

AccessMenu: Enhancing Usability of Online Restaurant Menus for Screen Reader Users

1 Nithiya Venkatraman*

2 Old Dominion University

3 Department of Computer Science

4 Norfolk, Virginia, USA

5 nvenk001@odu.edu

6 Yash Prakash

7 Old Dominion University

8 Department of Computer Science

9 Norfolk, Virginia, USA

10 yprak001@odu.edu

11 Akshay Kolgar Nayak*

12 Old Dominion University

13 Department of Computer Science

14 Norfolk, Virginia, USA

15 anaya001@odu.edu

16 Hae-Na Lee

17 Michigan State University

18 Department of Computer Science and

19 Engineering

20 East Lansing, Michigan, USA

21 leehaena@msu.edu

22 Suyog Dahal

23 Old Dominion University

24 Department of Computer Science

25 Norfolk, Virginia, USA

26 sdaha005@odu.edu

27 Vikas Ashok

28 Old Dominion University

29 Department of Computer Science

30 Norfolk, Virginia, USA

31 vganjigu@odu.edu

ABSTRACT

Online food ordering has become commonplace due to its convenience. The wide variety of culinary choices, combined with fast and economical door-delivery services, encourages more people to order food online. To facilitate this process, food vendors, including restaurants, often provide full menus on their websites, typically in visual formats such as images or PDFs. While this is convenient for sighted users, blind and visually impaired (BVI) individuals face significant challenges accessing these visual menus with their screen reader assistive technology. An interview study with 12 BVI screen reader users revealed that present assistive tools do not adequately satisfy the needs of these users, with issues ranging from text-ordering errors, to inaccurate inferences (e.g., incorrectly categorizing a Caesar salad with anchovies as vegetarian), to misinterpretation of symbols and legends. Moreover, the users expressed a need for a screen reader-tailored interface to access the information in menus. To address these access barriers and users' needs, we present AccessMenu, a browser extension that automatically detects visual menus in restaurant websites, uses multi-modal large language models to extract and analyze the menu content, and re-renders it in a conveniently navigable HTML format accessible with screen readers. AccessMenu also enables BVI users to issue natural language queries, allowing them to efficiently distill specific information from the menus. In a user evaluation with 10 blind participants, AccessMenu significantly outperformed a state-of-the-art solution in usability and task workload, by providing convenient menu navigation and query-based menu filtering capabilities.

*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

CCS CONCEPTS

• Human-centered computing → Accessibility technologies;
Empirical studies in accessibility.

KEYWORDS

blind, screen reader, visual impairment, restaurant menu, usability, accessibility, large language models

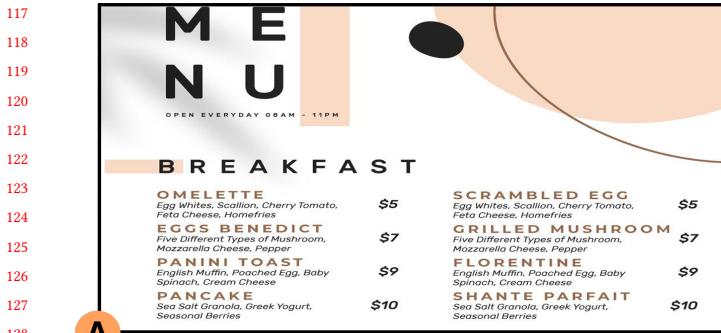
ACM Reference Format:

Nithiya Venkatraman, Akshay Kolgar Nayak, Suyog Dahal, Yash Prakash, Hae-Na Lee, and Vikas Ashok. 2018. AccessMenu: Enhancing Usability of Online Restaurant Menus for Screen Reader Users. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation emai (Conference acronym 'XX)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

The landscape of the restaurant industry has witnessed a profound transformation in the last decade with the proliferation of online food ordering platforms. The convenience offered by online platforms has revolutionized the way consumers engage with restaurants and other dining establishments. According to recent statistics, the global online food delivery market has experienced exponential growth of \$294 billion¹. In the fast-moving and busy world, ordering food online from restaurants has become more efficient and convenient for people all over the world. To facilitate convenient online ordering, food establishments present digital online menus on their websites, so that customers can obtain an overview of available dishes along with associated information such as price, ingredients, customization options, and sometimes even pictures. While these menus significantly elevate the food-ordering experience for sighted customers, the menus pose significant access challenges for blind and visually impaired (BVI) customers, particularly those who interact with digital content using screen reader assistive technology (e.g., JAWS [72], NVDA [7], VoiceOver [10]).

¹https://en.wikipedia.org/wiki/Online_food_ordering

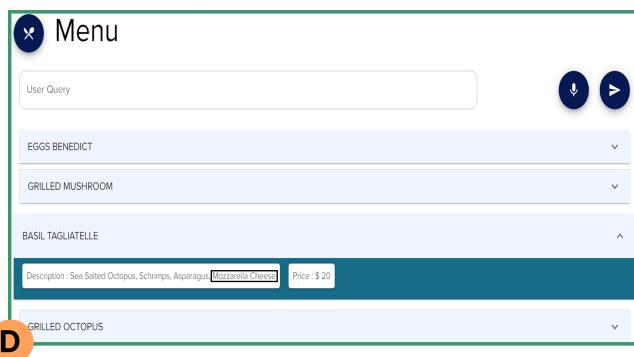


Original Restaurant Menu

4410851 6 - OCR Results	
OMELETTE	\$5
Egg Whites, Scallion, Cherry Tomato, Feta Cheese, Home fries	
EGGS BENEDICT	\$7
Five Different Types of Mushroom, Mozzarella Cheese, Pepper	
PANINI TOAST	\$9
English Muffin, Poached Egg, Baby Spinach, Cream Cheese	
PANCAKE	

175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231

OCR Output with JAWS



AccessMenu Response Interface



AccessMenu Proxy Interface

Figure 1: (A) The original restaurant menu. (B) Output from JAWS Convenient OCR. (C) The AccessMenu interface, where the red box highlights the natural language query field, the yellow box indicates the voice-input button, and the green box is the submit button. (D) The updated AccessMenu interface showing only the subset of menu items matching the user's query.

A screen reader narrates content and assists BVI users in navigating the web using keyboard shortcuts. However, this navigation is predominantly one-dimensional, requiring users to methodically traverse elements on a web page to find the desired element or information. While text elements are simply read out using text-to-speech, visual elements such as images are handled via either website-provided alternative textual descriptions, i.e., alt text, or automatic AI-generated texts describing the visual content [5, 6, 31, 41]. While such texts are suitable for simple images that can be fully described using captions, they are impractical and inadequate in case of complex two-dimensional document images such as restaurant menus (see Figure 1A). Consequently, blind screen reader users typically rely on AI-driven assistive tools [22, 37, 53, 75, 81] to access information in image documents such as menus.

However, in an interview study with 12 BVI participants, we found that the present assistive tools are inadequate in their ability to address the BVI users' challenges and needs with regards to interaction with visually complex and information-rich documents such as restaurant menus. The participants stated that interacting with assistive-tools' OCR outputs was cumbersome and mentally taxing with a screen reader, as the information layout in the outputs often did not *logically* match the screen reader narration order. For example, in the information extracted from the menu in Figure 1A, the screen reader reads out "Scrambled Egg" after "Omelette", instead

of reading out the ingredients and price of "Omelette". The participants also reported hallucinations in AI applications due to lack of contextual awareness (e.g., answering 'vegetarian' to a question about a Caesar salad in the menu that clearly lists anchovies as one of the ingredients) and legend misinterpretations (e.g., mistaking vegetarian icon for the vegan icon). Additionally, a majority of the participants expressed a need for an alternative BVI-friendly user interface for quick-and-easy perusal of menu items, and they also provided design ideas regarding this user interface.

Informed by the findings of our interview study, we designed and developed AccessMenu, a browser extension that automatically builds a semantics-based menu model (list of menu items and properties of each item including customization options for each item, item category, etc.) from online restaurant menu images using multimodal large language models (MLLMs), and then leverages this menu model to provide an alternative menu interface tailored for both convenient screen reader navigation and natural language query access. As illustrated in Figure 1C, AccessMenu presents users with a *proxy* menu interface, where the linearly-organized menu items are conveniently navigable using basic screen reader shortcuts. Moreover, as seen in the figure, AccessMenu enables the user to quickly and accurately access specific information about menu items via natural language queries (e.g., *suggest meat-based appetizers*). Key technical innovations of this work include the development of a robust multimodal pipeline capable of extracting

and contextualizing menu information, and a seamless integration of this pipeline into a real-time browser extension. This system effectively bridges the gap between visual menu content and accessibility requirements, offering BVI users a faster and more intuitive interaction with online menus.

An evaluation of AccessMenu in a user study with 10 blind participants showed significant improvements in the usability, task workload and overall user experience with online restaurant menus, compared to a state-of-the-art OCR-based solution. A majority of the participants also stated that the AccessMenu interface would motivate them to be more active in ordering food online and consider a plethora of food options and make an informed dining choice. In sum, this paper makes the following contributions.

- The findings of an interview study detailing the usability issues faced by BVI screen reader users while accessing restaurant menus online for ordering food.
- The design and evaluation of a novel AccessMenu browser extension that presents inaccessible/unusable restaurant menus via alternative BVI-friendly interfaces using MLLMs.

2 RELATED LITERATURE

2.1 Web Interaction Using Screen Readers

Extensive research has been conducted to explore the complexities and challenges involved in interacting with web content using a screen reader [4, 9, 12, 13, 16, 43, 44, 64, 80]. An early work by Lazar et al. [43] identified significant access barriers that persist despite established accessibility guidelines [5], including poorly designed page layouts, technical conflicts between screen readers and web applications, and the absence of alt text for images. Borodin et al. [16] further investigated the strategies adopted by screen reader users to circumvent the web-interaction issues, and found that the users typically resorted to increasing the speech rate and using the headings' hotkeys to efficiently navigate the web content. More importantly, they observed that the screen reader vocabulary of most users was limited to a handful of basic screen reader shortcuts. Similar to these seminal works, more recent research efforts in this area have also investigated and uncovered numerous accessibility and usability issues of screen reader users in different web interaction scenarios [33, 70, 79]. While such investigations of issues are generic across the web, and therefore applicable to a certain extent to online restaurant menus, they do not capture the unique domain-specific issues that screen reader users face when interacting with online restaurant menus. We address this knowledge gap via an interview study with 12 blind participants, aiming to uncover their pain points, needs, and preferences while navigating and interacting with restaurant menus online.

Prior research has also explored solutions to overcome the numerous accessibility and usability challenges for screen reader users [11, 26, 27, 30, 36, 45, 59, 61, 63–65, 73, 77, 82, 83]. These solutions include automatic captioning of visual content [25, 48, 50, 60], web automation [11, 66, 67, 84], natural language assistants [12, 20, 28, 54], and even alternative third-party navigation devices [14, 15, 27, 30, 36, 46, 61, 64, 65]. While these solutions do significantly enhance usability of web screen reading in general, they are currently limited in their ability to address the specific issues that arise when interacting with online restaurant menus.

The arrangement of content in a typical menu is highly visual, with the document layout itself used to implicitly convey the semantics associated with the listed menu items. Moreover, many restaurant menus are in PDF or image formats so it is not possible to use the aforementioned generic web-based solutions to address the interaction problems. To fill this void, we present AccessMenu, a solution that specifically focuses on enhancing usability of visual-rich online documents, particularly restaurant menus.

2.2 Visual Document Understanding

Visual document understanding (VDU) tasks (e.g., visual question answering) involves the interpretation and analysis of a wide range of digital documents, including but not limited to forms, tables, reports, and academic papers [8, 51, 92]. The techniques employed in VDU can be broadly classified into two primary categories. The first category focuses on accomplishing the VDU tasks by aligning images with annotations sourced from external optical character recognition (OCR) systems [34, 35, 78, 88], whereas the second category comprises approaches that process document images directly, without relying on external OCR tools [42, 47, 55].

LayoutLMv2 [88], a notable example of the first category, leverages OCR to extract text and bounding boxes from visually-rich documents, combining text, layout, and image data for enhanced document understanding. By integrating OCR output during pre-training with spatial-aware self-attention, LayoutLMv2 captures document context more effectively. A contemporary example of the second category, the OCR-free Donut model [42] simplifies VDU tasks by eliminating dependency on OCR engines, directly mapping document images to structured outputs using a ‘transformer-only’ architecture. Through pre-training with custom curated synthetic data (SynthDoG [42]) and fine-tuning across diverse VDU tasks, Donut has demonstrated strong performance and has also been generalized across multiple languages and document types.

More recently, large language models (LLMs) such as LMDX [62], BLIP [49], LLaVA [56], MiniGPT-4 [91], and mPLUG-Owl [90] have demonstrated significant capabilities in accomplishing VDU tasks in visually-rich documents via minimal instructions [19, 86]. However, despite the impressive zero-shot reasoning capabilities demonstrated by multimodal LLMs, studies have shown that these LLMs face challenges in comprehending text-rich images [57]. Recent studies have also explored the effectiveness of MLLMs in Visual Question Answering (VQA), an important VDU task that involves accurately responding to questions based on the visual information in documents such as receipts, forms, and research papers [39, 40, 58]. Current VQA solutions employ an assortment of natural language processing and computer vision techniques to accurately answer posed questions [38, 76, 85]. However, none of the existing multimodal LLMs have been previously investigated for their efficacy in handling unique-style documents such as restaurant menus. In this paper, we conduct an in-depth investigation of MLLMs like GPT-4 [89], Claude [21], and LLaMA 3 [24] for information extraction and reasoning tasks on documents such as restaurant menus. Additionally, we explore the efficacy of MLLMs in comprehending menu-related queries and reasoning logically over menu content to generate valid responses.

349

3 UNCOVERING USABILITY ISSUES

350

We conducted an Institutional Review Board (IRB)-approved semi-structured interview study with 12 blind participants to uncover their current interaction challenges and needs while accessing online restaurant menus.

354

3.1 Participants

355

We recruited 12 blind screen reader users (6 female, 6 male), with an average age of 49.41 years (Median = 49, SD = 16.59, Range = 31–68). The inclusion criteria required the participants to be proficient in web screen reading and familiar with restaurant websites. All participants stated that they order food through phone at least once every week. Also, none of the participants had residual vision good enough to visually interact with digital content using screen magnifiers. The participants did not have any additional impairments, such as motor or hearing difficulties, that could affect their ability to complete study tasks effectively.

367

3.2 Interview Design and Procedure

368

The interviews were semi-structured with seed questions pertaining to the following topics:

369

- **Food ordering habits.** E.g., How often do you order food? How do you order food? How do you choose restaurants for ordering food?
- **Experience with restaurant menus.** E.g., What assistive technologies do you use to access menus online? What issues do you typically face while accessing these menus? How do you tackle these issues?
- **Needs and preferences.** E.g., Do you have any design suggestions for making these menus more screen reader friendly? What kind of additional support do you think you will need to better access online menus?

370

The interviews were conducted remotely via Zoom conferencing software². At the beginning of the study, informed consent was obtained remotely via the DocuSign service [23]. The experimenter then engaged the participant in conversations about the topics, starting with the seed questions. During the interview, the participant was also encouraged to explain responses through illustrations on actual restaurant websites. Each interview lasted about 45 to 60 minutes. Each of the participants was compensated with a \$25 Amazon gift card.

391

3.3 Data Collection and Analysis

392

With the participants' permission, all interviews were audio-recorded and also screen-captured (for capturing illustrations). We did not retain any personal or identifiable information besides the basic demographic details. We analyzed the collected and transcribed qualitative data using the standard open coding technique followed by axial coding [68]; we iteratively went over the user responses and identified key insights and patterns that reoccurred in the data.

400

3.4 Findings

401

The notable themes that emerged from the qualitative analysis are presented next.

405

²<https://www.zoom.com/>

406

Access menu online but order food through phone. Most participants stated that they preferred ordering over a phone call after perusing the menu online either on their computers or smartphones. These participants stated that most restaurant employees do not have the time to patiently describe the menu over the phone, and they often put them on hold for long durations. Therefore, they prefer to be 'more-or-less decided' before calling the restaurant. As for not ordering online, the participants stated that most restaurant websites are not usable and sometimes not accessible, which previously caused them to make mistakes such as ordering extra portions of food and ordering unintended dishes.

Most restaurant menus require additional assistive tools besides screen reader for access to its information. Almost all participants mentioned that they often had to rely on additional tools, predominantly OCR software to access content in menus, since these menus were mostly in image or PDF formats, both of which are not conducive to screen reader-based interaction. A few participants, who were adept at using screen readers, also mentioned using AI assistants often to query information in the menus.

Current assistive tools do not provide sufficient support to interact with menus. Most participants mentioned that the OCR outputs of present assistive tools (e.g., JAWS Convenient OCR, ABYY FineReader) often contained errors or inconsistencies. Moreover, they also stated that mentally parsing OCR output based on audio alone was cognitively taxing, as the screen reader narration order of the OCR output did not often match their expected 'logical' order implicitly conveyed through visual cues. This was best expressed by the participant P8: "*The OCR output is often a mess. I need to figure out which part is linked to which other part. Suppose I hear Appetizers from the screen reader, I am naturally expecting the next thing to be the name of an appetizer, instead I hear eleven dollars, and now I need to figure out which dish costs eleven dollars*". Some of the participants who used other LLM-based assistive tools such as ChatGPT, mentioned that these tools often provided incorrect or confusing responses to their queries. For example, P4 stated: "*I am careful when picking food, because I don't want fish or meat in what I order. I once asked ChatGPT to list vegetarian dishes in a restaurant menu, and its response contained many dishes which my friend said had fish sauce or seafood ingredients in them. Sometimes, I think it also gets confused between vegetarian and vegan, as it only mentions vegan dishes when I ask for all vegetarian dishes*".

Ask friends or family members for obtaining specific information. Nearly two-thirds of the participants stated that they often 'jointly' explored the menus with their sighted companions. The participants further stated that this joint interaction mostly entailed question-answering, where they asked their sighted companions a variety of questions or 'doubts' regarding the menu.

Need for an alternative interface to access menus. All participants specified a need for a 'new' interface to peruse menus using a screen reader. Seven participants mentioned that linear organization of menu items was more convenient for screen reader navigation. Two participants further suggested the idea of a popup

interface that could present the menu items in a linear arrangement, preferably as a list. One of these participants, P5, asked: “*Is it possible to put the menu items in one single list within a popup window? I can then go through menu linearly without missing anything*”. Four other participants suggested including an assistant in the interface for quickly querying information in the menu. One of these participants P2 stated: “*I would rather just ask the AI to give me all the gluten-free menu items instead of going over all the items myself and filtering them out one-by-one.*”

Summary. The interview study revealed several pain points and needs of blind screen-reader users when they interact with online menu documents. From the study observations, it is clear that an alternative non-visual interface is needed that enables users to conveniently navigate the menu items, while also providing an option to query specific information in the menu. Specifically, the interface must enable convenient perusal of menu items, with the items arranged in a simple linear list. All information about a menu item should also be available at one place in the alternative interface, i.e., where the menu item is listed, irrespective of how the information is scattered in the original menu. Guided by these findings, we designed and developed the AccessMenu prototype interface which is described next.

4 SYSTEM DESIGN

4.1 AccessMenu Overview

Figure 2 presents the operational workflow of AccessMenu, embodied as a browser extension, that generates an alternative screen reader-friendly interface to peruse menu items. On any restaurant’s webpage that contains the menu, users can access the AccessMenu’s alternative menu interface using the ‘Ctrl+J’ keyboard shortcut. Specifically, this hotkey triggers the following sequence of operations in the background: (i) Extract the menu items from the image menu by instructing a multimodal large language model (MLLM) with a custom crafted prompt; and (ii) Use the MLLM output to re-render the information of menu items in the AccessMenu’s conveniently-navigable linear menu interface. The AccessMenu’s interface also enables the user to issue natural language queries (e.g., *list only the gluten-free items in the menu*) to obtain specific information about the menu in the interface.

For extracting the menu items using an MLLM, we adopted the Chain-of-Thought (CoT) prompting strategy [87], where we carefully handcrafted ‘reasons’ or ‘thoughts’ to ensure that the MLLM accounted for the unique aspects of information presentation in menus, for instance, use of icons or symbols (e.g., a leaf) next to items with the legend describing these icons/symbols (e.g., vegan) placed somewhere else in the menu. In the prompt, we also included instructions for the MLLM to generate the output or the ‘menu model’ as a collection of JSON objects (i.e., one object per menu item) to ensure consistency and prevent potential ‘phantom information’ arising from model hallucinations. For supporting users’ natural language queries, we again crafted a custom CoT prompt with guardrails and few-shot examples [19] to ensure that the MLLM strictly based its responses on the extracted menu model. The MLLM output in this case too was in JSON format, to facilitate

convenient rendering of the query responses in the AccessMenu’s interface. The details of AccessMenu are provided next.

4.2 Extraction of Menu Data Items

When a user presses the ‘Ctrl+J’ keyboard shortcut to access the AccessMenu’s interface, AccessMenu first captures a series of menu images and sends them to a backend server. This raw input of menu images are diverse and complex comprising a mix of textual, graphical, and decorative elements. These images are then used as input contextual information in a custom ‘prompt’ for instructing an MLLM to accurately extract menu items. We specifically employed Chain-of-Thought (CoT) prompting [87], given its suitability for this task. A snippet of our custom prompt template is shown below.

CoT Prompt Template for Menu Item Extraction

Menu: [INSERT MENU IMAGES HERE]

Task: Extract structured menu info from an image with proper categorization, icon detection, and JSON formatting.

Steps:

- (1) Extract all visible text from the image and identify menu headers, item names, prices, and descriptions.
- (2) Detect visual cues such as icons (e.g., a red chili) and style differences (bold titles, colored texts) and cross-reference with any provided legend.
- (3) If no legend is present, infer icon meanings using common conventions (e.g., a red chili icon indicates spiciness).
- (4) Apply a rigid [JSON schema] to enforce consistent structure.
- (5) Filter out extraneous elements like watermarks, disclaimers, and decorative texts.

Examples:

- **Input:** Image of a restaurant menu
- **Raw Text Output:** [Lunch Specials, Spicy Chicken Burger \$8.99, Caesar Salad \$6.99, ...]
- **Reasoning Steps:** [Detected section header “Lunch Specials” ... parsed item “Spicy Chicken Burger” with price “\$8.99” ... inferred red chili icon implies “spicy” for “Spicy Chicken Burger” ... filtered out decorative footer text ... structured data using the designated [JSON schema]]

• **Final Output:**

```
{
  "menu_items": [
    {
      "name": "Spicy Chicken Burger",
      "description": "Grilled chicken ...",
      "icons": ["spicy"],
      "price": "\$8.99"
    },
    ...
  ]
}
```

As shown above, the prompt comprises different components: (i) Menu snapshots; (ii) Task description; (iii) Sequence of reasoning steps; and (iv) Demonstrative examples. Notice how the reasoning

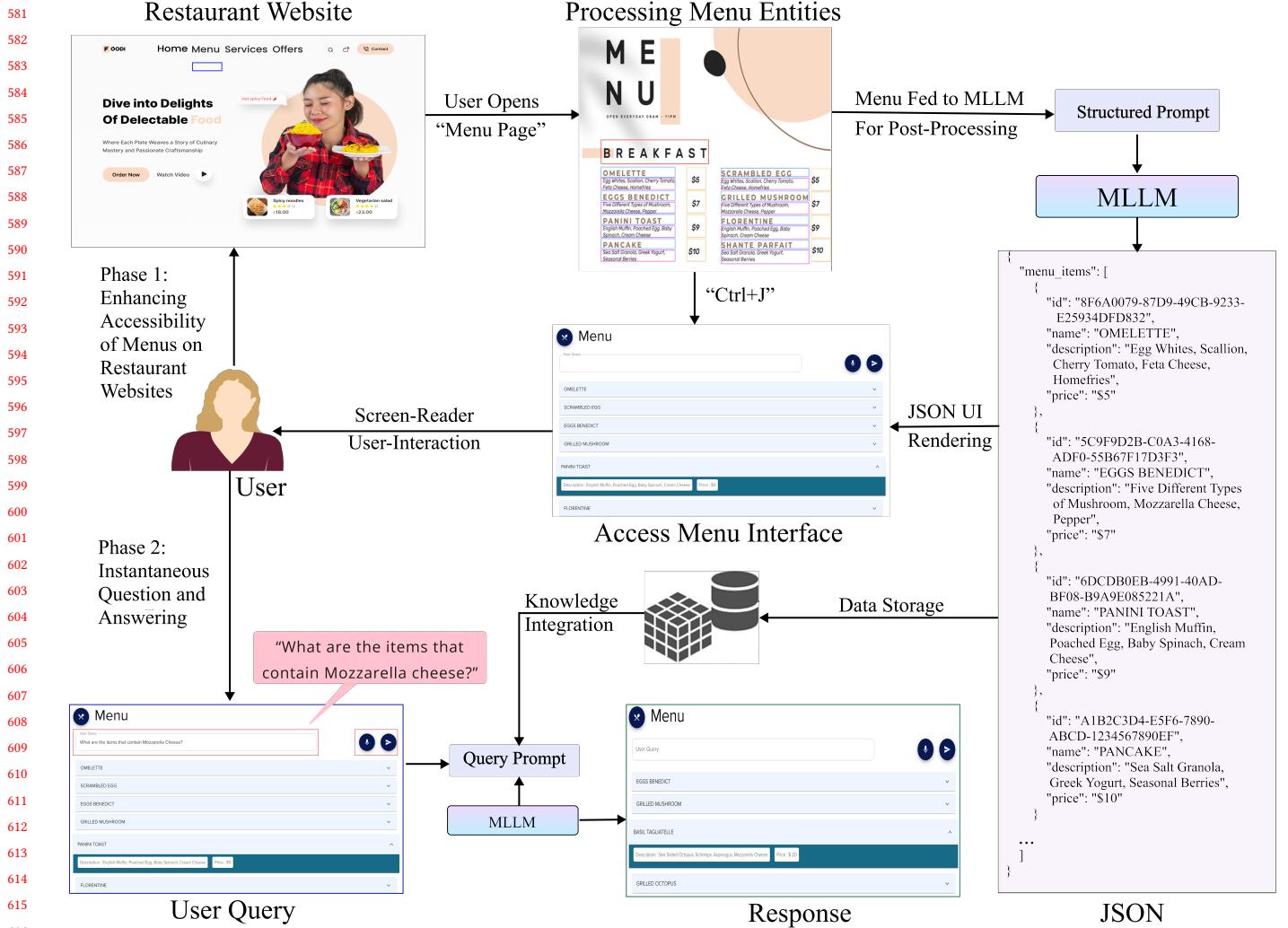


Figure 2: Architecture of AccessMenu consists of two phases: (a) Phase 1 - Information Extraction and (b) Phase 2 - Question Answering.

steps are designed to accommodate the unique aspects of restaurant menus such as spatial relationships, legends, icons, and symbols. As shown, the prompt also instructs the MLLM to structure the output, i.e., the extracted list of menu items, as a collection of JSON objects adhering to a pre-defined fixed schema for ensuring structural and processing consistency. Lastly, the prompt also ensures that the MLLM filters out extraneous elements such as watermarks, disclaimers, or promotional text. The final JSON output, which we henceforth refer to as ‘menu model’, is stored in the back-end server memory for the session, serving as the knowledge base for responding to subsequent user queries (refer to Section 4.3).

4.2.1 Evaluation. To evaluate the extraction performance of our approach, we collected a diverse dataset of 50 menus. The dataset was curated based on criteria such as cuisine type, menu format (e.g., à la carte, set menus), and geographical location to ensure diversity in both content and presentation styles. The dataset included menus

of different layouts, with varied information structures, ranging from simple item listings to complex hierarchical representations featuring categories, subcategories, and additional legend information (e.g., spice level indicators, vegetarian legends). For each of the 50 selected menus, annotators manually created a ground truth dataset in JSON format.

We evaluated the extraction performance of three different MLLMs: GPT-4o-mini [2], Claude-3-5-Sonnet [1], and Llama 3.2-90B-Vision [3]. To measure the models’ ability to capture the various components and relationships within the data presented in the menu, we used three main metrics: (i) **Entity F1 Score (EF1)**: Captured the MLLMs’ ability to extract individual menu elements, such as the names of dishes, their descriptions, prices, and any relevant legends (e.g., dietary symbols like vegan or gluten-free); (ii) **Relationship F1 Score (RF1)**: Assessed the model’s capacity to understand and extract associations between different entities. For example, it measured

697 how well the model linked menu items to their corresponding legends (such as indicating spice levels) or associated prices with the
 698 correct items; and (iii) **Structural F1 Score (SF1)**: Evaluated the
 699 model's ability to maintain the hierarchical organization of the
 700 menu, such as distinguishing between main sections like appetizers,
 701 main courses, and desserts, and further recognizing subsections
 702 or groupings within each category. GPT-4o-mini outperformed the
 703 other models in extracting information from menu images, achiev-
 704 ing an Entity F1 Score of 0.80, a Relationship F1 Score of 0.73, and
 705 a Structural F1 Score of 0.84. In comparison, Claude-3.5-Sonnet ob-
 706 tained an Entity F1 Score of 0.62, a Relationship F1 Score of 0.43, and
 707 a Structural F1 Score of 0.79, while Llama 3.2-90B-Vision achieved
 708 an Entity F1 Score of 0.79, a Relationship F1 Score of 0.61, and a
 709 Structural F1 Score of 0.78. We therefore integrated GPT-4o-mini
 710 in AccessMenu. Note however that AccessMenu follows a modular
 711 architecture, allowing individual components, including the MLLM,
 712 to be easily replaced with a better one if needed in the future.
 713

714 4.3 Processing Contextual User Queries

715 To handle menu-related user queries, we again crafted a similar
 716 ‘Chain-of-Thought’ prompt [87] with few-shot examples, that in-
 717 structed the LLM (GPT-4o-mini [2]) to comprehend and reason
 718 over the extracted ‘Menu Model’ (refer to Section 4.2) to generate
 719 the expected response. As in case of menu extraction, the prompt
 720 included different components: (i) Task description; (ii) User query;
 721 (iii) Menu model providing the context; (iv) Sequence of reasoning
 722 steps to be considered for generating the output; and (v) Few shot
 723 examples covering a variety of queries. The few shot examples
 724 also included ‘negative’ queries, i.e., queries unrelated to the menu
 725 content, to mitigate the impact of model hallucinations. The design
 726 of other menu-related few-shot examples were influenced by the
 727 participants’ feedback in the earlier interview study. These few-
 728 shot queries ranged from simple filtering (e.g., “List all vegetarian
 729 items.”) and single-hop reasoning (e.g., “What are the desserts under
 730 \$10?”) to more complex multi-hop reasoning (e.g., “What are the
 731 gluten-free appetizers with a drink under \$20?”), logical and arith-
 732 metic queries (e.g., “Find me a combination of a main dish and a
 733 dessert for less than \$30, with the main dish being vegetarian.”), and
 734 suggestive queries (e.g., “What’s a good vegan meal with a drink for
 735 under \$25?”). Lastly, the prompt instructed the MLLM to structure
 736 the output in JSON, which AccessMenu then parsed to render the
 737 response in its menu interface.

738
 739 4.3.1 **Evaluation.** To evaluate the quality of responses generated
 740 by our method for user queries, we conducted a study using five
 741 restaurant menus wherein we invited 10 research volunteers to
 742 interact with the menus and pose various menu-related questions.
 743 Each volunteer was provided with 10 minutes per menu, allowing
 744 them to explore and query each of the five menus within a 50-
 745 minute study window. The MLLMs’ responses to these questions
 746 were then evaluated against the ground truth using the F1 Score,
 747 with a final score of 0.83, indicating high similarity (generated vs.
 748 ground truth) and strong overall performance. The inaccuracies
 749 were primarily due to the model’s difficulty in understanding re-
 750 lationships between items placed far apart in the menu. In some
 751 other cases, ambiguity in the user’s phrasing contributed to errors
 752 in the responses. For example, when a user asked, “Give me the
 753

754 healthiest main course dishes,” without specifying the criteria for
 755 “healthy” (e.g., low calorie or vegetarian), the system’s response
 756 varied from what the user expected. In such instances, AccessMenu
 757 occasionally produced a response that did not fully align with the
 758 user’s intent in the query.
 759

760 4.4 User Interface

761 The AccessMenu’s menu interface comprises a query form, a submit
 762 button and a voice-input button at the top followed by a list of
 763 extracted menu items displayed as an accordion. The accordion is
 764 made up of vertically stacked headers representing the names of
 765 items from the menu, which, upon activation (using the ENTER
 766 key), expand to reveal further details about each item (refer Figure 1).
 767 The interface was carefully crafted to enable easy navigation using
 768 simple TAB, ENTER, and ARROW shortcuts. Moreover, the web
 769 elements were optimized for accessibility, employing tab-index and
 770 ARIA (Accessible Rich Internet Applications) attributes [74]. The
 771 tab-index attribute specifies the sequence where elements would
 772 gain keyboard focus, whereas the aria label offers users additional
 773 information about the web element. By default, upon the interface
 774 activation, the initial focus is set to the ‘Query’ form, allowing users
 775 to smoothly transition between the form, the control buttons, and
 776 the accordion via TAB/SHIFT+TAB or ARROW shortcuts. When
 777 the user poses a query, the interface (if needed) simply refreshes
 778 the list of menu items in the accordion based on the MLLM output.
 779 Responses to factual or invalid queries (e.g., what is the price of
 780 omelette? Why is omelette so expensive?) on the other hand are
 781 simply voiced out.
 782

783 4.5 Implementation Details

784 We implemented AccessMenu as a web browser extension, adhering
 785 to the open-source guidelines provided by Google for Chrome
 786 extensions³. When AccessMenu is activated, a service worker ini-
 787 tializes and listens for specific browser events, such as the loading
 788 or closing of a page. Once a menu webpage is loaded, content
 789 scripts are dynamically injected into the page. These JavaScript
 790 files interact with the parent extension code and have access to the
 791 webpage’s DOM, allowing AccessMenu to modify and enhance the
 792 page as needed. To capture menu snapshots, AccessMenu leverages
 793 the services of a Selenium driver [29]. These menu snapshots that
 794 serve as the preliminary input for subsequent extraction process
 795 are sent to the backend server via a POST request. The backend
 796 server was built using Django Rest Framework⁴ and Python mod-
 797 ules were used for all inter-module communication. Integration
 798 of the MLLM into AccessMenu was done using the LangChain
 799 framework⁵, which is known to seamlessly orchestrate query pro-
 800 cessing and response generation. Additionally, the backend was
 801 containerized using Docker⁶ to ensure a consistent environment
 802 across different systems, simplify dependency management, and
 803 enable seamless deployment.
 804

³<https://developer.chrome.com/docs/extensions/mv3/devguide/>

⁴<https://www.djangoproject.org/>

⁵https://python.langchain.com/docs/get_started/introduction

⁶<https://www.docker.com/>

813

5 EVALUATION

814

We conducted an IRB-approved user study with 10 blind screen reader users to assess the effectiveness of AccessMenu and compare it with the status quo OCR-based assistive tool.

815

816

817

818

819

5.1 Participants

820

We enlisted 10 participants with visual impairments (6 female, 4 male), averaging 47.3 years old (Median = 47, SD = 12.7, Range = 23–66), recruited through email lists and snowball sampling. To preserve external validity, we ensured that there was no overlap between the participant groups in this study and the previous interview study. The inclusion criteria required the participants to be proficient in web browsing using the JAWS screen reader, as the study was conducted on the Windows OS platform with JAWS installed as the primary screen reader. Moreover, familiarity with the JAWS Convenient OCR feature [71] was essential, as this was the study baseline condition for assessing AccessMenu. All participants reported accessing online restaurant websites at least once a week. Participant demographics are detailed in Table 1. No participant reported having other difficulties (e.g., hearing, motor control) that could possibly affect their ability to perform study tasks.

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

5.2 Design

838

In a within-subject experimental setup, the participants were asked to freely explore the contents of a restaurant menu under the following two conditions:

839

840

841

842

843

844

845

846

847

848

- **OCR** – The participants were allowed to interact with the textual output generated by JAWS Convenient OCR [71] to access the menu content.
- **AccessMenu** – The participants were allowed to interact with the AccessMenu’s alternative user interface to access the menu content.

We chose this free-form exploration task to emulate real-world scenarios where people typically start perusing menus freely without any specific focus. The participants were also asked to think-aloud during menu exploration. To minimize any learning effects, we ensured that menus from the same restaurant were not used more than once when performing the tasks under different conditions. Instead, we selected menus from two different restaurants for the two conditions. The assignment of restaurant menus to conditions and the ordering of conditions were counterbalanced across the study participants using the well-known Latin-square method [17]. A maximum of five minutes was allotted for each task.

849

850

851

852

853

854

855

856

857

858

859

860

861

862

5.3 Apparatus

863

The study was conducted using a Windows-based Lenovo ThinkPad laptop equipped with all the required software, including the Google Chrome browser, the AccessMenu Chrome extension, and the JAWS screen reader with JAWS Convenient OCR installed. An external QWERTY desktop keyboard was plugged in since all participants mentioned that they were familiar with the standard keyboard during the recruitment process.

864

865

866

867

868

869

870

5.4 Procedure

The experimenter began the study by obtaining the participant’s informed consent and explaining the objectives of the study to the participant. The experimenter then allowed the participant to get familiar with AccessMenu and also configure the screen reader parameters according to their preferences. This was done to ensure that the participant’s comfort level with the study apparatus was more-or-less similar to that with their own computers at home. The experimenter then asked the participant to complete the study tasks according to the predetermined counterbalanced order. After each task, the experimenter administered the SUS and NASA-TLX questionnaires [18, 32] to obtain feedback regarding the usability and task workload respectively for the corresponding study condition. All conversations were in English and the participants were compensated \$25 for their time. Each study lasted about 45 minutes.

5.5 Data Collection and Analysis

Other than the SUS and NASA-TLX responses, we also recorded the participants’ think-aloud utterances while doing the tasks as well as the number of items covered in each task. The experimenter also noted down any peculiar screen reader behavior from the participants while doing the tasks. We analyzed the SUS and TLX responses using standard descriptive and inferential statistical methods. Qualitative data was transcribed and analyzed using open coding and axial coding [69] to identify key insights and themes recurring in the data. We detail our findings next.

5.6 Results

5.6.1 Average number of items covered. All participants fully used the allotted 5 minutes exploring the menu in all the tasks. On average, the participants perused 14.6 items (Median = 15, Minimum = 7, Maximum = 28) under the OCR condition and 30.5 items (Median = 31.5, Minimum = 21, Maximum = 36) under the AccessMenu condition. This difference was found to be statistically significant (Wilcoxon signed rank test, $Z = 2.76$, $W = 0$, $p = 0.005$). Qualitative analysis of the participants’ think aloud responses and experimenter’s notes revealed the causes underlying this significant difference in the number of items covered between conditions. In the OCR condition, the participants were frequently complaining of getting confused by the screen reader output, and therefore they spent extra time going back-and-forth listening to the same content multiple times in order to not only comprehend it but also discover boundaries between the different menu items. Such an issue was not observed in the AccessMenu condition, where the participants went through the list one-by-one in a linear fashion. Also, in the OCR condition, the participants spent time searching for desired information regarding a group of items (e.g., “gluten-free options”), whereas in the AccessMenu condition, they avoided this overhead by simply asking the AccessMenu to filter the menu via a natural language command.

5.6.2 Accuracy of query responses. Overall, 108 queries were issued by the participants during the study, with an average of 10.8 commands ($\sigma = 1.03$) per participant. A manual inspection of the generated AccessMenu responses to these questions revealed a precision of approximately 0.71 and a recall of 0.85, resulting in an

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

ID	Age	Gender	Age of Vision Loss	Occupation	Preferred Screen Reader	Proficiency
P1	43	M	Cannot remember	Unemployed	JAWS	Intermediate
P2	66	F	Since birth	Self-employed	JAWS	Advanced
P3	38	F	Cannot remember	Student	JAWS	Beginner
P4	53	M	Age 3	Self-employed	JAWS	Intermediate
P5	37	F	Since birth	Social worker	JAWS	Intermediate
P6	63	M	Since birth	Corporate	JAWS	Expert
P7	59	F	Cannot remember	Teacher	NVDA	Advanced
P8	40	M	Age 5	Unemployed	JAWS	Intermediate
P9	51	F	Since birth	Corporate	JAWS	Expert
P10	23	F	Since birth	Student	NVDA	Beginner

Table 1: Demographics of blind participants in the evaluation study. All information was self-reported.

F1 score of 0.77. Error analysis revealed that the majority of inaccuracies (82.3%) were caused by ambiguities in user queries and issues in transcribing complex menu items from voice input. Ambiguous questions often led to filtering errors; for instance, when a user asked for “light snacks,” the system struggled to interpret “light” as it could refer to either low-calorie items or small portions. Additionally, voice transcription errors occurred with complex menu item names (e.g., Wagyu with Béarnaise Sauce) leading to inaccuracies in response generation.

5.6.3 Usability and Task Workload. The System Usability Scale (SUS) questionnaire [18] asks participants to respond to alternating positive and negative Likert items on a scale from 1 to 5, where 1 indicating strong disagreement, 3 representing neutrality, and 5 representing strong agreement. These responses are combined into a single usability score between 0 to 100, with higher scores reflecting better usability. As shown in Figure 3a, the AccessMenu condition received significantly higher SUS ratings (Average (μ) = 69.25, Standard Deviation (σ) = 16.36) compared to the screen reader OCR condition (Average (μ) = 46.25, Standard Deviation (σ)

= 11.25), as determined by a one-way ANOVA ($F = 12.08, p < 0.005, \eta^2 = 0.40$). The relatively high effect size suggests a strong influence of the condition on SUS ratings.

An in-depth examination of the System Usability Scale (SUS) responses illuminated the specific items that contributed more significantly to the observed variations in usability scores between the conditions. In particular, responses to statement 1 (I would like to use this system regularly), statement 3 (I found the system simple to use), statement 8 (I found the system unnecessarily complex), and statement 9 (I felt confident while using the system) displayed the most noticeable differences. The AccessMenu condition received consistently positive feedback on these items, while the screen reader OCR condition received unfavorable feedback. Although the responses to other SUS items followed a similar pattern, the differences were relatively less pronounced.

The NASA Task Load Index (NASA-TLX) questionnaire [32] is typically used to assess participants’ perceived workload while doing the tasks. NASA-TLX scores also range from 0 to 100, but lower ratings indicate better performance, i.e., reduced taskload. We observed that there was a significant impact of the study conditions on the NASA-TLX scores (ANOVA test; $F = 161.26, p < 0.005$). Specifically, the TLX scores for the AccessMenu condition (Average (μ) = 48.03, Standard Deviation (σ) = 6.03) were significantly lower than those for the screen reader OCR condition (Average (μ) = 77.93, Standard Deviation (σ) = 3.67), suggesting a substantial reduction in perceived workload when using the proposed system (Figure 3b). A deeper inspection of the individual ratings revealed that responses to the Mental Demand, Effort, and Frustration subscales contributed relatively more to the difference in TLX scores between the conditions than those to the other subscales.

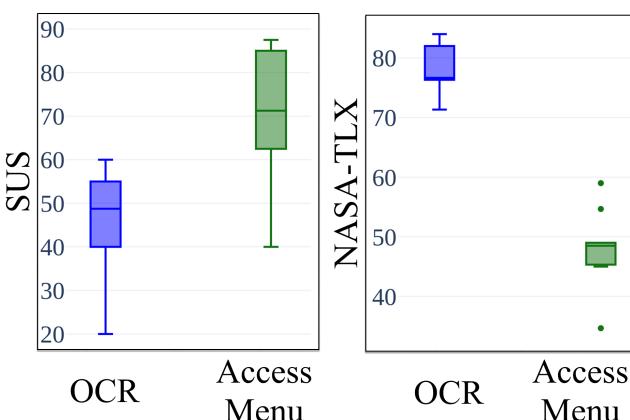


Figure 3: (a) System usability scale (SUS) and (b) NASA Task Load Index (NASA-TLX) for the two study conditions.

5.6.4 Qualitative Feedback. Qualitative analysis of the participants’ feedback revealed the following notable themes:

Exploring menu items with screen reader OCR was cumbersome. All participants reported experiencing fatigue and frustration when interacting with restaurant menus using the default screen reader OCR feature. The primary challenges contributing to

1045 this experience included *mentally linking scattered pieces of information, difficulty in locating specific content, and a need for memorization*. As P2 explained, “You have to listen back-and-forth and figure out on your own, what are all related to each other and what are not. The entire structure is complex without a proper order, so it is challenging to search for specific items I am interested in.”
 1046 Towards this, four other participants also expressed a preference
 1047 for a system feature that will enable them to maintain a “favorite
 1048 list”, which the system can then use to automatically filter the items
 1049 in the menu.
 1050

1051 **AccessMenu was perceived to be easy to learn and use.** A
 1052 majority of the participants (8) attributed their high usability ratings
 1053 for AccessMenu to its simplicity and very short learning curve.
 1054 They noted that memorizing a few shortcuts to access and navigate
 1055 the AccessMenu’s interface was a reasonable trade-off, given the
 1056 substantial advantage of reduced navigation effort with the content.
 1057 **Challenges in menu search and platform-wide filtering.** Nearly
 1058 half of the participants pointed out the lack of personalization in
 1059 menu searches. Specifically, they asked if they had to go through the
 1060 same search process of issuing the same queries when they accessed
 1061 different menus while comparing restaurants. As P9 stated, “I want
 1062 my previous search to carry over when I move to a different menu
 1063 instead of typing it again.” Additionally, a few participants emphasized
 1064 the need for filtering at a platform-wide level rather than
 1065 being limited to individual restaurant menus. They explained that
 1066 while filtering menu items within a restaurant is helpful, the ability
 1067 to search and filter across the menus of multiple restaurants, e.g.,
 1068 by relying on Google Maps platform, would be significantly more
 1069 beneficial. Towards this, P3 noted, “It would be helpful if I could just
 1070 search for a type of food or dietary preference across all available
 1071 restaurants, rather than going through each one separately.”
 1072

6 DISCUSSION

6.1 Limitations

1073 A notable limitation of our evaluation study was that the selection
 1074 of restaurant menus was confined to those with high extraction ac-
 1075 curacy. While this strategy aimed to reduce confounding variables,
 1076 it inadvertently restricted our capacity to assess the AccessMenu
 1077 ‘in-the-wild’, i.e., in the presence of extraction inaccuracies. Future
 1078 work should explore how blind participants respond and adapt to
 1079 potential extraction errors, providing insights into the system’s
 1080 effectiveness under less ideal conditions.
 1081

1082 Another limitation is that AccessMenu can presently support
 1083 only restaurant menus in English. In real-world scenarios, menu lan-
 1084 guages often vary based on geographical location, reflecting local
 1085 linguistic preferences. Extending our method to support multilin-
 1086 gual restaurant menus is a promising direction for future research,
 1087 allowing for broader applicability of our work.
 1088

1089 The third limitation relates to the inherent latency associated
 1090 with large language models such as GPT-4o-mini. While none of
 1091 the participants reported any noticeable latency issues when using
 1092 AccessMenu, this may not fully capture real-world scenarios where
 1093 delays could potentially impact user experience. In future work, we
 1094 aim to optimize the deployment process to mitigate any potential
 1095 latency concerns, ensuring AccessMenu operates seamlessly and
 1096 efficiently across various use cases.
 1097

1098 Lastly, AccessMenu is currently designed exclusively for desktop
 1099 environments. Given the widespread use of smartphones and the
 1100 growing trend of mobile-based activities, enabling efficient non-
 1101 visual web interactions with restaurant menus on smart mobile
 1102 devices is essential, which is also in the scope of our future research.
 1103

6.2 Platform-Wide Menu Filtering

1104 Our user study highlighted that while filtering items within a restau-
 1105 rant menu is helpful for blind users, they would immensely benefit
 1106 from the ability to filter items from multiple menus across different
 1107 restaurants, i.e., platform-wide level filtering by leveraging services
 1108 such as Google Maps, Uber Eats, Grubhub, and DoorDash. Recogniz-
 1109 ing the increased adoption of LLM agents [52] in online platforms,
 1110 we plan to develop a custom LLM agent that would enable blind
 1111 users to issue filter queries at a platform-level, e.g., in Google Maps,
 1112 and the agent would respond by providing an assimilated list of
 1113 items extracted from multiple menus. The user would be able to
 1114 therefore compare the items ‘in-one-place’ before deciding on the
 1115 restaurant for ordering food.
 1116

6.3 Personalized Query-Based Menu Filtering

1117 In our user study, we identified a need for personalization. Specifi-
 1118 cally, the participants wanted AccessMenu to carry over their prior
 1119 search queries when navigating different restaurant menus. In the
 1120 current system design, user queries are not stored, thereby requir-
 1121 ing users to reissue the same query for each new menu, hindering
 1122 efficient comparison. To address this, we plan to incorporate a per-
 1123 sonalization feature that will store user queries and automatically
 1124 apply them (to the best extent possible) on other menus accessed
 1125 in the same browsing session. In addition to queries, we also plan
 1126 to store and apply user preferences such as allergen-related filters.
 1127 For example, if a user requests to exclude items with specific aller-
 1128 gens, the system will remember this preference and curate other
 1129 subsequently accessed menus accordingly, providing a tailored,
 1130 user-centric experience.
 1131

7 CONCLUSION

1132 In this paper, we introduced AccessMenu, an intelligent browser
 1133 extension designed to enhance the usability of online restaurant
 1134 menus for blind and visually impaired (BVI) users who rely on
 1135 screen readers. The design of AccessMenu was based on the find-
 1136 ings of an interview study with 12 participants, which illuminated
 1137 the various pain points and needs of blind screen reader users
 1138 regarding online restaurant menus. AccessMenu provides an alter-
 1139 native screen reader-friendly interface to conveniently peruse
 1140 information in the menus. The AccessMenu interface also supports
 1141 natural language queries, enabling users to swiftly retrieve relevant
 1142 information without the need to manually scan the entire menu.
 1143 The findings from user evaluations showed that AccessMenu sig-
 1144 nificantly improved usability, surpassing the capabilities of status
 1145 quo OCR-based solutions.
 1146

ACKNOWLEDGMENTS

1147 We sincerely thank Amanda Kelly Dcosta, Aditya Vishal, Sri Harshitha
 1148 Gattu and Sushmitha Halli Sudhakara for their invaluable contribu-
 1149 tions to this work.
 1150

1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160

REFERENCES

- [1] [n. d.]. claude. ([n. d.]). <https://www.anthropic.com/news/clause-3-5-sonnet>
- [2] [n. d.]. gpt4. ([n. d.]). <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>
- [3] [n. d.]. llama. ([n. d.]). <https://www.llama.com/>
- [4] [n. d.]. w3org. ([n. d.]). <https://www.w3.org/TR/2006/WDaria-roadmap-20060926/>
- [5] Nov. 2008. WCAG. (Nov, 2008). <https://www.w3.org/WAI/WCAG22/quickref/?versions=2.1#text-alternatives/>
- [6] Nov. 2008. webaim. (Nov, 2008). <https://webaim.org/techniques/alttext/>
- [7] NV Access. 2018. NVDA screen-reader.
- [8] Kriti Aggarwal, Aditi Khandelwal, Kumar Tanmay, Owais Khan Mohammed, Qiang Liu, Monojit Choudhury, Hardik Chauhan, Subhojit Som, Vishrav Chaudhary, and Saurabh Tiwary. 2023. DUBLIN: Visual Document Understanding By Language-Image Network. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*. 693–706.
- [9] Patrizia Andronico, Marina Buzzi, Carlos Castillo, and Barbara Leporini. 2006. Improving search engine interfaces for blind users: a case study. *Universal Access in the Information Society* 5 (2006), 23–40.
- [10] Inc Apple. 2023. VoiceOver. https://www.apple.com/voiceover/info/guide/_1121.html
- [11] Vikas Ashok, Syed Masum Billah, Yevgen Borodin, and IV Ramakrishnan. 2019. Auto-suggesting browsing actions for personalized web screen reading. In *Proceedings of the 27th ACM Conference on User Modeling, Adaptation and Personalization*. 252–260.
- [12] Vikas Ashok, Yury Puzis, Yevgen Borodin, and IV Ramakrishnan. 2017. Web screen reading automation assistance using semantic abstraction. In *Proceedings of the 22nd International Conference on Intelligent User Interfaces*. 407–418.
- [13] Shirley Ann Becker. 2009. Web Accessibility and Compliance Issues. In *Encyclopedia of Information Science and Technology, Second Edition*. IGI Global, 4047–4052.
- [14] Syed Masum Billah, Vikas Ashok, Donald E Porter, and IV Ramakrishnan. 2017. Speed-dial: A surrogate mouse for non-visual web browsing. In *Proceedings of the 19th International ACM SIGACCESS Conference on Computers and Accessibility*. 110–119.
- [15] Syed Masum Billah, Vikas Ashok, Donald E Porter, and IV Ramakrishnan. 2018. SteeringWheel: a locality-preserving magnification interface for low vision web browsing. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–13.
- [16] Yevgen Borodin, Jeffrey P Bigham, Glenn Dausch, and IV Ramakrishnan. 2010. More than meets the eye: a survey of screen-reader browsing strategies. In *Proceedings of the 2010 International Cross Disciplinary Conference on Web Accessibility (W4A)*. 1–10.
- [17] James V Bradley. 1958. Complete counterbalancing of immediate sequential effects in a Latin square design. *J. Amer. Statist. Assoc.* 53, 282 (1958), 525–528.
- [18] John Brooke. 1996. SUS: A 'quick and dirty' usability scale. Usability Evaluation in Industry. PW Jordan, B Thomas, BA Weerdmeester and AL McClelland.
- [19] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [20] Giovanni Campagna, Silei Xu, Rakesh Ramesh, Michael Fischer, and Monica S Lam. 2018. Controlling fine-grain sharing in natural language with a virtual assistant. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 2, 3 (2018), 1–28.
- [21] Loredana Caruccio, Stefano Cirillo, Giuseppe Polese, Giandomenico Solimando, Shanmugam Sundaramurthy, and Genoveffa Tortora. 2024. Claude 2.0 large language model: Tackling a real-world classification problem with a new iterative prompt engineering approach. *Intelligent Systems with Applications* 21 (2024), 200336.
- [22] Yihao Ding, Siwen Luo, Hyunsuk Chung, and Soyeon Caren Han. 2023. VQA: A new dataset for real-world VQA on PDF documents. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 585–601.
- [23] DocuSign. 2025. Join More Than 1 Billion Users Who Trust DocuSign. <https://www.docusign.com> Accessed: [Your Access Date].
- [24] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Scheltens, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [25] Yansong Feng and Mirella Lapata. 2012. Automatic caption generation for news images. *IEEE transactions on pattern analysis and machine intelligence* 35, 4 (2012), 797–812.
- [26] Javedul Ferdous, Hae-Na Lee, Sampath Jayaratna, and Vikas Ashok. 2022. In-Support: Proxy Interface for Enabling Efficient Non-Visual Interaction with Web Data Records. In *27th International Conference on Intelligent User Interfaces*. 49–62.
- [27] Javedul Ferdous, Sami Uddin, and Vikas Ashok. 2021. Semantic table-of-contents for efficient web screen reading. In *Proceedings of the 36th Annual ACM Symposium on Applied Computing*. 1941–1949.
- [28] Fernando A Mikic Fonte, Martín Llamas Nistal, Martín Llamas Nistal, and Manuel Caeiro Rodriguez. 2016. NLAST: A natural language assistant for students. In *2016 IEEE global engineering education conference (EDUCON)*. IEEE, 709–713.
- [29] Boni Garcia, Mario Munoz-Organero, Carlos Alario-Hoyos, and Carlos Delgado Kloos. 2021. Automated driver management for selenium WebDriver. *Empirical Software Engineering* 26 (2021), 1–51.
- [30] Cole Gleason, Amy Pavel, Emma McCamey, Christina Low, Patrick Carrington, Kria M Kitani, and Jeffrey P Bigham. 2020. Twitter A11y: A browser extension to make Twitter images accessible. In *Proceedings of the 2020 chi conference on human factors in computing systems*. 1–12.
- [31] Darren Guinness, Edward Cutrell, and Meredith Ringel Morris. 2018. Caption crawler: Enabling reusable alternative text descriptions using reverse image search. In *Proceedings of the 2018 chi conference on human factors in computing systems*. 1–11.
- [32] Sandra G Hart and Lowell E Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of empirical and theoretical research. In *Advances in psychology*. Vol. 52. Elsevier, 139–183.
- [33] Sayed Kamrul Hasan and Terje Gjøsæter. 2021. Screen Reader Accessibility Study of Interactive Maps. In *International Conference on human-computer interaction*. Springer, 232–249.
- [34] Teakgyu Hong, Donghyun Kim, Mingi Ji, Wonseok Hwang, Daehyun Nam, and Sungrae Park. 2022. Bros: A pre-trained language model focusing on text and layout for better key information extraction from documents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 10767–10775.
- [35] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *Proceedings of the 30th ACM International Conference on Multimedia*. 4083–4091.
- [36] Mina Huh, Yi-Hao Peng, and Amy Pavel. 2023. GenAssist: Making image generation accessible. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–17.
- [37] Wonseok Hwang, Seongheyon Kim, Minjoon Seo, Jinyeong Yim, Seunghyun Park, Sungrae Park, Junyeop Lee, Bado Lee, and Hwalsuk Lee. 2019. Post-ocr parsing: building simple and robust parser via bio tagging. In *Workshop on Document Intelligence at NeurIPS 2019*.
- [38] Md Farhan Ishamm, Md Sakib Hossain Shovon, Muhammad Firoz Mridha, and Nilanjan Dey. 2024. From image to language: A critical analysis of visual question answering (vqa) approaches, challenges, and opportunities. *Information Fusion* 16 (2024), 102270.
- [39] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5648–5656.
- [40] Kushal Kafle, Robik Shrestha, Scott Cohen, Brian Price, and Christopher Kanan. 2020. Answering questions about data visualizations using efficient bimodal fusion. In *Proceedings of the IEEE/CVF Winter conference on applications of computer vision*. 1498–1507.
- [41] Muiz Ahmed Khan, Pias Paul, Mahmudur Rashid, Mainul Hossain, and Md Atiqur Rahman Ahad. 2020. An AI-based visual aid with integrated reading assistant for the completely blind. *IEEE Transactions on Human-Machine Systems* 50, 6 (2020), 507–517.
- [42] Geewook Kim, Teakgyu Hong, Moonbin Yim, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. 2021. Donut: Document understanding transformer without ocr. *arXiv preprint arXiv:2111.15664* 7, 15 (2021), 2.
- [43] Jonathan Lazar, Aaron Allen, Jason Kleinman, and Chris Malarkey. 2007. What frustrates screen reader users on the web: A study of 100 blind users. *International Journal of human-computer interaction* 22, 3 (2007), 247–269.
- [44] Jonathan Lazar, Abiodun Olalere, and Brian Wentz. 2012. Investigating the accessibility and usability of job application web sites for blind users. *Journal of Usability Studies* 7, 2 (2012).
- [45] Hae-Na Lee and Vikas Ashok. 2022. Customizable Tabular Access to Web Data Records for Convenient Low-vision Screen Magnifier Interaction. *ACM Transactions on Accessible Computing (TACCESS)* 15, 2 (2022), 1–22.
- [46] Hae-Na Lee, Vikas Ashok, and IV Ramakrishnan. 2020. Rotate-and-Press: A Non-visual Alternative to Point-and-Click? In *International Conference on Human-Computer Interaction*. Springer, 291–305.
- [47] Kenton Lee, Mandar Joshi, Iulia Radu Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenbach, Urveshi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *International Conference on Machine Learning*. PMLR, 18893–18912.
- [48] Maurizio Leotta, Fabrizio Mori, and Marina Ribaudo. 2023. Evaluating the effectiveness of automatic image captioning for web accessibility. *Universal access in the information society* 22, 4 (2023), 1293–1313.
- [49] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*. PMLR, 19730–19742.

1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275

- 1277 [50] Sheng Li, Zhiqiang Tao, Kang Li, and Yun Fu. 2019. Visual to text: Survey of image
1278 and video captioning. *IEEE Transactions on Emerging Topics in Computational
1279 Intelligence* 3, 4 (2019), 297–312.
- 1280 [51] Xin Li, Yunfei Wu, Xinghua Jiang, Zhihao Guo, Mingming Gong, Haoyu Cao,
1281 Yinsong Liu, Deqiang Jiang, and Xing Sun. 2024. Enhancing visual document
1282 understanding with contrastive learning in large visual-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15546–15555.
- 1283 [52] Yuanchun Li, Hao Wen, Weijun Wang, Xiangyu Li, Yizhen Yuan, Guohong Liu,
1284 Jiacheng Liu, Wenxing Xu, Xiang Wang, Yi Sun, et al. 2024. Personal llm agents:
1285 Insights and survey about the capability, efficiency and security. *arXiv preprint arXiv:2401.05459* (2024).
- 1286 [53] Christos Liambas and Miltiadis Saratzidis. 2016. Autonomous OCR dictating
1287 system for blind people. In *2016 IEEE Global Humanitarian Technology Conference
1288 (GHTC)*. IEEE, 172–179.
- 1289 [54] Kate Lister, Tim Coughlan, Francisco Iniesto, Nick Freear, and Peter Devine.
1290 2020. Accessible conversational user interfaces: considerations for design. In *Proceedings of the 17th international web for all conference*. 1–11.
- 1291 [55] Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee,
1292 Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Martin Eisenchlos. 2022.
1293 Matcha: Enhancing visual language pretraining with math reasoning and chart
1294 derendering. *arXiv preprint arXiv:2212.09662* (2022).
- 1295 [56] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual
1296 instruction tuning. *Advances in neural information processing systems* 36 (2024).
- 1297 [57] Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Cheng-lin Liu, Lian-
1298 wen Jin, and Xiang Bai. 2023. On the hidden mystery of ocr in large multimodal
1299 models. *arXiv preprint arXiv:2305.07895* (2023).
- 1300 [58] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset
1301 for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on
1302 applications of computer vision*. 2200–2209.
- 1303 [59] Valentyn Melnyk, Vikas Ashok, Valentyn Melnyk, Yury Puzis, Yevgen Borodin,
1304 Andrii Soviak, and IV Ramakrishnan. 2015. Look ma, no aria: generic accessible
1305 interfaces for web widgets. In *Proceedings of the 12th International Web for All
1306 Conference*. 1–4.
- 1307 [60] Yue Ming, Nannan Hu, Chunxiao Fan, Fan Feng, Jiangwan Zhou, and Hui Yu.
1308 2022. Visuals to text: A comprehensive review on automatic image captioning.
1309 *IEEE/CAA Journal of Automatica Sinica* 9, 8 (2022), 1339–1365.
- 1310 [61] Juan Nino, Sherezada Ochoa, Jocelyne Kiss, Geoffreyen Edwards, Ernesto
1311 Morales, James Hutson, Frédérique Poncet, and Walter Wittich. 2024. Assistive
1312 Technologies for Internet Navigation: A Review of Screen Reader Solutions for
1313 the Blind and Visually Impaired. *International Journal of Recent Engineering
1314 Science* 11, 6 (2024).
- 1315 [62] Vincent Perot, Kai Kang, Florian Luisier, Guolong Su, Xiaoyu Sun, Ramya Sree
1316 Boppana, Zilong Wang, Zifeng Wang, Jiaqi Mu, Hao Zhang, et al. 2023. Lmdx:
1317 Language model-based document information extraction and localization. *arXiv
1318 preprint arXiv:2309.10952* (2023).
- 1319 [63] Yash Prakash, Akshay Kolgar Nayak, Shoib Mohammed Alyaan, Pathan Aseef
1320 Khan, Hae-Na Lee, and Vikas Ashok. 2024. Improving Usability of Data Charts
1321 in Multimodal Documents for Low Vision Users. In *Proceedings of the 26th
1322 International Conference on Multimodal Interaction*. 498–507.
- 1323 [64] Yash Prakash, Akshay Kolgar Nayak, Mohan Sunkara, Sampath Jayarathna, Hae-
1324 Na Lee, and Vikas Ashok. 2024. All in One Place: Ensuring Usable Access to
1325 Online Shopping Items for Blind Users. *Proceedings of the ACM on Human-
1326 Computer Interaction* 8, EICS (2024), 1–25.
- 1327 [65] Yash Prakash, Mohan Sunkara, Hae-Na Lee, Sampath Jayarathna, and Vikas
1328 Ashok. 2023. AutoDesc: Facilitating Convenient Perusal of Web Data Items for
1329 Blind Users. In *Proceedings of the 28th International Conference on Intelligent User
1330 Interfaces*.
- 1331 [66] Yury Puzis. 2012. Accessible web automation interface: a user study. In *Proceedings
1332 of the 14th international ACM SIGACCESS conference on Computers and
1333 accessibility*. 291–292.
- 1334 [67] Yury Puzis, Yevgen Borodin, Faisal Ahmed, and IV Ramakrishnan. 2012. An intuitive
1335 accessible web automation user interface. In *Proceedings of the International
1336 Cross-Disciplinary Conference on Web Accessibility*. 1–4.
- 1337 [68] Johnny Saldaña. 2015. *The coding manual for qualitative researchers*. Sage.
- 1338 [69] Johnny Saldaña. 2021. The coding manual for qualitative researchers. (2021).
- 1339 [70] Anastasia Schaadhardt, Alexis Hiniker, and Jacob O Wobbrock. 2021. Understanding
1340 blind screen-reader users' experiences of digital artboards. In *Proceedings of the
1341 2021 CHI Conference on Human Factors in Computing Systems*. 1–19.
- 1342 [71] Freedom Scientific. 2020. JAWS OCR, What It Is and How It Works! <https://www.freedomscientific.com/webinars/jaws-ocr-what-it-is-and-how-it-works/> Accessed: Mar. 7, 2025.
- [72] Freedom Scientific. 2020. JAWS® – Freedom Scientific. <http://www.freedomscientific.com/products/software/jaws>.
- [73] Woosuk Seo and Hyunggu Jung. 2022. Challenges and opportunities to improve
the accessibility of YouTube for people with visual impairments as content
creators. *Universal Access in the Information Society* 21, 3 (2022), 767–770.
- [74] Weishi Shi, Heather Moses, Qi Yu, Samuel Malachowsky, and Daniel E Krutz.
2023. All: Supporting experiential accessibility education and inclusive software
development. *ACM Transactions on Software Engineering and Methodology* 33, 2
(2023), 1–30.
- [75] Shalini Sonth and Jagadish S Kallimani. 2017. OCR based facilitator for the
visually challenged. In *2017 International Conference on Electrical, Electronics,
Communication, Computer, and Optimization Techniques (ICEECCOT)*. IEEE, 1–7.
- [76] Yash Srivastava, Vaishnav Murali, Shiv Ram Dubey, and Snehasis Mukherjee.
2021. Visual question answering using deep learning: A survey and performance
analysis. In *Computer Vision and Image Processing: 5th International Conference,
CVIP 2020, Prayagraj, India, December 4–6, 2020, Revised Selected Papers, Part II 5*.
Springer, 75–86.
- [77] Mohan Sunkara, Yash Prakash, Hae-Na Lee, Sampath Jayarathna, and Vikas
Ashok. 2023. Enabling Customization of Discussion Forums for Blind Users.
Proceedings of the ACM on Human-Computer Interaction 7, ELCS (2023), 1–20.
- [78] Zineng Tang, Ziyi Yang, Guoxin Wang, Yuwei Fang, Yang Liu, Chenguang Zhu,
Michael Zeng, Cha Zhang, and Mohit Bansal. 2023. Unifying vision, text, and layout
for universal document processing. In *Proceedings of the IEEE/CVF Conference
on Computer Vision and Pattern Recognition*. 19254–19264.
- [79] Sachin Tanwar and PVM Rao. 2024. Inequality in User Experience: Can Mobile
User Interfaces that Help Sighted Users Create Barriers for Visually Challenged
People?. In *International Conference on Computers Helping People with Special
Needs*. Springer, 19–30.
- [80] Mary Frances Theofanos and Janice Redish. 2003. Bridging the gap: between
accessibility and usability. *interactions* 10, 6 (2003), 36–51.
- [81] Rubén Tito, Minesh Mathew, CV Jawahar, Ernest Valveny, and Dimosthenis
Karatzas. 2021. Icdar 2021 competition on document visual question answering.
In *Document Analysis and Recognition—ICDAR 2021: 16th International Conference,
Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part IV 16*. Springer,
635–649.
- [82] Tejal Tiwary and Rajendra Prasad Mahapatra. 2023. Enhancement in web ac-
cessibility for visually impaired people using hybrid deep belief network-bald
eagle search. *Multimedia Tools and Applications* (2023), 1–22.
- [83] Utku Uckun, Rohan Tumkur Suresh, Md Javedul Ferdous, Xiaojun Bi, IV Ra-
makrishnan, and Vikas Ashok. 2022. Taming User-Interface Heterogeneity with
Uniform Overlays for Blind Users. (2022).
- [84] Helmut Vieritz, Farzan Yazdi, Daniel Schilberg, Peter Göhner, and Sabina Jeschke.
2011. User-centered design of accessible web and automation systems. In *Information
Quality in e-Health: 7th Conference of the Workgroup Human-Computer
Interaction and Usability Engineering of the Austrian Computer Society, USAB 2011,
Graz, Austria, November 25–26, 2011. Proceedings 7*. Springer, 367–378.
- [85] Peng Wang, Qi Wu, Chunhua Shen, and Anton Van den Hengel. 2017. The
vqa-machine: Learning how to use existing vision algorithms to answer new
questions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern
Recognition*. 1173–1182.
- [86] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian
Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models
are zero-shot learners. *arXiv preprint arXiv:2109.01652* (2021).
- [87] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi,
Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning
in large language models. *Advances in neural information processing systems*
35 (2022), 24824–24837.
- [88] Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan
Liu, Dinei Florencio, Cha Zhang, Wanxiang Che, et al. 2020. Layoutlmv2: Multi-
modal pre-training for visually-rich document understanding. *arXiv preprint arXiv:2012.14740* (2020).
- [89] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng
Liu, and Lijuan Wang. 2023. The dawn of lmms: Preliminary explorations with
gpt-4v (vision). *arXiv preprint arXiv:2309.17421* 9, 1 (2023), 1.
- [90] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiaobo Ye, Ming Yan, Yiyang Zhou, Junyang
Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023. mpplug-owl: Modular-
ization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178* (2023).
- [91] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023.
Minigpt-4: Enhancing vision-language understanding with advanced large lan-
guage models. *arXiv preprint arXiv:2304.10592* (2023).
- [92] Wang Zhu, Alekh Agarwal, Mandar Joshi, Robin Jia, Jesse Thomason, and Kristina
Toutanova. 2023. Efficient End-to-End Visual Document Understanding with
Rationale Distillation. *arXiv preprint arXiv:2311.09612* (2023).

1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392