

Not All Visitors are Bilingual: A Measurement Study of the Multilingual Web from an Accessibility Perspective

Masudul Hasan Masud Bhuiyan

masudul.bhuiyan@cispa.de

CISPA Helmholtz Center for Information Security
Saarbrücken, Germany

Yasir Zaki

yasir.zaki@nyu.edu

New York University Abu Dhabi
Abu Dhabi, UAE

Matteo Varvello

matteo.varvello@nokia.com

Nokia Bell Labs
Holmdel, NJ, USA

Cristian-Alexandru Staicu

staicu@cispa.de

CISPA Helmholtz Center for Information Security
Saarbrücken, Germany

ABSTRACT

English is the predominant language on the web, powering nearly half of the world’s top ten million websites. Support for multilingual content is nevertheless growing, with many websites increasingly combining English with regional or native languages in both visible content and hidden metadata. This multilingualism introduces significant barriers for users with visual impairments, as assistive technologies like screen readers frequently lack robust support for non-Latin scripts and misrender or mispronounce non-English text, compounding accessibility challenges across diverse linguistic contexts. Yet, large-scale studies of this issue have been limited by the lack of comprehensive datasets on multilingual web content. To address this gap, we introduce LangCrUX, the first large-scale dataset of 120,000 popular websites across 12 languages that primarily use non-Latin scripts. Leveraging this dataset, we conduct a systematic analysis of multilingual web accessibility and uncover widespread neglect of accessibility hints. We find that these hints often fail to reflect the language diversity of visible content, reducing the effectiveness of screen readers and limiting web accessibility. We finally propose Kizuki, a language-aware automated accessibility testing extension to account for the limited utility of language-inconsistent accessibility hints.

1 INTRODUCTION

An estimated 2.2 billion people live with some form of visual impairment, and 90% of them reside in low- and middle-income countries [44]. As internet adoption grows in these regions, an increasing share of web content is created in languages other than English. Many of these languages use non-Latin writing systems such as Devanagari, Bengali, Arabic, or Thai and are often presented alongside English in mixed-language interfaces [27–29]. Popular screen readers like JAWS [2] and NVDA [7] still exhibit limited support for non-Latin scripts and often perform poorly when confronted

with mixed languages [14, 24, 39, 40]. These tools may misrender non-English words or produce unintelligible output, as has been documented with languages like Nepali [37].

The root of the problem lies in the inadequate usage of *text alternatives* [12, 23] metadata such as lang attributes or image alt text, which screen readers use to process content appropriately. When such metadata is absent, incorrect, or inconsistent with the visible text, it creates a mismatch between the displayed content and what the assistive tool offers. This issue is compounded for scripts that require complex shaping or language-specific pronunciation models. For example, Apple’s VoiceOver [11], the default screen reader on macOS and iOS devices, does not provide any support for languages such as Urdu, Amharic, or Burmese [11]. This lack of linguistic inclusivity is at odds with the principles laid out by the W3C Web Accessibility Initiative, which advocates for equal digital access regardless of language or ability [43].

Web accessibility research only marginally studies the intersection of multilingualism and assistive technology [16, 17]. A key bottleneck in this area is the lack of data. Widely used datasets like Tranco [30] focus on popularity but do not provide any insight into the linguistic composition of websites. This limits the researchers’ ability to measure the prevalence of multilingual content and to assess the accessibility landscape for speakers of less commonly supported languages. To bridge this gap, we introduce LangCrUX, a large-scale dataset of 120,000 popular websites across 12 languages that primarily use non-Latin scripts. We construct LangCrUX by selecting high-traffic websites from the Chrome User Experience Report (CrUX) dataset, verifying language use via automated script detection and manual sampling, and through Puppeteer-based crawls routed via country-specific VPNs to capture localized versions of each site.

We leverage LangCrUX to conduct the first large-scale analysis of multilingual web accessibility, focusing on how assistive technologies interact with real-world, multilingual web content. Our findings reveal that nearly 40% of websites in

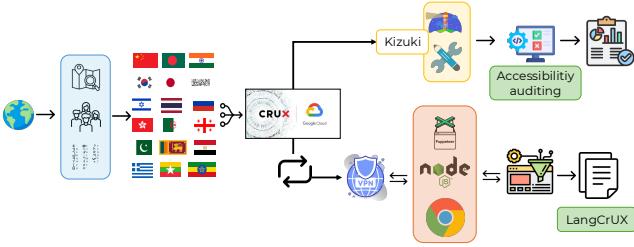


Figure 1: Overview of our methodology for constructing and analyzing the LangCrUX dataset.

Bangladesh and India lack any accessibility text in the native language, despite having predominantly native-language visible content. More broadly, language-aware accessibility remains a widespread and under-addressed issue: many websites use non-descriptive, placeholder, or untranslated text (e.g., “file1,” “button,” or generic English terms) in critical accessibility elements like image alt text. These mismatches significantly reduce the utility of screen readers, which rely on metadata to convey meaningful information to users. We observe that current automated accessibility testing tools fail to detect these language inconsistencies, as they typically only evaluate the presence of accessibility hints. To address this gap, we propose and develop **Kizuki**¹, a testing extension that identifies such mismatches and evaluates metadata based on alignment with the surrounding linguistic context, offering a more inclusive measure of accessibility.

2 METHODOLOGY

Our study investigates the accessibility of websites in non-Latin-script languages, focusing on linguistic diversity in underrepresented language communities. We collect **LangCrUX**, a dataset comprising websites with a high proportion of non-Latin-script content, combining data from the Chrome User Experience Report (CrUX) [8], custom web crawls using Puppeteer, and language metadata verification. Figure 1 visualizes our methodology, including dataset construction, language selection, and automated accessibility analysis. LangCrUX is released as an open-source dataset on GitHub.²

Language and Country Selection Criteria: We begin with a pool of 26 widely spoken non-Latin-script languages, including Hindi, Bangla, Modern Standard Arabic, Tamil, Telugu, Mandarin Chinese, Urdu, Amharic, Russian, Marathi, and others [6]. Our goal is to identify a diverse, representative set of languages underrepresented in web accessibility research, particularly those using non-Latin writing systems. Language selection is guided by three main factors: script type (non-Latin), size of the global speaker base, and geographic and linguistic diversity. To ensure representativeness

¹named after the Japanese word for “awareness”

²<https://anonymous.4open.science/r/LangCrux-F68F>

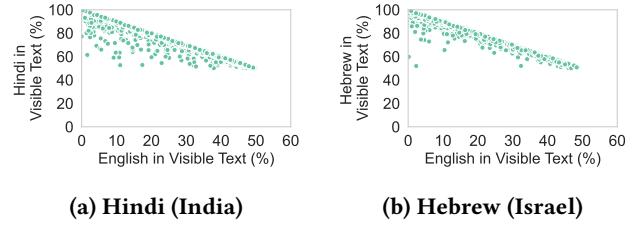


Figure 2: Native language distribution in visible text for India and Israel in LangCrUX

and sufficient data, we apply two strict inclusion criteria: 1) at least 10,000 websites with 50% or more visible textual content in the target language, and 2) inclusion in the CrUX dataset – which provides user experience metrics from Chrome users – with sufficient traffic and performance data. Applying these filters results in a final set of 12 language-country pairs, each uniquely representing a distinct non-Latin-script language with a verifiable web presence. Because the initial pool alone did not yield enough languages that met these thresholds, we expanded the selection to include additional ones such as Hebrew, Sinhala, Greek, and Burmese. These were added to increase script and regional diversity while still satisfying our inclusion criteria.

The selected countries, their corresponding languages, and approximate global speaker populations are China (Mandarin Chinese, 1.2 billion), India (Hindi, 609 million), Algeria (Modern Standard Arabic, 335 million), Bangladesh (Bangla, 284 million), Russia (Russian, 253 million), Japan (Japanese, 126 million), Egypt (Egyptian Arabic, 119 million), Hong Kong (Cantonese, 85.5 million), South Korea (Korean, 82 million), Thailand (Thai, 71 million), Greece (Greek, 13.5 million), and Israel (Hebrew, 9 million). Collectively, these 12 languages are spoken by over 3.19 billion people, representing about 39.5% of the global population. For languages spoken in multiple countries, such as Modern Standard Arabic, used in Algeria, Saudi Arabia, and Morocco, we select the country with the highest population of native speakers, in this case, Algeria. In multilingual countries like India, we include all major non-Latin-script languages with substantial speaker populations. However, only Hindi meets our data threshold; other widely spoken languages, such as Tamil and Telugu, do not meet the 10,000-website requirement and are excluded. Similar exclusions apply to Sinhala (Sri Lanka) and Georgian (Georgia), where websites with sufficient native-language content fell below the threshold despite initial inclusion.

Website Selection: We use Google CrUX to identify and rank websites by real-world usage metrics. For each selected language-country pair, we extract the top 10,000 websites based on CrUX rankings, which reflect user engagement, load performance, and interaction quality. To validate language presence, we use a Unicode-based heuristic that matches visible text content against script-specific character ranges

(e.g., Devanagari for Hindi, Hangul for Korean, and Cyrillic for Russian). For overlapping scripts, such as Arabic and Urdu, we include additional language-specific characters to improve precision. A website is retained if at least 50% of its visible textual content is in the target language. Websites that do not meet this threshold are excluded and replaced with the next-ranking candidate from the CrUX list. In cases where 10,000 qualifying websites could not be found among the top-ranked entries, we extended our search to lower-ranked websites within the CrUX database to fulfill the quota. Figure 2 illustrates the distribution of visible content by language for two representative cases: Hindi websites in India and Hebrew websites in Israel.

Data Collection: To extract accessibility-related features from the selected websites, we develop a web crawler using Puppeteer [9], which simulates web browsing conditions in a Chromium environment. Each website is visited programmatically, allowing us to capture network-level metadata, page structure, and accessibility indicators such as alternative text, ARIA (Accessible Rich Internet Applications) attributes (which enhance the semantics of web elements for assistive technologies [1]), and declared language tags.

To capture the localized experience of users in each country, we route all browser traffic through VPN servers physically hosted in the corresponding country. This step is critical for collecting region-specific versions of websites, as many sites dynamically serve content—including language settings, layout, or accessibility features—based on the user’s IP location. Without VPN-based localization, web crawlers risk accessing global or English-dominant versions of websites that do not accurately reflect the intended user experience of native speakers. We use a combination of commercial VPN services, including ProtonVPN [10] and Hotspot Shield [22], to achieve broad geographic coverage. Since not all VPN providers have servers in every target country, we select the provider on a per-country basis to ensure reliable and consistent access from within national borders. Compared to crawling from generic cloud-hosted IPs, this approach offers a significantly more realistic vantage point, reducing the risk of content variation, redirection, or censorship artifacts that may otherwise skew accessibility analysis.

Accessibility Element Selection Criteria: To identify accessibility elements where natural language plays a critical role, we follow a structured process based on the Lighthouse accessibility auditing framework [19]. Lighthouse evaluates a range of accessibility checks, each associated with specific HTML elements and best practices. Our goal is to select those elements for which the presence, clarity, and appropriateness of natural language directly influence accessibility outcomes.

button-name	document-title	image-alt
frame-title	summary-name	label
input-image-alt	select-name	link-name
input-button-name	svg-img-alt	object-alt

Table 1: Web elements requiring natural language.

We begin by examining the set of Lighthouse accessibility tests, which internally relies on the Axe-core accessibility engine [5]. For each test, we identify the corresponding HTML element it targets (e.g., ``, `<button>`, `<input>`). We then analyze the test rationale by referencing the associated rule definitions from Axe-core, which provide detailed explanations and criteria for each audit. From these specifications, we assess whether natural language content is integral to the test. Specifically, we ask whether the accessibility of the element depends on human-readable text, for example, whether a screen reader user would rely on the clarity and relevance of that text to understand the element’s purpose. If natural language is central to the function or evaluation of the element, we include it in our set. Following this process, we identify the twelve elements listed in Table 1 as language-sensitive accessibility features.

These elements span a range of interface components, including images, forms, buttons, and navigation elements, all of which rely on meaningful textual descriptions to be accessible to users with visual impairments. We explicitly exclude certain tests, such as `video-caption`. Although captions are inherently language-dependent and critical for accessibility, accurately evaluating them at scale poses challenges. In many cases, captions are not embedded in the HTML but are provided through separate files (e.g., VTT or SRT) or dynamically loaded via JavaScript. These may be inaccessible to crawlers unless the video is played or fully rendered in the browser. Furthermore, identifying whether a video has accurate, synchronized, and complete captions often requires manual inspection or audio-visual comparison—steps outside the scope of automated large-scale analysis. Due to these limitations, we omit `video-caption` checks to maintain consistency and reproducibility in our methodology.

Limitations: Our methodology has several limitations. First, reliance on CrUX limits our scope to websites with measurable Chrome traffic, potentially excluding low-traffic or highly localized websites. Next, while using a VPN allows accessing region-specific versions of websites, some websites may detect VPN use and return generic or restricted versions. In such cases, we replace the affected websites with the next eligible candidate. Additionally, Puppeteer’s simulated browsing environment may not fully reflect user experiences. Finally, language detection (both automated and manual) can be challenged by embedded content or non-standard scripts, though our 50% content threshold and manual verification aim to reduce this risk.

3 IS MULTILINGUAL WEB ACCESSIBLE?

Table 2 provides statistics on the quality and presence of accessibility text across twelve HTML elements. For each element, we report the median, standard deviation, and average percentage of websites where the accessibility attribute is missing or empty. We also include metrics like text length (in characters) and word count to assess richness and verbosity. These measurements allow us to compare how frequently accessibility features are implemented and how informative they are when present. In the following, we analyze these patterns by examining missing and empty values, evaluating text length and word count, filtering uninformative content, and characterizing the language distribution of informative accessibility text.

Prevalence of Missing and Empty Accessibility Texts:

Missing accessibility text is a widespread issue across HTML elements. Several elements exhibit extremely high average missing rates, including `label` (98.5%), `svg-img-alt` (96.6%), `link-name` (95.9%), `input-button-name` (93.9%), and `summary-name` (90.47%). The `image-alt` attribute stands out for its relatively lower average missing rate (17.12%) but exhibits the highest percentage of empty values (25.39%). Although our findings show that it is possible to pass the Lighthouse audit for image accessibility by setting the `alt` attribute to an empty string, such text does not convey meaningful information to users (Appendix D). Compared to non-multilingual websites, the missing percentage is slightly higher (15.19% vs. 17.12%) and the empty percentage is notably higher (16.36% vs. 25.39%), indicating a greater tendency to include but not meaningfully populate the attribute in multilingual contexts.

Attributes like `button-name` and `link-name`, which are essential for identifying interactive components, also show high missing rates (61.92% and 95.96%, respectively). Similarly, attributes associated with form controls, such as `input-button-name` and `select-name`, are frequently absent or left empty, compromising the clarity of form interactions. A likely reason for these high rates is that, in many cases, screen readers fall back to reading visible HTML text (such as the inner text of a button or link) when accessibility attributes like `aria-label`, `label`, or `alt` are missing. This fallback behavior reduces the perceived need for developers to explicitly include accessibility metadata, especially when the element already includes visible text.

Text Length and Word Count Analysis: Table 2 also captures the descriptive quality of accessibility text through text length and word count metrics. Among the elements, `link-name` has a relatively high average text length (~27 characters) and word count (~5), compared to `summary-name`, `select-name`, `label`, and `button-name`, which show lower average word counts (1.17, 2.31, 1.67, and 3.86, respectively).

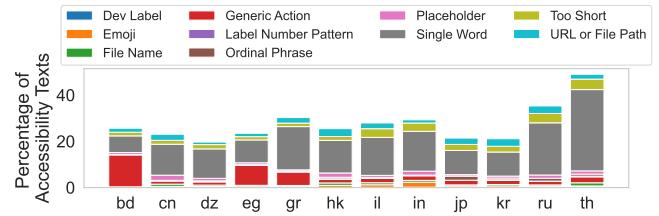


Figure 3: Distribution of filtered accessibility texts by discard reason across countries.

This suggests that link descriptions tend to be more detailed and contextually informative. However, for some other elements, shorter texts are often acceptable; for example, buttons labeled “Login,” “Send,” or “Submit” typically provide sufficient clarity with just one or two words.

For elements like `image-alt`, which require contextual descriptions to convey the meaning of an image, we find an average word count of only ~4. This shortness, especially when combined with the high empty rate noted earlier, suggests a broader trend of developers including minimal or superficial alt text, possibly to satisfy automated checks rather than to genuinely enhance accessibility. Finally, the table reveals substantial outliers, e.g., `image-alt` has a maximum of 261,864 characters and 12,306 words, while `link-name` 5,228 characters and 518 words. These extreme but rare cases likely indicate instances where extraneous content, such as metadata, boilerplate text, or full paragraphs, has been mistakenly inserted into accessibility attributes, potentially overwhelming assistive technologies and undermining user experience (see Appendix E for examples).

Filtering Uninformative Accessibility Text: To assess the quality of accessibility text, we apply a filtering step to discard uninformative or placeholder texts. This is essential because the presence of an `alt` or `aria-label` attribute does not guarantee usefulness. Labels such as `button`, `file1`, or `image1` may satisfy automated checks but provide no semantic value to screen reader users. We define a set of heuristics to classify accessibility texts into eleven categories, distinguishing between useful and discardable content. These include short strings, file paths, placeholders, developer labels, etc. Appendix H shows the full list of filtering rules. Appendix G shows the breakdown by HTML element.

Figure 3 shows the percentage distribution of filtered accessibility texts across countries and categories. One of the most common issues is the use of generic single words. For example, in Thailand, over 33% of accessibility texts are single-word labels. High rates are also observed in Russia (22.2%), Greece (18.03%), and India (17.1%). In contrast, countries like Bangladesh (6.9%) and Egypt (10.5%) show lower proportions. A small but non-negligible portion of texts are too short to convey meaning. In Russia, 4.26% of labels fall into this category, followed by Thailand (4.24%), Israel (4.03%), and India (3.6%). Some websites also use raw URLs or file

Element	Missing (%)			Empty (%)			Text Length			Word Count		
	Median	Std Dev	Mean	Median	Std Dev	Mean	Median	Std Dev	Mean	Median	Std Dev	Mean
button-name	71.43	37.25	61.92	0.00	4.18	0.36	14	25.01	21.35	3	4.72	3.83
frame-title	87.50	30.09	75.81	0.00	3.33	0.21	13	11.57	17.45	1	2.39	2.54
image-alt	1.89	28.86	17.12	7.46	32.40	25.39	14	1332.15	22.97	2	8.27	3.67
input-button-name	100.00	22.62	93.90	0.00	4.03	0.19	12	13.12	14.26	2	2.70	2.83
input-image-alt	0.00	47.17	35.07	0.00	21.27	4.85	3	6.92	5.66	1	1.26	1.41
label	100.00	10.01	98.55	0.00	1.27	0.02	8	6.49	9.28	1	1.04	1.67
link-name	100.00	11.98	95.96	0.00	0.98	0.04	22	27.01	26.64	3	4.39	4.67
object-alt	100.00	23.30	94.19	0.00	5.09	0.26	9	17.64	14.26	1	3.48	2.49
select-name	100.00	28.78	89.84	0.00	2.00	0.05	10	23.03	12.94	2	2.26	2.30
summary-name	100.00	25.84	90.47	0.00	2.98	0.17	5	3.68	5.69	1	0.59	1.18
svg-img-alt	100.00	15.15	96.66	0.00	2.90	0.15	13	5.38	11.98	2	0.74	1.88

Table 2: Accessibility element statistics showing median, standard deviation, and mean values for missing and empty percentages, as well as descriptive richness (text length and word count).

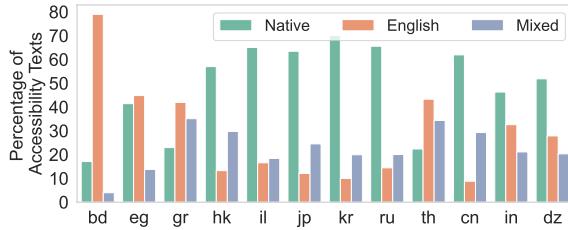


Figure 4: Language distribution of filtered accessibility texts across countries. Only texts classified as potentially useful are included.

paths as labels. This affects 3.8% of labels in Hong Kong, 3.5% in South Korea, and 3.17% in Russia. These findings show that even when accessibility text is present, it often lacks descriptive content. This reinforces the need to go beyond presence-based metrics and assess the actual semantic value of accessibility labels.

Language Distribution of Informative Accessibility Text: After removing uninformative and placeholder accessibility text, we reanalyze the remaining content to understand the language distribution of texts that are potentially useful. This filtered set reflects more intentional and meaningful uses of alt, label, and other accessibility attributes. Figure 4 shows the proportion of accessibility texts written in native languages, English, or a mix of both, across the 12 analyzed countries. A prominent pattern is the heavy reliance on English, even in countries where it is not the primary language. In Bangladesh, 79% of informative accessibility texts are in English, the highest among all countries analyzed. Egypt, Thailand, and Greece also show a strong tendency to default to English. This suggests that developers may rely on English for accessibility metadata due to limited localization tools or being unaware of screen reader needs in native languages.

Another important pattern is the use of mixed-language accessibility hints, where a single alt attribute contains both the native language and English. This occurs frequently in

Greece (35%), Thailand (34%), and Hong Kong (30%), and in over 20% of websites in China, Russia, Japan, and India. While sometimes intended to aid multilingual users, such mixing often confuses screen readers, which typically do not handle language switching within a single label, resulting in mispronunciations or reduced clarity.

4 LANGUAGE-AWARE ACCESSIBILITY

Mismatch Between Visible and Accessibility Text: Figure 5 compares native language usage in visible versus accessibility text across the 12 analyzed countries. While the visible content of many websites is multilingual or predominantly in the native language, the associated accessibility metadata, such as alt text, aria-labels, and form labels, is typically written in English. This mismatch is especially noticeable in countries like India and Bangladesh, where over 40% of websites have less than 10% of their accessibility text in the native language. Thailand, China, and Hong Kong also show similar trends, with more than a quarter of their websites falling into this category. In contrast, countries like Japan and Israel have significantly lower rates of mismatch, with fewer than 10% of websites showing such disparities. For blind users who rely on screen readers, this language discrepancy introduces an additional barrier, forcing them to navigate a bilingual interface where visible content and assistive text do not align. Appendix F provides a detailed view of the mismatch by visualizing the distribution of visible versus accessibility text across all countries in our dataset.

For example, in Bangladesh, <https://teachers.gov.bd> is a widely used government education portal. Although more than 98% of its visible content is in Bangla, only one of the 79 images with alt attributes uses Bangla. Comparable patterns are also seen in India, Thailand, and China. The Indian website <https://cmhelpline.mp.gov.in> has a Hindi version where

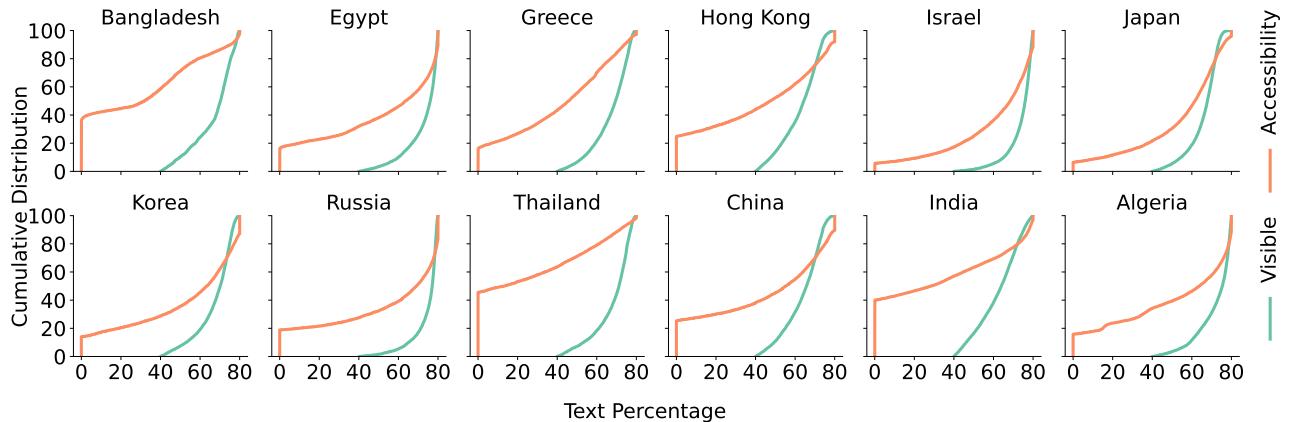


Figure 5: CDFs of native language usage in visible vs. accessibility text across countries. Most websites with visible content in native languages still rely on English in accessibility metadata.

the interface is almost entirely in Hindi, but all accessibility text is in English. The Thai news website <https://www.khaosod.co.th> contains over 92% of its visible content in Thai, but its accessibility labels are mostly in English. The Chinese provincial government site <https://kjt.shaanxi.gov.cn> is almost fully in Chinese, yet its accessibility texts are entirely in English. Mismatch examples are provided in Appendix I.

Adding Language Awareness to Lighthouse: Automated testing tools such as Lighthouse do not consider the language of accessibility text when evaluating compliance. As a result, alt attributes are marked as present regardless of whether their content matches the language of the surrounding interface. To address this limitation, we introduce **Kizuki**, a Lighthouse extension that incorporates language awareness in accessibility evaluation. Specifically, we extend the audit for image alt text to verify whether the description is written in the same language as the page’s visible content.

We evaluate Kizuki on 10,000 websites from Bangladesh and Thailand, two countries where language mismatch between visible content and accessibility metadata is particularly common. For fairness, we exclude websites that fail the original Lighthouse test due to missing alt attributes. Figure 6 shows the resulting shift in accessibility scores. Without considering language, 43% of websites received a Lighthouse score above 90 (considered “good” [3, 4]), and 5.6% achieved a perfect score. After applying Kizuki’s language-aware check, these numbers dropped significantly: only 15.8% of websites scored above 90, and just 1.8% retained a perfect score.

5 RELATED WORK

Several studies have explored the challenges of multilingual web accessibility. Vázquez et al. [32, 33, 35, 41] conducted user studies highlighting the limitations of screen readers in handling multilingual interfaces. Casalegno [42]

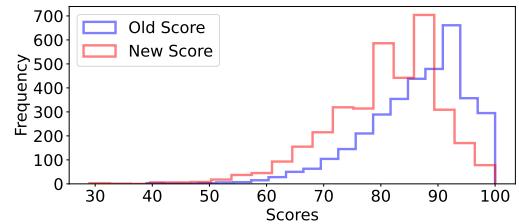


Figure 6: Accessibility score distribution before and after applying Kizuki’s language-aware alt text evaluation on websites from Bangladesh and Thailand.

reported similar findings, emphasizing the cognitive strain users face when navigating mixed-language content. García-Garcinuño et al. [16] confirmed that language mismatches in screen reader output persist even when markup is correctly annotated. Vázquez et al. [34, 36] also showed that accessibility is often excluded from localization workflows, leading to untranslated alt text and inconsistent metadata. Several works highlight problems such as mispronunciation, broken accents, or robotic voices in the context of multilingual content for screen readers [24, 39, 40]. Sankhi et al. [37] documented similar barriers in Nepal, while Raghavendra et al. [31] pointed to the lack of robust multilingual speech synthesis systems in Indian languages, emphasizing infrastructural challenges for regional screen reader development.

Researchers also explored automated solutions to tackle accessibility issues. Several approaches are proposed for alt text generation, including human-curated [18, 46], image-search-based [20, 26], and AI-based methods [13, 15, 45]. While human and search-based approaches emphasize contextual accuracy and reusability, AI-driven techniques offer scalable alternatives, including generating captions [21, 25, 38] or even producing the image itself with embedded descriptions [13]. However, these methods still rely on high-quality training data and often require human oversight to ensure contextual relevance and inclusivity.

6 CONCLUSION

Multilingual web content is increasing in prevalence, but accessibility support lags behind, especially for non-Latin scripts. This paper introduces LangCrUX, the first large-scale dataset of 120,000 popular websites across 12 languages that primarily use non-Latin scripts. Analysis of the LangCrUX dataset reveals widespread issues and motivates language-aware accessibility improvements such as Kizuki, a Google Lighthouse extension we developed, which incorporates language consistency checks into accessibility testing. We hope that our dataset, which we will open source, will spark the community's interest to further measure the implications of an increasingly multilingual web.

REFERENCES

- [1] [n. d.]. ARIA - Accessibility | MDN – developer.mozilla.org. <https://developer.mozilla.org/en-US/docs/Web/Accessibility/ARIA>.
- [2] [n. d.]. JAWS® – Freedom Scientific – freedomscientific.com. <https://www.freedomscientific.com/products/software/jaws/>.
- [3] [n. d.]. Lighthouse performance scoring – graphite.dev. <https://graphite.dev/guides/lighthouse-scoring>.
- [4] [n. d.]. Lighthouse performance scoring | Chrome for Developers – developer.chrome.com. <https://developer.chrome.com/docs/lighthouse/performance/performance-scoring>.
- [5] [n. d.]. List of axe 4.7 rules | Deque University | Deque Systems – dequeuniversity.com. <https://dequeuniversity.com/rules/axe/4.7>.
- [6] [n. d.]. List of languages by total number of speakers - Wikipedia – en.wikipedia.org. https://en.wikipedia.org/wiki/List_of_languages_by_total_number_of_speakers.
- [7] [n. d.]. NV Access – nvaccess.org. <https://www.nvaccess.org/>.
- [8] [n. d.]. Overview of CrUX | Chrome UX Report. <https://developer.chrome.com/docs/crux>
- [9] [n. d.]. Puppeteer | Puppeteer – pptr.dev. <https://pptr.dev/>.
- [10] [n. d.]. The best VPN for speed and security – protonvpn.com. <https://protonvpn.com/>.
- [11] [n. d.]. VoiceOver User Guide for Mac – support.apple.com. <https://support.apple.com/en-gb/guide/voiceover/welcome/mac>.
- [12] [n. d.]. Web Content Accessibility Guidelines (WCAG) 2.1 – w3.org. <https://www.w3.org/TR/WCAG21/#non-text-content>.
- [13] Nouar Aldahoul, Joseph Hong, Matteo Varvello, and Yasir Zaki. 2023. Exploring the Potential of Generative AI for the World Wide Web. doi:10.48550/arXiv.2310.17370 arXiv:2310.17370 [cs].
- [14] Peter Auer. 2013. *Code-switching in conversation: Language, interaction and identity*. Routledge.
- [15] Maitraye Das, Alexander J Fiannaca, Meredith Ringel Morris, Shaun K Kane, and Cynthia L Bennett. 2024. From provenance to aberrations: Image creator and screen reader user perspectives on alt text for AI-generated images. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–21.
- [16] Álvaro García Garcinúñ and Jesús Torres-del Rey. 2024. Multilingual accessibility in human-screen reader interaction with web content: an exploratory study. *Tradumática* 22 (2024), 0426–449.
- [17] Xurxe Toivo García. [n. d.]. The troubled state of screen readers in multilingual situations – uxdesign.cc. <https://uxdesign.cc/the-troubled-state-of-screen-readers-in-multilingual-situations-f6a9da4ecdf3>. [Accessed 13-05-2025].
- [18] Gleason, 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, ASSETS ’19*. New York, NY, USA, 367–376.
- [19] Google Developers. [n. d.]. Lighthouse. <https://developer.chrome.com/docs/lighthouse> Chrome for Developers.
- [20] 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, CHI ’18*. New York, NY, USA, 1–11.
- [21] Margot Hanley, Solon Barocas, Karen Levy, Shiri Azenkot, and Helen Nissenbaum. 2021. Computer vision and conflicting values: Describing people with automated alt text. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 543–554.
- [22] Hotspotshield. [n. d.]. Hotspot Shield: Fastest VPN for Streaming, Gaming & More – hotspotshield.com. <https://www.hotspotshield.com/>.
- [23] Sathish Kumar. [n. d.]. Understanding WCAG SC 1.1.1 Non-text Content • DigitalA11Y – digitala11y.com. <https://www.digitala11y.com/understanding-sc-1-1-1-non-text-content/>.
- [24] Ted McCarthy, Joyojeet Pal, Tanvi Marballi, and Edward Cutrell. 2012. An analysis of screen reader use in India. In *Proceedings of the Fifth International Conference on Information and Communication Technologies and Development*. ACM, Atlanta Georgia USA, 149–158. doi:10.1145/2160673.2160694
- [25] Yunseo Moon, Hyunmin Lee, SeungYoung Oh, and Hyunggu Jung. 2024. SaGol: using MiniGPT-4 to generate alt text for improving image accessibility. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*. 8745–8748.
- [26] Sujeath Pareddy, Anhong Guo, and Jeffrey P. Bigham. 2019. X-ray: Screenshot accessibility via embedded metadata. In *Proceedings of the 21st International ACM SIGACCESS Conference on Computers and Accessibility, ASSETS ’19*. New York, NY, USA, 389–395.
- [27] Daniel Pimienta. [n. d.]. Reliably exploring the presence of languages on the Internet. (n. d.).
- [28] Daniel Pimienta. 2022. Resource: Indicators on the presence of languages in Internet. In *Proceedings of the 1st Annual Meeting of the ELRA/ISCA Special Interest Group on Under-Resourced Languages*. 83–91.
- [29] Daniel Pimienta, Álvaro Blanco, and Gilvan Müller de Oliveira. 2023. The method behind the unprecedented production of indicators of the presence of languages in the Internet. *Frontiers in Research Metrics and Analytics* 8 (2023), 1149347.
- [30] Victor Le Pochat, Tom Van Goethem, Samaneh Tajalizadehkhoob, Maciej Korczyński, and Wouter Joosen. 2018. Tranco: A research-oriented top sites ranking hardened against manipulation. *arXiv preprint arXiv:1806.01156* (2018).
- [31] E Veera Raghavendra and Kishore Prahallad. 2010. A multilingual screen reader in Indian languages. In *2010 National Conference On Communications (NCC)*. IEEE, 1–5.
- [32] Silvia Rodríguez Vázquez. 2014. Introducing web accessibility to localization students: implications for a universal web. In *Proceedings of the 16th international ACM SIGACCESS conference on Computers & accessibility - ASSETS ’14*. ACM Press, Rochester, New York, USA, 333–334. doi:10.1145/2661334.2661414
- [33] Silvia Rodríguez Vázquez. 2015. Exploring Current Accessibility Challenges in the Multilingual Web for Visually-Impaired Users. In *Proceedings of the 24th International Conference on World Wide Web*. ACM, Florence Italy, 871–873. doi:10.1145/2740908.2743010
- [34] Silvia Rodríguez Vázquez. 2016. *Assuring accessibility during web localisation: an empirical investigation on the achievement of appropriate text*

- alternatives for images*. Ph.D. Dissertation. Universidad de Salamanca. doi:10.14201/gredos.132982
- [35] Silvia Rodríguez Vázquez. 2016. Measuring the impact of automated evaluation tools on alternative text quality: a web translation study. In *Proceedings of the 13th International Web for All Conference*. ACM, Montreal Canada, 1–10. doi:10.1145/2899475.2899484
- [36] Silvia Rodríguez Vázquez and Sharon O’Brien. 2017. Bringing Accessibility into the Multilingual Web Production Chain: Perceptions from the Localization Industry. In *Universal Access in Human–Computer Interaction. Design and Development Approaches and Methods*, Margherita Antona and Constantine Stephanidis (Eds.). Vol. 10277. Springer International Publishing, Cham, 238–257. doi:10.1007/978-3-319-58706-6_20 Series Title: Lecture Notes in Computer Science.
- [37] Prakash Sankhi and Frode Eika Sandnes. 2022. A glimpse into smartphone screen reader use among blind teenagers in rural Nepal. *Disability and Rehabilitation: Assistive Technology* 17, 8 (Nov. 2022), 875–881. doi:10.1080/17483107.2020.1818298 Publisher: Taylor & Francis _eprint: <https://doi.org/10.1080/17483107.2020.1818298>
- [38] Yixian Shen, Hang Zhang, Yanxin Shen, Lun Wang, Chuanqi Shi, Shaoshuai Du, and Yiyi Tao. 2024. AltGen: AI-Driven Alt Text Generation for Enhancing EPUB Accessibility. *arXiv preprint arXiv:2501.00113* (2024).
- [39] Aditya Vashistha and Richard Anderson. 2016. Technology use and non-use by low-income blind people in India. *ACM SIGACCESS Accessibility and Computing* 116 (2016), 10–21.
- [40] Aditya Vashistha, Erin Brady, William Thies, and Edward Cutrell. 2014. Educational Content Creation and Sharing by Low-Income Visually Impaired People in India. In *Proceedings of the Fifth ACM Symposium on Computing for Development*. ACM, San Jose California USA, 63–72. doi:10.1145/2674377.2674385
- [41] Silvia Rodríguez Vázquez. 2015. Unlocking the potential of web localizers as contributors to image accessibility: what do evaluation tools have to offer? In *Proceedings of the 12th International Web for All Conference*. ACM, Florence Italy, 1–4. doi:10.1145/2745555.2746662
- [42] Silvia Rodríguez Vázquez and Lucía Morado Vázquez. [n. d.]. Usability of Partially Localised Websites in Switzerland: A Study with Screen Reader Users. ([n. d.]).
- [43] W3C. 2023. Success Criterion 3.1.2: Language of Parts. <https://www.w3.org/WAI/WCAG21/Understanding/language-of-parts.html>
- [44] World Health Organization. 2021. World Report on Vision. <https://www.who.int>
- [45] Shaomei Wu, Jeffrey Wieland, Omid Farivar, and Julie Schiller. 2017. Automatic Alt-text: Computer-generated Image Descriptions for Blind Users on a Social Network Service. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing (CSCW ’17)*. Association for Computing Machinery, New York, NY, USA, 1180–1192. doi:10.1145/2998181.2998364
- [46] Zhang, Mingrui Ray Zhang, Mingyuan Zhong, and Jacob O. Wobbrock. 2022. Ga11y: An automated gif annotation system for visually impaired users. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems, CHI ’22*. New York, NY, USA.

A ETHICS

This work does not raise any ethical issues.

B DATA AND TOOL AVAILABILITY

We have open-sourced both the LangCrUX dataset and Kizuki. Kizuki includes detailed documentation and a README file explaining how to use it and how to extend it with custom

accessibility tests. We have also created an interactive website for LangCrUX, where users can explore the dataset in greater detail, including language distribution across individual websites, with sampling and filtering options.

The dataset and testing tool are available at:

<https://anonymous.4open.science/r/LangCrux-F68F/>

The interactive website is available at:

<https://anonymous.4open.science/w/LangCrux-F68F/>

C DISTRIBUTION OF WEBSITE RANKINGS

Figure 7 displays the distribution of website rankings across different countries in LangCrUX. The rankings are based on CrUX (Chrome User Experience) data. Most countries have the majority of their websites ranked within the top 50,000 rank, indicating high visibility and usage. However, a notable exception is India, where the rankings tend to be significantly lower, often reaching into the 1 million range.

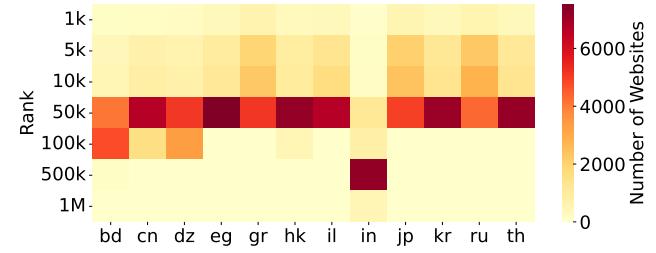


Figure 7: Distribution of website rankings across countries in LangCrUX. Each cell indicates the number of websites within a specific global rank range for a given country. Most countries have websites concentrated in the top 50,000 ranks, while India shows a broader distribution extending toward the 1 million rank range.

D EVALUATION OF ACCESSIBILITY ELEMENTS USING LIGHTHOUSE

Table 3 summarizes the results of our Lighthouse tests for individual accessibility elements discussed in Section §2. To understand how Lighthouse responds to different accessibility conditions, we created isolated test pages, each containing only a single target element. For each element, we evaluated three scenarios: the element being completely missing, present but with an empty value, and present with content in a different language. The table reports whether Lighthouse flagged the element as failing or passing in each scenario.

Table 3: Lighthouse test results for individual accessibility elements under different conditions. A ✓ indicates the test passed, while a ✗ indicates the test failed.

Accessibility Rule	Missing Element	Empty Value	Incorrect Language
button-name	✗	✓	✓
document-title	✓	✗	✓
frame-title	✗	✗	✓
image-alt	✗	✓	✓
input-button-name	✓	✗	✓
input-image-alt	✗	✗	✓
label	✓	✓	✓
link-name	✗	✗	✓
object-alt	✗	✗	✓
select-name	✗	✗	✓
summary-name	✓	✓	✓
svg-img-alt	✓	✓	✓

E EXAMPLES OF EXTREME ACCESSIBILITY TEXT VALUES

Table 4 shows examples of image alt texts that exceed 1000 characters. These values were extracted from real-world web-pages from Bangladesh, India, Japan, Greece, and Thailand, and are shown along with the corresponding source URLs. The examples illustrate cases where accessibility attributes contain unusually long descriptive content, often including entire paragraphs or embedded metadata.

F COUNTRY-LEVEL SCATTER PLOTS

To complement the main analysis, we include country-specific scatter plots, shown in Figure 8, displaying the distribution of websites based on the percentage of native-language content in visible versus accessibility text. Each point represents one site, with the x-axis representing the share of visible content in the native language, and the y-axis representing the share of accessibility metadata in the same language. These plots offer a more detailed view of the language mismatch patterns discussed in Section §3.

As an illustration, consider the scatter plot for Thai websites, shown in Figure 8k. Points near the bottom of the plot represent websites where there is almost no accessibility metadata in Thai, regardless of the language used in the visible content. A dense cluster in the bottom right corner indicates websites where the visible content is almost entirely in Thai, but the accessibility text includes little or no Thai. In contrast, points in the top right corner represent websites where both the visible and accessibility content are predominantly in Thai, indicating consistent language use across both types of content.

G ELEMENT-LEVEL FILTERING ANALYSIS

Figure 9 shows the distribution of uninformative accessibility text by HTML element. Generic action labels are especially common in `<button>` (14.2%) and `<input>` buttons (13.5%),

indicating vague, non-descriptive usage. Single-word labels are the most prevalent issue overall, notably in `<label>` (24.4%), `<image-alt>` (17.1%), and `<select>` (15.3%) elements. These short, generic strings often fail to provide meaningful context. `<summary>` elements show both high generic action (42.9%) and single-word rates (40.5%), highlighting minimal semantic value. These trends suggest a need for deeper evaluation of accessibility text beyond its presence.

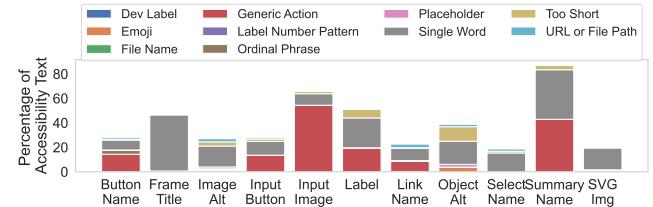


Figure 9: Breakdown of filtered uninformative accessibility text by HTML element and category.

H FILTERING UNINFORMATIVE ACCESSIBILITY TEXT

To assess the informativeness of accessibility text (e.g., alt, aria-label, or label attributes), we apply a rule-based filtering pipeline to discard generic or low-quality entries. Below, we outline each discard category, its rationale, and representative examples.

- **Emoji:** Emoji are discarded because screen readers often fail to interpret them reliably or skip them altogether, making them unsuitable for conveying meaningful accessible content.
Example: [emoji character removed for compatibility]
- **Too Short:** Texts below a language-specific character threshold are considered too short to be useful. For CJK (Chinese, Japanese, Korean) scripts, the limit is 1 character; for others, it is 3 characters.
Example: "go", "图"
- **File Name:** Strings that appear to be image or asset file names (e.g., ending in .jpg, .png, .svg) are removed.
Example: "banner_img123.jpg"
- **URL or File Path:** URLs or file system paths are excluded as they are not meaningful for screen reader users.
Example: "https://example.com/image.png", "/assets/img/logo.svg"
- **Generic Action:** Common UI actions (e.g., "close", "search") in multiple languages are filtered if used alone without context.
Example: "search", "닫기" (Korean for "close")

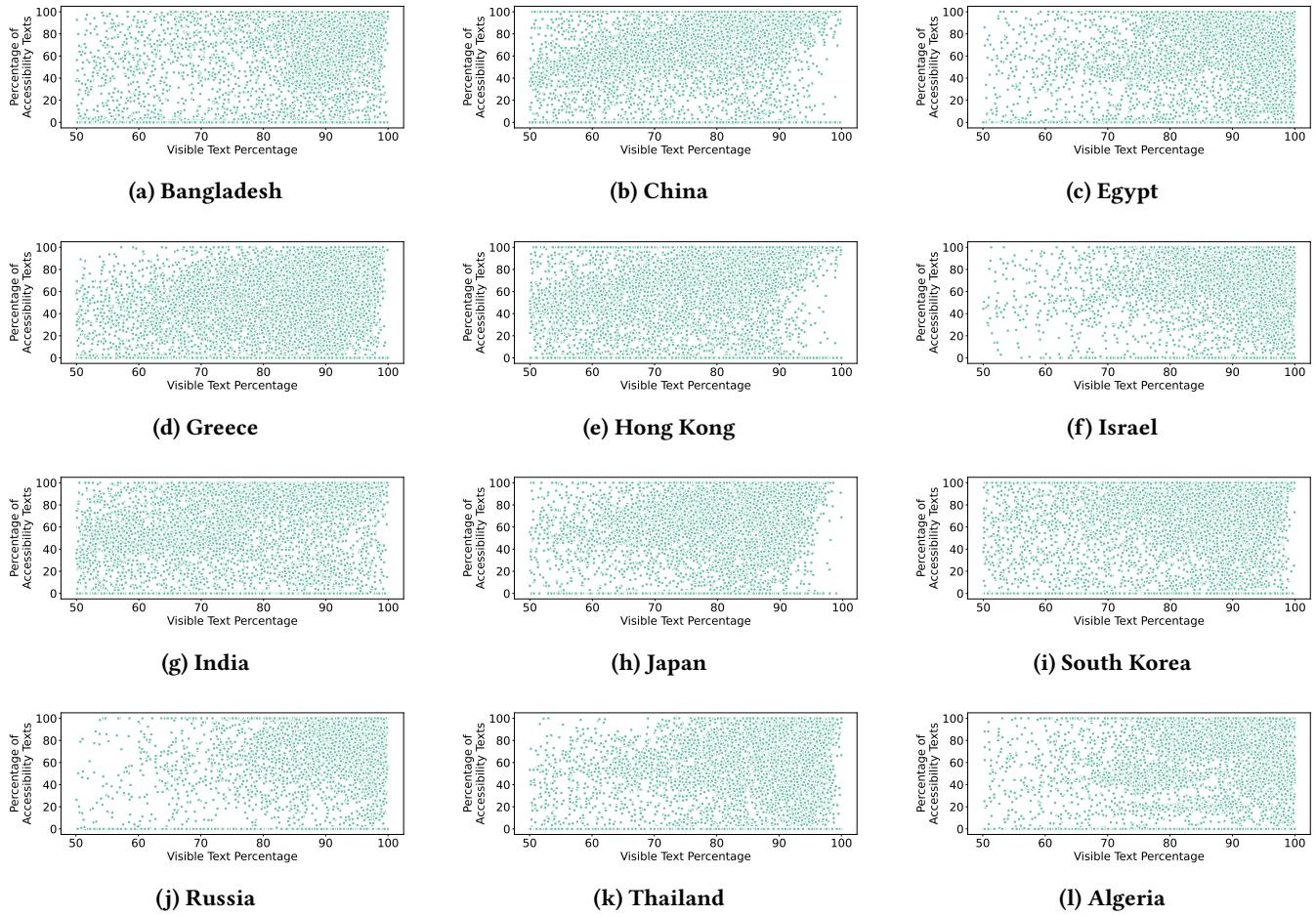


Figure 8: Scatter plot showing the percentage of native language usage in visible text versus accessibility text for websites from 12 different countries. Each point represents a single website. The plots highlight the degree of alignment or mismatch between the languages used in visible content and the corresponding accessibility attributes.

- **Placeholder:** Generic placeholders for images or UI components, such as "image", "icon", or "button", are removed. These include translations in various languages.
Example: "icon", "图像" (Chinese for "image")
- **Dev Label:** Developer-generated IDs or component labels (e.g., "navbar-toggle", "carousel1") are excluded.
Example: "btn-submit", "nav_menu"
- **Label Number Pattern:** Common patterns like "image 1", "button 2", etc., are discarded as they typically lack descriptive value.
Example: "slide 3", "figure 5"
- **Single Word:** For non-CJK scripts, single-word entries are filtered unless they appear to carry descriptive meaning.
Example: "photo", "submit"

- **Mixed Alnum:** Strings with alphanumeric IDs are typically programmatic or internal references and are removed.
Example: "img123", "icon2"
- **Ordinal Phrase:** Numeric phrases like "3 of 5" are common in pagination or sliders and provide limited context.
Example: "2 of 10", "1 of 3"

Text that does not match any of the above patterns is retained and considered **useful** for accessibility analysis. This filtering process helps distinguish meaningful metadata from boilerplate or autogenerated labels, enabling more accurate assessments of multilingual accessibility practices.

I EXAMPLES OF MISMATCH BETWEEN VISIBLE AND ACCESSIBILITY TEXT

Table 5 illustrates examples of accessibility mismatches on six websites across Bangladesh, India, Thailand, Egypt, China, and Hong Kong. Each cell consists of an image from the

website, the corresponding URL, and the alt text or accessibility description associated with that image. The examples highlight cases where website content is presented in the native language, yet the accessibility descriptions, such as alt texts, are written in English. This language inconsistency creates confusion for screen reader users and demonstrates the need for language-aligned accessibility practices.

Table 4: Examples of image alt texts exceeding 1000 characters in length, collected from websites in Bangladesh, India, Thailand, Greece, and Japan. Each row shows the corresponding alt text and the source URL where the alt text was extracted.

 <p>Link: https://developer.nvidia.cn/zh-cn/blog/nvidia-nim-operator-2-0-boasts-ai-deployment-with-nvidia-nemo-microservices-support/</p> <p>Alt text: “The image depicts a stack diagram highlighting NVIDIA NIM Operator, a Kubernetes Operator that is designed to facilitate the deployment, management, and scaling of NVIDIA NIM microservices on Kubernetes clusters.”</p>	 <p>Link: https://goodmoneybygsb.com</p> <p>Alt text: “A hand is holding a smartphone displaying the "GOOD MONEY" app by GSB. The screen shows the app's logo, featuring a white "G" on a gradient background transitioning from pink to orange.”</p>
 <p>Link: https://www.jagonews24.com/photo/bangladesh/news/12915</p> <p>Alt text: “Three people were killed in a head-on collision between an ambulance and an engine-powered van in Jhikargachha, Jessore. Another person was injured in the incident. Photo: Milan Rahman”</p>	 <p>Link: https://www.ajnet.me</p> <p>Alt text: “Vladimir Putin - Khalifa Haftar meeting in Moscow Vladimir Putin - Khalifa Haftar meeting in Moscow- - MOSCOW, RUSSIA - MAY 10: (-EDITORIAL USE ONLY - MANDATORY CREDIT - 'KREMLIN PRESS SERVICE / HANDOUT' - NO MARKETING NO ADVERTISING CAMPAIGNS - DISTRIBUTED AS A SERVICE TO CLIENTS--) President of Russia Vladimir Putin (R) meets with Khalifa Haftar, the leader of the armed forces in the east of the country in Moscow, Russia on May 10, 2025. DATE 11/05/2025 SIZE x Country Rusya SOURCE Anadolu/Kremlin Press Service.”</p>
 <p>Link: https://hindi.cricketaddictor.com/cricket-news/preity-zintas-team-saint-lucia-kings-won-the-1st-title-of-cpl-2024-under-the-captaincy-of-faf-du-plessis-7288602/</p> <p>Alt text: “Preity Zinta's team Saint Lucia Kings won the 1st title of CPL 2024 under the captaincy of Faf du Plessis”</p>	 <p>Link: https://www.aia.com.hk</p> <p>Alt text: “Lovely Asian girl and mother celebrating birthday with friends / family by having a virtual birthday party at home. Enjoying birthday celebration in front of the laptop during pandemic. Practicing social distancing. Birthday lifestyle theme.”</p>

Table 5: Illustrative examples of accessibility mismatches on six websites across Bangladesh, India, Thailand, Egypt, China, and Hong Kong. Each cell consists of the image from the website, the corresponding URL, and the alt text or accessibility description associated with that image. The examples highlight cases where website content is presented in the native language, yet the accessibility descriptions, such as alt texts, are written in English.