

WebMind

¹nd Prayas Raj

*Department of Computer Science and Engineering
Indian Institute of Information Technology, Raichur
Raichur, India
cs21b1022@iiitr.ac.in*

²st Shivanshu Gupta

*Department of Computer Science and Engineering
Indian Institute of Information Technology, Raichur
Raichur, India
shivanshugupta@gmail.com*

Abstract—With the rapid advancements in artificial intelligence, AI models have developed strong reasoning and decision-making abilities. However, their capability to execute real-world tasks autonomously remains limited. This project aims to bridge this gap by enabling AI models to interact with and control web-based applications on behalf of users. By leveraging AI’s reasoning power, the system can perform browser-based tasks such as navigating websites, extracting information, and executing user-defined operations. This approach enhances automation, reduces human effort, and expands the practical applications of AI in everyday digital interactions.

Index Terms—AI automation, web interaction, browser control, task execution, intelligent agents.

I. INTRODUCTION

The rapid advancements in artificial intelligence have significantly improved its ability to reason and make decisions. However, AI still lacks the capability to autonomously execute tasks within digital environments, especially for everyday browser-based interactions. Our project aims to bridge this gap by developing a system that enables AI models to control web applications and perform tasks on behalf of users. This will allow users to automate repetitive actions, enhance productivity, and simplify complex workflows.

A. Impact and Benefits

By allowing AI to interact with and manipulate web-based user interfaces, this system has the potential to significantly reduce manual effort. Users can delegate routine tasks such as information retrieval, form filling, and online transactions to the AI, freeing up time for more critical decision-making. This solution also holds value for individuals with accessibility needs, helping them navigate digital platforms more efficiently.

B. Existing Research and Comparative Analysis

To develop an effective AI-driven automation system, we will analyze and discuss existing research papers that focus on key tasks such as object detection, UI element recognition, and relationship mapping between interface components. We will compare various models based on their accuracy, efficiency, and applicability to real-world scenarios.

C. Proposed Solution

Building on our research, we will design and implement our own AI-driven automation system. This will involve a structured approach to UI element detection, relationship

recognition, and task execution. We will detail the internal workings of our solution, including its architecture, model selection, and optimization strategies.

D. Performance Evaluation and Market Comparison

Finally, we will evaluate our system’s performance against existing automation tools, such as OpenAI’s Operator, Runner H, and AgentE. By benchmarking our approach against these solutions, we aim to highlight its strengths, identify potential areas for improvement, and demonstrate its practical viability in real-world applications.

II. IMPACT AND BENEFITS

Our AI-powered web automation software has numerous applications across different domains. By enabling AI models to interact with and control browser-based applications, it enhances productivity, reduces manual effort, and opens up new possibilities for automation. Below are some key use cases, their target audience, and their impact.

A. Personal Task Automation

Target Audience: General users, busy professionals, students

Usage: Users can automate tasks like online shopping, scheduling appointments, filling out forms, and gathering information from multiple websites.

Impact: Saves time and effort in repetitive digital tasks, allowing users to focus on more important activities.

B. AI-Powered Virtual Assistants

Target Audience: General users, business professionals, digital marketers

Usage: The software can act as a personal assistant, handling tasks like checking emails, finding schedules, booking flights, or setting reminders.

Impact: Enhances convenience by automating daily activities and streamlining workflow.

C. Web Scraping and Data Extraction

Target Audience: Researchers, journalists, financial analysts, market researchers

Usage: AI can extract and analyze data from multiple sources, such as collecting news, stock prices, product comparisons, or social media trends.

Impact: Enables quick and accurate data gathering for research and decision-making.

D. Accessibility for Disabled Users

Target Audience: Visually impaired and physically challenged individuals

Usage: AI can navigate websites, read content aloud, and interact with forms or buttons for users who face difficulties using traditional input devices.

Impact: Increases digital accessibility, making the internet more inclusive for disabled individuals.

E. Automated Customer Support & Chatbot Enhancements

Target Audience: Businesses, e-commerce platforms, customer service teams

Usage: AI can navigate websites to find relevant answers, handle customer queries, and assist chatbots in providing faster responses.

Impact: Reduces human workload in customer service, leading to quicker responses and improved user experience.

F. Social Media Management

Target Audience: Social media managers, influencers, marketing teams

Usage: AI can schedule posts, reply to messages, monitor trends, and even generate reports based on user engagement.

Impact: Saves time and effort in managing multiple social media accounts, improving efficiency.

G. Automated Testing for Web Applications

Target Audience: Software testers, QA engineers, developers

Usage: The AI can simulate user interactions, test UI elements, and verify web application functionality without manual intervention.

Impact: Speeds up software testing, improves accuracy, and reduces human errors in quality assurance.

H. Fraud Detection and Security Monitoring

Target Audience: Cybersecurity professionals, financial institutions, online service providers

Usage: AI can monitor web activity, detect unusual patterns, and flag potential fraud or cyber threats.

Impact: Enhances digital security by providing proactive monitoring and fraud prevention.

I. Financial and Investment Automation

Target Audience: Investors, traders, financial advisors

Usage: AI can track stock markets, execute trades, and provide real-time insights based on web data.

Impact: Helps investors make informed decisions faster, reducing manual monitoring efforts.

J. Educational Assistance

Target Audience: Students, teachers, online learners

Usage: AI can fetch learning materials, summarize articles, assist with online assignments, and manage study schedules.

Impact: Makes learning more efficient and accessible by reducing time spent searching for resources.

K. E-Governance and Public Services

Target Audience: Government agencies, citizens seeking online services

Usage: AI can assist in automating form submissions, tracking application statuses, and fetching government-related information.

Impact: Improves accessibility to government services and reduces paperwork for citizens.

L. E-Commerce and Price Comparison

Target Audience: Online shoppers, businesses, product researchers

Usage: AI can compare product prices, find discounts, and alert users about the best deals.

Impact: Helps users save money and businesses optimize pricing strategies.

III. PROJECT BOUNDARIES AND SCOPE

This project is focused on enabling AI models to control and interact with web-based applications. However, to maintain clarity and feasibility, certain boundaries and constraints have been defined.

A. Scope of the Project

- **Browser-Based Tasks Only:** This software is exclusively designed to handle tasks within a web browser. It will not be capable of controlling or automating native desktop applications.
- **Simple Web Applications:** The project is primarily meant for automating interactions with standard web applications, such as form filling, data retrieval, and navigation. Complex activities like playing online games or dynamically surfing entertainment platforms are beyond the scope of this project.
- **Extensions of the Core Solution:** While various potential use cases have been mentioned earlier, they are merely extensions of the core technology. Our main focus remains on developing the fundamental interaction capabilities, leaving additional functionalities such as voice-based input to the open-source community for further enhancements.
- **Privacy Considerations:** The privacy of user data depends on the deployment choice. If a cloud-based LLM, such as ChatGPT, is used, the data privacy is subject to the policies of the service provider. Users are advised to share only publicly available information when interacting with cloud-based models. Alternatively, a self-hosted LLM offers greater control over privacy and security.

B. Aspects Not Covered in This Project

The following advanced capabilities are currently not part of this project but may be explored in future research:

- **Reinforcement Learning:** The ability to remember past tasks and improve future performance based on user interactions is not implemented in this version.

- **Cross-Device Operation:** The software does not support running on a server or multiple devices to enable remote control of desktops via mobile phones.
- **Task Reminders and Scheduling:** Features such as reminding users of deadlines or scheduled activities are not included in the initial implementation.
- **Heuristic-Based Adjustments:** While AI-driven automation is the focus, the software does not currently integrate heuristics such as automatically retrying failed actions with alternative methods.
- **No Interaction with Encrypted or Highly Dynamic Content:** The software will not be designed to interact with content that requires decryption, such as DRM-protected media, banking portals with dynamic token-based authentication, or highly volatile AJAX-based pages where elements constantly change.
- **Not a Full AI Agent, but a Task Executor:** The system will not function as a fully autonomous AI agent that makes independent decisions beyond executing specific user-defined tasks. It will always require user input or predefined workflows to operate.
- **No Direct API-Based Automation:** Many web services provide APIs for automation, but this project is focused on browser-based interaction, not direct API calls to services like Gmail API or Twitter API.
- **No Real-Time Multi-User Collaboration:** The system is designed for individual use, and it does not support multiple users interacting with the same session simultaneously. Collaborative AI agents controlling the same browser across multiple users or locations are out of scope.
- **Not Meant for Bots or Mass Web Interaction:** This software is not intended for large-scale web crawling, scraping, or bot-driven automation that could violate website terms of service. It focuses on assisting users with personal and professional tasks within ethical boundaries.
- **No Deep Customization Without Developer Input:** While the system may allow some level of configuration, advanced custom workflows, integrations, or AI model modifications would require developer intervention rather than a no-code interface for end-users.
- **Limited to Standard Web Browsers:** The system is built to work within standard web browsers (Chrome, Firefox, Edge, etc.) and does not support automation in specialized browsing environments like embedded web views in games or proprietary software.

By defining these boundaries, we ensure a focused and structured approach to development while leaving room for future enhancements and community-driven improvements.

IV. OBJECTIVES

The primary objective of this project is to develop an AI-driven software solution capable of interacting with and controlling web browsers to automate tasks and extract information. The following key objectives define the scope of the project:

- 1) **Develop an AI-Powered Web Automation Software:** Design and implement a software system that enables AI models to control web browsers and facilitate interaction with web pages for information extraction and task execution.
- 2) **Accurate UI Element Detection:** Develop mechanisms to accurately detect and classify UI elements (buttons, input fields, links, tables, etc.) using a combination of techniques, including:
 - Parsing DOM (Document Object Model) content to extract structured data.
 - Employing computer vision techniques to analyze UI components visually.
- 3) **Understanding UI Element Relationships:** Implement algorithms to determine the relationships between UI elements, such as:
 - Identifying parent-child relationships in DOM structures.
 - Recognizing contextual groupings (e.g., input field and its corresponding label).
 - Detecting logical flows between elements for smooth automation.
- 4) **Leveraging Large Language Models (LLMs) for Decision Making:** Integrate the power of LLMs to analyze extracted UI information and decide on appropriate follow-up actions, such as:
 - Interpreting user intent and converting it into actionable browser interactions.
 - Handling ambiguous UI structures by reasoning about possible actions.
 - Generating step-by-step execution plans dynamically.
- 5) **Executing AI-Decided Operations through Python Automation:** Develop a robust Python-based execution pipeline that:
 - Receives operation instructions from the LLM.
 - Uses automation libraries (e.g., Selenium, Playwright, or Puppeteer) to interact with web pages.
 - Handles exceptions and dynamically adjusts to unexpected scenarios.
- 6) **Implementing a User-Friendly Interface:** Design an intuitive UI that allows users to interact with the software easily, including:
 - A dashboard to monitor AI-driven interactions.
 - Controls to adjust automation parameters.
 - Logs and explanations of executed actions for transparency.
- 7) **Ensuring Privacy and Security:** Address privacy concerns by:
 - Providing an option for users to self-host the AI model to avoid sending sensitive data to cloud-based services.
 - Implementing secure handling of user credentials and private browsing sessions.

- Avoiding unauthorized interactions that may violate website terms of service.

8) **Optimizing for Performance and Scalability:** Ensure that the system is lightweight, efficient, and capable of handling:

- Large-scale automation without excessive resource consumption.
- Seamless execution across different web browsers.
- Adaptability to new web layouts and technologies over time.

V. LITERATURE SURVEY

A. Research Paper 1: AutoWebGLM: A Large Language Model-based Web Navigating Agent (Lai et al., 2024)

This paper introduces AutoWebGLM, a web navigation agent based on ChatGLM3-6B, which utilizes HTML simplification, curriculum training, and reinforcement learning to autonomously navigate and interact with webpages for complex tasks [?]. AutoWebGLM addresses the challenges of complex HTML, versatile webpage actions, and open-domain task difficulty by first simplifying HTML to retain vital information, then using a hybrid human-AI method to create curriculum training data, and finally employing reinforcement learning with rejection sampling to enhance webpage comprehension, browser operation, and task decomposition. The research also introduces AutoWebBench, a bilingual benchmark for evaluating real-world web navigation tasks, demonstrating AutoWebGLM’s potential to outperform even GPT-4 in challenging web environments.

B. Research Paper 2: Steward: An LLM-Powered Web Automation Tool for Cost-Effective and Scalable End-to-End Web Interactions

This paper presents Steward, a novel LLM-powered web automation tool integrated with the Playwright framework, designed to provide a cost-effective and scalable solution for automating web interactions based on natural language instructions [?]. Steward aims to overcome the limitations of traditional frameworks that require manual coding by enabling natural language-driven interaction with websites. The system is designed to be fully autonomous, requiring only a high-level natural language goal to perform operations on websites until the desired end state is reached. Evaluations demonstrate Steward’s efficiency in terms of speed and cost, achieving a median runtime of 8.52 to 10.14 seconds and a cost of \$0.028 per action, with further reductions through a caching mechanism. It also shows a 40% task completion rate on real websites, highlighting its potential as a practical web automation solution.

C. Research Paper 3: The Potential of LLMs in Automating Software Testing: From Generation to Reporting

This paper explores an agent-oriented framework leveraging LLMs to automate various aspects of software testing, including generating unit tests and controlling a browser-like interface for test execution and reporting [?], [?]. The proposed

framework, named Kashef, integrates LLMs to generate unit tests dynamically, visualize call graphs, and automate the execution and reporting of tests, aiming to reduce human intervention and enhance testing efficiency. Evaluations across multiple applications in Python and Java demonstrate the system’s high test coverage and efficient operation, underscoring the potential of LLM-powered agents to streamline software testing workflows while addressing challenges in scalability and accuracy. Kashef employs a multi-agent architecture comprising a Test Engineer and an HTML Interpreter (both LLM-powered), alongside a Code Executor.

D. Research Paper 4: Large Language Model-Brained GUI Agents: A Survey (Zhang et al., 2024)

This survey paper provides a comprehensive overview of the rapidly evolving field of LLM-powered GUI agents, including those for web navigation, discussing their historical evolution, core components, and advanced techniques [?]. The paper explores how the advent of LLMs, particularly multimodal models, has ushered in a new era of GUI automation, enabling agents to interpret complex GUI elements and autonomously execute actions based on natural language instructions. It addresses critical research questions such as existing GUI agent frameworks, data collection and utilization for training, the development of large action models, and evaluation metrics. By consolidating foundational knowledge and state-of-the-art developments, the survey aims to guide researchers and practitioners in overcoming challenges and unlocking the full potential of LLM-brained GUI agents across web, mobile, and desktop domains.

E. Research Paper 5: Utilizing LLMs to Infer Clickability of HTML Elements for Enhanced Deep Reinforcement Learning-Based Web Application Testing

This paper proposes a method to enhance deep reinforcement learning for web application testing by using LLMs to infer whether HTML elements are clickable, incorporating this inference into the DRL state to improve testing efficiency [?]. The research addresses the challenges of traditional GUI testing by integrating the extensive knowledge of LLMs, which are not limited to web apps, to aid in obtaining states in DRL when direct collection is difficult. The hypothesis is that the probabilistic nature of LLM inference can be compensated by the robustness of DRL models. By using LLMs to infer the clickability of HTML elements and incorporating this into the DRL state, the proposed method aims to improve the efficiency and accuracy of automated black-box GUI testing for web applications.

F. Comparative Analysis

The five research papers summarized above showcase a variety of architectural designs and implementation strategies for browser control software powered by LLMs.

Feature	AutoWebGLM	Steward	Software Testing Paper	GUI Agents Survey	Clickability Inference Paper
Control Mechanism	Monolithic (Fine-tuned LLM)	LLM + Framework	Distributed (Multi-Agent)	Varies (Direct LLM to Complex Systems)	LLM as Component in DRL
Browser Interaction Framework	Not explicitly mentioned	Playwright	Selenium	Discusses systems using Selenium & Playwright	Selenium
Webpage Understanding	HTML Simplification Algorithm	Implied through Playwright integration	HTML Interpreter Agent	DOM Parsing, Computer Vision, Multimodal LLMs	LLM infers clickability from HTML
LLM Used	ChatGLM3-6B	Not specified in abstract	GPT-3.5, GPT-4, CodeLlama, Llama2	Wide range, including GPT-4 & open-source LLMs	Not detailed in abstract
Fine-tuning/Training Data	Curriculum training, Reinforcement Learning	High accuracy without training for top elements	Not detailed in abstract	Discusses data collection and utilization	Not detailed in abstract

VI. NOVELTY

Our project introduces several innovative aspects that distinguish it from existing solutions in the field of AI-driven web automation:

A. Open-Source Architecture

Unlike many proprietary solutions in the market, our system is designed with an open-source architecture. This approach fosters transparency, community collaboration, and continuous improvement. The open-source nature allows developers to:

- Access and modify the source code freely
- Contribute to the project's development
- Create custom extensions and plugins
- Share improvements with the community

B. Self-Hostable Design

A key differentiator of our system is its self-hostable architecture. Users have the flexibility to:

- Deploy the system on their own infrastructure
- Maintain complete control over their data and operations
- Customize the system according to their specific needs
- Avoid dependency on external cloud services

C. Privacy-First Approach

Our system prioritizes user privacy through multiple mechanisms:

- **Self-Hosted LLM Option:** Users can deploy their own LLM instance, ensuring sensitive data never leaves their infrastructure
- **Local Processing:** Critical operations are performed locally, minimizing data exposure
- **Configurable Privacy Settings:** Users can fine-tune privacy parameters based on their requirements
- **Transparent Data Handling:** Clear documentation of data flow and storage practices

D. Adaptive Learning System

Unlike traditional automation tools, our system incorporates:

- Dynamic adaptation to different website layouts and structures
- Learning from user interactions to improve future performance
- Context-aware decision making for complex web interactions
- Intelligent error recovery and alternative action planning

E. Multi-Modal Understanding

Our solution combines multiple approaches for robust web interaction:

- DOM-based structural analysis
- Visual element recognition
- Natural language understanding
- Contextual relationship mapping

F. User-Centric Design

The system is built with a strong focus on user experience:

- Natural language interface for task specification
- Intuitive control mechanisms
- Detailed operation logging and explanation
- Customizable automation workflows

G. Extensible Framework

The architecture supports future enhancements through:

- Modular plugin system
- Standardized API interfaces
- Custom action definition capabilities
- Integration with external tools and services

These novel aspects collectively position our solution as a comprehensive, user-friendly, and privacy-conscious approach to AI-driven web automation, addressing key limitations in existing solutions while providing unique advantages for end users.

REFERENCES

- [1] Browser automation: Going beyond the limits of LLMs — CLOVA, accessed March 25, 2025, <https://clova.ai/en/tech-blog/browser-automation-going-beyond-the-limits-of-llms>
- [2] Large Language Model-Brained GUI Agents: A Survey - arXiv, accessed March 25, 2025, <https://arxiv.org/html/2411.18279v3>
- [3] Steward: Natural Language Web Automation - arXiv, accessed March 25, 2025, <https://arxiv.org/html/2409.15441v1>
- [4] (PDF) The Potential of LLMs in Automating Software Testing: From ..., accessed March 25, 2025, https://www.researchgate.net/publication/387670287_The_Potential_of_LLMs_in_Automating_Software_Testing_From_Generation_to_Reporting
- [5] The Potential of LLMs in Automating Software Testing: From Generation to Reporting - arXiv, accessed March 25, 2025, <https://arxiv.org/html/2501.00217v1>
- [6] Skyvern-AI/skyvern: Automate browser-based workflows with LLMs and Computer Vision - GitHub, accessed March 25, 2025, <https://github.com/Skyvern-AI/skyvern>
- [7] www.scitepress.org, accessed March 25, 2025, <https://www.scitepress.org/Papers/2025/132488/132488.pdf>
- [8] Develop Browser Agents: Integrating LLMs, Playwright, Browser-Use and Web-UI — by KailashPathak — Feb, 2025, accessed March 25, 2025, <https://kailash-pathak.medium.com/develop-intelligent-browser-agents-integrating-llms-playwright-browser-use-and-web->
- [9] AutoWebGLM: A Large Language Model-based Web Navigating Agent - arXiv, accessed March 25, 2025, <https://arxiv.org/html/2404.03648v2>
- [10] [2404.03648] AutoWebGLM: A Large Language Model-based Web Navigating Agent - arXiv, accessed March 25, 2025, <https://arxiv.org/abs/2404.03648>
- [11] arxiv.org, accessed March 25, 2025, <https://arxiv.org/abs/2409.15441>
- [12] Large Language Model-Brained GUI Agents: A Survey - arXiv, accessed March 25, 2025, <https://arxiv.org/html/2411.18279v1>
- [13] [2411.18279] Large Language Model-Brained GUI Agents: A Survey - arXiv, accessed March 25, 2025, <https://arxiv.org/abs/2411.18279>
- [14] Paper page - Large Language Model-Brained GUI Agents: A Survey - Hugging Face, accessed March 25, 2025, <https://huggingface.co/papers/2411.18279>
- [15] Should you let ChatGPT control your browser? - WithSecure™ Labs, accessed March 25, 2025, <https://labs.withsecure.com/publications/browser-agents-llm-prompt-injection>