

GUI Agents with Foundation Models: A Comprehensive Survey

Shuai Wang¹, Weiwen Liu¹, Jingxuan Chen¹, Yuqi Zhou², Weinan Gan¹,
Xingshan Zeng¹, Yuhan Che¹, Shuai Yu¹, Xinlong Hao¹, Kun Shao¹,
Bin Wang¹, Chuhan Wu¹, Yasheng Wang¹, Ruiming Tang¹, Jianye Hao¹

¹Huawei Noah’s Ark Lab ²Renmin University of China

{wangshuai231, liuweiben8}@huawei.com

Abstract

Recent advances in foundation models, particularly Large Language Models (LLMs) and Multimodal Large Language Models (MLLMs), have facilitated the development of intelligent agents capable of performing complex tasks. By leveraging the ability of (M)LLMs to process and interpret Graphical User Interfaces (GUIs), these agents can autonomously execute user instructions, simulating human-like interactions such as clicking and typing. This survey consolidates recent research on (M)LLM-based GUI agents, highlighting key innovations in data resources, frameworks, and applications. We begin by reviewing representative datasets and benchmarks, followed by an overview of a generalized, unified framework that encapsulates the essential components of prior studies, supported by a detailed taxonomy. Additionally, we explore relevant commercial applications. Drawing insights from existing work, we identify key challenges and propose future research directions. We hope this survey will inspire further advancements in the field of (M)LLM-based GUI agents.

1 Introduction

Graphical User Interfaces (GUIs) are the primary medium through which humans interact with digital devices. From mobile phones to websites, people engage with GUIs daily, and well-designed GUI agents can significantly enhance the user experience. Thus, research on GUI agents has been extensive. However, traditional rule-based and reinforcement learning-based methods struggle with tasks requiring human-like interactions [Liu *et al.*, 2018], limiting their applicability.

Recent advancements in Large Language Models (LLMs) and Multimodal Large Language Models (MLLMs) have significantly enhanced their capabilities in language understanding and cognitive processing [Achiam *et al.*, 2024; Touvron *et al.*, 2023; Yang *et al.*, 2024a]. With improved natural language comprehension and enhanced reasoning abilities, (M)LLM-based agents can now effectively interpret and utilize human language, formulate detailed plans, and execute complex tasks. These breakthroughs provide new opportuni-

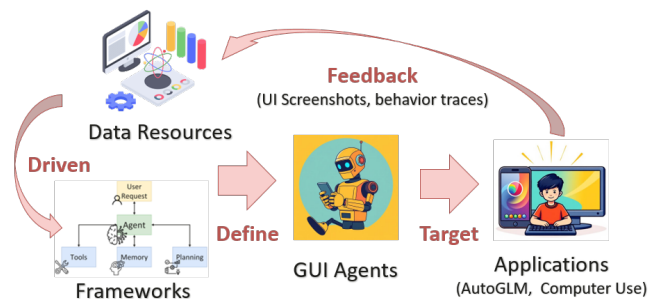


Figure 1: The foundational aspects and goals of GUI agents.

ties for researchers to address challenges previously considered highly difficult, such as automating tasks within GUIs.

As shown in Figure 2, recent studies on GUI agents illustrate a shift from simple Transformer-based models to (M)LLM-based agentic frameworks. Their capabilities have expanded from single-modality interactions to multimodal processing, making them increasingly relevant to commercial applications. Given these advancements, we believe it is timely to systematically analyze the development trends of GUI agents, particularly from an application perspective.

This paper aims to provide a structured overview of the latest and influential work in the field of GUI agents. As depicted in Figure 1, we focus on the foundational aspects and goals of GUI agents. Data resources, such as user instructions, User Interface (UI) screenshots, and behavior traces, drive the design of GUI agents [Rawles *et al.*, 2023; Lu *et al.*, 2024a]. Frameworks define the underlying algorithms and models that enable intelligent decision-making [Li *et al.*, 2024b; Wang *et al.*, 2024a; Zhu *et al.*, 2024]. Applications represent the optimized and practical goals [Lai *et al.*, 2024; Liu *et al.*, 2024]. The current state of these aspects reflects the maturity of the field and highlights future research priorities.

To this end, we organize this survey around three key areas: **Data Resources**, **Frameworks**, and **Applications**. The main contributions of this paper are: 1) a comprehensive summary of existing research and a detailed review of current data sources, providing a useful guide for newcomers to the field; 2) a unified and generalized GUI agent framework with clearly defined and categorized functional components to facilitate a structured review; 3) an analysis of trends in both

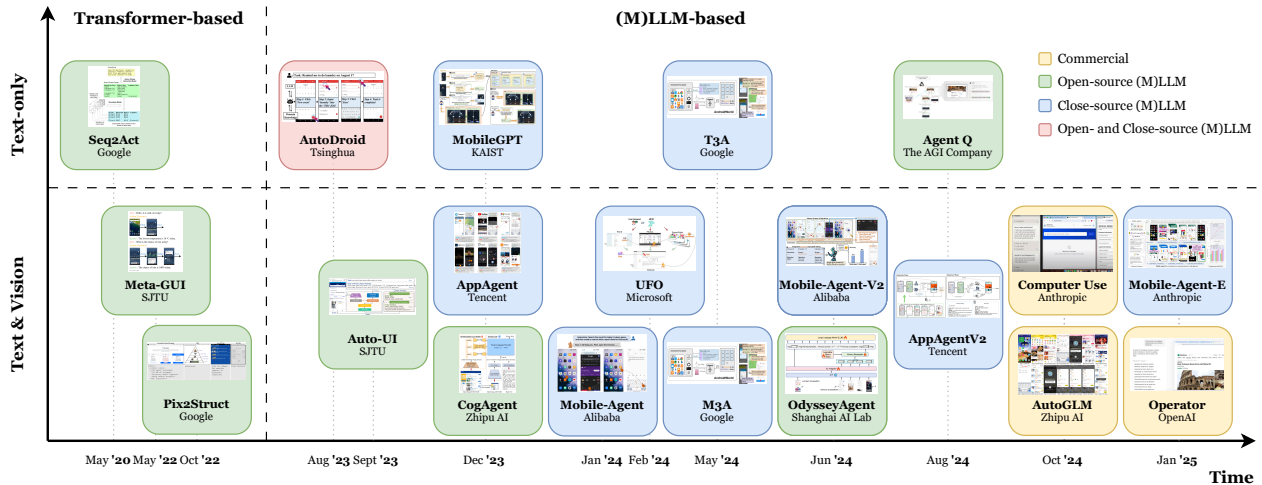


Figure 2: Illustration of the growth trend in the field of GUI agents with foundation models.

research and commercial applications of GUI agents.

2 GUI Agent Data Resources

Recent research has focused on developing datasets and benchmarks to train and evaluate the capabilities of (M)LLM-based GUI agents. A variety of datasets are available for training GUI agents. These agents employ different approaches to interact with environments. Additionally, multiple methods have been proposed for evaluation.

Dataset: Common datasets for training GUI agents typically contain natural language instructions that describe task goals, along with demonstration trajectories that include screenshots and action pairs. A pioneering work in this area is PIXELHELP [Li *et al.*, 2020], which introduces a new class of problems focused on translating natural language instructions into actions on mobile user interfaces. In recent years, Android in the Wild [Rawles *et al.*, 2023] has created a dataset featuring a variety of single-step and multi-step tasks. Aimed at advancing GUI navigation agent research, Android-In-The-Zoo [Zhang *et al.*, 2024b] introduces a benchmark dataset with chained action reasoning annotations.

Insight-UI [Shen *et al.*, 2024] automatically constructs a GUI pre-training dataset that simulates multiple platforms across 312,000 domains. To assess model performance both within and beyond the scope of training data, AndroidControl [Li *et al.*, 2024a] includes demonstrations of daily tasks along with both high- and low-level human-generated instructions. The scope of mobile control datasets is further extended from single-application to cross-application scenarios by GUI-Odyssey [Lu *et al.*, 2024a].

Most of the aforementioned datasets are primarily limited to English and image-based tasks. However, UGIF Dataset [Venkatesh *et al.*, 2024] covers eight languages, Mobile3M [Wu *et al.*, 2024] focuses on Chinese, and GUI-WORLD [Chen *et al.*, 2024a] includes video annotations, expanding the dataset landscape for broader multilingual and multimodal research.

Environment: GUI agents require environments for task execution, which can be broadly categorized into three types. The first category is static environments, where the environment remains fixed as it was when developed. Agents in this category operate within predefined datasets without the ability to create new states.

In contrast, the second and third categories involve dynamic environments, where new outcomes can emerge during agent execution. The key distinction between these categories lies in whether the dynamic environment is simulated or realistic. Simulations of real-world environments require additional implementation but are often cleaner and free of distractions, such as pop-ups and advertisements. WebArena [Zhou *et al.*, 2023] implements a versatile website covering e-commerce, social forums, collaborative software development, and content management. Similarly, GUI Testing Arena [Zhao *et al.*, 2024] provides a standardized environment for testing GUI agents, including defect injection.

For realistic environments, agents interact directly with web or mobile platforms as human users do, better reflecting real-world conditions. SPA-Bench [Chen *et al.*, 2024b] encompasses tasks that involve both system and third-party mobile applications, supporting single-app and cross-app scenarios in both English and Chinese.

Evaluation: Another critical component of GUI agent datasets is the evaluation of agent performance. The most common and important metric is success rate, which measures how effectively an agent completes tasks. Additional metrics, such as efficiency, are sometimes considered as well.

Evaluation methods are often closely tied to the environment type. In static environments, action matching is a widely used method that compares an agent’s executed action sequence with a human demonstration (e.g., Rawles *et al.* [2023], Li *et al.* [2024a]). However, a major limitation of action matching is its inability to account for multiple successful execution paths, leading to false negatives when evaluating agent performance.

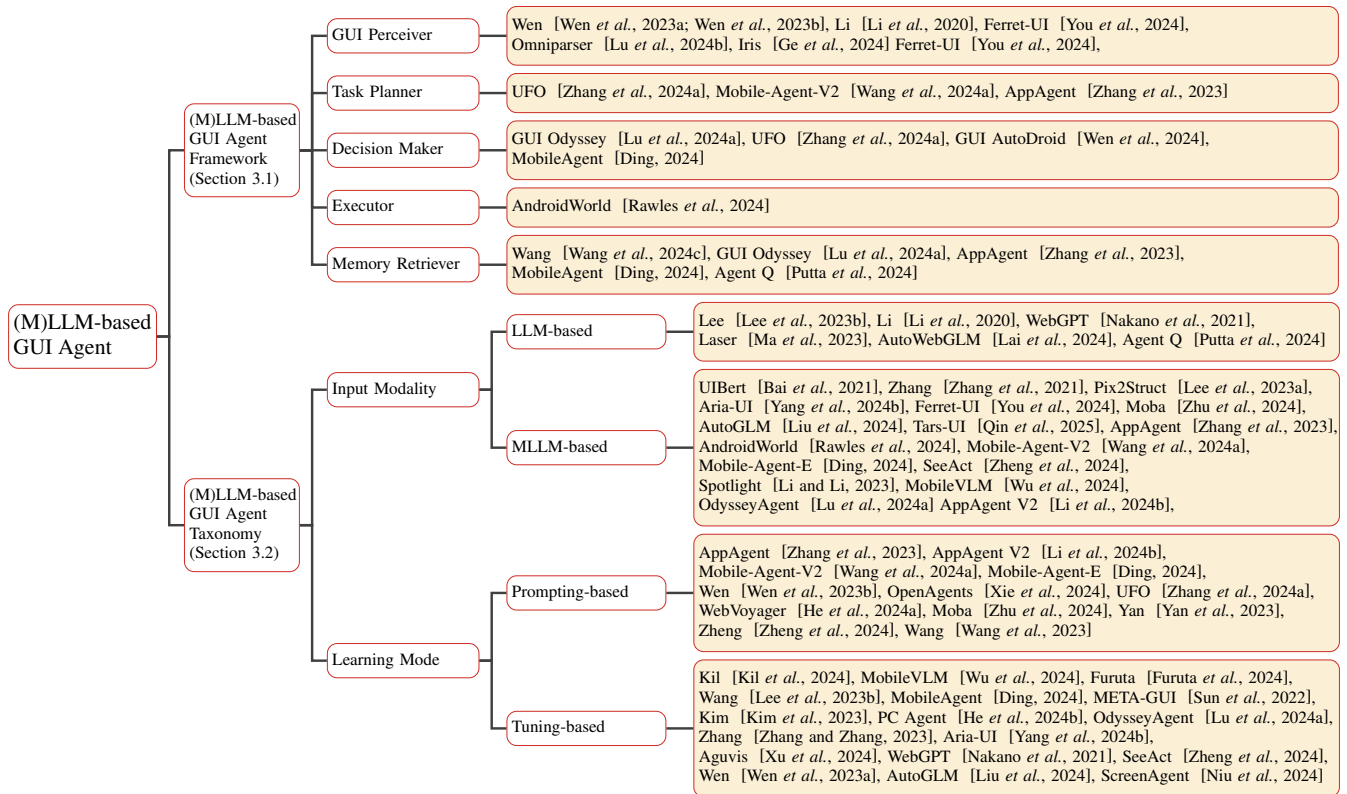


Figure 3: A comprehensive taxonomy of (M)LLM-based GUI Agents: frameworks, modality, and learning paradigms.

Evaluating dynamic environments, whether simulated or realistic, presents additional challenges due to their uncertain conditions. Evaluation methods can range from fully human-dependent to semi-automated and fully automated approaches. Human evaluations require manual verification, making them non-reusable. In AppAgent [Li et al., 2024b] and MobileAgent [Ding, 2024], human evaluators assess whether each agent-executed task was successful. Semi-automated evaluations involve human-developed validation logic that can be reused for different execution trajectories of the same task. For example, WebArena [Zhou et al., 2023] and AndroidWorld [Rawles et al., 2024] incorporate handcrafted validation functions for task completion. Fully automated evaluations eliminate human involvement by relying on models for success detection. SPA-Bench [Chen et al., 2024b], for instance, employs MLLMs for evaluating task completion. Although reducing human labor is crucial for large-scale evaluation, balancing efficiency with accuracy remains a key research challenge.

3 (M)LLM-based GUI Agent

With the human-like capabilities of (M)LLMs, GUI agents aim to handle various tasks to meet users’ needs. Organizing the frameworks of GUI agents and designing methods to optimize their performance is crucial to unlocking the full potential of (M)LLMs. As shown in Figure 3, we summarize a generalized **Framework** and discuss its components in relation to existing works in Section 3.1. Building on this founda-

tion, we then review recent influential **Methods** for constructing and optimizing GUI agents, categorizing them with an exhaustive taxonomy in Section 3.2.

3.1 (M)LLM-based GUI Agent Framework

The goal of GUI agents is to automatically control a device to complete tasks defined by the user. Typically, GUI agents take a user’s query and the device’s UI status as inputs and generate a series of human-like actions to achieve the tasks.

As shown in Figure 4, we present a generalized (M)LLM-based GUI agent framework, consisting of five components: GUI Perceiver, Task Planner, Decision Maker, Memory Retriever, and Executor. Many variations of this framework exist. For instance, Wang et al. [2024a] proposes a multi-agent GUI control framework comprising a planning agent, a decision agent, and a reflection agent to tackle navigation challenges in mobile device operations. This approach shares functional similarities with our proposed framework. A follow-up study by Wang et al. [2025] further disentangles high-level planning from low-level action decisions by employing dedicated agents and introduces memory-based self-evolution to enhance performance.

GUI Perceiver: To effectively complete a device task, a GUI agent should accurately interpret user input and detect changes in the device’s UI. Although language models excel in understanding user intent [Touvron et al., 2023; Achiam et al., 2024], navigating device UIs requires a reliable visual

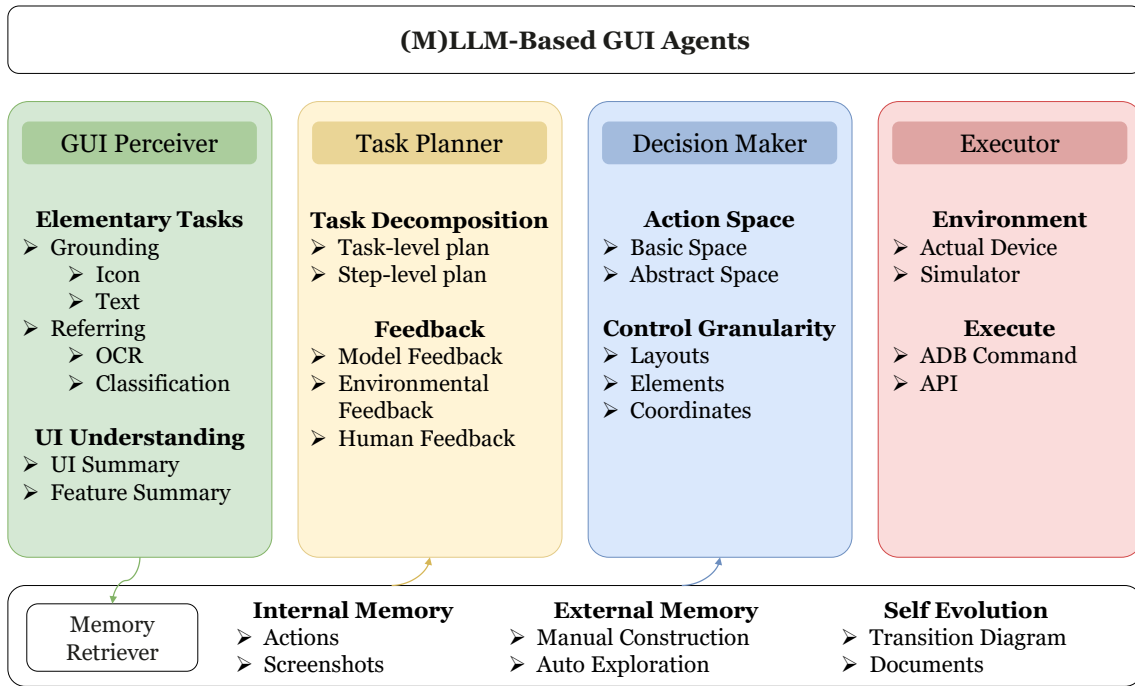


Figure 4: (M)LLM-based GUI agents: the generalized framework and key technologies.

perception model to understand GUIs.

A GUI Perceiver appears explicitly or implicitly in GUI agent frameworks. For agents based on single-modal LLMs [Wen *et al.*, 2023a; Wen *et al.*, 2023b; Li *et al.*, 2020], a GUI Perceiver is usually an explicit module of the frameworks. However, for agents with multi-modal LLMs [Hong *et al.*, 2024; Zhang *et al.*, 2023; Wang *et al.*, 2024b], UI perception is seen as a capability of the model itself.

UI perception is also an important problem in GUI agent research, some work [You *et al.*, 2024; Zhang *et al.*, 2021; Lu *et al.*, 2024b] focuses on understanding and processing UIs, rather than building the agent. For example, Pix2struct [Lee *et al.*, 2023a] employs a ViT-based image-encoder-text-decoder architecture, which pre-trains on Screenshot-HTML data pairs and fine-tunes for specific tasks. This method has shown strong performance in web-based visual comprehension tasks. Screen2words [Wang *et al.*, 2021] is a novel approach that encapsulates a UI screen into a coherent language representation, which is based on a transformer encoder-decoder architecture to process UIs and generate the representation. To address the defects of purely vision-based screen parsing methods, Ge *et al.* [2024] introduces Iris, a visual agent for GUI understanding, addressing challenges related to architectural limitations for heterogeneous GUI information and annotation bias in GUI training via two innovations: An information-sensitive architecture to prioritize high-density UI regions via edge detection, and a dual-learning strategy that refines visual/functional knowledge iteratively using unlabeled data, reducing annotation dependence.

Task Planner: The GUI agent should effectively decompose complex tasks, often employing a Chain-of-Thought (CoT) approach. Due to the complexity of tasks, recent stud-

ies [Zhang *et al.*, 2024a; Wang *et al.*, 2024a] introduce an additional module to support more detailed planning.

Throughout the GUI agent’s process, plans may adapt dynamically based on decision feedback, typically achieved through a ReAct-style. For instance, Zhang *et al.* [2023] uses on-screen observations to enhance the CoT for improved decision-making, while Wang *et al.* [2024a] develops a reflection agent that provides feedback to refine plans.

Decision Maker: A Decision Maker provides the next operation(s) to control a device. Most studies [Lu *et al.*, 2024a; Zhang *et al.*, 2024a; Wen *et al.*, 2024] define a set of UI-related actions—such as click, text, and scroll—as a basic action space. In a more complicated case, Ding [2024] encapsulates a sequence of actions to create Standard Operating Procedures(SOPs) to guide further operations.

As the power of GUI agents improves, the granularity of operations becomes more refined. Recent work has progressed from element-level operations [Zhang *et al.*, 2023; Wang *et al.*, 2024b] to coordinate-level controls [Wang *et al.*, 2024a; Hong *et al.*, 2024].

Executor: An Executor maps outputs to the relevant environments. While most studies use Android Debug Bridge (ADB) to control real devices [Li *et al.*, 2024b; Wang *et al.*, 2024a], Rawles *et al.* [2024] develops a simulator to access additional UI-related information.

Memory Retriever: A Memory Retriever is designed as an additional source of information to help agents perform tasks more effectively [Wang *et al.*, 2024c].

GUI agents’ memory is typically divided into internal and external categories. Internal memory [Lu *et al.*, 2024a] consists of prior actions, screenshots, and system states during execution, while external memory [Zhang *et al.*, 2023; Ding, 2024] includes knowledge and rules related to the UI

or task, providing additional inputs for the agent.

3.2 (M)LLM-based GUI Agent Taxonomy

Consequently, this paper classifies existing work with the difference of input modality and learning mode in Figure 3.

GUI Agents with Different Input modality

LLM-based GUI Agents: With the limited multimodal capability, earlier GUI agents [Lee *et al.*, 2023b; Li *et al.*, 2020; Ma *et al.*, 2023; Lai *et al.*, 2024; Putta *et al.*, 2024; Nakano *et al.*, 2021; Nakano *et al.*, 2021] often require a GUI perceiver to convert GUI screens into text-based inputs.

So, parsing and grounding the GUI screens is the first step. For instance, Li *et al.* [2020] transforms the screen into a series of object descriptions and applies a transformer-based action mapping. The problem definitions and datasets have spurred further research. You *et al.* [2024] proposes a series of referring and grounding tasks, which provide valuable insights into the pre-training of GUIs. Lu *et al.* [2024b] proposes a screen parsing framework incorporating the local semantics of functionality with interactable region detection for better UI understanding and element grounding.

Afterward, LLMs are used as the core of the agents. Wen *et al.* [2024] further converts GUI screenshots into a simplified HTML representation for compatibility with the LLMs. By combining GUI representation with app-specific knowledge, they build Auto-Droid, a GUI agent based on online GPT and on-device Vicuna. In the field of web automation, LASER [Ma *et al.*, 2023] navigates web environments purely through text, treating web navigation as state-space exploration to enable flexible state transitions and error recovery. Similarly, AutoWebGLM [Lai *et al.*, 2024] processes HTML text data without visual inputs, refining webpage structures to preserve key information for ChatGLM3-6B. Agent Q [Putta *et al.*, 2024] further extends this paradigm by relying solely on HTML DOM text for reasoning and decision-making, emphasizing language models for planning and action execution. WebGPT [Nakano *et al.*, 2021], a fine-tuned GPT-3 model, uses text-based web browsing (processing HTML content) to collect information via commands like searching and clicking. It generates answers supported by references and is optimized using human feedback and rejection sampling.

MLLM-based GUI Agents: Recent studies [Wang *et al.*, 2024a; Bai *et al.*, 2021; Zhang *et al.*, 2023; Kim *et al.*, 2023] utilize the multimodal capabilities of advanced (M)LLMs to improve GUI comprehension and task execution.

Leveraging the visual understanding capabilities of MLLMs, recent studies [Wang *et al.*, 2024a; Li and Li, 2023; Bai *et al.*, 2021; Zhu *et al.*, 2024; Qin *et al.*, 2025] explore end-to-end frameworks for GUI device control. For example, Spotlight [Li and Li, 2023] proposes a Vision-Language model framework, pre-trained on Web/mobile data and fine-tuned for UI tasks. This model greatly improves the ability to understand UIs. By combining screenshots with a user focus as input, Spotlight outperforms previous methods on multiple UI understanding tasks, showing verified gains in downstream tasks. Likewise, VUT [Li *et al.*, 2021] is proposed for GUI understanding and multi-modal UI input modeling, using two Transformers: one for encoding and fusing

image, structural, and language inputs, and the other for linking three task heads to complete five distinct UI modeling tasks and learn downstream multiple tasks end-to-end. Experiments show that VUT’s multi-task learning framework can achieve state-of-the-art (SOTA) performance on UI modeling tasks. UIbert [Bai *et al.*, 2021] focuses on heterogeneous GUI features and considers that the multi-modal information in the GUI is self-aligned. UIbert is a transformer-based joint image-text model, which is pre-trained in large-scale unlabeled GUI data to learn the feature representation of UI elements. Zhu *et al.* [2024] presents a two-level agent structure for executing complex and dynamic GUI tasks. Moba’s Global Agent handles high-level planning, while the Local Agent selects actions for sub-tasks, streamlining the decision-making process with improved efficiency. UI-TARS [Qin *et al.*, 2025] navigates interfaces through screenshots, enabling human-like interactions via keyboard and mouse. Leveraging a large-scale GUI dataset, it achieves context-aware UI understanding and precise captioning.

To enhance performance, some studies [Zhang *et al.*, 2023; Rawles *et al.*, 2024] utilize additional invisible metadata. For instance, AndroidWorld [Rawles *et al.*, 2024] establishes a fully functional Android environment with real-world tasks, serving as a benchmark for evaluating GUI agents. They propose M3A, a zero-shot prompting agent that uses Set-of-Marks as input. Experiments with M3A variants assess how different input modalities—text, screenshots, and accessibility trees—affect GUI agent performance. Yang *et al.* [2024b] proposes a framework incorporating dynamic action history with both textual and interleaved text-image formats, which allows it to ground elements more effectively for dynamic, multi-step scenarios.

GUI Agents with Different Learning Mode

Prompting-based GUI Agents: Prompting is an effective approach to building agents with minimal extra computational overhead. Given the diversity of GUIs and tasks, numerous studies [Zhang *et al.*, 2023; Li *et al.*, 2024b; Wang *et al.*, 2024a; Wen *et al.*, 2023b; Xie *et al.*, 2024; Zhang *et al.*, 2024a; He *et al.*, 2024a] use prompting to create GUI agents, adopting CoT or ReAct styles.

Recent studies use prompting to build and simulate the functions of GUI agent components. For example, Yan *et al.* [2023] introduces MM-Navigator, which utilizes GPT-4V for zero-shot GUI understanding and navigation. For the first time, this work demonstrates the significant potential of LLMs, particularly GPT-4V, for zero-shot GUI tasks. Manual evaluations show that MM-Navigator achieves impressive performance in generating reasonable action descriptions and single-step instructions for iOS tasks. Additionally, Wen *et al.* [2023b] presents DroidBot-GPT, which summarizes the app’s status, past actions, and tasks into a prompt, using ChatGPT to choose the next action. Beyond mobile applications, prompting-based approaches have also been widely adopted in web-based GUI agents. Zheng *et al.* [2024] proposes SeeAct, a GPT-4V-based generalist web agent. With screenshots as input, SeeAct generates action descriptions and converts them into executable actions with designed action grounding techniques. OpenAgents [Xie *et al.*, 2024] leverages prompts

Table 1: Overview of (M)LLM-Based GUI Agents.

Model Name	Category	GUI Perceiver	Learning Method	Base Model	Scenarios
Prompting-based					
PaLM [Wang <i>et al.</i> , 2023]	Single Step	HTML	Few-shot prompting	PaLM	Mobile
MM-Navigator [Yan <i>et al.</i> , 2023]	Single Step	Screenshot	Zero-shot prompting	GPT-4V	Mobile
MemōDroid [Lee <i>et al.</i> , 2023b]	End-to-End	HTML	Few-shot prompting	ChatGPT/GPT-4V	Mobile/Desktop
AutoTask [Pan <i>et al.</i> , 2023]	End-to-End	Screenshot/API	Zero-shot prompting	GPT-4V	Mobile
AppAgent [Zhang <i>et al.</i> , 2023]	End-to-End	Screenshot	Exploration-based/In-context learning	GPT4V	Mobile
DroidBot-GPT [Wen <i>et al.</i> , 2023b]	End-to-End	Screenshot	Zero-shot prompting	ChatGPT	Mobile
Mobile-Agent-V2 [Wang <i>et al.</i> , 2024a]	End-to-End	Screenshot	Zero-shot prompting	GPT4V	Mobile
SeeAct [Zheng <i>et al.</i> , 2024]	End-to-End	Screenshot/HTML	Few-shot prompting	GPT-4V	Web
Mobile-Agent-E [Wang <i>et al.</i> , 2025]	End-to-End	Screenshot	Zero-shot prompting	GPT-4o/Claude-3.5-Sonnet/Gemini-1.5-pro	Mobile
Learning-based					
Spotlight [Li and Li, 2023]	UI modeling	Screenshot	Pretrain/SFT	ViT	Mobile/Web
Pix2Struct [Lee <i>et al.</i> , 2023a]	UI modeling	Screenshot	Pretrain/SFT	ViT	Web
VUT [Li <i>et al.</i> , 2021]	UI modeling	Screenshot	SFT	Transformer	Mobile/Web
Screen Recognition [Zhang <i>et al.</i> , 2021]	UI modeling	Screenshot	SFT	Faster R-CNN	Mobile
Screen2Words [Wang <i>et al.</i> , 2021]	UI modeling	Screenshot	SFT	Transformer	Mobile
Aria-UI [Yang <i>et al.</i> , 2024b]	UI modeling	Screenshot	Pretrain/SFT	Aria	Mobile/Web/Desktop
Ferret-UI [You <i>et al.</i> , 2024]	UI modeling	Screenshot	Pretrain/SFT	Ferret	Mobile
AutoDroid [Wen <i>et al.</i> , 2024]	End-to-End	HTML	Exploration-based/SFT	Vicuna-7B	Mobile
Seq2Act [Li <i>et al.</i> , 2020]	End-to-End	Texts	Supervised learning	Transformer	Mobile
Meta-GUI [Sun <i>et al.</i> , 2022]	End-to-End	Screenshot/XML	Supervised learning	Transformer	Mobile
Agent Q [Putta <i>et al.</i> , 2024]	End-to-End	Screenshot/DOM	RL/BC Training	Transformer	Web
WebGUM [Furuta <i>et al.</i> , 2024]	End-to-End	Screenshot/HTML	SFT	Flan-T5	Web
CogAgent [Hong <i>et al.</i> , 2024]	End-to-End	Screenshot	SFT	CogVLM	Mobile/Desktop
MobileVLM [Wu <i>et al.</i> , 2024]	End-to-End	XML/Screenshot	Pretrain/SFT	Qwen-VL-Chat	Mobile
WebGPT [Nakano <i>et al.</i> , 2021]	End-to-End	Texts	SFT	GPT-3	Web
AutoGLM [Liu <i>et al.</i> , 2024]	End-to-End	Screenshot/HTML	Pretrain/SFT/RL	ChatGLM	Mobile/Web
OdysseyAgent [Lu <i>et al.</i> , 2024a]	End-to-End	Screenshot	SFT	Qwen-VL	Mobile

to guide browser extensions in executing tasks such as web navigation and form filling, operating purely on the reasoning capabilities of LLMs without additional training. Similarly, WebVoyager [He *et al.*, 2024a] integrates visual and textual information from screenshots and web pages, using prompts to interpret UI elements and execute interactions like clicking and typing. UFO [Zhang *et al.*, 2024a] dynamically generates task plans and executes actions through prompting, allowing it to generalize across diverse web tasks without requiring task-specific adaptations.

Some studies enable the GUI agent to fully leverage external knowledge through prompting to complete GUI tasks.

AppAgent [Zhang *et al.*, 2023] proposes a multi-modal agent framework to simulate human-like mobile phone operations. The framework is divided into two phases: Exploration, where agents explore applications and document their operations, and Deployment, where these documents guide the agent in observing, thinking, acting, and summarizing tasks. This is the first work to claim human-like GUI automation capabilities. AppAgent V2 [Li *et al.*, 2024b] further improves GUI parsing, document generation, and prompt integration by incorporating optical character recognition (OCR) and detection tools, moving beyond the limitations of off-the-shelf parsers for UI element identification. Wang *et al.* [2023] uses a pure in-context learning method to implement interaction between LLMs and mobile UIs. The method divides the conversations between agents and users into four categories from the originator and designs a series of structural CoT prompting to adapt an LLM to execute mobile UI tasks. MobileGPT [Lee *et al.*, 2023b] emulates the cognitive processes of human use of applications to enhance the LLM-based agent with a human-like app memory. MobileGPT uses a random explorer to explore and generate screen-related sub-tasks on many apps and save them as app memory. During the execution, the related memory is recalled to complete tasks.

SFT-based GUI Agents: Supervised fine-tuning (SFT) allows (M)LLMs to adapt to specific domains and perform customized tasks with high efficiency. Recent studies on GUI agents [Wen *et al.*, 2023a; Furuta *et al.*, 2024; Niu *et al.*, 2024; He *et al.*, 2024b; Kil *et al.*, 2024] demonstrate the benefits of SFT for GUI agents to process new modal inputs, learn specific procedures, or execute specialized tasks.

For instance, Furuta *et al.* [2024] proposes WebGUM for web navigation. WebGUM is jointly fine-tuned with an instruction-optimized language model and a vision encoder, incorporating temporal and local perceptual capabilities. The evaluation results on MiniWoB show that WebGUM outperforms GPT-4-based agents. Zhang and Zhang [2023] introduces Auto-UI, a multimodal solution combining an image-language encoder-decoder architecture with a Chain of Actions policy, fine-tuned on the AitW dataset. This Chain of Actions captures intermediate previous action histories and future action plans. Yang *et al.* [2024b] proposes a data-centric pipeline to generate high-quality generalization data from publicly available data. This data is used to fine-tune the VLM for diverse instructions in various environments. Xu *et al.* [2024] introduces a two-stage training paradigm for AGU-VIS. In the first stage, the agent learns visual representations of GUI components through self-supervised learning. In the second stage, it fine-tunes interactive tasks using reinforcement learning, enabling efficient autonomous GUI interaction. On computer-based environments, ScreenAgent [Niu *et al.*, 2024] fine-tunes the ScreenAgent dataset, mapping screenshots to action sequences. It operates via VNC, following a planning-acting-reflecting framework inspired by Kolb’s experiential learning. PC-Agent [He *et al.*, 2024b] employs a multi-agent architecture, fine-tuning a planning agent on cognitive trajectories collected via PC Tracker, enabling it to model human cognitive patterns. Additionally, Kil *et al.* [2024] fine-tunes DeBERTa for element ranking and

Flan-T5 for action prediction, incorporating visual signals to enhance web navigation.

In summary, we provide a systematic overview of recent influential research on (M)LLM-based GUI agents. We address their goal formulations, input perceptions, and learning paradigms, as shown in Table 1

4 Industrial Applications of (M)LLM-Based GUI Agents

GUI agents have been widely used in industrial settings, such as mobile assistants and search agents, demonstrating significant commercial value and potential.

Google Assistant for Android: By saying phrases like “Hey Google, start a run on Example App,” users can use Google Assistant for Android to launch apps, perform tasks, and access content. App Actions, powered by built-in intents (BIIs), enhance app functionality by integrating with Google Assistant. This enables users to navigate apps and access features through voice queries, which the Assistant interprets to display the desired screen or widget.

Apple Intelligence: Apple Intelligence is the suite of AI-powered features and services developed by Apple. This includes technologies such as machine learning, natural language processing, and computer vision that power features like Siri, facial recognition, and photo organization. Apple also integrates AI into its hardware and software ecosystem to improve device performance and user experience. Their focus on privacy means that much of this AI processing happens on-device, ensuring that user data remains secure.

New Bing: Microsoft’s search engine is designed to offer users a more intuitive, efficient, and comprehensive search experience. Leveraging cutting-edge artificial intelligence and machine learning technologies, New Bing goes beyond traditional keyword searches to understand the context and intent behind user queries. With New Bing as an example, the LLM-based deep search engine is also an important form of GUI agents.

Anthropic Computer Use: Anthropic’s “Computer Use” feature enables Claude to interact with tools and manipulate a desktop environment. By understanding and executing commands, Computer-Using Agent(CUA) can perform the necessary actions to complete tasks, much like a human.

OpenAI Operator: OpenAI recently introduced Operator, an AI agent capable of autonomously performing tasks using its own browser. This agent leverages the CUA model, which combines GPT-4o’s vision capabilities with advanced reasoning through reinforcement learning. Operator can interpret screenshots and interact with GUIs—such as buttons, menus, and text fields—just as humans do. This development marks a significant advancement in AI capabilities, enabling more efficient and autonomous interactions with digital interfaces.

Microsoft Copilot: An AI tool in Microsoft 365 apps for productivity with GPT-based suggestions, task automation, and content generation. Enhances workflows, creativity, and decision-making with real-time insights.

AutoGLM: AutoGLM [Liu *et al.*, 2024] is designed for autonomous mission completion via GUIs on platforms like phones and the web. Its Android capability allows it to understand user instructions autonomously without manual input, enabling it to handle complex tasks such as ordering takeout, editing comments, shopping, and summarizing articles.

MagicOS 9.0 YOYO: An advanced assistant with four main features: natural language and vision processing, user behavior learning, intent recognition and decision-making, and seamless app integration. It understands user habits to autonomously fulfill requests, such as ordering coffee through voice commands, by navigating apps and services.

5 Challenges

Due to the rapid development of this field, we summarize several key research questions that require urgent attention:

Personalized GUI Agents: Due to the personal nature of user devices, GUI agents inherently interact with personalized information. As an example, users may commute from home to work during weekdays, while walking to their favorite restaurants and cafes on weekends. The integration of personalized information would clearly enhance the user experience with GUI agents. As the capabilities of (M)LLMs continue to improve, personalized GUI agents have become a priority. Effectively collecting and utilizing personal information to deliver a more intelligent experience for users is an essential topic for future research and applications.

Security of GUI Agents: GUI devices play a crucial role in modern life, making the idea of allowing GUI agents to take control a significant concern for users. For instance, improper operations in financial apps could lead to substantial financial losses, while inappropriate comments on social media apps could damage one’s reputation and privacy. Ensuring that GUI agents are not only highly efficient and capable of generalizing but also uphold user-specific security and provide transparency about their actions is an urgent research challenge. This is a critical issue, as it directly impacts the viability of applying GUI agents in real-world scenarios.

Inference Efficiency: Humans are highly sensitive to GUI response time, which significantly impacts the user experience. Current (M)LLM-based GUI agents still face notable drawbacks with inference latency. Additionally, communication delay is also an important consideration in real-world applications. As a result, efficient device-cloud collaboration strategies and effective device-side (M)LLM research will become critical areas of focus in the future.

6 Conclusion

In this paper, we provide a comprehensive review of the rapidly evolving field of (M)LLM-based GUI Agents. The review is organized into three main perspectives: Data Resources, Frameworks, and Applications. Additionally, we present a detailed taxonomy that connects existing research and highlights key techniques. We also discuss several challenges and propose potential future directions for GUI Agents that leverage foundation models.

References

- [Achiam *et al.*, 2024] Josh Achiam, Steven Adler, et al. Gpt-4 technical report, 2024.
- [Bai *et al.*, 2021] Chen Bai, Xiaoyu Zang, Yan Xu, Sriniwas Sunkara, Abhinav Rastogi, and Jiешan Chen. Uibert: Learning generic multimodal representations for ui understanding, 2021.
- [Chen *et al.*, 2024a] Dongping Chen, Yue Huang, Siyuan Wu, et al. Gui-world: A dataset for gui-oriented multimodal llm-based agents. *arXiv preprint arXiv:2406.10819*, 2024.
- [Chen *et al.*, 2024b] Jingxuan Chen, Derek Yuen, Bin Xie, et al. Spa-bench: A comprehensive benchmark for smartphone agent evaluation. In *NeurIPS 2024 Workshop on Open-World Agents*, 2024.
- [Ding, 2024] Tinghe Ding. Mobileagent: enhancing mobile control via human-machine interaction and sop integration. *arXiv preprint arXiv:2401.04124*, 2024.
- [Furuta *et al.*, 2024] Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, et al. Multimodal web navigation with instruction-finetuned foundation models. In *ICLR*, 2024.
- [Ge *et al.*, 2024] Zhiqi Ge, Juncheng Li, Xinglei Pang, et al. Iris: Breaking gui complexity with adaptive focus and self-refining. *arXiv preprint arXiv:2412.10342*, 2024.
- [He *et al.*, 2024a] Hongliang He, Wenlin Yao, Kaixin Ma, et al. Webvoyager: Building an end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*, 2024.
- [He *et al.*, 2024b] Yanheng He, Jiahe Jin, Shijie Xia, et al. Pc agent: While you sleep, ai works—a cognitive journey into digital world. *arXiv preprint arXiv:2412.17589*, 2024.
- [Hong *et al.*, 2024] Wenyi Hong, Weihang Wang, Qingsong Lv, et al. Cogagent: A visual language model for gui agents. In *CVPR*, pages 14281–14290, 2024.
- [Kil *et al.*, 2024] Jihyung Kil, Chan Hee Song, Boyuan Zheng, Xiang Deng, Yu Su, and Wei-Lun Chao. Dual-view visual contextualization for web navigation. In *CVPR*, pages 14445–14454, 2024.
- [Kim *et al.*, 2023] Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. In *NIPS*, pages 39648–39677, 2023.
- [Lai *et al.*, 2024] Hanyu Lai, Xiao Liu, Iat Long Iong, et al. Autowebglm: A large language model-based web navigating agent. In *SIGKDD*, pages 5295–5306, 2024.
- [Lee *et al.*, 2023a] Kenton Lee, Mandar Joshi, Iulia Raluca Turc, et al. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *ICML*, pages 18893–18912, 2023.
- [Lee *et al.*, 2023b] Sunjae Lee, Junyoung Choi, Jungjae Lee, et al. Explore, select, derive, and recall: Augmenting llm with human-like memory for mobile task automation. *arXiv preprint arXiv:2312.03003*, 2023.
- [Li and Li, 2023] Gang Li and Yang Li. Spotlight: Mobile ui understanding using vision-language models with a focus. In *ICLR*, 2023.
- [Li *et al.*, 2020] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. Mapping natural language instructions to mobile ui action sequences. In *ACL*, pages 8198–8210, 2020.
- [Li *et al.*, 2021] Yang Li, Gang Li, Xin Zhou, Mostafa Dehghani, and Alexey Gritsenko. Vut: Versatile ui transformer for multi-modal multi-task user interface modeling. *arXiv preprint arXiv:2112.05692*, 2021.
- [Li *et al.*, 2024a] Wei Li, William Bishop, Alice Li, et al. On the effects of data scale on computer control agents. *arXiv preprint arXiv:2406.03679*, 2024.
- [Li *et al.*, 2024b] Yanda Li, Chi Zhang, Wanqi Yang, et al. Appagent v2: Advanced agent for flexible mobile interactions. *arXiv preprint arXiv:2408.11824*, 2024.
- [Liu *et al.*, 2018] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. In *ICLR*, 2018.
- [Liu *et al.*, 2024] Xiao Liu, Bo Qin, Dongzhu Liang, et al. Autoglm: Autonomous foundation agents for guis. *arXiv preprint arXiv:2411.00820*, 2024.
- [Lu *et al.*, 2024a] Quanfeng Lu, Wenqi Shao, Zitao Liu, et al. Gui odyssey: A comprehensive dataset for cross-app gui navigation on mobile devices. *arXiv preprint arXiv:2406.08451*, 2024.
- [Lu *et al.*, 2024b] Yadong Lu, Jianwei Yang, Yelong Shen, and Ahmed Awadallah. Omniparser for pure vision based gui agent. *arXiv preprint arXiv:2408.00203*, 2024.
- [Ma *et al.*, 2023] Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, and Dong Yu. Laser: Llm agent with state-space exploration for web navigation. In *NeurIPS Workshop*, 2023.
- [Nakano *et al.*, 2021] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- [Niu *et al.*, 2024] Runliang Niu, Jindong Li, Shiqi Wang, et al. Screenagent: A vision language model-driven computer control agent. *arXiv preprint arXiv:2402.07945*, 2024.
- [Pan *et al.*, 2023] Lihang Pan, Bowen Wang, Chun Yu, Yuxuan Chen, Xiangyu Zhang, and Yuanchun Shi. Autotask: Executing arbitrary voice commands by exploring and learning from mobile gui. *arXiv preprint arXiv:2312.16062*, 2023.
- [Putta *et al.*, 2024] Pranav Putta, Edmund Mills, Naman Garg, et al. Agent q: Advanced reasoning and learning for autonomous ai agents. *arXiv preprint arXiv:2408.07199*, 2024.
- [Qin *et al.*, 2025] Yujia Qin, Yining Ye, Junjie Fang, et al. Ui-tars: Pioneering automated gui interaction with native agents. *arXiv preprint arXiv:2501.12326*, 2025.

- [Rawles *et al.*, 2023] Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy P Lillicrap. Androidinthewild: A large-scale dataset for android device control. In *NIPS Datasets and Benchmarks Track*, 2023.
- [Rawles *et al.*, 2024] Christopher Rawles, Sarah Clinckemaillie, Yifan Chang, et al. Androidworld: A dynamic benchmarking environment for autonomous agents. *arXiv preprint arXiv:2405.14573*, 2024.
- [Shen *et al.*, 2024] Huawei Shen, Chang Liu, Gengluo Li, et al. Falcon-ui: Understanding gui before following user instructions. *arXiv preprint arXiv:2412.09362*, 2024.
- [Sun *et al.*, 2022] Liangtai Sun, Xingyu Chen, Lu Chen, Tianle Dai, Zichen Zhu, and Kai Yu. Meta-gui: Towards multi-modal conversational agents on mobile gui. In *EMNLP*, pages 6699–6712, 2022.
- [Touvron *et al.*, 2023] Hugo Touvron, Thibaut Lavril, Gautier Izacard, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [Venkatesh *et al.*, 2024] Sagar Gubbi Venkatesh, Partha Talukdar, and Srini Narayanan. Ugif-dataset: A new dataset for cross-lingual, cross-modal sequential actions on the ui. In *Findings of NAACL*, pages 1390–1399, 2024.
- [Wang *et al.*, 2021] Bryan Wang, Gang Li, Xin Zhou, Zhouong Chen, Tovi Grossman, and Yang Li. Screen2words: Automatic mobile ui summarization with multimodal learning. In *UIST*, pages 498–510, 2021.
- [Wang *et al.*, 2023] Bryan Wang, Gang Li, and Yang Li. Enabling conversational interaction with mobile ui using large language models. In *CHI*, pages 1–17, 2023.
- [Wang *et al.*, 2024a] Junyang Wang, Haiyang Xu, Haitao Jia, et al. Mobile-agent-v2: Mobile device operation assistant with effective navigation via multi-agent collaboration. In *NIPS*, 2024.
- [Wang *et al.*, 2024b] Junyang Wang, Haiyang Xu, Jiabo Ye, et al. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception. *arXiv preprint arXiv:2401.16158*, 2024.
- [Wang *et al.*, 2024c] Lei Wang, Chen Ma, Xueyang Feng, et al. A survey on large language model based autonomous agents. *FCS*, 18(6):186345, 2024.
- [Wang *et al.*, 2025] Zhenhailong Wang, Haiyang Xu, Junyang Wang, et al. Mobile-agent-e: Self-evolving mobile assistant for complex tasks. *arXiv preprint arXiv:2501.11733*, 2025.
- [Wen *et al.*, 2023a] Hao Wen, Yuanchun Li, Guohong Liu, et al. Empowering LLM to use Smartphone for Intelligent Task Automation. *arXiv preprint arXiv:2308.15272*, 2023.
- [Wen *et al.*, 2023b] Hao Wen, Hongming Wang, Jiaxuan Liu, and Yuanchun Li. Droidbot-gpt: Gpt-powered ui automation for android. *arXiv preprint arXiv:2304.07061*, 2023.
- [Wen *et al.*, 2024] Hao Wen, Yuanchun Li, Guohong Liu, et al. Autodroid: Llm-powered task automation in android. In *MobiCom*, pages 543–557, 2024.
- [Wu *et al.*, 2024] Qinzhuo Wu, Weikai Xu, Wei Liu, et al. Mobilevlm: A vision-language model for better intra-and inter-ui understanding. In *EMNLP*, pages 10231–10251, 2024.
- [Xie *et al.*, 2024] Tianbao Xie, Fan Zhou, et al. Openagents: An open platform for language agents in the wild. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.
- [Xu *et al.*, 2024] Yiheng Xu, Zekun Wang, Junli Wang, et al. Aguis: Unified pure vision agents for autonomous gui interaction. *arXiv preprint arXiv:2412.04454*, 2024.
- [Yan *et al.*, 2023] An Yan, Zhengyuan Yang, Wanrong Zhu, et al. Gpt-4v in wonderland: Large multimodal models for zero-shot smartphone gui navigation. *arXiv preprint arXiv:2311.07562*, 2023.
- [Yang *et al.*, 2024a] An Yang, Baosong Yang, et al. Qwen2 technical report, 2024.
- [Yang *et al.*, 2024b] Yuhao Yang, Yue Wang, Dongxu Li, et al. Aria-ui: Visual grounding for gui instructions. *arXiv preprint arXiv:2412.16256*, 2024.
- [You *et al.*, 2024] Keen You, Haotian Zhang, Eldon Schoop, et al. Ferret-ui: Grounded mobile ui understanding with multimodal llms. In *ECCV*, pages 240–255, 2024.
- [Zhang and Zhang, 2023] Zhuosheng Zhang and Aston Zhang. You only look at screens: Multimodal chain-of-action agents. *arXiv preprint arXiv:2309.11436*, 2023.
- [Zhang *et al.*, 2021] Xiaoyi Zhang, Lilian De Greef, Amanda Swearngin, et al. Screen recognition: Creating accessibility metadata for mobile applications from pixels. In *CHI*, pages 1–15, 2021.
- [Zhang *et al.*, 2023] Chi Zhang, Zhao Yang, Jiaxuan Liu, et al. AppAgent: Multimodal Agents as Smartphone Users. *arXiv preprint arXiv:2312.13771*, 2023.
- [Zhang *et al.*, 2024a] Chaoyun Zhang, Liqun Li, Shilin He, et al. Ufo: A ui-focused agent for windows os interaction. *arXiv preprint arXiv:2402.07939*, 2024.
- [Zhang *et al.*, 2024b] Jiwen Zhang, Jihao Wu, Yihua Teng, et al. Android in the zoo: Chain-of-action-thought for gui agents. *arXiv preprint arXiv:2403.02713*, 2024.
- [Zhao *et al.*, 2024] Kangjia Zhao, Jiahui Song, Leigang Sha, et al. Gui testing arena: A unified benchmark for advancing autonomous gui testing agent. *arXiv preprint arXiv:2412.18426*, 2024.
- [Zheng *et al.*, 2024] 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.
- [Zhou *et al.*, 2023] Shuyan Zhou, Frank F Xu, Hao Zhu, et al. Webarena: A realistic web environment for building autonomous agents. In *ICLR*, 2023.
- [Zhu *et al.*, 2024] Zichen Zhu, Hao Tang, Yansi Li, et al. Moba: A two-level agent system for efficient mobile task automation. *arXiv preprint arXiv:2410.13757*, 2024.