THE IMPACT OF ELEMENT ORDERING ON LM AGENT PERFORMANCE

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

030

Paper under double-blind review

Abstract

There has been a surge of interest in language model agents that can navigate virtual environments such as the web or desktop. To navigate such environments, agents benefit from information on the various elements (e.g., buttons, text, or images) present. However, it remains unclear which element attributes have the greatest impact on agent performance, especially in environments that only provide a graphical representation (i.e., pixels). Here we find that the *ordering* in which elements are presented to the language model is surprisingly impactful—randomizing element ordering in webpages compromises average agent performance to a degree comparable to removing all visible text from webpages. While web agents benefit from the semantic hierarchical ordering of elements available via the browser, agents that parse elements directly from pixels do not have access to any such ordering. Here we endeavor to derive effective orderings and investigate the impact of various element ordering methods in web and desktop environments. We find that dimensionality reduction provides a viable ordering for pixel-only environments. We train a UI element detection model to derive elements from pixels and apply our findings to an agent benchmark—OmniACT—where we only have access to pixels. Our method completes more than two times as many tasks on average relative to the previous state-of-the-art.

1 INTRODUCTION

There has been growing interest in using language model (LM) agents to autonomously navigate virtual environments. Autonomous web agents (Zhou et al., 2023; Kim et al., 2023; Zheng et al., 2024; Gur et al., 2024; He et al., 2024) have become a particularly popular area of research. Typically, a web agent takes as input a task prompt from a user, observes a text and visual representation of the environment, and then outputs one or more actions to execute the task in the environment. Recently, research interest has expanded to include agents that can navigate mobile (Rawles et al., 2023; Yan et al., 2023) and desktop (Xie et al., 2024; Kapoor et al., 2024; Bonatti et al., 2024) environments as well.

At a high level, a virtual environment consists of numerous elements—some are interactive (e.g. 040 buttons or widgets), while others are not (e.g. plain text). To allow for human navigation, these 041 elements are usually represented in the pixel space via a Graphical User Interface (GUI). In contrast, 042 agents often rely on distinct state representations to navigate virtual environments. The exact format 043 of a state representation varies between environments and approaches. In web environments, common 044 text representations include the HTML or accessibility tree (Zhou et al., 2023; Koh et al., 2024). For visual representations, a popular approach is to label UI elements with bounding boxes and numeric identifiers (Koh et al., 2024; He et al., 2024), known as Set-of-Mark (Yang et al., 2023). In either case, 046 the state representation is derived from the underlying Document Object Model (DOM) (Zhou et al., 047 2023; Koh et al., 2024; He et al., 2024). However, many environments lack a descriptive DOM and 048 only provide pixel information, which we show is insufficient for existing agents (see Section 4.2). To construct an effective state representation from only pixels, it is important to answer the following basic questions about these representations: What aspects of a state representation are most important 051 to an agent? How can we derive these important aspects with only pixels? 052

In almost all implementations of the agent's state representation, there exists a list of interactable or non-interactable elements which the agent uses to determine the next action (Koh et al., 2024; Zhou



089 090

091

092

094

095

096

Figure 1: We are motivated by the goal of enabling agents to act on environments where an underlying DOM does not exist. Instead, the agent must determine its next action using only the environment's graphical representations. In Step 1, we first detect a list of unordered UI elements using an object detection model and identify them with bounding boxes. In Step 2, we convert these UI elements to their text representation. In Step 3, we order the elements via 2D-to-1D dimensionality reduction. Due to the sequential nature of a language model, elements are always presented in a specific order to the LM Agent. Finding an effective ordering is non-trivial, yet can significantly affect agent performance. Elements that are visually close together are often functionally associated with each other. t-SNE's ability to retain local structure allows it to generate an effective ordering.

098

100 et al., 2023; He et al., 2024; Dupont, 2024; Yan et al., 2023; Ishan0102, 2024; Kapoor et al., 2024; 101 Xie et al., 2024). Elements are characterized by various attributes such as visual appearance, text 102 descriptions, or type labels. Because the state representation is the input to an LM, this list of elements 103 is always presented in a specific ordering. For example, the default method to derive elements from a 104 webpage performs a pre-order traversal of the DOM tree (World Wide Web Consortium, 2013). We 105 analyze various attributes of a popular state representation for agents and find element ordering to be the single most impactful attribute to agent performance. We find that the ordering of elements 106 can dramatically affect the performance of an agent, resulting in differences of up to 49% relative 107 performance.



Figure 2: Many software applications lack informative accessibility trees or DOMs. The accessibility tree for a popular game development engine (Godot, left) contains only the exit, minimize, and full screen buttons. For a presentation slide (Google Slides, right), no interactive elements (e.g. title 122 and subtitle boxes) are present in the DOM. Language model agents rely on this information to 123 navigate applications, and agent performance can accordingly be compromised in scenarios where it 124 is incomplete.

125 126

120

121

127 This can prove problematic as many environments lack obvious methods to both derive and order elements. For example, many mobile and desktop applications (see Figure 2) do not properly expose 128 interactable elements (Chen et al., 2020; Ross et al., 2020; Zhang et al., 2021). In such environments, 129 pixels may be the only source of information available. Previous approaches to deriving interactable 130 elements from pixels either leverage off-the-shelf segmentation models (Yan et al., 2023; Kapoor 131 et al., 2024) or build custom models that target accessibility features (Wu et al., 2023). In our 132 approach, we leverage common crawl (Common Crawl, 2023) to train an object detection model (Ren 133 et al., 2016) that detects interactable UI elements specifically for agents. To the best of our knowledge, 134 the elements detected through these approaches are ordered arbitrarily (e.g. based on confidence 135 scores); visually, the ordering is effectively random. Our experiments indicate that a random ordering 136 consistently results in the lowest performance across multiple scenarios.

137 Here we propose and evaluate strategies for deriving effective element orderings in scenarios where 138 a hierarchical ordering based on the GUI design is not explicitly provided by the environment. 139 Across multiple agent benchmarks, we find that ordering elements via a 2D-to-1D dimensionality 140 reduction (Van der Maaten and Hinton, 2008) reliably yields improvements to agent performance 141 relative to other baselines. We experiment on the VisualWebArena (Koh et al., 2024) and OmniACT 142 (Kapoor et al., 2024) benchmarks and achieve new state-of-the-art performance on OmniACT.

- 143 Out contributions are as follows. 144
 - We conduct a thorough ablation of VisualWebArena's state representation for agents by including or removing each element attribute individually. Despite advancements in vision language models, we find that a text representation is still necessary for web and desktop agents. We find that element ordering is, perhaps surprisingly, more impactful than any other attribute in the text representation.
 - We demonstrate that ordering via dimensionality reduction consistently provides performance improvements over random ordering. Additionally, we find that ordering via dimensionality reduction performs better than a simple position-based ordering in most scenarios.
 - We achieve a new state-of-the-art result on OmniACT, an agent benchmark that considers the scenario of operating on pixels. Our approach more than doubles the expected average task success rate compared to the previous state-of-the-art.

156 157

145

146

147

148

149

150

151

152

153

RELATED WORK 2

¹⁵⁸ 159

Agent Benchmarks. World-of-bits provided the first environment for evaluating web GUI navigation using an agent (Shi et al., 2017). Over time, more realistic web (Zhou et al., 2023; Koh 161 et al., 2024; Yao et al., 2023), desktop (Kapoor et al., 2024; Xie et al., 2024), and mobile (Rawles

et al., 2023) agent benchmarks have been created. Kim et al. (2023) provided one of the first LM agent approaches, successfully navigating World-of-bits. However, existing agents are still unable to properly navigate more realistic benchmarks, completing only 15% of web tasks (Zhou et al., 2023) and 12% of desktop tasks.

166

188

189

190

191

192 193

194 195

196

197

199

200

201

202

203

204

205 206

207

208

167 Agents With Direct Access to Elements. Despite work on multimodal agents (Koh et al., 2024), 168 existing techniques in navigating web and desktop environments still rely heavily on ground-truth 169 text representations. Zhou et al. (2023); Koh et al. (2024) both utilize the accessibility tree and its 170 elements as their state representation. Koh et al. (2024); He et al. (2024) consider approaches where interactable elements are also represented in an image via Set-of-Mark (Yang et al., 2023) bounding 171 boxes and labels. Xie et al. (2024) provides an agent that navigates desktop applications by observing 172 a filtered down version of the accessibility tree. All of these approaches require access to either a 173 webpage's underlying DOM or an accessibility tree to derive elements; however, our focus is on 174 environments that only give access to their graphical representations which is significantly more 175 challenging. 176

177 Agents With Access to Only a Graphical Representation. There have been several approaches— 178 primarily focused on desktop and mobile environments—to directly navigating a GUI via its pixels. 179 Kapoor et al. (2024) and Yan et al. (2023) focus on navigating desktop and mobile applications respectively. Both leverage an off-the-shelf-segmentation model-Segment Anything (Kirillov et al., 181 2023)—to find icons in the image. These icons are then either represented in text (Kapoor et al., 2024) 182 or labeled with Set-of-Mark (Yang et al., 2023) bounding boxes and labels in the image (Yan et al., 2023). We instead train an object detection model that detects interactable UI elements directly. While 183 previous UI element detection models are trained to detect accessibility features (Wu et al., 2023), our model is trained specifically to detect interactable elements that would be useful to an agent. All 185 three of our approaches use Optical Character Recognition (OCR) modules such as EasyOCR (AI, 186 2020) to extract text from pixel information. 187

Agent Input Ablations. While most agent studies include some ablations, few focus on detailed analysis of an agent's input. To our knowledge, Huq et al. (2023) is the only other study that directly studies this. Their study focuses on broader components to an input prompt such as the selection of few-shot examples used, while we focus on specific element attributes such as element ordering.

3 PROBLEM DEFINITION

We define the environment state \mathcal{E} as a set of elements $\mathcal{E} = \{e_1, e_2, \dots, e_n\}$, where each element e_j is a tuple $\langle i_j, C_j, A_j, S_j \rangle$ defined by the following parameters:

- $i_j \in \{0, 1\}$ denotes the interactability of element e_j . An element with $i_j = 1$ is interactable, while an element with $i_j = 0$ is not.
- $C_j = \{ \langle x_1, y_1 \rangle, \langle x_2, y_2 \rangle \}$ is the set of pixel coordinates that form the bounding box around e_j . $\langle x_1, y_1 \rangle$ is the top left coordinate and $\langle x_2, y_2 \rangle$ is the bottom right coordinate.
- $A_j = \{a_1, a_2, \dots, a_m\}$ is the set of potential actions that can be taken on that element. For example, the potential actions for a search bar might be {click,type}; a non-interactive text element has the action set \emptyset .
- S_j represents the set of other environment-specific attributes for element e_j . These attributes can include image captions, type labels (e.g., Button, Text Field), or accessibility information.

While certain environments may provide full access to this environment state, here we are focus on environments where only the pixel information, \mathcal{P} , is available. We then must predict elements from the pixel information to construct a state representation for the agent. In other words, we must find a function $g: \mathcal{P} \to \mathcal{E}$

Elements can be represented in both visual and text modalities. For images, the most common approach is to overlay bounding boxes with numeric identifiers around each interactable element (Koh et al., 2024; He et al., 2024; Yan et al., 2023) in a manner inspired by Set-of-Mark Prompting

(Yang et al., 2023). In text, a common approach is to represent each element as a string containing its index, coordinates, and other attributes, such as "[1] [x,y] [Description]". Because LMs operate on sequential data, elements must be given an ordering; in most approaches, this is implicitly defined by the method used to identify the elements.

Ordering. The ordering is defined as a permutation σ of the indices $\{1, 2, ..., n\}$ where $\{e_{\sigma(1)}, e_{\sigma(2)}, ..., e_{\sigma(n)}\}$ represents a specific sequence of elements. An ordering function is a function $f: \mathcal{E} \to \sigma$ that takes the environment as an input and yields a specific ordering σ .

4 WHICH ASPECTS OF AGENT STATE REPRESENTATIONS ARE MOST IMPACTFUL?

Here we describe a series of ablation experiments designed to examine which aspects of an LM state representation are most impactful to the performance of LM agents. In particular, we experiment on the VisualWebArena (VWA) (Koh et al., 2024) benchmark and ablate attributes of the state representation of the state-of-the-art agent (proposed in the same paper). We pick this representation in particular due to both its strong performance on VWA, as well as its similarity to common practices seen in other agent research (He et al., 2024) and open source projects (Ishan0102, 2024; Dupont, 2024). Our experiments in turn ablate the impact of (1) multimodal (image and text) aspects of the state representation, and (2) individual element attributes within the text component alone.

236 237 4

221

222

223 224 225

226

227 228

229

230

231

232

233

234

235

4.1 VISUALWEBARENA

238 VisualWebArena focuses on multimodal tasks in the web and provides a self-hostable environment 239 for language agents to navigate (Koh et al., 2024). Agents operating on VisualWebArena have full 240 access to the DOM. The current best approach (Koh et al., 2024) utilizes a multimodal representation 241 where elements are parsed in a pre-order traversal of the DOM tree (World Wide Web Consortium, 242 2013). Each element j has attributes $S_i = \{id, tag, text\}$ where id is a unique numeric identifier for 243 interactable elements and \varnothing otherwise, tag is the HTML tag (e.g. BTN or IMG), and text is the alt 244 text and captions for images and the HTML text otherwise. In the text representation, an example of 245 an element would be "[1] [IMG] [alt text, caption]". In the image representation, each element is labeled with bounding boxes and numeric labels. 246

247 The original agent in (Koh et al., 2024) only achieves 15% success rate across all tasks. Since our 248 goal is not to improve agent performance on VisualWebArena but rather to understand the importance 249 of attributes in the state representation, we examine a subset of tasks to reduce costs.⁰ Specifically, 250 we explore tasks marked as "easy" within tasks that the original agent completed successfully. Due to 251 variance associated with stochastic LM outputs, our reproduction of these originally-successful tasks 252 yields a success rate around $74.07\% \pm 5.56\%$. The exact list of tasks can be found in Appendix A.8. We reuse the action space from the original agent which consists of executing high-level actions 253 (e.g., click, hover) on individual elements—see Appendix A.6 for more details. 254

255 256

4.2 Ablation Setup and Findings

257 The agent state representation we explore is multimodal and consists of image and text information. 258 The image consists of a screenshot of the webpage along with Set-of-Mark annotations, while the text 259 consists of a DOM-ordered list of elements with the attributes outlined above. Our ablation protocol 260 consists of removing individual attributes from the image or text representation and measuring task 261 success rate—we say an attribute has high "impact" if its removal leads to substantial reduction in task 262 success rate. To provide evidence that these findings may be robust across different LM backbones, 263 we explore both GPT-4V as used in the original agent, and Gemini 1.5 Pro. Some experiments were 264 not run on GPT-4V due to high associated costs, though we found ordering to be consistent between these two LM backbones on all ablations where we ran both. 265

In Table 1, we report results ablating aspects of the multimodal representation. In Table 2, we report the impact of ablating various attributes in the text representation specifically. Across all experiments, we consider the pre-order traversal of the DOM tree as the ground truth element ordering, and define

⁰A full run of the state-of-the-art agent on VisualWebArena can cost up to \$800 with GPT-4V.

Table 1: Ablating the multimodal aspects of state representation in VisualWebArena. \checkmark indicates ground truth obtained from the HTML. x indicates removal of the attribute. We find that removing the text representation can dramatically harm agent performance.

Observations			Success Rate (†)		
Screenshot	Set-of-Mark	Text Representation	Gemini 1.5	GPT-4V	
\checkmark	\checkmark	\checkmark	64.20%	74.07%	
Х	х	\checkmark	46.30%	38.89%	
\checkmark	х	\checkmark	50.00%	-	
\checkmark	\checkmark	X	3.70%	38.89%	

the "removal" of ordering information as substituting an ordering σ_{rand} picked uniformly at random from all possible permutations. We summarize a few key findings from both sets of ablations below.

Text Representation is Still Necessary. While adding a visual representation clearly improves performance, we find that it alone is insufficient even with Set-of-Mark labels. This contradicts previous findings on agents for mobile applications which found that a screenshot with Set-of-Mark labels achieves similar performance with or without text (Yan et al., 2023). We speculate that this is due to the substantial difference in viewport sizes between mobile and desktop environments. Specifically, the average mobile device has a viewport size of 360x800 while the average desktop has a viewport size of 1920x1080 (Statcounter Global Stats, 2024). Additionally, larger viewport sizes have been shown to improve agent performance in desktop environments (Xie et al., 2024). We speculate that this may be because current agents almost never understand when to change the screen view (e.g. by scrolling).

Removing Ordering Information Harms Performance More Than Removing Any Other Attribute. Although most element attributes are important, we find that ordering is the single most important attribute for agent performance. Random ordering results in a similar performance drop to removing all HTML text descriptions.

Captions Impact Performance More Than Alt Text. Removing captions causes a greater decrease to performance than removing alt text. From our experience, captions almost always provide more information than alt text. In fact, captions frequently include the alt text directly in its description.

Table 2: Ablating attributes of the text component of the VisualWebArena state representation. All results include the screenshot with Set-of-Mark bounding boxes and labels. TAG is the HTML tag. CAPTIONS are image captions generated using BLIP-2-T5XL(Li et al., 2023). TEXT_{Alt}, TEXT_{Interact}, and TEXT_{Static} are the alt text, text for interactable elements, and text for noninteractable elements respectively. ORDER is element ordering. \checkmark indicates ground truth obtained from the HTML. x indicates removal of the attribute. x Element Ordering indicates a random shuffling of the elements. - denotes experiments that were not run due to cost.

		Success F	Rate (†)				
TAG	CAPTIONS	TEXT_{Alt}	$TEXT_{Interact}$	TEXT_{Static}	ORDER	Gemini 1.5	GPT4-V
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	64.03%	74.07%
х	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	61.11%	61.11%
\checkmark	х	\checkmark	\checkmark	\checkmark	\checkmark	46.30%	-
\checkmark	\checkmark	x	\checkmark	\checkmark	\checkmark	68.15%	66.67%
\checkmark	\checkmark	\checkmark	X	\checkmark	\checkmark	53.70%	-
\checkmark	\checkmark	\checkmark	\checkmark	х	\checkmark	57.40%	-
\checkmark	\checkmark	\checkmark	Х	Х	\checkmark	35.18%	-
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	Х	37.04%	44.44%

324 5 EXPERIMENTAL SETUP

For the remainder of this paper, we leverage the insights gained from our state representation ablations
 on VisualWebArena to tackle a more challenging task: enabling LM agents to act in environments
 that only expose pixel information.

Most applications are built on top of an underlying hierarchical representation. For example, a 330 webpage is modeled by the DOM which is hierarchical. When exposed, this hierarchy can be used 331 to determine a strong element ordering. However, the availability and quality of an underlying 332 hierarchical representation can vary greatly between environments. For example, Chen et al. (2020) 333 found that 77% of mobile applications had missing labels and Ross et al. (2020) found that 53% 334 of image buttons had missing labels and were incorrectly sized. Additionally, Ross et al. (2020) 335 found that 8% of applications lacked interactable element information altogether. In such scenarios, 336 we may only have access to the application's pixel information. We continue to experiment with VisualWebArena and also experiment with the OmniACT benchmark as a scenario where we only 337 have access to pixel information. Details on the LM agent backbones used in our experiments can be 338 found in Appendix A.7. 339

340

341 5.1 OmniACT

OmniACT provides both web and desktop environments for agents to benchmark on. OmniACT contains 177 application screenshots overall and 2021 tasks in the test set. Agents are tasked with generating pyautogui code that can navigate the application screenshot. We consider OmniACT as a setting where only a pixel information is available.

To detect UI elements $\{e_1, e_2, \ldots, e_n\}$ when given only pixel information \mathcal{P} we train an object 347 detection model (Ren et al., 2016) to detect interactable UI elements in the screenshot and use 348 EasyOCR (AI, 2020) to extract visible text. In other words, the function $q: \mathcal{P} \to \mathcal{E}$ is defined by the 349 trained object detection model. We add visible text and captions to each interactable UI element. We 350 gather a dataset by finding 674,416 interactable elements over 1468 Common Crawl webpages. We 351 selected our webpages based on top websites from Similarweb (2024). Despite the domain shift from 352 webpages to desktop applications, we found that our object detection model worked reasonably well 353 on the OmniACT benchmark in the end-to-end agent setting. We publicly release this model along 354 with our paper. Training details can be found in the Appendix A.1.

OmniACT provides partial human annotations for each screenshot; multiple, but not all, interactable UI elements are annotated with bounding boxes. The original intent of these bounding boxes is for evaluation only. As a result, there are significantly less UI elements annotated compared to possible UI elements in the application. We experiment with elements derived from these human annotated bounding boxes to understand a) impacts to ordering performance in an easier setting and b) the potential performance that can be gained by improving UI element detection.

5.2 METRICS

362

372

374

376

363 Our primary metric is task success rate which is the standard for agent evaluations (Koh et al., 2024; 364 Zhou et al., 2023; He et al., 2024). VisualWebArena provides an evaluation framework for task success. Task success criteria include achieving an expected final webpage state or receiving a desired 366 response from the language agent. OmniACT does not provide task success rate directly, instead 367 introducing the sequence score and action score metrics. Sequence score evaluates if the output 368 contains the correct high level action (e.g. click or type), but does not check if the action element 369 or parameter (e.g. click [1] or type [parameter]) are correct. Action score evaluates if 370 the output contains both the correct high level action and the correct element or parameter. Thus, for 371 OmniACT we focus our evaluations on the action score as it is more similar to task success rate.

- 373 5.3 ORDERING METHODS
- In addition to random ordering, we experiment with two different ordering methods.
- **Random.** We pick the ordering σ_{rand} uniformly at random from all possible permutations. This provides a baseline performance.

Raster. Elements $\{e_1, e_2, \dots, e_n\}$ are ordered in a left-to-right raster scan. We define a raster scan as an ordering σ_{raster} where i < j iff $\lfloor \frac{y_i}{8} \rfloor < \lfloor \frac{y_j}{8} \rfloor$ and $x_i < x_j$. We chose to discretize the scan to prevent jumps in ordering from minor pixel variations. This method mimics the natural way English speakers read text and images from left-to-right and top-to-bottom.

t-SNE. We apply dimensionality reduction techniques to better capture the spatial relationships. Using t-SNE (Van der Maaten and Hinton, 2008), we reduce the dimensionality with the function $g: \langle x, y \rangle \rightarrow z$. The set of values $Z = \{z_1, z_2, \dots, z_m\}$ is used to determine the ordering σ_{tsne} . We use the scikit-learn (Pedregosa et al., 2011) implementation of t-SNE and keep the default parameters.

Our intuition for choosing t-SNE stems from visually and qualitatively inspecting our agent trajectories. We observed that agents often use adjacently ordered elements as context clues to determine the correct action. Furthermore, agents have difficulty reasoning about functionally associated elements that are separated from each other in the ordering. t-SNE generates an effective ordering as elements that are visually close together (i.e. nearby in 2D pixel space) are often functionally associated with each other as well. When reducing dimensions, t-SNE retains local structure which increases the odds that functionally associated elements are adjacent in the induced 1D ordering.

394 395

382

5.4 ACTION SPACE

396

397

412 413 414

415 416

417

418 419

420

We use the same high-level action space as 399 described in OmniACT. Unlike OmniACT, 400 we do not have the model directly output pyautogui code. Instead, we use high 401 level actions that map to pyautogui code. 402 For example, for element e_1 defined by 403 $\langle 1, \{\langle 50, 50 \rangle, \langle 100, 100 \rangle\}, \{ click \}, \emptyset \rangle$ 404 would output click the model 405 [1] which would be converted to 406 pyautoqui.click(75, 75). This 407 prevents the model from having to reason 408 about pixel coordinates directly. Each 409 element's unique identifier reflects the 410 position of the element in the ordering. 411 The full action space is in Table 3.

Table 3: The set of possible actions in OmniACT.

Action	Description
Click	Perform a single click on an element.
Double Click	Perform a double click on an element.
Right Click	Perform a right click on an element.
Move/Hover	Move the cursor over an element.
Drag	Click and drag an element to a new position.
Scroll	Scroll up or down the page.
Horizontal Scroll	Scroll left or right on the page.
Press	Press a key on the keyboard.
Keyboard Hotkey	Use a keyboard shortcut or hotkey.
Write	Type text using the keyboard.

6 RESULTS

Our main findings on the impact of ordering are in Table 4. We utilize our various findings to improve upon OmniACT; our experiments against their baseline are in Table 5.

6.1 IMPACT OF ORDERING

421
 422
 423
 424
 426
 426
 427
 428
 429
 429
 429
 420
 420
 421
 421
 422
 423
 424
 425
 425
 424
 425
 425
 425
 421
 421
 422
 423
 424
 424
 425
 425
 424
 425
 425
 425
 424
 425
 425
 425
 426
 427
 428
 428
 429
 429
 420
 420
 421
 421
 422
 423
 424
 425
 425
 424
 425
 425
 425
 426
 427
 428
 428
 429
 429
 420
 420
 421
 421
 422
 423
 424
 424
 425
 425
 424
 425
 425
 426
 427
 428
 428
 429
 429
 420
 420
 420
 421
 421
 422
 422
 423
 424
 425
 425
 424
 425
 425
 426
 427
 428
 428
 429
 429
 420
 420
 421
 421
 422
 421
 422
 423
 424
 425
 424
 425
 425
 425
 426
 427

t-SNE Best For Larger Models And More Challenging Tasks. Navigating by using detected
elements is a harder task than navigating by using human annotated bounding boxes; not only are
there more elements—on average double the amount—there is the possibility that the correct element
is missing from the detected elements. We see that when elements are derived from the DOM and
our UI detection model, t-SNE ordering generally outperforms raster ordering. Additionally, more
powerful models see an increased benefit from t-SNE ordering with Gemini-1.5 and GPT-4v seeing
larger improvements than LLama3.

455

456

457

458 459

460 461

462

464

467

468

469

470

471

472

473

474

475

476 477

478 479

480

481

484

485

Table 4: Performance of different ordering methods across various models and information scenarios. 433 The baseline approach for VisualWebArena is the same as their paper. Human annotations are 434 from OmniACT's annotation files. The Faster-RCNN model is trained to detect interactable UI 435 elements from CommonCrawl webpages. VisualWebArena is evaluated on success rate. OmniACT 436 is evaluated on unweighted action score (i.e. each task is weighted equally). We use Llama3-70B 437 for VisualWebArena and Llama3-8B for OmniACT due to its large size². GPT-4v is evaluated on 438 a 100 random tasks for OmniACT; the exact list is in Appendix A.8. Gemini-1.5 and GPT-4v are 439 multimodal while Llama3 is text only. 440

Expe	Success Rate				
Element Source	Benchmark	Ordering Method	Gemini-1.5 (†)	GPT-4v (†)	Llama3 (†
Ground Truth (DOM)	VWA	Pre-order	64.03%	74.07%	27.79%
		Random	37.04%	37.04%	20.37%
		Raster	38.88%	53.70%	29.63%
		t-SNE	44.44%	61.11%	24.07%
Human	OmniACT	Random	57.29%	61.52%*	28.67%
Annotations		Raster	61.04%	65.88%*	33.65%
		t-SNE	59.17%	62.11%*	31.99%
Detected	OmniACT	Random	39.59%	44.63%*	18.88%
(Faster-RCNN)		Raster	45.21%	47.38%*	21.58%
		t-SNE	47.16%	49.18%*	24.61%

Raster Ordering Performs Best With Human Annotations. Raster ordering performs the best with human annotated elements. Unfortunately, these annotations are fewer and partially guaranteed to contain important information. Additionally, high quality human annotations are difficult to scale across applications.

6.2 STATE-OF-THE-ART PERFORMANCE ON OMNIACT

We achieve a new state-of-the-art performance on OmniACT. Due to cost, we look to our previous experiments to pick the best combination of features for our approach. We observe that multimodal representations are still helpful. We find that t-SNE ordering improves performance the best in 463 OmniACT when elements are detected by our model. Koh et al. (2024) states that high level actions are easier for an LM to reason with. 465

466 We analyze the differences between our best approach and OmniACT's baseline. These are as follows.

- Element Source: OmniACT obtains elements by searching for icons with Segment Anything (Kirillov et al., 2023) and text with EasyOCR. Unfortunately, there are no shared artifacts for their icon detection system. We obtain UI elements through an object detection model and text with EasyOCR (AI, 2020).
- Ordering: It is unclear how elements are ordered in OmniACT. Considering how most approaches don't pay specific attention to ordering, we assume the ordering in OmniACT is effectively random. We order our elements using our t-SNE ordering.
 - Action Space: OmniACT directly outputs pyautogui code as their actions. We consider a higher level action space that maps to pyautogui code.
 - Intractability OmniACT lists out each element, but does not indicate which element is interactable. We specify whether elements are interactable or not.
 - Multimodal Representation: OmniACT evaluates their full test set using a text only representation³. We experiment with a multimodal representation.

We apply our findings and achieve more than one-fold increase over the existing best action score. 482 Our results can be found in Table 5. 483

²We use the Groq API for our Llama3 models.

³OmniACT evaluates the impact of adding a visual representation on smaller subset; however, this subset is not shared and varies empirically from the full test set.

487 Table 5: \checkmark and x indicates whether the feature is available when building a	representation for the
488 model. Ours indicates a high level actions such as click [1]. Code indicat	es that pyautogui code
is directly generated. bold is best and <i>italics</i> are second best. While sequence	e score only checks for
the correct high level action (e.g. click), action score checks for both the	correct action and the
correct element or parameter (e.g. click [1]). Thus, action score is the n	nost equivalent to task
success rate. * is as reported in Kapoor et al. (2024)	

Model	State	Actions	Screenshot	Action Score (↑)	Sequence Score (\uparrow)
GPT-4	OmniACT	Code	x	11.60*	32.75*
Gemini 1.5	Ours	Code	X	16.53	21.67
Gemini 1.5	Ours	High Level	X	22.29	29.42
Gemini 1.5	Ours	High Level	\checkmark	22.86	28.91
Llama8b	Ours	High Level	X	18.64	26.22
GPT-4v	Ours	High Level	\checkmark	23.34	30.47
		0		I	

501 502

508

FUTURE WORK 7

504 **Further Improvements to Ordering** We provided two simple methods to apply ordering when 505 a default ordering is not given. However, both approaches still fall short when compared to the hierarchical ordering derived from the DOM. We hope that future research can introduce more 506 sophisticated methods to find ordering with only pixel information. 507

Image Only Ordering We focused on the impact of element ordering in a text representation 509 (although element labels re-ordered accordingly in the visual representation as well). The impact 510 of element ordering may or may not generalize to an image only representation. Unfortunately, our 511 results indicated that a visual representation alone was insufficient for web and desktop environments 512 which prevents us from conducting this experiment. However, Yan et al. (2023) found that a visual 513 representation in mobile environments was able to achieve comparable performance with and without 514 a text representation. In the future, we hope to experiment with various ordering methods on an 515 image only representation.

516

517 Expanded Scenarios and Benchmarks In this paper, we explored two benchmarks— 518 VisualWebArena and OmniACT—as web and desktop scenarios. In the future, we hope to explore other benchmarks and settings with our approach. For example, Xie et al. (2024) uses the OS level 519 accessibility tree for desktop agent navigation. We hope to compare our approach against theirs and 520 believe that combining both approaches may lead to further improvements. Additionally, mobile 521 environments often have only pixel-level information (Chen et al., 2020; Ross et al., 2020); we hope 522 to apply our approach to a mobile benchmark such as Rawles et al. (2023). 523

- 524 525
- 526 527

529

530

531

533 534

536

537

8 CONCLUSION

We conducted thorough ablations to show that element ordering has a significant impact on the performance of agents. We provided a method of ordering elements through dimensionality reduction 528 and showed that it performed best in realistic environments. We trained a UI element detection model on Common Crawl data and publicly share the model. We demonstrated an end-to-end method which allows a LM agent to act on environments that provide only pixel information. Using this method, we were able to achieve a new state-of-the-art performance on OmniACT. 532

REFERENCES

- J. AI. Easyocr: A simple and robust ocr library. https://github.com/JaidedAI/EasyOCR, 2020. Accessed: 2024-05-07.
- R. Bonatti, D. Zhao, F. Bonacci, D. Dupont, S. Abdali, Y. Li, Y. Lu, J. Wagle, K. Koishida, A. Bucker, 538 L. Jang, and Z. Hui. Windows agent arena: Evaluating multi-modal os agents at scale, 2024. URL https://arxiv.org/abs/2409.08264.

540 541 542 543	J. Chen, C. Chen, Z. Xing, X. Xu, L. Zhu, G. Li, and J. Wang. Unblind your apps: predicting natural- language labels for mobile gui components by deep learning. In <i>Proceedings of the ACM/IEEE</i> 42nd International Conference on Software Engineering, ICSE '20. ACM, June 2020. doi: 10.1145/3377811.3380327. URL http://dx.doi.org/10.1145/3377811.3380327.
544 545	Common Crawl. Common crawl. https://commoncrawl.org, 2023. Accessed: 2023-05-19.
546 547 548	D. Dupont. Gpt-4v-act. https://github.com/ddupont808/GPT-4V-Act, 2024. Accessed: 2024-05-19.
549 550	I. Gur, H. Furuta, A. Huang, M. Safdari, Y. Matsuo, D. Eck, and A. Faust. A real-world webagent with planning, long context understanding, and program synthesis, 2024.
551 552 553	H. He, W. Yao, K. Ma, W. Yu, Y. Dai, H. Zhang, Z. Lan, and D. Yu. Webvoyager: Building an end-to-end web agent with large multimodal models, 2024.
554 555 556 557	F. Huq, J. P. Bigham, and N. Martelaro. What's important here?: Opportunities and challenges of LLM in retrieving information from web interface. In <i>RO-FoMo:Robustness of Few-shot and Zero-shot Learning in Large Foundation Models</i> , 2023. URL https://openreview.net/forum?id=Jd8mD3SU8j.
558 559 560	Ishan0102. vimgpt. https://github.com/ishan0102/vimGPT, 2024. Accessed: 2024-05-19.
561 562 563	R. Kapoor, Y. P. Butala, M. Russak, J. Y. Koh, K. Kamble, W. Alshikh, and R. Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web, 2024.
564 565	G. Kim, P. Baldi, and S. McAleer. Language models can solve computer tasks, 2023.
566 567	A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, WY. Lo, P. Dollár, and R. Girshick. Segment anything, 2023.
568 569 570	J. Y. Koh, R. Lo, L. Jang, V. Duvvur, M. C. Lim, PY. Huang, G. Neubig, S. Zhou, R. Salakhutdinov, and D. Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks, 2024.
571 572	J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023.
573 574 575 576	TY. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In <i>European conference on computer vision</i> , pages 740–755. Springer, 2014.
577 578 579 580	F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten- hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. <i>Journal of Machine Learning Research</i> , 12:2825–2830, 2011.
581 582 583	C. Rawles, A. Li, D. Rodriguez, O. Riva, and T. Lillicrap. Android in the wild: A large-scale dataset for android device control, 2023.
584 585	S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks, 2016.
586 587 588 589	A. S. Ross, X. Zhang, J. Fogarty, and J. O. Wobbrock. An epidemiology-inspired large-scale analysis of android app accessibility. ACM Trans. Access. Comput., 13(1), apr 2020. ISSN 1936-7228. doi: 10.1145/3348797. URL https://doi.org/10.1145/3348797.
590 591 592 593	T. Shi, A. Karpathy, L. Fan, J. Hernandez, and P. Liang. World of bits: An open-domain platform for web-based agents. In D. Precup and Y. W. Teh, editors, <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 3135–3144. PMLR, 06–11 Aug 2017. URL https://proceedings.mlr.press/v70/shi17a.html.

594 595 596	Similarweb. Top websites ranking, 2024. URL https://www.similarweb.com/ top-websites/. Accessed: 2024-05-21.
597 598 599	Statcounter Global Stats. Mobile screen resolution stats worldwide, 2024. URL https://gs.statcounter.com/screen-resolution-stats/mobile/worldwide. Accessed: 2024-05-21.
600 601	L. Van der Maaten and G. Hinton. Visualizing data using t-sne. <i>Journal of machine learning research</i> , 9(11), 2008.
602 603 604	J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
605 606	World Wide Web Consortium. Selectors api level 1. W3C Recommendation, Feb 2013. Available online: https://www.w3.org/TR/selectors-api/#findelements.
608 609	J. Wu, S. Wang, S. Shen, YH. Peng, J. Nichols, and J. P. Bigham. Webui: A dataset for enhancing visual ui understanding with web semantics, 2023.
610 611	Y. Wu, A. Kirillov, F. Massa, WY. Lo, and R. Girshick. Detectron2. https://github.com/facebookresearch/detectron2, 2019.
613 614 615	T. Xie, D. Zhang, J. Chen, X. Li, S. Zhao, R. Cao, T. J. Hua, Z. Cheng, D. Shin, F. Lei, Y. Liu, Y. Xu, S. Zhou, S. Savarese, C. Xiong, V. Zhong, and T. Yu. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments, 2024.
616 617 618	A. Yan, Z. Yang, W. Zhu, K. Lin, L. Li, J. Wang, J. Yang, Y. Zhong, J. McAuley, J. Gao, Z. Liu, and L. Wang. Gpt-4v in wonderland: Large multimodal models for zero-shot smartphone gui navigation, 2023.
619 620 621	J. Yang, H. Zhang, F. Li, X. Zou, C. Li, and J. Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v, 2023.
622 623	S. Yao, H. Chen, J. Yang, and K. Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents, 2023.
624 625 626 627	X. Zhang, L. de Greef, A. Swearngin, S. White, K. Murray, L. Yu, Q. Shan, J. Nichols, J. Wu, C. Fleizach, A. Everitt, and J. P. Bigham. Screen recognition: Creating accessibility metadata for mobile applications from pixels, 2021.
628 629	B. Zheng, B. Gou, J. Kil, H. Sun, and Y. Su. Gpt-4v(ision) is a generalist web agent, if grounded, 2024.
630 631 632	S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, T. Ou, Y. Bisk, D. Fried, U. Alon, and G. Neubig. Webarena: A realistic web environment for building autonomous agents, 2023.
633	
634	
636	
637	
638	
639	
640	
641	
642	
643	
644	
645	
646	
647	

648 A APPENDIX

A.1 FASTER-RCNN TRAINING DETAILS

We trained our model using the detectron2 Wu et al. (2019) implementation of faster-rcnn. We did not change much from the default implementation and recognize that there are significant improvements that could be made to the model.

Table 6: Key hyperparameters for the Faster-RCNN model.

Hyperparameter	Value
Base Learning Rate	0.00025
Number of Classes	1
Iterations	200000
Optimizer	SGD
Backbone	Resnet-50 (ImageNet Pretrained)
ResNet Depth	50
Images per Batch	16
Objects per Image	128
Devices	8

We share the remaining hyperparameters in a config file. We also share the model artifacts and dataset in our github.

A.2 FASTER-RCNN EVALUATION

We evaluated our detection results using the standard COCO evaluation metrics Lin et al. (2014), which include mean Average Precision (mAP) across IoU thresholds from 0.50 to 0.95, as well as performance breakdowns across different object scales.

Table 7: Object detection performance metrics across different dataset splits using out Faster R-CNN Ren et al. (2016) model. AP represents Average Precision, with subscripts denoting IoU thresholds. AP_S, AP_M, and AP_L represent performance on small ($< 32^2$ pixels), medium (32^2 to 96^2 pixels), and large ($> 96^2$ pixels) objects respectively. All values are percentages (%).

Split	AP _{50:95}	AP_{50}	AP ₇₅	AP_S	AP_M	AP_L	AR ₁₀₀
Train Validation	68.94 69.40	77.53 77.51	73.08 73.18	63.61 62.94	73.95 75.04	62.55 62.16	73.40 73.80
Test	66.87	75.47	70.96	64.77	70.16	62.10	71.30

A.3 DATASET DETAILS

We built our dataset by automatically annotating Common Crawl webpages. We split the dataset into a randomized 70/20/10 (train/dev/test) split.





Figure 5: Distribution of the number of bounding boxes per image. This histogram shows the density of annotations across our dataset.

A.4 EFFECTS OF ELEMENT COUNT ON ORDERING

We notice that the relative impact of element ordering becomes more significant as interfaces grow in complexity (Table 9). This affects both raster and t-SNE orderings similarly.

Table 9: Difference in task success rate compared to random baseline for Raster and t-SNE across different ranges of elements in the interface. The split is set to the median number of elements observed across VisualWebArena and OmniACT datasets.

Element Cour	nt Raster	t-SNE
0-25	+2.8%	+3.2%
25+	+6.3%	+6.5%

A.5 T-SNE HYPERPARAMETERS

We evaluated t-SNE across various perplexity values. In Table 10 we can see that varying perplexity values has a neglible effect to overall performance.

Table 10: Action scores achieved by the model across different t-SNE perplexity values.

Perplexity	10	20	30	40	50
Action Score	18.25	17.67	17.95	18.37	18.47

A.6 VISUALWEBARENA AGENT ACTION SPACE

We use the same action space as VisualWebArena for all VisualWebArena experiments.

813 814

812

815

816

826 827

828

829 830

846

849

850

817 We use the same action space as described 818 in VisualWebArena. VisualWebArena uses 819 high level actions that act directly on el-820 ements rather than pixel coordinates. In-821 teractable elements possess a unique id 822 identifier while non-interactable elements do not. The id identifier reflects the po-823 sition of the element in the ordering. The 824 full action space is in Table 11. 825

Action	Description
click [id]	Click on element id.
hover [id]	Hover on element id.
type [id] [text]	Type text on element id.
press [key_comb]	Press a key combination.
new_tab	Open a new tab.
tab_focus [index]	Focus on the i-th tab.
tab_close	Close current tab.
goto [url]	Open url.
go_back	Click the back button.
go_forward	Click the forward button.
scroll [up down]	Scroll up or down the page.
stop [answer]	End the task with an optional output.

Table 11: The set of possible actions in VisualWebArena.

A.7 LM AGENT HYPERPARAMETERS AND SETTINGS

We set our temperature, top-p, and input 831 token limits based on existing works Zhou 832 et al. (2023); Koh et al. (2024); Kapoor et al. 833 (2024). Prompts for all three backbones 834 contain few-shot examples and use chain-835 of-thought prompting Wei et al. (2023) as 836 described in Koh et al. (2024). In GPT-837 4v and Llama3 each example is a different 838 message; Gemini-1.5's context length al-839 lowed us to input all examples as a single prompt. We detail our LM Agent hyperpa-840 rameters in Table 12 841

Setting	Language Model Backbone					
	GPT-4v	Gemini-1.5	Llama3			
Input Token Limit	3840	900000	3840			
Temperature	1.0	1.0	1.0			
Тор-р	0.9	0.9	0.9			

Table 12: Settings for different LM agent backbones.

A.8 VISUALWEBARENA AND OMNIACT SUBSET

We experimented with subsets of VisualWebArena and OmniACT to save on costs. We list them herefor reproducibility.

For all VisualWebArena experiments we used the following:

```
851[13, 15, 50, 129, 164, 167, 0, 77, 86, 89, 98, 99, 100, 101, 105,852130, 131, 142, 143, 146, 150, 189, 16, 29, 37, 38, 39, 47, 49, 52,85353, 56, 60, 61, 62, 69, 73, 76, 77, 81, 148, 173, 193, 196, 201,854212, 216, 231, 273, 314, 315, 322, 445]
```

855 For GPT-4v ordering ablations on OmniACT we used the following:

856 [4, 58, 115, 147, 156, 162, 165, 178, 179, 194, 204, 218, 235, 857 240, 248, 297, 353, 374, 391, 392, 395, 404, 409, 419, 434, 462, 858 487, 492, 517, 533, 556, 573, 598, 658, 667, 673, 678, 719, 795, 859 827, 896, 910, 944, 961, 975, 1018, 1025, 1038, 1084, 1093, 1101, 860 1103, 1128, 1130, 1138, 1142, 1147, 1181, 1192, 1219, 1252, 1284, 861 1291, 1353, 1427, 1442, 1448, 1514, 1521, 1538, 1580, 1590, 1594, 1600, 1606, 1622, 1636, 1641, 1665, 1684, 1694, 1696, 1710, 1711, 862 1719, 1726, 1731, 1740, 1743, 1845, 1877, 1883, 1918, 1924, 1951, 863 1960, 1993, 1994, 1997, 2011]

A.9 PROMPT

864

866	Listing 1: Language Model Prompt
867	Vey are an autonomous intelligent agent tasked with newigating dealter
868	and web applications. You will be given tasks that can be
869 870	accomplished by various actions that will be mapped to pyautogui code
871	
872	Here's the information you'll have:
873	The user's objective: This is the task you're trying to complete.
874	The current desktop screenshot: This is a screenshot of the desktop or
875	webpage, with each interactable element assigned a unique numerical
876	id. Each bounding box and its respective id shares the same color.
877	current screenshot with their text content if any in the format [id]
878	[tagType] [text content]. tagType is the type of the element. text
870	content is the text content of the element. For example, [1234] [
880	BUTTON] ['Add to Cart'] means that there is a button with id 1234 and
881	text content 'Add to Cart' on the current web page. [] [StaticText]
882	interactable
883	
884	The actions you can perform fall into two categories:
885	
886	Mouse Actions:
887	click [10]: Inis action clicks on an element with a specific 10.
888	specific id.
889	right_click [id]: This action right clicks on an element with a specific
890	id.
891	hover [id]: Hover over an element with id.
892	Keyboard Actions:
893	type [content]: Use this to type content. Be sure to use other commands
894	to click before or press enter after if necessary.
895	press [key_comb]: Simulates the pressing of a key combination on the
896	hotkey [key1] [key2]: Simulates the pressing of a multiple key
897	combinations on the keyboard. For example, hotkey [Ctrl] [Alt] [
898	Delete] will press Ctrl+Alt+Delete.
899	The because full it is now important to fallow the fallowing walks.
900	To be successful, it is very important to follow the following rules:
901	observation. Everything is possible. You MUST issue actions.
902	2. You can issue a sequence of actions separated by newlines.
903	3. You should follow the examples from past messages to reason step by
904	step and then issue the next actions.
905	5. Generate the actions in the correct format. Start with a "In summary.
906	the next actions I will perform are" phrase, followed by the actions
907	inside ```. Each action should be split by a newline. There should be
908	no text inside `` except for the actions. For example, "In summary,
909	the actions I will perform are ''Click [1234] type [sample text]
910	press [enter] .
911	Here are a few examples:
912	Example 1:
913	(Example 1)
914	[nyamhre 1]
915	Example 2:
916	
917	{Example 2}

xample 3:						
Example 3}						
haga wara t	be examples	Nou make	a prodictio	a given th	o obcontration	
nose were t	lie examples	. NOW Make	a predictio.	i given ti	le observation	•
BSERVATION:	:					
Observatior	ן ב					
	<pre>kample 3: Example 3} nose were t BSERVATION: Dbservation</pre>	<pre>kample 3: Example 3} hose were the examples BSERVATION: Dbservation}</pre>	<pre>kample 3: Example 3} hose were the examples. Now make BSERVATION: Dbservation}</pre>	<pre>kample 3: Cxample 3} nose were the examples. Now make a prediction SSERVATION: Dbservation}</pre>	<pre>kample 3: 2xample 3) hose were the examples. Now make a prediction given th 3SERVATION: 2bservation)</pre>	<pre>cample 3; ixample 3; nose were the examples. Now make a prediction given the observation SSERVATION: Dbservation}</pre>