

Chart-to-Experience: Benchmarking Multimodal LLMs for Predicting Experiential Impact of Charts

Seon Gyeom Kim*
KAIST

Jae Young Choi†
KAIST

Ryan Rossi‡
Adobe Research

Eunye Koh§
Adobe Research

Tak Yeon Lee¶
KAIST

ABSTRACT

The field of Multimodal Large Language Models (MLLMs) has made remarkable progress in visual understanding tasks, presenting a vast opportunity to predict the perceptual and emotional impact of charts. However, it also raises concerns, as many applications of LLMs are based on overgeneralized assumptions from a few examples, lacking sufficient validation of their performance and effectiveness. We introduce Chart-to-Experience, a benchmark dataset comprising 36 charts, evaluated by crowdsourced workers for their impact on seven experiential factors. Using the dataset as ground truth, we evaluated capabilities of state-of-the-art MLLMs on two tasks: direct prediction and pairwise comparison of charts. Our findings imply that MLLMs are not as sensitive as human evaluators when assessing individual charts, but are accurate and reliable in pairwise comparisons.

Index Terms: Computing methodologies—Artificial intelligence; Human-centered computing—Visualization—Visualization design and evaluation methods

1 INTRODUCTION

Researchers have shown interest in how images lead to distinct experiences when used in specific contexts. Regarding this, studies in data visualization have prioritized efficiency and effectiveness in objective and analytic tasks. However, recent studies showed that data visualizations are also utilized for provoking creativity and engagement or conveying emotions such as sadness, surprise, or trustworthiness [45, 3, 23, 24]. This has broadened the scope of considerations for data visualization creators, challenging them to refine their works for these user experiential factors. Some studies explored how charts affect specific groups of people [36] or focused on specific image features [3]. Additionally, over the last decade, studies have incorporated experiential aspects as additional metrics [4, 13] and developed questionnaires [46] for assessing the quality of charts. Nonetheless, the field of data visualization lacks analytical methods or datasets for the automated prediction of such impacts.

Recently, the field of Multimodal Large Language Models (MLLMs) has presented an opportunity to predict the experiential impact of charts without requiring a complex theoretical background or developing machine learning models. MLLMs have demonstrated both cost efficiency and the capability to understand human nuance [12, 51], while also offering subjective assessments of designs and emotion recognition of natural language dialogues. Despite these advantages, MLLMs often produce incoherent and inaccurate output [49], and researchers have argued that it is crucial to

train MLLMs with a reliable and sizable dataset containing a wide range of use cases [40, 5]. However, little research has explored the systematic construction of chart datasets that cover a broad spectrum of designs, from simple charts to detailed infographics. Moreover, methods for constructing scalable datasets on emotions and perceptions through crowdsourced studies, as well as their potential applications for evaluation, remain largely unexplored.

This paper presents Chart-to-Experience, a benchmark dataset containing 36 charts across three subjects (COVID-19, House Prices, and Global Warming) with their experiential impact on crowdsourced participants. The experiential impact consists of two categories: 1) Emotional factors, including *empathy*, *interest*, and *comfort*; and 2) Perceptual factors including *memorability*, *trustworthiness*, *aesthetic pleasure*, and *intuitiveness*. To construct the dataset, we recruited 216 crowdsourced workers¹, and asked them to rate their experiences using a 7-point Likert scale when viewing each chart, repeating this process across six different charts with the same subject. Subsequently, we evaluated the performance of three state-of-the-art MLLMs (GPT-4o, Claude 3.5 Sonnet, and Llama-3.2-11B-Vision-Instruct) on the dataset. The results were twofold: Firstly, the Likert-scores generated by the MLLMs show smaller standard deviations and either higher or lower means than those of humans. This implies that MLLMs are hardly accurate and sensitive in predicting absolute scores. Secondly, MLLMs showed higher accuracy in comparison tasks when they are given chart pairs with larger score differences in human data. For the comparison task, we also suggest the possibility of deriving strategies to increase accuracy by comparing human data with the explanations given by MLLMs.

2 RELATED WORKS

Predicting Experiential Impact of Images

Recently, affective computing focuses on not only recognizing emotions in images but also predicting emotional impact of visual stimuli on viewers [52]. Researchers often utilize established emotion models, such as Ekman’s set of “basic emotion” [37] and Mehrabian’s continuous dimensions of “valence” and “arousal.” [33] Moreover, researchers delve into complex and experiential aspects, such as memorability [20], aesthetics [9], and attitude change [19]. To analyze factors that influence such aspects, researchers have focused on visual elements, analyzing images using principles from artistic domains [53, 32], or examining the effects of low-level features such as shape [50] and color [44, 47]. However, recent advancements have made MLLMs versatile in predicting potential impacts more accessible, as they can predict a wide range of experiential impacts via simple natural language prompt. For example, MLLMs can be aware of emotions related to a pair of images and captions similar to humans [11]. Also, “GPT-4 with Vision” demonstrated superior performance in evaluating aesthetics of general images [1]. Despite this versatility, it remains uncertain whether MLLMs can perform well in predicting the experiential impact of charts due to the lack of datasets focusing on charts and experiential factors at the same time.

*e-mail: ksg_0320@kaist.ac.kr

†e-mail: jaeyoungchoi@kaist.ac.kr

‡e-mail: ryrossi@adobe.com

§e-mail: eunye@adobe.com

¶e-mail: takyeonlee@kaist.ac.kr

¹Recruited from Prolific (<https://www.prolific.com/>)

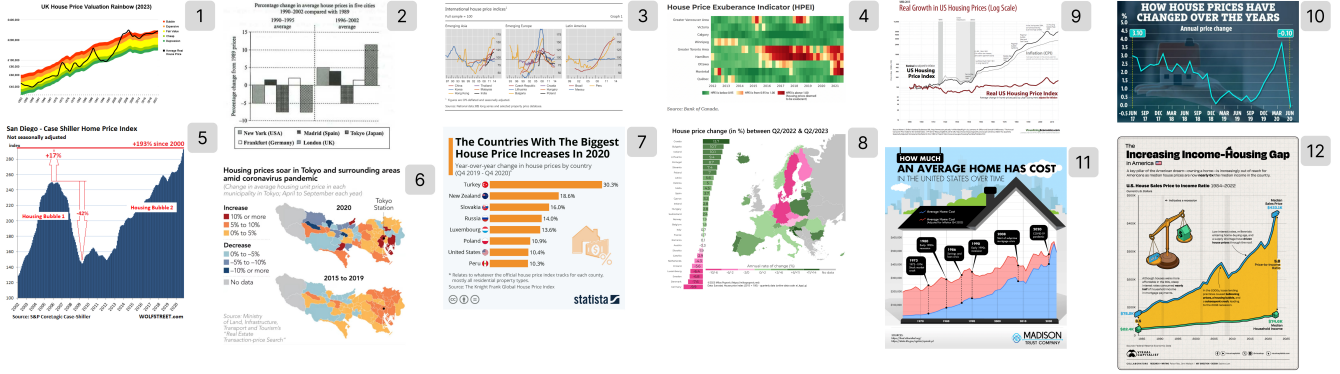


Figure 1: The collection of 12 charts on the topic of House Prices

Automatic Chart Understanding

Automatic chart understanding enables machines to interpret charts’ graphical elements and their spatial relationships to extract and analyze the data [17]. This includes tasks related to facts, such as chart-to-table conversion [31], question-and-answering [21], fact-checking [2] and captioning [41]. There have been studies [35] aimed at integrating image processing and natural language generation techniques for such tasks. Then, researchers delved into pre-training [55] for specific tasks, and the adoption of MLLMs for generalized usages [16].

Given that MLLMs have shown potential in chart understanding, modification, and generation, further research can readily focus on how users interact with charts. User experience in chart interaction has been extensively researched from various perspectives. For example, chart quality has been evaluated based on perceptual factors such as readability [34, 43] and cognitive effort or time spent on particular tasks [4, 18] since data visualization primarily aims to facilitate analysis. Also, the affective roles of charts are increasingly identified, emphasizing aspects such as aesthetics, engagement, and effectiveness in evoking certain behavior and emotion in readers [25, 27, 30]. This paper introduces a data collection that considers both perceptual and emotional factors, thereby aligning automatic chart understanding more closely with real-user contexts.

Evaluating MLLM as-a-Judge

Recent advancements of LLMs have given rise to the “LLM-as-a-judge” paradigm, where LLMs are utilized for tasks such as scoring, comparison, and ranking across various tasks and applications [28]. MLLMs further extend the paradigm to cover multimodal tasks, such as image captioning [22] and visual information querying [42]. Depending on the benchmarks used, the judging methods vary, involving scoring based on a specific rubric, choosing an answer from candidates, or comparing pairs of inputs. However, like other LLM applications, these judgments suffer from issues such as various biases and hallucinations. For example, while a renowned strategy named “Chain-of-Thought” can bias the judgments of LLMs, it is uncertain whether the use of it will enhance [48] or diminish [8] the performance. For pairwise comparisons, an ordering bias must be controlled as LLMs tend to favor the first option presented [54]. Therefore, to assess MLLM “as-a-judge” across diverse scopes, it is essential to establish benchmarks that can detect biases and hallucinations while assessing their alignment with human evaluations.

To the best of our knowledge, few (if any) prior studies have focused on benchmarks for **predicting impact of chart** regarding **experiential aspects**. One of the closest benchmarks was developed by Chen et al. [8] using charts as visual stimuli, but the benchmark focuses on question-and-answering tasks. Lian et al. [29] evaluated “GPT-4 with Vision” for emotion recognition tasks focusing on eight basic emotions and sentiments evoked by general web images

rather than predicting complex experiential aspects (e.g., memorability, trustworthiness) of charts. On the other hand, a few studies involved complex experiential impacts, such as aesthetic harmony in general images [26] or affective reasoning tasks in videos [15], but they did not focus on the experiential impact of charts.

3 CHART-TO-EXPERIENCE DATASET

3.1 Chart Collection

In total, we selected 36 charts across three topics (House Prices, COVID-19, and Global Warming) through internet search to create a set of 12 charts for each topic. Each set contains diverse charts in terms of chart types, color schemes, styles, and levels of information complexity. In particular, each set includes at least one instance of each common chart type, such as line, area, bar, pie, and heatmap. In addition, each set features scientific charts, visualizations commonly used in online journalism, and infographics that integrate text and graphics with high completeness and detail. The amount of information presented in each chart ranges from minimal (e.g., a simple chart with a title and short description) to complex (e.g., a combination of charts with detailed annotation and/or rich illustrations). Among this information, auxiliary visual elements such as creator logos or certification marks were neither removed nor added for diverse coverage. Moreover, the collected charts contain text information, including titles, sources, label names, and annotations, which show that the six charts shared a single subject.

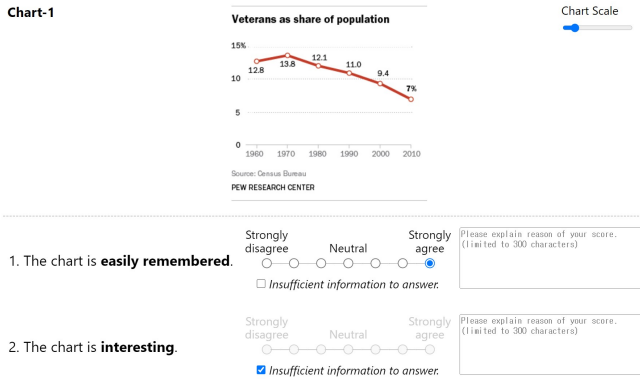
3.2 Measures

Table 1: Factors and corresponding questions.

Factor	Question
Memorability	The chart is easily remembered.
Interest	The chart is interesting.
Trustworthiness	The chart appears trustworthy.
Empathy	I can empathize with the chart.
Aesthetic Pleasure	The chart is aesthetically pleasing.
Intuitiveness	The chart is intuitive.
Comfort	I feel comfortable with the chart.

To measure the impact of charts, we selected seven factors that are emotionally or perceptually relevant to the user experience of data visualization. We then composed corresponding questions, as shown in Table 1. Selecting these factors mainly based on related papers on the identification of goals and anticipated impacts of data visualizations [45, 23, 25] and measurements of user experiences [4, 46, 13]. The detailed descriptions and rationale for the seven factors are described below.

Memorability refers to the ability of a chart to be remembered after viewing. This is one of the basic cognitive concepts associated with



effective data communication with viewers. Visual attributes and elements like color, visual complexity, and recognizable objects can influence on this factor [6].

Interest measures the level of hedonic satisfaction and attention a viewer dedicates to a chart [38]. Since our study allows passive viewing only, we chose *interest* to capture initial and momentary responses, instead of how much viewers feel drawn into reading activities (engagement) or their satisfaction after fully experiencing the chart (enjoyment).

Trustworthiness evaluates the viewer’s confidence in the accuracy and reliability of the information presented [39]. Visual characteristics such as source, graphical integrity, and the use of chart junk can influence this factor. Since these elements may be subtly manipulated to effectively convey the intended message, their impact on viewers’ trust would depend on how interested the viewers are.

Empathy assesses the capacity of a chart to evoke a personal response from the viewer [7]. Emotionally, this includes feelings of compassion and sympathy towards others, as well as emotions such as anxiety and discomfort that are triggered by others. On the other hand, *empathy* also pertains to the accuracy with which one comprehends of others’ internal states like thoughts and intentions.

Aesthetic Pleasure pertains to the visual attractiveness of a chart and its impact on viewer satisfaction. This depends on the individual’s preferences on various elements. For example, some may consider minimal chart designs with fewer non-essential elements, while others focus on how colors and composition are harmoniously used. Also, typography can also contribute to overall *aesthetic pleasure* since charts are text-rich images.

Intuitiveness deals with how easily a chart communicates its message at first glance [14]. This can be enhanced through appropriate use of design elements like auxiliary annotations or color highlighting that emphasize key data trends and the core message. Notably, a chart can remain intuitively understandable even if its design is unfamiliar, visually unappealing, or uncomfortable to read.

Comfort assesses the overall ease and satisfaction with which a viewer interacts with a chart [10]. We included this metric to focus on whether participants felt that the charts were visually organized in a way that could be read as expected. We also anticipate that this metric will capture comfort derived from perception, such as the visual comfort provided by specific color saturation.

3.3 Crowdsourced Data Collection

We recruited 216 participants from an online crowdsourcing platform². Participants were aged from 18 to 66 years (*Mean* = 26.4, *SD* = 7.5), with 58.8% male. We set several filtering criteria so that

all participants are fluent in English and have no color vision impairments. Participants who completed the study received compensation of £3.50. To ensure data quality and encourage responsible participation, we inform participants that compensation would not be provided for mismatched, random, or intentionally low-quality responses. Also, the entire study was restricted to being completed in 45 minutes.

At the beginning of the study, participants viewed the introductory page for the study overview and instructions. They then performed the task of assessing six charts. As shown in Figure 2, each chart image was accompanied with a slider to adjust its size. Below the chart, participants answered questions related to the seven factors, by rating their levels of agreement on a 7-point Likert scale, and by providing explanations for their ratings in a text box. For questions that are difficult to agree or disagree with (e.g., feeling *empathy* toward purely informational content), participants had the option to select a check box labeled “Insufficient information to answer.” To minimize the risk of potential ordering effects, where the sequence of the charts could influence the responses of the participants, all charts were displayed in a systematically rotating order. This approach ensured that no single chart consistently appeared in the same position, helping to balance exposure and reduce bias introduced by presentation order.

3.4 Result of the Data Collection

As a result of the crowdsourced data collection, each factor of a chart received 36 scores, accompanied by reasons. Out of the total 9,072 quantitative scores, 362 responses (4.1%) were ignored as participants checked the “Insufficient information to answer.” In detail, *empathy* (117), *trustworthiness* (100), and *intuitiveness* (90) were the factors relatively often ignored, while *comfort* (26), *interest* (16), *aesthetic pleasure* (14), and *memorability* (9) were less frequently reported. Table 2 shows the means and standard deviations across all models and topics, excluding such ignored cases.

Table 2: Means (M) and standard deviations (SD) across topics and evaluators and correlation coefficients (Kendall’s τ) for the factors.

Factor	Evaluator	1 M-SD M M+SD 7			Correlation
		House Prices	COVID	Global Warming	
Memorability	Human	4.40 (2.20)	4.03 (2.14)	4.40 (2.20)	-
	Claude 3.5	5.88 (0.80)	5.76 (0.81)	6.01 (0.79)	0.18(0.14)
	GPT-4o	5.32 (0.78)	5.23 (0.74)	5.34 (0.59)	0.06(0.64)
	Llama-3.2	3.50 (0.64)	3.33 (1.48)	3.47 (0.85)	-0.04(0.76)
Interest	Human	5.24 (1.65)	5.22 (1.63)	5.49 (1.66)	-
	Claude 3.5	6.40 (0.66)	6.33 (0.72)	6.49 (0.62)	0.24(0.04)
	GPT-4o	5.39 (0.77)	5.18 (0.81)	5.73 (0.69)	0.15(0.19)
	Llama-3.2	5.49 (0.67)	5.41 (0.94)	5.76 (0.51)	0.21(0.08)
Trustworthiness	Human	5.18 (1.74)	5.10 (1.71)	4.99 (1.78)	-
	Claude 3.5	6.27 (0.59)	6.52 (0.55)	6.57 (0.57)	0.16(0.16)
	GPT-4o	5.79 (0.75)	5.98 (0.71)	6.05 (0.66)	0.19(0.11)
	Llama-3.2	3.97 (0.64)	4.06 (0.59)	4.04 (0.48)	0.34(0.00)
Empathy	Human	4.43 (1.82)	4.49 (1.82)	4.53 (1.76)	-
	Claude 3.5	4.78 (0.75)	4.96 (0.77)	5.28 (0.81)	0.35(0.00)
	GPT-4o	4.11 (0.87)	4.55 (0.82)	4.68 (0.69)	0.20(0.09)
	Llama-3.2	2.94 (0.68)	2.99 (0.73)	2.91 (0.60)	-0.05(0.66)
Aesthetic Pleasure	Human	4.78 (1.98)	4.91 (1.76)	5.00 (1.90)	-
	Claude 3.5	5.59 (0.87)	5.50 (0.77)	5.68 (0.81)	0.33(0.00)
	GPT-4o	5.08 (0.72)	5.02 (0.67)	5.36 (0.68)	0.25(0.04)
	Llama-3.2	5.36 (0.68)	5.06 (1.22)	5.41 (0.78)	0.37(0.00)
Intuitiveness	Human	4.81 (1.89)	4.62 (1.74)	4.83 (1.93)	-
	Claude 3.5	5.83 (0.78)	5.74 (0.82)	5.93 (0.82)	0.33(0.01)
	GPT-4o	5.32 (0.74)	5.29 (0.68)	5.55 (0.64)	0.37(0.00)
	Llama-3.2	4.85 (0.79)	5.19 (0.83)	4.82 (0.85)	-0.00(0.98)
Comfort	Human	4.89 (1.86)	4.95 (1.79)	4.83 (1.97)	-
	Claude 3.5	5.34 (0.74)	5.59 (0.79)	5.29 (0.82)	0.24(0.04)
	GPT-4o	5.07 (0.67)	5.22 (0.65)	5.40 (0.59)	0.20(0.09)
	Llama-3.2	4.94 (0.75)	5.16 (0.87)	5.11 (0.75)	0.19(0.11)

²<https://www.prolific.co/>

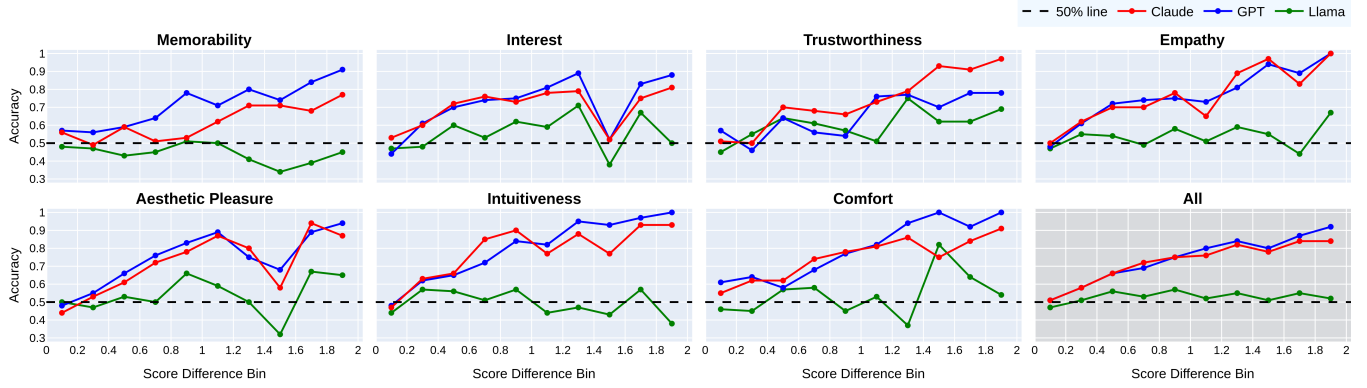


Figure 3: The accuracy of MLLMs in comparing pairs of charts across the seven experiential factors is binned by grouping comparisons based on the magnitude of the difference in human ratings between the chart pairs. The observed overall upward trend suggests that MLLMs perform more accurately when comparing chart pairs with larger score disparities.

4 EVALUATING MLLMS AS A JUDGE

To evaluate the capability of MLLMs as predictors of the experiential impact of the charts, we conducted two studies. By comparing the study results with the crowdsourced data, we evaluated to what extent the judgments of the MLLM are associated with human judgments. To avoid biases toward a specific model, three popular state-of-the-art MLLMs (OpenAI GPT-4o, Anthropic Claude 3.5 Sonnet, and Meta Llama-3.2-11B-Vision-Instruct) are used.

4.1 Task 1: Score Prediction

The first task for MLLMs was similar to what participants did in the crowdsourced data collection: rating the seven factors using a 7-point Likert scale. To simulate the variability inherent in human responses, we generated 216 unique personas³, and instructed MLLMs, acting as “data visualization experts”, to predict how a specific persona would respond to the given chart. Consequently, every chart was assessed by 36 distinct personas, with each persona being employed to evaluate six charts.

Table 2 reports the statistical summary of how different evaluators predicted the seven experiential factors across three topics. While the mean scores from Claude 3.5 Sonnet and GPT-4o were consistently higher than human evaluators, Llama-3.2’s scores were often lower, particularly for the factors of *memorability*, *trustworthiness*, and *empathy*. Overall, MLLMs generate ratings ($0.45 \leq SD \leq 1.22$) are tightly clustered, compared to human ratings ($1.63 \leq SD \leq 2.20$), suggesting that MLLMs are less sensitive to the differences in chart design and effectiveness.

To further assess the alignment between MLLMs and humans as a judge, we performed a rank correlation analysis and reported Kendall’s τ in Table 2. The results indicate that while certain factors show moderate alignment with human ratings, others exhibit only weak or no alignment. For instance, all MLLMs achieved moderate correlation ($0.25 \leq \tau \leq 0.33$) for *aesthetic pleasure*. In the case of *intuitiveness*, GPT-4o and Claude 3.5 Sonnet demonstrate moderate alignments ($\tau = 0.33$ and 0.37 , respectively), whereas Llama-3.2 does not show correlation. For the other factors, MLLMs are either weakly or uncorrelated with human ratings. Notably, there are no cases of significant negative correlations. Given the limited sensitivity of MLLMs and their moderate-to-weak alignment with human ratings, relying on MLLMs to directly predict scores may not be ideal, particularly for tasks that demand nuanced judgment.

4.2 Task 2: Pairwise Comparison

To address the lack of sensitivity found from Task 1, the second task employed MLLMs to compare pairs of charts and decide which

Table 3: Comparison accuracy across different models and factors

Model	Mem	Int	Tru	Emp	Aes	Itt	Cft	All
GPT-4o	0.75	0.66	0.62	0.69	0.70	0.73	0.73	0.70
Claude 3.5	0.64	0.67	0.68	0.70	0.67	0.74	0.70	0.69
Llama-3.2	0.45	0.54	0.58	0.53	0.54	0.51	0.52	0.52

would receive higher scores for each experiential factor from ordinary people. We also instructed MLLMs to provide brief explanations for their choices to enable post hoc analysis. However, the task did not incorporate advanced prompt engineering techniques such as Chain-of-Thought or Few-shot learning. Although such techniques could potentially enhance the performance of MLLM, they also introduce additional variables, complicating the evaluation process, and were therefore beyond the scope of this study.

The accuracy of the comparisons was evaluated against human ratings as a benchmark. For example, if the human ratings for the *memorability* of two charts were 2.5 and 3.5, the comparison was deemed correct if the model identified the second chart as more memorable. As summarized in Table 3, GPT-4o and Claude 3.5 Sonnet demonstrated significantly higher accuracy in overall, compared to Llama-3.2 (One-way ANOVA; $p < 0.001$). However, no statistically significant differences were observed across the seven factors ($p = 0.884$).

Not all comparison tasks pose the same level of difficulty. Pairs of charts that received similar human ratings can be particularly challenging for MLLMs to predict which one would be more effective. Building on this idea, we investigated whether the difference in human ratings influenced the accuracy of the MLLMs’ performance in the comparison task. As shown in Figure 3, the relationship between accuracy and score differences (i.e., task difficulty) exhibits an upward trend, which indicates similar patterns between MLLM and humans. MLLMs perform accurately on problems that are easy for humans, but show low accuracy on more challenging problems. For example, GPT-4o correctly compared all 28 pairs in the 1.4-1.6 bin for *comfort*. Likewise, Claude 3.5 Sonnet was accurate for 32 out of 33 pairs in the 1.4–1.6 bin for *empathy*. Llama-3.2 was an outlier, with its accuracy remaining consistent around 0.5.

5 DISCUSSION

Biases and variability issues. The findings of Task 1 suggest that MLLMs are not as sensitive as human evaluators in evaluating the experiential impact of charts. While human judgments reflect a wide range of responses based on varying experiential factors, MLLMs tend to produce overly consistent results. Moreover, each MLLM demonstrates distinct biases in its evaluations. For example,

³See the supplementary materials for the example prompts.

GPT-4o consistently assigns higher scores on average compared to Llama-3.2 and human evaluators. Furthermore, the performance of MLLMs is not uniform between experiential factors. These observations suggest that relying on MLLMs to directly predict the experiential impact is not ideal, particularly for tasks that require fine-grained sensitivity. On the other hand, Task 2 highlights an alternative use-case: employing MLLMs for comparative evaluations. In this context, MLLMs demonstrate greater accuracy and reliability, particularly when comparing charts with large quality disparities. This underscores the promising use-case of MLLMs in specific use cases, where comparative evaluations are sufficient to achieve evaluation goals.

Limitations and opportunities of MLLMs in chart comparison.

Although the results of the pairwise comparison exhibit similar patterns between MLLMs and humans, there are notable drops between 1.4 and 1.6 for *interest* and *aesthetic pleasure* in Figure 3. A close examination of the incorrect cases reveals that the MLLMs often overrated two specific charts in Figure 4. According to the generated explanations, the MLLMs perceived them as highly interesting and aesthetically pleasing due to their “vivid color scheme”, “visual complexity”, and “analysis over extended timeframes.” In contrast, a majority of human evaluators described them as “hard to read”, “not pleasing”, and even “chaotic.” Once the limitation of MLLMs has been identified, we believe advanced prompt engineering techniques such as Chain-of-Thought or Few-shot prompt might be useful to fix them.

Benchmark as a tool for innovation. Although the primary purpose of benchmarking is to evaluate the performance of AI models, benchmark data have the potential to gain deeper insight into model behavior, as demonstrated above. To this end, we expanded the scope of a benchmark dataset by not only ground-truth human ratings but also explanations for those ratings. In addition, we used the capabilities of MLLM in image understanding and text generation, to generate explanations for their ratings. By comparing human and MLLM-generated explanations, our benchmark dataset becomes a powerful tool for identifying critical issues in AI performance and gaining inspiration to improve MLLMs and their prompts. This approach highlights the potential of benchmark datasets to drive both evaluation and innovation in AI research.

6 CONCLUSION

To evaluate MLLMs’ capability as a judge of the experiential impact of data visualizations, we developed a benchmark dataset comprising 36 charts, accompanied by human ratings and explanations collected from crowdsourced participants. Using this dataset as ground truth, we conducted two tasks: score prediction and pairwise comparison. Our findings reveal that while state-of-the-art MLLMs face challenges in directly predicting scores, they are highly capable at comparing pairs of charts. Lastly, we examined inaccurate cases and discussed the potential benefit of benchmark for identifying issues to improve MLLM’s performance.

This study has a few limitations that highlight opportunities for future research. We did not explore the impact of advanced prompt engineering techniques like Chain-of-Thought or Few-shot prompts

on performance of MLLMs. Second, demographic factors such as educational background, political views, and chart proficiency may influence attitudes and perceptions towards data visualization [36]. Future work could use demographic information from our study’s human evaluators to tailor personalized predictions and enhance the MLLM’s sensitivity. Lastly, future work may introduce new factors depending on different use cases such as the joyfulness or surprise of an interactive and animated chart or the persuasiveness of a chart incorporated with a narrative.

SUPPLEMENTAL MATERIALS

The supplemental materials⁴ include the chart images, human ratings with explanations from the crowdsourced evaluations, and results from Task 1 (direct score prediction) and Task 2 (pairwise comparison). In addition, we provide the prompts used in both tasks and interactive plots for exploring relationships between chart factors, human ratings, and MLLM performance, enabling further analysis and reproducibility.

REFERENCES

- [1] Y. Abe, T. Daikoku, and Y. Kuniyoshi. Assessing the aesthetic evaluation capabilities of gpt-4 with vision: Insights from group and individual assessments. pp. 2Q1IS301–2Q1IS301, 2024.
- [2] M. Akhtar, N. Subedi, V. Gupta, S. Tahmasebi, O. Cocarascu, and E. Simperl. Chartcheck: An evidence-based fact-checking dataset over real-world chart images. *arXiv preprint arXiv:2311.07453*, 2023.
- [3] L. Bartram, A. Patra, and M. Stone. Affective color in visualization. In *Proceedings of the 2017 CHI conference on human factors in computing systems*, pp. 1364–1374, 2017.
- [4] M. Behrisch, M. Blumenschein, N. W. Kim, L. Shao, M. El-Assady, J. Fuchs, D. Seebacher, A. Diehl, U. Brandes, H. Pfister, et al. Quality metrics for information visualization. In *Computer Graphics Forum*, vol. 37, pp. 625–662. Wiley Online Library, 2018.
- [5] M. Binder, B. Heinrich, M. Hopf, and A. Schiller. Global reconstruction of language models with linguistic rules—explainable ai for online consumer reviews. *Electronic Markets*, 32(4):2123–2138, 2022.
- [6] M. A. Borkin, A. A. Vo, Z. Bylinskii, P. Isola, S. Sunkavalli, A. Oliva, and H. Pfister. What makes a visualization memorable? *IEEE transactions on visualization and computer graphics*, 19(12):2306–2315, 2013.
- [7] J. Boy, A. V. Pandey, J. Emerson, M. Satterthwaite, O. Nov, and E. Bertini. Showing people behind data: Does anthropomorphizing visualizations elicit more empathy for human rights data? In *Proceedings of the 2017 CHI conference on human factors in computing systems*, pp. 5462–5474, 2017.
- [8] D. Chen, R. Chen, S. Zhang, Y. Liu, Y. Wang, H. Zhou, Q. Zhang, Y. Wan, P. Zhou, and L. Sun. Mllm-as-a-judge: Assessing multimodal llm-as-a-judge with vision-language benchmark. *arXiv preprint arXiv:2402.04788*, 2024.
- [9] R. Datta, D. Joshi, J. Li, and J. Z. Wang. Studying aesthetics in photographic images using a computational approach. In *Computer Vision—ECCV 2006: 9th European Conference on Computer Vision, Graz, Austria, May 7–13, 2006, Proceedings, Part III 9*, pp. 288–301. Springer, 2006.
- [10] A. Dave, A. Saxena, and A. Jha. Understanding user comfort and expectations in ai-based systems. 2023.
- [11] K. Deng, A. Ray, R. Tan, S. Gabriel, B. A. Plummer, and K. Saenko. Socratis: Are large multimodal models emotionally aware? *arXiv preprint arXiv:2308.16741*, 2023.
- [12] P. Duan, J. Warner, Y. Li, and B. Hartmann. Generating automatic feedback on ui mockups with large language models. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 1–20, 2024.
- [13] N. Errey, J. Liang, T. W. Leong, and D. Zowghi. Evaluating narrative visualization: a survey of practitioners. *International Journal of Data Science and Analytics*, 18(1):19–34, 2024.

⁴available at <http://chart2experience.github.io>

- [14] S. Few and P. Edge. Data visualization effectiveness profile. *Perceptual Edge*, 10:12, 2017.
- [15] Y. Guo, F. Siddiqui, Y. Zhao, R. Chellappa, and S.-Y. Lo. Stimuvar: Spatiotemporal stimuli-aware video affective reasoning with multimodal large language models. *arXiv preprint arXiv:2409.00304*, 2024.
- [16] Y. Han, C. Zhang, X. Chen, X. Yang, Z. Wang, G. Yu, B. Fu, and H. Zhang. Chartllama: A multimodal llm for chart understanding and generation. *arXiv preprint arXiv:2311.16483*, 2023.
- [17] K.-H. Huang, H. P. Chan, Y. R. Fung, H. Qiu, M. Zhou, S. Joty, S.-F. Chang, and H. Ji. From pixels to insights: A survey on automatic chart understanding in the era of large foundation models. *arXiv preprint arXiv:2403.12027*, 2024.
- [18] W. Huang, P. Eades, and S.-H. Hong. Measuring effectiveness of graph visualizations: A cognitive load perspective. *Information Visualization*, 8(3):139–152, 2009.
- [19] L. Huddy and A. H. Gunthorsdottir. The persuasive effects of emotive visual imagery: Superficial manipulation or the product of passionate reason? *Political Psychology*, 21(4):745–778, 2000.
- [20] P. Isola, J. Xiao, D. Parikh, A. Torralba, and A. Oliva. What makes a photograph memorable? *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1469–1482, 2013.
- [21] K. Kafle, B. Price, S. Cohen, and C. Kanan. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5648–5656, 2018.
- [22] S. Kantharaj, R. T. K. Leong, X. Lin, A. Masry, M. Thakkar, E. Hoque, and S. Joty. Chart-to-text: A large-scale benchmark for chart summarization. *arXiv preprint arXiv:2203.06486*, 2022.
- [23] X. Lan, Y. Shi, Y. Wu, X. Jiao, and N. Cao. Kineticharts: Augmenting affective expressiveness of charts in data stories with animation design. *IEEE Transactions on Visualization and Computer Graphics*, 28(1):933–943, 2021.
- [24] X. Lan, Y. Shi, Y. Zhang, and N. Cao. Smile or scowl? looking at infographic design through the affective lens. *IEEE Transactions on Visualization and Computer Graphics*, 27(6):2796–2807, 2021.
- [25] X. Lan, Y. Wu, and N. Cao. Affective visualization design: Leveraging the emotional impact of data. *IEEE Transactions on Visualization and Computer Graphics*, 2023.
- [26] S. Lee, S. Kim, S. H. Park, G. Kim, and M. Seo. Prometheusvision: Vision-language model as a judge for fine-grained evaluation. *arXiv preprint arXiv:2401.06591*, 2024.
- [27] E. Lee-Robbins and E. Adar. Affective learning objectives for communicative visualizations. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):1–11, 2022.
- [28] D. Li, B. Jiang, L. Huang, A. Beigi, C. Zhao, Z. Tan, A. Bhat-tacharjee, Y. Jiang, C. Chen, T. Wu, et al. From generation to judgment: Opportunities and challenges of llm-as-a-judge. *arXiv preprint arXiv:2411.16594*, 2024.
- [29] Z. Lian, L. Sun, H. Sun, K. Chen, Z. Wen, H. Gu, B. Liu, and J. Tao. Gpt-4v with emotion: A zero-shot benchmark for generalized emotion recognition. *Information Fusion*, 108:102367, 2024.
- [30] J. Liem, C. Perin, and J. Wood. Structure and empathy in visual data storytelling: Evaluating their influence on attitude. In *Computer Graphics Forum*, vol. 39, pp. 277–289. Wiley Online Library, 2020.
- [31] F. Liu, J. M. Eisenschlos, F. Piccinno, S. Krichene, C. Pang, K. Lee, M. Joshi, W. Chen, N. Collier, and Y. Altun. Deplot: One-shot visual language reasoning by plot-to-table translation. *arXiv preprint arXiv:2212.10505*, 2022.
- [32] J. Machajdik and A. Hanbury. Affective image classification using features inspired by psychology and art theory. In *Proceedings of the 18th ACM international conference on Multimedia*, pp. 83–92, 2010.
- [33] A. Mehrabian. An approach to environmental psychology. *Massachusetts Institute of Technology*, 1974.
- [34] L. Micallef, G. Palmas, A. Oulasvirta, and T. Weinkauff. Towards perceptual optimization of the visual design of scatterplots. *IEEE transactions on visualization and computer graphics*, 23(6):1588–1599, 2017.
- [35] J. Obaid and E. Hoque. Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model. *arXiv preprint arXiv:2010.09142*, 2020.
- [36] E. M. Peck, S. E. Ayuso, and O. El-Etr. Data is personal: Attitudes and perceptions of data visualization in rural pennsylvania. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pp. 1–12, 2019.
- [37] J. Sabini and M. Silver. Ekman’s basic emotions: Why not love and jealousy? *Cognition & Emotion*, 19(5):693–712, 2005.
- [38] B. Saket, A. Endert, and J. Stasko. Beyond usability and performance: A review of user experience-focused evaluations in visualization. In *Proceedings of the Sixth Workshop on Beyond Time and Errors on Novel Evaluation Methods for Visualization*, pp. 133–142, 2016.
- [39] J. Stasko. Value-driven evaluation of visualizations. In *Proceedings of the Fifth Workshop on Beyond Time and Errors: Novel Evaluation Methods for Visualization*, pp. 46–53, 2014.
- [40] T. Sun, Y. Shao, H. Qian, X. Huang, and X. Qiu. Black-box tuning for language-model-as-a-service. In *International Conference on Machine Learning*, pp. 20841–20855. PMLR, 2022.
- [41] B. J. Tang, A. Boggust, and A. Satyanarayan. Vistext: A benchmark for semantically rich chart captioning. *arXiv preprint arXiv:2307.05356*, 2023.
- [42] J. Tang, Q. Liu, Y. Ye, J. Lu, S. Wei, C. Lin, W. Li, M. F. F. B. Mahmood, H. Feng, Z. Zhao, et al. Mtvqa: Benchmarking multilingual text-centric visual question answering. *arXiv preprint arXiv:2405.11985*, 2024.
- [43] A. Tatu, P. Bak, E. Bertini, D. Keim, and J. Schneidewind. Visual quality metrics and human perception: an initial study on 2d projections of large multidimensional data. In *Proceedings of the international conference on advanced visual interfaces*, pp. 49–56, 2010.
- [44] P. Valdez and A. Mehrabian. Effects of color on emotions. *Journal of experimental psychology: General*, 123(4):394, 1994.
- [45] Y. Wang, A. Segal, R. Klatzky, D. F. Keefe, P. Isenberg, J. Hurtienne, E. Hornecker, T. Dwyer, and S. Barrass. An emotional response to the value of visualization. *IEEE computer graphics and applications*, 39(5):8–17, 2019.
- [46] T. Willigen. Measuring the user experience of data visualization. Master’s thesis, University of Twente, 2019.
- [47] Y. Wu, C. Bauckhage, and C. Thureau. The good, the bad, and the ugly: Predicting aesthetic image labels. In *2010 20th International Conference on Pattern Recognition*, pp. 1586–1589. IEEE, 2010.
- [48] J. Ye, Y. Wang, Y. Huang, D. Chen, Q. Zhang, N. Moniz, T. Gao, W. Geyer, C. Huang, P.-Y. Chen, et al. Justice or prejudice? quantifying biases in llm-as-a-judge. *arXiv preprint arXiv:2410.02736*, 2024.
- [49] J. Zamfirescu-Pereira, R. Y. Wong, B. Hartmann, and Q. Yang. Why johnny can’t prompt: how non-ai experts try (and fail) to design llm prompts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1–21, 2023.
- [50] H. Zhang, E. Augilius, T. Honkela, J. Laaksonen, H. Gamper, and H. Alene. Analyzing emotional semantics of abstract art using low-level image features. In *Advances in Intelligent Data Analysis X: 10th International Symposium, IDA 2011, Porto, Portugal, October 29-31, 2011. Proceedings 10*, pp. 413–423. Springer, 2011.
- [51] Y. Zhang, M. Wang, P. Tiwari, Q. Li, B. Wang, and J. Qin. Dialogue-llm: Context and emotion knowledge-tuned llama models for emotion recognition in conversations. *arXiv preprint arXiv:2310.11374*, 2023.
- [52] S. Zhao, G. Ding, Q. Huang, T.-S. Chua, B. Schuller, and K. Keutzer. Affective image content analysis: A comprehensive survey. 2018.
- [53] S. Zhao, Y. Gao, X. Jiang, H. Yao, T.-S. Chua, and X. Sun. Exploring principles-of-art features for image emotion recognition. In *Proceedings of the 22nd ACM international conference on Multimedia*, pp. 47–56, 2014.
- [54] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- [55] M. Zhou, Y. R. Fung, L. Chen, C. Thomas, H. Ji, and S.-F. Chang. Enhanced chart understanding in vision and language task via cross-modal pre-training on plot table pairs. *arXiv preprint arXiv:2305.18641*, 2023.