

---

# MYPCBENCH: A Benchmark for Personally Intelligent Computer-Use Agents

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Current benchmarks for computer-use agents evaluate models in impersonal environments. This leaves a gap between evaluation and deployment, since personal assistants are expected to work across a user’s whole digital life, including  
2       their context, historical data, and logged-in accounts. The gap is widest on web  
3       tasks, where live-web evaluations cannot exercise sites that require logging in or  
4       personal information, the kind of site a real personal assistant has to drive. We  
5       introduce MYPCBENCH, which tests computer-use agents as personal assistants  
6       on a Linux desktop populated with 17 simulated real-world web applications and  
7       a full desktop stack, all seeded for one canonical persona, Michael Scott from  
8       *The Office*<sup>1</sup>. We define 184 tasks in this environment, each inspired by a real  
9       request drawn from the OpenClaw community, and benchmark six frontier and  
10      open-weight models under each provider’s native CUA agent. Claude Opus 4.6  
11      reaches 55.4% perfect, the only model above 50%; the gap to other models concentrates on tasks that span many applications and on long trajectories, the regimes  
12      where personalization stresses an assistant the most. We release the environment,  
13      the task set with rubrics, the agent harness, and the rubric-grading judge at [TBD].  
14  
15  
16

## 17 1 Introduction

18    A person’s computer is not a blank slate. Bank transactions, calendar events, email threads, travel  
19    bookings, and work chats accumulate across many applications, together forming the record of a  
20    user’s working and personal life. Current benchmarks for computer-use agents ignore this. Tasks  
21    run against empty desktops, generic application states, and minimally seeded databases. In most  
22    tasks, the agent is told exactly which application to open and the exact workflow to complete, with  
23    no realistic user data behind the application. An agent that can place a delivery order but cannot find  
24    which restaurant the user actually orders from every Friday has not demonstrated useful capabilities  
25    as a personal assistant. As LLM-based assistants for personal computers move from research demos  
26    to consumer products (e.g. OpenClaw [16] and Claude CoWork [2]), evaluation has to keep up: it  
27    should test whether these systems are actually personal, and whether personalization improves or  
28    regresses with each new model release.

29    Existing agent benchmarks span the web [26, 11, 5, 7], full desktops [19, 4, 22], enterprise plat-  
30    forms [6], and mobile devices [18]. Most are *simulated* so that grading is deterministic and repro-  
31    ducible. The cost of that reproducibility is impersonality. Each application carries only the data the  
32    current task literally needs, and there is no user history behind it. Web benchmarks in particular  
33    do not evaluate any site that requires logging in or variable personal information, which rules out a  
34    large fraction of what real users ask their assistants to do.

---

<sup>1</sup>*The Office* (US), an American mockumentary sitcom developed by Greg Daniels and aired on NBC, 2005–2013. <https://www.imdb.com/title/tt0386676/>.

# MyPCBench

MyPCBench is a reproducible Linux desktop benchmark seeded from a single canonical persona. The environment provides seventeen pre-logged-in web applications, the full LibreOffice suite, and a 184-task evaluation set.

17

WEB APPS

184

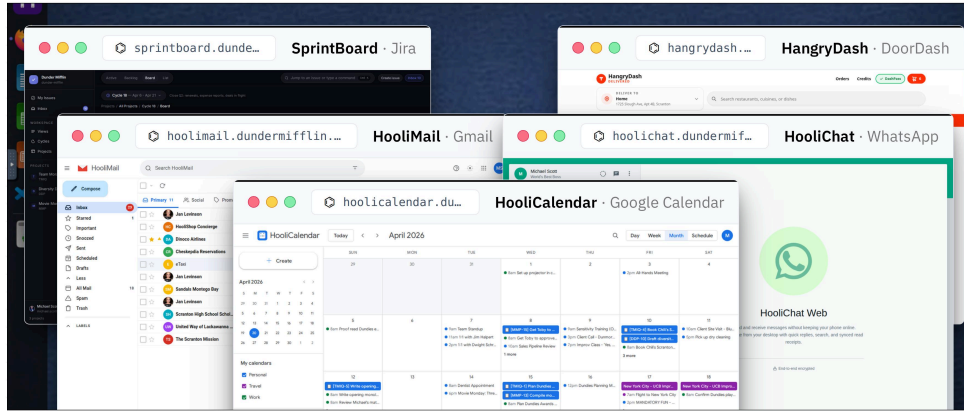
TASKS

6

DOMAINS

17K

RECORDS



<b>MS</b>	<b>Michael G. Scott</b> Regional Manager · Dunder Mifflin · Scranton, PA	<b>1,429</b> BANK TXNS	<b>1,162</b> EMAILS	<b>656</b> CAL. EVENTS	<b>2,586</b> CHAT MSGS	<b>426</b> ORDERS	<b>10,776</b> WEB VISITS	<b>35</b> BOOKMARKS
-----------	--	---------------------------	------------------------	---------------------------	---------------------------	----------------------	-----------------------------	------------------------

Figure 1: **Overview of MYPCBENCH.** A reproducible Linux-desktop benchmark for personally intelligent computer-use agents, seeded end-to-end from a single canonical persona (Michael Scott). The image hosts 17 pre-logged-in web apps mirroring real consumer products plus the full LibreOffice suite; the persona’s records (1,429 bank txns, 1,162 emails, 656 calendar events, 2,586 chats, 426 orders, 10,776 web visits, 35 bookmarks; bottom strip) are cross-linked so that one trip leaves correlated records across every app that would plausibly book it.

35 Simulated benchmarks have traded personalization for reproducibility because no benchmark  
36 has previously seeded a coherent user identity at the scale of a user’s full personal computer.  
37 MYPCBENCH closes this gap. A single persona specification and a deterministic multi-application  
38 generator together produce an environment that is personal, consistent across applications, and re-  
39 producible.

40 Our canonical persona is Michael Scott, the regional manager of a paper company in Scranton,  
41 Pennsylvania. Michael’s desktop is seeded with 1,429 bank transactions, 1,162 emails, 656 calendar  
42 events, 2,586 chat and workplace messages, 115 rideshare requests, 300 food-delivery orders, 126  
43 retail orders, 28 restaurant reservations, and a Firefox profile with 35 bookmarks and 10,776 page-  
44 history visits, distributed across 17 pre-logged-in web applications and the surrounding desktop  
45 stack. The 17 apps and 184 tasks were chosen by manually inspecting the OpenClaw Discord,  
46 the largest personalized-LLM-agent community we are aware of, so that MYPCBENCH reflects the  
47 types of requests users actually issue to a personal assistant. Our contributions are:

- 48 1. A reproducible, cross-app-consistent desktop environment for evaluating personalized agents: 17  
49 custom web apps and a full Linux desktop (Firefox, LibreOffice, file manager), deterministically  
50 populated from one persona seed and packaged as a Docker container.
- 51 2. 184 tasks inspired by real OpenClaw personal-assistant requests, each with a natural-language  
52 rubric, plus an agent harness that drives the standard CUA ReAct [24] loop against the environ-  
53 ment and a rubric-grading LLM-as-a-judge.
- 54 3. Benchmarking of six closed- and open-weight models (Claude Opus / Sonnet 4.6, GPT-5.4 / mini,  
55 Qwen 3.5 35B-A3B / 9B) under each provider’s native CUA agent, with a failure taxonomy and  
56 two scaling analyses (task length and personal-data load).

57 We release the environment image, the 184-task evaluation set with its rubrics, the agent harness,  
58 and the LLM-as-a-judge configuration at [TBD]. Our headline finding is that even the strongest  
59 current frontier agent (Claude Opus 4.6) perfects only 55.4% of MYPCBENCH tasks and only 36%  
60 of tasks that span 7 or more applications. GPT-5.4, Qwen 3.5 35B-A3B, and Qwen 3.5 9B all reach  
61 0% perfect on the same 7+-app slice; the personalization regime stays wide open for future work.

## 62 2 Related Work

63 **Web and desktop agent benchmarks.** The first web benchmarks began with purely synthetic environ-  
64 ments (MiniWoB++ [13], WebShop [23]) and progressed to synthetic realistic websites (We-  
65 bArena [26], VisualWebArena [11]) and then to live Internet evaluations (Mind2Web [5], Web-  
66 Voyager [7], Online-Mind2Web [21]). Desktop benchmarks such as OSWorld [19] extend evalua-  
67 tion to Linux desktops with manually handcrafted reward verifiers; Windows Agent Arena [4] and  
68 MacOSWorld [22] cover the other major operating systems. Static grounding benchmarks such as  
69 ScreenSpot [12] and OmniAct [10] evaluate an LLM’s ability to ground actions on desktop, browser,  
70 and mobile interfaces. LLM agents with personas have been evaluated in limited work settings:  
71 TheAgentCompany [20] places agents in the role of an employee at a simulated software company,  
72 and WorkArena [6] drives ServiceNow workflows. Generative Agents [15] studies a society of  
73 agents in a text-based environment and how they interact under separate personas.

74 **Addressing the personalization gap.** What every benchmark above shares is an impersonal environ-  
75 ment. Desktop evaluations seed only what the task literally needs, so the agent never has to  
76 read across a real user’s data or preferences. Web evaluations live on sites with no logged-in profile,  
77 so any task that requires personal information or pages behind a login is excluded by construction.  
78 You cannot, for example, evaluate calling an Uber, paying a friend back on Zelle, or fulfilling a  
79 user’s usual DoorDash order. MYPCBENCH is a Linux-desktop benchmark with a fixed VM image  
80 and a deterministic snapshot reset (like OSWorld), but the desktop is seeded end-to-end with one  
81 user’s data across every application, not just the data each task touches. The closest existing bench-  
82 mark in spirit is TheAgentCompany, which simulates a software company populated by multiple  
83 coworker personas; MYPCBENCH differs by pinning a single user identity and by spanning the  
84 consumer-application surface that personal computers actually run (banking, travel, food delivery,  
85 calendar, messaging, etc.) rather than internal company tools. The result is that the same agent loop  
86 OSWorld-style evaluations already use can finally be pointed at tasks that require knowing who the  
87 user is.

## 88 3 MYPCBENCH

### 89 3.1 Environment

90 We release MYPCBENCH as a reproducible, open-source Linux desktop through a Docker image  
91 that runs a real QEMU/KVM Ubuntu 24.04 VM with GNOME Shell. The VM hosts 17 pre-logged-  
92 in websites (each modelled on a real-world analogue), the full LibreOffice suite (Writer, Calc, Im-  
93 press), and a Firefox profile pre-loaded with a realistic browsing history and bookmark set. Two  
94 of the web apps, HooliWork (Slack) and HooliChat (WhatsApp), are also exposed as native desk-  
95 top apps, mirroring how a real user would have them open. The home directory is populated with  
96 files relating to Michael’s personal and work life. Figure 2 shows screenshots of all applications.  
97 MYPCBENCH is built around three properties for evaluating personalized agents:

98 **(1) Cross-app consistency.** Any trip, dinner, or client deal leaves correlated records in every appli-  
99 cation that would plausibly record it. Michael’s Philadelphia trip generates a Cheskepdia (Airbnb)  
100 booking, two Gringotts (Chase) charges, a HooliCalendar (Google Calendar) block, two Dinoco  
101 (Delta) boarding passes, browsing history for “Radisson Blu Warwick”, three Travel-folder emails,  
102 and HooliChat (WhatsApp) messages referencing the trip. The seed pipeline writes those records  
103 together so they line up at boot, and runtime cross-app effects keep them in sync. For example, a  
104 HangryDash (DoorDash) order placed at runtime posts a charge to Gringotts and drops a confirma-  
105 tion in HooliMail (Gmail).

106 **(2) Persona coherence.** The user is a specific person, not a generic account. Friends, co-workers,  
107 routines, and preferences are interleaved characteristics of a user’s data, and our environment reflects

# MyPCBench Environment Suite

Each cell depicts one of the seventeen web applications hosted within MyPCBench, paired with the real-world service it mirrors. Every tile is a live screenshot of the corresponding pre-logged-in application.

**17** WEB APPS    **6** DOMAINS    **184** TASKS

## EXAMPLE TASKS

Reconcile last week's Cooper's Seafood charge against the calendar block — confirm who the dinner was with.  
TableFind · Gringotts · HooliCalendar · HooliChat · HooliMail

Total my upcoming NYC trip — flight, hotel, ride to AVP, and any incidentals — and email me the breakdown.  
Dinoco · Cheskepdia · eTaxi · Gringotts · HooliMail

List every Threat Level Midnight purchase I've made on HooliShop and DM Dwight the receipts.  
HooliShop · Gringotts · HooliChat · HooliMail

Reschedule the Benihana team dinner to next Friday and notify everyone who was originally invited on the work channel.  
TableFind · HooliWork · HooliCalendar · Gringotts

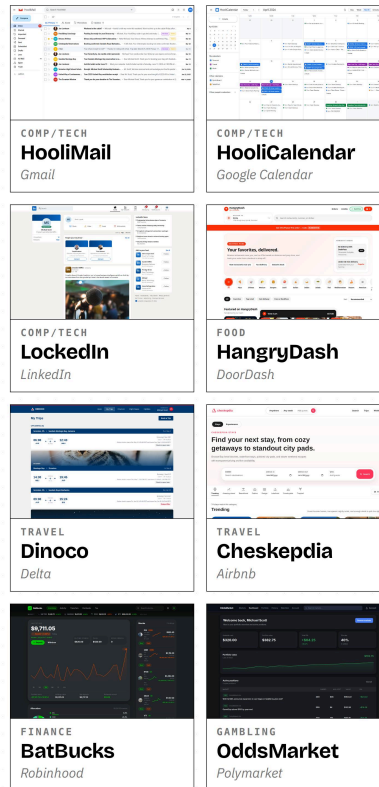


Figure 2: **MYPCBENCH environment suite**. The 17 pre-logged-in web apps span six SimilarWeb top-level domains (Computers/Tech, Finance, Travel, Food, Ecommerce, Gambling). The four example tasks (top-left) each require threading 3–5 of these apps together against Michael’s seeded history.

108 that. Because the persona is Michael Scott, frontier coding agents can draw on *The Office*<sup>1</sup> canon to  
109 populate it with coherent, realistic data.

110 **(3) Real-world fidelity.** Each web application is a local clone with the security and reproducibility  
111 constraints of a fixed VM, but its UI, navigation, and supported flows match the real-world analogue.  
112 Every feature touched by a benchmark task was exercised end-to-end by a human author against the  
113 live VM during the QA pass (§4.1).

## 114 3.2 Environment Creation And Infrastructure

115 **Synthetic website generation.** We built 17 clones of real consumer products using Claude  
116 Code [1], each a full Next.js build rather than a static mock, following prior work on coding-agent  
117 web cloning [25]. Gringotts supports transfers, bill pay, Zelle, and statement downloads; Dinoco  
118 Airlines generates boarding passes with QR codes; eTaxi uses OSRM for realistic routing over 700+  
119 Scranton-area locations; TableFind exposes a reservation inventory of 3,360 slots with hold-and-  
120 release semantics. Across the canonical Michael Scott seed, the 17 applications expose 185 distinct

121 database tables and roughly 18,000 rows of state, of which around 17,000 are the user-facing records  
122 itemised in Figure 1 (transactions, emails, events, messages, orders, browsing history).

123 **Persona generation.** The persona is specified as a JSON document covering identity, financial  
124 profile, social network, travel history, work context, routines, preferences, and recent and upcoming  
125 life events. A deterministic Python pipeline populates every part of the desktop from this  
126 spec: SQLite databases for the 17 web apps with cross-consistent references, a Firefox profile  
127 with bookmarks/history/cookies/form-fields, and a filesystem of meeting notes, expense reports,  
128 trip itineraries, boarding-pass PDFs, and resume drafts.

129 **Infrastructure.** The default resource budget for a single virtual machine is 8 CPU cores and 16 GB  
130 of RAM. Boot-to-ready is about 90 seconds, and a base snapshot is captured after first boot and used  
131 to reset between tasks, avoiding state leakage. Adding personas or websites uses the same template  
132 (Appendix I).

## 133 4 Tasks and Evaluation Setup

### 134 4.1 Task Suite

Type	# Tasks (%)	Representative instruction
Bounded action	64 (35%)	<i>Zelle Pam a hundred bucks. She covered me last weekend. Check HooliChat first to make sure Pam Beesly is on my contacts, then put a memo on the transfer.</i>
Multi-step orchestration	48 (26%)	<i>The Threat Level Midnight Fan Club has been dormant. Peek at the group chat, scroll my LockedIn contacts for Dunder Mifflin folks to recruit, draft them an invitation email, and book a watch party on my calendar for next month.</i>
Cross-source reconciliation	25 (14%)	<i>I've got the Jamaica trip AND the Barbados trip booked about four weeks apart. Given my credit-card balance, can I actually afford both, or am I about to max out?</i>
Aggregation & reporting	23 (12%)	<i>How much am I sending via Zelle each month, and who's getting the money? Check the most recent two complete calendar months and rank the recipients in a LibreOffice Calc spreadsheet.</i>
Personal lookup	13 (7%)	<i>What's my current FlyMiles loyalty tier on Dinoco Airlines, and how many miles do I have in the bank?</i>
Pattern inference	11 (6%)	<i>What do I usually tip on food delivery, in dollars and as a percent? I want to set a smart default so I'm not thinking about it every order.</i>

Table 1: The six behavioural task types in MYPCBENCH, with counts and a representative instruction for each. The split separates analysis tasks (lookup, aggregation, inference, reconciliation; 72 tasks) from action tasks (bounded vs. orchestrated; 112 tasks). 68% of all tasks are multi-application. Full definitions in Appendix C.

135 MYPCBENCH includes 184 tasks, each one inspired by a real use case or request from the Open-  
136 Claw community. To ensure the task set reflects the true distribution of requests issued to personal-  
137 computer assistants, the authors manually sieved through 2,749 anonymized and paraphrased use-  
138 cases from the OpenClaw Discord, the largest community of users running personalized LLM agents  
139 on their own desktops to our knowledge. We dropped requests that (i) were near-duplicates of an  
140 already-kept request, (ii) were infeasible inside any deterministic VM (e.g. “call my mom”), or (iii)  
141 required an app outside the 17 we host. The remaining requests were rewritten by Claude Code so  
142 the named entities (people, restaurants, dates, accounts) match Michael Scott’s seeded data; in the  
143 same pass the coding agent generated a per-task rubric in the Odysseys [8] format. Both the rewrite  
144 and the rubric were then audited by the authors (§4.1). The final task set is stored as JSON, with  
145 each task carrying both its natural-language instruction and its rubric.

146 **Quality assurance.** Because coding agents generate the initial task drafts and the application  
147 clones, we manually verify both. Each task was reviewed by at least two authors through a cus-  
148 tom web interface (Appendix H, Figure 8). Reviewers ran each task end-to-end on the live VM and  
149 confirmed that (a) every named entity exists in the seeded environment, (b) the expected answer is  
150 obtainable from the environment alone, (c) each rubric criterion is individually checkable from a  
151 step-level screenshot, and (d) the task is not a near-duplicate of another in the suite. All 184 tasks  
152 survived this round.

153 **Domain coverage.** We manually map each application to a top-level SimilarWeb<sup>2</sup> category by  
154 inspecting its real-world analogue, mirroring the categorization scheme used by Odysseys [8]. The  
155 17 apps span six SimilarWeb top-level categories and fourteen distinct subcategories. Computers,  
156 Electronics & Technology covers HooliMail, HooliCalendar, HooliWork, SprintBoard, HooliChat,  
157 and LockedIn. Finance covers Gringotts, BatBucks, and SpeedTax. Travel & Tourism covers Dinoco  
158 Airlines, Cheskepdia, and eTaxi. Food & Drink covers HangryDash and TableFind, Ecommerce &  
159 Shopping covers HooliShop and Kwik-E-Mart, and Gambling covers OddsMarket. The full per-app  
160 subcategory mapping, including Banking, Investing, Air Travel, and Restaurants & Delivery, is in  
161 Appendix A.

162 **Apps per task.** Tasks span from one to nineteen co-touched applications, and 68% are multi-  
163 application; 40% span at least two SimilarWeb top-level categories. The multi-application regime  
164 is what tests personalization, since the agent has to reconcile data across the persona’s environment  
165 rather than drive a single tool in isolation. Figure 5 (Appendix B) gives the apps-per-task distribu-  
166 tion, the per-domain task coverage, and the behavioural task-type split.

167 **Task types.** Independently of the SimilarWeb domain, every task is also assigned a behavioural  
168 *type* that captures what the agent must *do* with the persona’s data, independent of which apps are  
169 involved (Table 1). We arrived at the type taxonomy by reading every task instruction and clustering  
170 by the primary capability under test. The resulting six categories cleanly split analysis tasks (lookup,  
171 aggregation, inference, reconciliation) from action tasks (bounded vs. orchestrated).

## 172 4.2 Agent Harness

173 The harness lets us point standard CUA agents at the MYPCBENCH environment with as little adap-  
174 tation as possible. We model MYPCBENCH as a partially observable Markov decision process [9]  
175  $\mathcal{E} = (\mathcal{S}, \mathcal{A}, \Omega, T)$ . At each step  $t$  the agent receives an observation  $o_t \in \Omega$  from the guest VM  
176 and emits an action  $a_t \in \mathcal{A}$ , which the harness executes through the OSWorld [19] VNC bridge. Our  
177 harness is a thin extension of the OSWorld runner. It boots the MYPCBENCH Docker image, re-  
178 stores a fresh QEMU snapshot before every task so each run begins from an identical desktop state,  
179 attaches to the VNC display, and drives the standard agent loop until the agent emits DONE or FAIL  
180 or exhausts the step budget.

181 **Observation space.** At each step the agent receives a 1280×800 screenshot of the full Linux  
182 desktop, augmented with the running tool-call history. Screenshots are passed unmodified to vision-  
183 language model APIs.

184 **Action space.** The action space is the unmodified OSWorld pyautogui surface (click, type, key,  
185 scroll, drag, wait, screenshot, done, fail). Frontier computer-use APIs (Claude Computer Use, Ope-  
186 nAI CUA) each define their own action vocabularies, and we map them onto this surface through  
187 a unified translation layer (for example, Claude’s `computer.click` becomes `click`, and CUA’s  
188 `drag path` becomes `drag`). Anthropic’s Computer Use bundle additionally ships with a native bash  
189 tool and a `str_replace_based_edit_tool`, and we expose those through the same VNC channel  
190 when running Claude. The OpenAI CUA and Qwen CUA APIs do not expose equivalent tools; we  
191 therefore do not provide them to those agents, since adding tools they were not trained against would  
192 be both unfair and out of distribution. The full per-action table, the per-provider mapping, and the  
193 system prompts are in Appendix K.

## 194 4.3 Grading

195 We grade each trajectory against each rubric using an LLM-as-a-Judge, following the scheme pro-  
196 posed by Odysseys [8]. Every task ships with its own set of rubrics: a list of natural-language criteria  
197  $\{r_1, \dots, r_N\}$  authored alongside the task and audited during the QA pass. Across the suite, rubrics  
198 range from 3 to 13 items, with a mean of 6.5 per task and 1,191 rubric items in total. The judge runs  
199 once per rubric item over the *full* trajectory. It is shown the task instruction for context, that one  
200 rubric item, the agent’s complete action history (the running tool-call log), and every screenshot the

---

<sup>2</sup><https://www.similarweb.com/category/>

Model	Perfect $\uparrow$	Rubric score $\uparrow$	Avg steps	Traj. Eff. $\uparrow$
<i>API (closed-weight)</i>				
Claude Opus 4.6	<b>55.4</b>	<b>81.8</b>	46.5	<b>3.61</b>
Claude Sonnet 4.6	39.1	65.4	45.8	3.03
GPT-5.4	23.9	50.5	41.2	2.26
GPT-5.4 mini	18.5	47.5	<b>30.8</b>	2.31
<i>Open weights</i>				
Qwen 3.5 35B-A3B	10.3	39.0	53.2	1.38
Qwen 3.5 9B	4.3	20.2	42.8	0.78

Table 2: Main results on the 184-task suite under each provider’s native CUA agent (gemini-3.1-flash-lite-preview judge, 100-step budget, shared persona context). *Perfect* is the fraction of tasks for which every rubric in the task passed and is our headline metric. *Rubric score* is the per-task average of rubric pass rates (each rubric weighted equally) and gives partial credit. *Traj. Eff.* is rubric score per agent step (Odysseys) in percent. Claude Opus 4.6 leads on every metric, with a perfect rate more than  $2\times$  any non-Claude model and  $13\times$  Qwen 9B. Best per column in **bold**.

201 trajectory produced. The judge prompt caps the screenshot list at the most recent 200 as a defensive  
202 limit on context length; in practice no run in this paper hit that cap because we use a 100-step  
203 budget and capture one screenshot per step. The full action history is always provided in text form  
204 regardless of the cap. The judge returns “Status: success” or “Status: failure” for that item. A rubric  
205 item is considered satisfied iff the judge returns “success”, and we denote this  $s_r \in \{0, 1\}$ . We write  
206  $s_i$  for the per-task average  $(\sum_r s_{i,r})/N_i$  on task  $i$ . We use gemini-3.1-flash-lite-preview as  
207 the judge model throughout; we use a single judge for all runs (Appendix M discusses the resulting  
208 correlation in judge errors across rubrics on the same trajectory).

209 We report three metrics per model. The *rubric score*  $\bar{s} = \frac{1}{T} \sum_i s_i$  is the per-task average of rubric  
210 pass rates (each rubric weighted equally) averaged across tasks; it credits partial completion. The  
211 stricter *perfect rate*  $\frac{1}{T} \sum_i \mathbb{1}[s_i = 1]$  requires every rubric in a task to pass. *Trajectory Efficiency*,  
212 also from Odysseys, measures how much rubric score the agent extracts per step,

$$\text{Traj. Eff.} = \frac{1}{T} \sum_{i=1}^T \frac{s_i}{n_i},$$

213 where  $n_i$  is the number of agent steps on task  $i$ . Trajectory Efficiency penalises agents that arrive at  
214 a correct outcome only after substantial wasted effort, and we report it scaled by 100 (percent of the  
215 rubric satisfied per agent step) for readability. The full judge prompt is in Appendix J.

## 216 5 Experiments and Analysis

### 217 5.1 Main Results

218 We evaluate six models on the full 184-task suite using each provider’s native computer-use (CUA)  
219 agent (§4.2). The four frontier closed-weight models are Claude Opus 4.6 and Claude Sonnet 4.6 [3],  
220 GPT-5.4 and GPT-5.4 mini [14]. The two open-weight models are Qwen 3.5 [17] 35B-A3B and 9B,  
221 picked for contrasting scale within the same family. Every run uses a 100-step budget and the shared  
222 persona-and-environment context reproduced in Appendix K, and all runs are graded by the same  
223 Gemini judge (§4.3). Table 2 reports the three metrics from §4.3 alongside the average and median  
224 step counts each agent actually consumed.

225 Closed-weight frontier agents lead by a wide margin. Claude Opus 4.6 reaches **55.4%** perfect at  
226 81.8% rubric score, the only model above 50% perfect and more than twice the perfect rate of  
227 any non-Claude model. Within the API tier the Opus, Sonnet, GPT-5.4, GPT-5.4 mini ordering is  
228 preserved across both metrics, with Sonnet sitting roughly halfway between Opus and the GPT-5.4  
229 family on both. The open-weight Qwen 3.5 models trail every closed-weight model. Within the  
230 Qwen family the larger 35B-A3B more than doubles 9B on perfect rate (10.3% vs. 4.3%) and on  
231 rubric score (39.0% vs. 20.2%); both still fall well below GPT-5.4 mini.

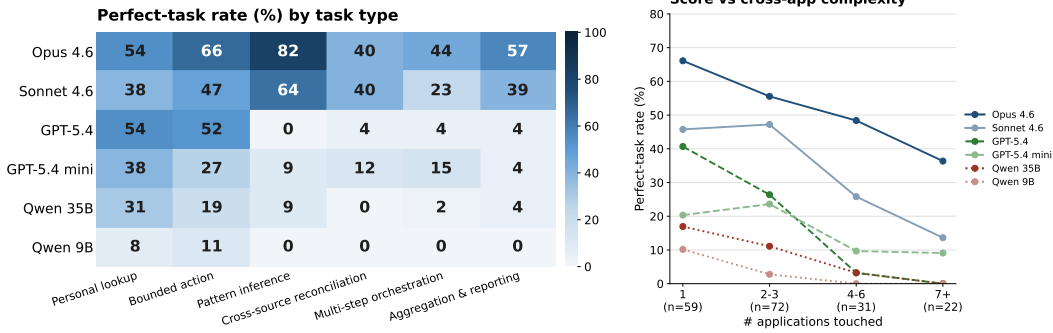


Figure 3: **Left:** per-task-type perfect-task rate (%), models  $\times$  task types. Pattern inference is the largest single-category gap (Opus 82%, GPT-5.4 0%); on the three analysis categories (aggregation, multi-step orchestration, cross-source reconciliation) the GPT-5.4 family and Qwen all stay below 16% perfect. **Right:** perfect-task rate versus the number of distinct applications a task touches. At 7+ apps, GPT-5.4, Qwen 35B, and Qwen 9B all reach 0% perfect, and only the Claude tier and GPT-5.4 mini perfect any 7+-app task at all.

232 Trajectory Efficiency adds a step-budget view. Opus extracts **3.61** rubric points per step, almost  $5 \times$   
 233 Qwen 3.5 9B (0.78). GPT-5.4 mini consumes the fewest steps per task (30.8) but converts them  
 234 less productively than the Claude tier, so its Traj. Eff. (2.31) still trails Sonnet (3.03) by 0.7 points.  
 235 Average step count therefore does not by itself predict efficiency: low step counts can reflect either  
 236 tight execution (Sonnet, 45.8 steps, Eff. 3.03) or premature stopping (Qwen 9B, 42.8 steps, Eff.  
 237 0.78), and the two are distinguished by the failure-mode breakdown in §5.4.

## 238 5.2 Performance by Task Type

239 The aggregate gap in Table 2 hides large differences across the six task types. Figure 3-left shows  
 240 the per-type perfect-task rate for every model.

241 Two categories localise the gap. *Bounded action* and *personal lookup* are the only types where every  
 242 model in the API tier clears 38% perfect; GPT-5.4 matches Sonnet on personal lookup (both 54%)  
 243 and overtakes it on bounded action (52% vs. 47%). The remaining four categories all require either  
 244 reasoning over persona history or coordinating writes across multiple apps, and the gap to Opus  
 245 widens accordingly. The largest single-category gap is *pattern inference*, where Opus reaches 82%  
 246 perfect and GPT-5.4 reaches 0% perfect on the same 11 tasks. These tasks ask the agent to infer an  
 247 unstated rule from many records (“what do I usually tip?”), and the rubric only credits answers that  
 248 match the rule the seeded history supports. On the three remaining analysis categories (aggregation,  
 249 multi-step orchestration, cross-source reconciliation), the GPT-5.4 family and both Qwen models  
 250 stay below 16% perfect, and Qwen 9B records zero perfect tasks across the three combined.

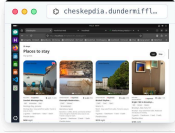
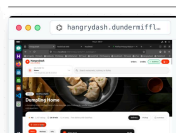
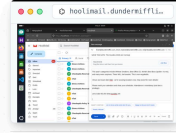
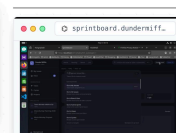
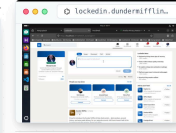
## 251 5.3 Performance Scaling by Steps and Apps

252 The two scaling axes are (a) the number of distinct applications a task touches and (b) the number of  
 253 agent steps the trajectory consumes. The Figure 3-right curve covers the apps axis (per-bin numbers  
 254 in Table 6, Appendix D); a step-budget scaling law (cumulative perfect-task rate at each step budget)  
 255 is in Figure 6 (Appendix D).

256 **Apps and steps both stress horizon.** Every model degrades as apps-touched grows, and the drop  
 257 is sharper on perfect rate than on rubric score. From single-app to 7+-app bins, the perfect rate  
 258 falls from 66%  $\rightarrow$  36% for Opus, 50%  $\rightarrow$  14% for Sonnet, 41%  $\rightarrow$  0% for GPT-5.4, and 23%  $\rightarrow$  9% for  
 259 GPT-5.4 mini; the Qwen 35B model and Qwen 9B both end at 0% in the 7+-app bin (Figure 3-  
 260 right; per-bin numbers and rubric-score equivalents in Table 6). The step axis (Figure 6) sorts the  
 261 cumulative perfect-rate curves into three regimes: Opus is still adding 7pp between steps 60 and  
 262 100, the GPT-5.4 family flattens by step 60 ( $< 2$ pp in the last 40 steps), and Qwen saturates by step  
 263 25. The two axes interact with the task-type findings: 68% of MYPCBENCH tasks span 2+ apps

“Run the full upcoming-Dundies lifecycle plan. Open the categories doc, build a Cheskepdia venue shortlist, scope catering on TableFind and HangryDash, send save-the-dates, coordinate in HooliChat, block the day on HooliCalendar, stand up a SprintBoard task list, and post a teaser on LockedIn.”

**99** **10** **9/9**  
STEPS APPS RUBRICS

<b>01</b> STEPS 4–14 <i>venue search</i>	<b>02</b> STEPS 25–32 <i>catering search</i>	<b>03</b> STEPS 34–48 <i>compose &amp; send</i>	<b>04</b> STEPS 69–89 <i>create task</i>	<b>05</b> STEPS 90–95 <i>compose post</i>
				
<b>Cheskepdia</b> venue shortlist	<b>HangryDash</b> catering scope	<b>HooliMail</b> save-the-date emails	<b>SprintBoard</b> Dundies task list	<b>LockedIn</b> teaser to network

**RUBRICS PASSED 9/9**

- |  |  |   |
|--|--|---|
| <input checked="" type="checkbox"/> R1<br>Read the Dundies categories doc.           | <input checked="" type="checkbox"/> R2<br>Build a Cheskepdia shortlist of ≥2 venues. | <input checked="" type="checkbox"/> R3<br>Scope catering on TableFind and HangryDash. |
| <input checked="" type="checkbox"/> R4<br>Send save-the-date HooliMails to the five. | <input checked="" type="checkbox"/> R5<br>Post a coordination message in HooliChat.  | <input checked="" type="checkbox"/> R6<br>Block the day on HooliCalendar.             |
| <input checked="" type="checkbox"/> R7<br>Stand up SprintBoard with ≥5 tasks.        | <input checked="" type="checkbox"/> R8<br>Post a teaser on LockedIn.                 | <input checked="" type="checkbox"/> R9<br>All artifacts persist across reset.         |

**OUTCOME** All nine artifacts survive a fresh harness boot (emails, HooliChat post, calendar block, SprintBoard tasks, LockedIn teaser).

**PASS · 9/9**

Figure 4: **One full successful Opus trajectory: Dundies-lifecycle plan on long\_horizon-f050** (99 steps, 10 apps, 9/9 rubrics). Cells are real screenshots from the steps where the agent is actively driving each app; the bottom strip enumerates the nine rubric criteria the judge marked passed.

264 and 35% touch 4+, so the steep cross-app slope for the GPT and Qwen families directly produces  
 265 the low perfect rates on the analysis categories of §5.2.

266 **5.4 Personalization-Specific Failures**

267 We read every failed-rubric judge explanation and tagged it with up to five categories: skipped re-  
 268 quired app, premature DONE, surface error as terminal, partial artifact, and hallucinated persona data  
 269 (definitions and per-model counts in Appendix F). Skipped required apps (578 hits) and premature  
 270 DONE (532) account for most of the loss; surface-error abandonment (123), partial artifact (90), and  
 271 hallucinated persona data (36) follow. The three model families concentrate in different modes: GPT  
 272 in premature DONE (307 of 532 hits), Qwen in persona-data hallucination (24 of 36), and Claude in  
 273 console-script shortcuts that hit app REST endpoints with its native bash tool instead of driving the  
 274 visible UI. Mean step count on zero-score trajectories splits the families the same way: GPT-5.4  
 275 (17.7), GPT-5.4 mini (19.7), and Sonnet (29.7) abandon early, while Opus (52.1), Qwen 35B (52.6),  
 276 and Qwen 9B (39.2) keep going until they hit the step budget. Figure 4 traces a clean Opus run  
 277 end-to-end for contrast; per-family pass/fail vignettes are in Appendix L.

278 **6 Conclusion**

279 MYPCBENCH seeds the empty desktop of computer-use evaluation with a coherent persona: 17  
 280 cross-consistent web apps, 17,000 personal records, and 184 OpenClaw tasks. Claude Opus 4.6  
 281 leads at 55.4% perfect; the gap concentrates on multi-app, long-horizon, and history-dependent  
 282 tasks. Environment, tasks, harness, and judge are released at [TBD].

## 283 References

- 284 [1] Anthropic. Claude code. [https://docs.anthropic.com/en/docs/claude-code/](https://docs.anthropic.com/en/docs/claude-code/overview)  
285 [overview](https://docs.anthropic.com/en/docs/claude-code/overview), 2026.
- 286 [2] Anthropic. Claude cowork. <https://claude.com/product/cowork>, 2026.
- 287 [3] Anthropic. Claude opus 4.6. <https://www.anthropic.com/news/claude-opus-4-6>,  
288 2026.
- 289 [4] Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li,  
290 Yadong Lu, Justin Wagle, Kazuhito Koishida, Arthur Bucker, et al. Windows agent arena:  
291 Evaluating multi-modal os agents at scale. *ICML*, 2025.
- 292 [5] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and  
293 Yu Su. Mind2web: Towards a generalist agent for the web. *NeurIPS*, 2023.
- 294 [6] Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H Laradji, Manuel Del Verme,  
295 Tom Marty, Léo Boisvert, Megh Thakkar, Quentin Cappart, David Vazquez, et al. Workarena:  
296 How capable are web agents at solving common knowledge work tasks? *ICML*, 2024.
- 297 [7] Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong  
298 Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal  
299 models. *ACL*, 2024.
- 300 [8] Lawrence Keunho Jang, Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. Odysseys:  
301 Benchmarking web agents on realistic long horizon tasks. *arXiv preprint arXiv:2604.24964*,  
302 2026.
- 303 [9] Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting  
304 in partially observable stochastic domains. *Artificial Intelligence*, 101:99–134, 1998.
- 305 [10] Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Al-  
306 shikh, and Ruslan Salakhutdinov. Omniaact: A dataset and benchmark for enabling multimodal  
307 generalist autonomous agents for desktop and web. *ECCV*, 2024.
- 308 [11] Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Lim, Po-Yu Huang, Graham  
309 Neubig, Shuyan Zhou, Russ Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating  
310 multimodal agents on realistic visual web tasks. *ACL*, 2024.
- 311 [12] Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang,  
312 and Tat-Seng Chua. Screenspot-pro: Gui grounding for professional high-resolution computer  
313 use. *MM*, 2025.
- 314 [13] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement  
315 learning on web interfaces using workflow-guided exploration. *ICLR*, 2018.
- 316 [14] OpenAI. Gpt-5.4. <https://developers.openai.com/api/docs/models/gpt-5.4>,  
317 2026.
- 318 [15] Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and  
319 Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior. *UIST*,  
320 2023.
- 321 [16] Peter Steinberger. Openclaw. <https://openclaw.ai/>, 2026.
- 322 [17] Qwen Team. Qwen3.5-35b-a3b. <https://huggingface.co/Qwen/Qwen3.5-35B-A3B>,  
323 2026.
- 324 [18] Christopher Rawles, Sarah Clinckemaiillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Mary-  
325 beth Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, et al. Androidworld:  
326 A dynamic benchmarking environment for autonomous agents. *ICLR*, 2025.
- 327 [19] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh J  
328 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal  
329 agents for open-ended tasks in real computer environments. *NeurIPS*, 2024.

- 330 [20] Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Z.  
331 Wang, Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Mar-  
332 tin, Zhe Su, Leander Maben, Raj Mehta, Wayne Chi, Lawrence Jang, Yiqing Xie, Shuyan  
333 Zhou, and Graham Neubig. Theagentcompany: Benchmarking llm agents on consequential  
334 real world tasks. *NeurIPS*, 2025.
- 335 [21] Tianci Xue, Weijian Qi, Tianneng Shi, Chan Hee Song, Boyu Gou, Dawn Song, Huan Sun,  
336 and Yu Su. An illusion of progress? assessing the current state of web agents. *COLM*, 2025.
- 337 [22] Pei Yang, Hai Ci, and Mike Zheng Shou. Macosworld: A multilingual interactive benchmark  
338 for gui agents. *NeurIPS*, 2025.
- 339 [23] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable  
340 real-world web interaction with grounded language agents. *NeurIPS*, 2022.
- 341 [24] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan  
342 Cao. React: Synergizing reasoning and acting in language models. *ICLR*, 2023.
- 343 [25] Shuyan Zhou. Webarena-infinity: Generating browser environments with verifiable tasks at  
344 scale. *shuyanzhou.com*, 2026.
- 345 [26] Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,  
346 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for  
347 building autonomous agents. *NeurIPS*, 2024.

## 348 **A Application Details**

349 Table 3 summarises the 17 web applications hosted within the MYPCBENCH environment image,  
350 the real-world service each one mirrors, and the SimilarWeb top-level category and subcategory  
351 inherited from that analogue. Live screenshots of every application are in Figure 2 (§3).

App	Analogue	SimilarWeb category	cate-	Subcategory	Description
HooliMail	Gmail	Computers, Electronics & Tech	Electronics & Tech	Email	Gmail-style web client over a local Maildir; 667 seeded messages across Inbox, Sent, and four labelled folders.
HooliCalendar	Google Calendar	Computers, Electronics & Tech	Electronics & Tech	Productivity	Google-Calendar-style scheduling app exposing 637 personal and work events with recurrence and attendee lists.
HooliWork	Slack	Computers, Electronics & Tech	Electronics & Tech	Programming & Developer Software	Slack-style team messenger with the persona's branch channels, DMs, and read state.
SprintBoard	Jira	Computers, Electronics & Tech	Electronics & Tech	Programming & Developer Software	Jira-style issue tracker holding the persona's running sprints, tickets, and assignees.
HooliChat	WhatsApp	Computers, Electronics & Tech	Electronics & Tech	Telecommunications	WhatsApp-style messenger with one-to-one and group threads spanning friends, family, and co-workers.
LockedIn	LinkedIn	Computers, Electronics & Tech	Electronics & Tech	Social Media Networks	LinkedIn-style professional network exposing the persona's profile, feed, and connections.
HangryDash	DoorDash	Food & Drink	Food & Drink	Restaurants & Delivery	DoorDash-style food delivery surface with order history and active carts.
TableFind	OpenTable	Food & Drink	Food & Drink	Restaurants & Delivery	OpenTable-style restaurant reservation app with past bookings and search.
Kwik-E-Mart	Instacart	Ecommerce & Shopping	Ecommerce & Shopping	Marketplace	Instacart-style grocery delivery with multi-store carts and an order log.
HooliShop	Amazon	Ecommerce & Shopping	Ecommerce & Shopping	Marketplace	Amazon-style retail front with personalised recommendations, orders, and a cart.
Dinoco	Delta	Travel & Tourism	Travel & Tourism	Air Travel	Delta-style airline app exposing flights, loyalty, baggage, and check-in.
Cheskepdia	Airbnb	Travel & Tourism	Travel & Tourism	Accommodation & Hotels	Airbnb-style stays / experiences booking surface with trips, wishlists, and search.
eTaxi	Uber	Travel & Tourism	Travel & Tourism	Ground Transportation	Uber-style ride-hail app with ride history, saved places, and active requests.
SpeedTax	TurboTax	Finance	Finance	Accounting & Auditing	TurboTax-style tax preparation surface with prior-year returns, W-2s, 1099s, and a current-year draft.
Gringotts	Chase Bank	Finance	Finance	Banking, Credit & Lending	Chase-style bank dashboard covering checking, savings, credit card, and a 1,370-row transaction log.
BatBucks	Robinhood	Finance	Finance	Investing	Robinhood-style brokerage with a holdings view, watchlist, and position-level history.
OddsMarket	Polymarket	Gambling	Gambling	Other	Polymarket-style prediction-market exchange with the persona's open positions and watchlist.

Table 3: The 17 web applications hosted within the MYPCBENCH environment image, with the SimilarWeb top-level category and subcategory each one inherits from its real-world analogue. Screenshots of every app are in Figure 2.

## 352 B Task-Distribution Plots

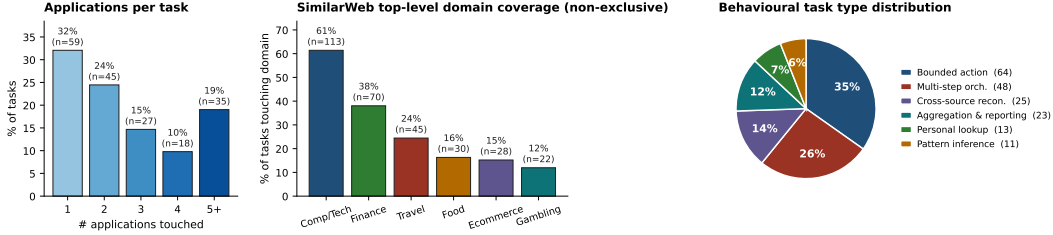


Figure 5: **Left:** distribution of tasks by the number of distinct applications they touch. **Middle:** fraction of tasks that touch at least one application in each SimilarWeb top-level category (non-exclusive; a single multi-app task can contribute to several bars). **Right:** behavioural task-type split (exclusive 1-of-6 categorisation per task).

## 353 C Task-Type Definitions

Type	n (%)	Definition	Representative instruction
Bounded action	64 (35%)	Execute one concrete write action, or a tightly-scoped sequence within a single application, grounded in the persona’s current state.	<i>Zelle Pam a hundred bucks. She covered me last weekend. Check HooliChat first to make sure Pam Beesly is on my contacts, then put a memo on the transfer.</i>
Multi-step orchestration	48 (26%)	Chain reads and writes across multiple applications, typically producing artefacts (LibreOffice docs, decks) plus coordinated side-effects (calendar blocks, sent messages, board updates).	<i>The Threat Level Midnight Fan Club has been dormant. Peek at the group chat, scroll my LockedIn contacts for Dunder Mifflin folks to recruit, draft them an invitation email, and book a watch party on my calendar for next month.</i>
Cross-source reconciliation	25 (14%)	Reconcile a claim, assumption, or hypothetical against the persona’s data by pulling from several apps and quantifying the gap. Covers both contradiction-finding and counterfactual feasibility.	<i>I’ve got the Jamaica trip AND the Barbados trip booked about four weeks apart. Given my credit-card balance, can I actually afford both, or am I about to max out?</i>
Aggregation & reporting	23 (12%)	Compute a total, distribution, ranking, chart, or rollup from many persona records, typically delivered into a LibreOffice document.	<i>How much am I sending via Zelle each month, and who’s getting the money? Check the most recent two complete calendar months and rank the recipients in a LibreOffice Calc spreadsheet.</i>
Personal lookup	13 (7%)	Surface a specific named value, file, or record from the persona’s environment. Single-fact retrieval, no rollup, no inference.	<i>What’s my current FlyMiles loyalty tier on Dinoco Airlines, and how many miles do I have in the bank?</i>
Pattern inference	11 (6%)	Infer an unstated habit, preference, or stylistic pattern from historical persona data, never explicitly stored anywhere in the environment.	<i>What do I usually tip on food delivery, in dollars and as a percent? I want to set a smart default so I’m not thinking about it every order.</i>

Table 4: The six behavioural task types in MYPCBENCH: definition, count, and a representative instruction. Counts cover all 184 tasks. (Main paper Table 1 reproduces only the counts and example instructions.)

## 354 D Per-Task-Type, Cross-App, and Step-Budget Scaling

355 Tables 5 and 6 are the raw numbers behind Figure 3. Figure 6 reads the step axis as a scaling law:  
 356 at each step budget  $X$  on the horizontal axis, the curve plots the fraction of the 184 tasks the model  
 357 graded perfect with  $\leq X$  agent steps consumed.

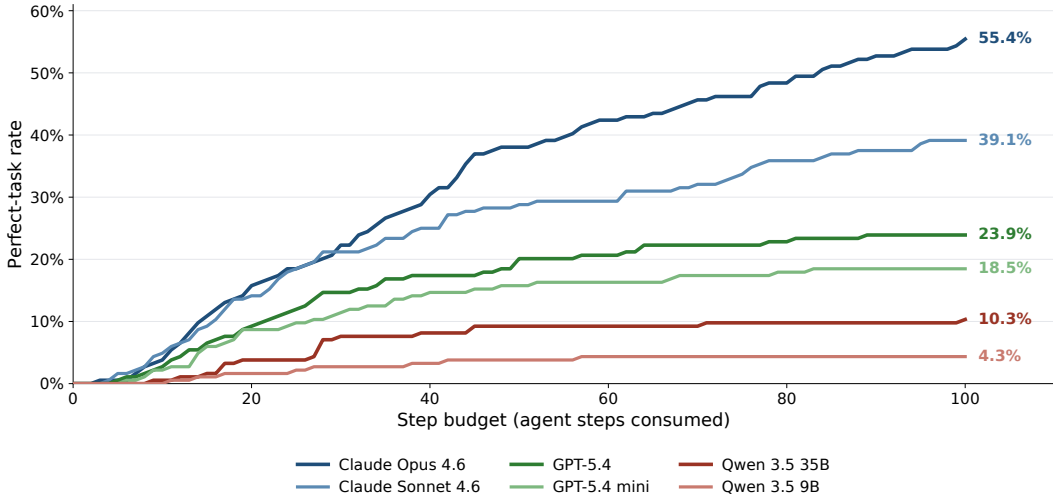


Figure 6: **Step-budget scaling law.** For each model, the curve at step budget  $X$  is the fraction of the 184 tasks the model graded perfect with  $\leq X$  agent steps consumed. Curve shape, not just height, separates the families: Claude Opus is still climbing at the 100-step cap; GPT-5.4 saturates near 60 steps; Qwen flatlines below 25 steps.

Task type	n	Opus	Sonnet	GPT-5.4	GPT-5.4 mini	Qwen 35B	Qwen 9B	Mean
Personal lookup	13	<b>83.8</b>	66.0	71.9	69.0	58.1	32.3	63.5
Bounded action	64	<b>85.3</b>	70.4	73.2	60.6	48.2	31.0	61.5
Pattern inference	11	<b>94.7</b>	77.3	41.3	36.3	28.0	18.5	49.4
Cross-source reconciliation	25	<b>76.1</b>	61.8	41.1	41.4	27.1	13.1	43.4
Multi-step orchestration	48	<b>76.6</b>	59.5	31.5	39.2	37.5	11.4	42.6
Aggregation & reporting	23	<b>81.9</b>	61.9	29.4	28.6	23.6	10.3	39.3

Table 5: Rubric score (%) by task type. Rows ordered by descending cross-model average. *Mean* is the unweighted average across the six models. Best per row in **bold**.

Apps touched	n	Opus	Sonnet	GPT-5.4	GPT-5.4 mini	Qwen 35B	Qwen 9B
1	59	87.4	69.9	59.7	51.6	44.7	31.8
2-3	72	82.4	66.8	57.3	54.2	37.7	16.0
4-6	31	79.8	61.8	31.7	35.9	34.1	15.9
7+	22	67.9	54.1	29.9	30.9	34.5	9.1
$\Delta (1 \rightarrow 7+)$		-19.5	-15.8	-29.8	-20.7	-10.2	-22.7

Table 6: Rubric score (%) versus the number of distinct applications a task touches. 68% of MYPCBENCH tasks are multi-app.

## 358 E Family-Signature Plots

359 Figure 7 breaks failed-rubric hits out by failure mode and groups them by family (left), and shows  
 360 a per-model error budget (right) with the share of zero-score tasks that terminated under 20 steps  
 361 (premature termination) versus those that hit the 99-step budget without success.

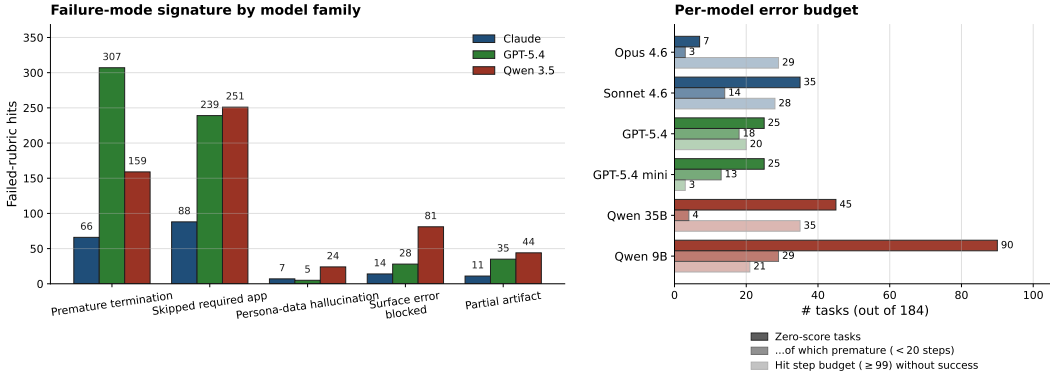


Figure 7: **Left:** failed-rubric hits by failure mode, grouped by family (Claude=Opus+Sonnet; GPT-5.4=GPT-5.4+mini; Qwen=35B-A3B+9B). Each family’s bar set traces a distinct trend. **Right:** per-model error budget. Top dark bar: zero-score tasks; middle: subset that terminated under 20 steps; lightest: trajectories that hit the 99-step budget without success.

## 362 F Detailed Failure Modes

Failure mode	Hits	Description
Premature termination	532	Agent emits DONE before the trajectory satisfies the remaining rubric items.
Skipped required app	578	Multi-app task closed after only some apps are visited; unvisited app’s rubrics fail.
Surface error as terminal	123	Agent hits a captcha, console error, slow page, or modal and quits instead of recovering.
Partial artifact	90	Artifact started but not saved/exported (e.g., spreadsheet opened but never saved).
Hallucinated persona data	36	Agent fabricates a value instead of reading the seeded source.

Table 7: Failure-mode counts on rubric items the judge marked failed, aggregated across all six models. Per-family breakouts are in Figure 7.

363 This section expands the family-level failure trends summarised in §5.4, the compounding effect  
 364 that makes perfect rate fall faster than rubric score, and the two failure modes that are catalogued  
 365 only briefly in the main paper.

366 **Why the perfect rate falls faster than the rubric score.** A failed rubric is rarely an isolated  
 367 event. Skipped-app failures co-occur with premature-DONE on the same trajectory (the agent quits  
 368 because it considers itself done after the apps it did open), and surface errors trigger partial-artifact  
 369 failures (an opened spreadsheet that is never saved). Because perfect-task rate requires *every* rubric  
 370 to pass, even one such co-occurring failure zeroes the task. This compounding is why the perfect-  
 371 rate gap (Opus 55.4%, GPT-5.4 23.9%, Qwen 9B 4.3%) is wider than the rubric-score gap (81.8 /  
 372 50.5 / 20.2).

373 **Per-family breakdown.** Three trends fall out of the per-model counts in Table 7 and Figure 7.

374 • **The GPT family stops too early.** 307 of the 532 premature-DONE hits are in the GPT family  
 375 (GPT-5.4 mini 170, GPT-5.4 137), at least 3.6× either Claude model on its own. This lines up  
 376 with the early-stopping step regime in §5.4 and with the family’s flat step-budget curve in §5.3.

- **The Qwen family makes things up.** 24 of the 36 persona-data-hallucination hits are in Qwen (9B 14, 35B 10), versus 7 in Claude and 5 in GPT. Qwen 9B alone accounts for 39% of all hallucination hits and an additional 58 surface-error abandonments (47% of all surface-error hits).
- **The Claude family takes UI shortcuts.** The Anthropic Computer Use API exposes a native bash tool. Claude trajectories use it to query app REST endpoints (e.g., 55 curl calls to the SprintBoard backend on hard\_app-f026) instead of driving the visible UI. Rubrics that require a user-visible side-effect (move a card, save a file from the menu) then fail even though the agent retrieved the correct value. We do not observe this in the GPT or Qwen runs because the harness exposes no equivalent tool to those agents (§4.2).

**Console-script shortcuts (Claude-specific).** A failure pattern unique to the two Claude models is the console-script shortcut: the agent opens a JavaScript console (or, given Anthropic’s native bash tool, hits the app’s REST endpoint directly with curl) and reads the persona’s data without driving the visible UI. When the rubric only requires that the agent *know* the value, this satisfies it. When the rubric requires a user-visible side-effect (drag a card, open the project, save a file from the menu), the script reads the data and DONEs the task without producing the artifact. hard\_app-f026 is canonical: the judge notes that the agent “investigates SprintBoard throughout the trajectory using API calls in the browser console. . . but it never actually opens the three SprintBoard projects.” Opus issues 55 curl http://localhost:3013/api/projects . . . calls on that trajectory.

**Skipped-app concrete examples.** On 228 failed rubrics, the agent finishes a multi-app task without ever opening one of the named apps. long\_horizon-f066: the agent archives a HooliChat conversation and never opens OddsMarket. long\_horizon-f074: visits nine apps, never opens TableFind to make the required reservation (Figure 10, row 04). aggregation-f010: searches HooliChat extensively but never reads the Dundies categories file in ~/Documents. These are exactly the failures predicted by the cross-app scaling collapse in §5.3.

## 401 G Persona Specification and Event Chains

402 The canonical Michael Scott persona is stored as a single JSON document  
403 (personas/michael\_scott.json in the released code) with fifteen top-level sections that  
404 the generator reads in dependency order. Every seeded record in the 17-app environment can be  
405 traced back to one of these sections.

### 406 Top-level schema.

407

```

408 {
409   "identity":      { name, age, city, address, employer, role,
410                   salary, email, phone, gender, bio, ... },
411   "contacts":     [ { name, relationship, frequency, email, phone,
412                   shared_activities, apps_present_in,
413                   message_personality, birthday, address }, ... ],
414   "financial":    { checking_balance, savings_balance, credit_limit,
415                   credit_used, monthly_income_net,
416                   recurring_charges },
417   "investments":  { cash_balance, holdings, order_history, dividends },
418   "prediction_markets": { balance, total_invested, net_pnl,
419                   active_positions, watchlist },
420   "routines":     { commute, exercise, meals, improv_class },
421   "trips":        [ { destination, dates, hotel, flights, ... }, ... ],
422   "work":         { projects: [ ... ] },
423   "tax_info":     { tax_year, w2, freelance_1099, deductions,
424                   state_code },
425   "planted_contradictions": [ ... ],
426   "planted_dependencies": [ ... ],
427   "browsing_patterns": { research_threads, routine_browsing,
428                   humor_searches },
429   "shopping":     { online_orders, wishlist },
430   "app_overrides": { hoolishop, lockedin, batbucks, speedtax,
431                   hoolichat, etaxi, hangrydash, tablefind },
432   "cross_app_events": [ ... ]
433 }
434

```

436 The four sections that drive cross-app consistency are `cross_app_events` (cross-app side-effects:  
437 a trip seeds rows in six apps), `planted_contradictions` (deliberate red herrings that test whether  
438 agents read all sources), `planted_dependencies` (records the rubric needs the agent to chain  
439 through), and `app_overrides` (per-app tuning, e.g. a Cheskepdia booking that a later HangryDash  
440 record references).

441 **Annotated event chain.** A single `cross_app_events` entry produces correlated rows across ev-  
442 ery application that would plausibly record the event. Below is the canonical Cooper's Seafood  
443 House dinner-plan event taken verbatim from the seed.

```
444 {  
445   {  
446     "type": "dinner_plan",  
447     "description": "Romantic dinner at Cooper's Seafood House for Holly",  
448     "date": "2026-03-28",  
449     "time": "7:30pm",  
450     "apps": ["tablefind", "hoolichat", "gringotts", "hoolimail",  
451             "hoolicalendar"],  
452     "generates": {  
453       "hoolichat_mention": {  
454         "contact": "Jim Halpert",  
455         "context": "Jim. JIM. I need your help. I'm taking Holly to Cooper's  
456           tonight. What do I wear? Should I bring flowers? Is it too much if  
457           I also bring a boombox? Please respond immediately."  
458       },  
459       "browser_history": [  
460         "coopers seafood house scranton reviews",  
461         "romantic restaurants scranton",  
462         "how to be charming at dinner wikihow",  
463         "what wine goes with steak date night"  
464       ],  
465       "calendar_event": true  
466     }  
467   }  
468 }
```

469 The seeders fan out: `seed_webapps.py` writes the TableFind reservation and the Gringotts charge;  
470 `seed_calendar.py` writes the HooliCalendar block; `seed_browser.py` writes the Firefox his-  
471 tory rows; the HooliChat seeder writes the message thread. Because every seeder reads from the  
472 same event record, the entire chain stays internally consistent: TableFind, Gringotts, HooliCalendar,  
473 Firefox history, and HooliChat all reference the same date, time, restaurant identifier, and contact.

## 474 H Task-Review Interface

475 We built a single-page web reviewer (Figure 8) for the quality-assurance pass described in Sec-  
476 tion 4.1. The interface lists every task in the suite grouped by primary application, surfaces the in-  
477 struction, difficulty, and the apps the task touches inline, and exposes per-task review state (unsorted  
478 / unsolved / fix-needed / complete / app-fix-needed) with one-key keyboard shortcuts. Selecting a  
479 task expands a side-pane with the rubric items and a deep-link to the corresponding live application  
480 URL inside the VM, so a reviewer can run the task end-to-end without leaving the page.

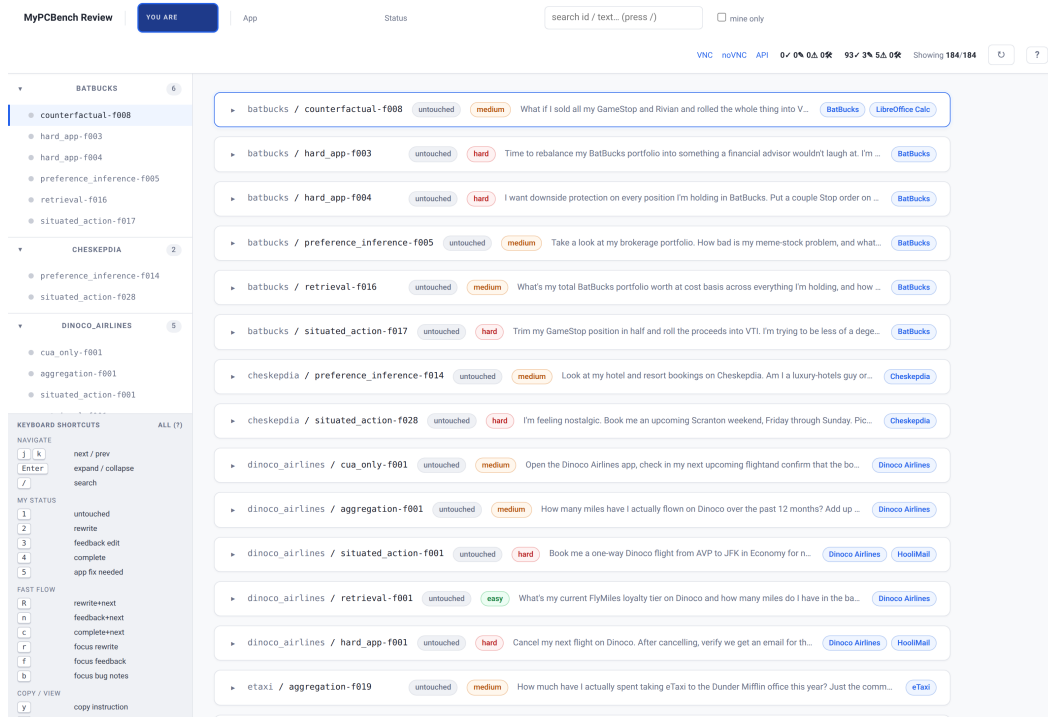


Figure 8: The MYPCBENCH task-review interface used during quality assurance. The left pane is grouped by primary application; each row shows the task identifier, review state, difficulty, the verbatim instruction preview, and pills for the apps the task touches. Reviewer-identifying fields have been removed for double-blind review.

## 481 I Data Generation Pipeline

482 The pipeline turns the persona JSON in Appendix G into a fully populated Linux desktop image. It  
 483 is a single Python entry point (`generator/generate.py`) that calls one seeder per data surface in  
 484 dependency order. Every seeder consumes the same persona document, so adding a new persona is  
 485 a one-file change.

### 486 Seeder modules.

- 487 • `persona_registry.py`: resolves the active persona (default Michael Scott) and exposes identity  
 488 helpers used by every downstream seeder.
- 489 • `persona_data.py` and `enrich_personas.py`: expand the JSON into derived fields (e.g.  
 490 `paystub` line items inferred from `salary` and `tax_info`).
- 491 • `seed_webapps.py`: writes the SQLite databases for the 17 Next.js apps. Honours  
 492 `cross_app_events` so a single trip leaves correlated rows across Cheskepdia, Dinoco, eTaxi,  
 493 Gringotts, HooliMail, and HooliCalendar.
- 494 • `seed_email.py`: builds the on-disk Maildir (`/home/user/Maildir`) and the HooliMail SQLite  
 495 mirror.
- 496 • `seed_calendar.py`: writes HooliCalendar events plus the `.ics` files that the LibreOffice  
 497 toolchain reads.
- 498 • `seed_browser.py` and `inject_firefox_cookies.py`: populate the Firefox profile, history,  
 499 bookmarks, form-fields, and the per-app session cookies that keep every web app pre-logged-in.
- 500 • `seed_filesystem.py`: writes the persona's documents (meeting notes, expense reports,  
 501 trip itineraries, boarding-pass PDFs, resume drafts) into `/home/user/Documents` and  
 502 `/home/user/Downloads`.
- 503 • `document_personas.py`: emits a per-persona Markdown summary used by reviewers and by  
 504 the LLM-as-a-judge as background context.

505 **Determinism.** Every seeder is seeded from a deterministic RNG keyed on the persona name and  
506 the seeder identifier, so identical inputs produce byte-identical outputs across runs and machines.  
507 The default reference time is the persona’s `tax_info.tax_year` end-of-day, overridable through  
508 the `MYPCBENCH_REFERENCE_TIME` environment variable for time-travel debugging.

509 **From persona to image.** After all seeders run, a single `docker build` bakes the populated home  
510 directory, the Firefox profile, and the 17 Next.js apps into the released environment image. The first  
511 boot of the QEMU guest captures a base snapshot; every subsequent task starts from this snapshot,  
512 so the agent always sees the same initial state. Adding a new persona requires (i) writing a JSON  
513 document conforming to the schema in Appendix G and (ii) rerunning `generator/generate.py`  
514 with the new persona slug; no seeder code has to change.

## 515 J Grading and Rubric Prompts

516 This appendix reproduces, verbatim from the released code, the prompts used to grade every task.  
517 Each task’s rubric is a list of natural-language criteria authored alongside the task and audited during  
518 the QA pass (§4.3). The judge runs once per rubric item over the agent’s full trajectory; we use  
519 `gemini-3.1-flash-lite-preview` as the judge model.

### 520 Judge system prompt.

521

```
522 You are an expert evaluator of desktop-agent trajectories.  
523  
524  
525 You will receive:  
526 - The user task (for context).  
527 - ONE specific rubric item with a criterion and (optional) verification description.  
528 - The agent’s full action history (one line per step).  
529 - Every screenshot from the trajectory, in chronological order.  
530  
531 Your goal is to decide whether this single rubric item is satisfied by the trajectory.  
532  
533 Evaluation rules:  
534 - Judge ONLY the one rubric item you are given; ignore all other implicit requirements.  
535 - Ground your judgment in what the screenshots and actions actually show. Do not invent state.  
536 - Filtering / sorting / form requirements must be applied AND confirmed (visible) to count as satisfied.  
537 - If the agent was blocked (captcha, access denied, crash, etc.) and therefore could not satisfy the  
538 ↪ rubric, report failure.  
539 - If a later step UNDONE the rubric (e.g. user-visible state was correct, then was overwritten with wrong  
540 ↪ data), report failure.  
541  
542 Respond in exactly this format:  
543  
544 Thoughts: <your reasoning, citing specific steps/screenshots>  
545 Status: "success" or "failure"
```

547 **Judge user prompt (one call per rubric).** The user message is instantiated from the template  
548 below with the task instruction, the single rubric item being evaluated, and the agent’s compacted  
549 action history; up to the most recent 200 screenshots from the trajectory are attached in chronological  
550 order in the same message.

```
551 User Task (context only): {task_instruction}  
552  
553  
554 Evaluate ONLY this rubric item:  
555 Rubric ID: {rubric_id}  
556 Requirement: {rubric_criterion}  
557  
558 Full Action History:  
559 {action_history}  
560  
561 Screenshots attached below: {n_screenshots} (trajectory had {n_steps} total step(s)).  
562  
563 Decide whether the rubric ({rubric_id}) is satisfied. Use the required 'Thoughts:' / 'Status:' format.
```

565 **Aggregation.** Letting  $s_r \in \{0, 1\}$  denote whether the judge returned “success” for rubric  $r$ , the  
 566 two reported metrics are

$$\text{rubric score} = \frac{1}{N} \sum_{r=1}^N s_r, \quad \text{perfect} = \mathbb{1}[\forall r : s_r = 1].$$

## 567 K Agent Harness

568 This appendix documents the release of the agent harness: its interface, action space (Table 8), step  
 569 budgets, snapshot reset, and the actual system prompts used for each evaluated model family.

570 **Harness interface.** The harness spins up the Docker image (Ubuntu 24.04, XFCE over VNC on  
 571 port 5901, supervisord managing the 17 Next.js services), waits for desktop-ready (typically  $\sim 90$  s),  
 572 and enters the step loop: capture screenshot over VNC; construct the agent message (system prompt  
 573 + task instruction + screenshot + accumulated history); dispatch to the model’s native API; parse  
 574 the returned action; execute on the guest via VNC or shell; repeat until the agent emits done/DONE  
 575 or the step budget is reached. A fresh base snapshot is restored between tasks so each task sees an  
 576 identical initial state.

577 **Action space.** Table 8 enumerates every action exposed to an evaluated agent. The top block is  
 578 the unmodified OSWorld pyautogui surface and is mapped onto every provider’s CUA vocabulary.  
 579 The bottom block lists the two tools that ship natively with Anthropic’s Computer Use API; the  
 580 harness exposes them when running Claude but not for the OpenAI or Qwen CUA agents.

Action	Parameters	Available to
click, right_click	double_click, $x, y$ pixel coordinates	all CUA agents
type	text string	all CUA agents
key	key combination (e.g., ctrl+c)	all CUA agents
scroll	$x, y$ , scroll amount	all CUA agents
drag	start $(x, y)$ , end $(x, y)$	all CUA agents
wait	duration (seconds)	all CUA agents
screenshot	—	all CUA agents
done / fail	—	all CUA agents
bash	shell command string	Claude (native)
str_replace_based_edit_tool	view / create / replace / insert	Claude (native)

Table 8: The MYPCBENCH action space.

581 **Shared persona-and-environment context (all agents).** Every evaluated agent, regardless of  
 582 model family, receives the following persona-and-environment block appended to its system  
 583 prompt. This is the shared frame that keeps the agents grounded on the correct persona, appli-  
 584 cations, and conventions; it is reproduced verbatim from the released code (agents/prompts.py,  
 585 MYPCBENCH\_CONTEXT).

```

586 ## Persona
587 - Name: Michael Scott
588 - Email: 'michael.scott@dundermifflin.com'
589 - Linux user: 'user' (sudo password: '{CLIENT_PASSWORD}')
590
591 ## Environment
592 - Ubuntu 24.04 GNOME desktop. Browser: Firefox (pre-logged-in to every
593 web app via the bookmarks toolbar).
594 - Pinned to dock: HooliChat, HooliWork, Firefox, LibreOffice Writer/Calc/Impress, VS Code.
595 - '/home/user/Documents/', '/home/user/Downloads/', '/home/user/Maildir/' hold persona files.
596 - Python 3.12 and the LibreOffice CLI are available in the VM.
597
598 ## Web apps
599
600
601
602
```

```
603 Each is served at 'http://localhost:PORT', pre-authenticated as the persona.
```

```
604
605 | Port | App | Domain |
606 |-----|-----|-----|
607 | 3001 | Gringotts | personal banking: accounts, transactions, transfers |
608 | 3002 | BatBucks | stock / crypto trading: portfolio, orders |
609 | 3003 | OddsMarket | prediction markets: bets, positions |
610 | 3004 | HooliChat | direct + group messaging |
611 | 3005 | HooliWork | workplace channels |
612 | 3006 | eTaxi | ride hailing: trips, drivers |
613 | 3007 | HangryDash | food delivery: orders, restaurants |
614 | 3008 | TableFind | restaurant reservations |
615 | 3009 | Kwik-E-Mart | grocery orders, inventory |
616 | 3010 | HooliShop | e-commerce: orders, carts, products |
617 | 3011 | Dinoco Airlines | flight bookings, itineraries |
618 | 3012 | Cheskepdia | short-term rental bookings |
619 | 3013 | SprintBoard | project tasks, sprints |
620 | 3014 | LockedIn | professional networking, jobs, connections |
621 | 3015 | SpeedTax | tax returns, filings |
622 | 3016 | HooliMail | email inbox |
623 | 3017 | HooliCalendar | events, invitations |
```

```
624
625 ## Output
```

```
626
627 - Place your final answer (numbers, text, file paths) as plain text in
628 your last assistant turn before any stop signal.
```

630 **Task-completion discipline (all agents).** A single shared block on *when not* to terminate, also  
631 appended to every system prompt.

```
632 ## Task completion discipline
```

```
633
634 - Do NOT emit 'DONE', 'terminate', or any stop signal until you have actually completed the task. A task
635 ↪ is only complete when you have produced the specific output the user asked for AND verified it looks
636 ↪ correct.
637
638 - Use all available steps - plan, act, observe, iterate. Don't bail out early just because the first
639 ↪ approach didn't work.
640
641 - If something fails, try a different approach (different coordinates, different app, different bash
642 ↪ command). Never give up on the first error.
643
644 - Always write your final answer (numbers, text, file contents) before terminating - the grader reads
645 ↪ your last response to check correctness.
646
647 - Only emit 'FAIL' if the task is genuinely impossible (required data literally does not exist). Never
648 ↪ use 'FAIL' as a shortcut when the task is just hard.
```

647 **Claude Computer Use system prompt.** Prepended to the Claude Computer Use request before  
648 the shared persona context. Used by Claude Opus 4.6 and Claude Sonnet 4.6.

```
649 You are an agent on a Linux workstation. Your tools are 'computer', 'bash', and '
650 ↪ str_replace_based_edit_tool'.
```

```
651
652 Stop signal: state your final answer in plain text, then emit "'DONE'" (or "'FAIL'" / '[INFEASIBLE]'
653 ↪ if impossible).
```

656 **OpenAI CUA operator prompt.** Injected as the text portion of the first user message for OpenAI  
657 computer-use agents (GPT-5.4, GPT-5.4 mini); the shared persona context is appended.

```
658 You are an agent on a Linux desktop. Your tools are 'computer' and 'bash'.
```

```
659
660 Stop signal: state your final answer in plain text, then emit "'DONE'" (or "'FAIL'" / '[INFEASIBLE]'
661 ↪ if the task is impossible).
```

664 **Vision-only screenshot-to-pyautogui prompt.** Used for the open-weight agents (Qwen 3.5 35B-  
665 A3B and Qwen 3.5 9B); ported from OSWorld [19] with MYPCBENCH-specific additions in the  
666 shared persona context. The agent emits pyautogui code, which the harness executes on the guest.

```
667 You are an agent which follows my instruction and perform desktop computer tasks as instructed.
668 You have good knowledge of computers and good internet connection and assume your code will run on a
669 ↪ computer for controlling the mouse and keyboard.
```

```

671 For each step, you will get an observation of an image, which is the screenshot of the computer screen
672 ↪ and you will predict the action of the computer based on the image.
673
674 You are required to use 'pyautogui' to perform the action grounded to the observation, but DONOT use the
675 ↪ 'pyautogui.locateCenterOnScreen' function to locate the element you want to operate with since we
676 ↪ have no image of the element you want to operate with. DONOT USE 'pyautogui.screenshot()' to make
677 ↪ screenshot.
678 Return one line or multiple lines of python code to perform the action each time, be time efficient.
679 ↪ When predicting multiple lines of code, make some small sleep like 'time.sleep(0.5)'; interval so
680 ↪ that the machine could take; Each time you need to predict a complete code, no variables or function
681 ↪ can be shared from history
682 You need to to specify the coordinates of by yourself based on your observation of current observation,
683 ↪ but you should be careful to ensure that the coordinates are correct.
684 You ONLY need to return the code inside a code block, like this:
685 ```python
686 # your code here
687 ```
688 Specially, it is also allowed to return the following special code:
689 When you think you have to wait for some time, return ```WAIT```;
690 When you think the task can not be done, return ```FAIL``` , don't easily say ```FAIL``` , try your best
691 ↪ to do the task;
692 When you think the task is done, return ```DONE```.
693
694 My computer's password is '{CLIENT_PASSWORD}', feel free to use it when you need sudo rights.
695 First give the current screenshot and previous things we did a short reflection, then RETURN ME THE CODE
696 ↪ OR SPECIAL CODE I ASKED FOR. NEVER EVER RETURN ME ANYTHING ELSE.

```

## 698 L Example Trajectories

699 Figures 9 and 10 pair one passing and one failing trajectory from each of the three evaluated model
700 families (Claude, GPT, Qwen) at higher fidelity than the main-paper Figure 4. Each vignette row
701 reproduces the verbatim task instruction issued to the agent, three screenshots from the actual tra-
702 jectory (one early, one mid, one near the end), and a short observed-behaviour note explaining
703 how the trajectory arrives at the judge verdict shown in the pill. The pairing surfaces intra-family
704 contrast: the same model family can both complete a task cleanly and exhibit one of the failure
705 modes catalogued in Table 7. The six selected runs cover aggregation-f001, hard\_app-f011,
706 retrieval-f009, long\_horizon-f074, and situated\_action-f036, drawn from three of the
707 failure modes and the family-level trends discussed in §5.4. The Qwen pair places the larger MoE
708 model (Qwen 3.5 35B-A3B) and the smaller dense model (Qwen 3.5 9B) on the same aggregation
709 task to surface scale-driven differences within the family.

Each row is one trajectory drawn from the runs in §5.1: three screenshots span an early, middle, and near-final step. The verdict pill states the rubric outcome assigned by the LLM-as-a-judge. This page shows the Claude pair and the GPT pass; the corresponding GPT and Qwen rows continue on the next figure.

PASS VERDICT  
FAIL VERDICT

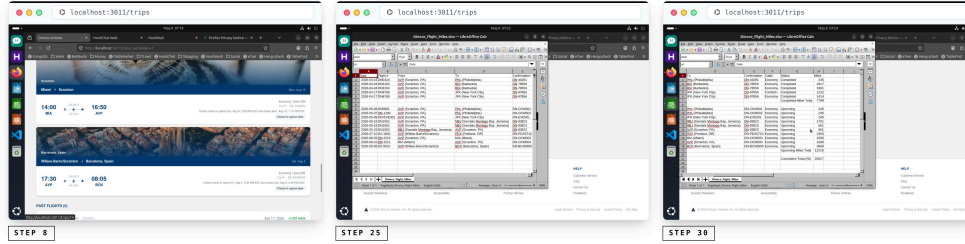
### 01 Claude Opus 4.6

Claude family · aggregation-f001

PASS 4 / 5 RUBRICS

**TASK INSTRUCTION**

How many miles have I actually flown on Dinoco over the past 12 months? Add up the miles from every completed Dinoco flight in that window. Give me a spreadsheet of each trip, and the miles for each one with a final row of the cumulative miles. Also include the upcoming trips and their mileage.



**OBSERVED BEHAVIOUR**

Opens Dinoco and reads the live flight history (step 8), transcribes each completed and upcoming flight into a Calc spreadsheet (step 25), and finishes with a per-trip breakdown and a final cumulative-miles row (step 30). All five rubrics are scored 1 by the judge.

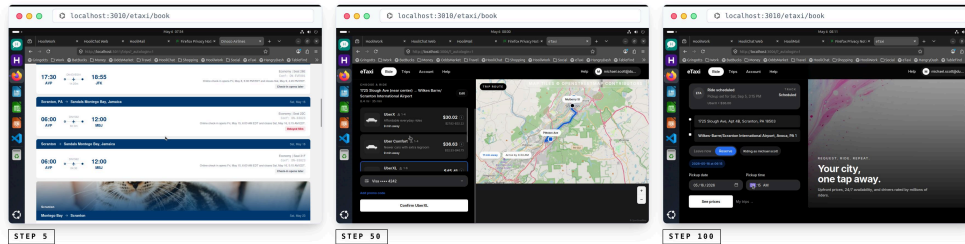
### 02 Claude Opus 4.6

Claude family · hard\_app-f011

FAIL 5 OF 9 RUBRICS MISSED

**TASK INSTRUCTION**

Book me a round-trip eTaxi for my Jamaica trip. I need a ride from my place to AVP in the wee hours of departure day and a return ride back from AVP when I land. Go with the cheapest option that has a driver wait under five minutes. While you're in eTaxi, add two new saved locations for me: the Dundies venue (1901 Mulberry St, Scranton, PA 18510) and the improv academy (its address is on the calendar event location).



**OBSERVED BEHAVIOUR**

Confirms the Jamaica flight dates in Dinoco (step 5) and advances the outbound eTaxi flow to a vehicle-selection card (step 50), but at the 100-step budget the booking is still unconfirmed and the return leg and saved locations were never completed (step 100). Misses R2 (cheapest under-5-min ETA), R3 (return ride), R6 (verification), R8 and R9 (booking confirmations).

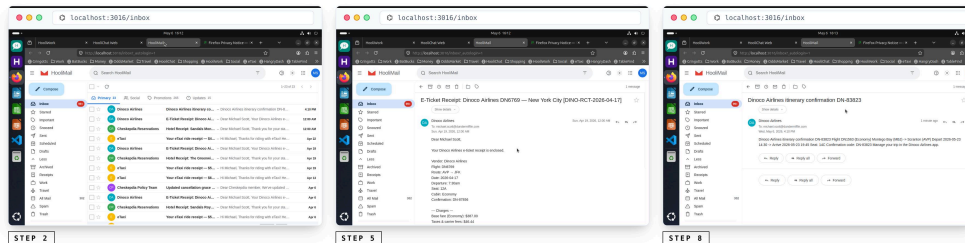
### 03 GPT-5.4

GPT family · retrieval-f009

PASS 3 / 3 RUBRICS

**TASK INSTRUCTION**

For my recent NYC trip, pull up the hotel confirmation, the flight confirmation, and the check-in date. I need to send them to someone. I stayed at the Greenwich and flew Dinoco.



**OBSERVED BEHAVIOUR**

Searches HooliMail and lands on the Receipts folder (step 2), opens the Dinoco DN6769 NYC e-ticket for AVP-JFK (step 5), and surfaces the matching itinerary confirmation DN-83823 (step 8). All four rubrics, including hotel and flight confirmation numbers and the check-in date, are scored 1 by the judge.

Figure 9: Trajectory vignettes, part 1 of 2. Rows: Claude Opus 4.6 PASS on aggregation-f001, Claude Opus 4.6 FAIL on hard\_app-f011, GPT-5.4 PASS on retrieval-f009. Each row shows the verbatim task instruction, three screenshots (early, mid, near-final step), and an observed-behaviour note linked to the rubric outcome assigned by the LLM-as-a-judge.

Continued from the previous figure. The GPT fail uses the same GPT-5.4 backbone as the GPT pass. The Qwen pair contrasts Qwen 3.5 35B-A3B (the larger MoE model) with Qwen 3.5 9B (the smaller dense model) on a single aggregation task to surface scale-driven differences within the family.

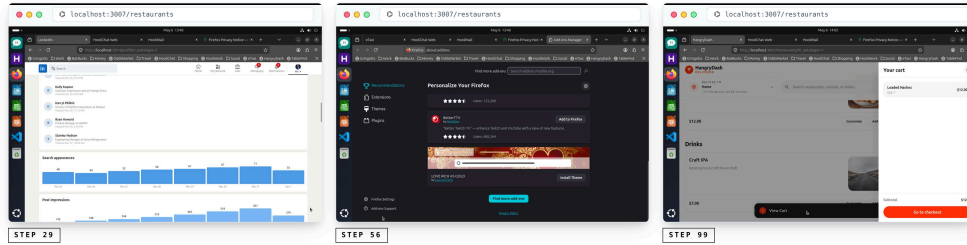
■ PASS VERDICT  
■ FAIL VERDICT

**04 GPT-5.4**  
GPT Family · long\_horizon-f074

FAIL 5 OF 16 RUBRICS MISSED

TASK INSTRUCTION

Scope a side-consultancy: 'Somehow I Manage Consulting'. Build a real go-to-market packet, including a TableFind dinner reservation for the kickoff meeting and a cross-app summary drawing from BatBucks, OddsMarket, SpeedTax, LockedIn, and the Files folder.



OBSERVED BEHAVIOUR

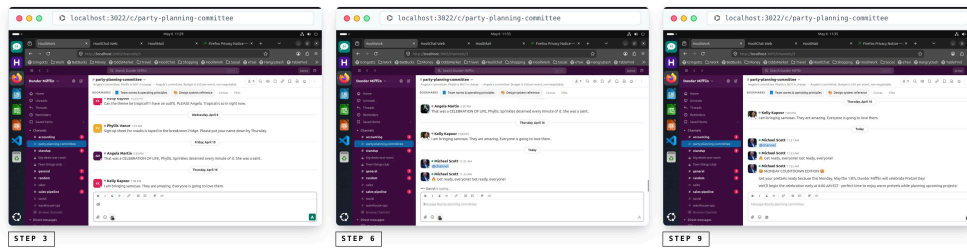
Browses LockedIn for company connections (step 29), drifts into Firefox add-on listings instead of TableFind (step 56), and at step 99 has a \$12 Loaded Nachos cart open in HangryDash. Across the full 100 steps the agent never interacts with TableFind, so the dinner reservation, cross-app summary, and verification rubrics (R3, R7, R8, R9, R10) are all scored 0 by the judge.

**05 Qwen 3.5 35B-A3B**  
Qwen Family · situated\_action-f036

PASS 4 / 4 RUBRICS

TASK INSTRUCTION

Post a Pretzel Day countdown for next Monday at 8:00 AM EST in the #party-planning-committee channel at work. Build the hype.



OBSERVED BEHAVIOUR

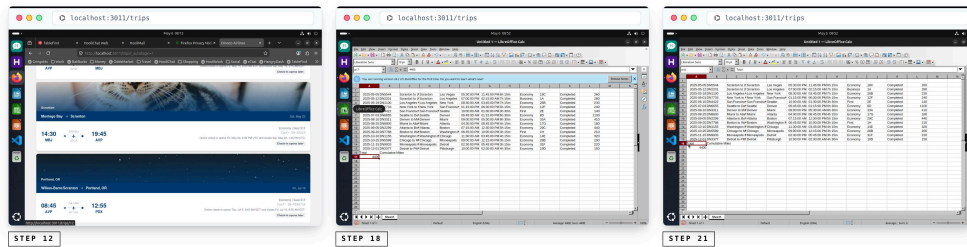
Opens HooliWork (step 3), drafts the Pretzel Day countdown message in the channel composer (step 6), and posts the message to #party-planning-committee with the correct date and timezone (step 9). All four rubrics, including channel selection, persona-correct date, and post persistence, are scored 1 by the judge.

**06 Qwen 3.5 9B**  
Qwen Family · aggregation-f001

FAIL RUBRICS R1, R2, R5

TASK INSTRUCTION

How many miles have I actually flown on Dinoco over the past 12 months? Add up the miles from every completed Dinoco flight in that window. Give me a spreadsheet of each trip, and the miles for each one with a final row of the cumulative miles. Also include the upcoming trips and their mileage.



OBSERVED BEHAVIOUR

Briefly opens Dinoco's flight history (step 12) before switching to LibreOffice Calc through the startup tip dialog (step 18) and filling the sheet with rows that do not match the live history — Las Vegas, Chicago, Atlanta, Greensboro (step 21). The judge scores rubrics R1 (use the live history), R2 (list upcoming flights), and R5 (single coherent answer) as failures; R3 and R4 (table format and final-row total) still pass.

Figure 10: Trajectory vignettes, part 2 of 2 (continued from Figure 9). Rows: GPT-5.4 FAIL on long\_horizon-f074, Qwen 3.5 35B-A3B PASS on situated\_action-f036, Qwen 3.5 9B FAIL on aggregation-f001. Same layout convention as Figure 9.

## 710 **M Limitations**

711 MYPCBENCH commits to one canonical persona (Michael Scott) and one Linux/GNOME/Firefox  
712 software stack; quantitative patterns specific to other demographics, locales, or stacks may not trans-  
713 fer. Grading uses a single Gemini judge, and the judge’s errors are correlated across rubrics on the  
714 same trajectory; absolute failure-mode counts (Table 7) should be read as a structural breakdown  
715 across three models, not as a precise prevalence estimate. The seeded persona is intentionally low-  
716 sensitivity (a public fictional character), so behaviours that emerge only on agents reasoning over  
717 genuinely sensitive personal data are not exercised by this benchmark; we view that as an explicit  
718 out-of-scope choice.

## 719 **N Broader Impact**

720 A personally intelligent agent is, by construction, an agent that can act on a user’s full digital life.  
721 We see two broad classes of impact. On the upside, a benchmark that explicitly tests cross-app,  
722 cross-history personalization should make it harder to ship assistants that look competent on stock-  
723 state demos but fail the moment they meet real personal data, and the failure-mode catalogue gives  
724 developers a concrete checklist (premature DONE, surface-error abandonment, skipped apps, hallu-  
725 cinated persona values, console-script shortcuts) to drive against. On the downside, the same skills  
726 that drive a clean Dundies-lifecycle plan against Michael Scott’s seeded desktop are the skills re-  
727 quired to drive an agent against a real user’s logged-in accounts; numbers on this benchmark should  
728 not be read as a clearance to deploy CUA agents on production accounts. We mitigate the immediate  
729 dual-use surface by (i) seeding only synthetic data tied to a public fictional persona, so the released  
730 image contains no real PII or real correspondence, (ii) hosting every web application locally inside  
731 the QEMU guest, so credentials and form values cannot be exfiltrated to the live web during evalu-  
732 ation, and (iii) recommending offline benchmarking against the released image as the intended use,  
733 and only that.

## 734 **O Release Artifacts**

735 The following artifacts are located at the URL given in §1: (i) the environment image (Docker +  
736 QEMU snapshot), (ii) the set of task evaluations 184, (iii) the rubrics per-task, (iv) the agent harness  
737 that connects standard CUA agents to the environment, and (v) the configuration of the rubric-  
738 grading judge. We do *not* release the agent trajectories or per-rubric judge outputs produced by the  
739 runs in this paper; any future work using the same harness and judge can reproduce them on the  
740 released image.

743 **1. Claims**744 Question: Do the main claims made in the abstract and introduction accurately reflect the  
745 paper’s contributions and scope?

746 Answer: [Yes]

747 Justification: The abstract and §1 claim three contributions: (i) a reproducible cross-  
748 consistent personalized desktop environment, (ii) a 184-task evaluation set with rubrics  
749 and an agent harness, and (iii) a benchmarking of six closed- and open-weight models with  
750 a failure taxonomy and two scaling analyses. Each of these is delivered in the body of the  
751 paper (§3 for the environment, §4 and §4.2 for the task set and harness, and §5 for the  
752 benchmarking, with results in Tables 2 and 5 and Figures 3 and 4).

753 Guidelines:

- 754 • The answer [N/A] means that the abstract and introduction do not include the claims
- 
- 755 made in the paper.
- 
- 756 • The abstract and/or introduction should clearly state the claims made, including the
- 
- 757 contributions made in the paper and important assumptions and limitations. A [No] or
- 
- 758 [N/A] answer to this question will not be perceived well by the reviewers.
- 
- 759 • The claims made should match theoretical and experimental results, and reflect how
- 
- 760 much the results can be expected to generalize to other settings.
- 
- 761 • It is fine to include aspirational goals as motivation as long as it is clear that these
- 
- 762 goals are not attained by the paper.

763 **2. Limitations**

764 Question: Does the paper discuss the limitations of the work performed by the authors?

765 Answer: [Yes]

766 Justification: A dedicated Limitations appendix (Appendix M) discusses the single-persona  
767 / single-stack scope, the single-judge grading setup with correlated errors, and the deliber-  
768 ate decision to use a low-sensitivity public fictional persona rather than real personal data.

769 Guidelines:

- 770 • The answer [N/A] means that the paper has no limitation while the answer [No] means
- 
- 771 that the paper has limitations, but those are not discussed in the paper.
- 
- 772 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 
- 773 • The paper should point out any strong assumptions and how robust the results are to
- 
- 774 violations of these assumptions (e.g., independence assumptions, noiseless settings,
- 
- 775 model well-specification, asymptotic approximations only holding locally). The au-
- 
- 776 thors should reflect on how these assumptions might be violated in practice and what
- 
- 777 the implications would be.
- 
- 778 • The authors should reflect on the scope of the claims made, e.g., if the approach was
- 
- 779 only tested on a few datasets or with a few runs. In general, empirical results often
- 
- 780 depend on implicit assumptions, which should be articulated.
- 
- 781 • The authors should reflect on the factors that influence the performance of the ap-
- 
- 782 proach. For example, a facial recognition algorithm may perform poorly when image
- 
- 783 resolution is low or images are taken in low lighting. Or a speech-to-text system might
- 
- 784 not be used reliably to provide closed captions for online lectures because it fails to
- 
- 785 handle technical jargon.
- 
- 786 • The authors should discuss the computational efficiency of the proposed algorithms
- 
- 787 and how they scale with dataset size.
- 
- 788 • If applicable, the authors should discuss possible limitations of their approach to ad-
- 
- 789 dress problems of privacy and fairness.
- 
- 790 • While the authors might fear that complete honesty about limitations might be used by
- 
- 791 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
- 
- 792 limitations that aren’t acknowledged in the paper. The authors should use their best

793 judgment and recognize that individual actions in favor of transparency play an impor-  
794 tant role in developing norms that preserve the integrity of the community. Reviewers  
795 will be specifically instructed to not penalize honesty concerning limitations.

### 796 3. Theory assumptions and proofs

797 Question: For each theoretical result, does the paper provide the full set of assumptions and  
798 a complete (and correct) proof?

799 Answer: [N/A]

800 Justification: The paper introduces a benchmark and an empirical evaluation; it contains no  
801 theoretical results requiring formal assumptions or proofs.

802 Guidelines:

- 803 • The answer [N/A] means that the paper does not include theoretical results.
- 804 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
805 referenced.
- 806 • All assumptions should be clearly stated or referenced in the statement of any theo-  
807 rems.
- 808 • The proofs can either appear in the main paper or the supplemental material, but if  
809 they appear in the supplemental material, the authors are encouraged to provide a  
810 short proof sketch to provide intuition.
- 811 • Inversely, any informal proof provided in the core of the paper should be comple-  
812 mented by formal proofs provided in appendix or supplemental material.
- 813 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 814 4. Experimental result reproducibility

815 Question: Does the paper fully disclose all the information needed to reproduce the main  
816 experimental results of the paper to the extent that it affects the main claims and/or conclu-  
817 sions of the paper (regardless of whether the code and data are provided or not)?

818 Answer: [Yes]

819 Justification: The environment is shipped as a Docker + QEMU snapshot with a determin-  
820 istic snapshot reset between tasks (§3); the harness is documented at the action-space level  
821 in Table 8 and Appendix K; the judge model, prompt, and per-rubric protocol are repro-  
822 duced verbatim in Appendix J; and Appendix O lists the five release artefacts that together  
823 let any third party reproduce every number in the paper using the same persona seed and  
824 step budget.

825 Guidelines:

- 826 • The answer [N/A] means that the paper does not include experiments.
- 827 • If the paper includes experiments, a [No] answer to this question will not be per-  
828 ceived well by the reviewers: Making the paper reproducible is important, regardless  
829 of whether the code and data are provided or not.
- 830 • If the contribution is a dataset and/or model, the authors should describe the steps  
831 taken to make their results reproducible or verifiable.
- 832 • Depending on the contribution, reproducibility can be accomplished in various ways.  
833 For example, if the contribution is a novel architecture, describing the architecture  
834 fully might suffice, or if the contribution is a specific model and empirical evaluation,  
835 it may be necessary to either make it possible for others to replicate the model with  
836 the same dataset, or provide access to the model. In general, releasing code and data  
837 is often one good way to accomplish this, but reproducibility can also be provided via  
838 detailed instructions for how to replicate the results, access to a hosted model (e.g., in  
839 the case of a large language model), releasing of a model checkpoint, or other means  
840 that are appropriate to the research performed.
- 841 • While NeurIPS does not require releasing code, the conference does require all sub-  
842 missions to provide some reasonable avenue for reproducibility, which may depend  
843 on the nature of the contribution. For example  
844 (a) If the contribution is primarily a new algorithm, the paper should make it clear  
845 how to reproduce that algorithm.

- 846 (b) If the contribution is primarily a new model architecture, the paper should describe  
847 the architecture clearly and fully.
- 848 (c) If the contribution is a new model (e.g., a large language model), then there should  
849 either be a way to access this model for reproducing the results or a way to re-  
850 produce the model (e.g., with an open-source dataset or instructions for how to  
851 construct the dataset).
- 852 (d) We recognize that reproducibility may be tricky in some cases, in which case au-  
853 thors are welcome to describe the particular way they provide for reproducibility.  
854 In the case of closed-source models, it may be that access to the model is limited in  
855 some way (e.g., to registered users), but it should be possible for other researchers  
856 to have some path to reproducing or verifying the results.

## 857 5. Open access to data and code

858 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
859 tions to faithfully reproduce the main experimental results, as described in supplemental  
860 material?

861 Answer: [Yes]

862 Justification: We open-source the environment image, the 184-task evaluation set, the per-  
863 task rubrics, the agent harness, and the rubric-grading judge configuration at the URL given  
864 in §1; Appendix O enumerates the five release artefacts and explicitly notes which artefacts  
865 (per-trajectory rubric outputs from the runs in this paper) are not released. At submission  
866 time the URL is anonymised in keeping with the double-blind policy.

867 Guidelines:

- 868 • The answer [N/A] means that paper does not include experiments requiring code.
- 869 • Please see the NeurIPS code and data submission guidelines ([https://neurips.  
870 cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 871 • While we encourage the release of code and data, we understand that this might not  
872 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
873 including code, unless this is central to the contribution (e.g., for a new open-source  
874 benchmark).
- 875 • The instructions should contain the exact command and environment needed to run to  
876 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
877 //neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 878 • The authors should provide instructions on data access and preparation, including how  
879 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 880 • The authors should provide scripts to reproduce all experimental results for the new  
881 proposed method and baselines. If only a subset of experiments are reproducible, they  
882 should state which ones are omitted from the script and why.
- 883 • At submission time, to preserve anonymity, the authors should release anonymized  
884 versions (if applicable).
- 885 • Providing as much information as possible in supplemental material (appended to the  
886 paper) is recommended, but including URLs to data and code is permitted.

## 887 6. Experimental setting/details

888 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
889 rameters, how they were chosen, type of optimizer) necessary to understand the results?

890 Answer: [Yes]

891 Justification: §4.2 specifies the harness interface, the unified action space, and the 100-step  
892 budget; §4.3 specifies the judge model (`gemini-3.1-flash-lite-preview`), the per-  
893 rubric grading protocol, and the three reported metrics; Appendix K reproduces the system  
894 prompts and the per-provider action mapping verbatim; Appendix J reproduces the judge  
895 prompt verbatim. There is no training: the paper benchmarks pre-trained agent stacks at  
896 inference time only.

897 Guidelines:

- 898 • The answer [N/A] means that the paper does not include experiments.

- 899           • The experimental setting should be presented in the core of the paper to a level of  
900 detail that is necessary to appreciate the results and make sense of them.  
901           • The full details can be provided either with the code, in appendix, or as supplemental  
902 material.

## 903 7. Experiment statistical significance

904 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
905 information about the statistical significance of the experiments?

906 Answer: [No]

907 Justification: Following standard practice in computer-use benchmarking (OSWorld,  
908 Online-Mind2Web, Odysseys), we run each agent on each of the 184 tasks once and report  
909 the per-task aggregate, because end-to-end CUA evaluation is dominated by API and  
910 VM-rollout cost. Headline gaps in this paper (e.g., Opus 55.4% vs. Qwen 9B 4.3% perfect)  
911 are large multiples of any plausible single-run variance, but we do not report formal error  
912 bars and acknowledge this in the Limitations section.

913 Guidelines:

- 914           • The answer [N/A] means that the paper does not include experiments.
- 915           • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
916 intervals, or statistical significance tests, at least for the experiments that support  
917 the main claims of the paper.
- 918           • The factors of variability that the error bars are capturing should be clearly stated (for  
919 example, train/test split, initialization, random drawing of some parameter, or overall  
920 run with given experimental conditions).
- 921           • The method for calculating the error bars should be explained (closed form formula,  
922 call to a library function, bootstrap, etc.)
- 923           • The assumptions made should be given (e.g., Normally distributed errors).
- 924           • It should be clear whether the error bar is the standard deviation or the standard error  
925 of the mean.
- 926           • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably  
927 report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of  
928 Normality of errors is not verified.
- 929           • For asymmetric distributions, the authors should be careful not to show in tables or figures  
930 symmetric error bars that would yield results that are out of range (e.g., negative  
931 error rates).
- 932           • If error bars are reported in tables or plots, the authors should explain in the text how  
933 they were calculated and reference the corresponding figures or tables in the text.

## 934 8. Experiments compute resources

935 Question: For each experiment, does the paper provide sufficient information on the computer  
936 resources (type of compute workers, memory, time of execution) needed to reproduce  
937 the experiments?

938 Answer: [Yes]

939 Justification: §3.2 reports the per-VM resource budget (8 CPU cores, 16 GB RAM, ~90 s  
940 boot-to-ready). Each agent runs are dispatched against the model provider’s hosted CUA  
941 API (no local GPUs required) at the 100-step budget per task. Across 6 models × 184 tasks  
942 at ≤100 steps each, total wall time per full run was approximately 2–3 days on a single  
943 host with four parallel workers; the released harness logs each run’s wall time alongside its  
944 trajectory.

945 Guidelines:

- 946           • The answer [N/A] means that the paper does not include experiments.
- 947           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
948 or cloud provider, including relevant memory and storage.
- 949           • The paper should provide the amount of compute required for each of the individual  
950 experimental runs as well as estimate the total compute.

- 951 • The paper should disclose whether the full research project required more compute  
952 than the experiments reported in the paper (e.g., preliminary or failed experiments  
953 that didn't make it into the paper).

#### 954 9. Code of ethics

955 Question: Does the research conducted in the paper conform, in every respect, with the  
956 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

957 Answer: [Yes]

958 Justification: We have reviewed the NeurIPS Code of Ethics. The persona is a public  
959 fictional character; the seeded data is wholly synthetic and contains no real PII, no real  
960 correspondence, and no real financial data; every web app is a local clone running inside a  
961 sandboxed QEMU guest with no network egress at evaluation time.

962 Guidelines:

- 963 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of  
964 Ethics.
- 965 • If the authors answer [No], they should explain the special circumstances that require  
966 a deviation from the Code of Ethics.
- 967 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
968 eration due to laws or regulations in their jurisdiction).

#### 969 10. Broader impacts

970 Question: Does the paper discuss both potential positive societal impacts and negative  
971 societal impacts of the work performed?

972 Answer: [Yes]

973 Justification: A dedicated Broader Impact appendix (Appendix N) discusses both the up-  
974 side (a benchmark that makes it harder to ship assistants that fail on real personal data; a  
975 concrete failure-mode catalogue for developers) and the dual-use downside (the same skills  
976 generalize to driving agents against real logged-in accounts), along with three concrete mit-  
977 igations baked into the release.

978 Guidelines:

- 979 • The answer [N/A] means that there is no societal impact of the work performed.
- 980 • If the authors answer [N/A] or [No], they should explain why their work has no soci-  
981 etal impact or why the paper does not address societal impact.
- 982 • Examples of negative societal impacts include potential malicious or unintended uses  
983 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
984 (e.g., deployment of technologies that could make decisions that unfairly impact spe-  
985 cific groups), privacy considerations, and security considerations.
- 986 • The conference expects that many papers will be foundational research and not tied  
987 to particular applications, let alone deployments. However, if there is a direct path to  
988 any negative applications, the authors should point it out. For example, it is legitimate  
989 to point out that an improvement in the quality of generative models could be used to  
990 generate Deepfakes for disinformation. On the other hand, it is not needed to point out  
991 that a generic algorithm for optimizing neural networks could enable people to train  
992 models that generate Deepfakes faster.
- 993 • The authors should consider possible harms that could arise when the technology is  
994 being used as intended and functioning correctly, harms that could arise when the  
995 technology is being used as intended but gives incorrect results, and harms following  
996 from (intentional or unintentional) misuse of the technology.
- 997 • If there are negative societal impacts, the authors could also discuss possible mitiga-  
998 tion strategies (e.g., gated release of models, providing defenses in addition to attacks,  
999 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
1000 feedback over time, improving the efficiency and accessibility of ML).

#### 1001 11. Safeguards

1002 Question: Does the paper describe safeguards that have been put in place for responsible  
1003 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
1004 image generators, or scraped datasets)?

1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057

Answer: [Yes]

Justification: The Broader Impact appendix (Appendix N) names three release-level safeguards: (i) all seeded data is synthetic and tied to a public fictional persona (no real PII), (ii) every web app is a local clone hosted inside the QEMU guest with no live-web traffic, and (iii) the recommended use is offline benchmarking against the released image, not pointing CUA agents at production accounts.

Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We build on OSWorld (Apache-2.0) for the harness backbone, Odysseys (CC-BY-4.0) for the rubric-judge scheme, OpenClaw for the source request distribution, and standard model APIs (Anthropic Claude, OpenAI GPT-5.4, Qwen 3.5) for the evaluated agents; each is cited at first use. The web applications inside the benchmark are independent clones built from publicly-visible UI references; they are functionally similar to but not derived from the real services they mirror, and bundle no proprietary art assets. Specific versions and citations are listed in the references.

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The five new assets we release (environment image, 184-task evaluation set, per-task rubrics, agent harness, judge configuration; enumerated in Appendix O) ship with a README at the release URL that documents the persona seed, the per-app database schemas, the task JSON schema, the harness action space, and the exact judge prompt and

1058 decoding parameters. The release URL itself is anonymised at submission time per the  
1059 double-blind policy.

1060 Guidelines:

- 1061 • The answer [N/A] means that the paper does not release new assets.
- 1062 • Researchers should communicate the details of the dataset/code/model as part of their  
1063 submissions via structured templates. This includes details about training, license,  
1064 limitations, etc.
- 1065 • The paper should discuss whether and how consent was obtained from people whose  
1066 asset is used.
- 1067 • At submission time, remember to anonymize your assets (if applicable). You can  
1068 either create an anonymized URL or include an anonymized zip file.

#### 1069 14. Crowdsourcing and research with human subjects

1070 Question: For crowdsourcing experiments and research with human subjects, does the pa-  
1071 per include the full text of instructions given to participants and screenshots, if applicable,  
1072 as well as details about compensation (if any)?

1073 Answer: [Yes]

1074 Justification: The only human-subjects component is the QA pass, performed by paper au-  
1075 thors using the in-house task-review interface. The interface, instructions, and review states  
1076 are documented in Appendix H (with a screenshot in Figure 8). No external annotators or  
1077 crowdworkers were employed, so no compensation question arises.

1078 Guidelines:

- 1079 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1080 with human subjects.
- 1081 • Including this information in the supplemental material is fine, but if the main contri-  
1082 bution of the paper involves human subjects, then as much detail as possible should  
1083 be included in the main paper.
- 1084 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-  
1085 tion, or other labor should be paid at least the minimum wage in the country of the  
1086 data collector.

#### 1087 15. Institutional review board (IRB) approvals or equivalent for research with human 1088 subjects

1089 Question: Does the paper describe potential risks incurred by study participants, whether  
1090 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1091 approvals (or an equivalent approval/review based on the requirements of your country or  
1092 institution) were obtained?

1093 Answer: [N/A]

1094 Justification: The QA pass was performed by paper authors only, on synthetic data tied to  
1095 a public fictional persona. There were no external study participants, no exposure to real  
1096 personal data, and no risks of the type that would normally require IRB review.

1097 Guidelines:

- 1098 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1099 with human subjects.
- 1100 • Depending on the country in which research is conducted, IRB approval (or equiva-  
1101 lent) may be required for any human subjects research. If you obtained IRB approval,  
1102 you should clearly state this in the paper.
- 1103 • We recognize that the procedures for this may vary significantly between institutions  
1104 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1105 guidelines for their institution.
- 1106 • For initial submissions, do not include any information that would break anonymity  
1107 (if applicable), such as the institution conducting the review.

#### 1108 16. Declaration of LLM usage

1109 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1110 non-standard component of the core methods in this research? Note that if the LLM is used  
1111 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
1112 scientific rigor, or originality of the research, declaration is not required.

1113 Answer: [Yes]

1114 Justification: LLMs are used in three core, declared roles. (1) Claude Code generates the  
1115 17 web-app clones and adapts the OpenClaw use cases to Michael Scott’s seeded data; both  
1116 passes are author-supervised and every output was human-verified through the QA inter-  
1117 face (§4, Appendix H). (2) The evaluated CUA agents (Claude Opus / Sonnet 4.6, GPT-5.4 /  
1118 mini, Qwen 3.5 35B-A3B / 9B) are the systems-under-test and are documented per-provider  
1119 in §5.1 and Appendix K. (3) The grading judge is gemini-3.1-flash-lite-preview  
1120 run with the per-rubric protocol of Odysseys; the prompt is reproduced verbatim in Ap-  
1121 pendix J.

1122 Guidelines:

- 1123 • The answer [N/A] means that the core method development in this research does not  
1124 involve LLMs as any important, original, or non-standard components.
- 1125 • Please refer to our LLM policy in the NeurIPS handbook for what should or should  
1126 not be described.