# GuideWeb: A Benchmark for Automatic In-App Guide Generation on Real-World Web UIs

**Chengguang Gan**[1]    **Yoshihiro Tsujii**[1]    **Yunhao Liang**[2]
**Tatsunori Mori**[4]    **Shiwen Ni**[3]    **Hiroki Itoh**[1]

[1]Techtouch, Inc.
[2]University of Chinese Academy of Sciences
[3]Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences
[4]Yokohama National University

{chengguang.gan, yoshihiro.tsujii, hiroki.itoh}@techtouch.co.jp
liangyunhao22@mails.ucas.ac.cn    tmori@ynu.ac.jp    sw.ni@siat.ac.cn

## Abstract

Digital Adoption Platform (DAP) provide web-based overlays that deliver operation guidance and contextual hints to help users navigate complex websites. Although modern DAP tools enable non-experts to author such guidance, maintaining these guides remains labor-intensive because website layouts and functionalities evolve continuously, which requires repeated manual updates and re-annotation. In this work, we introduce **GuideWeb**, a new benchmark for automatic in-app guide generation on real-world web UIs. GuideWeb formulates the task as producing page-level guidance by selecting **guide target elements** grounded in the webpage and generating concise guide text aligned with user intent. We also propose a comprehensive evaluation suite that jointly measures the accuracy of guide target element selection and the quality of generated intents and guide texts. Experiments show that our proposed **GuideWeb Agent** achieves **30.79%** accuracy in guide target element prediction, while obtaining BLEU scores of **44.94** for intent generation and **21.34** for guide-text generation. Existing baselines perform substantially worse, which highlights that automatic guide generation remains challenging and that further advances are necessary before such systems can be reliably deployed in real-world settings.

## 1   Introduction

In recent years, rapid progress in large language models (LLMs) (OpenAI, 2025; Anthropic, 2025; Comanici et al., 2025) has accelerated research on LLM-based agents that plan and act to achieve user goals (Huang et al., 2024). Among these, web agents (Ning et al., 2025) have attracted growing attention for their ability to automate interactions with websites and complete complex tasks. While substantial effort has been devoted to benchmarks and methods for action-oriented web agents, an important and widely deployed form of web assistance remains underexplored: in-app guidance for
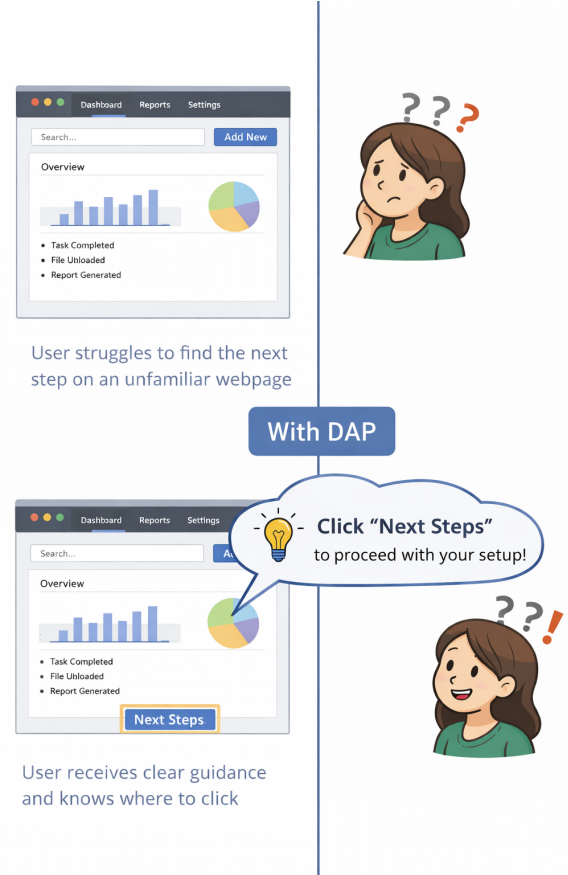


Figure 1: Illustration of a digital adoption platform (DAP) as an in-app overlay for unfamiliar web UIs. **Top:** without guidance, the user cannot confidently identify the correct guide target element to proceed. **Bottom:** with a DAP overlay, the webpage is augmented with contextual hints and step-by-step instructions, enabling the user to complete the action efficiently.

helping users understand and operate unfamiliar web interfaces, as commonly provided by digital adoption platforms (DAPs).

DAPs are typically implemented as browser-based overlays that augment existing web services with contextual hints and step-by-step instructions. As illustrated in Figure 1, users often struggle to identify the correct UI elements to interact with
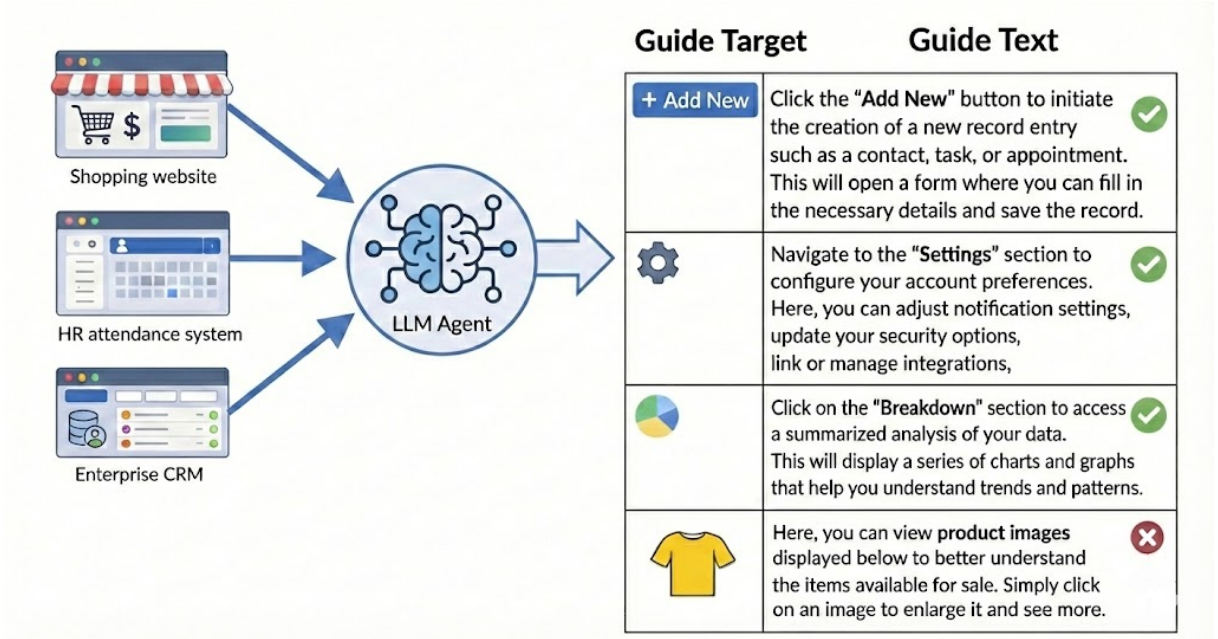
Figure 2: Overview of GUIDEWEB. Given the main page of a real-world website, an LLM-based agent identifies guide targets, namely interactive UI elements whose usage may benefit from guidance, and generates corresponding guide text grounded in visible on-page content. The examples on the right illustrate both correct guides and typical failure cases, where the agent produces low-utility guidance for elements that are already self-explanatory.

when encountering new or complex webpages, especially in enterprise systems with dense functionality and limited visual cues. This difficulty increases learning cost, leads to inefficient usage, and results in higher customer support burden. By annotating key workflows directly on top of webpages, DAPs enable users to quickly locate and understand essential operations, improving productivity and reducing support cost.

Despite their practical importance, existing DAP solutions rely heavily on manual authoring. Human experts must identify which UI elements require guidance and write corresponding instructional text. This process is costly and brittle, as web layouts and functionalities evolve continuously. Layout changes require re-locating guide targets, and newly introduced features require additional guide creation, while outdated guides can even mislead users. These challenges motivate automatic in-app guide generation with LLM-based web agents.

We address this gap by introducing GUIDEWEB, a benchmark for automatic web guide annotation and generation grounded in real-world web UIs, as illustrated in Figure 2. GUIDEWEB focuses on main pages of diverse websites and formulates guide generation as a two-stage task: identifying which interactive UI elements should be guided

and generating concise guide text aligned with user intent. This formulation reflects a key challenge of in-app guidance: not all interactive elements require explanation, and effective guides must be both selective and informative. To support systematic evaluation, we propose a comprehensive evaluation protocol that measures both guide target selection accuracy and guide text quality.

In addition, we train a lightweight and efficient GUIDEWEB AGENT tailored to this benchmark. Experimental results demonstrate that the trained agent substantially outperforms existing baselines across major evaluation metrics, highlighting both the difficulty of the task and the effectiveness of task-specific modeling for automatic in-app guide generation.

## 2 Related Work

**Web Agents and Web Interaction Benchmarks.** Recent advances in large language models have led to rapid progress in web agents that autonomously operate websites to accomplish user goals. Early benchmarks focus on controlled and templated environments, such as Liu et al. (2018), which enables reproducible evaluation over synthetic web interfaces. More recent efforts emphasize realism and scale. Yao et al. (2022) introduces end-to-end shopping tasks grounded in structured product pages,

while Deng et al. (2023) collects diverse real-world trajectories covering a wide range of websites and user intents. Zhou et al. (2023) further advances this direction by hosting interactive websites that require multi-step reasoning and form-based interactions across domains. Multimodal extensions such as Koh et al. (2024) incorporate visual perception to better model real browsing behavior. Collectively, these benchmarks study how an agent perceives a webpage and executes actions to complete a task. In contrast, our work addresses a different but complementary problem: instead of automating user actions, we focus on generating in-app guidance artifacts that help end users understand how to operate unfamiliar web interfaces.

**In-App Guidance and UI-Level Assistance.** Beyond action execution, prior research in human-computer interaction and intelligent assistance has explored how to provide contextual guidance over existing interfaces. Systems such as Zhong et al. (2021) generate visual tutorials aligned with UI states, while approaches like Li et al. (2017) infer task procedures from user demonstrations and UI signals. Other studies investigate interactive walk-throughs, tooltips, and overlay-based explanations to reduce learning cost for complex software. However, these works typically focus on mobile applications, scripted tutorials, or automation by demonstration, and they do not formalize the problem of web-scale guide generation as a benchmarked learning task. Moreover, existing web agent benchmarks do not evaluate whether an agent can identify which UI elements actually require guidance, nor whether the generated explanations are useful to users. GuideWeb fills this gap by introducing a benchmark grounded in real-world web UIs that explicitly decomposes automatic guide generation into guide target identification and guide text generation, together with a comprehensive evaluation protocol for both stages.

## 3 GuideWeb: A Benchmark for Automatic In-App Guide Generation on Real-World Web UIs

We introduce GUIDEWEB, a benchmark designed to study automatic in-app guide generation on real-world web user interfaces. GUIDEWEB focuses on the main pages of diverse websites and aims to evaluate whether an agent can identify interaction points that require guidance and generate appropriate guide text grounded in visible on-page content.

Unlike prior web agent benchmarks that emphasize long-horizon task completion, GUIDEWEB targets a complementary but practically important problem: assisting users in understanding and operating unfamiliar web interfaces through lightweight, contextual guidance.

### 3.1 Task Definition and Output Schema

**Input Representation.** Given the main page of a website, we represent the input as the raw HTML source

$$x \in \mathcal{X}, \tag{1}$$

from which we construct a DOM tree and extract a set of interactive elements

$$\mathcal{E}(x) = \{e_1, e_2, \ldots, e_N\}. \tag{2}$$

Each element $e \in \mathcal{E}(x)$ is associated with observable attributes

$$\phi(e) = (\texttt{tag}(e), \texttt{visible\_text}(e), \texttt{xpath}(e)). \tag{3}$$

**Output and Task.** The goal is to generate a structured guide annotation

$$y = (g, \mathcal{A}), \tag{4}$$

where $g \in \{0, 1\}$ indicates whether the page requires any guidance, and $\mathcal{A}$ is a set of element-grounded guide annotations.

Formally, a GUIDEWEB system is a mapping

$$f : \mathcal{X} \to \{0, 1\} \times \mathcal{Y}, \quad f(x) = (g, \mathcal{A}). \tag{5}$$

**Stage 1: Guide Target Identification.** Conditioned on the page $x$ and its interactive set $\mathcal{E}(x)$, the model selects a subset of elements that require guidance:

$$\mathcal{E}^+ = S(x) \subseteq \mathcal{E}(x). \tag{6}$$

Elements in $\mathcal{E}^+$ are referred to as *guide targets*. If the model predicts $g = 0$, we define $\mathcal{E}^+ = \emptyset$ and $\mathcal{A} = \emptyset$.

**Stage 2: Element-Grounded Guide Generation.** For each guide target $e \in \mathcal{E}^+$, the model generates a structured annotation

$$a(e) = (i(e), t(e), s(e), \phi(e)), \tag{7}$$

where $i(e)$ is a natural-language intent, $t(e)$ is an action type (e.g., search, navigation, selection), and $s(e)$ is the guide text describing how to use the

| Field | Description |
|---|---|
| `site` | Unique identifier of the website. |
| `html_file` | Raw HTML of the main page (`page.html`). |
| `needs_guides` | Boolean indicating whether the page requires guidance. |
| `page_category` | Coarse category of the webpage (e.g., landing, listing). |
| `annotations` | List of guide annotations for selected guide targets. |
|   `intent` | Natural language description of user intent. |
|   `action_type` | High-level action type (e.g., search, navigation). |
|   `guide_text` | Generated textual guidance for the target element. |
|   `tag` | HTML tag of the target element. |
|   `visible_text` | Visible text associated with the element (if any). |
|   `xpath` | XPath locating the target element in the DOM. |

Table 1: Simplified output schema of GUIDEWEB annotations.

element. The tuple explicitly carries $\phi(e)$ so that the annotation is grounded in the original page $x$.

The final output set is

$$\mathcal{A} = \{a(e) \mid e \in \mathcal{E}^+\}. \tag{8}$$

This formulation jointly captures *where* guidance is needed (target selection via $S(x)$) and *what* guidance should be provided (element-conditioned generation of $a(e)$), while grounding all outputs in the input webpage through `xpath` and visible content.

**Output Schema.** Each annotated webpage in GUIDEWEB is stored in a structured JSON format. The schema consists of page-level fields and a list of guide-level annotations, summarized in Table 1.

As summarized in Table 1, each webpage in GUIDEWEB is stored in a dedicated subdirectory. The directory contains the raw HTML file of the webpage main page, as well as a JSON file that stores all annotation-related information. At the page level, the JSON file records whether the webpage requires guidance and includes coarse-grained category information. At the guide level, the file contains a collection of annotated guide targets, each corresponding to an interactive element that requires guidance.

For each guide target, we store the element's visible text, namely the textual content presented to users on the webpage. We note that some interactive elements are purely icon-based and do not expose visible text; in such cases, the element is referenced solely by its XPath. The XPath field enables precise localization of the target element within the DOM, which is essential for downstream DAP systems to reliably attach guidance overlays to the correct UI components.

In addition to grounding information, we annotate each guide target with an action type (e.g., search, login), which captures the high-level functional category of the interaction. The `tag` field records the raw HTML tag type of each guide target element as it appears in the DOM. For example, many interactive controls in real-world websites are implemented using anchor tags (`<a>`), even when they function visually as buttons, tabs, or menu items. Therefore, `tag` reflects the structural implementation of an element rather than its semantic action type, and is used primarily for grounding guide annotations to concrete DOM nodes.

Finally, we associate each guide target with a natural-language intent description. The intent serves two purposes. First, it encourages the LLM to explicitly reason about why a particular guide is needed and what user goal it supports, thereby improving the accuracy of guide target identification. Second, the intent enables downstream DAP systems to proactively confirm whether the user shares the same goal, facilitating faster and more targeted assistance during interaction.

```
{
    "intent": "How do I log in to my account?",
    "action_type": "login",
    "guide_text": "Choose "Login" to go to the sign-in page.",
    "tag": "a",
    "visible_text": "Login",
    "xpath": "//a[@href='https://app.benevity.com/login']"
}
```

Figure 3: A single guide annotation in GUIDEWEB, containing intent, action type, guide text, and DOM grounding fields.

As illustrated in Figure 3, all of the above information is unified within a single JSON file, al-
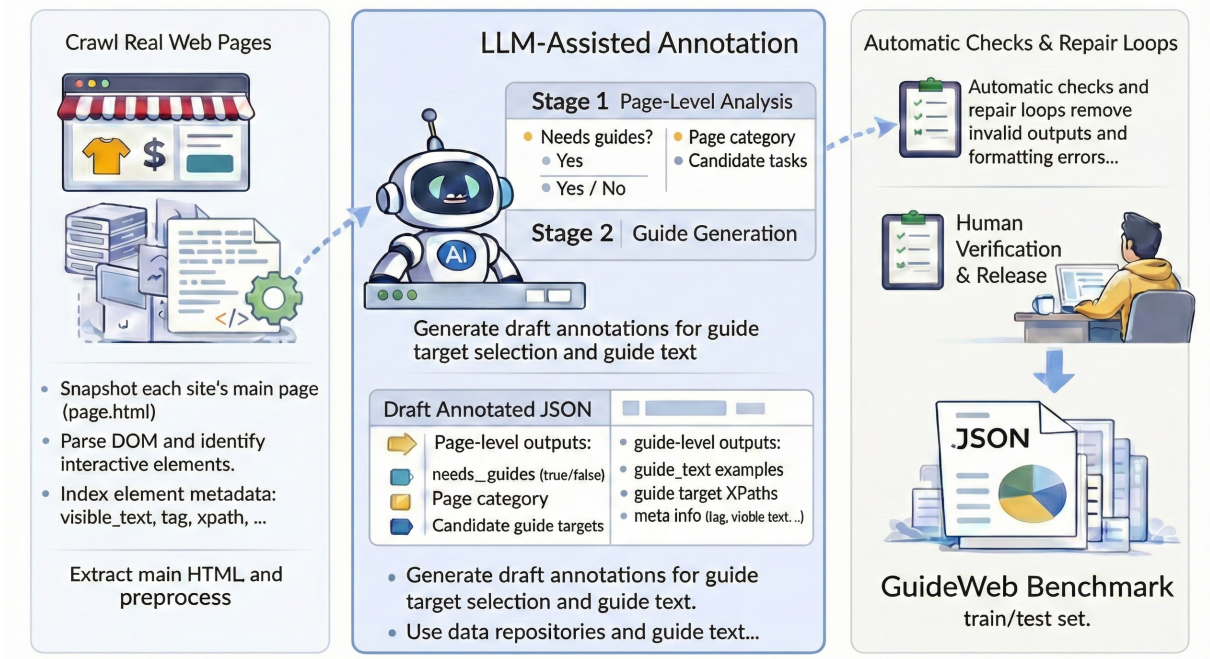
Figure 4: Overview of the GUIDEWEB construction pipeline with LLM-assisted annotation and human verification.

lowing models to generate, store, and consume annotations in a consistent and structured format.

## 3.2 Dataset Construction Pipeline

During dataset construction, we adopt a hybrid annotation strategy that combines LLM-assisted labeling with human verification and correction. The overall pipeline of constructing GUIDEWEB is illustrated in Figure 4. This design allows us to efficiently scale the annotation process while maintaining high annotation quality.

The websites included in GUIDEWEB are sampled from the Cisco Umbrella Popularity List,[1] which ranks domains based on aggregated passive DNS usage observed across Umbrella's global network. Unlike browser-centric rankings such as Alexa, which rely primarily on HTTP traffic collected from browser plugins, the Umbrella list reflects Internet-wide activity across diverse protocols and applications. This property makes it particularly suitable for collecting a broad and realistic set of web services.

We begin with the Top 1M domains provided by the Umbrella list[2] and randomly sample candidate domains. For each selected domain, we use an automated browser to crawl and snapshot its main landing page. We retain only websites whose main pages satisfy minimum structural requirements, such as containing a sufficient number of interactive elements (e.g., clickable controls, input fields, and forms). This filtering step ensures that each retained page presents meaningful opportunities for in-app guidance generation. After filtering, we obtain a final set of 1,000 real-world websites spanning diverse categories, including e-commerce platforms, enterprise systems, and general service websites.

For each selected website, the main page serves as the annotation target. We first store the raw HTML snapshot of the page and parse its DOM structure to identify interactive elements. Each interactive element is indexed using a unified `elements` index, which records its raw HTML tag type (`tag`), visible text content (`visible_text`), and a precise XPath (`xpath`) for reliable localization. Notably, some interactive elements, such as icon-only buttons, do not expose visible text and are therefore indexed using their structural attributes alone. This element indexing step forms the foundation for grounding all subsequent guide annotations to concrete DOM nodes.

The processed and indexed HTML is then provided to an LLM for assisted annotation. At the first stage, the LLM is instructed to perform page-level analysis, including determining whether the page contains interactive elements that merit guidance and assigning a coarse-grained page category.

---

[1] https://umbrella-static.s3-us-west-1.amazonaws.com/index.html

[2] http://s3-us-west-1.amazonaws.com/umbrella-static/top-1m.csv.zip

| Statistic | Value | Statistic | Value |
|---|---|---|---|
| Total crawled websites | 1,000 | Websites removed after verification | 4 |
| Valid annotated websites | 996 | Websites requiring guides (needs_guides=true) | 980 |
| Websites without guides | 16 | Ratio requiring guides | 98.4% |
| Average guides per page | 3.09 | Pages reaching guide cap (5 guides) | 546 (54.8%) |

Table 2: Overall statistics of the GUIDEWEB benchmark (guide-level annotations only).

If the page is judged to require guidance, the LLM proceeds to identify a subset of interactive elements that should be annotated as guide targets. At the second stage, the LLM generates corresponding guide text for each selected target, along with an explicit user intent describing the underlying user goal. In this formulation, the guide text serves as an answer-like explanation, while the intent captures the implicit question a user may have when interacting with the page. All annotations are generated in a predefined JSON schema specified in the prompt. If the model output does not conform to the required format, an automatic regeneration and repair loop is triggered until a valid structured output is obtained.

The LLM-generated annotations are subsequently reviewed and corrected by human annotators with higher education backgrounds. Each annotation is verified against the original webpage to ensure correctness, clarity, and practical usefulness. Through this verification process, we remove four pages with severe structural errors or annotation ambiguities. The final benchmark therefore consists of 996 valid annotated webpages. We split these samples into training and test sets with a ratio of 7.5:2.5, which are used consistently across all experiments.

### 3.3 Dataset Statistics

Table 2 reports high-level benchmark statistics. Starting from 1,000 crawled domains, we obtain 996 validated main-page snapshots after human verification. The vast majority of pages are deemed guide-worthy (needs_guides=true), reflecting that real-world homepages typically expose multiple interactive entry points that benefit from lightweight in-app guidance. On average, each page contains 3.09 guide annotations, and more than half of the pages hit the per-page annotation budget of five guides, indicating that even a conservative cap is frequently saturated in practice. This design choice reflects the observation that densely annotating a page with excessive guidance is neither necessary nor desirable, as it may hinder us-

| Action type | # Guides |
|---|---|
| search | 728 |
| navigate | 520 |
| login | 412 |
| contact_support | 307 |
| other (incl. <50 types) | 296 |
| signup | 199 |
| subscribe_newsletter | 112 |
| start_trial | 107 |
| checkout | 93 |
| pricing | 86 |
| filter_sort | 79 |
| download_install | 75 |
| settings_profile | 64 |

Table 3: Distribution of guide annotations by action type (guide-level counts). Rare types with fewer than 50 guides are merged into other.

ability and visual clarity; the limit of five guides is therefore introduced in an exploratory setting to prioritize the most valuable, frequently used, and user-critical interactions.

Table 3 summarizes the distribution of guide annotations by action type. The benchmark is dominated by information-seeking and navigation behaviors (e.g., search and navigate), followed by common account and support workflows (e.g., login, signup, contact_support). Long-tail categories are aggregated into other to stabilize analysis and avoid over-interpreting sparse types. A more detailed breakdown of page categories and their relationship with guide necessity is provided in Appendix A.

## 4  GuideWeb Agent

In addition to the benchmark, we train a lightweight GUIDEWEB AGENT on the GuideWeb training set for automatic in-app guide generation. To address the challenge of long and noisy webpage inputs, we incorporate a *Shorter* mechanism that removes irrelevant HTML content and compresses the remaining structure before inference. This design substantially reduces input length, leading to faster inference and lower computational cost, while preserving the information necessary for guide target
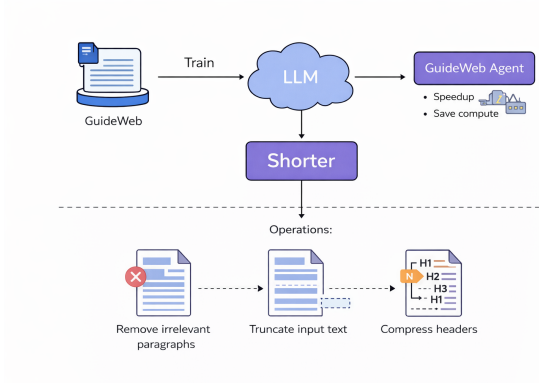
Figure 5: Overview of the GUIDEWEB AGENT with the *Shorter* mechanism. The agent is trained on the GUIDEWEB benchmark to perform automatic in-app guide generation. During inference, the *Shorter* module reduces input length by removing irrelevant content, truncating long text, and compressing headers, enabling faster inference and reduced computational cost while preserving guide-relevant information.

identification and guide text generation.

Empirical results in Section 6 show that the *Shorter* mechanism does not degrade performance and in some cases improves it. These findings suggest that, for web guide generation, the information required by the model is largely concentrated in interactive elements, their visible texts, and nearby contextual descriptions. In contrast, large portions of raw webpage text and auxiliary HTML structures contribute little to the task. This observation supports the use of compact, interaction-focused representations for efficient and effective in-app guide generation.

## 5 Experiments

### 5.1 Evaluation Metrics

We evaluate GUIDEWEB from three complementary aspects: (1) **guide target selection** (which UI elements should be annotated), (2) **text generation quality** (intent and guide text), and (3) **structured grounding accuracy** (field-level exact matches). While the underlying measures are standard, we tailor their evaluation targets to the GUIDEWEB setting.

**Guide target selection (P/R/F1).** For each page, let $G$ be the set of gold guide targets and $\hat{G}$ be the predicted set, where each target corresponds to a unique DOM element. We report precision, recall, and F1 over set overlap (definitions in Appendix B). Intuitively, precision measures how many predicted targets are truly guide-worthy, recall measures how

many gold targets are recovered, and F1 summarizes the trade-off. This metric captures *selective* guide authoring, penalizing redundant or low-utility annotations.

**Intent and guide-text generation (BLEU, ROUGE-L).** For matched guide targets, we compare generated intent strings and guide texts against references using BLEU and ROUGE-L (standard definitions). These metrics quantify lexical overlap and sequence-level similarity, respectively, and reflect whether the agent produces concise, faithful descriptions aligned with the intended interaction.

**Structured field correctness (Exact-match F1).** We additionally evaluate whether predictions are grounded in the webpage structure by computing exact-match F1 for key fields: action_type, tag, visible_text, and xpath. Here, a field value is counted as correct only if it exactly matches the reference. This provides fine-grained diagnostic signals: xpath tests precise localization, whereas tag and visible_text reflect semantic and surface alignment.

### 5.2 Baselines

We evaluate GUIDEWEB using a diverse set of representative large language models as baselines, covering both proprietary and open-source systems. Specifically, we include GPT-5 (Singh et al., 2025), Claude Sonnet 4.5 (Anthropic, 2025), and Gemini 2.5 Pro (Comanici et al., 2025) as strong closed-source LLM baselines, as well as Qwen3-8B (Yang et al., 2025) and LLaMA 3.1-8B (Grattafiori et al., 2024) as competitive open-source counterparts.

All baseline models are prompted to perform guide target identification and guide text generation directly from webpage content under a unified prompt template, without task-specific fine-tuning. This setting reflects a realistic zero-shot or instruction-following deployment scenario and allows us to isolate the inherent capability of general-purpose LLMs on the GUIDEWEB benchmark.

## 6 Results

As shown in Table6, we evaluate all baseline models and the proposed GuideWeb Agent using the comprehensive metric suite introduced in Section 5.1. We first focus on guide target element selection. The GuideWeb Agent achieves the highest F1 score among all models, indicating a substantially better balance between precision and recall.

Table 4: Main results on GuideWeb. Left: guide target element selection. Right: intent and guide-text generation quality.

| Model | (a) Guide target element selection | | | | (b) Text generation quality | | | |
| | P | R | F1 | Match/Pred. | Intent BLEU | Intent ROUGE-L | Guide BLEU | Guide ROUGE-L |
|---|---|---|---|---|---|---|---|---|
| GPT-5 | 15.29 | 54.69 | 23.90 | 356/2328 | 13.84 | 24.99 | 4.30 | 11.69 |
| Claude-Sonnet-4.5 | 15.99 | **58.99** | 25.16 | 384/2402 | 15.88 | 25.99 | 1.48 | 8.98 |
| Gemini-2.5-Pro | 14.43 | 54.84 | 22.85 | 357/2474 | 16.53 | 30.54 | 2.53 | 13.43 |
| Qwen3-8B | 9.11 | 12.33 | 10.48 | 72/790 | 3.69 | 11.54 | 1.93 | 11.47 |
| LLaMA3.1-8B | 10.98 | 2.77 | 4.42 | 18/164 | 4.13 | 13.57 | 2.05 | 11.45 |
| **GuideWeb Agent** | **29.99** | 31.64 | **30.79** | 206/687 | **44.94** | **52.89** | **21.34** | **28.44** |

Although several baseline models exhibit much higher recall than the trained GuideWeb Agent, this does not imply superior performance. Instead, these models achieve high recall by producing a significantly larger number of predictions, effectively relying on over-generation to increase the chance of matching gold guide targets.

This behavior is clearly reflected in the *Match/Predicted* statistics. Compared to the GuideWeb Agent, baseline models generate approximately three to four times more predicted guide targets, yet they only identify about $1.5\times$ more correct matches. Such a strategy inflates recall at the cost of precision and indicates that these models fail to selectively identify guide-worthy elements. Rather than understanding which interactive elements truly require guidance, they tend to mark many elements indiscriminately, resulting in redundant and low-utility guides.

In contrast, for both user intent generation and guide text generation, the trained GuideWeb Agent consistently outperforms all baselines by a large margin across BLEU and ROUGE-L metrics. This demonstrates that task-specific training enables the agent to better align generated intents and guide texts with the underlying page semantics and user needs. Overall, these results highlight that existing general-purpose models struggle to solve the GuideWeb task effectively, and that dedicated modeling and training are crucial for accurate and practical in-app guide generation.

## 7 Ablation of the Shorter Mechanism

Table 5 reports an ablation study on the proposed *Shorter* mechanism using Qwen3-8B as the backbone model. Without Shorter, the model processes raw webpage inputs with substantially longer HTML context, resulting in limited recall and overall F1 performance. After enabling the Shorter mechanism, which selectively preserves

Table 5: Ablation of the Shorter mechanism.

| Qwen3-8B | | | | |
|---|---|---|---|---|
| **Setting** | **P** | **R** | **F1** | **Match/Pred.** |
| w/o Shorter | 9.11 | 12.33 | 10.48 | 72/790 |
| Shorter | **9.52** | **29.67** | **14.42** | 181/1901 |

interactive elements and nearby informative text, recall improves markedly from 12.33 to 29.67, leading to a consistent F1 gain.

Notably, this improvement is achieved without increasing model capacity or modifying training objectives. Instead, Shorter reduces redundant and low-utility HTML content, allowing the model to focus on semantically relevant signals for guide target selection. These results indicate that effective context reduction is critical for GuideWeb tasks, and that long-form raw HTML is not only inefficient but can actively hinder accurate guide prediction.

## 8 Conclusion and Future Work

We introduce GUIDEWEB, the first benchmark for automatic in-app guide generation on real-world web UIs. GuideWeb formalizes guide generation as selecting guide-worthy interactive elements and generating concise, user-aligned guide text. We design a comprehensive evaluation suite and propose a lightweight GUIDEWEB AGENT with a context-shortening mechanism that significantly outperforms strong general-purpose baselines.

Our results show that existing models tend to over-generate guide targets to boost recall, while failing to selectively identify meaningful guidance. In contrast, task-specific training and input structuring enable more precise and informative guide generation. Future work includes extending GuideWeb to multi-step workflows and exploring tighter integration with real-world DAP systems.

## 9 Limitations

This work focuses exclusively on main pages of web applications and does not address multi-page or stateful workflows that require long-horizon planning. In addition, although GuideWeb leverages human verification to ensure annotation quality, the benchmark remains limited in scale compared to fully automated web interaction datasets. Finally, our GuideWeb Agent is trained in an offline setting and does not adapt to user feedback or evolving webpage layouts at inference time. We leave interactive learning, broader page coverage, and online adaptation to future work.

## References

Anthropic. 2025. Claude sonnet 4.5 system card. Technical report, Anthropic.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36:28091–28114.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang,

Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. 2024. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*.

Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Russ Salakhutdinov, and Daniel Fried. 2024. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 881–905.

Toby Jia-Jun Li, Amos Azaria, and Brad A Myers. 2017. Sugilite: creating multimodal smartphone automation by demonstration. In *Proceedings of the 2017 CHI conference on human factors in computing systems*, pages 6038–6049.

Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. 2018. Reinforcement learning on web interfaces using workflow-guided exploration. In *International Conference on Learning Representations (ICLR)*.

Liangbo Ning, Ziran Liang, Zhuohang Jiang, Haohao Qu, Yujuan Ding, Wenqi Fan, Xiao-yong Wei, Shanru Lin, Hui Liu, Philip S Yu, et al. 2025. A survey of webagents: Towards next-generation ai agents for web automation with large foundation models. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pages 6140–6150.

OpenAI. 2025. Gpt-5 system card. Technical report, OpenAI.

Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, Akshay Nathan, Alan Luo, Alec Helyar, Aleksander Madry, Aleksandr Efremov, Aleksandra Spyra, Alex Baker-Whitcomb, Alex Beutel, Alex Karpenko, Alex

Makelov, Alex Neitz, Alex Wei, Alexandra Barr, Alexandre Kirchmeyer, Alexey Ivanov, Alexi Christakis, Alistair Gillespie, Allison Tam, Ally Bennett, Alvin Wan, Alyssa Huang, Amy McDonald Sandjideh, Amy Yang, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrei Gheorghe, Andres Garcia Garcia, Andrew Braunstein, Andrew Liu, Andrew Schmidt, Andrey Mereskin, Andrey Mishchenko, Andy Applebaum, Andy Rogerson, Ann Rajan, Annie Wei, Anoop Kotha, Anubha Srivastava, Anushree Agrawal, Arun Vijayvergiya, Ashley Tyra, Ashvin Nair, Avi Nayak, Ben Eggers, Bessie Ji, Beth Hoover, Bill Chen, Blair Chen, Boaz Barak, Borys Minaiev, Botao Hao, Bowen Baker, Brad Lightcap, Brandon McKinzie, Brandon Wang, Brendan Quinn, Brian Fioca, Brian Hsu, Brian Yang, Brian Yu, Brian Zhang, Brittany Brenner, Callie Riggins Zetino, Cameron Raymond, Camillo Lugaresi, Carolina Paz, Cary Hudson, Cedric Whitney, Chak Li, Charles Chen, Charlotte Cole, Chelsea Voss, Chen Ding, Chen Shen, Chengdu Huang, Chris Colby, Chris Hallacy, Chris Koch, Chris Lu, Christina Kaplan, Christina Kim, CJ Minott-Henriques, Cliff Frey, Cody Yu, Coley Czarnecki, Colin Reid, Colin Wei, Cory Decareaux, Cristina Scheau, Cyril Zhang, Cyrus Forbes, Da Tang, Dakota Goldberg, Dan Roberts, Dana Palmie, Daniel Kappler, Daniel Levine, Daniel Wright, Dave Leo, David Lin, David Robinson, Declan Grabb, Derek Chen, Derek Lim, Derek Salama, Dibya Bhattacharjee, Dimitris Tsipras, Dinghua Li, Dingli Yu, DJ Strouse, Drew Williams, Dylan Hunn, Ed Bayes, Edwin Arbus, Ekin Akyurek, Elaine Ya Le, Elana Widmann, Eli Yani, Elizabeth Proehl, Enis Sert, Enoch Cheung, Eri Schwartz, Eric Han, Eric Jiang, Eric Mitchell, Eric Sigler, Eric Wallace, Erik Ritter, Erin Kavanaugh, Evan Mays, Evgenii Nikishin, Fangyuan Li, Felipe Petroski Such, Filipe de Avila Belbute Peres, Filippo Raso, Florent Bekerman, Foivos Tsimpourlas, Fotis Chantzis, Francis Song, Francis Zhang, Gaby Raila, Garrett McGrath, Gary Briggs, Gary Yang, Giambattista Parascandolo, Gildas Chabot, Grace Kim, Grace Zhao, Gregory Valiant, Guillaume Leclerc, Hadi Salman, Hanson Wang, Hao Sheng, Haoming Jiang, Haoyu Wang, Haozhun Jin, Harshit Sikchi, Heather Schmidt, Henry Aspegren, Honglin Chen, Huida Qiu, Hunter Lightman, Ian Covert, Ian Kivlichan, Ian Silber, Ian Sohl, Ibrahim Hammoud, Ignasi Clavera, Ikai Lan, Ilge Akkaya, Ilya Kostrikov, Irina Kofman, Isak Etinger, Ishaan Singal, Jackie Hehir, Jacob Huh, Jacqueline Pan, Jake Wilczynski, Jakub Pachocki, James Lee, James Quinn, Jamie Kiros, Janvi Kalra, Jasmyn Samaroo, Jason Wang, Jason Wolfe, Jay Chen, Jay Wang, Jean Harb, Jeffrey Han, Jeffrey Wang, Jennifer Zhao, Jeremy Chen, Jerene Yang, Jerry Tworek, Jesse Chand, Jessica Landon, Jessica Liang, Ji Lin, Jiancheng Liu, Jianfeng Wang, Jie Tang, Jihan Yin, Joanne Jang, Joel Morris, Joey Flynn, Johannes Ferstad, Johannes Heidecke, John Fishbein, John Hallman, Jonah Grant, Jonathan Chien, Jonathan Gordon, Jongsoo Park, Jordan Liss, Jos Kraaijeveld, Joseph Guay, Joseph Mo, Josh Lawson, Josh McGrath, Joshua Vendrow, Joy Jiao, Julian Lee, Julie Steele, Julie Wang, Junhua Mao, Kai Chen, Kai Hayashi, Kai Xiao, Kamyar Salahi, Kan Wu, Karan Sekhri, Karan Sharma, Karan Singhal, Karen Li, Kenny Nguyen, Keren Gu-Lemberg, Kevin King, Kevin Liu, Kevin Stone, Kevin Yu, Kristen Ying, Kristian Georgiev, Kristie Lim, Kushal Tirumala, Kyle Miller, Lama Ahmad, Larry Lv, Laura Clare, Laurance Fauconnet, Lauren Itow, Lauren Yang, Laurentia Romaniuk, Leah Anise, Lee Byron, Leher Pathak, Leon Maksin, Leyan Lo, Leyton Ho, Li Jing, Liang Wu, Liang Xiong, Lien Mamitsuka, Lin Yang, Lindsay McCallum, Lindsey Held, Liz Bourgeois, Logan Engstrom, Lorenz Kuhn, Louis Feuvrier, Lu Zhang, Lucas Switzer, Lukas Kondraciuk, Lukasz Kaiser, Manas Joglekar, Mandeep Singh, Mandip Shah, Manuka Stratta, Marcus Williams, Mark Chen, Mark Sun, Marselus Cayton, Martin Li, Marvin Zhang, Marwan Aljubeh, Matt Nichols, Matthew Haines, Max Schwarzer, Mayank Gupta, Meghan Shah, Melody Huang, Meng Dong, Mengqing Wang, Mia Glaese, Micah Carroll, Michael Lampe, Michael Malek, Michael Sharman, Michael Zhang, Michele Wang, Michelle Pokrass, Mihai Florian, Mikhail Pavlov, Miles Wang, Ming Chen, Mingxuan Wang, Minnia Feng, Mo Bavarian, Molly Lin, Moose Abdool, Mostafa Rohaninejad, Nacho Soto, Natalie Staudacher, Natan LaFontaine, Nathan Marwell, Nelson Liu, Nick Preston, Nick Turley, Nicklas Ansman, Nicole Blades, Nikil Pancha, Nikita Mikhaylin, Niko Felix, Nikunj Handa, Nishant Rai, Nitish Keskar, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Oona Gleeson, Pamela Mishkin, Patryk Lesiewicz, Paul Baltescu, Pavel Belov, Peter Zhokhov, Philip Pronin, Phillip Guo, Phoebe Thacker, Qi Liu, Qiming Yuan, Qinghua Liu, Rachel Dias, Rachel Puckett, Rahul Arora, Ravi Teja Mullapudi, Raz Gaon, Reah Miyara, Rennie Song, Rishabh Aggarwal, RJ Marsan, Robel Yemiru, Robert Xiong, Rohan Kshirsagar, Rohan Nuttall, Roman Tsiupa, Ronen Eldan, Rose Wang, Roshan James, Roy Ziv, Rui Shu, Ruslan Nigmatullin, Saachi Jain, Saam Talaie, Sam Altman, Sam Arnesen, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Sarah Yoo, Savannah Heon, Scott Ethersmith, Sean Grove, Sean Taylor, Sebastien Bubeck, Sever Banesiu, Shaokyi Amdo, Shengjia Zhao, Sherwin Wu, Shibani Santurkar, Shiyu Zhao, Shraman Ray Chaudhuri, Shreyas Krishnaswamy, Shuaiqi, Xia, Shuyang Cheng, Shyamal Anadkat, Simón Posada Fishman, Simon Tobin, Siyuan Fu, Somay Jain, Song Mei, Sonya Egoian, Spencer Kim, Spug Golden, SQ Mah, Steph Lin, Stephen Imm, Steve Sharpe, Steve Yadlowsky, Sulman Choudhry, Sungwon Eum, Suvansh Sanjeev, Tabarak Khan, Tal Stramer, Tao Wang, Tao Xin, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Degry, Thomas Shadwell, Tianfu Fu, Tianshi Gao, Timur Garipov, Tina Sriskandarajah, Toki Sherbakov, Tomer Kaftan, Tomo Hiratsuka, Tongzhou Wang, Tony Song, Tony Zhao, Troy Peterson, Val Kharitonov, Victoria Chernova, Vineet Kosaraju, Vishal Kuo, Vitchyr Pong, Vivek Verma, Vlad Petrov, Wanning Jiang, Weixing Zhang, Wenda Zhou, Wenlei Xie, Wenting Zhan, Wes McCabe, Will DePue, Will Ellsworth, Wulfie Bain, Wyatt Thompson, Xiangning Chen, Xiangyu Qi, Xin

Xiang, Xinwei Shi, Yann Dubois, Yaodong Yu, Yara Khakbaz, Yifan Wu, Yilei Qian, Yin Tat Lee, Yinbo Chen, Yizhen Zhang, Yizhong Xiong, Yonglong Tian, Young Cha, Yu Bai, Yu Yang, Yuan Yuan, Yuanzhi Li, Yufeng Zhang, Yuguang Yang, Yujia Jin, Yun Jiang, Yunyun Wang, Yushi Wang, Yutian Liu, Zach Stubenvoll, Zehao Dou, Zheng Wu, and Zhigang Wang. 2025. Openai gpt-5 system card. *Preprint*, arXiv:2601.03267.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. 2025. Qwen3 technical report. *Preprint*, arXiv:2505.09388.

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757.

Mingyuan Zhong, Gang Li, Peggy Chi, and Yang Li. 2021. Helpviz: Automatic generation of contextual visual mobile tutorials from text-based instructions. In *The 34th Annual ACM Symposium on User Interface Software and Technology*, pages 1144–1153.

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

## A  Page Category Analysis

This appendix provides additional quantitative analyses of page categories in the GUIDEWEB benchmark. Figure 6 presents the distribution of page categories, highlighting the dominance of landing pages alongside a diverse set of real-world entry points. Figure 7 further breaks down each category by whether pages require in-app guides, showing that only a small fraction of pages are labeled as not guide-worthy, and that such cases are concentrated in a few content-oriented categories.

## B  Metric Definitions

**Guide target selection.**  For a given webpage, let $G$ be the set of gold guide targets and $\hat{G}$ the predicted targets, where each target corresponds to a unique DOM element. We define precision, recall, and F1 as follows:

$$\text{Precision} \;=\; \frac{|G \cap \hat{G}|}{|\hat{G}|}. \qquad (9)$$

$$\text{Recall} \;=\; \frac{|G \cap \hat{G}|}{|G|}. \qquad (10)$$

$$\text{F1} \;=\; \frac{2\,\text{Precision}\,\text{Recall}}{\text{Precision} + \text{Recall}}. \qquad (11)$$

**Field-level exact-match F1.**  For each structured field (e.g., `action_type`, `tag`, `visible_text`, `xpath`), we compute exact-match precision, recall, and F1 over field values aggregated across annotations. Let $S$ denote the (multi)set of gold field values and $\hat{S}$ the (multi)set of predicted values.[3] We define:

$$\text{Precision} \;=\; \frac{|S \cap \hat{S}|}{|\hat{S}|}. \qquad (12)$$

$$\text{Recall} \;=\; \frac{|S \cap \hat{S}|}{|S|}. \qquad (13)$$

$$\text{F1} \;=\; \frac{2\,\text{Precision}\,\text{Recall}}{\text{Precision} + \text{Recall}}. \qquad (14)$$

## C  Implementation Details

**Baseline Inference Settings.**  All baseline models are evaluated in an inference-only setting without any task-specific fine-tuning. We use the original model parameters and default decoding configurations provided by each model implementation. To ensure fair comparison across models and to cover the majority of real-world webpages, we set the maximum input length to 130K tokens for all baselines. This context window is sufficient to accommodate full HTML source code for most main pages in the GUIDEWEB benchmark after preprocessing, avoiding aggressive truncation that may remove critical interactive elements. All baselines are prompted using a unified instruction template and generate guide target predictions and guide text in a single forward pass.

**GuideWeb Agent Training.**  The GUIDEWEB AGENT is trained using the GUIDEWEB training split with full-parameter fine-tuning. Training is conducted on an NVIDIA GB10 GPU with

---

[3]Our implementation handles duplicates consistently between $S$ and $\hat{S}$.
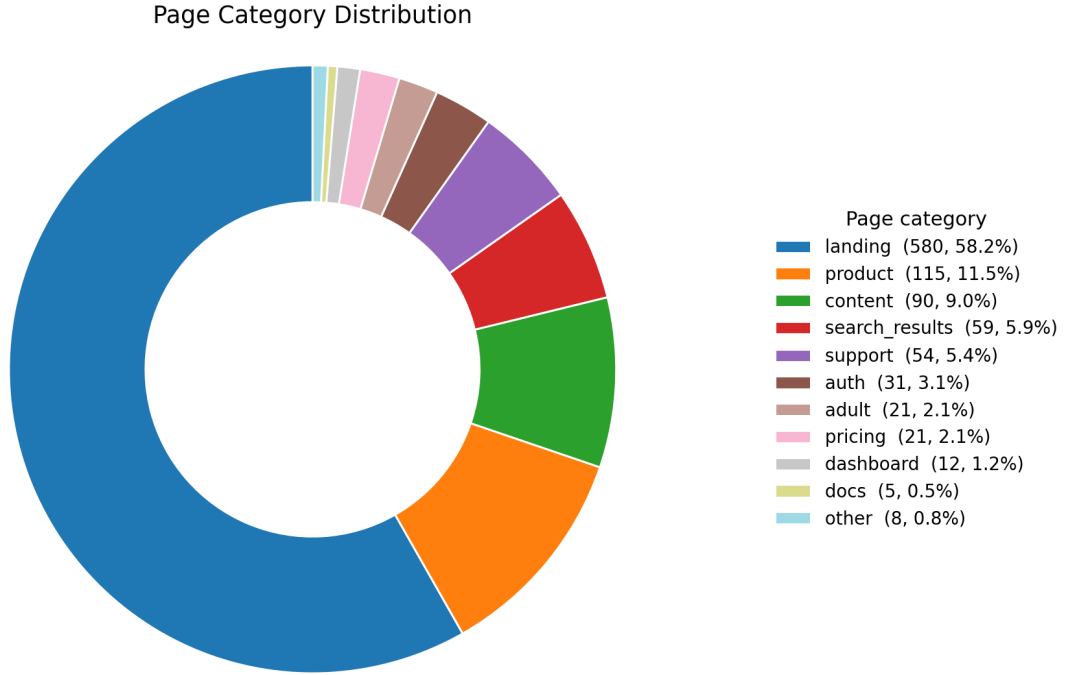
Page Category Distribution



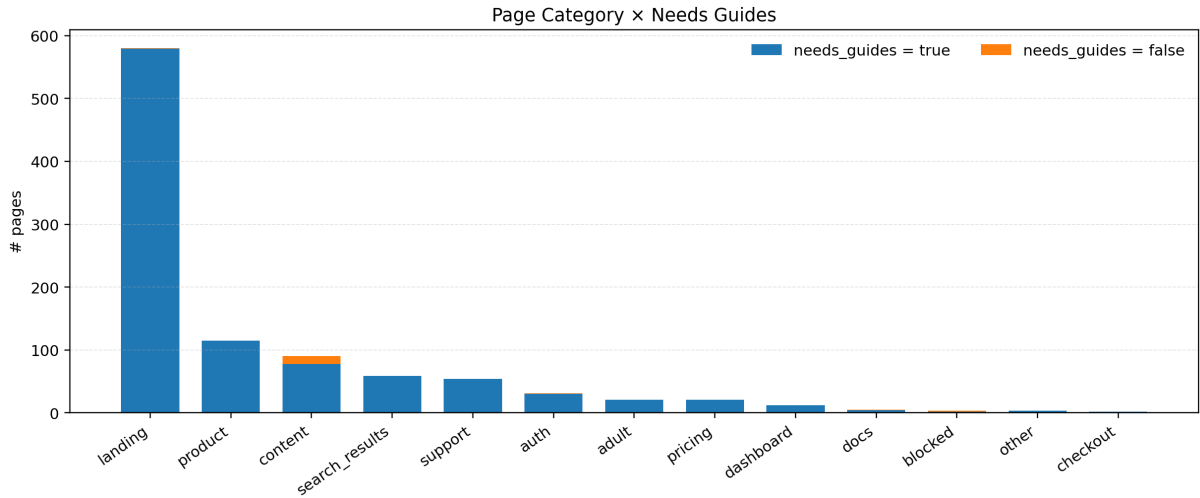Figure 6: Page category distribution in GUIDEWEB.



Figure 7: Page category by whether in-app guides are needed in GUIDEWEB.

120 GB of memory, which allows stable optimization with long-context inputs and structured output constraints. We adopt a supervised fine-tuning (SFT) paradigm, where the model is trained to jointly predict guide targets and corresponding guide text conditioned on preprocessed webpage representations.

To reduce unnecessary computation caused by redundant or non-informative HTML content, we employ a *shorter* mechanism that selectively retains interactive elements, visible text, and local contextual cues surrounding candidate targets. This mechanism significantly shortens the effective input length while preserving task-relevant information, resulting in faster training and inference without degrading accuracy, as confirmed by our experimental results.

**Training Hyperparameters.** We fine-tune the GUIDEWEB AGENT with full-parameter supervised learning using AdamW. We train for **3 epochs** with a learning rate of $1 \times 10^{-5}$, warmup ratio **0.03**, and weight decay **0.0**. The per-device batch size is **1** with gradient accumulation of **4** (effective batch

size **4**). We log every **10** steps and save checkpoints every **2000** steps. For numerical stability and memory efficiency, we load the base model in **FP16** and train with **BF16** mixed precision, enabling **gradient checkpointing**. We set the maximum prompt length to **3950** tokens and the maximum generation length to **2000** tokens, resulting in a maximum sequence length of **6014** tokens (including a small buffer). All runs use a fixed random seed (**13**).

**Input Construction and Shorter Configuration.** We render each webpage using a fixed viewport of **1280×720**. For visible text, we keep up to **800** text blocks, truncating each block to at most **400** characters. We additionally include a short global excerpt of the full-page text using an excerpt ratio of **0.02**, bounded to **[100, 200]** characters. For structural cues, we keep headings (enabled) with at most **20** headings and at most **40** characters per heading. For interactive candidates, we extract up to **2000** interactive elements, truncate each element's visible text to **120** characters, and enforce an overall interactive-text budget of **6500** characters. To represent targets robustly, we provide XPaths using the `stable_then_abs` mode, limiting XPath textual fields to **80** characters and attribute fields to **200** characters. Model outputs are required to end with a fixed end marker (`</JSON>`).

**Filtering Long Samples.** To avoid near-limit truncation and unstable supervision, we filter overly long samples during training. Specifically, we drop samples whose estimated total tokens exceed **5200**, or whose prompt/output lengths are within **98%** of the configured maximums. Dropped instances are recorded in `dropped_samples.jsonl` for traceability.