

# FedGUI: Benchmarking Federated GUI Agents across Heterogeneous Platforms, Devices, and Operating Systems

Wenhao Wang<sup>1,3,4,\*</sup>, Haoting Shi<sup>2,4,\*</sup>, Mengying Yuan<sup>5</sup>, Yiquan Lin<sup>1</sup>, Panrong Tong<sup>3</sup>, Hanzhang Zhou<sup>3</sup>, Guangyi Liu<sup>1</sup>, Pengxiang Zhao<sup>1</sup>, Yue Wang<sup>3,†</sup>, Siheng Chen<sup>2,4,†</sup>,

<sup>1</sup> Zhejiang University   <sup>2</sup> Shanghai Jiao Tong University   <sup>3</sup> Tongyi Lab

<sup>4</sup> Multi-Agent Governance & Intelligence Crew (MAGIC)   <sup>5</sup> Wuhan University

## Abstract

Training GUI agents with traditional centralized methods faces significant cost and scalability challenges. Federated learning (FL) offers a promising solution, yet its potential is hindered by the lack of benchmarks that capture real-world, cross-platform heterogeneity. To bridge this gap, we introduce FedGUI, the first comprehensive benchmark for developing and evaluating federated GUI agents across mobile, web, and desktop platforms. FedGUI provides a suite of six curated datasets to systematically study four crucial types of heterogeneity: cross-platform, cross-device, cross-OS, and cross-source. Extensive experiments reveal several key insights: First, we show that cross-platform collaboration improves performance, extending prior mobile-only federated learning to diverse GUI environments; Second, we demonstrate the presence of distinct heterogeneity dimensions and identify platform and OS as the most influential factors. FedGUI provides a vital foundation for the community to build more scalable and privacy-preserving GUI agents for real-world deployment. Code and data are publicly available at <https://github.com/wwh0411/FedGUI>.

## 1 Introduction

Recent advances in vision–language models (VLMs) have enabled the emergence of GUI agents that can perceive graphical user interfaces (GUI) and execute user instructions through sequential interactions (Zhou et al., 2025; Liu et al., 2025b; Lian et al., 2026). Traditional approaches to GUI agents largely rely on centralized data collection and manual labeling. While effective, such paradigm suffers from high data collection costs and limited scalability (Sun et al., 2024). Meanwhile, the widespread and frequent use of GUI devices naturally generates abundant supervisory signals, which could serve as a low-cost data source for training GUI agents

(Berkovitch et al., 2025). However, this real-world, large-scale data remains largely underutilized, as it cannot be publicly shared due to user privacy concerns (Zhang et al., 2024). This necessitates a distributed learning paradigm where each client trains locally on its own data without direct transmission (He et al., 2024).

Initial research has explored this via federated learning (FL) for privacy-preserving collaborative training (Kuang et al., 2023). FedMABench (Wang et al., 2025c) is the first benchmark designed for federated mobile agents, but it is limited to collaboration among Android users and overlooks the significant potential of incorporating users from web and desktop environments to further enhance performance. Moreover, FedMABench does not account for broader forms of heterogeneity across devices, operating systems (OSs), and data sources.

These limitations give rise to two fundamental challenges: (1) **RQ1**: How to enable cross-platform collaboration of GUI agent training, and does the expanded collaboration from distinct platforms improve performance? (2) **RQ2**: How to quantitatively characterize and measure the real-world heterogeneity spanning diverse platforms, OSs, devices, and data sources.

To address the challenges outlined above, we introduce **FedGUI**, a comprehensive benchmark designed for distributed GUI agents across diverse platforms and devices. FedGUI is characterized by three key features: (1) **Diversity**. FedGUI covers a wide range of real-world GUI environments, including over 900 mobile applications, 40+ desktop applications, and 200+ websites. It supports both multi-step and cross-application tasks, enabling the evaluation of agents partitioned with different levels of complexity. (2) **Comprehensiveness**. FedGUI integrates seven representative FL algorithms and supports more than 20 base models, including state-of-the-art open-source VLMs and proprietary models. In addition, FedGUI provides a

\*Equal contributions. †Corresponding authors.

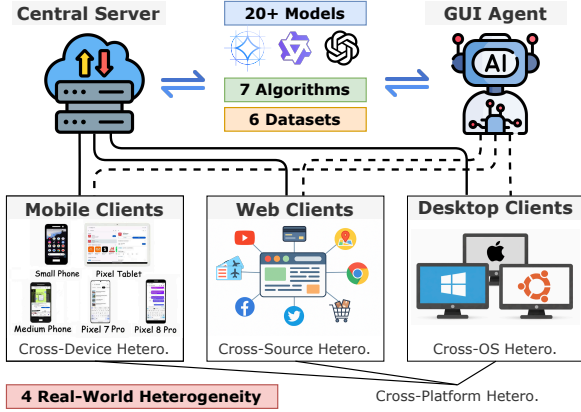


Figure 1: Overview of FedGUI. FedGUI provides a federated framework that coordinates a central server and heterogeneous clients across mobile, web, and desktop platforms to train a generalized cross-platform GUI agent. Hetero. is short for heterogeneity.

comprehensive set of evaluation metrics that jointly assess task performance and system efficiency. (3) **Heterogeneity.** FedGUI models four representative real-world heterogeneity settings, simulating the complexity of user collaboration across different platforms, devices, and operating systems, thereby reflecting realistic deployment scenarios.

Specifically, FedGUI constructs six datasets from eight data sources to study four different heterogeneities. (1) *FedGUI-Platform* is constructed using samples collected from three platforms (i.e., mobile, web and desktop), aiming to model **cross-platform heterogeneity**. (2) We further evaluate **cross-device heterogeneity** within the same mobile platform. Samples are collected from five different devices, forming the *FedGUI-Device* dataset. (3) To study **cross-OS heterogeneity**, we construct *FedGUI-OS*, which comprises samples collected from desktop environments on Ubuntu, macOS, and Windows. (4) To isolate and analyze the impact of data sources, we additionally construct more controlled datasets within each platform. *FedGUI-Web* and *FedGUI-Mobile* are designed to measure **cross-source heterogeneity**, where data originate from different underlying datasets. (5) Finally, *FedGUI-Full* is introduced to jointly capture both cross-platform and cross-source heterogeneity, providing a comprehensive benchmark for complex realistic distributed settings.

Based on FedGUI, we conduct an exhaustive empirical study to investigate and compare performance across three key axes: federated algorithm, heterogeneity level, and base model. Our extensive experiments yield several key observations:

| Benchmark     | Distributed Training | Cross-Platform Evaluation |     |         |
|---------------|----------------------|---------------------------|-----|---------|
|               |                      | Mobile                    | Web | Desktop |
| Mobile-Bench  | ✗                    | ✓                         | ✗   | ✗       |
| MMBench-GUI   | ✗                    | ✓                         | ✓   | ✓       |
| FedMABench    | ✓                    | ✓                         | ✗   | ✗       |
| FedGUI (Ours) | ✓                    | ✓                         | ✓   | ✓       |

Table 1: Comparison of related benchmarks.

| Dimension     | FedMABench  | FedGUI (Ours)                   |
|---------------|-------------|---------------------------------|
| Platform      | Mobile Only | Mobile, Web, Desktop            |
| Supported OS  | Android     | Android, Ubuntu, mac, Win       |
| Data Sources  | AC, AitW    | 9 Diverse Datasets              |
| Heterogeneity | User-based  | 4 Types (Plat., Dev., OS, Src.) |
| Eval Metric   | TF-IDF Sim. | Action Type, Grounding, SR      |

Table 2: Detailed benchmark comparison between FedGUI and FedMABench (Wang et al., 2025c).

(1) Increasing user participation in federated learning boosts model performance, even when these users contribute from highly diverse platforms and devices. (2) Platform-level heterogeneity poses a more significant challenge to model performance than within-platform heterogeneity. (3) Adaptive FL algorithms (e.g., FedYogi (Reddi et al., 2020)) outperform other baselines, demonstrating particular robustness in cross-platform settings. To summarize, our contributions are:

1. We present FedGUI, a comprehensive and unified benchmark for training cross-platform GUI agents through federated learning.
2. We construct six datasets targeted for four type of real-world heterogeneity. The datasets covers diverse applications and websites.
3. Extensive experiments thoroughly investigate federated GUI agents across platforms and devices, revealing insightful discoveries.

## 2 Related Work (Summarized in Table 1)

### 2.1 Centralized GUI Agents

**Single-Platform Agents.** Previous research has explored GUI agents for individual platforms, with mobile agents (Ye et al., 2025; Zhang et al., 2026b) focusing on app navigation, web agents (Ning et al., 2025; Chen et al., 2025) specializing in browser-based tasks, and desktop agents (Xie et al., 2024) targeting workflow automation on computing environments. However, these single-platform approaches exhibit limited cross-platform generalization, as their specialized capabilities are confined to corresponding environments and fail to transfer effectively across different interface modalities.

| Dataset Name    | Heterogeneity           | Source Datasets    | N# Clients | N# Subsets | N# Epi. | N# Steps |
|-----------------|-------------------------|--------------------|------------|------------|---------|----------|
| FedGUI-Platform | Cross-Platform          | AC, GA, AS         | 15         | 3          | 3,000   | 23,157   |
| FedGUI-Device   | Cross-Device            | GO                 | 5          | 4          | 2,500   | 38,064   |
| FedGUI-OS       | Cross-OS                | AS, OA-Mac, OA-Win | 3          | 5          | 1,800   | 5,813    |
| FedGUI-Web      | Cross-Source            | M2W, GA-W, OA-W    | 3          | 5          | 1,800   | 9,975    |
| FedGUI-Mobile   | Cross-Source            | AC, AitW, GO       | 3          | 5          | 6,000   | 59,328   |
| FedGUI-Full     | Cross-Platform & Source | All Nine Sources*  | 9-36       | 7          | 5,400   | 33,341   |

\* **Source Datasets:** **AC:** AndroidControl; **GA:** GUIAct; **AS:** AgentSynth; **GO:** GUI Odyssey; **OA-W/Mac/Win:** OmniAct-Web/MacOS/Windows; **M2W:** Mind2Web; **AitW:** Android-in-the-Wild.

\* **Evaluation Datasets:** Prior to constructing the six training datasets, we create a dedicated test set of 100 episodes for each of the nine sources. These test sets serve as a fixed benchmark for all subsequent evaluations.

Table 3: Statistic summary of FedGUI’s datasets. The table details the heterogeneity type, source composition, and scale for each of the six sets. N# is short for "the number of" and Epi. represents "episode"

**Cross-Platform Agents.** To address the generalization limitations, recent work has shifted toward cross-platform GUI agents capable of operating across diverse interface environments. Representative efforts span both models (Qin et al., 2025) and benchmarks (Wang et al., 2025d), demonstrating unified reasoning across mobile, web, and desktop platforms. Despite their promising performance, these approaches remain constrained by centralized data collection pipelines that are costly to scale and fail to capture the full spectrum of data diversity encountered in real-world distributed scenarios.

## 2.2 Distributed GUI Agents

**Federated Learning.** FL (Wen et al., 2023) offers a decentralized paradigm in which models are trained collaboratively without transmitting raw data. By performing local optimization and periodic server-side aggregation, FL enables large-scale learning under privacy constraints while naturally capturing user-specific behaviors.

**Federated GUI Agents.** FedMABench (Wang et al., 2025c) represents the first benchmark for federated mobile agents. However, as detailed in Table 2, its scope is fundamentally limited: it only considers collaboration among mobile device users, ignoring the broader context where users from different platforms can contribute to a unified agent system. In contrast, FedGUI presents the first comprehensive benchmark that addresses heterogeneous data sources, device types, and operating systems in federated agent training, benchmarking across a substantially broader spectrum of scenarios.

## 3 FedGUI

### 3.1 System Overview.

FedGUI presents a comprehensive benchmark accompanied by six datasets that emphasize multi-

dimensional heterogeneity and diversity in real-world mobile, web, and desktop GUI interactions. As illustrated in Figure 1, FedGUI follows the typical federated learning protocol and offers a research-friendly infrastructure that integrates 7 representative federated learning algorithms and supports over 20 mainstream VLMs.

To systematically characterize the complex non-IID patterns inherent in GUI interaction data, we categorize heterogeneity along four key dimensions: cross-platform, cross-device, cross-source, and cross-OS. Based on these dimensions, we construct six datasets with 29 subsets that span varying degrees of data skew and client scales, with the statistics summarized in Table 3.

### 3.2 Data Collection (Details in Appendix C)

**Data Source.** We collect a diverse set of GUI interaction datasets covering mobile, web, and desktop platforms, including AndroidControl (AC) (Li et al., 2024), Android-in-the-Wild (AitW) (Rawles et al., 2023), GUI Odyssey (GO) (Lu et al., 2025), GUIAct (GA) (Chen et al., 2025), Mind2Web (M2W) (Deng et al., 2023), OmniAct (OA) (Kapoor et al., 2024), and AgentSynth (Xie et al., 2025). These datasets differ in task complexity, platform, device and sources, enabling the study of heterogeneous data in federated settings.

**Unified Action Space.** We design a unified action space to standardize user interactions across mobile, web, and desktop platforms. Specifically, we identify six basic actions shared across all platforms (e.g., CLICK, TYPE), while mapping platform-specific actions into two distinct domains defined in the system prompt. This unified action space enables consistent policy learning and parameter aggregation across platforms; details are provided in Appendix, Table 16.

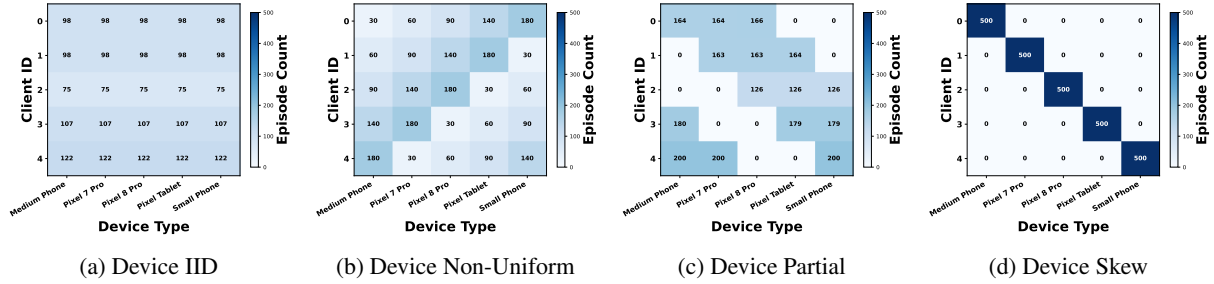


Figure 2: Distributions of the 5 different Android device across five clients in *FedGUI-Device*. Through 4 types of cross-device heterogeneity subsets, diverse patterns across clients are simulated.

**Data Processing.** We implement systematic data cleaning to remove episodes with missing images, redundant or semantically unnecessary actions, and ambiguous action parameters (e.g., undefined spatial ranges). To model realistic non-IID scenarios, we construct various federated partitioning configurations spanning platform-specific settings, varying client counts, and extreme distribution skew.

### 3.3 Data Description & Visualization

#### 3.3.1 FedGUI-Platform

**Cross-Platform Heterogeneity.** Data collected from different platforms, such as mobile devices and web browsers, inherently exhibit significant distributional shifts, as characterized by differences in visual attributes, application ecosystems, task structures, and action spaces. This phenomenon, termed cross-platform heterogeneity, poses a substantial challenge to federated learning, as clients contributing platform-specific data can degrade the performance and convergence of the global model. Therefore, a systematic study and quantification of this heterogeneity is of paramount importance.

**Description.** To enable a systematic investigation of cross-platform heterogeneity, we construct the *FedGUI-Platform* dataset. To systematically control the degree of heterogeneity, we propose three carefully curated data partitions: (1) In the Platform IID subset, each client holds a uniform mixture of data from all platforms. (2) A moderate degree of heterogeneity is introduced in Platform Partial, where each client is denied access to data from one particular platform. (3) In stark contrast, Platform Skew exclusively assign each client data from a single platform, thereby maximizing statistical heterogeneity across the network.

#### 3.3.2 FedGUI-Device

**Cross-Device Heterogeneity.** In real-world scenarios, participating clients provide data collected

from heterogeneous devices, including smartphones, laptops, and tablets. Even among clients operating on identical Android OS, substantial device-level variations persist, such as distinct visual patterns in image resolution, screen size, aspect ratio, and background rendering. To minimize the confounding effects introduced by platform and OS differences, we restrict our analysis to Android devices and define cross-device heterogeneity as the fine-grained variability that arises solely from device-specific characteristics.

**Description.** Leveraging the data collection methodology of GUI Odyssey (Lu et al., 2025), which employs multiple device simulators, we construct *FedGUI-Device* by labeling each episode with its device type to model cross-device heterogeneity. Specifically, we consider the following settings with increasing level of heterogeneity: (1) Device IID, where all clients observe data from all devices with identical distributions; (2) Device Non-Uniform, where all clients observe all devices but with heterogeneous proportions. (3) Device Partial, where each client observes data from a subset of devices and (4) Device Skew, where each client receives data from only one device.

As visualized in Figure 2, these variants induce progressively stronger cross-source heterogeneity while keeping platform and device factors fixed.

#### 3.3.3 FedGUI-OS

**Cross-OS Heterogeneity.** Even when underlying hardware and application functionalities are comparable, the diversity of client operating systems, such as macOS and Windows, introduces substantial discrepancies in system architecture, GUI design, window management, and interaction conventions. These differences lead to distinct visual layouts, system behaviors, and event-handling mechanisms, resulting in heterogeneous data distributions across clients. We define cross-OS hetero-

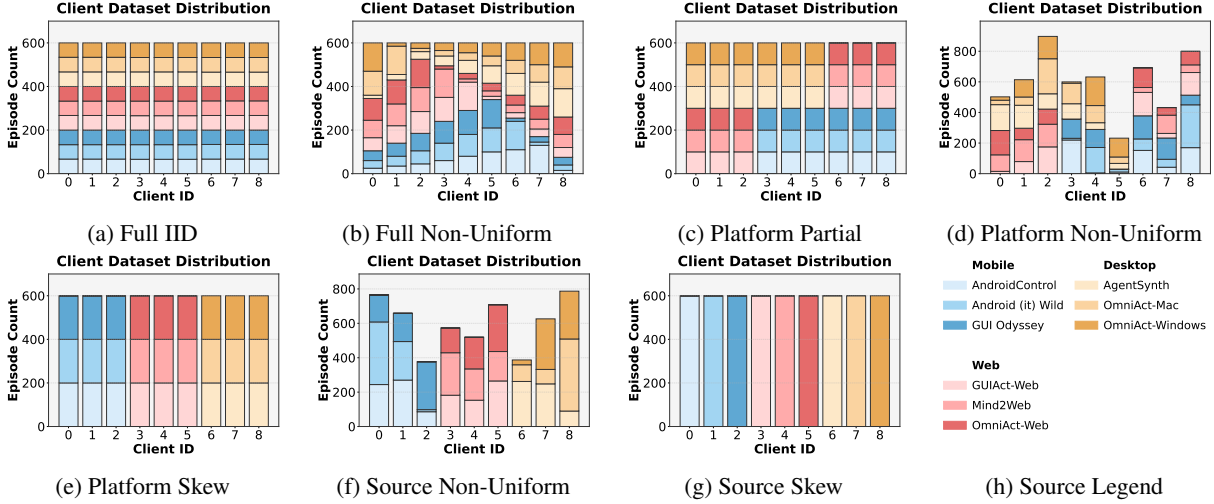


Figure 3: Distributions of episode counts across clients in *FedGUI-Full*. The 7 subsets demonstrate various degrees of cross-source and cross-platform heterogeneity, ranging from IID coverage to highly skewed distributions.

generality as the variability arising from differences among operating systems, independent of device-level or platform-level factors.

**Description.** We construct *FedGUI-OS* to facilitate the exploration of cross-OS heterogeneity. For Ubuntu, we collect data using AgentSynth, which is sourced from desktop simulations built upon OS-World (Xie et al., 2024). For macOS and Windows, we leverage data from OmniAct and manually categorize the samples according to their operating systems. The client participation protocols follow those of *FedGUI-Device*, including IID, non-uniform, skewed, and partial-observation settings. In addition, we introduce a fully random variant to provide an additional homogeneous evaluation.

### 3.3.4 FedGUI-Mobile & FedGUI-Web

**Cross-Source Heterogeneity.** Cross-source heterogeneity refers to the statistical discrepancies in data that arise from different acquisition methods and sources. Such heterogeneity can be pronounced even when data is collected on the same platform and operating system. For instance, the trajectory data from AC (collected via large-scale outsourcing) and GO (generated by automated synthesis on simulators) exhibit different distributions.

**Description.** To investigate cross-source heterogeneity at a granular level, we constructed two platform-specific datasets: *FedGUI-Mobile* and *FedGUI-Web*. These datasets are designed to highlight the nuanced discrepancies stemming from their distinct data acquisition methods. *FedGUI-Mobile* comprises three sources, AC, AitW, and

GO, originating from different task distributions and collection pipelines. In contrast, *FedGUI-Web* is built from OA-W, GA, and M2W, covering diverse web interaction tasks and website domains. Both datasets follow the same construction principle as *FedGUI-OS*, while focusing specifically on source-level differences in GUI environments.

### 3.3.5 FedGUI-Full

*FedGUI-Full* represents the most realistic and challenging setting in our benchmark. It constitutes all data sources and simultaneously incorporates cross-platform and cross-source heterogeneity, while further scaling the number of participating clients. To systematically characterize diverse degrees and forms of heterogeneity, we construct seven partition variants by jointly controlling source coverage, platform purity, and quantity balance across clients.

As visualized in Figure 3, we consider: (1) Full IID, where each client observes all sources and platforms with equal proportions; (2) Full Non-Uniform, where all sources are observed but with uneven data distributions; (3) Platform Partial, where each client observes only two platforms; (4) Platform Non-Uniform, which further introduces quantity imbalance across platforms; (5) Platform Skew, where each client is restricted to a single platform; (6) Source Non-Uniform, combining platform skew with source-level imbalance; and (7) Source Skew, where each client receives data from only one source and one platform. Together, these variants span a broad spectrum of realistic federated heterogeneity patterns.

| <b>Supported Base Models</b> , Details in Appendix D.3              |  |
|---|--|
| OpenAI:   | GPT-5, GPT-4o, GPT-4o-mini                             |
| Tongyi:   | Qwen3-VL-2B/4B/8B,<br>Qwen2.5-VL-3B/7B, Qwen2-VL-2B/7B |
| Google:   | Gemma-3-4B/12B, Gemini-Pro-2.5                         |
| Intern:   | InternVL2-1B/2B/4B/8B                                  |
| DeepSeek:   | DeepSeekVL2, DeepSeekVL2-tiny/small                    |
| Seed:   | UI-TARS-1.5-7B, Doubao-1.5-Vision-pro                  |
| <b>Integrated FL Algorithms</b> , Details in Appendix D.2           |  |
| FedAvg, FedProx, SCAFFOLD, FedAvgM,<br>FedAdam, FedYogi, FedAdagrad |  |

Table 4: Summary of supported base models and integrated federated learning algorithms in FedGUI.

### 3.4 Framework Description

FedGUI is built upon the prevalent *ms-swift* framework (Zhao et al., 2024), which we extend for federated learning by modularizing the core pipeline to seamlessly integrate federated components. This architecture establishes FedGUI as a unified and extensible benchmark for evaluating GUI agents across diverse platforms, supporting a wide array of VLMs and federated algorithms while enabling efficient experimentation.

**Model Support & Algorithm Integration.** As summarized in Table 4, FedGUI accommodates diverse GUI-capable foundation models including both proprietary and open-source VLMs, enabling evaluation across different model scales and architectural designs. In addition, we integrate commonly used federated optimization algorithms, ensuring fair and reproducible comparisons under identical training configurations.

## 4 Experiments

To address the key research questions, we first conduct experiments on *FedGUI-Full*, the most comprehensive and realistic dataset, in Section 4.2. Subsequently, Section 4.3 compares the performance of various FL algorithms, while Section 4.4 explores additional forms of heterogeneity across their corresponding datasets. We also analyze the impact of backbone models in Section 4.6 and assess system efficiency in Section 4.7. Supplementary results are detailed in Appendix B.

### 4.1 Basic Setups (Details in Appendix D)

**Base Model.** We adopt Qwen2-VL-7B (Wang et al., 2024) as the main base model in our experiments. To enable efficient adaptation under resource constraints on client-side devices, we apply Low-Rank Adaptation (LoRA) (Hu et al., 2021).

**Training Configuration.** All models are trained for 30 federated rounds. At each round, we uniformly sample 10% of the training data and randomly select 3 clients to participate, simulating realistic user availability and client dropout scenarios (Jiang et al., 2024).

**Evaluation Metrics.** We evaluate model performance at the action level using three metrics. *Type* measures whether the predicted action type matches the ground truth. *Grounding Accuracy (Ground)* evaluates spatial correctness for grounding actions (e.g., CLICK), which measures whether the predicted target GUI element matches the gold element; *Success Rate (SR)* requires both correct action type and correct content grounding, including spatial alignment for coordinate-based actions and semantic similarity for text-based actions.

### 4.2 Verification of Cross-Platform Federated Collaboration

**Setups.** To verify whether cross-platform collaboration improves federated learning performance, we first conduct experiments on *FedGUI-Full* across different distribution subsets using FedAvg (McMahan et al., 2017). The comparisons are organized along two dimensions: (1) comparison between collaboration on **homogeneous** and **heterogeneous** distributions with increasing heterogeneity levels; (2) comparison of collaboration among clients from **a single platform** versus clients from **all platforms (Source Skew)**.

**Results.** From Table 5, we draw the following key conclusions: (1) Single-domain models suffer from catastrophic failure on unseen platforms, but federated learning (even under *Source Skew*) significantly restores baseline utility. This underscores that federated collaboration is essential to bridge the domain isolation that renders localized GUI agents non-functional across platforms. (2) Overall performance degrades consistently as the data distribution shifts from IID coverage to increasingly heterogeneous settings, confirming that cross-platform heterogeneity fundamentally challenges federated GUI learning. (3) A clear platform sensitivity hierarchy emerges: desktop exhibits the highest sensitivity, followed by web, then mobile. This suggests asymmetric knowledge transfer, where mobile interactions transfer more effectively to other platforms than desktop does.

| Distribution         | AndroidControl |        |       | GUIAct-Web |        |       | AgentSynth |        |       | Average |        |       |
|----------------------|----------------|--------|-------|------------|--------|-------|------------|--------|-------|---------|--------|-------|
|                      | Type           | Ground | SR    | Type       | Ground | SR    | Type       | Ground | SR    | Type    | Ground | SR    |
| Central              | 48.41          | 45.47  | 59.89 | 68.48      | 51.48  | 48.73 | 70.54      | 76.74  | 54.91 | 71.28   | 55.17  | 49.32 |
| Only Mobile          | 72.53          | 29.14  | 40.21 | 46.74      | 18.79  | 23.91 | 43.98      | 19.55  | 8.58  | 54.42   | 22.49  | 24.23 |
| Only Web             | 54.78          | 4.49   | 5.16  | 66.12      | 45.09  | 43.66 | 48.96      | 22.08  | 11.89 | 56.62   | 23.89  | 20.24 |
| Only Desktop         | 53.57          | 6.46   | 10.17 | 34.78      | 7.89   | 2.90  | 72.06      | 74.77  | 54.63 | 53.47   | 29.71  | 22.57 |
| Full IID             | 67.83          | 22.79  | 30.86 | 64.13      | 44.79  | 43.84 | 69.85      | 55.77  | 40.53 | 67.27   | 41.12  | 38.39 |
| Full Non-Uniform     | 67.07          | 19.57  | 30.96 | 59.78      | 41.86  | 37.50 | 68.74      | 57.55  | 41.91 | 65.20   | 39.66  | 36.79 |
| Platform Partial     | 68.44          | 18.71  | 30.50 | 65.76      | 38.55  | 40.12 | 70.40      | 51.71  | 39.54 | 68.20   | 36.32  | 36.72 |
| Platform Non-Uniform | 70.11          | 21.20  | 30.02 | 66.49      | 39.87  | 39.64 | 69.02      | 56.85  | 40.17 | 68.54   | 39.31  | 36.61 |
| Platform Skew        | 70.86          | 18.44  | 33.84 | 60.87      | 39.75  | 40.58 | 64.87      | 37.18  | 26.14 | 65.53   | 31.79  | 33.52 |
| Source Non-Uniform   | 72.08          | 20.67  | 34.29 | 59.60      | 32.76  | 36.05 | 66.80      | 40.22  | 28.63 | 66.16   | 31.22  | 32.99 |
| Source Skew          | 61.91          | 6.91   | 16.24 | 56.52      | 42.08  | 34.24 | 69.29      | 58.07  | 42.05 | 62.57   | 35.69  | 30.84 |

Table 5: Results of FedAvg on *FedGUI-Full* under different data distributions. Only X denotes that the participating clients are restricted to the X domain in the *Source Skew* subset. Overall, the performance of the GUI agent degrades as the level of heterogeneity increases, while contributions from other platforms can help improve performance.

| Algorithm  | Platform IID   |            |            |         | Platform Partial |            |            |         | Platform Skew  |            |            |         |
|------------|----------------|------------|------------|---------|------------------|------------|------------|---------|----------------|------------|------------|---------|
|            | AndroidControl | GUIAct-Web | AgentSynth | Average | AndroidControl   | GUIAct-Web | AgentSynth | Average | AndroidControl | GUIAct-Web | AgentSynth | Average |
| SR         |                |            |            |         |                  |            |            |         |                |            |            |         |
| Central    | 48.10          | 53.26      | 60.72      | 54.03   | -                | -          | -          | -       | -              | -          | -          | -       |
| Local      | 27.77          | 35.51      | 28.63      | 30.64   | 28.07            | 33.15      | 24.62      | 28.61   | 33.84          | 10.51      | 5.53       | 16.63   |
| FedAvg     | 35.05          | 43.12      | 33.06      | 37.08   | 32.78            | 42.21      | 35.27      | 36.75   | 27.16          | 43.84      | 34.16      | 35.05   |
| FedProx    | 32.63          | 42.93      | 37.22      | 37.59   | 31.63            | 41.85      | 36.21      | 36.56   | 27.77          | 43.84      | 34.72      | 35.44   |
| SCAFFOLD   | 33.69          | 43.48      | 34.02      | 37.06   | 32.12            | 44.20      | 33.33      | 36.55   | 26.86          | 43.48      | 33.89      | 34.74   |
| FedYogi    | 35.81          | 44.38      | 52.28      | 44.16   | 35.81            | 43.12      | 48.13      | 42.35   | 34.75          | 48.55      | 50.90      | 44.73   |
| FedAdam    | 37.94          | 43.84      | 53.53      | 45.10   | 35.05            | 42.93      | 48.82      | 42.27   | 33.84          | 44.38      | 50.62      | 42.95   |
| FedAvgM    | 33.54          | 42.93      | 34.30      | 36.92   | 34.33            | 42.03      | 33.59      | 36.65   | 28.22          | 44.93      | 34.85      | 36.00   |
| FedAdagrad | 35.81          | 43.12      | 46.75      | 41.89   | 35.36            | 41.30      | 33.61      | 36.76   | 26.71          | 26.99      | 35.13      | 29.61   |

Table 6: Success Rate (SR, %) results on *FedGUI-Platform*. Each algorithm is evaluated under IID and non-IID settings. Results are shown for three benchmarks (AC: AndroidControl, GA-W: GUIAct-Web, AS: AgentSynth) and their average. Algorithms exhibit diverse performance patterns, while optimizer-based methods are more robust under skewed distributions.

### 4.3 Comparison of FL Algorithms

**Setups.** To examine their performance under diverse data distributions, we further compare different FL algorithms in three representative subsets on both *FedGUI-Platform* and *FedGUI-Full*, with results shown in Table 6 and Table 14, respectively. We also provide centralized and local baselines under identical configurations for reference.

We conclude the following findings: (1) Federated learning exhibits a strong salvage effect under cross-platform heterogeneity: while local models trained on a single platform fail to generalize to others, FL algorithms substantially recover cross-platform usability, elevating unusable models to a functional level. (2) Adaptive optimizers (e.g., FedAdam) achieve higher average performance under platform-skewed distributions, suggesting stronger robustness to platform-induced interference compared to correction-based methods (e.g., FedProx). (3) Despite these gains, a substantial gap to centralized training remains, indicating con-

siderable headroom for future work and exposing an inherent trade-off in federated settings between global generalization and platform-specific expertise under extreme heterogeneity.

### 4.4 Investigation of Cross-Device, Cross-OS, and Cross-Source Heterogeneity

**Setups.** Beyond cross-platform heterogeneity, we further investigate more fine-grained types of heterogeneity: cross-device, cross-OS, and cross-source variations on the corresponding datasets described in Section 3.3. For cross-device heterogeneity, we evaluate all algorithms across all five devices on *FedGUI-Device*, with complete results provided in Table 11. For cross-OS and cross-source experiments, we evaluate models on the corresponding source test sets. For example, on *FedGUI-OS* shown in Figure 4, we test on samples from AS (Ubuntu), OA-Mac (macOS), and OA-Win (Windows). All models are fine-tuned for 30 rounds under three data distribution settings: IID,

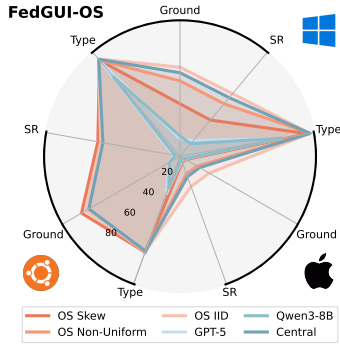
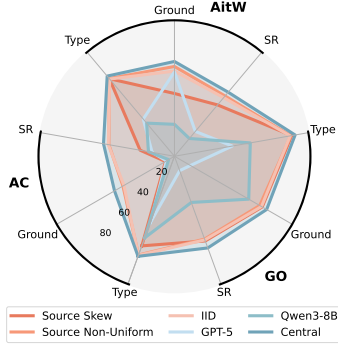
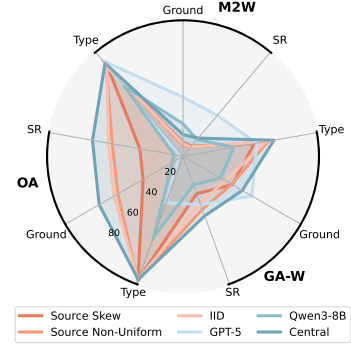


Figure 4: Results on *FedGUI-OS*, with evaluations conducted on Windows, Ubuntu, and macOS. Our Federated GUI agents demonstrate superior cross-OS generalization compared to strong VLMs like GPT-5.



(a) FedGUI-Mobile.



(b) FedGUI-Web.

Figure 5: Radar-chart comparison on *FedGUI-Mobile* and *FedGUI-Web*. Results for models trained with FedAvg on their respective source-specific subsets are plotted in red. The diverse data distributions in our dataset effectively reveal varied performance patterns, highlighting its capability to measure robustness against real-world distribution shifts.

| Setting     | Edu   | Fin   | Game  | Travel | Shop  | Ent   | Lit   | Ref   | Trans | Avg.         |
|-------------|-------|-------|-------|--------|-------|-------|-------|-------|-------|--------------|
| App IID     | 71.54 | 60.71 | 48.48 | 43.75  | 53.66 | 34.33 | 51.61 | 39.58 | 25.45 | <b>47.68</b> |
| App Partial | 70.59 | 62.12 | 43.55 | 41.26  | 53.12 | 32.89 | 50.98 | 39.12 | 23.56 | 46.35        |
| App Skew    | 69.23 | 64.28 | 42.42 | 40.00  | 52.65 | 34.33 | 50.00 | 38.52 | 21.81 | 45.92        |

Table 7: Cross-application heterogeneity results on web environments using Qwen2-VL-7B. Performances across homogeneous and heterogeneous distributions show the model’s robustness across diverse application domains.

non-uniform, and skewed.

**Results.** From Figure 4 and 5, we observe that: (1) The gap between skewed and IID settings (where the light-red curves diverge from and cover the dark-red curves) further evidences the presence of the identified heterogeneities in FedGUI. (2) Federated fine-tuning consistently outperforms strong VLM baselines (e.g., GPT-5), highlighting the value of aggregating in-domain GUI interaction knowledge beyond zero-shot generalization. (3) In some cases (e.g., *FedGUI-OS*), FedAvg achieve results close to centralized learning, demonstrating the practical value of federated GUI agents. (4) Cross-OS heterogeneity appears more critical than source or device heterogeneity, as OS Skew causes substantially larger performance drops compared to the modest variations observed for the same algorithm across device distributions in Table 11.

#### 4.5 Cross-Application Heterogeneity.

**Setups.** To explore how cross-application heterogeneity affects federated training performance, we conduct experiments on *FedGUI-Web*, covering 10 application categories ranging from education to transportation. Heterogeneity increases from balanced App IID, to subset-based App Partial, and finally single-app dominated App Skew.

**Results.** Key observations from Table 7 are as follows: (1) Although VLM agents can handle a wide variety of web tasks, performance is still affected by extreme skew in website applications. (2) Compared to the larger degradation under cross-platform heterogeneity (Table 5), cross-application shifts have a much smaller impact, suggesting that platform-level GUI differences are more challenging than application-level variations.

#### 4.6 Comparison of Model Backbones

**Setups.** To analyze the impact of model backbones on cross-platform federated GUI training, we conduct experiments on over 20 VLMs, covering both open-source and proprietary models, including Qwen, Gemma, and GPT (OpenAI, 2023). For open-source models, we evaluate both their base performance and federated fine-tuned counterparts on *FedGUI-Full* as a representative setting.

**Results.** As visualized in Figure 6 and 8: (1) The inconsistent performance of base VLMs across domains necessitates in-domain adaptation for executable GUI behavior. (2) After federated training, performance becomes more structured and shows a clearer positive correlation with model scale, particularly within the Qwen3 family (Bai et al., 2025). (3) Collaboration on distributed GUI data narrows

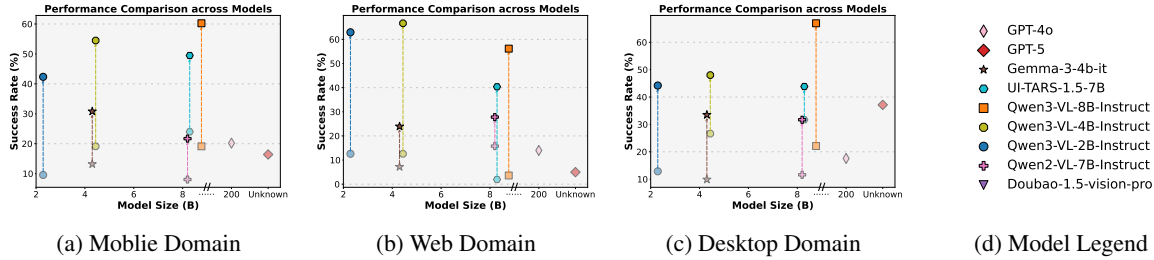


Figure 6: Model performance comparison across mobile, web, and desktop domains. Semi-transparent and solid markers denote base and FedAvg fine-tuned models, respectively, with dashed arrows indicating federated gains. Notably, after federated fine-tuning, several compact open-source models outperform large proprietary models.

| Approach                                 | Overhead                  |
|--|---------------------------|
| Central + Unpacked data (5k Epi.)        | $\approx 20$ GB           |
| Central + Unpacked data (All)            | $\approx 159$ GB          |
| FL + Full model (Qwen3-VL-8B)            | $16.50$ GB $\times$ round |
| FL + LoRA adapter (rank=8, $\alpha=32$ ) | $84.32$ MB $\times$ round |

Table 8: Communication overhead of FedGUI using different approaches. Transmitting only LoRA adapters achieves the highest communication efficiency.

| Base Model    | GPU Memory(MB) | Time per Round(s) |
|---------------|----------------|-------------------|
| SmoVLM-500M   | 6,056          | 04:47             |
| Intern2-VL-1B | 10,498         | 06:56             |
| Qwen3-VL-2B   | 14,536         | 10:24             |
| Gemma-4B      | 17,650         | 18:10             |
| Qwen2-VL-7B   | 21,610         | 18:24             |
| Qwen3-VL-8B   | 30,490         | 31:03             |

Table 9: Computational statistics for training VLMs using *FedGUI-Full*, showing its efficiency in both GPU memory usage and training time.

the gap between small open-source models and large proprietary models, enabling compact VLMs to achieve competitive GUI performance.

#### 4.7 Efficiency Evaluation

**Communication.** Assuming equal episode sizes, we compare three settings: centralized training, federated full fine-tuning, and federated training with LoRA. As shown in Table 8, LoRA-based FL achieves the lowest communication cost.

**Computation.** Table 9 reports the computational cost of *FedGUI-Full* across a diverse set of VLMs, spanning from lightweight to medium-scale models. We measure the peak GPU memory usage and the average wall-clock time per federated round. These results reinforce the training efficiency of FedGUI, which can be fully executed on a single RTX 4090 GPU. Notably, 2B models are compact enough to be deployed on standard mobile devices with 16GB of RAM.

| Model                | Distribution | 📱     | 🌐     | 💻     | 🏢     |
|----------------------|--------------|-------|-------|-------|-------|
| Qwen3-VL-2B-Instruct | Base         | 9.56  | 12.59 | 12.86 | 11.67 |
|                      | Full IID     | 42.34 | 62.93 | 44.20 | 49.82 |
|                      | Source Skew  | 32.92 | 58.32 | 43.32 | 44.85 |
| SmoVLM-500M          | Base         | 0.95  | 1.02  | 0.85  | 0.94  |
|                      | Full IID     | 15.15 | 22.28 | 20.60 | 19.34 |
|                      | Source Skew  | 5.31  | 7.42  | 20.47 | 11.07 |

Table 10: Success Rate (%) results using mobile-trainable models show that even models deployable on-device remain effective under FedGUI, demonstrating the feasibility of practical edge-side agents.

**Deployability on Edge Devices.** To better reflect realistic mobile and web scenarios, we evaluate the performance of sub-2B models that are deployable on high-end smartphones or edge devices. We specifically test Qwen3-VL-2B-Instruct and SmoVLM-500M-Instruct (Marafioti et al., 2025). As shown in Table 10, federated training consistently improves sub-2B models across all distributions, making them hardware-feasible for real-world mobile and edge deployment.

## 5 Conclusion

In this paper, we introduce **FedGUI**, the first unified benchmark for cross-platform federated GUI agents that systematically models heterogeneity across mobile, web, and desktop environments. Through six carefully curated datasets, FedGUI exposes realistic non-IID challenges largely overlooked by existing benchmarks. Extensive experiments across 7 algorithms and 20+ models show that such heterogeneity fundamentally degrades performance, while cross-platform collaboration exhibits a strong salvage effect even under severe isolation.

Overall, FedGUI bridges federated learning research and practical GUI agent evaluation, providing a solid foundation for robust and privacy-preserving agent development.

## Limitations

Despite its comprehensive design, coverage of cross-platform heterogeneity, and the informative insights provided, FedGUI still faces one major challenge. Our benchmark relies on publicly available GUI datasets rather than real user-owned private data collected in live federated deployments. While real user data would better reflect authentic interaction patterns, privacy, ethical, and reproducibility constraints make such data difficult to obtain and release at scale. Therefore, to better benefit the research community and facilitate more accessible experimentation, we incorporate diverse publicly available data sources, which, by virtue of their scale and heterogeneity, can reasonably simulate real-world distributed user data.

## References

- Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*.
- Omri Berkovitch, Sapir Caduri, Noam Kahlon, Anatoly Efros, Avi Caciularu, and Ido Dagan. 2025. Identifying user goals from ui trajectories. In *Companion Proceedings of the ACM on Web Conference 2025*, pages 2381–2390.
- Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H. Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. 2016. [Practical secure aggregation for federated learning on user-held data](#). *Preprint*, arXiv:1611.04482.
- Yuxiang Chai, Siyuan Huang, Yazhe Niu, Han Xiao, Liang Liu, Dingyu Zhang, Peng Gao, Shuai Ren, and Hongsheng Li. 2024. [AMEX: Android Multi-annotation Expo Dataset for Mobile GUI Agents](#). *Preprint*, arXiv:2407.17490.
- Wentong Chen, Junbo Cui, Jinyi Hu, Yujia Qin, Junjie Fang, Yue Zhao, Chongyi Wang, Jun Liu, Guirong Chen, Yupeng Huo, and 1 others. 2025. Guicourse: From general vision language model to versatile gui agent. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 21936–21959.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, and 1 others. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36:28091–28114.
- Cynthia Dwork. 2006. Differential privacy. In *International colloquium on automata, languages, and programming*, pages 1–12. Springer.
- Feng He, Tianqing Zhu, Dayong Ye, Bo Liu, Wanlei Zhou, and Philip S. Yu. 2024. [The Emerged Security and Privacy of LLM Agent: A Survey with Case Studies](#). *Preprint*, arXiv:2407.19354.
- Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. 2019. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2021. Lora: Low-rank adaptation of large language models. In *ICLR*.
- Zhifeng Jiang, Wei Wang, and Ruichuan Chen. 2024. Dordis: Efficient federated learning with dropout-resilient differential privacy. In *Proceedings of the Nineteenth European Conference on Computer Systems*, pages 472–488.
- Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem AlShikh, and Ruslan Salakhutdinov. 2024. Omniaact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web. In *European Conference on Computer Vision*, pages 161–178. Springer.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. 2020. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR.
- Quyu Kong, Xu Zhang, Zhenyu Yang, Nolan Gao, Chen Liu, Panrong Tong, Chenglin Cai, Hanzhang Zhou, Jianan Zhang, Liangyu Chen, Zhidan Liu, Steven Hoi, and Yue Wang. 2025. [Mobileworld: Benchmarking autonomous mobile agents in agent-user interactive, and mcp-augmented environments](#). *Preprint*, arXiv:2512.19432.
- Weirui Kuang, Bingchen Qian, Zitao Li, Daoyuan Chen, Dawei Gao, Xuchen Pan, Yuexiang Xie, Yaliang Li, Bolin Ding, and Jingren Zhou. 2023. Federatedscope-llm: A comprehensive package for fine-tuning large language models in federated learning. *arXiv preprint arXiv:2309.00363*.
- Sunjae Lee, Junyoung Choi, Jungjae Lee, Munim Hasan Wasi, Hojun Choi, Steve Ko, Sangeun Oh, and Insik Shin. 2024. [Mobilegpt: Augmenting llm with human-like app memory for mobile task automation](#). In *Proceedings of the 30th Annual International Conference*

- on *Mobile Computing and Networking*, ACM MobiCom '24, page 1119–1133, New York, NY, USA. Association for Computing Machinery.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 2020. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450.
- Wei Li, William E Bishop, Alice Li, Christopher Rawles, Folawiyo Campbell-Ajala, Divya Tyamagundlu, and Oriana Riva. 2024. On the effects of data scale on ui control agents. *Advances in Neural Information Processing Systems*, 37:92130–92154.
- Shuquan Lian, Yuhang Wu, Jia Ma, Yifan Ding, Zihan Song, Bingqi Chen, Xiawu Zheng, Hui Li, and Rongrong Ji. 2026. [Ui-agile: Advancing gui agents with effective reinforcement learning and precise inference-time grounding](#). *Preprint*, arXiv:2507.22025.
- Guangyi Liu, Pengxiang Zhao, Yaozhen Liang, Qinyi Luo, Shunye Tang, Yuxiang Chai, Weifeng Lin, Han Xiao, WenHao Wang, Siheng Chen, Zhengxi Lu, Gao Wu, Hao Wang, Liang Liu, and Yong Liu. 2026. [Memgui-bench: Benchmarking memory of mobile gui agents in dynamic environments](#). *Preprint*, arXiv:2602.06075.
- Guangyi Liu, Pengxiang Zhao, Liang Liu, Zhiming Chen, Yuxiang Chai, Shuai Ren, Hao Wang, Shibo He, and Wenchao Meng. 2025a. Learnact: Few-shot mobile gui agent with a unified demonstration benchmark. *arXiv preprint arXiv:2504.13805*.
- Guangyi Liu, Pengxiang Zhao, Liang Liu, Yaxuan Guo, Han Xiao, Weifeng Lin, Yuxiang Chai, Yue Han, Shuai Ren, Hao Wang, Xiaoyu Liang, Wenhao Wang, Tianze Wu, Linghao Li, Hao Wang, Guanqing Xiong, Yong Liu, and Hongsheng Li. 2025b. [Llm-powered gui agents in phone automation: Surveying progress and prospects](#). *Preprint*, arXiv:2504.19838.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024. [Deepseek-vl: Towards real-world vision-language understanding](#). *Preprint*, arXiv:2403.05525.
- Quanfeng Lu, Wenqi Shao, Zitao Liu, Lingxiao Du, Fangqing Meng, Boxuan Li, Botong Chen, Siyuan Huang, Kaipeng Zhang, and Ping Luo. 2025. [Guidyssey: A comprehensive dataset for cross-app gui navigation on mobile devices](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22404–22414.
- Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril Zakkas, Loubna Ben Allal, Anton Lozhkov and Noumane Tazi, Vaibhav Srivastav, Joshua Lochner, Hugo Larcher, Mathieu Morlon, Lewis Tunstall, Leandro von Werra, and Thomas Wolf. 2025. [Smolvlm: Redefining small and efficient multimodal models](#). *arXiv preprint arXiv:2504.05299*.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR.
- Liangbo Ning, Ziran Liang, Zhuohang Jiang, Haohao Qu, Yujuan Ding, Wenqi Fan, Xiao-yong Wei, Shanru Lin, Hui Liu, Philip S Yu, and 1 others. 2025. [A survey of webagents: Towards next-generation ai agents for web automation with large foundation models](#). In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pages 6140–6150.
- OpenAI. 2023. [Gpt-4: A large-scale multimodal model](#). *arXiv preprint arXiv:2303.08774*.
- Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, and 1 others. 2025. [Uitars: Pioneering automated gui interaction with native agents](#). *arXiv preprint arXiv:2501.12326*.
- Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. 2023. [Android in the Wild: A Large-Scale Dataset for Android Device Control](#). *Preprint*, arXiv:2307.10088.
- Sashank J Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. 2020. Adaptive federated optimization. In *International Conference on Learning Representations*.
- Qiushi Sun, Kanzhi Cheng, Zichen Ding, Chuanyang Jin, Yan Wang, Fangzhi Xu, Zhenyu Wu, Chengyou Jia, Liheng Chen, Zhoumianze Liu, Ben Kao, Guohao Li, Junxian He, Yu Qiao, and Zhiyong Wu. 2024. [OS-Genesis: Automating GUI Agent Trajectory Construction via Reverse Task Synthesis](#). *Preprint*, arXiv:2412.19723.
- Gemma Team. 2025. [Gemma 3](#).
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. [Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution](#). *arXiv preprint arXiv:2409.12191*.
- Wenhao Wang, Peizhi Niu, Zhao Xu, Zhaoyu Chen, Jian Du, Yaxin Du, Xianghe Pang, Keduan Huang, Yanfeng Wang, Qiang Yan, and Siheng Chen. 2025a. [Mcp-flow: Facilitating llm agents to master real-world, diverse and scaling mcp tools](#). *Preprint*, arXiv:2510.24284.

- Wenhao Wang, Zijie Yu, William Liu, Rui Ye, Tian Jin, Siheng Chen, and Yanfeng Wang. 2025b. [Fedmobileagent: Training mobile agents using decentralized self-sourced data from diverse users](#). *Preprint*, arXiv:2502.02982.
- WenHao Wang, Zijie Yu, Rui Ye, Jianqing Zhang, Guangyi Liu, Liang Liu, Siheng Chen, and Yanfeng Wang. 2025c. [FedMABench: Benchmarking mobile GUI agents on decentralized heterogeneous user data](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 26398–26419, Suzhou, China. Association for Computational Linguistics.
- Xuehui Wang, Zhenyu Wu, JingJing Xie, Zichen Ding, Bowen Yang, Zehao Li, Zhaoyang Liu, Qingyun Li, Xuan Dong, Zhe Chen, Weiyun Wang, Xiangyu Zhao, Jixuan Chen, Haodong Duan, Tianbao Xie, Shiqian Su, Chenyu Yang, Yue Yu, Yuan Huang, and 8 others. 2025d. [Mmbench-gui: Hierarchical multi-platform evaluation framework for gui agents](#). *arXiv preprint arXiv:2507.19478*.
- Jie Wen, Zhixia Zhang, Yang Lan, Zhihua Cui, Jianghui Cai, and Wensheng Zhang. 2023. A survey on federated learning: challenges and applications. *International journal of machine learning and cybernetics*, 14(2):513–535.
- Jingxu Xie, Dylan Xu, Xuandong Zhao, and Dawn Song. 2025. [Agentsynth: Scalable task generation for generalist computer-use agents](#). *arXiv preprint arXiv:2506.14205*.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. 2024. [Os-world: Benchmarking multimodal agents for open-ended tasks in real computer environments](#). *Preprint*, arXiv:2404.07972.
- Jiabo Ye, Xi Zhang, Haiyang Xu, Haowei Liu, Junyang Wang, Zhaoqing Zhu, Ziwei Zheng, Feiyu Gao, Junjie Cao, Zhengxi Lu, and 1 others. 2025. [Mobile-agent-v3: Fundamental agents for gui automation](#). *arXiv preprint arXiv:2508.15144*.
- Xinyu Zhang, Huiyu Xu, Zhongjie Ba, Zhibo Wang, Yuan Hong, Jian Liu, Zhan Qin, and Kui Ren. 2024. [PrivacyAsst: Safeguarding User Privacy in Tool-Using Large Language Model Agents](#). *IEEE Transactions on Dependable and Secure Computing*, pages 1–16.
- Yuzhe Zhang, Feiran Liu, Yi Shan, Xinyi Huang, Xin Yang, Yueqi Zhu, Xuxin Cheng, Cao Liu, Ke Zeng, Terry Jingchen Zhang, and Wenyuan Jiang. 2026a. [Silo-bench: A scalable environment for evaluating distributed coordination in multi-agent llm systems](#). *Preprint*, arXiv:2603.01045.
- Yuzhe Zhang, Xianwei Xue, Xingyong Wu, Mengke Chen, Chen Liu, Xinran He, Run Shao, Feiran Liu, Huanmin Xu, Qiutong Pan, and Haiwei Wang. 2026b. [Don't act blindly: Robust gui automation via action-effect verification and self-correction](#). *Preprint*, arXiv:2604.05477.
- Pengxiang Zhao, Guangyi Liu, Yaozhen Liang, Weiqing He, Zhengxi Lu, Yuehao Huang, Yaxuan Guo, Kexin Zhang, Hao Wang, Liang Liu, and Yong Liu. 2025. [Mas-bench: A unified benchmark for shortcut-augmented hybrid mobile gui agents](#). *Preprint*, arXiv:2509.06477.
- Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang, Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. 2024. [Swift:a scalable lightweight infrastructure for fine-tuning](#). *Preprint*, arXiv:2408.05517.
- Hanzhang Zhou, Xu Zhang, Panrong Tong, Jianan Zhang, Liangyu Chen, Quyu Kong, Chenglin Cai, Chen Liu, Yue Wang, Jingren Zhou, and Steven Hoi. 2025. [Mai-ui technical report: Real-world centric foundation gui agents](#). *Preprint*, arXiv:2512.22047.

## A Future Directions

### A.1 Cross-Platform Heterogeneity-Aware Optimization

While our experiments demonstrate that cross-platform collaboration can improve overall performance, the substantial heterogeneity across mobile, web, and desktop environments remains a significant challenge for federated GUI agents. Future work should explore specialized aggregation mechanisms that can dynamically weight client contributions based on platform-specific characteristics and task relevance. For instance, developing platform-aware federated optimization algorithms that can identify and leverage complementary knowledge across different GUI modalities (e.g., touch-based interactions on mobile vs. keyboard-mouse interactions on desktop) could substantially improve convergence speed and final performance.

Additionally, investigating hierarchical federated learning architectures that first aggregate updates within the same platform before cross-platform aggregation may help mitigate the negative impact of extreme heterogeneity. The promising method based on our empirical results are design optimizer based algorithms based on FedYogi (Reddi et al., 2020).

### A.2 Real-World Data Collection and Continuous Learning

The current benchmark relies on curated datasets from existing sources, which may not fully capture the complexity and diversity of real-world user

| Algorithm      | P7P                   | P8P          | Sm.          | Med.         | Tab          | Avg.         | Algorithm      | P7P                       | P8P          | Sm.          | Med.         | Tab          | Avg.         |
|----------------|-----------------------|--------------|--------------|--------------|--------------|--------------|----------------|---------------------------|--------------|--------------|--------------|--------------|--------------|
| Central        | 78.95                 | 76.74        | 77.71        | 78.81        | 63.64        | 75.17        | Central        | 78.95                     | 76.74        | 77.71        | 78.81        | 63.64        | 75.17        |
| <b>Homo.</b>   | <b>Device IID</b>     |              |              |              |              |              | <b>Hetero.</b> | <b>Device Non-Uniform</b> |              |              |              |              |              |
| Local          | 56.15                 | 60.91        | 54.29        | 57.25        | 46.46        | 55.01        | Local          | 53.85                     | 63.64        | 45.14        | 61.83        | 49.49        | 54.79        |
| FedAvg         | 66.15                 | 70.00        | 63.71        | 70.99        | 60.61        | 66.29        | FedAvg         | 64.62                     | 67.27        | 60.76        | 64.12        | 54.55        | 62.26        |
| FedProx        | 67.69                 | 68.18        | 62.86        | 69.47        | 61.62        | 65.96        | FedProx        | 63.08                     | 66.36        | 62.29        | 64.12        | 57.58        | 62.69        |
| FedAvgM        | 67.69                 | 70.91        | 61.71        | 67.18        | 60.61        | 65.62        | FedAvgM        | 66.15                     | 66.36        | 58.86        | 64.89        | 56.57        | 62.57        |
| SCAFFOLD       | 66.15                 | 69.09        | 61.14        | 70.23        | 56.57        | 64.64        | SCAFFOLD       | 64.62                     | 67.27        | 61.71        | 66.41        | 54.55        | 62.91        |
| FedAdagrad     | 66.15                 | 73.64        | 65.71        | 72.52        | 64.65        | 68.53        | FedAdagrad     | 68.46                     | 71.82        | 68.57        | 68.70        | 59.60        | 67.43        |
| FedYogi        | <b>71.54</b>          | 75.45        | <b>76.57</b> | <b>75.57</b> | <b>68.69</b> | <b>73.56</b> | FedYogi        | 70.77                     | <b>76.36</b> | <b>70.86</b> | 74.05        | 61.62        | 70.73        |
| FedAdam        | 66.92                 | <b>77.27</b> | 70.86        | 71.76        | 60.61        | 69.48        | FedAdam        | <b>73.08</b>              | 72.73        | 69.71        | <b>75.57</b> | <b>63.64</b> | <b>70.95</b> |
| <b>Hetero.</b> | <b>Device Partial</b> |              |              |              |              |              | <b>Hetero.</b> | <b>Device Skew</b>        |              |              |              |              |              |
| Local          | 58.46                 | 60.00        | 45.14        | 58.02        | 45.45        | 53.41        | Local          | 59.23                     | 67.27        | 46.29        | 51.91        | 42.42        | 53.42        |
| FedAvg         | 59.88                 | 69.09        | 61.14        | 62.60        | 56.57        | 61.86        | FedAvg         | 52.38                     | 66.36        | 55.43        | 60.31        | 54.55        | 57.81        |
| FedProx        | 64.62                 | 68.18        | 61.71        | 64.89        | 52.53        | 62.39        | FedProx        | 68.46                     | 68.18        | 57.14        | 64.12        | 52.53        | 62.09        |
| FedAvgM        | 64.62                 | 66.36        | 64.57        | 64.12        | 52.53        | 62.44        | FedAvgM        | 66.92                     | 68.18        | 57.14        | 63.36        | 55.56        | 62.23        |
| SCAFFOLD       | 66.15                 | 70.91        | 61.71        | 67.18        | 54.55        | 64.10        | SCAFFOLD       | 68.46                     | 67.27        | 55.43        | 62.60        | 54.55        | 61.66        |
| FedAdagrad     | 67.69                 | <b>71.82</b> | 70.29        | 65.65        | 57.58        | 66.61        | FedAdagrad     | 66.15                     | 71.82        | 61.14        | 66.41        | 60.61        | 65.23        |
| FedYogi        | 74.62                 | 70.91        | 73.14        | <b>74.81</b> | 60.61        | 70.82        | FedYogi        | 73.08                     | <b>79.09</b> | <b>70.29</b> | <b>71.76</b> | <b>66.67</b> | <b>72.18</b> |
| FedAdam        | <b>75.38</b>          | <b>71.82</b> | <b>73.71</b> | 72.52        | <b>64.65</b> | <b>71.62</b> | FedAdam        | <b>74.62</b>              | 78.18        | 66.86        | <b>71.76</b> | 65.66        | 71.42        |

Table 11: Experimental results on *FedGUI-Device*. We evaluate four types of cross-device distributions across five devices: Pixel 7 Pro (P7P), Pixel 8 Pro (P8P), Small Phone (Sm.), Medium Phone (Med.), Pixel Tablet (Tab) and their average. Accuracy is generally higher in **Device IID** settings, while heterogeneous scenarios (**Device Non-Uniform**, **Device Partial**, and **Device Skew**) demonstrate the varying robustness of FL algorithms to cross-device variability.

interactions across different geographical regions, application versions, and usage patterns. Future research should investigate scalable frameworks for continuous data collection from real users in production environments (Kong et al., 2025; Zhang et al., 2026a), while addressing the inherent challenges of data quality control, annotation efficiency, and temporal distribution shifts. Specifically, developing semi-supervised or self-supervised learning techniques that can leverage abundant unlabeled GUI interaction traces would significantly reduce the annotation burden.

### A.3 Privacy-Preserving Mechanisms for Federated GUI Agents

Although federated learning provides a foundation for privacy preservation by keeping raw data local, GUI interaction traces often contain highly sensitive information including personal identities, financial transactions, health records, and behavioral patterns that could be exposed through model updates or inferred from agent actions. Future work should integrate advanced privacy-enhancing technologies specifically designed for GUI agent training. Differential privacy (DP) (Dwork, 2006) mechanisms and secure aggregation (Bonawitz et al., 2016) could be adapted to account for the sequential and high-dimensional nature of GUI interaction data, ensuring formal privacy guarantees without

severely degrading model utility.

## B Supplementary Experiments & Results

### B.1 Client Scalability under Heterogeneity.

**Setups.** We conduct ablation studies on *FedGUI-Full* to analyze the effect of scaling client number under different heterogeneity settings. By varying the number of participating clients (9, 18, 27, and 36) and applying multiple heterogeneous partitions (Full IID, Platform Skew, and Source Skew), we simulate federated training scenarios with different degrees of fragmentation and client participation.

**Results.** From the ablation results in Table 12, we draw the following conclusions: (1) As the client number scales from 9 to 36, performance consistently declines across most settings, indicating that finer partitioning leads to data sparsity and hampers effective GUI feature learning. (2) Regardless of client scales, Skew settings consistently underperform the IID setting, confirming that data heterogeneity remains the primary obstacle that cannot be mitigated by changing participation scale. (3) Increasing the number of participating clients does not compensate for performance loss under heterogeneity, suggesting that scaling alone is insufficient to address the challenges of non-IID data in federated GUI learning.

| Client Count | Full IID |       |       |       | Platform Skew |       |       |       | Source Skew |       |       |       |
|--------------|----------|-------|-------|-------|---------------|-------|-------|-------|-------------|-------|-------|-------|
|              | 📱        | 🌐     | 💻     | 📊     | 📱             | 🌐     | 💻     | 📊     | 📱           | 🌐     | 💻     | 📊     |
| Client 9     | 30.86    | 43.84 | 40.53 | 38.39 | 16.24         | 34.24 | 42.05 | 30.84 | 33.84       | 40.58 | 26.14 | 33.52 |
| Client 18    | 21.40    | 32.79 | 32.79 | 28.99 | 22.46         | 31.16 | 31.16 | 28.26 | 11.53       | 15.58 | 15.58 | 14.23 |
| Client 27    | 22.61    | 32.07 | 32.07 | 28.92 | 23.67         | 33.51 | 33.51 | 30.23 | 14.11       | 20.65 | 20.65 | 18.47 |
| Client 36    | 18.97    | 29.71 | 29.71 | 26.13 | 14.87         | 24.82 | 24.82 | 21.50 | 13.35       | 24.82 | 24.82 | 21.00 |

Table 12: Client scalability evaluation across different client numbers (9, 18, 27, 36). Each setting is evaluated under Full IID, Platform Skew and Source Skew distributions. Results are shown for three benchmarks (AndroidControl: 📱, GUIAct-Web: 🌐, AgentSynth: 💻) and their average (📊). Model performance degrades with a larger number of clients, owing to the escalating difficulty of collaboration.

| Sample Count | Full IID |       |       |       | Source Skew |       |       |       |
|--------------|----------|-------|-------|-------|-------------|-------|-------|-------|
|              | 📱        | 🌐     | 💻     | 📊     | 📱           | 🌐     | 💻     | 📊     |
| Sample 3     | 22.61    | 28.62 | 32.07 | 27.77 | 18.61       | 26.62 | 26.32 | 23.85 |
| Sample 6     | 22.91    | 28.72 | 32.68 | 28.10 | 19.52       | 27.27 | 28.65 | 25.15 |
| Sample 9     | 23.51    | 28.98 | 32.32 | 28.27 | 21.51       | 28.98 | 32.32 | 27.60 |
| Sample 12    | 23.32    | 29.34 | 33.38 | 28.68 | 22.22       | 28.34 | 33.68 | 28.08 |

Table 13: Client sampling scalability evaluation. Success rate (%) across different number of sampled clients (3, 6, 9, 12) under Full IID and Source Skew distributions. Results cover Mobile, Web, and OS platforms with their average performance.

## B.2 Supplementary Comparison of FL Algorithms

**Setups.** We further extend our evaluation to *FedGUI-Mobile*, and *FedGUI-Web* to make a more comprehensive view of result concluded in Section 4.3. For the mobile domain, we use AC, AitW, and GO each 2000 episodes. For the web domain, we evaluate on M2W (500 episodes), GA-W (500 episodes), and OA-W (1,000 episodes). In both settings, we simulate a federated environment with 15 clients and a participation rate of 3 clients per round, enabling a consistent comparison of federated algorithms under diverse UI distributions.

**Results.** Key observations from Tables 15 include: (1) Decentralized training is more feasible in the mobile domain, where IID performance reaches within 10–15% of the centralized baseline, whereas the web domain exhibits a larger gap, particularly on M2W, due to higher variability in layouts, dynamic content, and interaction flows; (2) data heterogeneity significantly degrades performance in both domains, with the web setting suffering nearly 50% average degradation relative to IID, indicating that source skew poses a more severe challenge for cross-website GUI tasks than for mobile ones; and (3) adaptive methods such as FedAdam, FedAdagrad, and FedYogi consistently outperform FedAvg across both domains, demonstrating that per-parameter adaptation and gradient

scaling are crucial for handling high-variance GUI interaction sequences.

## C Dataset Details

### C.1 Data Composition & Example

To illustrate the structure of our dataset and the composition of a data episode, we present representative examples from each domain in this section.

As shown in Figure 9, each episode consists of three components, following typical GUI agent research (Chai et al., 2024; Lee et al., 2024; Liu et al., 2025a, 2026; Zhao et al., 2025): (1) an instruction, i.e., a natural-language description of the task to be completed; (2) a sequence of screenshots captured from the beginning to the end of the task; and (3) a corresponding sequence of actions, aligned one-to-one with the screenshots, specifying the user interactions that lead to the next state. Notably, an episode may involve cross-application or cross-website interactions, reflecting realistic multi-step workflows in GUI environments. All actions are defined within a unified action space of 17 action types, including both basic interactions and domain-specific custom actions, as detailed in Table 16.

### C.2 Data Collection & Processing Details

To ensure the quality and consistency of the federated GUI datasets across heterogeneous platforms, we implement a rigorous data processing pipeline

| Algorithm  | Full IID |       |       |       | Platform Skew |       |       |       | Source Skew |       |       |       |
|------------|----------|-------|-------|-------|---------------|-------|-------|-------|-------------|-------|-------|-------|
|            | 📱        | 🌐     | 💻     | 📄     | 📱             | 🌐     | 💻     | 📄     | 📱           | 🌐     | 💻     | 📄     |
| Central    | 44.31    | 48.73 | 54.91 | 49.32 | -             | -     | -     | -     | -           | -     | -     | -     |
| Local      | 31.11    | 67.22 | 33.20 | 33.33 | 37.10         | 45.15 | 5.97  | 23.96 | 30.81       | 26.46 | 5.29  | 19.94 |
| FedAvg     | 30.80    | 69.85 | 40.53 | 38.39 | 16.24         | 69.29 | 42.05 | 30.84 | 33.84       | 64.87 | 26.14 | 33.52 |
| FedProx    | 31.26    | 70.12 | 41.63 | 38.97 | 35.66         | 64.73 | 26.69 | 34.37 | 16.24       | 68.74 | 41.22 | 29.48 |
| SCAFFOLD   | 30.50    | 69.85 | 41.49 | 38.43 | 34.90         | 64.87 | 26.28 | 33.86 | 16.39       | 69.85 | 41.36 | 30.06 |
| FedYogi    | 37.78    | 73.31 | 58.64 | 48.14 | 43.40         | 69.85 | 39.28 | 43.62 | 21.24       | 71.09 | 54.63 | 39.00 |
| FedAdam    | 39.45    | 72.61 | 56.57 | 47.95 | 42.49         | 69.85 | 39.83 | 43.20 | 21.55       | 70.82 | 55.19 | 40.19 |
| FedAvgM    | 31.26    | 69.43 | 40.66 | 38.41 | 35.81         | 63.76 | 25.73 | 34.46 | 16.54       | 69.02 | 40.25 | 29.32 |
| FedAdagrad | 32.47    | 68.46 | 48.82 | 42.01 | 39.15         | 67.50 | 33.33 | 39.86 | 21.24       | 69.85 | 54.36 | 36.49 |

Table 14: Performance comparison of FL algorithms on *FedGUI-Full*. The table reports Success Rate (SR, %), as it is the most informative metric among the three evaluated. These results, in conjunction with those in Table 6, reaffirm the significant performance variation of FL algorithms under different data distributions.

| Algorithm  | 📱 Mobile Evaluation |       |       |       | Algorithm  | 🌐 Web Evaluation |       |       |       |
|------------|---------------------|-------|-------|-------|------------|------------------|-------|-------|-------|
|            | AC                  | AitW  | GO    | Avg.  |            | M2W              | GA-W  | OA-W  | Avg.  |
| Central    | 55.08               | 58.02 | 69.78 | 60.96 | Central    | 17.73            | 52.72 | 67.74 | 46.06 |
| FedAvg     | 42.19               | 52.19 | 55.29 | 49.89 | FedAvg     | 5.18             | 32.79 | 53.77 | 30.58 |
| FedProx    | 42.49               | 51.56 | 54.53 | 49.53 | FedProx    | 4.58             | 33.33 | 54.53 | 30.81 |
| FedYogi    | 48.41               | 53.54 | 64.35 | 55.43 | FedYogi    | 3.59             | 39.13 | 54.91 | 32.54 |
| FedAdam    | 48.25               | 52.19 | 61.54 | 53.99 | FedAdam    | 3.98             | 39.67 | 62.08 | 35.24 |
| SCAFFOLD   | 41.88               | 52.19 | 55.56 | 49.88 | SCAFFOLD   | 4.38             | 33.88 | 54.15 | 30.80 |
| FedAvgM    | 43.25               | 52.08 | 55.08 | 50.14 | FedAvgM    | 5.18             | 34.06 | 54.15 | 31.13 |
| FedAdagrad | 46.89               | 47.92 | 59.13 | 51.31 | FedAdagrad | 5.38             | 35.33 | 60.19 | 33.63 |
| FedAvg     | 23.98               | 47.60 | 60.10 | 43.89 | FedAvg     | 6.77             | 13.04 | 29.81 | 16.54 |
| FedProx    | 24.43               | 46.98 | 60.30 | 43.90 | FedProx    | 6.18             | 12.50 | 28.68 | 15.79 |
| FedYogi    | 29.59               | 50.83 | 65.32 | 48.58 | FedYogi    | 5.58             | 14.31 | 30.75 | 16.88 |
| FedAdam    | 30.35               | 49.48 | 65.52 | 48.45 | FedAdam    | 5.38             | 11.59 | 40.38 | 19.12 |
| SCAFFOLD   | 24.43               | 61.30 | 59.89 | 48.54 | SCAFFOLD   | 5.38             | 12.68 | 30.57 | 16.21 |
| FedAvgM    | 24.58               | 47.08 | 59.75 | 43.80 | FedAvgM    | 5.58             | 12.86 | 29.06 | 15.83 |
| FedAdagrad | 27.92               | 50.17 | 70.77 | 49.62 | FedAdagrad | 7.57             | 17.21 | 31.51 | 18.76 |

Table 15: Supplementary experiments results on *FedGUI-Mobile* (left) and *FedGUI-Web* (right). Performance of each algorithm, measured as Success Rate (SR, %), is evaluated under IID and Skewed data distributions and compared against the Centralized baseline. The table presents results on individual datasets within the mobile and web domains, along with their respective averages. W

as follows: We first perform systematic data cleaning to remove low-quality episodes, including samples with missing visual observations, ambiguous or undefined action parameters, and redundant interaction steps (Wang et al., 2025a). Semantically unnecessary actions are filtered to ensure concise and meaningful interaction trajectories.

To improve cross-platform learnability, we normalize raw GUI actions into a compact and standardized representation. Platform-specific redundancies are removed by extracting core interaction primitives, e.g., collapsing compound operations such as `select_from_to` followed by `copy` into a single `COPY` action. Actions with equivalent semantics are further unified into a common format (e.g., `press_enter` and `press_button` are mapped to `HOTKEY`), ensuring semantic consistency across mo-

bile, web, and desktop environments.

We adopt a unified and extensible action space to support GUI interaction across mobile, web, and desktop platforms. As summarized in Table 16, the action space consists of 17 discrete action types, each optionally associated with structured attributes such as screen coordinates, text inputs, or key sequences.

The action space is organized into three categories: (1) Basic Actions, which capture platform-agnostic operations, including clicking, typing, scrolling, waiting, and task termination; (2) Custom Actions for Mobile, which extend the basic set with smartphone-specific interactions such as long-press gestures, application launching, and system-level navigation; and (3) Custom Actions for Web & Desktop, which incorporate mouse-

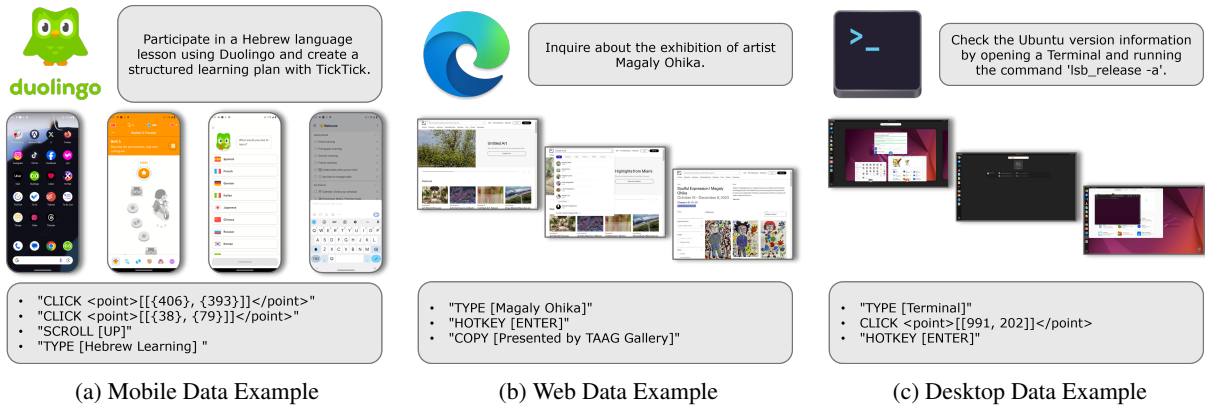
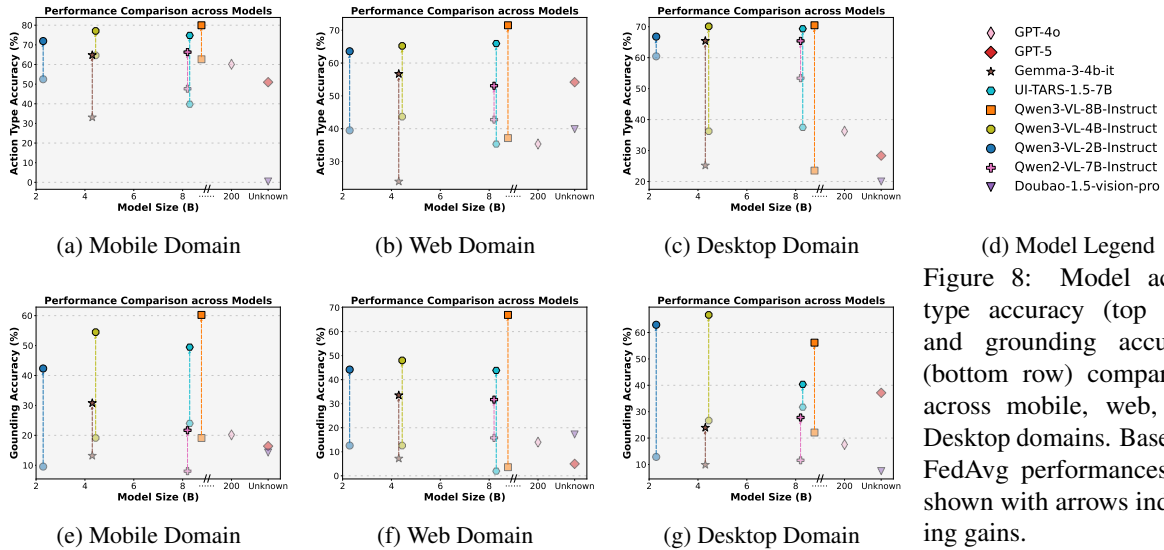


Figure 9: Example data episodes for training GUI agents across mobile, web, and desktop platforms in FedGUI.

and keyboard-centric operations including double-clicking, right-clicking, cursor movement, hotkeys, and clipboard interactions. This unified design enables consistent modeling and evaluation of heterogeneous GUI tasks while preserving platform-specific interaction semantics, and all predicted actions are constrained to this predefined action space during both training and evaluation.

### C.3 Source Datasets

We consolidate multiple publicly available GUI interaction datasets spanning mobile, web, and desktop environments. These datasets provide diverse task formulations, action spaces, and interface modalities, enabling comprehensive evaluation of federated GUI agents under heterogeneous data distributions. Figure 10 provides a visual overview of the specific interface characteristics across these benchmarks.

AitW is a large-scale mobile device control dataset containing screenshots, human interaction

trajectories, and natural language instructions collected across more than 800 Android devices and applications. It supports multi-step interaction tasks and emphasizes generalization across apps and device configurations.

AndroidControl focuses on common Android application tasks with hierarchical task formulations, including both high-level goals and fine-grained low-level instructions. The dataset contains approximately 15K episodes and is widely adopted for studying data scaling effects and instruction-following behaviors in mobile agents.

GUI Odyssey is a mobile GUI navigation dataset designed for cross-application reasoning. It covers a wide range of device types and hundreds of mobile applications, with step-wise annotations that explicitly capture intermediate reasoning and navigation decisions.

**GUIAct-Web (GA-W).** GUIAct-Web is a web-based GUI action dataset comprising both single-

| Index                                       | Action Type   | Attribute  | Description  |
|---|---------------|------------|--|
| <b>Basic Actions</b>                        |               |            |  |
| 1   | CLICK         | (x, y)     | Click at specific screen coordinates.  |
| 2   | TYPE          | input_text | Type the given input text into an active text field.                           |
| 3   | SCROLL        | direction  | Scroll in a specific direction (up, down, left, or right).                     |
| 4   | COMPLETE      | –          | Mark the task as completed.  |
| 5   | IMPOSSIBLE    | –          | Indicate that the current task cannot be completed.                            |
| 6   | WAIT          | –          | Wait for a period (e.g., loading or animation).                                |
| <b>Custom Actions for Mobile</b>            |               |            |  |
| 7   | LONG_PRESS    | (x, y)     | Long press at the given coordinates.   |
| 8   | OPEN_APP      | app_name   | Open a specific application by name.   |
| 9   | NAVIGATE_BACK | –          | Go back to the previous page or screen.  |
| 10  | NAVIGATE_HOME | –          | Return to the home screen.   |
| 11  | PRESS_RECENT  | –          | Open the recent apps menu.   |
| 12  | PRESS_ENTER   | –          | Simulate pressing the Enter key.   |
| <b>Custom Actions for Web &amp; Desktop</b> |               |            |  |
| 13  | DOUBLE_CLICK  | (x, y)     | Perform a double click at the specified coordinates.                           |
| 14  | RIGHT_CLICK   | (x, y)     | Right click at the given position.   |
| 15  | MOVETO        | (x, y)     | Move the cursor to a target coordinate without clicking.                       |
| 16  | HOTKEY        | keys       | Execute a keyboard shortcut (e.g., ENTER, ESC, DOWN, SPACE, RIGHT, LEFT, TAB). |
| 17  | COPY          | text       | Copy the given text to clipboard.  |

Table 16: Action space used in FedGUI. The table summarizes the unified set of 17 executable actions, consisting of platform-agnostic **basic actions** and platform-specific custom actions for **mobile** and **web & desktop** environments. Each action is defined with its required attributes and semantic description.

step and multi-step instructions. It includes fundamental operations such as clicking, scrolling, and text input, and is designed to benchmark web interaction understanding and action prediction.

**Mind2Web (M2W).** Mind2Web targets generalist web agents and is collected from 137 websites across 31 domains. The dataset contains complex multi-step tasks and places particular emphasis on cross-site and cross-domain generalization.

**AgentSynth (AS).** AgentSynth is a synthetic dataset generated via a scalable task and trajectory synthesis pipeline, primarily collected in Linux environments. It enables low-cost construction of diverse GUI tasks and serves as a complementary data source for training and evaluating generalist agents.

**OmniAct-Web (OA-W).** OmniAct-Web is the web subset of the OmniAct benchmark. It supports end-to-end evaluation of web interaction capabilities and facilitates systematic comparison between web and desktop GUI tasks under a unified action abstraction.

**OmniAct-MacOS (OA-Mac).** OmniAct-MacOS is a desktop interaction dataset specifically collected within the macOS environment. It captures unique visual cues and interface paradigms of the macOS environment, enabling the evaluation

of cross-platform agent adaptation and execution.

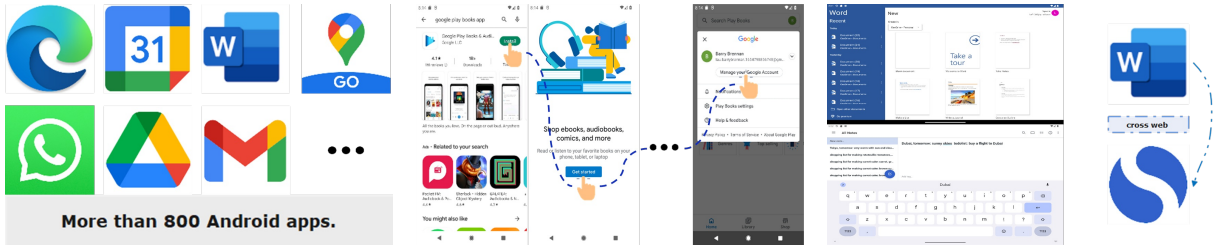
**OmniAct-Windows (OA-Win).** OmniAct-Windows is the Windows-specific subset of OmniAct, focusing on productivity tasks in Microsoft Windows. It provides executable scripts paired with GUI states, enabling rigorous evaluation of task automation and correctness across Windows software.

## D Experiment Details

### D.1 Evaluation & Metric Details

**Standardized Evaluation Sets.** To ensure a robust and comprehensive assessment, we curate a standardized evaluation set by sampling 100 interaction episodes from the test partition of each source dataset. The sampling process is designed to capture diversity across multiple axes, including varying trajectory lengths and distinct task categories. This balanced selection ensures that the evaluation is not biased toward specific simple tasks and reflects the complexity of real-world GUI environments.

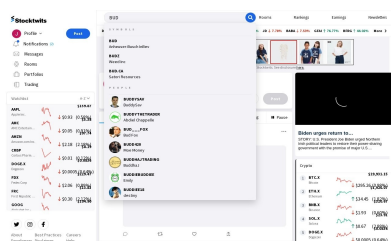
We implement a hierarchical evaluation protocol to dissect the agent’s performance at different granularities. Following the implementation in our evaluation pipeline, we define three primary metrics:



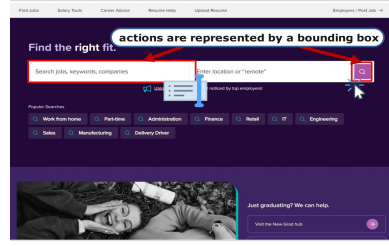
(a) AndroidControl (AC) contains over 800 Android applications.

(b) Android-in-the-Wild (AitW) is characterized by its procedural, multi-step tasks.

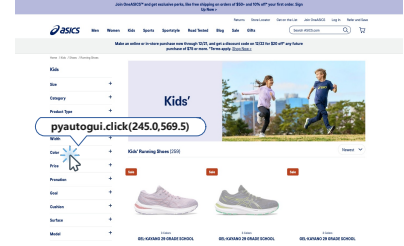
(c) GUI Odyssey (GO) is a benchmark for cross-app GUI navigation.



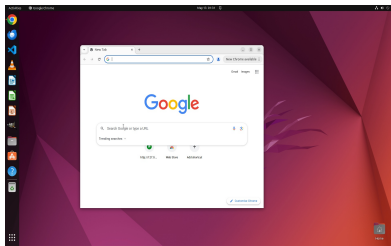
(d) Mind2Web (M2W) covers diverse domains, websites, and user tasks.



(e) GUIAct-Web (GA-W), whose actions are represented by bounding boxes.



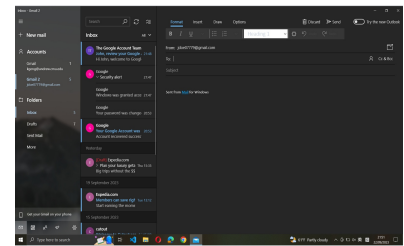
(f) OmniAct-Web (OA-W) provides executable PyAutoGUI scripts.



(g) AgentSynth (AS) is collected on Linux Ubuntu environments.



(h) OmniAct-MacOS (OA-Mac) is collected on macOS environments.



(i) OmniAct-Windows (OA-Win) is collected on Windows environments.

Figure 10: Overview of the datasets used in our experiments across mobile (AC, AitW, GO), web (M2W, GA-W, OA-W), and desktop (AS, OA-Mac, OA-Win) environments. These datasets cover diverse task formulations, interface modalities, and interaction patterns, enabling evaluation of federated GUI agents under various heterogeneity.

- **Action Type Accuracy (Type):** This measures the consistency between the predicted and ground-truth interaction intent. A prediction is marked as a successful type match if the first token of the generated action string aligns with the target category in our unified action space.
- **Grounding Accuracy (Ground):** For coordinate-based actions such as CLICK and DOUBLE\_CLICK, we adopt a spatial proximity criterion. A grounding attempt is considered successful if the Euclidean distance between the predicted coordinates and the ground-truth coordinates is within a specific threshold. In our implementation, this threshold is set to 14% of the screen diagonal. This relative ratio allows the benchmark to accommodate various UI element scales and screen resolutions

across heterogeneous federated clients.

- **Success Rate (SR):** This metric represents the end-to-end execution accuracy, requiring both the action type and its associated parameters to be correct. For text-centric actions like TYPE, OPEN\_APP, or COPY, we evaluate semantic alignment using a **Similarity Score**. This score integrates F1-token matching and character-level intersection to support multilingual scenarios. As defined in our pipeline, a prediction is deemed successful if the similarity score exceeds 0.5, ensuring the agent captures the core semantic intent even with minor variations in the generated text.

## D.2 Training Details

This section summarizes the key configurations used in our experiments, ensuring reproducibility

and clarity.

The framework provides adaptive default hyperparameters for multiple federated optimization algorithms. FedYogi (Reddi et al., 2020) uses momentum parameters ( $\beta_1 = 0.9, \beta_2 = 0.999$ ) with learning rate  $\eta = 10^{-3}$  and stabilization constant  $\tau = 10^{-6}$ . FedAvgM (Hsu et al., 2019) applies a 0.9/0.1 weighting between historical and current models. FedProx (Li et al., 2020) incorporates proximal regularization with coefficient  $\mu = 0.2$  via a  $\|w - w^t\|^2$  penalty term. SCAFFOLD (Karimireddy et al., 2020) maintains a server learning rate of  $\eta_s = 1.0$  with client-side control variate compensation, while FedAdam and FedAdam (Reddi et al., 2020) share base parameters ( $\beta_1 = 0.9, \beta_2 = 0.999$ ) with adaptive learning rate scaling. All algorithm-specific coefficients are exposed through a unified parameter interface.

**General Parameters.** Our implementation is based on the `ms-swift` library (Zhao et al., 2024) and adopts parameter-efficient fine-tuning. We employ Low-Rank Adaptation (LoRA) with a rank of 8 and an alpha scaling factor of 32, together with a dropout rate of 0.05 to mitigate overfitting. The maximum sequence length is set to 4,096 tokens. We use a batch size of 1 with a gradient accumulation step of 4, and the learning rate is fixed to  $5 \times 10^{-5}$  across all experiments.

**Hardware Configuration.** Training is conducted on two NVIDIA GeForce RTX 3090 GPUs with CUDA version 12.4. Under this configuration, training a single client for one communication round with 10 episodes requires approximately 2 minutes.

### D.3 Model Details.

We employ two categories of models in our experiments: open-source vision-language models (VLMs) that support local training and federated optimization, and API-based closed-source models used exclusively for annotation and evaluation.

**Open-Source VLMs.** This category forms the backbone of our federated training framework. We include a diverse set of state-of-the-art open-source VLMs spanning multiple model families and parameter scales, including Qwen (Qwen3-VL-2B/4B/8B, Qwen2.5-VL-3B/7B, Qwen2-VL-2B/7B) (Wang et al., 2024), InternVL (InternVL2-1B/2B/4B/8B) (Chen et al., 2024), Gemma (Gemma-3-4B) (Team,

| Configuration         | Value                               |
|-----------------------|-------------------------------------|
| Model Architecture    | Qwen2VLForConditional<br>Generation |
| Model Type            | qwen2_v1                            |
| Torch Dtype           | bf16                                |
| Tokenizer             | Qwen2TokenizerFast                  |
| Tokenization Strategy | Greedy search                       |

Table 17: Model and tokenizer configurations.

2025), DeepSeek (DeepSeekVL2, DeepSeekVL2-tiny/small) (Lu et al., 2024), and Seed models (UI-TARS-1.5-7B) (Qin et al., 2025). These models support supervised fine-tuning and parameter-efficient adaptation (e.g., LoRA), enabling systematic evaluation of federated learning algorithms on heterogeneous GUI interaction data across mobile, web, and desktop platforms.

We additionally incorporate commercial vision-language APIs, including GPT-5, GPT-4o, and Doubao-1.5-Vision. Due to their closed-source and API-only nature, these models do not support supervised fine-tuning within our framework and are therefore exclusively used as annotation and reference models, providing strong off-the-shelf baselines for GUI understanding and action grounding.

**Model & Tokenization Configuration.** Our primary experiments are conducted using the Qwen2-VL model. The detailed model and tokenizer configuration is provided in Table 17.

### D.4 Prompt Format

Following FedMobileAgent (Wang et al., 2025b), we employ a structured prompt format, shown in Figure 11 and 12, to support executable decision-making. The prompt integrates a high-level task goal, the current visual context (i.e., the screen screenshot), and an explicit action set comprising both basic and domain-specific actions. This design grounds model decisions in the task objective while constraining outputs to valid, screen-aware GUI interactions.

## Full System Prompt (Part I)

You are a foundational action model capable of automating tasks across various digital environments, including desktop systems like Windows, macOS, and Linux, as well as mobile platforms such as Android and iOS. You also excel in web browser environments. You will interact with digital devices in a human-like manner: by reading screenshots, analyzing them, and taking appropriate actions. Your expertise covers two types of digital tasks:

- **Grounding:** Given a screenshot and a description, you assist users in locating elements mentioned. Sometimes, you must infer which elements best fit the description when they aren't explicitly stated.
- **Executable Language Grounding:** With a screenshot and task instruction, your goal is to determine the executable actions needed to complete the task.

You are now operating in Executable Language Grounding mode. Your goal is to help users accomplish tasks by suggesting executable actions that best fit their needs. Your skill set includes both basic and custom actions:

Basic actions are standardized and available across all platforms. They provide essential functionality and are defined with a specific format, ensuring consistency and reliability.

- **Basic Action 1: CLICK**
  - purpose: Click at the specified position.
  - format: CLICK <point>[[x-axis, y-axis]]</point>
  - example usage: CLICK <point>[[101, 872]]</point>
- **Basic Action 2: TYPE**
  - purpose: Enter specified text at the designated location.
  - format: TYPE [input text]
  - example usage: TYPE [Shanghai shopping mall]
- **Basic Action 3: SCROLL**
  - purpose: Scroll in the specified direction.
  - format: SCROLL [direction (UP/DOWN/LEFT/RIGHT)]
  - example usage: SCROLL [UP]
- **Basic Action 4: COMPLETE**
  - purpose: Indicate the task is finished.
  - format: COMPLETE
  - example usage: COMPLETE
- **Basic Action 5: IMPOSSIBLE**
  - purpose: Indicate the task is infeasible to reach.
  - format: IMPOSSIBLE
  - example usage: IMPOSSIBLE
- **Basic Action 6: WAIT**
  - purpose: Wait for the screen to load.
  - format: WAIT
  - example usage: WAIT

### 2. Custom Actions for Mobile Platforms

Custom actions are unique to each user's platform and environment. They allow for flexibility and adaptability, enabling the model to support new and unseen actions defined by users.

- **Custom Action 7: LONG\_PRESS**
  - purpose: Long press at the specified position.
  - format: LONG\_PRESS <point>[[x-axis, y-axis]]</point>
  - example usage: LONG\_PRESS <point>[[272, 341]]</point>
- **Custom Action 8: NAVIGATE\_BACK**
  - purpose: Press a back button to navigate to the previous screen.
  - format: NAVIGATE\_BACK
- **Custom Action 9: NAVIGATE\_HOME**
  - purpose: Press a home button to navigate to the home page.
  - format: NAVIGATE\_HOME
- **Custom Action 10: OPEN\_APP**
  - purpose: Open the specified application.
  - format: OPEN\_APP [app\_name]
  - example usage: OPEN\_APP [Google Chrome]
- **Custom Action 11: PRESS\_RECENT**
  - purpose: Press the recent button to view or switch between recently used applications.
  - format: PRESS\_RECENT

Figure 11: The full system prompt for training and evaluation (Part I).

## Full System Prompt (Part II)

### 3. Custom Actions for Web and Desktop Platforms

- **Custom Action 12: DOUBLE\_CLICK**
  - purpose: Double click at the specified position.
  - format: DOUBLE\_CLICK <point>[[x-axis, y-axis]]</point>
- **Custom Action 13: RIGHT\_CLICK**
  - purpose: Right click at the specified position.
  - format: RIGHT\_CLICK <point>[[x-axis, y-axis]]</point>
- **Custom Action 14: MOVETO**
  - purpose: Move the object to the specified position.
  - format: MOVETO <point>[[x-axis, y-axis]]</point>
- **Custom Action 15: HOTKEY**
  - purpose: Use the hot key.
  - format: HOTKEY [keys]
  - example usage: HOTKEY [CTRL+ALT]
- **Custom Action 16: COPY**
  - purpose: Copy a sentence to answer user questions.
  - format: COPY [text with answer]
  - example usage: COPY [Wednesday]

In most cases, task instructions are high-level and abstract. Carefully read the instruction and action history, then perform reasoning to determine the most appropriate next action. Ensure you strictly generate one section: Actions.

**Actions:** Specify the actual actions you will take based on your reasoning.

Figure 12: The full system prompt for training and evaluation (Part II).