# Automatically Detecting Online Deceptive Patterns in Real-time

**Asmit Nayak, Shirley Zhang**, **Yash Wani**, **Rishabh Khandelwal, Kassem Fawaz**

University of Wisconsin – Madison

## Abstract

Deceptive patterns (DPs) in digital interfaces manipulate users into making unintended decisions, exploiting cognitive biases and psychological vulnerabilities. These patterns have become ubiquitous across various digital platforms. While efforts to mitigate DPs have emerged from legal and technical perspectives, a significant gap in usable solutions that empower users to identify and make informed decisions about DPs in real-time remains. In this work, we introduce *AutoBot*, an automated, deceptive pattern detector that analyzes websites' visual appearances using machine learning techniques to identify and notify users of DPs in real-time. *AutoBot* employs a two-staged pipeline that processes website screenshots, identifying interactable elements and extracting textual features without relying on HTML structure. By leveraging a custom language model, *AutoBot* understands the context surrounding these elements to determine the presence of deceptive patterns. We implement *AutoBot* as a lightweight Chrome browser extension that performs all analyses locally, minimizing latency and preserving user privacy. Through extensive evaluation, we demonstrate *AutoBot*'s effectiveness in enhancing users' ability to navigate digital environments safely while providing a valuable tool for regulators to assess and enforce compliance with DP regulations.

## 1 Introduction

Deceptive patterns (DPs) are design elements that manipulate users into making unintended decisions while interacting with interfaces on applications or websites. These patterns exploit cognitive biases and psychological vulnerabilities to influence user behavior, often in ways that benefit the service provider at the user's expense. They have become increasingly prevalent across digital environments, significantly impacting user experiences on social media, mobile devices, cookie consent banners, and even gaming applications [33]. Prior works have shown that DPs can result in financial loss [42], privacy breaches [8], or exploitation of vulnerable groups, including children [46]. A typical example of a DP is the intentionally convoluted subscription cancellation process, where users must navigate through multiple obscure options to terminate a service, as seen on platforms like www.dailymail.co.uk.

Efforts to mitigate the effect of DPs have emerged from the policy and the technical sides. Regarding policy, regulations like CPRA [15] and the GDPR [4] have released guidance on deceptive patterns. On the technical front, researchers have explored classifying the existing types of deceptive patterns [11, 21, 33, 48, 44, 7, 16, 39, 32], assessing their effectiveness [30], and compiling different deceptive pattern cases [12]. Despite these ongoing efforts, deceptive patterns continue to pose significant challenges for both users and regulators. Users frequently encounter DPs in their digital interactions. For instance, sponsored advertisements are strategically placed at the top of search results, creating a false impression that they are the most relevant or popular items. Furthermore, regulators often lack the tools to assess and enforce compliance with regulations concerning deceptive patterns effectively. This technological gap leaves users vulnerable to sophisticated DPs.

In this work, we propose a new paradigm to bridge this gap and improve the usability of online websites by developing a solution to detect deceptive patterns and automatically warn users in real time. Our solution makes users aware of deceptive patterns as they browse and enables regulators to assess and evaluate online services for deceptive

---

*Equal Contribution

patterns. Achieving this objective requires us to overcome several challenges. First, the lack of a standardized taxonomy and usable datasets complicates classification efforts. Second, the vast diversity of web designs and the ingenuity of designers in creating new deceptive techniques have led to a trade-off between scalability and accuracy in detection methods. For instance, existing solutions like disabling third-party cookies or using filter lists have proven inadequate in preventing user tracking, as highlighted by Chen et al. [16]. Third, user-friendly solutions must operate with minimal infrastructure to be practical for widespread adoption.

We first create a standardized taxonomy of deceptive patterns to realize our objectives, building upon existing taxonomies. Leveraging this standardized taxonomy, we then built *AutoBot*, an automated, deceptive pattern detector that analyzes the website's visual appearances, uses machine learning tools to identify deceptive patterns, and brings deceptive patterns to users' attention in real-time. Specifically, *AutoBot* relies on an invariant behavior of deceptive patterns: the visual representation and how users perceive them. Using this invariant behavior, *AutoBot* applies a two-staged pipeline that identifies the deceptive patterns present in the website's screenshot.

First, *AutoBot* analyzes the website screenshot, identifying interactable elements using visual analysis and extracting textual features without relying on the HTML structure. We note that here, we rely on the insight that deceptive patterns work by manipulating what the users perceive; thereby, analyzing their visual appearance provides a comprehensive signal to be able to identify them. Next, *AutoBot* leverages a custom language model to understand the context surrounding these elements, including color, font size, and text, to determine the presence of deceptive patterns. We created an annotated dataset for deceptive patterns and designed a new training paradigm to fine-tune language models. We also note here that empirically, we have found that off-the-shelf language models like Gemini and ChatGPT often fail to detect deceptive patterns from screenshots, as shown in Figure 2.

We showcase the usability of *AutoBot* by building a Chrome browser extension that (1) detects the deceptive patterns in real-time by performing the complete analysis locally and (2) annotates the current webpage to warn the users about the deceptive patterns. We also characterize the performance of *AutoBot* on several devices and find that the latency is less than one second on modern devices, which minimally impacts the usability of the websites.

In this work, we put forth *AutoBot*, a automated solution to detect deceptive patterns on the web and make the following contribution:

1. We compile an annotated large-scale deceptive design dataset (D3) based on carefully curated taxonomy. We also introduce a novel pipeline to generate a diverse dataset of synthetically generated websites.
2. We show that our system, *AutoBot*, can detect deceptive designs present on the web and classify them per our taxonomy, with high recall resulting in most deceptive patterns being detected.
3. We leverage recent advances in optimizing LLM inference on consumer hardware to run our system, *AutoBot*, on device with real-time inference.

## 2 Background and Related Works

Deceptive patterns refer to interface design choices that manipulate or deceive users into making decisions they might not otherwise make [11]. In this paper, we focus on deceptive patterns found on websites. Examples of such patterns include hidden costs, forced continuity, misdirection in e-commerce websites, and privacy-invasive default settings in social media platforms. In the rest of this section, we survey existing works on detecting deceptive patterns on the web and present a taxonomy to categorize them.

### 2.1 Detecting Deceptive Patterns

Given their effect on users' online privacy, security, and safety, researchers developed several mechanisms to detect and measure the prevalence of online deceptive patterns [11, 21, 33, 48, 44, 7, 16, 39, 32].

#### 2.1.1 Measurements of Deceptive Patterns

Early efforts to identify online deceptive patterns relied on manual exploration. Brignull et al. [11] manually explored the web to compile a "Hall of Shame:" a list of websites with deceptive patterns. Gray et al. [21] expanded this corpus
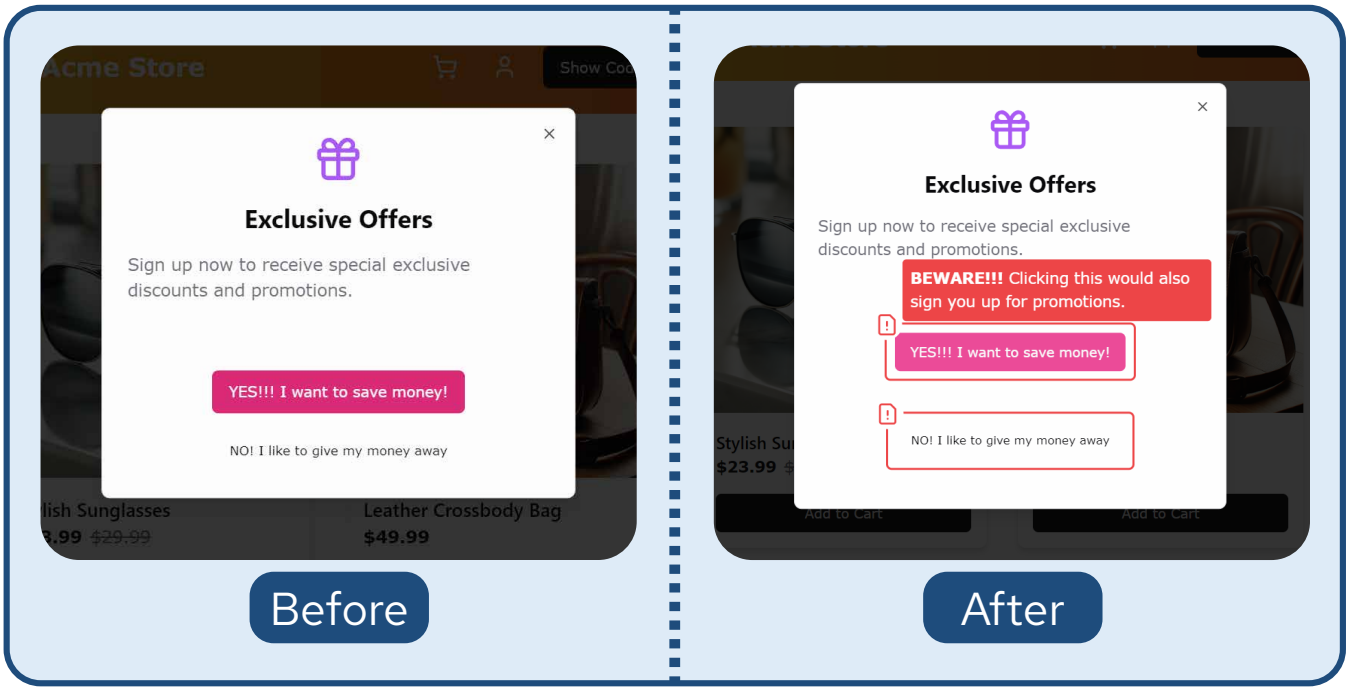
**Figure 1:** Here we show how *AutoBot* highlights and informs the users about deceptive patterns on a site.

by performing keyword searches on social media for posts highlighting websites with deceptive patterns, which were then manually validated. While highly accurate, this methodology is limited to domain experts and lacks scalability due to manual validation requirements.

To expedite the manual exploration process, Mathur et al. [33] proposed a clustering-based pipeline to group similar websites based on their text content, which is then manually inspected. Although this process accelerates data collection, it still lacks scalability. Attempts to automate the manual inspection process include Tung et al.'s naive Bayes classifier [48], Soe et al.'s gradient-boosted tree to flag texts showcasing deceptive patterns [44], and Adorna et al.'s combined naive Bayes classifier and VGG-19 model to identify deceptive designs in cookie notices [7]. However, these works are often domain-specific and struggle to generalize to new domains. Additionally, most works rely solely on text-based classifiers or rules, limiting their ability to detect deceptive patterns, such as those based on colors or trick wordings.

Recent works by Chen et al., Mansur et al., and Raju et al. [16, 32, 39] focus on automated detection of deceptive patterns in mobile apps. However, Chen et al. [16] and Mansur et al. [32] used predefined heuristics, limiting detection to simpler deceptive patterns like disguised ads. Raju et al. [39] employed rule-based source code analysis to detect patterns like ads, forced action, and nagging [39].

### 2.1.2 User facing Solutions

While most previous works in the field of deceptive patterns focused on creating corpora of websites and apps with deceptive patterns, few have addressed user-facing solutions. To the best of our knowledge, only two works have focused on this aspect: UIGuard by Chen et al. [16] and Ariadne by Adorna et al.[7]. However, both solutions are limited by their text-only classifiers and are confined to specific domains - mobile apps for UIGuard and cookie notices for Ariadne. In the broader domain of Privacy, Safety, and Security (PSS), researchers have developed tools to help users navigate certain deceptive designs on websites. For instance, works by Khandelwal et al. [25, 26] enable users to find and adjust privacy settings and disable non-essential cookies automatically. Similarly, OptOutEasy by Kumar et al. [9] automatically finds the opt-out links from privacy policies and surfaces them to the users.

In this work, we propose *AutoBot*, a automated deceptive pattern detector that integrates both the visual signal and the text signal to identify deceptive patterns. This allows *AutoBot* to make users aware of the deceptive patterns as they are browsing without impacting the usability of the websites.

### 2.1.3 UI Element Detection

Detecting deceptive patterns requires gathering contextual clues from web pages, necessitating a Web Element Detector to identify buttons, checkboxes, radio buttons, and switches [45]. Existing works like the Element Detector presented by Wu et al. [50]. do not detect elements such as checkboxes and switches. Detectors like UISketch [43] and Roboflow-Universe [3] perform poorly on real-world data. UISketch, designed to detect UI elements on the web, was primarily trained on hand-drawn images to identify UI elements based on their sketches. This training approach limits the model's applicability to real-world web elements. Testing on ten popular websites revealed that the model detector failed to identify the most required elements.

The ScreenRecognition model presented by Zhang et al. [52] shows promise for this task but was trained on mobile screenshots and is closed-source, preventing its use or testing for web generalization. The dataset is also private, precluding its use for training our detector.

UEID, a tool developed for GUI Element Detection by Xing et al. [51], performs well by merging text and visual detectors. However, UEID with an average runtime of 4.8s will not fit into a user-friendly solution on the browser level. The RICO dataset [29] they used contains mobile UIs and thus could not be used to train our detector for webpages.

In this work, we build a real-time UI element detector using `YOLOv10`. In prior work, we also tackled the dataset limitation by generating a synthetic dataset of 62K websites.

### 2.1.4 Limitations

Existing techniques for automated detection of deceptive patterns are often limited to tasks solved via rule-based heuristic methods [18, 16, 39]. Furthermore, these works usually extract text and icons but do not extract any information regarding the color of the text or if the text belongs to a button. Due to the lack of this information, these approaches can overlook deceptive patterns like nudge. *AutoBot* addresses these limitations by building on a more comprehensive taxonomy and leveraging recent advancements in Large Language Models (LLMs) to automatically detect a wide range of deceptive patterns on the web.

## 2.2 Taxonomy for Deceptive Patterns

There have been several attempts at categorizing online deceptive patterns [12, 17, 10, 21, 33, 28, 32, 18]. Below, we discuss the existing taxonomies and build upon them to propose a more comprehensive and standardized taxonomy that enables *AutoBot* to identify deceptive patterns on websites.

**Existing Taxonomies:** Brignull et al. [12] developed the first taxonomy specifying different types of deceptive patterns in 2010. Conti et al. [17] expanded on this taxonomy, including malicious interface designs. Bösch et al. [10] introduced a similar taxonomy called "privacy dark patterns", which included more privacy-centric categories such as 'Forced Registration' and 'Hidden Legalese Stipulations.'

Gray et al. 2018 created a unified corpus to detect deceptive designs in user interfaces, building on previous taxonomies and categories [21]. Since then, various works have further adapted the taxonomy for specific domains. Lewis et al. [28] codified dark patterns for mobile apps and games, while Mathur et al. [33] adapted the taxonomy to focus on deceptive patterns present in shopping websites. Mansur et al. [32] extended the taxonomy to detect deceptive patterns in mobile and web apps. More recently, Curley et al. [18] created a taxonomy based on how dark patterns can be detected in manual, automated, or semi-automated fashion.

**Limitations:** The problem of detecting deceptive patterns requires a taxonomy that comprehensively covers the related practices. Existing taxonomies, however, overlap in some patterns, do not describe the patterns consistently, and often miss some patterns described in other taxonomies. For instance, the taxonomy of Chen et al. [16] lacks the categorization for nudge, confirmshaming and jargon, and the taxonomy by Curley et al. [18] does not include categorization for nudge, fake-scarcity/fake-urgency, pre-selection, and disguised ads.

**Proposed Taxonomy:** We address this problem by establishing a consistent and comprehensive taxonomy of detectable deceptive patterns. We instantiated the new combined taxonomy with the Brignull et al. [12] taxonomy from 2010. We expanded the combined taxonomy with the taxonomies from Conti et al. [17], Bösch et al. [10],

| Category | Sub-Type | Description | Examples |
|---|---|---|---|
| **Interface Interference** | Confirmshaming | Guilt-tripping users into making a specific choice | "No, I prefer to pay more" |
| | Fake-Scarcity/ Fake-Urgency | Create a false sense of urgency/scarcity to pressure users into making a choice | "Only 3 left in stock" |
| | Nudge | Nudge a user towards a specific choice. | "Accept All" in bright colors, while "Reject" is hard to notice |
| | Hard to Cancel | Design choices that visually interfere to make cancellation hard for users. | The unsubscribe button is tiny and hard to click. |
| **Forced-Action** | Forced-Action | Design tactic forcing users to complete a specific task to proceed. | Pop-up ads |
| **Obstruction** | Hard to Cancel | Making the cancellation process unnecessarily complicated | User need to call customer service to cancel membership. |
| | Pre-Selection | Choices given to users are already selected. | The checkbox of "Sign up for news and updates" is checked by default. |
| | Visual Interference | Misleading design elements that distract or mislead users from important information. | The term of use for service users are signing up for is in tiny font and cannot be clearly seen on sites. |
| | Jargon | The use of non-user-friendly language to prevent users from understanding important information. | "By affirming this selection, you consent to the perpetuation of automatic pecuniary transactions at designated intervals." |
| **Sneaking** | Hidden Subscription | Users are not clearly informed that they are signing up for a service. | "By signing up for this email, you are agreeing to news and information from us" |
| | Hidden Costs | Users are not clearly informed about all the costs associated with a service. | During checkout, unexpected fees such as "handling charges"appear. |
| | Disguised Ads | Visually misleading ads embedded into page content. | Prominent "Download" button at the top of a page that redirects the user to unrelated adware. |
| | Trickwording | The use of non-user-friendly language to trick users into making certain choices | "Newsletter subscription by default, tick here to unsubscribe" |
| **Non-Deceptive** | | Common, user-friendly design element that does not show any deceptive pattern. | Any text or web element that does not exhibit any deceptive behavior. |

**Table 1:** Curated taxonomy of Deceptive Patterns. The table shows the Category, the Sub-Type, the description of the deceptive pattern, and an example.

Gray et al. [21], Mathur et al. [33], Curley et al. [18], Mansur et al. [32], and Chen et al. [16]. In particular, we incrementally added each pattern from existing taxonomies to the combined taxonomy. If the pattern already existed in our taxonomy, we consolidated the terms by deferring to the later work. Otherwise, we added the pattern to the taxonomy.

The outcome is a taxonomy that combines all patterns in existing taxonomies and uses standardized language to

**Figure 2:** Dialog with Gemini and GPT-4V, showing they cannot identify correct deceptive patterns from a screenshot.

describe them. Table 1 shows the final taxonomy, comprising four high-level categories and 13 low-level categories. *AutoBot* utilizes a language model to label UI elements from a webpage using this taxonomy. We note here that we only consider the patterns that can be detected using a screenshot. For instance, *nagging* is a pattern where a banner is repeatedly shown to the user, even after they close it. Detection of such patterns requires temporal analysis of the webpages, which is out of scope for this work.

# 3   System Overview

We propose *AutoBot* as an user-facing solution to alert users about online deceptive patterns. In essence, *AutoBot* is a browser extension that highlights UI elements that include a deceptive pattern on a webpage the user is visiting. It achieves its objectives in real-time by running our classification models on-device, enhancing user privacy and simplifying deployment. The underlying design approach of *AutoBot* involves analyzing the rendered webpage, overcoming challenges related to non-standard and widely varying implementations. *AutoBot* analyzes a screenshot of the webpage to identify and highlight deceptive patterns as per the taxonomy from Table 1.

Recent advances in multimodal foundational models offer a venue for analyzing screenshots and highlighting patterns with proper prompting. To this end, we empirically evaluated how accurately the state-of-the-art models, including GPT4-V [35] and Gemini 1.5 [20], detect online deceptive patterns. Our experiments showed that while these models can detect deceptive patterns, they often hallucinate and give false positive answers. For example, Figure 2 shows how Gemini and GPT-4V struggle to identify the deceptive patterns in screenshots. While Gemini 1.5 hallucinates about the lack of "More Options", which does exist in the screenshot, GPT-4V falsely identifies the "Apply" button as a nudge because it is brightly colored compared to its surroundings. Another significant shortcoming in current multimodal models is the challenges in deploying them. These models require developers to utilize external (non-private) API services or employ large, slow, computationally expensive models locally.
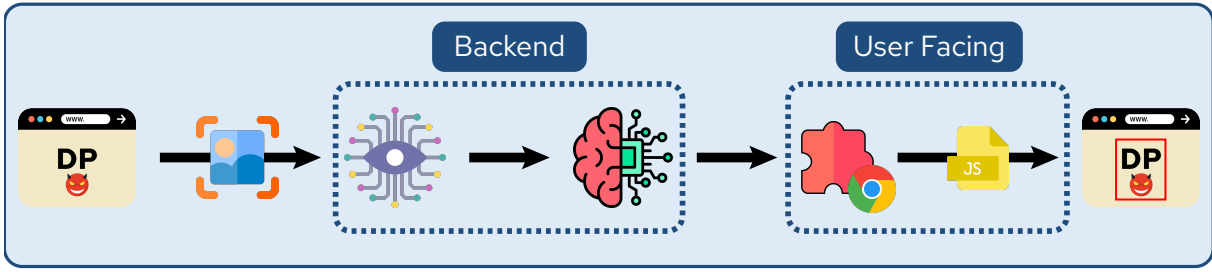
**Figure 3:** An overview of *AutoBot*'s components. *Backend* generates the list of deceptive patterns in the website, and *User Facing* highlights them for the user.

*AutoBot* addresses these shortcomings through a pipelined multi-modal approach to better understand a website. Figure 3 shows the outline of the high-level system overview of *AutoBot*. In particular, *AutoBot* addresses three challenges. First, as we aim to provide service in browser extension, the model size must be small to run locally and provide real-time response on the user's browser. Additionally, while minimizing the model's size, *AutoBot* needs to provide comparatively accurate evaluations of the webpage. Finally, *AutoBot* should perform its task without affecting the website's usability.

To overcome the challenges, *AutoBot* interleaves advanced vision models with recent advances in language models to understand the visual layout of a page and reason about the text and UI elements.

**Vision Module:** The *Vision Module* (Section 4) takes the screenshot of the webpage as input from the *User Facing Interface* and performs OCR and web-element detection on it to create a tabular representation of the website including features such as text content, bounding box, font color, and font size.

**On-device Language Model:** Using the tabular representation provided by *Vision Module*, the *Language Module* (Section 5) identifies deceptive designs on the website by cohesively interpreting the interconnected elements within the same file.

**User-Facing Component:** The user-facing component consists of the user interface for the browser extension (Section 6), *AutoBot*, that uses the output from the *Language Module* to automatically highlight all the deceptive patterns on the site.

## 4 Vision Module

Prior works analyzing websites [25, 26] have primarily relied on HTML analysis. This method faces significant challenges due to the dynamic nature of websites. Websites increasingly employ *JavaScript* frameworks that modify the Document Object Model (DOM) on the fly, rendering static HTML analysis insufficient. Furthermore, the diversity in the coding practices and obfuscation techniques provide additional challenges in HTML-based analyses.

To overcome these challenges, we adopt a novel approach focusing on the invariant aspect of websites: *the user experience*. By leveraging the visual signals, our methodology relies on how the users perceive and interact with websites. This approach offers several advantages: (1) It is resilient to change in underlying technologies as it captures the actual rendered content. (2) It allows us to analyze the same information that the user encounters, providing a more accurate representation of the potential deceptive patterns.

The Vision Module is the foundation of our analysis pipeline, performing three key tasks. First, as it receives the screenshot from the user-facing component (See Section 6), it performs an Optical Character Recognition (OCR) operation on the image, extracting the text and generating a CSV file consisting of the text and their bounding box location and their associated features such as position, font size, and color. We note that color and font can be crucial in determining deceptive patterns, like nudging. Next, the *Web-UI Element Detector* generates a list of the detected web elements and their corresponding bounding boxes. Finally, the *Vision Module* merges the two outputs to create the comprehensive textual representation of the website, referred to as the *TextMap*. An example of the process is in Figure 4.
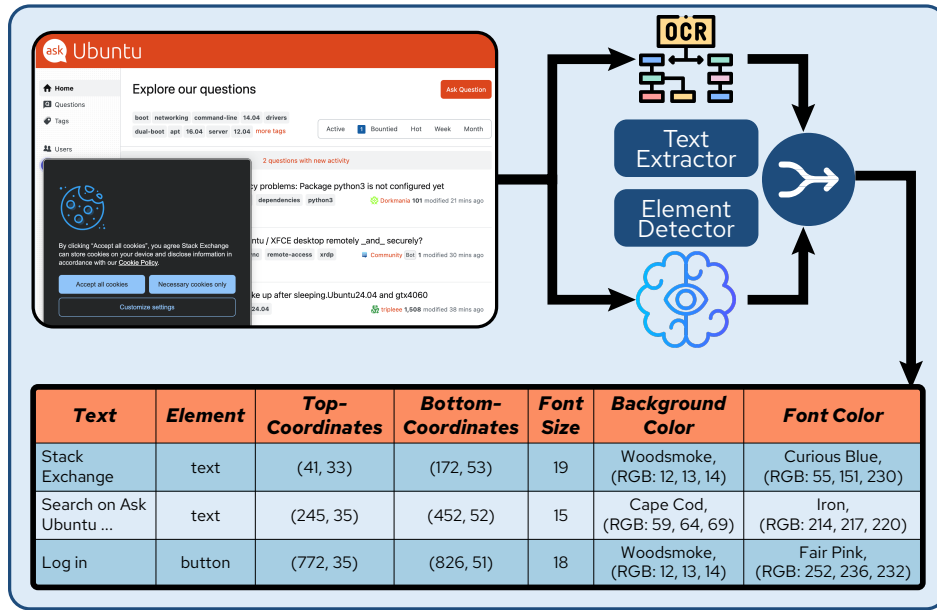
**Figure 4:** Overview of flow of the *Vision Module*. The screenshot is processed by the Text Extractor and the Element Detector, whose outputs are then merged to create the TextMap.

## 4.1 Text Extraction

The *Text Extractor* module starts by performing Optical Character Recognition (OCR) on the website screenshot. For this task, we employ the Google Vision OCR API [19] as it has high accuracy in extracting text from images of varying scales and resolutions. The API returns detected texts that our OCR module concatenates based on proximity to form coherent text blocks. These blocks are a list of bounding boxes with their respective text content. Leveraging these bounding boxes, we derive additional textual features, including font size, font color, and background color, as illustrated in Figure 4. We calculate font size by subtracting the bottom y-coordinate from the top y-coordinate of the bounding box. To determine color information, we utilize the *extcolors* [14] package, extracting the most prominent color (background) and the second most prominent color (font).

This approach addresses a key limitation identified by Soe et al. [44] in the previous ML methods: the lack of representation of the UI richness that a user perceives, such as text placement and contrast between the font and background colors. Our *Text Extractor* module captures these detailed text features, along with the relative location of text elements in the website, providing a comprehensive representation of the textual content experienced by the users.

To further contextualize the extracted text, we employ the *Web-UI Element Detector* to identify the web elements located within each text block. This provides additional information about the role and significance of each textual element within the larger structure of the webpage.

## 4.2 Web-UI Element Detector

In order to allow the *Language Module* to understand the context better, we trained a *Web-UI Element Detector* to identify UI elements like buttons, unchecked/checked checkboxes, unchecked/checked switches, and unchecked/checked radio buttons. This provides context to the extracted text and enables a more comprehensive understanding of the webpage's structure. For instance, distinguishing between a clickable button and a static text block can be significant, as a seemingly simple text might be a deceptive call to action when recognized as a button. Similarly, identified checkbox states (checked or unchecked) can reveal pre-selected options that users might overlook.
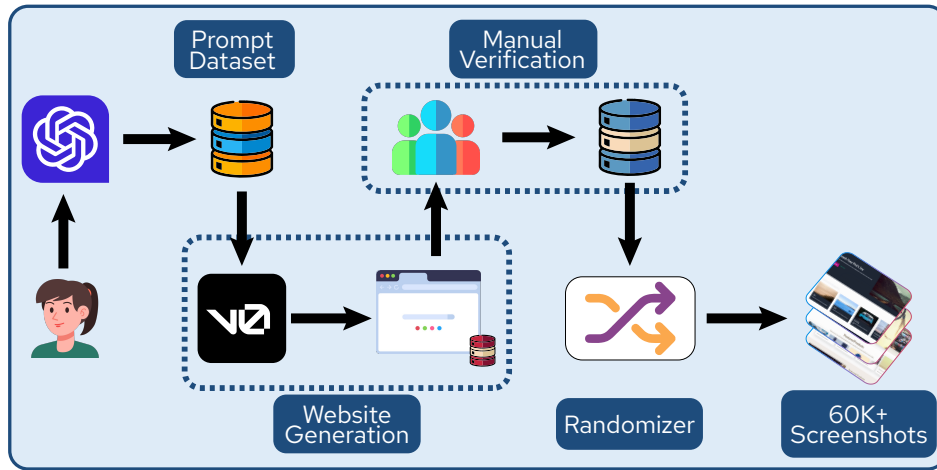
**Figure 5:** Pipeline of Generating Web-UI Element Dataset.

### 4.2.1 Choosing YOLO over CNNs

Unlike previous works [25, 32, 16] that used Convolution Neural Networks (CNNs) to detect and classify elements and icons in the screenshots, we opted for using You Only Look Once (YOLO) model, specifically the `YOLOv10` model. While CNNs are effective for object detection tasks, YOLO models offer multiple advantages over CNNs that make them a perfect fit for our goals.

**Real-Time Detection** YOLO models are designed for real-time object detection tasks, which is the primary requirement for our work. Unlike CNN-based object detection models involving complex multi-stage feature extraction and classifications processes, YOLO models perform a single pass to extract the bounding boxes and class predictions. Previous YOLO models would generate multiple class predictions with similar bounding boxes, requiring the use of the non-maximal suppression (NMS) post-processing step; however, `YOLOv10`'s efficient and optimized architecture can bypass the NMS stage while maintaining similar accuracy, significantly reducing the computational overhead [47].

**Balancing Efficiency with Accuracy** While CNN-based detectors, like Faster R-CNN, provide predictions with high accuracy, they are computationally intensive and require considerable computational power and time [40, 41]. YOLO models, however, are optimized for speed and real-time inference, striking a balance between accuracy and efficiency [47]. For our task of locally running a model to detect Web-UI elements with constrained resources, we opted for a smaller model with slightly lower accuracy over a resource-intensive model.

**Model Size** The `YOLOv10` model is comparatively lightweight, around 40MB, allowing for a seamless local in-browser deployment for users without adversely impacting their resources.

### 4.2.2 Dataset Creation

We develop a novel approach to create a diverse and representative dataset to address the lack of suitable Web-UI element dataset for training our *Web-UI Element Detector*. Instead of relying on manual scraping and labeling [52, 43, 51], we leverage AI-powered tools to generate a synthetic dataset of over 62K website interfaces that reflect the current web landscape. A key advantage of our synthetic approach is its precise control over element positioning and labeling.

Our dataset generation process centered around `v0`[1], an AI-powered generative user interface system by Vercel[2], in conjunction with Astro[3], a JavaScript framework. To scale the generation of diverse prompts for `v0`, we employ `GPT-4` and create a wide range of prompts. We then submit these prompts to `v0` and collect the three designs generated for each prompt. We specifically instruct GPT-4 to generate prompts that include the UI elements. Refer

---

[1]https://v0.dev/
[2]https://vercel.com/
[3]https://astro.build

| Class | Precision | Recall | F1-Score | # Elements |
|-------|-----------|--------|----------|------------|
| **button** | 0.942 | 0.882 | 0.911 | 5226 |
| **checked checkbox** | 0.845 | 0.947 | 0.893 | 113 |
| **unchecked checkbox** | 0.979 | 0.76 | 0.856 | 246 |
| **checked radio button** | 0.854 | 0.934 | 0.892 | 76 |
| **unchecked radio button** | 0.859 | 0.886 | 0.872 | 132 |
| **checked switch** | 0.959 | 0.981 | 0.97 | 52 |
| **unchecked switch** | 0.914 | 0.942 | 0.928 | 34 |
| **Total** | 0.907 | 0.905 | 0.906 | 5879 |

**Table 2:** Performance of the visual classifier on the test dataset

to Appendix A.1 Since shadcn[4] components, used by v0, are inherently customizable. We were able to leverage this feature to track the exact location, size and type of each UI element. This level of control allowed us to generate accurate bounding box annotations without manual labeling.

To ensure diversity in our dataset, we also incorporated components from six popular UI libraries, including Material UI [5] and Bootstrap [6] [1, 2, 37]. This integration, facilitated by Astro's ability to combine components from various JavaScript libraries, ensured our dataset encompassed a wide range of UI design patterns and styles. Additionally, we randomized images across all webpages using the Unsplash API [7] and Picsum [8], adding visual diversity to our synthetic websites.

**Manual Verification:** 3 Authors manually verified all generations made by v0 to ensure the pages were rendered and had content. With a failure rate of less than 2% this process was relatively fast as any major errors were caught by the compiler and fixed by the authors.

As the world of web design evolves, this pipeline can still be used to generate datasets that represent the design language of the web by simply adding popular UI libraries to the `Astro` framework and generating newer web layouts using a generative model like `v0`.

### 4.2.3  Training and Performance of YOLOv10

Using the above methodology, we generated over 60K screenshots along with their ground truth labels. These websites were then randomly split into training sets consisting of 85% of the images and the remaining 15% in the testing set.

During training, we observed that the `YOLOv10` model was able to specialize in detecting a subset of UI Elements like a button and a checkbox. Hence, in order to boost the performance, we adopted a step-by-step approach in training three different `YOLOv10` models, each specializing in detecting 2-3 elements each, which we combined to make an ensemble of models. This ensemble of models achieves near-perfect accuracy on the seven classes on the testing dataset.

To measure the accuracy of the ensemble of models on real-life websites, we evaluated the performance of the models on the deceptive pattern websites dataset curated by Mathur et al. [33]. One of the authors annotated over 1.3K website screenshots, which were then used to evaluate the models' performance. Table 2 shows the models' F1-scores for each of the 7 classes. We note that to compute the precision and recall, we used an IoU value of 0.5 and a confidence threshold of 0.3. We chose a lower IoU value because while training was done on a synthetic dataset, the evaluation bounding boxes are human-annotated and will not perfectly match with the predicted bounding box, and a high IoU score would negatively affect model performance.

---

[4]https://ui.shadcn.com/

[5]https://mui.com/material-ui/

[6]https://getbootstrap.com/

[7]https://unsplash.com

[8]https://picsum.photos

# 5   Language Module

The language module identifies the deceptive patterns by describing the UI elements on a webpage. In particular, it takes as input the type of each UI element, its coordinates, its state, and the surrounding text.

## 5.1   Overview

Previous approaches, such as manual analysis, traditional classifiers, and heuristic or rule-based methods, as discussed in Section 2, can be limited in locally identifying deceptive patterns accurately over diverse webpages. Furthermore, as discussed earlier in Section 3, using out-of-the-box LLMs and Vision Language Models (VLMs) exhibit different shortcomings. *AutoBot* addresses these shortcomings in the following ways:

**High False Positive Rate:** We empirically observed that LLMs like Gemini [20] and GPT4 [35], when prompted to detect deceptive patterns on screenshots of websites, had a high false positive rate and would often miss the important content on the site Figure 2. Additionally, these LLMs would frequently classify the same deceptive design with different names due to multiple taxonomies with diverse definitions [21]. To solve this issue, we created our taxonomy as stated in Section 2.2. With this taxonomy's help, we could ground the LLMs to persistently classify the deceptive designs with set labels. However, the LLMs would often miss important parts of the website screenshot, leading to misclassification and, ultimately, a high false positive rate.

*Solution:* While LLMs can solve specific tasks, they often fail when planning and solving complex tasks [24]. Additionally, through empirical experimentation, we observed that LLMs perform poorly in finding deceptive designs when given a screenshot. To properly leverage the knowledge in LLMs, we devised a setup wherein we provide a textual representation of the site (Section 4) as a CSV file, along with a finetuned prompt to allow the LLM to understand how texts on the site are presented to the user.

**Slow Response Time:** Another drawback of using LLMs is their slow response time. LLMs normally consist of trillions of parameters [5], which lowers the inference time.

*Solution* To achieve real-time results, we focused on distilling the knowledge to smaller LLMs like T5. We detail our entire methodology later in this section.

**Localization Issues** LLMs and VLMs, when given a screenshot, cannot provide the bounding boxes for the deceptive design, even if they can detect them correctly.

*Solution* To accomplish this goal, we incorporate the bounding boxes in the textual representation of the site sent to the LLM. Since each row of the CSV file represents a unique element in the site, we can quickly look up that element's position when the LLM detects a deceptive design.

**Data Privacy Concerns** One concern with using external LLM services is that once the data leaves the device, it is unknown what happens. The API service has complete access to the queried content and the querying party's IP address. This potentially allows these third-party services to track users as they browse through the web.

*Solution* The best possible solution is to confine the end-to-end analysis to a user's local browsers. Our work proposes a solution that can be deployed locally in the browser with real-time inference.

Tackling these four challenges, we designed a language model that runs completely on-device and can provide real-time results. Given recent advancements in compressing and running optimized ML models on the browser, we were able to run a T5 [38] base model (hereby, referred to as T5) in the browser using the `Transformer.js` [23] package. However, we cannot use the T5 model out of the box due to its knowledge limitation. To overcome this limitation, we created our deceptive design dataset to finetune T5 further. We followed a similar methodology as that by Hsieh et al. [22], where they extracted rationales from LLMs and then used them as an additional subtask when finetuning smaller models.

## 5.2   Dataset Creation

To create the dataset to train and evaluate T5, we crawled the sites in the deceptive pattern dataset from Mathur et al. [33]. Out of the 1400 sites mentioned, only 555 were still online. We visited these websites and captured their

| Category | Subtype | # Samples |
|---|---|---|
| non-deceptive | not-applicable | 374 |
| forced-action | forced-action | 116 |
| interface-interference | nudge | 102 |
| | fake-scarcity-fake-urgency | 101 |
| | confirmshaming | 16 |
| obstruction | visual-interference | 4 |
| | pre-selection | 10 |
| sneaking | trick-wording | 149 |
| | hidden-costs | 26 |
| | hidden-subscription | 88 |
| | disguised-ads | 2 |

**Table 3:** Distribution of Samples in D3 Dataset.

screenshots. We then passed these screenshots through the *Vision Module* (Section 4) to retrieve their TextMaps in the form of a CSV file as shown in Figure 4.

Prior works have shown that LLMs can generate high-quality reasoning steps when carefully prompted [49]. As detailed by Hsieh et al. [22], extracting rationales from LLMs allows us to leverage an LLM's ability to reason to train smaller models in a data-efficient manner [22]. Based on this characteristic of LLMs, we employ an LLM-assisted human labeling methodology to create the ground truth dataset.

For each site, we provide the taxonomy in Section 2.2 and use chain-of-thought (CoT) prompting [49] and few-shot prompting methods on Gemini 1.5 Flash model [20] to generate three additional columns: Deceptive Category, Deceptive Subtype, and Reasoning, for each row of the TextMap. Using prompting techniques from previous works that encourage the LLM to reason about its choices [34], we were able to not only generate proper labels and the associated reasoning but also minimize hallucinations, as seen earlier when providing just images to LLMs. These classifications were then verified by two authors who are experts in this field. The authors also annotated a standard set of 20 websites and exhibited a near-perfect agreement on classifying deceptive patterns ($\kappa = 0.9$) [27]. We manually collected and annotated over 13K samples from over 500 websites to create our golden Deceptive Design Dataset (D3 dataset).

However, in the D3 dataset, more than 95% (12K+) of the samples were non-deceptive, with only about 600 being deceptive. We randomly sampled 370 non-deceptive examples to balance the dataset and removed the rest. The complete distribution of the samples is shown in Table 3.

## 5.3 Distilling T5

To distill the knowledge from LLMs into smaller models like T5, we used the paradigm suggested by Hsieh et al. [22]. Here, we formally define the training dataset, $D_{\text{train}}$, as:

$$d_i \in D_{\text{train}} \subset \mathcal{D} \tag{1}$$

$$d_i = (x_i, y_i, z_i, r_i) \tag{2}$$

Where $D_{\text{train}}$ is a subset of the D3 dataset ($\mathcal{D}$) consisting of a website's textual representation along with its deceptive pattern classification and reasoning. $x_i$ represents the web-element's text to classify, $y_i$ represents the deceptive design category, $z_i$ represents the deceptive design subtype, and $r_i$ represents the associated reasoning for the category and reasoning.

Using the above framework, we train the smaller T5 model, $f$, on a multi-task problem. The objective of the model is to produce a concatenated label for both category and subtype into a classification label ($\hat{y}_i$) and the reasoning behind that label ($\hat{r}_i$) and to minimize the prediction loss.

**Splitting Dataset** We split the D3 dataset into an 80% training set and a 20% testing set grouped by the sites the samples were extracted from. This ensures that samples from the same site are not present in training and testing sets.

**Baseline Approach** Based on this initial methodology, we train the smaller T5 model to predict the deceptive design category, subtype as the label, and reasoning as the rationale. We use the task-specific prefix `[classify]` to generate the label and `[reason]` for the reasoning.

We define the model and loss functions as follows:

$$f(x,t) = \begin{cases} (\hat{y}_i \oplus \hat{z}_i), & \text{if } t = [\texttt{category}] \\ \hat{r}_i, & \text{if } t = [\texttt{reason}] \end{cases} \tag{3}$$

$$\mathcal{L} = \mathcal{L}_{\text{label}} + \alpha \mathcal{L}_{\text{reason}} \tag{4}$$

Here, $\mathcal{L}_{\text{label}}$ and $\mathcal{L}_{\text{reason}}$ are the label prediction loss and reason generation loss, defined as:

$$\mathcal{L}_{\text{label}} = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i, t_{\text{category}}), y_i \oplus z_i) \tag{5}$$

$$\mathcal{L}_{\text{reason}} = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i, t_{\text{reason}}), r_i) \tag{6}$$

Where $\ell$ is the cross entropy loss between the predicted and target tokens, and $\oplus$ is a string concatenation function.

Each web element that is to be classified, $x_i$, is provided alongside its neighboring web elements, $x_{i-1 \to 0}$ to $x_{i+1 \to N}$. The model's performance is measured as the exact match of its `[category]` outputs with the ground truth data. An example of the model's sample input and expected output is shown below.

---

**Input Sample**

Input: `[category]`: Line 14,Preferences,checked checkbox,...</s>Line 10,"MAGIC We use cookies to personalise ...</s>Line 11,COMING SOON ,text,...</s>
Output: *obstruction,pre-selection*

Input: `[reason]`: Line 14,Preferences,checked checkbox,...</s>Line 10,"MAGIC We use cookies to personalise ...</s>Line 11,COMING SOON ,text,...</s>
Output: *Cookie banner option is pre-selected to indicate users to allow extra cookies.*

---

After training T5 on the ground truth dataset for 50 epochs (~4000 steps), we observed the training accuracy to saturate to ~48% and train loss around 4.3. Here, accuracy refers to the correct prediction of both category and sub-types.

**Our Approach:** To overcome the low accuracy in the baseline approach, we split the labeling task into two separate tasks: category and subtype, introducing an additional "task prefix" `[subtype]` for the new task. We redefine the model and loss function to:

$$f(x,t) = \begin{cases} \hat{y}_i, & \text{if } t = [\texttt{category}] \\ \hat{z}_i, & \text{if } t = [\texttt{subtype}] \\ \hat{r}_i, & \text{if } t = [\texttt{reason}] \end{cases} \tag{7}$$

$$\mathcal{L} = \alpha(\mathcal{L}_{\text{category}} + \mathcal{L}_{\text{subtype}}) + (1-\alpha)\mathcal{L}_{\text{reason}} \tag{8}$$

Where $\alpha$ is a tuning factor.

By separating `label` into `category` and `subtype` in addition to the `reason` tasks, the model can learn the relation between the category and the subtype and how they both relate to the reasoning. A sample of the new inputs and expected outputs is shown below.

---

**Input Sample**

**Input:** `[category]:` Line 14,Preferences,checked checkbox,...</s>Line 10,"MAGIC We use cookies to personalise ...</s>Line 11,COMING SOON ,text,...</s>
**Output:** *obstruction*

**Input:** `[subtype]:` Line 14,Preferences,checked checkbox,...</s>Line 10,"MAGIC We use cookies to personalise ...</s>Line 11,COMING SOON ,text,...</s>
**Output:** *pre-selection*

**Input:** `[reason]:` Line 14,Preferences,checked checkbox,...</s>Line 10,"MAGIC We use cookies to personalise ...</s>Line 11,COMING SOON ,text,...</s>
**Output:** *Cookie banner option is pre-selected to indicate users to allow extra cookies.*

---

Upon training the T5 model on the new loss function and tasks using the same ground truth dataset and the same number of epochs, we observe the test accuracy of detecting deceptive and non-deceptive patterns saturate to $\sim 95\%$ and the training loss saturating to $\sim 0.18$.

**Further Improving Accuracy** Language models cannot directly understand natural language; in order to process natural language, text is converted into tokens. The way prediction loss is calculated is as follows:

$$\ell(a',a) = \frac{1}{N}\sum_{i=0}^{N}(a_i - a'_i) \tag{9}$$

where, $a'$ is the predicted output, $a$ is the target output, and $a_i$ is the $i^{th}$ token of the $a$.

In our training dataset, the target labels are often multi-tokens, due to which noise is often introduced in loss calculations. To reduce this noise, we updated the labels to have single tokens. For instance, the following original labels:

| Label | Tokens |
|---|---|
| forced-action | $[5241, 18, 4787, 1]$ |
| nudge | $[3, 29, 13164, 1]$ |
| pre-selection | $[554, 18, 7, 15, 12252, 1]$ |
| hidden-subscription | $[5697, 18, 7304, 11830, 1]$ |

are updated to the following single token alternatives respectively:

| Label | Tokens |
|---|---|
| obligation | $[10472, 1]$ |
| push | $[3292, 1]$ |
| set | $[356, 1]$ |
| conceal | $[23808, 1]$ |

## 5.4 Performance of Distilled T5

Using this new methodology, of a single token, we again finetune the T5 model, achieving a final training accuracy of $\sim 97\%$. In Table 4, we observe that on the testing set, the T5 model can predict deceptive patterns with a high recall; the model can capture almost all the deceptive patterns in the site. When compared with the baseline approach and the predictions by Gemini 1.5, we observe that the recall of our model is the highest.

| Pattern Type | Our Approach | | Baseline T5 | | Baseline Gemini | |
|---|---|---|---|---|---|---|
| | Precision | Recall | Precision | Recall | Precision | Recall |
| **non-deceptive** | 0.96 | 0.68 | - | - | 1.00 | 1.00 |
| **obstruction** | 1.00 | 0.86 | - | - | 1.00 | 1.00 |
| **sneaking** | 0.66 | 0.91 | - | - | 0.91 | 0.90 |
| **interface-interference** | 0.75 | 0.72 | - | - | 0.94 | 0.91 |
| **forced-action** | 0.70 | 0.95 | - | - | 0.76 | 0.90 |
| not-applicable | 0.97 | 0.69 | 0.42 | 1.00 | 1.00 | 1.00 |
| hidden-costs | 0.72 | 0.64 | - | - | 0.93 | 1.00 |
| fake-scarcity-fake-urgency | 0.71 | 1.00 | - | - | 0.99 | 0.96 |
| forced-action | 0.72 | 0.84 | 0.93 | 0.38 | 0.76 | 0.90 |
| pre-selection | 1.00 | 1.00 | - | - | 1.00 | 1.00 |
| confirmshaming | 1.00 | 1.00 | - | - | 1.00 | 0.94 |
| nudge | 0.67 | 0.59 | - | - | 0.88 | 0.84 |
| trick-wording | 0.48 | 0.88 | - | - | 0.88 | 0.99 |
| **Deceptive** | 0.83 | 0.98 | 1.00 | 0.20 | 0.90 | 0.91 |

**Table 4:** Comparison of Precision and Recall across different Approaches and Pattern Types on DP-200

# 6   User-Facing Component

As discussed in Section 3, the *user-facing component* of *AutoBot* is a browser extension for Google Chrome, designed to integrate seamlessly into a user's browsing experience. To achieve this, it takes a two-step approach: (1) it sends a screenshot to the on-device *Vision module*; (2) it then parses the results from the *Language Module* to highlight deceptive patterns for the user as shown in Figure 1.

The extension overlays boxes around identified DPs by injecting an empty `<div>` tag into the DOM tree to bring the deceptive patterns to the users' attention. This div serves as a visual indicator, with a red border around the deceptive element (as shown in Figure 1), ensuring that users can quickly identify areas they should pay attention to.

We note here that the user experience is designed with passivity, allowing users to browse without distractions or unnecessary engagement. This passive design choice aligns with findings from works like [31, 6], emphasizing that user experiences maximizing functionality without competing for attention are more effective. Moreover, we allow users to change the message they see for any deceptive pattern, enhancing the overall user experience.

Our system integrates seamlessly with the user's browsing experience without burdening the user with unnecessary engagement.

# 7   Evaluations

We perform multiple experiments to evaluate the accuracy and usability of *AutoBot*, as well as to showcase its utility on a large-scale dataset. We seek to answer these questions:
  **Q1.** *What is the end-to-end performance of AutoBot?*
  **Q2.** *Can AutoBot run in real-time in a browser?*
  **Q3.** *What is the impact of AutoBot on the user experience?*

## 7.1   End-to-End Evaluation

We performed a manual end-to-end evaluation of *AutoBot* on 200 websites to determine its performance in helping users identify deceptive patterns in the web.

**Dataset Creation:** We captured single screenshots from the top-200 websites from the Tranco [36] list[9] to create our evaluation dataset for *AutoBot*. Next, we pass the 200 screenshots through the *Vision Module*. As mentioned in Section 4, the *Vision Module* produces the TextMap of a given image. Next, using the screenshots, one of the authors manually annotated the deceptive pattern present in the TextMaps, creating the ground truth dataset for evaluation. The distribution of deceptive patterns present in the dataset is given in Table 5

| Category | Subtype | # Samples |
|---|---|---|
| non-deceptive | not-applicable | 2378 |
| forced-action | forced-action | 26 |
| interface-interference | nudge | 22 |
| | fake-scarcity-fake-urgency | 3 |
| obstruction | visual-interference | 2 |
| sneaking | trick-wording | 10 |
| | hidden-costs | 10 |

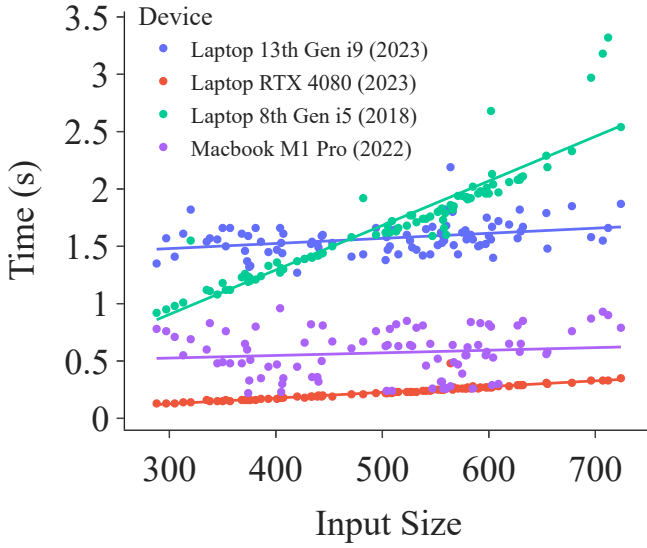**Table 5:** Distribution of deceptive patterns in top-200 Websites from Tranco List.

**Findings:** We apply the *Language Module* to the output of the *Vision Module* for the 200 websites and compare the model's prediction to the manually annotated dataset. We find that the model was able to detect deceptive patterns with a precision of 95% and a recall of 96% with an F1-score of 0.95. Detailed model performance results are shown in Table 6.

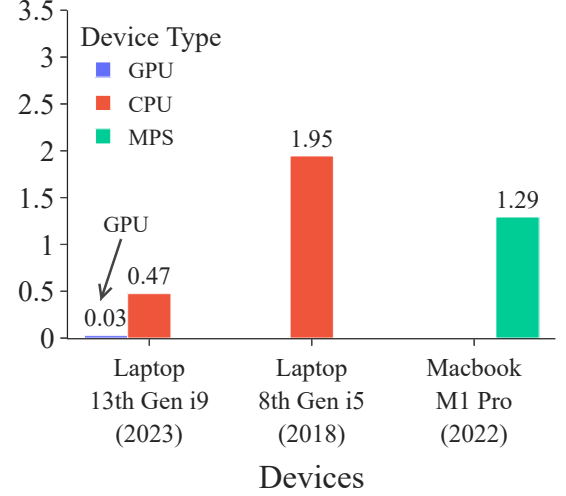| Pattern Type | Our Approach | | |
|---|---|---|---|
| | Precision | Recall | F1-Score |
| **non-deceptive** | 0.97 | 0.85 | 0.91 |
| **forced-action** | 0.74 | 0.96 | 0.83 |
| **interface-interference** | 0.78 | 0.72 | 0.75 |
| **obstruction** | 1.00 | 0.67 | 0.80 |
| **sneaking** | 0.95 | 0.90 | 0.92 |
| not-applicable | 1.00 | 0.88 | 0.94 |
| hidden-costs | 1.00 | 0.70 | 0.82 |
| fake-scarcity-fake-urgency | 0.75 | 1.00 | 0.86 |
| forced-action | 0.78 | 0.81 | 0.79 |
| visual-interference | 1.00 | 1.00 | 1.00 |
| nudge | 0.74 | 0.77 | 0.76 |
| trick-wording | 0.67 | 1.00 | 0.80 |
| **Deceptive** | 0.95 | 0.99 | 0.97 |

**Table 6:** Comparison of Precision and Recall across different Approaches and Pattern types

As mentioned in Section 4, the TextMaps generated by the *Vision Module* consists of the text extracted from OCR merged with the element detection by the YOLO models. The OCR performs with near-perfect accuracy, capturing all the text in the screenshot. The YOLO model is then used to predict the element types and states, selected or unselected. Such prediction assists the *Language Module* in detecting deceptive patterns like pre-selection, which the model can detect with high accuracy, as shown in Table 4. Prediction for UI elements like buttons helps the model understand the context of a text and what a user is experiencing.

---

[9]https://tranco-list.eu/list/QGJ74/1000000

**(a)** T5's latency when running on the three machines.

**(b)** YOLO's latency when running on the three machines

**Figure 6:** `YOLOv10` model inference time running on the 3 machines

We note here that the system's overall performance depends on the prediction of the T5 model. Even if certain misclassifications are present in the TextMap generated by the *Vision Module*, we find that the T5 model can detect the deceptive patterns accurately.

**Error Analysis:** In the end-to-end evaluation, we observe that the *AutoBot* fails to identify 22 deceptive pattern categories/subtypes correctly. Investigating the error cases further, we find that in 18 instances, *AutoBot* misclassified the category or sub-category of the deceptive pattern (while correctly identifying that the pattern is deceptive). For example, there are instances where *AutoBot* incorrectly identifies *forced-action* as *nudge*. We also find that on Wikipedia, *AutoBot* incorrectly tags a non-deceptive pattern as deceptive. The classification text contains money or cost-related text, causing the model to classify the text as *trick-wording* incorrectly. We also note that we only have one false negative in our evaluation case.

**Impact of Errors:** We note that the impact of the misclassifications due to category or sub-category mismatch is minimal on the users as the users will still be notified of a deceptive pattern. For false positives, the user gets notified for a pattern where none exists. Upon further inspection, users can safely ignore the notification, causing minimal distraction. For false negatives - the user impact can be severe. In such cases, users might get into a sense of false security and get deceived by the deceptive pattern because of *AutoBot*'s error. We emphasize that high recall of *AutoBot* ensures that such cases will be minimal.

## 7.2 System Level Evaluation

To assess the real-time performance of *AutoBot*, we conducted latency evaluations across four distinct machine configurations. We specifically measured the inference time for both the core Vision Module (YOLOv10) and the Language Module (T5) under two conditions: utilizing only the CPU and leveraging GPU acceleration, where available.

**Measurement Setup:** The latency of each module's execution was measured on three different systems. For the Vision Module, we focused on the YOLOv10 object detection model. In the Language Module, we measured the inference time of the fine-tuned T5 model when processing our evaluation dataset. We utilized the `YOLOv10m` model (approximately 40MB in size) and a fine-tuned T5 model (approximately 750MB).

**Results:** As shown in Figure 6a and Figure 6b, we observe that on modern machines, T5 model inference typically completes in under 2 seconds when using the CPU and in less than 0.5 seconds when utilizing a GPU. The YOLOv10 model processes images in approximately 1.2 seconds on CPU and under 30 milliseconds with GPU acceleration.

**Takeaway:** Our latency evaluations demonstrate that *AutoBot* can achieve near real-time performance on modern

17

computing devices. Leveraging GPU acceleration, core components like the T5 model and `YOLOv10` module exhibit inference times well under a second, ensuring a responsive user experience. This highlights the feasibility of deploying *AutoBot* for practical use, enabling seamless integration into real-world applications without significant performance bottlenecks.

## 7.3 User Study

We perform a user-based evaluation to study the effects of our system's highlighting feature on a website's usability. We recruited 151 participants from Prolific who reside in the United States. We paid each participant $2 to complete the study, with a median time of around 7 minutes. We did not collect demographic data and asked Prolific to distribute the survey evenly across the demographics. The IRB at our institution determined that the proposed activity is not research involving human subjects as defined by DHHS and FDA regulations.

### 7.3.1 Study Design

We conducted a within-subject study to measure if highlighting deceptive patterns on websites through our extension, *AutoBot*, affected the website's usability. We informed the participants that the objective of the survey was to test the usability of a web interface. Next, participants were asked to visit two custom-made websites, one with and one without highlighting, and fill out a System Usability Scale questionnaire [13] after each. In both these websites, the participants were asked to perform one of the four tasks: sign up, download, do shopping, or read a news article. Once the participants visit a website and hover over the highlighted text, a banner cautioning them about the deceptive pattern is shown. Note that the participants interacted with highlighted websites directly without installing a browser extension.

These custom-made websites were created by the authors based on various real websites, available at *UXP²* *Dark Patterns*[10], caught using deceptive patterns. An example of such a website is shown in Figure 1. After visiting each website, the participant is asked to complete a System Usability Scale questionnaire [13]. After finishing all the tasks, participants are asked to fill out a post-study questionnaire consisting of three questions:

Q1. *Did you find the hints related to the deceptive patterns useful?*
Q2. *Did you feel that the highlighted box encouraged you to notice the deceptive patterns?*
Q3. *Did you feel that the highlighted box encouraged you to change your choice?*

Finally, there was an open-ended question asking for general comments and feedback at the end of the survey.

### 7.3.2 Findings

We measure the effect of highlighting on the user's usability of a site using the SUS metric, and Figure 7 shows the SUS score for both original and highlighted websites. In this case, we consider the non-hypothesis to be that the website's usability is unaffected after highlighting. Based on the statistical test, the *p*-value for comparing two groups is 0.106, indicating no statistically significant difference between the SUS score for (un)highlighted websites. Therefore, the null hypothesis stands that usability is unaffected. We also note that the mean SUS score for highlighted is higher than that of unhighlighted by 2 points.

Additionally, we asked the participants three questions mentioned above. Their responses are shown in Figure 11. We see 56.3% of the participants consider the hint corresponding to the highlight to be useful in terms of recognizing deceptive patterns, 65.5% of them consider the highlighted box itself helpful, and 38% of them would change their behavior after introducing the highlighted deceptive patterns.

Reflecting on the user study's results, we claim that the design of highlighting deceptive patterns on webpages will not affect website usability and is effectively helping users recognize deceptive patterns.
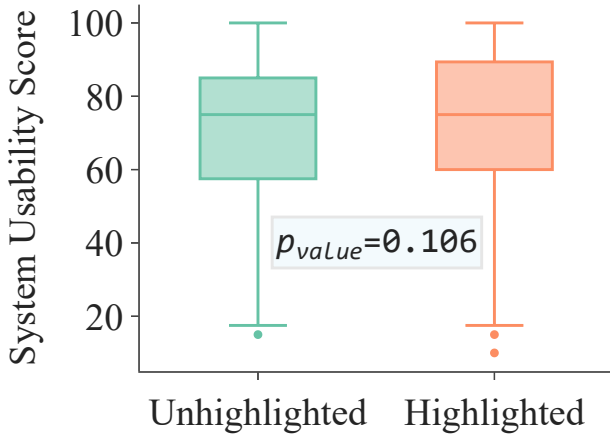
---

[10]`https://darkpatterns.uxp2.com/`

**Figure 7:** The results from the usability study show how highlighting affected the website's user experience. We find that usability is not affected by the highlights ($p = 0.106$).
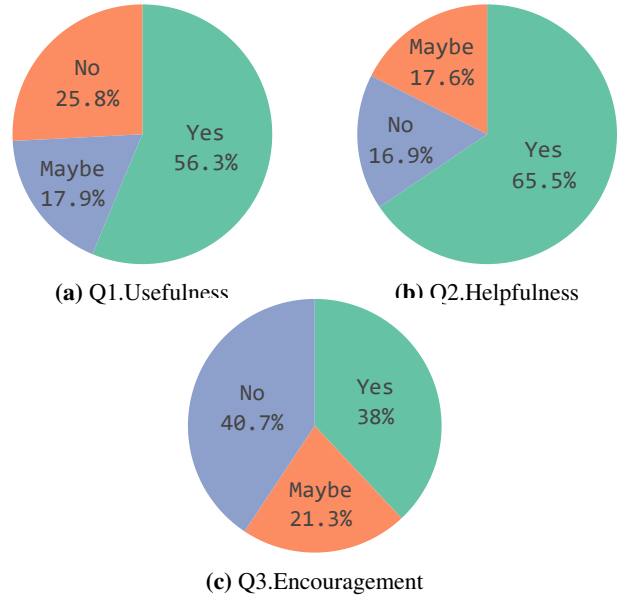


**(a)** Q1.Usefulness



**(b)** Q2.Helpfulness



**(c)** Q3.Encouragement

**Figure 8:** The distribution of participant responses to qualitative Q1-3 asked after the user study.

# 8 Discussion

We discuss the limitations of *AutoBot*, the associated deployment aspects, and how it can advance in multimodal foundational models.

**Limitations:** The main limitations of *AutoBot* stem from the underlying machine learning models regarding real-time performance and dependency on training data. Hardware constraints on user devices can limit *AutoBot*'s real-time performance. Although our implementation uses packages with optimized processing time for inference, the performance is slower when running on older devices without GPU support. Additionally, *AutoBot* exhibits a trade-off between latency and performance. While we can utilize larger, resource-intensive models, such as CNN instead of YOLO or Gemma 2B instead of T5, to get better accuracy, the hardware limitation will largely increase the inference time. Furthermore, limitations in the training data also impact *AutoBot*'s performance. We tried to ensure diversity within the training dataset; however, the models may exhibit reduced performance on certain websites. For instance, the lack of dark-themed websites in YOLO's training dataset affects its performance on high-contrast websites.

**Deployment Aspects:** *AutoBot* employs a modular design comprising of a pipeline of models. This design allows developers the flexibility to adjust *AutoBot* per their needs. In particular, the architecture of *AutoBot* can be updated over time: developers can customize their T5 models to substitute ours, adapting to new deceptive patterns or refined taxonomies. Moreover, developers can expose different interfaces for users. These interfaces can include passive alerts and active interventions with different highlighting options. Regulators can directly instrument the ML backend to enforce regulations like GDPR in real-time at a very low computational cost, making this a feasible solution for a long-standing problem.

**Multimodal Models:** Looking forward, *AutoBot* can enable multimodal foundational models customized for deceptive patterns and, more generally, online privacy, security, and safety (PSS). A developer can swap in models to use *AutoBot* at scale to provide a high-quality dataset of websites annotated for PSS tasks. For example, developers can include language models that target cookie notices. Collecting such a dataset can enable finetuning general-purpose, end-to-end multimodal models to generate a specialized model that can solve several PSS tasks on a webpage with minimal prompting and engineering.

**Assisting Regulators:** *AutoBot* offers a significant advancement for regulators who lack efficient tools to assess website compliance with regulations. By automatically detecting deceptive patterns, *AutoBot* provides a scalable solution for monitoring and enforcing rules like GDPR. This empowers regulators to proactively identify non-compliant

practices, reducing the burden on manual audits and enabling more effective enforcement.

## 9 Conclusion

In this paper, we introduce *AutoBot*, a browser extension that automatically detects the deceptive patterns on websites the user is visiting. *AutoBot* employs a pipelined multi-modal approach: first, it captures a screenshot of the website and processes it using OCR and YOLO to provide a text representation of the website. This representation is then analyzed by a distilled T5 model, trained on a specially curated dataset (D3), to determine whether deceptive patterns exist and their type. The user-facing component interprets the analysis result and renders the corresponding alert on the user's webpage. We conduct three sets of evaluations to verify that (1) *AutoBot* accurately detects deceptive patterns with high precision, (2) *AutoBot* operates in real-time within the browser, and (3) the website's usability is not affected by *AutoBot*.

## References

[1] Best 19 React UI Component Libraries in 2024.

[2] The best React UI library | Smile.

[3] WebForm elements detection Dataset > Overview.

[4] Regulation (EU) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/ec (general data protection regulation). `https://eur-lex.europa.eu/eli/reg/2016/679/oj`, 2016. Accessed: 2024-09-02.

[5] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[6] admin. Passive and Active Design Approaches, June 2017.

[7] J. H. Adorna, A. J. Dantis, R. Feria, L. L. Figueroa, and R. Solamo. Developing a browser extension for the automated detection of deceptive patterns in cookie banners. In *Proceedings of the Workshop on Computation: Theory and Practice (WCTP 2023)*, volume 20, page 101. Springer Nature, 2024.

[8] S. aka Steve Spiker. Tweet by spike aka steve spiker on x. `https://x.com/spjika/status/1686492710910427137`, 2023. Accessed: 2024-09-02.

[9] V. Bannihatti Kumar, R. Iyengar, N. Nisal, Y. Feng, H. Habib, P. Story, S. Cherivirala, M. Hagan, L. Cranor, S. Wilson, F. Schaub, and N. Sadeh. Finding a Choice in a Haystack: Automatic Extraction of Opt-Out Statements from Privacy Policy Text. In *Proceedings of The Web Conference 2020*, pages 1943–1954, Taipei Taiwan, Apr. 2020. ACM.

[10] C. Bösch, B. Erb, F. Kargl, H. Kopp, and S. Pfattheicher. Tales from the dark side: Privacy dark strategies and privacy dark patterns. *Proceedings on Privacy Enhancing Technologies*, 2016.

[11] H. Brignull. Deceptive patterns. `https://www.deceptive.design/`. Accessed: 2024-09-2.

[12] H. Brignull. Types of deceptive design. *Deceptive Design*, 2010.

[13] J. Brooke. Sus: A quick and dirty usability scale. *Usability Evaluation in INdustry/Taylor and Francis*, 1996.

[14] T. Cairns. extcolors: Extract colors from an image using k-means clustering. `https://pypi.org/project/extcolors/`, 2021. Accessed: 2024-09-02.

[15] California Privacy Protection Agency. California Privacy Rights Act (CPRA). `https://thecpra.org/`, 2024. Accessed: 2024-09-04.

[16] J. Chen, J. Sun, S. Feng, Z. Xing, Q. Lu, X. Xu, and C. Chen. Unveiling the tricks: Automated detection of dark patterns in mobile applications. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–20, 2023.

[17] G. Conti and E. Sobiesk. Malicious interface design: exploiting the user. In *Proceedings of the 19th international conference on World wide web*, pages 271–280, 2010.

[18] A. Curley, D. O'Sullivan, D. Gordon, B. Tierney, and I. Stavrakakis. The design of a framework for the detection of web-based dark patterns. 2021.

[19] Google Cloud. Google cloud vision ocr. `https://cloud.google.com/vision/docs/ocr`, 2024. Accessed: 2024-09-02.

[20] Google DeepMind. Google gemini. `https://gemini.google.com/app`, 2024. Accessed: 2024-09-02.

[21] C. M. Gray, Y. Kou, B. Battles, J. Hoggatt, and A. L. Toombs. The dark (patterns) side of ux design. In *Proceedings of the 2018 CHI conference on human factors in computing systems*, pages 1–14, 2018.

[22] C.-Y. Hsieh, C.-L. Li, C.-K. Yeh, H. Nakhost, Y. Fujii, A. Ratner, R. Krishna, C.-Y. Lee, and T. Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. *arXiv preprint arXiv:2305.02301*, 2023.

[23] Hugging Face. Transformers.js: Javascript library for running transformer models. `https://huggingface.co/docs/transformers.js/en/index`, 2024. Accessed: 2024-09-02.

[24] S. Kambhampati, K. Valmeekam, L. Guan, K. Stechly, M. Verma, S. Bhambri, L. Saldyt, and A. Murthy. Llms can't plan, but can help planning in llm-modulo frameworks. *arXiv preprint arXiv:2402.01817*, 2024.

[25] R. Khandelwal, T. Linden, H. Harkous, and K. Fawaz. {PriSEC}: A Privacy Settings Enforcement Controller. pages 465–482, 2021.

[26] R. Khandelwal, A. Nayak, H. Harkous, and K. Fawaz. Automated Cookie Notice Analysis and Enforcement.

[27] J. R. Landis and G. G. Koch. The measurement of observer agreement for categorical data. *biometrics*, pages 159–174, 1977.

[28] C. Lewis. Gameful Patterns. In C. Lewis, editor, *Irresistible Apps: Motivational Design Patterns for Apps, Games, and Web-based Communities*, pages 33–50. Apress, Berkeley, CA, 2014.

[29] T. F. Liu, M. Craft, J. Situ, E. Yumer, R. Mech, and R. Kumar. Learning design semantics for mobile apps. In *The 31st Annual ACM Symposium on User Interface Software and Technology*, UIST '18, pages 569–579, New York, NY, USA, 2018. ACM.

[30] J. Luguri and L. J. Strahilevitz. Shining a Light on Dark Patterns. *Journal of Legal Analysis*, 13(1):43–109, 03 2021.

[31] J. Madrid and M. C. Hout. Examining the effects of passive and active strategies on behavior during hybrid visual memory search: evidence from eye tracking. *Cognitive Research: Principles and Implications*, 4(1):39, Sept. 2019.

[32] S. H. Mansur, S. Salma, D. Awofisayo, and K. Moran. Aidui: Toward automated recognition of dark patterns in user interfaces. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 1958–1970. IEEE, 2023.

[33] A. Mathur, G. Acar, M. J. Friedman, E. Lucherini, J. Mayer, M. Chetty, and A. Narayanan. Dark patterns at scale: Findings from a crawl of 11k shopping websites. *Proceedings of the ACM on human-computer interaction*, 3(CSCW):1–32, 2019.

[34] A. Nayak, R. Khandelwal, E. Fernandes, and K. Fawaz. Experimental security analysis of sensitive data access by browser extensions. In *Proceedings of the ACM on Web Conference 2024*, pages 1283–1294, 2024.

[35] OpenAI. Gpt-4 technical report. Technical report, OpenAI, 2023. Accessed: 2024-09-02.

[36] V. L. Pochat, T. Van Goethem, S. Tajalizadehkhoob, M. Korczyński, and W. Joosen. Tranco: A research-oriented top sites ranking hardened against manipulation. *arXiv preprint arXiv:1806.01156*, 2018.

[37] Priya. 25+ Best React UI Component Libraries / Frameworks for 2024, Feb. 2019.

[38] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

[39] S. H. Raju, S. F. Waris, S. Adinarayna, V. C. Jadala, and G. S. Rao. Smart dark pattern detection: Making aware of misleading patterns through the intended app. In *Sentimental Analysis and Deep Learning: Proceedings of ICSADL 2021*, pages 933–947. Springer, 2022.

[40] J. Redmon. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.

[41] S. Ren. Faster r-cnn: Towards real-time object detection with region proposal networks. *arXiv preprint arXiv:1506.01497*, 2015.

[42] Schneider Wallace. Dark patterns: Making online subscriptions harder to cancel draws lawsuits, settlements, and government scrutiny, 2023. Accessed: 2024-09-02.

[43] V. P. Sermuga Pandian, S. Suleri, and P. D. M. Jarke. Uisketch: a large-scale dataset of ui element sketches. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2021.

[44] T. H. Soe, O. E. Nordberg, F. Guribye, and M. Slavkovik. Circumvention by design-dark patterns in cookie consent for online news outlets. In *Proceedings of the 11th nordic conference on human-computer interaction: Shaping experiences, shaping society*, pages 1–12, 2020.

[45] T. H. Soe, C. T. Santos, and M. Slavkovik. Automated detection of dark patterns in cookie banners: how to do it poorly and why it is hard to do it any other way. *arXiv preprint arXiv:2204.11836*, 2022.

[46] L. Stein. Tweet by luke stein on x. `https://x.com/lukestein/status/1150014732486742016`, 2019. Accessed: 2024-09-02.

[47] A. Wang, H. Chen, L. Liu, K. Chen, Z. Lin, J. Han, and G. Ding. Yolov10: Real-time end-to-end object detection. *arXiv preprint arXiv:2405.14458*, 2024.

[48] D. Wang, J. Anthony Ribando, D. Mo, and N. Tung. insite. a chrome extension that protects consumers from marketing tricks when they shop online., 2019.

[49] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

[50] J. Wu, S. Wang, S. Shen, Y.-H. Peng, J. Nichols, and J. Bigham. Webui: A dataset for enhancing visual ui understanding with web semantics. *ACM Conference on Human Factors in Computing Systems (CHI)*, 2023.

[51] M. Xie, S. Feng, Z. Xing, J. Chen, and C. Chen. Uied: a hybrid tool for gui element detection. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ESEC/FSE 2020, page 1655–1659, New York, NY, USA, 2020. Association for Computing Machinery.

[52] X. Zhang, L. De Greef, A. Swearngin, S. White, K. Murray, L. Yu, Q. Shan, J. Nichols, J. Wu, C. Fleizach, et al. Screen recognition: Creating accessibility metadata for mobile applications from pixels. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–15, 2021.

# A Appendix

## A.1 Website Generation Prompts

To generate website using `v0` we used prompts from GPT4. Some example prompts are shown below:

1. Create a virtual learning platform with courses, webinars, and interactive tools for students of all ages.

2. Create an online gourmet food shop featuring high-quality ingredients, kitchen tools, and gourmet recipes.

3. Design a warm and inviting UI for a pet care blog that radiates friendliness and approachability. The primary color palette should include soft, earthy tones with playful accents. Use a clean, easy-to-read sans serif font and incorporate elements like paw prints or pet silhouettes to enhance the thematic appeal. The homepage must prominently feature an engaging welcome message, and sections for various pets like dogs, cats, birds, and more, encouraging user navigation. Include a dynamic sidebar with widgets for pet care tips, a search bar, and featured posts. Dropdown menus should be intuitive, providing categories such as nutrition, training, health, and grooming. Visuals are key: integrate heart-warming pet images and infographics to explain care practices. End with footer links to contact details, social media, and a cute, animated pet mascot offering useful tips periodically.

## A.2 User Study Results

We created custom websites for the user study, measuring the effect of highlighting on usability of a website. We show some example screenshots of the websites used in Figures 9 and 10. In Figure 11, we show the response distribution received to our qualitative question asked to participants after the user study.
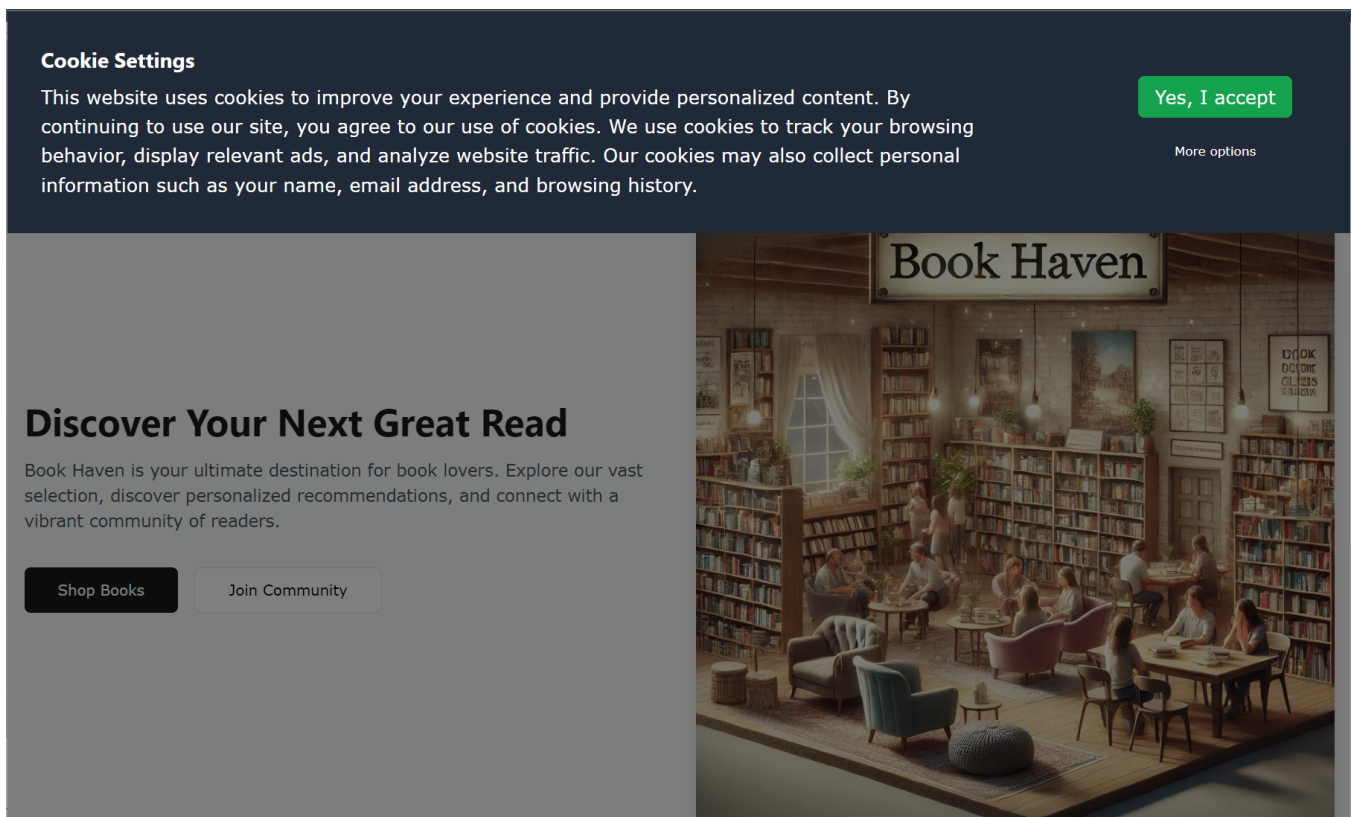
**Figure 9:** A website showing a book store's homepage with a cookie notice, nudging users to accept cookies.
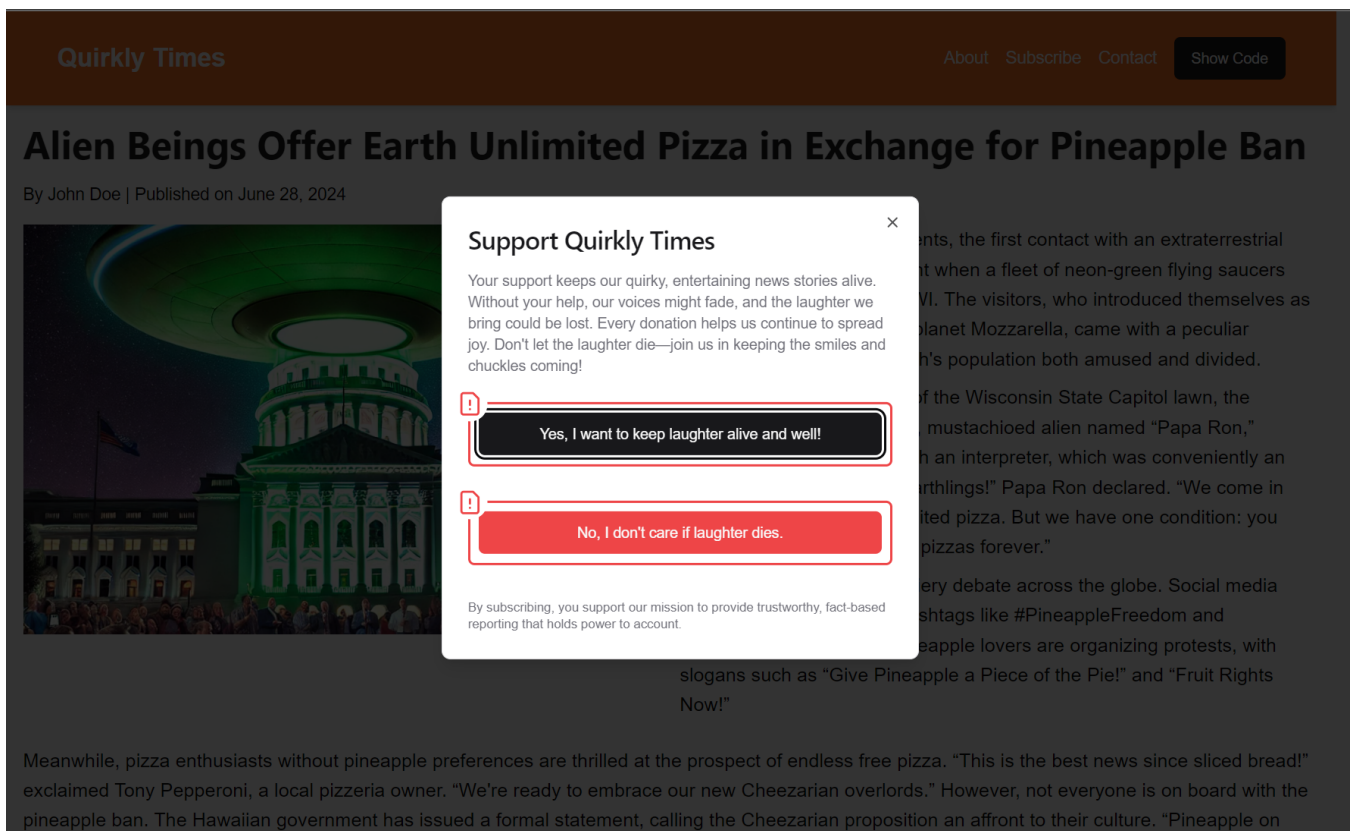


**Figure 10:** A news organization's website asking users to donate. The deceptive patterns are highlighted in this example with a red box around them.

**(a)** Q1. Usefulness        **(b)** Q2. Helpfulness        **(c)** Q3. Encouragement
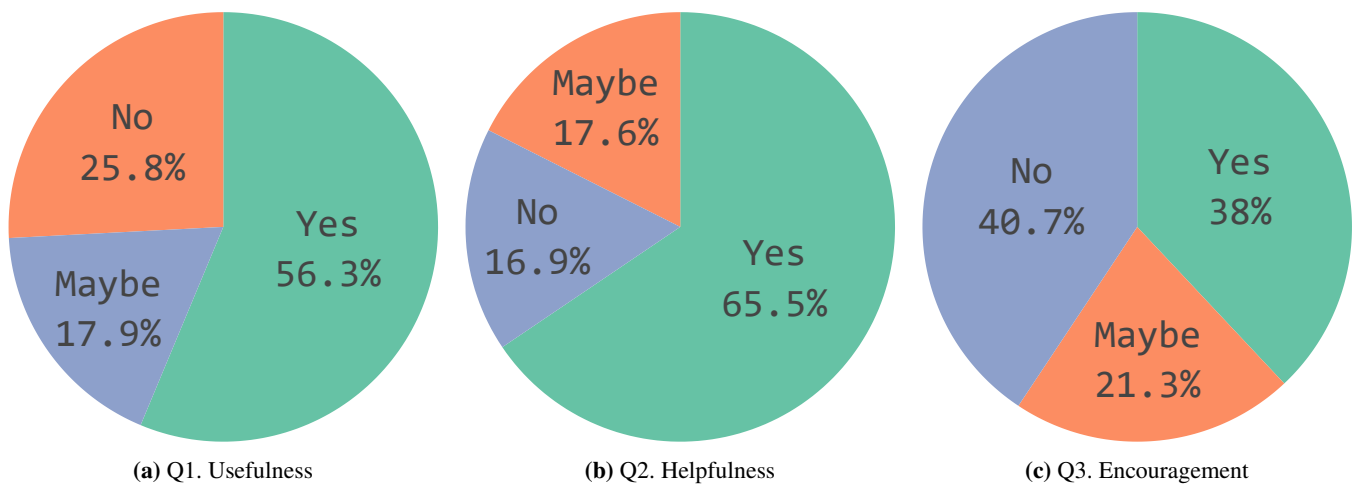
**Figure 11:** The distribution of participant responses to qualitative Q1-3 asked after the user study.