

Do MLLMs Capture How Interfaces Guide User Behavior? A Benchmark for Multimodal UI/UX Design Understanding

Anonymous ACL submission

Abstract

User interface (UI) design goes beyond visuals to shape user experience (UX), underscoring the shift toward UI/UX as a unified concept. While recent studies have explored UI evaluation using Multimodal Large Language Models (MLLMs), they largely focus on surface-level features, overlooking how design choices influence user behavior at scale. To fill this gap, we introduce WisierUI-Bench, a novel benchmark for *multimodal understanding of how UI/UX design affects user behavior*, built on 300 real-world UI image pairs from industry A/B tests, with empirically validated winners that induced more user actions. For future design progress in practice, post-hoc understanding of why such winners succeed with mass users is also required; we support this via expert-curated key interpretations for each instance. Experiments across multiple MLLMs on WisierUI-Bench for two main tasks, (1) predicting the more effective UI image between an A/B-tested pair, and (2) explaining it post-hoc in alignment with expert interpretations, show that models exhibit limited understanding of the behavioral impact of UI/UX design. We believe our work will foster research on leveraging MLLMs for visual design in user behavior contexts.

1 Introduction

User interface (UI) design is central to application development. While visual aesthetics lay its foundation, the fundamental goal is to steer user behavior during interaction, such as encouraging sign-ups or purchases, and thereby boost service revenue (Fogg, 2002). This behavioral focus expands to *UI/UX design*, uniting visual form with user experience (UX) considerations (Norman, 1988).

Practitioners validate such decisions through large-scale A/B tests that randomly assign users to alternative UI variants and measure which design more effectively drives the desired action (Kohavi et al., 2009). For example, as shown in Figure

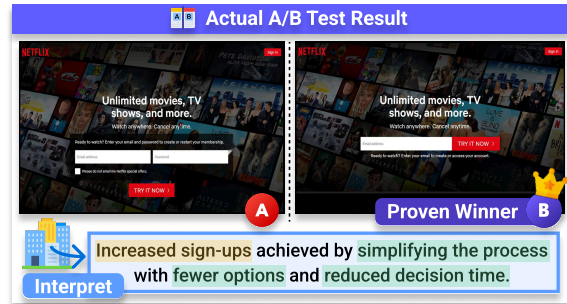


Figure 1: A real-world, behavior-aware design decision example from WisierUI-Bench, grounded in A/B test results, illustrating how UI changes steer user actions.

1, version B achieves higher sign-ups overall, interpreted as reduced user friction from fewer options. Such interpretation reflects an abductive reasoning process (Hobbs et al., 1993; Bhagavatula et al., 2020), inferring the most plausible explanations for aggregate outcomes that are used to guide subsequent design decisions (Quin et al., 2024). Given that predicting and explaining these results is nontrivial even for humans, hinging on complex cognitive processes underlying user experience, assessing *multimodal understanding of UI/UX design* becomes a crucial challenge, requiring not only identifying visual elements but also reasoning about how they influence user behavior.

Recent efforts have explored how to evaluate UI design quality with Multimodal Large Language Models (MLLMs); however, they all fall short in capturing the behavioral impact of UI/UX design, as summarized in Table 1. While some focus solely on basic visual, heuristic-based attributes (Yang and Li, 2024; Wu et al., 2024), others rely on handwritten critiques of individual UIs, lacking validation against real user data and offering limited coverage of cognitive aspects (Duan et al., 2024).

To address these limitations, we introduce **WisierUI-Bench** (Pairwise UI/UX Design Understanding Benchmark), a novel benchmark for eval-

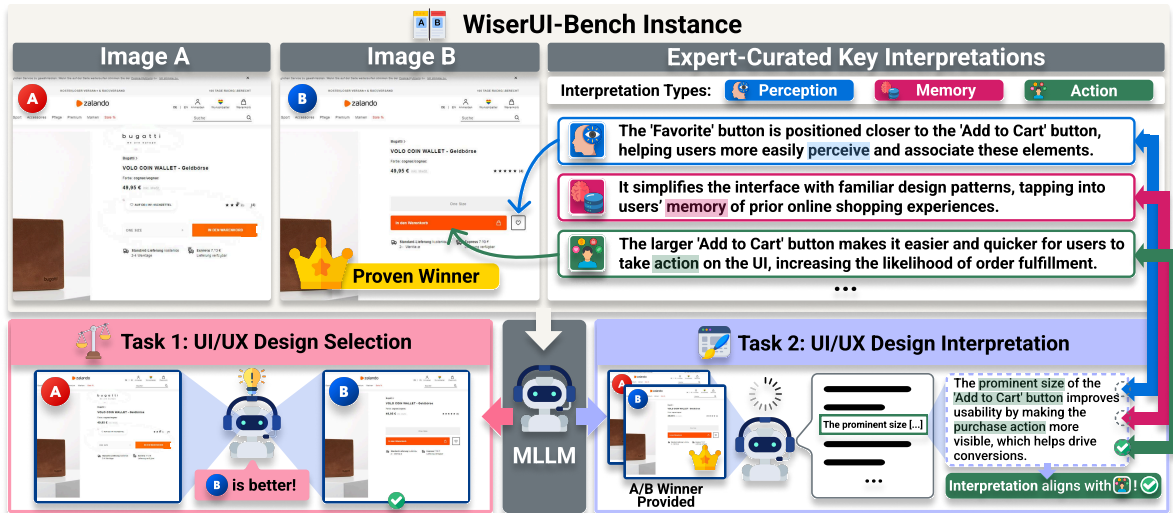


Figure 2: Overview of WiserUI-Bench and two main tasks. Each instance contains a UI image pair with verified A/B test winner and expert-curated key interpretations explaining it. These span the three cognitive UX dimensions; the example includes all three, though most cover fewer. MLLMs are evaluated on (1) selecting the more effective UI/UX design by predicting the verified winner, and (2) explaining the effectiveness of a given winner, measured by whether the model captures each expert interpretation.

uating a model’s multimodal understanding of how UI/UX design affects user behavior. WiserUI-Bench comprises 300 real-world UI image pairs with winning variants verified by A/B tests from actual companies, spanning varied contexts and sourced from reliable platforms. Each pair includes expert-curated key interpretations, 684 in total, written post-hoc with knowledge of the winner to explain its effectiveness, covering core behavioral drivers. An example instance is illustrated in the upper part of Figure 2. Grounded in large-scale real user behavior data, our benchmark captures authentic usage and preference patterns, providing a reliable basis for evaluation. Each interpretation is tagged with one of the three cognitive UX dimensions articulated by Norman (2007): perception (what users notice), memory (what users remember), and action (how users act).

Our benchmark supports two core tasks, illustrated in the lower part of Figure 2: (1) selecting the more effective UI/UX design from a given pair by predicting the A/B test-verified winner, and (2) interpreting the UI/UX effectiveness of the given winner by assessing alignment between model-generated and expert-curated interpretations. This task setup captures the essential capabilities for UI/UX assessment, not only to identify which design better guides user behavior, but also to reliably explain why in a post-hoc manner.

Experiments on a broad range of proprietary and open-source MLLMs on WiserUI-Bench show

that existing models struggle with visual reasoning in user behavior contexts. Models perform near-random on selection, and their post-hoc interpretations remain insufficient to reach expert level. We expect that WiserUI-Bench and our novel UI/UX-focused tasks will spur further research on this critical and underexplored, user-centered area, advancing behavior-aware reasoning about visual design in MLLMs.

In summary, our main contributions are:

- We introduce a novel and challenging task of multimodal UI/UX understanding that centers on the behavioral impact of design.
- We release WiserUI-Bench, a benchmark of 300 real-world UI pairs uniquely constructed from verified results of large-scale industry A/B tests, along with 684 carefully curated expert key interpretations.
- Through comprehensive experiments and analyses across diverse MLLMs, we uncover their limited ability to understand the behavioral impact of UI/UX design, highlighting a critical gap in current visual reasoning.

2 Related Work

2.1 Visual Reasoning

Visual reasoning benchmarks for Multimodal Large Language Models (MLLMs) evaluate models’ ability to process visual information and reason at a

Benchmark	# Samples	Unit	Source Quality	Source Type	Large-Scale User Validation	Expert Interpretation	Evaluation Objective
Yang and Li (2024)	382	Single	Legacy UI	Mobile	✗	✗	Guideline Violation Detection
BetterWeb (Wu et al., 2024)	892	Pair	Synthetic	Mobile	✗	✗	Basic Visual Quality
UICrit (Duan et al., 2024)	983	Single	Legacy UI	Mobile	✗	✓	UI Design Critique
WiserUI-Bench (Ours)	300	Pair	Actual Production	Mobile + Web	✓ (Real A/B Test)	✓	UI/UX Design Understanding (User Behavior-Focused)

Table 1: Comparison between WiserUI-Bench and existing UI evaluation benchmarks.

cognitive level. While some target general tasks (Yue et al., 2024; Liu et al., 2024b; Fu et al., 2024; Li et al., 2024b), others focus on specific domains like mathematics or science (Hao et al., 2025; Lu et al., 2024; Xu et al., 2025). However, few benchmarks address reasoning about human behavior or preference patterns in visual contexts. Moreover, although some studies explore visual reasoning across multiple images (Cheng et al., 2025; Zhao et al., 2024), such cases remain rare. Our task incorporates both aspects, reasoning over paired UI images and aligning with behavior patterns, as captured by actual A/B test outcomes.

2.2 UI Evaluation

Recent advances have driven progress in evaluating UI quality with MLLMs. Yang and Li (2024) propose a benchmark for detecting guideline violations in individual screens, while Wu et al. (2024) generates synthetic UI pairs via perturbations such as color noise or spacing changes. However, both are limited to heuristic analysis of surface-level appearance. Duan et al. (2024) evaluates single screens using expert-critique annotations, but it lacks validation with real-world user behavior and relies on RICO (Deka et al., 2017), an outdated, mobile-only dataset that limits its relevance to today’s production-level designs. Prior work also suggests that pairwise evaluation yields more reliable signals, particularly for subjective tasks like ours that target behavior patterns captured by real-world A/B tests (Abel et al., 2020; Bawabe et al., 2021). In this context, WiserUI-Bench offers a stronger foundation by grounding design effectiveness in large-scale user behavior data and requiring models to engage in deeper, nuanced reasoning over both visuals and UX in a pairwise format.

2.3 Simulation of Human Behavior

Efforts to simulate human behavior with AI are ongoing. Park et al. (2023) places LLM-driven

characters in a sandbox town, producing socially plausible interactions. Lu et al. (2025a) moves closer to real behavior by using actual platform logs to predict next user actions from prior context, yet remains text-only. In the UI domain, Lu et al. (2025b) and Wang et al. (2025a) deploy persona-based agents to browse interfaces and examine how designs may steer user behavior; the former checks realism only via qualitative interviews, while the latter compares against an actual A/B test but only under a single scenario. By contrast, our work grounds evaluation in numerous real-user A/B test results and introduces diverse multimodal tasks that reason over UI variant image pairs to predict and explain user behavior.

3 WiserUI-Bench

We propose WiserUI-Bench, a new benchmark for evaluating multimodal understanding of UI/UX design, a key factor in crafting interfaces that influence user behavior.

3.1 Benchmark Construction

Data Sources To construct a benchmark that faithfully reflects real-world UI/UX design decisions, we collect UI images from industry-validated A/B testing data. Specifically, we aggregate data from widely recognized A/B testing and analysis platforms¹, which showcase large-scale A/B tests conducted by companies across various industries and regions. These sources provide both UI variants and verified outcomes based on actual user behavior data. By grounding our benchmark in such real, rigorously tested cases, we ensure a credible basis for evaluating UI/UX design understanding in practical scenarios.

Data Curation Pipeline As illustrated in Figure 4, our curation process involves multiple stages to

¹<https://vwo.com/success-stories/>, <https://goodui.org/leaks/>, <https://abtest.design/>

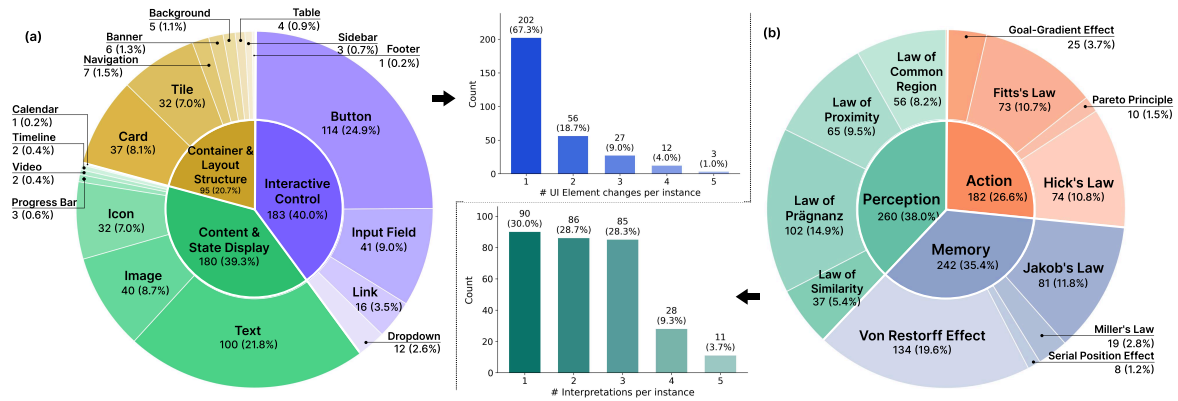


Figure 3: Distribution of (a) UI change element types and (b) expert interpretations for WisierUI-Bench.

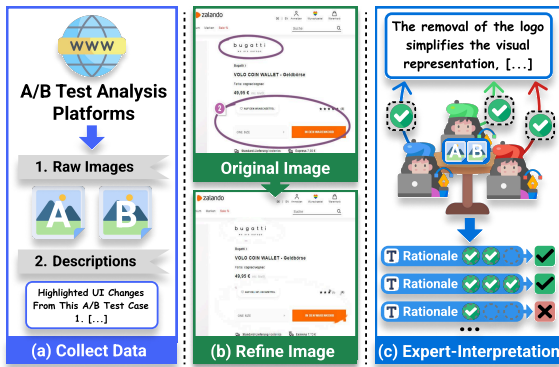


Figure 4: Benchmark construction pipeline. (a) Raw A/B test data were collected from reliable platforms, (b) then UI images were refined for clarity. (c) Finally, key interpretations explaining the test outcomes were curated by UI/UX experts.

ensure high-quality annotations. First, we extract raw UI images and accompanying textual descriptions from our selected platform sources. Next, we refine these images by eliminating any visual markers that were externally added, such as arrows or circles, to preserve a clean, unbiased evaluation of visual reasoning.

Finally, we construct key interpretations explaining the winning variant’s superior performance observed in A/B test results, curated by three experienced UI/UX experts aware of the outcomes. Such expert involvement is essential, as interpreting user behavior often requires deep domain expertise. To identify core behavioral drivers, only interpretations independently identified by at least two experts were retained. This resulted in final 684 (562 with full, 122 with two-way agreement) out of 845 initial candidates. The process was guided by textual descriptions from data sources and well-established principles, the Nielsen Norman 10 Usability Heuristics (Nielsen, 1994). Each was then

tagged with the most relevant law from a curated set of 12 UX laws. Since user behavior is often shaped by multiple factors, a single instance may include multiple interpretations. See Appendix A.1 for details.

3.2 Benchmark Analysis

Our pipeline produces a benchmark of 300 real-world UI image pairs, each with a definitive winner identified through A/B test results and annotated with expert-curated key interpretations of the core design factors driving user behavior, totaling 684 entries. We include detailed statistics in Appx. A.2.

Ensuring contextual diversity is essential, as effective UI/UX strategies vary widely across use cases (Oinas-Kukkonen and Harjumaa, 2018). For example, page types correspond to different stages of the user journey, requiring distinct strategies. To capture this, WisierUI-Bench spans 11 page types, 13 industry domains, and two device types.

Moreover, we annotate UI differences using element–attribute associations, specifying the UI element being modified and the corresponding attribute change (e.g., a *button* with a *position* change), to capture real-world UI variants with concurrent modifications. On average, each instance contains 1.53 element-level and 1.99 attribute-level annotations, spanning 19 element types and 14 attribute types. We further group the element types by their primary functional roles (Figure 3(a)): *Interactive Control*, which directly elicits intentional user actions; *Content & State Display*, which primarily conveys information, semantics, or system state; and *Container & Layout Structure*, which organizes spatial layout and guides attention at a structural level.

WisierUI-Bench also exhibits diversity in the interpretations used to explain user behavior, as

# Cases	Overall (300)				IC (100)		CSD (79)		CLS (49)		Mixed (72)	
	FA	SA	AA	CA	AA	CA	AA	CA	AA	CA	AA	CA
Random	50.00	50.00	50.00	25.00	50.00	25.00	50.00	25.00	50.00	25.00	50.00	25.00
o1	16.56	97.78	57.17	15.56	55.67	16.00	55.70	11.81	56.46	11.56	61.34	21.76
GPT-4o	31.89	88.33	60.11	30.11	59.67	31.00	61.60	29.96	51.70	17.69	64.81	37.50
Claude 3.5 Sonnet	26.11	86.67	56.39	24.22	56.67	22.67	57.59	<u>29.96</u>	51.70	12.24	57.87	28.24
Qwen2.5-VL-32B	33.67	81.44	57.56	31.00	58.67	<u>31.67</u>	56.54	28.69	58.84	28.57	56.25	34.26
Qwen2.5-VL-7B	23.00	77.89	50.44	12.89	48.17	8.67	56.12	19.41	43.88	9.52	51.85	13.89
InternVL-2.5-38B	55.67	60.00	<u>57.83</u>	34.56	56.17	34.33	59.70	36.29	55.10	28.57	59.95	<u>37.04</u>
InternVL-2.5-8B	58.33	45.67	52.00	21.33	52.17	23.00	48.73	15.61	<u>57.14</u>	<u>25.85</u>	51.85	22.22
LLaVA-NeXT-7B	54.44	40.56	47.50	10.78	48.50	<u>12.67</u>	43.67	8.44	46.60	8.16	50.93	12.50
LLaVA-OneVision-7B	19.78	81.56	50.67	10.44	50.00	11.00	54.43	<u>13.50</u>	47.28	3.40	49.77	11.11

Table 2: Results of the UI/UX design selection task on WisierUI-Bench, shown overall and by UI change element types. Best and second-best model scores for each metric are marked in **bold** and underlined, respectively. For CA, the best and second-best UI element type performances within each model are highlighted in blue and yellow, respectively. All metrics are reported as percentages (%).

shown in Figure 3(b), with an average of 2.28 interpretations per instance. The tagged UX laws are grouped into three cognitive dimensions of UX, perception, memory, and action, following Norman (2007), with four laws assigned to each.

4 Task Description

Proper evaluation of a model’s multimodal understanding of UI/UX design requires assessing not only its ability to predict which design more effectively guides user behavior, but also whether it can explain *why* in a post-hoc manner, grounding its reasoning in visual and behavioral cues. To this end, we define two main tasks, selection and interpretation, as illustrated in Figure 2, both essential for demonstrating genuine understanding. Success requires identifying key visual differences between two UI variants, as well as inferring the underlying behavioral goal from context (e.g., promoting sign-ups) and how the interface supports that goal.

Task 1: UI/UX Design Selection The first task assesses whether the model can identify which of two UI images is more effective at influencing user behavior, with the correct answer determined by large-scale A/B testing run by real companies. Although the task is framed as a binary choice, solving it effectively requires nuanced visual reasoning beyond surface-level comparison.

Task 2: UI/UX Design Interpretation The second task evaluates a model’s post-hoc interpretation ability on the observed results. Unlike the first task, it presents the winning UI/UX design upfront and asks the model to explain the factors

contributing to its effectiveness in guiding user behavior. The model’s explanation is then compared to expert-curated key interpretations, with the degree of alignment serving as a measure of the quality of its expert-level interpretive understanding.

5 Experiments

Models We evaluate a diverse set of MLLMs, spanning both *proprietary* and *open-source* models. On the proprietary side, we include o1 (OpenAI et al., 2024b), GPT-4o (OpenAI et al., 2024a), Claude 3.5 Sonnet (Anthropic, 2024), all noted for their strong visual reasoning. For open-source models, we evaluate widely adopted models such as Qwen2.5-VL (7B, 32B) (Bai et al., 2025), InternVL-2.5 (8B, 38B) (Chen et al., 2025), LLaVA-NeXT-7B (Liu et al., 2024a) and LLaVA-OneVision-7B (Li et al., 2024a). Refer to Appendix B.1 for details.

5.1 UI/UX Design Selection

5.1.1 Evaluation Metrics

We report model accuracies across input-order variations of each UI image pair, reflecting the multi-image input setting. Specifically, we report: (1) *First Accuracy (FA)*, when the more effective UI/UX design is presented first in the input; (2) *Second Accuracy (SA)*, when it is presented second; and (3) *Average Accuracy (AA)*, the mean of FA and SA, reflecting overall correctness. Finally, we compute (4) *Consistent Accuracy (CA)*, which checks whether the model selects the same, correct UI across both input orders of each pair. This metric isolates the model’s core visual reasoning ability by ensuring that its choices are based on content

Model	Method	FA	SA	AA	CA
GPT-4o	Zero-Shot	31.89	88.33	60.11	30.11
	CoCoT	33.00	83.33	58.17	29.67
	Self-Refine	32.44	82.67	57.56	29.89
	DDCoT	36.56	86.89	61.72	34.78
	MAD (R1)	51.22	67.44	59.33	39.00
	MAD (R3)	34.44	80.89	57.67	29.89
Claude 3.5 Sonnet	Zero-Shot	26.11	86.67	56.39	24.22
	CoCoT	29.78	81.78	55.78	27.22
	Self-Refine	25.67	76.44	51.06	20.44
	DDCoT	26.22	86.67	56.44	23.56
	MAD (R1)	39.00	72.89	55.94	30.00
	MAD (R3)	23.33	85.44	54.39	21.44

Table 3: Results of the UI/UX design selection task across reasoning strategies on WisierUI-Bench, figures in percentages (%).

rather than position bias, a well-known issue in pairwise comparison settings (Wang et al., 2025b; Zhang et al., 2023). (Random guessing yields a CA of 25%: since there’s a 50% chance of being correct in each order.) All evaluations are conducted over three independent runs, and we report the averaged results. Experiment details, including the prompts, are provided in Appendix B.2.

5.1.2 Results

Overall Performance Table 2 summarizes model performance on the UI/UX design selection task. Most AA scores hover only slightly above the 50% random baseline for a binary choice. Moreover, they are inflated by strong position bias, especially in proprietary models and the Qwen2.5-VL family, which tend to pick the second image. When we examine CA, the order-invariant metric, performance drops back to near-random levels. Interestingly, Qwen2.5-VL-32B and InternVL-2.5-38B score slightly higher than proprietary models on CA, but still fall short of being adequate. Contrary to expectations, o1 shows the strongest bias and a lower CA than GPT-4o, suggesting that self-reasoning of o1 does not transfer well to this task. Smaller (7B/8B) models perform significantly worse than larger ones, confirming the scaling trend. Overall, existing MLLMs struggle to consistently predict the more effective UI/UX design in real-world UI pairs.

Performance Across UI Element Types The UI element-wise results in Table 2 show performance variations across change type categories, while preserving overall trends: GPT-4o and InternVL-2.5-38B generally achieve the highest AA and CA, respectively. Focusing on CA, Mixed cases,

Model	Interpretation Recall	Instance-level Recall
o1	64.18	78.33
GPT-4o	50.15	66.67
Claude 3.5 Sonnet	63.16	78.67
Qwen2.5-VL-32B	60.23	73.00
Qwen2.5-VL-7B	43.71	61.67
InternVL-2.5-38B	45.32	59.33
InternVL-2.5-8B	35.09	51.67
LLaVA-NeXT-7B	11.99	21.67
LLaVA-OneVision-7B	21.78	33.33

Table 4: Results of the UI/UX design interpretation task on WisierUI-Bench, figures in percentages (%).

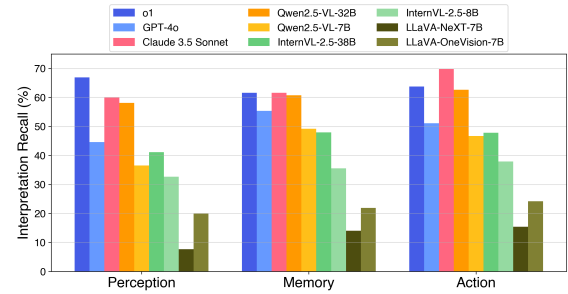


Figure 5: Results of the UI/UX design interpretation task by each UX dimension on WisierUI-Bench.

where changes span multiple functional UI element groups, perform best across most models, suggesting that aggregated multi-functional changes provide stronger behavioral cues than single-group modifications. Interactive Control (IC) and Content & State Display (CSD) exhibit comparable CA and rank below Mixed, indicating relatively moderate sensitivity to action-oriented (e.g., buttons) and content-level changes (e.g., images or text). By contrast, Container & Layout Structure (CLS) is the most challenging category, with sharply reduced CA across most models, reflecting limited sensitivity to structural layout changes, such as card or tile-level differences. Notably, Qwen2.5-VL-32B and InternVL-2.5-8B perform relatively better on CLS than other models, including GPT-4o.

Effect of Reasoning Strategies We additionally evaluate various reasoning strategies on GPT-4o and Claude 3.5 Sonnet to assess their effectiveness, including Contrastive Chain-of-Thought (CoCoT) (Zhang et al., 2024), Self-Refine (Madaan et al., 2023), Duty-Distinct Chain-of-Thought (DDCoT) (Zheng et al., 2023), and Multi-Agent Debate (MAD) (Liang et al., 2024; Du et al., 2023). For MAD, we stop and moderate the debate after the first (R1) and third (R3) rounds. As shown in Ta-

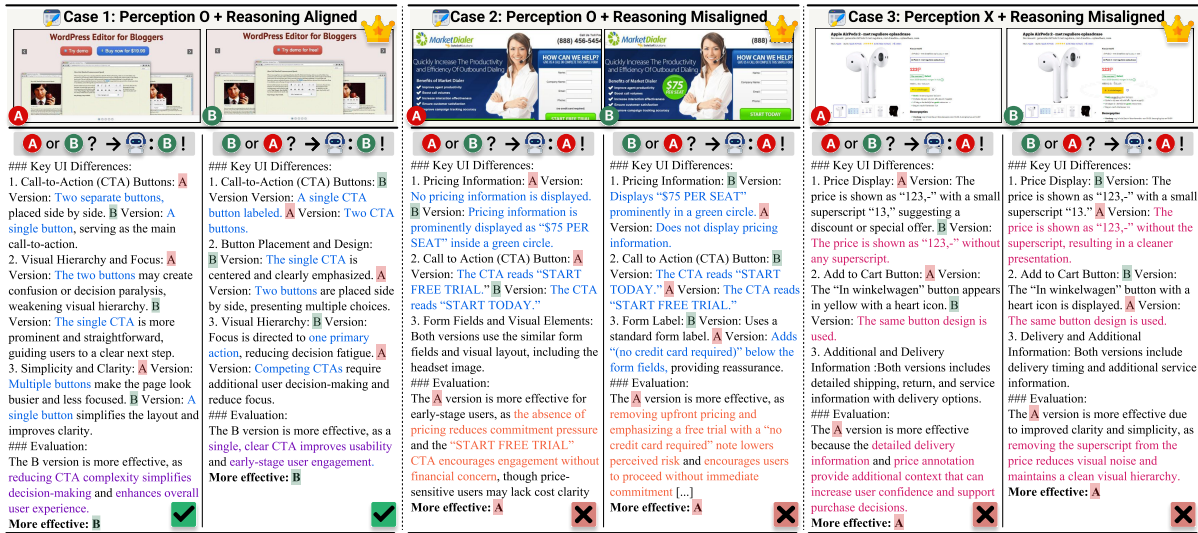


Figure 6: Example responses from InternVL-2.5-38B on the UI/UX design selection task, showing model outputs under different input image orders for each case. Blue marks correct UI change perception, while red denotes incorrect or hallucinated perception. Purple indicates user behavior pattern-aligned reasoning, whereas misaligned reasoning in orange.

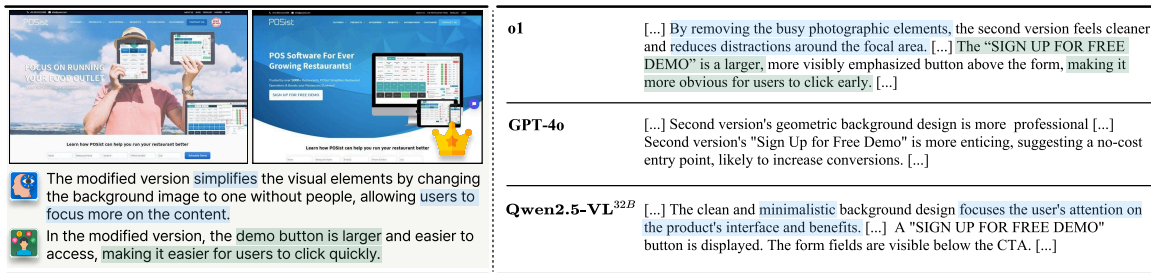


Figure 7: Example responses from o1, GPT-4o, and Qwen2.5-VL-32B on the UI/UX design interpretation task. Green and blue highlight points capturing the expert interpretations fully.

ble 3, AA scores remain similar to or below the zero-shot setting, while CA scores often surpass it, especially for MAD (R1) across models, DD-CoT on GPT-4o, and CoCoT on Claude 3.5 Sonnet. This gain is mainly due to reduced position bias, enabling more consistent decisions regardless of input order. MAD (R1) performs well by generating and moderating conflicting perspectives, which suits this human behavior pattern-driven task, benefiting from multi-perspective reasoning. However, excessive internal debate may hurt performance, as seen in the drop with MAD (R3).

5.2 UI/UX Design Interpretation

5.2.1 Evaluation Metrics

To judge how well a model’s generated interpretation explains the observed outcome, we compare it against expert-curated key interpretations for each instance. Specifically, we instruct the model to produce an unconstrained free-form explanation

enumerating all relevant factors, and each expert interpretation is independently evaluated in a binary manner, whether it is fully captured by the explanation or not. We report two metrics: *Interpretation Recall*, measuring the proportion of expert interpretations captured, and *Instance-level Recall*, which counts an instance as correct if at least one expert interpretation is captured. Binary inclusion judgments are made using a GPT-4o-based evaluator (Liu et al., 2023), validated against human judgments on 1,000 random samples (83.0% accuracy, Cohen’s kappa 0.66). The prompts used, along with the results of BLEU, ROUGE, METEOR, and BERTScore (Papineni et al., 2002; Lin, 2004; Banerjee and Lavie, 2005; Zhang et al., 2020) are reported in the Appendix B.3.

5.2.2 Results

Overall Performance Table 4 reports performance on the UI/UX design interpretation task.

o1 and Claude 3.5 Sonnet exhibit relatively strong post-hoc interpretation capability when the more effective variant is provided, despite having noticeably underperformed on selection in Task 1. By contrast, GPT-4o and InternVL-2.5-38B perform poorly under this setting, while Qwen2.5-VL-32B performs reasonably well. Among smaller models, Qwen2.5-VL-7B shows relatively strong alignment, whereas other 7B/8B models underperform. The consistent gap between Interpretation and Instance-level Recall scores indicates that while models can often recover parts of the expert interpretations, they rarely capture them in full.

Performance Across UX Aspects We analyze Interpretation Recall scores by UX dimension, examining which types of UX interpretations are correctly recovered, as a true understanding of UI/UX requires capturing the full range of cognitive aspects—perception, memory, and action. As shown in Figure 5, o1 leads on Perception, indicating stronger ability to explain which visual cues are likely noticed by users. Claude 3.5 Sonnet performs best on Action, suggesting greater sensitivity to interaction-oriented factors. Notably, model performance converges more on Memory than on the others. The persistent sub-70% ceiling across all dimensions underscores the need for deeper, more holistic modeling of user cognition in UI/UX understanding and interpretation.

6 Discussion

Case Study For the UI/UX design selection task, we analyze the success and failure modes of the strong-performing model, InternVL-2.5-38B, as shown in Figure 6, focusing on cases where the model makes consistent decisions under both input image orders. In Case 1, the model correctly perceives the UI change as the presence of a CTA button and selects the more effective UI/UX, providing a reason that behaviorally aligns with the user behavior pattern implied by the A/B test outcome, hypothesizing that fewer CTAs better focus user attention. In Case 2, although the model correctly identifies differences such as pricing badges and CTA text, it fails to select the winning variant identified by the A/B test, instead offering misaligned reasoning that higher visual saliency negatively impacts user behavior by inducing commitment pressure. In Case 3, the model fails at perception level, hallucinating non-existent changes instead of recognizing the actual button color difference, which

ultimately leads to an incorrect final decision.

Figure 7 shows a case from the UI/UX design interpretation task. The UI on the right was verified by A/B test results as the more effective design, with experts attributing its effectiveness to a simplified background and an enlarged demo button, along with their associated effects, factors that models must explicitly capture for an accurate interpretation. o1 succeeds by directly highlighting these key points. GPT-4o correctly identifies the relevant UI elements, but provides an incorrect or vague explanation for their impact. Qwen2.5-VL-32B correctly notes the minimalist design, but stops at identifying the button, missing the key detail of its enlargement and resulting in an incomplete explanation.

7 Conclusion

We introduce a new task that evaluates multimodal understanding of how design drives user actions in the critical field of UI/UX. To support this, we present WiserUI-Bench, a behavior-focused benchmark uniquely built from diverse, real-world A/B-tested UI variants and expert-curated key interpretations, providing a reliable basis for assessment. Our two tasks, selection and interpretation, go beyond surface-level recognition, requiring models to reason about the behavioral implications of visual design, a cognitively challenging problem. Extensive experiments across both proprietary and open-source MLLMs reveal they exhibit only shallow levels of behavior-aware visual reasoning. Achieving stronger performance on this significant task, particularly across the analyzed UI elements or UX dimensions, may require further training on large-scale, behaviorally grounded UI/UX data that reflects real-world user dynamics. We envision WiserUI-Bench and our work as a step toward multimodal systems capable of not only understanding but also generating visual interfaces grounded in user behavior contexts, aligning with human cognition to optimize user experience. More broadly, we hope to catalyze future research toward deeper UI/UX and human behavior literacy in AI.

8 Limitations

Given that user experience patterns vary across cultural contexts and societal norms, cultural biases inevitably exist in our benchmark (Karreman et al., 2016). Also, interactive UIs were not extensively explored in our work. Although the number of

300 UI image pairs may seem limited, publicly available real-world A/B test results are extremely rare, as companies rarely disclose them. Moreover, independently conducting production-level UI experiments at a comparable level and scale was practically infeasible. However, as the first benchmark in this area, we focused on reliability, diversity, and annotation quality, covering a wide range of real-world UI changes and providing expert-curated key interpretations.

9 Ethical Considerations

To ensure expert-level annotation quality, we recruited three annotators with professional UI/UX backgrounds who are currently employed at global IT companies specializing in UX and AI. We clearly informed them of the academic nature of the task, and they voluntarily participated in the study. We provided each annotator with monetary compensation equivalent to approximately \$225 USD, reflecting fair market rates for domain-specific labor, and issued payments within 48 hours of annotation completion. We instructed annotators to focus solely on task-relevant elements and ensured that they were not exposed to or asked to label any sensitive, offensive, or triggering content.

Our benchmark does not include any private, user-identifiable, or confidential information. All materials were collected from publicly released sources originally disclosed by the originating platforms or companies. Nevertheless, to avoid downstream commercial use beyond the original disclosure scope, we plan to release the benchmark under a non-commercial research license.

References

Edward Abel, Ixent Galpin, Norman Paton, and John Keane. 2020. [Pairwise comparisons or constrained optimization? a usability evaluation of techniques for eliciting decision priorities](#). *International Transactions in Operational Research*, 29.

Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. [Qwen2.5-vl technical report](#). *Preprint*, arXiv:2502.13923.

Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Sarah Bawabe, Laura Wilson, Tongyu Zhou, Ezra Marks, and Jeff Huang. 2021. [The ux factor: Using comparative peer review to evaluate designs through user preferences](#). *Proc. ACM Hum.-Comput. Interact.*, 5(CSCW2).

Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Han-nah Rashkin, Doug Downey, Scott Wen tau Yih, and Yejin Choi. 2020. [Abductive commonsense reasoning](#). *Preprint*, arXiv:1908.05739.

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang, Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo, Jiahao Wang, Tan Jiang, Bo Wang, and 23 others. 2025. [Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling](#). *Preprint*, arXiv:2412.05271.

Ziming Cheng, Binrui Xu, Lisheng Gong, Zuhe Song, Tianshuo Zhou, Shiqi Zhong, Siyu Ren, Mingxiang Chen, Xiangchao Meng, Yuxin Zhang, Yanlin Li, Lei Ren, Wei Chen, Zhiyuan Huang, Mingjie Zhan, Xiaojie Wang, and Fangxiang Feng. 2025. [Evaluating mllms with multimodal multi-image reasoning benchmark](#). *Preprint*, arXiv:2506.04280.

Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschan, Daniel Afegan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. 2017. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th annual ACM symposium on user interface software and technology*, pages 845–854.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multi-agent debate. In *Forty-first International Conference on Machine Learning*.

Peitong Duan, Chin-Yi Cheng, Gang Li, Bjoern Hartmann, and Yang Li. 2024. Uicrit: Enhancing automated design evaluation with a ui critique dataset. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pages 1–17.

Hermann Ebbinghaus, Henry A Ruger, and Clara E Bussenius. 1913. Memory: A contribution to experimental psychology.

Paul M Fitts and James R Peterson. 1964. Information capacity of discrete motor responses. *Journal of experimental psychology*, 67(2):103.

631	Brian J Fogg. 2002. Persuasive technology: using computers to change what we think and do. <i>Ubiquity</i> , 2002(December):2.	686
632		687
633		688
634	Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. 2024. <i>Mme: A comprehensive evaluation benchmark for multimodal large language models</i> . <i>Preprint</i> , arXiv:2306.13394.	689
635		690
636		691
637		692
638		693
639		694
640	Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. 2025. <i>Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark</i> . <i>Preprint</i> , arXiv:2501.05444.	695
641		696
642		697
643		698
644		699
645	William E Hick. 1952. On the rate of gain of information. <i>Quarterly Journal of experimental psychology</i> , 4(1):11–26.	700
646		701
647		702
648	Jerry R Hobbs, Mark E Stickel, Douglas E Appelt, and Paul Martin. 1993. Interpretation as abduction. <i>Artificial intelligence</i> , 63(1-2):69–142.	703
649		704
650		705
651	Clark L Hull. 1932. The goal-gradient hypothesis and maze learning. <i>Psychological review</i> , 39(1):25.	706
652		707
653	Joyce Karreman, Pietro Romeo, and Qian Li. 2016. Cross-cultural hci and ux design: a comparison of chinese and western user interfaces. <i>Unpublished work</i> .	708
654		709
655		710
656		711
657	Kurt Koffka. 1922. Perception: an introduction to the gestalt-theorie. <i>Psychological bulletin</i> , 19(10):531.	712
658		713
659	Ron Kohavi, Roger Longbotham, Dan Sommerfield, and Randal M Henne. 2009. Controlled experiments on the web: survey and practical guide. <i>Data mining and knowledge discovery</i> , 18:140–181.	714
660		715
661		716
662		717
663	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. <i>Llava-onevision: Easy visual task transfer</i> . <i>Preprint</i> , arXiv:2408.03326.	718
664		719
665		720
666		721
667		722
668	Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. 2024b. Seed-bench: Benchmarking multimodal large language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13299–13308.	723
669		724
670		725
671		726
672		727
673		728
674	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. <i>Encouraging divergent thinking in large language models through multi-agent debate</i> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 17889–17904, Miami, Florida, USA. Association for Computational Linguistics.	729
675		730
676		731
677		732
678		733
679		734
680		735
681		736
682	Chin-Yew Lin. 2004. <i>ROUGE: A package for automatic evaluation of summaries</i> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	737
683		738
684		739
685		740
		741
	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 26296–26306.	
	Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruochen Xu, and Chenguang Zhu. 2023. <i>G-eval: NLG evaluation using gpt-4 with better human alignment</i> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	
	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, and 1 others. 2024b. <i>Mmbench: Is your multi-modal model an all-around player?</i> In <i>European conference on computer vision</i> , pages 216–233. Springer.	
	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. <i>Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts</i> . In <i>International Conference on Learning Representations (ICLR)</i> .	
	Yuxuan Lu, Jing Huang, Yan Han, Bingsheng Yao, Sisong Bei, Jiri Gesi, Yaochen Xie, Zheshen Wang, Qi He, and Dakuo Wang. 2025a. <i>Prompting is not all you need! evaluating llm agent simulation methodologies with real-world online customer behavior data</i> . <i>Preprint</i> , arXiv:2503.20749.	
	Yuxuan Lu, Bingsheng Yao, Hansu Gu, Jing Huang, Jessie Wang, Yang Li, Jiri Gesi, Qi He, Toby Jia-Jun Li, and Dakuo Wang. 2025b. <i>Uxagent: A system for simulating usability testing of web design with llm agents</i> . <i>Preprint</i> , arXiv:2504.09407.	
	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, and 1 others. 2023. <i>Self-refine: Iterative refinement with self-feedback</i> . <i>Advances in Neural Information Processing Systems</i> , 36:46534–46594.	
	George A Miller. 1956. The magical number seven, plus or minus two: Some limits on our capacity for processing information. <i>Psychological review</i> , 63(2):81.	
	Jakob Nielsen. 1994. Enhancing the explanatory power of usability heuristics. In <i>Proceedings of the SIGCHI conference on Human Factors in Computing Systems</i> , pages 152–158.	
	Jakob Nielsen. 2000. End of web design. <i>Alert-box [en línea] Disponible en: www.useit.com/alertbox/20000723.html</i> , 23(06):2000.	
	Don Norman. 2007. <i>Emotional design: Why we love (or hate) everyday things</i> . Basic books.	
	Donald A Norman. 1988. <i>The psychology of everyday things</i> . Basic books.	

742	Harri Oinas-Kukkonen and Marja Harjumaa. 2018. Persuasive systems design: key issues, process model and system features 1. In <i>Routledge handbook of policy design</i> , pages 87–105. Routledge.	<i>ACM Symposium on User Interface Software and Technology</i> , pages 1–16.	798
743			799
744			
745			
746	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and 401 others. 2024a. Gpt-4o system card . <i>Preprint</i> , arXiv:2410.21276.	Weiye Xu, Jiahao Wang, Weiyun Wang, Zhe Chen, Wengang Zhou, Aijun Yang, Lewei Lu, Houqiang Li, Xiaohua Wang, Xizhou Zhu, Wenhai Wang, Jifeng Dai, and Jinguo Zhu. 2025. Visulogic: A benchmark for evaluating visual reasoning in multi-modal large language models . <i>Preprint</i> , arXiv:2504.15279.	800
747			801
748			802
749			803
750			804
751			805
752			
753	OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, and 244 others. 2024b. Openai o1 system card . <i>Preprint</i> , arXiv:2412.16720.	Bo Yang and Shanping Li. 2024. Uisgpt: Automated mobile ui design smell detection with large language models. <i>Electronics</i> , 13(16):3127.	806
754			807
755			808
756			
757			
758			
759			
760	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, and 3 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of CVPR</i> .	809
761			810
762			811
763			812
764			813
765			814
766			815
767	Vilfredo Pareto. 1919. <i>Manuale di economia politica con una introduzione alla scienza sociale</i> , volume 13. Società editrice libraria.	Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. 2024. Cocot: Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs . <i>Preprint</i> , arXiv:2401.02582.	817
768			818
769			819
770	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th annual acm symposium on user interface software and technology</i> , pages 1–22.	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert . <i>Preprint</i> , arXiv:1904.09675.	822
771			823
772			824
773			825
774			
775			
776	Federico Quin, Danny Weyns, Matthias Galster, and Camila Costa Silva. 2024. A/b testing: A systematic literature review. <i>Journal of Systems and Software</i> , 211:112011.	Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan, William Yang Wang, and Linda Ruth Petzold. 2023. Gpt-4v(ision) as a generalist evaluator for vision-language tasks . <i>Preprint</i> , arXiv:2311.01361.	826
777			827
778			828
779			829
780	Hedwig Von Restorff. 1933. Über die wirkung von bereichsbildungen im spurenfeld. <i>Psychologische Forschung</i> , 18:299–342.	Bingchen Zhao, Yongshuo Zong, Letian Zhang, and Timothy Hospedales. 2024. Benchmarking multi-image understanding in vision and language models: Perception, knowledge, reasoning, and multi-hop reasoning . <i>Preprint</i> , arXiv:2406.12742.	831
781			832
782			833
783	Dakuo Wang, Ting-Yao Hsu, Yuxuan Lu, Hansu Gu, Limeng Cui, Yaochen Xie, William Headean, Bingsheng Yao, Akash Veeragouni, Jiapeng Liu, Sreyashi Nag, and Jessie Wang. 2025a. Agenta/b: Automated and scalable web a/btesting with interactive llm agents . <i>Preprint</i> , arXiv:2504.09723.	Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibe Yang. 2023. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. <i>Advances in Neural Information Processing Systems</i> , 36:5168–5191.	836
784			837
785			838
786			839
787			840
788			
789	Ziqi Wang, Hanlin Zhang, Xiner Li, Kuan-Hao Huang, Chi Han, Shuiwang Ji, Sham M. Kakade, Hao Peng, and Heng Ji. 2025b. Eliminating position bias of language models: A mechanistic approach . <i>Preprint</i> , arXiv:2407.01100.		
790			
791			
792			
793			
794	Jason Wu, Yi-Hao Peng, Xin Yue Amanda Li, Amanda Swearngin, Jeffrey P Bigham, and Jeffrey Nichols. 2024. Uiclip: a data-driven model for assessing user interface design. In <i>Proceedings of the 37th Annual</i>		
795			
796			
797			

841 A WiserUI-Bench Details

842 A.1 Curation Details

843 **Image Preprocess** During benchmark construction process, we observed that some UI images
844 from our data sources included manual indicators,
845 such as circled regions, numbered markers, and
846 arrows, highlighting specific UI elements related to
847 reasoning points from A/B test results (see Figure
848 17). While helpful for human interpretation, these
849 indicators introduced artificial attention cues that
850 could potentially bias model evaluation.
851

852 To mitigate this, we cleaned the dataset using in-
853 painting techniques. We first identified common in-
854 dicator patterns via edge detection and color-based
855 segmentation, including circular highlights drawn
856 in semi-transparent colors, numbered markers, and
857 arrows. We then applied automated inpainting fol-
858 lowed by manual refinement to remove these ele-
859 ments while preserving UI fidelity. This process
860 ensured a cleaner, more naturalistic dataset for fair
861 visual reasoning evaluation.

862 **Interpretation Annotation** As shown in Figure
863 20, three UI/UX experts are presented with a pair of
864 UI images, one designated as the winner based on
865 A/B test results, along with a corresponding source
866 website link. Experts were instructed to provide
867 key interpretations explaining why the winning ver-
868 sion supported user actions more effectively than
869 the alternative. We engaged experts for annota-
870 tion, as user behavior is often subconscious and
871 shaped by implicit context, making evaluation by
872 non-experts inherently unreliable. To distill salient
873 behavioral drivers, we retained only interpretations
874 independently identified by at least two experts, re-
875 sulting in 684 final interpretations (562 with unan-
876 imous agreement and 122 with two-expert agree-
877 ment) from an initial pool of 845 candidates.

878 Expert annotators were initially guided by the
879 original sources’ textual description which often
880 highlights key design intentions, as well as by a
881 robust theoretical heuristic widely used in usability
882 research, the Nielsen Norman 10 Usability Heuristics
883 (Nielsen, 1994), which served as broad concep-
884 tual criteria. The resulting interpretations were then
885 consolidated and mapped to 12 curated UX laws
886 that remain widely applied in practice. These UX
887 laws provide concrete, actionable guidance across
888 diverse aspects of user experience, which can be
889 categorized into three cognitive dimensions: Per-
890 ception, Memory, and Action. This classification

891 is inspired by the framework proposed by Norman
892 (2007), which introduced the concept of *Emotional*
893 *Design*. Emotional design focuses on creating prod-
894 ucts that evoke emotions and foster positive user
895 experience by engaging users at three cognitive lev-
896 els: visceral, behavioral, and reflective. The full list
897 of UX laws and their corresponding classifications
898 is presented in Table 5.

899 A.2 Statistics Details

900 Figure 18 presents the distribution of WiserUI-
901 Bench across page types, industry domains, device
902 types. For clarification on the page type classifica-
903 tion, the *homepage* refers to the main entry point of
904 a website, typically offering broad navigation and
905 general content. In contrast, a *landing page* is usu-
906 ally a standalone page optimized for a specific user
907 action, such as signing up or making a purchase.

908 Full distribution of UI change attribute types is
909 shown in Figure 19. For clarification, element-level
910 counts reflect the number of distinct UI compo-
911 nents affected, while attribute-level counts capture
912 the total number of modified attributes across all
913 components. For example, an instance annotated
914 as <Input Field: Padding, Color / Button: Color>
915 is counted as 2 element changes and 3 attribute
916 changes.

917 B Experiment Details

918 B.1 Setup

919 B.1.1 Model Configuration

920 For proprietary MLLMs, we accessed the models
921 via their official API endpoints as follows:

- 922 • o1 (OpenAI et al., 2024b): 922
923 o1-2024-12-17 923
- 924 • GPT-4o (OpenAI et al., 2024a): 924
925 gpt-4o-2024-08-06 925
- 926 • Claude 3.5 Sonnet (Anthropic, 2024): 926
927 claude-3-5-sonnet-20240620 927

928 For open-source models, we used publicly avail-
929 able weights from the Huggingface Hub². Specific
930 model versions are as follows:

- 931 • Qwen2.5-VL-32B (Bai et al., 2025): 931
932 Qwen/Qwen2.5-VL-32B-Instruct 932
- 933 • Qwen2.5-VL-7B (Bai et al., 2025): 933
934 Qwen/Qwen2.5-VL-7B-Instruct 934

²<https://huggingface.co/models>

- 935 • InternVL-2.5-38B (Chen et al., 2025):
936 OpenGVLab/InternVL2_5-38B
- 937 • InternVL-2.5-8B (Chen et al., 2025):
938 OpenGVLab/InternVL2_5-8B
- 939 • LLaVA-NeXT-7B (Liu et al., 2024a):
940 llava-hf/llava-v1.6-mistral-7b-hf
- 941 • LLaVA-OneVision-7B (Li et al., 2024a):
942 llava-hf/llava-onevision-qwen2-7b-ov-
943 hf

944 The temperature for all models was set to 0.2.

945 B.1.2 Resources

946 All experiments using open-source models were
947 conducted on a Linux machine (Ubuntu 20.04.6
948 LTS) equipped with 4 NVIDIA RTX A6000 GPUs,
949 each with 48GB VRAM. We used the following
950 software libraries and frameworks: transformers
951 version 4.53.2 and vllm version 0.9.2.

952 B.2 UI/UX Design Selection

953 B.2.1 Prompts

954 Figures 8 through 14 show the prompts used, in-
955 cluding those designed for diverse reasoning strate-
956 gies.

957 B.2.2 Token Counts

958 Token counts across reasoning strategies on propri-
959 etary models are presented in Table 6.

960 B.2.3 Standard Deviation

961 Table 7 presents the standard deviation across three
962 experiment runs for each model and reasoning
963 method. For metric abbreviation, *FS* and *SS* de-
964 note the standard deviation when the first or second
965 UI image input was the more effective one, while
966 *AS* and *CS* refer to the standard deviation of total
967 correct and both correct cases, respectively.

968 B.2.4 Additional Case Study

969 Figure 21 and 22 present the additional case stud-
970 ies across different models. Figure 23 shows an
971 example comparing various reasoning strategies
972 executed with GPT-4o.

973 B.3 UI/UX Design Interpretation

974 B.3.1 Prompts

975 Figure 15 presents the prompt used for models’
976 interpretation generation, and Figure 16 presents
977 the prompt used for its GPT-4o-based evaluation.

978 B.3.2 Results on Text Similarity Metrics

979 Alongside Interpretation Recall and Instance-level
980 Recall from GPT-4o-based evaluation, we report
981 BLEU, ROUGE, METEOR, and BERTScore as
982 well. (Papineni et al., 2002; Lin, 2004; Banerjee
983 and Lavie, 2005; Zhang et al., 2020) in Table 8.

984 B.3.3 Additional Case Study

985 Figure 24 presents an additional case study across
986 different models.

987 C Icon Attribution

988 Icons used in this paper are from Flaticon
989 (<https://www.flaticon.com>), and they are attributed
990 to the respective authors as required by Flaticon’s
991 license.

UX Dimension	UX Law	Description
Perception	Law of Common Region (Koffka, 1922)	Elements sharing a boundary or background are perceived as a unified group. This can be a visible border, a colored box, or a distinct region in the layout.
	Law of Proximity (Koffka, 1922)	Items positioned close together are seen as belonging to the same group. This mental shortcut helps users quickly identify relationships between adjacent elements.
	Law of Prägnanz (Koffka, 1922)	Humans naturally perceive and interpret complex images in their simplest forms. It reduces cognitive effort by helping us group and understand elements more quickly.
	Law of Similarity (Koffka, 1922)	People group together elements that look alike, assuming they serve a similar purpose. Consistency in color, shape, or size leads the user to see these items as related.
Memory	Von Restorff Effect (Von Restorff, 1933)	People remember a distinct element more easily than others that blend in. This is also known as the Isolation Effect, emphasizing our attention on anything unusual.
	Serial Position Effect (Ebbinghaus et al., 1913)	Users tend to remember the first and last items in a series more than those in the middle. This influences how they recall sequential information or actions.
	Miller's Law (Miller, 1956)	On average, people can hold only about seven items in their working memory at once. Going beyond this limit causes cognitive overload.
	Jakob's Law (Nielsen, 2000)	Users expect new interfaces to operate similarly to those they've used before. Leveraging established patterns lowers cognitive load and accelerates adoption.
Action	Hick's Law (Hick, 1952)	Decision time rises with the number and complexity of choices. Too many options can overwhelm users and slow their actions.
	Pareto Principle (Pareto, 1919)	Roughly 80% of outcomes stem from 20% of causes. Focusing on the most impactful features or tasks delivers the greatest improvement.
	Fitts's Law (Fitts and Peterson, 1964)	The time required to reach a target depends on its size and distance. Larger, closer elements are easier to access quickly and accurately.
	Goal-Gradient Effect (Hull, 1932)	Motivation increases as people get closer to finishing a task. Providing clear progress markers encourages users to continue and complete it faster.

Table 5: Overview of UX laws used for interpretation annotation, and their categorization by cognitive dimension (Perception, Memory, Action).

Models	Methods	Input Tokens	Output Tokens
o1	Zero-Shot	1292.5 (\pm 559.77)	1347.7 (\pm 489.76)
GPT-4o	Zero-Shot	1453.4 (\pm 634.41)	248.8 (\pm 57.98)
	CoCoT	1458.4 (\pm 634.41)	329.7 (\pm 63.29)
	Self-Refine	5240.4 (\pm 1935.80)	820.9 (\pm 142.43)
	DDCoT	1521.4 (\pm 634.41)	415.3 (\pm 73.42)
	MAD (R1)	4942.4 (\pm 1903.15)	483.1 (\pm 58.5)
	MAD (R3)	11501.5 (\pm 4435.48)	1158.1 (\pm 82.60)
Claude 3.5 Sonnet	Zero-Shot	1686.6 (\pm 1062.43)	401.4 (\pm 72.94)
	CoCoT	1692.5 (\pm 1062.58)	487.0 (\pm 82.28)
	Self-Refine	6230.7 (\pm 3245.69)	1201.4 (\pm 207.32)
	DDCoT	1754.5 (\pm 1064.29)	582.2 (\pm 97.13)
	MAD (R1)	5895.5 (\pm 3224.31)	834.1 (\pm 125.13)
	MAD (R3)	13853.4 (\pm 7512.63)	1936.5 (\pm 276.14)

Table 6: Overview of Token Counts ($\mu \pm \sigma$) in the UI/UX design selection task on WisierUI-Bench.

Model	Method	FS	SS	AS	CS
o1	Zero-Shot	1.528	2.082	3.606	0.577
GPT-4o	Zero-Shot	1.528	1.000	2.082	1.528
	CoCoT	4.000	2.646	1.732	1.732
	Self-Refine	2.517	1.000	3.215	0.577
	DDCoT	2.082	2.309	0.577	2.517
	MAD (R1)	3.215	3.055	6.245	2.646
	MAD (R3)	0.577	2.082	1.732	3.215
Claude 3.5 Sonnet	Zero-Shot	1.528	2.000	3.512	2.887
	CoCoT	3.055	3.512	0.577	0.577
	Self-Refine	2.646	3.512	4.619	1.528
	DDCoT	2.082	1.732	0.577	2.517
	MAD (R1)	1.000	1.528	2.309	4.359
	MAD (R3)	1.000	2.887	3.786	3.786
Qwen2.5-VL-32B	Zero-Shot	1.732	2.309	2.082	2.000
Qwen2.5-VL-7B	Zero-Shot	2.646	0.577	2.517	3.215
InternVL-2.5-38B	Zero-Shot	2.646	1.000	3.000	3.055
InternVL-2.5-8B	Zero-Shot	1.000	2.646	2.646	2.000
LLaVA-NeXT-7B	Zero-Shot	3.512	1.155	4.583	6.028
LLaVA-OneVision-7B	Zero-Shot	4.041	2.309	4.359	3.512

Table 7: Results on standard deviation across three experiment runs in the UI/UX design selection task for each model and reasoning strategy on WisierUI-Bench.

Model	BLEU-2	BLEU-4	ROUGE-L	METEOR	BERTScore
o1	0.0181	0.0041	0.0972	0.1630	0.8414
GPT-4o	0.0162	0.0040	0.0875	0.1557	0.8233
Claude 3.5 Sonnet	0.0141	0.0037	0.0654	0.1418	0.8076
Qwen2.5-VL-32B	0.0098	0.0025	0.0442	0.1016	0.8085
Qwen2.5-VL-7B	0.0180	0.0048	0.0789	0.1481	0.8331
InternVL-2.5-38B	0.0128	0.0034	0.0574	0.1225	0.8080
InternVL-2.5-8B	0.0111	0.0032	0.0510	0.1124	0.8082
LLaVA-NeXT-7B	0.0077	0.0018	0.0386	0.0908	0.8127
LLaVA-OneVision-7B	0.0157	0.0043	0.0718	0.1416	0.8255

Table 8: Results of text similarity metrics per expert interpretation for the UI/UX design interpretation task on WiserUI-Bench.

You are an expert in designing UI/UX for web/apps.

The two screenshots show two different versions of the same page.
Identify the key UI differences between the two versions, and then evaluate which variant is more effective UI/UX design that leads to better user experience and conversion.

You should end your answer with following the format (No bold, etc):
More effective: <First/Second>

Figure 8: Prompt used for the UI/UX design selection task with zero-shot strategy.

You are an expert in designing UI/UX for web/apps.

The two screenshots show two different versions of the same page.
Identify the key UI differences between the two versions, and then evaluate which variant is more effective UI/UX design that leads to better user experience and conversion.

Think step by step.

You should end your answer with the following format (No bold, etc):
More effective: <First/Second>

Figure 9: Prompt used for the UI/UX design selection task with CoCoT strategy.

You are an expert in designing UI/UX for web/apps.

The two screenshots show two different versions of the same page.
We want to identify the key UI differences between the two versions, and then evaluate which variant is more effective UI/UX design that leads to better user experience and conversion.

Review your previous answer and find problems with your answer.
Previous answer:
{previous_answer}

Figure 10: Prompt used for the UI/UX design selection task with Self-Refine strategy on the reviewing step. Prompt inputs are in **bold**.

You are an expert in designing UI/UX for web/apps.

The two screenshots show two different versions of the same page.
We want to identify the key UI differences between the two versions, and then evaluate which variant is more effective UI/UX design that leads to better user experience and conversion.

Your previous answer:
{previous_answer}

Feedback on your previous answer:
{feedback}

Based on the feedback for your previous answer, improve your answer.

You should end your answer with following the format (No bold, etc):
More effective: <First/Second>

Figure 11: Prompt used for the UI/UX design selection task with Self-Refine strategy on the improving step. Prompt inputs are in **bold**.

You are an expert in designing UI/UX for web/apps.

The two screenshots show two different versions of the same page.
We want to identify the key UI differences between the two versions and then evaluate which variant is more effective UI/UX design that leads to better user experience and conversion.

You should

- (1) Think step by step about the preliminary knowledge to answer the question, deconstruct the problem as completely as possible down to necessary sub-questions.
- (2) With the aim of helping answer the original question, try to answer the sub-questions.
- (3) Give your final answer according to the sub-questions and sub-answers.

You should end your answer with the following format (No bold, etc):
More effective: <First/Second>

Figure 12: Prompt used for the UI/UX design selection task with DDCoT strategy.

You are an expert in designing UI/UX for web/apps.

And you are a debater. Hello and welcome to the debate.

The two screenshots show two different versions of the same page.
Identifying the key UI differences between the two versions, you think the **{first/second}** variant is more effective UI/UX design that leads to better user experience and conversion.

It's not necessary to fully agree with each other's perspectives, as our objective is to find the correct answer.

The opponent proposed that:
{opponent_opinion}

Figure 13: Prompt used for the UI/UX design selection task with MAD strategy on the debate step arguing the first/second version better drives user actions. Prompt inputs are in **bold**.

You are an expert in designing UI/UX for web/apps.

The two screenshots show two different versions of the same page.

If the **SECOND** version has proven to be more effective UI/UX design that leads to better user experience and conversion, provide all possible and appropriate interpretations for its superiority.

Figure 15: Prompt used for the UI/UX design interpretation task.

You are an expert in designing UI/UX for web/apps.

And you are a moderator. There will be two debaters involved in a debate.

The two screenshots show two different versions of the same page.

They are discussing which variant is more effective UI/UX design that leads to better user experience and conversion.

At the end of each debate round, you will evaluate answers and decide which is correct.

Now the **{number}** round of debate for both sides has ended.

Argue from the side arguing for the first version:
{first_reason}

Argue from the side arguing for the second version:
{second_reason}

Then what is the correct answer?
Please give the final answer that you think is correct and summarize your reasons.

You should end your answer with following the format (No bold, etc):
More effective: <First/Second>

Figure 14: Prompt used for the UI/UX design selection task with MAD strategy on the moderation step. Prompt inputs are in **bold**.

Given two statements are interpretations why one UI is better than another in terms of guiding user behavior effectively.

Statement 1:
{model_reason}

Statement 2:
{reference_reason}

Your task is to assess whether the first statement contains the same reasoning and essential content as the second statement.

If it does, answer with 'Yes'. Otherwise, answer with 'No'.

Figure 16: Prompt used for GPT-4o-based evaluation on the UI/UX design interpretation task. Prompt inputs are in **bold**.

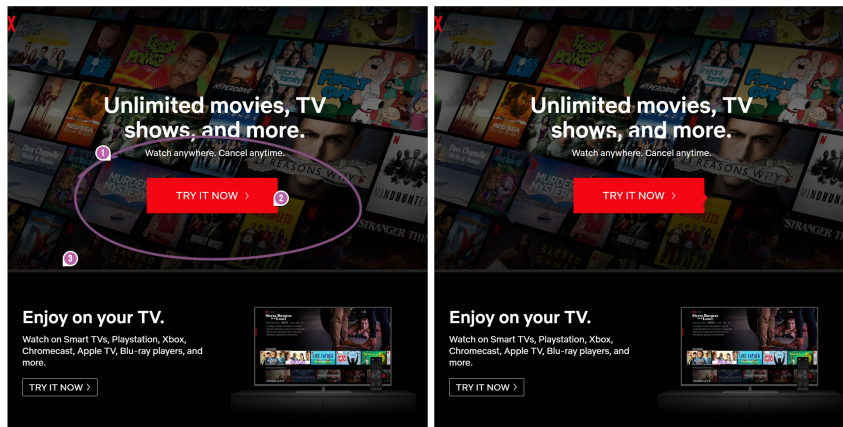


Figure 17: Example of visual indicator removal in WiserUI-Bench. Left: Original image with markers added by data source website. Right: After inpainting, with markers removed while preserving UI elements.

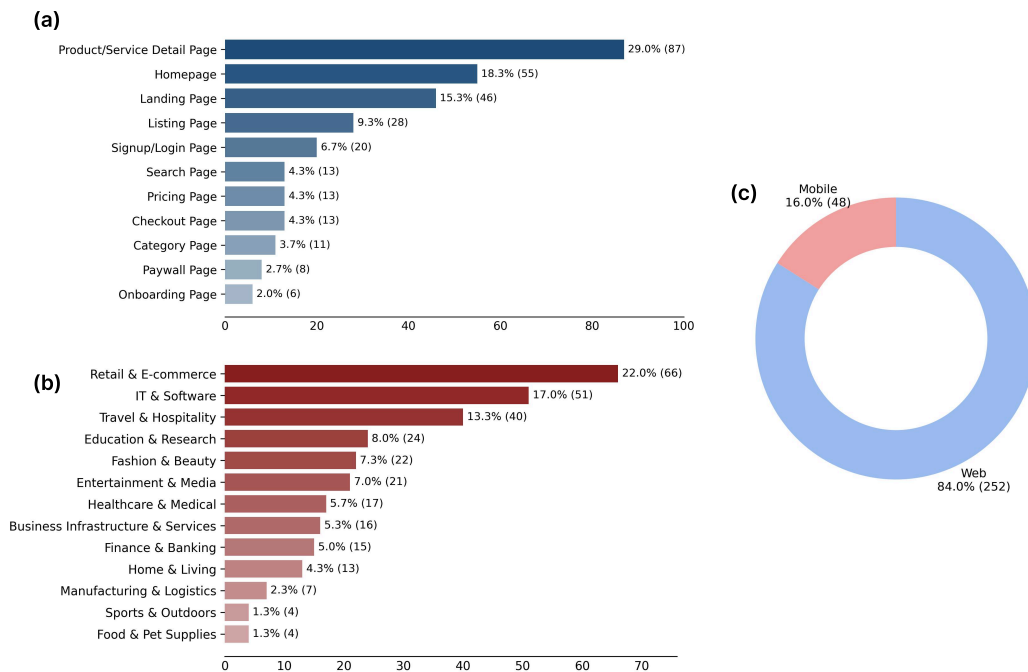


Figure 18: Detailed statistics of WiserUI-Bench on (a) page type, (b) industry domain, (c) device type

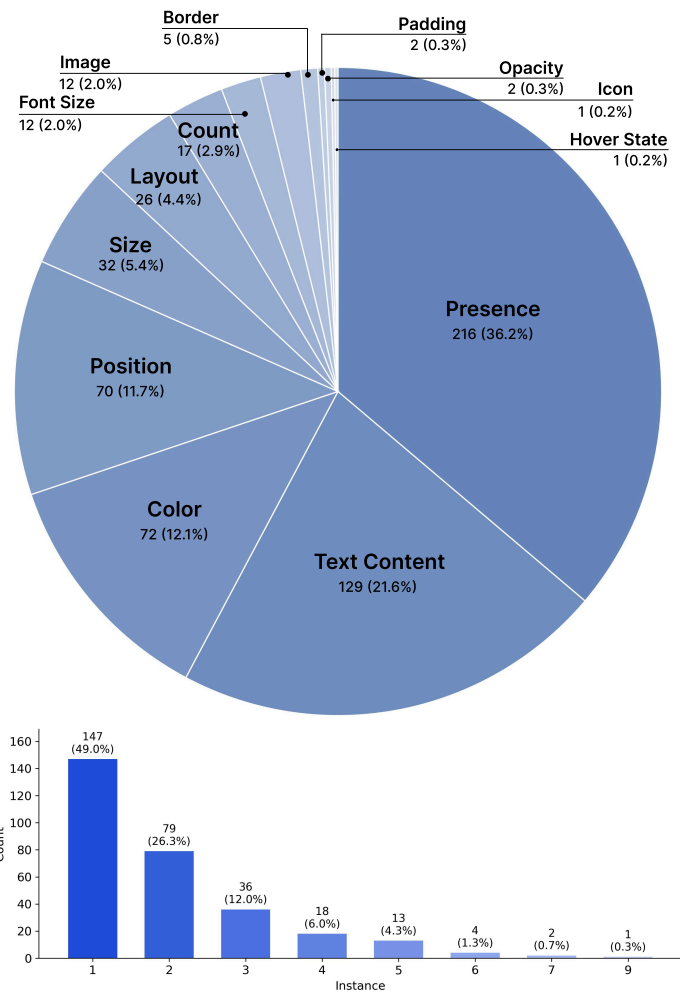


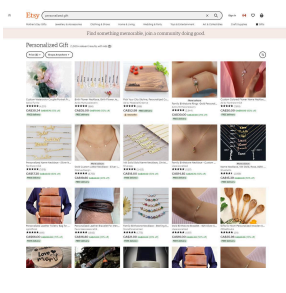
Figure 19: Detailed statistics of WiserUI-Bench on UI change attribute type

Your Objective

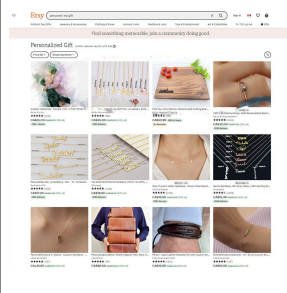
Review the pair of website images along with the information provided in the website link, and write key rationales explaining why the winning design more effectively guided user behavior, based on the provided criteria.

A/B Test Winner

A



B



[Click here to access source website](#)

Criteria

Nielsen Norman 10 Usability Heuristics

1. Visibility of System Status
2. Match Between the System and the Real World
3. User Control and Freedom
4. Consistency and Standards

⋮

12 Laws of UX (Practical)

1. Law of Common Region
Elements sharing a boundary or background are perceived as a unified group. This can be a visible border, a colored box, or a distinct
2. Law of Proximity
Items positioned close together are seen as belonging to the same group. This mental shortcut helps users quickly identify

⋮

Key Interpretations

1. The larger tiles in the right version stand out more, making them easier to remember and likely to attract user attention.

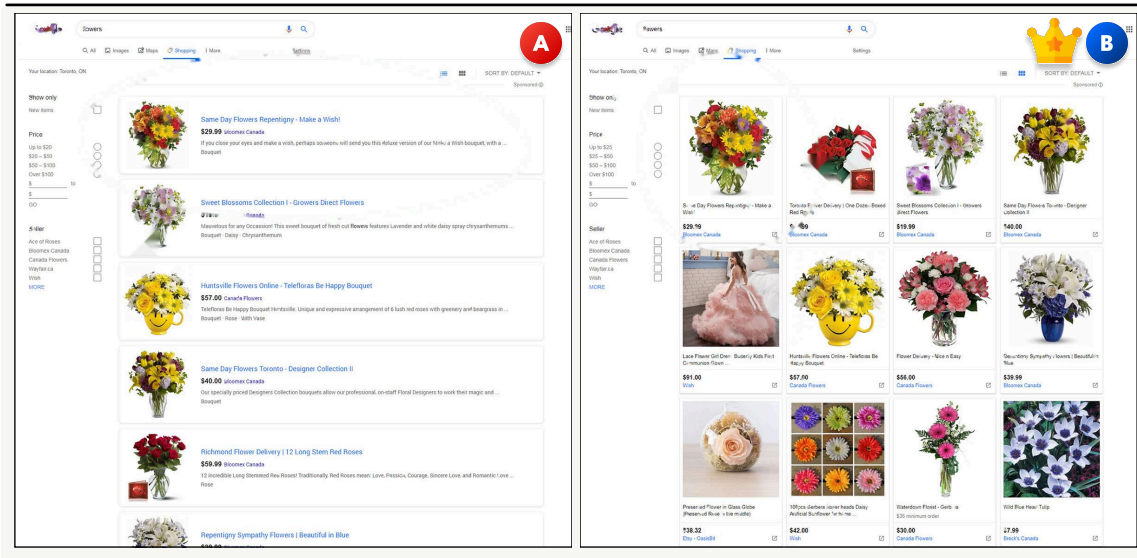
1 2 3 4 5 6 7 8 9 10 11 12

2. The larger tile size in the right version makes it easier to click, reducing the time and effort required for interaction.

1 2 3 4 5 6 7 8 9 10 11 12

[Save & Submit](#)

Figure 20: A screenshot of the interpretation annotation interface for UI/UX experts based on original data sources and guidelines. After writing the key interpretations, they also selected the most relevant UX law.



Claude 3.5 Sonnet

The key UI differences between the two versions are:

1. Layout: The first version uses a list view, while the second version uses a grid view.
2. Product cards: The second version has more compact product cards with less text description visible.
3. Image size: The second version features larger product images.
4. Information density: The second version displays more products on the screen at once.
5. Share button: The second version includes a share icon for each product.
6. View options: The second version offers both list and grid view options at the top right.
7. Product variety: The second version shows a wider variety of products, including non-flower items like dresses.

Evaluation: The second version is more effective in terms of UI/UX design and likely to lead to better user experience and conversion for the following reasons:

1. Visual appeal: The grid layout with larger images is more visually engaging and allows users to see more products at a glance.
2. Efficient browsing: Users can scan through more options quickly, potentially finding what they're looking for faster. [...]
7. Reduced cognitive load: Less text on each product card allows users to focus on the visual aspects, which are crucial for flower selection. While the first version provides more detailed information upfront, the second version's design is more likely to keep users engaged and browsing, potentially leading to higher conversion rates.

More effective: Second

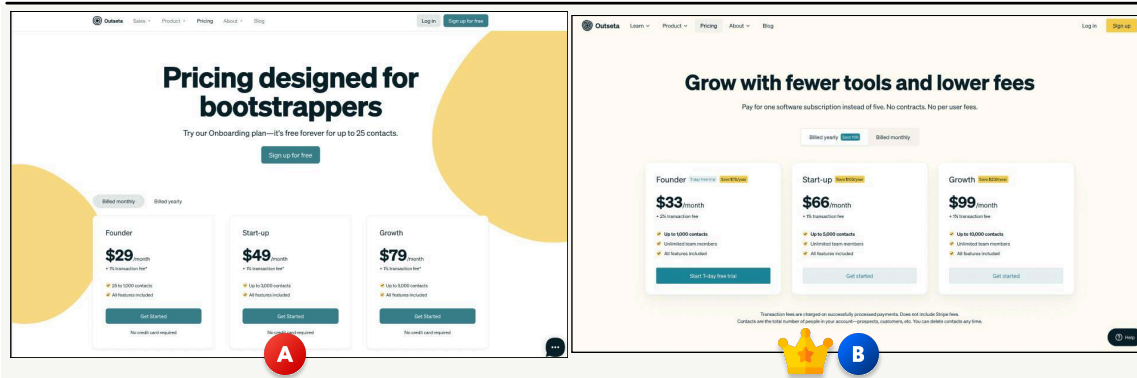
InternVL-2.5^{38B}

The second version of the page appears to be more effective in terms of UI/UX design, leading to a better user experience and conversion. Here are some reasons why:

1. Simplified Layout: - The second version has a cleaner and more streamlined layout. The removal of the "Refine Your Search" section and the sorting options simplifies the page, making it less cluttered and easier to navigate.
2. Prominent Product Display: - The products are more prominently displayed in the second version. The featured products are larger and more visible, which can attract users' attention more effectively.
3. Clear Call-to-Action: - The "More Details" buttons are more prominent and consistent in the second version. This makes it easier for users to find and click on them, leading to more engagement with the products. [...]
7. Better Mobile Responsiveness: - The simplified layout in the second version is likely more responsive on mobile devices, providing a better experience for users accessing the site on smaller screens.

Overall, the second version's design improvements contribute to a more user-friendly and effective shopping experience, which can lead to higher engagement and conversion rates.

Figure 21: A qualitative example from the UI/UX design selection task on WiserUI-Bench: Both Claude 3.5 Sonnet and InternVL-2.5-38B selected the A/B test winner.



GPT-4o

The two versions of the page have distinct differences in their UI design:

1. Color and Aesthetics: - The first version has a simple white background with a purple header, giving it a clean and minimalistic look. - The second version uses a gradient background with more vibrant colors and highlights, making it visually richer. 2. Layout and Structure: - The first version has a straightforward list format indicating key features like total investment, frequency, and next date, with clear labels and values. - The second version is organized into sections with boxes, providing a different hierarchy with clear separation for different types of information (deposits and investments). 3. Information Density: - The first version has less information density, focusing only on a few key settings. - The second version provides more detailed options and information, allowing for additional interaction (like setting deposits and investments). 4. Interactivity and Options: - The first version has a simpler interface with fewer interactive buttons. - The second version includes interactive buttons for deposits and investments, offering more direct actions on this screen. 5. Clarity and Focus: - The first version focuses on clarity and simplicity, potentially making it easier for fast comprehension without distraction. - The second version may engage users more with its richness but could also overwhelm some users with additional options and visuals.

Evaluation:- If the goal is simplicity, ease of use, and focusing on key features, the first version is more effective.- If providing more options and a vibrant interface is preferred to keep users engaged, the second version is more effective.

More effective: Second

LLaVA-NeXT^{7B}

Based on the provided screenshots, here are the key UI differences between the two versions:

1. Color Scheme: The first version uses a darker color scheme, while the second version uses a lighter one. The darker version may be more accessible for users with visual impairments, but the lighter version can appear more modern and clean. 2. Button Styling: The first version has more prominent and contrasting buttons, which can make them more noticeable and easier to tap on a touchscreen. The second version has buttons that blend more with the background, which might be less effective for touch targets. 3. Typography: The first version uses a larger font size for the main headline, which can make the message more prominent and easier to read. The second version uses a smaller font size for the main headline, which might be less effective for users with visual impairments. 4. Spacing and Layout: The first version has more spacing between elements, which can make the page feel less cluttered and easier to navigate. The second version has less spacing, which might make the page feel more cramped and harder to read. 5. Call to Action (CTA): The first version has a more prominent CTA button ("Get Started") with a contrasting color, which can make it more enticing for users to click. The second version has a less prominent CTA button with a color that blends more with the background, which might make it less effective at drawing attention. [...]

10. Responsiveness: The first version appears to be more responsive, with the layout adjusting well to different screen sizes. The second version does not showcase this feature, so it's unclear how well it adapts to different devices. In terms of which variant is more effective UI/UX design, the first version appears to be more effective due to its clearer information hierarchy, more prominent CTA, and overall better use of whitespace and design consistency. These factors can contribute to a better user experience and potentially higher conversion rates.

More effective: First version

Figure 22: A qualitative example from the UI/UX design selection task on WiserUI-Bench: GPT-4o selected the A/B test winner, whereas LLaVA-NeXT-7B did not.



Zero-Shot The screenshots show two versions of a webpage with the following differences:

1. Call to Action (CTA) Button Text: - First: "START FREE TRIAL" - Second: "START TODAY"
2. Price Information: - First: No price information displayed. - Second: Displays "\$75 PER SEAT" prominently.
3. Form Requirement Note: - First: Includes "(no credit card required)" next to the email field. - Second: No note regarding credit card requirement.

Evaluation: 1. CTA Effectiveness: - The first version with "START FREE TRIAL" may attract more users interested in trying the service without commitment, likely leading to higher conversion rates for initial user acquisition. 2. Pricing Visibility: - The second version provides clear pricing information ("\$75 PER SEAT"), which could deter or attract users based on their pricing expectations. Transparency in pricing is beneficial but might discourage users not ready to pay. [...] Overall, the first version is likely more effective in encouraging initial engagement due to the emphasis on a free trial and no credit card requirement, potentially leading to better user experience and initial conversion rates.

More effective: First

DDCoT To determine which variant offers a more effective UI/UX design, consider the following sub-questions:

1. Visual Hierarchy and Clarity: - How are elements arranged, and do they guide the user effectively?
2. Call-to-Action (CTA) Presence and Clarity: - Are the CTAs clear and motivating? Is their purpose obvious?
3. Messaging and Information: - Is the essential information presented clearly and concisely?
4. User Engagement: - Which design might users find easier and more compelling to interact with?

Sub-Question Analysis: 1. Visual Hierarchy and Clarity: - First: Clear division between benefits and contact form, good use of whitespace. - Second: Similar layout, but inclusion of pricing information may distract from the primary CTA. 2. Call-to-Action (CTA) Presence and Clarity: - First: CTA is "Start Free Trial" with no cost barrier, presumably more inviting. - Second: CTA "Start Today" presents cost (\$75 per seat), which could deter some users. 3. Messaging and Information: - First: Highlights free trial prominently, focusing on user benefits. - Second: Adds immediate pricing info, which could be informative but also distracting or off-putting. 4. User Engagement: - First: Free trial offer encourages users to engage without risk. - Second: Pricing can shift focus or create hesitation.

Final Decision: Taking these considerations into account, the first variant, with its emphasis on a free trial, is likely more effective for drawing initial user interest and engagement without creating cost barriers.

More effective: First

MAD (R1) In comparing the two screenshots, I advocate for the first version as a more effective UI/UX design for several reasons. The absence of the "\$75 per seat" pricing in the first version reduces cognitive load for users who might be deterred by cost information too early in the decision process. This version can lead to better initial engagement as users are not immediately discouraged by financial details, which can enhance conversion rates. Additionally, the call to action "START FREE TRIAL" presents a low-risk opportunity, [...]

In the comparison of these two versions, the second variant stands out as a more effective UI/UX design. Displaying the "\$75 per seat" pricing upfront, it aligns with transparency best practices, allowing users to make informed decisions early on. This approach attracts users who are genuinely interested and willing to invest, thereby improving the quality of leads. The clarity offered by upfront pricing reduces surprises later in the user journey, enhancing trust and reducing drop-off rates. The call to action, "START TODAY," suggests immediacy and commitment, which can effectively spur decisive users into action, [...]

Both arguments present valid points regarding their respective designs. [...] Ultimately, the choice between the two depends on the specific goals of the business and its target audience. If the intention is to cast a wide net and engage many users initially, the first version is more effective. On the other hand, if the priority is to attract quality leads who are ready for financial commitment, the second version may be more effective. Considering general practices for enhancing trust and transparency in UI/UX design, and fostering immediate, committed engagement:

More effective: Second

Figure 23: A qualitative example from the UI/UX design selection task on WisierUI-Bench, showcasing the effect of diverse reasoning strategies executed on GPT-4o.

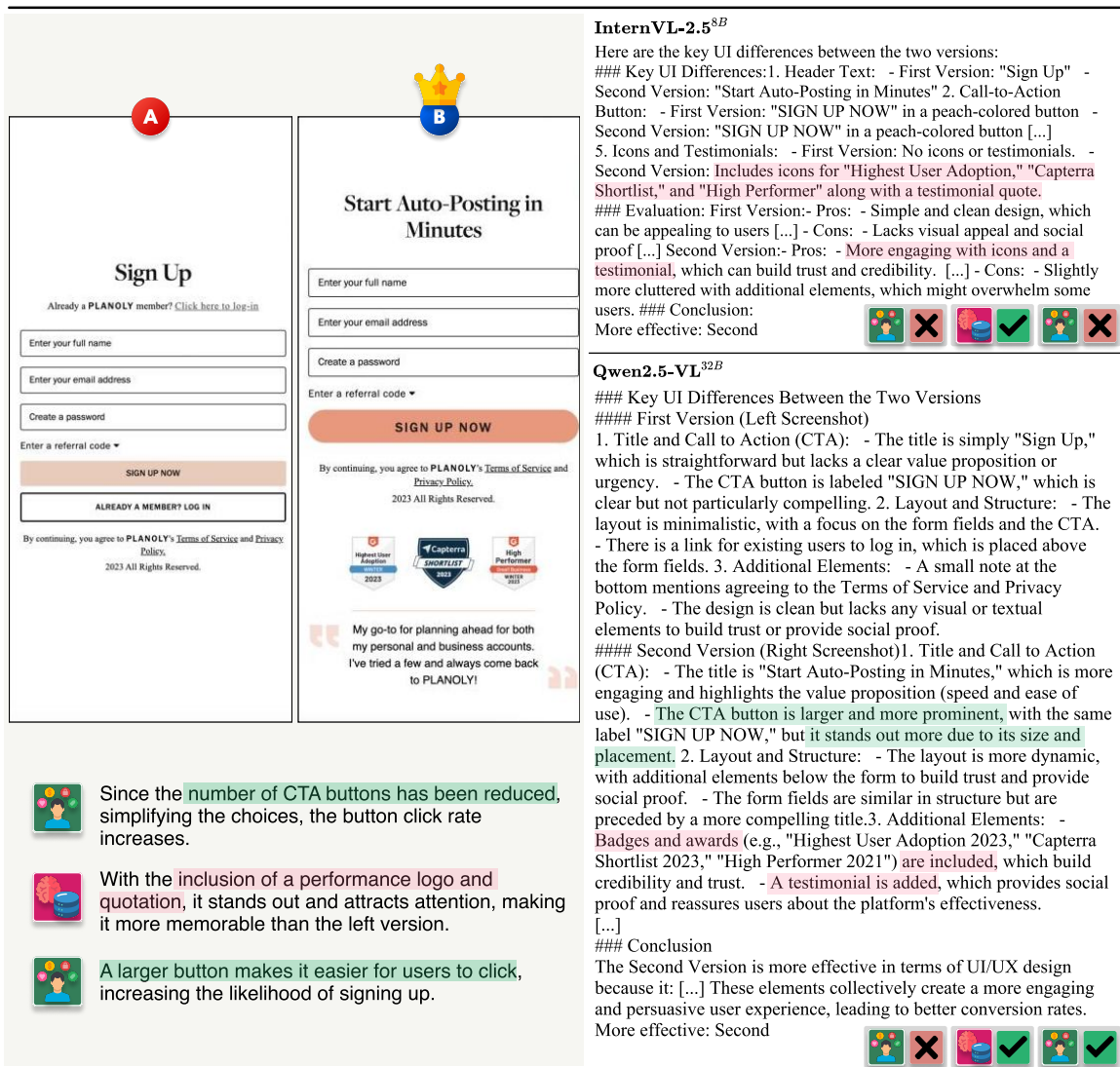


Figure 24: A qualitative example from the UI/UX design interpretation task on WisersUI-Bench: InternVL-2.5-8B captured only the memory part, while Qwen2.5-VL-32B identified two key interpretations. The highlights in the models' responses indicate the parts that align with the expert-curated key interpretations.