# AppCopilot: Toward General, Accurate, Long-Horizon, and Efficient Mobile Agent

*https://github.com/OpenBMB/AppCopilot*

Jingru Fan★†, Yufan Dang♠†, Jingyao Wu★, Huatao Li★, Runde Yang★, Xiyuan Yang★, Yuheng Wang★, Zhong Zhang♠, Yaxi Lu♠, Yankai Lin♣, Zhiyuan Liu♠, Dahai Li♦, Chen Qian★✉

★School of Artificial Intelligence, Shanghai Jiao Tong University

♠Department of Computer Science and Technology, Tsinghua University

♣Gaoling School of Artificial Intelligence, Renmin University of China

♦Modelbest Inc.

fanjingru510@sjtu.edu.cn, dangyf21@mails.tsinghua.edu.cn, qianc@sjtu.edu.cn

†Equal Contribution

✉Corresponding Author

# Abstract

With the raid evolution of large language models and multimodal foundation models, the mobile-agent landscape has proliferated without converging on the fundamental challenges. This paper identifies four core problems that must be solved for mobile agents to deliver practical, scalable impact: (1) generalization across tasks, modalities, apps, and devices; (2) accuracy, specifically precise on-screen interaction and click targeting; (3) long-horizon capability for sustained, multi-step goals; and (4) efficiency, specifically high-performance runtime on resource-constrained devices.

We present AppCopilot, a multimodal, multi-agent, general-purpose on-device assistant that operates across applications and constitutes a full-stack, closed-loop system from data to deployment. AppCopilot operationalizes this position through an end-to-end autonomous pipeline spanning data collection, training, deployment, high-quality and efficient inference, and PC/mobile application development. At the model layer, it integrates multimodal foundation models with robust Chinese–English support. At the reasoning and control layer, it combines chain-of-thought reasoning, hierarchical task planning and decomposition, and multi-agent collaboration. At the execution layer, it enables user personalization and experiential adaptation, voice interaction, function/tool calling, cross-app and cross-device orchestration, and comprehensive mobile app support. The system design incorporates profiling-driven optimization for latency, memory, and energy across heterogeneous hardware. Empirically, AppCopilot achieves significant improvements along all four dimensions: stronger generalization, higher-precision on-screen actions, more reliable long-horizon task completion, and faster, more resource-efficient runtime.

By articulating a cohesive position and a reference architecture that closes the loop from "data collection—training and deployment—high-quality, efficient inference—application development", this paper offers a concrete roadmap for general-purpose digital assistants and provides actionable guidance for both academic research and industrial adoption.

For updates and forthcoming releases, see the project page: https://github.com/OpenBMB/AppCopilot.

**Keywords:** Multimodal Foundation Model, Large Language Model, Autonomous Agent, Multi-Agent Collaboration, Mobile Agent, GUI Agent

# Contents

# 1
# Research Background

## 1.1 From Models to Agents

Language is at the core of human intelligence, serving as the fundamental bond that underpins complex human social collaboration systems [54]. From ancient oral communication to modern highly abstract written language, programming languages, and symbol systems, language is more than just an information transfer tool; it profoundly shapes human cognition, knowledge organization, and multi-agent cooperation paradigms. In high-complexity scenarios such as software development, legal proceedings, education, scientific research, and business negotiations, language acts as a "universal interface" for knowledge and intent, enabling individuals from different backgrounds to understand each other, share knowledge, and collaboratively achieve goals. Therefore, an efficient language interaction mechanism directly determines the capacity for interdisciplinary knowledge fusion and indirectly drives social innovation and technological progress. However, even in human societies, cross-professional collaboration faces numerous challenges [55]: differences in knowledge domains often lead to misunderstandings, the uncertainty of language can cause miscommunication, and asymmetries in time, attention, and intent among collaborators frequently result in information silos, communication delays, and even team failure. Under real-world constraints of limited resources and time, these "unstructured bottlenecks" become critical issues that hinder efficient knowledge collaboration and limit the quality of multi-domain solutions. These problems are not only prevalent in reality but have also become a core challenge in the design of intelligent systems.

Fortunately, the emergence of Large Language Models (LLMs) like ChatGPT [6] has brought unprecedented possibilities for overcoming these collaborative difficulties. By training on vast amounts of text, LLMs have demonstrated powerful language understanding, knowledge reasoning, and cross-domain generalization capabilities. They can handle traditional natural language processing tasks (e.g., Q&A, translation, summarization) and possess "expert-like" general abilities in tasks like code generation, mathematical reasoning, scientific inquiry, and document creation. Research from institutions like Stanford University highlights the significance of LLMs for achieving artificial general intelligence [5, 102]. China's "Next Generation Artificial In-

telligence Development Plan" (2017) explicitly lists natural language processing as a national key technology direction, reflecting a global consensus on the intelligent potential of language models. Furthermore, the "Agent" paradigm developed upon LLMs marks a leap from static Q&A tools to dynamic, autonomous, perception-enabled, and actionable systems. For example, systems like AutoGPT [100], BabyAGI [1], and AgentVerse [9] can now leverage the reasoning power of language models to interact with external environments and automatically plan and execute complex tasks. This shift not only enhances a model's adaptability in open environments but also expands its capabilities in task decomposition, tool use, and result verification, evolving language models from "passive responders" to "active executors" and laying the technical foundation for human-AI co-creation and collaborative intelligence.

However, when facing truly complex tasks, a single agent still has significant limitations. Complex tasks often involve multiple sub-goals, multi-dimensional information inputs, multi-stage execution steps, and an environment characterized by high dynamism and uncertainty. These task features are particularly prominent in natural language environments. For instance, personalized student tutoring in educational systems, multi-disciplinary joint diagnosis in healthcare, and multi-device coordination in industrial control cannot be completed by a single model alone within a reasonable timeframe. In this context, introducing a "Multi-Agent System" (MAS) becomes a crucial strategy to overcome a model's capacity bottlenecks. Multiple agents with heterogeneous capabilities and specialized roles can simulate human organizational collaboration, achieving greater adaptability and task efficiency through task decomposition, distributed execution, and language-based interaction. Multi-agent collaborative systems, by simulating the organization and communication mechanisms of human teams, are better equipped to handle issues of heterogeneous knowledge, diverse perspectives, and conflicting goals. Agents can negotiate role assignments, exchange intermediate results, and coordinate action plans based on a shared language, dynamically adjusting their strategies in response to environmental changes. This mechanism not only improves the system's scalability and robustness but also lays the foundation for achieving autonomous, efficient, and general task-oriented agency.

Against the backdrop of continuous development in agent technology, agents are significantly enhancing task automation and intelligence by understanding natural language instructions, integrating external tools, and building collaborative systems composed of multiple heterogeneous agents. However, current mainstream agent systems still largely focus on language interaction and task execution within the text modality. They face significant challenges when dealing with large-scale complex tasks that involve multi-modal information such as vision, spatial location, and actions. This is especially true in highly realistic operating environments like Graphical User Interfaces (GUIs), where agents need to simultaneously perceive dynamic visual states, understand spatial structures, track control functions, and adjust action sequences based on real-time feedback. This far exceeds the capabilities of traditional language-based systems.

To address these challenges, a new type of system, the mobile agent (also referred to as GUI agent), has recently emerged. These agents aim to break the application boundaries of traditional language models in multi-modal tasks. A **mobile agent** is an intelligent agent that can perceive

---

[1] https://github.com/yoheinakajima/babyagi

and understand graphical user interfaces and autonomously execute user instructions in various scenarios, including mobile, PC, and Web. By integrating the natural language understanding and planning capabilities of LLMs with the control recognition and visual perception abilities of image models, they can perform reasoning operations within a graphical interface without relying on structured interfaces. Specifically, a mobile agent can identify and understand interface controls like buttons, text boxes, and menu items, perceive their spatial location and state, and then generate an executable action plan based on the task intent, achieving precise manipulation of the GUI environment. Compared to traditional Application Programming Interface (API) tool calls, this method more closely aligns with human usage habits, offers good cross-platform adaptability, and maintains a high success rate and task stability even when interfaces are unavailable or unpredictable. The advent of mobile agents represents a breakthrough in the deep integration of language models and visual perception, laying a new foundation for the comprehensive development of agent systems from "can speak" to "can see" and "can act." However, as a new system still in its rapid development phase, its capabilities are not yet fully mature when facing real-world multi-modal complex tasks, and it is still constrained by various factors, including insufficient generalization, limited long-range capabilities, inadequate accuracy, and low efficiency.

Given this context, research into intelligent agent systems is gradually moving from rule-driven to model-driven and from isolated responses to multi-party collaboration, making more general, open, and flexible intelligent collaboration architectures a realistic necessity. In line with this trend, this project focuses on how to integrate the language understanding and reasoning capabilities of LLMs, the environmental perception abilities of multi-modal inputs, and the autonomous behavior mechanisms of agents. We aim to explore a new system paradigm for multi-end and cross-application collaboration. Specifically, this project proposes building a multi-modal, multi-agent-driven, general-purpose on-device intelligent system. Its goal is to expand the interaction boundaries of large models, endowing them with structured collaborative capabilities and environmental adaptability for complex real-world tasks. This project will review the current research status and development trends in LLMs, autonomous agents, multi-agent collaboration systems, and mobile agents, systematically analyze existing problems, propose corresponding technical solutions, and validate them through quantitative experiments. Finally, we will provide future development directions and organize and provide the models, scripts, and software resources used.

### 1.1.1 Large Language Models

In recent years, **Pretrained Language Models** have become the dominant technological paradigm in Natural Language Processing (NLP), profoundly driving the key process of artificial intelligence's evolution from perceptual intelligence to cognitive intelligence. This paradigm primarily consists of two stages: first, self-supervised pre-training on massive amounts of unlabeled text data to learn linguistic statistical structures and basic knowledge graphs; then, fine-tuning with limited supervision for specific tasks, enabling the model to adapt to particular downstream applications. Since the advent of the Transformer architecture [85], representative models like

the GPT series, BERT [19], T5 [67], and BART [38] have emerged, forming the current mainstream model family. These models can be structurally categorized into three types: autoregressive models (e.g., ChatGPT, InstructGPT [59], GPT-4, PaLM [14], DeepSeek [47]), autoencoding models (e.g., BERT, RoBERTa [49]), and sequence-to-sequence models (e.g., T5, BART), each excelling at different types of language tasks.

With the continuous increase in model parameter scale, pretrained language models have gradually evolved into LLMs. In late 2022, OpenAI's release of ChatGPT demonstrated unprecedented capabilities in language understanding and generation, human-machine dialogue, and code generation, marking the formal establishment of the autoregressive model-dominated LLM paradigm. Since then, LLM research has accelerated globally. Commercial closed-source models such as Claude and GPT-4 continue to push the boundaries of capability, while open-source models like LLaMA [83], Alpaca [80], GLM [21], and Qwen [2] have exploded, driving a new wave of competition with multi-modal large models (e.g., Gemini [82], V-JEPA [4], Sora [50]). Notably, in 2025, OpenAI officially released the open-weight GPT-OSS [56] series, signaling its attempt to gradually move toward open-source community collaboration while ensuring safety. Overall, model scale is widely regarded as a key factor influencing performance. Models with larger parameter scales often provide better performance, a phenomenon known as the "Scaling Law" [35]. As training data volume and model parameters continuously increase, foundation models begin to exhibit "Emergent Abilities" [91]—that is, when a model's scale reaches a certain threshold, it naturally demonstrates advanced intelligent characteristics like zero-shot learning, cross-task generalization, and logical reasoning [90].

Along with technological evolution, research on optimizing performance and aligning behavior for practical applications of LLMs has also developed rapidly. First, how to achieve more efficient deployment and fine-tuning while maintaining performance for massive model parameters has become a significant topic. To this end, the academic community has proposed Parameter-Efficient Tuning (PEFT) techniques [30], which achieve performance close to full fine-tuning through the tuning of a limited subset of the model's parameters. Methods like LoRA [30], Prefix Tuning [45], and Adapter [29] are widely adopted. Additionally, Instruction Tuning guides models to understand and execute natural language instructions, enabling stronger task generalization in zero-shot or few-shot conditions, laying the foundation for the general interactivity of LLMs. To make a model's generated content more aligned with human expectations, Reinforcement Learning from Human Feedback (RLHF) [27] has been introduced. This mechanism guides the model to follow human value preferences and reasoning logic during generation, significantly enhancing its controllability and social acceptability.

Beyond basic capabilities, LLMs have shown extraordinary potential in complex task reasoning. Traditional symbolic reasoning systems rely on explicit rules, which struggle to adapt to the uncertainty and semantic ambiguity of open environments. In contrast, LLMs can achieve flexible "pseudo-reasoning" processes based on natural language. Prompt Learning [6], as a core mechanism, allows users to guide the model to complete target tasks by constructing specific contextual prompts. Furthermore, Chain-of-Thought (CoT) reasoning [90] breaks down complex problems into manageable sub-problems by constructing a multi-step logical chain, significantly enhancing the model's interpretability and reasoning stability. This technique has been

widely applied in high-cognitive tasks like mathematical reasoning, causal inference, and decision planning, becoming an important path for models to advance toward higher intelligence.

However, despite the groundbreaking progress of LLMs on various tasks, they still have significant limitations, mainly in the following areas:

1. LLMs are essentially static response systems. Their interaction process is typically limited to a "single-turn input, single-turn output" mode. They lack dialogue memory, state maintenance, and dynamic behavior planning capabilities, making them unable to handle complex interactive tasks that require continuous decision-making and context preservation.

2. Traditional LLMs lack perception and control interfaces with external environments. They cannot acquire information, call tools, or perform operations in the physical world, which makes them inherently inadequate for tasks requiring real-time observation and execution feedback (e.g., web browsing, robot control).

3. LLMs are typically designed as general-purpose answerers, not as "agents" with autonomous goals, planning, and behavioral capabilities. They cannot decompose tasks based on long-term goals, manage resources, or coordinate strategies. They also lack the structured ability for multi-role, multi-step collaboration to complete a task.

It is based on these technical limitations that academia and industry are gradually exploring how to embed LLMs as the core intelligence within an "agent" framework. This would give them the ability to think, act, and collaborate, thereby overcoming the bottlenecks of traditional language models in perception, task planning, and environmental adaptation. This agent-ization not only expands the boundaries of language models but also provides a key opportunity for building a new generation of cognitive intelligent systems.

### 1.1.2   Autonomous Agents

Amid the rapid development of LLMs, the concept of the **autonomous agent** is widely considered a natural extension of their capabilities, and a crucial step in the models' transition from "Q&A tools" to "cognitive agents." Although LLMs excel at tasks like language understanding, summarization, Q&A, and code generation, their core remains a passive language system. They lack goal awareness, behavioral control, and the ability to execute complex decisions, making it difficult for them to independently complete full task workflows in open environments. To overcome this limitation, academia and industry have proposed the "autonomous agent" paradigm, which seeks to endow language models with external tools and behavioral planning capabilities, enabling them to complete multi-stage tasks in a more intelligent and autonomous way.

The basic idea of an autonomous agent is to use the large language model as the central "cognitive hub," combining it with peripheral modules such as task planning, tool interfaces, memory systems, and role-playing. Together, they form a closed-loop system with the ability to perceive, think, act, and self-reflect. This type of system can autonomously formulate plans, call external tools, perceive environmental feedback in real time, and adjust its process based on a task goal, thereby completing complex tasks with minimal human intervention. For example,

AutoGPT uses multi-round reasoning and task decomposition to automatically execute software construction; the ReAct framework [103] combines reasoning chains and action chains, allowing the model to think while acting; and Reflexion [76] introduces a self-reflection mechanism after task failure, enabling it to review and correct error paths. These advancements indicate that LLMs are transforming from "language generators" into "task agents" with perception and execution capabilities.

From an architectural perspective, autonomous agents are highly modular. The language model itself is responsible for semantic understanding and high-level reasoning, while peripheral modules supplement its capabilities in environmental interaction, information processing, and state maintenance. For example, a planning module is responsible for breaking down tasks into multi-step execution plans; a tool-use module allows the agent to access external tools like search engines, databases, compilers, etc.; a memory module provides context storage and historical state retrieval capabilities; and a role-playing module defines the agent's behavior style and professional domain. These modules work together to build a closed-loop structure around the language model, turning it from a passive conversational system into a composite agent with a degree of autonomous behavior.

Although this architecture has shown initial success in several applications, there are still significant bottlenecks. The current planning and actions of agents heavily rely on prompt engineering and manual strategy settings, making them difficult to generalize. Memory mechanisms are not yet efficiently persistent, facing issues of context invalidation and information redundancy. Tool interfaces are also mostly based on pre-set functions or API calls, lacking task-driven dynamic adaptation capabilities. Moreover, it is particularly noteworthy that when task scenarios are highly complex, information is heterogeneously distributed, or task structures are nonlinear, a single agent can hardly handle the full workflow. This is not only due to the inherent limitations of a single agent's perception and knowledge but also its lack of the ability to integrate multiple perspectives and perform collaborative reasoning.

**Related Technologies**

To support autonomous agents in completing complex tasks, their architecture is typically built around four key technical modules: Thinking and Planning, Tool Use, Memory Mechanisms, and Role-Playing. These modules, in synergy with the large language model, provide the agent with a complete capability chain, from understanding the problem and formulating a plan to calling resources and generating feedback. Figure 1.1 shows a typical autonomous agent architecture.

1. **Thinking and Planning Module**: To execute complex tasks, an autonomous agent must possess a deep understanding of the problem structure and a clear plan of action. In this process, "task decomposition" and "self-reflection" play crucial roles. Task decomposition techniques enable the agent to break down high-complexity goals into a series of logically clear, easy-to-execute sub-tasks, thereby more efficiently solving problems step-by-step through a "slow thinking" reasoning approach. For example, using a "Chain-of-Thought" prompting structure can guide the model to reason incrementally, while "Tree-of-Thought"

**Figure 1.1:** *Autonomous Agent Framework*

explores different paths to build a reasoning structure, effectively expanding the agent's decision space and ability to handle complex situations. Building on this, the self-reflection mechanism further enhances the model's execution robustness, allowing it to review its historical trajectory to identify and correct errors when tasks fail or deviate, thus continuously optimizing task strategies. This integrated mechanism of task planning and self-correction provides the agent with an execution process that is close to human cognitive paths.

2. **Tool-Use Module**: The ability to use tools is a key mark of the large language model's transition to an agent. It significantly enhances the model's capacity to interact with the real world, overcoming the inherent limitations of language models in fact-checking, precise calculation, and logical execution. Agent systems integrate various external tools, such as search engines, translation engines, knowledge graphs, databases, compilers, Web browsers, and calculators. This allows the model to not only generate text but also actually call resources to complete task execution. During a task, the agent can decide when to use which tool based on the context and generate call instructions in natural language, thereby achieving dynamic interaction with the external environment. The model can also adjust its strategy and provide content feedback based on the call results, realizing a closed-loop execution of decisions. In some tasks, agents can even self-iteratively call the same or multiple models to collaboratively solve complex problems, improving the system's robustness, security, and response quality. This "tool-augmented agent" architecture has been widely validated in practical applications. Systems like AutoGPT, Manus[2], Toolformer [71], and

---

[2] https://manus.im/

WebGPT [102] have demonstrated the immense potential of deeply integrating models with tools.

3. **Short-term and Long-term Memory Module**: An efficient and stable memory system is fundamental for autonomous agents to maintain context, manage long-term tasks, and continuously accumulate knowledge. Inspired by human cognitive structure, agents typically adopt a parallel architecture of short-term and long-term memory: the former relies on the language model's own context window to handle the current interaction state and temporary information; the latter uses a vector database to store and retrieve long-term knowledge and historical trajectories. In practical systems, long-term memory content is converted into high-dimensional semantic embeddings. Approximate Nearest Neighbor (ANN) retrieval techniques, such as Maximum Inner Product Search (MIPS), are used to achieve efficient information matching and recall. Common implementation methods include Locality-Sensitive Hashing (LSH), HNSW small-world graph structures, and FAISS vector quantization algorithms [115]. These methods strike a balance between recall rate and response efficiency while significantly enhancing large-scale knowledge retrieval and historical state restoration capabilities. By combining multi-source memory with dynamic update mechanisms, agents can maintain contextual coherence and behavioral consistency in multi-turn dialogues, long-term tasks, or cross-task scenarios.

4. **Role-Playing Module**: The role-playing capability provides autonomous agents with a highly personalized path for behavioral expression, allowing them to exhibit professional identity and style in specific tasks. By setting a clear role for an agent, such as a programmer, financial analyst, or educator, it not only augments the model's ability to activate knowledge in a specific domain but also regulates its language style, interaction methods, and task preferences. Common construction strategies include manually writing role profiles, using a language model to automatically generate role settings, and creating a persona through real-world data alignment. The manual method has the advantage of high precision and controllability, making it suitable for critical task deployments. The model-generated approach is better for batch construction and dynamic settings, although it has some uncertainty in accuracy. The data alignment method uses real user behavior data to create semantic embeddings, which greatly enhances the realism and task relevance of the role. In a multi-agent system, role-playing can further support task division and collaborative interaction among agents, forming the basis for a collective intelligence structure.

### 1.1.3 Multi-Agent Collaboration

With the continuous expansion of LLM capabilities, their performance in single-agent tasks has become quite mature. However, when faced with complex, dynamic, and heterogeneously distributed tasks, a single agent struggles to cover all processes. It is limited by a restricted cognitive perspective, tight task coupling, and a lack of behavioral flexibility. Therefore, to further enhance a system's task adaptability, scenario generalizability, and collaborative efficiency, research is shifting from building single, fully-functional agent systems to creating collaborative systems composed of multiple agents with different capabilities and specializations—that

is, **Multi-Agent Systems**. Through task division, negotiation, communication, and game theory among agents, it becomes possible to process complex tasks in parallel and complete them efficiently, which is a key path for moving LLMs toward higher levels of intelligence.

Multi-agent systems are not a new concept. Traditional research often used centralized or distributed control architectures, combined with methods such as reinforcement learning, game modeling, policy transfer, communication protocols, and hierarchical learning, to optimize individual agents' perception, reasoning, and decision-making capabilities, thus building systems that can collaborate effectively on specific tasks. However, these systems typically rely on external signal-based supervised training and lack a general language interaction mechanism, which limits their adaptability and generalization in open-ended tasks and dynamic collaborative environments.

The introduction of LLMs has significantly changed this paradigm. By using natural language interaction as a unified coordination mechanism, LLM-driven multi-agent systems have achieved a deep coupling of cognition and collaboration. On one hand, agents can use language for efficient knowledge sharing, task synchronization, and decision transmission, reducing collaboration costs. On the other hand, the generalization capabilities of LLMs enable the system to show stronger adaptability when facing unknown scenarios and goals. This new paradigm not only enhances the naturalness and scalability of collaboration but also reduces reliance on manual rules and supervised signals, providing a feasible path for building general, multi-modal, multi-task collective intelligence systems.

A typical LLM-driven multi-agent system framework is shown in Figure 1.2:



**Figure 1.2:** *Multi-agent System Framework*

**Types of Multi-Agent Systems**

Based on their group goals, LLM-driven multi-agent systems are typically divided into the following two types:

1. **Social Simulation [39]**: These systems primarily simulate social behaviors or group dynamics rather than targeting predefined tasks. They use language-driven agents to recreate real-world social interactions in a virtual environment to explore social science questions or validate theories. For example, the Stanford SmallVille [61] uses 25 agents in a virtual sandbox town to simulate social behavior, demonstrating the potential of LLMs to enhance agent autonomy and behavioral complexity. By retrieving an agent's experiences and a periodic "reflection" mechanism, agents are prompted to generate high-level abstract thoughts based on recent experiences, recursively creating details and continuously adapting to the environment. The AgentVerse framework, on the other hand, dynamically adjusts agent combinations and observes emergent group behaviors. Experiments show that a collaborative multi-agent combination is more efficient than a single agent and can spontaneously exhibit social behaviors like willingness, obedience, and destruction. These studies indicate that language-driven multi-agent systems have immense potential in simulating complex social interactions and decision-making.

2. **Task-Solving**: These systems have a clear end goal and emphasize the efficient completion of tasks through multi-agent collaboration, not just simple social simulation. For example, ChatDev [63] builds a virtual software development team with multiple roles, enabling design, coding, and testing to communicate through natural language, forming a complete software development workflow. This collaborative framework not only significantly improves development efficiency but can also find and fix code issues during the interaction, demonstrating the potential of LLM collaboration in software development and inspiring subsequent task-oriented research. MetaGPT [28] further encodes standardized operating procedures in the prompts, breaking down complex tasks into a pipeline of subtasks. By assigning roles and using programmatic planning, it reduces the occurrence of logical inconsistencies. AutoGen [94] provides a general multi-agent conversational orchestration framework that supports the creation, configuration, and management of multiple agents that can communicate and collaborate. This makes it easy to build complex task execution workflows, such as data analysis, code generation, and debugging, and achieves greater flexibility and scalability in open-ended tasks through programmable interaction logic. These frameworks prove that in complex tasks with clear goals, LLM-driven multi-agent collaboration can significantly improve efficiency and result quality.

### 1.1.4 Mobile Agents

As the demand for complex model interaction environments grows, the tool-use capabilities of large models are gradually extending to the visual modality. Compared to the text modality, visual tool learning is closer to real-world interaction, enabling complex task manipulation directly through the perception and understanding of visual information from graphical inter-

faces or real-world environments. In contrast to traditional scripted automation or rule-based interface control, a mobile agent can understand complex task intentions in both visual and linguistic modalities and automatically generate cross-application, multi-step interaction sequences. This brings new opportunities for fields like software testing, human-machine collaboration, and accessibility technology. However, when it comes to generalizing control over truly complex interfaces, current research is still in its early stages, facing significant challenges in interface perception accuracy, cross-platform adaptation, and security.

International research institutions have been at the forefront of this exploration. OpenAI, building on ChatGPT, launched a prototype for web and desktop operations called "Operator," demonstrating the language model's ability to automatically fill out forms and perform web searches through control recognition. Google DeepMind's AndroidEnv [84] uses reinforcement learning to train agents to automate operations on Android apps, and it uses deep networks to integrate visual information with environmental states, enhancing the feasibility of deployment on real devices. Microsoft Research's UFO (UI-Focused Agent) [108] combines Windows Automation interfaces with LLMs to perform multi-step operations and published a survey, "LLM-Brained Mobile Agents" [107], which systematically reviews the research path of mobile agents. Additionally, UC Irvine's RCI [36] (Recursive Criticize and Improve) method achieves a high task completion rate in web interface tasks like MiniWob with minimal instructions. The Aguvis [98] system, jointly developed by Salesforce and Columbia University, uses a purely visual approach for cross-platform GUI manipulation and shows good generalization. The Screen2Words model [86] is dedicated to converting mobile interfaces into readable text descriptions for accessibility.

Keeping pace with international research, China is also showing rapid development in the field of mobile agents. In industry, numerous mobile phone manufacturers like Huawei, Apple, Honor, OPPO, vivo, and Xiaomi, as well as large model companies like Alibaba, Zhipu, and Minabai, are actively deploying related technologies and products, continuously expanding the capabilities and application scenarios of on-device equipment. Tencent's AppAgent [109] focuses on typical mobile operations like app installation, page navigation, and information input, using a combination of visual and language modalities to achieve multi-step control. ByteDance has developed a cross-platform GUI operation model software, attempting to enable large models to achieve consistent interface perception and decision-making on various terminals such as mobile phones, PCs, and the Web [66]. The academic community has also produced numerous results. Tsinghua University proposed AutoWebGLM [37], which integrates the GLM model with web interaction data. It builds a diverse web scenario dataset through "human-computer collaboration + self-supervised learning" to enhance the model's adaptability in real networks. The Renmin University-Tsinghua team proposed the GUICourse framework [10], which builds a multimodal dataset with three subsets: GUIEnv, GUIAct, and GUIChat. It injects GUI functional knowledge and contextual dialogue capabilities into a pretrained vision-language model, significantly improving the model's generalization on multi-tasks. The Shanghai Jiao Tong University-OPPO team's MobileUse system [42] introduces a hierarchical reflection mechanism and active exploration strategy, significantly reducing perception and planning errors in the Android-World/AndroidLab environments. Microsoft and Nanjing University's collaborative GUI-Actor

model [93] introduces an attention mechanism at the end of the VLM architecture to achieve patch-level action alignment, effectively improving cross-resolution operation accuracy. Sense-Time [33] built the SpiritSight system, proposing the GUI-Lasagne ultra-large-scale dataset and a Universal Block Parsing (UBP) mechanism, which addresses the issue of component location ambiguity in high-resolution visual input, further enhancing the localization ability of GUI objects.



**Figure 1.3:** *Development Timeline of Mobile Agents [107]*

Despite rapid research progress, mobile agents still face numerous challenges: recognition deviations and action failures are common in dynamic interfaces, model deployment performance on edge devices does not yet meet real-time requirements, secure read/write mechanisms for user interface information are incomplete, and task transfer and generalization capabilities are still insufficient. Most current achievements are concentrated on specific applications or benchmark scenarios and have not yet truly achieved general intelligence for diverse GUI environments. However, these bottlenecks also create vast opportunities for future innovation.

Based on this background, this project aims to leverage the advantages of multi-agent collaboration to build a multi-modal, multi-agent-driven, cross-application and cross-device on-device assistant. Through a unified language interaction mechanism and dynamic role-playing, this system supports adaptive task decomposition and collaborative work in different application scenarios (e.g., text, voice input). It can accumulate experience and adjust strategies across tasks and enable collaborative interaction and data exchange between multiple different devices. This not only enhances the system's adaptability and scalability but also meets users' cross-scenario and personalized needs, driving intelligent assistants from single-task processing toward a more complex, cross-domain collaboration model.

## 1.2　Related Work

### 1.2.1　Representative Mobile Agents

Mobile agents are a rapidly developing branch of the agent market. A mobile agent is defined as a highly automated agent that can understand user instructions in natural language, analyze the GUI interface and its elements, and autonomously perform corresponding operations without relying on complex, platform-specific scripts or pre-defined workflows [106]. Such agent systems utilize LLMs as their core reasoning and cognitive engines to generate, plan, and execute actions in a flexible and adaptive manner. This section 1.2.1 introduces several mature mobile agent foundation models and technical products, analyzing their model capabilities and evaluation results across different dimensions.

**UI-TARS**

UI-TARS [66] is an open-source, purely vision-driven mobile agent released by ByteDance Seed in March 2025. It directly takes a screenshot as input and outputs low-level instructions, such as cursor operations, to directly interact with the interface. Unlike existing frameworks that rely on closed source models and complex prompt engineering and workflows, the core innovation of this model lies in its end-to-end architecture, which outperforms more complex designs. UI-TARS introduces enhanced perception capabilities, achieving context-aware understanding and precise descriptions of UI elements through a large-scale GUI screenshot dataset. Its unified action modeling standardizes cross-platform actions and enables precise positioning and interaction through a massive action trajectory dataset. Additionally, the model incorporates a multi-step thinking mode, including task decomposition, reflection, and task milestone recognition, allowing it to handle complex multi-step tasks. To address the scarcity of training data, UI-TARS employs self-iterative learning, collecting and filtering data on hundreds of virtual machines to achieve self-growth. In multiple GUI benchmarks, UI-TARS's performance surpasses several leading closed-source models, including GPT-4o and Claude.

**Qwen-VL**

Qwen-VL [3] consists of a collection of large-scale vision-language models developed by Alibaba Group. Based on the Qwen-LM, it is equipped with powerful visual capabilities through a vision encoder, input-output interfaces, a three-stage training process, and a multi-lingual, multi-modal cleaned dataset. In addition to traditional image description and question-answering functions, Qwen-VL also achieves advanced functions like visual localization and text recognition by aligning images, text descriptions, and bounding boxes. The models in this series, including Qwen-VL and Qwen-VL-Chat, outperform general-purpose models of the same scale in various visual benchmarks (such as image description, QA, and visual localization) and different settings (e.g., zero-shot, few-shot). The instruction-optimized dialogue model, Qwen-VL-Chat, demonstrates performance superior to existing vision-language chatbots in real dialogue scenarios. Qwen-VL can accept images, text, and bounding boxes as input and output text, and

bounding boxes. Its core functions cover general image understanding, multi-image dialogue, fine-grained OCR, and strong generalization ability.

**AgentCPM-GUI**

AgentCPM-GUI [113] is an open-source, on-device large agent model jointly developed by Tsinghua University's Natural Language Processing and Social Humanities Computing Laboratory (TsinghuaNLP), Renmin University of China, and ModelBest. This model can take a mobile phone screen image as input and automatically execute user-provided tasks, showing strong performance in Chinese GUI scenarios. For existing challenges including data noise, poor generalization, and insufficient support for non-English applications, AgentCPM-GUI introduces several key innovations. First, it constructs a high-quality mixed Chinese-English training dataset with targeted data collection and strict deduplication, significantly enhancing its cross-lingual and cross-application generalization capabilities. Second, the model uses a progressive training pipeline, including perception-centric pre-training, supervised fine-tuning on a collection of high-quality multilingual data and reinforcement fine-tuning with the GRPO algorithm, systematically strengthening its perception, imitation, and reasoning abilities. Additionally, AgentCPM-GUI features a lightweight design for mobile devices, using a compact JSON action format to compress the average output length and improve on-device execution efficiency. By pre-training on a large-scale Chinese-English Android dataset, AgentCPM-GUI significantly improves the recognition and localization accuracy of various GUI controls. As the first mobile agent finely optimized for Chinese apps, it supports over 30 mainstream mobile applications, including AutoNavi Map, Dianping, Bilibili, and Xiaohongshu. The model also introduces a reasoning mechanism through reinforcement fine-tuning, giving it preliminary planning capabilities for multi-step tasks.

**InternVL**

InternVL [11] is a family of large-scale vision-language models developed by Shanghai AI Lab and Tsinghua University, among others. It leverages massive web image-text data to progressively align with large text language models. Its powerful visual capabilities enable it to excel in tasks like image-level and pixel-level recognition and serve as an effective alternative to ViT-22B [17]. The latest version, InternVL3 [117], introduces several innovative designs. During pre-training, it jointly learns multi-modal and language capabilities from diverse multi-modal data and pure text corpora. This unified training paradigm effectively solves the complexity and alignment challenges common in traditional post-processing training pipelines. InternVL3 integrates Variable Visual Position Encoding (V2PE) to support extended multi-modal contexts and employs post-training techniques like supervised fine-tuning (SFT) and mixed preference optimization (MPO) for refinement. InternVL3 shows outstanding performance on multi-modal tasks compared with leading models such as ChatGPT-4o [58], Claude 3.5 Sonnet [1] and Gemini 2.5 Pro [15] while maintaining strong pure language capabilities.

**OS-Atlas**

OS-Atlas [96] is a mobile agent foundation model jointly proposed by Shanghai AI Laboratory, Shanghai Jiao Tong University, and other teams. It aims to solve the severe dependency of existing mobile agents on closed-source models like GPT-4o. This work builds a powerful foundational GUI action model through dual innovations in data and models. For data, OS-Atlas developed and open-sourced the first cross-platform GUI localization data synthesis toolkit and utilized it to construct the largest open-source, cross-platform corpus for GUI localization, filling the gap in desktop data. For the model, by solving action naming conflicts during the training process, OS-Atlas has become a general and highly accurate foundational action model for all GUIs. Extensive evaluations show that OS-Atlas outperforms previous closed-source models on 6 benchmarks covering desktop, mobile, and web environments.

**OS-Genesis**

OS-Genesis [78] is a data synthesis method jointly proposed by Shanghai AI Laboratory, the University of Hong Kong, Shanghai Jiao Tong University, and other teams. It solves the problem of data scarcity for training mobile agents by allowing agents to autonomously explore graphical user interfaces and retroactively generate high-quality tasks from their interactions. The core innovation of OS-Genesis lies in its inverse task synthesis method, which solves the bottleneck of collecting high-quality data for mobile agent training. Unlike the traditional "task-driven" model, OS-Genesis adopts an "interaction-driven" approach, where the agent first autonomously explores the GUI interface and performs operations, and then retrospectively converts these interaction processes into high-quality task instructions. This achieves end-to-end data synthesis without manual supervision. This mechanism not only greatly increases data diversity and bridges the gap between abstract instructions and dynamic GUIs, but experiments have also shown that an agent trained with OS-Genesis significantly outperforms traditional methods on benchmarks like AndroidWorld. This new data synthesis paradigm offers great potential for transforming general vision-language models into specialized mobile agents.

**AGUVIS**

AGUVIS [99] is a unified vision-based framework for GUI task automation proposed by the University of Hong Kong. It gets rid of the reliance on text representations, platform-specific action spaces, and closed-source models, achieving a completely vision-based autonomous mobile agent. To achieve this, AGUVIS constructed the large-scale AGUVIS DATA COLLECTION dataset, which includes cross-platform trajectories, multi-modal localization, and explicit reasoning paths captured through inner monologues, providing high-quality annotations for training. Second, it designed a novel two-stage training process that decouples GUI localization from planning and reasoning abilities, enhancing autonomous navigation through a structured thinking process. Ultimately, AGUVIS performs well in both offline and real online benchmarks, becoming the first fully autonomous visual mobile agent that does not rely on closed-source models. It achieves performance on par with top closed-source models on related datasets like AIME [62].

**GPT-OSS**

GPT-OSS [56] is a series of open-source reasoning models released by OpenAI's research team in August 2025, including two models with different parameter sizes: gpt-oss-120b and gpt-oss-20b. Thanks to its powerful reasoning and planning capabilities, the model can act as an autonomous agent. With post-training data on Agentic Tool Use, it gains the planning and reasoning abilities required for a foundation agent model, allowing it to take natural language instructions and break them down into executable subtasks. Combined with its multi-modal capabilities, GPT-OSS can not only handle pure text instructions but also analyze visual information from screenshots, documents, or web pages to guide its actions. For example, it can understand a user's intent from a provided webpage screenshot and automatically execute complex operations such as filling out forms, browsing websites, or generating reports. It integrates language understanding, visual analysis, and action planning.

**GPT-5**

The GPT-5 [57] series models, released by OpenAI's research team on August 7, 2025, represent the latest advancements in general artificial intelligence. Its core innovation is the construction of a unified system module rather than a single giant model. This system integrates an efficient, responsive "smart model" and a "deep reasoning model" (GPT-5 thinking) specifically for complex tasks, using a real-time router to dynamically distribute tasks based on their nature, complexity, and user intent. This architecture aims to optimize resource utilization, ensuring a quick response for most routine queries while providing deeper, more reliable solutions for difficult problems. At the same time, the model's multi-modal understanding of voice, video, and audio has also been significantly improved, reaching top performance on multi-modal understanding datasets like MMMU [105].

Furthermore, GPT-5's hallucination rate for open-ended factual questions is significantly lower than its predecessors, and it is more honest about its limitations when faced with impossible tasks or a lack of crucial tools, rather than generating deceptive responses. These technological advances collectively make GPT-5 a system that has achieved a leap forward in intelligence, reliability, and safety.

**Mobile Agent Capability Comparison**

Table 1.1 shows a horizontal comparison of several open-source mobile agents presented in Section 1.2.1 across various metrics. These include foundation model capabilities (e.g., whether the model is open source, whether it supports Chinese data during training and inference) and user-level module metrics (UI module, cross-device support, cross-app usage, and user feedback).

The results in the table show that current mobile agent foundation models already possess considerable multi-modal understanding capabilities. However, they still have numerous limitations in practical applications. Specifically, existing agents have deficiencies and shortcomings in generalization, accuracy, long-range capabilities, and efficiency, hindering their ability to man-

age complex and varied task scenarios. These shortcomings collectively restrict the seamless, efficient, and intelligent autonomous operation of existing mobile agents in the real world. The AppCopilot proposed in this paper aims to fill these research gaps by comprehensively enhancing these key capabilities.

**Table 1.1:** *Mobile Agent Capability Comparison*

| Mobile Agent | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| InternVL [117] | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| AGUVIS [99] | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| OS-ATLAS [96] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Qwen-VL [3] | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| UI-TARS [66] | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ |
| AgentCPM-GUI [113] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| GPT-OSS [56] | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ |
| GPT-5 [57] | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✓ |
| **AppCopilot (Ours)** | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

**1. Open-Source**: Whether the model's weights are publicly available.

**2. Chinese Data**: Whether the model was specifically trained on labeled Chinese-language data.

**3. Multi-modal**: Whether the model can process with multiple input types of data.

**4. Model Inference**: Whether the model can display model reasoning process while performing the operations.

**5. Task Decomposition**: Whether the agent can break down a complex task into a series of smaller, manageable sub-tasks.

**6. Cross-App**: Whether the agent can perform operations across multiple applications.

**7. Cross-Device**: Whether the agent can coordinate and perform tasks across different devices.

**8. Physical Device**: Whether the agent can complete tasks on a physical Android device.

**9. User Information**: Whether the agent can leverage and memorize user information.

**10. Experiential**: Whether the agent can retrieve and reuse past experiences on new, similar tasks.

**11. UI Module**: Whether the agent includes UI module demonstrating model's operations.

**12. Function Calling**: Whether the agent can autonomously call external tools and functions (e.g., API calls, web searches) to complete a task.

**13. User Feedback**: Whether the agent can receive and process user feedback to correct its actions or improve its behavior in real time. User feedback refers to the agent recognizing and returning 'STATUS: NEED FEEDBACK' when an action requires information from the user (e.g., personal details, passwords), waiting for the user's input.

### 1.2.2  Case Studies and Capability Evaluation

**Baseline Model Experimental Setup**

To make a horizontal comparison between existing mobile agents and our AppCopilot, we selected several open-source and closed-source foundation models as baselines for experiments in specific telecommunication service scenarios. Table 1.2 shows the six baseline models, including three open-source and three closed-source models.

**Table 1.2:** *Baseline Foundation Models*

| Type | Model Name |
|------|-----------|
| Open-Source Models | UI-TARS-1.5-7B [66] <br> AgentCPM-GUI [113] <br> Qwen-2.5-VL-7B [3] |
| Closed-Source Commercial Models | gpt-4o-2024-08-06 [58] <br> claude-opus-4-20250514 [1] <br> gemini-2.5-pro-preview-05-06 [15] |

To evaluate the performance of baseline models in their telecommunication capabilities and other general utility functions, the tests were divided into three categories: **Phone and SMS Scenarios**, **China Unicom APP Usage Scenarios**, and **Other APP Usage Scenarios**. The agents have permissions to operate both Android virtual machines and physical devices, and their action space primarily includes operations like tapping, long-pressing, swiping, and typing. Table 1.3 shows the action space settings for several open-source models (UI-TARS-1.5-7B [66], AgentCPM-GUI [113], Qwen-2.5-VL-7B [3]).

**Phone and SMS APP Scenario Tests**

The phone and SMS scenario tests were divided into three tasks: **Answering and Hanging Up Calls**, **Making Calls**, and **Sending SMS Messages**.

1. **Answering and Hanging Up Calls**
   In the task of answering and hanging up a call, the agent needs to tap the "Answer" button on the phone notification. The task is considered successful if the call is connected correctly.
   Figure 1.4 shows screenshots of key steps when the agent operates a physical device to complete the task "Wait to answer the call from 13951696300."

2. **Making Calls**
   In the task of making a call, the agent needs to accept the instruction "Make a call to [phone number]" and complete three sub-operations: "tap the phone icon, type [phone number] character by character, and tap the call button." The task is considered successful if the target phone number receives a call notification.
   Figure 1.5 shows screenshots of key steps when the agent operates a physical device to complete the task "Make a call to 13951696300."

**Table 1.3:** *Action Space of Several Open-Source Vision Models*

| Model | Action Space |
|---|---|
| UI-TARS-1.5-7B | click(start_box='<\|box_start\|>(x1,y1)<\|box_end\|>'): Clicks an element (represented by bounding box coordinates)<br>long_press(start_box='<\|box_start\|>(x1,y1)<\|box_end\|>'): Long-presses an element, with specified duration<br>type(content=""): Types text content (e.g., filling forms, searching)<br>scroll(direction= "down & up & right & left" ): Scrolls the screen, can move up/down/left/right<br>press_back(): Simulates a back operation (Android back key)<br>press_home(): Simulates returning to the home screen (Android Home key)<br>wait(): Pauses for a period (used to wait for pages to load, etc.) |
| AgentCPM-GUI | thought: A Chinese description of the agent's current thinking process and decision logic<br>POINT: Specifies the starting coordinate for a tap or swipe<br>TO: Indicates the swipe direction ("up", "down", etc.) or the end coordinate<br>DURATION: Duration of the action (in milliseconds)<br>PRESS: Simulates a key press ("HOME", "BACK", "ENTER")<br>TYPE: Text content to be typed |
| Qwen-2.5-VL-7B | click: Clicks an element<br>long_press: Long-presses an element<br>swipe: Swipes<br>type: Types text content<br>system_button: Simulates a system key (BACK, HOME, ENTER)<br>terminate: Marks the task as complete |



**(a)** *Instruction: Wait to answer the call from XXX*   **(b)** *Operation: Check the phone call pop-up at the top of the screen*   **(c)** *Operation: Tap the pop-up to enter the phone screen*   **(d)** *Operation: Tap the answer button to connect the call*

**Figure 1.4:** *Example of an agent completing the "Answer a Call" task. The first image shows the initial phone screen and the user's specific instruction, while the subsequent images show the agent's autonomous operations and the changes in the interface state.*

**(a)** *Instruction: Call the phone number "13951696300"*  **(b)** *Operation: Tap the green phone APP on the bottom left of the home screen*  **(c)** *Operation: Type the target phone number*  **(d)** *Operation: Tap the call button to connect the call*

**Figure 1.5:** *Task example: Make a call to the target phone number*

3. **Sending SMS Messages**

   In the task of sending an SMS message, the agent needs to accept the instruction "Send a message with content [SMS content] to [user's phone number]" and complete five sub-operations: "tap the messages icon, tap the 'start chat' or plus icon, type [user's phone number], type [SMS content], and tap send." The task is considered successful if the target phone number successfully receives the SMS. Figure 1.6 shows screenshots of key steps when the agent operates a physical device to complete the task "Send a message to 18012985692, telling him to come home for dinner a bit later today."

**China Unicom APP Scenario Tests**

The China Unicom APP usage scenario tests primarily include two sub-tasks: **Logging into the China Unicom APP** and **Scheduling a broadband service appointment within the China Unicom APP**.

1. **Logging into the China Unicom APP**

   In the task of logging into the China Unicom APP, the agent needs to complete sub-operations such as "swipe to find the China Unicom APP, tap to open the China Unicom APP, tap the login button within the APP, and type in the relevant login information." The task is considered successful if a "Login successful" prompt appears. Figure 1.7 shows screenshots of key steps when the agent operates a physical device to complete the task "Log into the China Unicom APP account."

2. **Scheduling a broadband service appointment within the China Unicom APP**

   To complete this task, the agent needs to perform sub-operations such as "tap to open the China Unicom APP, type 'broadband' in the search bar and search, type in the relevant

**(a)** *Instruction: Send a message to XXX with specific information*

**(b)** *Operation: Enter the Messages APP*

**(c)** *Operation: Type the target phone number and the message text*

**(d)** *Operation: Tap to send the message, which shows as successfully delivered*

**Figure 1.6:** *Task example: Send an SMS message to the target phone number*



**(a)** *Instruction: Log in to the China Unicom APP*

**(b)** *Operation: Open the China Unicom APP*

**(c)** *Operation: Tap the bottom navigation bar to enter the login page*

**(d)** *Operation: Type the login phone number to complete the login*

**Figure 1.7:** *Task example: Log in to the China Unicom APP*

appointment information, and tap 'Book Now'." The task is considered successful if an "Appointment successful" prompt appears. Figure 1.8 shows screenshots of key steps when the agent operates a physical device to complete the task "Schedule a broadband service appointment within the China Unicom APP."



**(a)** *Instruction: Schedule a broadband service appointment within China Unicom*

**(b)** *Operation: Open the China Unicom APP and enter the broadband service page*

**(c)** *Operation: Type in the relevant user information*

**(d)** *Operation: Submit the order, which shows as successfully submitted*

**Figure 1.8:** *Task example: Schedule a broadband service appointment within China Unicom*

**Other APP Scenario Tests**

The other APP usage scenario tests mainly examine the scenario of downloading the China Unicom APP via a third-party application, specifically including **downloading the China Unicom APP from Yingyongbao** and **searching for and downloading the China Unicom APP via a browser**[3].

1. **Downloading the China Unicom APP from Yingyongbao**
   For this sub-task, the agent needs to complete operations such as "open the Yingyongbao APP, type 'China Unicom' in the search bar and search, tap the download button for the China Unicom APP in the results list, and complete the download." The task is considered successful if the China Unicom APP is successfully installed on the phone.

2. **Searching for and downloading the China Unicom APP via a browser**
   For this sub-task, the agent needs to complete operations such as "open the mobile browser, type 'baidu.com' in the browser's search bar and search, enter the Baidu search page, type 'China Unicom APP' in the Baidu search bar and search the results, and tap the 'Official Download' button on the page to complete the download." The method for checking task success is the same as above.

---

[3] https://baidu.com

Figure 1.9 shows screenshots of key steps when the agent operates a physical device to complete the task "Download the China Unicom APP from Yingyongbao."



**(a)** *Instruction: Download the China Unicom APP from Yingyongbao*

**(b)** *Operation: Open Yingyongbao, type in 'China Unicom APP'*

**(c)** *Operation: Tap the search result: China Unicom APP*

**(d)** *Operation: Tap download, and the result page shows the APP download is complete*

**Figure 1.9:** *Task example: Download the China Unicom APP from Yingyongbao*

**Experimental Results Analysis**

Tables 1.4, 1.5 and 1.6 show the performance and overall score rates of the different baseline models (Table 1.2) on the six fundamental tasks.

**Table 1.4:** *Comparison of Score Rates for Different Foundation Models on Specific Tasks*

| Model Name | Score |
|---|---|
| gpt-4o-2024-08-06 [58] | 1/6 |
| claude-opus-4-20250514 [1] | 1/6 |
| gemini-2.5-pro-preview-05-06 [25] | 6/6 |
| UI-TARS-1.5-7B [66] | 1/6 |
| AgentCPM-GUI [113] | 4.5/6 |
| Qwen-2.5-VL-7B [3] | 4.5/6 |

The evaluation of the six models on mobile agent-related tasks reveals significant performance differences. GPT-4o, as a multi-modal model, performs poorly on GUI tasks, possibly because it has not been specifically optimized for such application scenarios. The model generally suffers from inaccurate icon localization, and its output coordinates are often integers in units of 10, which suggests that its visual localization capability is not well-trained. Claude-opus-4's localization capability is superior to GPT-4o's, but it has a clear defect in adhering to instruction formats, often deviating from the specified output format to describe image content instead. In contrast, Gemini-2.5-pro demonstrates excellent visual localization and reasoning capabilities,

with format adherence issues occurring only in a few cases. UI-TARS's visual localization capability evaluation result is average, but this result might be affected by reproducibility issues, as its GUI operation script is not open-source. AgentCPM-GUI and Qwen-2.5-VL-7B show good visual localization performance, but their reasoning capabilities are relatively average, and they occasionally exhibit hallucinations.

Current mobile agents generally take text and images as input and output text that includes a thought process and specific actions. Therefore, a good mobile agent requires solid text reasoning and visual understanding capabilities, and its foundation model must strictly adhere to the specified output format (e.g., JSON) and a valid action space (e.g., $\mathrm{click}$ <1,200>). Furthermore, strong visual localization capability is key to ensuring accurate action points. However, current mobile agents still have limitations in generalization, accuracy, long-range capability, and efficiency. These four limitations will be further explained in Section 2.

**Table 1.5:** *Comparison of AI Model and Application Function Performance (Tasks 1-3)*

| Agent | Answering Calls | | Making Calls | | Sending Messages | |
|---|---|---|---|---|---|---|
| | Status | Description | Status | Description | Status | Description |
| gpt-4o-2024-08-06 | ✓ | Succeeded, but with low localization accuracy | ✗ | Unable to accurately localize the dial pad | ✗ | Unable to accurately localize the message button |
| claude-opus-4-20250514 | ✓ | Succeeded, but with low localization accuracy | ✗ | Unable to accurately localize and tap the dial pad | ✗ | Unable to accurately localize and tap the message button |
| gemini-2.5-pro-preview-05-06 | ✓ | Accurate localization | ✓ | Successfully recognized and dialed the phone number | ✓ | Can navigate to the message sending interface, but does not follow the instruction format for output |
| UI-TARS-1.5-7B | Partial | Unable to answer accurately, inaccurate localization | ✗ | Unable to accurately localize and tap the dial pad | ✗ | Unable to accurately localize and tap the message button |
| AgentCPM-GUI | ✓ | Accurate localization | ✓ | Successfully recognized and dialed the phone number | Partial | Can open the 'Messages' button, but cannot find the specific sending page |
| Qwen-2.5-VL-7B | ✓ | Successfully answered the call | ✓ | Successfully dialed the phone number | Partial | Hallucinations in the message content |

**Table 1.6:** *Comparison of AI Model and Application Function Performance (Tasks 4-6)*

| Agent | China Unicom Login | | Scheduling Broadband Services | | Downloading Unicom Apps | |
|---|---|---|---|---|---|---|
| | Status | Description | Status | Description | Status | Description |
| gpt-4o-2024-08-06 | ✗ | Inaccurate localization, opened an adjacent APP | ✗ | Inaccurate localization | ✗ | Inaccurate localization |
| claude-opus-4-20250514 | ✗ | Inaccurate localization; follows instruction format for multi-step output | ✗ | Inaccurate localization; follows instruction format for multi-step output | ✗ | Inaccurate localization; follows instruction format for multi-step output |
| gemini-2.5-pro-preview-05-06 | ✓ | Completed; can correctly tap the APP and navigate to the login page | ✓ | Completed; can correctly tap the APP and navigate to the broadband service page | ✓ | Mostly completed, but may fail at the end due to not following the instruction format for output |
| UI-TARS-1.5-7B | ✗ | Inaccurate localization, opened an adjacent APP | ✗ | Inaccurate localization, opened an adjacent APP | ✗ | Inaccurate localization, opened an adjacent APP |
| AgentCPM-GUI | ✓ | Completed; can correctly tap the APP and navigate to the login page | ✓ | Completed; can correctly tap the APP and navigate to the broadband service page | ✗ | Incorrect reasoning, did not infer the need to open the 'Yingyongbao' app |
| Qwen-2.5-VL-7B | ✓ | Tapped the APP to enter the 'My' interface, then a login method selection popped up | ✓ | Can localize the recharge interface | ✗ | Failed to open Yingyongbao due to network issues |

# 2

# Existing Problems

## 2.1 Generalization: Insufficient Generalization Capability in Application Scenarios

With the rapid development of GUI agents, although significant progress has been made in various benchmark tests, there are still many shortcomings in open and variable real-world application scenarios, highlighting the core bottleneck of "insufficient generalization." This bottleneck is not caused by a single factor but stems from a systematic disconnect between current training data and the complexity of the real world across multiple dimensions. **This chapter will conduct an in-depth analysis of the key challenges constraining model generalization from three progressive dimensions: the regional ecological differences, internal structural diversity, and behavioral authenticity of data.**

### 2.1.1 Data Resource Scarcity in Chinese Scenarios

In the process of building GUI agents with generalization, task execution, and environmental adaptation capabilities, high-quality data resources have always been a fundamental support. Training data not only carries the model's ability to understand interface structures, user intentions, and operational logic but also directly determines its stability and robustness when performing tasks in real scenarios.

However, despite the significant progress in GUI agent research in recent years and the emergence of a series of datasets and evaluation benchmarks—such as those designed for grounding, which aim to associate natural language instructions with on-screen controls, with representative works including Mind2Web [18] and OS-ATLAS [96]; and other datasets focusing on **agent task execution**, representing tasks as sequences of "action-observation pairs," such as AITZ [110], AndroidControl [44], and GUI-Odyssey [52]—the existing training data as a whole still heavily relies on English semantic environments and English interface designs (as shown in Table 2.1). **In the Chinese context, there is still an extreme scarcity of GUI operation task data for real-world APPs, which restricts the application and promotion of models in Chinese**

**scenarios.**

On the one hand, the language prompts, labels, and content in Chinese APPs are predominantly in Chinese. Their interface structures, control styles, and interaction sequences also differ significantly from English applications, making direct transfer from English training data to a Chinese environment challenging. Chinese interfaces often adopt designs that conform to local user habits, such as complex nested menus, vertical text, and multi-level navigation paths, whereas English interfaces tend to be flatter with horizontal layouts. Furthermore, Chinese interfaces widely feature highly localized elements like ambiguous icons, Pinyin abbreviations, and local symbols, which are almost non-existent in English data. This makes it difficult for models to accurately identify, locate, and parse the semantics of controls, thereby affecting the accuracy of operation execution.

On the other hand, the limited existing Chinese GUI datasets mostly focus on low-level tasks such as static screenshots, OCR recognition, or control classification, lacking dynamic interaction data structured around "action-observation" pairs. Such data often only cover static information about interface elements or single-step operations, missing continuous operational sequences, clear user task objectives, and semantic annotations of user behavior. This makes it difficult for models to learn complete task planning and execution logic from them.

**Table 2.1:** *Overview Comparison of Major GUI Agent Datasets*

| Dataset Name | Main Task | Data Scale | Language | Platform | Operating System |
|---|---|---|---|---|---|
| Mind2Web [18] | Grounding / Agent | 2,350 tasks; 137 websites | English | Web | N/A |
| ScreenSpot-Pro [41] | Grounding | 1,581 instructions | English | Desktop, Mobile, Web | Windows, macOS, Linux |
| OS-ATLAS [96] | Grounding | Over 13 million elements; 2.3 million screenshots | English | Desktop, Mobile, Web | Windows, Linux, macOS, Android |
| GUICourse [10] | Grounding / Agent | 73,000 web tasks, 9,000 mobile tasks | English | Web, Mobile | Android |
| UGround [26] | Grounding | 10 million elements; 1.3 million screenshots | English | Web, Desktop, Mobile | Windows, Linux, macOS, Android, iOS |
| AITW [68] | Agent | 715,000 trajectories; 30,000 instructions | English | Mobile | Android |
| AITZ [110] | Agent | 18,643 screen-action pairs; 2,500 instructions | English | Mobile | Android |
| AndroidControl [44] | Agent | 15,283 demonstrations; 14,500 tasks; 833 apps | English | Mobile | Android |
| GUI-Odyssey [52] | Agent | 8,834 trajectories; 212 apps | English | Mobile | Android |
| AMEX [7] | Agent | Over 104,000 screenshots; 110 apps | English | Mobile | Android |
| AndroidWorld [69] | Agent | 116 programmatic tasks; 20 apps | English | Mobile | Android |
| E-ANT [88] | Agent | 40,000+ trajectories | Chinese | Mobile | Android |
| Mobile3M [95] | Agent | 20,138,332 actions | Chinese | Mobile | Android |
| CAGUI [113] | Grounding / Agent | 3,000 grounding data, 600 agent tasks | Chinese | Mobile | Android |

### 2.1.2 Lack of Data Diversity in Vertical Scenarios

While achieving significant progress in benchmark tests, existing GUI agents still face severe challenges in their generalization capabilities in vertical domains. **The lack of generalization mainly stems from systematic biases in training data across two dimensions: insufficient application coverage breadth and the homogenization of task types.**

Firstly, at the application level, most existing training datasets show a long-tail distribution, focusing primarily on a few mainstream applications—such as YouTube and Amazon Shopping in English contexts, and WeChat and Taobao in Chinese contexts. Taking the large-scale Chinese Android dataset Mobile3M [95] as an example, its data distribution (see Figure 2.1) also confirms this phenomenon, with the vast majority of samples coming from a few popular applications. Annotated data for a large number of specific domains (such as finance, communication, etc.) or emerging applications is extremely scarce. This makes it difficult for the model to learn diverse interface layouts, interaction patterns, and user intentions. When the model is deployed in a new,

**Figure 2.1:** *Data Proportions of Categories and Specific Applications in the Mobile3M Dataset*

unseen application environment, its performance often plummets, severely limiting its practical value.

Secondly, at the level of data content, a significant number of benchmark datasets exhibit a distinct predisposition towards specific task typologies. Notably, **Navigation Tasks** account for a disproportionately high percentage of the data [110] [52] [79]. The fundamental objective of these tasks is to direct an agent through a sequence of operations—such as clicks and swipes across various application interfaces—based on natural language commands, with the ultimate goal of arriving at a designated target page. This bias is especially evident in several major large-scale datasets. For example, in AITZ [110], one of the largest collections of real-world mobile interactions, most task trajectories follow a canonical navigational pattern, typically abstracted as $\text{open APP} \rightarrow \text{multi-step click/swipe} \rightarrow \text{complete objective}$.

Such a distributional bias in the training data induces model overfitting to navigation-specific instruction patterns and environmental cues. While the resulting agent may exhibit proficiency in simple navigational tasks, it fails to acquire the more sophisticated capabilities required for understanding and decomposing complex instructions. Consequently, a significant degradation in performance is observed when the agent is confronted with **Composite Tasks** that require a fusion of distinct operational steps, such as navigating to a page before performing information extraction. Moreover, existing datasets are markedly deficient in their support for **Cross-APP tasks**. Many real-world user workflows—for example, *find a restaurant in a map app, copy its address, then switch to WeChat and send it to a friend*—necessitate seamless inter-application control and data transfer. The current datasets, with their overwhelming focus on intra-application operations, are

insufficient for endowing agents with the generalization capabilities required for such real-world scenarios, thereby severely constraining their practical utility.

In summary, the lack of application breadth and task depth are intertwined, collectively creating a hard-to-break ceiling that limits the model's generalization ability. This dual dilemma indicates that simply increasing the amount of homogeneous data can no longer bring about fundamental improvements. How to design and construct a training dataset that is more balanced, diverse, and challenging in both the application and task dimensions is a significant current challenge.

### 2.1.3 Challenges in the Authenticity of Behavioral Data



**Figure 2.2:** *Schematic of Spatial Position Generalization for a Target Control*

In the development of generalizable GUI agents, the **behavioral authenticity** of training data emerges as a factor that is both decisive and exceptionally challenging. A robust foundation for generalization can only be established if the data fully encapsulates the inherent diversity and complexity of genuine user operation processes. This high-fidelity requirement manifests as two core challenges: **interaction execution generalization** at the micro-level, and **task strategy generalization** at the macro-level.

At the micro-level of interaction execution, even after an agent has successfully formulated a rational plan and identified the correct control for the next step, it must still address the challenge of spatial position generalization. As illustrated in Figure 2.2, contemporary GUI agents typically employ **pixel-level localization**, where an action is parameterized as a single coordinate point $(x, y)$ on the screen. This representation is fundamentally misaligned with the nature of GUI controls, whose interactive affordance spans a spatial region (i.e., a bounding box) rather than a discrete point. Training data often exacerbates this issue by recording only a single, specific location clicked by a human demonstrator (e.g., the geometric center), inducing a bias that causes the model to overfit to this point-based representation while failing to learn the concept of the broader interactive area where multiple points are viable.

This brittleness becomes particularly evident under conditions of UI dynamism, such as when controls are rearranged, resized, or partially occluded. For instance, a minor shift in a button's visual position on a different terminal device can cause a coordinate-based action to

result in an 'off-target' click, failing the interaction even if the control's semantic identity was correctly recognized. Therefore, to enhance interactional robustness, the agent must develop a structural understanding of a control's interactive region and acquire the ability to dynamically select a viable operation point within that valid area, rather than merely reproducing a previously observed coordinate.

Second, at the task strategy level, the model must handle the **diversity and non-determinism** of paths leading to the same goal. For example, for the intent *query last month' s electricity bill*, users may follow very different routes:

1. Navigate step-by-step through the Utility Bills module on the homepage;
2. Enter keywords in the search box for a direct jump;
3. Access the historical bills module via the "My" page; and so on.

Although these paths are functionally equivalent, their interaction logic, interface navigation structure, and reliance on contextual information are all different. Most current training datasets only cover one or two "golden paths," which makes the model more inclined to learn fixed operational flows rather than context-adaptive task planning strategies. Therefore, if the interface layout is adjusted or a recommended entry point is missing, the model may fail the task due to a path mismatch, reflecting its lack of a semantic-level generalized understanding of the task objective.

## 2.2   Accuracy: Weak Single-Step Control Accuracy

During task execution, the ultimate success of a GUI agent depends not only on its high-level strategic planning capabilities but is also directly limited by the accuracy of each basic control step. Currently, even if an agent can correctly understand the task intent, it still exhibits significant fragility at the single-step operation level, specifically manifesting in three core problems caused by **technical architecture, positioning granularity, and the model's own uncertainty**.

### 2.2.1   Error Accumulation Caused by Non-End-to-End Models

In terms of technical architecture, some current GUI agent systems adopt a modular, non-end-to-end design. This paradigm is centrally reflected in the current agent frameworks. Its core mechanism is to decompose complex tasks through a preset Workflow and chain together a series of relatively independent functional modules. These modules—such as the "Interface Understanding Module" for visual perception, the "Task Reasoning Module" for decision-making, and the "Action Execution Module" for specific operations—are orchestrated and communicate via Prompt Engineering or external scripts [66]. Although this divide-and-conquer design has certain advantages in development efficiency and functional decoupling, **its inherent flaw is that local errors in the reasoning process accumulate and amplify along the multi-stage chain, ultimately leading to unstable task performance and difficulty in holistic optimization towards a unified objective.**

**Figure 2.3:** *Positioning Deviations in Control Interaction*

Under this "pipeline" processing mechanism, any slight error from an upstream module can be inherited and amplified by downstream modules, eventually having a fatal impact on task completion. For example, the interface understanding module might misidentify "Confirm Submission" as "Confirm Proposal" due to blurry fonts or insufficient recognition accuracy. When this incorrect input is fed into the task reasoning module, it may lead the language model to mistakenly judge that the current page lacks a submission entry, thus generating erroneous instructions like "scroll down to find the submission button." Once such an error occurs, subsequent modules can hardly correct it and can only execute the instruction, causing the operation to deviate from the correct trajectory. Overall, the system's single-step decision accuracy approaches the product of the accuracies of each module, forming a typical cascading error accumulation mechanism that significantly weakens the system's stability and interaction success rate.

More critically, this modular design faces the challenge of "indirect optimization" during the training phase. Each submodule is often optimized independently around its own local objective. For example, the visual module aims to maximize OCR accuracy, while the reasoning module strives to minimize the language model's perplexity. However, these local optimization goals often deviate significantly from the global task success rate, which is what the system truly cares about. Taking the "Confirm Submission" recognition error as an example, from the visual module's perspective, this might just be an insignificant character-level mistake, but in the context of the global task, it could cause task failure. Because the current architecture lacks an end-to-end error attribution mechanism that propagates from the overall task success/failure result back to the perception and understanding modules, the system cannot effectively capture and correct these "high-impact" local errors. This disconnect between objective functions greatly limits the potential for synergistic optimization of the system's overall performance.

### 2.2.2 Excessively Pixel-Level Positioning Granularity Leads to Control Errors

At the action execution level, many current methods use pixel-level coordinates (x, y) as the final output to specify the click position. **However, this excessively fine positioning granularity is misaligned with the physical reality of GUI interaction, making it extremely sensitive to prediction errors and a direct cause of control failures.**

A fundamental discrepancy exists between the model's action parameterization and the physical reality of GUI interaction. While an interface control's interactive affordance is a spa-

tial region (i.e., a bounding box), the agent is typically trained to predict a discrete coordinate point (x,y) within this area. This granularity mismatch renders the system highly sensitive to prediction errors. Consequently, even a minor deviation of a few pixels in the predicted coordinates can cause the action to fall outside the intended boundary, nullifying the operation. This vulnerability is exacerbated in dense interfaces where adjacent controls with opposing functions, such as "Delete" and "Cancel," are present in Figure 2.3,. In such cases, a slight coordinate shift can trigger an unintended and potentially critical action. This problem is particularly pronounced in dynamic GUI environments characterized by responsive layouts, varying device resolutions, or asynchronous content loading, as these factors compromise the generalizability of a specific coordinate learned in a static context.

### 2.2.3 Generative Model Reliability is Affected by Sampling Fluctuations



**(a)** *Schematic of the generation process of a Large Language Model [70]*

**(b)** *The effect of temperature on the generation results in a Large Language Model [70]*

**Figure 2.4:** *Randomness in the Generation Process of Large Language Models*

In action generation, the core decision-making ability of modern GUI agents typically relies on a generative reasoning mechanism driven by Large Language Models (LLMs). **However, the inherent random sampling strategy in LLMs, while enhancing the diversity and flexibility of language generation, inevitably introduces non-deterministic fluctuations, becoming a significant factor that constrains stability and reproducibility.**

In practical applications, to avoid monolithic generation results, large language models naturally adopt decoding strategies with randomness (as shown in Figure 2.4a), such as temperature sampling or top-k/top-p sampling (see Figure 2.4b). This mechanism means that even with identical inputs (e.g., the same screen state and task instruction), the intermediate chain of thought and the final action decision generated by the model may differ across different runs. This uncertainty in the generation process can easily lead to so-called "accidental failures" or "behavioral jitter." For example, in repeated trials of the same task, the model may succeed in some runs but, due to sampling perturbations, occasionally generate a reasoning chain that deviates from the correct trajectory, leading to stochastic task failure. Another example is when there are multiple functionally equivalent options (e.g., button A and button B have the same function), the model might frequently switch between them in consecutive operations, increasing the instability of the path.

This behavioral inconsistency induced by the sampling mechanism presents a significant challenge for GUI agents in high-reliability scenarios (e.g., financial transactions and medical applications). On the one hand, the system cannot ensure deterministic mapping from identical inputs to consistent decisions, which undermines task flow controllability and complicates

debugging. On the other hand, such instability heightens user uncertainty, thereby eroding the foundation of trust in human–computer interaction.

## 2.3 Long-Horizon Capability: Limited Long-Range Task Orchestration

One of the core hallmarks of a general agent is its ability to understand, plan, and execute complex long-range tasks over extended time scales. However, current GUI agents are mostly confined to single-step or short-range operations, lacking effective planning and execution capabilities for long-term goals and complex processes. This deficiency in "long-range capability" is specifically reflected at multiple levels, **from causal reasoning within sequences to collaborative work across applications and even across devices**.

### 2.3.1 Lack of Supervised Reinforcement for Sequential Reasoning

In the process of interaction between a GUI agent and its environment, each time step follows the loop of "perceive-reason-act": first, it needs to perceive the current interface state; then, it reasons based on the high-level task goal and historical operation records to decide the most reasonable next action; finally, it translates this decision into a specific command and acts upon the environment. This process is not a one-time static calculation. Each action of the agent changes the state of the environment, and this new state becomes the input for the next "perceive-reason-act" loop. This series of interlocking loops constitutes a complete "decision chain" or "action trajectory." **Therefore, the ultimate success of a task depends not only on the accuracy of single-step operations but more so on the logical coherence and cumulative effectiveness of the entire reasoning chain.**

However, the training paradigm of many current GUI agents heavily relies on "single-step behavior supervision" based on imitation learning. Its core method is to train the model to imitate the next single-step action given the current screen state. The fundamental flaw of this method is that it "slices" a continuous task sequence along the time dimension, with the implicit assumption that each "observation-action pair" is independently and identically distributed. This oversimplified approach inevitably leads the model to focus only on the "local optimum," i.e., how to accurately reproduce the immediate step, while neglecting to model and reinforce the long-range logical dependencies between steps [112] [23].

Although methods based on single-step imitation learning have advantages such as simple implementation and clear supervision signals in model training, such "local optimum" strategies often perform poorly in actual multi-step tasks. The fundamental reason is the model's lack of ability to model the internal structure and causal relationships of the task [92] [75] [74]. In GUI scenarios, an agent's operations are not composed of a series of isolated actions but exhibit highly structured process characteristics, with significant temporal dependencies and logical prerequisites between steps. For example, the "Register" button only becomes clickable after the user checks the "Agree to User Agreement" checkbox; ignoring this precondition will directly lead to the failure of subsequent operations. The single-step imitation learning approach focuses on

learning "what to do" in a specific interface state, rather than understanding "why to do it." This causes the model to exhibit "locally reasonable but globally incorrect" behavior when executing multi-step tasks, which is essentially due to a lack of modeling capability for the causal structure between operations.

Furthermore, in long-sequence reasoning tasks, this training mechanism that ignores the global causal chain will lead to the "temporal credit assignment" problem: a seemingly harmless minor misoperation early on, such as incorrectly selecting a shipping address in an e-commerce scenario, may only manifest its impact as a task failure a dozen steps later at the payment stage, yet the model cannot reasonably attribute this failure to the preceding error [97] [42]. Therefore, without sequence-level supervision, the agent can hardly extract useful signals from task feedback and achieve self-correction, thus finding it difficult to complete long-range tasks.

### 2.3.2 Limited Complex Task Planning and Decomposition Capabilities

**When executing long-range tasks in the real world, in addition to requiring the agent to have reasoning capabilities across long sequences of steps, a more critical part is its ability for structured planning and dynamic decomposition of complex tasks** [32] [31] [111]. This means that the agent not only needs to interact specifically with the GUI interface and complete low-level operations but also needs to possess high-level planning capabilities at an abstract level. It should be able to accurately infer hidden intentions and constraints from high-level, semantically ambiguous instructions from the user, and autonomously break them down into a series of clear, executable sub-goals, thereby constructing a hierarchical and operationally feasible overall task plan.

However, at the current stage, GUI agents based on large language models, despite having strong common-sense reasoning and language generation capabilities, still face multiple challenges in terms of truly effective task planning. First, effective planning requires a deep analysis of the user's ambiguous intent. The agent needs to be able to take high-level, semantically vague instructions from the user and accurately decompose them into a series of well-defined, logically consistent, and operationally feasible sub-tasks, ultimately forming a hierarchical overall task plan. Second, a plan that appears logically perfect at the text level may be ineffective if it is detached from a deep understanding of the real GUI environment. Current agents, when planning, often produce "hallucinations" and generate impractical steps due to their inaccurate understanding of the GUI environment. Finally, the agent lacks adaptability and correction capabilities when facing the dynamism of the real world. Once a plan is hindered by interface changes or operational failures, agents generally lack effective failure diagnosis and dynamic replanning mechanisms, making it difficult to recover from setbacks and adjust their strategies.

### 2.3.3 Insufficient Coordination and Linkage for Cross-Application Tasks

In reality, user tasks are often not confined to a single application but unfold naturally across applications, forming a "task chain" that spans multiple system components. For example, a user might receive an address in Enterprise WeChat, copy it, switch to AutoNavi Map for navigation, and then take a screenshot of the estimated arrival time to share back to Enterprise WeChat.

Such cross-application operations are extremely common in real scenarios, reflecting the highly systemic and collaborative nature of user task flows.

However, **the design and training paradigm of most current GUI agents is still confined within a single APP, lacking the ability to model cross-application tasks globally**. Achieving efficient cross-application collaboration faces multiple technical bottlenecks:

1. **Loss and forgetting of context**: When the agent switches between applications, intermediate states, task progress, and short-term memory from the previous application (e.g., App A) are often not effectively retained and transferred. This causes the task chain to "break" midway, affecting the coherence and correctness of subsequent operations [43] [46].

2. **Drift and loss of sub-task goals**: Cross-application tasks usually involve multiple phased sub-goals. The agent must not only maintain a clear awareness of the current sub-task during execution but also dynamically adjust its operational strategy when the interface changes or the context is updated, to ensure the consistency and correctness of the task goals.

3. **Insufficient adaptability to heterogeneous interaction interfaces**: Different applications use their own independent UI design paradigms and interaction logic. The agent needs to have strong UI representation generalization capabilities to achieve seamless transfer and immediate response in highly heterogeneous environments with different interface structures, control styles, and interaction modes.

Therefore, breaking through the limitation of "single-application closed execution" and enabling the agent with cross-application perception, memory, and planning capabilities—evolving towards a ubiquitous assistant with system-level context management and cross-domain orchestration capabilities—is the key path to achieving automation of real-world long-range tasks.

### 2.3.4 Cross-Device Collaborative Mechanisms Not Yet Established

**The ultimate challenge for long-range tasks lies in extending the agent's capabilities from serving a "single user" to supporting the collaborative work of "multiple users"** [48] [61]. This is not just an extension of capability but a paradigm shift from a "personal intelligent assistant" to a "distributed collaborative network." In this vision, each user's device runs an autonomous agent representing their individual preferences and intentions. These agents can discover, communicate with, and collaborate with each other to jointly accomplish a complex goal that is beyond the capability of a single agent. The establishment of this collaborative paradigm faces three core challenges. First, in a multi-user environment, the task states of different users are naturally distributed across different devices and private contexts. The agent needs to perform reasoning and coordination with incomplete information, possessing the ability to handle partially observable environments. Second, communication and negotiation mechanisms need to be built between agents. Agents must express intentions, divide tasks, and provide feedback through structured, semantically aligned interaction protocols to ensure policy consistency in asynchronous and heterogeneous environments. Finally, when their respective goals conflict (e.g., resource competition, temporal inconsistency), the system needs an effective consensus

mechanism to achieve an overall optimal collaborative solution while satisfying individual constraints.

Therefore, the implementation of cross-device collaboration essentially requires the GUI agent system to upgrade from a single agent's perceive-decision-execute loop to a system-level agent architecture with multi-agent collaboration, distributed state modeling, and mechanism design capabilities, marking the entry of GUI agent research into the "swarm intelligence" stage.

## 2.4 Efficiency: Low Inference and Decision-Making Efficiency

**Although GUI agents show great potential in task automation, their bottlenecks in inference and decision-making efficiency severely constrain their practicality**. This challenge stems from multiple levels: first, the reliance on large-scale models leads to high computational costs and response latency, posing high requirements for on-device deployment; second, systems generally lack long-term memory for user personalization and mechanisms to learn from past successful experiences, resulting in redundant interactions and an inability to quickly handle repetitive tasks; finally, because they cannot directly call the internal APIs of applications, agents must rely on slow and fragile graphical interface simulation, which greatly compromises both efficiency and stability.

### 2.4.1 Model Parameter Size Affects Inference Efficiency



**(a)** *Relationship between Model Size and Inference Latency*

**(b)** *Relationship between Model Size and Computational Cost*

**Figure 2.5:** *Relationship between Model Size, Inference Latency, and Computational Cost*

The powerful capabilities of current mainstream GUI agents depend on their underlying pretrained language or vision-language models with large-scale parameters, typically ranging from billions to hundreds of billions. Although larger language models possess excellent generality and expressive power, their inference costs are extremely high. Specifically, when processing input sequences of the same length, the model's forward computation cost (measured in FLOPs) grows linearly with the parameter size, thereby significantly increasing the response latency in each interaction round (as shown in Figure 2.5) [**Chinchilla**].

This enormous computational burden makes it nearly impossible to directly deploy fully functional agents on resource-constrained edge devices such as mobile phones, tablets. Consequently, most current systems are forced to adopt a cloud-based deployment solution, where

interface information is sent to a server, inference is completed by powerful cloud computing resources, and the instructions are then sent back. Although this approach solves the computing power problem, it introduces three equally tricky new issues: first, the unavoidable network latency, which reduces the smoothness of interaction; second, the data privacy risks associated with uploading sensitive information like user screens to third-party servers; and third, the considerable operational costs required to maintain the cloud service.

Meanwhile, the other extreme—overly lightweight models—is not an ideal solution either. Although such models are easy to deploy and have fast response times, their limited parameter size and simplified structure cannot effectively support complex multimodal understanding capabilities. This leads to their severe inadequacy in parsing ambiguous, multi-step natural language instructions and generalizing to unseen interfaces and tasks, ultimately falling into the dilemma of being "deployable but not capable."

Therefore, exploring a viable path to migrate from large cloud-based models to lightweight on-device models has become imperative. The core of this path is not just about unilaterally applying compression techniques like Distillation, Quantization, and Pruning, but more about finding the optimal balance (sweet spot) between model performance and resource consumption. **The ultimate goal is to build an on-device model that is both "light" and "powerful"—it must meet the low-power, low-latency requirements of on-device deployment while retaining sufficient core reasoning and generalization capabilities. Achieving this goal will be the core challenge and key to enhancing the practicality and deployment flexibility of GUI agents.**

### 2.4.2 Lack of Memory and Retrieval Mechanisms for User Personalization

An efficient GUI agent should possess long-term memory capabilities similar to a human assistant [60, 104], able to continuously accumulate and retrieve users' personalized preferences and historical behaviors to reduce redundant interactions and improve overall collaborative efficiency. However, **most current systems are still at the stage of stateless or short-term memory, lacking mechanisms for structured storage, dynamic updating, and real-time retrieval of user habits**, leading to severe "short-term memory impairment" in multi-round interactions. Specifically, agents often cannot remember information that the user has repeatedly entered, such as frequently used shipping addresses, dining preferences, commuting routes, or commonly accessed APP module entry points. Each task execution starts almost from scratch, forcing the user to repeatedly provide contextual information, which significantly increases interaction costs and weakens the system's user-friendliness.

Therefore, how to build a secure, controllable, and efficient personalized memory system for GUI agents that supports user-level long-term preference modeling and context memory retrieval has become one of the key technical paths to enhancing their practicality and intelligence level.

### 2.4.3 Lack of Mechanisms to Leverage High-Value Historical Behaviors

An efficient GUI agent should also have the ability to summarize experience from its own historical behaviors and reuse efficient strategies [114] [64], thereby avoiding redundant "from-

scratch reasoning" for repetitive tasks. Current systems generally lack this capability. **Even when faced with a task that has been successfully executed multiple times, the agent tends to restart the complete perceive-reason-execute loop, repeating interface analysis and task planning, which leads to redundant reasoning.**

Taking "recharge 100 yuan for a mobile number" as an example, ideally, the agent should be able to recognize the consistency of this task with past records and directly invoke the operation sequence that has been verified as optimal in the past, rather than reconstructing the reasoning chain. This "no experience leveraged" behavior pattern not only wastes computational resources but also significantly increases the system's response latency.

Therefore, establishing a generalizable mechanism for mining and reusing historical behaviors, enabling the GUI agent to achieve "experience-based rapid decision-making" when facing task repetitions or structurally similar tasks, is a key direction for improving its reasoning efficiency and user experience.

### 2.4.4 Third-Party Systems Cannot Directly Access Internal APIs

The mainstream interaction method for current GUI agents is to simulate the human "look at the screen—tap the screen" behavior path to operate with the front-end graphical user interface (GUI). This method has strong universality, but it is essentially an indirect and fragile interaction mechanism. In contrast, directly calling the application's function interfaces would significantly improve task execution efficiency and robustness. However, limited by the closed nature of operating systems and application ecosystems, **GUI agents, as third-party systems, usually have difficulty obtaining access permissions to these internal interfaces, leading to the following problems:**

1. **Lengthy interaction paths:** For example, to get the "order status," an agent typically needs to complete several steps in sequence: "open APP → click 'My' → click 'My Orders' → locate the target order → read the status." However, if it could directly call an API (e.g., get_order_status(order_id)), the same task could be completed in milliseconds.
2. **Poor stability:** UI-based interaction is highly dependent on the stability of the interface structure, and mobile interfaces are frequently updated. Any minor change in a button's position or text could cause the operation to fail. In contrast, function call interfaces are usually published as a public contract by the system and have higher stability and backward compatibility.

In summary, how to break down the "interface barrier" between the agent and the application, and enable it to intelligently switch between "direct function call" and "UI operation simulation" modes, will be a key breakthrough direction for improving task execution efficiency and system robustness.

# 3

# Generalization: Enhancing Generalization Capability for Chinese Universal and Vertical Scenarios

The preceding chapters have provided an in-depth analysis showing that the generalization bottleneck of current GUI agents originates from a systematic mismatch between training data and the real world in three dimensions: regional and ecological differences, internal structural diversity, and behavioral authenticity. To address this core challenge, our team has designed and implemented a comprehensive data collection and annotation methodology targeting both Chinese-English universal scenarios and vertical application scenarios. **We have designed and executed a complete data strategy consisting of three core stages: "Universal Capability Development," "Specialized Capability Enhancement," and "Real-Behavior Alignment." The final output of this strategy is a hierarchical, multi-source hybrid training dataset, whose overall architecture is presented in Table 3.1.** To systematically elaborate on its construction process, the following three chapters will describe these key stages in sequence:

**Bilingual Chinese-English Universal Data Foundation**: This stage focuses on the universal data required for model training. We will describe how, by collecting open-source datasets and independently constructing Chinese universal data, we infuse the model with broad cross-lingual foundational capabilities to address the challenge of "regional and ecological differences."

**Specialized Data Construction for Vertical Domain Applications**: This stage aims to collect the specialized data necessary for model training. We will present refined production methods for domain-specific agent data targeting particular application scenarios, thereby enriching the internal structure of the training data.

**Collection of Real User Behavior Data**: This stage is key to ensuring the data remains consistently aligned with real-world conditions. We will explain how to standardize the collection, filtering, and processing of real user behavior data, thereby fundamentally addressing the

challenge of "behavioral authenticity."

**Table 3.1:** *Overall Composition of Training Data*

| Corresponding Chapter | Data Type | Language | Data Source |
|---|---|---|---|
| **1. Universal Data Construction** | Grounding Data | English | Open-source datasets |
| | Grounding Data | Chinese | Open-source datasets |
| | Agent Data | English | Open-source datasets |
| | Agent Data | Chinese | Proprietary dataset |
| **2. Specialized Data Enhancement** | Agent Data | Chinese | Proprietary dataset |

## 3.1 Construction of the Chinese-English Universal Data Foundation

To develop an intelligent agent for graphical user interfaces (GUI agents) that can deeply understand and proficiently operate both Chinese and English interfaces, our model training relies on two core categories of data: **Natural Language–GUI Grounding Data and Agent Task Data**.

Natural Language–GUI Grounding Data is designed to establish precise correspondences between UI elements and natural language descriptions, enabling the model to "comprehend" the screen by accurately aligning visual elements (such as buttons, icons, and text boxes) with functional language (e.g., "search button," "user avatar," "enter password"). This alignment equips the model with essential visual comprehension capabilities. To achieve this goal, our dataset must not only include basic element recognition tasks but also encompass multi-dimensional visual understanding capabilities ranging from low-level to high-level, including fine-grained GUI icon classification, interface control detection, and optical character recognition (OCR) for key content. This design enables the model to develop a hierarchical perception ability akin to that of humans, progressing from "seeing clearly" (OCR) to "recognizing" (control detection) and ultimately "understanding" (icon semantics). Agent Task Data, on the other hand, is represented in the form of *(instruction, action sequence)* pairs. It aims to train the model to complete specific tasks based on natural language instructions by executing sequences of operations such as clicking, swiping, and text input—thereby enabling task-level intelligent interaction.

We posit that the model's bilingual capability is fundamentally rooted in the bilingual coverage of its training data. Accordingly, we have established a core principle: both data categories must be supported by high-quality Chinese and English datasets. **To achieve this, our data construction approach is two-pronged: fully leveraging existing open-source resources while purposefully building high-quality proprietary Chinese datasets.**

### 3.1.1 Integration and Analysis of Open-Source Datasets

We conducted a comprehensive and in-depth survey and integration of mainstream GUI-related datasets available in the community. First, we seek to "stand on the shoulders of giants" to rapidly build the foundational abilities of the model. Moreover, we aim to analyze existing datasets to clarify their strengths and limitations.

For both Natural Language–GUI Grounding tasks and Agent tasks, we have selected and integrated the following industry-leading open-source datasets as part of our training data preparation:

**Table 3.2:** *Datasets Related to Natural Language–GUI Grounding Capability*

| Dataset | Scale | Language | Core Task | Annotation Depth | Platform |
|---------|-------|----------|-----------|------------------|----------|
| AITZ [110] | 18k | English | Instruction-to-element mapping | Coordinates / Chain-of-Action-and-Thought | Android |
| GUICourse [10] | 700k | English | Text/region localization | Bbox / Text | Web |
| OS-Atlas [96] | 5.5M | English | General element localization | Bbox / Instruction | Cross-platform |
| SeeClick [13] | 280k | English | Visual element localization | Bbox / Instruction | Cross-platform |
| Wave-UI [81] | 16k | Chinese-English | UI understanding | Functional/intent description | Web |
| AMEX [7] | 100k | English | Multi-level element localization | Coordinates / Element function description | Android |
| UGround-V1 [26] | 6.5M | English | Visual localization | Bbox / Textual referring expressions | Cross-platform |

**Table 3.3:** *Comparative Analysis of Mainstream Agent Task Datasets*

| Dataset | Scale | Language | Task Type | Average Steps | Task Diversity | Annotation Depth |
|---------|-------|----------|-----------|---------------|----------------|------------------|
| AITW [68] | 3M / 715k | English | Single/multi-step (single-app) | 2–16 | Hundreds of apps and websites | Low: OCR + icon annotation only |
| Android_Control [44] | 250k | English | Single/multi-step (single-app) | 4.8 | Extremely high: 833 apps, 15k+ tasks | Medium: high/low-level instructions; partial reasoning chains |
| GUI_Odyssey [52] | 90k | English | Multi-step (cross-app) | 15.3 | High: 212 apps, 1.4k+ combinations | High: semantic labels for decision reasoning and screen understanding |
| AITZ [110] | 13k | English | Single/multi-step (single-app) | - | Based on AITW, covering 70+ apps | Extremely high: complete CoAT reasoning annotation |
| AMEX [7] | 38k | English | Multi-step (single-app) | 12.8 | High: 192 apps, 3k+ complex instructions | Extremely high: triple-layer annotation (function, localization, instruction chain) |

Through the systematic integration of existing open-source datasets, we observe that model training in English contexts already enjoys a solid data foundation, particularly in terms of dataset scale and task diversity. However, a critical and increasingly evident issue has emerged: **high-quality Chinese data remains severely lacking**. Although certain datasets (e.g., Wave-UI) contain a small amount of Chinese data, in terms of overall volume, application coverage, and task complexity, such data is far from sufficient to support the development of robust generalization and reasoning capabilities in Chinese scenarios. This shortfall will directly constrain the model's generalization capability and real-world relevance in Chinese contexts.

### 3.1.2 Self-Built Universal Data Pipeline for Chinese Scenarios

To address the industry-wide shortage of high-quality Chinese Agent task data and to build the core competitiveness of our model in localized applications, we initiated a large-scale, high-standard self-built data project. **This project employed a standardized data collection pipeline to obtain a vast quantity of user instructions originating from real-world scenarios.**

**Systematic Application Selection**

The process began with a systematic selection of applications within the domestic app ecosystem. We considered factors such as market penetration, user activity level, and the irreplaceability of application scenarios, ultimately constructing an application matrix consisting of more than 30 mainstream Chinese Android apps. This matrix, covering the primary scenarios of Chinese users' digital lives (Table 3.4), includes nationally popular applications such as WeChat (social communication), Amap (mapping and navigation), Xiaohongshu (content community), Dianping (lifestyle services), and Bilibili (media and entertainment), as well as frequently used tools in multiple vertical domains. This ensures diversity and representativeness in our dataset.

**Table 3.4:** *Matrix of Mainstream Apps Covered in the Self-Built Chinese Dataset (36 Apps in Total)*

| | | | | | |
|---|---|---|---|---|---|
| Ele.me | Amap | Dianping | WeChat | Bilibili | Tencent Meeting |
| Xiaohongshu | Taobao | Browser | Alarm | Notes | Meituan |
| Ctrip | Baidu Maps | Tencent Video | Gallery | JD.com | QQ Music |
| NetEase Cloud Music | Cainiao | Kugou Music | Youku Video | Didi Chuxing | 5G Wide Vision |
| Tencent Maps | Douyin | Alipay | Kuwo Music | Feishu | Ximalaya |
| Tomato Novel | Weibo | Qunar | WeChat Reading | Meituan Waimai | iQIYI |

**Task Design Based on Real Use Cases**

To construct Agent task data with authentic semantics, diverse expressions, and large-scale coverage, we designed and implemented a structured task generation process based on the "instruction template and slot filling" paradigm. This process, grounded in expert knowledge modeling, enables systematic expansion of data scale while maintaining instruction authenticity and linguistic flexibility, forming one of the key pathways for building our Chinese agent data ecosystem.

**Definition of Instruction Patterns and Template Library Construction** In the first stage, human experts lead the abstraction and definition of "instruction patterns" that cover the core functionalities of each app. From a user perspective, we analyze the core usage scenarios of typical apps and summarize several high-frequency, mission-critical task types. For example, the typical task type for "Ele.me" is "ordering takeaway," whereas "Amap" focuses on scenarios such as "route planning" and "location search."

For each task type, we further construct multiple base instruction templates with rich linguistic diversity. These templates vary in syntactic structure, informational emphasis, and colloquial style, while reserving task parameters through "<slot>" placeholders. For example, for the "ordering takeaway" task in "Ele.me," the template set may be designed as follows:

- Order <quantity ><item >from <store name >on Ele.me, deliver to <delivery destination >
- Order <quantity ><item >from <brand/store >, deliver to <delivery destination >
- Order <quantity ><item >from the nearest <store name >, deliver to home
- Get one <item >, address is <delivery destination >

Through meticulous template design, we ensure that subsequent data maintains naturalness and diversity in linguistic expression while retaining structural consistency.

**Slot Filling for Instruction Templates** After template construction, the second stage involves users and annotators filling the slots based on real-life experience to generate complete instructions. This stage emphasizes contextual authenticity and clarity of intent, avoiding the semantic distortion caused by random word filling. For example, for the template "Order <quantity ><item >from <store name >on Ele.me, deliver to <delivery destination >," a user's filling might be:

User A (office worker ordering lunch):

- <store name >→ "Yunhaiyao Yunnan Cuisine"
- <quantity >→ "one"
- <item >→ "steam pot chicken"
- <delivery destination >→ "company front desk"
- **Generated instruction** → "Order one steam pot chicken from Yunhaiyao Yunnan Cuisine, deliver to the company front desk"

By leveraging this "expert template design + user semantic filling" mechanism, we achieve both structural and large-scale optimization of task instructions. On this basis, the platform can further apply enhancement strategies such as semantic perturbation and context reuse, ensuring comprehensive coverage across task types, expression styles, and contextual scenarios. This process not only guarantees the naturalness and diversity of linguistic expressions in task data but also greatly improves data generation efficiency and controllability. Compared to methods relying entirely on manual writing or purely rule-based automatic generation, this approach offers greater scalability and stronger assurance of semantic authenticity.

**Ultimately, through the above process, we supplemented and annotated approximately 50,000 high-quality Chinese task instructions covering multiple mainstream app scenarios. Along with the integrated open-source datasets (see Table 3.2 and Table 3.3), these data constitute one of the core training corpora for our model.**

## 3.2 Data Annotation and Expansion for Vertical Domain Applications

For GUI agents, the quality and task relevance of data directly determine their actual performance in specific domains. In vertical industries such as telecommunications, business scenarios are often complex and diverse, encompassing a large number of domain-specific terms and standardized operational workflows. Although we have constructed a general-purpose dataset with broad coverage, such data typically falls short of meeting the stringent requirements of vertical domains for high professionalism, precision, and reliability.

Therefore, obtaining and processing vertical domain data closely tied to real-world business is not only the starting point for building high-performance GUI agents but also a critical factor determining their ultimate effectiveness. Such domain-specific data not only provides the semantic foundation for the model to understand business logic but also serves as the core guarantee for making precise decisions and executing efficiently in actual operations.

To this end, we conducted in-depth analyses of representative applications in real business contexts, systematically mapping out their core business processes and high-frequency user needs. As an illustrative case, we selected China Unicom's app ecosystem, which is typical of telecommunications applications, and carried out structured task analysis. Based on this methodology, we organized expert teams to design task instructions and completed high-quality domain-specific data annotation work. **As a result, we constructed datasets closely aligned**

**with real business scenarios of vertical industries such as telecommunications, providing a solid foundation for building GUI agent capabilities that can be broadly applied across different application ecosystems.**

### 3.2.1   Systematic Application Mapping in Representative Ecosystems

To systematically cover diverse business scenarios, we first conducted a comprehensive review and categorization of application ecosystems in vertical domains. As an example, in the case of China Unicom, we grouped its diverse applications and services into logically coherent and functionally cohesive categories (see Table 3.5), laying a solid foundation for subsequent instruction creation and data annotation work. This methodology is equally applicable to other enterprises and domains with complex service ecosystems.

**Table 3.5:** *Core Application Function Categories in the China Unicom Ecosystem (as an illustrative case)*

| Function Category | Representative Applications / Functional Examples |
|---|---|
| **Communication & Messaging** | **Unicom Video Ringtone** (incoming call video display and interactive service) |
| **Network & Broadband** | **China Unicom App** (account management, plan subscriptions, balance inquiry), **Unicom Smart Home App** (intelligent networking) |
| **Shopping & Cultural Tourism** | **China Unicom App** (includes Unicom Mall for mobile phones, accessories, smart devices, and cultural tourism products) |
| **Entertainment & Content** | **5G Wide Vision** (HD video, VR content, online performances, and other immersive multimedia services) |
| **Cloud Storage & Utilities** | **Unicom Cloud Drive** (file storage, backup, and multi-terminal synchronization services) |

### 3.2.2   Instruction Determination and Augmentation for Applications

Instruction construction is the key bridge between user intent and application operations, aiming to produce a high-quality instruction set that covers diverse scenarios. So **we adopted a biphasic strategy of "expert knowledge definition + large language model augmentation."**

First, an expert team authored "seed instructions" with clear intent and complete elements for the core functions of representative applications, thereby forming the high-quality nucleus of the instruction set. For instance, for the *Unicom Smart Home App*, an instruction might be: "Restart the router in the living room"; or for the *China Unicom App*: "Purchase a Xiaomi phone with a budget of 2,000–3,000 yuan." Similar instruction construction principles can be applied to applications in other vertical industries, such as finance, healthcare, or e-commerce.

Next, to simulate the varied linguistic habits of real users, we designed precise system prompts to strictly regulate the large language model's behavior, applying **augmentation-based paraphrasing** to the seed instructions. This process follows the principles below:

- **Strict constraint on intent preservation:** Ensure that LLM-generated instructions are 100% consistent with the core intent of the seed instruction. For example, an augmented instruction aimed at "buying a phone" must not drift into "checking phone information."
- **Protection of core keywords:** Key entities in the instruction, such as the application name "China Unicom App," product name "Apple phone," or constraints like "4,000–7,000 yuan," are set as immutable "protected keywords" to prevent loss of critical task elements.
- **Injection of harmless procedural phrasing:** Without altering the core operation, the model may insert minor, natural, conversational flow words (e.g., "I'd like to⋯," "Please help me⋯," "Go inside⋯and do⋯") to improve naturalness.

- **Output format control:** Length and formatting are constrained to maintain conciseness and structural consistency, facilitating automated downstream processing.

By utilizing the generative abilities of large language models while enforcing precise rule constraints in the prompts, we avoided semantic drift and information loss, ultimately producing a high-quality, wide-coverage instruction dataset. This methodology, while exemplified using China Unicom, is broadly applicable to other vertical applications requiring precise and reliable GUI agent operations.

## 3.3 Real User Behavior Data Collection Strategy

### 3.3.1 Challenges in Data Collection

Through the processes described in the preceding sections, we have constructed a rich in-struction dataset covering mainstream Chinese applications (e.g., WeChat, Taobao) as well as vertical domains (e.g., the Unicom app suite). However, these data contain only the textual description of user intent—i.e., the "instruction." A complete, trainable sample requires not only the instruction but also the strictly corresponding "ground truth" user behavior trajectory, namely the *action sequence*. Efficiently and reliably acquiring these action sequences is the core challenge in transforming instruction data into usable training assets.

Two main industry approaches exist to address this challenge:

**Approach 1: LLM-based automatic annotation.** This approach uses advanced multi-modal large models (e.g., GPT-4V, Gemini-2.5-pro) to "understand" both the instruction and current UI screenshot, then autonomously generate the required action sequence to complete the task. While potentially enabling fully automated collection at minimal human cost, it suffers from inherent limitations:

- **Lack of behavioral authenticity:** The actions generated are based on the model's param-eterized knowledge rather than actual human behavior. They may contain "hallucinations" or non-human-like operational habits, and cannot fully replicate real decision-making in complex scenarios.
- **Uncertain accuracy and generalization:** When facing unfamiliar or domain-specific complex UIs, the model's operation success rate is unreliable—especially for long-sequence, fine-grained tasks.

**Approach 2: Human-in-the-loop behavior trajectory collection.** Here, trained annota-tors execute task instructions on real devices, with their full interaction trajectories recorded. Compared to automated generation, this approach offers clear advantages:

- **Native data with high fidelity:** All trajectories originate from real user operations, faith-fully reflecting decision-making processes and behavioral patterns in real interactive con-texts. Such "gold standard" data provide highly reliable supervision for model training.
- **Diversity in task completion strategies:** Different users may choose different paths to the same goal—via navigation menus, search entries, or history records—based on personal

habits, UI prompts, or prior experience. This natural strategy diversity enriches the dataset, enabling models to generalize across paths rather than rely on a fixed procedure.

- **Spatial generalization at the interaction level:** Human operations inherently tolerate spatial variation. For example, in click actions, users often click anywhere within a control's active area rather than a single pixel coordinate. Thus, collected data naturally cover multi-point distributions within interactive regions, helping models build structural understanding of "clickable areas" and improving robustness to layout changes or occlusions.
- **High accuracy and adaptability:** Humans are adept at handling complex layouts and diverse task logic, flexibly responding to abnormal states, loading delays, and other issues —yielding higher collection success rates and data quality.

In pursuit of data authenticity and accuracy, we conclude that high-quality ground truth data are the foundation for building reliable intelligent agents. Therefore, we ultimately adopt the "human behavior trajectory collection" approach, supplemented by engineering methods to address potential efficiency and standardization challenges.

### 3.3.2 Standardized Collection Tool Implementation



**(a)** *Interface for selecting an action*

**(b)** *Interface with action selected but not executed*

**Figure 3.1:** *Data collection APK interface examples*

To scale and standardize the human annotation approach, we developed a dedicated data collection APK. This tool is not a simple screen recorder but an integrated system combining task distribution, process guidance, and structured data capture. Its core design centers on an **"intent-before-action" collection workflow**. As shown in Fig. 3.1, before performing each task step, annotators must explicitly declare the type of action they intend to execute (e.g., "click," "swipe," "type") via a floating menu on the right side of the screen. The collection tool then enters the corresponding capture mode, guiding and precisely recording the action's parameters (e.g., click coordinates, swipe trajectory, input text).

This "declare-then-execute" mechanism fundamentally ensures process standardization:

**Strict action space constraint and ambiguity elimination:** By requiring annotators to first select an action primitive, we discretize and structure continuous, ambiguous physical operations at the source. For example, the system does not need to guess whether a touch was a "click" or "long press" based on duration—the intent is explicitly declared. This guarantees that each recorded action $A_t$ comes strictly from a predefined, unified action space $\mathcal{A}$, eliminating operational ambiguity. Table 3.6 details the standardized action space $\mathcal{A}$.

**Accurate capture of interaction parameters:** Once the action intent is declared, the APK launches a targeted parameter collection module. Selecting "click" triggers precise logging of the next touch's normalized coordinates; selecting "swipe" logs start point and direction/endpoint. Guided capture ensures parameter completeness and accuracy, avoiding the errors and omissions common in manual logging.

**Unified data structure generation:** Within the "intent–action–parameter" framework, the APK automatically records the entire process as highly consistent structured data. A complete sample $\mathcal{D}$ is defined as the instruction $I$ and its trajectory $\mathcal{T}$ in the tuple $\mathcal{D} = (I, \mathcal{T})$, where $\mathcal{T}$ is an alternating sequence of screen states $S_t$ and standardized actions $A_t$: $\mathcal{T} = \{S_0, A_0, S_1, A_1, \ldots, S_n, A_n\}$. This guarantees uniformity across all collected data, facilitating downstream processing.

**Table 3.6:** *Atomic operations and parameter specifications of standardized action space $\mathcal{A}$*

| Atomic Operation / Parameter | Description |
| --- | --- |
| POINT | Locates an operation coordinate, accepting an integer tuple $(x, y)$ normalized to $[0,1000]$. Defaults to click; can be combined with to or duration to perform swipe or long press. |
| to | Defines swipe direction or target coordinate. Can be set to a preset direction (e.g., "up", "down") or combined with POINT to specify a swipe trajectory. |
| TYPE | Inputs text into the current focus element, accepting a string as parameter. |
| PRESS | Triggers a system-level key, such as "HOME" (home screen), "BACK" (back), or "ENTER" (confirm/newline). |
| STATUS | Updates the current task execution status, such as "finish" or "impossible". |
| duration | Specifies the duration of the action (in milliseconds), used for waits or long presses. |

In summary, by employing our self-developed collection tool built on the "intent-before-action" workflow, we have successfully scaled the human annotation approach. This strategy not only ensures the authenticity of behavioral data but, through its mandatory intent declaration and guided parameter capture, fundamentally guarantees the standardization, consistency, and accuracy of the collected data—laying a solid and reliable foundation for subsequent model training.

# 4

# Accuracy: Enhancing Single-Step Manipulation Accuracy in GUI Mode-Driven Systems

## 4.1 End-to-End Inference Architecture Based on Multimodal Large Language Models

In order to fundamentally cope with the prevalent issues of "error accumulation" and "optimization fragmentation" in traditional modular technology pipelines, we introduce an end-to-end inference architecture grounded in a Multimodal Large Language Model (MLLM) as the core foundation of the GUI agent. This architecture abandons the "pipeline-style" design that relies on multiple independent submodules in series, and instead integrates visual perception, instruction comprehension, and action decision-making into a unified model, achieving tighter cross-modal collaboration and a more natural reasoning flow.

Our system is built upon the industry-leading open-source pretrained multimodal model MiniCPM-V, and is further fine-tuned and reinforced on downstream tasks specific to GUI agent scenarios, thereby enhancing its adaptability in interactive execution domains. This section introduces the fundamental structural design of the multimodal large model we adopt, and further explains how this architecture effectively supports the end-to-end task execution process, along with the systemic advantages it brings.

### 4.1.1 Fundamental Principles and Advantages of the Multimodal Large Language Model

The multimodal large language model integrates a vision encoder, a large language model core, and a vision-language merger to achieve unified modeling of heterogeneous visual-text information, providing a solid foundation for building intelligent and efficient end-to-end interactive systems.

50

**Vision Encoder** The vision encoder is responsible for converting raw visual inputs such as GUI screenshots into high-dimensional feature vectors, namely "visual tokens." Mainstream implementations are typically based on the Vision Transformer (ViT)[20], which partitions images into patches and applies a Transformer architecture to capture both local details and global layout information.

**Large Language Model Core** The large language model serves as the reasoning and decision-making center of the system. Our model is initialized from the pretrained language model MiniCPM-V, inheriting its powerful language understanding, logical reasoning, and knowledge generalization capabilities.

**Vision-Language Merger** Located at the interface between the vision and language modalities, this module is primarily used to compress the long sequence features output by the vision encoder, thereby reducing computational load, and to map visual features into a representation aligned with the embedding space of the language model. For example, in Qwen2.5-VL, a two-layer multilayer perceptron (MLP) is employed as the merger: spatially adjacent patch features are grouped and concatenated, then projected via the MLP to the text embedding dimension. This mechanism improves processing efficiency while preserving the visual-semantic information crucial for downstream tasks.

Supported by the above architecture, an MLLM-based end-to-end system demonstrates three key advantages over traditional approaches. First, in terms of **deep integration of heterogeneous information**, an MLLM can directly process raw, pixel-level low-level visual features rather than relying on structured text generated by intermediate modules, thereby avoiding semantic loss caused by information conversion. This allows the model to capture detailed associations among appearance, color, and icons of UI elements and their functional semantics, resulting in deeper understanding of GUI interfaces. Second, the architecture possesses **end-to-end global optimization capability**: as a continuous, differentiable neural network, it supports backpropagation from the final task objective, enabling unified optimization of all module parameters, which significantly enhances stability and efficiency in real manipulation tasks. Finally, regarding the **paradigm of capability inheritance and development**, the model not only inherits the strong generalization and zero-shot reasoning abilities of the pretrained LLM, but also supports task-oriented fine-tuning to further shape professional operational strategies for GUI agents. This enables the effective unification of generality and domain-specific expertise, with the ability to continuously adapt and expand in novel interfaces and complex environments.

### 4.1.2 GUI Agent with End-to-End Architecture

Based on the multimodal large language model, we are able to construct a truly end-to-end inference system. As shown in Figure 4.1, **the core design philosophy of this architecture is to simulate the most natural human interaction logic: a person perceives the GUI interface through "vision," integrates it with the "language"-based task intent, and finally forms an action decision also expressed in language. The adopted multimodal large language model**

**Figure 4.1:** *Input-output schematic of the GUI agent with end-to-end architecture*

**directly maps multimodal inputs (GUI screenshot + user text instruction) to structured text outputs.** This output is an explicit, language-based action command (e.g., click("Login button") or type("hello", on="Input field")), which directly corresponds to specific executable operations. Throughout this process, there are no manually designed intermediate modules or feature engineering; the entire information flow is seamlessly and naturally completed within the model.

This end-to-end "text generation" paradigm has its most significant advantage in greatly facilitating holistic model optimization. It supports backpropagation from the success or failure of the final task via a unified loss function, optimizing the full chain of parameters from visual perception to textual decision-making. Whether subtle deviations in perception or flaws in reasoning logic, both can be "captured" and corrected by the final task objective. This "result-driven model improvement" training paradigm is difficult for traditional methods to match. In addition, the architecture provides two further benefits: since the model always directly accesses the most raw and complete input information, it completely avoids the cascading amplification of errors across independent modules, significantly improving the robustness of the reasoning process; secondly, the overall system architecture is more concise, requiring no complex interface protocols or submodule logic, thereby greatly reducing development and maintenance costs.

In summary, the MLLM-driven end-to-end inference architecture not only improves accuracy and generalization at the model level but also achieves a return to human natural interaction logic in design philosophy. By building a deeply integrated chain from perception, understanding, to decision-making, it provides a solid technical foundation for creating truly intelligent and reliable GUI agents.

## 4.2   OCR- and OR-Based Region Localization Calibration Mechanism

At the action execution level of the GUI agent, many existing methods adopt pixel-level coordinates as the final output form to specify the click position. However, excessively fine-grained localization makes the model extremely sensitive to minor errors, often resulting in manipulation failures. To mitigate this issue, we introduce a control semantic-assisted localization mechanism based on OCR (Optical Character Recognition) and OR (Object Recognition). **This mechanism utilizes the OCR+OR module to identify semantic labels and bounding box regions of interface elements, which are then used to geometrically calibrate the model's**

**output coordinates.**

### 4.2.1    Widget Region Recognition and Semantic Parsing



**(a)** *Original interface image*        **(b)** *Visualization of OCR + OR recognition results*

**Figure 4.2:** *Example of widget region recognition: from raw interface to structured semantic region extraction*

**Table 4.1:** *Example of structured widget information output by the OCR + OR module*

| Widget Index | Widget Type | Content | Interactivity | Normalized Position (bbox) |
|---|---|---|---|---|
| 11 | Icon (Application) | YouTube | Yes | [0.74, 0.63, 0.90, 0.73] |
| 12 | Icon (Application) | Play Store | Yes | [0.10, 0.63, 0.26, 0.73] |
| 13 | Text | Thu, Aug 7 | No | [0.12, 0.14, 0.18, 0.17] |
| | | … | | |
| | | … | | |

A Graphical User Interface (GUI) serves as the central hub connecting users with the functional capabilities of an application. Through graphical elements such as buttons, icons, and text, it fulfills the dual role of information presentation and interaction control. Fundamentally, a GUI functions as a visual language, conveying operational intent and system state via the spatial layout of widgets, visual style, and textual content. Consequently, the task of widget region recognition and semantic parsing in GUIs must go beyond the mere detection of visual objects and recognition of textual content; it must also achieve a structured understanding of widget semantics and interaction intentions, enabling machines to emulate human-like GUI comprehension capabilities [22][16].

Traditional computer vision approaches often focus on "detection" and "recognition," e.g., locating a bounding rectangle in the interface or identifying the text contained within it. However, the real challenge lies in integrating these raw perceptual results into high-level semantic structures with functional meaning—bridging the "semantic gap" between visual pixels and interaction functionalities.

Although Optical Character Recognition (OCR) and generic Object Recognition (OR) technologies have matured and achieved wide success in practical applications, their isolated use still exhibits significant limitations in the context of complex GUI understanding.

**Semantic Blind Spots of OCR**    The primary goal of OCR is to convert text within images into structured string information. The conventional OCR pipeline consists of image preprocessing, text detection, character recognition, and result integration. While effective in tasks such as document digitization, this approach shows two notable shortcomings in GUIs with complex layouts and semantically interdependent widgets [22]:

First, **lack of layout and structural awareness**. GUI semantics are not only embedded in textual content but also in the spatial relationships and arrangements of widgets. For example, a "Username" label positioned next to an input box forms a logical unit, yet OCR merely extracts the text without capturing this semantic grouping.

Second, **lack of interaction-function understanding**. While OCR can recognize the word "Submit," it cannot determine whether the button containing it is interactive, nor can it infer the action triggered by clicking it. Therefore, OCR inherently lacks the ability to model semantic functions.

**Generalization Bottleneck of OR**    Object recognition focuses on extracting bounding box information for specific object categories from images. As deep learning models like YOLO and Faster R-CNN have progressed, widget detection accuracy has improved considerably. Nevertheless, in the context of GUIs, OR still struggles with **weak recognition of semantic specificity**. For instance, an OR model may classify multiple widgets as "Button" but cannot distinguish their functional differences—such as "Save," "Cancel," and "Delete" —which may share similar visual styles but differ entirely in purpose. This inability to capture instance-level semantic distinctions limits the support such models can provide for higher-level interaction decisions [8][16].

**Collaborative OCR + OR Parsing Paradigm**    Given the above, it is necessary to deeply integrate OCR and OR techniques into a GUI parsing framework that possesses both structural awareness and semantic understanding capabilities. The aim of this collaborative paradigm is not only to identify "what elements exist" but also to understand "what they mean," "whether they are interactive," and "what consequences interaction would produce."

We adopt Microsoft's open-source OmniParser model [53], which exemplifies this collaborative design philosophy. OmniParser combines object detection and text recognition with semantic description generation, thereby enabling an end-to-end pipeline from visual perception to semantic modeling. Its core idea is to "tokenize" a GUI screenshot into a structured, semantic, DOM-like representation suitable for downstream reasoning and decision-making by language models. The framework comprises three parallel modules:

**A. Interactable Region Detection** A fine-tuned YOLOv8 model detects bounding boxes for interactive widgets such as buttons, icons, and text boxes. This module answers the fundamental questions: "What interactive elements are present on the interface?" and "Where are

they located?"

**B. Semantic Description Generation** A pretrained vision-language model generates image-to-text descriptions for widget regions, including not only visual attributes but also anticipated functional roles. For example, a gear icon is described as "Settings menu entry" rather than simply "gear icon."

**C. Optical Character Recognition** The OCR module extracts all textual information from the interface—including labels, menus, and tooltips—providing linguistic cues for constructing a semantic structure.

OmniParser's output preserves the spatial layout and visual characteristics of interface elements while integrating textual semantics and interaction properties. This structured output facilitates semantic understanding and action planning. As shown in Figure 4.2, in a typical smartphone interface, OmniParser can identify semantically meaningful widgets such as "Gmail" and "YouTube" icons, along with their precise locations and interactivity labels. The structured output, illustrated in Table 4.1, specifies widget type, semantic content, interactivity, and normalized spatial position (bbox). This representation provides stable and precise target regions for subsequent action coordinate calibration, mitigating uncertainty inherent in pixel-level targeting.

### 4.2.2 Action Coordinate Calibration Based on Widget Regions



**Android phone interface**

**Touchable area of the *"Phone"* widget**

**Bounding box of the *"Phone"* widget**

**(a)** *Mobile interface*          **(b)** *Interactive widget regions and bounding boxes*

**Figure 4.3:** *Example of interactive regions and bounding boxes corresponding to widgets*

To improve the interaction accuracy and robustness of GUI agents operating in complex visual environments, we design an action coordinate calibration strategy that combines static spatial constraints with dynamic historical memory mechanisms, leveraging the spatial and semantic information of widget regions.

Each interactive widget is associated with a bounding box (bbox) that approximately delineates its operational area (Figure 4.3).

**Static Bounding Box Constraint Correction**: For a model-generated click point $p =$

$(x, y)$, we first check whether it falls within any widget bounding box in the set $\mathcal{B} = \{b_1, b_2, \ldots, b_n\}$:

$$\exists b_i \in \mathcal{B}, \quad p \in b_i$$

If $p$ is inside a widget region, the click is considered valid and the original coordinate is retained. Otherwise, we treat it as a potential misalignment requiring correction. We compute the minimum perpendicular distance from $p$ to each bounding box $b_i$ and relocate $p$ to the center of the closest widget $b^*$:

$$b^* = \operatorname{argmin}_{b_i \in \mathcal{B}} d(p, b_i), \quad p = \operatorname{center}(b^*)$$

where $d(p, b_i)$ is the minimum Euclidean distance from $p$ to rectangle $b_i$:

$$d(p, b_i) = \sqrt{\max(x_1 - x, 0, x - x_2)^2 + \max(y_1 - y, 0, y - y_2)^2}$$

with $b_i = [x_1, y_1, x_2, y_2]$ representing bbox coordinates (top-left and bottom-right corners).

The rationale is to leverage detected widget bounding boxes to monitor and correct GUI agent click commands in real time, ensuring operational validity and preventing failures due to coordinate drift. This is grounded in the behavioral assumption that a user's intention is typically to interact with a specific widget on the interface; hence, a valid click should fall within an interactive region. When a click lies outside all recognized widget regions, it is reasonable to infer a targeting error and adjust it to the nearest widget center.

**Dynamic Historical Correction Avoidance**: For each frame $I$, the system computes a perceptual hash $h(I)$ and records the set of previously corrected click coordinates that failed to produce interaction $\mathcal{F}(h(I))$. Before applying a correction, if:

$$p \in \mathcal{F}(h(I))$$

then the system bypasses the correction and executes the original $p$, avoiding repeated erroneous adjustments. This mechanism is based on the assumption that a valid interaction should cause a state change in the interface; if a click at a given location has historically failed to elicit any change, it is deemed ineffective, and repeating it should be avoided to prevent infinite loops.

## 4.3  Multi-Agent Collaborative Decision-Making Strategy

Multi-agent voting among large language model (LLM) agents is a key ensemble decision-making strategy [40][34]. Its central premise is that aggregating the "opinions" of multiple autonomous LLM agents can produce a decision that is more accurate and robust than any individual model. This approach mitigates hallucinations and biases common to single models, while leveraging "collective intelligence" to address complex reasoning and decision-making tasks.

Various forms of implementation exist, ranging from basic majority voting [40], to weighted voting based on historical performance or confidence scores [65], to ranked-choice voting capable of expressing complete preference orderings [116].

In our system, we employ a multi-agent ensemble strategy with majority voting to enhance decision-making robustness. At each decision step, multiple independent agents receive the exact same environmental state and independently propose an action. The system then aggregates these proposals via voting to select a consensus-based action for execution.

**Furthermore,** in our collaborative framework, each decision is decomposed into two stages to ensure consensus from "what to do" to "how to do it":

**Stage 1: Action Type Voting**  The system collects all proposed $\mathtt{action}$ outputs and conducts majority voting on their **action type**. Action types include: POINT (click coordinate), PRESS (press a key), TYPE (enter text), STATUS (status check), CLEAR (clear field), etc. The type with the highest vote count becomes the primary action type for that step.

**Table 4.2:** *Multi-agent action aggregation strategy*

| Action Type | Parameter Aggregation Strategy |
| --- | --- |
| (x,y) coordinates | **Geometric centroid voting**: compute the centroid of all coordinates, select the closest candidate |
| PRESS/STATUS/CLEAR | **Discrete count voting**: count occurrences, choose the most frequent value |
| text | **Text frequency voting**: count suggested texts, choose the most common |
| duration | **Nearest-to-mean selection**: compute the mean, choose the original duration closest to it |
| thought | **Source tracing**: adopt the reasoning chain from the agent whose action was ultimately selected |

**Stage 2: Action Parameter Aggregation**  To consolidate outputs from parallel decision-making agents into a unified and effective instruction, we design a parameter aggregation strategy tailored to the data characteristics of each action type (Table 4.2). Instead of a one-size-fits-all majority rule, the method applies different voting or fusion mechanisms: for continuous spatial data like $(x, y)$ coordinates, geometric centroid voting identifies a consensus location; for discrete actions like PRESS or candidate texts, frequency counting selects the most common choice; and for the decision rationale ($\mathtt{thought}$), we preserve coherence by tracing back to the originating agent. This differentiated aggregation scheme condenses collective agent intelligence into a single, logically consistent action, thereby enhancing decision accuracy and interpretability.

Through this mechanism, we build a careful and robust decision-making system that synthesizes multiple agents' inputs while fine-tuning parameter handling, significantly improving overall task execution success rates.

# 5

# Long-Horizon Capability: Enhancing Planning through Complex Task Decomposition

## 5.1 Supervised Reinforcement Algorithms for Sequential Reasoning

The interaction between mobile agents and their environment is inherently a multi-turn, sequential decision-making process. In such long-horizon tasks, the ultimate success or failure of an agent does not depend on a single correct action in isolation, but rather on the logical coherence and cumulative effectiveness of the entire action sequence. However, as discussed previously, traditional training paradigms exhibit inherent limitations in addressing the core challenges of sequential reasoning, particularly in credit assignment over time and modeling long-horizon causal dependencies.

To overcome these limitations, we introduce the training approach of **Reinforcement Learning** (RL). Specifically, instead of providing the model with a "ground truth" answer for each step, we allow it to engage in autonomous exploration through interaction with the environment, thereby generating complete action sequences. Subsequently, the system provides a fine-grained reward signal based on the overall task outcome. By iteratively optimizing its policy to maximize long-horizon cumulative rewards, the model is compelled to learn the causal chain between single-step behaviors and final outcomes. Consequently, the learning objective naturally shifts from achieving "local optima" to attaining "global optima".

In practice, we adopt the Group Relative Policy Optimization (GRPO) algorithm for Reinforcement Finetuning (RFT). GRPO [73] is a simplified and improved variant of the classical Proximal Policy Optimization (PPO) algorithm [72]. Its key idea is to replace the complex value function (Value Critic) network in PPO with relative comparisons among a group of candidate generations. This design not only simplifies the training architecture but also provides more

stable and efficient gradients through direct relative reward signals.

The algorithm proceeds in four core steps:

1. **Sampling**: Given an input query $q$, the old version of the current policy $\pi_{\theta_{old}}$ is used for sampling, generating $N$ candidate responses (completions), denoted as the set $\{o_1, o_2, \ldots, o_N\}$.

2. **Evaluation**: For each generated response $o_i$, one or more external reward functions are applied to evaluate it, producing a scalar task reward $r_i$. Thus, we obtain the reward set corresponding to the response set: $\mathbf{r} = \{r_1, r_2, \ldots, r_N\}$.

3. **Advantage Calculation**: The GRPO algorithm does not rely on a learned value network to estimate the expected return of a state. Instead, it normalizes the reward of each response within its group to compute its relative advantage. Specifically, the advantage value $\hat{A}_i$ of response $o_i$ is obtained by calculating the Z-score of its reward:

$$\hat{A}_i = \frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r}) + \epsilon} \tag{5.1}$$

   where $\text{mean}(\mathbf{r})$ and $\text{std}(\mathbf{r})$ represent the mean and standard deviation of the reward set $\mathbf{r}$, and $\epsilon$ is a small constant introduced for numerical stability. The advantage value $\hat{A}_i$ directly quantifies the relative quality of response $o_i$ compared to the average level generated by the current policy.

4. **Policy Update**: After computing the advantages, we update the policy network using an objective similar to PPO. This objective includes a Clipped Surrogate Objective and a KL Divergence penalty. The goal of the policy update is to maximize the following function $J_{\text{GRPO}}(\theta)$:

$$J_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim P(Q), \{o_i\} \sim \pi_{\theta_{old}}} \left[ \sum_{i=1}^{N} \mathcal{L}_{\text{CLIP}}(\theta, o_i) - \beta D_{\text{KL}}(\pi_{\theta_{old}} \| \pi_\theta) \right] \tag{5.2}$$

   where $\mathcal{L}_{\text{CLIP}}$ is the clipped surrogate objective function, defined as:

$$\mathcal{L}_{\text{CLIP}}(\theta, o_i) = \min \left( \rho_i(\theta)\hat{A}_i, \quad \text{clip}(\rho_i(\theta), 1 - \varepsilon, 1 + \varepsilon)\hat{A}_i \right) \tag{5.3}$$

   Here, $\rho_i(\theta) = \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}$ denotes the probability ratio between the new and old policies for response $o_i$. $D_{\text{KL}}(\pi_{\theta_{old}} \| \pi_\theta)$ represents the KL divergence between the old and new policies, which measures the magnitude of the policy update. The hyperparameter $\beta$ controls the penalty strength to ensure the stability of the training process.

In summary, we adopt the GRPO algorithm for reinforcement learning fine-tuning of the model, aiming to address the challenges of long-horizon dependencies and credit assignment faced by single-step supervised learning in sequential tasks. This end-to-end reinforcement learning paradigm enables the agent to better establish causal relationships between individual actions and final outcomes, thereby enhancing its strategic planning capability and overall task completion rate in complex GUI tasks.

## 5.2   Long-Horizon Task Planning and Decomposition Mechanism

Many real-world user tasks often require execution across different operating environments, which can be categorized into two forms: **cross-application** collaboration, i.e., coordinating multiple applications on a single device, and **cross-device** collaboration, i.e., performing cooperative operations across heterogeneous devices such as smartphones and computers.

Although these two scenarios differ in their manifestation, they share a common intrinsic requirement for task planning and decomposition: a high-level, complex user intent must be systematically planned and orchestrated into a logically coherent workflow composed of multiple atomic operations. Each operation within the workflow must then be precisely mapped to a specific "execution endpoint" (i.e., a concrete application or physical device) for completion.

**Accordingly, we propose a unified hierarchical planning and execution framework. This framework leverages a central high-level planning agent, $\pi_{\mathbf{plan}}$, to decompose and allocate high-level tasks to different execution endpoints (applications or devices). Furthermore, it incorporates an efficient state synchronization mechanism to ensure that all subtasks are executed accurately and sequentially according to their dependency structure.**

### 5.2.1   Task decomposition planning mechanism

To transform ambiguous natural language commands from users into structured tasks that machines can understand and execute, we design a task decomposition and planning mechanism led by a **High-level Planning Agent** $\pi_{\text{plan}}$. This agent, built on a large language model, takes a complex high-level task command $T$ as its input, with the core objective of deeply parsing user intent and ultimately outputting a structured **Task Dependency Graph** $G_T$, laying a solid foundation for subsequent task allocation and execution.

Specifically, when $\pi_{\text{plan}}$ receives a command, it does not adopt a traditional linear decomposition approach, as this fails to effectively capture parallel relationships among subtasks. Instead, we innovatively represent the task as a **Directed Acyclic Graph (DAG)** to more precisely capture its internal structure. The task dependency graph $G_T$ is defined as:

$$G_T = (V, E), \quad V = \{st_1, st_2, \ldots, st_n\}, \quad E \subseteq V \times V \tag{5.4}$$

Here, $V$ is the set of $n$ nodes corresponding to sub-tasks, where each sub-task $st_i$ is a fine-grained operation executable independently on a single endpoint. The directed edge set $E$ describes the execution dependencies among subtasks. An edge from $st_i$ to $st_j$, $(st_i, st_j) \in E$, specifies that the completion of $st_i$ is a prerequisite for the execution of $st_j$, formally:

$$(st_i, st_j) \in E \implies st_i \prec st_j \tag{5.5}$$

This graph-based dependency modeling provides a rigorous and flexible mathematical foundation for handling complex task structures. Unlike linear sequences, DAGs naturally represent both *parallel* and *sequential* relationships, thereby enabling parallel scheduling and improved execution efficiency.

The mechanism proceeds through three major stages: task decomposition, task allocation, and task state synchronization, as described below.

**Task Decomposition:** When the system receives a complex task command $T$, the high-level planning agent $\pi_{\text{plan}}$, powered by the LLM, first parses and decomposes it. For example, for the task "buy a suitable gift for Lili", the planning agent decomposes the goal into multiple atomic operations. The set $V$ of subtasks enumerates all the required atomic operations:

$$
\begin{aligned}
V = \{ st_1 : \ & \text{Check the recent history of 5G fashion videos on Lili's phone,} \\
st_2 : \ & \text{Analyze the viewing history and generate Lili's preference tags,} \\
st_3 : \ & \text{Send Lili's preference tags to my phone,} \\
st_4 : \ & \text{Receive and parse the preference information on my phone,} \\
st_5 : \ & \text{Search on Taobao and purchase a gift that matches her preferences} \},
\end{aligned}
\tag{5.6}
$$

**Task Allocation:** In parallel with constructing the task dependency graph, the planning agent must also determine the execution endpoint for each subtask. The proposed framework is general: an endpoint may be either a physical device or a specific application on a device.

In the more general **cross-device** collaboration scenario, suppose the system has $k$ available devices, denoted as the set:

$$
D = \{d_1, d_2, \ldots, d_k\}.
$$

We define an allocation function $A : V \rightarrow D$ to map each subtask $st_i$ to a specific execution device $d_j$.

It is worth noting that **cross-application** collaboration can be regarded as a special case of this framework. When $k = 1$, the execution endpoint set is no longer the set of physical devices $D$, but rather the set of applications available on that single device. In this case, the target domain of the allocation function $A$ changes from $D$ to APP.

$$
\text{APP} = \{\text{APP}_1, \text{APP}_2, \ldots, \text{APP}_m\},
$$

Thus, in both cross-device and cross-application scenarios, the final output of the planning agent can be uniformly represented as a tuple $(G_T, A)$, where $G_T$ encodes the structural task dependencies, and $A$ specifies the execution endpoint (a concrete device or application) for each subtask.

For the above cross-device task "buy a suitable gift for Lili", the allocation result after task decomposition can be formalized as:

$$
A = \{st_1 \mapsto d_{\text{lily\_phone}}, \ st_2 \mapsto d_{\text{lily\_phone}}, \ st_3 \mapsto d_{\text{lily\_phone}}, \ st_4 \mapsto d_{\text{my\_phone}}, \ st_5 \mapsto d_{\text{my\_phone}}\},
$$
$$
E = \{(st_1, st_2), (st_2, st_3), (st_3, st_4), (st_4, st_5)\}.
$$
$$
\tag{5.7}
$$

Here, $A : V \rightarrow D$ is the device allocation function mapping each subtask to its execution device, while $E$ is the dependency set where $(st_i, st_j) \in E$ indicates that $st_i$ must be completed before $st_j$ can begin.

**Task State Synchronization:**   To ensure correct execution of cross-device tasks, we design a synchronization protocol: each subtask $st_i$ is assigned a state $S(st_i) \in \{\text{Pending}, \text{Running}, \text{Completed}, \text{Failed}\}$, and any subtask assigned to device $d_j$ can only start if all its predecessors have been completed:

$$\forall st_p \in \text{Pred}(st_i), \quad S(st_p) = \text{Completed} \tag{5.8}$$

where $\text{Pred}(st_i)$ denotes the set of all predecessor nodes pointing to $st_i$ in $G_T$. For example, in the case of the gift-purchasing task, my phone $d_{\text{my\_phone}}$ can only perform the subtask "search on Taobao and purchase a gift" after it has received the confirmation signal from Lili's phone $d_{\text{lily\_phone}}$ that the "preference information has been transmitted".

Through this centralized planning–distributed execution architecture, both complex cross-device collaboration tasks and long-horizon cross-application tasks can be decomposed into goal-oriented subtask sequences executed by low-level execution strategies on individual devices. At the same time, the state synchronization mechanism ensures orderliness and reliability in cross-device execution, thereby guaranteeing the overall task completion.

## 5.3   Cross-Application Task Transition Mechanism

**Table 5.1:** *Types and Core Dimensions of Cross-Application Tasks*

| Task Type | Requires Memory Maintenance | Proactive APP Switching | Linear Control Flow |
|---|---|---|---|
| Passively triggered system linkage tasks | ✗ | ✗ | ✓ |
| Proactively triggered data-passing tasks | ✓ | ✓ | ✓ |
| Proactively triggered collaborative multi-tasks | ✓ | ✓ | ✗ |

We define a **cross-application** task as one where, in order to accomplish a user's final intent, the execution path must actively or passively invoke, interact with, or reference the functionalities or data of two or more independent applications. The essence lies in the user's "final intent" rather than the operations themselves; cross-application behaviors are merely the objective manifestation of the technical path toward fulfilling that intent.

Cross-application tasks introduce unique challenges absent in single-application tasks, which can be deconstructed along several key dimensions:

- **Memory Maintenance**: Whether the agent needs to actively maintain internal memory across applications during execution, e.g., temporarily storing text fragments, screenshots, or other intermediate results.
- **Triggering Mechanism**: Whether cross-application switching is explicitly performed by the agent through multiple atomic operations (proactive triggering), or automatically executed by the system based on a single action (passive triggering).
- **Control Flow**: Whether the task execution process is linear and sequential, or non-linear, requiring concurrency or alternating handling.

Based on the above three dimensions and their combinations, we categorize the wide variety of cross-application tasks into representative types. These types essentially form a progressive ladder of difficulty, imposing increasing demands on the agent's planning, perception, and execution capabilities. The classifications and corresponding solutions are described as follows:

**Passively Triggered System Linkage Tasks:** This category of tasks, in our dimensional analysis, is characterized by no need for memory maintenance, passively triggered app transitions, and a linear control flow. Its core feature lies in the fact that the agent only needs to perform a single atomic action—such as clicking a link or button—after which the operating system's built-in deep linking mechanism automatically handles the application switch and context transfer. The agent itself does not need to proactively retain information or plan complex navigation steps. For example, when clicking the "Pay with Alipay" button in an e-commerce app, the system automatically transitions to the Alipay app for verification.

The main challenge here lies in accurately recognizing the key "trigger" UI element (e.g., the "Pay with Alipay" button in an e-commerce app or the "Share" button in a video app), rather than in complex process planning. Our solution focuses on enhancing the agent's perception and recognition ability. Specifically, we **augmented the training data with a large number of cross-application linkage samples**. This targeted fine-tuning enables the agent to reliably identify the correct switching button at the right time, effectively simplifying complex system-level interactions into a single successful recognition and click.

**Proactively triggered data-passing tasks:** These tasks require memory maintenance, proactive application switching, and linear control flow. The agent must retrieve information (e.g., text, images, files) from one app, store it temporarily in memory, and then perform a sequence of operations (e.g., switching apps, locating an input field, pasting) to transfer the information to another app. For example, copying an address from a browser app and then opening a maps app to paste the address and start navigation.

Such tasks demand explicit step-by-step planning and memory capabilities, making them core application scenarios of the **task decomposition and planning mechanism** described earlier. In this context, the high-level planning agent $\pi_{\text{plan}}$ decomposes the user's intent into a linear task graph $G_T$ consisting of multiple atomic operations. This structured task graph serves as an execution blueprint, precisely guiding the agent through a sequence of actions—including application transitions—while leveraging contextual information to transfer necessary data across steps.

**Proactively triggered collaborative multi-tasks:** These are the most complex type of tasks, characterized by extensive memory maintenance, proactive application switching, and non-linear control flow. The agent must simultaneously track the states of two or more apps and dynamically perform operations in one app based on information from another. The process often requires repeated switching and cross-referencing across applications. For example, in split-screen mode, a user may follow a tutorial in a video app while simultaneously operating and recording in a notes app.

Due to their high non-linearity, dynamic perception demands, and complex decision logic, these tasks often involve deep human-agent collaboration. They currently lie beyond the scope of our present work and are regarded as an important direction for future research.

## 5.4  Cross-Device Collaboration Mechanism



**Figure 5.1:** *Cross-device task classification quadrant chart based on synchronization mode and role relationship*

We define a cross-device task as one where, in order to accomplish a user's final intent, the execution path must be jointly carried out by two or more physically distinct computing devices (e.g., smartphone, computer, or devices belonging to different users), and task states or data must be synchronized across these endpoints. **The essence of "cross-device" lies in inter-device collaboration and data exchange, not merely the independent execution of the same functionality on multiple devices.**

Similar to cross-application tasks, cross-device tasks face inherent challenges due to contextual isolation and state differences among physical devices. We analyze these challenges along two key dimensions:

- **Data & State Synchronicity:** Whether data exchange between devices requires immediate response via real-time synchronization, or whether delayed, non-blocking asynchronous synchronization is acceptable.

- **Role Relationship:** Whether devices operate under a master-slave paradigm with asymmetric functionality, or in a peer-to-peer mode with functional parity.

By combining these two dimensions, all cross-device tasks can be mapped clearly into the four quadrants shown in Figure 5.1.

Regardless of whether tasks fall into the "task delivery" (Type I) in master-slave mode or "information sharing" (Type III) in peer-to-peer mode, they can essentially be decomposed into a set of well-defined, bounded subtasks. The previously described task decomposition framework, with its generation of $G_T$ and state synchronization mechanisms, effectively manages subtask dependencies and execution flows across devices, ensuring successful task completion.

However, for tasks requiring real-time synchronization (Types II and IV), the technical challenge shifts from macro-level "task planning" to micro-level "system engineering". For example, bottlenecks in "remote control" tasks lie in establishing low-latency, high-throughput event streams rather than generating semantic steps. Similarly, "real-time collaboration" tasks demand complex distributed consistency protocols to resolve conflicts in concurrent modifications. Addressing these problems requires a fundamentally different communication and state management infrastructure than our current framework. Consequently, our present work focuses on asynchronous collaboration scenarios, leaving real-time synchronization tasks as an important direction for future research.

# 6

# Efficiency: Three-layer Optimization of Reasoning Efficiency across Model–Data–Pattern

## 6.1 Efficient Model Selection for Edge Deployment

**Table 6.1:** *Comparison of large language models with different parameter sizes for GUI scenarios*

| Model | Developer | Open-source | Parameters (B) | Multimodal Input Support |
|---|---|---|---|---|
| GPT-4o | OpenAI | ✗ | N/A | ✓ |
| Claude 3.5 Sonnet | Anthropic | ✗ | N/A | ✓ |
| Gemini 2.5 Pro | Google | ✗ | N/A | ✓ |
| Qwen2-VL-72B | Alibaba | ✓ | 72B | ✓ |
| CogAgent | THUDM | ✓ | 18B | ✓ |
| Llama 4 | Meta | ✓ | 17B | ✓ |
| Ferret-UI 2 | Apple | ✓ | 8.4B | ✓ |
| MiniCPM-V | OpenBMB | ✓ | 8B | ✓ |
| Qwen2.5-VL-7B | Alibaba | ✓ | 7B | ✓ |
| Phi-3-vision | Microsoft | ✓ | 4.2B | ✓ |
| Llama 3-3B | Meta | ✓ | 3B | ✗ |
| Gemma 3-1B | Google | ✓ | 1B | ✗ |

Current powerful mobile agents generally rely on large-scale pretrained models, but this brings severe challenges of high inference cost and significant response latency. Such computational burdens make it difficult to deploy these models directly on resource-constrained edge devices such as smartphones. Meanwhile, existing cloud-based solutions introduce a new set of problems, including network latency, data privacy concerns, and operational costs. Therefore, exploring a technological path for migrating from cloud-based large models to efficient and lightweight edge-side models has become critical to the advancement of this field.

Between the ultimate vision of fully on-device deployment and the substantial overhead of current cloud-based large models, we must identify a balance between **model capability** and **inference efficiency**. To this end, after careful evaluation, we **selected a state-of-the-art open-source model with 8B (8 billion) parameters as the backbone of our mobile agent**.

This decision was made because a 8B-scale model occupies a strategic "sweet spot", as shown in Table 6.1. On the one hand, compared with giant models with tens or hundreds of billions of parameters, a 8B open-source model significantly reduces inference cost and latency, enabling the construction of an economical and efficient cloud service. On the other hand, compared with smaller models (e.g., 1B–3B scale), the 8B model crosses the critical threshold of complex reasoning and multimodal capabilities. Smaller models often struggle to reliably execute multi-step complex instructions, whereas the 8B model demonstrates sufficiently strong logical planning and tool-use abilities, making it a solid backbone for a reliable mobile agent.

More importantly, the 8B-scale model is at a critical inflection point—it paves the way toward future fully on-device deployment. It already has the potential for local deployment on high-end edge devices, such as flagship smartphones, laptops, or dedicated AI hardware. This enables our research not only to optimize current cloud-based services but also to provide a practical technical route toward realizing a fully on-device mobile agent with no reliance on network connectivity, enhanced privacy protection, and zero latency.

In summary, adopting the 8B model as the backbone represents the best practice under current technological conditions, balancing core model capabilities with deployment and operational costs. This approach not only allows us to deliver high-performance, low-latency agent services today, but also constitutes a forward-looking strategic layout for achieving fully autonomous on-device intelligence in the future.

## 6.2 Personalized Information Memory and Retrieval Mechanism

In the mobile agent system, to achieve deeply personalized task execution for users, **we design and implement a mechanism for personalized information memory and retrieval, enabling the agent to automatically accumulate user information and preferences across multi-round task interactions, and seamlessly invoke them in subsequent executions.** This mechanism effectively solves the problem of users repeatedly inputting similar information (e.g., phone numbers, addresses, ID numbers) in repetitive scenarios, thereby significantly reducing interaction costs and improving the naturalness and efficiency of task execution. Compared with the traditional "instant recognition–instant usage" paradigm, our mechanism achieves long-term personalized information management across tasks and rounds through persistence and dynamic updating, endowing the agent with stronger contextual continuity and user adaptation capabilities.

### 6.2.1 Information Collection and Structured Representation

The mechanism abstracts a user's personal information profile as an information repository, denoted as $P$. Each entry in the repository is represented as a triplet $(r, f, v)$, where $r$ denotes

the relational entity (e.g., "mother"), $f$ denotes the information field (e.g., "phone number"), and $v$ is the concrete value associated with the entity and field.

The collection process begins with parsing the user's natural language instruction $I$. The parsing engine identifies one or more target information slots $(r_t, f_t)$ from $I$. If the query over $P$ yields no corresponding value (i.e., $\neg\exists v : (r_t, f_t, v) \in P$), the system activates the Optical Character Recognition (OCR) module. This module parses newly appearing on-screen text and validates it using the regular expression $\text{RegEx}(f_t)$ associated with the field $f_t$. Once a text fragment $v_{\text{new}}$ passes validation, the system executes an update operation $\text{Update}(P, (r_t, f_t, v_{\text{new}}))$, storing the entry in the personal information repository.

### 6.2.2 Information Change Detection and Logging

To ensure the timeliness of repository $P$, the mechanism continuously monitors stored entries. For an existing entry $(r, f, v_{\text{old}}) \in P$, the system compares it against newly recognized data $v_{\text{new}}$ obtained through OCR in subsequent interactions. When $v_{\text{new}}$ meets two conditions—$v_{\text{new}} \neq v_{\text{old}}$ and $v_{\text{new}}$ passes $\text{RegEx}(f)$ validation—the system automatically triggers an update, replacing $v_{\text{old}}$ in the repository with $v_{\text{new}}$.

To ensure transparency and traceability, each successful update appends an audit record with a timestamp and change details to a separate change log repository $L_c$. This ensures that the historical evolution of all personal information fields can be precisely traced.

### 6.2.3 Seamless Invocation of Personal Information

The ultimate value of this mechanism lies in enabling "seamless invocation" of personalized information. When the user issues an instruction $I$ requiring stored information, the system, after parsing the target slot $(r_t, f_t)$, performs a retrieval operation from $P$:

$$v_{\text{retrieved}} = \text{Retrieve}(P, (r_t, f_t))$$

The retrieved value $v_{\text{retrieved}}$ is injected into the agent's context $C$. The model then generates the final execution sequence $A = \text{Model}(I, C)$ based on both the original instruction and the enhanced context. Since all required parameters are already complete, the system skips any manual intervention and directly executes the task, thereby achieving smooth and efficient automation.

Architecturally, the mechanism separates the high-frequency real-time query repository $P$ from the historical change log repository $L_c$. This design ensures low latency and high efficiency for daily retrieval operations, while the independent logging system guarantees data interpretability and security auditing, striking a balance between performance and robustness.

## 6.3 Efficient Reuse of Historical Behavioral Data

**To optimize the performance of agents in executing repetitive operational tasks, we propose a novel experience reuse framework.** This framework systematically accumulates historical behavioral data and leverages semantic retrieval based on LLMs to enable precise

replay of previously executed tasks. The core mechanism consists of three modules: structured accumulation of experiential data, semantic-driven retrieval, and log-based task replay. The framework aims to significantly reduce redundant reasoning overhead and execution latency in long-horizon or high-frequency tasks. Once a matching experience is successfully retrieved, the system bypasses complex planning and perception modules, entering a **"zero-reasoning"** execution mode, thereby greatly enhancing efficiency and stability.

### 6.3.1 Structured Accumulation of Experience Data

The foundation of experience reuse lies in the systematic and structured recording of historical task execution processes. We design an automated logging and accumulation mechanism to provide high-quality data sources for subsequent experience retrieval.

Each independent task execution is regarded as a potential experience. The system automatically captures and records the complete execution trajectory. A structured experience entry $E$ can be formalized as a tuple:

$$E = (q, S, A, \mathcal{T})$$

where: $q$ denotes the original user query, expressed as a natural language string. $S$ represents the final execution status, with $S \in \{\text{success}, \text{failure}, \text{in-progress}\}$. $A$ is an ordered action sequence, $A = \langle a_1, a_2, \ldots, a_n \rangle$, where each atomic action $a_i$ is itself a structured entry containing the UI screenshot and the executed operation (e.g., click, input). $\mathcal{T}$ includes task metadata such as start/end timestamps and execution duration.

Only when a task is successfully completed ($S = \text{success}$) is the record $E$ deemed valuable for reuse and archived into the experience pool $\mathcal{P}$. The pool $\mathcal{P} = \{E_1, E_2, \ldots, E_m\}$ grows dynamically as more tasks are executed, ensuring continuous evolution and expansion of coverage. This accumulation mechanism provides a robust data foundation for achieving high recall in experience matching.

### 6.3.2 Semantic-driven Experience Retrieval

The effectiveness of retrieval directly determines the performance of the entire reuse framework. Traditional keyword- or template-based matching struggles with the complexity and diversity of natural language. To overcome this, we introduce a semantic matching strategy powered by LLMs, enabling high-accuracy and robust task recognition.

When a new task query $q_{\text{current}}$ arrives, the system does not immediately proceed with the standard reasoning pipeline. Instead, it first attempts retrieval from the experience pool $\mathcal{P}$ using a semantic matching function $\mathcal{M}$. The goal of $\mathcal{M}$ is to identify a semantically equivalent query from the set of historical queries $Q_{\mathcal{P}} = q_i \mid (q_i, S_i, A_i, \mathcal{T}_i) \in \mathcal{P}$.

Formally:

$$E_{\text{matched}} = \mathcal{M}(q_{\text{current}}, \mathcal{P})$$

Specifically, $\mathcal{M}$ leverages an LLM $\mathcal{L}$ to evaluate the semantic similarity between $q_{\text{current}}$

and each $q_i \in Q_{\mathcal{P}}$. By constructing prompts that incorporate contextual information, the model considers multiple dimensions such as operation intent and target entities. If the model determines that a semantically equivalent query $q_k$ exists, it returns the corresponding experience $E_k$; otherwise, it returns $\emptyset$:

$$\mathcal{M}(q_{\text{current}}, \mathcal{P}) = \begin{cases} E_k & \text{if } \exists E_k \in \mathcal{P} \text{ s.t. } q_k \text{ matched } q_{\text{current}} \\ \emptyset & \text{otherwise} \end{cases}$$

The performance of this mechanism is evaluated by the Experience Hit Rate $H$, defined as the ratio of tasks successfully matched to reusable experiences $N_{\text{hit}}$ over the total number of tasks $N_{\text{total}}$:

$$H = \frac{N_{\text{hit}}}{N_{\text{total}}}$$

A higher hit rate is directly associated with overall system efficiency improvement.

### 6.3.3 Log-based Task Replay Execution

Once a matching experience $E_{\text{matched}}$ is retrieved, the system bypasses perception, planning, and decision-making modules, and activates the log replay mechanism.

The replay module directly reads the recorded action sequence $A_{\text{matched}} = \langle a_1, a_2, \ldots, a_n \rangle$ and executes each atomic action $a_i$ strictly in the original order. Since this process involves no real-time reasoning or visual analysis, execution speed far exceeds that of the standard pipeline.

We define the efficiency gain $\eta$ as follows. Let $T_{\text{std}}$ be the average execution time of the standard reasoning pipeline, and $T_{\text{replay}}$ be the average execution time of log replay:

$$\eta = \frac{T_{\text{std}} - T_{\text{replay}}}{T_{\text{std}}} = 1 - \frac{T_{\text{replay}}}{T_{\text{std}}}$$

According to our preliminary evaluation, in scenarios with stable UI layouts and fixed task structures, this mechanism can achieve efficiency gains of $40\% \sim 60\%$. Furthermore, log replay offers determinism and interpretability, providing strong support for task debugging, auditing, and the construction of more advanced memory systems.

## 6.4 Control Optimization via Function Calling

The function calling paradigm (Function Calling[1]) enables LLMs to automatically invoke external tool functions or application programming interfaces (APIs) through structured data formats. Although GUI-based control provides broad generality, it inherently suffers from limitations in execution efficiency, stability, and completeness of information acquisition. To build a more efficient and robust agent, **we introduce an API-based function calling paradigm and integrate it with GUI control, forming a hybrid control framework**.

---

[1] https://platform.openai.com/docs/guides/function-calling

The core idea of this framework is to expand the agent's action space. We define the total action space $A_{\text{total}}$ as the union of the GUI atomic operation space $A_{\text{GUI}}$ and the function calling action space $A_{\text{Function Call}}$:

$$A_{\text{total}} = A_{\text{GUI}} \cup A_{\text{Function Call}} \tag{6.1}$$

Here, $A_{\text{GUI}}$ includes operations that directly interact with UI elements such as 'click(element)' and 'input(element, text)' (see Table 3.6). Each action $a_i \in A_{\text{Function Call}}$ is a structured function call represented as a tuple:

$$a_i = (f_i, P_i), \quad P_i = \{p_k \mapsto v_k\}_{k=1}^m \tag{6.2}$$

where $f_i$ is a predefined function name, and $P_i$ is the set of key–value parameter pairs. In this framework, to ensure accuracy and reduce model dependency, function call actions are actively driven and selected by the user rather than determined by the model. Based on intent $T$, the user matches a suitable function and parameter list $a_{\text{match}} = (f_{\text{match}}, P_{\text{match}}) \in A_{\text{Function Call}}$ from the predefined domain-specific API toolkit. For each predefined parameter $p_k \in P_{\text{match}}$, the user may provide a custom value or use the system's default value $p_k^{(0)}$. Finally, the user executes the function call in a structured format on the terminal, completing the task at lower cost and shorter latency, with GUI operations serving as a general fallback. Moreover, the function-calling optimization is implemented as an independent, minimal, realizable module, enhancing control efficiency.

**Implementation and Application of Domain-specific API Toolkits**   To validate the effectiveness of this hybrid framework, we developed and integrated two domain-specific API toolkits to extend the agent's capabilities.

First, we incorporated an *email service toolkit $Tool_{\text{email}}$*, which encapsulates low-level mail transfer protocols into high-level semantic interfaces. For example, a core function $\text{send\_email(to,}$ $\text{subject, body, attachments)}$ abstracts what would otherwise be a complex GUI sequence— launching an app, navigating the interface, entering text, selecting files, and transferring files across devices—into a single deterministic function call. This dramatically improves automation efficiency in scenarios such as batch invoice sending or system alert notifications. Moreover, through its parameterized attachment interface, it resolves the challenge of cross-application file transfer on Android devices.

Second, leveraging existing open-source libraries, we integrated an *automated social media interaction toolkit $Tool_{\text{social}}$*, exemplified by the Bilibili platform. This toolkit provides atomic function calls such as $\text{like(video\_id)}$, $\text{coin(video\_id, count)}$, and $\text{search(keyword)}$. By invoking these tools, the agent can conduct deeper social interactions and information retrieval tasks. For instance, it can automatically search for content creators based on user preferences, or perform keyword-driven searches (e.g., "LLM pretraining"), which are difficult to achieve efficiently in a single step through pure GUI control.

**Methodological Analysis and Discussion**   Integrating function calling into a GUI-agent framework fundamentally represents a trade-off between *generality* and *efficiency*.

Function calling methods demonstrate superior accuracy and speed in predefined tasks, with determinism and structured data advantages unmatched by pure GUI methods. However, this advantage comes at the cost of generality. The core bottleneck lies in dependency on external APIs: the agent's functional boundary is constrained to the set of encapsulated interfaces, leaving it unable to handle unseen tasks and vulnerable to changes in third-party APIs.

Therefore, we argue that the optimal agent architecture is not a single pathway but an **organic integration of both**. Our proposed hybrid control framework embodies such a pragmatic solution: GUI control provides a universal fallback ensuring broad applicability, while function calling delivers performance leaps in domains where APIs are available. **This design allows the agent to operate any interface like a human while achieving superhuman efficiency and reliability in specific tasks like a program.**

# 7

# Experimental Results and Analysis

The experimental results and analysis are divided into two main parts: **Basic Capability Evaluation** and **Scenario Capability Evaluation**. The evaluation architecture is illustrated in Figure 7.1. The basic capability evaluation includes **Model Grounding Capability**, **Model General Capability**, and **Post-Training Model Capability**. The scenario capability evaluation covers **Basic Scenarios**, **Complex Scenarios**, and **Real Scenarios**. The Grounding capability evaluation primarily supports the core metric of accuracy, while the model general capability and post-training model capability mainly support generalization and long-horizon capability, further contributing to accuracy.

## 7.1 Basic Capability Evaluation

Basic capabilities form the technical foundation for addressing four core challenges. This section evaluates grounding capability, general model capability, and post-training optimization effects, providing fundamental evidence to validate core challenges of **generalization, accuracy, long-horizon capability, and efficiency**.

### 7.1.1 Grounding Capability Evaluation

**Task Definition**

This evaluation supports the core capability of **accuracy**, focusing on vision-language alignment and understanding through three sub-tasks:

- **Fun2Point Task**: Requires the model to locate coordinates of target UI components based on functional descriptions;
- **Text2Point Task**: Requires the model to locate positions of given text strings in target UIs;

**Figure 7.1:** *Experimental Evaluation Architecture*

- **Bbox2Text Task**: Requires the model to output text content within specified bounding box coordinates.

These tasks evaluate precise GUI understanding from functional semantics, text localization, and visual content extraction perspectives. Fun2Point and Text2Point focus on GUI localization, while Bbox2Text evaluates OCR capability.

**Evaluation Dataset**

The evaluation uses the CAGUI [113] benchmark—the first large-scale Chinese Android GUI benchmark for natural language-to-GUI mapping. It contains real-app screenshots with XML metadata, featuring detailed annotations (bounding boxes, text, component types). For Text2Point and Bbox2Text tasks, annotations are directly extracted from XML; Fun2Point uses VLM Qwen2.5-VL-72B [3] to generate functional descriptions. The benchmark includes 1,500 samples covering typical Chinese mobile scenarios.

**Metrics and Methodology**

**Accuracy** is the core metric. Outputs are standardized: Bounding box models (e.g., InternVL series [11][12][117]) use IoU (threshold=0.5); Coordinate-output models (e.g., AppCopilot) use spatial tolerance; GPT-4o [58] uses OmniParser [53] for layout parsing. Standardization ensures fair comparison.

**Baseline Models**

Baselines include closed-source (GPT-5-mini, GPT-5 [57], GPT-5 with grounding, GPT-4o, GPT-4o with grounding) and open-source models (Qwen2.5-VL-7B [3], InternVL2.5-8B/26B [12], OS-Genesis-7B [78], UI-TARS-7B [66], OS-Atlas-7B [96], Aguvis-7B [99]).

**Results and Analysis**

As shown in Table 7.1, AppCopilot outperforms all baselines with 71.3% average accuracy. It achieves 79.1% (Fun2Point), 76.5% (Text2Point), and 58.2% (Bbox2Text)—surpassing others by >53% in Bbox2Text. This demonstrates exceptional visual-text alignment in complex mobile UIs (small text, overlapping elements). Closed-source models (e.g., GPT-5) underperform in localization tasks; open-source models struggle in text extraction.

**Table 7.1:** *GUI Localization Accuracy on CAGUI Benchmark (Fun2Point/Text2Point/Bbox2Text). Bold and underline indicate best/second-best.*

| Model | Fun2Point | Text2Point | Bbox2Text | Avg |
|---|---|---|---|---|
| *Closed-Source* | | | | |
| GPT-5-mini | 28.5 | 27.5 | 23.9 | 26.6 |
| GPT-5 | 33.7 | 25.7 | 30.3 | 29.9 |
| GPT-5 w/ grounding | 47.3 | 47.9 | 30.3 | 41.8 |
| GPT-4o | 22.1 | 19.9 | 14.3 | 18.8 |
| GPT-4o w/ grounding | 44.3 | 44.0 | 14.3 | 34.2 |
| *Open-Source* | | | | |
| Qwen2.5-VL-7B | 59.8 | 59.3 | <u>50.0</u> | <u>56.4</u> |
| InternVL2.5-8B | 17.2 | 24.2 | 45.9 | 29.1 |
| InternVL2.5-26B | 14.8 | 16.6 | 36.3 | 22.6 |
| OS-Genesis-7B | 8.3 | 5.8 | 4.0 | 6.0 |
| UI-TARS-7B | 56.8 | <u>66.7</u> | 1.4 | 41.6 |
| OS-Atlas-7B | 53.6 | 60.7 | 0.4 | 38.2 |
| Aguvis-7B | <u>60.8</u> | **76.5** | 0.2 | 45.8 |
| **AppCopilot** | **79.1** | **76.5** | **58.2** | **71.3** |

### 7.1.2 Model General Capability Evaluation

**Task Definition**

This evaluation supports **long-horizon capability** and **generalization**, further enhancing **accuracy**. It assesses action prediction in multi-step scenarios, requiring multimodal reasoning to generate precise actions (click, swipe, text input, etc.) based on user instructions, history, and screenshots.

**Evaluation Datasets**

Five benchmarks: Four public (AndroidControl [44] [Low/High], GUI-Odyssey [52], AITZ [110]) and CAGUI (600 tasks, 4,516 images across 8 Chinese app domains).

**Metrics and Methodology**

Core metrics: **Type Match (TM)** and **Exact Match (EM)**. TM checks action type correctness; EM requires full parameter match (coordinates, text). Outputs are standardized: coordinates normalized to [0,1000], time units to milliseconds.

**Baseline Models**

Closed-source: GPT-5, GPT-5-mini, GPT-4o, Gemini 2.0 [25], Claude [1]. Open-source: Qwen2.5-VL-7B, UI-TARS-7B, OS-Genesis-7B, OS-Atlas-7B, Aguvis-7B, OdysseyAgent.

**Results and Analysis**

**Accuracy, Generalization, and long-horizon Capability Established**    Table 7.2 shows App-Copilot leading all benchmarks: AndroidControl-Low (94.4% TM, 90.2% EM), AndroidControl-

High (77.7% TM, 69.2% EM), GUI-Odyssey (90.9% TM, 75.0% EM), AITZ (85.7% TM, 76.4% EM), and CAGUI (96.9% TM, 91.3% EM).

Closed-source models (e.g., GPT-5) underperform in Chinese scenarios. GPT-OSS [56] was excluded due to text-only limitations. Open-source UI-TARS-7B approaches but lags in Chinese tasks.

**This validates AppCopilot's cross-lingual generalizability and complex task handling.**

**Table 7.2:** *Step-Level Action Prediction on Five Benchmarks (TM/EM). Bold/underline: best/second-best. *OS-Atlas uses different train/test split.*

| Model | AC-Low | | AC-High | | Odyssey | | AITZ | | CAGUI | |
|---|---|---|---|---|---|---|---|---|---|---|
| | TM | EM | TM | EM | TM | EM | TM | EM | TM | EM |
| *Closed-Source* | | | | | | | | | | |
| GPT-5-mini | 16.8 | 16.8 | 15.6 | 15.5 | - | - | 16.2 | 11.7 | 12.6 | 9.7 |
| GPT-5 | 69.8 | 69.1 | 56.3 | 55.2 | - | - | 56.0 | 31.3 | 70.6 | 33.0 |
| GPT-4o | - | 19.5 | - | 20.8 | - | 20.4 | 70.0 | 35.3 | 3.7 | 3.7 |
| Gemini 2.0 | - | 28.5 | - | 60.2 | - | 3.3 | - | - | - | - |
| Claude | - | 19.4 | - | 12.5 | 60.9 | - | - | - | - | - |
| *Open-Source* | | | | | | | | | | |
| Qwen2.5-VL-7B | 94.1 | 85.0 | 75.1 | 62.9 | 59.5 | 46.3 | 78.4 | 54.6 | 74.2 | 55.2 |
| UI-TARS-7B | **95.2** | **91.8** | **81.6** | **74.4** | 86.1 | 67.9 | <u>80.4</u> | <u>65.8</u> | <u>88.6</u> | <u>70.3</u> |
| OS-Genesis-7B | 90.7 | 74.2 | 65.9 | 44.4 | 11.7 | 3.6 | 20.0 | 8.5 | 38.1 | 14.5 |
| OS-Atlas-7B | 73.0 | 67.3 | 70.4 | 56.5 | 91.8* | 76.8* | 74.1 | 58.5 | 81.5 | 55.9 |
| Aguvis-7B | 93.9 | 89.4 | 65.6 | 54.2 | 26.7 | 13.5 | 35.7 | 19.0 | 67.4 | 38.2 |
| OdysseyAgent | 65.1 | 39.2 | 58.8 | 32.7 | <u>90.8</u> | <u>73.7</u> | 59.2 | 31.6 | 67.6 | 25.4 |
| AppCopilot | <u>94.4</u> | <u>90.2</u> | <u>77.7</u> | <u>69.2</u> | **90.9** | **75.0** | **85.7** | **76.4** | **96.9** | **91.3** |

**Efficiency Capability Established** Figure 7.2 shows Compact Schema reduces instruction length by >50% via action space optimization. As Table 7.3 demonstrates, instructions use minimal JSON: e.g., "POINT":[x,y] for coordinates, "TYPE":"交话费" for task type. Average length: 9.7 tokens.

Figure 7.3 visualizes the action space in a phone UI. It maps raw operations to optimized sequences under interface constraints, enabling efficient execution.



**Figure 7.2:** *Output Length Distribution*

**Table 7.3:** *Compact Schema Example Outputs*

| Example Outputs: |
| --- |
| {"POINT":[87,445],"STATUS":"start"} |
| {"POINT":[123,456],"to":"left"} |
| {"TYPE":" 交话费"} |
| {"PRESS":"BACK"} |
| {"STATUS":"finish"} |



**Figure 7.3:** *Action Space*

**Case Analysis**

**Task: Open YouTube, Search, and Play Video**

**Execution Process:**

**Chinese Task:** Figure 7.4: Trajectory for " 打开 YouTube，搜寻并播放 '音乐家的无聊人生' 频道视频": (a) Instruction → (b) Open folder → (c) Launch YouTube → (d) Tap search → (e) Input text → (f) Identify channel → (g) Play video. Validates Chinese instruction understanding and GUI operation stability.

**English Task:** Figure 7.5: Trajectory for "Open YouTube, search for Mr.Beast, play videos": (a) Instruction → (b) Open folder → (c) Launch YouTube → (d) Tap search → (e) Input "mrbeast" → (f) Identify channel → (g) Play video. Confirms cross-lingual capability and consistent "instruction→interaction→goal" pipeline.

**(a)** *Instruction*      **(b)** *Open folder*      **(c)** *Launch YouTube*      **(d)** *Tap search box*

**(e)** *Input text*      **(f)** *Identify channel*

**Figure 7.4:** *Chinese Task:* 打开 *YouTube*，搜寻并播放'音乐家的无聊人生'频道视频

**(a)** *Instruction*      **(b)** *Open folder*      **(c)** *Launch YouTube*      **(d)** *Tap search box*



**(e)** *Input "mrbeast"*      **(f)** *Identify channel*

**Figure 7.5:** *English Task: Open YouTube, Search, and Play "Mr.Beast"*

### 7.1.3 Post-Training Model Capability Evaluation

**Results and Analysis**

**Enhanced Accuracy, Generalization, and long-horizon Capability**    The post-training model evaluation serves as advanced validation for core capabilities like **long-horizon capability** and **generalization**, while significantly strengthening the foundational metric of **accuracy**. As shown in Table 7.4, performance comparisons before and after reinforcement fine-tuning are presented across all benchmarks. On datasets like AndroidControl-Low, GUI-Odyssey, and AITZ, reinforcement fine-tuning delivers substantial improvements, particularly in exact match accuracy. This outcome can be further explained by the reward curve during model training: as shown in Figure 7.6, the reward curve exhibits a steady upward trend during reinforcement fine-tuning, indicating that this phase effectively enhances the model's effectiveness in managing long-horizon reasoning and complex decision-making tasks, leading to superior performance on datasets requiring dynamic decision-making.

However, on benchmarks like AndroidControl-High and CAGUI, the model with only supervised fine-tuning already demonstrates competitive performance, even slightly outperforming in some cases. As shown in Figure 7.7, the loss curve during supervised fine-tuning reveals that these benchmarks provide sufficiently large and diverse training datasets. The model's loss decreases rapidly and stabilizes early, indicating it has fully learned patterns within the data. In such scenarios, mere imitation learning achieves strong generalization, and additional reinforcement learning yields limited marginal gains.



**Figure 7.6:** *Reinforcement Fine-Tuning Reward Curve*

**Table 7.4:** *Ablation Study: AppCopilot Performance Before/After Reinforcement Fine-Tuning*

| Model | AC-Low | | AC-High | | Odyssey | | AITZ | | CAGUI | |
|---|---|---|---|---|---|---|---|---|---|---|
| | TM | EM | TM | EM | TM | EM | TM | EM | TM | EM |
| AppCopilot-SFT | 87.6 | 83.1 | **78.6** | **69.5** | 86.1 | 66.7 | 79.0 | 61.1 | **96.9** | **91.5** |
| AppCopilot-RFT | **94.4** | **90.2** | 77.7 | 69.2 | **90.9** | **75.0** | **85.7** | **76.4** | **96.9** | 91.3 |

**Figure 7.7:** *Supervised Fine-Tuning Loss Curve*

## 7.2 Scenario Capability Evaluation

### 7.2.1 Basic Scenario Task Evaluation Results and Analysis

**Basic Communication Tasks**

**Task 1.1: Send SMS to grandson: "Are you coming back for dinner tonight?" / send message to mom, and tell her I will go back home**

**Task Objective:** This task demonstrates capabilities in **cross-lingual intent recognition, task decomposition, and contextual understanding**.

(a) *Instruction: Send SMS to grandson...*

(b) *Action: Tap SMS icon*

(c) *Action: Tap search box*

(d) *Action: Input "Grandson"*

(e) *Action: Select contact "Grandson"*

(f) *Action: Send SMS*

**Figure 7.8:** *Instruction: Send SMS to grandson: "Are you coming back for dinner tonight?"*

**(a)** *Instruction: send message...*  **(b)** *Action: Tap SMS icon*  **(c)** *Action: Tap search box*  **(d)** *Action: Input "mom"*



**(e)** *Action: Select contact "mom"*  **(f)** *Action: Send SMS*

**Figure 7.9:** *Instruction: send message to mom, and tell her I go back home*

**Execution Process:**

As shown in Fig. 7.8, during initial execution, the Agent completes **multi-level parsing** of the instruction. It accurately captures the core action intent "send SMS" and clearly distinguishes two key elements: the recipient "Grandson" and the message content "Are you coming back for dinner tonight?" **This structural decomposition capability allows the Agent to clarify element boundaries and logical relationships through intent recognition, preventing confusion between recipient and content.**

During entity association and localization, the Agent directly searches for the "Grandson" contact and triggers interaction. This process demonstrates precise **intent recognition and contextual understanding**. The core of task planning lies in mapping abstract kinship terms in natural language to specific contact entities, requiring both awareness of contact data structures and commonsense contextual understanding.

Additionally, the Agent accurately populates the message content without deviation, reflecting high-fidelity information capture and restoration capabilities. In the transformation from natural language to interactive content, the Agent avoids information omissions, semantic distortions, or formatting errors. The SMS content perfectly matches the instruction description, indicating high-fidelity information transformation capabilities crucial for executing **complex text instructions**.

The task successfully completes sending specific content to a specific recipient. Through the complete demonstration of the "send specific SMS to grandson" task, it showcases the Agent's capabilities in intent recognition, task decomposition, and contextual understanding. The system can accurately split and map natural language contacts and message content to specific interaction steps, maintaining execution accuracy and consistency even with incomplete information.

In the execution of the English instruction "send message to mom, and tell her I will go back home" (Fig. 7.9), the Agent similarly demonstrates precise parsing capabilities, accurately identifying the core action "send message", recipient "mom", and content "I will go back home", while fully inputting the English content. This verifies the system's language-agnostic parsing capability.

Through bilingual scenario demonstrations, this task comprehensively showcases the Agent's capabilities in intent recognition, task decomposition, and contextual understanding. The system can accurately map natural language contacts and message content to interaction steps, maintaining execution accuracy and consistency in multilingual environments.

**Task 1.2: Call Mom**

**Task Objective:** This task highlights capabilities in **ambiguous semantic parsing** and **user module invocation**.

**(a)** *Instruction: Call Mom*  **(b)** *Action: Open Phone*  **(c)** *Action: Navigate to "Calls" page*  **(d)** *Action: Tap search box*



**(e)** *Action: Input "Mommy"*  **(f)** *Action: Find "Mother"*  **(g)** *Action: Tap call button*

**Figure 7.10:** *Instruction: Call Mom (First Execution)*

**(a)** *Instruction: Call Mom*    **(b)** *Action: Open Phone*    **(c)** *Action: Input contact number*    **(d)** *Action: Initiate call*

**Figure 7.11:** *Instruction: Call Mom (Second Execution)*

**Execution Process:**

As shown in Fig. 7.10, during the first execution without user module parameters, the Agent relies entirely on real-time semantic parsing and interface interaction logic. During intent recognition, the Agent accurately identifies the core action "make call" and generalizes the colloquial "Mom" to the more universal search keyword "Mom m y". This process focuses on **dynamic matching of ambiguous entities** rather than precise positioning. The Agent completes entity association through interface search functions, essentially an **autonomous exploration** behavior without prior knowledge. During execution, the Agent strictly follows standard GUI interaction processes: open phone app → enter search interface → input keyword "Mommy" → match contact results → finally locate the contact corresponding to "Mother". This approach relies on the synonym expansion capability of the natural language understanding module and the element recognition accuracy of the GUI control module, adapting to possible variations in contact naming while providing visual verification opportunities.

As shown in Fig. 7.11, during the second execution with the **user module** enabled, the Agent's processing logic changes. The pre-stored "Mother" entity information in the user module becomes the core basis for **intent recognition** and **task decomposition and planning**. During intent recognition, the Agent directly links "Mom" in the instruction with the "Mom" entity recorded in the user module. During execution, the Agent retrieves the phone number corresponding to "Mother" by calling the user module interface and passes the number to the dialing function to complete the call. This path offers significant efficiency advantages, eliminating intermediate steps like search and matching, greatly reducing end-to-end response time. At the same time, since it relies on verified entity relationships in the user module, accuracy is more strongly guaranteed, **avoiding matching errors caused by semantic ambiguity**.

The system successfully completes the task of calling a person with a specific title. Through the two execution processes demonstrated above, it showcases capabilities in ambiguous semantic parsing, dynamic entity matching, and user module invocation.

**Transaction Service Tasks**

**Task 1.3: Recharge 50 yuan to this device's number via China Unicom**

**Task Objective:** This task highlights capabilities in **human-AI collaboration**.



**(a)** *Instruction: Recharge via China Unicom...*  **(b)** *Action: Open China Unicom*  **(c)** *Action: Tap "Recharge"*



**(d)** *Action: Select 50 yuan*  **(e)** *Action: Tap "Pay"*

**Figure 7.12:** *Instruction: Recharge 50 yuan to this device's number via China Unicom*

**Execution Process:**

As shown in Fig. 7.12, during initial instruction execution, the Agent accurately identifies and launches the China Unicom APP, directly locating the "Recharge" core function module. This process relies on its feature learning in vertical application scenarios, establishing direct mapping from instructions to functional entry points by recognizing UI components.

After entering the recharge interface, the Agent automatically matches "this device's number" by reading device-bound user data, transforming the abstract "this device" expression into

a specific phone number and filling it into the input box. For the "50 yuan" amount parameter, the Agent correctly selects the corresponding numerical value. The entire process shows no number mismatches or amount selection errors, validating its precise parsing capability for entity information and numerical parameters.

When the process advances to the payment stage, the Agent successfully navigates to the payment interface and actively pauses. This design is not accidental but incorporates reinforcement training for mechanisms requiring user feedback at critical nodes, aiming to enhance the Agent's control over practical application security and understanding of human-AI interaction boundaries. Since recharges involve sensitive financial transactions, users typically need to confirm transaction information accuracy and choose payment methods independently. Therefore, we emphasize in training that the Agent must actively stop the automation process at the payment interface rather than directly executing the payment action.

The task successfully completes the full process from application launch to payment interface, ultimately pausing at the pending payment interface awaiting user confirmation, demonstrating capabilities in vertical scenario function localization, precise entity parameter parsing, and human-AI collaboration.

**Task 1.4: Please select and purchase appropriate benefits in China Unicom based on my video viewing habits**

**Task Objective:** This task demonstrates capabilities in **user preference analysis, task decomposition, and personalized decision-making**.

**(a)** *Instruction: Please select benefits...*

**(b)** *Action: Open China Unicom*

**(c)** *Action: Close ad*

**(d)** *Action: Enter benefits interface*

**(e)** *Action: Tap search box*

**(f)** *Action: Select corresponding benefit*

**(g)** *Action: Tap "Pay"*

**(h)** *Action: Confirm payment*

**Figure 7.13:** *Instruction: Please select and purchase appropriate benefits in China Unicom based on my video viewing habits*

**Execution Process:**

As shown in Fig. 7.13, during initial instruction execution, the Agent demonstrates precise extraction of **user preferences**. The abstract requirement "video viewing habits" is resolved by invoking the user module to locate specific "Tencent Video" preference data. This process requires both stored behavioral records and the Agent's ability to interpret data, extracting core preferences from scattered behavioral data. This transformation of habits into specific service types represents the concrete application of user profiling.

After obtaining the "Tencent Video" preference, the Agent completes the entire screening and purchase process within the China Unicom APP. This choice reflects its accurate understanding of the instruction boundary to provide services within the Unicom ecosystem. The Agent achieves full-chain automation of data invocation, preference analysis, benefit screening, and purchase guidance. In practical operation, the Agent first **autonomously skips interface ads**, demonstrating generalization capabilities in unknown interaction scenarios and task execution efficiency. It then navigates to the benefits supermarket module, locates Tencent Video-related benefits through search, and guides to the purchase interface. This process involves multiple interface transitions and element interactions, yet the Agent maintains process continuity without path deviation.

This task demonstrates the core value of **Agent personalized service**. Unlike deterministic instructions like "recharge 50 yuan", "select appropriate benefits" requires the Agent to make subjective judgments based on user data. This judgment both **conforms to user preferences** and **fits the service scenario**. This capability breaks through the limitations of traditional GUI mechanical operations, achieving **active decision-making** based on user data—essentially the integration of user understanding and business understanding capabilities.

The task is successfully completed by invoking user module preference data for personalized recommendation and purchase, demonstrating capabilities in user preference analysis, business scenario boundary control, and personalized decision-making.

**Entertainment Content Tasks**

**Task 1.5: Browse short videos in Unicom Video Ringtone and like favorites**

**Task Objective:** This task demonstrates capabilities in **user preference analysis, multimodal information interaction, and autonomous execution**.

**(a)** *Instruction: Browse short videos...*

**(b)** *Action: Open Unicom Video Ringtone*

**(c)** *Action: Swipe to next video*

**(d)** *Action: Swipe to next video*

**(e)** *Action: Swipe to next video*

**(f)** *Action: Like*

**Figure 7.14:** *Instruction: Browse short videos in Unicom Video Ringtone and like favorites*

**Execution Process:**

As shown in Fig. 7.14, during initial instruction execution, the Agent demonstrates the ability to concretize the ambiguous requirement "favorites". By invoking the "scenery" preference tag from the user module, it **transforms abstract requirements into executable judgment criteria**. This process relies not only on stored preference data but also requires semantic understanding of preference tags.

During video content recognition, since Unicom Video Ringtone's short videos are dynamically loaded in a streaming manner, the Agent needs to analyze frame content in real-time to determine compliance with "scenery" preference. This heavily depends on the model's **visual capabilities**. In the application, the Agent swipes the screen to load new videos, executing a play-identify-judge loop until finding scenery-aligned content. When identified, it proactively triggers the like operation. This dynamic decision-making relies on collaboration between the GUI control module and content understanding module.

The task highlights the core value of **active service**. Beyond simply "browsing videos", the Agent incorporates active judgment and feedback based on preferences. This **active service** capability depends on the Agent's deep understanding of user needs.

The task is successfully completed by combining the "scenery" preference tag to achieve personalized browsing and proactive liking, demonstrating capabilities in user preference understanding, real-time content recognition, dynamic interaction decision-making, and active service awareness.

**Task 1.6: Watch Episode 10 of "Falling Flowers Meet Jun Again" on 5G Kuan Shijie and save to Unicom Cloud Drive**

**Task Objective:** This task demonstrates capabilities in **error tolerance** and **complex task decomposition**.

**(a)** *Instruction: Watch on 5G Kuan Shijie...*

**(b)** *Action: Open 5G Kuan Shijie*

**(c)** *Action: Tap search box*

**(d)** *Action: Input drama name*

**(e)** *Action: Tap search*

**(f)** *Action: Tap drama*

**(g)** *Action: Tap episode selector*

**(h)** *Action: Tap Episode 10*

**(i)** *Action: Tap share*

**(j)** *Action: Save to cloud drive*

**Figure 7.15:** *Instruction: Watch Episode 10 of "Falling Flowers Meet Jun Again" on 5G Kuan Shijie and save to Unicom Cloud Drive*

**Execution Process:**

As shown in Fig. 7.15, during initial instruction execution, the Agent accurately launches the 5G Kuan Shijie application and locates the drama through search functionality. This relies on understanding video application interaction logic. Precise drama name recognition prevents search failures.

During episode lookup error handling, the Agent demonstrates **active decision-making**. When only the first six episodes are displayed by default, the Agent proactively triggers the "All Episodes" list expansion. This behavior reflects problem-solving awareness—recognizing insufficient information and invoking alternative paths.

After successfully playing Episode 10, the Agent completes the subsequent "save to cloud drive" task. This involves functional linkage between 5G Kuan Shijie and Unicom Cloud Drive. The Agent locates the share-to-cloud function and confirms the storage operation.

The task validates the Agent's capability to recognize and operate complex interface elements. As a video application, 5G Kuan Shijie contains rich multimedia elements. Experimental results show the Agent can stably recognize various interface elements.

The task is successfully completed, demonstrating capabilities in contextual interaction understanding, error handling decision-making, cross-module collaboration, and interface adaptation robustness.

**Account and Information Query Tasks**

**Task 1.7: Help me generate an account sheet for July's recharge records and send it to "Little Assistant" on WeChat**

**Task Objective:** This task demonstrates capabilities in **generalization, long-horizon context modeling, complex task decomposition, and planning**.

**(a)** *Instruction: Generate account sheet...*

**(b)** *Action: Open China Unicom*

**(c)** *Action: Tap "Recharge"*

**(d)** *Action: Tap records, find July*

**(e)** *Action: Generate account sheet*

**(f)** *Action: Return to desktop*

**(g)** *Action: Open WeChat*

**(h)** *Action: Open chat with "Little Assistant"*

**(i)** *Action: Tap dialog box*

**(j)** *Action: Paste account sheet*

**(k)** *Action: Send account sheet*

**Figure 7.16:** *Instruction: Generate account sheet for July's recharge records and send to "Little Assistant" on WeChat*

**Execution Process:**

As shown in Fig. 7.16, this task chains four core processes: **data retrieval, information extraction, text generation, and cross-APP sending**. The entire instruction involves two applications with multiple critical steps.

The Agent progresses in logical order: first finds recharge records, locates "July records" through swiping, generates the account text, then switches to WeChat for sending. This **long-horizon coherence** stems from global task understanding.

During record retrieval, the Agent's swiping behavior adapts to real interaction scenarios. Unicom APP typically displays records in reverse chronological order. The Agent simulates human operation while understanding temporal filtering logic.

After finding July records, the Agent extracts key recharge information through **OCR recognition and structured information extraction** technologies.

After generating the account sheet, the Agent switches from China Unicom APP to WeChat. This involves system-level operations. Smooth execution demonstrates global control in multi-app environments.

**Tool Efficiency Tasks**

**Task 1.8: View the e-invoice for the last recharge in China Unicom and save as PDF to Personal Cloud in Unicom Cloud Drive**

**Task Objective:** This task demonstrates capabilities in **precise intent parsing, complex task decomposition, and planning**.

**(a)** *Instruction: View e-invoice...*

**(b)** *Action: Open China Unicom*

**(c)** *Action: Tap "Recharge"*

**(d)** *Action: Tap "E-invoice"*

**(e)** *Action: Tap "Issued Invoices"*

**(f)** *Action: View last recharge's e-invoice*

**(g)** *Action: Tap "Save PDF"*

**(h)** *Action: Tap "Save to Cloud Drive"*

**(i)** *Action: Select "Personal Cloud"*

**(j)** *Action: Tap "Confirm"*

**Figure 7.17:** *Instruction: View e-invoice for last recharge in China Unicom and save as PDF to Personal Cloud in Unicom Cloud Drive*

**Execution Process:**

As shown in Fig. 7.17, during initial instruction execution, the Agent navigates from homepage to "Recharge" module, then to "E-invoice" module. This long-horizon operation contains multiple key steps executed in logical order without path deviation or step omission. This coherence stems from deep understanding of Unicom APP business logic.

The execution addresses the pain point of invoice management. Traditional management requires manual lookup, download, and app switching. The Agent integrates multiple steps into one automated flow. Precise location of "last" record ensures users get the latest standardized invoice file.

The task is successfully executed, demonstrating capabilities in precise intent parsing, long-horizon process control, cross-functional module collaboration, and file format handling—comprehensively validating comprehensive service capabilities in complex business scenarios.

### 7.2.2 Complex Scenario Task Evaluation Results and Analysis

**Aging-Friendly and Vision-Impaired Tasks**

**Task 2.1: Listen to "A Record of a Mortal's Journey to Immortality" on Qimao Free Novel**

**Task Objective:** This task demonstrates **precise content localization** and **voice accessibility support**, highlighting intelligent interaction's core value in accessibility scenarios.

**(a)** *Instruction: Listen to novel...*

**(b)** *Action: Open Qimao Free Novel*

**(c)** *Action: Skip ad page*

**(d)** *Action: Tap search box*

**(e)** *Action: Search novel name*

**(f)** *Action: Tap confirm*

**(g)** *Action: Enter page and listen*

**Figure 7.18:** *Instruction: Listen to "A Record of a Mortal's Journey to Immortality" on Qimao Free Novel*

**Execution Process:**

As shown in Fig. 7.18, for elderly or visually impaired users, manually searching novels in complex interfaces presents challenges. After launching Qimao Novel APP, the Agent locates the search box, inputs the novel name, and quickly locks the target work. This approach minimizes visual recognition and manual input barriers.

**Activation and adaptation of the voice module** is core to this task. After entering the reading interface, the Agent activates **voice reading** rather than displaying text. For visually impaired users, this removes visual information barriers; for elderly users, it alleviates visual fatigue. During reading, the Agent automatically adapts to chapter transitions and controls speed/pauses.

The task is successfully completed, demonstrating precise content localization and voice accessibility support—providing equal access to digital content for elderly and visually impaired groups.

**long-horizon Tasks**

**Task 2.2: Check the highest-rated nearby restaurant on Dianping**

**Task Objective:** This task demonstrates capabilities in **task decomposition/planning** and **contextual understanding**.

**(a)** *Instruction: Check restaurant...*

**(b)** *Action: Open Dianping*

**(c)** *Action: Enter food page*



**(d)** *Action: Swipe page*

**(e)** *Action: View recommendation posts*

**Figure 7.19:** *Instruction: Check highest-rated nearby restaurant on Dianping (First Execution)*

**Execution Process:**

Fig. 7.19 shows the first execution choosing the "homepage recommendation post" path. Dianping's homepage typically aggregates personalized recommendations based on location and history. The Agent prioritizes using platform recommendations rather than starting from scratch. This path leverages platform algorithms to reduce operational steps while meeting needs.

Fig. 7.20 shows the second execution using the "search bar filtering" path, demonstrating **active control** over requirement boundaries. The core logic is **precise requirement decomposition**: first search "restaurant" for full coverage, then set "Nearby" distance filter, finally switch to "Highest Rating" sorting. Each step directly corresponds to core instruction conditions.

The path difference reflects the Agent's dynamic balance between **"efficiency" and "accuracy"**. The recommendation path prioritizes efficiency when content matches requirements; the search path prioritizes accuracy when recommendation reliability is uncertain. This balancing depends on real-time judgment of requirement clarity and recommendation reliability.

The task is successfully completed, demonstrating multi-strategy adaptation capabilities in complex application scenarios.

**Cross-APP Tasks**

**Task 2.3: Send my itinerary from TravelSky to WeChat contact Manager Wu**

**Task Objective:** This task validates capabilities in **intent recognition**, **complex task decomposition**, and **planning**.

**Execution Process:**

As shown in Fig. 7.21, the Agent navigates to the itinerary module and extracts core elements through **structured parsing** of the page. This automatically identifies and extracts key data (departure, destination, time, flight number), ensuring clear, complete information transmission.

After extraction, the Agent switches from TravelSky to WeChat—involving system-level operations. The Agent maintains state coherence: after completing extraction in TravelSky, it switches to WeChat via system operations, locates "Manager Wu" through search/contacts, enters chat interface, pastes integrated itinerary, and triggers sending. This cross-end execution relies on deep cognition of multi-app interface logic.

The task is successfully completed, demonstrating cross-APP collaboration, information extraction/integration, and long-horizon task control capabilities—validating full-chain intelligent service in complex scenarios.

**Cross-Device Tasks**

**Task 2.4: Buy a gift for Lili based on her 5G Kuan Shijie history**

**Task Objective:** This task demonstrates **multi-agent collaboration**, **user preference analysis**, and **error tolerance**.

**Execution Process:**

As shown in Fig. 7.22, Lili's device stores 5G Kuan Shijie viewing history data, while the

**(a)** *Instruction: Check restaurant...*

**(b)** *Action: Open Dianping*

**(c)** *Action: Tap search box*

**(d)** *Action: Search restaurant*

**(e)** *Action: Tap confirm*

**(f)** *Action: Tap sort method*

**(g)** *Action: Select "Nearby"*

**(h)** *Action: Tap confirm*

**(i)** *Action: Tap sort method*

**(j)** *Action: Select "Highest Rating"*

**(k)** *Action: View best-rated restaurant*

**Figure 7.20:** *Instruction: Check highest-rated nearby restaurant on Dianping (Second Execution)*

**(a)** *Instruction: Send itinerary...*

**(b)** *Action: Open TravelSky*

**(c)** *Action: Enter itinerary page*

**(d)** *Action: Extract info, return to desktop*

**(e)** *Action: Open WeChat*

**(f)** *Action: Open chat with Manager Wu*

**(g)** *Action: Tap dialog box*

**(h)** *Action: Paste itinerary*

**(i)** *Action: Send message*

**Figure 7.21:** *Instruction: Send my itinerary from TravelSky to WeChat contact Manager Wu*

**(a)** *Instruction: Buy gift...*



**(b)** *Action: Open 5G Kuan Shijie*



**(c)** *Action: Open history*



**(d)** *Action: Extract keywords from recent video*



**(e)** *(Cross-device to my phone)*



**(f)** *Action: Open Taobao*



**(g)** *Action: Tap search box*



**(h)** *Action: Input gift info*



**(i)** *Action: Tap search*



**(j)** *Action: Tap gift*



**(k)** *Action: Add to cart*

**Figure 7.22:** *Instruction: Buy gift for Lili based on her 5G Kuan Shijie history*

user's device completes gift purchasing. This demonstrates **cross-device multi-agent collaboration**, **user preference extraction**, and **cross-application decision-making**.

After authorization verification, the Agent locates the history module and extracts key information from the most recent video. But raw video lists can't directly guide gift selection. Here, the Agent extracts IP keywords from "Crayon Shin-chan" to infer potential interests. This transcends simple data transfer by achieving a leap from data to preference to demand through content understanding.

On the user's device, the Agent receives "Crayon Shin-chan" keywords and launches Taobao. It locates relevant gifts through search, maintaining process coherence across multiple operations.

Critically, the task highlights the core value of cross-device service—breaking device barriers to achieve precise data-to-service docking. Traditional scenarios require manual preference inquiry and product search; the Agent automates the entire process from data collection to product recommendation.

**Task 2.5: Purchase gifts based on Lili's and Fanfan's 5G Kuan Shijie history**

**Task Objective:** This task extends agent capabilities **from serving individual users to multi-user** collaborative paradigms.

**(a)** *Instruction: Purchase gifts...*

**(b)** *Action: Open 5G Kuan Shijie*

**(c)** *Action: Open history*

**(d)** *Action: Extract keywords*

**(e)** *(Cross-device to Lili's phone)*

**(f)** *Action: Open 5G Kuan Shijie*

**(g)** *Action: Open history*

**(h)** *Action: Extract keywords*

**(i)** *(Cross-device to my phone)*

**(j)** *Action: Open Taobao*

**(k)** *Action: Tap search box*

**(l)** *Action: Input gift info*

**(m)** *Action: Tap search*

**(n)** *Action: Select gift*

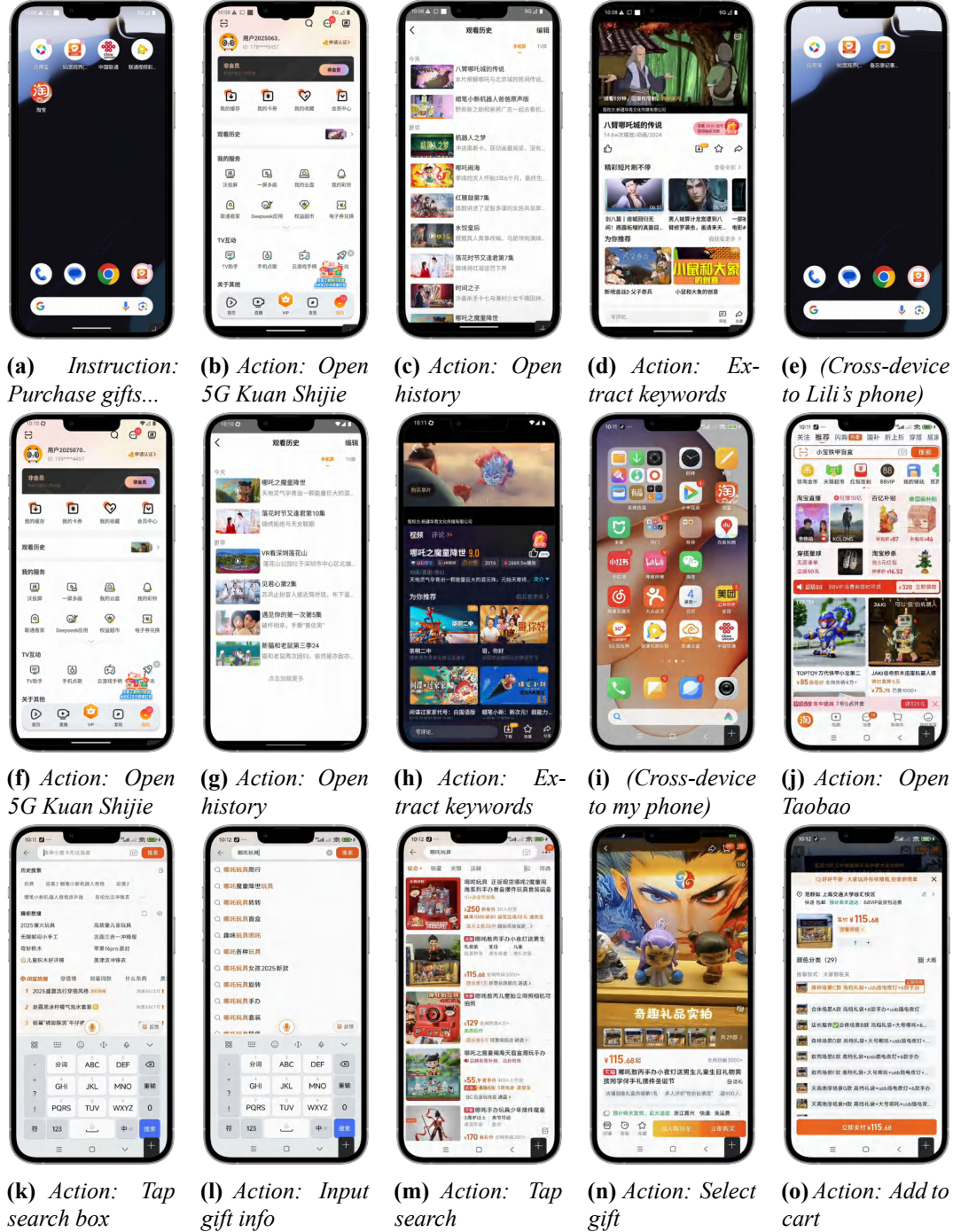**(o)** *Action: Add to cart*

**Figure 7.23:** *Instruction: Purchase gifts based on Lili's and Fanfan's 5G Kuan Shijie history*

**Execution Process:**

As shown in Fig. 7.23, Lili's and Fanfan's devices store viewing history data, while the user's device completes gift purchasing. This extends agent capabilities **from individual users to multi-user collaborative operations**. This isn't simple technical stacking but a paradigm shift from personal assistants to distributed collaboration networks.

Each agent focuses on parsing its own 5G Kuan Shijie history to **extract individual preferences**. The user's agent integrates preference data from both ends to drive targeted gift selection. In this architecture, each agent has independent computational space and decision boundaries, preserving core attributes representing individual intent while breaking limitations through collaboration.

The multi-device task confirms the feasibility of a **distributed intelligence system**. First, addressing reasoning and coordination challenges with incomplete information: Lili's and Fanfan's agents only know their own task status; the user's agent cannot directly access raw data on other devices. Through **data desensitization and intent recognition mechanisms**, agents collaborate accurately despite information gaps.

Second, addressing communication and negotiation mechanisms: agents achieve precise intent transmission through unified protocols despite heterogeneous systems. The successful execution validates that the **mobile agent system has upgraded from single-agent to a system-level architecture with multi-agent collaboration, distributed state modeling, and mechanism design capabilities**. This upgrade's core value enables intelligent services to break single-user boundaries and complete complex cross-domain long-horizon tasks through multiple autonomous agents collaborating—moving toward realizing the "theoretically expandable to massive terminals" vision of collective intelligence.

### 7.2.3 Real Scenario Task Evaluation Results and Analysis

In the real scenario capability evaluation phase, based on the technical framework of the AppCopilot system, our team independently developed an Android mobile application using Java as the development language.

This study conducted a system performance evaluation of the developed Android mobile application in real-world scenarios. The research designed a **complete chain of daily user behaviors** as the evaluation scenario:

The user wants to call a friend for dinner but finds the mobile phone has no credit, so they recharge; then opens Dianping to search for the highest-rated nearby restaurant; then needs to drive to the restaurant, so opens Amap for navigation; finally, after the meal, bookmarks the restaurant. This involves the following four subtasks:

1. **Task 3.1: Recharge credit on China Unicom**: Simulates users completing online credit recharge operations when discovering insufficient balance during calls. Specific instructions and process screenshots are shown in Fig 7.24;

2. **Task 3.2: Open Dianping to search for the highest-rated nearby restaurant**: Implements restaurant information retrieval and sorting (by descending rating) through a third-party platform (Dianping). Specific instructions and process screenshots are shown in Fig

7.25;

3. **Task 3.3: Navigate to the restaurant using Amap**: Calls map navigation service (Amap) to search for the destination and enter the navigation interface. Specific instructions and process screenshots are shown in Fig 7.26;

4. **Task 3.4: Bookmark the restaurant on Dianping**: Completes the post-consumption bookmarking operation on the third-party platform (Dianping). Specific instructions and process screenshots are shown in Fig 7.27.
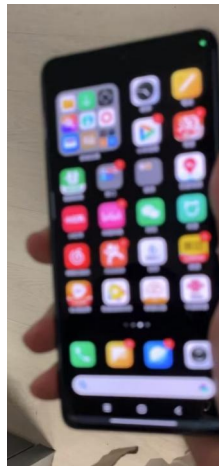
**Task Objectives:** This evaluation scheme aims to verify the system's generalization, accuracy, and efficiency in **real complex scenarios**. It simulates high-frequency daily user scenarios (recharge, restaurant search, navigation, bookmarking), demonstrating the practicality of the mobile application intelligent assistant for real needs.

**Task Execution Process:** Currently, there is a growing demand for seamless multi-application collaborative operations. This evaluation, through the complete closed-loop of credit recharge → restaurant search → navigation → bookmarking, reflects the potential of intelligent assistants in enhancing life efficiency and optimizing service experience. For example, automatic recharge reminders when credit is low, precise local service recommendations, and integrated navigation can significantly reduce user operation costs, especially benefiting groups unfamiliar with digital interactions (e.g., the elderly), promoting inclusive development of digital services.

As shown below, the Android mobile application developed based on AppCopilot can accurately recognize and jump to the corresponding APP to complete tasks in all four subtasks, demonstrating the system's generalization, accuracy, and efficiency in real complex scenarios.

**(a)** *Operation: Voice input command*



**(b)** *Operation: Open China Unicom APP*



**(c)** *Operation: Click to recharge credit*



**(d)** *Operation: Select recharge amount*



**(e)** *Operation: Enter payment page*



**(f)** *Operation: Confirm payment*

**Figure 7.24:** *Instruction: Recharge 20 yuan for this device number on China Unicom*

**(a)** *Operation: Voice input command*



**(b)** *Operation: Open Dianping*



**(c)** *Operation: Click "Highest Rating First"*



**(d)** *Operation: Locate the restaurant*

**Figure 7.25:** *Instruction: Search for the highest-rated nearby restaurant on Dianping*



**(a)** *Operation: Voice input command*



**(b)** *Operation: Enter Amap*



**(c)** *Operation: Search for Lao Xing Xian*



**(d)** *Operation: Enter navigation page*

**Figure 7.26:** *Instruction: Open Amap to navigate to Lao Xing Xian*

**(a)** *Operation: Open Dianping*    **(b)** *Operation: Search for Lao Xing Xian*    **(c)** *Operation: Enter the page*    **(d)** *Operation: Click bookmark*

**Figure 7.27:** *Instruction: Bookmark Lao Xing Xian shop on Dianping*

### 7.2.4 Scenario Summary

The above presents the main scenario evaluation tasks. Guided by the four core capabilities of accuracy, generalization, long-horizon capability, and efficiency, and combined with scenario characteristics, the evaluation plan is systematically designed into three major categories: Basic scenarios, complex interaction scenarios, and real-life work scenarios. The evaluation results are summarized in Table 7.5. The testing covers Basic scenarios, common general scenarios, and practical work-life scenarios. All scenarios were validated through standardized and detailed procedures. The results not only confirm the effectiveness of the solutions proposed in this study for the four core capability issues but also visually demonstrate that the project achieves precise task execution in both Basic scenarios and complex scenarios. Especially in long-horizon tasks, relying on the collaborative optimization of multiple mechanisms, the model exhibits stable and reliable execution performance, while the overall task execution efficiency is significantly improved.

**Table 7.5:** *Evaluation Metrics and Task Scenario Correlation Table*

| Evaluation Metric | Task Scenario | Associated Tasks |
|---|---|---|
| **Accuracy** | **Basic Scenarios** | |
| |   - Basic Communication Tasks | Task 1.1, Task 1.2 |
| |   - Transaction Service Tasks | Task 1.3, Task 1.4 |
| |   - Entertainment Content Tasks | Task 1.5, Task 1.6 |
| |   - Account and Information Query Tasks | Task 1.7 |
| |   - Tool Efficiency Tasks | Task 1.8 |
| | **Complex Scenarios** | |
| |   - Aging-Friendly and Vision-Impaired Tasks | Task 2.1 |
| |   - long-horizon Tasks | Task 2.2 |
| |   - Cross-APP Tasks | Task 2.3 |
| |   - Cross-Device Tasks | Task 2.4, Task 2.5 |
| | **Real Scenarios** | Task 3.1, Task 3.2, Task 3.3, Task 3.4 |
| **Generalization** | **Basic Scenarios** | |
| |   - Transaction Service Tasks | Task 1.4 |
| |   - Entertainment Content Tasks | Task 1.5, Task 1.6 |
| |   - Account and Information Query Tasks | Task 1.7 |
| |   - Tool Efficiency Tasks | Task 1.8 |
| | **Complex Scenarios** | |
| |   - Aging-Friendly and Vision-Impaired Tasks | Task 2.1 |
| |   - long-horizon Tasks | Task 2.2 |
| |   - Cross-APP Tasks | Task 2.3 |
| |   - Cross-Device Tasks | Task 2.4, Task 2.5 |
| | **Real Scenarios** | Task 3.2, Task 3.3, Task 3.4 |
| **long-horizon Capability** | **Basic Scenarios** | |
| |   - Entertainment Content Tasks | Task 1.6 |
| |   - Account and Information Query Tasks | Task 1.7 |
| |   - Tool Efficiency Tasks | Task 1.8 |
| | **Complex Scenarios** | |
| |   - long-horizon Tasks | Task 2.2 |
| |   - Cross-APP Tasks | Task 2.3 |
| |   - Cross-Device Tasks | Task 2.4, Task 2.5 |
| | **Real Scenarios** | Task 3.1, Task 3.2, Task 3.4 |
| **Efficiency** | **Basic Scenarios** | |
| |   - Basic Communication Tasks | Task 1.2 |
| |   - Transaction Service Tasks | Task 1.4 |
| |   - Entertainment Content Tasks | Task 1.5, Task 1.6 |
| |   - Account and Information Query Tasks | Task 1.7 |
| |   - Tool Efficiency Tasks | Task 1.8 |
| | **Complex Scenarios** | |
| |   - long-horizon Tasks | Task 2.2 |
| |   - Cross-APP Tasks | Task 2.3 |
| |   - Cross-Device Tasks | Task 2.4, Task 2.5 |
| | **Real Scenarios** | Task 3.2, Task 3.4 |

# 8
# Future Directions

With the continuous refinement of mobile agents in real-world business and complex terminal environments, substantial breakthroughs have been achieved across many technical aspects. Looking ahead, the further evolution of mobile agents will center on four core capabilities: generalization, accuracy, long-horizon capability, and efficiency. Continuous optimization along these dimensions will be crucial for enabling large-scale adoption, stable deployment, and innovative applications of mobile agents.



**Figure 8.1:** *Human-Agents-Physical Device Collaboration Architecture Diagram*

To intuitively illustrate its technical implementation logic, Figure8.1 establishes an interaction framework entitled "Human Layer - Agents Layer - Physical Device Layer" applicable to cross-device and cross-APP scenarios, which explicitly decomposes the complete chain through which human needs are transformed into physical world operations under the mediation of multi-agents.

Specifically, the Human Layer, as the initiator of needs, does not depend on a single carrier or scenario. Instead, it collaborates via cross-APP interactions (e.g., life service and device control applications) and cross-device terminals (e.g., mobile phones, tablets, and computers) to precisely transmit users' personalized needs, scenario-specific instructions, and associated information to the intermediate layer functioning as the "core transformation hub." This layer transcends the functional limitations of traditional single-agent systems and employs a multi-agent architecture inherently suited for collaboration, thereby effectively addressing complex needs in cross-scenario contexts. Ultimately, upon receiving tasks allocated by the multi-agents, the Physical Device Layer executes the instructions through cross-device coordination.

This framework not only resolves the functional constraints and adaptation bottlenecks inherent in traditional single-device and single-agent scenarios but also enhances the flexibility and efficiency of interactions between humans and physical devices through inter-layer collaborative design, providing concrete implementation support for the four core evolutionary directions (generalization, accuracy, long-horizon capability, and efficiency) elaborated earlier.

## 8.1 Generalization: Data Scale Expansion and Training Paradigm Innovation

The generalization ability of an agent determines its adaptability in unknown or rapidly changing environments. Although mobile agents have already demonstrated strong performance in real-world settings through multimodal and multi-scenario training, the depth and breadth of current data remain insufficient to sustain long-term development. Future optimization can focus on the following three aspects.

### 8.1.1 Expanding Heterogeneous Data Collection and Coverage

To equip the platform with "immunity" across different terminals, it is essential to broaden the coverage of data sources. This includes not only traditional desktop and mobile operating systems but also bilingual interfaces, emerging applications, and complex workflows. Currently, long-tail applications and mixed-language interfaces often lack sufficient annotation and interaction records, which limits model generalization. By building systematic annotation pipelines, enabling automated data collection, and leveraging device-side logs, the platform can more comprehensively capture user behaviors and interface features, enabling more robust reasoning when faced with unseen UIs. Moreover, heterogeneous data collection should not only cover more device types (e.g., different screen sizes and input modalities) but also consider diverse cultural and linguistic environments, thereby enhancing global adaptability.

### 8.1.2 Integrating Self-Supervised, Transfer, and Incremental Learning

Merely increasing data quantity is not sufficient to enhance generalization. More efficient training paradigms are required to maximize data utility. Self-supervised learning creates supervision signals from unlabeled data and has opened new possibilities across many domains. For rapid adaptation to new tasks, transfer learning allows knowledge acquired in previous tasks to be applied to related problems, enabling faster performance gains with limited data. Furthermore, incremental learning enables models to update continuously without retraining from scratch, helping them adapt to new environments while retaining old knowledge. Research shows that incremental learning is efficient, flexible, and scalable. Mobile agents should adopt these paradigms by leveraging self-supervision to reduce reliance on manual labeling, transfer learning to accelerate adaptation to new environments, and incremental learning to cope with shifting data distributions.

### 8.1.3 Establishing Open Data and Collaborative Learning Mechanisms

Sustained improvements in generalization require openness and collaboration. The open data movement initially aimed to make more data freely available to promote transparency and innovation. Building on this, researchers should explore Open Data 2.0 paradigms, where data owners retain sovereignty while enabling collaboration with model developers through decentralized data agents and federated learning. In federated learning, devices train models locally and upload only parameter updates, ensuring privacy while allowing performance improvements through aggregated learning. This paradigm has been proven effective in enabling collaborative innovation without leaking sensitive data. By promoting open data sharing mechanisms and collaborative training, mobile agents can expand into new domains more quickly, avoid redundant efforts, and ensure data security and compliance.

Overall, the key to enhancing an agent's generalization capability lies in the synergy between data scale expansion and training paradigm innovation. On the one hand, continuously broadening the scope of heterogeneous data collection—particularly targeting long-tail scenarios such as language-mixed interfaces, emerging applications, and complex workflows—enables the model to accumulate experience across more dimensions. On the other hand, novel training strategies such as self-supervised learning, transfer learning, and incremental learning transform data into powerful general-purpose capabilities. At the same time, collaborative mechanisms such as open data sharing and federated learning allow cross-enterprise knowledge transfer while preserving privacy. Together, these three aspects complement each other, laying a solid foundation for mobile agents to maintain continuous adaptability and cross-domain generalization in dynamic environments.

## 8.2 Accuracy: Dual-Track Integration of API and GUI with Enhanced Robustness

In enterprise-level deployment and large-scale adoption, operational accuracy serves as the core threshold for assessing platform value. Existing mobile agents have demonstrated that in scenarios with stable APIs and clear semantics, high task correctness can be achieved. However, in real-world business contexts involving long-tail applications and complex workflows, a single interface alone cannot cover all operations. Therefore, future systems must enhance robustness along two parallel tracks: on one hand, strengthening the stability and self-healing capacity of API invocation chains; on the other hand, leveraging GUI automation to ensure correctness when APIs are absent or unstable. By integrating GUI-based and API-based modes into a unified foundation model for both training and inference, and enabling hybrid scheduling, systems can achieve seamless cross-mode switching and on-demand complementarity, thereby maintaining consistent execution effectiveness and long-term accuracy across diverse business scenarios.

### 8.2.1 Robustness and Reliability in API Design

Robust interface design is the cornerstone of improving system accuracy. APIs with clear and consistent endpoint design can significantly reduce user operation errors, particularly in RESTful architecture contexts where standardized semantic conventions allow developers to quickly understand service logic and reduce integration difficulty. Moreover, interface versioning is essential for maintaining backward compatibility. Studies show that approximately 45% of enterprises face compatibility issues due to poor API management. Semantic versioning (Semantic Versioning) not only helps developers anticipate the impact of interface updates but also reduces the risk of failures during system evolution. Furthermore, the robustness and self-healing capacity of APIs form the core support for platform reliability. Interfaces equipped with comprehensive error-code systems and unified exception-handling strategies can provide clear feedback during failures, assisting developers in quickly locating and resolving issues. Related surveys indicate that 90% of front-line engineers rely primarily on detailed error information and recovery suggestions provided by APIs for troubleshooting. Consequently, future mobile agents should continuously strengthen the integration of interface abstraction, anomaly detection, self-recovery strategies, task-chain visualization, and log tracing within the API module, thereby enabling rapid integration and interpretable feedback to meet the stability requirements of complex enterprise workflows.

### 8.2.2 Generality and Redundancy Safeguards in GUI Automation

In many legacy systems, closed platforms, or emerging applications, the absence of stable or open APIs is the norm. GUI automation demonstrates strong generality and controllability in such contexts. By simulating mouse, keyboard, and touch gestures, it achieves precise control over front-end interfaces without relying on internal system permissions or invocation capabilities. Research shows that modern mobile agents can leverage visual recognition, OCR technologies, and object tracking to accomplish task transfer and widget identification across

multi-resolution, multilingual interfaces—particularly in scenarios with complex pop-ups, asynchronous loading, and multimodal interference. More importantly, GUI automation provides natural advantages of "zero-intrusion and low-dependency", enabling rapid deployment and iterative updates without altering business system source code. In stages where APIs are unavailable, frequently fail, or require rapid prototyping, the GUI module can serve as a high-accuracy alternative to ensure task continuity. Therefore, mobile agents in system design should always preserve an automated GUI control pathway as a redundant fallback channel when API robustness proves insufficient.

### 8.2.3 Dual-Track Collaboration and Dynamic Decision-Making

In the pursuit of high accuracy and robustness, future systems should not treat API and GUI modes as two isolated capability tracks. Instead, they must establish a flexible coordination mechanism that enables adaptive path selection under varying environmental conditions. Specifically, GUI-based and API-based modes should be integrated into a unified foundation model for joint training and collaborative inference, thereby bridging the gap between the two modes and achieving end-to-end seamless integration from planning to execution to verification. Integrating both modes into the same foundation model first requires unifying the action representations and semantic instructions of the two operation types at the training stage: abstracting API calls into functional actions, reducing GUI operations to visual primitives, and aligning semantics and learning strategies under a common plan–execute–verify framework. By jointly collecting and constructing API logs and GUI playback trajectories, the model can establish semantic mappings across modes and mitigate semantic gaps and perceptual noise during mode switching through contrastive learning and data augmentation.

At the inference stage, the foundation model no longer relies on a single mode but dynamically chooses strategies based on real-time availability, confidence, and execution costs—for example, prioritizing APIs when available to maximize efficiency, while using GUI operations for result verification, or rapidly falling back to GUI pathways upon API failure to prevent task interruption. Through such stepwise coordination and seamless switching mechanisms, systems can not only ensure accuracy but also significantly enhance robustness, achieving full coverage from stable API scenarios to real-world complex interfaces and addressing version changes, environmental variations, and unpredictable long-tail requirements in enterprise-level deployments.

Accuracy, as the foundational guarantee of intelligent agent systems in real-world applications, should evolve through dual-track synergy. On one side, API modules must continue strengthening semantic clarity and consistency in interface abstraction, while adopting semantic versioning, structured error feedback, and automatic rollback mechanisms to construct stable invocation chains with fault tolerance and maintainability. On the other side, GUI automation must enhance adaptability to multimodal inputs, cross-platform environments, and dynamic interface structures, improving widget recognition and workflow reconstruction for broad generalization, ensuring operability even in interface-lacking or heterogeneous environments. More critically, future systems should not only establish dynamic fusion scheduling between the two modes but

also unify GUI-based and API-based operations into a single foundation model for mixed training and collaborative inference. This will allow models to adaptively choose optimal execution paths across planning, execution, and verification, achieving seamless switching between APIs and GUIs. By integrating contextual states, historical execution traces, and API availability into the decision process, systems can prioritize APIs when stable to improve efficiency, while falling back to GUI pathways under API anomalies or low confidence to ensure continuity. This results in a closed-loop architecture spanning from interface abstraction to UI operations, from stable scenarios to critical exceptions, and from static orchestration to runtime reconstruction. Only through this evolution can mobile agents secure long-term accuracy and reliability amid business complexity and environmental uncertainty, ultimately becoming a trustworthy hub for universal intelligent interaction.

## 8.3 Long-Horizon Capability: Personalized Memory, Cross-Application Collaboration, and Ecosystem Synergy

For complex sequential tasks, conditional branching, and cross-application scheduling, long-horizon capability is a critical milestone for next-generation intelligent agents to move toward practical deployment. In this regard, limited rule sets or shallow sequence modeling are far from sufficient to cope with the diversity and uncertainty of real-world environments. To enable mobile agents to execute dozens or even hundreds of composite instructions, seamlessly switch across applications, and make decisions based on users' historical preferences and real-time intentions, systems must simultaneously enhance their "long-horizon capabilities" on the following three levels.

### 8.3.1 Deep User Data and Personalized Long-Term Memory

True long-horizon intelligence requires not only continuous task execution but also an understanding of users' long-term preferences and behavioral patterns, as well as the ability to make appropriate decisions in dynamically changing contexts. Recent studies suggest that AI models must be equipped with Long-Term Memory (LTM) to store and manage authentic user interaction data in order to achieve self-evolution during reasoning[77]. LTM not only captures statistical patterns of long-tail individual data but also promotes continuous self-evolution by supporting diverse experiences across different environments and agents. From this perspective, future mobile agents should build personalized databases by continuously collecting users' operation histories, preference feedback, and abnormal behaviors, and integrate them with the language model's contextual understanding and user intent modeling to support task decomposition and intelligent recommendations. By combining traceable long-term memory with short-term interaction feedback, agents can maintain coherence across multi-session operations, anticipate user needs in advance, and provide decisions better aligned with individual habits. Meanwhile, the construction of long-term memory must comply with privacy protection and data governance requirements, adopting mechanisms such as on-device storage and federated learning to ensure that user data is used solely for local personalization and safeguarded from unauthorized leakage.

### 8.3.2 Cross-Application Collaboration and Multi-Agent Self-Evolution

Long-horizon tasks often involve multiple applications, complex conditional branches, and lengthy operation chains, which pose tremendous challenges for a single agent. Academic research highlights that current "procedural" agents face significant limitations in handling cross-application instructions, mainly due to task complexity, application heterogeneity, and error propagation across multi-step executions. To address these challenges, frameworks such as MobileSteward[51] propose application-oriented multi-agent systems: dynamically recruiting specialized agents for different applications, explicitly associating tasks via scheduling graphs, delegating subtasks to application-specific StaffAgents, and mitigating error propagation through result evaluation. Moreover, "memory-based self-evolution" mechanisms are used to summarize successful experiences and iteratively improve execution capabilities. This line of thought suggests that long-horizon tasks require cross-application, multi-model collaboration, along with memory reuse mechanisms for continual improvement. Future mobile agents can draw inspiration from these strategies by dynamically planning cross-application execution paths according to task content, while leveraging accumulated experiences to achieve robust execution of long-chain tasks.

### 8.3.3 Ecosystem Co-Construction and System-Level Interface Openness

Model innovation alone cannot fully resolve the challenges of cross-application and long-horizon tasks; it must be supported by deep collaboration with operating systems and hardware vendors. Accessibility services, system-level scheduling interfaces, and notification systems provided by mobile and desktop platforms form the foundation for cross-application navigation, element recognition, and smooth operation. By collaborating with OS developers, device manufacturers, and third-party applications to open and standardize these underlying capabilities, intelligent agents can more reliably recognize interface elements, invoke system permissions, and achieve seamless data transfer and task coordination across applications. Furthermore, ecosystem synergy helps unify security standards and privacy protocols, driving collaborative advances in algorithmic safety, data protection, and user experience. For instance, within a co-built ecosystem, unified privacy permission management and transparent interaction logs can be designed to allow users to clearly understand how their data is used and retain full control. Such a trustworthy collaborative ecosystem will become a key enabler for the large-scale deployment of long-horizon intelligent agents.

In summary, the enhancement of mobile agents' long-horizon capabilities relies not only on the accumulation of user data and the establishment of personalized memory systems but also on the development of cross-application collaboration frameworks and the co-construction of ecosystem-level interfaces. By advancing simultaneously on these dimensions, systems can truly meet the demands of complex long-horizon tasks and take a crucial step toward practical deployment.

## 8.4   Efficiency:  Knowledge Density Enhancement and Lightweight Deployment on Devices

The inference efficiency and deployability of AI systems directly constrain the cost and user experience of large-scale applications. With the increasing model capabilities and the demand for multi-end collaboration, the efficiency of mobile agents must not only address model "slimming" but also incorporate dynamic scheduling between edge and cloud inference as well as fast task response. In the future, this capability can be advanced in the following aspects.

### 8.4.1   Model Compression and Knowledge Density Enhancement

Deep neural networks typically contain a large number of redundant weights, resulting in high computational overhead and energy consumption, making them difficult to deploy on resource-constrained mobile devices. To address this, researchers have developed various model compression techniques, such as quantization, pruning, and knowledge distillation, which significantly reduce model size without notably compromising accuracy. An experiment on convolutional models demonstrated that structured pruning reduced model size by 75%, dynamic quantization achieved a 95% reduction in parameter count, and combining pruning with quantization reduced model size overall by 89.7%, while improving accuracy by 3.8%. Ultimately, this enabled 92.5% accuracy with 20ms inference latency on edge devices[24]. These results suggest that proper pruning and quantization can substantially lower computational and storage demands while maintaining or even enhancing performance. Moreover, the study highlights that knowledge distillation transfers the knowledge of a large teacher model to a smaller student model, thereby achieving significant parameter compression while incurring almost no loss in accuracy. Therefore, for LLMs powering mobile agents, structural pruning can eliminate redundant connections, weight quantization can reduce precision requirements, and knowledge distillation can produce smaller yet equally effective student models. Together, these techniques increase the knowledge density per parameter and enable efficient on-device inference.

### 8.4.2   Dynamic Inference and Early Exit

Beyond model compression, dynamic inference strategies can also reduce inference costs. Early-exit techniques insert multiple exits in the intermediate layers of a network, allowing inference to stop early once sufficient confidence is achieved, thus lowering computational demand. A recent study proposed a class-exclusion-based early-exit mechanism: the model leverages intermediate features to eliminate irrelevant classes and continues inference only on the remaining ones; if a single candidate class remains at a certain layer, the inference stops immediately. Experiments showed that this early-exit mechanism significantly reduces inference costs, making it particularly suitable for mobile and embedded devices [87]. In mobile agents, similar dynamic routing or early-exit strategies can be applied to simpler interfaces or tasks. By adapting inference depth to input difficulty, unnecessary computation is avoided, achieving a flexible and streamlined inference process.

### 8.4.3 Edge-Cloud Collaboration and Dynamic Scheduling

Relying solely on cloud inference incurs high latency and communication costs, while performing all inference locally on devices is limited by compute capacity. To balance this trade-off, researchers have proposed edge-cloud collaborative architectures: inference is dynamically scheduled to run either on the device or in the cloud depending on task complexity and resource availability. A study on large language model services found that handling all requests in the cloud resulted in massive network transmission and energy overhead, while on-device deployment consumed less energy but failed to handle complex tasks. By adopting collaborative edge-cloud scheduling, overall energy consumption was reduced by more than 50% while still meeting user latency requirements[101]. Mobile agents can adopt a similar framework by implementing elastic local-cloud inference: tasks requiring high real-time responsiveness and low computational complexity are prioritized for local inference, while more complex tasks or those exceeding device capacity are dynamically migrated to the cloud or supported by cloud microservices. Task scheduling algorithms can then intelligently allocate resources between edge and cloud to balance cost and performance.

### 8.4.4 Caching of Frequent Tasks and Local Fine-Tuning

In everyday operations, user interactions are often concentrated around a limited set of high-frequency tasks, such as launching common applications, text input, and system settings. In these cases, caching inference results or fine-tuning lightweight local models can greatly improve responsiveness. Recent research in edge AI has shown that caching frequently used data and models across edge devices and cloud infrastructure can reduce both storage and communication overhead[89]. Similarly, leveraging local computation and preprocessing can reduce the volume of data transmitted to the cloud, thereby lowering latency. For mobile agents, local caching of UI elements and action sequences associated with high-frequency applications, combined with lightweight fine-tuning based on user preferences, can enable "instant responses". For long-tail scenarios, dynamic loading or cloud inference ensures generalization capability.

Enhancing efficiency is critical for the large-scale deployment of mobile agents. By leveraging pruning, quantization, and distillation to increase knowledge density, adopting early-exit and other dynamic inference mechanisms to reduce unnecessary computation, and incorporating edge-cloud collaboration and task caching, it becomes possible to reduce energy consumption and response times while preserving accuracy. These strategies, together with the capabilities of generalization, accuracy, and long-horizon reasoning discussed earlier, will jointly form the core competitiveness of future mobile agents, enabling them to deliver low-latency and highly reliable intelligent services in complex and dynamic environments.

## Conclusion

In summary, mobile agents are entering a new era of ecosystem development in intelligent automation, cross-platform operation, and continual learning. Importantly, these abilities should not be viewed as a mere summary of existing achievements, but rather as a vision for future

evolution. Through coordinated advancement across multiple dimensions and collaborative development with industry ecosystems, mobile agents are expected to become the next-generation foundational infrastructure for human-computer interaction and business automation. In this process, the four core capabilities—generalization, accuracy, long-horizon capability, and efficiency—will intertwine to define platform competitiveness, continuously driving new architectures and modes of collaboration to adapt to increasingly diverse and complex application environments.

# References

[1] Anthropic. *Introducing computer use, a new Claude 3.5 Sonnet, and Claude 3.5 Haiku.* 2024. URL: https://www.anthropic.com/news/3-5-models-and-computer-use.

[2] Jinze Bai et al. *Qwen technical report.* 2023. URL: https://arxiv.org/abs/2309.16609.

[3] Jinze Bai et al. *Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond.* 2023. URL: https://arxiv.org/abs/2308.12966.

[4] Adrien Bardes et al. *V-jepa: Latent video prediction for visual representation learning.* 2024. URL: https://openreview.net/forum?id=WFYbBOEOtv.

[5] Rishi Bommasani et al. *On the opportunities and risks of foundation models.* 2021. URL: https://arxiv.org/abs/2108.07258.

[6] Tom Brown et al. *Language models are few-shot learners.* 2020. URL: https://arxiv.org/abs/2005.14165.

[7] Yuxiang Chai et al. *AMEX: Android Multi-annotation Expo Dataset for Mobile GUI Agents.* 2025. URL: https://arxiv.org/abs/2407.17490.

[8] Jieshan Chen et al. *Object detection for graphical user interface: old fashioned or deep learning or a combination?* 2020. URL: http://dx.doi.org/10.1145/3368089.3409691.

[9] Weize Chen et al. *Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents.* 2024. URL: https://proceedings.iclr.cc/paper_files/paper/2024/hash/578e65cdee35d00c708d4c64bce32971-Abstract-Conference.html.

[10] Wentong Chen et al. *GUICourse: From General Vision Language Models to Versatile GUI Agents.* 2025. URL: https://arxiv.org/abs/2406.11317.

[11] Zhe Chen et al. *InternVL: Scaling up Vision Foundation Models and Aligning for Generic Visual-Linguistic Tasks.* 2024. URL: https://arxiv.org/abs/2312.14238.

[12] Zhe Chen et al. *Expanding Performance Boundaries of Open-Source Multimodal Models with Model, Data, and Test-Time Scaling.* 2025. URL: https://arxiv.org/abs/2412.05271.

[13] Kanzhi Cheng et al. *SeeClick: Harnessing GUI Grounding for Advanced Visual GUI Agents*. 2024. URL: https://arxiv.org/abs/2401.10935.

[14] Aakanksha Chowdhery et al. *Palm: Scaling language modeling with pathways*. 2023. URL: https://jmlr.org/papers/v24/22-1144.html.

[15] Gheorghe Comanici et al. *Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities*. 2025. URL: https://arxiv.org/abs/2507.06261.

[16] Seyed Shayan Daneshvar and Shaowei Wang. *GUI Element Detection Using SOTA YOLO Deep Learning Models*. 2024. URL: https://arxiv.org/abs/2408.03507.

[17] Mostafa Dehghani et al. *Scaling Vision Transformers to 22 Billion Parameters*. 2023. URL: https://arxiv.org/abs/2302.05442.

[18] Xiang Deng et al. *Mind2Web: Towards a Generalist Agent for the Web*. 2023. URL: https://arxiv.org/abs/2306.06070.

[19] Jacob Devlin et al. *Bert: Pre-training of deep bidirectional transformers for language understanding*. 2019. URL: https://aclanthology.org/N19-1423/.

[20] Alexey Dosovitskiy et al. *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*. 2021. URL: https://arxiv.org/abs/2010.11929.

[21] Zhengxiao Du et al. *Glm: General language model pretraining with autoregressive blank infilling*. 2022. URL: https://aclanthology.org/2022.acl-long.26/.

[22] El Hassane Ettifouri et al. *Visual Grounding Methods for Efficient Interaction with Desktop Graphical User Interfaces*. 2025. URL: https://arxiv.org/abs/2407.01558.

[23] Dylan J. Foster, Adam Block, and Dipendra Misra. *Is Behavior Cloning All You Need? Understanding Horizon in Imitation Learning*. 2024. URL: https://arxiv.org/abs/2407.15007.

[24] Samer Francy and Raghubir Singh. *Edge ai: Evaluation of model compression techniques for convolutional neural networks*. 2024. URL: https://arxiv.org/abs/2409.02134.

[25] Google DeepMind. *Introducing Gemini 2.0: Our New AI Model for the Agentic Era*. 2024. URL: https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/.

[26] Boyu Gou et al. *Navigating the Digital World as Humans Do: Universal Visual Grounding for GUI Agents*. 2025. URL: https://openreview.net/forum?id=kxnoqaisCT.

[27] Shane Griffith et al. *Policy shaping: Integrating human feedback with reinforcement learning*. 2013. URL: https://papers.nips.cc/paper_files/paper/2013/hash/e034fb6b66aacc1d48f445ddfb08da98-Abstract.html.

[28] Sirui Hong et al. *MetaGPT: Meta programming for a multi-agent collaborative framework*. 2024. URL: https://openreview.net/forum?id=VtmBAGCN7o.

[29] Neil Houlsby et al. *Parameter-efficient transfer learning for NLP*. PMLR, 2019. URL: https://proceedings.mlr.press/v97/houlsby19a.html.

[30] Edward J Hu et al. *Lora: Low-rank adaptation of large language models.* 2022. URL: https://openreview.net/forum?id=nZeVKeeFYf9.

[31] Jing Huang et al. *ScaleTrack: Scaling and back-tracking Automated GUI Agents*. 2025. URL: https://arxiv.org/abs/2505.00416.

[32] Xu Huang et al. *Understanding the planning of LLM agents: A survey*. 2024. URL: https://arxiv.org/abs/2402.02716.

[33] Zhiyuan Huang et al. *Spiritsight agent: Advanced gui agent with one look*. 2025. URL: https://openreview.net/forum?id=jY2ow7jRdZ.

[34] Lars Benedikt Kaesberg et al. *Voting or Consensus? Decision-Making in Multi-Agent Debate*. 2025. URL: http://dx.doi.org/10.18653/v1/2025.findings-acl.606.

[35] Jared Kaplan et al. *Scaling laws for neural language models*. 2020. URL: https://arxiv.org/abs/2001.08361.

[36] Geunwoo Kim, Pierre Baldi, and Stephen McAleer. *Language models can solve computer tasks*. 2023. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/7cc1005ec73cfbaac9fa21192b622507-Abstract-Conference.html.

[37] Hanyu Lai et al. *Autowebglm: A large language model-based web navigating agent*. 2024. URL: https://dl.acm.org/doi/10.1145/3637528.3671620.

[38] Mike Lewis et al. *BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension*. 2020. URL: https://aclanthology.org/2020.acl-main.703/.

[39] Guohao Li et al. *Camel: Communicative agents for" mind" exploration of large language model society*. 2023. URL: https://openreview.net/forum?id=3IyL2XWDkG.

[40] Junyou Li et al. *More Agents Is All You Need*. 2024. URL: https://arxiv.org/abs/2402.05120.

[41] Kaixin Li et al. *ScreenSpot-Pro: GUI Grounding for Professional High-Resolution Computer Use*. 2025. URL: https://arxiv.org/abs/2504.07981.

[42] Ning Li et al. *MobileUse: A GUI Agent with Hierarchical Reflection for Autonomous Mobile Operation*. 2025. URL: https://arxiv.org/abs/2507.16853.

[43] Tianle Li et al. *Long-context LLMs Struggle with Long In-context Learning*. 2024. URL: https://arxiv.org/abs/2404.02060.

[44]  Wei Li et al. *On the Effects of Data Scale on UI Control Agents*. 2024. URL: https: //arxiv.org/abs/2406.03679.

[45]  Xiang Lisa Li and Percy Liang. *Prefix-tuning: Optimizing continuous prompts for generation*. 2021. URL: https://aclanthology.org/2021.acl-long.353/.

[46]  Zhuowan Li et al. *Retrieval Augmented Generation or Long-Context LLMs? A Comprehensive Study and Hybrid Approach*. 2024. URL: https://arxiv.org/abs/2407.16833.

[47]  Aixin Liu et al. *Deepseek-v3 technical report*. 2024. URL: https://arxiv.org/abs/2412.19437.

[48]  Wei Liu et al. *Autonomous Agents for Collaborative Task under Information Asymmetry*. 2024. URL: https://arxiv.org/abs/2406.14928.

[49]  Yinhan Liu et al. *Roberta: A robustly optimized bert pretraining approach*. 2020. URL: https://openreview.net/forum?id=SyxS0T4tvS.

[50]  Yixin Liu et al. *Sora: A review on background, technology, limitations, and opportunities of large vision models*. 2024. URL: https://arxiv.org/abs/2402.17177.

[51]  Yuxuan Liu et al. *MobileSteward: Integrating Multiple App-Oriented Agents with Self-Evolution to Automate Cross-App Instructions*. 2025. URL: https://dl.acm.org/doi/10.1145/3690624.3709171.

[52]  Quanfeng Lu et al. *GUIOdyssey: A Comprehensive Dataset for Cross-App GUI Navigation on Mobile Devices*. 2025. URL: https://arxiv.org/abs/2406.08451.

[53]  Yadong Lu et al. *OmniParser for Pure Vision Based GUI Agent*. 2024. URL: https: //arxiv.org/abs/2408.00203.

[54]  James R Martin and Peter R White. *The language of evaluation*. 2003. URL: https: //link.springer.com/book/10.1057/9780230511910.

[55]  Douglas W Maynard. *Language, interaction, and social problems*. 1988. URL: https: //www.researchgate.net/publication/240083988_Language_Interaction_and_ Social_Problems.

[56]  OpenAI. *gpt-oss-120b and gpt-oss-20b Model Card*. 2025. URL: https://cdn.openai. com/pdf/419b6906-9da6-406c-a19d-1bb078ac7637/oai_gpt-oss_model_card. pdf.

[57]  OpenAI. *Introducing GPT-5*. 2025. URL: https://openai.com/index/introducing-gpt-5/.

[58]  OpenAI et al. *GPT-4o System Card*. 2024. URL: https://arxiv.org/abs/2410.21276.

[59]  Long Ouyang et al. *Training language models to follow instructions with human feedback*. 2022. URL: https://arxiv.org/abs/2203.02155.

[60] Charles Packer et al. *MemGPT: Towards LLMs as Operating Systems*. 2024. URL: https://arxiv.org/abs/2310.08560.

[61] Joon Sung Park et al. *Generative agents: Interactive simulacra of human behavior*. 2023. URL: https://dl.acm.org/doi/fullHtml/10.1145/3586183.3606763.

[62] Bhrij Patel et al. *AIME: AI System Optimization via Multiple LLM Evaluators*. 2024. URL: https://arxiv.org/abs/2410.03131.

[63] Chen Qian et al. *Chatdev: Communicative agents for software development*. 2024. URL: https://aclanthology.org/2024.acl-long.810/.

[64] Chen Qian et al. *Experiential Co-Learning of Software-Developing Agents*. 2024. URL: https://arxiv.org/abs/2312.17025.

[65] Yiyue Qian et al. *Enhancing LLM-as-a-judge via multi-agent collaboration*. 2025. URL: https://www.amazon.science/publications/enhancing-llm-as-a-judge-via-multi-agent-collaboration.

[66] Yujia Qin et al. *UI-TARS: Pioneering Automated GUI Interaction with Native Agents*. 2025. URL: https://arxiv.org/abs/2501.12326.

[67] Colin Raffel et al. *Exploring the limits of transfer learning with a unified text-to-text transformer*. 2020. URL: https://jmlr.org/papers/v21/20-074.html.

[68] Christopher Rawles et al. *Android in the Wild: A Large-Scale Dataset for Android Device Control*. 2023. URL: https://arxiv.org/abs/2307.10088.

[69] Christopher Rawles et al. *AndroidWorld: A Dynamic Benchmarking Environment for Autonomous Agents*. 2025. URL: https://arxiv.org/abs/2405.14573.

[70] Ivan Reznikov. *How Exactly an LLM Generates Text*. 2023. URL: https://www.linkedin.com/pulse/how-exactly-llm-generates-text-ivan-reznikov/.

[71] Timo Schick et al. *Toolformer: Language models can teach themselves to use tools*. 2023. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/d842425e4bf79ba039352da0f658a906-Abstract-Conference.html.

[72] John Schulman et al. *Proximal Policy Optimization Algorithms*. 2017. URL: https://arxiv.org/abs/1707.06347.

[73] Zhihong Shao et al. *DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models*. 2024. URL: https://arxiv.org/abs/2402.03300.

[74] Tianlin Tim Shi et al. *World of bits: an open-domain platform for web-based agents*. Sydney, NSW, Australia, 2017. URL: https://proceedings.mlr.press/v70/shi17a.html.

[75] Yucheng Shi et al. *MobileGUI-RL: Advancing Mobile GUI Agent through Reinforcement Learning in Online Environment*. 2025. URL: https://arxiv.org/abs/2507.05720.

[76] Noah Shinn et al. *Reflexion: Language agents with verbal reinforcement learning*. 2023. URL: https://dl.acm.org/doi/10.5555/3666122.3666499.

[77] Theodore Sumers et al. *Cognitive architectures for language agents*. 2024. URL: https://openreview.net/forum?id=1i6ZCvflQJ.

[78] Qiushi Sun et al. *OS-Genesis: Automating GUI Agent Trajectory Construction via Reverse Task Synthesis*. 2025. URL: https://arxiv.org/abs/2412.19723.

[79] Yuchen Sun et al. *GUI-Xplore: Empowering Generalizable GUI Agents with One Exploration*. 2025. URL: https://arxiv.org/abs/2503.17709.

[80] Rohan Taori et al. *Alpaca: A strong, replicable instruction-following model*. 2023. URL: https://crfm.stanford.edu/2023/03/13/alpaca.html.

[81] AgentSea Team. *Wave-UI: A Large-Scale Dataset for GUI Grounding*. 2024. URL: https://huggingface.co/datasets/agentsea/Wave-UI.

[82] Gemini Team et al. *Gemini: a family of highly capable multimodal models*. 2023. URL: https://arxiv.org/abs/2312.11805.

[83] Hugo Touvron et al. *Llama: Open and efficient foundation language models*. 2023. URL: https://arxiv.org/abs/2302.13971.

[84] Daniel Toyama et al. *Androidenv: A reinforcement learning platform for android*. 2021. URL: https://arxiv.org/abs/2105.13231.

[85] Ashish Vaswani et al. *Attention is all you need*. 2017. URL: https://proceedings.neurips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.

[86] Bryan Wang et al. *Screen2words: Automatic mobile ui summarization with multimodal learning*. 2021. URL: https://dl.acm.org/doi/10.1145/3472749.3474765.

[87] Jingcun Wang, Bing Li, and Grace Li Zhang. *Early-exit with class exclusion for efficient inference of neural networks*. IEEE, 2024. URL: https://ieeexplore.ieee.org/document/10595861.

[88] Ke Wang et al. *E-ANT: A Large-Scale Dataset for Efficient Automatic GUI NavigaTion*. 2024. URL: https://arxiv.org/abs/2406.14250.

[89] Xubin Wang and Weijia Jia. *Optimizing edge AI: a comprehensive survey on data, model, and system strategies*. 2025. URL: https://arxiv.org/abs/2501.03265.

[90] Jason Wei et al. *Chain-of-thought prompting elicits reasoning in large language models*. 2022. URL: https://dl.acm.org/doi/10.5555/3600270.3602070.

[91] Jason Wei et al. *Emergent abilities of large language models*. 2022. URL: https://arxiv.org/abs/2206.07682.

[92]  Chuan Wen et al. *Fighting Copycat Agents in Behavioral Cloning from Observation Histories*. 2020. URL: https://arxiv.org/abs/2010.14876.

[93]  Qianhui Wu et al. *GUI-Actor: Coordinate-Free Visual Grounding for GUI Agents*. 2025. URL: https://arxiv.org/abs/2506.03143.

[94]  Qingyun Wu et al. *Autogen: Enabling next-gen LLM applications via multi-agent conversations*. 2024. URL: https://openreview.net/forum?id=BAakY1hNKS.

[95]  Qinzhuo Wu et al. *MobileVLM: A Vision-Language Model for Better Intra- and Inter-UI Understanding*. 2024. URL: https://arxiv.org/abs/2409.14818.

[96]  Zhiyong Wu et al. *OS-ATLAS: A Foundation Action Model for Generalist GUI Agents*. 2024. URL: https://arxiv.org/abs/2410.23218.

[97]  Yuquan Xie et al. *Mirage-1: Augmenting and Updating GUI Agent with Hierarchical Multimodal Skills*. 2025. URL: https://arxiv.org/abs/2506.10387.

[98]  Yiheng Xu et al. *Aguvis: Unified pure vision agents for autonomous gui interaction*. 2024. URL: https://openreview.net/forum?id=FHtHH4ulEQ.

[99]  Yiheng Xu et al. *Aguvis: Unified Pure Vision Agents for Autonomous GUI Interaction*. 2025. URL: https://arxiv.org/abs/2412.04454.

[100]  Hui Yang, Sifu Yue, and Yunzhong He. *Auto-gpt for online decision making: Benchmarks and additional opinions*. 2023. URL: https://arxiv.org/abs/2306.02224.

[101]  Zheming Yang et al. *Perllm: Personalized inference scheduling with edge-cloud collaboration for diverse llm services*. 2024. URL: https://arxiv.org/abs/2405.14636.

[102]  Shunyu Yao et al. *Webshop: Towards scalable real-world web interaction with grounded language agents*. 2022. URL: https://proceedings.neurips.cc/paper_files/paper/2022/hash/82ad13ec01f9fe44c01cb91814fd7b8c-Abstract-Conference.html.

[103]  Shunyu Yao et al. *React: Synergizing reasoning and acting in language models*. 2023. URL: https://openreview.net/forum?id=WE_vluYUL-X.

[104]  Hongli Yu et al. *MemAgent: Reshaping Long-Context LLM with Multi-Conv RL-based Memory Agent*. 2025. URL: https://arxiv.org/abs/2507.02259.

[105]  Xiang Yue et al. *MMMU: A Massive Multi-discipline Multimodal Understanding and Reasoning Benchmark for Expert AGI*. 2024. URL: https://arxiv.org/abs/2311.16502.

[106]  Chaoyun Zhang et al. *Large Language Model-Brained GUI Agents: A Survey*. 2025. URL: https://arxiv.org/abs/2411.18279.

[107]  Chaoyun Zhang et al. *Large language model-brained gui agents: A survey*. 2025. URL: https://openreview.net/forum?id=xChvYjvXTp.

[108]  Chaoyun Zhang et al. *Ufo: A ui-focused agent for windows os interaction*. 2025. URL: https://aclanthology.org/2025.naacl-long.26/.

[109]  Chi Zhang et al. *Appagent: Multimodal agents as smartphone users*. 2025. URL: https://dl.acm.org/doi/full/10.1145/3706598.3713600.

[110]  Jiwen Zhang et al. *Android in the Zoo: Chain-of-Action-Thought for GUI Agents*. 2024. URL: https://arxiv.org/abs/2403.02713.

[111]  Shaoqing Zhang et al. *Dynamic Planning for LLM-based Graphical User Interface Automation*. Ed. by Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen. Miami, Florida, USA, 2024. URL: https://aclanthology.org/2024.findings-emnlp.70/.

[112]  Thomas T. Zhang et al. *Imitation Learning in Continuous Action Spaces: Mitigating Compounding Error without Interaction*. 2025. URL: https://arxiv.org/abs/2507.09061.

[113]  Zhong Zhang et al. *AgentCPM-GUI: Building Mobile-Use Agents with Reinforcement Fine-Tuning*. 2025. URL: https://arxiv.org/abs/2506.01391.

[114]  Andrew Zhao et al. *Expel: Llm agents are experiential learners*. 2024. URL: https://ojs.aaai.org/index.php/AAAI/article/view/29936.

[115]  Wayne Xin Zhao et al. *Dense text retrieval based on pretrained language models: A survey*. 2024. URL: https://dl.acm.org/doi/10.1145/3637870.

[116]  Xiutian Zhao, Ke Wang, and Wei Peng. *An Electoral Approach to Diversify LLM-based Multi-Agent Collective Decision-Making*. 2024. URL: https://arxiv.org/abs/2410.15168.

[117]  Jinguo Zhu et al. *InternVL3: Exploring Advanced Training and Test-Time Recipes for Open-Source Multimodal Models*. 2025. URL: https://arxiv.org/abs/2504.10479.