

ANDROIDWORLD: A DYNAMIC BENCHMARKING ENVIRONMENT FOR AUTONOMOUS AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

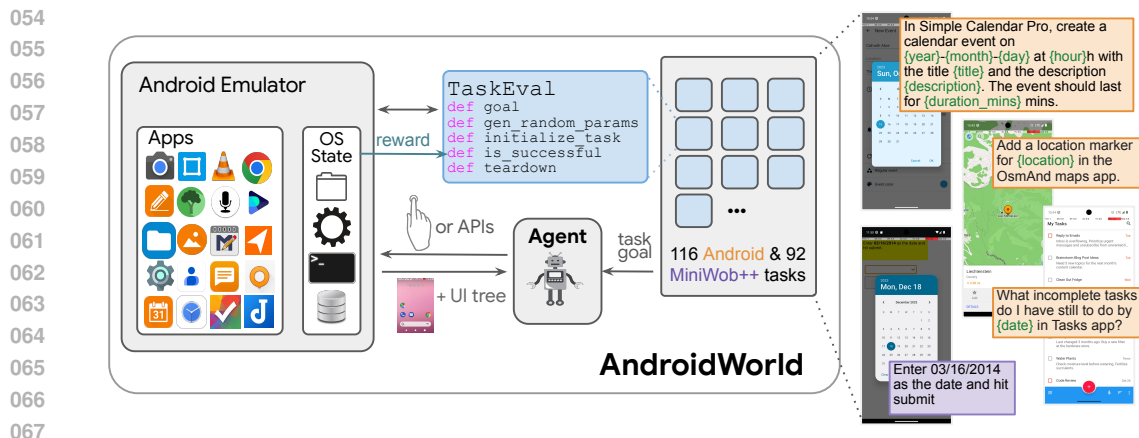
Autonomous agents that execute human tasks by controlling computers can enhance human productivity and application accessibility. However, progress in this field will be driven by realistic and reproducible benchmarks. We present ANDROIDWORLD, a fully functional Android environment that provides reward signals for 116 programmatic tasks across 20 real-world Android apps. Unlike existing interactive environments, which provide a static test set, ANDROIDWORLD dynamically constructs tasks that are parameterized and expressed in natural language in unlimited ways, thus enabling testing on a much larger and more realistic suite of tasks. To ensure reproducibility, each task includes dedicated initialization, success-checking, and tear-down logic, which modifies and inspects the device’s system state.

We experiment with baseline agents to test ANDROIDWORLD and provide initial results on the benchmark. Our best agent can complete 30.6% of ANDROIDWORLD’s tasks, leaving ample room for future work. Furthermore, we adapt a popular desktop web agent to work on Android, which we find to be less effective on mobile, suggesting future research is needed to achieve universal, cross-platform agents. Finally, we also conduct a robustness analysis, showing that task variations can significantly affect agent performance, demonstrating that without such testing, agent performance metrics may not fully reflect practical challenges.

1 INTRODUCTION

Autonomous agents that interpret natural language instructions and operate computing devices can provide enormous value to users by automating repetitive tasks, augmenting human intelligence, and accomplishing complex workflows. However, a key research challenge remains the realistic evaluation of these agents in real-world settings. Despite growing enthusiasm for building autonomous agents (Rawles et al., 2023; Deng et al., 2023; Zheng et al., 2024a; Koh et al., 2024; Kim et al., 2024; He et al., 2024; Gravitas, 2023; Age, 2024; Web, 2024; Wu et al., 2023; Xie et al., 2023), most existing approaches for evaluation compare an agent’s actions at each step to a previously collected human demonstration (Deng et al., 2023; Rawles et al., 2023; Yang et al., 2023b; Zhang & Zhang, 2023; Lù et al., 2024; Zhang et al., 2024c; Yan et al., 2023; Li et al., 2024). Measuring performance in this way can be misleading because when performing tasks online in real environments agents can take multiple paths to solve tasks, environments may behave non-deterministically, and agents can dynamically learn from mistakes to correct their actions (Shinn et al., 2023; Liu et al., 2018b; Li et al., 2023b; Pan et al., 2024). For this reason, online evaluation of agents in realistic environments able to reward task outcome provides a gold standard for evaluation. While there is an emerging body of work to address this need across different environments (Zhou et al., 2023; Koh et al., 2024; Drouin et al., 2024; Lee et al., 2024; Xie et al., 2024; Bonatti et al., 2024), there is no comprehensive solution for mobile platforms, such as Android, which are used by billions of users and therefore represent environments in which automation agents may be very productively employed. We introduce ANDROIDWORLD to address this.

At its core, ANDROIDWORLD offers a reliable means of obtaining reward signals for tasks performed by agents in realistic mobile environments. Reward signals are quantitative metrics that indicate functional correctness of a task, i.e. is the stated goal achieved? For example, for the task “Send a text message to Jane confirming I’ll be there,” a positive reward indicates that the relevant



068 Figure 1: ANDROIDWORLD is an environment for building and testing autonomous agents.

069
070
071 message has been sent. Unlike simulated environments (Tassa et al., 2018; Shridhar et al., 2020) or
072 games (Mnih et al., 2013; Silver et al., 2016; Vinyals et al., 2019; Wang et al., 2023b; Tan et al.,
073 2024; Toyama et al., 2021), real-world apps and websites do not inherently offer explicit reward sig-
074 nals. While human (Rawles et al., 2023; Zheng et al., 2024a; Pan et al., 2024; Kinniment et al., 2023)
075 or LLM-based (Chiang et al., 2024; Zheng et al., 2023; Liu et al., 2023; Du et al., 2023; Ma et al.,
076 2023; Pan et al., 2024; He et al., 2024) judges can be employed to reward the outcome of a task,
077 these approaches scale poorly or are not fully reliable, respectively. Alternatively, environments
078 for autonomous agents which provide automated ground-truth rewards for complex workflows have
079 been developed (Yao et al., 2023; Zhou et al., 2023; Koh et al., 2024; Xie et al., 2024; Bonatti
080 et al., 2024). We find two problems with these environments. First, they are constrained to desktop
081 computing environments, overlooking the mobile domain, which is of paramount importance given
082 the ubiquity and diversity of mobile devices in the real world. Secondly, they are limited in their
083 real-world diversity and scale. Crucially, unlike in real-world scenarios where conditions and task
084 inputs vary widely, these environments support only static test specifications, meaning that when
085 task parameters deviate, the reward signal is likely to break.

085 We seek to develop a comprehensive benchmark that addresses the limitations of the existing
086 approaches above for evaluating automation agents in mobile environments. ANDROIDWORLD does
087 this by spanning 20 Android apps on a total of 116 programmatic tasks to provide ground truth-
088 rewards. Unlike existing test environments (with MiniWoB++ (Shi et al., 2017) being a notable ex-
089 ception), each task in ANDROIDWORLD is dynamically instantiated using randomly-generated param-
090 eters, challenging agents with millions of unique task goals and conditions. While MiniWoB++
091 consists of simple, synthetic websites, ANDROIDWORLD leverages actual Android applications. A
092 main challenge that ANDROIDWORLD must address is how to ensure that reward signals are durable
093 when using real-world applications and varying task parameters dynamically. ANDROIDWORLD’s
094 key insight is to leverage the extensive and consistent state management capabilities of the Android
095 operating system, using the same mechanisms that the apps themselves utilize to store and update
096 data.

096 In addition to providing a comprehensive benchmark, ANDROIDWORLD is lightweight, requiring
097 only 2 GB of memory and 8 GB of disk space, and is designed with convenience in mind. It
098 connects agents to the Android OS by leveraging the Python library AndroidEnv (Toyama et al.,
099 2021) to connect to the freely available Android Emulator.¹ In addition to the 116 Android tasks,
100 we extend ANDROIDWORLD with web tasks by integrating the MiniWoB++ (Shi et al., 2017; Liu
101 et al., 2018a) benchmark into it.

102 To demonstrate ANDROIDWORLD’s usefulness as a benchmark, we build and release a multi-modal
103 agent, M3A (Multimodal Autonomous Agent for Android), and establish state-of-the-art results on
104 ANDROIDWORLD. We analyze M3A’s performance using both multimodal and text-only input,
105 and we observe that while multimodal perception can improve performance in some cases, it gen-
106

107 ¹The Android Emulator is packaged as part of Android Studio, which can be downloaded from
https://developer.android.com/studio

erally does not outperform the text-only approach. On ANDROIDWORLD, M3A achieves a 30.6% success rate, which surpasses that of a web agent adapted for Android but remains significantly lower than the human success rate of 80.0%. In pursuit of building robust UI control agents, our study includes comprehensive tests under varied real-world conditions, demonstrating significant performance variations primarily driven by changes in intent parameters.

Overall, we make the following contributions: (i) the creation of a new, highly diverse and realistic mobile UI control agent environment; (ii) establishment of benchmark performance with a state-of-the-art multimodal agent, and (iii) a careful analysis demonstrating the need to evaluate agents across variable task parameters and conditions due to the inherent stochasticity in both models and environments.

2 RELATED WORK

Table 1 compares existing evaluation environments for autonomous UI agents.

2.1 INTERACTIVE EVALUATION ENVIRONMENTS

Effective evaluation of autonomous agents requires benchmarks that mimic real-world scenarios, but also interactive environments that provide reward signals upon successful task completion (Rawles et al., 2023; Deng et al., 2023; Abramson et al., 2022; Ruan et al., 2023; Chen et al., 2021). Many existing benchmarking environments target web browsing. MiniWoB++ (Shi et al., 2017; Liu et al., 2018b) consists of small, synthetic HTML pages with parameterizable tasks which allow for unlimited task variability. WebShop (Yao et al., 2023) provides a simulated e-commerce environment, whereas WebArena (Zhou et al., 2023) and VisualWebArena (Koh et al., 2024) consist of simulated websites across up to six domains. WorkArena (Drouin et al., 2024) consists of 29 tasks for enterprise software. GAIA (Mialon et al., 2023) is a static dataset that tests an agent’s ability to interact with live web environments. MMInA (Zhang et al., 2024e) is a multihop and multimodal benchmark designed to evaluate agents for compositional Internet tasks.

Towards building generalist agents that control operating systems (OSes), OSWorld (Xie et al., 2024) and WindowsAgentArena (Bonatti et al., 2024) provide a test suite of tasks for desktop computer use cases and custom execution-based evaluation scripts across across 9 and 11 apps, respectively. In the mobile domain, existing benchmarks are limited and do not capture the diversity of real-world mobile interactions, containing low-complexity tasks or on a limited number of applications. B-MoCA’s (Lee et al., 2024) evaluation is based on 6 simple tasks (e.g., “Call 911”, “turn on airplane mode”) across 4 apps², validated using regular expressions. Mobile-Env (Zhang et al., 2024b) offers task reproducibility limited to 13 task templates for a single app (WikiHow).

While ANDROIDWORLD shares the mobile platform focus of B-MoCA and Mobile-Env, it is more comparable to OSWorld (and WindowsAgentArena, which builds on top of OSWorld) in terms of task complexity and the diversity of interactions it supports. ANDROIDWORLD enhances OSWorld’s approach by dynamically constructing the start states of an agent’s run and varying the task parameters in unlimited ways, thus allowing for a new type of evaluation under varying real-world conditions.

Other studies leverage human evaluation (Rawles et al., 2023; Zheng et al., 2024a; Bishop et al., 2024) for tasks where automatic evaluation is not available. Lastly, emerging research (Pan et al., 2024; He et al., 2024; Xing et al., 2024) explores the potential of multimodal models to generalize agent evaluations to new settings, though this area requires further research to achieve accuracy comparable to manually-coded rewards.

AndroidEnv (Toyama et al., 2021) provides a mechanism to manage communication with the Android emulator, similar to Playwright and Selenium for web environments. While ANDROIDWORLD leverages this functionality, it diverges in its reward system. AndroidEnv’s approach requires modifying application source code and implementing task-specific logging statements, making it well-suited for gaming environments with easily verifiable success criteria. In contrast, ANDROIDWORLD implements a non-invasive reward mechanism, allowing it to create a benchmark suite

²Based on what reported in the Experiments Section of the B-MoCA manuscript as of October 1st, 2024.

Table 1: Comparison of different datasets and environments for benchmarking computer agents.

	Env?	# of apps or websites	# task templates	Avg # task instances	Reward method	Platform
GAIA	✗	n/a	466	1	text-match	None
MIND2WEB	✗	137	2350	1	None	Desktop Web
WEBLINX	✗	155	2337	1	None	Desktop Web
WEBVOYAGER	✗	15	643	1	LLM judge	Desktop Web
PIXELHELP	✗	4	187	1	None	Android
METAGUI	✗	6	1125	1	None	Android
MOTIF	✗	125	4707	1	None	Android (Apps+Web)
AITW	✗	357+	30378	1	None	Android (Apps+Web)
ANDROIDCONTROL	✗	833	15283	1	None	Android (Apps+Web)
OMNIACT	✗	60+	9802	1	None	Desktop (Apps+Web)
ANDROIDARENA	✗	13	221	1	Action match/LLM	Android (Apps+Web)
LLAMATOUCH	✗	57	496	1	Screen match	Android (Apps+Web)
MINIWOB++	✓	1	114	∞	HTML/JS state	Web (synthetic)
WEBSHOP	✓	1	12k	1	product attrs match	Desktop Web
WEBARENA	✓	6	241	3.3	url/text-match	Desktop Web
VISUALWEBARENA	✓	4	314	2.9	url/text/image-match	Desktop Web
WORKARENA	✓	1	29	622.4	cloud state	Desktop Web
MOBILE-ENV	✓	1	13	11.5	regex	Android (Apps)
B-MOCA	✓	4	6	1.9	regex	Android (Apps+Web)
MMINA	✓	14	1050	1	text-match	Desktop web
OSWORLD	✓	9	369	1	device/cloud state	Desktop (Apps+Web)
WINDOWSAGENTARENA	✓	11	154	1	device state	Desktop (Apps+Web)
ANDROIDWORLD	✓	20	116	∞	device state	Android (Apps+Web)

for apps where source code access is unavailable and to reuse validation components across different apps. This approach enables ANDROIDWORLD to cover a broader range of real-world mobile tasks.

2.2 STATIC DATASETS FOR UI AUTOMATION

Datasets derived from human interactions provide proxy metrics that correlate with real-world agent performance (Li et al., 2020; Burns et al., 2021; Deng et al., 2023; Rawles et al., 2023). On mobile platforms, AitW (Rawles et al., 2023), AndroidControl (Li et al., 2024), PixelHelp (Li et al., 2020), AndroidArena (Xing et al., 2024), LlamaTouch (Zhang et al., 2024d), UGIF (Venkatesh et al., 2022), and MoTIF (Burns et al., 2021) consist of demonstrations across Android apps and mobile websites, with screens often represented via accessibility trees. In contrast, desktop web environments typically utilize the DOM for representing website content, with Mind2Web (Deng et al., 2023), OmniAct (Kapoor et al., 2024) and others, across various desktop websites. Mobile-based datasets frequently involve more complex actions, such as scrolling, which are not as useful in DOM-based desktop interactions where the entire action space is readily accessible. Additionally, API-centric datasets like API-Bank (Li et al., 2023a), ToolTalk (Farn & Shin, 2023), and ToolBench (Xu et al., 2023) assess agents’ capabilities to manipulate computer systems via APIs.

2.3 INTERACTIVE AGENTS

Prior to today’s foundation models, traditional approaches to developing user interface-operating agents primarily used reinforcement learning and behavioral cloning to simulate interactions like mouse clicks and keyboard typing (Liu et al., 2018b; Li et al., 2020; Shvo et al., 2021; Gur et al., 2022a; Humphreys et al., 2022). More recent work tends to leverage off-the-shelf foundational models (Team, 2023a;b; Touvron et al., 2023) with in-context learning (ICL) and fine-tuning applied to mobile (Rawles et al., 2023; Hong et al., 2023; Wang et al., 2023a; Yan et al., 2023; Zhang & Zhang, 2023; Bishop et al., 2024; Zhang et al., 2023), desktop web (Zheng et al., 2024a; Deng et al., 2023; Zhou et al., 2023; Koh et al., 2024; Cheng et al., 2024; Lai et al., 2024), and desktop OS (Wu et al., 2024; Zhang et al., 2024a; Xie et al., 2024). Recent work explores agents that reflect on system state (Shinn et al., 2023; Yao et al., 2022; Madaan et al., 2024) by leveraging exploration, self-evaluation, and retry-capabilities enabling continual learning and adaptation (Li et al., 2023b; Yang et al., 2023b; Pan et al., 2024; Wu et al., 2024; Gao et al., 2023; Murty et al., 2024).

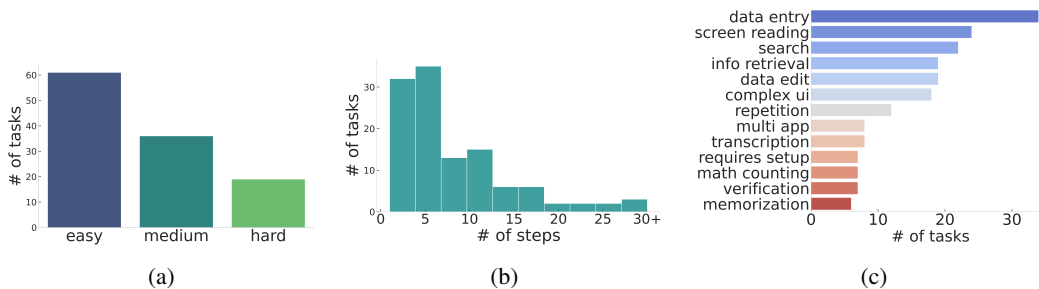


Figure 2: Annotators performed the tasks assigned to them, assigned a difficulty level (2a) and estimated the number of steps required to complete each task (2b), using the action space available to an agent. For each task, they selected relevant category tags from a predefined list (2c).

3 ANDROIDWORLD

3.1 ANDROID FOR AUTONOMOUS AGENTS

Android is an ideal environment for developing autonomous agents. It is the most widely-used OS globally³ and is highly flexible for research, while providing an open world of the Web⁴ and over 2M apps for agents to operate in. Using emulation, an Android environment is easy to deploy, does not require specialized hardware, and can be run on a laptop. Android Virtual Devices or emulator images are well suited for research as they are self-contained, easy to distribute, and configurable.

Compared to desktop environments, mobile environments, like Android, pose unique research challenges for computer control agents. On one hand, mobile UIs tend to be simpler than their desktop counterparts because of their smaller screen size. On the other hand, the action space on mobile devices is more complicated and more actions can be required to complete tasks. Precise gestures are needed to fully operate the UI, such as when navigating a carousel widget, long-pressing on a widget, or performing multi-finger gestures to zoom in. Since it is an OS, Android is a fully open environment compared to web-browser-only environments. Android’s flexibility is also reflected in its action space; in addition to UI actions (click, scroll, type, etc.), Android provides function-calling APIs, such as sending a text message, for example, which allow computer control agents to utilize a broader action space.

3.2 THE OBSERVATION AND ACTION SPACE

ANDROIDWORLD provides an interface for agents to receive observations and execute actions on Android. It uses AndroidEnv (Toyama et al., 2021) and the Android Device Bridge to facilitate interaction between Android and the agent. The observation space consists of a full-resolution screenshot and a UI tree representation developed for accessibility purposes. The action space is similar to that which humans use, consisting of gestures (i.e., tapping, long-press, and swiping), typing, and navigation buttons (i.e., go home and go back). In addition to these naturalistic actions, ANDROIDWORLD exposes a limited set of function calling APIs, such as `send_text_message`, to help agents accomplish goals. Appendix B provides more details on the observation format and action space.

3.3 REPRODUCIBLE AND PARAMETERIZED TASKS

ANDROIDWORLD consists of a suite of 116 tasks, spread across 20 diverse applications (see Appendix C for more details). These tasks simulate practical, everyday activities, including note-taking, scheduling appointments, communicating through messaging, and interacting with system utilities. The suite consists of open-source apps and built-in Android system apps, such as Settings and Contacts. As rated by humans, the tasks vary in difficulty, duration, and categories (Figure 2).

³<https://gs.statcounter.com/os-market-share>

⁴Mobile is the most popular platform for accessing the web; <https://gs.statcounter.com/platform-market-share/desktop-mobile/worldwide/>

Table 2: Selected tasks with code describing validation logic.

Task	Validation code
In Simple Calendar Pro, create a calendar event on $\{event.year\}$ - $\{event.month\}$ - $\{event.day\}$ at $\{event.hour\}$ h with the title ' $\{event.title\}$ ' and the description ' $\{event.description\}$ '. The event should last for $\{event.duration\}$ mins.	<code>event_exists(event)</code>
Send a text message to $\{phone_number\}$ with message: $\{message\}$.	<code>message_exists(phone_number, message, messaging_db)</code>
Create a new drawing in Simple Draw Pro. Name it $\{file_name\}$. Save it in the Pictures folder.	<code>file_exists(file_path)</code>
Create a timer with $\{hours\}$ hours, $\{minutes\}$ minutes, and $\{seconds\}$ seconds. Do not start the timer.	<code>timer_displays(time, ui_hierarchy)</code>
Create a new note in Markor named $\{file_name\}$ with the following text: $\{text\}$. Share the entire content of the note with the phone number $\{number\}$ via SMS.	<code>(file_exists(file_name, content=text) + message_exists(phone_number, message)) / 2.0</code>
Turn on WiFi and open $\{app_name\}$.	<code>(wifi_enabled() + app_launched(app_name)) / 2.0</code>

To achieve a high degree of reproducibility in real-world scenarios, ANDROIDWORLD precisely controls the OS and app states in several ways. The Android OS is fixed, consisting of a Pixel 6 emulator running Android 13 with a fixed time on the date October 15th, 2023. All applications in ANDROIDWORLD are fully-functional and consists of both open-source apps and OS-level apps included with Android. For the open-source apps, ANDROIDWORLD maintains a constant environment by installing a fixed version of each app, acquired from F-Droid.⁵ OS-level apps' versions are determined by the Android OS, which is also fixed. To maintain a reproducible environment, ANDROIDWORLD utilizes apps that do not require login/authentication and can store their application data on device.

In addition to managing the states of apps and operating systems, ANDROIDWORLD precisely defines and controls the state during task execution. Each task has its own unique setup, reward determination logic, and teardown procedures (see Appendix C.2 and C.3 for more details), ensuring a fully reproducible suite of tasks.

Automatic task *parameterization* is a critical mechanism, unique to ANDROIDWORLD, to evaluate agents on a much larger and more realistic suite of tasks than current benchmarks support. Achieving this requires significantly more effort than randomly generating new task parameters because it involves developing evaluation logic that remains valid across different task instantiations. It is exactly through its careful state management that in addition to reproducibility AndroidWorld ensures that the reward mechanisms function correctly. Task parameters, initialized randomly at the start of each task based on a controlled random seed, dictate the initial state and influence reward outcomes. Similarly to MiniWoB++ (Shi et al., 2017; Liu et al., 2018a), ANDROIDWORLD consists of a practically infinite set of varying initial conditions and success criteria.

This approach provides more granular analyses of agents' adaptability — a vital attribute for real-world deployment. Beyond testing agent robustness, the dynamic construction of tasks supports the use of online learning methodologies, particularly reinforcement learning (Shi et al., 2017; Liu et al., 2018a; Humphreys et al., 2022; Gur et al., 2022a). It also simplifies the generation of distinct train/test datasets, facilitating supervised learning experiments (Humphreys et al., 2022; Shaw et al., 2023; Furuta et al., 2023).

3.4 DURABLE REWARDS FROM SYSTEM STATE

ANDROIDWORLD provides reward signals by managing application state using the Android Debug Bridge (adb). With the adb tool ANDROIDWORLD has complete access to system resources including the file system, application databases, and system settings. Determining reward signals

⁵<https://f-droid.org/>

324 from system state has several benefits. It is highly accurate because an application’s state can be
 325 quickly inspected and manipulated using the same mechanisms that the app itself utilizes. Using
 326 the underlying system state is much more durable than matching superficial UI changes. Addition-
 327 ally, it facilitates easy re-use across disparate apps, which tend to use the same underlying caching
 328 mechanisms. For instance, logic for checking existence of a specific file is used across many un-
 329 related applications, including those for file management, note-taking, and media playback. For
 330 applications leveraging SQLite databases, a common pattern, ANDROIDWORLD implements evalu-
 331 ators that verify the existence of new and deleted rows. Table 2 shows examples of the validators in
 332 ANDROIDWORLD. For more examples see Table 5.

333 3.5 TASK COMPOSABILITY

335 In addition to facilitating accurate and reusable evaluations, inferring a task’s success from system
 336 state makes it easy to create *composite* tasks by combining together existing tasks. For example, the
 337 task “Create a calendar event with details and text the details to contact” was created by combining
 338 together two existing tasks for creating a calendar event and for sending a text message, which is
 339 possible because each task initialization and success detection logic is hermetic. Composite tasks
 340 tend to be more challenging because of their complexity, although they provide partial rewards based
 341 on completion of sub tasks, to help facilitate hill climbing. The last two rows of Table 2 show the
 342 validation code for composite tasks.

343 3.6 INTEGRATING MINIWOB++

345 We implement MiniWoB++ in the ANDROIDWORLD framework and term it MobileMiniWoB++.
 346 Each MobileMiniWoB++ task is instantiated using the standard ANDROIDWORLD interface, in-
 347 heriting from `TaskEval` base class, and contains methods like `initialize_state` and
 348 `is_successful`. Since MiniWoB++ leverages JavaScript for task configuration and success de-
 349 tection, we built a WebView app to communicate between Python and the app. For instance, the
 350 `is_successful` method of each task retrieves the reward value from the WebView app via an
 351 Android intent.

352 MobileMiniWoB++ introduces modifications in both observations and actions compared to the origi-
 353 nal benchmark. For example, HTML5 `<input>` elements are rendered with native Android UI wid-
 354 gets like the date-picker (see Figure 4), enhancing the realism of the tasks. MobileMiniWoB++ uses
 355 the same observation space as the Android tasks (accessibility tree and screenshot). Notably, it does
 356 not include the DOM as in the original implementation. The action space from ANDROIDWORLD
 357 is retained. We manually review and test each task to ensure they are solvable. We excluded twelve
 358 of the original tasks that failed to render correctly on Android, presented compatibility issues with
 359 the touch interface, or required near real-time interaction, which poses challenges on emulators.
 360 Overall, ANDROIDWORLD supports 92 MiniWoB++ tasks. See Appendix B.3 for more details.

361 4 ANDROIDWORLD AS A COMPUTER-CONTROL BENCHMARK

362 To test ANDROIDWORLD’s applicability for autonomous agents, we develop and test a state-of-the-
 363 art agent and its variants across all 20 apps and 116 tasks, as well as on MobileMiniWoB++.

364 4.1 COMPUTER CONTROL AGENTS

365 4.1.1 M3A

366 We develop a multimodal autonomous agent for Android, M3A. It is zero-shot, integrating ReAct-
 367 style (Yao et al., 2022) and Reflexion-style (Shinn et al., 2023) prompting to consume user instruc-
 368 tions and screen content, reason, take actions, and update its decision-making based on the outcome
 369 of its actions.

370 In the first stage, M3A generates an action, represented in JSON, and reasoning for that action.
 371 To generate this output, the agent is provided with a list of available action types, guidelines for
 372 operating the phone, and a list of UI elements derived from the Android accessibility tree’s leaf
 373 nodes. The agent receives the current screenshot and a Set-of-Mark (SoM) (Yang et al., 2023a)

Table 3: Success Rates (SR) on ANDROIDWORLD and MobileMiniWoB++.

Agent	Input	Base model	SR _{ANDROIDWORLD}	SR _{MobileMiniWoB++}
<i>Human</i>	<i>screen</i>	<i>N/A</i>	<i>80.0</i>	<i>100.0</i>
M3A	al1y tree	GPT-4 Turbo	30.6	59.7
M3A	al1y tree	Gemini 1.5 Pro	19.4	57.4
M3A	SoM (screen + al1y tree)	GPT-4 Turbo	25.4	67.7
M3A	SoM (screen + al1y tree)	Gemini 1.5 Pro	22.8	40.3
SeeAct (Zheng et al., 2024a)	SoM (screen + al1y tree)	GPT-4 Turbo	15.5	66.1

annotated screenshot, which includes bounding boxes with numeric labels on the top-left corner for each UI element (see screenshot in Figure 5). The agent attempts to execute outputted action by referencing the specific mark (if applicable). In addition to the multimodal agent, we have developed a text-only variant that consumes the screen represented using the accessibility tree and selects the relevant action in JSON format.

After executing an action, M3A reflects on its effect by observing any state changes that may have occurred. During this stage, the agent is provided with available action types, general operating guidelines, the actual action taken, and its reasoning, as well as before-and-after UI states, represented by UI element representations and screenshots with SoM annotations. We request the agent to provide a concise summary of this step, including the intended action, success or failure, potential reasons for failure, and recommendations for subsequent actions. This summary will serve as the action history and be used for future action selection. See Appendix D for more details on the agent.

4.1.2 SEEACT BASELINE

We implement a baseline agent based on SeeAct (Zheng et al., 2024a), which was originally designed for GPT-4V for web navigation. Specifically, we implement the best-performing variant, SeeAct_{choice}, which grounds actions via textual choices. We implement SeeAct for the Android environment to evaluate how an existing model that performs well on web tasks (Deng et al., 2023) can be adapted and applied to Android.

To accommodate the Android environment, we adapt SeeAct in several ways. Firstly, we augment the action space from the original SeeAct implementation to support actions needed for mobile, including scroll, long press, navigate home and back, and open app actions. Secondly, in lieu of the DOM, which is not available for Android apps, we utilize the accessibility tree to construct candidate UI actions. Due to the lack of the DOM representation, we do not use the bespoke ranker model from the original implementation. However, we observe that after applying a filtering heuristic to remove non-interactable elements, the majority of screens contains less than 50 candidate elements. See Appendix D.5 for more details on the implementation.

4.2 EXPERIMENTAL RESULTS

We evaluate M3A and SeeAct on ANDROIDWORLD and MobileMiniWoB++. We set the seed to 30 and provide a task-specific step budget. We use Gemini 1.5 Pro and GPT-4 Turbo as base models. For MobileMiniWoB++, we evaluate on a subset of 62 tasks, consistent with recent studies (Zheng et al., 2024b; Kim et al., 2024; Gur et al., 2022b).

Table 3 presents the success rates (SR) for the agents and human performance on both task suites. Although the agents have far from human performance, they demonstrate out-of-the-box capabilities in operating mobile UIs, exhibiting basic understanding and control capabilities of UIs. They can perform a variety of actions, including long-press, scrolling to search for information, and revising their plan if actions do not work out. The best performance is obtained for M3A when using GPT-4. On ANDROIDWORLD the SoM-based variant is less performant, while on MobileMiniWoB++ it performs best. A similar result was obtained in recent work in the context of computer agents for desktop applications (Xie et al., 2024). We posit SoM plays a more critical role in MobileMiniWoB++ tasks due to the often incomplete accessibility tree, compared to that of native Android apps.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

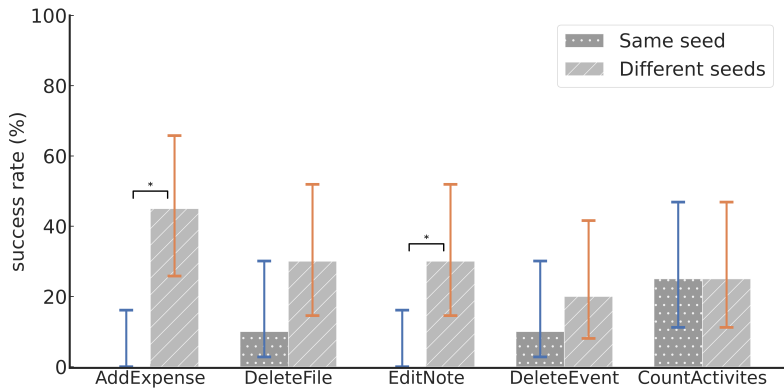


Figure 3: Success rate variation across tasks due to the parametrization built into ANDROIDWORLD. Using a fixed seed, the agent appears completely incapable of solving some tasks due to “bad luck” with the seed. In contrast, under different task parameterizations, we observe the agent is capable of solving the tasks fairly often. Wilson binomial proportion confidence intervals (95%) are shown for the different seed group (orange) and the same seed group (blue). The different seed group has higher variance than the same seed group. Significant differences, with p -value < 0.05 , are indicated by “*”.

4.3 ANALYSIS

Agents have difficulty understanding mobile UIs, often failing to detect visual cues that are essential for task completion (see Figure 6a). Additionally, agents struggle with certain UI patterns and affordances, and when they make reasoning mistakes (see Figure 6b), they often lack the capability to explore and adapt as humans do (see Figure 6c). Moreover, agents sometimes struggle with tasks that simply involve confirming system states, e.g., confirming the WiFi is turned on, suggesting challenges in both task and screen understanding.

The agents struggle with grounding, particularly when executing precise interactions, such as manipulating text (see Figure 7) or operating sliders, and they are often unable to recover from mistyping errors. In addition, for tasks that demand memory, such as performing transcriptions across apps, multiplying numbers, or scrolling, the agents struggle as they are unable to “remember” content.

SeeAct performs less effectively than M3A on the ANDROIDWORLD task suite and similarly on MobileMiniWoB++, reflecting its optimization for web rather than mobile environments. It struggles with mobile-specific actions like long-presses and swipes, and often fails to select appropriate actions due to not incorporating screen elements during action generation. Memory-intensive tasks are particularly challenging, as SeeAct only caches actions without remembering outcomes, leading to repetitive, ineffective behaviors such as endless scrolling. This lack of quick error recovery often results in task termination once maximum steps are reached.

Finally, we note that large foundation models significantly increase latency, taking three times longer than humans on average to complete tasks. On average, M3A takes 3.9 minutes to complete a task, with the text-only version taking 2.5 minutes.

4.4 AGENT ROBUSTNESS UNDER RANDOM SEEDS

We evaluate agent robustness under two conditions: (1) identical tasks with the same parameters and (2) tasks with different parameter combinations, which change the initial state and task definition. Due to computational constraints, we perform this analysis on a representative subset of ANDROIDWORLD tasks (listed in Appendix D.4). We use the strongest agent, M3A using the accessibility tree and GPT-4, for this analysis. Our results are shown in Figure 3.

In the baseline experiment with a constant seed, the agent fails the add and edit tasks, and rarely solves the two delete tasks. For the add and edit tasks, the agent struggles with UI operations, while for the delete tasks, it often gets confused before the step budget is consumed. Surprisingly, the

486 agent’s performance varies even with a fixed seed, suggesting the model’s non-determinism affects
487 agent reliability.

488 Agent performance varies much more using different seeds, with statistically significant differences
489 for the add expense and edit note tasks. The high intra-task variation indicates the model’s sensitivity
490 to the seed and its task parameters, stemming from a lack of robustness from the underlying model.
491

492 Research in RL environments (Henderson et al., 2018; Raffin et al., 2021; Colas et al., 2018) has
493 noted similar sensitivity to seeds. Consistent with these studies, we observe agent performance is
494 best represented by the mean across random seeds. Notably, the observation of non-zero rewards
495 under some seeds points to potential enhancements through RL-like mechanisms in future work.

496 This experiment underscores the importance of extensively testing agents under varied task param-
497 eters to ensure they can handle real-world variability, a capability that ANDROIDWORLD supports
498 effectively.

500 5 LIMITATIONS

502 ANDROIDWORLD currently supports tasks from open-source Android apps with at least 1 million
503 downloads and from built-in Android system apps. While testing on more trending apps is desirable,
504 we found open-source apps to be equally realistic and, in some cases, more challenging than apps
505 with larger user bases. Trending apps tend to have UIs which are heavily optimized for smooth user
506 experience, offering more UI functionality and shortcuts. Testing on less-optimized UIs makes the
507 agent’s task harder. In an example failure we show in Figure 6c, the agent needed to delete all notes
508 in a list and failed by repeatedly searching for a “delete-all” button. Instead, an agent with stronger
509 reasoning capabilities would have probably searched once for that functionality, realized it was not
510 available, and deleted the notes one by one.

512 6 CONCLUSION

514 We introduced ANDROIDWORLD, a realistic and robust agent environment for Android that enables
515 the development and evaluation of autonomous agents across a wide range of tasks and apps. AN-
516 DROIDWORLD provides a reproducible task suite consisting of 116 tasks across 20 apps, with each
517 task dynamically generated using random parameters to challenge agents with millions of unique
518 goals. By releasing ANDROIDWORLD and establishing benchmark performance with M3A, we aim
519 to accelerate research and development in this area, ultimately leading to the creation of computer
520 control agents capable of operating effectively in real-world environments. Further, the dynamic
521 nature of ANDROIDWORLD opens up new research opportunities for online learning algorithms in
522 computer control agents.

524 REFERENCES

526 AgentGPT, 2024. <https://agentgpt.reworkd.ai/>.

528 WebLlama, 2024. <https://webllama.github.io/>.

529 Josh Abramson, Arun Ahuja, Federico Carnevale, Petko Georgiev, Alex Goldin, Alden Hung, Jes-
530 sica Landon, Timothy Lillicrap, Alistair Muldal, Blake Richards, Adam Santoro, Tamara von
531 Glehn, Greg Wayne, Nathaniel Wong, and Chen Yan. Evaluating multimodal interactive agents,
532 2022.

534 William E Bishop, Alice Li, Christopher Rawles, and Oriana Riva. Latent state estimation helps ui
535 agents to reason, 2024.

537 Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong
538 Lu, Justin Wagle, Kazuhito Koishida, Arthur Bucker, Lawrence Jang, and Zack Hui. Windows
539 Agent Arena: Evaluating Multi-Modal OS Agents at Scale, 2024. URL <https://arxiv.org/abs/2409.08264>.

- 540 Andrea Burns, Deniz Arsan, Sanjna Agrawal, Ranjitha Kumar, Kate Saenko, and Bryan A. Plummer.
541 Mobile app tasks with iterative feedback (motif): Addressing task feasibility in interactive
542 visual environments. *CoRR*, abs/2104.08560, 2021. URL [https://arxiv.org/abs/
543 2104.08560](https://arxiv.org/abs/2104.08560).
- 544 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
545 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
546 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
547 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
548 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios
549 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex
550 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
551 Christopher Hesse, Andrew N Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,
552 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob
553 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating
554 large language models trained on code. July 2021.
- 555 Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong
556 Wu. Seeclck: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint
557 arXiv:2401.10935*, 2024.
- 558 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,
559 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena:
560 An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*,
561 2024.
- 562 Cédric Colas, Olivier Sigaud, and Pierre-Yves Oudeyer. How many random seeds? statistical power
563 analysis in deep reinforcement learning experiments. *arXiv preprint arXiv:1806.08295*, 2018.
- 564 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and
565 Yu Su. Mind2Web: Towards a generalist agent for the web, 2023.
- 566 Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H Laradji, Manuel Del Verme, Tom
567 Marty, Léo Boisvert, Megh Thakkar, Quentin Cappart, David Vazquez, et al. Workarena:
568 How capable are web agents at solving common knowledge work tasks? *arXiv preprint
569 arXiv:2403.07718*, 2024.
- 570 Yuqing Du, Ksenia Konyushkova, Misha Denil, Akhil Raju, Jessica Landon, Felix Hill, Nando
571 de Freitas, and Serkan Cabi. Vision-Language models as success detectors. March 2023.
- 572 Nicholas Farn and Richard Shin. Tooltalk: Evaluating tool-usage in a conversational setting. *arXiv
573 preprint arXiv:2311.10775*, 2023.
- 574 Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane
575 Gu, and Izzeddin Gur. Multimodal web navigation with Instruction-Finetuned foundation models.
576 May 2023.
- 577 Difei Gao, Lei Ji, Zechen Bai, Mingyu Ouyang, Peiran Li, Dongxing Mao, Qinchen Wu, Weichen
578 Zhang, Peiyi Wang, Xiangwu Guo, et al. Assistgui: Task-oriented desktop graphical user interface
579 automation. *arXiv preprint arXiv:2312.13108*, 2023.
- 580 Significant Gravitas. AutoGPT. <https://agpt.co>, 2023. <https://agpt.co>.
- 581 Izzeddin Gur, Natasha Jaques, Yingjie Miao, Jongwook Choi, Manoj Tiwari, Honglak Lee, and
582 Aleksandra Faust. Environment generation for zero-shot compositional reinforcement learning,
583 2022a.
- 584 Izzeddin Gur, Ofir Nachum, Yingjie Miao, Mustafa Safdari, Austin Huang, Aakanksha Chowdhery,
585 Sharan Narang, Noah Fiedel, and Aleksandra Faust. Understanding html with large language
586 models. *arXiv preprint arXiv:2210.03945*, 2022b.
- 587 Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan,
588 and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models.
589 *arXiv preprint arXiv:2401.13919*, 2024.

- 594 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
595 Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial*
596 *intelligence*, volume 32, 2018.
- 597 Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
598 Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. *arXiv*
599 *preprint arXiv:2312.08914*, 2023.
- 601 Peter C Humphreys, David Raposo, Tobias Pohlen, Gregory Thornton, Rachita Chhaparia, Alistair
602 Muldal, Josh Abramson, Petko Georgiev, Adam Santoro, and Timothy Lillicrap. A data-driven
603 approach for learning to control computers. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,
604 Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Con-*
605 *ference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp.
606 9466–9482. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/v162/](https://proceedings.mlr.press/v162/humphreys22a.html)
607 [humphreys22a.html](https://proceedings.mlr.press/v162/humphreys22a.html).
- 608 Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh,
609 and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist
610 autonomous agents for desktop and web. *arXiv preprint arXiv:2402.17553*, 2024.
- 611 Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks.
612 *Advances in Neural Information Processing Systems*, 36, 2024.
- 613 Megan Kinniment, Lucas Jun Koba Sato, Haoxing Du, Brian Goodrich, Max Hasin, Lawrence Chan,
614 Luke Harold Miles, Tao R Lin, Hjalmar Wijk, Joel Burget, et al. Evaluating language-model
615 agents on realistic autonomous tasks. *arXiv preprint arXiv:2312.11671*, 2023.
- 616 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham
617 Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating
618 multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*, 2024.
- 619 Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen
620 Zhang, Xiaohan Zhang, Yuxiao Dong, et al. Autowebglm: Bootstrap and reinforce a large lan-
621 guage model-based web navigating agent. *arXiv preprint arXiv:2404.03648*, 2024.
- 622 Juyong Lee, Taywon Min, Minyong An, Changyeon Kim, and Kimin Lee. Benchmarking mo-
623 bile device control agents across diverse configurations. In *ICLR 2024 Workshop on Generative*
624 *Models for Decision Making*, 2024.
- 625 Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei
626 Huang, and Yongbin Li. API-Bank: A comprehensive benchmark for Tool-Augmented LLMs.
627 April 2023a.
- 628 Tao Li, Gang Li, Zhiwei Deng, Bryan Wang, and Yang Li. A Zero-Shot language agent for computer
629 control with structured reflection. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings*
630 *of the Association for Computational Linguistics: EMNLP 2023*, pp. 11261–11274, Singapore,
631 December 2023b. Association for Computational Linguistics.
- 632 Wei Li, William Bishop, Alice Li, Chris Rawles, Folawiyo Campbell-Ajala, Divya Tyamagundlu,
633 and Oriana Riva. On the effects of data scale on computer control agents. In *Advances in Neural*
634 *Information Processing Systems (NeurIPS 2024)*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2406.03679)
635 [2406.03679](https://arxiv.org/abs/2406.03679).
- 636 Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. Mapping natural language
637 instructions to mobile UI action sequences. In *Proc. of the 58th Annual Meeting of the Associa-*
638 *tion for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pp. 8198–8210. Associa-
639 tion for Computational Linguistics, 2020. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/2020.acl-main.729/)
640 [2020.acl-main.729/](https://www.aclweb.org/anthology/2020.acl-main.729/).
- 641 Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, and Percy Liang. Reinforcement learning on
642 web interfaces using workflow-guided exploration. In *6th International Conference on Learn-*
643 *ing Representations (ICLR '18)*, 2018a. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=ryTp3f-0-)
644 [ryTp3f-0-](https://openreview.net/forum?id=ryTp3f-0-).

- 648 Thomas F. Liu, Mark Craft, Jason Situ, Ersin Yumer, Radomir Mech, and Ranjitha Kumar. Learning
649 design semantics for mobile apps. In *Proc. of the 31st Annual ACM Symposium on User Interface*
650 *Software and Technology*, UIST '18, pp. 569–579, New York, NY, USA, 2018b. Association for
651 Computing Machinery. ISBN 9781450359481. doi: 10.1145/3242587.3242650. URL <https://doi.org/10.1145/3242587.3242650>.
- 653 Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-Eval: NLG
654 evaluation using GPT-4 with better human alignment, 2023.
- 656 Xing Han Lù, Zdeněk Kasner, and Siva Reddy. WebLINX: Real-World website navigation with
657 Multi-Turn dialogue. February 2024.
- 658 Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayara-
659 man, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via
660 coding large language models, 2023.
- 662 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
663 Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, et al. Self-refine: Iterative refinement
664 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- 665 Grégoire Mialon, Clémentine Fourier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas
666 Scialom. GAIA: a benchmark for general AI assistants. November 2023.
- 667 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wier-
668 stra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.
- 670 Shikhar Murty, Christopher Manning, Peter Shaw, Mandar Joshi, and Kenton Lee. BAGEL: Boot-
671 strapping agents by guiding exploration with language, 2024. URL <https://arxiv.org/abs/2403.08140>.
- 673 Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. Autonomous
674 evaluation and refinement of digital agents. April 2024.
- 676 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dor-
677 mann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine*
678 *Learning Research*, 22(268):1–8, 2021.
- 679 Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the
680 wild: A large-scale dataset for android device control. *arXiv preprint arXiv:2307.10088*, 2023.
- 682 Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann
683 Dubois, Chris J Maddison, and Tatsunori Hashimoto. Identifying the risks of LM agents with
684 an LM-Emulated sandbox. September 2023.
- 685 Peter Shaw, Mandar Joshi, James Cohan, Jonathan Berant, Panupong Pasupat, Hexiang Hu, Urvashi
686 Khandelwal, Kenton Lee, and Kristina Toutanova. From pixels to UI actions: Learning to follow
687 instructions via graphical user interfaces. May 2023.
- 688 Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An
689 open-domain platform for web-based agents. In Doina Precup and Yee Whye Teh (eds.), *Proc. of*
690 *the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine*
691 *Learning Research*, pp. 3135–3144. PMLR, 06–11 Aug 2017. URL <http://proceedings.mlr.press/v70/shi17a.html>.
- 694 Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic
695 memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023.
- 696 Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi,
697 Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions
698 for everyday tasks, 2020.
- 700 Maayan Shvo, Zhiming Hu, Rodrigo Toro Icarte, Iqbal Mohamed, Allan Jepson, and Sheila A.
701 McIlraith. Appbuddy: Learning to accomplish tasks in mobile apps via reinforcement learning,
2021. URL <https://arxiv.org/abs/2106.00133>.

- 702 David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche,
703 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman,
704 Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine
705 Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the Game of Go
706 with Deep Neural Networks and Tree Search. *Nature*, 529(7587):484–489, January 2016. doi:
707 10.1038/nature16961.
- 708 Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia,
709 Jiechuan Jiang, Longtao Zheng, Xinrun Xu, Yifei Bi, Pengjie Gu, Xinrun Wang, Börje F. Karlsson,
710 Bo An, and Zongqing Lu. Towards General Computer Control: A Multimodal Agent For
711 Red Dead Redemption II As A Case Study. *arXiv preprint arXiv:2403.03186*, 2024.
- 712 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Bud-
713 den, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Ried-
714 miller. DeepMind control suite. January 2018.
- 715 Gemini Team. Gemini: A family of highly capable multimodal models, 2023a.
- 716 OpenAI Team. GPT-4 technical report, 2023b.
- 717 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
718 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
719 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 720 Daniel Toyama, Philippe Hamel, Anita Gergely, Gheorghe Comanici, Amelia Glaese, Zafarali
721 Ahmed, Tyler Jackson, Shibl Mourad, and Doina Precup. Androidenv: A reinforcement learning
722 platform for android, 2021. URL <https://arxiv.org/abs/2105.13231>.
- 723 Sagar Gubbi Venkatesh, Partha Talukdar, and Srinu Narayanan. Ugif: Ui grounded instruction fol-
724 lowing, 2022. URL <https://arxiv.org/abs/2211.07615>.
- 725 Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Juny-
726 oung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan
727 Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P Agapiou,
728 Max Jaderberg, Alexander S Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David
729 Budden, Yury Sulsky, James Molloy, Tom L Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff,
730 Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom
731 Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver.
732 Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782):
733 350–354, November 2019.
- 734 Bryan Wang, Gang Li, and Yang Li. Enabling conversational interaction with mobile ui using
735 large language models. In *Proc. of the 2023 CHI Conference on Human Factors in Computing
736 Systems*, CHI '23. Association for Computing Machinery, 2023a. ISBN 9781450394215. doi:
737 10.1145/3544548.3580895. URL <https://doi.org/10.1145/3544548.3580895>.
- 738 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
739 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models.
740 *arXiv preprint arXiv:2305.16291*, 2023b.
- 741 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun
742 Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and
743 Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework.
744 2023.
- 745 Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao
746 Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement.
747 *arXiv preprint arXiv:2402.07456*, 2024.
- 748 Tianbao Xie, Fan Zhou, Zhoujun Cheng, Peng Shi, Luoxuan Weng, Yitao Liu, Toh Jing Hua, Jun-
749 ning Zhao, Qian Liu, Che Liu, Leo Z. Liu, Yiheng Xu, Hongjin Su, Dongchan Shin, Caiming
750 Xiong, and Tao Yu. Openagents: An open platform for language agents in the wild, 2023.

- 756 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing
757 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal
758 agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*,
759 2024.
- 760 Mingzhe Xing, Rongkai Zhang, Hui Xue, Qi Chen, Fan Yang, and Zhen Xiao. Understanding the
761 weakness of large language model agents within a complex android environment. In *Proceedings*
762 *of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 6061–
763 6072, 2024.
- 764 Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. On the tool
765 manipulation capability of open-source large language models. May 2023.
- 766 An Yan, Zhengyuan Yang, Wanrong Zhu, Kevin Lin, Linjie Li, Jianfeng Wang, Jianwei Yang, Yiwu
767 Zhong, Julian McAuley, Jianfeng Gao, Zicheng Liu, and Lijuan Wang. GPT-4V in wonderland:
768 Large multimodal models for Zero-Shot smartphone GUI navigation. November 2023.
- 769 Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark
770 prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*,
771 2023a.
- 772 Zhao Yang, Jiakuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent:
773 Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023b.
- 774 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
775 ReAct: Synergizing reasoning and acting in language models. October 2022.
- 776 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable
777 real-world web interaction with grounded language agents, 2023.
- 778 Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei
779 Lin, Saravan Rajmohan, et al. Ufo: A ui-focused agent for windows os interaction. *arXiv preprint*
780 *arXiv:2402.07939*, 2024a.
- 781 Danyang Zhang, Zhennan Shen, Rui Xie, Situo Zhang, Tianbao Xie, Zihan Zhao, Siyuan Chen,
782 Lu Chen, Hongshen Xu, Ruisheng Cao, and Kai Yu. Mobile-env: Building qualified evaluation
783 benchmarks for llm-gui interaction, 2024b. URL [https://arxiv.org/abs/2305.](https://arxiv.org/abs/2305.08144)
784 08144.
- 785 Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and
786 Duyu Tang. Android in the zoo: Chain-of-action-thought for gui agents. *arXiv preprint*
787 *arXiv:2403.02713*, 2024c.
- 788 Li Zhang, Shihe Wang, Xianqing Jia, Zhihan Zheng, Yunhe Yan, Longxi Gao, Yuanchun Li, and
789 Mengwei Xu. Llamatouch: A faithful and scalable testbed for mobile ui task automation, 2024d.
790 URL <https://arxiv.org/abs/2404.16054>.
- 791 Zhuosheng Zhang and Aston Zhang. You only look at screens: Multimodal chain-of-action agents,
792 2023.
- 793 Zhuosheng Zhang, Yao Yao, Aston Zhang, Xiangru Tang, Xinbei Ma, Zhiwei He, Yiming Wang,
794 Mark Gerstein, Rui Wang, Gongshen Liu, et al. Igniting language intelligence: The hitchhiker’s
795 guide from chain-of-thought reasoning to language agents. *arXiv preprint arXiv:2311.11797*,
796 2023.
- 797 Ziniu Zhang, Shulin Tian, Liangyu Chen, and Ziwei Liu. MMInA: Benchmarking multihop multi-
798 modal internet agents. *arXiv preprint arXiv:2404.09992*, 2024e.
- 799 Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web
800 agent, if grounded. *arXiv preprint arXiv:2401.01614*, 2024a.
- 801 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
802 Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica.
803 Judging LLM-as-a-Judge with MT-Bench and chatbot arena. June 2023.

810 Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar
811 prompting with memory for computer control, 2024b.
812

813 Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
814 Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic
815 web environment for building autonomous agents, 2023.
816

817 APPENDIX A ETHICAL CONSIDERATIONS

818 **Malicious use** There is a risk that malicious actors could engineer agents to bypass security mea-
819 sures like CAPTCHAs or engage in activities like spamming. Additionally, they could alter prompts
820 or screen outputs to further harmful objectives.
821

822 **Societal impact** Automation agents may transform societal norms, disrupt employment, and mod-
823 ify human behavior. While they can enhance efficiency, this improvement could pose risks if ex-
824 ploited by malevolent forces.
825

826 APPENDIX B ANDROIDWORLD ENVIRONMENT

827 B.1 OBSERVATION SPACE

828 In ANDROIDWORLD, the Android screen is represented using a `State` class, which includes the
829 following attributes:
830

- 831 • **Pixels:** An RGB array representing the current screen capture of the device. The screenshot
832 resolution is $2400 \times 1080 \times 3$.
- 833 • **Accessibility tree:** A raw representation of the accessibility tree.⁶ This UI tree provides a
834 detailed snapshot of all UI elements currently displayed on the screen. We utilize an acces-
835 sibility forwarding app from AndroidEnv (Toyama et al., 2021), which leverages gRPC to
836 transmit the accessibility tree data efficiently to the device.
- 837 • **UI elements:** A list of processed UI elements extracted from the children of the accessibil-
838 ity tree. Each `UIElement` contains attributes such as text, content description, bounding
839 boxes, and various state flags (e.g., clickable, scrollable, focused).
840

841 Since Android observations and actions are asynchronous, changes resulting from actions may take
842 some time to manifest. Therefore, instead of using an RL-based interface, which assumes a tight
843 coupling between actions and observations, we design an interface for the agent tailored for asyn-
844 chronous interaction. This interface implements a `get_state` method responsible for capturing
845 the current state of the environment, typically after executing an action. This method includes an
846 optional `wait_to_stabilize` flag, which, when enabled, employs heuristics to ensure the UI el-
847 ements are not in a transient state, thus providing a stable and accurate snapshot of the environment.
848

849 B.2 ACTION SPACE

850 Actions are stored using a Python dataclass (shown below) and are executed using `adb`. Each action
851 type corresponds to specific ADB commands that interact with the Android environment. For click-
852 based actions (click, double tap, long press), we use ADB to simulate touch events at specified
853 coordinates on the screen. For text input actions, we first focus on the text input field and then use
854 ADB to type the desired text and press the enter button. Navigation actions (home, back) involve
855 sending corresponding key events to the device. Scrolling and swiping actions, which are essentially
856 inverse operations, are both implemented by generating and issuing swipe commands through ADB
857 to simulate these gestures. To launch applications, we use ADB to start the desired app.
858

859 ⁶Represented using all current windows; [https://developer.android.com/reference/
860 android/view/accessibility/AccessibilityWindowInfo](https://developer.android.com/reference/android/view/accessibility/AccessibilityWindowInfo)
861


```

864 1 ACTION_TYPES = {
865 2     "CLICK": "click",
866 3     "DOUBLE_TAP": "double_tap",
867 4     "SCROLL": "scroll",
868 5     "SWIPE": "swipe",
869 6     "INPUT_TEXT": "input_text",
870 7     "NAVIGATE_HOME": "navigate_home",
871 8     "NAVIGATE_BACK": "navigate_back",
872 9     "KEYBOARD_ENTER": "keyboard_enter",
873 10    "OPEN_APP": "open_app",
874 11    "STATUS": "status",
875 12    "WAIT": "wait",
876 13    "LONG_PRESS": "long_press",
877 14    "ANSWER": "answer",
878 15    "UNKNOWN": "unknown"
879 16 }
880 17
881 18 @dataclasses.dataclass()
882 19 class JSONAction:
883 20     """Represents a parsed JSON action.
884 21
885 22     # Example
886 23     result_json = {'action_type': 'click', 'x': %d, 'y': %d}
887 24     action = JSONAction(**result_json)
888 25
889 26     Attributes:
890 27     action_type: The action type.
891 28     index: The index to click, if action is a click. Either an index or a <x, y>
892 29     should be provided. See x, y attributes below.
893 30     x: The x position to click, if the action is a click.
894 31     y: The y position to click, if the action is a click.
895 32     text: The text to type, if action is type.
896 33     direction: The direction to scroll, if action is scroll.
897 34     goal_status: If the status is a 'status' type, indicates the status of the goal.
898 35     app_name: The app name to launch, if the action type is 'open_app'.
899 36     """
900 37     action_type: str
901 38     index: int = None
902 39     x: int = None
903 40     y: int = None
904 41     text: str = None
905 42     direction: str = None
906 43     goal_status: str = None
907 44     app_name: str = None

```

Listing 1: Pseudo-code representation of the action space.

B.3 MOBILEMINIWOB++

Authors manually completed all tasks in MiniWoB++ to verify solvability on a mobile interface. MobileMiniWoB++ differs from MiniWoB++ due to the touch-based interface, which required different approaches for certain tasks. For instance, highlighting text from the `highlight-text` tasks involves using Android’s long-press and cursor-moving functionalities. HTML5 `<input>` elements are natively rendered with native Android UI widgets like the date-picker (see Figure 4).

Our implementation of MiniWoB++ contains 92 tasks in total. We exclude the following tasks: `chase-circle` (requires near-realtime movement, unachievable by humans on emulators), `moving-items` (too hard to click in emulator), `drag-cube` (drags will scroll the screen, moving the task out of view), `drag-items-grid` (elements are not interactable on Android), `drag-items` (elements are not interactable on Android), `drag-shapes` (drags will scroll the screen, moving the task out of view), `drag-sort-numbers` (elements are not interactable on Android), `text-editor` (cannot underline everything, weird glitch), `number-checkboxes` (not correctly rendered: only three columns), `use-slider-2` (slider implementation not working), `use-spinner` (slider implementation not working), and `click-menu` (the menu responsiveness breaks and the task does not behave as intended).

APPENDIX C ANDROIDWORLD BENCHMARK DETAILS

C.1 APP SELECTION

Our selection of apps (summarized in Table 4) was guided by three main factors: use case, popularity, and the need for consistency and reproducibility.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937

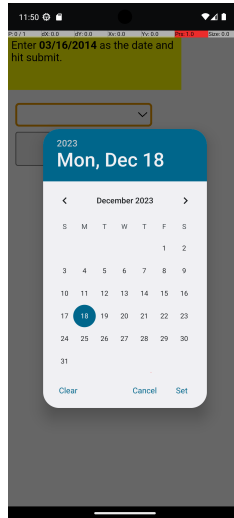


Figure 4: Native Android UI widget rendering for HTML5 `<input>` element.

938
939
940
941
942
943
944
945
946
947
948
949
950
951

Use case and categories We analyzed popular app categories in app stores, focusing on productivity, communication, and multimedia. Selected apps had to meet criteria such as not requiring a login and storing data locally on the device. Additionally, we considered apps from categories that the authors commonly used, ensuring the selection was representative of real-world Android usage.

Popularity We used download statistics from the Google Play Store to gauge app popularity, selecting apps with over 1 million downloads. Most of the selected apps exceeded this threshold. Less popular apps were also included if they featured common UI patterns and affordances, ensuring they are indicative of typical Android app usage. For instance, Simple Calendar Pro, though less downloaded, has a UI comparable to the widely-used Google Calendar app.

Consistency and reproducibility All apps were sourced from F-Droid, an open-source Android app repository. This allowed us to manage app versions precisely by selecting and distributing specific APKs. We use the newest version of each app at the time of download.

952
953

C.2 TASK CLASSIFICATION AND GENERATION

954
955
956
957

We categorize tasks into two types: those with side-effects and those without. Tasks with side-effects are those that modify the internal state of the device or applications, such as turning off Wi-Fi or creating a calendar event. These tasks are implemented as distinct Python classes, each with its own parameter generation, initialization, evaluation, and teardown methods.

958
959
960
961
962

Below we show an example of the task evaluation for a `SendSms` task, which involves sending and validating a text message. The pseudocode illustrates the task initialization, success check, and parameter generation methods. Each task has its own random parameter generation method and success logic.

963
964
965
966
967
968
969
970
971
972

```

1 class SendSms(TaskEval):
2     """Task sending and validating a text message has been sent.
3
4     It checks the SMS telephony database, which is located at:
5     /data/data/com.android.providers.telephony/databases/mmssms.db."""
6
7     template = (
8         "Send a text message using Simple SMS Messenger to "
9         "{number} with message: {message}"
10    )
11
12    def initialize_task(self, env: interface.AsyncEnv) -> None:
13        """Sets up the initial state of the task."""
14        super().initialize_task(env)
15        clear_sms_database(env.base_env)
16
17    def is_successful(self, env: interface.AsyncEnv) -> float:

```

Table 4: List of ANDROIDWORLD apps and number of tasks for each one.

App name	Description	# tasks
Simple Calendar Pro	A calendar app for creating, deleting, and managing events and appointments.	17
Settings	The Android system settings app for managing device settings such as Bluetooth, Wi-Fi, and brightness.	15
Markor	A note-taking app for creating, editing, deleting, and managing notes and folders.	14
Broccoli - Recipe App	A recipe management app for adding, deleting, and organizing recipes.	13
Pro Expense	An expense tracking app for adding, deleting, and managing expenses.	9
Simple SMS Messenger	An SMS app for sending, replying to, and resending text messages.	7
OpenTracks	A sport tracking app for recording and analyzing activities, durations, and distances.	6
Tasks	A task management app for tracking tasks, due dates, and priorities.	6
Clock	An app with stopwatch and timer functionality.	4
Joplin	A note-taking app.	4
Retro Music	A music player app.	4
Simple Gallery Pro	An app for viewing images.	4
Camera	An app for taking photos and videos.	3
Chrome	A web browser app.	3
Contacts	An app for managing contact information.	3
OsmAnd	A maps and navigation app with support for adding location markers, favorites, and saving tracks.	3
VLC	A media player app for playing media files.	3
Audio Recorder	An app for recording and saving audio clips.	2
Files	A file manager app for the Android filesystem, used for deleting and moving files.	2
Simple Draw Pro	A drawing app for creating and saving drawings.	1

```

18 """Checks if the SMS was sent successfully."""
19 super().is_successful(env)
20 messages = get_messages(env.base_env)
21 return check_message_exists(
22     phone_number=self.params["number"],
23     body=self.params["message"],
24 )
25
26 def teardown(self, env: interface.AsyncEnv) -> None:
27     """Clears the SMS database."""
28     super().teardown(env)
29     clear_sms_database(env.base_env)
30
31 @classmethod
32 def generate_random_params(cls) -> dict[str, Any]:
33     number = generate_random_number()
34     message = generate_random_message()
35     return {
36         "number": number,
37         "message": message,
38     }

```

C.3 INFORMATION RETRIEVAL TASKS

Tasks without side-effects are Information Retrieval tasks, requiring the agent to answer a question based on the device or app’s current state. For these tasks, instead of a Python class, we create a protobuf structure to specify the prompt, parameter values, and initialization and validation logic. We decided to use a structured data format with the belief that it would allow us to define new information retrieval tasks by simply adding new entries, making it easier to scale up the number of tasks without needing to write and maintain Python classes for each one.

Initialization is defined per app, including only the state relevant to the prompt’s answer and exclusion conditions for generating random states. This ensures that no random state contains information that could alter the expected answer. The initial state and prompt are parameterized using random values from the specified task parameters. For validation, we define the expected answer format within the prompt and use a few supported functions (“count”, “sum”, “identity”) to generate the answer from the initial state.

1026 Once an app and its specific logic are programmed, new tasks can be generated using an LLM to
 1027 generate the task’s protobuf. The process is not automatic and requires human review. Common
 1028 issues with LLM-generated tasks include missing fields, hallucinated fields, incompatible param-
 1029 eter generation, insufficient parameter usage, and non-specific task prompts. We observed that the
 1030 complexity of the proto structure correlates with an increase in generated task issues. Despite these
 1031 challenges, we found that editing LLM-generated protobufs can be more efficient than writing a
 1032 complete task from scratch.

1033 Below we show a simplified version of the task definition for the
 1034 SimpleCalendarEventsOnDate task which involves checking which events are on a
 1035 certain date. It specifies the relevant event, the exclusion conditions for any noisy event, how to
 1036 determine success, and possible parameter values to be chosen at random that will be used to fill
 1037 out the task definition.

```

1038 1 tasks {
1039 2   name: "SimpleCalendarEventsOnDate"
1040 3   prompt: "What events do I have {date} in Simple Calendar Pro? Answer with the titles only. If there are
      multiple titles, format your answer as a comma separated list."
1041 4   complexity: 1
1042 5   relevant_state {
1043 6     // Defines information for the goal events.
1044 7     state: {
1045 8       calendar {
1046 9         events {
1047 10          start_date: "{date}"
1048 11          start_time: "{time}"
1049 12          duration: "{duration}"
1050 13          title: "{title}"
1051 14        }
1052 15      }
1053 16    }
1054 17    // Non-goal events.
1055 18    exclusion_conditions {
1056 19      field: "start_date"
1057 20      operation: EQUAL_TO
1058 21      value: "{date}"
1059 22    }
1060 23  }
1061 24  success_criteria {
1062 25    expectations {
1063 26      field_transformation {
1064 27        operation: IDENTITY
1065 28        field_name: "title"
1066 29      }
1067 30      match_type: STRING_MATCH
1068 31    }
1069 32  }
1070 33  }
1071 34  task_params {
1072 35    name: "time"
1073 36    possible_values: "11:00am"
1074 37    // ...
1075 38  }
1076 39  }
1077 40  task_params {
1078 41    name: "date"
1079 42    possible_values: "October 15 2023"
1080 43    //...
1081 44  }
1082 45  }
1083 46  task_params {
1084 47    name: "duration"
1085 48    possible_values: "30 m"
1086 49    // ...
1087 50  }
1088 51  }
1089 52  task_params {
1090 53    name: "title"
1091 54    possible_values: "Data Dive"
1092 55    // ...
1093 56  }
1094 57  }

```

1071

1072

1073 C.4 HUMANS FOR TASK ANALYSIS

1074

1075 During development, we recruited six volunteers with proficient programming skills to analyze task
 1076 difficulty, duration, and category. Each human was assigned an equal portion of tasks and tasked
 1077 with identifying bugs during this annotation phase. This process resulted in the discovery and reso-
 1078 lution of over 30 bugs.

1079 To evaluate human performance, we enlisted two software engineers to complete the tasks using an
 Android emulator. Participants were provided with task descriptions and attempted to achieve the

goals based on their interpretations. Each participant had one attempt per task. The majority of errors stemmed from misinterpretations or minor errors, such as entering an incorrect file extension. Other errors occurred when participants encountered unfamiliar user interfaces, impeding their ability to solve the tasks on their first attempt.

In both exercises, we informed participants about the intended use of the collected data. Participants were not required to enter any personal information in the tested tasks.

C.5 TASK EXAMPLES

Table 5 lists some additional examples of tasks and highlights which task attributes can be parameterized in unlimited ways.

Table 5: Examples of ANDROIDWORLD tasks. We list the task nickname, the task template indicating which task attributes can be parameterized, the initialization logic that is executed before the task starts and pseudo code describing the success evaluation.

Task nickname	Task template	Initialization logic	Success evaluation code
VlcCreatePlaylist	Create a playlist in VLC, titled "{playlist_name}" with the following files, in order: {files}	Create new mpeg files: files + "noise" files that should not be added. Add them to VLC videos folder.	<code>execute_sql(vlc_query) == files</code>
RecipeAddMultiple RecipesFromImage	Add the recipes from recipes.jpg in Simple Gallery Pro to the recipe app.	Write a receipt file with recipes to Simple Gallery.	<code>sql.rows.exist(expected_recipes)</code>
MarkorEditNote	Edit {file_name} in Markor. {file_operation}.	Generate file with starting content, along with "noise" files not relevant to goal. <i>Note: file_operation can be to add a footer, header, or update note content.</i>	<code>file.exists(file_name, content=expected_content)</code>
ExpenseAddSingle	Add the following expenses into pro expense: {expense_csv}	Add to the app's SQLite database the expense that should be deleted, along with "noise" expenses that should not be deleted.	<code>sql.rows.exist(expense_obj)</code>
SimpleCalendarDelete EventsOnRelativeDay	In Simple Calendar Pro, delete all events scheduled for this {day_of_week}.	add to the app's SQLite database calendar events on specified day, along with "noise" events that should not be deleted.	<code>!sql.rows.exist(expected_events)</code>
FilesDeleteFile	Delete the file {file_name} from the Android filesystem located in the {subfolder} folder within the sdk_gphone_x86_64 storage area.	Generate specified file, along with "noise" files that should not be deleted.	<code>!file.exists(file_name)</code>
SportsTrackerActivities CountForWeek	How many {category} activities did I do this week in the OpenTracks app? Express your answer as a single integer.	add to the app's SQLite database activities for the specified category, along with "noise" activities.	<code>int(agent_response) == expected_count</code>

APPENDIX D ANDROIDWORLD AGENT DETAILS

D.1 M3A OBSERVATIONS

ANDROIDWORLD consumes the raw screen pixels, the screen shot with Set-of-Mark (SoM) (Yang et al., 2023a) annotations, and a list of UI elements on screen.

```

1 Here is a list of descriptions for some UI elements on the current screen:
2
1128 UIElement0: UIElement(text="VLC", content_description=None, class_name="android.widget.EditText",
1129 bbox_pixels=BoundingBox(x_min=98, x_max=886, y_min=146, y_max=311), ...)
1130 UIElement1: UIElement(text=None, content_description="Clear search box", class_name="android.widget.
    ImageButton",
1131 bbox_pixels=BoundingBox(x_min=886, x_max=1023, y_min=160, y_max=297), ...)
1132 UIElement2: UIElement(text="15:11", content_description="15:11", class_name="android.widget.TextView",
1133 bbox_pixels=BoundingBox(x_min=50, x_max=148, y_min=1, y_max=128), ...)
9 ... More elements listed ...
10
11 ... Guidelines on action selection emitted ...
12
```

1134 13 Now output an action from the above list in the correct JSON format, following the reason why you do that.
 1135 Your answer should look like:
 1136 14
 1137 15 Reason: ...
 1138 16 Action: {"action_type":...}

Listing 2: The prompt format pertaining to screen representation with UI elements.

1141 D.2 M3A ACTIONS

1142 For the SoM prompting, the screen is annotated based on the UI elements extracted from the acces-
 1143 sibility tree, which form the agent’s action space. Figure 5 shows one example.
 1144

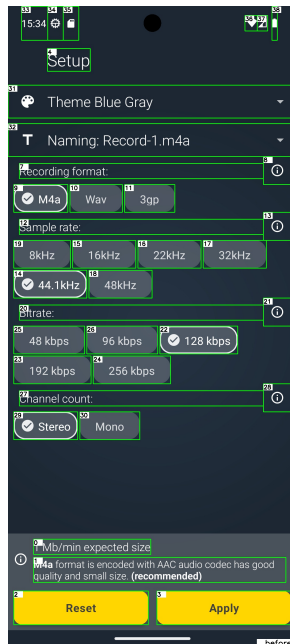


Figure 5: Set-of-marks overlaid on an Android screen.

1169 D.3 ERROR ANALYSIS

1170 We analyze M3A errors based on broader categories we observe during evaluation.

1171 **Perceptual errors** Perceptual errors are caused when the model fails to recognize crucial elements
 1172 on the screen necessary for task completion.

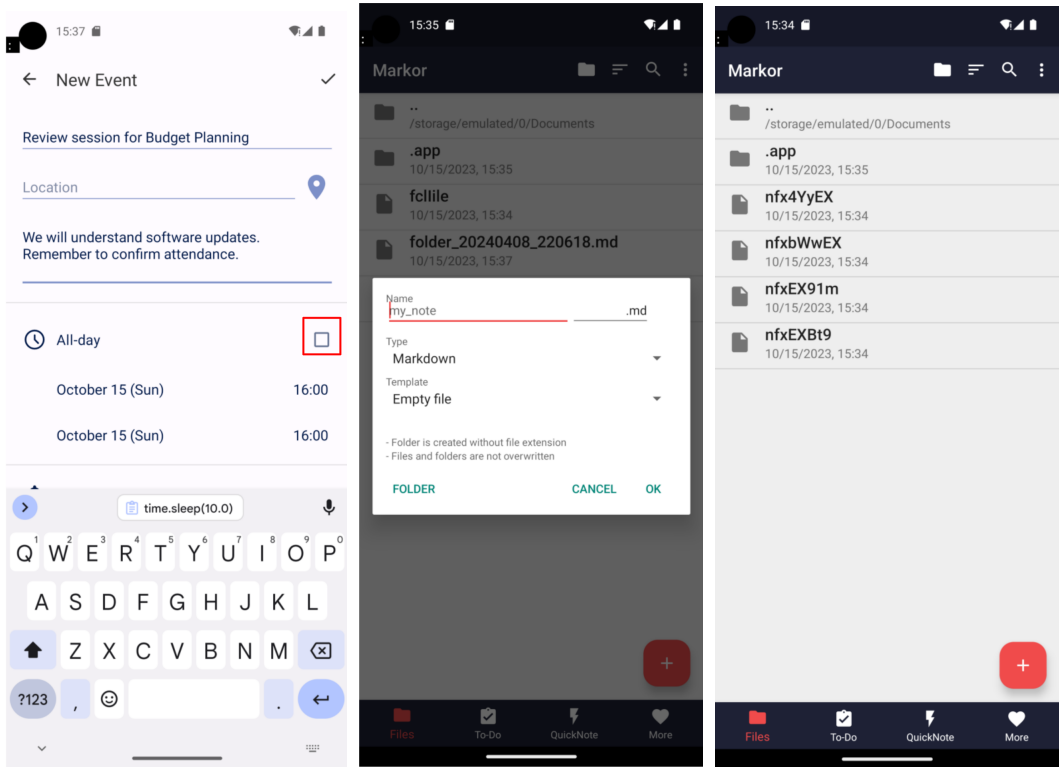
1173 For the task below, the model does not recognize that the “All-day” checkbox is currently not
 1174 checked (see Figure 6a):

1175 *In Simple Calendar Pro, create a recurring calendar event titled 'Review session
 1176 for Budget Planning' starting on 2023-10-15 at 14h. The event recurs weekly,
 1177 forever, and lasts for 60 minutes each occurrence. The event description should
 1178 be 'We will understand software updates. Remember to confirm attendance.'*

1183 **Reasoning errors** Reasoning errors occur when the model misinterprets the task requirements or
 1184 the current state, leading to incorrect actions.

1185 For example, in the task below, the model mistakenly believes the note name has already been
 1186 entered, so it types the note text into the “Name” field and cannot recover from this error (see
 1187 Figure 6b):

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241



(a) Perceptual error. Red square highlights issue.
 (b) Reasoning error. The agent's next action is to start entering the note's contents, which is incorrect because it needs to enter the note's name first.
 (c) Missing knowledge error. To delete all notes, the agent mistakenly looks for an option to delete all the notes at once, rather than trying to do it note-by-note.

Figure 6: Screens illustrating perceptual, reasoning and missing knowledge errors.

Create a new note in Markor named ipsum.mNcKd36pDq.txt with the following text: ...

Here is the agent's reasoning trace:

The user wants to create a new note with a specific name and text. The name of the note has already been entered, so all that needs to be done is to enter the text and click on the OK button.

Missing knowledge errors Missing knowledge errors occur when the model lacks the necessary understanding of the application's UI to complete a task efficiently.

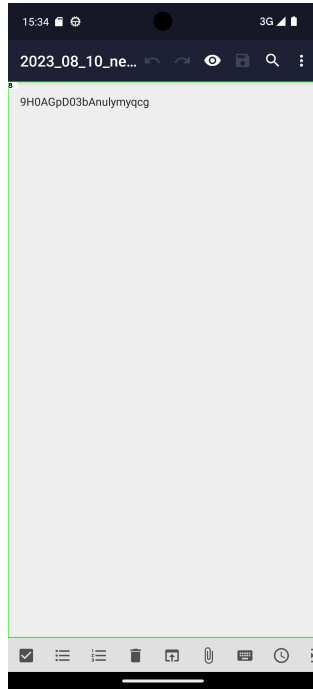
For example, in the task below, the agent does not know how to delete all files at once. It looks for an option to do that instead of deleting each file one by one (see Figure 6c):

Delete all files in the current directory.

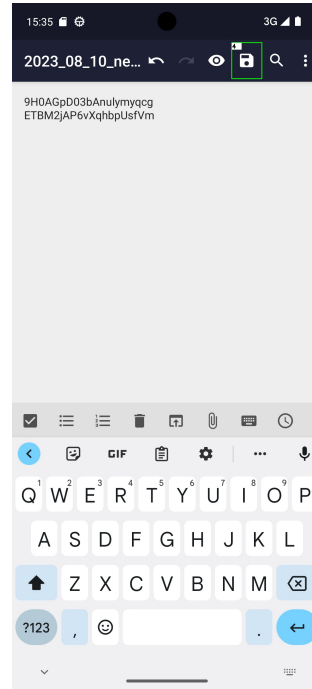
Grounding errors Grounding errors occur when the model fails to correctly interact with the UI elements based on their spatial or contextual positioning.

For the task below, the agent needs to update the Markor note by prepending text to the existing text. Figure 7 illustrates the errors the agent makes. It clicks the entire text field area, highlighted in green, which automatically places the cursor after the current text, resulting in the new text being appended after the current content.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295



(a) Error on initial click.



(b) Error with text entered and saving.

Figure 7: Screens illustrating grounding errors.

Update the Markor note ‘2023_08_10_neat_wolf.txt’ by adding the following text, along with a new blank line before the existing content: ”ETBM2jAP6vXqhbpuSfVm”, and rename it to ‘sure_ocean_uRnI.txt’.

Then, in the next screen, the text has been entered after the existing content, and the agent clicks the save button.

D.4 AGENT ROBUSTNESS EXPERIMENTS

We ran the agent on the following tasks (the nicknames shown in the figures in parentheses):

- MarkorEditNote (EditNote)
- ExpenseAddSingle (AddExpense)
- SimpleCalendarDeleteEventsOnRelativeDay (DeleteEvent)
- FilesDeleteFile (DeleteFile)
- SportsTrackerActivitiesCountForWeek (CountActivities)

More details about these tasks can be found in Table 5.

D.5 SEEACT DETAILS

We modify the SeeAct prompt (Zheng et al., 2024a) to reflect that the environment is Android by inputting elements from the accessibility tree and supporting additional actions (e.g., scrolling). Below we include the updated prompt. We annotate the system, user, and assistant roles that are each provided to the OpenAI API.

1
2 > Role: SYSTEM

1296
1297 3 Imagine that you are imitating humans operating an Android device for a task step by step. At each stage, you
1298 can see the Android screen like humans by a screenshot and know the previous actions before the current
1299 step decided by yourself through recorded history. You need to decide on the first following action to
1300 take. You can tap on an element, long-press an element, swipe, input text, open an app, or use the
1301 keyboard enter, home, or back key. (For your understanding, they are like 'adb shell input tap', 'adb
1302 shell input swipe', 'adb shell input text', 'adb shell am start -n', and 'adb shell input keyevent').
1303 One next step means one operation within these actions. Unlike humans, for typing (e.g., in text areas,
1304 text boxes), you should try directly typing the input or selecting the choice, bypassing the need for an
1305 initial click. You should not attempt to create accounts, log in or do the final submission. Terminate
1306 when you deem the task complete or if it requires potentially harmful actions.

4
1303 5 > Role: USER
1304 6 You are asked to complete the following task: <GOAL>

7
1305 8 Previous Actions:
1306 9 <PREVIOUS ACTIONS>

10
1307 11 The screenshot below shows the Android screen you see. Follow the following guidance to think step by step
1308 before outlining the next action step at the current stage:

12
1308 13 (Current Screen Identification)
1309 14 Firstly, think about what the current screen is.

15
1310 16 (Previous Action Analysis)
1311 17 Secondly, combined with the screenshot, analyze each step of the previous action history and their intention
1312 one by one. Particularly, pay more attention to the last step, which may be more related to what you
1313 should do now as the next step. Specifically, if the last action involved a INPUT TEXT, always evaluate
1314 whether it necessitates a confirmation step, because typically a single INPUT TEXT action does not make
1315 effect. (often, simply pressing 'Enter', assuming the default element involved in the last action,
1316 unless other clear elements are present for operation).

18
1315 19 (Screenshot Details Analysis)
1316 20 Closely examine the screenshot to check the status of every part of the screen to understand what you can
1317 operate with and what has been set or completed. You should closely examine the screenshot details to
1318 see what steps have been completed by previous actions even though you are given the textual previous
1319 actions. Because the textual history may not clearly and sufficiently record some effects of previous
1320 actions, you should closely evaluate the status of every part of the screen to understand what you have
1321 done.

21
1319 22 (Next Action Based on Android screen and Analysis)
1320 23 Then, based on your analysis, in conjunction with human phone operation habits and the logic of app design,
1321 decide on the following action. And clearly outline which element on the Android screen users will
1322 operate with as the first next target element, its detailed location, and the corresponding operation.

24
1322 25 To be successful, it is important to follow the following rules:
1323 26 1. You should only issue a valid action given the current observation.
1324 27 2. You should only issue one action at a time
1325 28 3. For handling the select dropdown elements on a screen, it's not necessary for you to provide completely
1326 accurate options right now. The full list of options for these elements will be supplied later.

29
1326 30 > Role: ASSISTANT
1327 31 <AGENT RESPONSE TO ABOVE>

32
1328 33 > Role: USER
1329 34 (Reiteration)
1330 35 First, reiterate your next target element, its detailed location, and the corresponding operation.

36
1330 37 (Multichoice Question)
1331 38 Below is a multi-choice question, where the choices are elements on the screen. All elements are arranged in
1332 the order based on their height on the screen, from top to bottom (and from left to right). This
1333 arrangement can be used to locate them. From the screenshot, find out where and what each one is on the
1334 screen, taking into account both their text content and details. Then, determine whether one matches
1335 your target element. Please examine the choices one by one. Choose the matching one. If multiple options
1336 match your answer, choose the most likely one by re-examining the screenshot, the choices, and your
1337 further reasoning. If you would like to perform a swipe action, you can optionally select the choice
1338 where you will swipe.

39
1336 40 A. "Home" icon
1337 41 B. "Phone" icon
1338 42 C. "Messages" icon
1339 43 D. "Chrome" icon
1340 44 E. "Search" icon
1341 45 ...

46
1339 46 If none of these elements match your target element, please select Z. None of the other options match the
1340 correct element.

47
1341 48 (Final Answer)
1342 49 Finally, conclude your answer using the format below. Ensure your answer is strictly adhering to the format
1343 provided below. Please do not leave any explanation in your answers of the final standardized format
1344 part, and this final part should be clear and certain. The element choice, action, and value should be
1345 in three separate lines.

50
1344 51 Format:
1345 52
1346 53 ELEMENT: The uppercase letter of your choice. (No need for TERMINATE, KEYBOARD ENTER, WAIT, ANSWER, OPEN APP,
1347 NAVIGATE HOME, NAVIGATE BACK; and optional for SWIPE.)

54
1347 54 ACTION: Choose an action from {CLICK, INPUT TEXT, LONG PRESS, NAVIGATE BACK, TERMINATE, KEYBOARD ENTER, SWIPE,
1348 WAIT, ANSWER, OPEN APP, NAVIGATE HOME}.

55
1348 55
1349 56
1349 57 VALUE: Provide additional input based on ACTION.
1350 58

1350
59 The VALUE means:
1351 60 If ACTION == INPUT TEXT, specify the text to be typed.
61 If ACTION == SWIPE, specify the direction: up, down, left, right.
1352 62 If ACTION == OPEN APP, provide the name of the app to be opened.
1353 63 If ACTION == ANSWER, specify the text of your answer to respond directly to a question or request for
information.
1354 64 For CLICK, LONG PRESS, KEYBOARD ENTER, NAVIGATE HOME, NAVIGATE BACK, WAIT, and TERMINATE, write "None".
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403