

Web UI Testing

Shamit Fatin*

January 26, 2025

1 Related Works

The automation of web user interface (UI) tasks and testing has garnered significant attention in recent years due to the growing complexity of modern web applications. Single Page Applications (SPAs), dynamic content, and responsive designs have introduced challenges for traditional approaches to web automation and testing. Addressing these challenges requires advancements in benchmarks, agent-based systems, and frameworks that leverage emerging technologies, such as large language models (LLMs), to enhance adaptability and scalability. However, despite significant progress, existing approaches often fall short in addressing critical aspects, leaving room for further exploration and improvement.

Benchmarks have been pivotal in evaluating the performance of web agents and testing frameworks by providing standardized datasets for assessing their capabilities. For instance, the BrowserGym framework integrates multiple benchmarks, including MiniWoB++ [8] and VisualWebArena [6], to enable a unified evaluation of agents across diverse web interaction tasks [2]. While this framework provides valuable tools for debugging and reproducibility, its reliance on simulated environments limits its applicability to real-world scenarios. The Mind2Web dataset addresses some of these shortcomings by introducing over 2,000 tasks across 137 real-world websites spanning 31 domains [3]. However, the dataset primarily focuses on static web environments, neglecting the dynamic interactions and evolving web technologies that are increasingly prevalent. Additionally, both benchmarks face challenges in providing sufficient diversity in tasks, which restricts their ability to test generalization effectively. This lack of comprehensiveness in benchmarks leaves open the need for datasets that better capture the complexities of modern web applications, including dynamic state management, real-time updates, and cross-device interactions.

Recent advancements in agent-based systems have leveraged LLMs and multimodal approaches to improve web navigation and interaction. Agents such as those described in the SeeAct framework utilize GPT-4V to process visual and structural inputs, such as screenshots and HTML, to generate action plans for task completion [11]. While these agents exhibit promising capabilities in multimodal environments, they struggle with grounding their actions to precise UI elements, often leading to errors in complex or ambiguous scenarios. Multimodal web agents, such as those combining Vision Transformers (ViT) with Flan-T5, have achieved state-of-the-art performance in benchmarks like MiniWoB++ and Mind2Web by integrating visual and structural representations [4]. However, these systems rely heavily on benchmark-specific optimizations, raising concerns about their ability to generalize to unseen tasks and websites. Similarly, the WebDreameer framework introduces a model-based planning approach, where LLMs simulate and evaluate potential actions before execution, resulting in safer and more effective navigation [5]. Despite these advancements, these systems often exhibit limitations in their long-term planning capabilities and their ability to recover from errors in cascading task sequences. Furthermore, the computational cost of multimodal models and their reliance on extensive fine-tuning for specific tasks highlight significant barriers to their practical deployment.

*BUET, CSE, Dhaka, Bangladesh. Email: shamit187@gmail.com

Testing frameworks have also seen significant innovation, with researchers exploring the use of LLMs to enhance the robustness and accessibility of automated UI testing. One approach introduces a pseudo-language for test creation, enabling users to describe actions in natural language, which are then translated into executable test scripts [1]. While this method lowers the barrier for non-technical users, its dependence on pseudo-languages introduces an additional layer of abstraction, which may not capture the nuances of complex UI behaviors. Similarly, the framework’s resilient locators, while adaptive, are limited in their effectiveness when faced with significant UI redesigns or dynamic elements. Another notable contribution leverages LLMs for automating test case generation, error detection, and A/B comparisons, demonstrating the scalability and cost-effectiveness of AI-driven testing methods [9]. However, the reliance on historical datasets for training these models creates challenges when encountering novel or unanticipated UI behaviors. Test-Agent, a multimodal framework, combines LLMs with computer vision and reinforcement learning to automate mobile app testing, showcasing adaptability across platforms [7]. Despite its efficiency, the framework heavily depends on the accuracy of visual perception modules, which can be unreliable in complex or visually cluttered interfaces. Additionally, hybrid approaches integrating traditional methods with LLMs for repairing broken tests have shown promise in improving the accuracy and reliability of test maintenance [10]. However, these methods often involve significant manual intervention for verification, undermining their promise of full automation.

The integration of benchmarks, agent-based systems, and testing frameworks reveals a complementary relationship that has the potential to address the challenges of modern web automation. While benchmarks provide the foundation for evaluating agents and frameworks, they often fail to capture the real-world complexities of dynamic and evolving web environments. Similarly, while agent-based systems and testing frameworks bring advanced capabilities, they are constrained by limited generalization, scalability challenges, and reliance on static datasets. These gaps highlight significant opportunities for future research, including the development of benchmarks that capture real-world dynamism, agents capable of fine-grained grounding and recovery, and testing frameworks that can seamlessly integrate with adaptive agents. Addressing these challenges will be critical to advancing the state-of-the-art in web UI automation and testing.

References

- [1] Maroun Ayli, Youssef Bakouny, Nader Jalloul, and Rima Kilany. Enhancing the resiliency of automated web tests with natural language. In *Proceedings of the 2024 4th International Conference on Artificial Intelligence, Automation and Algorithms*, pages 63–69, 2024.
- [2] De Chezelles, Thibault Le Sellier, Maxime Gasse, Alexandre Lacoste, Alexandre Drouin, Massimo Caccia, Léo Boisvert, Megh Thakkar, Tom Marty, Rim Assouel, et al. The browsergym ecosystem for web agent research. *arXiv preprint arXiv:2412.05467*, 2024.
- [3] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36, 2024.
- [4] Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models. *arXiv preprint arXiv:2305.11854*, 2023.
- [5] Yu Gu, Boyuan Zheng, Boyu Gou, Kai Zhang, Cheng Chang, Sanjari Srivastava, Yanan Xie, Peng Qi, Huan Sun, and Yu Su. Is your llm secretly a world model of the internet? model-based planning for web agents. *arXiv preprint arXiv:2411.06559*, 2024.
- [6] Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*, 2024.

- [7] Youwei Li, Yangyang Li, and Yangzhao Yang. Test-agent: A multimodal app automation testing framework based on the large language model. In *2024 IEEE 4th International Conference on Digital Twins and Parallel Intelligence (DTPI)*, pages 609–614. IEEE, 2024.
- [8] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. In *International Conference on Learning Representations (ICLR)*, 2018.
- [9] Fei Wang, Kamakshi Kodur, Michael Micheletti, Shu-Wei Cheng, Yogalakshmi Sadasivam, Yue Hu, and Zening Li. Large language model driven automated software application testing. *Technical Disclosure Commons*, (March 26, 2024) https://www.tdcommons.org/dpubs_series/6815, 2024.
- [10] Zhuolin Xu, Qiushi Li, and Shin Hwei Tan. Guiding chatgpt to fix web ui tests via explanation-consistency checking. *arXiv preprint arXiv:2312.05778*, 2023.
- [11] Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*, 2024.