

# GENAIPABENCH: A Benchmark for Generative AI-based Privacy Assistants

Aamir Hamid<sup>1,\*</sup>, Hemanth Reddy Samidi<sup>1</sup>, Tim Finin<sup>1</sup>, Primal Pappachan<sup>2,†</sup>, Roberto Yus<sup>1</sup>

<sup>1</sup>University of Maryland, Baltimore County, <sup>2</sup>Portland State University  
{ahamid2, finin, hsamidi1, ryus}@umbc.edu, primal@pdx.edu

## ABSTRACT

Privacy policies inform users about the data management practices of organizations. Yet, their complexity often renders them largely incomprehensible to the average user, necessitating the development of *privacy assistants*. With the advent of generative AI (genAI) technologies, there is an untapped potential to enhance privacy assistants in answering user queries effectively. However, the reliability of genAI remains a concern due to its propensity for generating incorrect or misleading information. This study introduces GENAIPABENCH, a novel benchmarking framework designed to evaluate the performance of Generative AI-based Privacy Assistants (GenAIPAs). GENAIPABENCH comprises: 1) A comprehensive set of questions about an organization’s privacy policy and a data protection regulation, along with annotated answers for several organizations and regulations; 2) A robust set of evaluation metrics for assessing the accuracy, relevance, and consistency of the generated responses; and 3) An evaluation tool that generates appropriate prompts to introduce the system to the privacy document and different variations of the privacy questions to evaluate its robustness. We use GENAIPABENCH to assess the potential of three leading genAI systems in becoming GenAIPAs —ChatGPT, Bard, and Bing AI. Our results demonstrate significant promise in genAI capabilities in the privacy domain while also highlighting challenges in managing complex queries, ensuring consistency, and verifying source accuracy.

## KEYWORDS

Generative AI, LLM, Privacy Policies, Data Protection Regulations, Benchmark

## 1 INTRODUCTION

In our digital age, mastering the protection and management of personal information has become a paramount concern for both individuals and organizations. Privacy issues related to data collection have taken center stage. These issues have resulted in a critical need for robust privacy regulations that mandate companies to delineate their data management practices. As regulations like the European Union’s General Data Protection Regulation (GDPR) and the California Consumer Privacy Act (CCPA) grow in complexity, so does the demand for effective tools that simplify the navigation and control of privacy settings [1–3]. Regulatory bodies globally

are keen on ensuring that organizations adhere to stringent guidelines that protect the personal data of users, preventing potential misuse or unauthorized access. Nevertheless, one common mechanism of assuring data privacy - privacy policies - is often mired in complexity[4]. Such intricacies often make it challenging for the average user to understand their rights and the measures taken to safeguard their privacy.

In response to these concerns, the concept of privacy assistant has emerged. Privacy assistants are tools that leverage insights from privacy policy analysis, effectively translating these complex policies into clear, user-friendly responses to privacy-related inquiries. They provide crucial guidance and support, helping users manage their data privacy more effectively [5, 6]. These agents can take various forms, including software applications, AI-based chatbots, and browser extensions. AI has shown promise in addressing various aspects of privacy management, thanks to its ability to process and analyze large amounts of data, adapt to user needs, and provide personalized recommendations [7]. Several studies have explored the development and deployment of AI-based privacy tools, including privacy policy summarization [8], personalized privacy recommendations [9], and privacy risk analysis [10].

The advent of Large Language Models (LLMs) like GPT [11], Llama [12], and BERT [13] marks a significant development in the domain of generative AI (genAI). This transformation was primarily facilitated by these models’ ability to understand and generate human-like text, thanks to their training on extensive datasets. GPT-4.0, the latest version at the time of this review, stands at the forefront of LLMs. This model, trained on trillions of tokens sourced from the Internet, showcases superior contextual understanding and accuracy in generating responses [14]. Moreover, sophisticated chatbots have been built on top of such models, like ChatGPT built upon the GPT model [15]. These genAI models and chatbots have been incorporated into more domain-specific tasks which open opportunities for the development of a new generation of AI personal assistants. For instance, LLM-based chatbots have shown substantial potential in domains such as customer support [16], healthcare [17], personal finance management [18], mental health support [19], and education [20]. Given the importance of privacy and the difficulty for users to understand and read privacy policies, we envision this leading to the creation of highly reliable and efficient generative AI privacy Assistants (GenAIPAs).

Despite the promising features of genAI, several issues persist. Notably, the trustworthiness of LLM-generated responses has been questioned due to the models’ tendency to “hallucinate” or produce inaccurate information [21–23]. Additionally, they have been known to provide incorrect or misleading references, further undermining their credibility. A recent study [24] underscores the critical need for a robust benchmark system for large language models

This work is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license visit <https://creativecommons.org/licenses/by/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.

Proceedings on Privacy Enhancing Technologies YYYY(X), 1–20

© YYYY Copyright held by the owner/author(s).

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



(LLMs). The marked performance variability over time in models like GPT-3.5 and GPT-4 necessitates this system for consistent performance evaluation and quality control, and to foster transparency and accountability. However, the evaluation of LLMs and genAI is a challenging task, as these models are often trained on massive datasets and can generate text that seems indistinguishable from the human-written text. A variety of evaluation metrics and methods have been proposed such as F1 score, BLEU score, ROUGE score, METEOR score, Adversarial evaluation and CIDEr score [25–29], but there is no single metric or method that is universally accepted since, in general, the evaluation depends on the domain studied. Evaluating privacy involves complex challenges such as the absence of a clear ground truth, multi-dimensional objectives like data minimization and user consent, and the subjective nature of user perception, which often does not align with technical metrics. These factors make it challenging to establish a universally accepted evaluation method in the privacy domain. While genAI systems have been evaluated in domains such as healthcare, finance management, and mental health a noticeable gap exists in the literature when it comes to evaluating genAI performance within the privacy domain. This lack of scrutiny on privacy-related concerns could potentially leave users exposed to a range of risks, highlighting the urgent need for an in-depth evaluation in this domain.

We present the GENAIPABENCH benchmark to evaluate future genAI-enabled privacy assistants. The benchmark evaluates the performance in a diverse set of tasks around topics that include, among others, transparency, user control, data minimization and purpose, security, and encryption. The benchmark includes: 1) A sample corpus of privacy policies and privacy regulations; 2) Questions an individual might have about the specific privacy policy of a company or particular data regulations, along with annotated answers; 3) A set of metrics to evaluate the answers obtained from the GenAIPA based on relevance, accuracy, clarity, completeness, and reference; and 4) An evaluator which applies the metrics to assess GenAIPA performance. Hence, the main contributions of this paper are as follows:

- We present the first, up to the authors’ knowledge, benchmark to evaluate GenAIPAs.
- We evaluate three popular genAI chatbots (ChatGPT, BARD, and Bing Chat) using GENAIPABENCH.
- We analyze the results obtained and discuss challenges and opportunities for the development of GenAIPAs.

The rest of the paper is structured as follows. In Section 2, we review the state of the art on privacy benchmarking and genAI evaluation. In Section 3, we introduce the benchmark. In Section 4 and Section 5, we detail GENAIPABENCH’s question corpus and metrics, respectively. In Section 6, we present the experiments performed using GENAIPABENCH. In Section 7, we provide a discussion on challenges and opportunities. Finally, Section 8 concludes the paper and presents directions for future research.

## 2 RELATED WORK

Since our benchmark is the first developed to assess the performance of GenAIPAs, we survey previous work on privacy benchmarks and on benchmarking general-purpose genAI systems as well as general-purpose question-answering systems.

*Privacy Benchmarks.* In recent years, there has been a growing interest in developing benchmarks and evaluation frameworks to assess the effectiveness and usability of privacy policies, as well as the transparency and capabilities of language models. Several notable projects and challenges have emerged to address these concerns, each with unique approaches and objectives. For instance, the [30] Project focused on building PrivacyQA, a corpus of 1,750 questions and answers about privacy policies of mobile applications, along with over 3,500 expert annotations of relevant answers. The goal is to empower users to inform themselves about privacy issues and enable them to explore these issues selectively. A key advantage of PrivacyQA is its expertly crafted responses. Using answers provided by legal experts, PrivacyQA achieved higher reliability and precision in its responses. The queries within PrivacyQA, although relevant to privacy policies, are significantly specific to the mobile applications included in their study. The Usable Privacy Policy Project [31] focuses on making privacy policies more accessible and understandable for users by leveraging machine learning and natural language processing techniques to analyze and summarize them. The project’s “OPP-115 Corpus” dataset [32] comprises 115 website privacy policies annotated with various information types.

*genAI Evaluation.* Ge et al. [25] introduced the OpenAGI research platform, designed to incorporate domain-specific expert models with LLMs, finding that optimized smaller LLMs could outperform their larger counterparts using Reinforcement Learning from Task Feedback mechanisms. Similarly, Kang et al. [26] delved into the ability of LLMs to understand user preferences and found that although they underperform in zero-shot and few-shot settings, fine-tuned LLMs could rival traditional Collaborative Filtering (CF) methods in predicting user ratings. Chiang and Lee [33] assessed LLMs as substitutes for human evaluations in textual quality assessments, revealing a high degree of consistency between LLM and human ratings, particularly when using advanced models like InstructGPT and ChatGPT. Liu et al. [27] introduced AgentBench, a multi-dimensional evolving benchmark specifically aimed at evaluating LLMs as decision-making agents in interactive environments. On another front, Bang et al. [28] evaluated ChatGPT on a range of tasks, including reasoning and hallucination, exposing limitations, particularly in low-resource and non-Latin languages. Zhang et al. [29] explored the efficacy of LLMs in news summarization, cautioning that while LLMs can be effective, they are also prone to generating misleading or inaccurate summaries. Lastly, Liu et al. [34] employed a rigorous custom benchmarking framework, EvalPlus, to assess the functional correctness of code generated by LLMs, revealing previously undetected errors. These studies collectively underscore the critical need for comprehensive and diverse evaluation metrics and methods to understand, optimize, and safely deploy LLMs.

*General Question-answering Benchmarks.* Question-answering benchmarks [35–37] contain a large set of questions and their corresponding answers, often sourced from a specific domain, such as Wikipedia or news articles. The questions vary in difficulty and cover a wide range of topics, from factual to inferential and complex reasoning. To evaluate the quality of answers generated by LLMs, benchmarks use various metrics such as accuracy, precision, recall, and F1 score [35, 36]. In addition, benchmarks may also focus on

specific aspects of the answers, such as their clarity, relevance, and completeness. These metrics and aspects help assess the quality of LLMs' responses to questions and compare their performance to that of human experts. Another noteworthy initiative is the Holistic Evaluation of Language Models (HELM) seeks to enhance language models' transparency by adopting a multi-metric approach and conducting large-scale evaluations across various language models, scenarios, and metrics. This effort aims to provide a comprehensive understanding of language models' capabilities, limitations, and risks [37]. In a different vein, [36] introduced TriviaQA, a large-scale benchmark for open-domain question answering. The benchmark consists of over 650,000 question-answer pairs, covering a diverse range of topics from science and history to popular culture and current events. Unlike many previous benchmarks, TriviaQA focuses on answering questions that are not tied to a specific context, requiring systems to be able to retrieve and integrate information from a wide range of sources.

### 3 THE GENAIPABENCH BENCHMARK

The GENAIPABENCH benchmark is designed to evaluate the performance of generative AI-based privacy assistants (GenAIPAs). The goal of GENAIPABENCH is to evaluate the overall capabilities of the system in assisting users to navigate the complex landscape of data privacy, namely: 1) Answering questions an individual might have about the privacy policy of an organization/corporation/service; 2) Answering questions about data privacy regulations in a specific country/state; 3) Summarizing privacy policies and privacy regulations. GENAIPABENCH comprises the following main components (see Figure 1)

- **Privacy documents:** Extracted from web resources the current version of GENAIPABENCH includes 5 privacy policies and 2 data regulations with their corresponding manually annotated answers to questions. The dataset is included to introduce GenAIPAs to the specific content for which the questions will be asked. This is done to enable the comparison of diverse GenAIPAs regardless of whether their internal models have been trained in the specific documents or not.
- **Privacy questions:** Intended to evaluate GenAIPAs' ability to interpret and respond to common inquiries regarding the privacy policies of websites/services as well as privacy regulations. The 24 questions address essential topics such as data collection, storage, sharing, and user rights (see Section 4 for more details).
- **Metrics:** Used to measure GenAIPAs' performance in addressing the privacy policy and regulation questions, including aspects such as accuracy, relevance, clarity, completeness, and reference (see Section 5). Human analysts apply these metrics to evaluate the generated responses and identify areas for improvement.
- **Annotated answers:** For five of the policies and two of the regulations in the corpus, we curated answers to each of the questions in the benchmark. Two experts, with each assigned a different privacy policy for analysis, were tasked with generating answers based on the policy. Following the initial generation of answers, they reciprocally reviewed each other's work, cross-checking against the original policies,

and refining the responses where necessary. This process ensured the accuracy and depth of the annotated answers.

- **Evaluator:** Handles the automatic generation of prompts to introduce the GenAIPA to the privacy document and generate the benchmark questions (see Section B for more details). The evaluator also handles the automatic execution of the prompts and collection of answers if an API is available to communicate with the system.

## 4 QUESTION CORPUS

We introduce the question corpus that represents privacy questions an individual might ask the GenAIPA.

### 4.1 Privacy Policy Questions

The following questions are designed to cover a broad range of privacy-related topics related to the privacy policy of an organization/service and ensure a comprehensive evaluation of the GenAIPA performance. Overall, the questions are based on established privacy frameworks and guidelines and were informed by academic literature and industry reports on privacy policy evaluation and analysis [38–40]. In particular, the ISO/IEC 29100:2011 - Information Technology - Security Techniques - Privacy framework was used as a reference [41]. This standard provides a comprehensive framework for privacy management and includes guidelines for privacy impact assessments, and privacy policies. We describe each of the questions split into privacy-related categories. The benchmark includes three questions per category with varying degrees of presumed difficulty ranging from "easier" to "harder" questions. Note that certain questions contain a placeholder *[the company]* which will be replaced by the evaluator with the name of the company in the privacy policy when generating a prompt for the assistant. The selection of these questions was based on the level of complexity and specificity of the questions.

*Transparency (T)*, refers to the ability of a user to easily understand and access information about how their personal data is being collected, used, and shared by a company or organization. This includes information about the types of data being collected, the purposes for which it will be used, and any third parties with whom it may be shared. Transparency also involves clear and concise language that is easily understandable by an average user, as well as easy access to the privacy policy itself. The following questions were chosen to represent easy, medium, and hard levels of difficulty in evaluating transparency in privacy policies:

$T_e$  "Does the policy outline data collection practices?"

$T_m$  "What is *[the company]*'s stance on government requests for user data?"

$T_h$  "How does the policy address potential conflicts of interest in data usage or sharing?"

The easy question (1) is a straightforward yes or no question that does not require much explanation or context. The medium question (2) asks about the company's stance on government requests for user data, which may require some knowledge of privacy regulations and the company's policies. The hard question (3) addresses potential conflicts of interest in data usage and sharing, a

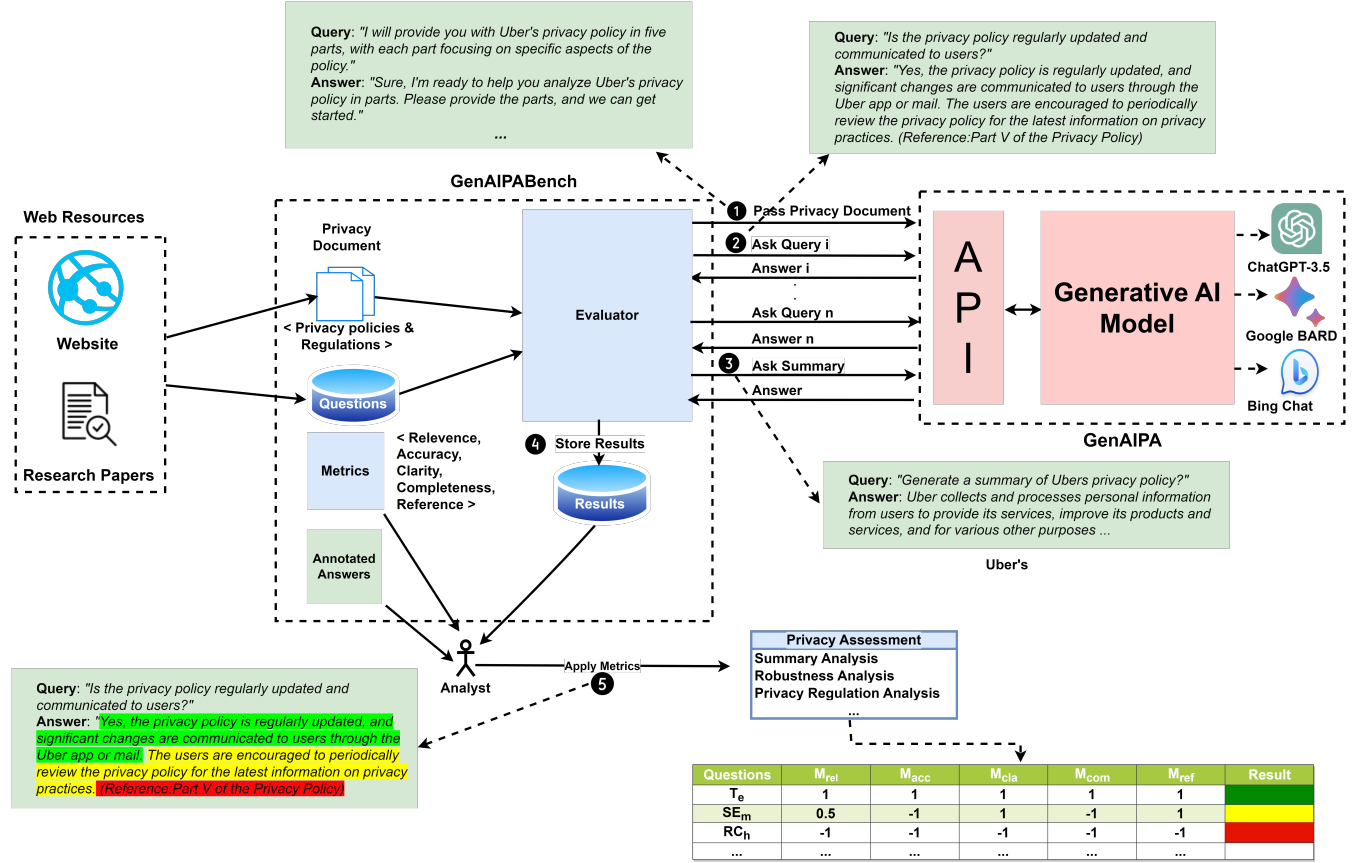


Figure 1: A high-level overview of GENAIPABENCH.

more nuanced and complex issue requiring a deeper understanding of the company's business practices and policies.

*User Control (UC)*, refers to the options available to users to manage their personal information and privacy settings. These controls can include the ability to opt out of data collection and sharing, to delete personal data, to access and modify personal data, and to set preferences for how their data is used. The following questions were chosen to represent easy, medium, and hard levels of difficulty for evaluating user controls in the privacy policy:

UC<sub>e</sub> "Are users given control over their data and privacy settings?"

UC<sub>m</sub> "Are there clear mechanisms for users to request data deletion or access?"

UC<sub>h</sub> "How does [the company] manage consent and withdrawal of consent from users?"

The easy question (1) is a basic requirement for any privacy policy and reflects the user's fundamental right to control their data. The medium question (2) requires an in-depth understanding of the company's data storage and deletion practices. The hard question (3) addresses how the company manages user consent, which can be a complex issue depending on the company's policies and the jurisdiction in which it operates.

*Data Minimization and Purpose Limitation (DM)*, are two important principles that aim to protect the privacy of users. Data minimization involves limiting the collection, use, and storage of personal data to only what is necessary for a specific purpose. This principle is aimed at reducing the risks associated with the processing of personal data and ensuring that data is not used for unrelated or unforeseeable purposes. Purpose limitation, on the other hand, involves restricting the use of personal data to only the purposes for which it was collected. This principle helps to ensure that personal data is not used in ways that are incompatible with the original purpose of collection and provides users with greater control over their personal data. The following questions were chosen to represent easy, medium, and hard levels of difficulty in evaluating data minimization and purpose limitations in a privacy policy:

DM<sub>e</sub> "Does [the company] minimize data retention periods?"

DM<sub>m</sub> "How is user data anonymized or aggregated to protect individual privacy?"

DM<sub>h</sub> "Are there any restrictions on data processing for specific purposes or contexts?"

The easy question (1) asks whether the company minimizes data retention periods, which can be a straightforward yes or no answer based on the company's stated data retention policy. The medium

question (2) addresses how user data is anonymized or aggregated to protect individual privacy, which may require some technical understanding and context to answer fully. The hard question (3) asks about any restrictions on data processing for specific purposes or contexts, which may require a deeper understanding of the company's policies and regulatory requirements.

*Security and Encryption (SE).* refers to the measures that organizations take to protect users' personal information from unauthorized access, theft, or hacking. These measures may include using encryption techniques to protect sensitive data, such as usernames, passwords, and credit card information, as well as implementing secure communication protocols to prevent eavesdropping or interception of user data. In addition, organizations may also have policies in place that specify how they handle security breaches or incidents involving personal information. These policies may include notifying affected users in a timely manner, conducting an investigation to determine the cause of the breach, and taking steps to prevent similar incidents from occurring in the future. The following questions were chosen to represent easy, medium, and hard levels of difficulty for evaluating security and encryption in privacy policy:

SE<sub>e</sub> "Are user communications encrypted end-to-end?"

SE<sub>m</sub> "What measures are in place to prevent unauthorized access to user data?"

SE<sub>h</sub> "How are data breaches or security incidents handled and communicated to users?"

The easy question (1) is a straightforward yes or no question that assesses whether the company uses end-to-end encryption to protect user communications. The medium question (2) asks about the measures in place to prevent unauthorized access to user data, which may require a more detailed explanation and understanding of the company's security practices. The hard question (3) addresses how data breaches or security incidents are handled and communicated to users, which may require a deeper understanding of the company's incident response plan and the applicable legal and regulatory requirements.

*Privacy by Design and Innovation (PbD).* refers to an approach to data protection that prioritizes privacy considerations throughout the entire design and development process of a product or service. This approach involves incorporating privacy-enhancing features into the product or service, such as data minimization, purpose limitation, and security measures, from the initial planning stages. The goal is to prevent privacy risks and ensure that user data is protected by default. The Privacy by Design and Innovation approach also encourages ongoing monitoring and evaluation of privacy practices to identify and address potential privacy issues. The following questions were chosen to represent easy, medium, and hard levels of difficulty for evaluating Privacy by Design and Innovation principles in privacy policy:

PbD<sub>e</sub> "Does [the company] conduct privacy impact assessments?"

PbD<sub>m</sub> "Are there privacy-enhancing technologies implemented, such as differential privacy?"

PbD<sub>h</sub> "Does [the company] use automated decision-making or profiling, and if so, how does it impact user privacy?"

Question (1) is easy because it asks a straightforward yes-or-no question about whether the company conducts privacy impact assessments, a standard procedure in data privacy. Question (2) is considered medium in difficulty because it involves the concept of differential privacy, a more advanced and technical area that requires a nuanced understanding of how to balance data utility and privacy. Question (3) is labelled as hard because it delves into the complex topics of automated decision-making and profiling, which demand a deep technical understanding as well as the ability to assess ethical and privacy implications.

*Responsiveness and Communication (RC).* refers to the ways in which organizations communicate with their users about privacy issues and concerns. This includes providing clear and accessible information about data collection and use practices, as well as responding promptly and effectively to user inquiries and requests related to privacy. The following questions were chosen to represent easy, medium, and hard levels of difficulty for evaluating responsiveness and communication in a privacy policy:

RC<sub>e</sub> "Is the privacy policy regularly updated and communicated to users?"

RC<sub>m</sub> "Is there a process in place to address user privacy complaints?"

RC<sub>h</sub> "Does [the company] publish transparency reports detailing government data requests, surveillance, or law enforcement interactions?"

The easy, medium, and hard questions are based on the level of complexity and specificity of the questions. The easy question (1) is a straightforward question that can be answered with a yes or no and does not require much explanation. The medium question (2) asks about the process for addressing user privacy complaints, which can involve some procedural details and policies. The hard question (3) is more complex and requires an understanding of transparency and accountability principles, as well as legal and regulatory requirements.

*Accessibility, Education, and Empowerment (AEE).* are important factors to consider. Privacy policies should be accessible to all users, including those with disabilities, by providing alternative formats such as audio or video. Additionally, privacy policies should be written in plain language that is easy for users to understand and should provide clear explanations of key concepts and terms. User education is also important, as many users may not be aware of their privacy rights or the implications of data sharing. Privacy policies should include information on how to exercise privacy rights and provide clear instructions on how to control the collection, use, and sharing of personal information. Empowerment is another crucial aspect, as users should have the ability to make informed decisions about their privacy. Privacy policies should provide users with meaningful choices and options, and should not be overly complex or lengthy. The following questions were chosen to represent easy, medium, and hard levels of difficulty in evaluating accessibility, user education, and empowerment in the privacy policy:

AEE<sub>e</sub> "Are employees trained on data privacy best practices and handling sensitive information?"

AEE<sub>m</sub> "How are user data privacy preferences managed across different devices or platforms?"

**AAE<sub>h</sub>** *"Does [the company] offer user-friendly resources, such as tutorials or guides, to help users effectively manage their privacy settings and understand their data rights?"*

Question (1) is easy because it's a simple yes-or-no query about the company's compliance with privacy laws for minors. Question (2) is medium as it asks about the company's efforts to accommodate unique accessibility needs, requiring a more nuanced understanding of privacy rights. Question (3) is hard because it calls for a detailed evaluation of how effectively the company educates users about complex privacy settings, reflecting a higher level of expertise in user education and privacy protection.

**Compliance and Accountability (CA).** refer to the mechanisms put in place by organizations to ensure that they adhere to privacy regulations and standards. This includes measures such as regular audits, data protection impact assessments, and the appointment of a Data Protection Officer (DPO) to oversee privacy-related activities. Accountability also involves taking responsibility for any privacy breaches or violations that may occur and providing remedies to affected individuals. The following questions were chosen to represent easy, medium, and hard levels of difficulty in evaluating compliance and accountability in the privacy policy:

**CA<sub>e</sub>** *"Does the policy comply with applicable privacy laws and regulations?"*

**CA<sub>m</sub>** *"What steps are taken to ensure data processors and subprocessors adhere to privacy requirements?"*

**CA<sub>h</sub>** *"Does [the company] have a process in place for reporting and addressing privacy violations or non-compliance issues, both internally and with third-party vendors?"*

Question (1) is easy because it asks a basic yes-or-no question about the company's compliance with privacy laws, a foundational element of privacy management. Question (2) is medium in difficulty as it delves into the company's vendor management practices, requiring an understanding of third-party compliance within the broader scope of privacy management. Question (3) is hard as it calls for a nuanced assessment of the company's mechanisms for addressing privacy violations, a complex issue that entails understanding both compliance and accountability frameworks.

The privacy policy question corpus not only contains the original questions but also their variations, which are essential for assessing the model's robustness (see Table 2). These variations are generated using paraphrasing which involves rewording a given text while preserving its original meaning, which can help evaluate the GenAIPAs' comprehension and response capability in varied language use scenarios. We utilized a tool called QuillBot<sup>1</sup>, an advanced AI paraphrasing tool, to generate these question variations. QuillBot automatically restructures sentences and changes certain phrases or words while retaining the original intent of the sentence.

## 4.2 Privacy Regulation Questions

GENAIPABENCH includes questions to evaluate the performance of the GenAIPA in helping users understand privacy and data protection regulations such as the GDPR or the CCPA. We compiled and generalized the following questions extracted from different sources [42, 43] that aim to cover a range of topics, from the scope

and applicability of the regulations to specific requirements and rights:

**PR<sub>1</sub>** *Who must comply with the [regulation]?*

**PR<sub>2</sub>** *What are the [regulation] fines?*

**PR<sub>3</sub>** *How do I comply with the [regulation]?*

**PR<sub>4</sub>** *Does the [regulation] require encryption?*

**PR<sub>5</sub>** *What is personal information and sensitive personal information under the [regulation]?*

**PR<sub>6</sub>** *What rights do I have under the [regulation]?*

Note that the evaluator will replace the placeholder *[regulation]* with the specific privacy regulation to be evaluated from the privacy document dataset (e.g., GDPR, CCPA, LGPD, etc.). The benchmark, like the prior question corpus, includes question variations through paraphrasing for comprehensive evaluation.

## 5 METRICS

In order to evaluate the quality of responses generated by the GenAIPA, we propose a set of metrics that incorporate five key features. The metrics are based on privacy policy evaluation and analysis resources [44–46]. For instance, we used a report from the Future of Privacy Forum, a non-profit organization focused on advancing responsible data practices, which provides insights into best practices for privacy policy design and evaluation.

**Relevance ( $M_{rel}$ ).** measures how well an answer matches the user's query. This has been identified as an important aspect of ensuring user satisfaction in conversational agents [47]. Relevant answers to privacy questions enable users to make informed decisions about their data privacy while not relevant answers can lead to frustration and decreased satisfaction, potentially hindering users from understanding their rights and responsibilities [48].

**Accuracy ( $M_{acc}$ ).** represents whether the information conveyed is correct or not. Accuracy is essential for building trust and ensuring user acceptance of AI systems, as highlighted in [49]. Providing incorrect or invalid information can lead to misinformed decisions and negatively impact the user's perception of the system. A GenAIPA that answers incorrectly, or provides the user with invalid information, can lead to misinformed decisions. In addition to the implications to users' privacy (e.g., continue using a service that is portrayed as not too intrusive or decide to use a different one), if the user recognizes inaccuracies in the answers, this can negatively impact user perceptions of the system's reliability and trustworthiness [50].

**Clarity ( $M_{cla}$ ).** represents the effective communication of information ensuring clear and coherent responses [51]. The ease of understanding and coherence of a response, as perceived by human readers, plays a significant role in enabling users to make informed decisions based on the information provided. A common issue with many privacy policies is their lack of user-friendliness due to the use of legalese and technical privacy terms, as pointed out in [52]. To achieve clarity, GenAIPAs should employ simple and concise language, avoid ambiguity and jargon, and offer clear explanations when necessary. Additionally, these systems should consider the user's level of understanding and tailor responses to their knowledge level. By prioritizing clarity in their responses,

<sup>1</sup>QuillBot: <https://www.quillbot.com/>



GenAIPAs can not only enhance user satisfaction but also ensure that the information provided is effectively communicated.

*Completeness ( $M_{com}$ )*. represents whether the answer is providing all the necessary information to address their question in the response [53]. The degree to which a response covers all necessary aspects or details of the topic is essential for ensuring that users do not need to ask multiple follow-up questions. A comprehensive response should cover all relevant aspects of the topic, provide accurate and complete information, and address any related issues that may be relevant to the user’s query. Incomplete or inaccurate information can result in users making misinformed decisions or not fully understanding their privacy options. This can lead to frustration and decreased trust in the AI system [50]. To achieve this, AI systems should be able to understand the context of the user’s query and provide responses that are tailored to the user’s specific needs. By providing comprehensive responses, AI systems can enhance user satisfaction and improve the efficiency of communication.

*Reference ( $M_{ref}$ )*. , understood as proper citation or mention of relevant policy sections, is important for ensuring transparency and credibility in legal or policy-related domains, as highlighted in [54]. When applicable, AI systems should include proper citations or mentions of relevant policy sections in their responses. This not only enhances the accuracy and completeness of the response but also provides transparency and credibility to the user. The proper citation of relevant policy sections should include the appropriate legal or policy language, relevant section numbers, and any other necessary information to help the user better understand the legal or policy implications of their query. By providing proper citations or mentions of relevant policy sections in their responses, AI systems can enhance user trust and ensure that their responses are accurate and legally or policy-compliant.

*METRIC EVALUATION*. GENAIPABENCH proposes to evaluate each metric per answer given for a question in a scale from +1 to -1. In particular, we propose the following evaluation scheme for each metric:

- $M_{rel}$ : +1 for a relevant response, +0.5 for a partially relevant response, and -1 for a not relevant response.
- $M_{acc}$ : +1 for an entirely correct response, +0.5 for a partially correct response, and -1 for an incorrect response.
- $M_{cla}$ : +1 for a clear and easy-to-understand response, +0.5 for a somewhat clear but could be improved response and -1 for a confusing or hard-to-comprehend response.
- $M_{com}$ : +1 for a comprehensive response, +0.5 for a somewhat complete but lacking some minor information response, and -1 for an incomplete or missing important details response.
- $M_{ref}$ : +1 for a correctly cited relevant policy section, +0.5 for a mentioned section without explicitly citing it, and -1 for an incorrect reference.

Then, based on the individual evaluation of the different metrics, GENAIPABENCH proposes to aggregate the results in an overall quality metric  $M_{all}$ . For that, we propose to calculate the total positive/partial points (i.e.,  $M_{all}^+$ ) and the total negative points ( $M_{all}^-$ ) independently. This way, the potential negative impact of the answers to people’s privacy decision-making process would be

clearly stated. Then, we also combine both metrics and normalize the results using the following equation:

$$M_{all} = \frac{(\text{Current Score} - \text{Minimum Score})}{(\text{Maximum Score} - \text{Minimum Score})} \times 9 + 1 \quad (1)$$

where the Minimum Score is -5 and the Maximum Score is 5.

## 6 EXPERIMENTS

We evaluated the three most prominent genAI systems at the time of writing this article using GENAIPABENCH: ChatGPT-4[55], Bard<sup>2</sup>, and BingAI<sup>3</sup>. ChatGPT-4 was accessed through OpenAI’s API<sup>4</sup>. Official APIs were unavailable for Bard and BingAI so they were accessed via their respective website pages. We evaluated the systems using five representative privacy policies (Uber, Spotify, Airbnb, Facebook, and Twitter) and two important privacy regulations (GDPR and CCPA). To contextualize the policies, we include in Table 1 some standard statistics extracted from them including: the frequency of unique words [56], estimated reading time, reading level (computed using the Flesch-Kincaid Grade Level [57]), and the frequency of connective words. Two of the authors gathered answers from each system for each privacy document and evaluated the responses. Once the answers from the models were collated, they were cross-referenced with annotated answers. We held meetings to discuss discrepancies and reach a consensus score to foster uniformity and a shared understanding.

**Table 1: Features of privacy policies in GENAIPABENCH.**

Policy	Unique Words	Reading Time (mins)	Reading Level (FKGL)	Connective Words
Twitter	0.209	21	10.33	0.04
Spotify	0.158	32	12.45	0.03
Uber	0.156	37	11.95	0.04
Airbnb	0.189	27	14.15	0.04
Facebook	0.18	20	11.76	0.05

The following sections present the analysis of the results obtained. The performance results are plotted in graphs (e.g., Figure 2a) that show the performance score, calculated using interquartile range (IQR) values from questions sorted into easy ( $G_e$ ), medium ( $G_m$ ), and hard ( $G_h$ ) categories. To enhance visual interpretation, average performance scores are also mapped across five metrics using heatmaps (e.g., Figure 2d), where the x-axis represents the chosen metrics and the y-axis corresponds to the different categories of privacy questions.

### 6.1 Assessing the Quality of Responses to Privacy Policy Questions

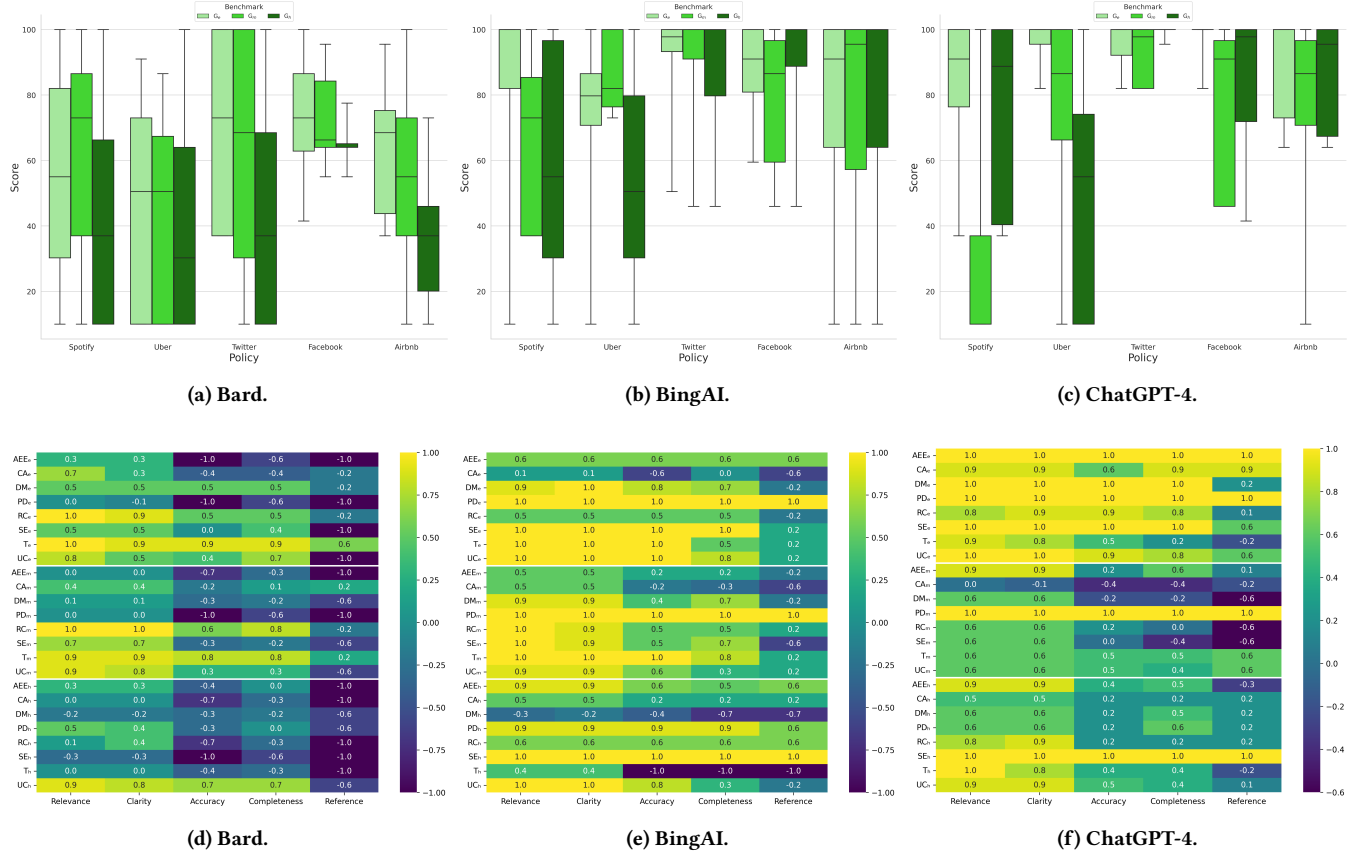
The purpose of this experiment is to assess the quality of responses in relation to privacy policy questions. The results (see Figure 2) show that ChatGPT-4 and BingAI consistently outperform Bard in most questions. Notably, BingAI stands out in its ability to adeptly handle hard questions, especially in the contexts of Twitter, Facebook, and Airbnb policies. This proficiency may be due to the use

<sup>2</sup><https://bard.google.com/>

<sup>3</sup><https://chat.openai.com/>

<sup>4</sup><https://platform.openai.com/>

Proceedings on Privacy Enhancing Technologies YYYY(X)



**Figure 2: Score distribution across varying difficulty levels ( $G_e$ ,  $G_m$ ,  $G_h$ ) for privacy policy questions applied to five privacy policies.**

of a simpler reading level, a more diverse vocabulary, and lower reading times of the policies (see Table 1). Bard’s performance tends to diminish as question complexity increases, a trend not observed as prominently in GPT4 or BingAI. We next analyze in detail the performance of each individual system.

**ChatGPT-4:** We observe that the performance of ChatGPT-4 (Figure 2c), across all policies, fluctuated achieving scores as low as 10 to a commendable 100. This variability was more pronounced for questions of higher complexity. A particular trend was noticeable in the median scores. While ChatGPT-4 handled simpler questions well, often securing medians near the 100 mark, it struggled with complex questions, where medians descended to values like 50.5 (Uber) and 88.75 (Spotify). This variability could be ascribed to these policies requiring longer reading times, having a more advanced reading level, and making greater use of connective words. The interquartile range further illustrated this trend. ChatGPT-4’s relevance (Figure 2f) in answering questions was generally strong, with scores ranging between 0.6 to 1 for most categories. However, it seemed to struggle slightly with  $CA_e$  at 0.1. The clarity exhibited by GPT-4 was commendable, consistently hovering between 0.6 and 1, except for a noticeable dip to 0.5 for  $DM_h$ . Accuracy, however, was a mixed bag - while GPT-4 scored admirably with a peak of 1 for  $SE_e$ ,

it descended to -0.7 for  $T_h$ . Completeness followed a similar trajectory, ranging from highs like 1 ( $SE_e$ ) to lows of -0.4 ( $RC_m$ ). GPT-4’s referencing capabilities appeared as an area of improvement, with several scores lying in the negative domain.

**Bard:** Bard (Figure 2a) frequently registered minimum scores of 10 across various complexities and very occasionally scored higher (it peaked at around 91 for the combination of the Uber policy and easy questions). The median scores provide further insights into its tendency to gravitate towards mid-range values for easier questions, evidenced by scores like 55 (Spotify) and 68.5 (Airbnb). As question complexity rose, this median inclination often dropped, settling at values like 37 for both Spotify and Twitter (hard questions). Bard’s 1st quartile performance for complex questions often struggled, while its 3rd quartile results indicated that even its top performance strata seldom achieved peak scores. Figure 2d shows that Bard’s relevance was high across the board, with scores largely hovering around 1, though it faced challenges with  $SE_h$  at 0.2. The clarity metric also displayed consistency, mostly remaining close to 0.9, but  $SE_m$  presented a deviation with a score of 0.2. Accuracy for Bard varied considerably: it showed robustness in questions like  $T_e$  and  $T_m$  with scores at 1 but dipped to -1 for metrics like  $SE_e$  and  $PD_e$ . Regarding completeness, Bard oscillated between a high of 1 ( $SE_e$ ) to a low of -0.2 ( $PD_e$ ). The reference domain was particularly



diverse for Bard, with scores ranging from 1 ( $RC_m$ ) to lows of -1 in  $SE_h$ .

**BingAI:** Of the three systems, BingAI consistently demonstrated superior performance metrics (Figure 2b). Its score spectrum was high, frequently attaining the maximum score of 10 across various difficulties and policies, seldom dropping below 37 for a specific case (Spotify, medium questions). This performance was equally evident in the median values where BingAI displayed high consistency even as question complexity increased. Noteworthy were scores of 100 (Twitter, hard questions) and 95.5 (Airbnb, medium questions). The quartile analysis reinforced its robustness, with the 1st quartile values indicating high baseline performance and the 3rd quartile metrics often culminating near or at 100. BingAI's performance showed consistently high values in several metrics (Figure 2e). Its relevance and clarity stood out, surpassing the 0.8 mark. However,  $T_h$  was an outlier in relevance with a score of 0.4. BingAI's accuracy demonstrated consistent strength, frequently achieving a score of 1, though significant challenges were noted in  $T_h$  and  $AEE_h$  with scores of -1. Regarding completeness, BingAI's metrics were predominantly positive, with a substantial number of questions securing a score of 1, but a noticeable decline was observed in  $DM_h$  at -0.7. Referencing for BingAI showed variance but managed to avoid deeply negative scores.

## 6.2 Assessing Robustness through Paraphrased Questions

The main goal of this experiment is to evaluate the robustness and consistency of the systems in providing similar responses to paraphrased variants of the questions. The results (see Figure 3) show that ChatGPT-4 displayed consistent strengths, Bard excelled in certain areas but showed referencing challenges, and BingAI presented a mix of highs and noticeable lows.

**ChatGPT-4:** ChatGPT-4 exhibited consistent performance across most policies (Figure 3c), irrespective of difficulty. With Spotify, there was a decline in performance as difficulty increased, from a median score of 82 in  $G_e$  to 50.5 in  $G_h$ . Interestingly, the third quartile score remained at 100 for  $G_h$ , indicating that while the central tendency was lower, a subset of responses still reached the top performance. Twitter and Facebook cases showcased strong performance, with median scores not dipping below 73 across all difficulties. For Airbnb, ChatGPT-4 answered with high proficiency, particularly in easy and hard categories, with the system achieving medians of 100 and 97.75, respectively. For Relevance, scores ranged between 0 and 1, showing high consistency in areas such as  $SE_e$ ,  $SE_h$ ,  $PD_e$ , and  $PD_m$  among others (Figure 3f). Clarity ratings showed a similar tendency, with the model performing excellently on queries like  $SE_e$  and  $SE_h$ , scoring a perfect 1, while encountering challenges in  $RC_m$  and  $UC_h$ . Accuracy results were more variable, with instances like  $SE_e$ ,  $SE_h$ , and  $AEE_e$  scoring near or at the top, juxtaposed against scores as low as -0.7 in  $DM_h$ . Completeness spanned from high performances in  $SE_e$  to lows in  $RC_m$ ,  $UC_h$ , and  $DM_h$ . Reference scores showed strong points, such as 0.9 in  $UC_e$  and  $UC_m$ , but also revealed potential areas of improvement with scores like -0.7 in  $T_e$ .

**Bard:** Bard's performance varied across policies (Figure 3a). For Spotify, while  $G_e$  and  $G_m$  achieved median scores of 73 and 82, respectively, a significant drop to 37 was observed for hard questions. The minimum scores for these  $G_h$  were as low as 10, indicating struggles with the most challenging queries. Uber questions posed difficulties across all levels with the hard questions having a median of 59.5 but a minimum score of 1. Both Twitter and Facebook had mid-range median scores, with hard questions in Twitter yielding a consistent median and third quartile, both at 55. Airbnb responses were relatively stable, with scores fluctuating between 61.75 to 64. Analyzing Bard's performance across metrics (see Figure 3d), we observe that relevance ranged from scores as high as 1 for  $UC_e$ ,  $UC_m$ , and  $UC_h$  to as low as 0.2 for  $T_h$  and  $PD_e$ . Clarity was similarly distributed, with certain questions like  $UC_e$  receiving high scores of 0.9, while others, such as  $DM_h$ , only achieved a score of 0.1. Accuracy proved to be a challenging area, with the lowest score being -1 for several questions, including  $SE_h$ ,  $PD_e$ , and  $AEE_e$ . However, the model managed to score 0.7 for  $UC_e$ . Completeness ranged from a notable 0.9 for  $SE_e$  to less promising results like -0.3 for  $SE_h$ . The Reference metric had its highs and lows, with the highest score being 1 for  $T_e$  and several instances of -1, indicating an inconsistency in this domain.

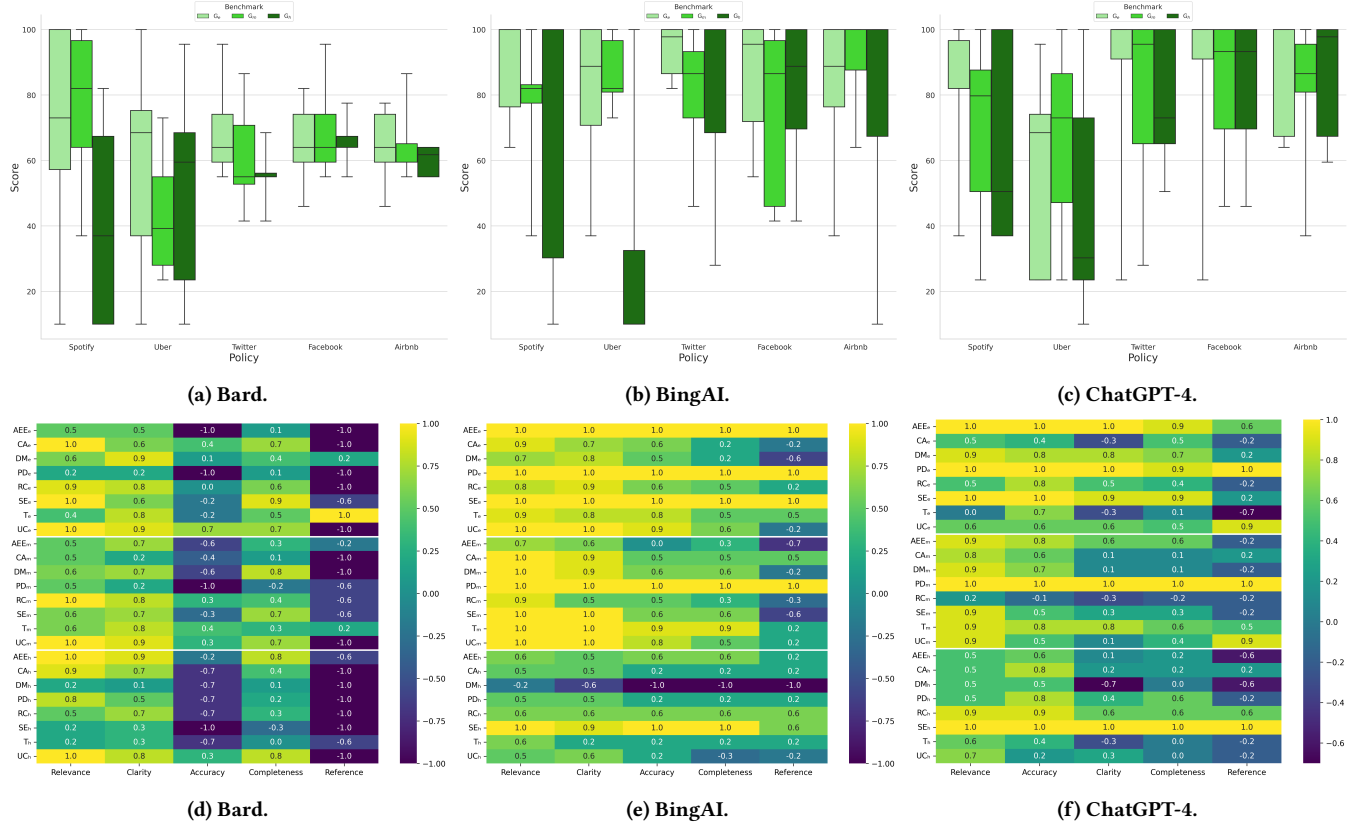
**BingAI:** BingAI exhibited a mix of outstanding and lackluster performances (Figure 3b). For Spotify, it achieved perfect medians of 10 for  $G_e$  and  $G_h$ , but the range in  $G_h$  was wide, from 10 to 100. The Uber policy was challenging, especially in  $G_h$ , with a median of just 1 and a narrow range, indicating a uniform struggle. Twitter and Facebook policies saw robust results, with medians consistently above 86.5. For Airbnb questions, BingAI's performance was notable, particularly in  $G_m$  and  $G_h$  categories, where the system reached a perfect median score of 100. BingAI demonstrated great performance in Relevance, particularly for questions like  $T_e$ ,  $PD_m$ ,  $UC_e$ , and  $SE_e$ , all scoring a perfect 1, but also showed weaker areas with scores like -0.2 for  $DM_h$  (Figure 3e). Clarity maintained a consistent trend, with scores predominantly leaning toward the higher end. For Accuracy, BingAI had top-performing scores in areas like  $AEE_e$ ,  $PD_m$ , and  $SE_h$ , but faltered in others, achieving a score of -1 for  $DM_h$ . In terms of Completeness, it exhibited excellence in  $SE_e$  and  $SE_h$ , both scoring 1, but saw a drop in areas like  $DM_h$ . The Reference scores varied, ranging from 1 in  $PD_e$  and  $PD_m$  to lows of -1 in areas such as  $DM_h$ .

## 6.3 Assessing the Ability to Recall Learned Privacy Policy Knowledge

The purpose of this experiment is to assess the performance of the systems when the privacy policy is not given explicitly and hence the system has to rely on the information it obtained when it was trained. The main question it seeks to answer is how well the system retains and recalls privacy policies in which it was trained.

**ChatGPT-4:** We observe that, considering the Spotify policy, ChatGPT-4's performance ranges between 23.5 and 100 scores in the  $G_e$  category (Figure 4c). Despite this variability, a strong median of 84.25 indicates its overall competence. The consistency was further emphasized by the narrow interquartile range (80.875 to 92.125). In the Uber policy, all three categories saw the model reaching its zenith with maximum scores of 100. For Twitter and Facebook, the

Proceedings on Privacy Enhancing Technologies YYYY(X)



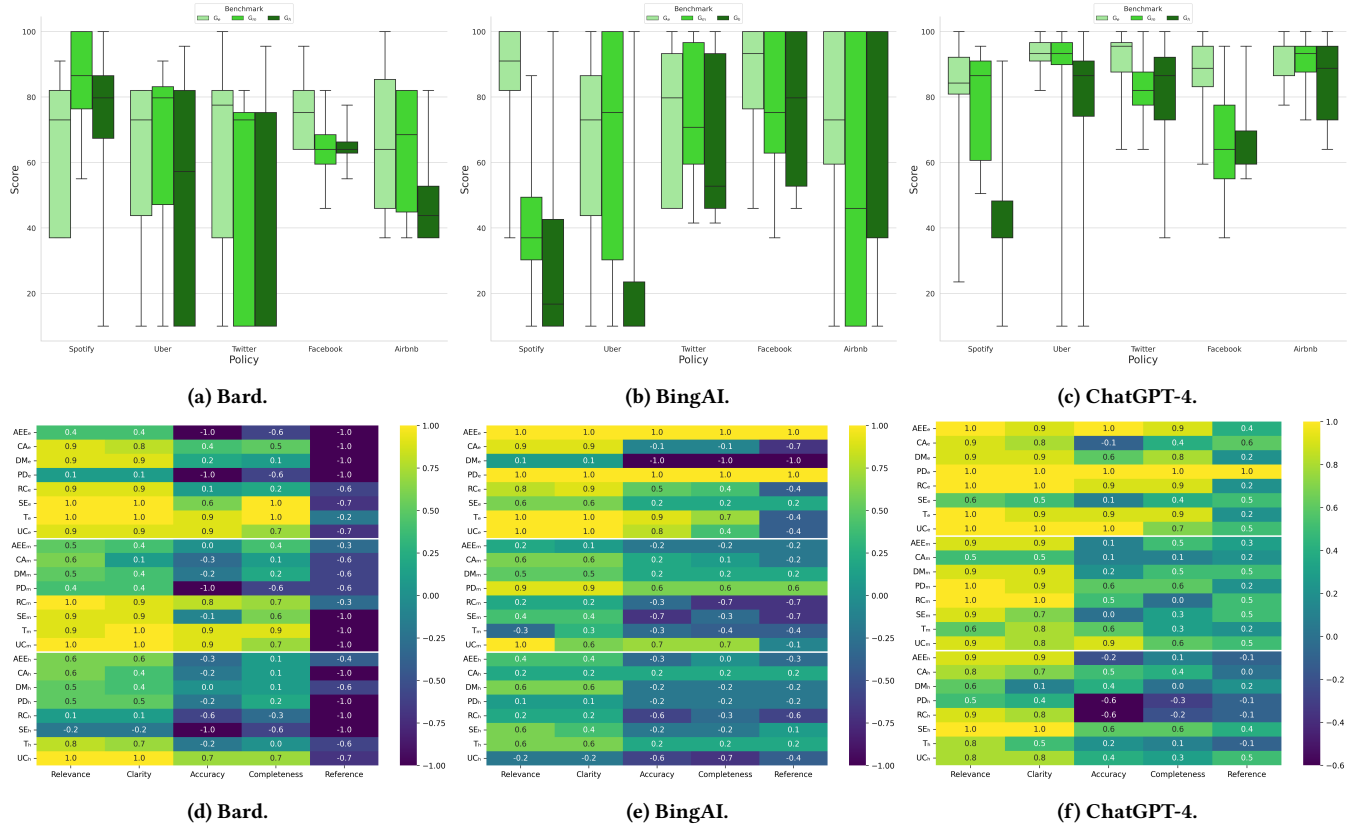
**Figure 3: Score distribution across varying difficulty levels ( $G_e$ ,  $G_m$ ,  $G_h$ ) for paraphrased privacy policy questions applied to five privacy policies.**

medians (95.5 and 88.75, respectively for  $G_e$ ) were strong, and the compact interquartile ranges again indicated reliable performances. Airbnb’s policy mirrored a similar trend with a median above 88.75. ChatGPT-4 predominantly had scores close to 1 in Relevance across different question complexities, with only a slight dip to 0.5 for  $PD_h$  (Figure 4f). Clarity remained fairly consistent, with many of its scores ranging between 0.8 to 1, but there was a notable drop to 0.1 for  $DM_h$ . In terms of Accuracy, while GPT-4 generally performed well in simpler and medium questions, there was a clear reduction in its performance in harder questions, dropping as low as -0.6 in the  $PD_h$  and  $RC_h$ . Completeness scores demonstrated a similar trend with higher scores in easier categories and diminishing results in the harder questions, the lowest being -0.3 for  $PD_h$ . The Reference, however, remained relatively low throughout, with a peak score of 1 for  $PD_e$  and a dip to -0.1 in several harder questions.

**Bard:** Bard displayed wider variability than ChatGPT-4 (Figure 4a). In the Spotify policy, the  $G_e$  category witnessed a spread from 37 to 91, suggesting more variance in its responses. The broader interquartile range (37 to 82) compared to ChatGPT-4 underlined this. Uber’s  $G_h$  difficulty indicated considerable inconsistency, with the lowest score being 10 and Q1 also at 10, suggesting that 25% of responses were at the floor of the scoring metric. Twitter’s  $G_h$  category further echoed this inconsistency, with both minimum and Q1 at 10. However, Facebook and Airbnb policies in the  $G_m$  difficulty

showed tighter interquartile ranges, hinting at better consistency. Bard maintained a high Relevance, predominantly fluctuating between 0.4 to 1 in  $G_e$  and  $G_m$  questions, but saw a drastic decline for  $SE_h$ , scoring -0.2 (Figure 4d). Its Clarity mostly mirrored ChatGPT-4’s pattern, though it had a steeper drop in  $G_h$  questions, reaching as low as -0.2 in the  $SE_h$  category. Accuracy exhibited significant variability, with scores ranging from a high of 0.9 in  $G_e$  like  $T_e$  and  $UC_e$ , to a troubling -1 in harder ones like  $SE_h$ ,  $PD_m$ , and  $PD_e$ . Completeness varied considerably as well, with scores peaking at 1 for  $T_e$  and plummeting to -0.6 in  $SE_e$  and  $PD_e$ . Reference scores were particularly notable for Bard due to their consistent negative values, dropping as low as -1 for multiple questions, suggesting possible issues with citation or source integrity.

**BingAI:** BingAI showcased a peculiar trend (Figure 4b). For Spotify’s  $G_e$  category, it ranged from 37 to a perfect 100, with a commendable median of 91. Yet, the  $G_m$  difficulty revealed stark contrasts, spanning from 1 to 86.5, with a median dropping to 37. This drastic disparity between  $G_e$  and  $G_m$  was further underscored by the interquartile range shift from 82-100 in  $G_e$  to a much broader 30.25-49.375 in  $G_m$ . Similarly, Uber’s  $G_h$  difficulty reflected a pronounced inconsistency with both the minimum and 25% of scores (Q1) languishing at 10, while the upper quartile (Q3) stretched to 23.5. Notably, in Airbnb’s  $G_h$  category, BingAI achieved a 10 median indicating that over half of its responses received the maximum



**Figure 4: Score distribution across varying difficulty levels ( $G_e$ ,  $G_m$ ,  $G_h$ ) for privacy policy questions applied to previously learned five privacy policies.**

score, though its minimum at 1 demonstrates the presence of some extreme outliers. BingAI’s performance in Relevance started strong, reaching 1 in categories like  $T_e$ ,  $UC_e$ , and  $PD_e$ , but faltered for  $G_m$  questions like  $T_m$ , which scored -0.3 as shown in Figure 4e. Clarity remained relatively stable, with many scores hovering around the 0.6 to 1 range. However, its accuracy was inconsistent, dropping to -1 for  $DM_e$  but redeeming itself with scores like 1 in  $PD_e$ . Completeness scores were highly variable, from a full score of 1 for  $PD_e$  to a concerning -1 for  $DM_e$ . As for the Referencing, it scored negatively for most of the questions.

In summary, the evaluation of the three Large Language Models, ChatGPT-4, Bard, and BingAI, revealed intricate patterns of strengths and challenges across the performance criteria. GPT-4 consistently showed high proficiency across policies, with its concentrated scores emphasizing reliability and consistency but also revealing potential weaknesses in referencing. While exhibiting proficiency in specific domains, Bard displayed broader variabilities and pronounced inconsistencies, especially in challenging contexts, with particular challenges in referencing and marked variability in accuracy and completeness for more complex questions. BingAI’s performance was a blend of exemplary moments counterbalanced by stark inconsistencies across all criteria.

## 6.4 Assessing the Quality of Responses to Privacy Regulation Questions

This experiment aims to examine the quality of responses generated by the systems for questions concerning the CCPA and GDPR data protection regulations. Figure 5 shows the results obtained after executing the privacy regulation benchmark for both data protection regulations. Both ChatGPT-4 and BingAI excelled in answering privacy regulation queries, with ChatGPT-4 consistently achieving top scores across every metric. While Bard demonstrated good performance, it consistently struggled to provide accurate references, placing it behind the other two models.

For all six questions (i.e.,  $PR_1$ , till  $PR_6$ ), ChatGPT-4 and BingAI responses were accurate, relevant, comprehensive, and included correct references to regulation details. BingAI’s scores took a hit due to its tendency to refer to online articles for its information rather than directly citing the articles from the GDPR and CCPA, a practice that ChatGPT-4 consistently followed. On the other hand, Bard’s responses for questions  $PR_1$ ,  $PR_2$  and  $PR_5$  scored 0.5 for completeness as they lacked some details. Also, the responses scored -1 across all questions wrt the reference metric for both CCPA and GDPR.

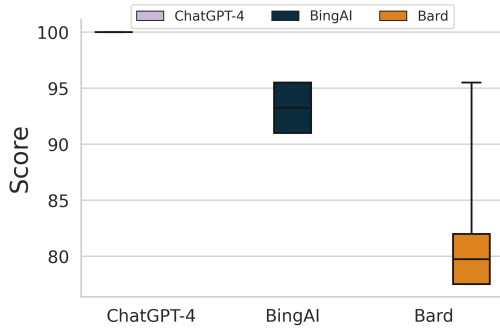


Figure 5: Scores for privacy regulation questions.

## 7 DISCUSSION

While, up to the author’s knowledge, no specific Generative AI-based Privacy Assistant (GenAIPA) has been proposed yet, our experiments indicate that current general-purpose genAI models can be a good starting point. Bard, BingAI, and ChatGPT-4 demonstrated commendable capabilities when confronted with GENAIPABENCH questions. Of course, there are also challenges the systems encountered which require further refinement and exploration.

When addressing questions related to an organization’s privacy policies, all systems obtained a fairly high score for low to medium-difficulty questions. In particular, BingAI emerged as the most consistent performer for those, demonstrating superior outcomes across most metrics. Interestingly, when paraphrased versions of the questions were used, BingAI performed worse than the other systems. This inconsistency highlights that some systems might expect users to express their questions in certain ways which would be an issue given the difference in perception about privacy among the general public [58].

The performance of all the systems declined when confronted with more complex questions, revealing limitations in handling situations requiring advanced reasoning or specialized knowledge. This is especially concerning given that these are the questions for which the general public might need more help. Of particular concern was a disconnect between the relevance and clarity of generated responses and their factual accuracy and completeness. Responses that were substantially incorrect were often presented coherently and relevantly, posing the risk of misleading users. Furthermore, we observed frequent issues wrt references that often point to outdated or incorrect data from the model’s training set, rather than the most recent privacy policy information (which was provided to the systems).

The three systems showed a strong understanding of the two privacy regulations evaluated. This might be due to the fact that there has been more discussion about data privacy regulations online than about specific privacy policies. This means that the underlying models of the three systems have potentially been trained on more information relevant to the privacy regulations. However, while the performance on the privacy regulations benchmark was excellent, we observe again the challenge of proper reference (especially in

the case of Bard). While this might not be concerning for the average user, it could be a problem for, for instance, small businesses using these systems to understand how to adapt their operations to the regulations in the countries where they operate.

In summary, we identified several critical challenges that GenAIPAs must overcome to become reliable at assisting users with their privacy questions. First, current genAI systems often fail to recognize and correct errors in their responses, posing a risk of disseminating inaccurate information on vital privacy matters. Second, they lack transparency in their reasoning, leaving users uncertain about the reliability of the provided answers, particularly when questions involve nuanced or ambiguous scenarios. Lastly, they demonstrated inconsistencies when faced with repeated queries and showed issues in recalling accurate source material, undermining their overall reliability and stability. Addressing these challenges is essential for future GenAIPAs to improve not only their accuracy but also to gain users’ trust through transparent and consistent performance. This highlights the need for models specialized in the privacy domain fine-tuned, particularly in handling complex questions, to maintain response consistency and accuracy. Such models should also pay special attention to maintaining their knowledge updated which is essential given the continuously changing landscape of privacy policies and regulations. These findings align with previous analysis and existing literature in other domains, confirming that while large language models like ChatGPT, Bard, and BingAI excel in general language tasks, their performance can vary significantly when applied to specialized domains. Hence, our analysis is, up to the author’s knowledge, the first to assess that this is the case in the domain of data privacy.

## 8 CONCLUSION AND FUTURE WORK

The emergence of generative AI and their ability to summarize text and answer questions generating human-like text presents an opportunity to develop more sophisticated privacy assistants (GenAIPAs). Due to the implications for individuals of receiving wrong information that might impact their privacy, it is required to evaluate such systems properly. In this paper we have presented a benchmark, GENAIPABENCH, to evaluate future GenAIPAs which includes questions about privacy policies and data privacy regulations, evaluation metrics, and annotated privacy documents. Our evaluation of popular genAI technology, including ChatGPT, Bard, and BingAI, shows promise for the technology but highlights that significant work remains to enhance their capabilities in handling complex queries, ensuring accuracy, maintaining response consistency and citing proper sources. We plan to continue expanding GENAIPABENCH with more annotated answers for a larger number of privacy documents. We also aim to develop the infrastructure to perform a periodic evaluation of current and future versions of genAI and GenAIPA systems.

## REFERENCES

- [1] P. Voigt and A. von dem Bussche, “The eu general data protection regulation (gdpr): A practical guide,” *Springer*, vol. 2, no. 1, pp. 1–16, 2017.
- [2] J. Greenberg and J. Maier, “California consumer privacy act (ccpa): Compliance guide,” *Business Law Today*, vol. 30, no. 3, pp. 1–11, 2020.
- [3] D. Solove, *Nothing to hide: The false tradeoff between privacy and security*. Yale University Press, 2013.

- [4] J. A. Obar and A. Oeldorf-Hirsch, "The biggest lie on the internet: Ignoring the privacy policies and terms of service policies of social networking services," in *TPRC 44: The 44th Research Conference on Communication, Information and Internet Policy*, 2018.
- [5] M. Langheinrich, "Privacy and mobile devices," *ACM SIGMOBILE Mobile Computing and Communications Review*, vol. 5, no. 1, pp. 34–44, 2001.
- [6] M. Ackerman, L. Cranor, and J. Reagle, "Privacy policies that people can understand," in *Proceedings of the SIGCHI conference on Human factors in computing systems*, pp. 415–422, ACM, 2001.
- [7] A. Cavoukian, "Privacy by design: The 7 foundational principles," in *2010 33rd International Conference on Privacy and Data Protection*, pp. 2–58, IEEE, 2010.
- [8] S. Wilson, S. Komanduri, G. Norcie, A. Acquisti, P. Leon, and L. Cranor, "Summarizing privacy policies with crowdsourcing and natural language processing," in *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pp. 2363–2374, ACM, 2016.
- [9] B. Knijnenburg and A. Kobsa, "Personalized privacy assistants for the internet of things: enabling user control over privacy in smart environments," in *Adjunct Proceedings of the 2013 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, pp. 1603–1608, ACM, 2013.
- [10] Y. Zhang, Y. Chen, and N. Li, "Privacy risk analysis for mobile applications," *IEEE Transactions on Dependable and Secure Computing*, vol. 15, no. 6, pp. 968–981, 2016.
- [11] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," 2018.
- [12] Y. Belinkov, I. Dagan, S. Shieber, and A. Subramanian, "Lama: Language-agnostic model agnosticism," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 215–225, 2020.
- [13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [14] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Gpt-3: Language models are few-shot learners," *arXiv preprint arXiv:2005.14165*, 2021.
- [15] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," *OpenAI Blog*, vol. 1, no. 8, p. 9, 2019.
- [16] J. Gao, M. Galley, and L. Li, *Neural Approaches to Conversational AI: Question Answering, Task-oriented Dialogues and Social Chatbots*. Now Foundations and Trends, 2019.
- [17] T. W. Bickmore and R. W. Picard, "Establishing and maintaining long-term human-computer relationships," *ACM Transactions on Computer-Human Interaction (TOCHI)*, vol. 12, no. 2, pp. 293–327, 2005.
- [18] B. Li, X. Wu, L. Qin, and J. Huang, "Alice: A conversational agent for financial planning," in *International Conference on Web Intelligence*, pp. 1163–1167, 2017.
- [19] K. K. Fitzpatrick, A. Darcy, and M. Vierhile, "Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): A randomized controlled trial," *JMIR Mental Health*, vol. 4, no. 2, p. e19, 2017.
- [20] R. Winkler, M. Söllner, and S. Neuweiler, "Evaluating the engagement with conversational agents: Experiments in education and health," in *International Conference on Design Science Research in Information Systems and Technology*, pp. 102–114, 2018.
- [21] A. Wang, K. Cho, and M. Lewis, "Truthfulqa: Measuring how models mimic human falsehoods," *arXiv preprint arXiv:2109.07958*, 2021.
- [22] T. Schick, A. Lauscher, and I. Gurevych, "'it's not a bug, it's a feature': Unwanted model outputs as bugs in ai systems," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2849–2859, 2021.
- [23] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, "On the dangers of stochastic parrots: Can language models be too big?," *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1354–1360, 2021.
- [24] L. Chen, M. Zaharia, and J. Zou, "How is chatgpt's behavior changing over time?," 2023.
- [25] Y. Ge, W. Hua, K. Mei, J. Ji, J. Tan, S. Xu, Z. Li, and Y. Zhang, "Openagi: When llm meets domain experts," 2023.
- [26] W.-C. Kang, J. Ni, N. Mehta, M. Sathiamoorthy, L. Hong, E. Chi, and D. Z. Cheng, "Do llms understand user preferences? evaluating llms on user rating prediction," 2023.
- [27] X. Liu, H. Yu, H. Zhang, Y. Xu, X. Lei, H. Lai, Y. Gu, H. Ding, K. Men, K. Yang, S. Zhang, X. Deng, A. Zeng, Z. Du, C. Zhang, S. Shen, T. Zhang, Y. Su, H. Sun, M. Huang, Y. Dong, and J. Tang, "Agentbench: Evaluating llms as agents," 2023.
- [28] Y. Bang, S. Cahyawijaya, N. Lee, W. Dai, D. Su, B. Wilie, H. Lovenia, Z. Ji, T. Yu, W. Chung, Q. V. Do, Y. Xu, and P. Fung, "A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity," 2023.
- [29] T. Zhang, F. Ladhak, E. Durmus, P. Liang, K. McKeown, and T. B. Hashimoto, "Benchmarking large language models for news summarization," 2023.
- [30] A. Ravichander, A. W. Black, S. Wilson, T. B. Norton, and N. M. Sadeh, "Question answering for privacy policies: Combining computational and legal perspectives," *ArXiv*, vol. abs/1911.00841, 2019.
- [31] N. Sadeh, A. Acquisti, T. D. Breaux, L. F. Cranor, A. M. McDonald, J. R. Reidenberg, N. A. Smith, F. Liu, N. C. Russell, F. Schaub, et al., "The usable privacy policy project," in *Technical report, Technical Report, CMU-ISR-13-119*, Carnegie Mellon University, 2013.
- [32] A. Wilson, F. Schaub, A. A. Dara, F. Liu, S. Cherivirala, P. G. Leon, M. S. Andersen, S. Zimmeck, K. M. Sathyendra, N. C. Russell, et al., "The creation and analysis of a website privacy policy corpus," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pp. 1330–1340, Association for Computational Linguistics, 2016.
- [33] C.-H. Chiang and H. yi Lee, "Can large language models be an alternative to human evaluations?," 2023.
- [34] J. Liu, C. S. Xia, Y. Wang, and L. Zhang, "Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation," 2023.
- [35] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "Squad: 100,000+ questions for machine comprehension of text," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, 2016.
- [36] M. Joshi, E. Choi, D. Weld, and L. Zettlemoyer, "Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1601–1611, 2017.
- [37] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C. D. Manning, C. Ré, D. Acosta-Navas, D. A. Hudson, E. Zelikman, E. Durmus, F. Ladhak, F. Rong, H. Ren, H. Yao, J. Wang, K. Santhanam, L. Orr, L. Zheng, M. Yuksekgonul, M. Suzgun, N. Kim, N. Guha, N. Chatterji, O. Khattab, P. Henderston, Q. Huang, R. Chi, S. M. Xie, S. Santurkar, S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang, and Y. Koreeda, "Holistic evaluation of language models," 2022.
- [38] Information and Privacy Commissioner of Ontario, "7 foundational principles of privacy by design," n.d.
- [39] I. Pollach, "What's wrong with online privacy policies?," *Commun. ACM*, vol. 50, pp. 103–108, 09 2007.
- [40] O. of the Australian Information Commissioner, "Office of the Australian information commissioner - oaic," 2023. Accessed: May 3, 2023.
- [41] "ISO/IEC 29100:2011 - Information technology - Security techniques - Privacy framework," 2011.
- [42] GDPR.eu, "Gdpr faqs," 2021.
- [43] C. A. General, "Ccpa faqs," 2021.
- [44] Future of Privacy Forum, "Best Practices for Consumer-Facing Privacy Notices and Consent Forms," June 2020.
- [45] K. Martin, "Ethical implications and accountability of algorithms," *Journal of Business Ethics*, vol. 160, 12 2019.
- [46] K. A. Bamberger and D. K. Mulligan, "Privacy on the books and on the ground," *Stanford Law Review*, vol. 63, p. 247, 2011.
- [47] T. W. Bickmore, L. M. Pfeifer, D. Schulman, and L. Yin, "Maintaining continuity in longitudinal, relational agents for chronic disease self-care," *Journal of Medical Systems*, vol. 42, no. 5, p. 91, 2018.
- [48] E. Luger and A. Sellen, "Like having a really bad pa: the gulf between user expectation and experience of conversational agents," *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, 2016.
- [49] Q. V. Liao, Y. Gao, Y. Wu, and Y. Zhang, "Evaluating the effectiveness of human-machine collaboration in human-in-the-loop text classification," *Proceedings of the 2019 Conference on Human Information Interaction and Retrieval*, 2019.
- [50] H. Choi, J. Park, and Y. Jung, "The role of privacy fatigue in online privacy behavior," *Computers in Human Behavior*, vol. 81, pp. 42–51, 2018.
- [51] H. P. Grice, "Logic and conversation," *Speech acts*, 1975.
- [52] C. Jensen and C. Potts, "Privacy policies as decision-making tools: an evaluation of online privacy notices," in *Proceedings of the SIGCHI conference on Human factors in computing systems*, pp. 471–478, ACM, 2004.
- [53] N. M. Radziwill and M. C. Benton, "Evaluating quality of chatbots and intelligent conversational agents," *arXiv preprint arXiv:1704.04579*, 2017.
- [54] J. Savelka and K. D. Ashley, "Extracting case law sentences for argumentation about the gdpr," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, ACL*, 2016.
- [55] OpenAI, "Gpt-4 technical report," 2023.
- [56] E. H. Hiebert, "Unique words require unique instruction,"
- [57] R. Flesch, "Flesch-kincaid readability test. retrieved october," 2007.
- [58] W. M. Steijn and A. Vedder, "Privacy under construction: A developmental perspective on privacy perception," *Science, Technology, & Human Values*, vol. 40, no. 4, pp. 615–637, 2015.

Category	Question	Difficulty
Transparency	Does the policy avoid using technical phrases or legalese that might mislead users?	Easy
	What is the business' position on demands from the government for user data?	Medium
	How does the privacy policy manage any conflicts of interest with regard to the application or exchange of data?	Hard
User Control	Do users have access to their data and privacy settings?	Easy
	Are there straightforward methods for users to request data access or deletion?	Medium
	How does the business handle user approval and revocation of approval?	Hard
Data Minimization and Purpose Limitation	Does the business reduce storage times?	Easy
	How is user data aggregated or masked to safeguard personal information?	Medium
	Are there any limitations on how data may be processed for particular situations or purposes?	Hard
Security and Encryption	Do user communications employ end-to-end encryption?	Easy
	How are user data being accessed without authorization prevented?	Medium
	How are security-related incidents or thefts of information handled and informed to users?	Hard
Privacy by Design and Innovation	Does the organization carry out privacy impact analyses?	Easy
	Have any privacy-enhancing technologies, such as differential privacy, been implemented?	Medium
	How does the use of automated decision-making or profiling by the firm affect the privacy of its customers?	Hard
Responsiveness and Communication	Is the privacy policy updated frequently and made available to users?	Easy
	Does a procedure exist to handle user privacy complaints?	Medium
	Does the business release transparency reports describing interactions with law enforcement or government agencies or data requests for surveillance?	Hard
Accessibility, User Education and Empowerment	Do employees receive training on handling sensitive information and best practices for data privacy?	Easy
	How are user choices for data privacy maintained across many platforms or devices?	Medium
	Does the company provide easy-to-use resources, such as tutorials or guides, to help users successfully manage their privacy settings and comprehend their data rights?	Hard
Compliance and Accountability	Does this policy conform with all relevant privacy laws and regulations?	Easy
	What measures are made to guarantee that data processors and subprocessors respect privacy standards?	Medium
	Is there a procedure in place at the business for reporting and dealing with privacy violations or non-compliance issues, both internally and with third-party vendors?	Hard

**Table 2: Paraphrased privacy policy questions.**

## Appendix A PARAPHRASED QUESTIONS

Table 2 lists out the paraphrased version of the privacy policy questions defined in Section 4. Quillbot<sup>5</sup> was used to automatically generate these questions.

<sup>5</sup><https://www.quillbot.com/>

## Appendix B EVALUATOR

The GENAIPABENCH includes a component evaluator whose goal is to communicate with the GenAIPA sharing the privacy documents and questions and collecting answers and summaries (see Algorithm 1). This is accomplished in a two-step process: response generation (see Step 1 of Algorithm 1) and summary generation (See Step 2 of Algorithm 1). In the response generation stage, an evaluator has two ways to introduce GenAIPA to the privacy document. In the first instance, the privacy document is partitioned into smaller sections, accommodating the input capacity of GenAIPA.



Simultaneously, an empty list intended for holding conversation sets is initialized. Next, a new conversation is initiated. After the initiation of a conversation with GenAIPA, the policy sections are systematically introduced to the GenAIPA. For the initiated conversation, the question list is shuffled, each question within the list is utilized as a prompt for the GenAIPA. The responses to each question are collected, organized into question-answer pairs, and added to thier respective conversation set. This entire procedure is repeated thrice for each privacy policy, resulting in a collection of three distinct conversation sets per policy. In the second instance, the privacy policy corpus is leveraged to ascertain whether the GenAIPA model was trained on a previous version of these policies. Using the evaluator, queries are then directed to GenAIPA to verify the version and date of the privacy policy with which it is familiar, assuming the system has access to the given policy. GenAIPA is thus requested to disclose the version and date of the privacy policy it last accessed. The disclosed version and date are subsequently used as a reference to retrieve the identical version of the privacy policy via the *Wayback Machine*<sup>6</sup>, an online digital archive that allows access to historical versions of internet content. This retrieved version serves as the gold standard in our analysis, providing accurate answers to the privacy questions that will be later asked to the GenAIPA. The ensuing step involves prompting GenAIPA with questions related to the privacy policy, using the evaluator which resembles step 1 of Algorithm1 to generate responses. The final phase of the procedure involves the evaluator comparing the GenAIPA responses with the annotated answers that were previously identified by analysts through the detailed examination of the retrieved policies from the Wayback Machine.

---

**Algorithm 1** GenAIPA Response and Summary Generation

---

```

1: procedure GENANDEVALRESP(Privacy Document PD, Questions Q, Runs r)
2:   Step 1: Response Generation
3:    $A \leftarrow \emptyset$ 
4:   for  $i = 1$  to  $r$  do
5:      $A_i \leftarrow \emptyset$ 
6:     IntroducePrivacyDocument(PD)
7:      $Q' \leftarrow \text{ShuffleQuestions}(Q)$ 
8:     for each  $q$  in  $Q'$  do
9:        $A_i \leftarrow \text{SendQuestion}(q)$ 
10:     $A \leftarrow \{q, A_i\}$ 
11:   Step 2: Summary Generation
12:    $C \leftarrow \text{SplitPolicy}(PD)$ 
13:   Initialize empty list  $D$ 
14:   for  $c \in C$  do
15:      $d \leftarrow \text{Prompt}(c)$ 
16:     Add  $d$  to  $D$ 
17:    $S \leftarrow \text{GenerateSummary}(D)$ 
18:   Return  $S$ 

```

---

In the summary generation step, the procedure begins by Utilizing the evaluator to query GenAIPA to craft a succinct summary of the privacy policy (see Algorithm 2). This involves creating an

efficient summary by feeding segments of the policy to GenAIPA. The evaluator initiates the process by partitioning the input policy into several segments, each supplied to GenAIPA individually. After all, segments have been processed, GenAIPA is queried to create a thorough summary. This synthesized summary is then utilized as the input for the ensuing stage. Once the summary is generated, it is used as input for GenAIPA by the evaluator. Subsequently, our curated set of privacy-related questions (see Section 4), are posed as prompts by the evaluator to GenAIPA in order to assess the depth and accuracy of the information contained in the summary. In response, GenAIPA crafts answers based on the summary. This procedure aligns with the response generation of the evaluator tool. For every question, the evaluator generates three different responses and evaluates the quality of the information derived from the summary by using the response evaluation (see Algorithm1).

---

**Algorithm 2** GenAIPA Response Evaluation

---

```

1: procedure EVALUATERESPONSE(Scores S provided by analyst)
2:    $P \leftarrow \emptyset$ 
3:   for  $i = 1$  to  $|S[1]|$  do
4:      $\bar{S}_i \leftarrow \frac{1}{|S|} \sum_{score \in S} score[i]$ 
5:     Categorize  $\bar{S}_i$  as Green, Yellow, or Red
6:      $P \leftarrow P \cup \text{Category}$ 
7:   return  $P$ 

```

---

In the GenAIPA Response Evaluation process (see Algorithm 2), the analyst scrutinizes the generated responses based on five key features: Relevance, Accuracy, Clarity, Completeness, and Reference to policy sections. Scores for each feature are provided as input S. The procedure starts by initializing an empty set P and then iterates through each score  $S_i$  in S, accumulating them for subsequent categorization. This evaluation is performed for each set of scores, representing multiple runs, and the average scores for each run are determined. Ultimately, an overall average score is calculated across all runs. This average score is then categorized into one of three groups: Green, Yellow, or Red, based on predefined criteria (see Section 5. The category is stored in set P, which upon completion of all iterations, contains the categorized average scores for all runs. The primary goal of this evaluation is to offer insights into GenAIPA’s capabilities in generating privacy policy-related responses and identify areas of potential improvement.

## Appendix C ADDITIONAL EXPERIMENTS

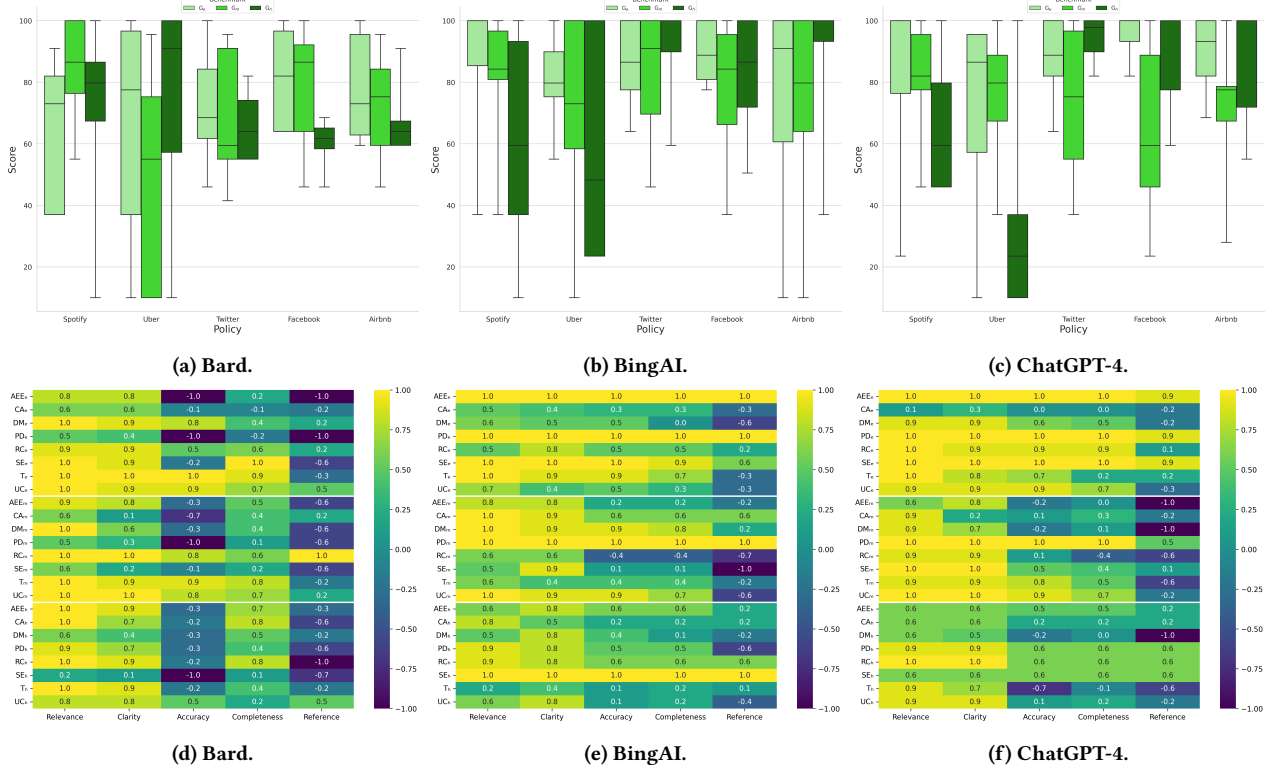
This section includes additional experiments and results from our evaluation of ChatGPT-4, Bard, and BingAI using GENAIPABENCH.

### C.1 Assessing the Quality of Privacy Policy Summaries

This experiment aims to examine the quality of the summary generated for privacy policies. The results (see Figure 6) show that Bard was consistently good across tasks, BingAI performed best in medium-difficulty challenges, and ChatGPT-4, while generally skilled, sometimes struggled with complex questions. Both Bard and BingAI had issues with citing references.

<sup>6</sup><https://archive.org/web/>

Proceedings on Privacy Enhancing Technologies YYYY(X)



**Figure 6: Score distribution across varying difficulty levels ( $G_e$ ,  $G_m$ ,  $G_h$ ) for original privacy benchmark questions applied to three distinct privacy policies.**

**Bard:** Bard’s performance across the various policies and difficulty levels exhibited a balanced range (Figure 6a). For  $G_e$ , Bard’s minimum scores, particularly with the Spotify policy, were competitive at 37, although its maximum scores did not consistently reach the peak value of 100 as seen in some other models. For  $G_m$ , Bard showcased a broad range of scores, with a notable low of 46 in policies like Spotify, Twitter, and Facebook. While its maximum scores in this bracket mostly touched 100, it did not achieve this for Twitter. For  $G_h$ , Bard managed to reach the top score across multiple policies, a testament to its ability to handle intricate tasks. The median scores, especially in the hard category, pointed to a balanced performance without extreme lows for the majority. Bard consistently scored highly across most metrics (Figure 6d). In terms of relevance, it consistently showed scores close to 1 for a majority of the questions, suggesting that its outputs were highly relevant to the user’s queries. Its clarity scores were generally between 0.7 and 1, demonstrating that its outputs were clear and comprehensible. However, Bard performed poorly in terms of accuracy, especially for questions like  $T_h$ ,  $DM_m$ , and  $AEE_e$  where the scores dipped to -0.2 or lower. Regarding completeness, the model consistently hovered between 0.4 and 1, suggesting a generally satisfactory, though not always complete, response to queries. Bard performed also poorly in terms of reference with negative scores for many of the questions.

**BingAI:** BingAI’s performance leaned towards higher scores in many instances (Figure 6b). For  $G_e$ , it matched Bard’s minimum score for the Spotify policy at 37 but showed a propensity to achieve the highest score of 100 across several policies. For  $G_m$ , BingAI achieved a consistent upper bound of 100 across most policies and a higher third quartile (Q3) in policies such as Spotify and Twitter. Within  $G_h$ , despite variations in the minimum score, BingAI consistently achieved a score of 100, underscoring its good performance in complex scenarios. In terms of specific metrics, BingAI’s performance was slightly more varied than Bard’s (Figure 6e). Its relevance scores ranged from 0.2 to 1, meaning that while it often provided relevant answers, there were instances where its responses could be off the mark. In terms of clarity, BingAI’s scores predominantly ranged between 0.4 and 1, which, while commendable, also suggests room for improvement. The model’s accuracy ranged from 0.1 to 1, with several low scores interspersed among generally high scores. Completeness was one area where BingAI lagged behind Bard, with scores mostly falling between 0.1 and 0.9. In terms of reference, BingAI had negative values for several questions, similar to Bard, but also exhibited positive scores for some.

**ChatGPT-4:** ChatGPT-4’s performance was marked by peaks of excellence as well as certain areas of inconsistency (Figure 6c). For  $G_e$ , it started at a slightly lower minimum for the Spotify policy at 23.5 but managed to match its competitors by frequently reaching

the maximum score of 100. For  $G_m$ , it maintained a narrow range in minimum scores, with the Spotify policy again being a common low at 37. However, for  $G_h$ , ChatGPT-4 displayed some potential areas of improvement. Although it achieved the maximum score across several policies, its median scores for certain policies like Spotify and Uber lagged behind its counterparts, hinting at challenges in maintaining consistent performance in more complex tasks. ChatGPT-4’s performance was more evenly spread across the metrics (Figure 6f), without pronounced highs or lows. Its relevance scores ranged between 0.1 and 1, suggesting varying levels of match between the model’s outputs and user queries. Clarity scores were predominantly high, with most values between 0.6 and 1, indicating its outputs were largely comprehensible. Accuracy was a mixed bag, with scores as low as -0.7 for  $T_h$  and as high as 1 for  $SE_e$ . In terms of completeness, the model hovered between -0.1 and 1, showing an inconsistent ability to provide complete answers. Lastly, its reference scores were a mix of negative, zero, and positive values, indicating an unpredictable citing behaviour.

## C.2 Result Analysis

Tables 3 through 14 offer an in-depth analysis of question performance across multiple experiments. Each table dissects the results on a per-question basis for each evaluated system. The rows have been color-coded to offer a quick, visual assessment of each question’s performance. In particular, the color scale is computed as follows based on the overall quality metric  $M_{all}$  (see Section 5):

- Green: High-quality answer ( $7 < M_{all} \leq 10$ ).
- Yellow: Medium-quality answer which might contain errors/omissions ( $4 < M_{all} \leq 7$ ).
- Red: Low-quality answers which are inaccurate/unreliable ( $0 < M_{all} \leq 4$ ).

Question	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	7.3	7.75	8.65	8.2	8.65	8.11
$UC_e$	7.75	6.4	9.1	10	8.2	8.29
$DM_e$	10	3.7	8.2	9.55	8.2	7.93
$SE_e$	10	9.55	10	10	10	9.91
$PD_e$	10	9.55	10	10	10	9.91
$RC_e$	10	9.55	8.2	10	6.85	8.92
$AEE_e$	10	9.55	10	10	10	9.91
$CA_e$	2.35	1	6.4	8.65	10	5.68
$T_m$	7.75	7.75	7.75	8.65	6.85	7.75
$UC_m$	7.75	7.75	10	9.55	7.75	8.56
$DM_m$	8.2	3.7	5.5	4.6	7.75	5.95
$SE_m$	10	8.65	9.55	4.6	8.2	8.2
$PD_m$	8.2	9.55	10	10	10	9.55
$RC_m$	9.55	3.7	7.3	4.6	6.4	6.31
$AEE_m$	4.6	8.2	5.5	2.35	7.75	5.68
$CA_m$	9.55	10	3.7	7.3	2.8	6.67
$T_h$	4.6	3.7	8.65	5.95	5.5	5.68
$UC_h$	7.3	3.7	9.55	7.75	7.75	7.21
$DM_h$	4.6	1	8.2	7.75	5.5	5.41
$SE_h$	10	1	10	10	10	8.2
$PD_h$	10	3.7	10	10	10	8.74
$RC_h$	4.6	10	10	10	10	8.92
$AEE_h$	7.3	1	10	10	10	7.66
$CA_h$	4.6	1	9.1	10	10	6.94

**Table 3: Scores obtained by ChatGPT-4 for Various Privacy Policies: Assessing the Quality of Privacy Policy Summaries**

Question	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	8.2	7.3	2.35	2.35	6.4	5.32
$UC_e$	9.55	2.35	10	10	10	8.38
$DM_e$	8.2	6.4	9.1	9.1	10	8.56
$SE_e$	8.2	7.3	10	10	10	9.1
$PD_e$	10	9.55	10	10	10	9.91
$RC_e$	8.2	2.35	9.1	10	6.85	7.3
$AEE_e$	10	7.75	10	10	10	9.55
$CA_e$	3.7	2.35	10	9.1	6.4	6.31
$T_m$	9.1	6.85	10	9.1	8.65	8.74
$UC_m$	8.65	3.7	9.55	10	8.2	8.02
$DM_m$	5.5	7.75	7.3	4.6	9.55	6.94
$SE_m$	7.75	5.05	9.55	4.6	8.65	7.12
$PD_m$	10	10	10	10	10	10
$RC_m$	2.35	2.35	2.8	7.75	9.55	4.96
$AEE_m$	3.7	8.2	10	10	7.75	7.93
$CA_m$	8.2	10	4.15	9.55	3.7	7.12
$T_h$	10	1	7.3	4.6	6.85	5.95
$UC_h$	6.4	6.4	5.05	4.6	9.55	6.4
$DM_h$	3.7	2.35	5.5	8.65	5.95	5.23
$SE_h$	10	10	10	10	10	10
$PD_h$	10	3.7	6.85	10	6.4	7.39
$RC_h$	3.7	10	10	10	10	8.74
$AEE_h$	3.7	2.35	7.3	7.75	10	6.22
$CA_h$	3.7	2.35	10	10	10	7.21

**Table 4: Scores obtained by ChatGPT-4 for Various Privacy Policies: Assessing Robustness through Paraphrased Questions**

Proceedings on Privacy Enhancing Technologies YYYY(X)

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	3.7	10	8.2	8.2	7.3	7.48
$UC_e$	7.3	10	9.55	10	10	9.37
$DM_e$	8.2	8.2	10	10	10	9.28
$SE_e$	10	10	8.2	10	10	9.64
$PD_e$	10	10	10	10	10	10
$RC_e$	7.75	8.2	10	10	7.3	8.65
$AEE_e$	10	10	10	10	10	10
$CA_e$	10	10	10	10	6.4	9.28
$T_m$	1	10	10	9.55	9.55	8.02
$UC_m$	1	9.55	10	9.55	9.55	7.93
$DM_m$	1	4.6	8.2	4.6	10	5.68
$SE_m$	1	7.3	8.2	4.6	7.3	5.68
$PD_m$	10	10	10	10	10	10
$RC_m$	1	7.75	10	4.6	7.75	6.22
$AEE_m$	3.7	10	9.55	10	6.4	7.93
$CA_m$	3.7	1	8.2	8.65	1	4.51
$T_h$	10	7.3	10	4.15	6.85	7.66
$UC_h$	7.75	7.75	10	5.5	9.1	8.02
$DM_h$	10	1	10	9.55	6.4	7.39
$SE_h$	10	10	10	10	10	10
$PD_h$	10	1	10	10	6.4	7.48
$RC_h$	4.15	3.7	10	10	10	7.57
$AEE_h$	3.7	7.3	9.55	7.75	10	7.66
$CA_h$	3.7	1	10	10	10	6.94

Table 5: Scores obtained by ChatGPT-4 for Various Privacy Policies: Assessing the Quality of Responses to Privacy Policy Questions

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	9.55	9.55	9.55	8.65	7.75	9.01
$UC_e$	9.1	9.1	9.1	9.55	9.55	9.28
$DM_e$	8.2	9.55	6.4	9.1	9.55	8.56
$SE_e$	2.35	9.1	9.55	7.3	8.65	7.39
$PD_e$	10	10	10	10	10	10
$RC_e$	8.65	8.2	9.55	9.55	9.55	9.1
$AEE_e$	8.2	10	10	8.65	9.55	9.28
$CA_e$	7.75	9.1	7.75	5.95	8.65	7.84
$T_m$	9.1	9.55	9.1	3.7	7.3	7.75
$UC_m$	9.1	9.1	7.75	9.1	9.1	8.83
$DM_m$	9.55	9.55	6.4	5.95	9.55	8.2
$SE_m$	5.05	8.65	8.65	6.85	9.1	7.66
$PD_m$	8.2	10	10	4.15	10	8.47
$RC_m$	9.1	9.1	7.75	7.3	7.75	8.2
$AEE_m$	6.4	10	7.75	5.95	9.55	7.93
$CA_m$	5.05	1	8.65	9.55	9.55	6.76
$T_h$	3.7	9.1	3.7	9.55	8.2	6.85
$UC_h$	9.1	9.1	9.1	6.4	6.4	8.02
$DM_h$	1	9.1	8.2	5.5	9.55	6.67
$SE_h$	8.2	10	9.55	5.95	10	8.74
$PD_h$	3.7	1	9.1	5.95	7.3	5.41
$RC_h$	3.7	8.2	5.95	5.95	7.3	6.22
$AEE_h$	3.7	7.75	7.75	5.95	9.55	6.94
$CA_h$	3.7	6.4	10	8.65	9.55	7.66

Table 6: Scores obtained by ChatGPT-4 for Various Privacy Policies: Assessing the Ability to Recall Learned Privacy Policy Knowledge

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	8.2	6.85	9.1	9.55	8.2	8.38
$UC_e$	3.7	7.75	7.75	7.75	7.75	6.94
$DM_e$	10	5.5	7.75	7.75	1	6.4
$SE_e$	10	7.75	10	10	10	9.55
$PD_e$	10	10	10	10	10	10
$RC_e$	10	8.2	8.2	8.2	10	8.92
$AEE_e$	10	10	10	10	10	10
$CA_e$	8.65	8.65	6.4	8.2	1	6.58
$T_m$	8.2	6.85	10	8.65	1	6.94
$UC_m$	7.75	6.85	8.2	9.55	8.2	8.11
$DM_m$	9.55	7.75	10	7.3	10	8.92
$SE_m$	8.2	2.8	7.75	3.7	7.75	6.04
$PD_m$	10	10	10	10	10	10
$RC_m$	8.65	1	4.6	4.6	7.3	5.23
$AEE_m$	3.7	10	10	8.2	3.7	7.12
$CA_m$	10	10	4.6	9.55	10	8.83
$T_h$	1	2.35	8.65	10	10	6.4
$UC_h$	9.1	2.35	9.1	5.5	7.3	6.67
$DM_h$	3.7	10	5.95	5.05	10	6.94
$SE_h$	10	10	10	10	10	10
$PD_h$	8.2	7.3	10	7.75	3.7	7.39
$RC_h$	3.7	10	10	9.55	10	8.65
$AEE_h$	10	2.35	10	7.75	10	8.02
$CA_h$	3.7	2.35	10	10	10	7.21

Table 7: Scores obtained by BingAI for Various Privacy Policies: Assessing the Quality of Privacy Policy Summaries

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	10	8.2	8.2	6.85	10	8.65
$UC_e$	7.75	7.3	9.55	10	7.75	8.47
$DM_e$	7.3	3.7	8.65	7.3	7.75	6.94
$SE_e$	10	10	10	10	10	10
$PD_e$	10	10	10	10	10	10
$RC_e$	10	6.4	10	9.1	3.7	7.84
$AEE_e$	10	10	10	10	10	10
$CA_e$	6.4	9.55	8.65	5.5	7.3	7.48
$T_m$	8.2	8.2	9.1	10	10	9.1
$UC_m$	7.75	7.3	9.1	9.55	9.55	8.65
$DM_m$	8.2	8.2	10	4.15	10	8.11
$SE_m$	8.2	8.2	8.2	4.6	10	7.84
$PD_m$	10	10	10	10	10	10
$RC_m$	7.75	9.55	7.75	4.6	6.4	7.21
$AEE_m$	3.7	7.75	5.95	7.75	6.4	6.31
$CA_m$	8.65	10	4.6	9.55	10	8.56
$T_h$	10	1	2.8	10	10	6.76
$UC_h$	10	1	8.2	4.15	7.75	6.22
$DM_h$	1	1	2.8	4.6	1	2.08
$SE_h$	10	10	10	7.75	10	9.55
$PD_h$	10	1	10	10	3.7	6.94
$RC_h$	1	10	10	10	10	8.2
$AEE_h$	10	1	10	7.75	10	7.75
$CA_h$	3.7	1	10	10	10	6.94

Table 8: Scores obtained by BingAI for Various Privacy Policies: Assessing Robustness through Paraphrased Questions

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	8.2	8.2	9.55	8.2	10	8.83
$UC_e$	10	7.75	9.55	10	8.2	9.1
$DM_e$	10	7.3	8.65	7.75	8.2	8.38
$SE_e$	8.2	8.2	10	10	10	9.28
$PD_e$	10	10	10	10	10	10
$RC_e$	10	6.4	10	8.2	1	7.12
$AEE_e$	10	1	10	10	10	8.2
$CA_e$	1	10	5.05	5.95	1	4.6
$T_m$	8.2	8.2	9.55	9.55	10	9.1
$UC_m$	3.7	7.75	10	9.55	9.55	8.11
$DM_m$	6.4	7.3	10	5.95	10	7.93
$SE_m$	8.2	7.3	7.75	5.95	9.55	7.75
$PD_m$	10	10	10	10	10	10
$RC_m$	9.55	10	10	4.6	7.3	8.29
$AEE_m$	3.7	8.2	10	10	1	6.58
$CA_m$	3.7	10	4.6	7.75	1	5.41
$T_h$	3.7	3.7	4.6	4.6	1	3.52
$UC_h$	7.3	7.3	8.2	9.55	8.2	8.11
$DM_h$	1	1	7.3	6.85	1	3.43
$SE_h$	10	10	10	10	10	10
$PD_h$	10	6.4	10	10	10	9.28
$RC_h$	1	10	10	10	10	8.2
$AEE_h$	9.55	3.7	10	10	10	8.65
$CA_h$	3.7	1	10	10	10	6.94

Table 9: Scores obtained by BingAI for Various Privacy Policies: Assessing the Quality of Responses to Privacy Policy Questions

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	8.2	8.2	9.1	9.1	7.3	8.38
$UC_e$	8.2	7.3	7.75	9.55	7.3	8.02
$DM_e$	3.7	1	4.6	4.6	1	2.98
$SE_e$	10	1	4.6	10	10	7.12
$PD_e$	10	10	10	10	10	10
$RC_e$	9.55	5.5	8.2	8.65	4.6	7.3
$AEE_e$	10	10	10	10	10	10
$CA_e$	8.65	7.3	4.6	4.6	6.4	6.31
$T_m$	3.7	7.75	6.4	3.7	1	4.51
$UC_m$	8.65	7.3	9.55	6.85	8.2	8.11
$DM_m$	3.7	1	10	10	10	6.94
$SE_m$	3.7	3.7	7.75	7.3	1	4.69
$PD_m$	3.7	10	10	10	10	8.74
$RC_m$	8.65	1	6.4	4.6	1	4.33
$AEE_m$	1	10	4.15	10	1	5.23
$CA_m$	1	10	4.6	7.75	10	6.67
$T_h$	10	1	4.6	10	10	7.12
$UC_h$	2.35	1	9.1	4.6	1	3.61
$DM_h$	1	10	4.6	4.6	10	6.04
$SE_h$	10	1	4.15	5.5	10	6.13
$PD_h$	1	1	10	10	3.7	5.14
$RC_h$	1	1	4.6	5.95	10	4.51
$AEE_h$	2.35	6.4	5.95	10	3.7	5.68
$CA_h$	1	1	10	10	10	6.4

Table 10: Scores obtained by BingAI for Various Privacy Policies: Assessing the Ability to Recall Learned Privacy Policy Knowledge

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	8.2	8.2	8.2	9.55	9.55	8.74
$UC_e$	9.1	7.3	9.1	10	10	9.1
$DM_e$	8.2	9.55	5.5	9.55	9.55	8.47
$SE_e$	8.2	10	6.4	6.4	5.95	7.39
$PD_e$	3.7	1	4.6	6.4	5.95	4.33
$RC_e$	3.7	10	10	10	7.75	8.29
$AEE_e$	3.7	3.7	6.4	6.4	6.4	5.32
$CA_e$	6.4	3.7	7.3	6.85	6.85	6.22
$T_m$	8.2	8.2	9.1	9.1	8.2	8.56
$UC_m$	7.75	7.3	9.55	10	9.55	8.83
$DM_m$	10	7.3	5.95	4.6	4.6	6.49
$SE_m$	7.3	1	4.15	9.1	7.3	5.77
$PD_m$	5.5	1	5.5	6.4	5.95	4.87
$RC_m$	10	9.55	9.1	9.55	9.1	9.46
$AEE_m$	10	3.7	5.5	6.4	7.75	6.67
$CA_m$	9.1	1	5.95	8.2	5.95	6.04
$T_h$	8.2	10	7.3	4.6	5.95	7.21
$UC_h$	8.65	7.3	8.2	6.85	9.1	8.02
$DM_h$	10	1	7.75	6.85	6.4	6.4
$SE_h$	1	1	6.85	6.4	6.4	4.33
$PD_h$	3.7	10	5.5	5.5	7.75	6.49
$RC_h$	7.75	8.2	5.95	6.4	5.95	6.85
$AEE_h$	8.65	10	5.5	5.95	6.4	7.3
$CA_h$	7.75	10	5.5	5.95	5.95	7.03

Table 11: Scores obtained by Bard for Various Privacy Policies: Assessing the Quality of Privacy Policy Summaries

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	10	10	6.85	4.6	7.3	7.75
$UC_e$	7.3	6.85	7.75	8.2	7.75	7.57
$DM_e$	10	4.15	9.55	5.95	7.75	7.48
$SE_e$	10	8.2	5.95	6.4	4.6	7.03
$PD_e$	1	2.35	5.5	5.95	5.95	4.15
$RC_e$	6.4	6.85	5.95	7.75	6.4	6.67
$AEE_e$	3.7	1	5.95	6.4	6.4	4.69
$CA_e$	7.3	7.3	7.3	7.3	5.95	7.03
$T_m$	8.2	2.8	8.65	9.55	8.65	7.57
$UC_m$	7.3	6.85	7.75	7.75	6.4	7.21
$DM_m$	8.2	3.7	5.5	6.4	5.95	5.95
$SE_m$	10	5.05	5.5	6.4	5.5	6.49
$PD_m$	3.7	2.8	4.15	5.95	5.95	4.51
$RC_m$	10	7.3	4.6	7.3	6.85	7.21
$AEE_m$	9.55	4.15	5.5	5.5	5.95	6.13
$CA_m$	3.7	2.35	6.85	5.95	5.95	4.96
$T_h$	1	5.5	5.5	6.4	5.5	4.78
$UC_h$	8.2	6.85	6.85	7.75	6.4	7.21
$DM_h$	1	1	5.5	7.75	6.4	4.33
$SE_h$	1	2.35	4.15	6.4	5.5	3.88
$PD_h$	3.7	6.4	5.5	5.5	5.5	5.32
$RC_h$	6.4	2.35	5.5	6.4	5.95	5.32
$AEE_h$	7.75	9.55	5.95	6.4	6.4	7.21
$CA_h$	3.7	6.85	5.5	6.4	6.4	5.77

Table 12: Scores obtained by Bard for Various Privacy Policies: Assessing Robustness through Paraphrased Questions

Proceedings on Privacy Enhancing Technologies YYYY(X)

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	8.2	9.1	10	10	9.55	9.37
$UC_e$	7.3	3.7	6.4	8.2	8.2	6.76
$DM_e$	10	1	10	8.2	6.4	7.12
$SE_e$	1	6.4	8.2	6.4	7.3	5.86
$PD_e$	1	1	3.7	5.95	3.7	3.07
$RC_e$	8.2	6.85	10	10	4.6	7.93
$AEE_e$	3.7	1	3.7	6.4	3.7	3.7
$CA_e$	3.7	8.65	3.7	4.15	7.3	5.5
$T_m$	8.2	6.4	10	9.1	10	8.74
$UC_m$	7.3	8.65	3.7	8.2	7.3	7.03
$DM_m$	10	5.05	1	6.4	1	4.69
$SE_m$	3.7	5.05	10	6.4	3.7	5.77
$PD_m$	1	1	3.7	6.4	3.7	3.16
$RC_m$	7.3	7.75	10	9.55	7.3	8.38
$AEE_m$	3.7	1	1	5.5	7.3	3.7
$CA_m$	10	1	10	6.85	3.7	6.31
$T_h$	1	5.05	6.4	6.4	1	3.97
$UC_h$	10	6.4	8.2	6.85	7.3	7.75
$DM_h$	1	1	10	7.75	1	4.15
$SE_h$	1	1	1	6.4	3.7	2.62
$PD_h$	1	10	3.7	5.5	7.3	5.5
$RC_h$	7.3	1	3.7	6.4	2.35	4.15
$AEE_h$	6.4	6.4	1	6.4	3.7	4.78
$CA_h$	6.4	1	1	6.4	3.7	3.7

**Table 13: Scores obtained by Bard for Various Privacy Policies: Assessing the Quality of Responses to Privacy Policy Questions**

Questions	Spotify	Uber	Twitter	Facebook	Airbnb	Average
$T_e$	8.2	8.2	8.2	9.55	10	8.83
$UC_e$	9.1	6.4	7.75	8.2	8.2	7.93
$DM_e$	8.2	4.6	7.75	8.2	3.7	6.49
$SE_e$	8.2	8.2	8.2	6.4	9.55	8.11
$PD_e$	3.7	1	1	6.4	4.6	3.34
$RC_e$	3.7	8.2	10	7.75	4.6	6.85
$AEE_e$	3.7	3.7	1	6.4	4.6	3.88
$CA_e$	6.4	8.2	4.6	7.3	8.2	6.94
$T_m$	8.2	8.2	8.2	6.85	8.2	7.93
$UC_m$	7.75	7.75	7.3	8.2	8.2	7.84
$DM_m$	10	8.2	1	5.95	3.7	5.77
$SE_m$	7.3	5.05	7.3	6.4	7.3	6.67
$PD_m$	5.5	3.7	1	6.4	4.6	4.24
$RC_m$	10	9.1	7.3	6.85	8.2	8.29
$AEE_m$	10	8.65	1	5.95	6.4	6.4
$CA_m$	9.1	1	8.2	4.6	4.15	5.41
$T_h$	8.2	3.7	9.55	5.5	3.7	6.13
$UC_h$	8.65	7.75	8.2	7.75	7.3	7.93
$DM_h$	10	1	7.3	7.3	3.7	5.86
$SE_h$	1	1	1	6.4	4.6	2.8
$PD_h$	3.7	8.2	1	6.4	8.2	5.5
$RC_h$	7.75	1	1	6.4	3.7	3.97
$AEE_h$	8.65	9.55	1	6.4	4.6	6.04
$CA_h$	7.75	8.2	1	5.95	4.15	5.41

**Table 14: Scores obtained by Bard for Various Privacy Policies: Assessing the Ability to Recall Learned Privacy Policy Knowledge**