

CodeA11y: Making AI Coding Assistants Useful for Accessible Web Development

Peya Mowar
pmowar@cs.cmu.edu
Carnegie Mellon University
Pittsburgh, PA, USA

Yi-Hao Peng
yihaop@cs.cmu.edu
Carnegie Mellon University
Pittsburgh, PA, USA

Jason Wu
jason_wu8@apple.com
Apple
Seattle, WA, USA

Aaron Steinfeld
steinfeld@cmu.edu
Carnegie Mellon University
Pittsburgh, PA, USA

Jeffrey P. Bigham
jbigham@cs.cmu.edu
Carnegie Mellon University
Pittsburgh, PA, USA

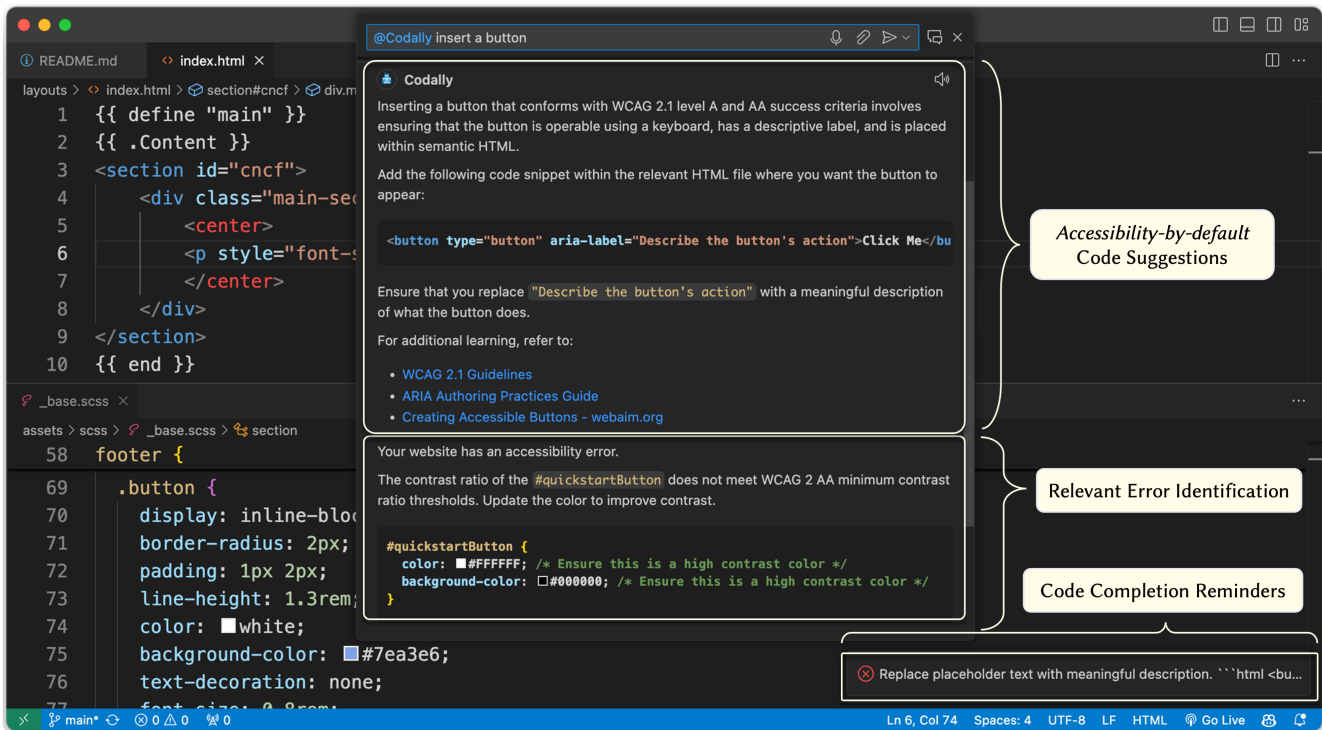


Figure 1: CodeA11y is a GitHub Copilot Extension for Accessible Web Development. CodeA11y addresses accessibility limitations of Copilot observed in our study with developers through three features: (1) accessibility-by-default code suggestions, (2) automatic identification of relevant accessibility errors, and (3) reminders to replace placeholders in generated code. Integrating these features directly into AI coding assistants would improve the accessibility of the user interfaces (UIs) developers create.

ABSTRACT

A persistent challenge in accessible computing is ensuring developers produce web UI code that supports assistive technologies. Despite numerous specialized accessibility tools, novice developers often remain unaware of them, leading to ~96% of web pages that contain accessibility violations. AI coding assistants, such as GitHub Copilot, could offer potential by generating accessibility-compliant code, but their impact remains uncertain [52]. Our formative study with 16 developers without accessibility training revealed three key issues in AI-assisted coding: failure to prompt AI for accessibility,

omitting crucial manual steps like replacing placeholder attributes, and the inability to verify compliance. To address these issues, we developed CodeA11y, a GitHub Copilot Extension, that suggests accessibility-compliant code and displays manual validation reminders. We evaluated it through a controlled study with another 20 novice developers. Our findings demonstrate its effectiveness in guiding novice developers by reinforcing accessibility practices throughout interactions, representing a significant step towards integrating accessibility into AI coding assistants.

CCS CONCEPTS

• **Human-centered computing** → **Accessibility design and evaluation methods; Interactive systems and tools**; • **Software and its engineering** → **Development frameworks and environments**.

KEYWORDS

AI Coding Assistants, Web Accessibility, Coding Agents, AI Agents

ACM Reference Format:

Peya Mowar, Yi-Hao Peng, Jason Wu, Aaron Steinfeld, and Jeffrey P. Bigham. 2025. CodeA11y: Making AI Coding Assistants Useful for Accessible Web Development. In *CHI Conference on Human Factors in Computing Systems (CHI '25)*, April 26–May 1, 2025, Yokohama, Japan. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3706598.3713335>

1 INTRODUCTION

Most websites contain extensive accessibility errors [87], despite decades of investment in standards and guidelines [13, 18], tools [8, 81, 82], advocacy [48, 57, 78], and policy. According to a recent analysis by WebAim [87], the homepages of the top million websites each contain 57 accessibility errors on average, including (but not limited to) missing alt-text for images [9, 28], inadequate color contrast [69], incorrect or missing labels for forms and links [36], and improper use of heading levels [12]. As a result, many people with disabilities will find it difficult to use these websites effectively and may not be able to use them at all.

Front-end web developers ultimately determine the accessibility (or inaccessibility) of the UI code that they write. Getting front-end developers to write more accessible code has proven exceptionally difficult. As Jonathan Lazar *et al.* wrote twenty years ago in 2004, “Since tools and guidelines are available to help designers and webmasters make their web sites accessible, it is unclear why so many sites remain inaccessible.” [40]. A survey of webmasters at the time indicated that they generally would like to make their web pages accessible but cited a number of reasons they do not: “lack of time, lack of training, lack of managerial support, lack of client support, inadequate software tools, and confusing accessibility guidelines.” Sixteen years later, Patel *et al.* reported remarkably similar results in their 2020 survey of 77 technology professionals [58]. Few developers had received formal accessibility training, implementing accessibility was considered confusing, and advocating for accessible development was in conflict with other business goals. Clearly, what we have done so far is not working.

We argue that AI coding assistants (*e.g.*, Github Copilot [1]) could offer an opportunity to make UI code more accessible. AI coding assistants are already widely adopted, which means that developers

do not need to be convinced to use them or to install a specialized tool for accessibility. They produce a wide variety of UI code and are capable enough to both reflect on the quality of arbitrary code and also prompt developers to fix what they are unable to do. This paper explores how AI coding assistants currently help developers create UI code, what problems remain, and presents a system called CodeA11y that shows that AI coding assistants can be made better at enabling developers to improve the accessibility of their UI code.

To explore this potential opportunity, we first conducted a user study (Section 3) with 16 developers not trained in accessibility to understand how current tools (GitHub Copilot) impact the production of accessible UI code. Our findings (Section 4) shows that while Copilot may potentially improve accessibility of UI code, three barriers prevent realization of those improvements: (1) developers may need to explicitly prompt the assistants for accessible code and thus not benefit if they fail to do so, (2) developers may overlook critical manual steps suggested by Copilot, such as replacing placeholders in alternative text for images, and (3) developers may not be able to verify if they fully implemented more complex accessibility enhancements properly. The formative study showed the potential of AI coding assistants to improve the accessibility of UI code, but revealed several gaps that led us to design goals (Section 5) for improving AI coding assistants to support accessibility.

We then built *CodeA11y* (Section 6, Figure 1), a GitHub Copilot Extension that addresses the observed gaps by consistently reinforcing accessible development practices throughout the conversational interaction. We evaluated CodeA11y (Section 7) with 20 developers, assessing its effectiveness in supporting accessible UI development and gathering insights for further refinement. We found that developers using CodeA11y are significantly more likely to produce accessible UI code. Finally, we reflect on the broader implications of integrating AI coding assistants into accessibility workflows, including the balance between automation and developer education, and the potential for AI tools to shape long-term developer behavior toward accessibility-conscious practices (Section 8).

The contributions of our paper are:

- We conducted a study with 16 developers that uncovered both benefits and limitations of current AI coding assistants for authoring accessible UI code.
- **CodeA11y**¹: a GitHub Copilot Extension that generates accessible UI code, identifies existing issues and reminds developers to perform manual validation.

2 RELATED WORK

Our research builds upon (i) *Assessing Web Accessibility*, (ii) *End-User Accessibility Repair*, and (iii) *Developer Tools for Accessibility*.

2.1 Assessing Web Accessibility

From the earliest attempts to set standards and guidelines, web accessibility has been shaped by a complex interplay of technical challenges, legal imperatives, and educational campaigns. Over the past 25 years, stakeholders have sought to improve digital inclusion by establishing foundational standards [13, 18], enforcing legal obligations [77, 90], and promoting a broader culture of accessibility awareness among developers [48, 57, 78]. Despite these

¹The source code for CodeA11y is available at <https://github.com/peyajm29/codea11y/>.

longstanding efforts, systemic accessibility issues persist. According to the 2024 WebAIM Million report [87], 95.9% of the top one million home pages contained detectable WCAG violations, averaging nearly 57 errors per page. These errors take many forms: low color contrast makes the interface difficult for individuals with color deficiency or low vision to read text; missing alternative text leaves users relying on screen readers without crucial visual context; and unlabeled form inputs or empty links and buttons hinder people who navigate with assistive technologies from completing basic tasks. Together, these accessibility issues not only limit user access to critical online resources such as healthcare, education, and employment but also result in significant legal risks and lost opportunities for businesses to engage diverse audiences. Addressing these pervasive issues requires systematic methods to identify, measure, and prioritize accessibility barriers, which is the first step toward achieving meaningful improvements.

Prior research has introduced methods blending automation and human evaluation to assess web accessibility. Hybrid approaches like SAMBA combine automated tools with expert reviews to measure the severity and impact of barriers, enhancing evaluation reliability [11]. Quantitative metrics, such as Failure Rate and Unified Web Evaluation Methodology, support large-scale monitoring and comparative analysis, enabling cost-effective insights [49, 85]. However, automated tools alone often detect less than half of WCAG violations and generate false positives, emphasizing the need for human interpretation [25, 86]. Recent progress with large pretrained models like Large Language Models (LLMs) [5, 24] and Large Multimodal Models (LMMs) [6, 44] offers a promising step forward, automating complex checks like non-text content evaluation and link purposes, achieving higher detection rates than traditional tools [20, 45]. Yet, these large models face challenges, including dependence on training data, limited contextual judgment, and the inability to simulate real user experiences. These limitations underscore the necessity of combining models with human oversight for reliable, user-centered evaluations [11, 20, 86].

Our work builds on these prior efforts and recent advancements by leveraging the capabilities of large pretrained models while addressing their limitations through a developer-centric approach. CodeA11y integrates LLM-powered accessibility assessments, tailored accessibility-aware system prompts, and a dedicated accessibility checker directly into GitHub Copilot—one of the most widely used coding assistants. Unlike standalone evaluation tools, CodeA11y actively supports developers throughout the coding process by reinforcing accessibility best practices, prompting critical manual validations, and embedding accessibility considerations into existing workflows.

2.2 End-user Accessibility Repair

In addition to detecting accessibility errors and measuring web accessibility, significant research has focused on fixing these problems. Since end-users are often the first to notice accessibility problems and have a strong incentive to address them, systems have been developed to help them report or fix these problems.

Collaborative, or social accessibility [76, 83], enabled these end-user contributions to be scaled through crowd-sourcing. AccessMonkey [10] and Accessibility Commons [37] were two examples

of repositories that store accessibility-related scripts and metadata, respectively. Other work has developed browser extensions that leverage crowd-sourced databases to automatically correct reading order, alt-text, color contrast, and interaction-related issues [31, 75].

One drawback of collaborative accessibility approaches is that they cannot fix problems for an “unseen” web page on-demand, so many projects aim to automatically detect and improve interfaces without the need for an external source of fixes. A large body of research has focused on making specific web media (e.g., images [26–28, 30, 41], design [43, 61, 65, 67], and videos [32, 59, 60, 63]) accessible through a combination of machine learning (ML) and user-provided fixes. Other work has focused on applying more general fixes across all websites.

Opportunity accessibility addressed a common accessibility problem of most websites: by default, content is often hard to see for people with visual impairments, and many users, especially older adults, do not know how to adjust or enable content zooming [7]. To this end, a browser script (`oppaccess.js`) was developed that automatically adjusted the browser’s content zoom to maximally enlarge content without introducing adverse side-effects (e.g., content overlap). While `oppaccess.js` primarily targeted zoom-related accessibility, recent work aimed to enable larger types of changes, by using LLMs to modify the source code of web pages based on user questions or directives [42].

Several efforts have been focused on improving access to desktop and mobile applications, which present additional challenges due to the unavailability of app source code (e.g., HTML). Prefab is an approach that allows graphical UIs to be modified at runtime by detecting existing UI widgets, then replacing them [22]. Interaction Proxies used these runtime modification strategies to “repair” Android apps by replacing inaccessible widgets with improved alternatives [92, 93]. The widget detection strategies used by these systems previously relied on a combination of heuristics and system metadata (e.g., the view hierarchy), which are incomplete or missing in the accessible apps. To this end, ML has been employed to better localize [15] and repair UI elements [14, 62, 89, 91].

In general, end-user solutions to repairing application accessibility are limited due to the lack of underlying code and knowledge of the semantics of the intended content.

2.3 Developer Tools for Accessibility

Ultimately, the best solution for ensuring an accessible experience lies with front-end developers. Many efforts have focused on building adequate tooling and support to help developers with ensuring that their UI code complies with accessibility standards.

Numerous automated accessibility testing tools have been created to help developers identify accessibility issues in their code: i) static analysis tools, such as IBM Equal Access Accessibility Checker [35] or Microsoft Accessibility Insights [51], scan the UI code’s compliance with predefined rules derived from accessibility guidelines; and ii) dynamic or runtime accessibility scanners, such as Chrome Devtools [29] or axe-Core Accessibility Engine [21], perform real-time testing on user interfaces to detect interaction issues not identifiable from the code structure. While these tools greatly reduce the manual effort required for accessibility testing, they are often criticized for their limited coverage. Thus, experts

often recommend manually testing with assistive technologies to uncover more complex interaction issues. Prior studies have created accessibility crawlers that either assist in developer testing [79, 80] or simulate how assistive technologies interact with UIs [72–74].

Similar to end-user accessibility repair, research has focused on generating fixes to remediate accessibility issues in the UI source code. Initial attempts developed heuristic-based algorithms for fixing specific issues, for instance, by replacing text or background color attributes [94]. More recent work has suggested that the code-understanding capabilities of LLMs allow them to suggest more targeted fixes. For example, a study demonstrated that prompting ChatGPT to fix identified WCAG compliance issues in source code could automatically resolve a significant number of them [55]. Researchers have sought to leverage this capability by employing a multi-agent LLM architecture to automatically identify and localize issues in source code and suggest potential code fixes [50].

While the approaches mentioned above focus on assessing UI accessibility of already-authored code (*i.e.*, fixing existing code), there is potential for more proactive approaches. For example, LLMs are often used by developers to generate UI source code from natural language descriptions or tab completions [1, 16, 33, 46, 71, 95], but LLMs frequently produce inaccessible code by default [4, 52], leading to inaccessible output when used by developers without sufficient awareness of accessibility knowledge. The primary focus of this paper is to design a more accessibility-aware coding assistant that both produces more accessible code without manual intervention (*e.g.*, specific user prompting) and gradually enables developers to implement and improve accessibility of automatically-generated code through IDE UI modifications (*e.g.*, reminder notifications).

3 FORMATIVE STUDY METHODS

We conducted a formative study to assess the implications of AI coding assistants on web accessibility. We recruited novice developers and tasked them with editing real-world websites using GitHub Copilot. Our goal was to better understand how the use of Copilot affected the accessibility of the user interface code they produced.

Tasks. The participants completed tasks in the codebases for two open-source websites, Kubernetes [68] and BBC News [54]. Both websites received over 2 million monthly visits worldwide [3] and belong to different categories in the IAB Content Taxonomy [2]. These websites were developed using different web development frameworks (Hugo and React, respectively). To choose the four specific tasks used in this formative study (Table 1), we sampled actual issues from each website’s repository. We chose issues for which accessibility needed to be considered to complete them correctly, but accessibility was not explicitly mentioned as a requirement in either the task description given to participants or on the issue description on the website’s code repository, as illustrated in Figure 2. Correctly performing the tasks required the consideration of several common web accessibility issues: color contrast, alternative text, link labels, and form labeling [87]. The goal was to mirror the kinds of specifications that developers often receive that do not explicitly mention accessibility.

Protocol. Our within-subjects user study had two conditions: (1) a control condition where participants received no AI assistance, and (2) a test condition where participants used GitHub Copilot. Each

participant was assigned to edit two distinct websites, each with two tasks. To counterbalance order effects, participants were evenly and randomly assigned to one of four user groups (Table 2), balanced by website order and control/test conditions. Further, to simulate real-world scenarios, we concealed the true purpose of the study from participants. Participants were informed that the study was about the usability of AI pair programmers in web development tasks but were not explicitly instructed to make their web components accessible. This allowed us to observe how developers naturally handle accessibility when it is not explicitly emphasized, reflecting typical developer behavior. The research protocol was reviewed and approved by the Institutional Review Board (IRB) at our university.

Participants. We employed convenience sampling and snowball sampling methods to recruit our participants. Our study was advertised on university bulletin boards, social media, and shared communication channels (Twitter, Slack, and mailing groups). Our recruitment criteria stipulated that participants must be over 18 years of age, live in the United States, and have self-assessed familiarity with web development. Further, we required the participants to be physically present on our university campus for the duration of the study. To avoid priming during participant recruitment, we did not stipulate awareness of web accessibility as an eligibility criterion. We chose university-specific avenues for recruiting CS students, that reflect a typical novice developer cohort.

Our study enlisted 16 participants (7 female and 9 male; ages ranged from 22 to 29). Almost all of our participants were students and had multiple years of coding experience. Most ($n=10$) had multi-year *industrial* programming experience (*e.g.*, full-time or intern experience in the company). Nearly all participants (except one) had previously used AI coding assistants. GitHub Copilot and OpenAI ChatGPT were the most popular ($n=10$). Others preferred Tabnine ($n=6$) and AWS CodeWhisperer ($n=2$). 12 participants had self-described substantial experience with HTML and CSS. 10 were proficient in JavaScript and 7 were proficient in React.js. Despite this expertise, the majority (14 participants) were unfamiliar with the Web Content Accessibility Guidelines (WCAG). Only 2 participants knew about these guidelines, but they had not actively engaged in creating accessible web user interfaces or received formal training on the subject (details are provided in Table 3).

Procedure. The study was conducted in person at our lab, where participants performed programming tasks on a MacBook Pro laptop equipped with IntelliJ IDEA with the GitHub Copilot plugin preinstalled. Before starting the study, we explained the study procedure to the participants and took their informed consent. The participants then watched a 5-minute instructional video explaining Copilot’s features, such as code autocompletion and the Copilot chat². Participants were assigned tasks related to two selected websites, with a total of four tasks to complete in 90 minutes. They were required to work on one website with and the other without GitHub Copilot. Further, they were allowed to access the web for task exploration or code documentation through traditional search engines like Google Search, but with generative results turned off. During the coding session, a researcher observed silently, offering help with tasks, tool usage, or debugging only if participants were stuck, and asked them to move on after 30 minutes, without

²<https://www.youtube.com/watch?v=jXp5D5ZnxGM>

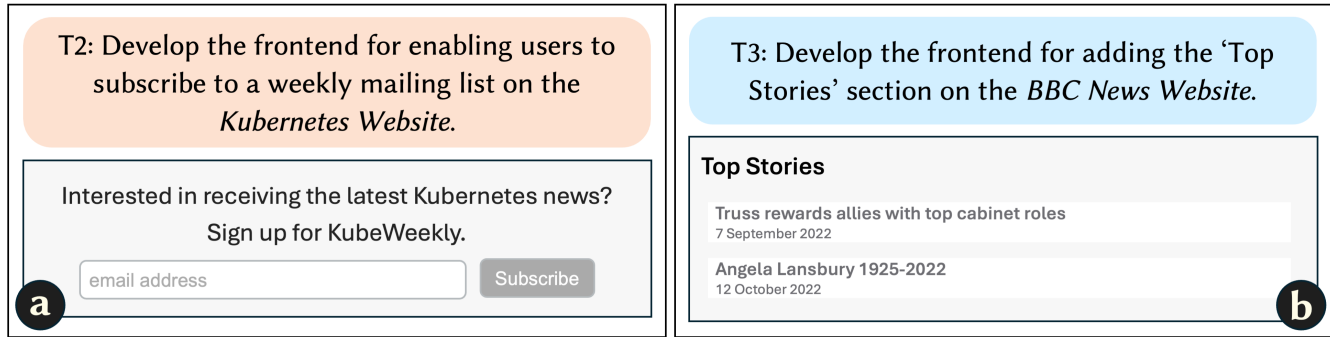


Figure 2: Examples of task descriptions and visual references given to our participants: (a) Task 2 was to implement a new contact form for subscribing to a mailing list, and (b) Task 3 was to add a ‘Top Stories’ section with linked articles. Successfully completing them required proper labeling of the form elements and links, but this was not explicitly stated in the instructions.

Table 1: Our formative study included four tasks. Each task was not primarily about accessibility but included an accessibility issue that was required to complete the task successfully. We adopt the scales of Unacceptable, Average and Good from prior work [66]. Uninformative attributes are those that merely reflect the field, such as ‘alt’ as alt-text or ‘click here’ as link description, without providing more meaningful or descriptive content [70]. Tasks are ranked from easy to difficult based on the time taken and success rates observed in our pilot studies.

Task	Difficulty	Accessibility Issue	Evaluation Criteria
(T1) Button Visibility	Easy	Color Contrast	<i>Unacceptable:</i> contrast ratio of < 4.5:1 for normal text and < 3:1 for large text <i>Average:</i> WCAG level AA in default state (contrast ratio of \geq 4.5:1 for normal text) <i>Good:</i> WCAG level AA in all states (default, hover, active, focus, etc.)
(T2) Form Element	Moderate	Form Labeling	<i>Unacceptable:</i> Missing form labels and keyboard navigation <i>Average:</i> One of form labels and keyboard navigation <i>Good:</i> Both form labeling and keyboard navigation
(T3) Add Section	Moderate	Link Labeling	<i>Unacceptable:</i> Missing link descriptions <i>Average:</i> Uninformative link descriptions <i>Good:</i> Descriptive link descriptions
(T4) Enhance Image for SEO	Difficult	Adding alt-text	<i>Unacceptable:</i> Missing or uninformative alt-text <i>Average:</i> Added alt-text with < 3 required descriptors [47] <i>Good:</i> Added alt-text with \geq 3 out of 4 required descriptors

Table 2: Participant User Groups: Each group is assigned specific order of tasks and testing conditions. Participants are evenly and randomly distributed among these groups.

#	Order 1, Testing Condition	Order 2, Testing Condition
1	Kubernetes, With AI Assistance	BBC News, No AI Assistance
2	Kubernetes, No AI Assistance	BBC News, With AI Assistance
3	BBC News, With AI Assistance	Kubernetes, No AI Assistance
4	BBC News, No AI Assistance	Kubernetes, With AI Assistance













giving any accessibility-related hints. Based on our observations from pilot studies, we set time limits ranging from 15 to 30 minutes per task. Finally, after completing the coding tasks, they were asked to participate in a 10-15 minute survey on their experience in AI-assisted programming and web accessibility, development

expertise, and open-ended feedback. In the end, the participants were compensated with a gift voucher worth 30 USD.

Data Collection and Analysis. We collected both quantitative and qualitative data for a mixed-method analysis. For quantitative data, we used an IntelliJ IDEA plugin [23] that tracked user actions – such as keyboard input (typing, backspace), IDE commands (copy, paste, undo), and interactions with GitHub Copilot (accepting suggestions, opening the Copilot Chat window) – and recorded their timestamps. Additionally, we employed the axe-Core Accessibility Engine 2 to gather accessibility violation metrics, including the type and count of WCAG failures, for each code submission, a method proven reliable in previous studies [56]. We also collected AI usage, programming languages and framework preferences, and expertise in web accessibility via a post-task survey.

On the qualitative side, we captured the entire study sessions through screen recordings, resulting in a total of 18.73 hours of video data. We complemented this with observational notes taken

Table 3: The distribution of participants’ opinions on AI-powered programming tools and their awareness of web accessibility. The percentages in the distribution column indicate the proportion of participants who either disagree (including both ‘strongly disagree’ and ‘disagree’) or agree (including both ‘strongly agree’ and ‘agree’) with the provided statements.

Statement	Distribution				
“I trust the accuracy of AI programming tools.”	13%				25%
“I am proficient in web accessibility.”	75%				19%
“I am familiar with the web accessibility standards, such as WCAG 2.0.”	88%				12%
“I am familiar with ARIA roles, states, and properties.”	69%				25%

■ Strongly Disagree
 ■ Disagree
 ■ Neutral
 ■ Agree
 ■ Strongly Agree

during the sessions, documenting verbal comments made by participants. The participants’ interactions with Copilot Chat were also recorded for further analysis between prompts and the final code. The analysis of this data was carried out using open coding and thematic analysis [19]. Some themes that emerged were: ‘visual enhancement’, ‘recalling syntax’, ‘feature request’, and ‘code understanding’. For accessibility evaluation, we manually inspected the websites created during the study and evaluated their accessibility on a qualitative scale of “Unacceptable,” “Average,” and “Good” adopted from prior work [66]. The criteria for these evaluations were developed per best practices identified in prior research published in CHI and ASSETS, detailed further in Table 1.

4 FORMATIVE FINDINGS

Our formative study revealed that while existing AI coding assistants can produce accessible code, developers still need accessibility expertise for effective use. Otherwise, (1) the accessibility introduced is likely to not be applied comprehensively, (2) the advanced features recommended by the assistant are unlikely to be implemented, (3) the accessibility errors introduced by the assistant are unlikely to be caught.

Developer Behavior. In the study, participants spent slightly more time on tasks without Copilot, averaging 30.84 minutes ($\sigma = 11.95$) compared to 28.94 minutes ($\sigma = 8.57$) with Copilot. Copilot also facilitated a greater volume of code edits (13.28 lines of code, $\sigma = 9.02$ vs 10.41 lines of code, $\sigma = 5.87$), indicating that AI-assisted workflows encouraged iterative coding practices. However, even with Copilot, participants spent approximately 39.84% of their task time (11.91 minutes, $\sigma = 8.00$) away from the IDE, browsing the web or checking the rendered HTML, highlighting the importance of traditional validation methods. The study also found fewer backspace key presses, an indicator of post-paste corrections, without Copilot ($\mu = 92.62$, $\sigma = 68.27$) than with the AI assistant ($\mu = 104.50$, $\sigma = 91.91$). Further, code pasting was slightly higher when participants solely browsed the web for exploration, averaging 12.68 times ($\sigma = 8.09$), compared to 11.43 times ($\sigma = 5.35$) with access to Copilot Chat. Participants dedicated about 7.39% (2.14 minutes, $\sigma = 1.72$) of their task time typing in the GitHub Copilot chat window, while, they also accepted Copilot’s code auto-complete suggestions around 5.44 times ($\sigma = 5.00$) on average.

AI Usage and Prompting Strategies. Participants mainly used the autocomplete feature only when they had a clear mental model of the desired code structure and sought to accelerate the code typing

process. In contrast, they heavily relied on the conversational interface for syntax inquiries, conceptual understanding, and the generation of code templates. We noticed that our participants wrote brief, task-oriented prompts that focused on immediate code solutions or specific interface modifications, often disregarding broader architectural considerations. Their prompting style was iterative and reactive, frequently requesting small incremental changes, fixes to previous outputs, or refinements to their vague prompts.

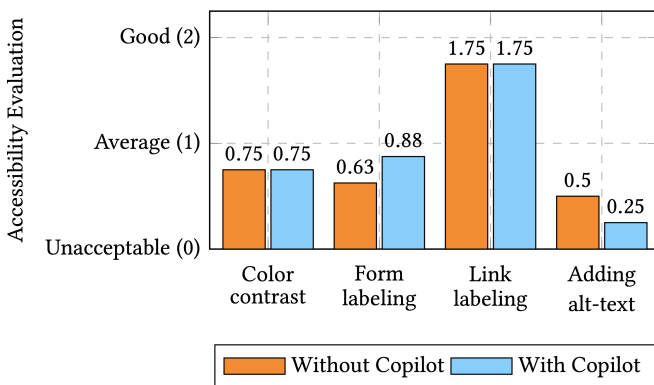
Furthermore, none of the participants, including the two who were familiar with web accessibility, prompted with accessibility in mind. Instead, our participants’ prompts centered around visual and functional attributes (e.g., ‘add a gray background to the subscription form’ (P4) or ‘add a grey patch’ (P1)). Consequently, the AI assistant’s suggestions often failed to incorporate accessibility best practices automatically. Occasionally, our participants prompted for enhancements that indirectly aligned with accessibility requirements, and Copilot provided relevant accessibility suggestions, as shown in Table 4. However, participants’ overreliance on AI assistance often led them to assume that Copilot’s code output was correct and complete. For instance, despite additional explanations from Copilot advising manual adjustments to image descriptions, participants directly pasted the code, resulting in code submissions with empty `alt` attributes.

Implications for Web Accessibility. Our study showed mixed results of AI assistants in considering accessibility issues with no statistically significant difference between the experimental conditions, as shown in Figure 3. Notably, Copilot could (sporadically) generate accessible components by utilizing patterns from other parts of a website. For example, it might automatically include proper labels for form fields, such as `<label for="email"> Email: </label>` in a signup form. However, there were also instances where Copilot inadvertently introduced new accessibility issues. For example, when adding new button components with hover effects, it failed to ensure adequate contrast between the hover color and background.

Further, the effectiveness of AI assistants was limited by the need for more sophisticated accessibility knowledge. Since our participants had limited awareness about these accessibility features, they would often ignore such suggestions by blindly accepting `alt = "" // Add your text here` or manually deleting the `<label>` tag. Some errors, such as providing blank alt-texts for informative images, were not even flagged by automated accessibility checkers because they interpret the image as decorative and consider this deliberate. This is particularly problematic as it

Table 4: The (partial) chat history revealed that directly copying Copilot’s code suggestions would be incomplete, as developers overlooked the additional recommended steps.

User:	“User Query”
Copilot:	Generated Code Associated Text Explanation
User 1:	“I want to label the images or add relevant image descriptors”
Copilot:	<code>alt=imgAlt // Add this line</code> Remember to replace imageAlt, <...> with your actual image attributes. In this example, imageAlt is the alt-text for the image, and <...>. You will need to replace these with the actual values.
User 2:	“Can you optimise the images on this article page for search engine optimization. Refactor the code accordingly by labeling the images or adding relevant image descriptions.”
Copilot:	<code>imgAlt=imageDescription // Add this line</code> Sure, to optimize images for SEO, you can add alt attributes to your img tags. The alt attribute provides alternative information for an image if a user for some reason cannot view it (because of slow connection, an error in the src attribute, or if the user uses a screen reader).

**Figure 3: Mean Accessibility Evaluation Scores by Tasks and Copilot Usage: Higher scores indicate success.**

implies that AI assistance might increase the risk of accessibility oversights, allowing critical errors to go unnoticed and uncorrected.

5 DESIGN REQUIREMENTS

Our formative study identified three limitations in novice developers’ interactions with AI assistants: (1) failing to prompt for accessibility considerations explicitly, (2) uncritically accepting incomplete code suggestions from Copilot, and (3) struggling to detect potential accessibility issues in their code. These shortcomings indicate possible directions to support accessibility through three design goals (G1-G3):

G1: Integrate System Prompts for Accessibility Awareness. Without explicit prompting, the AI assistant rarely produced accessibility-compliant code, reflecting the accessibility issues prevalent in its training data. However, it occasionally suggested accessibility features when participants indirectly prompted them, demonstrating its ability to recall accessibility practices from training data upon instruction. AI assistants should have a system prompt tuned towards following accessibility guidelines by default, for consistent generation of accessibility-compliant code, even when developers do

not mention accessibility specifically. Further, the system prompt should also direct the assistant to suggest accessibility-focused iterative refinements.

G2: Support Identification of Accessibility Issues. Due to their unfamiliarity with accessibility standards, our participants were unable to identify compliance issues in the existing and modified code. They primarily prompted changes to individual components (such as buttons and forms), hardly addressing broader page-level accessibility concerns (such as heading structure or landmark regions). AI assistants should not only automatically generate accessibility-compliant code, but also provide real-time feedback to detect and resolve accessibility violations within the codebase. In addition, AI assistants and automated accessibility checkers should work in tandem to ensure that incomplete or incorrect implementations of the AI-suggested code are always flagged by the latter.

G3: Encourage Developers to Complete AI-Generated Code. Our observations revealed that accessibility implementation in AI-assisted coding workflows commonly required critical manual intervention to complete and validate AI-generated code. This involved replacing placeholder attributes, such as labels and alt-texts, with meaningful values and verifying color contrast ratios. However, we found that participants blindly copy-pasted code and proceeded further if there were no apparent errors. This behavior of deferring thought to suggestions has also been documented in previous work [53]. To mitigate this, AI assistants should proactively remind developers to ensure that all necessary accessibility features – such as contrast ratios or keyboard navigation support – are fully implemented and verified.

6 CODEA11Y

Guided by the design goals identified through our user study, we built CodeA11y, a GitHub Copilot Extension for Visual Studio IDE. In this section, we present the interactions that it supports and its system architecture.

CodeA11y has three primary features (F1-F3, aligned to G1-G3, respectively): (F1) it produces user interface code that better complies with accessibility standards, (F2) it prompts the developer to resolve existing accessibility errors in their website, and (F3)

it reminds the developer to complete any AI-generated code that requires manual intervention. CodeA11y is integrated into Visual Studio Code as a GitHub Copilot Extension³, enabling CodeA11y to act as a chat participant within the GitHub Copilot Chat window panes. While we implemented this as an extension, it could be integrated directly into Copilot in the future.

Multi-Agent Architecture. CodeA11y has three LLM agents (Figure 4): Responder Agent, Correction Agent, and Reminder Agent. We provide their prompt instruction highlights in Table 5. These agents facilitate each of the above features (F1-F3) as follows:

- **Responder Agent** (for F1): This agent generates relevant code suggestions based on the developer’s prompt. It assumes that the developer is unfamiliar with accessibility standards and automatically generates accessible code. The prompt instruction for this agent is adapted from GitHub’s recommended user prompt for accessibility.⁴
- **Correction Agent** (for F2): This agent parses through accessibility error logs produced by an automated accessibility checker (axe DevTools Accessibility Linter⁵) to hint the developer at making additional accessibility fixes in the component or page being currently discussed in the chat context.
- **Reminder Agent** (for F3): This agent reviews the Responder Agent’s suggestions and identifies required manual steps for completing their implementation. It accordingly sends reminder notifications to the developer through the Visual Studio IDE infrastructure.

The agents are provided with several different sources of context:

- **Code Context:** the 100 lines of code centered around the cursor position in the active files.
- **Chat Context:** the current active chat window interaction.
- **Accessibility Linter Logs:** automated Axe DevTools Accessibility Linter error logs, refreshed periodically.
- **Project Context:** code context from the README and index files, which contain information about the web framework that is being used, information about project structure, and other key configuration details.

Due to the constraints in the context window, we optimized our prompts and filtered this context when it exceeded 4000 characters. The agents use GPT-4o as the back-end model from the same OpenAI GPT family of models that powers GitHub Copilot.

User Interaction. Developers invoke⁶ CodeA11y in the GitHub Copilot Chat window panes (includes Quick Chat and Chat View) using `@CodeA11y`. When a developer prompts CodeA11y, an internal `chat_context` state is established, storing the latest user prompts and agent responses. The `get_relevant_context()` function is called, which passes the source code and project context to the Responder Agent. The agent generates code suggestions, accessibility explanations, and links to additional resources and

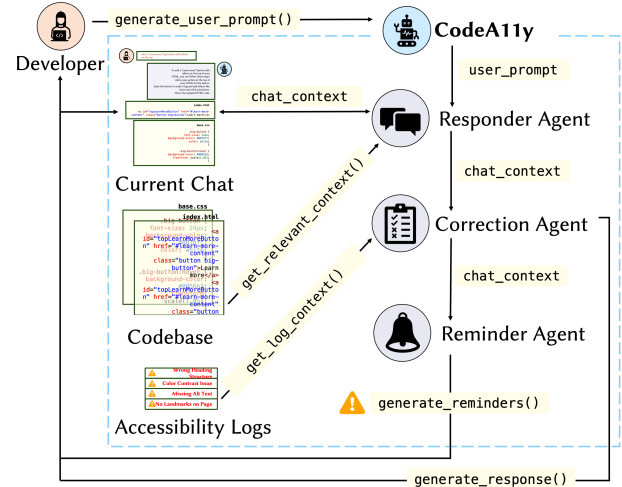


Figure 4: CodeA11y Architecture: Multi-agent workflow

updates `chat_context`. The `get_log_context()` function is called, which passes the accessibility linter logs to the Correction Agent. This agent suggests additional fixes and displays the responses in the chat pane. Lastly, the updated `chat_context` state is forwarded to the Reminder Agent, which generates and sends reminder notifications. Figure 5 illustrates a typical interaction between a developer and CodeA11y, showing how it compares to baseline assistants like GitHub Copilot.

7 USER EVALUATION

We conducted a within-subjects user study with 20 new participants to evaluate CodeA11y’s effectiveness in guiding novice developers toward adhering to accessibility standards, as compared to Copilot.

7.1 Methodology

We made the following revisions to our formative study protocol (Section 3). First, the experimental conditions were updated as follows: (1) the control condition involved using the baseline AI assistant (GitHub Copilot), and (2) the test condition where the participants used CodeA11y. Second, we changed the post-task survey to a brief semi-structured interview to get more nuanced insights about the usability of our system. We analyzed interview responses to better understand the factors shaping participants’ assistant preferences and their perceptions of any new coding practices introduced during the study. Third, we used Visual Studio Code as the IDE interface (which had advanced AI updates since the formative study – regarding both model performance and introduction of new features such as Inline Chat). Finally, we recruited 20 new participants for this subsequent study, with no prior exposure to the formative study, to evaluate CodeA11y’s performance on the same UI development tasks. These participants were the same demographic as our formative participants (students; multiyear programming experience; 6 female and 14 male; ages ranged from 22 to 30). Again, most participants were unfamiliar with the web accessibility standards (Table 6), but most (90%) had experience using AI programming tools. The IRB approved all our modifications.

³<https://docs.github.com/en/copilot/using-github-copilot/using-extensions-to-integrate-external-tools-with-copilot-chat>

⁴<https://github.blog/developer-skills/github/prompting-github-copilot-chat-to-become-your-personal-ai-assistant-for-accessibility/>

⁵<https://www.deque.com/axe/devtools/linter/>

⁶In the long term, the goal is for GitHub Copilot to invoke CodeA11y automatically during frontend development tasks.

Table 5: Prompt instructions for the three LLM agents in CodeA11y

Agent	Prompt Instruction Highlights
Responder Agent	<ul style="list-style-type: none"> • I am unfamiliar with accessibility and need to write code that conforms with WCAG 2.1 level AA criteria. • Be an accessibility coach that makes me account for all accessibility requirements. • Use reputable sources such as w3.org, webaim.org and provide links and references for additional learning. • Don't give placeholder variables but tell me where to give meaningful values. • Prioritise my current request and don't mention accessibility if I give a generic request like "Hi".
Correction Agent	<ul style="list-style-type: none"> • Review the accessibility checker log and provide feedback to fix errors relevant to current chat context. • If a log error relevant to current chat context occurs, provide a code snippet to fix it.
Reminder Agent	<ul style="list-style-type: none"> • Is there an additional step required by the developer to meet accessibility standards after pasting code? • Reminder should be single line. Be conservative in your response, if not needed, say "No reminders needed." • For example, remind the developer to replace the placeholder attributes with meaningful values or labels, or visually inspect element for colour contrast when needed.

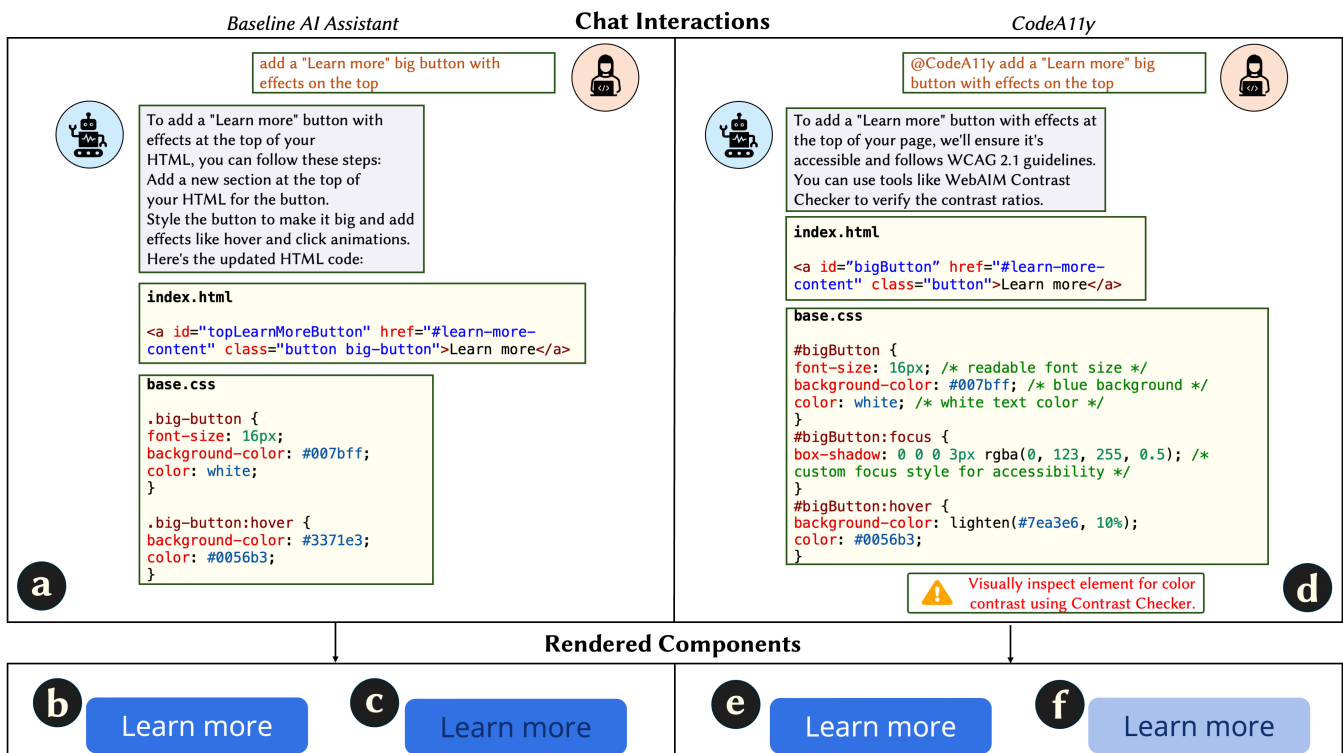






Figure 5: Contrasting responses for the same task across AI-assistants, showing differences in workflows. Developers had access to both the code and the rendered user interface. (a) and (d) represent conversations with the baseline assistant, and CodeA11y respectively. (b) and (e) show the buttons generated by each assistant in their default state. (c) and (f) display the buttons when hovered over, illustrating the differences in button color contrast.

Table 6: The distribution of opinions on AI-powered programming tools and their awareness of web accessibility based on the responses from participants in the evaluation study. The percentages in the distribution column indicate the proportion of participants who either disagree (including ‘strongly disagree’, ‘disagree’ and ‘slightly disagree’) or agree (including ‘strongly agree’, ‘agree’ and ‘slightly agree’) with the provided statements.

Statement	Distribution
“I trust the accuracy of AI programming tools.”	15%  70%
“I am proficient in web accessibility.”	70%  25%
“I am familiar with the web accessibility standards, such as WCAG 2.0.”	80%  15%
“I am familiar with ARIA roles, states, and properties.”	85%  10%




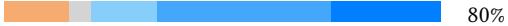
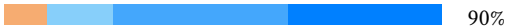



Table 7: The distribution of participants’ opinions on GitHub Copilot and CodeA11y, as well as their ease of completing tasks with these tools. The distribution column shows the count of responses from Strongly Disagree (1) to Strongly Agree (7).

Statement	Distribution
“I am satisfied with the code suggestions provided by”:	
GitHub Copilot	20%  75%
CodeA11y	15%  75%
“I found it easy to complete the coding tasks with”:	
GitHub Copilot	15%  80%
CodeA11y	10%  90%



To avoid biasing participants towards adhering to accessibility guidelines, we did not disclose the specific purpose of the CodeA11y plugin. For the duration of the study, we renamed the assistant “Codally” and described it as a general-purpose chat assistant for website editing. We assumed the interface would be intuitive, similar to widely used assistants, and therefore briefed participants only on basic AI assistant usage (e.g., Copilot), deliberately withholding explanations of error pop-ups to prevent influencing their behavior before the main study tasks. However, during the course of our study, we realized that VS Code was dismissing popup boxes created by our plugin more rapidly than expected - causing some participants to miss them. After 8 participants, we switched from floating popups to modals (which prevent the IDE’s auto-dismissal) due to a technical limitation. Both notification strategies do not require users to address errors, making them valid design choices. In our baseline comparison, we aggregate data from all users and include anecdotal observations of user behavior with each strategy. We acknowledge that such UI design choices may introduce variability and plan to investigate this further in future work.

7.2 Results

Here, we present the results of our subsequent evaluation study.

Accessibility Improvements. We implemented the accessibility assessments using the same measures outlined in our formative study (Table 1). Notably, our participants demonstrated a marked improvement in generating accessible web components and resolving accessibility issues with CodeA11y (Figure 6). CodeA11y

facilitated the automatic addition of form labels and ensured contrasting colors for button states, leading to statistically significant enhancements in accessibility outcomes. Specifically, participants performed better at adding form labels ($\mu = 1.5$, $\sigma = 0.85$) compared to GitHub Copilot ($\mu = 0.5$, $\sigma = 0.85$; $t = 2.63$, $p < 0.05$) and in ensuring contrasting button colors ($\mu = 1.3$, $\sigma = 0.67$ vs. $\mu = 0.7$, $\sigma = 0.82$; $t = 1.78$, $p < 0.05$). We also observed improvements in adding alt-texts with CodeA11y ($\mu = 0.7$, $\sigma = 0.95$ vs. $\mu = 0.1$, $\sigma = 0.32$; $t = 1.9$, $p < 0.05$). Though we did not find any statistical improvements in labeling links (perhaps because GitHub Copilot did a decent job at this task itself), all participants who used CodeA11y successfully completed this task ($\mu = 2$, $\sigma = 0$ vs. $\mu = 1.7$, $\sigma = 0.67$; $t = 1.9$, $p = 0.09$).

Developers’ Perspectives. Overall, participants reported no statistically significant difference in satisfaction ($\mu = 5.15$, $\sigma = 1.75$ vs. $\mu = 4.95$, $\sigma = 1.67$; $t = 0.37$, $p = 0.36$) and ease of use ($\mu = 5.8$, $\sigma = 1.47$ vs. $\mu = 5.4$, $\sigma = 1.67$; $t = 0.8$, $p = 0.21$) between CodeA11y and Github Copilot, respectively, as illustrated in Table 7.

During the post-study interviews, participants provided additional reasoning for their preferences. Most ($n=16$) participants did not have a specific preference between the two assistants, which is consistent with the conclusion of our statistical analysis. Others did indicate a preference ($n=3$) but provided reasoning that was based on the complexity of the task rather than assistant features, “I liked the first assistant (CodeA11y) better, maybe because of the tasks. The second one (GitHub Copilot) required me to understand the code, and the first directly gave me the code. That’s the difference.” (P18)

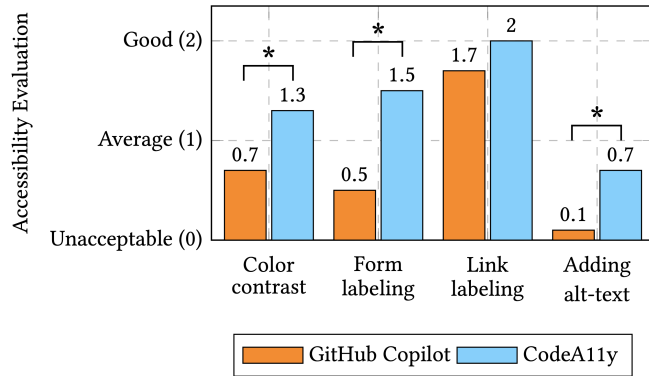


Figure 6: Mean Accessibility Evaluation Scores by Tasks and AI Assistant: Higher scores indicate success.

We asked our participants if they were introduced to any new coding practices by either of the assistants. To our surprise, only 4 participants mentioned accessibility, demonstrating CodeA11y’s effectiveness in “silently” improving the accessibility of our participants’ UI code. These participants noted that they had not these considerations before. However, some mentioned either not paying attention to them or subconsciously rejecting them, as they were primarily focused on completing the tasks, which they perceived to be unrelated to accessibility:

“I did not find any difference (between the assistants). When I was prompting CodeA11y, it was hinting at me to use alt texts, which was not happening in Copilot. It didn’t come to me by default, so that was good ... But I don’t think I implemented that.” (P27)

Still, they appreciated CodeA11y for emphasizing best practices for accessibility. For instance, P27 continued: *“I did see a few of the popups, and they did mention some interesting points like you need to consider the color of the button when you add a new button because if people are color blind, they might not be able to notice it.”* Further, P19, familiar with web accessibility but not proficient, realized that although he did not learn anything new about CodeA11y’s color contrast suggestion, he noticed a visible difference in user experience after accepting it. Our sole participant who claimed proficiency in accessibility valued support for a specific framework:

“I am familiar with accessibility coding practices but not in the React Native environment; I don’t know if I would have needed that help in HTML, but I liked that it tried to highlight accessibility practices in React Native.” (P21)

Other Observations. Participants frequently made minor edits to the AI-generated code for refining the visual appearance of web components. For instance, P18 remarked, *“CodeA11y did the job for me; I only had to change the property values.”* However, while focusing on visual adjustments, participants occasionally removed accessibility enhancements suggested by CodeA11y. Further, many participants still overlooked the manual validation steps required for implementing more advanced accessibility features. Participants appeared to consistently lack interest in the floating popups, so 50%

of participants using CodeA11y still added uninformative alt-texts. We observed a slightly different pattern with the modal reminders. Our participants initially paid attention to the modal interface, but then began to just close them. Although performance across the two reminder types is hard to assess objectively due to the small sample size, we observed higher means for the alt-text and form label tasks for the modal reminders. As a whole, the reminders proved somewhat effective: none of the participants using CodeA11y submitted empty alt-texts, which meant at least automated accessibility checkers would not consider the images decorative.

8 DISCUSSION

This paper has explored how AI coding assistants currently contribute to UI code that is accessible to people with disabilities. While these tools offer a new opportunity for achieving accessibility, we have revealed the remaining challenges and showed how they could be addressed with changes to the way the coding assistants operate.

Which comes first: Adoption or Awareness? Adoption is a perennial challenge faced by most accessibility technologies, even when the technology could lead to substantial improvements in user experience. One reason adoption is low is because awareness is low – people who could benefit from access technology do not know about it. For example, a survey found that only 10% of older adults knew what the term “accessibility” meant and therefore did not enable any useful settings [64, 88]. Similarly, developers benefit in many ways from tools that improve the accessibility of their code (e.g., linters, scanners), resulting in better-designed applications and reaching more users. Unfortunately, many developers are unaware of these tools or unwilling to adopt new practices that require changing their original workflow. For example, while AI assistants like Copilot are capable of generating accessible code, our formative study found that developers were unaware or unwilling to explicitly prompt it to do so.

One goal of our work was to investigate whether developers could increase the adoption of accessibility technology and development practices independently or in tandem with awareness. For example, prior work [7] found that by “opportunistically” zooming into web pages and configuring settings, users could automatically benefit from improved accessibility. Our motivation is similar: we aimed to improve code accessibility while introducing minimal changes to existing AI-assistant developer workflows *i.e.*, GitHub Copilot. According to the Visual Studio Marketplace, GitHub Copilot has been installed over 20 million times (as of the time of writing), suggesting that many developers are already familiar with the plugin’s interactions, tooling, and interface. Our results show that CodeA11y significantly improved code accessibility while maintaining a similar (slightly improved) ease of use to Copilot. This suggests that if GitHub Copilot included our set of features or were willing to use a similar plugin without requiring substantial deviation from existing workflows, millions of developers could start writing more accessible code immediately.

AI-Assisted but Developer-Completed. Even when developers adopt accessibility tools, additional expertise is required to maximize their utility. This is especially true for AI assistants, which are incapable of generating entirely correct or accessible code. Our work offers some insight into the manual effort needed to write

accessible code. Both our formative and evaluation studies underscore the necessity for developers to manually intervene with AI-generated code to effectively implement accessibility features. A recurrent challenge in AI-assisted coding is the generation of incomplete or boilerplate code, which often requires developers to take additional steps for completion and validation. Our findings reveal that novice developers tend to critically evaluate AI outputs in areas they prioritize, such as visual enhancements, while overlooking aspects they are less familiar with, like accessibility.

One of the tensions of this work is that while we aim to increase the adoption of AI-driven accessibility tools among users with little expertise or awareness, some degree of understanding is required to use these tools effectively. CodeA11y and other tools can employ strategies to scaffold this interaction, *e.g.*, asking a user “can you describe what’s in this image?” instead of asking them directly for “alternative text.” However, developers ultimately need to be willing to expend additional effort and manually implement the more challenging aspects of this work. Thus, while our work suggests that it is possible to “silently” improve the accessibility of developer-written code, it is ultimately not a replacement for better accessibility awareness, development practices, and education. Nevertheless, CodeA11y could help gradually improve awareness by slowly introducing and explaining accessibility concepts to users after they have found benefits from using the tool.

Limitations & Future Work. We describe several limitations in the current scope of the study and identify avenues for future work to build upon our findings:

First, the utility of the CodeA11y plugin was limited by the constraints placed by our target development environment (Visual Studio Code). Because CodeA11y was implemented as a Copilot plugin, we could only access a few APIs available to standard VS-Code IDE plugins. The Copilot plugin infrastructure was limiting because it restricted the source code that could be passed to the model (*i.e.*, context window length). Our implementation contained some mechanisms for heuristically determining the most relevant files but ultimately serves as a proof of concept of what would be possible in the future, well-integrated version (*e.g.*, built into Copilot). These factors affected the code generation of our system.

Second, although our user study provides statistical evidence that CodeA11y helps developers write more accessible website code, we acknowledge certain limitations of AI coding assistants that could affect their overall effectiveness and reliability. One well-documented issue, particularly with proactive AI assistants [17], is their potential to provide untimely or irrelevant guidance. For instance, the models may occasionally suggest fixes for problems that do not exist (*i.e.*, false positives), such as recommending changes to code that is already fully compliant with accessibility standards. While our study did not surface such occurrences—likely because CodeA11y’s suggestions were tied directly to verified issues from an accessibility checker, rather than the tool identifying issues on its own—this risk becomes more salient as AI assistants evolve to more proactively identify and address accessibility problems. Still, prior research suggests that developers often tolerate false positives more readily than false negatives [39], reasoning that overly cautious guidance from an assistant is less harmful than failing to flag genuine accessibility issues. Indeed, even standalone accessibility

checkers, which are also known for producing false positives [34], have been widely adopted due to their overall beneficial effect on UI quality. Moreover, the occasional presence of false positives does not necessarily negate the value of employing such tools. By raising awareness and prompting developers to consider accessibility from the outset, AI assistants can help cultivate a proactive mindset toward inclusive design. In this sense, the technology does not need to achieve perfect accuracy to have a net positive effect. As the underlying models and APIs improve, and assistants become better integrated with real-world workflows, their precision and utility in improving accessibility are likely to increase.

Third, while our study demonstrates the potential of CodeA11y to encourage developers to adopt accessibility practices, we acknowledge that the broader impact of AI coding assistants on long-term learning and behavior changes (*e.g.*, for accessibility awareness) remains underexplored in the current scope of the study. Research on how real-time AI tools can help developers internalize new practices, such as accessibility, or foster long-term behavioral changes could have strengthened the case for CodeA11y’s instructional components. For instance, prior work has shown that AI coding tools can enhance immediate task performance but may not consistently lead to deeper learning or sustained skill retention [38]. Similarly, studies on meta-cognitive demands in AI-assisted workflows emphasize the importance of tools promoting reflective learning and adaptive strategies, particularly as developers integrate them into their daily practices [84]. Although our findings suggest that CodeA11y has the potential to raise awareness of accessibility issues through direct integration with verified accessibility checks, further research is needed to understand whether such tools can foster a lasting developer’s commitment to accessibility or similar best practices. Additionally, exploring how these tools impact broader developer workflows, collaboration habits, and the ability to generalize learned behaviors across contexts would provide a more comprehensive view of their instructional value. By examining these dimensions, future studies could better elucidate the role of AI coding assistants in shaping not just productivity, but also the culture of inclusive and responsible software development. While CodeA11y focuses on improving accessibility, its approach could extend to other non-functional requirements, such as privacy and security. Investigating how specialized copilots could be seamlessly invoked within mainstream coding assistants for high-stakes scenarios—such as leveraging CodeA11y for front-end development tasks—represents a promising direction.

Finally, we see opportunities to iterate on and refine CodeA11y’s design. Our conservative approach adhered closely to Copilot’s existing interface to minimize friction during adoption. Future work could explore how developers respond to new features and interactions, identifying areas where innovation could enhance usability and functionality without compromising adoption. By addressing current limitations and exploring broader applications, tools like CodeA11y can refine how developers approach accessibility and other critical non-functional requirements. Beyond technical improvements, such advancements hold the potential to redefine AI’s role in shaping more inclusive, secure, and efficient coding practices.

9 CONCLUSION

Our work bridges decades of accessibility efforts with AI coding assistants, offering a novel solution to persistent web accessibility challenges. Through a formative user study, we identify shortcomings in how current AI-assisted development workflows handle accessibility implementation. Accordingly, we develop CodeA11y, a GitHub Copilot Extension, and demonstrate that novice developers using it are significantly more likely to create accessible interfaces. By focusing on integrating accessibility improvements seamlessly into everyday development workflows, this work marks a first step toward fostering accessibility-conscious practices in human-AI collaborative UI development.

REFERENCES

- [1] 2024. GitHub Copilot. <https://github.com/features/copilot>. [Accessed: 2024-04-22].
- [2] 2024. IAB Website Categories. <https://docs.webshrinker.com/v3/iab-website-categories.html#iab-categories>. Accessed: 2024-04-22.
- [3] 2024. SimilarWeb - Website Traffic & Market Intelligence. <http://similarweb.com>. Accessed: 2024-04-22.
- [4] Wajdi Aljedaani, Abdulrahman Habib, Ahmed Aljohani, Marcelo Eler, and Yunhe Feng. 2024. Does ChatGPT Generate Accessible Code? Investigating Accessibility Challenges in LLM-Generated Source Code. In *Proceedings of the 21st International Web for All Conference (Singapore, Singapore) (W4A '24)*. Association for Computing Machinery, New York, NY, USA, 165–176. <https://doi.org/10.1145/3677846.3677854>
- [5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609* (2023).
- [6] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966* (2023).
- [7] Jeffrey P Bigham. 2014. Making the web easier to see with opportunistic accessibility improvement. In *Proceedings of the 27th annual ACM symposium on User interface software and technology*, 117–122.
- [8] Jeffrey P Bigham, Jeremy T Brudvik, and Bernie Zhang. 2010. Accessibility by demonstration: enabling end users to guide developers to web accessibility solutions. In *Proceedings of the 12th international ACM SIGACCESS conference on Computers and accessibility*, 35–42.
- [9] Jeffrey P. Bigham, Ryan S. Kaminsky, Richard E. Ladner, Oscar M. Danielsson, and Gordon L. Hempton. 2006. WebInSight: making web images accessible. In *Proceedings of the 8th International ACM SIGACCESS Conference on Computers and Accessibility (Portland, Oregon, USA) (Assets '06)*. Association for Computing Machinery, New York, NY, USA, 181–188. <https://doi.org/10.1145/1168987.1169018>
- [10] Jeffrey P Bigham and Richard E Ladner. 2007. Accessmonkey: a collaborative scripting framework for web users and developers. In *Proceedings of the 2007 international cross-disciplinary conference on Web accessibility (W4A)*, 25–34.
- [11] Giorgio Brajnik and Raffaella Lomuscio. 2007. SAMBA: a semi-automatic method for measuring barriers of accessibility. In *Proceedings of the 9th International ACM SIGACCESS Conference on Computers and Accessibility*, 43–50.
- [12] Jeremy T. Brudvik, Jeffrey P. Bigham, Anna C. Cavender, and Richard E. Ladner. 2008. Hunting for headings: sighted labeling vs. automatic classification of headings. In *Proceedings of the 10th International ACM SIGACCESS Conference on Computers and Accessibility (Halifax, Nova Scotia, Canada) (Assets '08)*. Association for Computing Machinery, New York, NY, USA, 201–208. <https://doi.org/10.1145/1414471.1414508>
- [13] Ben Caldwell, Michael Cooper, Loretta Guarino Reid, Gregg Vanderheiden, Wendy Chisholm, John Slatin, and Jason White. 2008. Web content accessibility guidelines (WCAG) 2.0. *WWW Consortium (W3C) 290*, 1–34 (2008), 5–12.
- [14] Jieshan Chen, Chunyang Chen, Zhenchang Xing, Xiwei Xu, Liming Zhu, Guoqiang Li, and Jinshui Wang. 2020. Unblind your apps: Predicting natural-language labels for mobile gui components by deep learning. In *Proceedings of the ACM/IEEE 42nd international conference on software engineering*, 322–334.
- [15] Jieshan Chen, Mulong Xie, Zhenchang Xing, Chunyang Chen, Xiwei Xu, Liming Zhu, and Guoqiang Li. 2020. Object detection for graphical user interface: Old fashioned or deep learning or a combination?. In *proceedings of the 28th ACM joint meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 1202–1214.
- [16] Mark Chen, Jerry Twork, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).
- [17] Valerie Chen, Alan Zhu, Sebastian Zhao, Hussein Mozannar, David Sontag, and Ameet Talwalkar. 2024. Need Help? Designing Proactive AI Assistants for Programming. *arXiv preprint arXiv:2410.04596* (2024).
- [18] Wendy Chisholm, Gregg Vanderheiden, and Ian Jacobs. 2001. Web content accessibility guidelines 1.0. *Interactions* 8, 4 (2001), 35–54.
- [19] Victoria Clarke and Virginia Braun. 2017. Thematic analysis. *The journal of positive psychology* 12, 3 (2017), 297–298.
- [20] Giovanni Delnevo, Manuel Andruccioli, and Silvia Mirri. 2024. On the Interaction with Large Language Models for Web Accessibility: Implications and Challenges. In *2024 IEEE 21st Consumer Communications & Networking Conference (CCNC)*. IEEE, 1–6.
- [21] Deque Systems. 2024. axe: Accessibility Testing Tools and Software. <https://www.deque.com/axe/> Accessed: 2024-12-10.
- [22] Morgan Dixon and James Fogarty. 2010. Prefab: implementing advanced behaviors using pixel-based reverse engineering of interface structure. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 1525–1534.
- [23] dkandalov. 2024. activity-tracker. <https://github.com/dkandalov/activity-tracker>.
- [24] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [25] André P Freire, Renata PM Fortes, Marcelo AS Turine, and Debora MB Paiva. 2008. An evaluation of web accessibility metrics based on their attributes. In *Proceedings of the 26th annual ACM international conference on Design of communication*, 73–80.
- [26] Cole Gleason, Amy Pavel, Himalini Gururaj, Kris Kitani, and Jeffrey Bigham. 2020. Making GIFs Accessible. In *Proceedings of the 22nd International ACM SIGACCESS Conference on Computers and Accessibility*, 1–10.
- [27] Cole Gleason, Amy Pavel, Xingyu Liu, Patrick Carrington, Lydia B Chilton, and Jeffrey P Bigham. 2019. Making memes accessible. In *Proceedings of the 21st International ACM SIGACCESS Conference on Computers and Accessibility*, 367–376.
- [28] Cole Gleason, Amy Pavel, Emma McNamee, Christina Low, Patrick Carrington, Kris 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 (Honolulu, HI, USA) (CHI '20)*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3313831.3376728>
- [29] Google. 2024. Chrome DevTools Documentation. <https://developer.chrome.com/docs/devtools> Accessed: 2024-12-10.
- [30] 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.
- [31] Yun Huang, Brian Dobreski, Bijay Bhaskar Deo, Jiahang Xin, Natã Micael Barbosa, Yang Wang, and Jeffrey P Bigham. 2015. CAN: Composable accessibility infrastructure via data-driven crowdsourcing. In *Proceedings of the 12th International Web for All Conference*, 1–10.
- [32] Mina Huh, Saelyne Yang, Yi-Hao Peng, Xiang'Anthony' Chen, Young-Ho Kim, and Amy Pavel. 2023. AVscript: Accessible Video Editing with Audio-Visual Scripts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 1–17.
- [33] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2.5-coder technical report. *arXiv preprint arXiv:2409.12186* (2024).
- [34] Syed Fatiul Huq, Abdulaziz Alshayban, Ziyao He, and Sam Malek. 2023. #A11yDev: Understanding Contemporary Software Accessibility Practices from Twitter Conversations. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 1–18.
- [35] IBM. 2024. IBM Equal Access Toolkit. <https://www.ibm.com/able/toolkit> Accessed: 2024-12-10.
- [36] Muhammad Asiful Islam, Yevgen Borodin, and I. V. Ramakrishnan. 2010. Mixture model based label association techniques for web accessibility. In *Proceedings of the 23rd Annual ACM Symposium on User Interface Software and Technology (New York, New York, USA) (UIST '10)*. Association for Computing Machinery, New York, NY, USA, 67–76. <https://doi.org/10.1145/1866029.1866041>
- [37] Shinya Kawanaka, Yevgen Borodin, Jeffrey P Bigham, Darren Lunn, Hironobu Takagi, and Chieko Asakawa. 2008. Accessibility commons: a metadata infrastructure for web accessibility. In *Proceedings of the 10th international ACM SIGACCESS conference on Computers and accessibility*, 153–160.
- [38] Majeed Kazemitabaar, Justin Chow, Carl Ka To Ma, Barbara J Ericson, David Weintrop, and Tovi Grossman. 2023. Studying the effect of AI code generators on supporting novice learners in introductory programming. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 1–23.
- [39] Rafal Kocielnik, Saleema Amershi, and Paul N Bennett. 2019. Will you accept an imperfect ai? exploring designs for adjusting end-user expectations of ai systems. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 1–14.

- [40] Jonathan Lazar, Alfreda Dudley-Sponaugle, and Kisha-Dawn Greenidge. 2004. Improving web accessibility: a study of webmaster perceptions. *Computers in human behavior* 20, 2 (2004), 269–288.
- [41] Jaewook Lee, Yi-Hao Peng, Jaylin Herskovitz, and Anhong Guo. 2021. Image Explorer: Multi-layered touch exploration to make images accessible. In *Proceedings of the 23rd International ACM SIGACCESS Conference on Computers and Accessibility*. 1–4.
- [42] Amanda Li, Jason Wu, and Jeffrey P Bigham. 2023. Using llms to customize the ui of webpages. In *Adjunct Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–3.
- [43] Jingyi Li, Son Kim, Joshua A Miele, Maneesh Agrawala, and Sean Follmer. 2019. Editing spatial layouts through tactile templates for people with visual impairments. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–11.
- [44] Haotian Liu, Chanyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems* 36 (2024).
- [45] Juan-Miguel López-Gil and Juanan Pereira. 2024. Turning manual web accessibility success criteria into automatic: an LLM-based approach. *Universal Access in the Information Society* (2024), 1–16.
- [46] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173* (2024).
- [47] Kelly Mack, Edward Cutrell, Bongshin Lee, and Meredith Ringel Morris. 2021. Designing tools for high-quality alt text authoring. In *Proceedings of the 23rd International ACM SIGACCESS Conference on Computers and Accessibility*. 1–14.
- [48] Liliu Martin, Catherine Baker, Kristen Shinohara, and Yasmine N Elglaly. 2022. The Landscape of Accessibility Skill Set in the Software Industry Positions. In *Proceedings of the 24th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–4.
- [49] Beatriz Martins and Carlos Duarte. 2024. Large-scale study of web accessibility metrics. *Universal Access in the Information Society* 23, 1 (2024), 411–434.
- [50] Forough Mehralian, Titus Barik, Jeff Nichols, and Amanda Swearngin. 2024. Automated Code Fix Suggestions for Accessibility Issues in Mobile Apps. *arXiv preprint arXiv:2408.03827* (2024).
- [51] Microsoft. 2024. Accessibility Insights. <https://accessibilityinsights.io> Accessed: 2024-12-10.
- [52] Peya Mowar, Yi-Hao Peng, Aaron Steinfeld, and Jeffrey P Bigham. 2024. Tab to Autocomplete: The Effects of AI Coding Assistants on Web Accessibility. In *Proceedings of the 26th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–6.
- [53] Hussein Mozannar, Gagan Bansal, Adam Fourney, and Eric Horvitz. 2024. Reading between the lines: Modeling user behavior and costs in AI-assisted programming. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–16.
- [54] BBC News. 2024. BBC Home - Breaking News, World News, US News, Sports ... <https://www.bbc.com/>
- [55] Achraf Othman, Amira Dhouib, and Aljazi Nasser Al Jabor. 2023. Fostering websites accessibility: A case study on the use of the Large Language Models ChatGPT for automatic remediation. In *Proceedings of the 16th International Conference on Pervasive Technologies Related to Assistive Environments*. 707–713.
- [56] Luis P. Carvalho, Tiago Guerreiro, Shaun Lawson, and Kyle Montague. 2023. Towards real-time and large-scale web accessibility. In *Proceedings of the 25th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–9.
- [57] Maulishree Pandey and Tao Dong. 2023. Blending Accessibility in UI Framework Documentation to Build Awareness. In *Proceedings of the 25th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–12.
- [58] Rohan Patel, Pedro Breton, Catherine M. Baker, Yasmine N. El-Glaly, and Kristen Shinohara. 2020. Why Software is Not Accessible: Technology Professionals’ Perspectives and Challenges. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI EA '20). Association for Computing Machinery, New York, NY, USA, 1–9. <https://doi.org/10.1145/3334480.3383103>
- [59] Amy Pavel, Gabriel Reyes, and Jeffrey P Bigham. 2020. Rescribe: Authoring and automatically editing audio descriptions. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*. 747–759.
- [60] Yi-Hao Peng, Jeffrey P Bigham, and Amy Pavel. 2021. Slidecho: Flexible non-visual exploration of presentation videos. In *Proceedings of the 23rd International ACM SIGACCESS Conference on Computers and Accessibility*. 1–12.
- [61] Yi-Hao Peng, Peggy Chi, Anjali Kannan, Meredith Ringel Morris, and Irfan Essa. 2023. Slide Gestalt: Automatic Structure Extraction in Slide Decks for Non-Visual Access. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [62] Yi-Hao Peng, Faria Huq, Yue Jiang, Jason Wu, Xin Yue Li, Jeffrey P Bigham, and Amy Pavel. 2025. DreamStruct: Understanding Slides and User Interfaces via Synthetic Data Generation. In *European Conference on Computer Vision*. Springer, 466–485.
- [63] Yi-Hao Peng, JiWoong Jang, Jeffrey P Bigham, and Amy Pavel. 2021. Say it all: Feedback for improving non-visual presentation accessibility. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–12.
- [64] Yi-Hao Peng, Muh-Tarnng Lin, Yi Chen, Tzu-Chuan Chen, Pin Sung Ku, Paul Tael, Chin Guan Lim, and Mike Y Chen. 2019. PersonalTouch: Improving touchscreen usability by personalizing accessibility settings based on individual user’s touchscreen interaction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–11.
- [65] Yi-Hao Peng, Jason Wu, Jeffrey Bigham, and Amy Pavel. 2022. Diffscraper: Describing Visual Design Changes to Support Mixed-Ability Collaborative Presentation Authoring. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–13.
- [66] Athira Pillai, Kristen Shinohara, and Garreth W Tigwell. 2022. Website builders still contribute to inaccessible web design. In *Proceedings of the 24th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–4.
- [67] Venkatesh Potluri, Tadashi Grindeland, Jon E Froehlich, and Jennifer Mankoff. 2019. Ai-assisted ui design for blind and low-vision creators. In *the ASSETS'19 Workshop: AI Fairness for People with Disabilities*.
- [68] Kubernetes Project. 2024. Kubernetes. <https://kubernetes.io/>
- [69] Katharina Reinecke, David R. Flatla, and Christopher Brooks. 2016. Enabling Designers to Foresee Which Colors Users Cannot See. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) (CHI '16). Association for Computing Machinery, New York, NY, USA, 2693–2704. <https://doi.org/10.1145/2858036.2858077>
- [70] Anne Spencer Ross, Xiaoyi Zhang, James Fogarty, and Jacob O Wobbrock. 2018. Examining image-based button labeling for accessibility in Android apps through large-scale analysis. In *Proceedings of the 20th International ACM SIGACCESS Conference on Computers and Accessibility*. 119–130.
- [71] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950* (2023).
- [72] Navid Salehnamadi, Abdulaziz Alshayban, Jun-Wei Lin, Iftekhar Ahmed, Stacy Branham, and Sam Malek. 2021. Latte: Use-Case and Assistive-Service Driven Automated Accessibility Testing Framework for Android. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 274, 11 pages. <https://doi.org/10.1145/3411764.3445455>
- [73] Navid Salehnamadi, Ziyao He, and Sam Malek. 2023. Assistive-Technology Aided Manual Accessibility Testing in Mobile Apps, Powered by Record-and-Replay. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (CHI '23). Association for Computing Machinery, New York, NY, USA, Article 73, 20 pages. <https://doi.org/10.1145/3544548.3580679>
- [74] Navid Salehnamadi, Forough Mehralian, and Sam Malek. 2023. Groundhog: An Automated Accessibility Crawler for Mobile Apps. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering* (Rochester, MI, USA) (ASE '22). Association for Computing Machinery, New York, NY, USA, Article 50, 12 pages. <https://doi.org/10.1145/3551349.3556905>
- [75] Daisuke Sato, Masatomo Kobayashi, Hironobu Takagi, and Chieko Asakawa. 2009. What’s next? a visual editor for correcting reading order. In *Human-Computer Interaction—INTERACT 2009: 12th IFIP TC 13 International Conference, Uppsala, Sweden, August 24–28, 2009, Proceedings, Part I 12*. Springer, 364–377.
- [76] Daisuke Sato, Masatomo Kobayashi, Hironobu Takagi, and Chieko Asakawa. 2010. Social accessibility: the challenge of improving web accessibility through collaboration. In *Proceedings of the 2010 International Cross Disciplinary Conference on Web Accessibility (W4A)*. 1–2.
- [77] Brian Sierkowski. 2002. Achieving web accessibility. In *Proceedings of the 30th annual ACM SIGUCCS conference on User services*. 288–291.
- [78] David Sloan, Andy Heath, Fraser Hamilton, Brian Kelly, Helen Petrie, and Lawrie Phipps. 2006. Contextual web accessibility-maximizing the benefit of accessibility guidelines. In *Proceedings of the 2006 international cross-disciplinary workshop on Web accessibility (W4A): Building the mobile web: rediscovering accessibility?* 121–131.
- [79] Amanda Swearngin, Jason Wu, Xiaoyi Zhang, Esteban Gomez, Jen Coughenour, Rachel Stukenborg, Bhavya Garg, Greg Hughes, Adriana Hilliard, Jeffrey P Bigham, et al. 2024. Towards Automated Accessibility Report Generation for Mobile Apps. *ACM Transactions on Computer-Human Interaction* 31, 4 (2024), 1–44.
- [80] Maryam Taeb, Amanda Swearngin, Eldon Schoop, Ruijia Cheng, Yue Jiang, and Jeffrey Nichols. 2024. Axnav: Replaying accessibility tests from natural language. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–16.
- [81] Hironobu Takagi, Chieko Asakawa, Kentarou Fukuda, and Junji Maeda. 2003. Accessibility designer: visualizing usability for the blind. *SIGACCESS Access. Comput.* 77–78 (sep 2003), 177–184. <https://doi.org/10.1145/1029014.1028662>
- [82] Hironobu Takagi, Chieko Asakawa, Kentarou Fukuda, and Junji Maeda. 2003. Accessibility designer: visualizing usability for the blind. *ACM SIGACCESS accessibility and computing* 77-78 (2003), 177–184.

- [83] Hironobu Takagi, Shinya Kawanaka, Masatomo Kobayashi, Daisuke Sato, and Chieko Asakawa. 2009. Collaborative web accessibility improvement: challenges and possibilities. In *Proceedings of the 11th international ACM SIGACCESS conference on Computers and accessibility*. 195–202.
- [84] Lev Tankelevitch, Viktor Kewenig, Auste Simkute, Ava Elizabeth Scott, Advait Sarkar, Abigail Sellen, and Sean Rintel. 2024. The metacognitive demands and opportunities of generative AI. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–24.
- [85] Markel Vigo, Myriam Arrue, Giorgio Brajnik, Raffaella Lomuscio, and Julio Abascal. 2007. Quantitative metrics for measuring web accessibility. In *Proceedings of the 2007 international cross-disciplinary conference on Web accessibility (W4A)*. 99–107.
- [86] Markel Vigo, Justin Brown, and Vivienne Conway. 2013. Benchmarking web accessibility evaluation tools: measuring the harm of sole reliance on automated tests. In *Proceedings of the 10th international cross-disciplinary conference on web accessibility*. 1–10.
- [87] WebAIM. 2024. The WebAIM Million - The 2024 report on the accessibility of the top 1,000,000 home pages. <https://webaim.org/projects/million/>. Accessed: 2024-04-22.
- [88] Jason Wu, Gabriel Reyes, Sam C White, Xiaoyi Zhang, and Jeffrey P Bigham. 2021. When can accessibility help? An exploration of accessibility feature recommendation on mobile devices. In *Proceedings of the 18th international web for all conference*. 1–12.
- [89] Jason Wu, Siyan Wang, Siman Shen, Yi-Hao Peng, Jeffrey Nichols, and Jeffrey P Bigham. 2023. Webui: A dataset for enhancing visual ui understanding with web semantics. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [90] Yeliz Yesilada, Giorgio Brajnik, Markel Vigo, and Simon Harper. 2012. Understanding web accessibility and its drivers. In *Proceedings of the international cross-disciplinary conference on web accessibility*. 1–9.
- [91] Xiaoyi Zhang, Lilian De Greef, Amanda Swearngin, Samuel White, Kyle Murray, Lisa Yu, Qi Shan, Jeffrey Nichols, Jason Wu, Chris Fleizach, et al. 2021. Screen recognition: Creating accessibility metadata for mobile applications from pixels. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [92] Xiaoyi Zhang, Anne Spencer Ross, Anat Caspi, James Fogarty, and Jacob O Wobbrock. 2017. Interaction proxies for runtime repair and enhancement of mobile application accessibility. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 6024–6037.
- [93] Xiaoyi Zhang, Anne Spencer Ross, and James Fogarty. 2018. Robust annotation of mobile application interfaces in methods for accessibility repair and enhancement. In *Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology*. 609–621.
- [94] Yuxin Zhang, Sen Chen, Lingling Fan, Chunyang Chen, and Xiaohong Li. 2023. Automated and Context-Aware Repair of Color-Related Accessibility Issues for Android Apps. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (San Francisco, CA, USA) (ESEC/FSE 2023)*. Association for Computing Machinery, New York, NY, USA, 1255–1267. <https://doi.org/10.1145/3611643.3616329>
- [95] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568* (2023).