# TermSight: Making Service Contracts Approachable

ZIHENG HUANG, University of Illinois Urbana-Champaign, USA

TAL AUGUST, University of Illinois Urbana-Champaign, USA

HARI SUNDARAM, University of Illinois Urbana-Champaign, USA

Terms of Service (ToS) are ubiquitous, legally binding contracts that govern consumers' digital interactions. However, ToS are not designed to be read: they are filled with pages of ambiguous and complex legal terminology that burden potential users. We introduce TermSight, an intelligent reading interface designed to make ToS more approachable. TermSight offers visual summaries that highlight the relevance and power balance of information in a ToS. TermSight also categorizes and simplifies information within the ToS into concise plain-language summaries. To aid in reading the original text, TermSight offers contextualized definitions and scenarios for unfamiliar phrases. Our within-subjects evaluation of TermSight (N=20) revealed that TermSight significantly reduced the difficulty of reading ToS and increased participants' willingness to do so. We also observed emerging strategies that participants took when interacting with AI-powered features that highlight the diverse ways that TermSight assisted ToS reading.

CCS Concepts: • **Human-centered computing** → **Interactive systems and tools**.

Additional Key Words and Phrases: large language models, legal documents, reading support tools

## 1 INTRODUCTION

Terms of Service (ToS)—the ubiquitous, rarely-read but legally binding contracts that govern the use of digital services— are widely incomprehensible to potential consumers. Despite consumers' self-reported concern for the information presented in ToS [73, 76], and their eventual regret when they encounter issues related to this information (e.g., the selling of their data) [34, 77], almost no one reads ToS content fully [3, 34, 77]. There is good reason not to do so: ToS are not designed to be read. ToS are filled with legal jargon that can derail general-audience (also referred to as 'lay-audience') readers with confusing or under-defined terms (e.g., a ToS stating the service will collect certain 'information') [13, 45, 61, 94, 102]. ToS also contain more than one policy (e.g., a privacy policy) that in turn each contain reams of boilerplate or semi-boilerplate text. Within these pages of text, there might be information relevant to a reader (e.g., an arbitration clause or copyright policy) but such information is lost among the many pages of text [77, 89]. Legal scholars have argued that these ambiguous, jargoned, and exceptionally long contracts are designed not to cultivate a shared understanding, but to mitigate liability and establish unilateral control [11, 51, 52, 78].

Summarizing and simplifying text is a first step in making ToS more approachable, but resolves only part of the challenge these contracts pose. Prior interfaces have leveraged crowdsourcing [1, 2] and, more recently, AI-based approaches like large language models (LLMs) to extract and process text automatically from ToS [79]. However, these approaches primarily focus on presenting information from pre-defined categories (e.g., user privacy) or on single
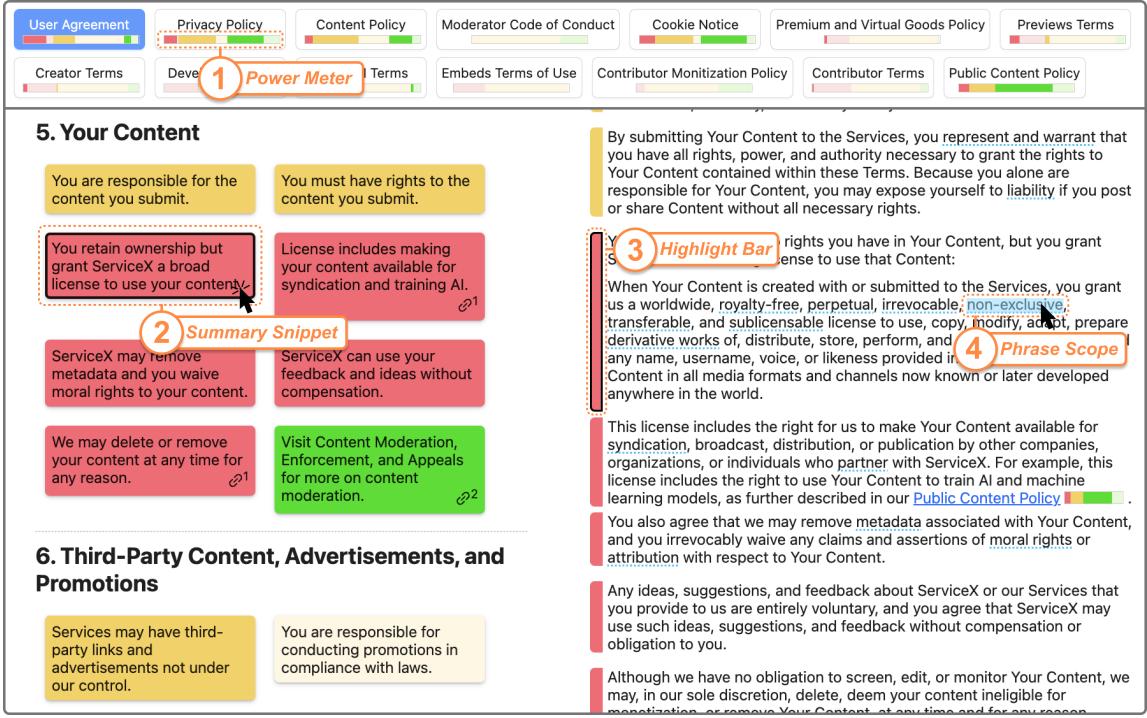
Fig. 1. The TermSight system provides guidance and reading support at multiple levels of ToS reading. At the contract level, TermSight visualizes the relevance and power balance of content in each policy (1). At the document level, TermSight chunks, summarizes, and categorizes content into 1-sentence plain-language summaries (Summary Snippets) highlighted with colors that reflect power and relevance (2). Readers can click the Summary Snippet (2) or Highlight Bar (3) to navigate between the summary snippets on the left and the original text on the right. At the phrase level, TermSight offers phrase definitions and scenarios for unfamiliar and vague phrases (4).

policies within a ToS (e.g., a privacy policy). In reality, ToS contain many policies that cover different information (§3.3.1), and what information a reader needs is dependent on their planned usage of the site (e.g., a seller vs. buyer on an e-commerce site).

We envision a new reading experience for ToS. Our goal is to enable readers to find information important to them anywhere within a ToS and scaffold their ability to understand and contextualize this information within their potential use of a service. We start by developing a more nuanced understanding of readers' information needs and challenges when reading ToS through a formative observational study (N=20). We found that participants cared about different information based on their envisioned usage of a service (e.g., a content creator will have different needs than a passive consumer on a social media website) and their own personal values (e.g., people had different levels of comfort when it came to privacy and data sharing). When participants sought to resolve these information needs, they encountered barriers at every level of granularity within a ToS. Within the full ToS contract, participants struggled to navigate nested policies and decide which policies to read. Within a specific policy of the ToS (e.g., the privacy policy), participants struggled to surface and interpret relevant information from extended, visually dense, and difficult text. When reading the text itself, participants struggled to interpret jargon and phrases that were ambiguous or vague. Interestingly, existing affordances (e.g., policy names, section headers with table of contents, and summaries) were often not helpful

and at times misleading. For example, participants often found that policy names and section headers were unclear and that overview summaries used deceptive friendly language to hide unfavorable clauses.

We address the specific challenges surfaced by our formative work with a new intelligent reading interface, TermSight. TermSight provides multi-level reading guidance and support through contract-level overviews (Power Meter), document-level plain-language summaries (Summary Snippets), and phrase-level definitions with scenarios of potential implications (Phrase Scope). To evaluate TermSight, we conducted a counterbalanced within-subject study (N=20) comparing TermSight and a baseline HTML reader during a timed reading task. Participants found reading ToS to be significantly easier with TermSight and were more willing to read and spend time on ToS. Participants reported that features of TermSight provided guidance and reading support at all granularities of ToS, resolving challenges surfaced in our formative study. We also observed interesting and unexpected user behaviors centered around trade-offs between TermSight's AI-generated features.

In this paper, we make the following contributions:

(1) We characterize the challenges readers face when trying to make sense of ToS contracts. Readers face challenges at all granularities of a ToS, from determining which policies to read to being able to resolve individual ambiguous terms within the text.
(2) We design and develop an intelligent reading interface, TermSight, for helping readers make sense of ToS contracts at the contract, document, and phrase levels.
(3) We collect empirical evidence from our user study (N=20) demonstrating the potential value of the interaction affordances of TermSight.

## 2 BACKGROUND AND RELATED WORK

### 2.1 Service Contracts and The Duty to Read

A Terms of Service (ToS) is a type of standard form contract consisting of a body of policies and conditions outlined by a service provider, dictating the rules and expectations for both the service provider and the consumer. These contracts have evolved to be more legally comprehensive compared to earlier ToS [4], encompassing a wide range of provisions such as copyright, privacy, returns, acceptable use, and service-specific terms. By clicking agree to a ToS, a legally binding contract is formed [28]. U.S. contract law is premised on the doctrine of the 'duty to read' where contracting parties are legally responsible for reading and understanding the contract [11]. However, the 'duty to read' was established in the era when contracts were short and based on mutual negotiation [51]. Modern ToS often have important terms buried within a large body of text and documents [89]. The language used to write ToS can also require advanced reading capability [4, 89] and contain unfamiliar [45, 94] or ambiguous terminology [13, 61, 102]. Due to the cognitive load associated with reading ToS, most consumers do not read modern ToS [7, 77], despite the binding nature of ToS and consumers' self-reports indicating concern for the information presented in ToS [73, 76]. ToS can also vary in content [30, 65] and structure [45] across services, leading to false consumer expectations. This widespread disengagement highlights what Kar and Radin described as a "paradigm slip" in contract law, where agreements once premised on shared understanding and mutual consent are now reduced to unilateral boilerplate "pseudo-contracts", eroding fundamental principles of contract law and personal autonomy [11, 51, 78]. Our formative study complements prior findings by aggregating and uncovering consumers' varying information needs and additional challenges encountered when navigating and reading policies within a ToS.

## 2.2 Augmenting Terms of Service

There are two major threads of work on improving the readability of ToS that we cover below: changing how ToS are designed or presented (§2.2.1) and providing summaries of ToS (§2.2.2).

*2.2.1 Stylistic Manipulations.* To help readers better understand the implications of using a service, prior work has investigated design recommendations for service providers to make documents within ToS more accessible [38, 54, 93]. Kay and Terry [54] proposed Textured Agreements that use typographic manipulation, pull quotes, vignettes, and iconic symbols to improve reader attention and comprehension on privacy policy compared to plain text. Habib et al. [38] investigated how icons and linked text could be designed to better convey privacy choices. Taber et al. [93] proposed crowdsourced sentiment highlighting of sentences in ToS, providing visual cues to guide readers toward sentences in ToS with negative sentiment. Design recommendations on styling user agreements have also been proposed [38, 54]. While helpful in navigating contracts, these design changes do not simplify the content itself and do not take into consideration the fact that readers have different information needs when reading a ToS.

*2.2.2 Summarization.* Work has also proposed helping consumers gain an overview of specific policies within a ToS, often the privacy policy (PP), without requiring reading the original document. Platform for Privacy Preferences (P3P) is an early attempt to standardize privacy policies by allowing service providers to submit privacy policies in a machine-readable format [83]. Tools such as Privacy Bird build atop the P3P format to provide warnings when a site's PP does not match consumers' privacy preferences [26]. However, P3P is rarely adopted by service providers [25] and no longer supported [72]. Work has also proposed the use of a single-page summary for end-user license agreements (EULAs) [34, 35], comic-based summaries for privacy policies [92], and Privacy Policy Nutrition Labels (PPNL) that label privacy policies in a table format [55, 84] to help users gain an overview of a policy.

Many earlier tools required buy-in from service providers to design better ToS. Other work, though, has focused on providing an overview summary of a ToS automatically or via crowdsourcing. For example, ToS;DR leverages volunteers to label and summarize terms in the privacy policy and the home page of ToS [2]. Works like CLAUDETTE automatically detect a predefined set of potentially unfair clauses from ToS [16, 17, 36, 62]. PrivacyCheck [74, 75, 104], PrivacyGuide [95], and Polisis [39] extract terms relevant to a list of privacy-related questions and criteria and generate a grading for each. Sathyendra et al. extract opt-out choices from privacy policies [8, 86]. Relatedly, some works evaluate the completeness of privacy policies [24] and the GDPR compliance of extracted clauses [64]. More recently, Pan et al. [79] investigated generating privacy nutrition labels from a privacy policy. Work has also attempted to generate an aggregated summary based on the extracted information from a privacy policy [56, 96]. TermSight extends this rich prior work by combining the benefit of plain language summaries with interactive affordances that enable readers to navigate the original ToS text, helping readers navigate and read the original body of policies within ToS beyond privacy policies.

## 2.3 Intelligent Reading Interfaces

TermSight takes inspiration from a rich body of work on intelligent reading interfaces that support readers in navigating documents and clarifying information. Prior work has explored helping readers navigate documents and locate relevant information in scholarly documents. For example, work has highlighted key information in research papers, such as objectives, results, and methods, to support skimming [32]. Also in research papers, ScholarPhi identifies position-sensitive definitions of terms and symbols, enabling readers to track their meaning across contexts [41]. Talk to Papers

allows readers to ask questions about research papers directly, while Qlarify supports clarification by letting readers expand text spans in paper abstracts to generate explanations [31, 106]. In medical research papers, Paper Plain provides a list of guiding questions to help readers navigate important information in a medical paper [6]. In the business domain, Marco allows readers to search and ask questions across collections of business documents [33]. In the news domain, Chen et al. [23] introduced Marvista, a reading interface that identifies the most summative portions of news articles based on a reader's time constraints or topic in the news article.

Compared to prior reading interfaces that have focused on research papers, news articles, or business documents, TermSight focuses on helping general audience readers make sense of legal documents, namely, ToS. Legal documents, and specifically ToS, present challenges different from prior reading contexts. In contrast to other knowledge-intensive settings (e.g., scientific or medical text), general audience readers in the legal setting are faced with language not designed to communicate information. Readers then struggle to (1) determine what text to read within and across policies, (2) define and contextualize vague language, and (3) navigate existing affordances that are themselves ambiguous or incomplete. Prior intelligent reading interfaces and context often assume that there is a single document [6, 23, 32], readers have consistent information needs [6], readers know what they need to read or ask [32, 33, 106], or readers can make sense of term definitions [6, 31, 41]. These existing affordances, while powerful, do not fully support readers' needs and challenges encountered when reading digital contracts such as ToS. TermSight takes inspiration from prior works to explore the design of reading interfaces to help readers surface and interpret relevant information from legal documents.

## 3 FORMATIVE STUDY

To better understand the challenges readers face when trying to make sense of the Term of Service and their information needs, we conducted a formative study to answer the following research questions:

**RQ1**: What information do readers want from a ToS?

**RQ2**: What are the challenges readers encounter when trying to get this information from ToS?

### 3.1 Methodology

*3.1.1 Study Procedure:* The entire study took 40 minutes and participants were compensated 10$. The study was approved by the Institutional Review Board (IRB) office at our organization. Participants initially completed a demographics form and reported their past experiences interacting with ToS. Following the form, participants were asked to imagine as if they were a first-time user of the assigned service registering for the service and were provided 15 minutes to go through the ToS. Participants were free to navigate to different policies within the ToS. While reading, participants were instructed to speak aloud any challenges they faced or thoughts about the ToS. After reading the ToS, we conducted semi-structured interviews with participants about their experiences and challenges with reading ToS. Lastly, we asked the participants to envision how computational supports could aid their reading. All interview sessions were transcribed. The semi-structured interview questions are listed in Appendix E.

*3.1.2 Analysis.* We conducted a reflexive thematic analysis on the interview transcripts to identify common challenges and themes. We followed the six phases of reflexive thematic analysis suggested by Braun and Clarke [18]. The lead researcher thoroughly explored the data, noted interesting features, systematically coded the data, and iteratively compared the codes to generate the initial themes. Then, the lead researcher discussed with other members of the research team to make sure the themes fit the data and further developed the themes.

*3.1.3   Materials.* Participants were randomly assigned one of the 10 ToS in the study. We selected ToS from 10 services spanning social media, e-commerce, video streaming, and internet services (Social Media: Reddit, Facebook, Instagram, Twitter; E-commerce: Amazon, eBay; Video: Youtube, Netflix; General: Google, Yahoo). These 10 services were selected from the top 20 visited sites on semrush.com (2024/05) where we selected the top 4 social media services, top 2 E-commerce services, top 2 video platforms, and top 2 internet services. We over-sampled social media services because they are the most represented service category in the top 10 visited services. In the study, participants were first directed to the main ToS and were allowed to navigate across any other policies.

*3.1.4   Participants.* We recruited 20 participants through Prolific. Each participant was randomly assigned one of the 10 selected ToS. All participants were over 18, fluent in English, and located in the US. 14 participants self-identified as female and 6 as male. 3 participants were between the ages of 18-25, 9 were between 26-35, 4 were between 36-45, 3 were between 46-55, and 1 was 67. 3 participants had never read or skimmed a ToS before and 17 participants had read or skimmed at least 1 ToS before. We did not require participants to have prior experience in reading the ToS. All the participants either had not read or did not remember reading the ToS they were assigned.

## 3.2   RQ1: What information do readers want from a ToS?

*3.2.1   **Participants have diverse information needs based on usage of the service and personal value***. All 20 participants described the need to know what control and rights the service had over the user, what the service was allowed to do, and how the user could be negatively impacted. Additionally, 12 participants highlighted the need to know *their* (as opposed to the service's) rights and control. Furthermore, 4 participants expressed wanting to know what they were allowed or not allowed to do on the platform. Participants also noted more specific and diverse information that they cared about such as data collection and usage (11), intellectual property right over user content (8), purchasing and returns (7), account deletion by the service (5), arbitration and liabilities (3), content moderation (2), and exposure to offensive content or misinformation (1).

10 participants explicitly mentioned how their information needs and priorities depended on usage and personal values. For example, P1 described that *"I don't think any information is irrelevant. I just think some parts are more of a priority of knowing. If I happen to become a person who wants to create graphic designs, I would definitely utilize the Copyright portion of the ToS to ensure that I'm being compliant with [the platform's] expectations."* On the other hand, P3, with a background in software, raised privacy concerns when using services for professional communication: *"For Microsoft Teams and Gmail, I was using them for a collaborative project with a small company. I wanted to ensure that whatever data or information was shared would remain proprietary to us."*

## 3.3   RQ2: What are the challenges readers encounter?

Our findings revealed that participants encountered challenges at the contract (§3.3.1), document (§3.3.2), and phrase (§3.3.4) level. In addition, participants highlighted how existing navigational and reading affordances were often ineffective, misleading, and failed to support meaningful comprehension (§3.3.3).

*3.3.1   **Contract Level: Participants struggled to navigate nested policies and decide which policy to read.***
11 participants mentioned that it's unclear to them what sub-policies are included when they are navigating the ToS because sub-policies are often hyperlinked throughout the documents. There often lack a centralized list of what policies are included as part of ToS. 7 participants described the challenge of navigating across these nested policies as getting *'lost in a whole tree of information'* (P18).

All 20 participants described relying on policy names to help decide whether to explore a policy. Yet, 17 participants described that the names of the policies alone did not provide a clear mental model of what was in the sub-policies and whether there was important information they should know about. For example, P14 mentioned that: *"When it's hyperlinked like 'our rules and policies' and 'privacy policy', it's difficult to determine if I need to know something from the link."* 8 participants explicitly expressed their reluctance to click into the linked polices due to the lack of awareness of relevant information and the challenge of finding relevant information from a hyperlinked document.

### 3.3.2 *Document Level: Participants lacked guidance and struggled to surface relevant information from extended, visually dense, and obfuscating text within a policy.* 17 out of 20 participants described their reading behavior as skimming the ToS to see if any information caught their eye. 6 participants noted that they were not sure what they should look for in the ToS and might have unintentionally missed important information. P12 further commented how it's not always obvious if the information was relevant at first sight: *"A lot of it seemed irrelevant until I started reading in a little bit more."*

18 participants highlighted the overwhelming length of the document, especially when considering the nested policies. These participants expressed concerns about missing important terms buried in the text. For example, P7 pointed out that *"When all the important information gets hidden behind walls of text, it doesn't really incentivize anyone to read it, and it gets hidden."* Reflecting P7's comment, 13 participants described the text as visually dense and difficult to parse when skimming. These participants further noted how important information can be positioned at the very bottom of an extended policy that was hard to discover. For example, P10 noticed how information about returns and cancellations of an E-commerce platform was positioned under the 13th section of the ToS with the title 'Additional Terms'.

Additionally, 3 participants described how the use of friendly language actually obfuscated the ToS and made skimming more challenging. For example, P15 noticed language such as *'We only use data to make [the service] a better place'* when reading a privacy policy. Yet, she later found out how user information can be used for targeted ads. As a result, P15 described that it's easy for lay people to misinterpret the language and miss important terms: *"If someone's just skimming the terms of service and is not skeptical at all. Then, you could read over that and think that means they're not using your data for targeted ads or stuff."*

### 3.3.3 *Document Level: Participants found existing affordances to be ineffective and misleading.* Some of the ToS participants read contained navigational and reading affordances like an interactive table of contents with section headers and overview summaries. However, 12 participants noted that, similar to policy names, section headings within policies were vague and not descriptive of the actual content. As a result, it was unclear what information was contained in a section by only looking at the section header. For example, when skimming the ToS, P10 skipped the 'Content' section but later realized that the section was about intellectual property rights over user-generated content: *"It just says 'content' which is pretty vague when they're talking about getting exclusive intellectual rights. I feel like it's misleading. So people would skip right past it. I did the first time."*

In addition to ineffective navigational affordances such as vague policy names (17/20) and section headers (12/20), 3 participants accessed a sub-policy with a summary at the top. All three participants (P4, P14, P15) described using the summaries to gain an overview at the start. Yet, after reading the policy, participants noticed that the summary was missing key details: *"The summary is kind of misleading. They gave you a more palatable version of [ToS]. If you scroll down, it will say we can terminate your account for no reason. They don't say that in the summary"* (P4). They also expressed concern that readers might only glance at the summary and mistakenly think it covers the entire policy.

*3.3.4* ***Phrase Level: Participants struggled to contextualize unfamiliar or ambiguous terminology***. When reading individual sections in a policy, 12 participants pointed out unfamiliar legal terminology (e.g., *Arbitration*, *Indemnity*, or *class lawsuit*). These participants described how the terminology hindered their understanding and motivation to read the section. For example, P8 described skipping sections that were full of unfamiliar terms. Additionally, participants noted how challenging it was to interpret the meaning of a word or phrase in the context of their situation: *"I can recognize these are legal terms, but I don't necessarily know what that means for me"* (P10). P6 also described how it's hard to envision scenarios in which things go wrong beforehand from the ToS.

8 participants noticed ambiguous or vague terminologies that made it challenging for them to understand the implications of signing. For example, P15 mentioned that *"I really think the big difficulty is the vagueness of the language and the constancy of exceptions that are vague which gives them a lot more leeway."* More specifically, 5 participants noticed vague terms around the service's information collection and sharing practices such as: *'we share you data with 3rd parties', 'track information about your interaction', and 'retain certain information about you'.* 4 participants noted vague terms around how the service can use user-generated content. For example, P3 and P14 pointed out that it's unclear what *"use, distribute, modify, and copy user content"* implicates and what are examples of what the service can do with user-generated content.

## 3.4  Design Guidelines:

We propose four design guidelines based on our observations that readers faced interrelated barriers at all levels of ToS reading:

> **DG1**: Contract level: Help readers form mental models of each policy of a ToS to determine important documents to interrogate.
> **DG2**: Document level: Help readers identify potentially relevant information within a policy.
> **DG3**: Phrase level: Help readers read and interpret original text with technical language, jargon, and vague phrases.
> **DG4**: All summaries should be tightly coupled with the text used to generate them. When reading any generated overviews, readers should be able to retrieve the original document text in at most one click.

## 4  THE TERMSIGHT SYSTEM:

We reify these design guidelines in TermSight, an augmented reading interface that supports user sense-making of ToS with features at each level of granularity of a ToS:

> **Power Meter:** A *contract-level* visualization of the relevance of information within a document and the distribution of power between the user and the platform that information represents.
> **Summary Snippets:** 1-sentence plain language summaries of information chunks at the *document-level.* Snippets are rendered with color and saturation that represent power and relevance to the user, similar to the Power Meter.
> **Phrase Scope:** In-situ *phrase-level* definitions and scenarios via a tooltip that contextualizes the meaning of the phrase while allowing users to ask clarification questions.

We adopted an iterative design approach in developing TermSight. 8 participants evaluated an early prototype of TermSight, and their feedback informed the final design. Additional details about the iterative design study, results, and design changes can be found in Appendix A.

### 4.1 Design Abstractions:

The design of TermSight is grounded in two design abstractions: 1) *Information Snippets*, and 2) *Power* and *Relevance* classification. These abstractions serve as foundational building blocks for TermSight's features that support contract and document-level sensemaking as illustrated in Figure 6.

*4.1.1 Information Snippets.* Our formative study reveals that the information readers cared about in ToS is often buried within extended and visually dense text. Consequently, prior works that support document understanding through document-level summaries [56, 96] or section-level summaries [6, 68] can fall short in helping readers surface pieces of information buried in large sections. To address this limitation, we use **Information Snippets** as a basic building block for TermSight. An information snippet is a continuous span of original text that shares the same topic and can be summarized in a single sentence of up to 12 words (detailed in Section 5.2). These information snippets form the basis of the Summary Snippets (§4.3.1). Information Snippets are non-overlapping and together reconstruct the original document's content (Examples in Figure 2). Figure 6 illustrates the process of how TermSight divides a paragraph into five Information Snippets, each paired with its Summary Snippet. The idea of designing interactions and sense-making around Information Snippets rather than entire sections was inspired by prior works that modularize and 'objectify' tools [9, 14], attributes [101], and AI agent's memory [44] to enable interaction flexibility, interaction specificity, and direct manipulation.
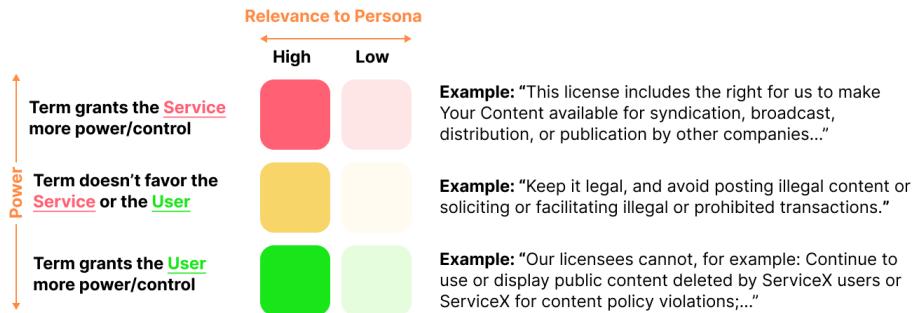


Fig. 2. Design conceptualization of Power and Relevance: Power refers to the degree to which a snippet grants control to the service provider or the user. Relevance refers to whether or not the snippet is relevant to the user's intended usage of the service and personal value. Examples of Information Snippets, TermSight classified under each category of power, are shown.

*4.1.2 Power and Relevance Classification.* We define **Power** and **Relevance** as two qualities of an information—and thus a summary—snippet. Power is the degree to which a snippet's text grants control to the Service Provider or the User (Categories: Service, Neutral, or User). Relevance is whether the snippet's text is directly relevant to the user's intended usage of the service or their values (Categories: High, Low). For example, a snippet related to selling would have low relevance for a user more focused on buying. As shown in Figure 2, we visually encode Power and Relevance in TermSight with a combination of hue (Red, Yellow, Green) and saturation (High, Low).
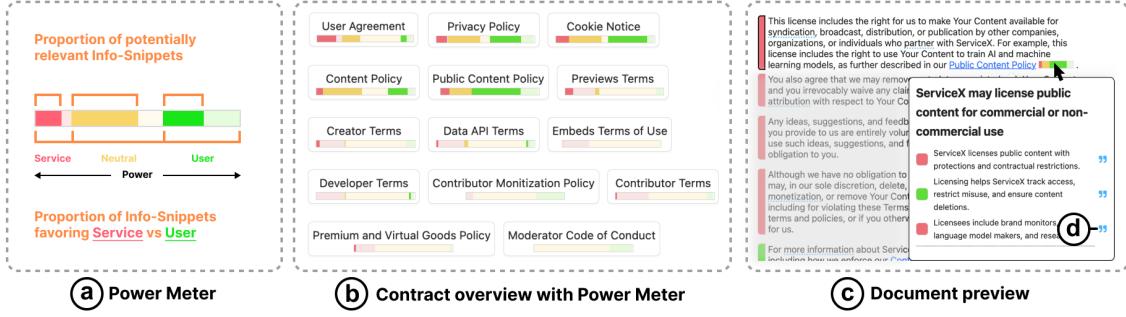
Fig. 3. Power Meter: a visualization of the distribution of the power and relevance of Information Snippets within a document (a). Users can gain an overview of the policies within the ToS (b). When hovering over the Power Meter, the list of Summary Snippets of the corresponding policy is shown (c). Clicking the blue backquote would allow users to view the referenced Information Snippet (d).

## 4.2   Power Meter: Policy level visualization of power and relevance

While prior works have focused on a single policy such as the Privacy Policy [74, 95] or the main ToS page [36, 62], Terms of Services include many policies (e.g., content policy, community guidelines, return policy, and more). To help users form mental models of each sub-policy, TermSight provides information scent and preview to help users gain an overview of the policies within ToS to decide which policies to read (DG1).

*4.2.1   Providing information scent of a policy.* In TermSight, the policies that are part of the ToS agreement are placed in the top navigation panel where users can navigate to different policies. Each policy in the top navigation panel or hyperlinked within the document is accompanied by a **Power Meter**: a horizontal bar visualization representing the distribution of Information Snippets in a document, with colors denoting power and relevance. Shown in Figure 3, Power Meter provides information scent [82] of the power distribution of Information Snippets within the policy and the fraction of Information Snippets that are potentially relevant to the user.

*4.2.2   Providing a preview of the document.* Hovering over a Power Meter reveals a preview popup containing a list of Summary Snippets, allowing users to skim and gain an overview of a policy's Information Snippets without having to be redirected to a new page (Figure 3c). Each Summary Snippet in the popup includes a clickable quote icon that allows users to view the referenced Information Snippet (DG4).

## 4.3   Summary Snippets: Plain language summaries of discrete information snippets

To help users surface and interrogate snippets of information buried in a wall of text (DG2) and make original text more accessible (DG4), TermSight provides Summary Snippets accompanied with Highlight Bars for seamless navigation between the summary and original text.

*4.3.1   Summary Snippets.* Unlike prior methods that rely on document-level overviews [56, 96] or section-level summaries [6, 68], TermSight supports interaction and sensemaking for finer-grained snippets of information within dense sections. TermSight features Summary Snippets - one-sentence plain language summaries of Information Snippets (§4.1.1). Each Summary Snippet is rendered with color based on the Power and Relevance of the corresponding Information Snippet (§4.1.2). To enable users to interrogate original text and verify AI-generated summaries (DG4), clicking on a Summary Snippet jumps a user to the corresponding Information Snippet in the document (Figure 4a).

*4.3.2 Highlight Bar.* To help users surface and identify Information Snippets while reading or skimming the original text (DG2), TermSight introduces Highlight Bars. Highlight bars are positioned to the left of each Information Snippet in the original text, visually segmenting dense sections into Information Snippets. For paragraphs containing multiple Information Snippets, additional line breaks are added after each snippet to add visual separation. Similar to the Summary Snippets, Highlight Bars are color-coded based on the Power and Relevance of the corresponding Information Snippet (§4.1.2). To support navigation from the original document to the summary, users can click on a Highlight Bar to jump to the corresponding Summary Snippet, as shown in Figure 4b.



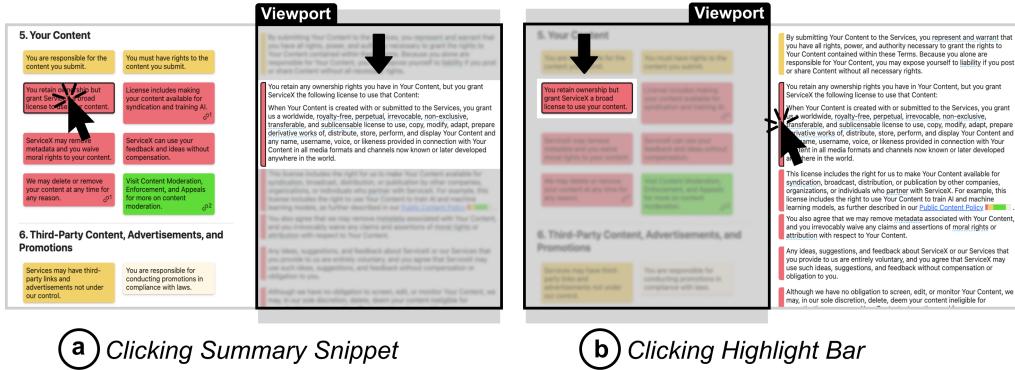(a) *Clicking Summary Snippet*   (b) *Clicking Highlight Bar*

Fig. 4. Viewport Scrolling: (a) By clicking the Summary Snippet, the viewport of the original passage will scroll to the referenced Information Snippet; (b) By clicking the Highlight bar, the left viewport will scroll to the corresponding Summary Snippet.
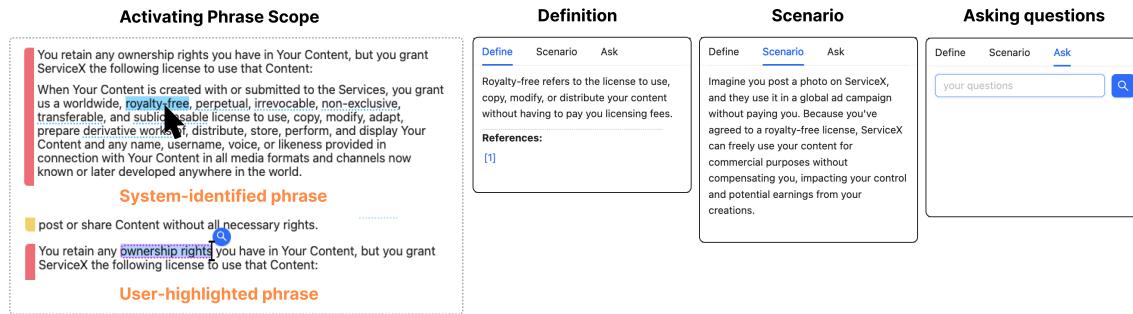


Fig. 5. Phrase Scope: A tooltip that provides term definition and scenario that contextualizes the meaning of the phrase while allowing users to ask clarification questions. Users can hover over the references to view the referenced text in the ToS.

## 4.4 Phrase Scope: Context-dependent definition, personalized scenario, and clarification question

To help users read original text with unfamiliar or vague terminology (DG3), TermSight extends prior work focusing on offering dictionary definitions (e.g., [6]) by introducing Phrase Scope: a set of features that support the identification and interpretation of unfamiliar or vague phrases (Figure 5). First, TermSight identifies words and phrases that might be unfamiliar or vague to general audience readers and underlines these phrases in blue. This aims to provide a starting point for users to explore potentially unfamiliar phrases without relying on users to identify them. Clicking on the

identified phrase opens Phrase Scope, a tooltip with the phrase's definition in context, a personalized scenario to envision potential implications, and an option to ask clarification questions (examples in Table 2). In addition to the suggested phrases, TermSight allows users to highlight any arbitrary span of text to request Phrase Scope.

## 5  IMPLEMENTATION DETAILS

TermSight is rendered as a web interface implemented with Next.js. Below, we explain the input, output, processing, and overall performance of TermSight.

### 5.1  System Input and Pre-processing:

TermSight assumes a set of HTML or markdown source files representing the policies included in the Terms of Service (ToS). For our user study, the research team manually copied the text from documents in ToS on the internet into markdown files with annotated hyperlinks and headers. We assumed clean source files because the focus of TermSight is not on developing web scraping and cleaning technologies but rather on investigating meaningful navigational and reading affordances for ToS. Each document is segmented into text **"chunks"**, composed of one or more paragraphs, which are further vectorized and stored in a vector database (detailed in Appendix B.1).

For the TermSight features that depend on user preferences, such as classifying the relevance of information snippets (§5.2) and generating personalized scenarios (§5.3.3), the system uses a text-based persona that includes users' intended usage of the service (e.g., content consumers vs. content creators) and their values or concerns (e.g., privacy, copyright, etc.). The user personas used for the study are further explained in Section 6.1.2 and displayed in Appendix F.1.
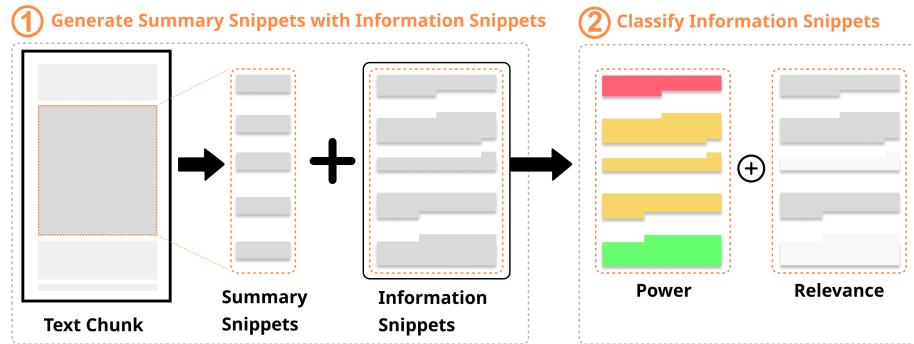


Fig. 6. A flowchart of the implementation for (1) obtaining Summary Snippets and Information Snippets and (2) classifying Power and Relevance for each Information Snippet.

### 5.2  Obtaining Summary Snippets and Information Snippets:

For each chunk of text obtained from the document pre-processing pipeline, we prompted GPT-4o to generate a list of 1-sentence summaries (Summary Snippets) each referenced to a span of the input text chunk (Information Snippet) as illustrated in Figure 6. While the length of the referenced text for each summary is unrestricted, each summary is constrained to a maximum of 12 words based on prior works [27, 91] and prompt engineering. Then, we prompted GPT-4o to classify each Information Snippet along two dimensions: Power and Relevance (to user persona). We detail our prompt design in Appendix B.3.

### 5.3 Phrase Scope

Below, we detail the implementation of each feature in Phrase Scope.

*5.3.1 Identifying unfamiliar or vague phrases.* Prior works have demonstrated the capabilities of LLM in personalized jargon identification [37] and detecting phrases that can potentially be elaborated [31]. Inspired by this work, we employed GPT-4o with two few-shot prompts to identify potentially unfamiliar and vague phrases within the document chunks produced by the pre-processing pipeline (Prompts in Appendix B.4.1). The input chunks are also referred to and used as the **"phrase context"** for each identified phrase.

*5.3.2 Generating phrase definitions or responses to user questions.* To generate in-context definitions for the identified phrases, we used a retrieval-augmented question answering approach [31]. We first query the vectorized document chunks (§5.1) to retrieve chunks that might define the phrase. The database query is framed as a question: *"What does {input phrase} refer to in the sentence: {phrase context}"*. Here, phrase context refers to the chunk of text containing the input phrase as defined in §5.3.1. Then, the query is embedded using OpenAI's text-embedding-3-small model, and the top 15 chunks from the vector database are retrieved based on cosine similarity. These chunks are collectively referred to as **"retrieved chunks"**. Using the phrase context and the retrieved chunks, we prompted GPT-4o to generate an in-context definition with references to the retrieved chunks. When users ask additional questions, the same retrieval-augmented question answering pipeline used for generating definitions is applied, with the only difference being the question asked. We detail our prompts in Appendix B.4.2.

*5.3.3 Generating scenarios for phrases.* Scenarios intend to complement phrase definitions, helping users interpret and envision the potential implications of a phrase in context. In addition, customizing scenarios based on user personas intends to ensure that the generated examples are relevant to the user's intended usage of the service (e.g., content creator or consumer) and personal value (e.g., privacy, copyright). Building on prior works that demonstrated the effectiveness of large language models (LLMs) in generating unintended consequences [99] and stories to teach legal concepts [50], we employed GPT-4o with zero-shot prompting (Figure 27) to generate the scenarios with the phrase context, phrase definition, and user persona.

### 5.4 Evaluation of System Output

Before reporting on our user study, we randomly sampled and conducted a manual evaluation of TermSight's core outputs. Our goal was not to assess advances in system performance, but to verify that the system can produce meaningful outputs that could support ToS reading. Across 116 sampled information snippets, we found 5 imperfect classifications of power or relevance due to the lack of full context in the input information snippet. Similarly, we reviewed 113 generated definitions and scenarios each. One scenario was found to be factually incorrect due to hallucination. All the definitions were correct except for 4 overly general definitions for service-specific phrases not explicitly defined anywhere in the ToS. More details are provided in Appendix C.

## 6 USER STUDY

We conducted a counterbalanced within-subject experiment with 20 participants to evaluate TermSight. We asked the following research questions:

> **RQ1**: How did participants perceive TermSight and its features?
> **RQ2**: How did participants read with TermSight and its features?

**RQ3**: How did TermSight influence comprehension and recall?

## 6.1 Study Design

*6.1.1 Study Conditions/Interface Variants.* To evaluate TermSight, we conducted a within-subject experiment where each participant used both TermSight and a baseline interface to read 2 services' ToS contracts, one with each interface variant. The order of the interface variant and service type were counterbalanced to reduce ordering effects. The baseline interface intended to mimic a standard ToS. The baseline interface had the same layout as TermSight without the Power Meter, Summary Snippet, Highlight Bar, and Phrase Scope. In place of the Summary Snippets on the left was a table of contents. Participants could click a section header in the table of contents to navigate to the corresponding section.

*6.1.2 Materials: ToS and User Persona.* For the user study, we used the ToS of one social media site (Reddit) and one e-commerce site (Poshmark) because social media and e-commerce platforms are among the most visited digital services and are representative of the standard long-form contract used in prior studies [34, 54, 77, 92, 93]. For each service, we manually collected policies that are linked or hyperlinked as part of the ToS. We did not include policies that are location-specific (e.g., a California Privacy Notice). For Reddit, we also did not include policies for advertisers, brand&press, or governmental agencies as they were not relevant to a lay user and to keep the number of policies for both services similar. We collected a total of 14 policies for Reddit and 15 for Poshmark. The service names for Reddit and Poshmark were anonymized as ServiceX and ServiceY. Features of TermSight rely on a user persona to determine relevance. For the user study, we designed two personas based on users' information needs identified in the formative study: a content consumer who posts personal content on social media sites (Figure 17) and a buyer who rarely posts reviews on E-commerce sites (Figure 18).



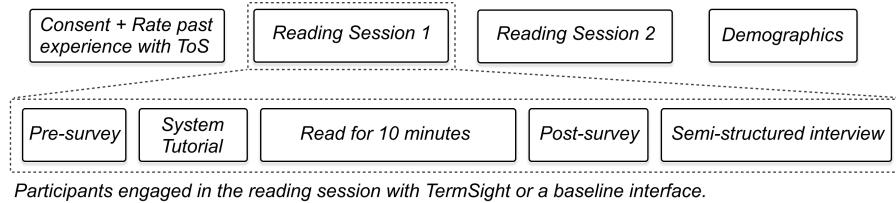*Participants engaged in the reading session with TermSight or a baseline interface.*

Fig. 7. The experiment procedure for the within-subject study. At the top, we show the procedure for the entire study. At the bottom, we show the process diagram for a reading session with either TermSight or Baseline. The orders for the reading session were counterbalanced in the study.

*6.1.3 Study Procedure.* The entire study took 90 minutes and participants were compensated 25$. The study was approved by the Institutional Review Board (IRB) office at our organization. The study was composed of 2 reading sessions, 45 minutes each (Figure 7). After obtaining consent and before the reading sessions, we asked participants about their past experiences with ToS. Within each reading session, participants first completed a pre-survey and reported their past experiences interacting with social media or e-commerce platforms. Then, an anonymized description of the service and user persona was given. Participants were given 10 minutes to read the ToS without having to speak aloud. Reminders were given at the 5 and 1-minute mark. After the reading session, participants completed seven 5-point Likert-style rating questions about their reading experiences followed by one free-response recall question and six

multiple-choice comprehension questions. We also conducted semi-structured interviews with participants about their experiences and reading strategies for 10 minutes. The same procedure was repeated for the second reading session. The study materials are included in Appendix F and the supplemental materials.
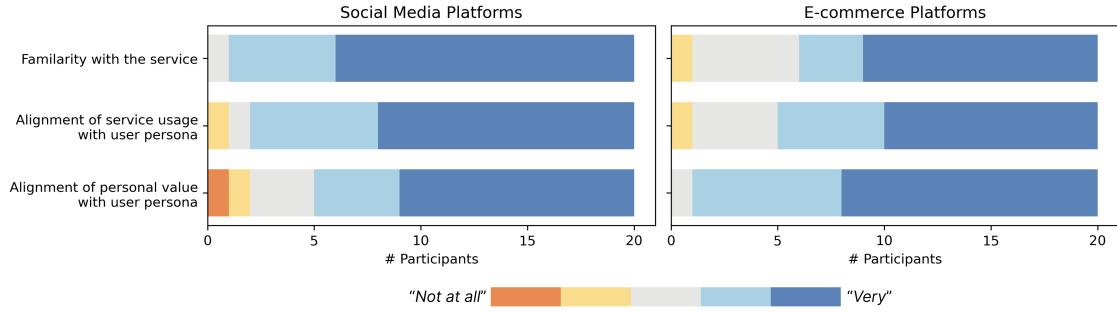


Fig. 8. Participants' self-ratings of their general familiarity with social media and e-commerce platforms and their alignment with the user persona given for each service. Most participants were familiar with both social media and e-commerce platforms in general and found the given persona to align with their personal usage of the service and personal value.

## 6.2 Participants and Setting

We recruited 22 participants through Prolific. The recruitment required participants to be over 18, fluent in English, located in the US, not have color deficiencies, and not be legal professionals. We did not require participants to have prior experience in reading ToS. Two participants encountered technical difficulties, and their data were removed from the analysis. Of the remaining 20 participants, 13 self-identify as female and 7 as male. 1 participant had a high school degree, 3 had an associate degree, 12 had a bachelor's degree, and 4 had a master's degree. 4 participants were between the ages of 18-25, 5 were between 26-35, 5 were between 36-45, and 6 were between 46-55. When asked about how many ToS they've read, 4 participants read none, 5 read between 1—3, 5 read between 4—6, 1 read between 7—9, and 5 read greater than 10. The median of participants' self-rated familiarity with ToS was 3.5 ($\sigma = 1.5$). Most participants were generally familiar with social media and e-commerce platforms and found their given user persona for both services to highly align with their personal usage of the service and their personal value (Figure 8). None of the participants had read or remembered the ToS for Reddit and Poshmark prior to the study.

## 6.3 Measures

*6.3.1 Perceived reading experiences.* We collected 5-point Likert-style ratings (1="Not at all," 5="Very") of participants' reading experiences (Appendix F.3). The experience measures included ease of reading, perceived understanding, and confidence in obtaining relevant information [6]. In addition, we included two 5-point Likert-style rating questions related to the ease of deciding which policy to read and what text to read within a policy, both of which are challenges identified in our formative study. Lastly, we included two 5-point Likert-style questions about their willingness to (1) read other ToS with the interface and (2) spend more time on their given ToS with the interface and get a link to the ToS after the study [5].

*6.3.2 Reading comprehension and recall.* For reading comprehension, we wrote 10 multiple-choice questions for each service and selected 6 questions that were relevant for the given user persona for each service, similar to prior works

on evaluating reading comprehension for ToS and privacy policy [6, 92, 97]. The questions are different for the two services to minimize learning effects. Comprehension was measured as the number of questions participants got right out of 6. Examples of comprehension questions and their answers are listed in Table 1.

For recall, participants were asked to write down information in ToS they found interesting, surprising, or any detail they remembered: *"What did you learn from reading the ToS? Recall one or more interesting things you learned from the ToS."* To analyze these free-form responses, we manually broke the response into single references to a clause in the ToS. Then, we labeled whether each reference was correct based on the ToS. Recall is measured as the number of correct facts participants wrote in the response [34, 93].

*6.3.3   Feature usage.* To understand users' interaction behavior, we logged all viewport scrolls and feature usage during the reading session along with timestamps. This allows us to measure the frequency of feature usage of TermSight and the baseline interface.

## 6.4   Analysis

In this section, we introduce our analysis framework for the user study. We used a causal framework (§6.4.1) to understand the effects of the treatment and a Bayesian analysis (§6.4.2) to estimate the treatment effect. Then, we discuss the thematic analysis (§6.4.3) used to analyze the qualitative responses from the semi-structured interviews.

*6.4.1   Structural Causal Framework.* We used a structural causal framework popularized by Pearl [81] to better understand the effects of the treatment. A structural causal model framework comprises exogenous variables $U$, endogenous variables $V$ and a set of functions $F$ that represent the physical process that results in the values of the endogenous variables. For example, the question $v_i = f(v, u)$ implies that the values of $v$ and $u$ completely determine the value of $v_i$ through the function $f$. The structural causal model can be represented as a Directed Acyclic Graph (DAG) where the nodes represent the variables and the edges represent the causal relationships among the variables. Thus, in the preceding example, $v$ and $u$ are the parents of $v_i$ in the DAG. The structural causal model framework is non-parametric, but frequently used with generalized linear models to estimate the functions $f$. The structural causal framework has connections with structural equation models [100] used in the Economic Sciences, though the latter were developed with linear relationships in mind. Figure 9 shows the DAG for the study.

*6.4.2   Bayesian Analysis.* We used Bayesian inference to estimate the treatment effect of TermSight on the outcomes of interest. While the use of Bayesian estimates is growing in HCI [53, 58], we briefly justify its use over traditional non-Bayesian methods. Standard statistical methods including t-tests of significance can be powerful, but they require us to test for the conditions under which a test is valid (*e.g.,* Normality for the t-test). Prior work in Biological and Psychological Sciences has shown that Null Hypothesis Statistical Testing (NHST) including $p$-values are problematic, and must be used with caution. They depend on the researchers' intentions, when they stop collecting samples, and pre-selection of hypothesis [21, 46, 47].

We believe that Bayesian techniques offer several advantages over non-Bayesian methods. First, as Kay et al. [53] point out, a Bayesian framework leads to an accumulation of knowledge within HCI, where the posterior of the parameters in a prior work can serve as the prior in the current experiment. Second, a Bayesian model is transparent—the researcher will foreground all the assumptions in the analysis via their model, allowing for a principled critique of the model. Third, use of a Bayesian framework shifts the discussion from "did it work" to "effect size of the intervention" [53]. While NHST techniques can compute the confidence interval, these intervals are susceptible to misinterpretation [10, 42] and
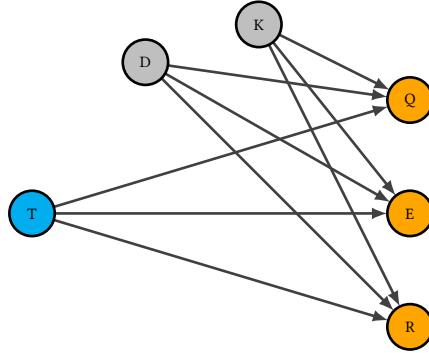
Fig. 9. The causal graph (a DAG) has three types of nodes: Treatment ($T$) colored blue, pre-treatment covariates $D$ (demographics), $K$ (knowledge of Terms of Service) colored gray, and measured outcomes $Q$ (comprehension), $R$ (recall of correct facts), and $E$ (experience of the interface) colored orange. Since we use randomized assignment, including the interface type (treatment vs. control), type of ToS (social media vs. e-commerce website), and order of presentation (first vs. second), the treatment $T$ is independent of the pre-treatment covariates $D$ and $K$. Thus, to estimate the causal effect of the treatment $T$ on the outcomes $Q$, $E$, and $R$, we can use a simple regression model, without conditioning on the pre-treatment covariates $D$ and $K$ since there is no backdoor path from the treatment $T$ to the outcome variables $Q$, $E$, and $R$.

importantly, underemphasized in favor of $p$-values. Finally, as McElreath [71] points out, a Bayesian model with use of maximum entropy priors for the parameters (*e.g.,* the normal distribution) is *the most conservative* given the evidence to estimate of the effect of the treatment.

*6.4.3   Thematic analysis.* To analyze the qualitative responses from the semi-structured interviews, we adopted the six phases of reflexive thematic analysis suggested by Braun and Clarke [18], following the same procedure as our formative study (§3.1.2).

## 7   FINDINGS

### 7.1   RQ1: How did participants perceive TermSight and its features?

When asked about which version of the interface they prefer, 19 participants preferred TermSight. One participant preferred the baseline interface, citing difficulty navigating TermSight's two scrollable columns on a small screen with a trackpad that lacks scroll functionality. Below we present our quantitative findings on the experience outcomes (§7.1.1), followed by qualitative findings organized around recurring themes about features of TermSight.

*7.1.1   Quantitative Findings: TermSight improved user experience measures.* Figure 10 shows participants' ratings of the 7 reading experience measures for both interfaces. Our Bayesian model analysis (model details in Appendix D.1) reveals a significant effect of the TermSight interface on the reading experience measures. The forest plots shown in Figure 19 in the Appendix demonstrate that the 94% Highest Posterior Density Interval (HPDI[1]) for each treatment–control pair for every user experience measure does not overlap with each other. This indicates a significant positive effect of the treatment on all 7 user experience measures.

---

[1]In Bayesian analysis, the 94% HPDI refers to that smallest interval (this interval is unique) of the posterior distribution that has 94% of the probability mass. It is common in Bayesian analysis to use intervals distinct (97% and 89% HPDI intervals are also used) from the typical 95% value to avoid the confusion of the frequentist confidence intervals.
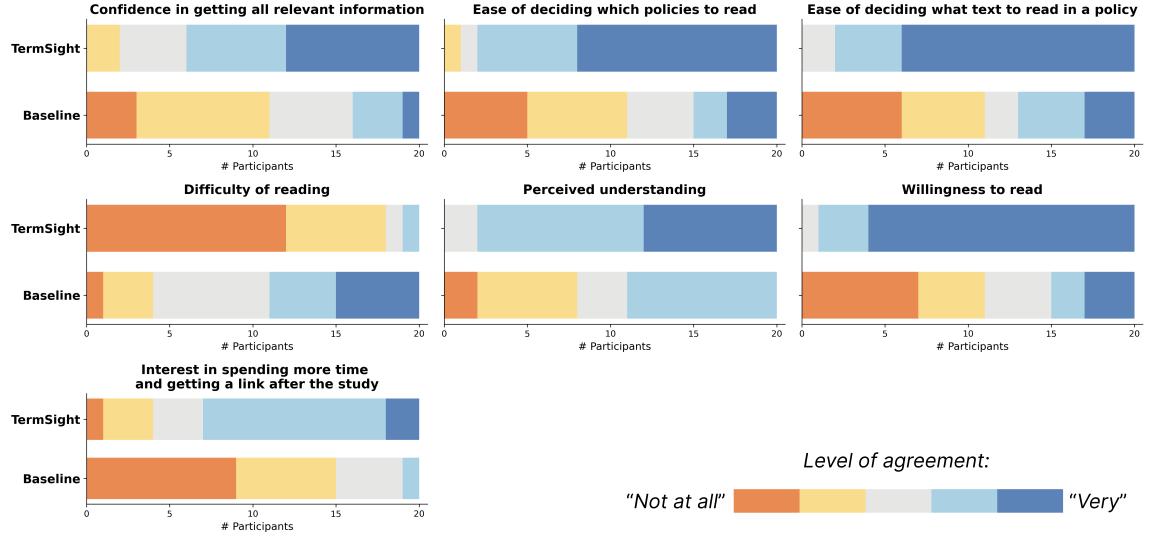
Fig. 10. 5-Point Likert Ratings of participants' reading experiences: Participants found it to be easier to find and understand relevant information. Participants were also more willing to spend more time on the ToS with TermSight and read other ToS with TermSight. The 7 experience questions are specified in Appendix F.3.

In Figure 11, we show the contrast between the treatment and control conditions averaged across 7 user experience measures as well as the posterior of the effect size. The posterior distribution (Figure 11b) of the effect size (Cohen's $d$) of the treatment on the user experience is centered around 2.7, with an HPDI of [2.3, 3.2], indicating a large effect size[2]. Because the HPDI doesn't overlap with the ROPE, the effect of TermSight on user experience is significant. Below, we further present recurring themes related to participants' perceptions and feedback on features of TermSight.



(a): Posterior distributions of the coefficients for the treatment and control conditions.

(b): Posterior distribution of the effect size of the treatment on the user experience.

Fig. 11. Figure 11a shows the posterior distributions of the coefficients for the treatment and control conditions. Note that the posterior distribution of the treatment condition is shifted to the right, indicating a positive effect of the treatment on the user experience. Figure 11b shows the posterior distribution of the effect size of the treatment on the user experience centered around 2.7, suggesting a large effect size. Additionally, the HPDI doesn't overlap with the ROPE, suggesting a significant effect of the treatment on user experience.

---

[2]This is the effect size measured not on the outcome Likert scale (1–5) but on latent scale of the model coefficients.

*7.1.2  Power Meter provided information scent of the policies and reduced navigation cost.* All 20 participants described that the Power Meter helped them gain an intuition of the content of a policy and provided guidance. Participants described how the saturated colors, especially red and green, *"piqued interest"*(P4). For example, P12 described how Power Meter provided guidance and supported decision-making on which policy to pay attention to.

> *"The color system was really guiding a lot of my decisions on which policies to read. The vibrant red was always what I was gonna try to get to 1st and the things that I cared the most about. So, finding which terms I'm agreeing to that I have the least amount of agency."* (P12)

P6 used the TermSight system first. However, he lamented the fact that when using the baseline interface without the Power Meter, deciding which policy to read becomes a guessing game.

> *"Policies like terms of service or privacy policy are kind of broad. Does that apply to something I would actually care about? I don't know. So I'd have to click each one and have to go through it. It was more like a guessing game."* (P6)

8 participants described how the document preview that appears when hovering over the Power Meter complemented and provided a more detailed overview of the policy to help them decide what to read. Additionally, 9 participants mentioned how the document preview allowed them to access information from linked pages or policies without *"clicking on the link and possibly losing where I am on"* (P1).

*7.1.3  Summary Snippet and Color provided guidance and helped surface power and relevance within documents.* 18 Participants described how the color of the Summary Snippets helped surface where power lies in the contract, which is challenging to "see" when presented in plain text with the baseline interface.

> *"The color coding gave me a better understanding of who's benefiting more, how is it benefiting them, and how it affects me. Versus [in the baseline interface] it feels like it only really benefits the company itself."* (P18)

Moreover, all participants described how the color and summary allowed them to make decisions and prioritize what to focus on strategically.

> *"I really liked having the summary with colors. I had a hierarchy of how I was going to read stuff. I made sure I read all the [saturated] red ones, and then the [saturated] green. If I had time, I would skim the [saturated] yellow and see if anything stuck out. It would take me a lot longer to filter for information otherwise."* (P6)

While participants described focusing on the more saturated colors, especially red and green, 6 participants also described skimming the less saturated ones to validate and make sure they didn't miss information.

*7.1.4  Summary Snippet simplified and broke down dense text to support sensemaking.* All participants highlighted how the Summary Snippets *"broke down the entire section into a few bullet-pointed summaries"* (P11), which helped them absorb the information and surface relevant pieces of information hidden in dense text. 9 participants further described how the Summary Snippets simplified the text and served as a scaffold to encourage reading the original text.

> *"Not only does it [TermSight] bring me up to things that I would normally miss. but then that would be an introduction or a guide, or a push towards reading the whole 3 paragraphs. ... It is really guiding me through this entire document, and it's simplifying it."* (P2)

Similarly, 4 participants mentioned how the highlight bars and line breaks added for each information snippet made it easier to consume the original text without having to *"mentally separate the concepts"* (P7). P20 described the experience of using Summary Snippets to make sense of the big ideas of the section while being able to dive deeper as forming *"a*

*web of thought and ideas: what everything kind of means together and piece by piece"*. Not only do the Summary Snippets serve as an entry point to the original text, 4 participants described how the Summary Snippets helped them check their understanding after reading the original text.

*7.1.5    Term definition and scenario help consume original text and envision implications.* 14 participants described that the underlined phrases attracted their attention, made them feel like these were phrases that *"should be known and understood"* (P18), and prompted them to check on their understanding of the phrase.

In addition, participants commented on how the definitions and scenarios made legal language more approachable (14/20).

> *"It really just made it a lot easier to understand, because most of the time they write it in legal terms that most people don't know. People don't speak legal language. So this broke it down easier for most people like the Layperson."* (P17)

12 participants highlighted how the scenarios helped interpret the meaning and envision potential implications of the phrase in the context of the user's intended usage of the service and was *"easier to resonate"*(P17).

> *"I like that it provides a real-world example of when and how [the phrase] would be relevant to someone. I think it actually does better explaining the concept than the definition. Because it's not just the definition. It's the definition in context."* (P11)

Moreover, one participant (P1) noted how the scenarios helped her come up with more questions, which she asked with the Ask function. Through this process, TermSight helped P1 to quickly obtain information from correlated policies.

## 7.2    RQ2: How did participants read with TermSight and its features?



Fig. 12.  Usage of features during the reading session. Each dot represents the number of times a feature is used by one participant during a reading session. Participants used most features of TermSight. In contrast to the baseline, participants when using TermSight more frequently clicked the Summary Snippets as opposed to the section titles to navigate to the original text.

*7.2.1    Feature Usage.* In the baseline interface, participants on average navigated to 7 policies ($\sigma = 5$). Participants' reading behaviors were similar to our findings from the formative study (§3.3.2). Participants described how they were

Fig. 13. Minute-by-minute usage of features of TermSight during the ten-minute reading task. Participants used these features throughout the reading session, rather than being limited to the beginning or end. The Summary Snippet is the most frequently used feature throughout.
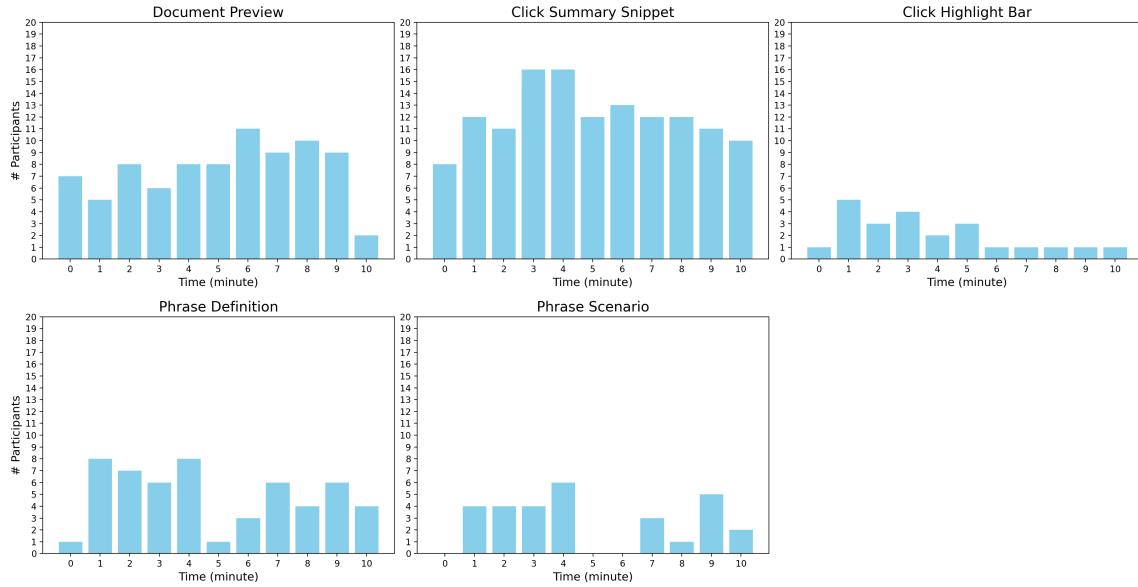
trying to skim everything from top to bottom. When they encounter sections that were hard to read, they tended to skip. Only 11 participants clicked the section heading in the table of contents at least once to navigate to targeted sections.

As shown in Figure 12 and 13, participants used most features of TermSight throughout the session, demonstrating that TermSight provided value throughout the session. When using TermSight, participants on average navigated to 6 policies ($\sigma = 3$). 19 participants used the document preview of Power Meter to gain an overview of documents in ToS for more than 1 second. 13% (17/135) of the total usage was for previewing documents hyperlinked inline, demonstrating that the Power Meter provided value also for inline hyperlinks.

Summary Snippet was the most frequently used feature, 19 participants clicked the Summary Snippet to navigate and read the original text at least once. On average, they clicked 17 Summary Snippets ($\sigma = 13$). In contrast, participants rarely clicked the section headings to navigate to the original text (4 clicks in total across 3 participants). This demonstrates that Summary Snippets were potentially more useful as navigational affordances to surface and interrogate relevant information from the original text compared to section headers. We also noticed how 10 participants clicked the Highlight Bar to read the summary snippet after reading the original text, suggesting how these participants were back-checking their understanding using the summary snippets.

For Phrase Scope, 16 participants viewed the term definitions and 12 participants viewed the term scenarios. Participants on average used term definition 3 times ($\sigma = 4.3$) and term scenario 1 time ($\sigma = 2.7$). Out of the 77 instances when Phrase Scope was accessed, 72 were suggested by TermSight. Participants rarely manually highlighted phrases to activate phrase scope, asked questions, or checked the references provided for the definitions. While only 1 participant asked questions using Phrase Scope, P1 asked 8 questions in a 10-minute reading session.

While TermSight provided multiple reading affordances, participants were in control of deciding whether or not to rely on these features. As opposed to a consistent preference for one interaction strategy over others, we observed diverse ways that participants relied or did not rely on TermSight (Figure 14). Below, we provide case studies of these non-exclusive emergent reading strategies.

*7.2.2 Summary driven reading.* 19 out of 20 participants at times relied on the summary snippets as a primary entry point for engaging with the text. One notable example of this behavior was P3, who mainly read the Summary Snippets and clicked 11 of them to interrogate the original text (Figure 14a). P3 described how colors were driving his attention to read. P3's decision to click the Summary Snippets and interrogate the original text relied on the color coding and evaluation of the information loss of the Summary Snippets.

> *"I mainly stay on the left side. I read them if it was like a saturated red or green. I clicked on them because I felt like there was more information for me [behind the Summary], and I want to learn more about it."* (P3)

*7.2.3 Original text driven reading .* At other times participants centered their interactions around the original text. 9 participants went through at least one policy focusing mainly on the original text. These participants described how they wanted to first read the actual text before relying on AI-generated summaries due to the potential imperfections of AI. These participants explained, *"I didn't want my opinion to be affected by AI"* (P5) and *"sometimes summaries don't have all of the information from my other experiences"* (P18).

When reading the Privacy Policy, P8 focused on the original text as shown by the frequent scrolling of the original text in Figure 14b. In addition, P8 clicked the Highlight Bar 11 times to refer back to the Summary Snippets. P8 described how he knew AI was not perfect from his prior professional experiences in AI. As a result, P8 took the colors and summary snippets as a guide rather than a guarantee.

> *"I'm not gonna just go and trust it without doing my due diligence. I would take it as a guide. I won't take it as a guarantee, because AI is not perfect. I think it [AI] would try to label it the best it can and I think it did, to be honest."* (P8)

*7.2.4 Calibration of trust.* 6 participants described that they went through a calibration process with AI-generated features such as the Summary and Color at the start of the reading session. For example, P6 described how his reading behavior was different during the calibration process and after. P6 described reading the original text on the right side to evaluate the colors and summaries at the beginning of the reading session (original text driven reading). Afterward, P6 relied on the summaries a lot more by scrolling and skimming the left panel and clicking on the summary snippets to dive to the original text as seen in Figure 14c (summary driven reading).

> *"So like early on in the reading, I kind of read the whole thing on the right first. Then, I looked at the summaries, and once I realized that it did do a good job summarizing it with colors. I trusted it enough to trust it for the rest of the time."* (P6)

*7.2.5 Switching strategy based on time and document.* ToS includes multiple documents, and in the study, there were 14-15 policies with varying lengths. We noticed how participants' reading behavior can vary based on document length (P5, P19). For example, P19 mainly read the summary for longer policies such as User Agreement and Privacy Policy (Summary driven reading). For shorter policies such as Buyer Policy, P19 only read the original text (Original text driven reading).

Fig. 14. Scrolling behavior of participants that employed each non-exclusive reading strategy. TermSight features two scrollable columns: Summary Snippets on the left (red lines) and Original Text on the right (blue lines). Clicking a Summary Snippet (orange X) auto-scrolls the Original Text panel to the referenced text, while clicking the Highlight Bar (light blue X) auto-scrolls the Summary panel to the corresponding Summary Snippet. (a) Summary-driven reading occurs when participants primarily read and scroll the Summary Snippets, occasionally clicking on them to interrogate the original text (19/20). (b) Original text-driven reading occurs when participants first read the original text and use the Summary Snippets as a supplement (9/20). (c) The calibration of trust mainly occurs at the beginning of the session where participants would compare and contrast the original text to the Summary Snippets with color (6/20).

We also noticed how time can change participants' reading strategy (P7, P8). For example, P7 spent 8.5 minutes and was very interrogative when reading the User Agreement policy. When she realized that she was running out of time, she skimmed through the summary snippets for 5 more policies without interrogating the original text.

### 7.3 RQ3: How did TermSight influence comprehension and recall?

Below we present the results of the comprehension quiz (§7.3.1) and recall task (§7.3.2).

*7.3.1 Comprehension.* For the comprehension quiz, out of 6 questions, participants on average scored 2.00 ($\sigma = 1.03$) when using TermSight and 2.15 ($\sigma = 1.23$) when using the baseline interface. Our Bayesian model analysis for the comprehension quiz (model details in Appendix D.2) reveals no significant differences in the comprehension scores across interface conditions. As shown in Figure 15, the posterior distributions for both conditions were nearly identical and centered around zero. When the question ID and service type are fixed, there are significant overlaps between the 94% HPDI for each treatment–control pair as shown by the forest plots in Figure 20 in the Appendix, suggesting no significant differences.
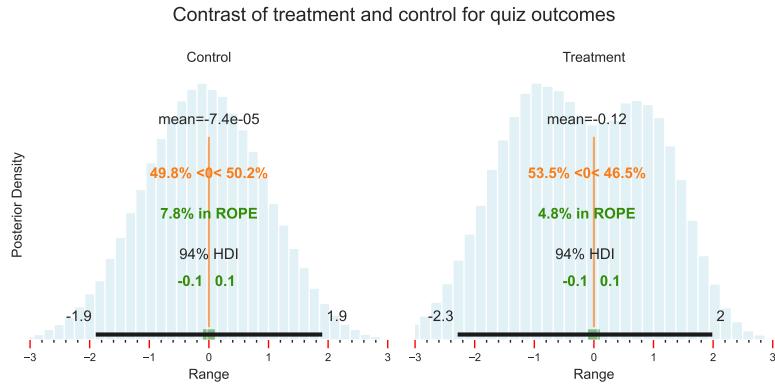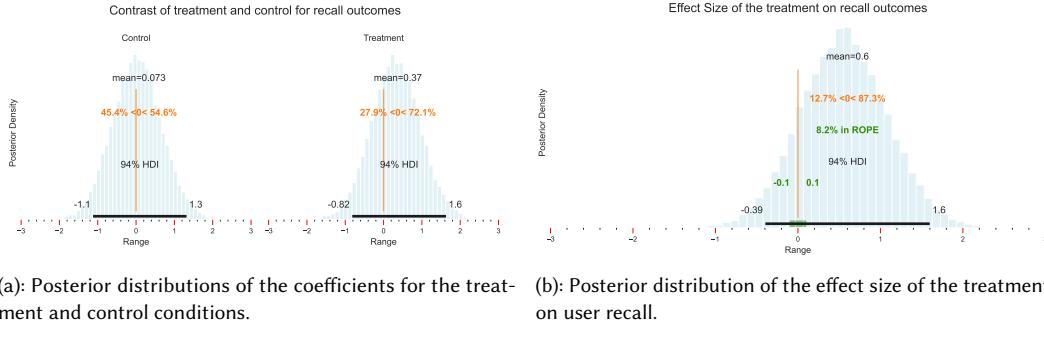


Fig. 15. The figure shows the posterior distributions of the coefficients for the treatment and control conditions. Both posterior distributions are nearly evenly split around zero, indicating no effect of the treatment on the comprehension quiz questions.

*7.3.2 Recall.* For the recall task, participants using TermSight recalled 39 correct facts and 3 incorrect facts in total. On the other hand, participants using the baseline interface totally recalled 26 correct facts and 7 incorrect facts. Our Bayesian model analysis for the recall task (model details in appendix D.3) reveals no significant differences in the recall scores across interface conditions. Figure 16b shows the posterior distribution of the effect size of the treatment on recall ($\mu = 0.6$, HPDI = [-0.39, 1.6]). Despite the mean being 0.6, suggesting a medium effect size, the HPDI interval has an overlap with the ROPE, implying that there is no significant effect of the treatment on the user recall.

## 8 DISCUSSION

In this paper, we set out to design a new reading experience that scaffolds readers' ability to find and understand information important to them within Terms of Service (ToS) contracts. Based on formative interviews and observations, we realized this vision in TermSight. Results from our evaluation of TermSight (N=20) suggest that TermSight made it easier to read, navigate, and contextualize information within ToS relevant to users. Especially interesting was how

(a): Posterior distributions of the coefficients for the treatment and control conditions.



(b): Posterior distribution of the effect size of the treatment on user recall.

Fig. 16. Figure 16a shows the posterior distributions of the coefficients for the treatment and control conditions. Specifically, the posterior distribution of the treatment condition is shifted to the right, indicating a positive influence of the treatment on Recall. Figure 16b shows the posterior distribution of the effect size of the treatment on the user recall centered around 0.6, suggesting a medium effect size. Yet, the HPDI has an overlap with the ROPE of [-0.1, +0.1], implying no significant effect of the treatment on the user recall.

TermSight afforded multiple different reading strategies (e.g., summary-driven reading, §7.2.2), all of which helped participants overcome the key difficulties of contract reading we observed in our formative study (§3). However, similar to prior works on simplifying ToS or other complex documents [6, 34, 57, 85, 92, 97], TermSight did not worsen or improve the comprehension or recall of the content. Below, we discuss the implications of our findings and the TermSight system.

## 8.1 Existing ToS are not designed to be read

Many existing ToS have static or interactive design elements, such as policy names, section headers, and overview summaries [97]. However, participants in our formative study found these features vague, misleading, and ineffective for surfacing and interpreting relevant information (§3.3). The HCI literature often refers to these ineffective and deceptive designs as dark patterns in the user interface that manipulate behavior or obscure information to serve the interests of the service providers [15, 59, 69].

TermSight offers one possible design to overcome these dark patterns by offering contract-level guidance (e.g., Power Meter and document previews), document-level guidance (e.g., color-coded Summary Snippets), and phrase-level guidance (e.g., phrase identification). Participants in our usability study (§7) affirmed these design goals in both self-report and behavior metrics. Participants reported that features of TermSight helped them decide which policy to read (Power Meter, Document Preview) and what text to read within a policy (Summary Snippet). Moreover, participants highlighted how TermSight made the original text more approachable with Summary Snippets that broke down and simplified dense sections and Phrase Scope that clarified unfamiliar phrases and helped envision implications.

Future systems could take these designs further by adapting them to other types of contracts, such as leasing contracts, medical contracts, or employment contracts. For example, Summary Snippets may help readers surface and interpret relevant or predatory clauses in leases, while phrase scope may help patients envision the implications of clauses in the medical contract through patient-specific scenarios. Lastly, while we designed TermSight as a contract reading interface, TermSight can be readily extended as a browser extension that highlights relevant and important clauses based on user-defined personas to provide in-situ legal awareness and warnings when encountering contracts.

## 8.2 Envisioning consequences to legal text

Prior work has explored providing term definitions to aid comprehension of academic text [6, 41]. Yet, participants in our study mentioned that definitions alone were not always sufficient for legal language. Participants noted that it can be difficult to contextualize abstract phrases in real-world scenarios or anticipate potential consequences—especially those they have not yet personally experienced. Our formative study further revealed how people's intended usage of the service, personal values, and worries can vary. While prior work explored predicting unintended consequences for scientific endeavors and AI [80, 99], TermSight explored the potential value of presenting personalized scenarios to illustrate the possible implications of legal phrases before a user agrees to a ToS. Some participants even noted how the scenario was more helpful than the definition (§7.1.5). By grounding abstract terms in concrete, relatable examples, these scenarios may help users better reflect on the potential implications behind legal language. Future work could investigate additional ways to help consumers reflect on unintended consequences, such as by showing multiple scenarios or relevant legal cases.

## 8.3 Design for calibration and trust

LLM output can be imperfect due to hallucinations or confusing output [43, 67, 70]. It is therefore important to help users form appropriate trust and reliance for fostering effective Human-AI collaboration [22, 87]. In this vein, prior reading interfaces integrating AI have introduced cognitive forcing functions [20], such as showing summaries only on demand [6] or at the end of the reading session [23], to encourage readers to first read the original document. Yet, it is unclear if these designs are always optimal for intelligent reading interfaces, since they can increase user frustration and may limit specific reading behaviors [19]. For example, these forcing functions, if applied to TermSight, might introduce significant friction for participants taking a summary-driven reading approach (§7.2.2).

TermSight took the approach of tightly linking the AI-generated summaries to the original text to allow users to compare and contrast the two. In TermSight, users can navigate between and compare the summary and the original text in one click. When given these affordances, participants took advantage of them. Some participants explicitly evaluated the Summary Snippets to calibrate their trust. Interestingly, one participant framed the evaluation of AI-generated output as his 'due diligence'. Even those who took a summary-driven reading approach frequently interrogated the referenced original text as opposed to solely reading the summaries. Future work can explore how to help users calibrate trust with the system and evaluate AI output. This could include providing explanations of how the reading interface works in general, providing examples of AI imperfections, and providing explanations for individual features [87].

## 8.4 Ethical Considerations and Social Implications

Features of TermSight (e.g., Summary Snippets, Classification of Power and Relevance, Term Definitions, and Scenarios) rely on the capabilities of LLMs. However, LLMs have been shown to produce imperfect outputs and factually incorrect hallucinations [43, 67, 70]. This introduces risks when interacting with LLM-generated text, as they are not part of the contract and are not legally binding. For example, incorrect summaries might cause readers to misinterpret the clauses in the contract that are binding. On one hand, LLMs have already shown strong capabilities in generating summaries [66, 105], term definitions [49, 103], and scenarios [29, 50, 88, 99]. Work is also actively investigating how to detect non-factual LLM output [40] and how to guide LLMs to generate factually correct information [98]. As a result, LLM output would ideally become more factually correct over time. On the other hand, for future deployment of systems like TermSight, we still believe a validation mechanism for LLM output is necessary. For example, future implementation

of TermSight's features (e.g., term definitions and scenarios) can be accompanied by additional sources of information when available (e.g., news articles, court cases, or information from other ToS). Furthermore, TermSight could facilitate end-user auditing [48, 60], allowing users to report inaccuracies. These user contributions could collectively help validate and refine model outputs.

## 8.5 Limitations and Future work

The generalizability of TermSight and our study findings is limited by the challenges associated with handling ToS documents and our sample of participants in the user study. TermSight currently relies on a set of precleaned Markdown or HTML files. Future works can investigate more advanced scraping and cleaning methods. In addition, our recruited participants were mostly college-educated and fluent in English. Marginalized communities, such as neurodivergent readers and individuals with limited English literacy, can face additional barriers when interpreting legal documents. For example, readers with ADHD might have a hard time staying focused when reading visually dense text such as ToS [12, 90]. Future works can investigate how to help marginalized populations make sense of legal contracts by evaluating TermSight with these populations of users or integrating additional supports such as translation.

## 9 CONCLUSION

ToS are service contracts outlining the rules and expectations for both the service provider and the consumer. Our formative study revealed ineffective and deceptive designs at all granularities of a ToS. To make ToS approachable for consumers, we designed and evaluated an intelligent reading interface: TermSight. Participants reported that TermSight reduced the difficulty of engaging with ToS and increased their willingness to do so. Participants highlighted that features of TermSight afforded them the ability to surface and interpret relevant information at all levels of ToS reading. They found the Summary Snippets with colors particularly helpful in breaking dense sections into snippets of simplified information, which helped them surface relevant information and power balance. The 'duty to read' originated when contracts were brief and based on mutual understanding and negotiation [11, 51]. Modern ToS contracts, however, prioritize unilateral control over mutual agreement by design. TermSight presents a vision and serves as a technology probe into making legal contracts more accessible to the general public, providing consumers with the tools to make more informed legal decisions.

## REFERENCES

[1] 2024. *PrivacySpy*. Retrieved Jan 10, 2025 from https://privacyspy.org
[2] 2024. *Terms of service; didn't read*. Retrieved April 2, 2024 from https://tosdr.org/
[3] Wasi Ahmad, Jianfeng Chi, Yuan Tian, and Kai-Wei Chang. 2020. PolicyQA: A Reading Comprehension Dataset for Privacy Policies. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. 743–749.
[4] Ryan Amos, Gunes Acar, Eli Lucherini, Mihir Kshirsagar, Arvind Narayanan, and Jonathan Mayer. 2021. Privacy policies over time: Curation and analysis of a million-document dataset. In *Proceedings of the Web Conference 2021*. 2165–2176.
[5] Tal August, Kyle Lo, Noah A Smith, and Katharina Reinecke. 2024. Know Your Audience: The benefits and pitfalls of generating plain language summaries beyond the" general" audience. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–26.
[6] Tal August, Lucy Lu Wang, Jonathan Bragg, Marti A Hearst, Andrew Head, and Kyle Lo. 2023. Paper plain: Making medical research papers approachable to healthcare consumers with natural language processing. *ACM Transactions on Computer-Human Interaction* 30, 5 (2023), 1–38.
[7] Yannis Bakos, Florencia Marotta-Wurgler, and David R Trossen. 2014. Does anyone read the fine print? Consumer attention to standard-form contracts. *The Journal of Legal Studies* 43, 1 (2014), 1–35.
[8] Vinayshekhar Bannihatti Kumar, Roger Iyengar, Namita Nisal, Yuanyuan Feng, Hana Habib, Peter Story, Sushain Cherivirala, Margaret Hagan, Lorrie Cranor, Shomir Wilson, et al. 2020. Finding a choice in a haystack: Automatic extraction of opt-out statements from privacy policy text. In *Proceedings of The Web Conference 2020*. 1943–1954.

[9] Benjamin B Bederson, James D Hollan, Allison Druin, Jason Stewart, David Rogers, and David Proft. 1996. Local tools: An alternative to tool palettes. In *Proceedings of the 9th annual ACM symposium on User interface software and technology*. 169–170.

[10] Sarah Belia, Fiona Fidler, Jennifer Williams, and Geoff Cumming. 2005. Researchers misunderstand confidence intervals and standard error bars. *Psychological methods* 10, 4 (2005), 389.

[11] Uri Benoliel and Shmuel I Becher. 2019. The duty to read the unreadable. *BCL Rev.* 60 (2019), 2255.

[12] Barbara Bental and Emanuel Tirosh. 2007. The relationship between attention, executive functions and reading domain abilities in attention deficit hyperactivity disorder and reading disorder: A comparative study. *Journal of Child Psychology and Psychiatry* 48, 5 (2007), 455–463.

[13] Jaspreet Bhatia, Travis D Breaux, Joel R Reidenberg, and Thomas B Norton. 2016. A theory of vagueness and privacy risk perception. In *2016 IEEE 24th International Requirements Engineering Conference (RE)*. IEEE, 26–35.

[14] Eric A Bier, Maureen C Stone, Ken Pier, William Buxton, and Tony D DeRose. 1993. Toolglass and magic lenses: the see-through interface. In *Proceedings of the 20th annual conference on Computer graphics and interactive techniques*. 73–80.

[15] Christoph Bösch, Benjamin Erb, Frank Kargl, Henning Kopp, and Stefan Pfattheicher. 2016. Tales from the dark side: Privacy dark strategies and privacy dark patterns. *Proceedings on Privacy Enhancing Technologies* (2016).

[16] Daniel Braun and Florian Matthes. 2021. NLP for consumer protection: Battling illegal clauses in German terms and conditions in online shopping. In *Proceedings of the 1st Workshop on NLP for Positive Impact*. 93–99.

[17] Daniel Braun, Elena Scepankova, Patrick Holl, and Florian Matthes. 2017. Satos: Assessing and summarising terms of services from german webshops. In *Proceedings of the 10th International Conference on Natural Language Generation*. 223–227.

[18] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative research in psychology* 3, 2 (2006), 77–101.

[19] Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z. Gajos. 2021. To Trust or to Think: Cognitive Forcing Functions Can Reduce Overreliance on AI in AI-assisted Decision-making. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 188 (April 2021), 21 pages. https://doi.org/10.1145/3449287

[20] Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z Gajos. 2021. To trust or to think: cognitive forcing functions can reduce overreliance on AI in AI-assisted decision-making. *Proceedings of the ACM on Human-computer Interaction* 5, CSCW1 (2021), 1–21.

[21] Katherine S Button, John PA Ioannidis, Claire Mokrysz, Brian A Nosek, Jonathan Flint, Emma SJ Robinson, and Marcus R Munafò. 2013. Power failure: why small sample size undermines the reliability of neuroscience. *Nature reviews neuroscience* 14, 5 (2013), 365–376.

[22] Shiye Cao, Anqi Liu, and Chien-Ming Huang. 2024. Designing for appropriate reliance: The roles of ai uncertainty presentation, initial user decision, and user demographics in AI-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction* 8, CSCW1 (2024), 1–32.

[23] Xiang "Anthony" Chen, Chien-Sheng Wu, Lidiya Murakhovs' ka, Philippe Laban, Tong Niu, Wenhao Liu, and Caiming Xiong. 2023. Marvista: exploring the design of a human-AI collaborative news reading tool. *ACM Transactions on Computer-Human Interaction* 30, 6 (2023), 1–27.

[24] Elisa Costante, Yuanhao Sun, Milan Petković, and Jerry Den Hartog. 2012. A machine learning solution to assess privacy policy completeness: (short paper). In *Proceedings of the 2012 ACM Workshop on Privacy in the Electronic Society*. 91–96.

[25] Lorrie Faith Cranor, Serge Egelman, Steve Sheng, Aleecia M McDonald, and Abdur Chowdhury. 2008. P3P deployment on websites. *Electronic Commerce Research and Applications* 7, 3 (2008), 274–293.

[26] Lorrie Faith Cranor, Praveen Guduru, and Manjula Arjula. 2006. User interfaces for privacy agents. *ACM Transactions on Computer-Human Interaction (TOCHI)* 13, 2 (2006), 135–178.

[27] Edward Cutrell and Zhiwei Guan. 2007. What are you looking for? An eye-tracking study of information usage in web search. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 407–416.

[28] Robert Lee Dickens. 2007. Finding common ground in the world of electronic contracts: the consistency of legal reasoning in clickwrap cases. *Marq. Intell. Prop. L. Rev.* 11 (2007), 379.

[29] Yi Feng, Mingyang Song, Jiaqi Wang, Zhuang Chen, Guanqun Bi, Minlie Huang, Liping Jing, and Jian Yu. 2024. SS-GEN: A Social Story Generation Framework with Large Language Models. *arXiv preprint arXiv:2406.15695* (2024).

[30] Casey Fiesler, Cliff Lampe, and Amy S Bruckman. 2016. Reality and perception of copyright terms of service for online content creation. In *Proceedings of the 19th ACM conference on computer-supported cooperative work & social computing*. 1450–1461.

[31] Raymond Fok, Joseph Chee Chang, Tal August, Amy X Zhang, and Daniel S Weld. 2024. Qlarify: Recursively Expandable Abstracts for Dynamic Information Retrieval over Scientific Papers. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*. 1–21.

[32] Raymond Fok, Andrew Head, Jonathan Bragg, Kyle Lo, Marti A Hearst, and Daniel S Weld. 2022. Scim: Intelligent Faceted Highlights for Interactive, Multi-Pass Skimming of Scientific Papers.(2022).

[33] Raymond Fok, Nedim Lipka, Tong Sun, and Alexa F Siu. 2024. Marco: Supporting Business Document Workflows via Collection-Centric Information Foraging with Large Language Models. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–20.

[34] Nathaniel Good, Rachna Dhamija, Jens Grossklags, David Thaw, Steven Aronowitz, Deirdre Mulligan, and Joseph Konstan. 2005. Stopping spyware at the gate: a user study of privacy, notice and spyware. In *Proceedings of the 2005 symposium on Usable privacy and security*. 43–52.

[35] Nathaniel S Good, Jens Grossklags, Deirdre K Mulligan, and Joseph A Konstan. 2007. Noticing notice: a large-scale experiment on the timing of software license agreements. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 607–616.

[36] Alfonso Guarino, Nicola Lettieri, Delfina Malandrino, and Rocco Zaccagnino. 2021. A machine learning-based approach to identify unlawful practices in online terms of service: analysis, implementation and evaluation. *Neural Computing and Applications* 33 (2021), 17569–17587.

[37] Yue Guo, Joseph Chee Chang, Maria Antoniak, Erin Bransom, Trevor Cohen, Lucy Lu Wang, and Tal August. 2024. Personalized jargon identification for enhanced interdisciplinary communication. In *Proceedings of the conference. Association for Computational Linguistics. North American Chapter*.

*Meeting*, Vol. 2024. 4535.

[38] Hana Habib, Yixin Zou, Yaxing Yao, Alessandro Acquisti, Lorrie Cranor, Joel Reidenberg, Norman Sadeh, and Florian Schaub. 2021. Toggles, dollar signs, and triangles: How to (in) effectively convey privacy choices with icons and link texts. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–25.

[39] Hamza Harkous, Kassem Fawaz, Rémi Lebret, Florian Schaub, Kang G Shin, and Karl Aberer. 2018. Polisis: Automated analysis and presentation of privacy policies using deep learning. In *27th USENIX Security Symposium (USENIX Security 18)*. 531–548.

[40] Jinwen He, Yujia Gong, Zijin Lin, Yue Zhao, Kai Chen, et al. 2024. Llm factoscope: Uncovering llms' factual discernment through measuring inner states. In *Findings of the Association for Computational Linguistics ACL 2024*. 10218–10230.

[41] Andrew Head, Kyle Lo, Dongyeop Kang, Raymond Fok, Sam Skjonsberg, Daniel S Weld, and Marti A Hearst. 2021. Augmenting scientific papers with just-in-time, position-sensitive definitions of terms and symbols. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–18.

[42] Rink Hoekstra, Richard D Morey, Jeffrey N Rouder, and Eric-Jan Wagenmakers. 2014. Robust misinterpretation of confidence intervals. *Psychonomic bulletin & review* 21, 5 (2014), 1157–1164.

[43] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2024. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems* (2024).

[44] Ziheng Huang, Sebastian Gutierrez, Hemanth Kamana, and Stephen MacNeil. 2023. Memory sandbox: Transparent and interactive memory management for conversational agents. In *Adjunct Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–3.

[45] Duha Ibdah, Nada Lachtar, Satya Meenakshi Raparthi, and Anys Bacha. 2021. "Why Should I Read the Privacy Policy, I Just Need the Service": A Study on Attitudes and Perceptions Toward Privacy Policies. *IEEE access* 9 (2021), 166465–166487.

[46] John P. A. Ioannidis. 2005. Why Most Published Research Findings Are False. *PLOS Medicine* 2, 8 (08 2005). https://doi.org/10.1371/journal.pmed.0020124

[47] John P. A. Ioannidis. 2012. The Importance of Potential Studies That Have Not Existed and Registration of Observational Data Sets. *JAMA* 308, 6 (08 2012), 575–576. https://doi.org/10.1001/jama.2012.8144 arXiv:https://jamanetwork.com/journals/jama/articlepdf/1309181/jvp120046_575_576.pdf

[48] Farnaz Jahanbakhsh, Amy X Zhang, Karrie Karahalios, and David R Karger. 2022. Our Browser Extension Lets Readers Change the Headlines on News Articles, and You Won't Believe What They Did! *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (2022), 1–33.

[49] James Jhirad, Edison Marrese-Taylor, and Yutaka Matsuo. 2023. Evaluating Large Language Models' Understanding of Financial Terminology via Definition Modeling. *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics: Student Research Workshop* (2023). https://doi.org/10.18653/v1/2023.ijcnlp-srw.12

[50] Hang Jiang, Xiajie Zhang, Robert Mahari, Daniel Kessler, Eric Ma, Tal August, Irene Li, Alex Pentland, Yoon Kim, Deb Roy, et al. 2024. Leveraging Large Language Models for Learning Complex Legal Concepts through Storytelling. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 7194–7219.

[51] Robin Bradley Kar and Margaret Jane Radin. 2019. Pseudo-contract and shared meaning analysis. *Harvard Law Review* 132, 4 (2019), 1135–1219.

[52] Michael Karanicolas. 2021. Too Long; Didn't Read: Finding Meaning in Platforms' Terms of Service Agreements. *U. Tol. L. Rev.* 52 (2021), 1.

[53] Matthew Kay, Gregory L. Nelson, and Eric B. Hekler. 2016. Researcher-Centered Design of Statistics: Why Bayesian Statistics Better Fit the Culture and Incentives of HCI. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (Santa Clara, California, USA) *(CHI '16)*. ACM, New York, NY, USA, 4521–4532. https://doi.org/10.1145/2858036.2858465

[54] Matthew Kay and Michael Terry. 2010. Textured agreements: re-envisioning electronic consent. In *Proceedings of the sixth symposium on usable privacy and security*. 1–13.

[55] Patrick Gage Kelley, Joanna Bresee, Lorrie Faith Cranor, and Robert W Reeder. 2009. A" nutrition label" for privacy. In *Proceedings of the 5th Symposium on Usable Privacy and Security*. 1–12.

[56] Moniba Keymanesh, Micha Elsner, and Srinivasan Sarthasarathy. 2020. Toward Domain-Guided Controllable Summarization of Privacy Policies.. In *NLLP@ KDD*. 18–24.

[57] Jana Korunovska, Bernadette Kamleitner, and Sarah Spiekermann. 2020. The Challenges and Impact of Privacy Policy Comprehension. ECIS.

[58] Vinay Koshy, Tanvi Bajpai, Eshwar Chandrasekharan, Hari Sundaram, and Karrie Karahalios. 2023. Measuring User-Moderator Alignment on r/ChangeMyView. *Proc. ACM Hum.-Comput. Interact.* 7, CSCW2, Article 286 (oct 2023), 36 pages. https://doi.org/10.1145/3610077

[59] Lin Kyi, Sushil Ammanaghatta Shivakumar, Cristiana Teixeira Santos, Franziska Roesner, Frederike Zufall, and Asia J Biega. 2023. Investigating deceptive design in GDPR's legitimate interest. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–16.

[60] Michelle S Lam, Mitchell L Gordon, Danaë Metaxa, Jeffrey T Hancock, James A Landay, and Michael S Bernstein. 2022. End-user audits: A system empowering communities to lead large-scale investigations of harmful algorithmic behavior. *proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (2022), 1–34.

[61] Logan Lebanoff and Fei Liu. 2018. Automatic Detection of Vague Words and Sentences in Privacy Policies. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. 3508–3517.

[62] Marco Lippi, Przemysław Pałka, Giuseppe Contissa, Francesca Lagioia, Hans-Wolfgang Micklitz, Giovanni Sartor, and Paolo Torroni. 2019. CLAUDETTE: an automated detector of potentially unfair clauses in online terms of service. *Artificial Intelligence and Law* 27 (2019), 117–139.

[63] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics* 12 (2024), 157–173.

[64] Shuang Liu, Baiyang Zhao, Renjie Guo, Guozhu Meng, Fan Zhang, and Meishan Zhang. 2021. Have you been properly notified? automatic compliance analysis of privacy policy text with gdpr article 13. In *Proceedings of the Web Conference 2021*. 2154–2164.

[65] Xingyu Liu, Annabel Sun, and Jason I Hong. 2021. Identifying Terms and Conditions Important to Consumers using Crowdsourcing. *arXiv preprint arXiv:2111.12182* (2021).

[66] Yixin Liu, Kejian Shi, Katherine He, Longtian Ye, Alexander Richard Fabbri, Pengfei Liu, Dragomir Radev, and Arman Cohan. 2024. On Learning to Summarize with Large Language Models as References. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. 8639–8656.

[67] Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2024. Factual consistency evaluation of summarization in the Era of large language models. *Expert Systems with Applications* 254 (2024), 124456.

[68] Laura Manor and Junyi Jessy Li. 2019. Plain English Summarization of Contracts. In *Proceedings of the Natural Legal Language Processing Workshop 2019*. 1–11.

[69] Arunesh Mathur, Mihir Kshirsagar, and Jonathan Mayer. 2021. What makes a dark pattern... dark? Design attributes, normative considerations, and measurement methods. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–18.

[70] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On Faithfulness and Factuality in Abstractive Summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 1906–1919.

[71] Richard McElreath. 2015. *Statistical rethinking: A Bayesian course with examples in R and Stan.* Chapman and Hall/CRC.

[72] Microsoft. 12/15/2016. *P3P is no longer supported.* https://learn.microsoft.com/en-us/previous-versions/windows/internet-explorer/ie-developer/compatibility/mt146424(v=vs.85)

[73] George R Milne and Mary J Culnan. 2004. Strategies for reducing online privacy risks: Why consumers read (or don't read) online privacy notices. *Journal of interactive marketing* 18, 3 (2004), 15–29.

[74] Razieh Nokhbeh Zaeem, Ahmad Ahbab, Josh Bestor, Hussam H Djadi, Sunny Kharel, Victor Lai, Nick Wang, and K Suzanne Barber. 2022. Privacycheck v3: empowering users with higher-level understanding of privacy policies. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 1593–1596.

[75] Razieh Nokhbeh Zaeem, Safa Anya, Alex Issa, Jake Nimergood, Isabelle Rogers, Vinay Shah, Ayush Srivastava, and K Suzanne Barber. 2020. PrivacyCheck v2: A tool that recaps privacy policies for you. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 3441–3444.

[76] Patricia A Norberg, Daniel R Horne, and David A Horne. 2007. The privacy paradox: Personal information disclosure intentions versus behaviors. *Journal of consumer affairs* 41, 1 (2007), 100–126.

[77] Jonathan A Obar and Anne Oeldorf-Hirsch. 2020. The biggest lie on the internet: Ignoring the privacy policies and terms of service policies of social networking services. *Information, Communication & Society* 23, 1 (2020), 128–147.

[78] Przemyslaw Palka. 2023. Terms of injustice. *W. Va. L. Rev.* 126 (2023), 133.

[79] Shidong Pan, Thong Hoang, Dawen Zhang, Zhenchang Xing, Xiwei Xu, Qinghua Lu, and Mark Staples. 2023. Toward the cure of privacy policy reading phobia: Automated generation of privacy nutrition labels from privacy policies. *arXiv preprint arXiv:2306.10923* (2023).

[80] Rock Yuren Pang, Sebastin Santy, René Just, and Katharina Reinecke. 2024. Blip: facilitating the exploration of undesirable consequences of digital technologies. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–18.

[81] Judea Pearl. 2009. *Causality.* Cambridge university press.

[82] Peter Pirolli and Stuart Card. 1999. Information foraging. *Psychological review* 106, 4 (1999), 643.

[83] Joseph Reagle and Lorrie Faith Cranor. 1999. The platform for privacy preferences. *Commun. ACM* 42, 2 (1999), 48–55.

[84] Daniel Reinhardt, Johannes Borchard, and Jörn Hurtienne. 2021. Visual interactive privacy policy: The better choice?. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–12.

[85] Eric P Robinson and Yicheng Zhu. 2020. Beyond "I agree": Users' understanding of web site terms of service. *Social media+ society* 6, 1 (2020), 2056305119897321.

[86] Kanthashree Mysore Sathyendra, Shomir Wilson, Florian Schaub, Sebastian Zimmeck, and Norman Sadeh. 2017. Identifying the provision of choices in privacy policy text. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2774–2779.

[87] Max Schemmer, Niklas Kuehl, Carina Benz, Andrea Bartos, and Gerhard Satzger. 2023. Appropriate reliance on AI advice: Conceptualization and the effect of explanations. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*. 410–422.

[88] Omar Shaikh, Valentino Emil Chai, Michele Gelfand, Diyi Yang, and Michael S Bernstein. 2024. Rehearsal: Simulating conflict to teach conflict resolution. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–20.

[89] Mukund Srinath, Shomir Wilson, and C Lee Giles. 2021. Privacy at Scale: Introducing the PrivaSeer Corpus of Web Privacy Policies. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 6829–6839.

[90] Pnina Stern and Lilach Shalev. 2013. The role of sustained attention and display medium in reading comprehension among adolescents with ADHD and without it. *Research in developmental disabilities* 34, 1 (2013), 431–439.

[91] Simon Sweeney and Fabio Crestani. 2006. Effective search results summary size and device screen size: Is there a relationship? *Information processing & management* 42, 4 (2006), 1056–1074.

[92] Madiha Tabassum, Abdulmajeed Alqhatani, Marran Aldossari, and Heather Richter Lipford. 2018. Increasing user attention with a comic-based policy. In *Proceedings of the 2018 chi conference on human factors in computing systems*. 1–6.

[93] Lee Taber, Paul May, Keane Yahn-Krafft, and Steve Whittaker. 2020. Beyond avoidance and passivity: Novel uis to make terms of service comprehensible. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–7.

[94] Jenny Tang, Hannah Shoemaker, Ada Lerner, and Eleanor Birrell. 2021. Defining privacy: How users interpret technical terms in privacy policies. *Proceedings on Privacy Enhancing Technologies* (2021).

[95] Welderufael B Tesfay, Peter Hofmann, Toru Nakamura, Shinsaku Kiyomoto, and Jetzabel Serna. 2018. I read but don't agree: Privacy policy benchmarking using machine learning and the EU GDPR. In *Companion Proceedings of the The Web Conference 2018*. 163–166.

[96] Noriko Tomuro, Steven Lytinen, and Kurt Hornsburg. 2016. Automatic summarization of privacy policies using ensemble learning. In *Proceedings of the Sixth ACM Conference on Data and Application Security and Privacy*. 133–135.

[97] Kim-Phuong L Vu, Vanessa Chambers, Fredrick P Garcia, Beth Creekmur, John Sulaitis, Deborah Nelson, Russell Pierce, and Robert W Proctor. 2007. How users read and comprehend privacy policies. In *Human Interface and the Management of Information. Interacting in Information Environments: Symposium on Human Interface 2007, Held as Part of HCI International 2007, Beijing, China, July 22-27, 2007, Proceedings, Part II*. Springer, 802–811.

[98] Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, et al. 2024. FreshLLMs: Refreshing Large Language Models with Search Engine Augmentation. In *Findings of the Association for Computational Linguistics ACL 2024*. 13697–13720.

[99] Zijie J Wang, Chinmay Kulkarni, Lauren Wilcox, Michael Terry, and Michael Madaio. 2024. Farsight: Fostering Responsible AI Awareness During AI Application Prototyping. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–40.

[100] Sewall Wright. 1934. The Method of Path Coefficients. *The Annals of Mathematical Statistics* 5, 3 (1934), 161–215. http://www.jstor.org/stable/2957502

[101] Haijun Xia, Bruno Araujo, Tovi Grossman, and Daniel Wigdor. 2016. Object-oriented drawing. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. 4610–4621.

[102] Johnathan Yerby and Ian Vaughn. 2022. Deliberately confusing language in terms of service and privacy policy agreements. *Issues in Information Systems* 23, 2 (2022).

[103] Yunting Yin and Steven Skiena. 2023. Word Definitions from Large Language Models. *arXiv preprint arXiv:2311.06362* (2023).

[104] Razieh Nokhbeh Zaeem, Rachel L German, and K Suzanne Barber. 2018. Privacycheck: Automatic summarization of privacy policies using data mining. *ACM Transactions on Internet Technology (TOIT)* 18, 4 (2018), 1–18.

[105] Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics* 12 (2024), 39–57.

[106] Tiancheng Zhao and Kyusong Lee. 2020. Talk to Papers: Bringing Neural Question Answering to Academic Search. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. 30–36.

## A   ITERATIVE DESIGN

A total of 8 participants were recruited through Prolific to evaluate an early prototype of TermSight. This prototype closely resembled the current version of TermSight but lacked a document preview when hovering over the Power Meter. Additionally, the prototype used three saturation levels for relevance, resulting in a total of nine colors (three colors for power X three saturation levels for relevance). Participants could also toggle between two layouts for the Summary Snippets: condensed layout and in-context layout. The condensed layout matched the current design shown in Figure 1. In contrast, the in-context layout placed each Summary Snippet directly next to its corresponding Information Snippet. After a brief interface tutorial, participants were given 10 minutes to read a ToS using the interface followed by a semi-structured interview about their experience. Each session lasted for 40 minutes. Below we present the main findings and changes made.

*General findings.*  All participants described that reading with the prototype was easier and more engaging compared to typical ToS, which they tend to skip entirely or skim superficially. In contrast to their usual approach of skimming and trying to search for keywords that catch their eyes, all participants described how the Summary Snippets directly highlighted relevant information otherwise buried in dense sections. Additionally, participants expressed interest in reading the original text after using the Summary Snippets to surface relevant information and valued the ability to navigate between the summary and the original text. 6 participants found the scenarios provided by Phrase Scope to be particularly helpful as the definitions alone can be difficult to contextualize and envision implications.

*Participants preferred condensed layout for quick navigation.*  7 participants preferred and primarily used the condensed layout of the Summary Snippets during the reading session. They noted that the purpose of the Summary Snippets was to support navigation, and the condensed layout made it easier to gain an overview and identify relevant information in a policy. In contrast, the in-context layout required significantly more scrolling, which participants found overwhelming. As a result, in the final version of TermSight, we used the condensed layout for the Summary Snippets (Figure 1).

*Reduce visual complexity of the color scheme.*  All participants found the colors to be intuitive and effective in highlighting power and relevance, helping them decide what to read. However, 3 participants pointed out that while the power dimension was intuitive, the relevance dimension with 3 saturation levels was challenging to differentiate as there would be 9 colors in total. To reduce visual complexity, the final version of TermSight included only two saturation levels (High vs. Low) for relevance (Figure 2).

*Complement Power Meter with Document Preview.*  Additionally, 4 participants suggested that while the Power Meter provided a general overview, they wanted a more concrete preview. This feedback led to the integration of a document preview feature in the final version of TermSight, which appears when users hover over the Power Meter (Figure 3).

## B   ADDITIONAL IMPLEMENTATION DETAILS OF TERMSIGHT

In this section, we provide additional implementation details of TermSight.

## B.1 Source Document and Pre-processing

TermSight takes a set of HTML or markdown source files as input. The source files should exclude irrelevant details (e.g., footers, navigation) while preserving the document's structure such as the section headings (e.g., h1, h2, h3, h4). Given a source file in HTML or markdown, the document is first segmented by headers (e.g., h1, h2, h3, h4). Within each section or subsection, the text is further chunked by newline separators (i.e., "\n") into segments of around 1,500 characters (approximately 250 words) using langchain's RecursiveCharacterTextSplitter[3], with no overlap between chunks. Importantly, paragraph structures are preserved, as the text splitter only splits at newline breaks. This chunking process is intended to accommodate for the limitations of large language models (LLMs) in processing and analyzing extended text [63] common in Terms of Service (ToS) documents. For consistency, the term **"chunk"** in subsequent sections refers to the output of this document pre-processing pipeline. The resulting text chunks are further vectorized using OpenAI's text-embedding-3-small model and stored in a vector database (Pinecone). These vectorized chunks support features of Phrase Scope where information retrieval is needed.

## B.2 Obtaining Summary Snippets and Information Snippets

The prompt used to obtain the Summary Snippets and Information Snippets is detailed in Figure 21. We constrained the summary snippets to be short because prior works have found that adding short summaries (10-20 words) under search results was more effective for navigational and information-seeking tasks compared to longer summaries, which can be harder to skim [27, 91]. This length constraint was further optimized through prompt engineering. We noticed that setting summaries shorter than 12 words led to excessive fragmentation (i.e., too many summary snippets), while longer summaries tried to cover too much information, making them visually dense and less skimmable.

## B.3 Classifying Information Snippets

In TermSight, two classifications were performed for each Information Snippet using GPT-4o with few-shot prompting (Figure 22 and 23). For the classification of Power, each Information Snippet was classified by the degree of control or agency it grants to the Service Provider or the User (Categories: Service, Neutral, User). For the classification of Relevance, Information Snippets were classified based on their relevance to the user persona (Categories: High, Low). User personas used for the study are included in Figure 17 and 18.

## B.4 Phrase Scope

*B.4.1 Identifying unfamiliar or vague phrases.* For each document chunk, we employ two few-shot prompts to GPT-4o to identify potentially unfamiliar (Figure 24) or vague phrases (Figure 25). Few-shot examples for identifying unfamiliar phrases were selected based on phrases participants found challenging in our formative study, which included legal jargon, copyright licenses, and privacy-related terminology. Few-shot examples for identifying vague phrases were informed by vague phrases participants identified in the formative study, as well as prior works analyzing vagueness in Terms of Service (ToS) documents [13, 61, 102]. Both prompts were applied to each document chunk, producing two sets of identified phrases. Since a phrase may be both unfamiliar and vague, the union of these two sets was taken as the final set of identified phrases for a given chunk of text.

*B.4.2 Generating phrase definitions and answers to user questions.* We used a retrieval-augmented question answering approach [31] to generate the phrase definitions. Figure 26 shows the prompt used to generate the definitions after

---

[3]https://python.langchain.com/v0.1/docs/modules/data_connection/document_transformers/recursive_text_splitter/

retrieving potentially relevant document chunks. The prompt explicitly instructs GPT-4o to produce a concise in-context definition of the input phrase and output a list of the retrieved chunks referenced in generating the definition. In the case when the definition is generated without referencing any of the retrieved chunks, it is still retained. TermSight intends to provide plausible definitions even when explicit definitions are absent from the document. This design decision is also supported by prior research demonstrating the capabilities of large language models in generating term definitions that are comparable to traditional dictionaries [49, 103]. When users ask additional questions, the same retrieval-augmented question answering pipeline used for generating definitions is applied, with the only difference being the question asked (Figure 28). If the answer is not specified in the retrieved chunks, the system generates a plausible answer based on the phrase context, similar to the definitions.

*B.4.3   Generating scenarios.* Our needfinding reveals the need to help readers not only interpret the definition of a phrase but to envision potential implications. As a result, we leveraged GPT-4o to generate customized scenarios of potential implications based on user persona. The prompt is specified in Figure 27.

## C   EVALUATION OF TERMSIGHT OUTPUT

We manually evaluated TermSight's feature pipelines and overall performance. Our goal was not to assess advances in system performance, but to verify that the system reliably produced meaningful outputs that could support ToS reading.

### C.1   Evaluating the Classification of Information Snippets

We sampled 116 Information Snippets with stratified sampling from all the sub-policies for each service included in the user study. For each Summary Snippet, we evaluated the classification of power and relevance. Of the 116 Information Snippets, 27 were classified by TermSight as favoring the user, 22 were classified as favoring the service, and 67 were classified as neutral. In addition, 59 were classified as of high relevance to the provided user persona and 57 were classified to be of low relevance.

Our evaluation revealed 3 instances where the power classifications were imperfect, mainly caused by the lack of context in the input information snippet. For example, the snippet (*'All Buy Now purchases in a ServiceY Show are final and binding'*) was classified by TermSight as favoring the service. In the generated explanation, the assumption was that the statement meant that the user can not return or refund the purchase. However, users can return and get a refund if there are problems with the purchase (e.g., an item doesn't match its description) as stated in the service's return policy, making this snippet more neutral. The return policy was not included in the snippet context, as a result, this neutral statement was considered more service-oriented.

For the classification of relevance, there were 2 instances, out of 116, where irrelevant information snippets were classified as relevant to the persona. In both cases, the snippets were more relevant to sellers, not buyers, even though the persona given was that of a buyer. There were 11 instances where the information snippets classified as relevant were indirectly relevant to the input user personas. For example, sellers have to pay a fee to the platform after selling an item. Despite this information being more relevant to sellers, TermSight classified it as relevant to buyers as it exposed the potential hidden fees that sellers might include as part of their listing price.

### C.2   Evaluating Term Definitions and Scenarios

We randomly sampled 113 generated definitions and scenarios using stratified sampling from all the sub-policies for each of the two ToS augmented for the user study. These were evaluated by members of the research team. Out of the

113 sampled phrases, 46 were identified by TermSight as vague phrases, 6 were identified as both jargon and vague phrases, and 61 were identified as jargon.

Our evaluation revealed that all the generated definitions were correct, but out of the 113 definitions, 4 were only general definitions of the phrase not specific to the ToS. On further inspection, we identified that this was because these vague phrases were not explicitly defined anywhere in the ToS. Additionally, their meanings are service-specific and cannot be extrapolated from common sense or inferred by large language models (LLMs). For example, for the phrase *'aggregated anonymized statistics'*, TermSight provided a general definition of what it might mean to aggregate and anonymize user data. However, the specific details—such as what user data are being aggregated—were not specified in the ToS. Rather than a limitation of TermSight, these imperfections highlight the ill-defined nature of the language used in ToS.

For the generated scenarios of phrases relevant to the user persona, we found one instance where the scenario was factually incorrect based on the input context. The input context and phrase stated that users do not gain ownership rights by downloading content from the service. However, the scenario claimed that users might lose ownership rights over the content they create by uploading it to the service. We also noticed that when the passages or phrases target a different audience (e.g., developers) than the user persona used to generate the scenarios (i.e., a lay user), the generated scenarios become less relevant or useful. We did not regenerate these imperfections to keep the user experience realistic to real-world settings where the LLM outputs are not guaranteed to be perfect. Participants in our user study were informed that AI output can be imperfect. This specific incorrect scenario was not accessed by any participants in the study.

## D    BAYESIAN ANALYSIS MODEL DETAILS

In this section, we provide additional details on the models used in the paper. All models used the DAG in Figure 9 to inform the model. Next, in appendix D.1, we describe the model used to estimate the effect of the treatment on the responses related to the experience of using the treatment interface. In appendix D.2, we describe the model used to estimate the effect of the treatment on the comprehension outcomes *i.e.,* the quiz questions. Finally, in appendix D.3, we describe the model used to estimate the effect of the treatment on the experiment subjects' ability to recall facts correctly. Figure 19 and Figure 20 show the forest plots of model coefficients for the experience and comprehension outcomes.

### D.1    Modeling Experience Outcomes

We asked the experiment subjects to rate their experience of using the treatment or the baseline interface on a Likert scale from 1 to 5 (L=5 outcome values). There were N subjects and we asked each subject K (K=7) questions. Since the likert scale is ordinal, we modeled the experience outcomes using an ordinal regression model. The DAG in Figure 9 indicates that the experimental condition alone affects the experience outcomes. There were three randomizations: whether the subject was shown the treatment or control interface, and whether the subject was shown the social media ToS or the shopping ToS, and whether the treatment was shown first or second. We modeled the experience outcomes using a cumulative `logit` model with a linear link function.

The model is given by:

$$Y_{i,j,k,l,m} \sim \text{OrderedLogistic}(\kappa, \phi_{i,j,k,l,m}) \quad \texttt{// ordinal outcome; cumulative logit link} \tag{1}$$

$$\phi_{i,j,k,l,m} = \alpha_i + c_{j,k} + s_l + o_m \quad \texttt{// linear link function} \tag{2}$$

$$\alpha_i \sim \text{Normal}(0,1), i \in \{1,\dots,N\} \quad \texttt{// prior on the intercept, one for each participant} \tag{3}$$

$$c_{j,k} \sim \text{Normal}(0,1), j \in \{1,2\}, k \in \{1,\dots,K\} \quad \texttt{// priors for each question and treatment/control} \tag{4}$$

$$s_l \sim \text{Normal}(0,1), l \in \{1,2\} \quad \texttt{// priors for each ToS case (shopping/social)} \tag{5}$$

$$o_m \sim \text{Normal}(0,1), m \in \{1,2\} \quad \texttt{// priors for each order of presentation (first/second)} \tag{6}$$

$$\kappa_i \sim \text{Normal}(0,1), i \in \{1,\dots,L-1\} \quad \texttt{// prior on the cutpoints} \tag{7}$$

In eq. (7), we further ensured that the samples are ordered *i.e.,* $\kappa_1 < \kappa_2 < \cdots < \kappa_{L-1}$. The priors are conservative. For example, on the logit scale, a coefficient that is normally distributed with mean 0 and standard deviation 1, the range [-3, 3] covers 99% of the distribution and evaluates using inverse logistic to outcome probability range of [0.04, 0.95].

## D.2   Modeling Comprehension Outcomes

Like earlier, there were N subjects, and we asked each subject K (K=6) questions to test their comprehension of the ToS. Thus, the outcome is a binary variable with 1 indicating a correct answer and 0 indicating an incorrect answer. The DAG in Figure 9 indicates that the experimental condition alone affects the comprehension outcomes. There were three randomizations: whether the subject was shown the treatment or control interface, whether the subject was shown the social media ToS or the shopping ToS, and whether the treatment was shown first or second. We modeled the comprehension outcomes using a logistic regression model. Notice that we created a joint variable $c_{j,k,l}$ for the coefficients corresponding to the question, service, and treatment/control. We did this because the questions were different across the two services used in the experiment. We included the treatment/control variable to be able to easily estimate the effect of the treatment on the comprehension outcomes per question.

The model is given by:

$$Y_{i,j,k,l,m} \sim \text{Binomial}(p_{i,j,k,l,m}) \quad \texttt{// binary outcome; logit link} \tag{8}$$

$$\text{Logistic}(p_{i,j,k,l,m}) = \alpha_i + c_{j,k,l} + o_m \quad \texttt{// linear link function} \tag{9}$$

$$\alpha_i \sim \text{Normal}(0,1), i \in \{1,\dots,N\} \quad \texttt{// prior on the intercept, one per participant} \tag{10}$$

$$c_{j,k,l} \sim \text{Normal}(0,1), j \in \{1,2\}, l \in \{1,2\}, k \in \{1,\dots,K\} \quad \texttt{// service, interface, question} \tag{11}$$

$$o_m \sim \text{Normal}(0,1), m \in \{1,2\} \quad \texttt{// priors for presentation order (first/second)} \tag{12}$$

As in the previous model, the priors are conservative. For example, on the logit scale, a coefficient that is normally distributed with mean 0 and standard deviation 1, the range [-3, 3] covers 99% of the distribution and evaluates using inverse logistic to outcome probability range of [0.04, 0.95].

## D.3   Modeling Recall Outcomes

There were N subjects, and we asked each subject a question on recall facts from the ToS. The outcome that we measure is the number of correctly recalled facts, making the outcome a count variable. Since many of the subjects could not

recall any facts correctly, we modeled the recall outcomes using a Zero Inflated Poisson regression model. The DAG in Figure 9 indicates that the experimental condition alone affects the comprehension outcomes. As before, there were three randomizations: whether the subject was shown the treatment or control interface, whether the subject was shown the social media ToS or the shopping ToS, and whether the treatment was shown first or second.

The model is given by:

$$Y_{i,j,k,l} \sim \text{ZeroInflatedPoisson}(\lambda_{i,j,k,l}, \phi) \quad \texttt{// binary outcome; logit link} \tag{13}$$

$$\text{Log}(\lambda_{i,j,k,l}) = \alpha_i + c_j + s_k + o_l \quad \texttt{// linear link function} \tag{14}$$

$$\alpha_i \sim \text{Normal}(0,1), i \in \{1, \ldots, N\} \quad \texttt{// prior on the intercept, one per participant} \tag{15}$$

$$c_j \sim \text{Normal}(0,1), j \in \{1,2\} \quad \texttt{// priors for treatment/control} \tag{16}$$

$$s_k \sim \text{Normal}(0,1), k \in \{1,2\} \quad \texttt{// priors for each ToS case (shopping/social)} \tag{17}$$

$$o_l \sim \text{Normal}(0,1), l \in \{1,2\} \quad \texttt{// priors for presentation order (first/second)} \tag{18}$$

$$\phi \sim \text{Beta}(2,2) \quad \texttt{// prior on the probability of zero inflation} \tag{19}$$

As in the previous model, the priors are conservative. For example, on the Log scale, a coefficient that is normally distributed with mean 0 and standard deviation 1, the range [-3, 3] covers 99% of the distribution and evaluates using inverse log to outcome range of [0.05, 21]. Notice that we model the number of correctly recalled facts, and thus is a conservative prior. Given that most subjects could not recall any facts, the zero inflation parameter is set to a weakly informative prior using a Beta distribution.

## E    STUDY MATERIALS: FORMATIVE STUDY

Questions asked before the reading session:

- How familiar are you with Terms of Service in general?
- Have you previously read or wished that you read the Terms of Service? What were your reasons for wanting or not wanting to read the Terms of Service?
- Have you used your assigned service before?
- Have you read the ToS for your assigned service before?

Semi-structured interview questions asked after the reading session:

- What were the challenges you faced when reading the Terms of Service?
- How did you go about reading the ToS?
- What information were you interested in in ToS? Both from your prior experience in reading ToS or this ToS?
- Imagine if you have a magic wand that can transform the Terms of Service in whatever ways you want. How would you transform the Terms of Service?

## F    STUDY MATERIALS: USER STUDY

### F.1    User Persona

Two personas were given to the participants during the user study for the social media service (Figure 17) and e-commerce service (Figure 18). The same personas were used for features of TermSight to classify relevance and generate

personalized scenarios. The personas were designed based on information participants in the formative study described caring about.

---

**User Persona for Social Media site: Content consumer who posts personal content**

```
Imagine you are a lay user of social media platforms. You are over 18 years old and located in the
United States.
Your usage of Social Media sites:
    • You spend most of your time on the platform scrolling through feeds, liking posts, chatting
      with other users, and sharing personal content such as photos.
Things you care about when using Social Media Sites:
    • You care about Privacy, particularly what data are being collected and how your data can be
      used and shared.
    • You care about what the service can do with user-generated content such as licenses over
      user content or advertising with user content.
    • You care about potential liabilities when using ServiceX.
```

Fig. 17.  User Persona for Social Media site: Content consumer who posts personal content

---

**User Persona for E-commerce site: Buyer who rarely posts reviews**

```
Imagine you are a lay user of E-commerce platforms. You are over 18 years old and located in the
United States.
Your usage of E-commerce sites:
    • You typically engage with the E-commerce platform to buy new or used items from other users.
    • You rarely post any reviews or content on the service.
Things you care about when using E-commerce sites:
    • You care about information related to making purchases, refunds, returns, user protection
      policies, termination, arbitration, and liabilities.
    • You also care about Privacy, particularly what data are being collected and how your data
      can be used and shared.
```

Fig. 18.  User Persona for E-commerce site: Buyer who rarely posts reviews

---

### F.2   Pre-survey Questions

Questions that were asked only once at the beginning of the interview:

- For how many online platforms have you read their Terms and Conditions (T&C) before? [None (0), Few (1-3), Some (4-6), Many (7-9), A lot (>10)]
- How familiar are you with Terms and Conditions (T&C) for online platforms? (5-point likert rating)

5-point Likert rating questions that were asked before each of the two reading sessions:

- How familiar are you with (E-commerce or Social Media) sites?
- How well does the above user persona align with your personal usage of (E-commerce or Social Media) sites?
- How well does the above user persona align with things you personally care about when using (E-commerce or Social Media) sites?

### F.3   5-point Likert Ratings of Reading Experience

Participants rated their reading experience after each of the two reading sessions.

- **E1:** I'm interested in spending more time reading the service's Terms and Conditions (T&C) with the current interface and wish to get a link after the study.
- **E2:** How hard did you have to work to read the T&C?
- **E3:** How easy was it for you to decide which sub-policies to read?
- **E4:** How easy was it for you to decide what text to read within a sub-policy?
- **E5:** How confident are you that you got all the relevant information from the T&C (including sub-policies)?
- **E6:** How much do you feel like you understood the T&C?
- **E7:** How much would you be willing to read the T&C of other services with this interface?

### F.4   Semi-structured Interview Questions

Below are questions that were asked after each of the two reading sessions, except the last question which was only asked after the second reading session where participants would compare and contrast both interfaces.

- Describe your experience reading the Terms of Service using the interface.
- How did you read the ToS using the interface?
- What were the challenges you faced when reading the Terms of Service? To what extent did the interface help you? (be specific about the features in the interface).
- Which interface do you prefer? Why? Compare and contrast. (Asked only after the second reading session)

| Question | Correct Answer | Relevant Passage |
| --- | --- | --- |
| How can ServiceX use photographs you post for advertisements or promotions? | ServiceX can use photos you post in ads without your permission and is not obligated to attribute you as the creator. | [Social Media] "...you grant us a worldwide, royalty-free, perpetual, irrevocable, non-exclusive, transferable, and sublicensable license to **use, copy, modify, adapt, prepare derivative works of, distribute, store, perform, and display Your Content** and any name, username, voice, or likeness provided in connection with Your Content in all media formats and channels now known or later developed anywhere in the world...you irrevocably **waive any claims and assertions of moral rights or attribution** with respect to Your Content." (User Agreement) |
| ServiceX licenses the content you post to third parties. Which of the following is a third-party licensee of user content on ServiceX allowed to do? | Train AI models with user content without user consent. | [Social Media] "For example, this license includes the right to **use Your Content to train AI** and machine learning models" (User Agreement) "**Our licensees cannot**: Combine third-party data with ServiceX public content to target ads; Perform background checks; Track, alert, monitor, or investigate sensitive events (for example, protests or rallies) or sensitive organizations (for example, unions or activist groups)..." (Public Content Policy) |
| You bought a pair of shoes from ServiceY. Before accepting the delivery, you tried the shoes on but found the shoes small for your foot. What can you do to return the item and guarantee a refund? | ServiceY does not accept "item doesn't fit" as a valid reason for return. | [E-commerce] "Trades, offline transactions, **items that do not fit**, or orders where the Buyer changed their mind are not covered by ServiceY Protect. These types of purchases are considered final sale and **not eligible for return**." (Return Policy) |
| According to ServiceY's authentication policy, when you buy an item over 500$, ServiceY will first authenticate the item before shipping it to you. If ServiceY decides that the item is fake, what would happen? | The item will be disposed. You will be refunded. | [E-commerce] "If we reasonably believe an item to be counterfeit, we will **provide a refund** for the full cost of the item and any accompanying shipping costs to the buyer ... if an item in question is deemed counterfeit, we shall ensure **disposal of the item at our discretion**." (Counterfeit Items Policy) |

Table 1. Example multiple choice comprehension questions: The questions were designed to be difficult and can't be directly answered by only reading the Summary Snippets. These questions can also require knowledge about information from multiple positions of the same policy or information across different policies.

| Phrase Type | Original Passage | Generated Definition | Generated Scenario |
|---|---|---|---|
| Jargon | "... When Your Content is created with or submitted to the Services, you grant us a worldwide, royalty-free, perpetual, irrevocable, non-exclusive, transferable, and **sublicensable** license to use, copy, modify, adapt, prepare derivative works of, distribute, store, perform, and display Your Content and any name, username, voice, or likeness provided ..." | Sublicensable refers to the ability to grant further licenses to third parties to use, copy, modify, or distribute your content. | Imagine posting a photo on ServiceX, and later discovering it's featured in a third-party app's advertisement. Due to the sublicensable license, ServiceX had the right to allow another company to use your photo, without needing your further permission, raising concerns about control over your content. |
| Jargon | "... The listed prices for Items do not include taxes, but the taxes will be displayed before a Buyer confirms the purchase. Use of **ServiceY Credit** (as defined below) may modify taxes that apply to a Buyer's purchase ... " | ServiceY Credit refers to non-redeemable promotional credits offered by ServiceY to be used exclusively for purchases on the Service. | Imagine Jane, a savvy shopper, receives $10 ServiceY Credit for a promotion. She buys a vintage lamp listed at $50. At checkout, ServiceY Credit reduces her total to $40, modifying applicable taxes. Jane saves money, but can't withdraw or transfer the Credit—it only applies to ServiceY purchases. |
| Vague | "To use certain features of our Services, you may be required to create a ServiceX account (an Account id) and provide us with a username, password, and **certain other information** about yourself as set forth in the Privacy Policy ..." | Certain other information refers to optional details such as a bio, gender, age, location, profile picture, or social link that you may provide when creating a ServiceX account. | Imagine signing up for ServiceX, and you're asked to provide a username, password, and 'certain other information' like your age and location. Later, you find out that ServiceX uses this data to tailor ads specifically for you, raising concerns about how much they know about you and potential privacy risks. |
| Vague | "... ServiceY reserves the right to discontinue providing Labels to any or all Users at any time and for **any reason** ..." | Any reason refers to ServiceY's discretion to stop providing shipping labels without needing to specify a particular cause or justification. | Imagine you sell vintage clothes online. ServiceY provides you with prepaid shipping labels. Suddenly, without explanation, they stop offering these labels to you. This means you'll need to cover shipping costs yourself, impacting your profits. This demonstrates ServiceY's right to discontinue services for 'any reason,' affecting your business. |

Table 2. Examples of the generated definitions and scenarios for potentially unfamiliar and vague phrases.
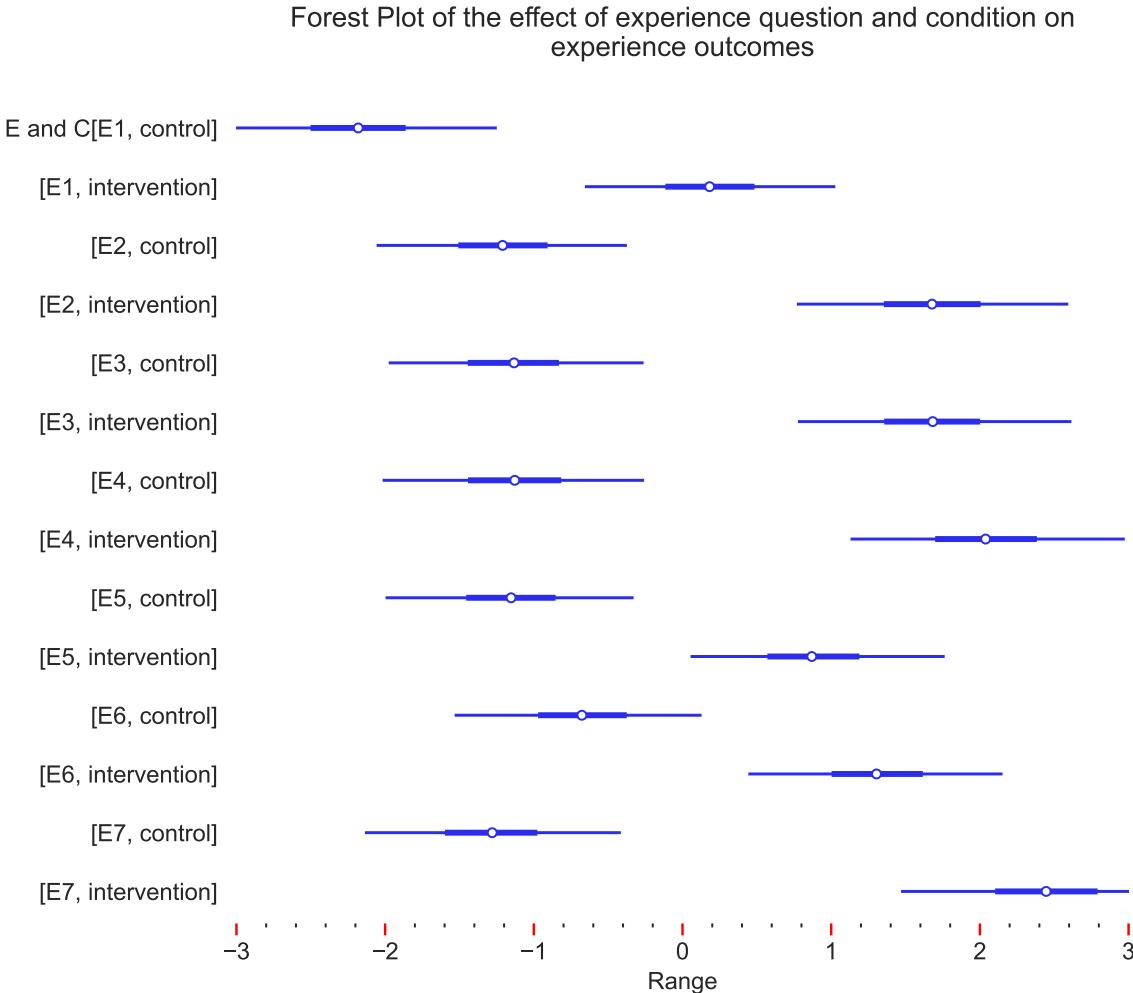
Fig. 19. The figure shows a forest plot of coefficients in the model, corresponding to treatment and control, for each of the experience questions. Each line of the plot shows a 94% High Density Interval (HPDI) for the coefficient. The inner, thicker line represents the 50% HPDI. The results show a significant effect of the TermSight interface on the user experience for every question. Since the 94% HPDI for each pair (control, treatment) for every question do not overlap with each other, we should expect a significant effect of the treatment on the user experience. E1–E7 refer to the 7 experience questions specified in Appendix F.3.
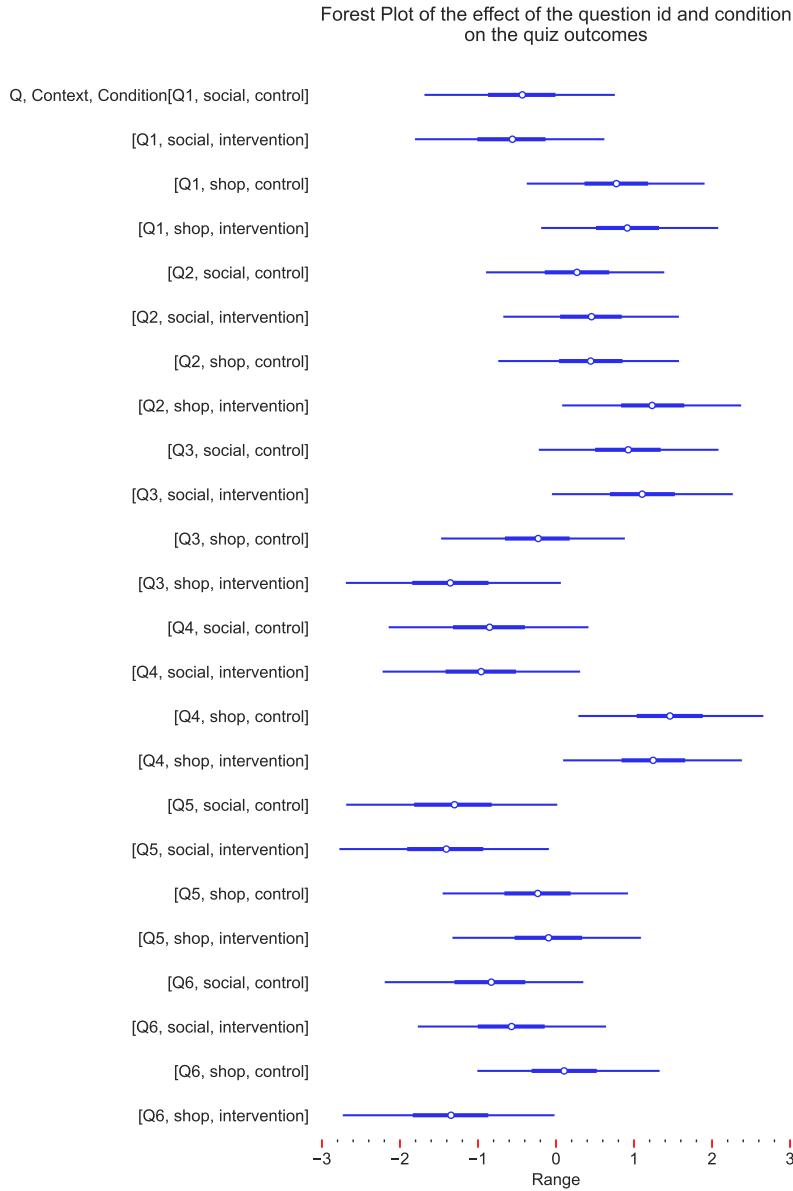
Fig. 20. The figure shows a forest plot of coefficients in the model, corresponding to treatment and control, for each of the experience questions. The $x$-axis is on the logistic scale, with +3 corresponding to a 0.95 probability value on the outcome scale (-3 corresponds to 0.05). Each line of the plot shows a 94% High Density Interval (HPDI) for the coefficient. The inner, thicker line represents the 50% HPDI. The 94% HPDI intervals for each treatment–control pair, when the question and service type are fixed, show significant overlap. This suggests no significant effect of the treatment on the user's comprehension. Note that the comprehension questions Q1–Q6 are different for the two service types and are included in the supplemental materials.

**Prompt for Obtaining Summary Snippets and Information Snippets**

```
Summarize the input section of the terms of service into concise bullet points (less than 12
words) in plain language. When adjacent paragraphs or sentences share a similar or related theme,
only output 1 single bullet point. For each bullet point summary, include the full-text reference
to the original passage in {} and don't use "..." to reduce text in the reference. When outputting
a reference, don't change anything in the original text such as spaces and newlines. There can be
multiple sentence or paragraph references to a single summary. The references to summary should
cover the original text

Example output format: {EXAMPLE OUTPUT}

Input: {INPUT TEXT CHUNK}
```

Fig. 21. Prompt for obtaining Summary Snippets and Information Snippets

**Prompt for Classifying Power**

```
Classify the input term from a Terms of Service agreement based on the power relationship and
benefit between the service and the user. Use the following categories (Bad, Neutral, Good):
- Service: The term grants the service provider disproportionate power or control over the user.
It may impose unfair restrictions, obligations, or liabilities on the user, or reduce the user's
rights and autonomy over their data or content.
- Neutral: The term outlines standard procedures, responsibilities, or conditions the user and
service have. For example, users take responsibility for the content they post. It neither
significantly favors the service provider nor the user, and does not substantially impact the
user's rights.
- User: The term empowers the user by offering clear protections, rights, or benefits, ensuring
transparency, and limiting the service provider's power.

Examples for each category:
Service:
- The service can delete specific content without prior notice and without a reason.
- The service can license user content to third parties.
- The service tracks your personal data for advertising
Neutral:
- Users are responsible for the content they post
- Users agree not to use the service for illegal purposes
- Blocking first party cookies may limit your ability to use the service
User:
- You can opt out of targeted advertising
- The service does not sell your personal data
- The service will not allow third parties to access your personal information without a legal
basis

Output format in JSON: {"Category": "Service/Neutral/User", "Explanation": "explanation of output"
}

Input: {INPUT INFORMATION SNIPPET}
```

Fig. 22. Prompt for classifying power balance of the Information Snippets

**Prompt for Classifying Relevance**

For the input term from a Terms of Service, assign a relevance rating (High/Low) of the input term with respect to the user persona. Output "Low" for low-relevance terms.

Relevance rating:
[High]: The term is directly relevant to the user's usage of the service or what the user cares about. The term applies to the user persona and is necessary for the user to know.
[Low]: The term is not relevant to the user's usage of the service or what the user cares about. The term doesn't apply to the user persona or is not necessary for the user to know.

User Persona: {INPUT USER PERSONA}

Output format in JSON: {"Relevance": "Low/High", "Explanation": "explanation of output" }

Input: {INPUT INFORMATION SNIPPET}

Fig. 23. Prompt for classifying the relevance of Information Snippets to the user persona

**Prompt for Identifying Unfamiliar Phrases**

You are a helpful assistant who extracts words or multi-word phrases in the input section of Terms of Service that a high schooler might not know the meaning of. Jargon refers to domain-specific terminologies that a lay user might not know about.

Example jargon:
- legal jargon: indemnity, arbitration, liability
- copyright licenses: sublicensable licenses, royalty-free licenses
- technical privacy terms: cookies, Ad identifiers, Authentication tokens

Return an empty array if the section does not contain jargon. The extracted word should exactly match the original input text with the same capitalization and sequence of words.

Output format in JSON: {"Jargon": []}

Input: {INPUT TEXT CHUNK}

Fig. 24. Prompt for identifying potentially unfamiliar phrases for lay users in a text chunk.

**Prompt for Identifying Vague Phrases**

```
You are a helpful legal assistant who extracts vague terms (can have multiple words in one term)
in the input section of Terms of Service. A vague term refers to information that is vaguely
abstracted without a clear definition provided in the section.

Example Vague terms: information, other, some, third parties, most, generally, personal data,
others, general, many, various, might, services, certain information

Return an empty array if the section does not contain vague terms. The extracted word should
exactly match the original input text with the same capitalization and sequence of words.

Output format in JSON: {"Vague": []}

Input: {INPUT TEXT CHUNK}
```

Fig. 25. Prompt for identifying vague phrases in a text chunk.

**Prompt for Generating Phrase Definitions**

```
Use information in the retrieved context to provide a definition of the user-selected phrase
or term. Avoid using long sentences. For example, if the user-selected term is "information",
define what the term "information" includes and refers to, such as: location data, interaction
data, profile data, etc. The output definition should be specific and straight to the point,
don't include language that doesn't contribute to the definition such as 'in the given context'.
Output the string list of reference ids (["ref1", ...]) used to generate the definition under
"References". If the definition of the phrase is not specified in the retrieved context, output a
definition of what the phrase might mean and output an empty array for "References".

Examples: {EXAMPLES}

Output format in JSON: {"Definition": "", "References": ["ref1", "ref2", "ref3"]}

Retrieved Context: {RETRIEVED CONTEXT}

Question: What does {INPUT PHRASE} refers to?
Context around the user-selected phrase: {PHRASE CONTEXT}
```

Fig. 26. Prompt for generating in context phrase definitions

**Prompt for Generating Scenarios**

```
Tell a concise what-if scenario or example in less than 50 words to demonstrate the meaning and
potential implications of the user-selected phrase based on the context around the user-selected
phrase. The scenario/example should be relevant to the below user persona using {an E-commerce
platform of used items / a Social Media platform}.

User Persona: {INPUT USER PERSONA}

Output format in JSON: {"Story": ""}

User selected phrase: {INPUT PHRASE}
Context around the user-selected phrase: {PHRASE CONTEXT}
Definition of user selected phrase: {GENERATED DEFINITION}
```

Fig. 27. Prompt for generating scenarios to contextualize the meaning and potential implications of a phrase

**Prompt for Generating Answers to User Questions**

```
You are an assistant for question-answering tasks. Use information in the retrieved context to
answer the user's question in less than 5 sentences. Output the string list of reference ids
(["ref1", ...]) used to generate the definition under "References". If the definition of the
phrase is not specified in the retrieved context, output a definition of what the phrase might
mean and output an empty array for "References".

Examples: {EXAMPLES}

Output format in JSON: {"Answer": "", "References": ["ref1", "ref2", "ref3"]}

Retrieved Context: {RETRIEVED CONTEXT}

Question: {USER QUESTION}
User selected phrase: {INPUT PHRASE}
Context around the user-selected phrase: {PHRASE CONTEXT}
```

Fig. 28. Prompt for generating answers to user questions