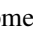# Do Multimodal Foundation Models Understand Enterprise Workflows? A Benchmark for Business Process Management Tasks

**Michael Wornow**     **Avanika Narayan**
**Ben Viggiano**     **Ishan S. Khare**     **Tathagat Verma**
**Tibor Thompson**     **Miguel Angel Fuentes Hernandez**     **Sudharsan Sundar**
**Chloe Trujillo**     **Krrish Chawla**     **Rongfei Lu**     **Justin Shen**
**Divya Nagaraj**     **Joshua Martinez**     **Vardhan Agrawal**     **Althea Hudson**
**Nigam H. Shah**     **Christopher Ré**
Stanford University

## Abstract

Existing ML benchmarks lack the depth and diversity of annotations needed for evaluating models on business process management (BPM) tasks. BPM is the practice of documenting, measuring, improving, and automating enterprise workflows. However, research has focused almost exclusively on one task – full end-to-end automation using agents based on multimodal foundation models (FMs) like GPT-4. This focus on automation ignores the reality of how most BPM tools are applied today – simply documenting the relevant workflow takes 60% of the time of the typical process optimization project. To address this gap we present 🔴🟡🔵 **WONDER-BREAD**, the first benchmark for evaluating multimodal FMs on BPM tasks beyond automation. Our contributions are: (1) a dataset containing 2928 documented workflow demonstrations; (2) 6 novel BPM tasks sourced from real-world applications ranging from workflow documentation to knowledge transfer to process improvement; and (3) an automated evaluation harness. Our benchmark shows that while state-of-the-art FMs can automatically generate documentation (e.g. recalling 88% of the steps taken in a video demonstration of a workflow), they struggle to re-apply that knowledge towards finer-grained validation of workflow completion ($F1 < 0.3$). We hope 🔴🟡🔵 **WONDERBREAD** encourages the development of more "human-centered" AI tooling for enterprise applications and furthers the exploration of multimodal FMs for the broader universe of BPM tasks. We publish our dataset and experiments here: ⭕ https://github.com/HazyResearch/wonderbread.

## 1   Introduction

Multimodal foundation models (FMs) such as GPT-4 [40] have the potential to revolutionize *business process management* (BPM), which is the discipline of measuring and improving enterprise workflows – e.g. a physician submitting a medication order. Typical BPM projects progress in four stages across the following *BPM tasks*: (1) *Documentation* – mapping the steps of an existing workflow; (2) *Knowledge Transfer* – ensuring a shared understanding of the documented workflow; (3) *Improvement* – identifying workflow inefficiencies and proposing fixes; and (4) *Automation* – writing software to execute the workflow without human involvement [49, 53]. FMs could be well-suited for these tasks due to their robust reasoning [65, 57, 2] and visual [8, 54, 67] understanding skills.

**However, existing ML benchmarks [69, 61, 17, 60] focus almost exclusively on one BPM task: end-to-end workflow automation** using agents based on multimodal FMs (see Table 1). This is despite the fact that **simply defining the relevant workflow takes 60% of the time** of the typical BPM project [23], and the BPM market is 4x larger than that of automation tools [47, 48, 27, 28].
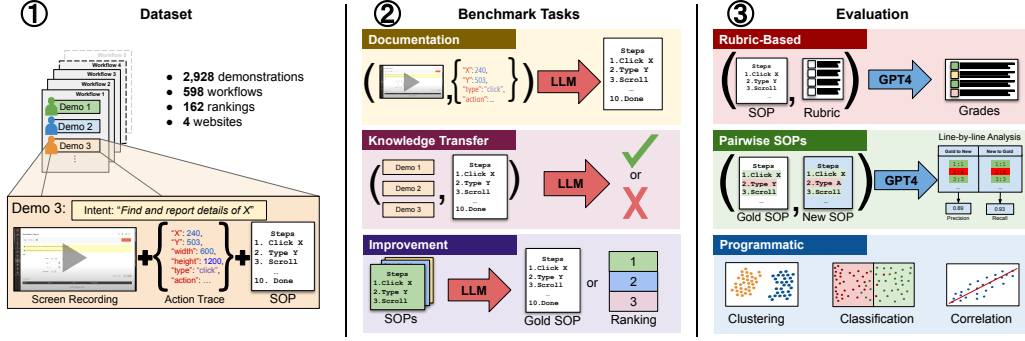
Figure 1: The three components of **WONDERBREAD**. (1) We curate 2928 human demonstrations across 598 web navigation tasks. Each demonstration includes an intent, a full screen recording, an action trace, and a written guide (SOP) describing the steps taken in the demonstration. (2) We create 6 BPM tasks that measure a model's ability to generate accurate documentation, assist in knowledge transfer, and improve workflows. (3) We provide automated evaluation pipelines for all tasks.

By **ignoring the most time-consuming aspects of BPM projects,** we overlook key opportunities to provide near-term value to enterprises. Several case studies have applied multimodal FMs to these broader BPM tasks and demonstrated better performance, easier set-up, and simpler maintenance than traditional BPM tools such as process mining [20, 49, 18, 53, 9, 24, 39]. While promising, however, these papers were largely anecdotal with small datasets ($< 50$ examples). **This motivates the creation of a large-scale benchmark and dataset specifically for BPM tasks.**

Unfortunately, no such dataset exists, and **current benchmarks designed around workflow automation cannot be readily repurposed** due to several limitations. First, their datasets either lack human demonstrations of workflows [69, 17] or do not contain sufficient annotation detail for BPM tasks [61, 15, 35, 30] – e.g. evaluating a model's ability to document a workflow requires reference documentation. Second, their evaluations typically only measure end-to-end workflow completion rates [69, 64, 17, 61] and thus do not consider the intermediate reasoning required for BPM tasks such as identifying inefficiencies within a successfully completed workflow. Third, they do not model real-world BPM use cases and instead focus on navigating websites or mobile apps – i.e. they are focused on workflow *execution* rather than *understanding* [15, 44, 64, 16, 30, 69, 33, 17, 66, 61, 35, 30].

Motivated by the overlooked potential for using multimodal FMs on a broader suite of BPM tasks, we thus introduce 🔴🔵 **WONDERBREAD**, a **WO**rkflow u**NDER**standing **B**enchma**R**k, **E**v**A**luation harness, and **D**ataset. Our contributions are as follows:

1. **Dataset:** We publish **2928 human demonstrations across 598 previously unannotated workflows** sourced from the WebArena benchmark [69]. Each workflow has an average of 4.9 independently collected demonstrations, and each demonstration contains a full screen recording, event log of all clicks/keystrokes/scrolls, and a manually written standard operating procedure ("SOP") – i.e. a step-by-step written guide which reflects the annotator's reasoning at each step of the workflow. For a subset of 162 workflows, we also have annotators rank all 5 demonstrations in order of perceived quality.

2. **Tasks:** Based on use cases drawn from the BPM literature around (1) Documentation, (2) Knowledge Transfer, and (3) Improvement, we define **6 novel BPM tasks** which require reasoning over multimodal data.

   (a) **Documentation:** Generate standard operating procedures (i.e. synthesize the steps of a workflow in writing) to fulfill quality control and audit requirements [5, 59].

   (b) **Knowledge Transfer:** Answer user queries about how workflows operate to simplify onboarding and reduce the 5.3 hours per week that knowledge workers spend waiting for information from colleagues.[42].

   (c) **Improvement:** Analyze workflows to identify inefficiencies and correct execution errors [19, 51].

3. **Evaluation:** We offer evaluation pipelines using automated metrics (e.g., F1, accuracy) and LLM-based evaluators with high correlation to human raters ($\rho > 0.8$). By focusing

on intermediate workflow steps, these evaluations provide a more comprehensive and transparent assessment of models than end-to-end workflow completion rates.

**Results.** We provide baseline results for three state-of-the-art multimodal FMs — GPT-4 [40], Claude 3 [4], and Gemini Pro [50]. Based on screen recordings, we find that models can generate accurate written documentation (F1 of 0.82) and determine whether a demonstration successfully achieved its desired goal (F1 of 0.90). While promising, increasing these numbers to enterprise-level accuracy (i.e. 0.99+) remains an open research challenge. We also identify more significant performance gaps. Models struggle with low-level error correction — for example, when prompted to classify whether a demonstration exactly followed a specific sequence of steps, the peak F1 achieved is 0.27. Models also score poorly when ranking multiple demonstrations of the same workflow on perceived quality and efficiency. We identify long context reasoning, lower-level process understanding, and human workflow preference alignment as key areas for future research.

Our dataset and code available at our Github repo: ⚙ https://github.com/HazyResearch/wonderbread

## 2 Background

We summarize traditional process mining approaches for BPM tasks, discuss recent work on applying multimodal FMs, and compare 🔴🔵 **WONDERBREAD** to existing multimodal FM benchmarks.

### 2.1 Process Mining

Process mining is the *de facto* tool currently used for most BPM tasks, acting as an organizational "X-Ray" [46] that enables large enterprises to identify, measure, and improve their workflows [52, 46, 6]. Techniques include statistical analysis of event logs, unsupervised machine learning, manual review of screen recordings, and user interviews [34, 46]. While interviews can provide an accurate picture of a workflow, they are costly and time-consuming; automated process mining tools are faster but significantly less accurate [1, 34]. Bridging the "semantic gap" between machine and human workflow understanding is an ongoing challenge [38, 34, 1] that 🔴🔵 **WONDERBREAD** aims to address.

### 2.2 Multimodal FMs

Foundation models (FMs) are large-scale ML models trained on vast datasets of unlabeled data which can be applied to a broad range of tasks with minimal adaptation [11]. Multimodal FMs such as GPT-4 combine natural language understanding with a vision model to process images and text jointly [67]. These models have shown promise in navigating graphical user interfaces and executing simple workflows [15, 63, 25, 26, 66, 60]. While the use of multimodal FMs for BPM tasks has been advocated [49], it has not yet been implemented. A failure mode of text-only FMs is the lack of an ability to "read between the lines" of human-generated textual summaries of workflows – e.g. when creating a process model from text, GPT-4 misses half the steps that a human would include [32, 24]. This motivates having multimodal FMs directly observe workflows, as done in our benchmark.

### 2.3 Benchmarks

A number of multimodal datasets have been published for end-to-end automation of websites [69, 17], mobile apps [44], and desktop applications [60, 61]. However, these datasets do not include step-by-step written guides (SOPs), nor do they evaluate on BPM tasks such as documentation, knowledge transfer, or process improvement [15, 44, 64, 16, 30, 69, 33, 17, 66, 61, 35, 30]. Several works have applied large language models to BPM tasks [20, 49, 18, 53, 9, 24, 39], but they conduct limited case studies (i.e. dozens of examples), rely on manual human evaluation, and do not consider multimodal inputs like screen recordings. Please see Table 1 for a detailed comparison with prior benchmarks.

## 3 Dataset

🔴🔵 **WONDERBREAD** includes 2928 human demonstrations across 598 distinct workflows. Each demonstration contains: **(1) Intent** – a short natural language description of the workflow's goal; **(2) Recording** – a full screen recording of the annotator performing the workflow; **(3) Action Trace** – a log of all actions taken (clicks, keystrokes, scrolls) and webpage states before/after each action; **(4) Key Frames** – images taken from the Recording at each action's timestamp; and **(5) SOP** – a written guide detailing all of the steps taken by the annotator. Please see Table 2 for complete definitions. Each workflow has demonstrations from at least 4 annotators to reflect the diversity of work habits present in an enterprise. The full dataset collection process is illustrated in Figure 2.

Table 1: Comparison of **WONDERBREAD** to existing benchmarks for workflows. For **Workflows**, "Env" stands for environment – $W$ is website, $M$ is mobile, and $D$ is desktop. For **Evaluation**, "Auto" means the benchmark contains evaluations for end-to-end workflow automation, "Doc" for documenting workflows, "KT" for knowledge transfer, and "Imp" for process improvement.

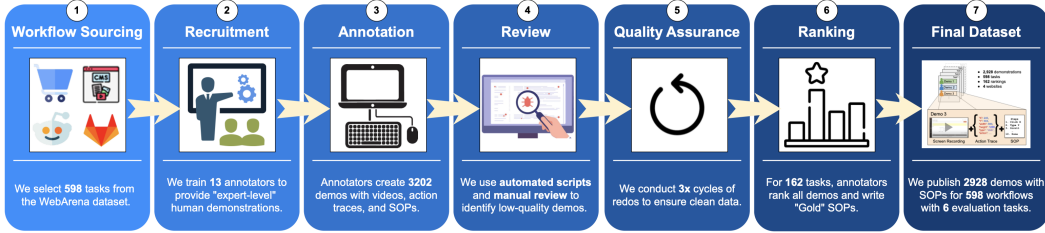| Benchmark | Workflows | | | Human Demonstrations | | | | | Evaluation | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | # Tasks | # Envs | Env Type | Action | Video | SOP | Ranking | Demos/Task | Auto | Doc | KT | Imp |
| AITW [44] | 30,378 | 357 | M | ✔ | ✔ | – | – | 23.5 | ✔ | – | – | – |
| Mind2Web [16] | 2,350 | 137 | W | ✔ | ✔ | – | – | 1 | ✔ | – | – | – |
| MoTIF [13] | 6,100 | 125 | M | ✔ | ✔ | – | – | 0.77 | – | – | – | – |
| WebArena [69] | 812 | 4 | W | ✔ | ✔ | – | – | 0.22 | ✔ | – | – | – |
| OmniAct [30] | 9,802 | 65 | D + W | ✔ | – | – | – | 1 | ✔ | – | – | – |
| WebShop [64] | 12,087 | 1 | W | ✔ | – | – | – | 0.13 | ✔ | – | – | – |
| VWA [33] | 910 | 3 | W | – | – | – | – | 0 | ✔ | – | – | – |
| WorkArena [17] | 23,150 | 5 | W | – | – | – | – | 0 | ✔ | – | – | – |
| WebLINX [35] | 2,337 | 155 | W | ✔ | ✔ | – | – | 1 | ✔ | – | – | – |
| OSWorld [61] | 369 | 13 | D + W | ✔ | ✔ | – | – | 1 | ✔ | – | – | – |
| Wonderbread | 598 | 4 | W | ✔ | ✔ | ✔ | ✔ | 4.9 | ✔ | ✔ | ✔ | ✔ |



Figure 2: The dataset collection process began by selecting 598 workflows from the WebArena dataset [69]. Thirteen annotators then recorded themselves demonstrating roughly 300 workflows each. After multiple rounds of QA, annotators ranked demonstrations for 162 workflows based on perceived quality. The final dataset contains 2928 demonstrations and 6 evaluation tasks.

We start with WebArena, a benchmark containing 812 workflows that require an agent to navigate open-source clones of an e-commerce, content management, forum, and developer tool website [69]. We filter this to 598 workflows by excluding workflows deemed impossible or inadequately specified.

We recruited 13 annotators to record themselves completing each workflow using a custom Python script. Existing workflow benchmarks often have low-quality demonstrations or inaccurate annotations [58], thus a key contribution of **WONDERBREAD** is the high quality of demonstrations achieved through several months of quality assurance. More details are provided in Appendix A.2.

In addition to demonstrations, we also curated 120 free response question-answer pairs to simulate inquiries that a BPM consultant might ask of a workflow. Examples are listed in Appendix A.4.
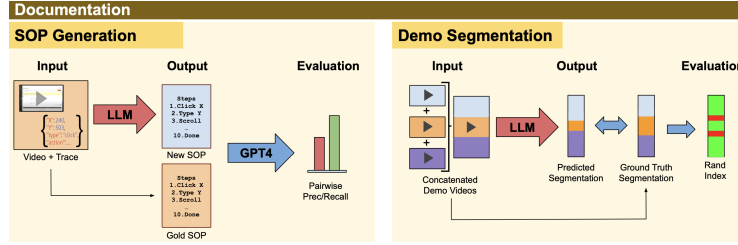
## 4 Benchmark

**WONDERBREAD** contains 6 tasks which cover three BPM applications not evaluated in prior benchmarks: automatically generating documentation from workflow demonstrations (**Documentation**), facilitating knowledge transfer (**Knowledge Transfer**), and identifying ways to improve inefficient workflows (**Improvement**). We provide a summary of each task below. Further details on the inputs, outputs, and evaluations are in Appendix B.

### 4.1 Documentation

Creating clear documentation of complex workflows is essential for operational continuity, compliance, and accountability [59, 5]. This can be achieved through Standard Operating Procedures ("SOP"), Process Definition Documents ("PDD"), or process maps. Our two documentation tasks – SOP Generation and Demo Segmentation – evaluate a model's ability to generate SOPs and accurately distill video recordings into discrete workflows.
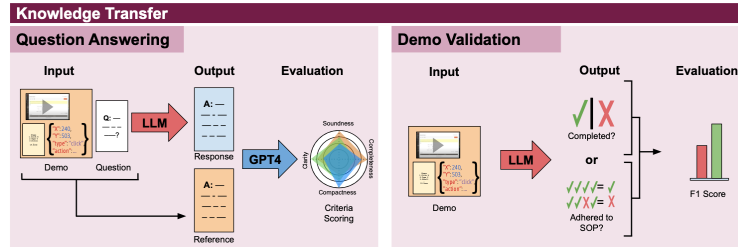
Table 2: Key terms and definitions

| Term | Definition | File Format |
|------|-----------|-------------|
| Task | One of the 6 evaluation tasks in our benchmark, as detailed in Section 4. | – |
| Workflow | A sequence of actions taken to complete a specific business goal. Also referred to as a process. A single workflow can have multiple demonstrations. | – |
| Demonstration | A single execution of a workflow. Each demonstration contains an Intent, Recording, Action Trace, Key Frames, and SOP. | Folder |
| Intent | A brief natural language specification of a workflow's goal, e.g. *"Cancel my last order"*. | .TXT |
| Recording | A video containing a full recording of the user's screen. | .MP4 |
| Action Trace | A log of all click, keystroke, and scroll actions (including associated elements and coordinates). | .JSON |
| Key Frames | Images taken from a Recording that are synced to events in the Action Trace. | .PNG(S) |
| SOP | A "Standard Operating Procedure" detailing (in writing) all of the steps taken in a demonstration. | .TXT |



Figure 3: Expected inputs, outputs, and evaluation settings for **Documentation** tasks.

**(A) SOP Generation**. Evaluation involves using GPT-4 to compare the generated SOP to an annotator-generated reference SOP, calculating precision (how many steps in the generated SOP are in the reference) and recall (how many steps in the reference are in the generated SOP). Each SOP step is evaluated atomically by GPT-4 for semantic equivalence. Details are in Appendix Section C.2.

**(B) Demo Segmentation**. We concatenate multiple workflow demonstrations into a single video and provide it to the model, which identifies the start and end of each workflow. This tests the model's ability to distinguish between sequential workflows. For evaluation, we calculate the adjusted rand index based on the model's assignment of each video frame to a workflow.



Figure 4: Expected inputs, outputs, and evaluation settings for **Knowledge Transfer** tasks.

## 4.2 Knowledge Transfer

The sharing of skills, know-how, and best practices within large organizations can be challenging [42]. By learning from workflow demonstrations, FMs could serve as a query-able repository of organizational knowledge for existing employees, and accelerate on-boarding of new hires by more quickly disseminating key information to trainees [22]. Our two Knowledge Transfer tasks – Question Answering and Demo Validation – assess whether a model can perform higher-level reasoning about the properties and correctness of a workflow.

**(A) Question Answering**. For questions about workflow demonstrations, the model generates a natural language answer, assessing its understanding of workflow semantics. We use GPT-4 to compare the generated answer to a reference answer for evaluation.

5

**(B) Demo Validation**. Given a demonstration, we predict whether (a) the workflow was successfully completed, or (b) the workflow followed the SOP exactly, with individual steps matching precisely. Since each demonstration in 🔴🟡 **WONDERBREAD** is "correct" by definition, we create synthetic negative examples by truncating recordings and shuffling frames. These binary classification tasks assess a model's ability to self-monitor and error-correct.
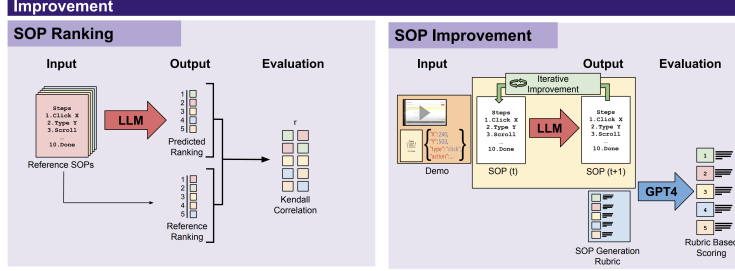


Figure 5: Expected inputs, outputs, and evaluation settings for **Improvement** tasks.

## 4.3 Improvement

The ability to continuously refine and enhance the workflows of an organization is crucial for reducing costs and staying ahead of competitors [19]. By focusing on the improvement of demonstrations and SOPs, we highlight the role of iterative learning and optimization in driving the evolution of workflows [51]. Our two Improvement tasks – SOP Ranking and SOP Improvement – evaluate whether a model can identify workflow inefficiencies and improve inaccurate documentation.

**(A) SOP Ranking**. The same end goal can often be achieved via many different sequences of actions. However, some sequences may be preferable to others as they are more efficient, robust, or avoid intermediate steps that could have undesirable side effects. Given a set of SOPs written by different annotators for the same workflow, this task requires the model to rank them in order of quality. This assesses a model's alignment with human perception of workflow quality. For evaluation, we measure the Kendall $\tau$ correlation between the generated ranking and a human annotator's ranking.

**(B) SOP Improvement**. Given a demonstration and a low-quality SOP, the model must generate an improved SOP that better aligns with the demonstration. The model will iterate to refine the SOP to a specified depth, assessing its ability to assist humans in documenting workflows. GPT-4 will evaluate the generated SOPs against a reference "gold" SOP.

## 4.4 Evaluation

We use programmatic metrics and LLM-based raters for our evaluations. Tasks involving clustering, classification, or ranking use metrics like adjusted rand index, F1, and correlation, respectively. Natural language tasks are evaluated using GPT-4-as-a-judge to assess input quality [14, 68]. Please see Appendix Table 6 for the specific metrics per task. Our LLM-based evaluations show high correlation with human raters ($\rho > 0.8$) (see Appendix Tables 8and 9).

## 5 Results

Our initial results show that current multimodal FMs, including GPT-4, Gemini, and Claude, excel at reasoning over short demonstration lengths but struggle with longer workflows. Our zero-shot evaluations focus on the out-of-the-box capabilities of these models across 162 workflows with rankings. Some models were excluded from specific tasks due to API budget and quota limitations.

### 5.1 Documentation

**SOP Generation.** *Description:* A model must generate a SOP that summarizes all of the actions taken in a video recording of a workflow. We ablate over different demonstration formats: only intent; intent with key frame screenshots; and intent with key frames plus a textual action log of clicks and keystrokes. *Results:* As shown in Table 3, GPT-4 performs best (F1-score of 0.82) with intent, keyframes, and action trace. Most model-demonstration pairs have higher recall than precision (avg. 0.06 points), indicating a tendency to hallucinate workflow steps. Upon qualitative review, we found that many hallucinated actions seemed reasonable but were not actually taken in the demonstration, e.g. adding *"Navigate to the shopping admin page"* even though the demonstration started on that page. Exact scores for each workflow and model are in Appendix Figure 9.

Table 3: **SOP Generation:** Accuracy of generated SOPs versus ground truth SOPs.

| Model | Intent | Keyframes | Trace | Precision | Recall | F1 | Avg. # of Steps |
|---|---|---|---|---|---|---|---|
| GPT-4 | ✓ | ✓ | ✓ | **0.80** | **0.88** | **0.82** | 10.26 |
| GPT-4 | ✓ | ✓ | | 0.69 | 0.79 | 0.71 | 10.32 |
| GPT-4 | ✓ | | | 0.48 | 0.59 | 0.49 | 13.10 |
| Claude 3 Sonnet | ✓ | ✓ | ✓ | 0.72 | 0.85 | 0.76 | 10.94 |
| Claude 3 Sonnet | ✓ | ✓ | | 0.67 | 0.78 | 0.70 | 11.35 |
| Claude 3 Sonnet | ✓ | | | 0.53 | 0.54 | 0.50 | 11.34 |
| Gemini Pro 1 | ✓ | ✓ | ✓ | 0.58 | 0.63 | 0.58 | 11.09 |
| Gemini Pro 1 | ✓ | ✓ | | 0.48 | 0.51 | 0.46 | 11.28 |
| Gemini Pro 1 | ✓ | | | 0.40 | 0.36 | 0.34 | 7.31 |
| Ground Truth | ✓ | ✓ | ✓ | 1 | 1 | 1 | 8.40 |

**Demo Segmentation.** *Description:* This task mimics what a video recording of a person's screen would capture during the typical workday, i.e. multiple workflows without clear boundaries. Concretely, the model receives $k$ concatenated demonstrations sampled from different workflows from our dataset, and must determine which frames belong to the same workflow. We set $k = 3$ and choose workflows that utilize the same website. *Results:* As shown in Table 4, segmenting a recording remains challenging. Providing additional information to the model in the form of an SOP and intent slightly increases performance for GPT-4 yet decreases performance for Gemini Pro 1.

Table 4: **Demo Segmentation:** Accuracy of clustering with $k = 3$ concatenated workflows.

| Model | Intent | SOP | Keyframes | V-Measure | Adj. RI |
|---|---|---|---|---|---|
| GPT-4 | ✓ | ✓ | ✓ | **0.85** | **0.88** |
| GPT-4 | | ✓ | ✓ | **0.85** | 0.87 |
| GPT-4 | | | ✓ | 0.80 | 0.86 |
| Gemini Pro 1 | ✓ | ✓ | ✓ | 0.54 | 0.65 |
| Gemini Pro 1 | | ✓ | ✓ | 0.53 | 0.65 |
| Gemini Pro 1 | | | ✓ | 0.58 | 0.69 |

## 5.2 Knowledge Transfer

**Question Answering.** *Description:* This task involves answering 120 free response questions about workflows, such as *"How would a user know the workflow is complete?"* and *"What is the purpose of this workflow?"*. These questions were drawn from the process mining literature [9, 20] and are provided in full in Appendix A.4. We use GPT-4-as-a-judge to evaluate model-generated answers by comparing to a reference answer from a human annotator. Following prior work [20], we have GPT-4 output four scores on a scale of 1 (bad) to 3 (good): completeness, soundness, clarity, and compactness. The Pearson correlation between GPT-4 and human raters was between 0.80 and 0.89 across all axes (see Appendix Table 8). *Results:* The average scores for each model are shown in Figure 6. All models perform well in "compactness" and "clarity" but score lower on "completeness." This may be due to their reliance on statistical patterns rather than contextual workflow understanding [36], leading to occasional omissions of relevant details despite scoring highly on the syntactic quality of their writing.
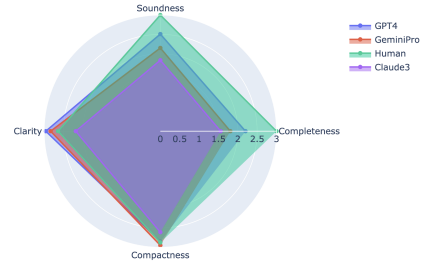


Figure 6: **Knowledge Transfer:** Scores across four axes – soundness, completeness, clarity, and compactness – across 120 free response questions for evaluating workflow understanding.

**Demo Validation.** *Description:* We consider two forms of validation: (a) workflow completion, where a demonstration is "correct" if the workflow's goal is achieved; and (b) workflow trajectory, where it is "correct" only if the goal is achieved following the specified SOP. "Correct" examples are sampled from our dataset, while "incorrect" examples are created by truncating, shuffling, or skipping states. *Results:* As shown in Table 5, GPT-4 is the best model. It can accurately determine whether a workflow completed its overall goal (peak F1 of 0.90) but struggles to validate that it followed the specific steps of an SOP (peak F1 of 0.27).

7

Table 5: **Demo Validation:** Accuracy on binary classification of whether a workflow was completed (*Completion*) or followed the exact steps outlined in the SOP (*Trajectory*).

| Model | Intent | Keyframes | SOP | Precision | Recall | F1 |
|---|---|---|---|---|---|---|
| **Completion** | | | | | | |
| GPT-4 | ✓ | ✓ | ✓ | **0.89** | **0.90** | **0.90** |
| GPT-4 | ✓ | ✓ | | 0.84 | 0.77 | 0.81 |
| Gemini Pro 1 | ✓ | ✓ | ✓ | 0.94 | 0.25 | 0.40 |
| Gemini Pro 1 | ✓ | ✓ | | 0.94 | 0.26 | 0.41 |
| Claude3 Sonnet | ✓ | ✓ | ✓ | 0.58 | 0.31 | 0.40 |
| Claude3 Sonnet | ✓ | ✓ | | 0.85 | 0.50 | 0.63 |
| **Trajectory** | | | | | | |
| GPT-4 | ✓ | ✓ | ✓ | 0.52 | **0.18** | **0.27** |
| Gemini Pro 1 | ✓ | ✓ | ✓ | **0.94** | 0.14 | 0.25 |

## 5.3 Improvement

**SOP Improvement.** *Description.* In this task we provide a model with a task recording and an SOP. The model is then tasked with subsequently improving the SOP given and SOP rubric. *Results.* As shown in Table 7a, current models are capable of improving the quality of their own SOPs (up to 1.4 points), conditioned upon a SOP rubric.

**SOP Ranking.** *Description:* In this task, we provide a model with SOPs from various annotators and have it rank them by quality. We then compare this ranking to a ground truth ranking by an annotator and measure the correlation between the model's and human's judgments. *Results:* As shown in Table 7b, current models struggle to rank SOPs based on perceived quality to human raters. The best model achieves a mean Kendall correlation of 0.05 with a standard deviation of 0.47, indicating essentially random rankings. Improving alignment between model and human judgment of workflow quality remains an area for further research.

| Model | Original SOP | Improved SOP |
|---|---|---|
| GPT-4 | 3.43 | **4.82** |
| Claude3 Sonnet | 3.43 | 4.26 |
| Gemini Pro 1 | 3.43 | 3.65 |

| Model | Spearman $\rho$ | Kendall $\tau$ |
|---|---|---|
| GPT-4 | **0.07 ± 0.58** | **0.06 ± 0.49** |
| Claude3 Sonnet | 0.06 ± 0.59 | 0.03 ± 0.50 |
| Gemini Pro 1 | 0.03 ± 0.58 | 0.03 ± 0.49 |

(a) **SOP Improvement:** Scores from 1 (bad) to 5 (good) for SOPs before/after model improvement.

(b) **SOP Ranking:** Corr. between model and human rankings of demonstrations for the same workflow.

Figure 7: Results for the two **Improvement** benchmark tasks.

## 6 Discussion

We discuss next steps, limitations, and the broader impacts of 🔴🔵 **WONDERBREAD** below.

**Improving Human-Model Alignment for BPM Tasks.** We find that out-of-the-box human and multimodal models alignment is low for SOP evaluation (see Section 5.3). Similar to how "human-model" alignment can be achieved for tasks like question-answering and instruction-following [55, 31], alignment also appears necessary for workflow understanding tasks. This might require fine-tuning models via supervised learning [56] or reinforcement learning on preference data [41].

**Expanding Multimodal Context Windows.** Even a 1-minute workflow can generate dozens of actions and key frames. Our results show that model accuracy on BPM tasks improves as more information is provided in the prompt. This might not be possible with longer workflows, leading to an incomplete representation for a workflow and lower downstream task performance. Longer context windows can help solve this problem and are a focal point of study in the community [29, 62].

**Low-Level Workflow Understanding**. Our results show that while multimodal FMs excel in high-level workflow analyses, they struggle with precise validation of individual steps (see Section 5.2). Enhancing this lower-level understanding may require supervised fine-tuning on GUIs as in [26, 7].

**Self-Improvement.** Our findings suggest that multimodal FMs can improve their outputs (i.e., SOPs) through multiple iterations of self-reflection (see Section 5.3). This highlights the potential of these

models to refine their outputs without human intervention [21, 3]. In the context of BPM tasks, this capability can help systems adapt to workflows as they change over time.

**Limitations.** There are several limitations to our work. First, dataset construction was constrained by our lack of access to real-world enterprise data due to privacy concerns. Second, the workflows in our dataset are taken from a limited set of 4 websites [69], and it is unclear how our results generalize to different environments with complex or longer workflows. Contemporaneous to our work, several datasets have been released which could be re-annotated following the process described in our paper [61, 35, 30], which we leave to future work. Third, our baseline results lack open-source models. Matching the performance of state-of-the-art proprietary models on these benchmarks with open source models remains an open research challenge.

**Societal Impact.** Our field's collective focus on automation contradicts recent advocacy for more *human-centered AI*, which aims to *augment* rather than *replace* human labor [43, 45, 12, 10]. While we intend for **WONDERBREAD** to serve as a counterpoint to this focus, we acknowledge that any AI tools aimed at improving productivity run the risk of automating jobs or replacing human labor.

# 7    Conclusion

We present **WONDERBREAD**, the first benchmark for evaluating multimodal models on common process mining tasks. It includes 2928 human demonstrations across videos, images, and text, along with step-by-step written guides (SOPs) and full action traces. We focus on applying these models to three BPM tasks that have been overlooked by existing ML benchmarks for workflow automation – documentation, knowledge transfer, and process improvement. **WONDERBREAD** features an automated evaluation harness with programmatic metrics and LLM-based assessments, providing baseline results for state-of-the-art multimodal models. Our work aims to inspire further efforts to support workers by *augmenting* rather than *replacing* human labor.

# References

[1] Simone Agostinelli, Andrea Marrella, and Massimo Mecella. Towards intelligent robotic process automation for bpmers. *arXiv preprint arXiv:2001.00804*, 2020.

[2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.

[3] Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang Li, Sheila Babayan, Kavya Kopparapu, Zachary Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Srinivasan, et al. Rest meets react: Self-improvement for multi-step reasoning llm agent. *arXiv preprint arXiv:2312.10003*, 2023.

[4] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.

[5] American Hospital Association. Assessing the regulatory burden on health systems, hospitals and post-acute care providers, 2017.

[6] Adriano Augusto, Raffaele Conforti, Marlon Dumas, Marcello La Rosa, Fabrizio Maria Maggi, Andrea Marrella, Massimo Mecella, and Allar Soo. Automated discovery of process models from event logs: Review and benchmark. *IEEE transactions on knowledge and data engineering*, 31(4):686–705, 2018.

[7] Gilles Baechler, Srinivas Sunkara, Maria Wang, Fedir Zubach, Hassan Mansoor, Vincent Etter, Victor Cărbune, Jason Lin, Jindong Chen, and Abhanshu Sharma. Screenai: A vision-language model for ui and infographics understanding. *arXiv preprint arXiv:2402.04615*, 2024.

[8] Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. Introducing our multimodal models, 2023.

[9] Alessandro Berti and Mahnaz Sadat Qafari. Leveraging large language models (llms) for process mining (technical report). *arXiv preprint arXiv:2307.12701*, 2023.

[10] Joseph R Biden. Executive order on the safe, secure, and trustworthy development and use of artificial intelligence. 2023.

[11] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

[12] Erik Brynjolfsson. The turing trap: The promise & peril of human-like artificial intelligence. *Daedalus*, 151(2):272–287, 2022.

[13] Andrea Burns, Deniz Arsan, Sanjna Agrawal, Ranjitha Kumar, Kate Saenko, and Bryan A. Plummer. Mobile app tasks with iterative feedback (motif): Addressing task feasibility in interactive visual environments, 2021.

[14] Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937*, 2023.

[15] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web, 2023.

[16] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web, 2023.

[17] Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H. Laradji, Manuel Del Verme, Tom Marty, Léo Boisvert, Megh Thakkar, Quentin Cappart, David Vazquez, Nicolas Chapados, and Alexandre Lacoste. Workarena: How capable are web agents at solving common knowledge work tasks?, 2024.

[18] Marlon Dumas, Fabiana Fournier, Lior Limonad, Andrea Marrella, Marco Montali, Jana-Rebecca Rehse, Rafael Accorsi, Diego Calvanese, Giuseppe De Giacomo, Dirk Fahland, et al. Ai-augmented business process management systems: a research manifesto. *ACM Transactions on Management Information Systems*, 14(1):1–19, 2023.

[19] Marlon Dumas, Marcello La Rosa, Jan Mendling, Hajo A Reijers, et al. *Fundamentals of business process management*, volume 2. Springer, 2018.

[20] Dirk Fahland, Fabian Fournier, Lior Limonad, Inna Skarbovsky, and Ava JE Swevels. How well can large language models explain business processes? *arXiv preprint arXiv:2401.12846*, 2024.

[21] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362, 2022.

[22] Keith Ferrazzi. Technology can save onboarding from itself. *Harvard Business Review*, March 2015.

[23] Fabian Friedrich, Jan Mendling, and Frank Puhlmann. Process model generation from natural language text. In *Advanced Information Systems Engineering: 23rd International Conference, CAiSE 2011, London, UK, June 20-24, 2011. Proceedings 23*, pages 482–496. Springer, 2011.

[24] Michael Grohs, Luka Abb, Nourhan Elsayed, and Jana-Rebecca Rehse. Large language models can accomplish business process management tasks. In *International Conference on Business Process Management*, pages 453–465. Springer, 2023.

[25] Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models, 2024.

[26] Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*, 2023.

[27] Mordor Intelligence. Business process management market - size, share and industry analysis, 2024.

[28] Mordor Intelligence. Robotic process automation market - size, share and industry analysis, 2024.

[29] Yixing Jiang, Jeremy Irvin, Ji Hun Wang, Muhammad Ahmed Chaudhry, Jonathan H Chen, and Andrew Y Ng. Many-shot in-context learning in multimodal foundation models. *arXiv preprint arXiv:2405.09798*, 2024.

[30] Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh, and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web, 2024.

[31] Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback. *arXiv preprint arXiv:2312.14925*, 2023.

[32] Nataliia Klievtsova, Janik-Vasily Benzin, Timotheus Kampik, Juergen Mangler, and Stefanie Rinderle-Ma. Conversational process modeling: Can generative ai empower domain experts in creating and redesigning process models?, 2024.

[33] Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks, 2024.

[34] Volodymyr Leno, Artem Polyvyanyy, Marlon Dumas, Marcello La Rosa, and Fabrizio Maria Maggi. Robotic process mining: vision and challenges. *Business & Information Systems Engineering*, 63:301–314, 2021.

[35] Xing Han Lù, Zdeněk Kasner, and Siva Reddy. Weblinx: Real-world website navigation with multi-turn dialogue. *arXiv preprint arXiv:2402.05930*, 2024.

[36] Amy Maitland, Ross Fowkes, and Stuart Maitland. Can chatgpt pass the mrcp (uk) written examinations? analysis of performance and errors using a clinical decision-reasoning framework. *BMJ open*, 14(3):e080558, 2024.

[37] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore, December 2023. Association for Computational Linguistics.

[38] Jorge Munoz-Gama, Niels Martin, Carlos Fernandez-Llatas, Owen A Johnson, Marcos Sepúlveda, Emmanuel Helm, Victor Galvez-Yanjari, Eric Rojas, Antonio Martinez-Millana, Davide Aloini, et al. Process mining for healthcare: Characteristics and challenges. *Journal of Biomedical Informatics*, 127:103994, 2022.

[39] Vinod Muthusamy, Yara Rizk, Kiran Kate, Praveen Venkateswaran, Vatche Isahagian, Ashu Gulati, and Parijat Dube. Towards large language model-based personal agents in the enterprise: Current trends and open problems. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6909–6921, 2023.

[40] R OpenAI. Gpt-4 technical report. *arXiv*, pages 2303–08774, 2023.

[41] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.

[42] Panopto. Workplace knowledge and productivity report, 2018.

[43] Lucia Rahilly, Melissa Valentine, Brooke Weddle, and Bryan Hancock. Human-centered ai: The power of putting people first, Dec 2023.

[44] Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control, 2023.

[45] Hope Reese. A human-centered approach to the ai revolution, 2022.

[46] Lars Reinkemeyer. Process mining in action. *Process Mining in Action Principles, Use Cases and Outloook*, 2020.

[47] Grand View Research. Business process management (bpm) market size, share report 2030, 2024.

[48] Grand View Research. Robotic process automation market size, share report 2030, 2024.

[49] Yara Rizk, Praveen Venkateswaran, Vatche Isahagian, Austin Narcomey, and Vinod Muthusamy. A case for business process-specific foundation models. In *International Conference on Business Process Management*, pages 44–56. Springer, 2023.

[50] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

[51] Wil Van Der Aalst and Wil van der Aalst. *Data science in action*. Springer, 2016.

[52] Wil MP Van der Aalst. Process mining in the large: a tutorial. *Business Intelligence: Third European Summer School, eBISS 2013, Dagstuhl Castle, Germany, July 7-12, 2013, Tutorial Lectures 3*, pages 33–76, 2014.

[53] Maxim Vidgof, Stefan Bachhofner, and Jan Mendling. Large language models for business process management: Opportunities and challenges. *arXiv preprint arXiv:2304.04309*, 2023.

[54] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.

[55] Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*, 2023.

[56] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.

[57] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

[58] Michael Wornow, Avanika Narayan, Krista Opsahl-Ong, Quinn McIntyre, Nigam H Shah, and Christopher Re. Automating the enterprise with foundation models. *arXiv preprint arXiv:2405.03710*, 2024.

[59] Danny TY Wu, Nikolas Smart, Elizabeth L Ciemins, Holly J Lanham, Curt Lindberg, and Kai Zheng. Using ehr audit trail logs to analyze clinical workflow: a case study from community-based ambulatory clinics. In *AMIA Annual Symposium Proceedings*, volume 2017, page 1820. American Medical Informatics Association, 2017.

[60] Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement. *arXiv preprint arXiv:2402.07456*, 2024.

[61] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*, 2024.

[62] Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin, Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oguz, et al. Effective long-context scaling of foundation models. *arXiv preprint arXiv:2309.16039*, 2023.

[63] An Yan, Zhengyuan Yang, Wanrong Zhu, Kevin Lin, Linjie Li, Jianfeng Wang, Jianwei Yang, Yiwu Zhong, Julian McAuley, Jianfeng Gao, et al. Gpt-4v in wonderland: Large multimodal models for zero-shot smartphone gui navigation. *arXiv preprint arXiv:2311.07562*, 2023.

[64] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents, 2023.

[65] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.

[66] Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, et al. Ufo: A ui-focused agent for windows os interaction. *arXiv preprint arXiv:2402.07939*, 2024.

[67] Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey, 2023.

[68] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.

[69] Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.

# A Dataset

## A.1 License & Availability

We license our code and dataset under the Apache 2.0 license. The authors bear all responsibility in case of violation of rights. Our code and data are available here: ⬡ https://github.com/HazyResearch/wonderbread.

Our dataset is based on the excellent WebArena benchmark [69], which also has an Apache 2.0 license and is available here: `https://github.com/web-arena-x/webarena`

## A.2 Dataset Curation

**1. Workflow Selection**. We begin with the WebArena [69] benchmark, which is a collection of 812 workflows instantiated from 187 workflow intents. For example, the template "Search for (term)" could have instantiations "Search for *jacket*" and "Search for *coat*". These 812 tasks require an agent to navigate fully functional open source clones of popular websites. In this dataset we use the e-commerce, content management system (Adobe Magneto), forum (PostMill), and developer tool (GitLab) sites provided by WebArena. We find that many workflows in WebArena are designed to be impossible, are *de facto* impossible, are underspecified, or have incorrect evaluations, and we purposely exclude these workflows from our dataset. We also ignore WebArena tasks that include multiple websites, leaving us with a total of 598 workflows.

**2. Annotator Recruitment and Training**. We enlisted 13 human annotators from a pool of approximately 60 applicants (all students at Stanford University) to participate in our data collection process. All selected annotators, who are undergraduate or graduate students at Stanford University with proficient computer literacy skills, were fully informed and consented to the publication of their complete demonstrations. They were also given the opportunity to review the entire codebase, experiments, and manuscript prior to submission. Annotators were aware that their full screen recordings would be made public and were advised to remove any personally identifiable information before recording. Prior to applying, they were informed that there would be no monetary compensation, as their participation would be on a voluntary basis for a research project.

An important distinction from the demonstrations contained in our dataset versus prior work is that our annotators were explicitly instructed not to perform "zero-shot" recordings, meaning annotators were told to rehearse each task before recording to ensure that the collected demonstrations were free of mistakes. More specifically, annotators were told to follow these principles:

- We are simulating **expert** users of the interface.
  - Do the optimal (i.e. most direct) way to complete each task.
  - Ensure that your demonstration contains no wasted clicks / typing.
  - Ensure that your demonstratoin has no mistakes – If you make a mistake while performing the demonstration, stop recording and re-record from scratch.
- We want a **clean** dataset
  - When you record, ensure that the selected interface within Google Chrome is visible.
  - Ensure you do not show any other applications.
  - Ensure you do not show personal information.

Therefore, the final dataset has a 100% task completion rate. In contrast, in the original WebArena benchmark [69], untrained human annotators could only complete 78% of tasks.

**3. Data Collection**. Each annotator utilized a custom Python script to record demonstrations of approximately 300 unique tasks. This script operated in the background while the annotator completed the demonstration, capturing and outputting four primary types of data: (1) a JSON trace detailing all user actions (clicks, keystrokes, and scrolls), including the precise HTML state of the website at the time of each action and attributes of the elements interacted with; (2) a video of the full screen recording of the annotator's computer; (3) a collection of screenshots corresponding to each recorded action; and (4) an initially blank Standard Operating Procedure (SOP) file.

Once the recording was complete, each annotator filled out the SOP file, creating a detailed, step-by-step list of the actions they performed. Annotators were directed to explain these steps with

the simplicity and clarity necessary for a five-year-old to follow. The annotators were instructed to provide the level of detail that a 5-year-old would need to complete the task. Finally, annotators assessed the difficulty of each task, classifying them as Easy, Medium, or Hard. On average, each annotator dedicated approximately 30 hours to this process, amounting to a collective total of nearly 300 man-hours of labeling over several months.

**4. Demonstration Ranking**. After completing the dataset collection process, we chose a subset of 162 tasks (all derived from different task templates) to form our collection of "Gold Tasks". Each annotator was then tasked with watching the demonstrations of approximately 15 "Gold Tasks", relatively ranking the demonstrations of the same task from 1 (best) to 5 (worst). The annotators then developed a more thorough SOP we call a "Gold SOP" based on the demonstration that received the top ranking. This process resulted in 162 tasks in our dataset containing demonstrations of ranked relatively quality, along with high quality "Gold SOPs" we use as the highest quality SOP representation of the "Gold Task"'s demonstrations. More details about this ranking procedure are included in Appendix A.5.

**5. Quality Assurance**. A key contribution of **WONDERBREAD** is high quality human task demonstrations. A review of existing benchmarks for web navigation tasks found consistently low quality demonstrations that have inaccurate annotations (e.g. misplaced bounding boxes for HTML elements) [58]. This made quality assurance a key concern while curating **WONDERBREAD**. We performed three rounds of quality assurance checks over the course of two months using a combination of automated scripts, manual review, and cross-referencing demonstrations across annotators. We had annotators redo any tasks that were of insufficient quality, and discarded any tasks that had less than 4 successful demonstrations. Additional details are available in the Appendix A.3.

**6. Workflow Understanding Questions**. To enable deeper evaluations of a model's workflow understanding, we also created a set of 11 free responses question templates, which are listed in Appendix A.4. These questions attempted to simulate actual inquiries that a BPM consultant might ask. Examples include *"Explain what the most common failure modes might be for a user performing this task"* and *"Why does the user click the "Commits" button in step #5?"*. We created 10 instances of all question templates, and an additional 10 instances for question template #2. This gives a total of 120 questions. We then had had a set of annotators write brief free-form answers based on the corresponding task.

## A.3   Quality Assurance

We ran a series of automated scripts to flag systematic errors, and had our annotators redo any tasks that were flagged. For example, we verify that all actions occur within Google Chrome and that major disagreements between annotators on each task are resolved. For example, we cross-reference task demonstrations across annotators and redo tasks where someone marked it as infeasible while someone else marked it as feasible. We also conduct manual review of all demonstrations corresponding to the 179 Gold tasks, as well as a random sampling of 300 other demonstrations across all tasks.

## A.4   Question Answering Dataset Questions

Listed below are the free response questions templates that we created for our Question Answering task, largely inspired by prior work in the process mining literature [20, 9].

1. Explain what the most common failure modes might be for a user performing this task.

2. Here are two demonstrations, one of which is more efficient than the other. Please describe ways to improve the less optimal workflow.

3. How would a user completing the task know that the workflow is completed?

4. What is the purpose of doing this workflow?

5. What if instead of X we wanted to do Y. How would you change this workflow to accomplish that?

6. Why does the user click the button X in step #Z?

7. Why does the user click the button X in screenshot #Y?

8. Why does the user type the string X in step #Z?

9. Why does the user type the string X in screenshot #Y?

10. Here are two workflows. Please identify the key differences between them.

11. Given the following three concatenated workflows, how would a user completing the workflow about X know that the workflow is completed?

## A.5 Factors for Quality of Gold SOPs

Listed below is the information given to annotators to aid them with writing high-quality Gold SOPs.

1. **Coverage of edge cases** – help the user complete the task by making note of ways in which the interface might change, and how to adapt:

   ○ e.g. If a task involves looking through a table of shipping orders to find a specific order, and your specific order just happens to be the first one, you should still make a note that the user might have to scroll / paginate through the results until they find the correct shipping order.

   ○ e.g. If you need to click a button at the bottom of a page, you should not assume that the user's browser window has the same size as yours, so you should let them know that they might need to scroll down if they can't see the button.

   ○ **Example**: Instead of "Click on the toggle labeled 'Enable Product'", you "should write "Look for the toggle labeled "Enable Product" which should be directly below the "Quantity" field. If the toggle is currently green, that means the product is currently enabled, which means you should click the toggle in order to disable the product. The toggle should change to a grey color to indicate the product is disabled. However, if the toggle is already greyed out, then do nothing since the product had already been disabled."

2. **Detailed localization of UI elements** – let the user know exactly where to find the element

   ○ e.g. "Click the 'Go to Result' button" is not sufficient. You must be extremely detailed in your specification of each element, i.e. its relative position on the screen, its proximity to other landmark elements, its color, what type of element it is, etc.

   ○ **Example**: Instead of "Click the 'Edit' link", you should write "Click on the blue "Edit" link at the far righthand side of the row corresponding to the "Configurable Product" we previously found."

3. **Generalizability** – the instructions should be written so that they could apply to any instantiation of the **Intent Template** corresponding to the task

   ○ e.g. The instructions should be written generally, providing task-specific information as asides.

   ○ **Example**: Instead of "Type "Out of Office" in the "What's your status?" input box.", you should write "Type the desired Gitlab status in the "What's your status?" input box. In this case, we should type "Out of Office""

4. **Explanations of each action** – briefly explain why we take each step (in the context of the next action, or the larger task)

   ○ e.g. What is the point of each individual action?

   ○ Example: Instead of "Click the "From" text field", you should write "Click the "From" text field to focus it."

   ○ Example: Instead of "Click on the toggle labeled 'Enable Product'", you should write "Click on the toggle labeled 'Enable Product' to disable the product."

# B Benchmark Tasks

For clarity, we define the following notation: Our dataset contains a set of workflow demonstrations $\mathcal{D}$. Each demonstration $d \in \mathcal{D}$ is defined as $d = (\text{W}, \text{SOP}, (s_1, a_1, s_2, a_2, ..., a_{n-1}, s_n))$ where W is the "Intent" or the natural language description of the workflow being done, SOP is a manually written step-by-step guide describing the steps taken in the demonstration, $s_i$ is the $i$th state of the webpage (i.e. a screenshot extracted from the screen recording of the demonstration), and $a_i$ is the action taken at state $s_i$ (i.e. a 'click', 'keystroke', or 'scroll' event extracted from the trace). There are multiple demonstrations $d$ for each workflow, so W is not unique. However, SOP, V, and $(s_1, a_1, s_2, a_2, ..., a_{n-1}, s_n))$ are unique across different demonstrations.

Table 6: Tasks in **WONDERBREAD**. Here, "Demo" can include some combination of an intent (W), a SOP, a recording (V), screenshot key frames $(s_1, ..., s_n)$, and/or action trace $(a_1, ..., a_{n-1})$ for that demonstration.

| Task | Input | Output | Eval | Multi-modal | Multiple Demos |
|------|-------|--------|------|:-----------:|:--------------:|
| **Documentation** | | | | | |
| SOP Generation | 1 Demo | SOP | LLM | ✔ | – |
| Demo Segmentation | 2+ Demos | Clustering | ARI | ✔ | ✔ |
| **Knowledge Transfer** | | | | | |
| Question Answering | Question & 1+ Demos | Free text | LLM | ✔ | ✔ |
| Demo Validation | 1 Demo with SOP | Binary label | F1 | ✔ | – |
| **Improvement** | | | | | |
| Demo Ranking | 3+ Demos | Ranking | Kendall $\tau$ | ✔ | ✔ |
| SOP Improvement | 1 Demo & SOP | SOP | LLM | ✔ | – |

## B.1 Documentation

These subtasks assess a model's ability to generate documentation for existing workflows.

1. **SOP Generation Description:** Given specified components of a workflow demonstration, the model is tasked with generating a new SOP that documents the steps of that workflow. This evaluates a model's ability to generate written documentation.

   **Input:** Given a demonstration $d = (\text{W}, \text{SOP}, \text{V}, (s_1, a_1, s_2, a_2, ..., a_{n-1}, s_n))$, we provide the model with either $(\text{W})$, $(\text{W}, (s_1, ..., s_n))$, or $(\text{W}, (s_1, a_1, ..., a_{n-1}, s_n))$.

   **Output:** An new SOP denoted as $s'$ describing the steps of demonstration $d$.

   **Evaluation:** Pairwise per-line comparison between $s$ and $s'$ that determines the precision and recall as described in Appendix Section C.2

2. **Demonstration Segmentation** Given multiple demonstrations from separate workflows concatenated into a single sequence, identify when each demonstration starts and ends. This evaluates the model's ability to disambiguate between different workflows occurring in sequence.

   **Input:** A concatenated sequence of $k$ demonstrations $\{d^i\}_{i=1}^k$, represented as either $(s_1^1, ..., s_n^1 || ... || s_1^k, ..., s_n^k)$ or $(s_1^1, a_1^1, ..., a_{n-1}^1, s_n^1 || ... || s_1^k, a_1^k, ..., a_{n-1}^k, s_n^k)$.

   **Output:** For each frame $s$ in the provided input, assign each of the frames to one of the $k$ demonstrations. This generates a clustering that maps frames to demonstrations. For example, given 20 frames from three demonstrations $(A, B, C)$, an output assignment clustering might map frames 1-5 to demonstration $A$, frames 6-10 to demonstration $C$, and frames 11-20 to demonstration $B$.

   **Evaluation:** Given the $k$ clusters of frames, measure the adjusted rand score.

## B.2 Knowledge Transfer

These subtasks assess a model's ability to apply knowledge of workflows in practical scenarios.

1. **Question Answering -** Given a question about one or more workflow demonstrations, generate a natural language answer.

   **Input:** A brief question (instantiated from one of the templates in Appendix A.4), and one or two demonstrations, where each demonstration is represented as either $(SOP)$ or $(s_1, a_1, ..., a_{n-1}, s_n)$.

   **Output:** A natural language answer to the question.

   **Evaluation:** Using GPT-4-as-a-judge, compare a human-written reference answer to the generated answer and determine a score for specified criteria on a scale from 1 (bad) to 3 (good). Specified criterion include **completeness** (the response fully answers the question), **soundness** (the response is logically consistent), **clarity** (the response is unambiguous), and **compactness** (the response is concise).

2. **Demonstration Validation -** Given a demonstration and SOP, determine whether (a) the workflow was successfully completed; and (b) whether the demonstration exactly followed the steps of the SOP. For (b), it is not sufficient to merely complete the workflow, but the steps taken to complete it must align with its corresponding SOP.

   **Input:** For (a) we create "positive" examples by sampling full sequences of $(s_1, a_1, ..., a_{n-1}, s_n)$ from our dataset, and create "negatives" by truncating some sequences by a random number of frames to get $(s_1, a_1, ..., s_{k-1}, s_k)$ where $k < n$. Given this sequence, we prompt the model to provide a binary assessment of whether the workflow was completed or not. For (b), we create "positives" by sampling full sequences $(s_1, a_1, ..., a_{n-1}, s_n)$ from our dataset, then and either (a) randomly shuffle or (b) randomly delete frames from this sequence to generate "negative" examples. We prompt the model with this sequence and the SOP, and have it output a binary assessment of whether the sequence exactly followed the SOP.

   **Output:** For (a), a binary assessment of whether the given sequence was truncated. For (b), a binary assessment of whether the given sequence exactly followed the steps of its associated SOP.

   **Evaluation:** Binary classification metrics (ie. Accuracy, F1-Score).

## B.3   Improvement

These subtasks evaluate a model's capacity to improve a given workflow's efficiency.

1. **SOP Ranking -** Given a set of SOPs written by different human annotators for the same workflow, rank the SOPs in order of quality.

   **Input:** A set of $k$ SOPS $\{SOP^i\}_{i=1}^k$ written by different annotators for the same workflow.

   **Output:** A ranking of the quality of the SOPs from $1...k$, where $1$ is best and $k$ is worst.

   **Evaluation:** Given a provided ground truth ranking, determine the Spearman correlation and Kendall's Tau between the predicted ranking and the ground truth.

2. **SOP Improvement -** Given a demonstration and low-quality SOP, and a rubric, generate an improved SOP that better captures what is shown in the demonstration.

   **Input:** One demonstration $d^1$, a low quality $SOP'^1$ generated by a human and an SOP generation rubric $r$.

   **Output:** An improved SOP $SOP'^2$ that better aligns with the provided rubric.

   **Evaluation:** LLM-based evaluation, where the model generates a rating of 1.0 - 5.0 conditioned upon a rubric.

# C  Evaluation

## C.1  Compute

We rely on the publicly available APIs for each of the multimodal FMs we benchmark in this report: GPT-4, Claude3 Sonnet, and Gemini Pro. Thus, we did not require any GPUs to run our benchmark. In terms of cost, the Gemini Pro 1 API was free to use, the Claude 3 API cost roughly $400 in credits, and the GPT-4 API cost roughly $1,000 in credits.

## C.2  LLM-Based Evaluation

### SOP Generation

The automated evaluation for the *SOP Generation* task utilized a pairwise per-step comparison operating over the generated new SOP and the reference high quality SOP. Through a series of iterative prompts, GPT-4 was tasked to identify if the intention of a step in the new SOP was encapsulated in any step of the reference SOP and vice versa. The record of which steps were not included in the alternative SOP were then utilized to calculate the per-step precision, recall, and F1-score.

The precision (P), recall (R), and F1-score (F1) are calculated as follows:

$$P = \frac{TP}{TP + FP}$$

$$R = \frac{TP}{TP + FN}$$

$$F1 = 2 \times \frac{P \times R}{P + R}$$

Where:

- $TP$ (True Positives) is the number of steps in the new SOP that correctly map steps in the reference SOP.
- $FP$ (False Positives) is the number of steps in the new SOP that do not map to any step in the reference SOP.
- $FN$ (False Negatives) is the number of steps in the reference SOP that do not map to any step in the new SOP.

For the *SOP Generation* task, we found that our LLM-based evaluator was able to achieve high correlation out of the box with human raters as shown in Appendix Table 9. We hypothesize that this is because the *SOP Generation* evaluation task is set up to only require the model to make a binary decision over an atomic fact, rather than assess the quality of an open-ended question as in the *Question Answering* task, as seen in other works on LLM-based evaluations [37].

### Question Answering

We rate each answer on a scale from 1 (bad) to 3 (good) on the following four criteria: **completeness** (the response fully answers the question), **soundness** (the response is logically consistent), **clarity** (the response is unambiguous), and **compactness** (the response is concise). Our original LLM-based evaluators had low correlation with human raters – an average Pearson correlation of 0.56 for scoring free reponses questions on a scale of 1 (low quality) to 3 (high) across the four axes of soundness, completeness, clarity, and compactness. We noticed that GPT-4 tended to be overly generous in its ratings. Adding a 3-shot example to our evaluation prompt (one for each possible score) and refining the prompt to "score harsher" helped increase the average correlation with human raters by 54% (to 0.86), as shown in Appendix Table 8.
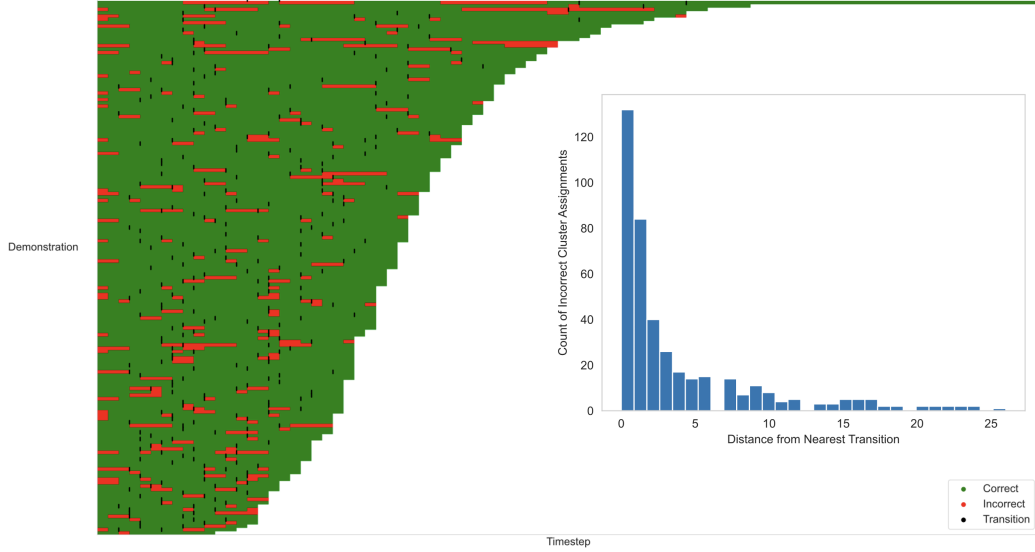
Figure 8: **Demo Segmentation:** Results from GPT-4 evaluated on $k = 3$ concatenated workflows when provided the workflow intent, SOP, and keyframes. (Outset) Each row represents a concatenated sequence of frames from 4 demonstrations. Green line segments are frames that were classified as belonging to the correct task. Red segments are incorrectly classified frames. Black markers indicate a transition between tasks in the ground truth sequence. (Inset) The distribution of distances between each incorrect frame prediction and its closest transition point. Its heavy right skew indicates that the transitions between workflows are where most errors occur, but that GPT-4 is typically able to recover within 3 frames into a workflow.

## D Additional Results

Table 7: **Knowledge Transfer:** Average scores across all 4 evaluation axes for question answering.

| Model | Completeness | Soundness | Clarity | Compactness | Average Score |
|---|---|---|---|---|---|
| Claude3 Sonnet | 1.56 | 1.83 | 2.18 | 2.61 | 2.05 |
| Gemini Pro 1 | 1.81 | 2.15 | 2.83 | 2.95 | 2.44 |
| GPT-4 | 2.20 | 2.51 | 2.96 | 2.85 | 2.63 |
| Human | 3.00 | 3.00 | 2.64 | 2.88 | 2.88 |

Table 8: **Knowledge Transfer:** Correlation between GPT-4 and human-based evaluation based on 60 randomly sampled question-answer pairs.

| Criteria | Pearson Corr. | Pearson p-value | Spearman Corr. | Spearman p-value |
|---|---|---|---|---|
| Completeness | 0.84 | 5.38e-09 | 0.86 | 1.12e-09 |
| Soundness | 0.92 | 1.51e-12 | 0.88 | 2.34e-10 |
| Clarity | 0.80 | 1.01e-07 | 0.80 | 1.01e-07 |
| Compactness | 0.89 | 2.07e-13 | 0.89 | 7.41e-11 |

Table 9: **SOP Generation:** Correlation between GPT-4 and human-based evaluation of the precision/recall of generated SOPs based on 30 randomly sampled examples.

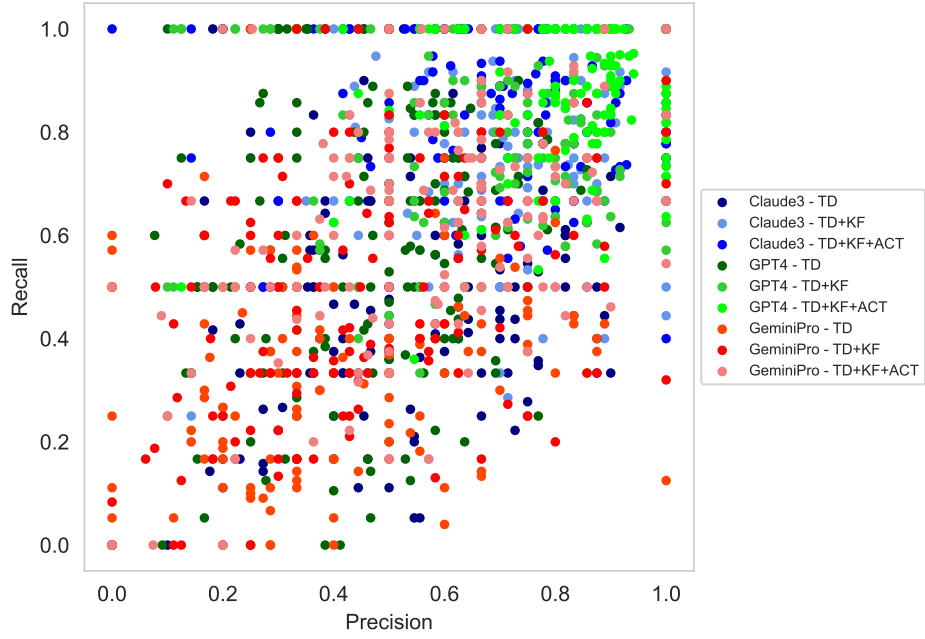| Criteria | Pearson Corr. | Pearson p-value | Spearman Corr. | Spearman p-value |
|---|---|---|---|---|
| Precision | 0.84 | 4.63e-09 | 0.85 | 2.80e-09 |
| Recall | 0.88 | 1.63e-10 | 0.82 | 3.97e-08 |

Figure 9: **SOP Generation:** Each point is an individual SOP. Higher and to the right is better. GPT-4 tends to excel at identifying all steps in a demonstration (i.e. higher recall) but hallucinates inaccurate or superfluous steps (i.e. lower precision).
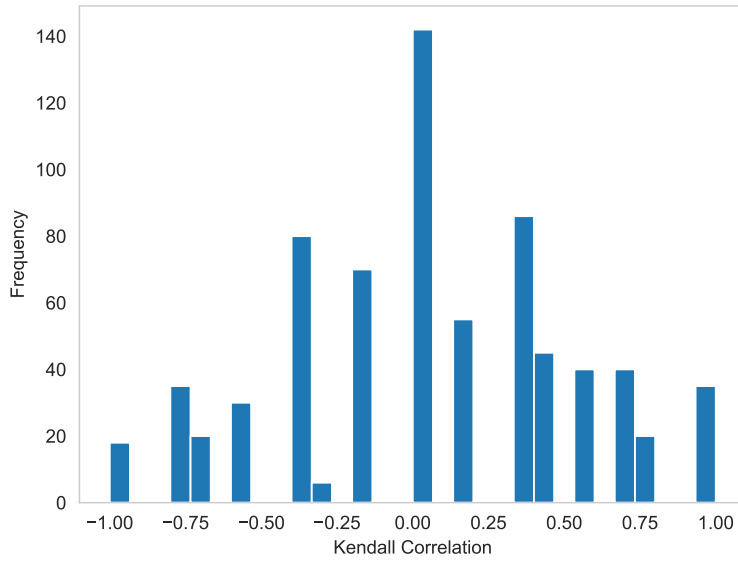


Figure 10: **SOP Ranking:** Ranking demos based solely on SOPs is essentially random
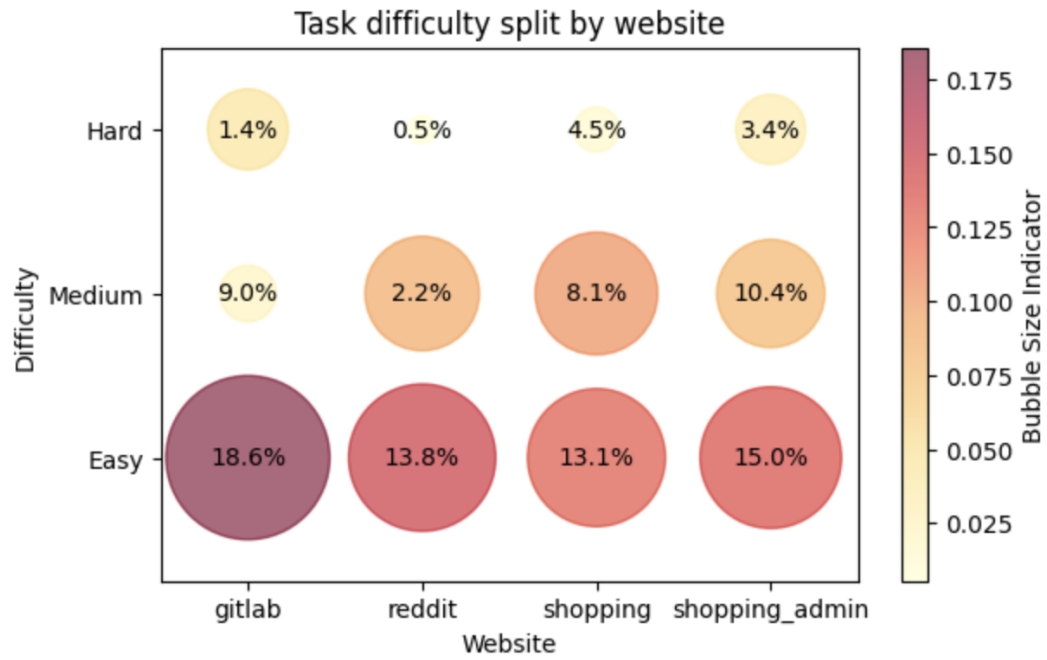
## D.1 Overall Dataset Stats



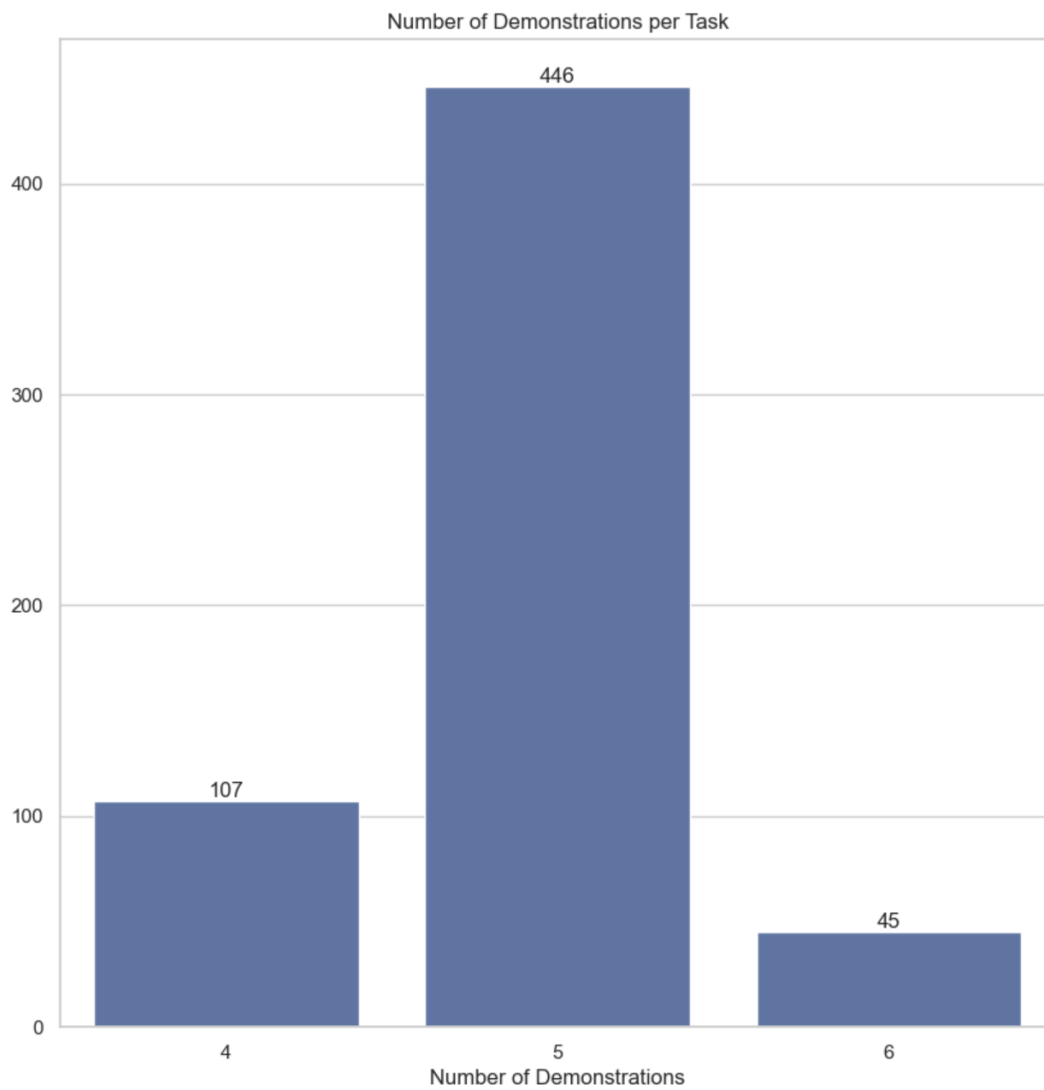Figure 11: Distribution of task difficulty across websites.
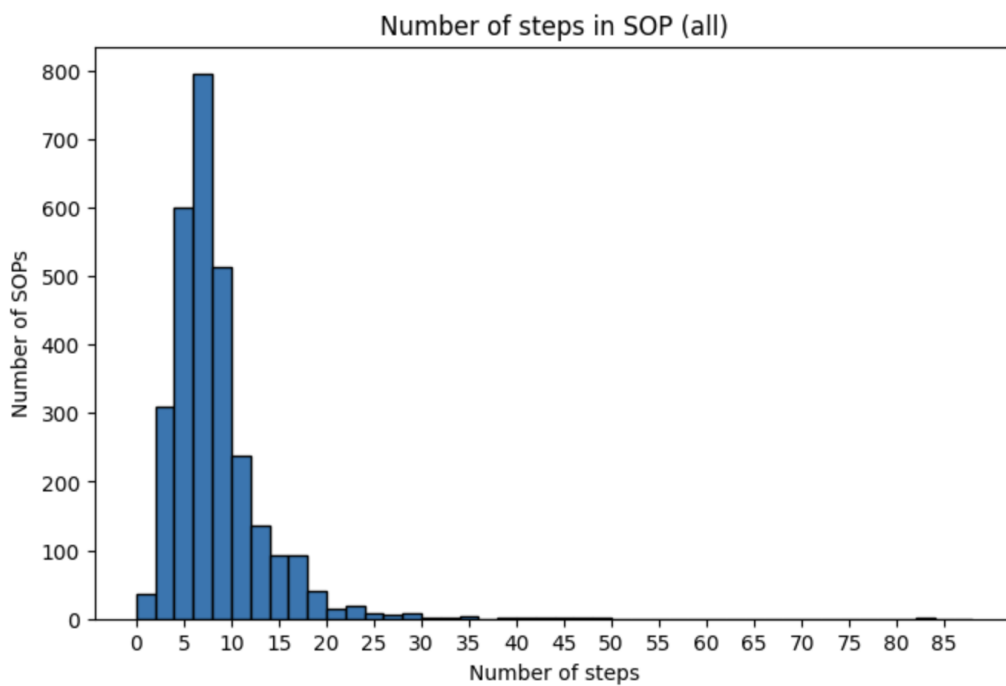
Figure 12: Number of demonstrations per task

Figure 13: Number of steps in SOP per demonstration



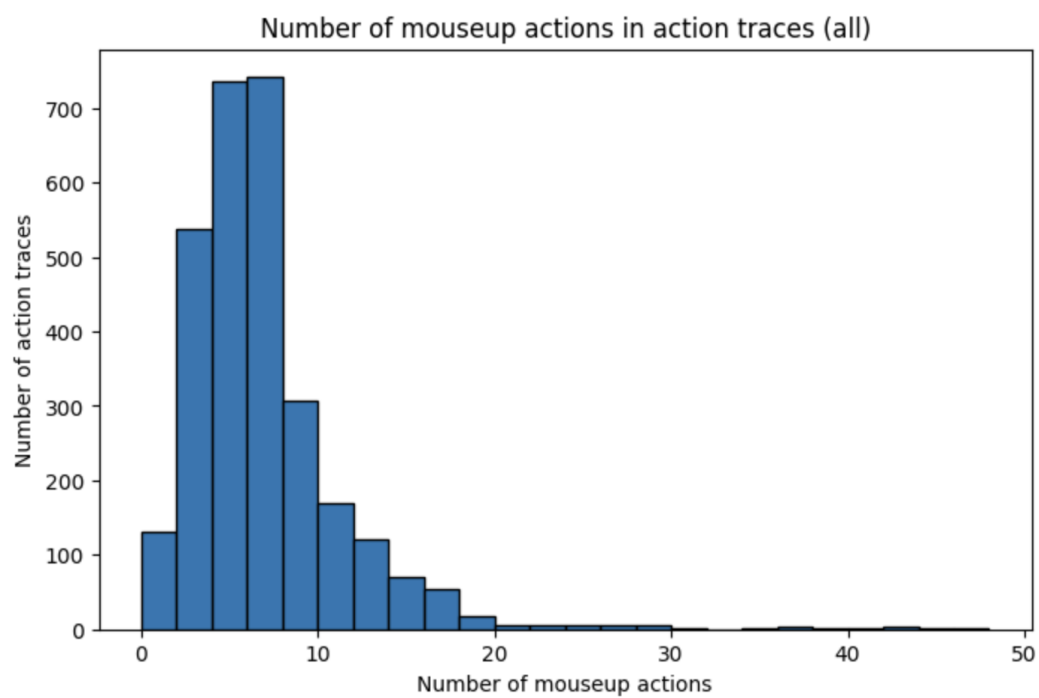Figure 14: Length of video recording (in seconds) per demonstration
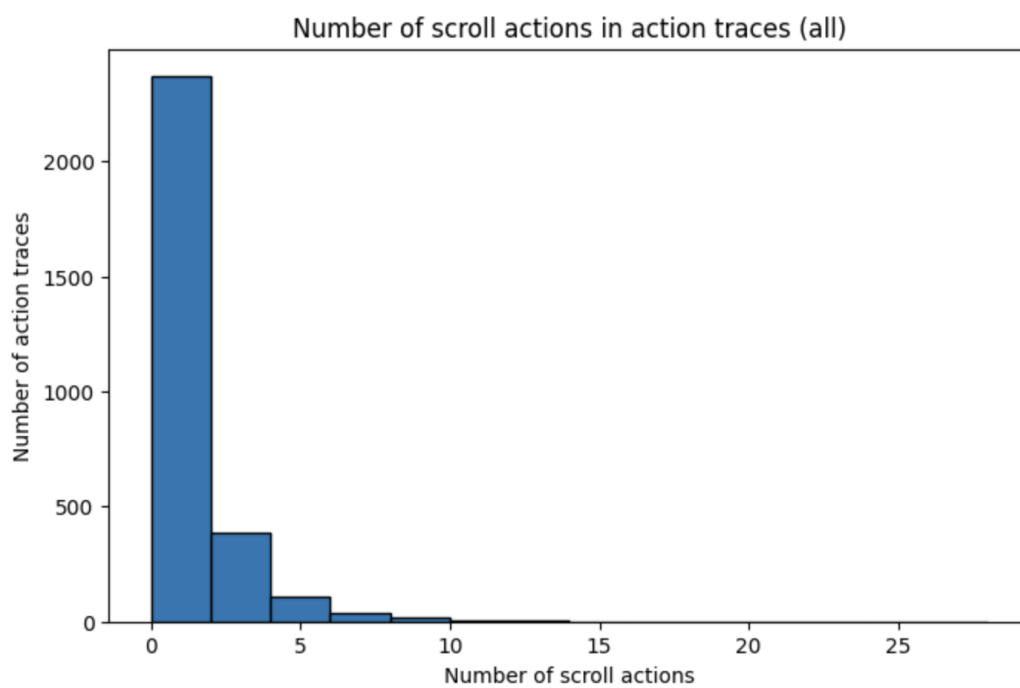
Figure 15: Number of clicks per demonstration

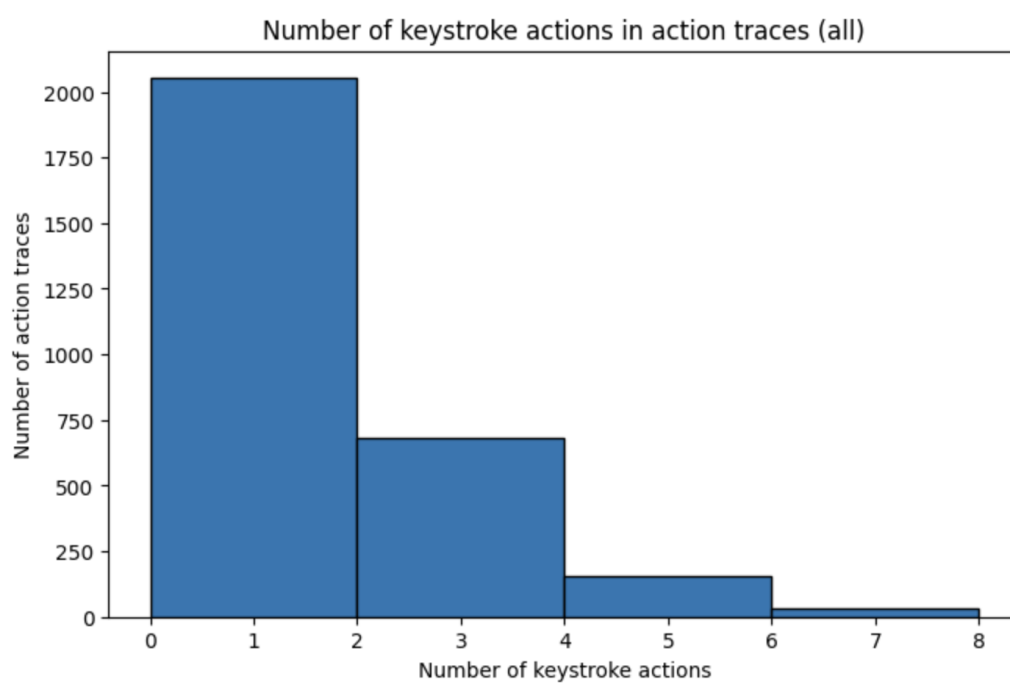

Figure 16: Number of scrolls per demonstration

Figure 17: Number of keystrokes per demonstration
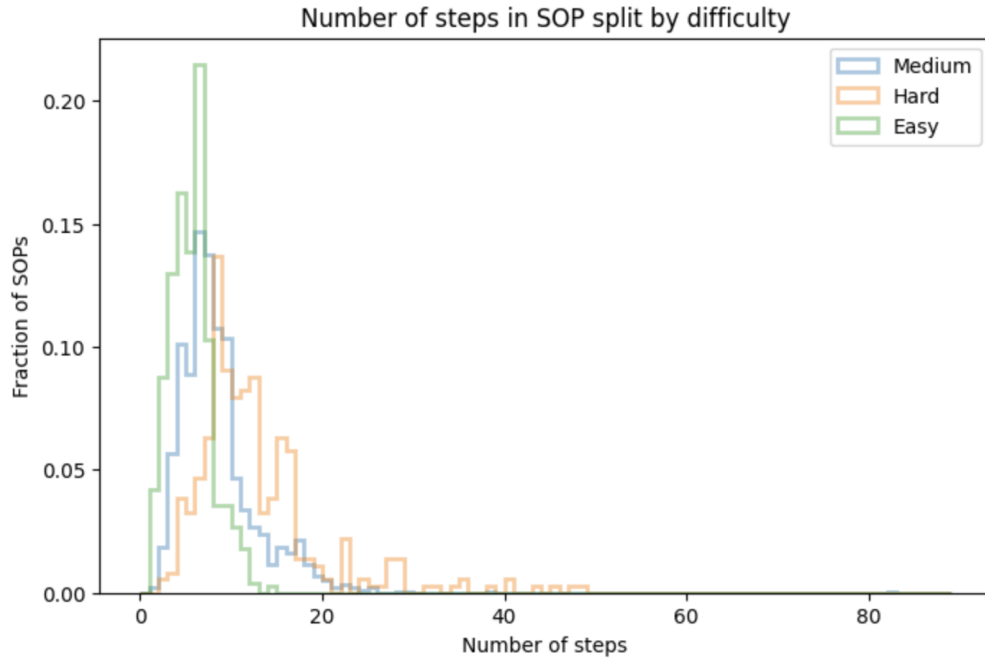
## D.2 Dataset Stats, Split By Difficulty



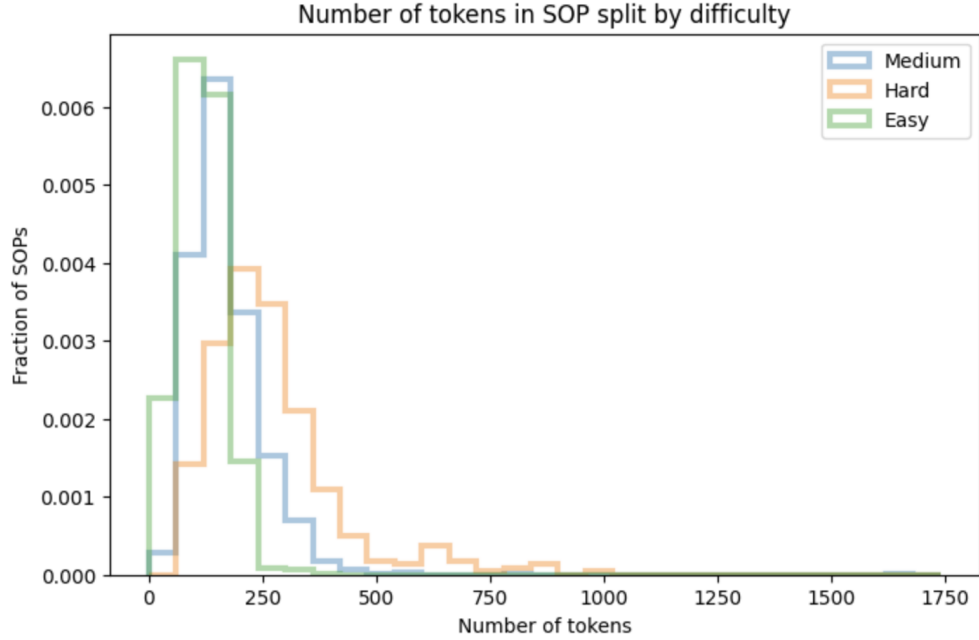Figure 18: Number of steps per SOP per demonstration, split by task difficulty



Figure 19: Number of tokens per SOP per demonstration, split by task difficulty
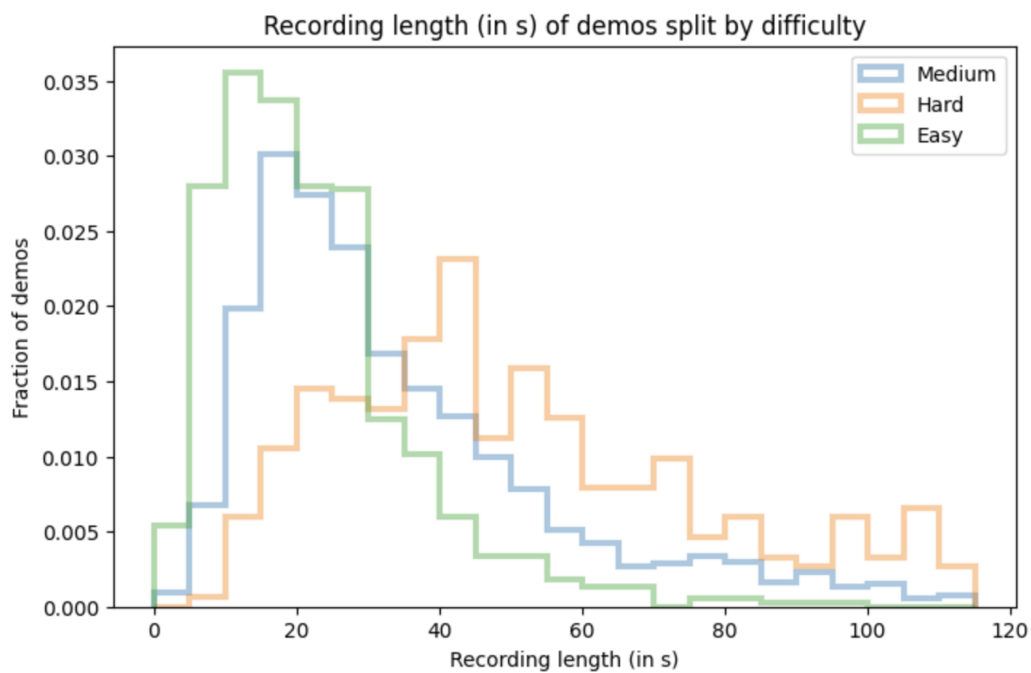
Figure 20: Length of video recording (in seconds) per demonstration, split by task difficulty
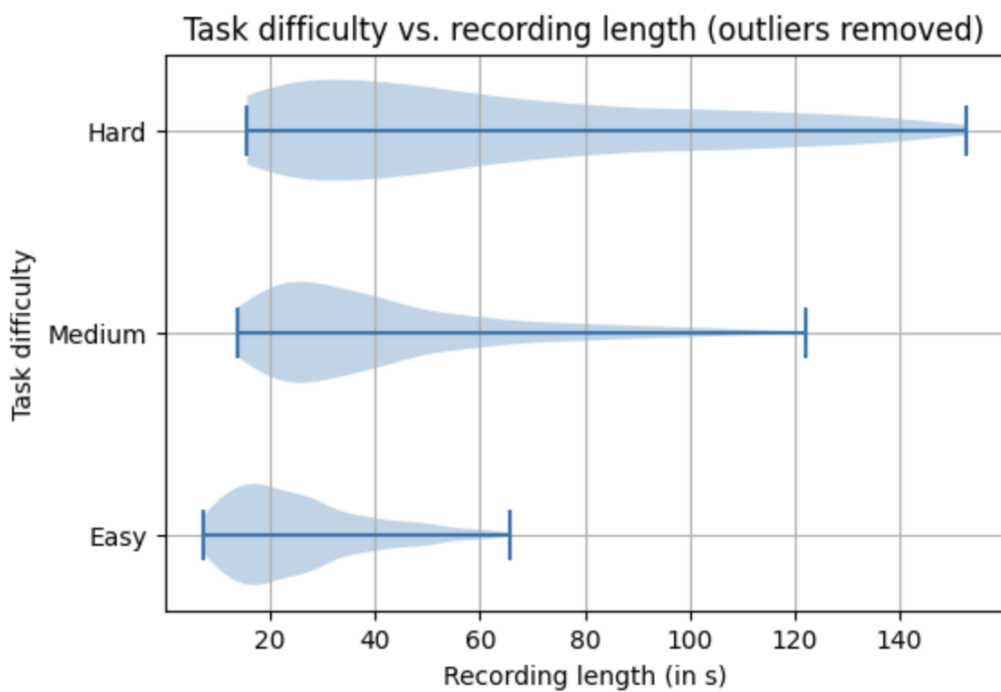


Figure 21: Length of video recording (in seconds) per demonstration, split by task difficulty

Figure 22: Number of actions per demonstration, split by task difficulty

| Difficulty | Min | Median | Max |
|---|---|---|---|
| Medium | 1 | 7 | 82 |
| Hard | 2 | 10 | 48 |
| Easy | 1 | 5 | 14 |

Table 10: Number of steps per SOP, split by task difficulty

| Difficulty | Min | Median | Max |
|---|---|---|---|
| Medium | 12 | 154 | 1631 |
| Hard | 62 | 240 | 976 |
| Easy | 18 | 114 | 382 |

Table 11: Number of tokens per SOP, split by task difficulty

## D.3 Dataset Stats, Split By Website

| Website | Min | Median | Max |
|---|---|---|---|
| shopping_admin | 30 | 163 | 704 |
| gitlab | 12 | 151 | 870 |
| shopping | 18 | 121 | 1631 |
| reddit | 43 | 148 | 382 |

Table 12: Number of tokens per SOP, split by website

| Website | Min | Median | Max |
|---|---|---|---|
| shopping_admin | 0 | 6 | 29 |
| gitlab | 0 | 6 | 44 |
| shopping | 1 | 4 | 47 |
| reddit | 2 | 6 | 23 |

Table 13: Number of mouseups per demonstration, split by website

| Website | Min | Median | Max |
|---|---|---|---|
| shopping_admin | 0 | 1 | 8 |
| gitlab | 0 | 1 | 7 |
| shopping | 0 | 0 | 7 |
| reddit | 0 | 1 | 8 |

Table 14: Number of keystrokes per demonstration, split by website

| Website | Min | Median | Max |
|---|---|---|---|
| shopping_admin | 0 | 0 | 6 |
| gitlab | 0 | 0 | 7 |
| shopping | 0 | 0 | 3 |
| reddit | 0 | 0 | 1 |

Table 15: Number of keypresses per demonstration, split by website

| Website | Min | Median | Max |
|---|---|---|---|
| shopping_admin | 0 | 1 | 13 |
| gitlab | 0 | 0 | 9 |
| shopping | 0 | 1 | 28 |
| reddit | 0 | 0 | 5 |

Table 16: Number of scrolls per demonstration, split by website

# E    Instructions for Annotators

The figures below contain the instructions and other training provided to the annotators.
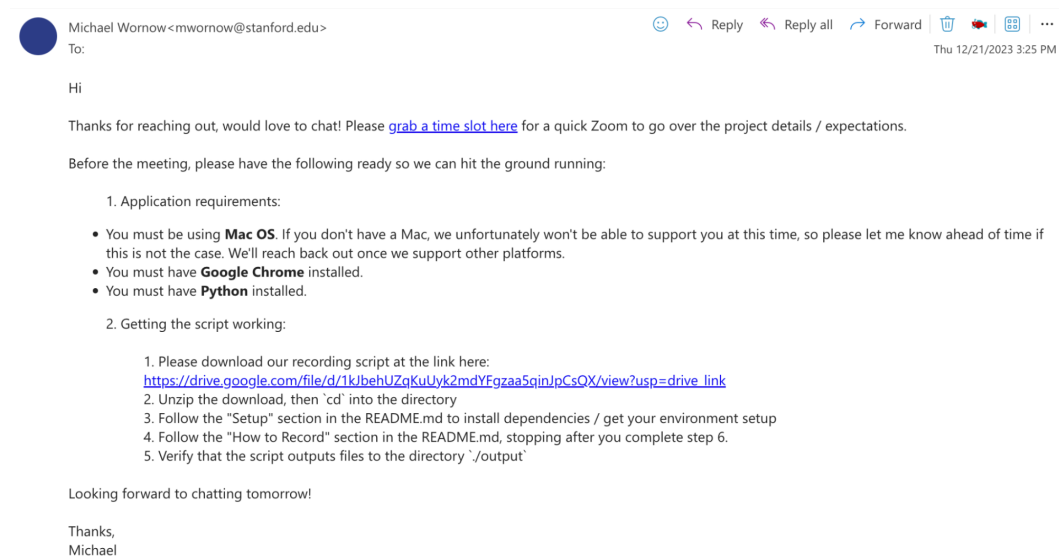


Figure 23: This is the initial onboarding email that asks annotators to set up an online meeting for training.
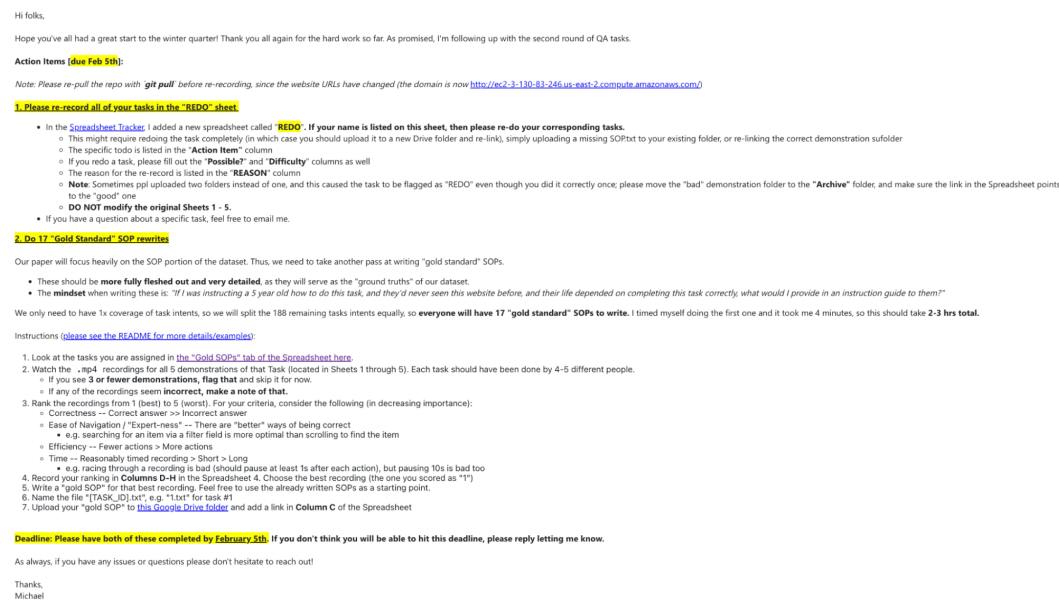


Figure 24: This email provides annotators instructions about how to re-record demonstrations after quality checks. Additionally, there are instructions to write Gold SOPs.

## Few Shot Demonstration Collection

**Goal:** Create human demonstrations for the tasks in WebArena.

**Principles:**

- We are simulating **experts**
  - Do the optimal (i.e. most direct) way to complete each task
  - No wasted clicks / typing.
  - No mistakes -- If you mis-step, then stop recording and re-record from scratch. You might need a couple rehearsals before you record for real.
- We want a **clean** dataset
  - When you record, **only have Google Chrome visible**.
    - Do not show any other Applications
    - No personal information
    - No messy desktop background
  - We recommend having Google Chrome in one monitor, and the rest of your stuff in another monitor.
- We want an **accurate** dataset, even though our **record.py** script has a slight lag.
  - Make sure that **record.py** has finished logging your action before you move onto the next action.
    - **Clicks** - usually no lag, so no need to worry
    - **Keystrokes** - usually no lag, so no need to worry
    - **Scrolls** - lags often b/c each 1px of scroll = 1 action, so make sure the script has "caught up" before you move onto the next action

**Disclaimers:**

- This only works on Mac currently.
- Please email mwornow@stanford.edu if you have any issues / questions / errors.

We will be publishing these recordings, so please make sure any personal information is hidden from view when recording. By participating in this project, you consent to having your recordings published.

## 📋 Tasks

All of the tasks are stored in `./tasks/`.

## 🔧 Setup

1. Enable the following Mac permissions:

   a. `System Preferences > Privacy & Security > Accessibility`, make sure VSCode and Terminal are enabled.

   b. `System Preferences > Privacy & Security > Screen Recording`, make sure VSCode and Terminal are enabled.

   c. `System Preferences > Privacy & Security > Input Monitoring`, make sure VSCode and Terminal are enabled.

2. Install **ffmpeg** with: `brew install ffmpeg`

3. Download this repo:

```
# Install repo
git clone https://github.com/Miking98/demonstration-collection.git
cd demonstration-collection/

# Create conda env + install dependencies
conda create -n demo_env python=3.10 -y
conda activate demo_env
pip3 install -r requirements.txt
pip3 install -e .
```

## ⚡ Quickstart

In a background terminal:

```
alias google-chrome="/Applications/Google\ Chrome.app/Contents/MacOS/Google\ Chrome"
google-chrome --remote-debugging-port=9222 --user-data-dir="/tmp/chrome_temp"
```

In another terminal:

```
conda activate demo_env
python record.py --is_webarena --name <TASK_ID>
```

Hit `Esc` to end the recording.

Then...

1. Fill out `SOP.txt`
2. Upload folder to the Google Drive
3. Fill out your task's row in the Google Spreadsheet tracker.

Figure 25: This screenshot provides instructions on how to get the environment and technology setup before recording demonstrations.

## 🌐 Websites

- Gitlab
  - URL: http://ec2-3-130-83-246.us-east-2.compute.amazonaws.com:8023
  - Username: byteblaze
  - Password: hello1234
- Shopping
  - URL: http://ec2-3-130-83-246.us-east-2.compute.amazonaws.com:7770
  - Username: emma.lopez@gmail.com
  - Password: Password.123
- Shopping Admin
  - URL: http://ec2-3-130-83-246.us-east-2.compute.amazonaws.com:7780/admin
  - Username: admin
  - Password: admin1234
- Reddit
  - URL: http://ec2-3-130-83-246.us-east-2.compute.amazonaws.com:9999
  - Username: MarvelsGrantMan136
  - Password: test1234

## 🙇 How to Record (long version)

The scripts below will automatically generate a trace for the task you demonstrate. We will keep track of each task in this Google Spreadsheet tracker and upload our recordings to this Google Drive folder.

1. Setup *Google Chrome* to run in debug mode:

```
# creates an alias for running Google Chrome
alias google-chrome="/Applications/Google\ Chrome.app/Contents/MacOS/Google\ Chrome"
# Opens up Chrome with port 9222 open for Python webdriver to attach to
google-chrome --remote-debugging-port=9222 --user-data-dir="/tmp/chrome_temp"
```

If you get a popup when you open this *Debugger Google Chrome* for the first time, uncheck both boxes and then click OK.

To simplify this in the future, you should add the line `alias google-chrome="/Applications/Google\ Chrome.app/Contents/MacOS/Google\ Chrome"` to your `~/.zshrc` file.

2. Login to all the relevant websites for your task in your *Debugger Google Chrome*. Credentials for each site are listed above. Basically, just go to each URL you care about in your *Debugger Google Chrome*, and enter the Username / Password.

For your initial setup, feel free to just do the **Gitlab** website, and then as you do other tasks you'll need to login to each website.

3. Make sure your *Debugger Google Chrome* only has one tab open before recording.

4. Run the **record.py** script in our repo:

```
python3 record.py --is_webarena --name <TASK_ID>
```

where `TASK_ID` is one of the `<TASK_ID>.json` files in `./tasks/`, e.g. an integer between [0, 811].

For your initial setup, use `python3 record.py --name 811` to do a **Gitlab** task.

5. Perform the task.

   - The task is defined in **<TASK_ID>.json** as the value for the key `intent`.
   - The **record.py** script will automatically jump to the proper start URL, and it will record the contents of your entire screen (so I'd recommend moving your Chrome window to a separate monitor to keep a clean background).
   - Try to not go too fast through the steps of the workflow.
   - Note that the **record.py** script takes a few seconds to get loaded, so wait until your console prints out `>>>>>>>> GOOD TO START RECORDING WORKFLOW <<<<<<<<<<` before you start taking actions (otherwise they will be ignored).

6. Hit `Esc` to end the recording. Outputs will be written to a folder in `./outputs/`

7. Within the empty file `SOP.txt` that was created in the `./outputs/` folder, manually write down a list of steps (i.e. an "SOP") for your task, based on what you did to accomplish the task. Examples can be found below.

8. Upload your task's output in the `./output/` folder to the Google Drive.

9. Fill out your task's row in the Google Spreadsheet tracker.

### Examples

Some good examples can be found here:

- Task 810 -- Google Drive link
- Task 811 -- Google Drive link

### Random Tips

1. Do not use the "Back" or "Forward" buttons on Google Chrome, or any other buttons of Google Chrome. Limit all interactions to the website itself.
2. Only Google Chrome should be visible and focused during your screen recording (i.e. no switching between Applications). Our script records your entire screen, so make sure it looks clean.
3. We are using a shared public website. This means other people are completing these tasks, and writing to the database. For some tasks, you may need to "reset" the website back to its initial start state before you record.

Figure 26: This screenshot provides instructions on the different websites and how to record.

📝 **How to Write an SOP**

When writing your SOP, try to be as explicit as possible. Each step should be a single, discrete action.

Imagine you were instructing a 4 year old to do this task. You can assume you are already at the proper starting page of the website.

For example:

**Bad**

```
1. Navigate to the GitLab group creation page by clicking on the "+" icon in the top navigation bar and selecting "N
```

**Good**

```
1. Navigate to the GitLab group creation page by clicking on the "+" icon in the top navigation bar
2. Select the "New group" from the dropdown menu.
```

Here are few examples SOPs:

```
Invite a collaborator to a Gitlab repository
1. Search for the repository in the search bar at the top of the page.
2. Click on the tag corresponding to the repository.
3. Click on the "Project Information" dropdown in the left sidebar of the project's main page.
4. Click on the "Members"  link in the left sidebar under the  Project Information  dropdown.
5. Click on the "Invite members" button located at the top right of the  Project members section
6. In the "Username or email address" field, type the collaborator's name
7. Click on the option corresponding to the collaborator in the dropdown menu.
8. Click the "Invite" button to send an invitation.
```

```
Find top-n bestsellers in a given period
1. Click on the "Reports" section on the left sidebar.
2. Under the "Product" section click on "Bestsellers" to open the Bestsellers Report.
3. Set the period of time for which you want to see the bestsellers.
4. In the "From" field enter start date in mm/dd/yy form.
5. In the "To" field enter end date in mm/dd/yy form.
6. Click on the "Show Report" button to generate the report.
7. Given the information, state the top-n bestsellers in the period in question.
```

Figure 27: This screenshot explains how to write a SOP.

## 🏆 How to Write a "Gold" SOP

A "gold" SOP should be significantly more fleshed out and detailed than the SOPs that have been written so far.

These will serve as the "ground truths" of our dataset.

The mindset you should be in when writing these is:

> "If I was instructing a 5 year old how to do this general task, and they'd never seen this website before, and their life depended on completing this task correctly, what information would I give them?"

Just like the normal SOPs, the "gold SOPs" must meet the following checklist:

1. **Task Description:** Include the task description at the top of the SOP.txt file. There should be 2 new lines between task description and the steps below.
2. **Element Specification:** Each element mentioned should have a...
   - Descriptive name (i.e., "Accounting tab under the Finances Section"
   - Descriptive location (i.e., "Accounting tab in the left hand side-bar")
3. **Action Specification:** Each step should...
   - Only reference the following actions: **Press, Delete, Click, Type, Scroll**; no hover!
   - Contain one discrete action (i.e., the step *"Search 'hi'"* is not a discrete action -- decompose it into three separate steps *"Click on the search bar"*, then *"Type 'hi'"*, then *"Press Enter"*)
   - Cover edge cases (i.e., if you don't see the button then scroll down)
   - Match the task (i.e. if the task is "Add Harry to the repo", then you should type "Harry" into the searchbar and not "John")
4. **Faithful to task recording:** The SOP should *exactly* match the steps in the `.mp4` recording.
5. **Generality:** Steps should be general, while also customizing to the specific task (i.e., "Click on the row corresponding to your order. In this case, that is the row with ID 000334")
6. **Formatting:**
   - Names of elements should be in quotes (i.e., click on the "Accounting" button)
   - Steps in SOP should be properly numbered (don't skip numbers or double count)
   - End each line with a period
   - Proper capitalization and spelling

For these "gold" SOPs, we want to put a special emphasis on:

1. **Coverage of edge cases** -- help the user complete the task by making note of ways in which the interface might change, and how to adapt
   - e.g. if a task involves looking through a table of results, and your Shipping Order just happens to be the first one, you should still make a note that the user might have to scroll / paginate through the results until they find the correct Shipping Order.
   - e.g. if you need to click the button at the bottom of the page, you should not assume that the user's browser window has the same size as yours, so you should let them know that they might need to scroll down if they can't see the button.
   - *Example:* Instead of *"Click on the toggle labeled 'Enable Product'"*, you *"should write "Look for the toggle labeled "Enable Product" which should be directly below the "Quantity" field. If the toggle is currently green, that means the product is currently enabled, which means you should click the toggle in order to disable the product. The toggle should change to a grey color to indicate the product is disabled. However, if the toggle is already greyed out, then do nothing since the product had already been disabled."*
2. **Detailed localization** of UI elements -- let the user know exactly where to find the element
   - e.g. "Click the 'Go to Result' button" is not sufficient. You must be extremely detailed in your specification of each element, i.e. its relative position on the screen, its proximity to other landmark elements, its color, what type of element it is, etc.
   - *Example:* Instead of *"Click the 'Edit' link"*, you should write *"Click on the blue "Edit" link at the far righthand side of the row corresponding to the "Configurable Product" we previously found."*
3. **Generalizability** -- write these instructions so that they could apply to any instantiation of the **Intent Template** corresponding to your task
   - e.g. Write your instructions generally, then provide your task-specific stuff as asides.
   - *Example:* Instead of *"Type "Out of Office" in the "What's your status?" input box."*, you should write *"Type the desired Gitlab status in the "What's your status?" input box. In this case, we should type "Out of Office""*
4. **Explanations** of each action -- briefly explain why we take each step (in the context of the next action, or the larger task)
   - e.g. What is the point of each individual action?
   - *Example:* Instead of *"Click the "From" text field"*, you should write *"Click the "From" text field **to focus it**."*
   - *Example:* Instead of *"Click on the toggle labeled 'Enable Product'"*, you should write *"Click on the toggle labeled 'Enable Product' **to disable the product**"*

Here is an example of a "gold SOP":

```
What is the top-1 best-selling product in 2022?

1. Click on the "Reports" button on the far lefthand sidebar. It has an icon which looks like a chart. It should be
2. In the popup menu that appears, click on the "Bestsellers" link to go to the "Bestsellers Report" page. The link :
3. Locate the field labeled "Period" and click on the dropdown menu to reveal our time options.
4. Click on the "Year" option to set the reporting period to Year.
5. Click on the "From" textbox to focus it. It should be located directly underneath the "Period" field.
6. Type in the first day of our desired time period, which in this case is "01/01/2022". Make sure the textbox is emp
7. Click on the "To" textbox to focus it. It should be located directly underneath the "From" field.
8. Type in the last day of our desired time period, which in this case is "12/31/2022". Make sure the textbox is emp
9. Click on the orange "Show Report" button, which can be found on the top right of the page, in order to generate o
10. The best-selling products will appear in a table at the bottom of the page. Scroll down if you cannot see the fu
```

Note the detailed localization for each UI element (*"'Reports' button on the far lefthand sidebar. It has an icon which looks like a chart..."*) the edge case coverage (*"Make sure the textbox is empty before you type into it."* and *"Scroll down if you cannot see the full table"*) and generalizability (*"Type in the first day of our desired time period, which in this case is "01/01/2022"."*).

Figure 28: This screenshot is the first part of instructions on how to write a gold SOP.

**Instructions:**

Follow these steps when writing your "gold SOPs":

1. Look at the tasks you are assigned in the Spreadsheet here.
2. Watch the `.mp4` recordings for all 5 demonstrations of that Task (located in Sheets 1 through 5). Each task should have been done by 4-5 different people.
    - If you see **3 or fewer demonstrations, flag that** and skip it for now.
    - If any of the recordings seem **incorrect, make a note of that.**
3. Rank the recordings from 1 (best) to 5 (worst). For your criteria, consider the following (in decreasing importance):
    - Correctness -- Correct answer >> Incorrect answer
    - Ease of Navigation / "Expert-ness" -- There are "better" ways of being correct
        - e.g. searching for an item via a filter field is more optimal than scrolling to find the item
    - Efficiency -- Fewer actions > More actions
    - Time -- Reasonably timed recording > Short > Long
        - e.g. racing through a recording is bad (should pause at least 1s after each action), but pausing 10s is bad too 4. Record your ranking in **Columns D-H** in the Spreadsheet 4. Choose the best recording (the one you scored as "1") 5. Write a "gold SOP" for that best recording. Feel free to use the already written SOPs as a starting point. 6. Name the file "[TASK_ID].txt", e.g. "1.txt" for task #1 6. Upload your "gold SOP" to this Google Drive folder and add a link in **Column C** of the Spreadsheet

# Output Format

Recordings are stored in two output formats:

1. `trace.json` - a JSON log of all actions and states
    - Actions
        - Clicks: (x,y) on screen, HTML element (if on a webpage), timestamp
        - Keystrokes: Key, HTML element (if on webpage), timestamp
    - States
        - Name of active desktop application (e.g. "Google Chrome", "Terminal", etc.)
        - Bounding boxes of accessible HTML elements (if on webpage)
2. `trace.mp4` - a full screen recordings of you completing this task
3. `screenshots/` - a folder containing screen shots of each state (taken from still frames in `trace.mp4` corresponding to the timestamps of each state)
4. `SOP.txt` - A manually generated summary of the steps of this process, based on the task description and screenshots

# Todos

Divergence from `record.py` in eclair:

1. `state.to_json()` auto-serializes `json_state` as string in eclair, but that's not true here.
2. In eclair, `execute_scripts` has been renamed `execute_js_scripts`, and now executes `proxy-select.js` (so that dropdowns appear in playwright/selenium). Here, we don't do `proxy-select.js`.
3. **Important:** Need to run `trace_cleanup.py` in eclair to fix (x,y) coords
4. Add `is_focused`, `is_checked`, etc. to elements

# Bugs

- Fix `window_position` and `window_size` bugs; also adjust `x,y` coords to account for browser chrome/position within `element_attributes` of `action`.
- Might want to expand `json_state` to be more inclusive, or to always include last element clicked

Figure 29: This screenshot is the second part of instructions on how to write a gold SOP.