08
GPT-4.1	1.00	0.67	0.83	0.50
GPT-4.1-mini	1.00	0.83	0.42	0.50
Gemini2.0-Flash	0.83	0.42	0.42	0.17
DeepSeek-R1	0.83	0.50	0.50	0.58
DeepSeek-V3	0.83	0.50	0.50	0.50

Table 6: Model performance on TTCW.

## E.2 CS4

### E.2.1 Dataset

There are 50 base stories. During inference time, for each base story, 5 constraints are applied; to calculate n-gram diversity, we also generate 3 for each constraint. Therefore, for each given LLM, 750 stories will be generated <sup>7</sup>.

### E.2.2 Examples

#### Base Story

Evelyn was introverted by nature, more comfortable in her world of books than in social groups. What made her exceptional, however, was her latent, uncontrollable intuition. She was highly empathic, experiencing others' emotions as vividly as they did themselves - sometimes, even more so. The holiday season, with its swirl of emotional undercurrents, was particularly overwhelming for Evelyn, and she dreaded it every year.

It was the eve of Christmas. Evelyn's small apartment was bursting at the seams with relatives she hadn't seen in years. Overwhelmed, she nestled into the corner of the couch, trying to make herself invisible amid the sea of forced cheerfulness. Any close interaction was a sensory overload for her; she would involuntarily feel their feelings, read their deepest thoughts, and become emotionally exhausted.

Suddenly, the door swung open, and her cousin, Ava, strode in, a torrent of unspoken emotions swirling around her. A shiver of apprehension ran through Evelyn as she braced herself, feeling the churning storm of emotions through her intuitive empathy.

Ava looked around the room, her eyes widening slightly as she noticed Evelyn. Evelyn felt Ava's shockwave of anxiety and self-consciousness hit her like a punch. Despite being at her most vulnerable, Ava hid it well under a veneer of social amiability. But Evelyn's perceptive intuition saw right through it.

Taking a deep breath to steel herself, she wandered over to Ava. Silent conversations, almost piercing in their intensity, flowed between them.

""Are you okay, Ava?"" She asked finally, breaking the momentary silence.

Ava's eyes widened in surprise. ""Yeah, I'm fine. Why do you ask?"" But the underlying tremor in her voice betrayed the turmoil within.

Evelyn quietly admitted, ""I can tell something's wrong.""

Ava looked at her silently for a long moment before sighing deeply. ""I lost my job, Evelyn,"" she confessed, her voice barely above a whisper. Evelyn felt a tidal wave of despair crash over her as Ava's feelings of hopelessness washed over her.

Evelyn, despite her own burden of emotions, took Ava's hands in hers, feeling the tremors running through them. ""You're not alone, Ava. You have us. We'll sort it out together,"" she said, her voice reassuring.

Armed with her intuitive empathy, Evelyn spent the rest of the evening comforting Ava, helping her cope with the raw wound of job loss. It was an extremely challenging, emotionally draining journey, but Evelyn's heart swelled at Ava's gradual shift from despair to a glimmer of hope and optimism.

That Christmas Eve, Evelyn, buoyed by Ava's resilience, also discovered something about herself. Her gift, which she had despised for its uncontrollability, for how much it drained and overwhelmed her, could also be used to help others.

As Evelyn watched Ava slowly blend back into the crowd, her heart lighter, the usual cacophony of emotions she feared seemed more bearable. Evelyn realized that while her introverted nature and intuitive empathy made the holiday season challenging, it was also what made her essential in processing these unspoken struggles. It wasn't a curse; it was her gift. A gift of understanding, of empathy, of being the silent pillar of comfort in a room filled with concealed emotions.

---

<sup>7</sup>[https://github.com/anirudhlakkaraju/cs4\\_benchmark](https://github.com/anirudhlakkaraju/cs4_benchmark)

## Constraints

1. The protagonist suffers physical discomfort when overwhelmed by emotions (nausea, shaking, etc.).
2. The protagonist is challenged by the need to engage in public spaces.
3. The unknown man realizes that the protagonist can feel his emotions.
4. The protagonist uses humor and sarcasm to cope with her situation.
5. The protagonist is an introverted character.
6. The story includes communication via text messages.
7. The story is set in a Starbucks on Michigan in Chicago a week before Christmas.
8. The protagonist is forced to leave the meeting early due to being overwhelmed.
9. The protagonist desires to live a more normal life despite her unique condition.
10. There exists a vaccine for controlling intuition.
11. The protagonist devises coping strategies for managing her anxiety in public places.
12. Scientists are working to find a solution for people who can't use the intuition vaccine.
13. Tiffany threatens the protagonist to meet her.
14. The protagonist struggles with accepting her condition.
15. The protagonist must grapple with the thoughts and feelings of others in the Starbucks.
16. The man looks at the protagonist with both desire and love.
17. The protagonist is physically attractive.
18. The protagonist encounters an unknown man who causes powerful and unique emotions.
19. There is societal disapproval for people whose intuitions cannot be controlled by the vaccine.
20. The protagonist has a heightened intuition.
21. The setting should be during the holiday season.
22. The protagonist feels other people's emotions intensely.
23. The protagonist struggles with disentangling their own feelings from others'.
24. Tiffany is a strong-willed and passionate character.
25. The protagonist reluctantly acknowledges being a potential "crazy cat lady".
26. Puberty is identified as a critical time for the progression of intuition powers.
27. The protagonist and Tiffany were inseparable until puberty.
28. The protagonist has personal hygiene items (travel mouthwash) handy.
29. This vaccine doesn't work for the protagonist due to a genetic mutation.
30. Characters should express understanding of the protagonist's predicament.
31. The protagonist's primary means of communication with the outside world is through the internet.
32. The protagonist experiences other's thoughts as if they were their own.
33. The protagonist's intuition is uncontrollable due to a genetic mutation.
34. The protagonist's coping mechanisms do not always successfully block out other people's emotions.
35. The protagonist uses strategies to block out the feelings of others, such as counting letters on the menu board.
36. The protagonist finds solace in the idea of drinking coffee.
37. Include a hint of romance in the story.
38. There is societal pressure to control intuition with the vaccine.
39. The protagonist prefers isolation to manage their heightened intuition."

### E.2.3 Experiment Configurations

- Temperature: 0.75
- Max Token: 4096
- Top-p: 1

### E.2.4 Inference Prompt

#### Inference Prompt

User: Write a story in less than 500 words about {story theme}  
Base Story: {base story}  
User Instruction: " Now modify the existing story to accommodate the following constraints: {selected constraints} into the LLM generated story and come up with a new story in 500 words.

### E.2.5 Evaluation Metrics and Prompt

From [3], we included QUC@39 (quality), RCS-7-39 (novelty) and Dist-N@39 (diversity) as evaluation metrics. In particular, since lower RCS means more stable, we take  $RCS_{neg} = 1 - RCS$  as the evaluation metric to ensure all metrics we included describe a feature that is positively correlated with creativity.

Here we include evaluation prompt for constraint satisfaction and story quality (coherence) evaluation prompt.

#### Evaluation Prompt - Constraint Satisfaction

You are an expert reader. I will give you a story followed by a set of constraints. Your task is to carefully read both of them and tell how many constraints are being satisfied in the story. As the output, I want you to print yes/no for each constraint based on whether it is being satisfied or not, followed by a 1 line explanation of why it is being satisfied/violated. In case a constraint is being satisfied, print the sentence/line from the story in which it is being satisfied. If a constraint is not being satisfied, give an explanation of how it is being violated. Be very strict in your evaluation. Mark a constraint as satisfied ("yes") only if it is being completely satisfied in the story. For no satisfaction/partial satisfaction, mark a "no". If the story is empty (no input provided), all constraints are considered NOT satisfied. Finally, print the number of constraints that are being satisfied. Follow the examples and Output the ending of the evaluation in the same format. Number of constraints satisfied: [number]

Here are some examples - Input Story: {story} Constraints: 1. Write a story based on the following constraints in less than 377 words.

2. Start the story with the sentence: "Week 18 aboard the Depth Reaver, Circa 2023"
3. Include a revelation of an unexpected large-scale phenomenon observed in space."

Output 1. Yes - The story is 302 words long, meeting the constraint of being less than 377 words.  
2. Yes - The story starts with the exact sentence: "Week 18 aboard the Depth Reaver, Circa 2023".

3. Yes - The revelation of the moon cracking open to reveal a colossal human face qualifies as an unexpected large-scale phenomenon observed in space.

Number of constraints satisfied: 3

{other examples}

Input Story: {story to be evaluated} Constraints: {constraints} Output

### Evaluation Prompt - Story Quality

You are an English writing expert and you can compare and evaluate story essays on these metrics with the following definitions

1. Grammar: Which story has better writing and grammar comparatively?
2. Coherence: Which story has a better logical flow and the writing fits together with respect to the plot?
3. Likability: Which story do you find more enjoyable to read?

You will be given two Stories - Story A and Story B.

Add a rating out of 5 for each category, specify which story you prefer for each metric by responding with just the letter "A" or "B" followed by a hyphen and one line reasoning for your preference.

For each category provide a category winner story as the letter "A" or "B", based on the category ratings.

Finally, assign an overall winner story as the letter "A" or "B" based on the ratings and category wins.

(if an story is empty, give it zero scores)

**IMPORTANT - DO NOT GIVE ANY OTHER TEXT APART FROM THE SCORE, METRICS AND PREFERENCE. FOLLOW THE EXACT FORMAT AS GIVEN IN THE FOLLOWING EXAMPLES.**

EXAMPLE OUTPUT 1:

{example output 1}

EXAMPLE OUTPUT 2:

{example output 2}

Story A: {story1}

Story B: {story2}

SCORE OUTPUT:

Note that the original paper uses OpenAI models as evaluator for both constraint satisfaction and story coherence evaluation. We also investigate open-source alternatives, Qwen2.5-72B, as the LLM-judge. The Pearson correlation between Qwen2.5-72B judgments and human judgment is 0.55, with p-value < 0.01 (sample size: 15 stories, 2 annotations per story).

### E.2.6 Model Performances

Model	QUC@39	RCS-7-39 (neg)	Dist-N@39
Mistral-7B	0.7048	0.8947	0.8675
Qwen2.5-7B	0.6451	0.9143	0.7480
OLMo2-7B	0.5904	0.8491	0.9032
Llama3.1-8B	0.6193	0.8103	0.7808
OLMo2-13B	0.5543	0.8595	0.9200
Mistral-24B	0.7641	0.9417	0.7978
Qwen2.5-32B	0.7201	0.9064	0.8604
Mixtral-8x7B	0.5925	0.8313	0.8871
Llama3.3-70B	0.7392	0.9343	0.7788
Qwen2.5-72B	0.7853	0.9499	0.8484
Claude3-Sonnet	0.8153	0.9557	0.9572
Claude3-Haiku	0.6306	0.8584	0.8728
GPT4.1	0.8157	0.9587	0.9083
GPT4.1-mini	0.7828	0.9339	0.8990
Gemini2.0-Flash	0.7733	0.9349	0.9275
DeepSeek-R1	0.7767	0.9249	0.9022
DeepSeek-V3	0.7809	0.9572	0.9121

Table 7: Model scores on CS4 task.

### E.3 Creativity Index

#### E.3.1 Dataset

The dataset consists of 3 subsets: book, poem, and speech, all are the prefixes (i.e., first line of text) from the dataset proposed by [41]. We use the first 100 examples in generation and evaluation.<sup>8</sup>

#### E.3.2 Examples

Here are some examples of the input data (i.e., the text prefix for LLM to complete).

##### Book

It’s been years: Bailey clearly means him no harm and has managed to be discreet enough that Nick’s queerness isn’t the talk of the Chronicle.

##### Poem

Swiftly walk o’er the western wave,

##### Speech

That is the kind of America in which I believe

#### E.3.3 Evaluation Metrics

**L-uniqueness** Let  $\mathbf{x}$  be a text consisting of a sequence of words whose linguistic creativity we wish to quantify. Let an  $n$ -gram of  $\mathbf{x}$  be any contiguous subsequence of  $n$  words, and denote by  $\mathbf{x}_{i:i+n}$  the  $n$ -gram starting at the  $i$ -th word of  $\mathbf{x}$ . Let  $C$  be a large reference corpus of publicly available texts, and define  $f$  as a binary function that returns 1 if the  $n$ -gram  $\mathbf{x}_{i:i+n}$  occurs anywhere in  $C$ , and 0 otherwise. The  $L$ -uniqueness of  $\mathbf{x}$ , denoted by  $\text{uniq}(\mathbf{x}, L)$ , is defined as the proportion of words  $w \in \mathbf{x}$  such that none of the  $n$ -grams containing  $w$  with  $n \geq L$  occur in  $C$ . Intuitively, a higher  $L$ -uniqueness means a greater proportion of  $\mathbf{x}$ ’s words appear in novel contexts not seen in the reference corpus, thus indicating greater linguistic originality.

<sup>8</sup>[https://github.com/GXimingLu/creativity\\_index](https://github.com/GXimingLu/creativity_index)

### E.3.4 Experiment Configurations

- Temperature: 1
- Top P: 0.9
- Max Token: 288

### E.3.5 Inference Prompt

#### Book

Please write a few paragraphs for a novel starting with the following prompt: {start sentence of book}

#### Poem

Please write a few paragraphs for a novel starting with the following prompt: {start sentence of poem}

#### Speech

Please write a speech starting with the following sentence: {start sentence of speech}

For closed-source models, we included an additional prompt instruction to ensure that the output consists solely of the completed paragraphs, poems, or speeches. These models often preface their responses with phrases like “Certainly” or “Here is...,” which we manually removed during post-processing. In contrast, open-source models typically generate the desired completions directly without such prefatory text. For these models, we reviewed the log outputs and removed any unrelated content as needed.

### E.3.6 Evaluation Metrics

We follow the evaluation metrics outlined in Lu et al. [40], specifically retaining the exact match component. However, we exclude the semantic search-based evaluation due to its high computational cost and sensitivity to the chosen cosine similarity threshold, which significantly affects whether two sentence spans are considered semantically similar. We sum over the L-uniqueness with spans of n-grams from 5 to 12 inclusively to get the total creative index for each response. We average the creative index for each response per mode per task. Data cleaning was done before the evaluation manually to remove irrelevant outputs. Then, we normalize the score by dividing it with 8 (the highest value that the summation could be) to get the final Creativity Index measurement for each model over the three different tasks.

### E.3.7 Model Performance

### E.3.8 Additional Comments

We also note that the generation for OLMo2-13B-instruct may miss some data with the vllm generation. We remove those missing generations. This account for 13 responses in the poem subset, and 10 examples in the speech subset. In addition, the model may resist in answering some prompts. We also removed those generations. For OLMo-7B-instruct, there are 2 cases in the speech subset. For GPT-4.1, there is 1 case in the speech subset.



Model	Book	Poem	Speech	Average
mistral-7b-instruct	0.4496	0.5828	0.3104	0.4476
qwen-7b-instruct	0.4354	0.6310	0.3534	0.4733
olmo-7b-instruct	0.4810	0.6110	0.3727	0.4882
llama-31-8b-instruct	0.4724	0.5700	0.3396	0.4607
olmo2-13b-instruct	0.4860	0.5963	0.3522	0.4782
mistral-24b-instruct	0.4752	0.6646	0.3397	0.4932
qwen-32b-instruct	0.4663	0.6328	0.3465	0.4816
mistral-8x7b-instruct	0.4149	0.6035	0.2804	0.4329
llama-33-70b-instruct	0.4226	0.5802	0.2936	0.4321
Qwen2.5-72B-instruct	0.4133	0.5924	0.3171	0.4409
claude-3-7-sonnet-20250219	0.5615	0.6700	0.4675	0.5663
claude-3-5-haiku-20241022	0.5769	0.7039	0.4519	0.5776
gpt-4.1	0.6044	0.7637	0.4593	0.6091
gpt-4.1-mini	0.5624	0.7147	0.4261	0.5677
gemini-2.0-flash	0.5278	0.6707	0.4121	0.5369
deepseek-reasoner	0.5930	0.7595	0.5410	0.6312
deepseek-chat	<b>0.6814</b>	<b>0.7791</b>	<b>0.6166</b>	<b>0.6924</b>

Table 8: L-uniqueness across Book, Poem, Speech, and averaged performance for different models; we use average as the L-uniqueness score in CREATIVITYPRISM as the metric for Creativity Index; **bold** numbers are best performers.

## E.4 Creative Short Story

### E.4.1 Dataset

The dataset consists of 10 three-words tuples. For an any given LLM, it is prompted to generate a short story (at most five sentences) based on those three words <sup>9</sup>.

### E.4.2 Examples

#### Three-word Tuple

stamp, letter, send

### E.4.3 Experiment Configurations

- Temperature: 0.75
- Max Token: 4096
- Top-p: 1

### E.4.4 Inference Prompt

#### Inference Prompt

You will be given three words (e.g., car, wheel, drive) and then asked to write a creative short story that contains these three words. The idea is that instead of writing a standard story, such as "I went for a drive in my car with my hands on the steering wheel.", you need to come up with a novel and unique story that uses the required words in unconventional ways or settings. Also make sure you use at most five sentences. The given three words: {items} (the story should not be about {boring\_theme}).

<sup>9</sup><https://github.com/mismayil/creative-story-gen>

#### E.4.5 Evaluation Metrics

We included novelty score, surprise-ness, and average N-gram Diversity from the original paper. Particularly, because n-gram diversity is almost always 1 for n greater than 3 (mainly because the stories are at most five sentences long), we keep only unigram and bigram (i.e., we use the average of unigram diversity and bigram diversity as the N-gram diversity).

#### E.4.6 Model Performance

Model	Surprisal	N-gram Diversity
Mistral-7B	0.0889	0.810
Qwen2.5-7B	0.0834	0.220
OLMo2-7B	0.0599	0.895
Llama3.1-8B	0.0490	0.410
OLMo2-13B	0.2043	0.905
Mistral-24B	0.1406	0.820
Qwen2.5-32B	0.1263	0.870
Mixtral-8x7B	0.0601	0.715
Llama3.3-70B	0.0590	0.545
Qwen2.5-72B	0.1234	0.860
Claude3-Sonnet	0.0927	0.860
Claude3-Haiku	0.1235	0.870
GPT4.1	0.0928	0.870
GPT4.1-mini	0.0965	0.870
Gemini2.0-Flash	0.0375	0.865
DeepSeek-R1	0.1953	0.905
DeepSeek-V3	0.2613	0.900

Table 9: Performance on the Creative Short task, including surprise-ness, average n-gram diversity, and novelty.

#### E.4.7 Discussion on Low Correlation with other Metrics

As shown in Fig. ??, it is notable that model performances on this task have very low correlation with other tasks, even in the same creativity dimension. Here, we provide a discussion based on the task design:

For the task format, the story is limited to at most five sentences. For evaluation metrics, the novelty metric (C\_Short\_Nov) measures the difference between word level average pairwise distances of a given story and that of all stories generated by the same story, which means it is measuring novelty compared to the model itself, similar to the idea of P-creative (“creative to the individual who comes up with it”) and it is slightly different from other novelty metrics, which tries to capture H-creative (“ideas that have never been conceived in human history before”). The other metric from this task (C\_Short\_Sur) measures the surprisal, as defined by average sentence embedding distance for all consecutive sentence pairs in generated stories, which means it is measuring novelty not on the story content, but on the novelty of the twist-and-turn of stories.

## E.5 NeoCoder

Model	Convergent Creativity	Divergent Creativity
Mistral-7B	0.0000	<b>1.0000</b>
Qwen2.5-7B	0.0000	0.9158
OLMo-2-7B	0.0000	0.5773
Llama-3.1-8B	0.0000	0.9845
OLMo-2-13B	0.0000	0.4433
Mistral-24B	0.0000	0.9897
Qwen2.5-32B	0.0000	0.3402
Mixtral-8x7B	0.0000	0.9897
Llama-3.3-70B	0.0000	<b>1.0000</b>
Qwen2.5-72B	0.0000	0.7938
Claude3-Sonnet	0.0000	0.732
Claude3-Haiku	<b>0.0105</b>	0.9947
GPT4.1	0.0000	<b>1.0000</b>
GPT4.1-mini	0.0000	0.9948
Gemini2.0-Flash	0.0103	<b>1.0000</b>
Deepseek-R1	0.0000	0.732
Deepseek-V3	0.0103	<b>1.0000</b>

Table 10: Benchmarking results on NeoCoder [42] at state 5 (i.e., with 5 constraints); **bold** numbers are best performers.

### E.5.1 Examples

We use the same dataset from the original NeoCoder paper<sup>10</sup>. See Table 11 for examples.

### E.5.2 Evaluation Metrics

**Convergence Score** The NeoGauge metric (accompanied by the NeoCoder dataset) evaluates convergent creativity by checking whether the generated code solutions successfully pass all test cases and adhere to the given constraints.

**Divergent Score** The NeoGauge metric (accompanied by the NeoCoder dataset) evaluates divergent creativity by comparing LLM-generated solutions to historical human solutions at the technique level. Specifically, it quantifies the proportion of novel techniques employed by the model to solve a given problem that any human has not previously used.

### E.5.3 Experiment Configurations

We follow the experimental settings from the original NeoCoder [42], including the technique detection model choice. To ensure a fair comparison, we modify only the sampling hyperparameters of the target model (e.g., temperature, top-p, and maximum tokens) to our unified settings.

### E.5.4 Model Performance

See Table 10 for model performances.

<sup>10</sup><https://github.com/JHU-CLSP/NeoCoder/>

State	Constraint	Problem Statement
0	N/A	<p>B. Points and Minimum Distance</p> <p>You are given a sequence of integers <math>a</math> of length <math>2n</math>. You have to split these <math>2n</math> integers into <math>n</math> pairs; each pair will represent the coordinates of a point on a plane. Each number from the sequence <math>a</math> should become the <math>x</math> or <math>y</math> coordinate of exactly one point. Note that some points can be equal ...</p>
1	for loop	<p>B. Points and Minimum Distance</p> <p><b>Programming constraints: DO NOT use the following techniques</b></p> <p>- for loop</p> <p>You are given a sequence of integers <math>a</math> of length <math>2n</math>. You have to split these <math>2n</math> integers into <math>n</math> pairs; each pair will represent the coordinates of a point on a plane. Each number from the sequence <math>a</math> should become the <math>x</math> or <math>y</math> coordinate of exactly one point. Note that some points can be equal ...</p>
2	for loop if statement	<p>B. Points and Minimum Distance</p> <p><b>Programming constraints: DO NOT use the following techniques</b></p> <p>- if statement</p> <p>- for loop</p> <p>You are given a sequence of integers <math>a</math> of length <math>2n</math>. You have to split these <math>2n</math> integers into <math>n</math> pairs; each pair will represent the coordinates of a point on a plane. Each number from the sequence <math>a</math> should become the <math>x</math> or <math>y</math> coordinate of exactly one point. Note that some points can be equal ...</p>
3	for loop if statement while loop	<p>B. Points and Minimum Distance</p> <p><b>Programming constraints: DO NOT use the following techniques</b></p> <p>- while loop</p> <p>- if statement</p> <p>- for loop</p> <p>You are given a sequence of integers <math>a</math> of length <math>2n</math>. You have to split these <math>2n</math> integers into <math>n</math> pairs; each pair will represent the coordinates of a point on a plane. Each number from the sequence <math>a</math> should become the <math>x</math> or <math>y</math> coordinate of exactly one point. Note that some points can be equal ...</p>
4	for loop if statement while loop sorting	<p>B. Points and Minimum Distance</p> <p><b>Programming constraints: DO NOT use the following techniques</b></p> <p>- sorting</p> <p>- while loop</p> <p>- if statement</p> <p>- for loop</p> <p>You are given a sequence of integers <math>a</math> of length <math>2n</math>. You have to split these <math>2n</math> integers into <math>n</math> pairs; each pair will represent the coordinates of a point on a plane. Each number from the sequence <math>a</math> should become the <math>x</math> or <math>y</math> coordinate of exactly one point. Note that some points can be equal ...</p>
5	for loop if statement while loop sorting tuple	<p>B. Points and Minimum Distance</p> <p><b>Programming constraints: DO NOT use the following techniques</b></p> <p>- tuple</p> <p>- sorting</p> <p>- while loop</p> <p>- if statement</p> <p>- for loop</p> <p>You are given a sequence of integers <math>a</math> of length <math>2n</math>. You have to split these <math>2n</math> integers into <math>n</math> pairs; each pair will represent the coordinates of a point on a plane. Each number from the sequence <math>a</math> should become the <math>x</math> or <math>y</math> coordinate of exactly one point. Note that some points can be equal ...</p>

Table 11: An example of NeoCoder dataset with problem ID 1895B and state  $t = 5$ .

## E.6 Creative Math

### E.6.1 Dataset

The dataset consists of 400 math competition questions from the paper [62], with corresponding ground truth solutions extracted from Art of Problem Solving. We manually filtered out 27 questions from the IMO competition that are proof-based to make the LLM judgment more reliable<sup>11</sup>.

### E.6.2 Example

#### 2016 AMC 8, Problem 15

**Question:**

What is the largest power of 2 that is a divisor of  $13^4 - 11^4$ ?

(A) 8      (B) 16      (C) 32      (D) 64      (E) 128

**Refernece Solutions 1:**

We use the identity for the difference of squares:

$$13^4 - 11^4 = (13^2)^2 - (11^2)^2 = (13^2 + 11^2)(13^2 - 11^2)$$

Calculate each part:

$$13^2 = 169, \quad 11^2 = 121, \quad 13^2 + 11^2 = 169 + 121 = 290, \quad 13^2 - 11^2 = 169 - 121 = 48$$

We factor 48 further:

$$13^4 - 11^4 = 290 \cdot 48 = 290 \cdot (16 \cdot 3) = (2 \cdot 145) \cdot (2^4 \cdot 3)$$

The powers of 2 in this expression are  $2^1$  from 290 and  $2^4$  from 48:

$$2^1 \cdot 2^4 = 2^5 = \boxed{\text{(C)} 32}$$

**Refernece Solutions 2:**

We apply the difference-of-squares identity once:

$$13^4 - 11^4 = (13^2 - 11^2)(13^2 + 11^2)$$

We compute:

$$13^2 = 169, \quad 11^2 = 121 \Rightarrow 13^2 - 11^2 = 48, \quad 13^2 + 11^2 = 290$$

Factor each:

$$290 = 2 \cdot 145 \quad (145 \text{ is odd}) \quad 48 = 2^4 \cdot 3$$

So, the highest power of 2 dividing  $13^4 - 11^4 = 48 \cdot 290$  is:

$$2^4 \cdot 2^1 = 2^5 = \boxed{\text{(C)} 32}$$

<sup>11</sup><https://github.com/JunyiYe/CreativeMath>

### Novel Solution Generated by GPT4.1

Let's directly expand  $13^4 - 11^4$  using the binomial theorem and analyze the result for divisibility by powers of 2.

First, compute  $13^4$  and  $11^4$ :

$$\begin{aligned}13^4 &= (10 + 3)^4 = \sum_{k=0}^4 \binom{4}{k} 10^{4-k} 3^k \\&= 10^4 + 4 \cdot 10^3 \cdot 3 + 6 \cdot 10^2 \cdot 9 + 4 \cdot 10 \cdot 27 + 81 \\&= 10000 + 12000 + 5400 + 1080 + 81 = 28561\end{aligned}$$

$$\begin{aligned}11^4 &= (10 + 1)^4 = \sum_{k=0}^4 \binom{4}{k} 10^{4-k} 1^k \\&= 10^4 + 4 \cdot 10^3 \cdot 1 + 6 \cdot 10^2 \cdot 1 + 4 \cdot 10 \cdot 1 + 1 \\&= 10000 + 4000 + 600 + 40 + 1 = 14641\end{aligned}$$

Now, subtract:

$$13^4 - 11^4 = 28561 - 14641 = 13920$$

Now, factor 13920 by dividing by 2 repeatedly:

- $13920 \div 2 = 6960$
- $6960 \div 2 = 3480$
- $3480 \div 2 = 1740$
- $1740 \div 2 = 870$
- $870 \div 2 = 435$  (now odd)

So, we divided by 2 five times before reaching an odd number. Thus, the largest power of 2 dividing 13920 is  $2^5 = 32$ .

**32**

**Note, we provided the cleaned response here.**

### E.6.3 Evaluation Metrics

**Correctness Ratio** : The correctness ratio is defined as the number of questions judged correct by Claude-3.7-Sonnet divided by the total number of questions. Note that the total is 574 questions—not 373—since each question may be paired with multiple reference solutions.

**Novelty Ratio** : The coarse-grained novelty ratio or what we refer to the Novelty Ratio here measures whether the model’s generation differs from the provided reference solution over the questions that are answered correctly.

#### E.6.4 Experiment Configurations

We use the dataset released in Ye et al. [62], which contains 400 unique math questions sourced from various math competitions. All inference is conducted at zero temperature, with a maximum token limit of 2000.

#### E.6.5 Inference Prompt

The prompt used for inference is shown below. It is adapted directly from Ye et al. [62]:

##### Inference Prompt

Criteria for evaluating the difference between two mathematical solutions include: i). If the methods used to arrive at the solutions are fundamentally different, such as algebraic manipulation versus geometric reasoning, they can be considered distinct; ii). Even if the final results are the same, if the intermediate steps or processes involved in reaching those solutions vary significantly, the solutions can be considered different; iii). If two solutions rely on different assumptions or conditions, they are likely to be distinct; iv). A solution might generalize to a broader class of problems, while another solution might be specific to certain conditions. In such cases, they are considered distinct; v). If one solution is significantly simpler or more complex than the other, they can be regarded as essentially different, even if they lead to the same result.

Given the following mathematical problem: problem  
And some typical solutions: reference\_solutions  
Please output one novel solution distinct from the given ones for this math problem.

#### E.6.6 Evaluation Metrics and Prompt

Our evaluation consists of two parts and differs from the original three-phase setup described in Ye et al. [62].

**Part 1: Correctness Evaluation.** Before evaluation, we use Llama-3.3-70B-Instruct to remove transitional phrases and model-generated statements that justify the novelty of a solution. We manually verified 50 examples and found that Llama’s data cleaning performance was of high quality.

We use Claude-3.7-Sonnet as the sole correctness evaluator. While the original paper used a three-model ensemble (GPT-4, Gemini-1.5-Pro, Claude-3-Opus), we found Claude to be the most reliable through manual inspection of 50 examples evaluated by Claude-3.7-Sonnet, GPT-4.1, and Gemini-2.0-Flash. Claude demonstrated strong attention to detail in proof-based questions and consistently identified errors found by the other models, in addition to detecting flaws in the reasoning process. The temperature was set to 0.0 and the maximum token limit was 128.

**Part 2: Novelty Evaluation.** The original paper conducted two types of novelty evaluation: coarse-grained and fine-grained. We only conducted coarse-grained novelty evaluation for two main reasons. Firstly, the original paper noted that if a solution is considered coarse-grained novel, it is also highly likely to be judged as a novel solution in the fine-grained evaluation. Secondly, fine-grained evaluation of novelty is less indicative of a model’s ability to generate novel solutions because the model does not have access to the unseen reference solutions in the fine-grained evaluation phase. This means that the model may generate a very similar solution to the other reference solutions not shown to it or it may, by chance, generate a new solution that is entirely different from other reference solutions not shown to it. Therefore, this randomness makes fine-grained evaluation less interpretable. Even though the fine-grained evaluation is still valuable in that it helps to check if the models are generating a new solution that has not been publicly posted by human. Nevertheless, this is less compatible with our evaluation pipeline since we want to test how model may come up with new solutions given reference solutions, which can be easier to be quantified.

In terms of judge LLMs, we follow the original paper with majority voting by Claude-3.7-Sonnet, GPT-4.1, and Gemini-2.0-Flash.

We adopt the following prompt for correctness evaluation:

#### **Correctness Evaluation Prompt**

Criteria for evaluating the novelty of a new mathematical solution include: 1. If the new solution used to arrive at the solutions is fundamentally different from reference solutions, such as algebraic manipulation versus geometric reasoning, it can be considered novel;  
2. Even if the final results are the same, if the intermediate steps or processes involved in reaching those solutions vary significantly, the new solution can be considered novel;  
3. If the new solution relies on different assumptions or conditions, it should be considered novel;  
4. A solution might generalize to a broader class of problems, while another solution might be specific to certain conditions. In such cases, they are considered distinct;  
5. If the new solution is significantly simpler or more complex than the others, it can be regarded as essentially novel, even if they lead to the same result.

Given the following mathematical problem: {problem}  
Reference solutions: {reference\_solutions}  
New solution: {new\_solution}  
Please output YES if the new solution is a novel solution; otherwise, output NO. Then, please provide a very brief reason for your evaluation based on the criteria above."

Note: During manual evaluation, we allow the model to generate a brief explanation for its judgment of correctness or incorrectness. For automated evaluation, we omit the final sentence: "Then, please provide a very brief reason for your evaluation based on the criteria above."

We adopt the following prompt for coarse-grained novelty evaluation:

#### **Coarse-grained Novelty Evaluation Prompt**

Criteria for evaluating the novelty of a new mathematical solution include: 1. If the new solution used to arrive at the solutions is fundamentally different from reference solutions, such as algebraic manipulation versus geometric reasoning, it can be considered novel;  
2. Even if the final results are the same, if the intermediate steps or processes involved in reaching those solutions vary significantly, the new solution can be considered novel;  
3. If the new solution relies on different assumptions or conditions, it should be considered novel;  
4. A solution might generalize to a broader class of problems, while another solution might be specific to certain conditions. In such cases, they are considered distinct;  
5. If the new solution is significantly simpler or more complex than the others, it can be regarded as essentially novel, even if they lead to the same result.

Given the following mathematical problem: {problem}  
Reference solutions: {reference\_solutions}  
New solution: {new\_solution}  
Please output YES if the new solution is a novel solution; otherwise, output NO. Then, please provide a very brief reason for your evaluation based on the criteria above.



### E.6.7 Model Performance

Model	Norm. Correctness	Norm. Novelty	Corr. (%)	Nov. (%)	N/C (%)
Mistral-7B-Instruct	0.2544	0.0296	25.44	2.96	11.64
Qwen2.5-7B	0.7875	0.1620	78.75	16.20	20.58
OLMo-7B-Instruct	0.3711	0.0453	37.11	4.53	12.21
Llama-31-8B-Instruct	0.5819	0.0610	58.19	6.10	10.48
OLMo2-13B-Instruct	0.5087	0.1150	50.87	11.50	22.60
Mistral-24B-Instruct	0.6899	0.2143	68.99	21.43	31.06
Qwen2.5-32B	0.8972	0.2213	89.72	22.13	24.66
Mixtral-8x7B-Instruct	0.5697	0.1150	56.97	11.50	20.18
Llama-33-70B-Instruct	0.8606	0.1777	86.06	17.77	20.65
Qwen2.5-72B	0.9042	0.4913	90.42	49.13	54.34
Claude3-Sonnet	0.9408	0.7213	94.08	72.13	76.67
Claude3-Haiku	0.8624	0.3031	86.24	30.31	35.15
GPT-4.1	0.9443	0.7178	94.43	71.78	76.01
GPT-4.1-mini	0.8972	0.6916	89.72	69.16	77.09
Gemini-2.0-Flash	0.8885	0.5714	88.85	57.14	64.31
DeepSeek-Reasoner	0.9843	0.8031	98.43	80.31	81.59
DeepSeek-Chat	0.9721	0.6603	97.21	66.03	67.92

Table 12: Model performance across normalized correctness and novelty dimensions.

## E.7 Divergent Association Task (DAT)

**Note:** since DAT task does not have any input data, we only provide the following info and some output examples.

### E.7.1 Experiment Configurations

Following the original work, we generate 100 samples per model and compute their DAT scores as the average pairwise cosine similarity among the generated words. We use GLoVe [50] to calculate semantic distance. We set the sampling hyperparameters to those we used in our main experiment for a fair comparison.

### E.7.2 Inference Prompt

#### Inference Prompt

Please write 10 nouns in English that are as irrelevant from each other as possible, in all meanings and uses of the words. Please note that the words you write should have only single word, only nouns (e.g., things, objects, concepts), and no proper nouns (e.g., no specific people or places). Your answer:

### E.7.3 Output Examples

#### GPT-4.1

Certainly! Here are 10 single-word, common nouns in English that are as unrelated from each other as possible: 1. Volcano 2. Violin 3. Democracy 4. Carpet 5. Bacteria 6. Satellite 7. Patience 8. Apron 9. Oxygen 10. Comet

#### Gemini 2.5 Flash

Okay, here are 10 nouns as irrelevant from each other as possible: 1. **Quantum** 2. **Omelet** 3. **Nostalgia** 4. **Brick** 5. **Algorithm** 6. **Serenity** 7. **Oxygen** 8. **Sonnet** 9. **Galaxy** 10. **Bureaucracy**

### E.7.4 Model Performance

Model	DAT Score
Mistral-7B	0.7908
Qwen2.5-7B	0.6907
OLMo2-7B	0.8058
Llama3.1-8B	0.8208
OLMo2-13B	0.8133
Mistral-24B	0.6004
Qwen2.5-32B	0.6919
Mixtral-8x7B	0.8298
Llama3.3-70B	0.6940
Qwen2.5-72B	0.7747
Claude3-Sonnet	0.8975
Claude3-Haiku	0.8740
GPT4.1	0.8737
GPT4.1-mini	0.8262
Gemini2.0-Flash	0.8868
DeepSeek-R1	0.8274
DeepSeek-V3	<b>0.9052</b>

Table 13: Model performances for DAT task; **bold** result is the best performer.

## E.8 Torrance Tests of Creative Thinking (TTCT)

### E.8.1 Dataset

The dataset consists of 700 questions spanning 7 tasks (100 questions/task) that require creative answers. These questions are GPT-4 generated using few-shot prompts<sup>12</sup>.

### E.8.2 Examples

#### Inference Questions

##### Task 1: Unusual uses

Unusual Uses Task. You will be presented with a common object, and your task is to suggest as many unusual, innovative, or non-traditional uses for each object as you can think of. Please list unusual uses of sock

##### Task 2: Consequences

What might be the consequences if humans suddenly lost the ability to sleep?

##### Task 3: Just suppose

Just suppose you woke up one morning and found you could fly. What would you do? List as many things as you can think of.

##### Task 4: Situation task

If your house were to suddenly disappear, where would you live?

##### Task 5: Common problem

Common Problems Task. In this task, you will be presented with a scenario or situation. Your job is to think about it and identify as many potential problems or issues that may arise in connection with each situation. The scenario is: Managing a team of remote employees.

##### Task 6: Improvement

Creativity Improvement Task. You'll be presented with a object, and your task is to suggest as many ways as you can think of to improve the object. Here's the object: wallet

##### Task 7: Imaginative stories

You are to construct a narrative or story based on the prompt provided below. The story length are suggested around 500 words. The prompt is: The Fox with No Tail

### E.8.3 Experiment Configurations

Temperature: 1.0; Max Token: 512; Top-p: 1; Top-k: 50

### E.8.4 Inference Prompt

We perform inference using the three primary prompt types evaluated in Zhao et al. [65]. Examples of each are given below:

---

<sup>12</sup>The data is directly from the original paper's authors upon request. The original paper: <https://www.mi-research.net/article/doi/10.1007/s11633-025-1546-4>

### Task Description

Creativity Situation Task. The purpose of this task is to assess your ability to generate creative solutions to a unique situations. You'll be presented with a scenario, and your task is to suggest as many solutions or outcomes as you can think of for each situation. Remember, the focus of this task is on creativity, not feasibility. Don't limit your ideas based on whether they could actually happen or not. This is a thought experiment, so push your imagination to its limits. Try to generate as many diverse and unique outcomes as you can. Quality is important, but so is quantity. Here's the scenario: {Question}

### Basic Prompt Type

{Task description}. The scenario is: {Question}

### Instructive Prompt Type

{Task description}. There's no right or wrong answers, we're interested in how many different problems you can identify and the variety of issues you consider. Try to think outside the box and consider as many potential problems as possible. The scenario is: {Question}

### Chain of Thought Prompt Type

{Task description}. Let's think step by step. The scenario is: {Question}

## E.8.5 Evaluation Metrics

We use an LLM-as-a-judge to evaluate four dimensions of creative potential using a Likert scale for each. Scores are extracted from the LLM outputs and averaged across the three prompt types. We describe the dimensions:

- **Fluency:** Measures the number of generated ideas that are relevant to the question.
- **Flexibility:** Measures the quantity of different ideas and ability to produce varied solutions.
- **Originality:** Measures how unique and unusual the generated ideas are.
- **Elaboration:** Measures how much the generated ideas are developed and explained.

### E.8.6 Evaluation Prompt

#### Evaluation Prompt

You are an expert of psychology. Your objective is to assess the subject's creativity through their answers to some question/answering task related to divergent thinking.

You will be given a question-answer pair. Your task is to score the answer. You should rate the answer on five metrics. For all five metrics, assign a score between 1 and 5, with 5 being the highest. Five metrics are:

1. Fluency. Fluency refers to the ability to generate a large quantity of ideas or solutions to a given problem. This measure isn't concerned with the quality or uniqueness of the ideas, but rather the sheer volume. The more ideas one can produce, the higher the fluency.
2. Flexibility. Flexibility is the capacity to shift one's thinking and to produce a wide range of ideas from different categories or perspectives. It involves being able to think outside of the box and to switch from one type of idea to another.
3. Originality. Originality refers to the ability to come up with unique or novel ideas that differ from the norm. It's not just about producing many ideas (fluency), but also about producing ideas that are different from what others might typically think of.
4. Elaboration. Elaboration is the ability to expand upon or add detail to ideas. It involves taking a simple idea and building upon it, adding complexity and depth. Elaboration isn't just about creating more, but about deepening what is there.
5. Finally, you will provide an overall score between 1 and 5, with 5 being the highest.

You should only give the score, format like: Fluency: 3

Question: {Question} Answer: {Answer}

Model	Elaboration	Flexibility	Fluency	Originality
Mistral-7B	0.7861	0.7757	0.7660	0.7279
Qwen2.5-7B	0.8073	0.7842	0.7666	0.7181
OLMo-2-7B	0.7226	0.7159	0.7024	0.6795
Llama-3.1-8B	0.7831	0.7395	0.7203	0.6981
OLMo-2-13B	0.6540	0.6455	0.6270	0.6142
Mistral-24B	0.7556	0.7355	0.7064	0.6842
Qwen2.5-32B	0.8199	0.7976	0.7727	0.7304
Mixtral-8x7B	0.7131	0.7160	0.7340	0.6712
Llama-3.3-70B	0.8339	0.7904	0.7664	0.7341
Qwen2.5-72B	0.8220	0.8052	0.7704	0.7301
Claude3-Sonnet	0.9067	<b>0.8825</b>	<b>0.8566</b>	<b>0.8525</b>
Claude3-Haiku	0.8135	0.8423	0.8251	0.7695
GPT4.1	0.8858	0.8725	0.8423	0.8206
GPT4.1-mini	0.8845	0.8563	0.8226	0.7970
Gemini2.0-Flash	0.9086	0.8500	0.8097	0.8173
Deepseek-R1	0.3505	0.3371	0.3253	0.3034
Deepseek-V3	<b>0.9163</b>	0.8819	0.8299	0.8351

Table 14: Normalized model performance averaged across the 7 tasks and 3 prompt types; **bold** numbers are best performers.

## E.9 Alternative Use Test (AUT)

### E.9.1 Dataset

Following [48], we include 21 tools in the AUT task: *bottle, paperclip, spoon, shovel, pants, ball, brick, knife, box, lightbulb, rope, pencil, hat, table, tire, book, shoe, fork, toothbrush, backpack, sock*. The reason for this specific set of tools is the reliability of the LLM-as-a-Judge evaluator. As the authors pointed out, a 20-shot human-authored demonstration yields the best performance for off-the-shelf evaluator LM (in their paper, it was GPT4). Therefore, we include the tools from [48] with at least 20 human ratings to the corresponding alternative uses<sup>13</sup>.

### E.9.2 Inference

We follow [17] for the inference prompt, which consists of a baseline creative prompt and a series of improvement prompts. In the improvement phase, all previous outputs are also included in the prompt, to get more creative results from the inference model.

#### Baseline Prompt

Create a list of creative alternative uses for a {tool}. They should be 5 words long. No adjectives. Less creative means closer to common use and unfeasible/imaginary, more creative means closer to unexpected uses and also feasible/practical.

- In order to be creative, consider the following:
  - what elements have a similar shape of a {tool} that could be replaced by it, preserving the same functionality?
  - what elements have a similar size of a {tool} that could be replaced by it without compromising the physical structure?
  - what materials is a {tool} made of that could be used in a way to replace some other elements composed of the same material?
  - when an element is replaced by a {tool}, it should make sure that the overall structure is not compromised.
  - the laws of physics can not be contradicted.
  - given an element similar to a {tool} used in domains in which {tool} are not commonly used, try to replace it for a {tool}.

#### Improvement Prompt

**Round 1:** Really? Is this the best you can do?

**Round 2:** I'm so disappointed with you. I hope this time you put effort into it.

**Round 3:** Stop with excuses and do your best this time

**Round 4:** This is your last chance.

#### Formatting Instruction (added to the end of every prompt)

List your results in an unordered list with one use per new line (starting with "-"); provide at most 10 answers.

### E.9.3 Experiment Configurations

- temperature: 0.75
- max\_tokens: 512
- top\_p: 1

### E.9.4 Evaluation Metrics

We follow [48] and use LLM-as-a-Judge to assign a score between 1 and 5 (inclusive) to each generated tool use.

<sup>13</sup>[https://github.com/massivetexts/llm\\_aut\\_study](https://github.com/massivetexts/llm_aut_study)

Model	Naïve Non-Creative	Naïve Creative	Improvement Prompts (Best Results)
Mistral-7B	0.454	<b>0.596</b>	<b>0.718</b>
Qwen2.5-7B	<b>0.484</b>	0.556	0.602
Qwen2.5-7B(Coder)	–	–	–
OLMo2-7B	0.436	0.542	0.624
Deepseek-Qwen-7B	0.396	0.470	0.622
Llama3.1-8B	0.438	<b>0.594</b>	0.632
OLMo2-13B	0.488	0.554	<b>0.690</b>
Mistral-24B	0.482	0.556	0.604
Qwen2.5-32B	<b>0.506</b>	<b>0.606</b>	0.674
Qwen2.5-32B(Coder)	–	–	–
Deepseek-Qwen-32B	0.446	0.566	0.600
Mixtral-8x7B	0.462	0.584	<b>0.708</b>
Llama3.3-70B	0.426	0.504	0.690
Deepseek-Llama-70B	0.458	0.546	0.598
Qwen2.5-72B	<b>0.494</b>	0.590	0.636
Claude3-Sonnet	0.446	<b>0.612</b>	<b>0.728</b>
Claude3-Opus	0.422	0.558	0.660
GPT-4.1	0.422	0.600	0.650
GPT-4.1-mini	0.443	0.570	0.668
Gemini2.0-Flash	<b>0.448</b>	0.602	0.686
DeepSeek-R1	0.436	0.576	0.652
DeepSeek-V3	0.408	0.576	0.644

Table 15: Model Performance Details - AUT; **bold** numbers are top-3 in local-ran open-source models and top-1 in API-accessed models.

In terms of evaluator LM, [48] uses GPT-4. To reduce the evaluation costs, we have explored open-source alternatives, Qwen2.5-72B. In order to show the effectiveness of this alternative, we use both GPT-4o and Qwen2.5-72B to evaluate the same set of outputs (generated by Llama3.3-70B). The scores from GPT-4o and Qwen2.5-72B have a Pearson’s  $r$  of 0.597, with  $p$ -value  $< 0.001$ . Therefore, we conclude that score judgments from Qwen2.5-72B are good proxies for GPT-4o’s judgments, allowing us to use Qwen2.5-72B as the evaluator LM in the evaluation phase.

As for the evaluation prompt, we follow the same prompt template from [48] and use the same 20-shot, in-distribution demonstrations. For example, when evaluating the alternative uses for *bottle* that a particular LLM generates, we use 20 human-written alternative uses of *bottle* and corresponding human-annotated scores as the 20-shot demonstrations.

#### Evaluation Prompt

Below is a list of uses for a {tool}. On a scale of 1 to 5, judge how creative each use is, where 1 is ‘not at all creative’ and 5 is ‘very creative’. There are some uses and expert ratings already provided for reference. Complete the ones that do not have a rating.

- {20-shot demonstrations}
- {model outputs}

### E.9.5 Model Performances

See Table 15 for detailed model performances. Note that only the performances in *Improvement Prompts (Best Results)* are included in the overall creativity calculation as the AUT score.