# Privacy and Data Protection in ChatGPT and Other AI Chatbots:
## Strategies for Securing User Information

Glorin Sebastian, Georgia Institute of Technology, USA*

(iD) https://orcid.org/0000-0003-2543-9127

## ABSTRACT

The evolution of artificial intelligence (AI) and machine learning (ML) has led to the development of sophisticated large language models (LLMs) that are used extensively in applications such as chatbots. This research investigates the critical issues of data protection and privacy enhancement in the context of LLM-based chatbots, with a focus on OpenAI's ChatGPT. It explores the dual challenges of safeguarding sensitive user information while ensuring the efficiency of machine learning models. It assesses existing privacy-enhancing technologies (PETs) and proposes innovative methods, such as differential privacy, federated learning, and data minimization techniques. The study also includes a survey of Chatbot users to measure their concerns related to data privacy with the use of these LLM-based applications. This study is meant to serve as a comprehensive guide for developers, policymakers, and researchers, contributing to the discourse on data protection in artificial intelligence.

## KEYWORDS

Artificial Intelligence, ChatGPT, Cybersecurity, Data Protection, Large Language Model

## 1. INTRODUCTION

"ChatGPT, developed by OpenAI in November 2022, is an AI chatbot that utilizes the Generative Pre-trained Transformer (GPT) model. OpenAI is an AI research and development company known for its innovative approaches in natural language processing. The GPT model, based on the Transformer architecture introduced by Vaswani et al. (2017), is trained on extensive datasets to generate contextually relevant and accurate responses to text-based inputs. However, as these systems become more sophisticated and widely used, concerns regarding user privacy and data protection have emerged. Large Language Models (LLMs) like ChatGPT aim to understand and generate human language, but their reliance on extensive datasets, which may contain sensitive information, raises privacy concerns. There is a risk of inadvertently capturing and exposing sensitive user data, particularly in the context of chatbots and virtual assistants where personal or confidential information is often disclosed. These concerns have been addressed in various research papers discussing the usage of LLM-based chatbots, such as those by Hariri (2023), Sebastian (2023), and Cao et al. (2023).

*Corresponding Author

However, these research papers have not delved deeply into the topic of data privacy risks in LLM chatbots, this paper addresses this research gap by reviewing the data privacy risks associated with LLM chatbots. Further, to mitigate these privacy concerns, it is essential to develop effective strategies and technologies that can safeguard user data while maintaining the utility of LLMs. This paper aims to address this need by examining current privacy concerns, exploring existing privacy-enhancing technologies (PETs), and proposing novel techniques to ensure robust data protection in LLMs like ChatGPT. The techniques include differential privacy, federated learning, data minimization, and secure multi-party computation. Additionally, this research explores legal and ethical frameworks that can guide the responsible development of AI systems, considering both the tremendous potential of LLMs and the importance of user privacy. The paper serves as a comprehensive guide for developers, policymakers, and researchers in this rapidly evolving field, contributing to the ongoing dialogue about data protection in AI and promoting the development of innovative technologies that prioritize user privacy."

## 1.1 Brief Overview of ChatGPT

ChatGPT is an advanced AI model developed by OpenAI, which utilizes the Generative Pretrained Transformer (GPT) series of models. GPT belongs to the category of large language models (LLMs), which are characterized by their extensive training on diverse and comprehensive linguistic datasets and their ability to generate human-like text that is contextually relevant and coherent. ChatGPT, specifically, is designed to engage in conversation with users, with applications ranging from virtual assistants and customer service bots to AI tutors and more. The power of ChatGPT lies in its capacity to understand and generate meaningful responses to a wide array of prompts, demonstrating a deep grasp of syntax, semantics, and even nuanced aspects of conversation such as humor and emotion. The ChatGPT is a closed model without information about its training dataset and how it is currently being trained. Preventing data leakage (training-test contamination) is one of the most fundamental principles of Machine learning because such leakage makes evaluation results unreliable (Aiyappa, Rachith, et al.,2023).

Training an LLM like ChatGPT involves two main steps: pre-training and fine-tuning (Zheng, Ou, et al.,2023). During pre-training, the model is exposed to a large corpus of Internet text to learn grammar, facts about the world, reasoning abilities, and unfortunately, some of the biases present in the training data. In the fine-tuning process, ChatGPT is further trained on a narrower dataset, generated with the help of human reviewers following specific guidelines provided by OpenAI. Despite its impressive capabilities, ChatGPT, like all AI systems, raises some important privacy and data protection issues. Since the model learns from vast amounts of data, there is a risk of it inadvertently learning and generating sensitive or personally identifiable information. Also, user interactions with ChatGPT could potentially expose personal data, either through the questions users ask or the context in which the system is deployed. Hence, it is critical to explore techniques and strategies to enhance privacy and data protection in ChatGPT and similar LLMs.

## 1.2 Importance of Privacy and Data Protection in AI Systems

The significance of privacy and data protection in AI systems cannot be overstated and encompasses ethical, legal, and user trust considerations. Preserving privacy and data protection is crucial for the ethical, responsible, and compliant development and deployment of AI systems, contributing to user trust and the overall success and acceptance of these technologies.

i) **Ethical Considerations:** Ethical principles dictate that the personal and sensitive data of users should be respected and safeguarded. AI systems, especially large language models (LLMs) that process vast amounts of data, may unintentionally learn and generate sensitive information. Any compromise of personal data can lead to detrimental consequences such as identity theft,
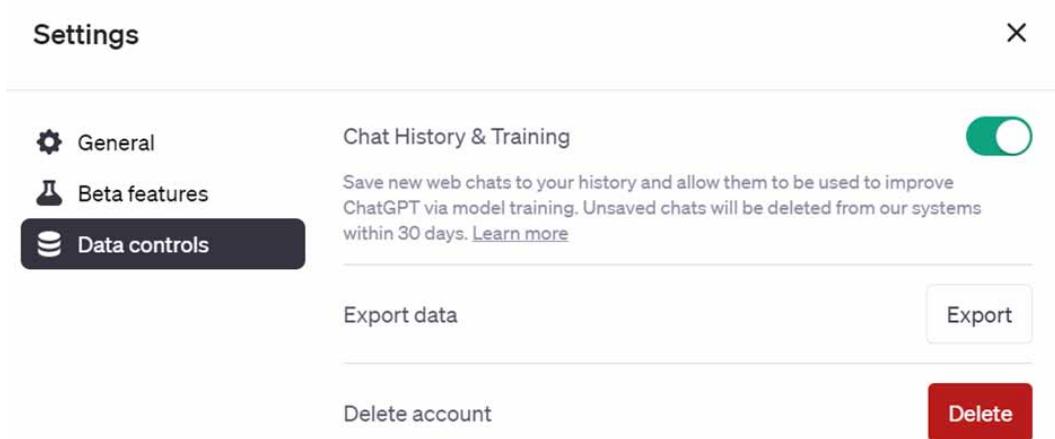
financial loss, or damage to one's reputation. Therefore, ensuring privacy and data protection in AI systems is an ethical imperative (Khoury, Raphaël, et al., 2023).

ii) **Legal Compliance:** Various jurisdictions have established data protection regulations, including the General Data Protection Regulation (GDPR) in the European Union and the California Consumer Privacy Act (CCPA) in the United States (Sherbini, Danya, 2023). These laws impose strict requirements on the collection, processing, and storage of personal data. GDPR guarantees rights such as access to data, rectification, and erasure, while CCPA empowers consumers with rights like knowledge, deletion, and opting out of personal information sale. If ChatGPT handles personal data from EU or California users, it must adhere to these provisions. Non-compliance can result in significant fines and legal consequences.

iii) **User Trust:** Trust plays a vital role in the adoption and continued use of AI systems. Research indicates that trust has a direct impact on both intention to use and actual utilization. Companies and policymakers should prioritize building trust and transparency in the development and deployment of chatbots (Choudhury, A. et al., 2023). Users need assurance that their personal data will be handled responsibly and securely. If an AI system is perceived as insecure or prone to privacy breaches, users may hesitate to utilize it, potentially limiting its overall utility and market acceptance.

iv) **Maintaining Data Integrity:** Adequate privacy protection and data security measures also contribute to maintaining data integrity. Without proper safeguards, data could be tampered with or manipulated, leading to unreliable outputs from the AI system. Real-world data inherently contain noise (Zongsheng Y., 2020). Considering this, data integrity may not guarantee fair evaluation but can further widen the gap between evaluation and real-world applicability (Lei Zhou et al., 2018).

v) **Mitigating Risks of Adversarial Attacks:** Adversarial attacks involve deliberate efforts to exploit vulnerabilities in AI models and manipulate their behavior by introducing carefully crafted input data. Such attacks can result in misleading outcomes, compromised decision-making, or even malicious actions. Mitigating the risks of adversarial attacks entails securing data through robust security measures like encryption and access controls. Ensuring high-quality data is vital for countering adversarial attacks. Adversarial training involves augmenting the training process with adversarial examples to enhance the resilience of AI models against attacks. Regularization techniques such as L1 or L2 regularization can be employed to reduce the susceptibility of AI models to adversarial attacks. These measures make it challenging for attackers to reverse-engineer the original model's behavior by introducing adversarial samples through defensive distillation. By securing data and ensuring privacy, the risks associated with such attacks can be mitigated.

## 1.3 Assessing ChatGPT's Data Collection, Storage and Handling of User Data

As per OpenAI's data usage policy, ChatGPT is trained on a large corpus of publicly available text from the internet during its pre-training phase, but it doesn't know specifics about which documents were part of its training set and cannot access any proprietary, classified, or confidential information. In the fine-tuning phase, the model is trained on a dataset generated with the help of human reviewers following specific guidelines provided by OpenAI. Regarding user interactions, OpenAI retains user interaction data for 30 days [1]. The data provided by users while interacting with models like ChatGPT is used to improve the models and is not used to personalize user experiences. More detailed and current information should be sourced directly from OpenAI's latest data usage policy OpenAI. (2021). Data collection and storage practices for AI models like ChatGPT and other large language models (LLMs) are critical to ensuring privacy and data protection. Here's an overview of these practices:

**Figure 1. Screenshot of the Data Controls that allow users to save web chats in history**



i) **Data Collection for Training:** ChatGPT and other LLMs are trained on a diverse range of internet text. The specifics about which documents were in the training set are not known to the model or to OpenAI, and the model cannot access any document or database during its operation.
ii) **Data Collection during User Interactions:** User interactions with these models may be stored and used to improve the system. The data collected include Log Data, Usage Data, Data is retained for 30 days (OpenAI., 2021).
iii) **Data Storage and Security:** OpenAI employs industry-standard security practices to protect the data, including encryption at rest and in transit. Our data will be shared with third parties exclusively upon obtaining your consent or in specific situations, such as legal obligations. OpenAI will guarantee that third parties involved in data processing adhere to comparable data handling practices and privacy standards. (OpenAI., 2021).

OpenAI's security page describes that it is compliant with legislation including GDPR and CCPA regulations. The option to create a custom Data Processing Agreement tailored to each organization or specific use case is available. The OpenAI API has undergone an assessment by an external security auditor, demonstrating SOC 2 Type 2 compliance. Additionally, it undergoes annual third-party penetration testing to proactively identify and address any security vulnerabilities that could be exploited by malicious individuals. (openai.com/security)

## 2. LITERATURE REVIEW

Literature Review included searching for research papers on the topic of privacy risks of ChatGPT. Refer to Table-1 below for details. While most researchers have mentioned the privacy risks of LLM-based Chatbots like ChatGPT, there has been no study that surveys ChatGPT users on their concerns about data privacy.

## 3. PRIVACY RISKS AND DATA LEAKAGE CONCERNS

LLM-based Chatbots like ChatGPT generate text based on extensive training data could lead to privacy concerns, particularly if it generates responses that appear to reveal sensitive information,

**Table 1. Literature review on privacy risks of LLM-based chatbots**

| S.No | Study | Summary |
|---|---|---|
| 1. | "Attention is all you need, Advances in neural information processing systems." (Vaswani, Ashish, et al., 2017) | proposes neural network architecture called the "Transformer" that uses self-attention mechanisms, without recurrent or convolutional layers. Transformer model achieves state-of-the-art results on several natural language processing tasks, outperforming traditional models that use sequential processing. |
| 2. | "Beyond the Safeguards: Exploring the Security Risks of ChatGPT." (Derner, Erik, and Kristina Batistič, 2023) | Investigates potential security risks associated with the language model ChatGPT. The paper goes explores the vulnerabilities and challenges that arise when deploying ChatGPT in real-world applications and discusses various attack vectors, information extraction, generation of harmful content, and adversarial manipulation of the model. |
| 3. | Unlocking the Potential of ChatGPT: A Comprehensive Exploration of its Applications, Advantages, Limitations, and Future Directions in Natural Language (Hariri, Walid., 2023) | Discusses practical applications for ChatGPT, including customer support, education, and creative writing. It highlights the advantages of ChatGPT, such as its ability to generate coherent and contextually relevant responses and its potential for personalization. The paper also states the limitations of ChatGPT, eg: occasional incorrect or nonsensical responses and difficulties in handling ambiguous queries. |
| 4. | Multi-step Jailbreaking Privacy Attacks on ChatGPT. arXiv preprint arXiv:2304.05197. (Li, Haoran, et al., 2023) | multi-step jailbreaking attacks that aim to extract sensitive user information through iterative interactions with the model. The authors demonstrate the ability to extract personal data such as names, and even credit card numbers by manipulating the model's responses through carefully crafted conversational inputs. |
| 5. | "Do ChatGPT and Other AI Chatbots Pose a Cybersecurity Risk?: An Exploratory Study" (Sebastian, G., 2023) | presents an exploratory study that examines potential vulnerabilities and threats that may arise from the use of AI chatbots. The author explores different attack vectors, including information extraction, social engineering, and malicious content generation. |
| 6. | "From ChatGPT to HackGPT: Meeting the Cybersecurity Threat of Generative AI" (Renaud, K., 2023) | talks on the cybersecurity risks including data privacy of generative AI such as ChatGPT and recommends replacing traditional, rule-based approach with "smarter" technology and training |

though it's important to note that the model itself doesn't have access to personal data. Additionally, concerns may arise from the potential misuse of user interactions with ChatGPT if privacy protocols for these interactions are not adequately enforced or understood. Below listed are some of the common privacy and Data leakage issues with these AI-based Chatbots. If any of these issues occur, it could significantly impact user trust in AI systems. Users need to be confident that their data is secure, and that the system will respect their privacy. Violations of privacy and data security can lead to users losing faith in AI systems, damaging the reputation and utility of these tools (Choudhury, Avishek, et al., 2023).

i)  **Unintended Sharing of Sensitive Information:** This occurs when a user unknowingly shares personal or sensitive data with the AI system (Sweeney, L., 2002). For instance, a user might share their credit card information, believing that the AI is secure. While non-PII AI models like ChatGPT do not have the ability to recall or store this information, given it stores non-PII information temporarily to improve performance for 30 days, the data could potentially be intercepted during transmission if the communication channel is not secure.

ii) **Data Leakage through Model Outputs:** Even though LLM models like ChatGPT do not know specifics about the data they were trained on, they can sometimes generate outputs that seem to refer to specific data or reveal sensitive information. However, these outputs are generated based on patterns learned during training and do not reflect access to any specific data sources or confidential databases. The AI could "hallucinate" specific, sensitive-looking details in responses. The model is not leaking real-world sensitive data that it learned during training—it's making things up based on the patterns it learned (Alkaissi, Hussam, et al., 2023)

iii) **Adversarial Attacks:** These attacks involve malicious actors attempting to manipulate or trick the AI into behaving in a certain way, usually for harmful purposes. For instance, an adversarial attack could involve inputting carefully crafted data designed to deceive the AI into generating inappropriate or harmful content. Some of the Cybersecurity attacks that LLM Chatbots could be subjected to include Evasion Attacks, Trojans Attacks, and Fake Review Attacks (Li, Jiazhao, et al., 2023) (Shi, Jiawen, et al., 2023).

iv) **Model Extraction:** This involves an attacker using the outputs of a machine learning model to create a copy of that model without access to the original training data. If successful, the attacker could use the extracted model for malicious purposes, potentially undermining the security and integrity of the original system (Li, Haoran, et al., 2023) (Tramèr, F., et al., 2023) (Shokri, R., et al., 2017) (Dwork, C. et al., 2006).

v) **Data Poisoning:** This is a type of attack where the attacker introduces harmful data into the model's training data with the aim of influencing its future predictions or behavior. It's a significant threat for systems that continually learn from their interactions with users (Li, Jiazhao, et al., 2023).

## 4. SURVEY RESULTS

The below Survey was conducted on MTurk participants who are familiar with using LLM-based Chatbots such as ChatGPT. The survey was conducted for a week between May 7 to May 12[th,] 2023, and received 177 responses. Out of the respondents, the majority were familiar with ChatGPT (86.2% rating their familiarity as 4 or 5) and expressed concerns about privacy and data protection (75.6% rated their concern as 4 or 5). 92.5% of the survey respondents mentioned they would be willing to sacrifice some performance or usability of the AI system for enhanced privacy and data protection. Almost all respondents (94.8%) agreed that AI systems should comply with data protection regulations. The majority of participants (78.8%) were aware of incidents involving the unintended sharing of sensitive information or data leakage. For further details of the survey please review Table-2 below and Figures 2-4.

## 5. MITIGATING PRIVACY RISKS AND STRENGTHENING DATA PROTECTION

Listed below are some techniques that can all contribute to privacy protection in the context of LLMs, but none of these is a silver bullet. The privacy risks associated with large language models like ChatGPT is a complex, multifaceted challenge that likely requires a combination of many techniques, as well as ongoing research and development.

i) **Data Anonymization and Aggregation:** Anonymization is a method of data protection where personally identifiable information fields within a data record are replaced by one or more artificial identifiers or pseudonyms. Aggregation, on the other hand, involves combining data in a way that the resulting dataset doesn't contain personally identifiable information (PII) (Mattas, Puranjay Savar., 2023).

**Table 2. Survey responses**

| Summary of survey results | Responses |
|---|---|
| **How familiar are you with ChatGPT and its applications? (Scale of 1 to 5, with 1 being "Not familiar at all" and 5 being "Extremely familiar")** | |
| - 2 | 06 (3.5%) |
| - 3 | 18 (10.4%) |
| - 4 | 93 (53.8%) |
| - 5 | 56 (32.4%) |
| **1. On a scale of 1 to 5, how concerned are you about privacy and data protection when using AI systems like ChatGPT? (1 "Not concerned at all" and 5 "Extremely concerned")** | |
| - 1 | 01 (0.6%) |
| - 2 | 08 (4.7%) |
| - 3 | 33 (19.2%) |
| - 4 | 93 (54.1%) |
| - 5 | 37 (21.5%) |
| **2. Do you believe that AI systems, such as ChatGPT, should comply with data protection regulations (e.g., GDPR, CCPA)? (Yes/No)** | |
| - Yes | 163 (94.8%) |
| - No | 09 (5.2%) |
| **3. Have you ever experienced or heard about unintended sharing of sensitive information or data leakage through AI-generated outputs? (Yes/No)** | |
| - Yes | 134 (78.8%) |
| - No | 036 (21.2%) |
| **4. How important do you think it is to implement data anonymization and aggregation techniques in AI systems like ChatGPT? (Scale of 1 to 5, with 1 being "Not important at all" and 5 being "Extremely important")** | |
| - 2 | 08 (4.6%) |
| - 3 | 23 (13.1%) |
| - 4 | 100 (57.1%) |
| - 5 | 44 (25.1%) |
| **5. Are you familiar with the concept of differential privacy? If yes, do you think it should be applied to AI systems like ChatGPT? (Yes/No/Not familiar with the concept)** | |
| - Yes | 145 (84.3%) |
| - No | 19 (11%) |
| - Not familiar with the concept | 08 (4.7%) |
| **6. How concerned are you about the potential security vulnerabilities in AI systems, such as adversarial attacks and data poisoning? (Scale of 1 to 5, with 1 being "Not concerned at all" and 5 being "Extremely concerned")** | |
| - 1 | 3 (1.8%) |
| - 2 | 5 (2.9%) |
| - 3 | 24 (14.1%) |
| - 4 | 102 (60%) |
| - 5 | 36 (21.1%) |

**Table 2. Continued**

| Summary of survey results | Responses |
|---|---|
| **7. Which countermeasure do you think is most effective in mitigating security vulnerabilities in AI systems like ChatGPT? (Select one: Adversarial training, Robustness testing, Data verification, Secure data sourcing, Rate-limiting, Blocking automated queries)** | |
| - Adversarial training (playing with someone who pretends to be a bad guy, so we can learn how to be safe) | 50 (28.7%) |
| - Robustness testing (checking if something, is strong enough to not break when we play with it.) | 29 (16.7%) |
| - Data verification (making sure the information we have, such as a list of our friends' names, is correct and true) | 47 (27%) |
| - Secure data sourcing (involves getting our information, like a secret message, from a safe and trusted place) | 32 (18.4%) |
| - Rate-limiting (eg having a rule that says you can only do something a certain number of times, such as having only three cookies a day) | 10 (5.7%) |
| - Blocking automated queries (This stops robots from asking too many questions or getting our secrets without permission) | 6 (3.4%) |
| **8. Would you be willing to sacrifice some performance or usability of an AI system like ChatGPT for enhanced privacy and data protection? (Yes/No)** | |
| - Yes | 161 (92.5%) |
| - No | 13 (7.5%) |
| **9. What additional recommendations or areas of research do you think should be explored to improve privacy and data protection in ChatGPT and similar AI systems?** | |
| - Explore ways to improve transparency and provide users with better visibility into how their data is being used by ChatGPT. | 33 (19.7%) |
| - Improving security protocols to better protect user data. This could include encrypting data while in transit and at rest, as well as developing systems to detect and report any potential malicious activity. | 57 (34.1%) |
| - Research into methods for preserving the privacy of sensitive data used as training data for AI systems, such as ChatGPT | 44 (26.3%) |
| - Develop an AI system that can detect and flag any potential data leaks or privacy violations. | 21 (12.6%) |
| - Need to incorporate stewardship requirements | 12 (7.2%) |

**Figure 2. Survey responses on if LLM based Chatbots should comply with data protection regulations**



3. Do you believe that AI systems, such as ChatGPT, should comply with data protection regulations (e.g., GDPR, CCPA)?
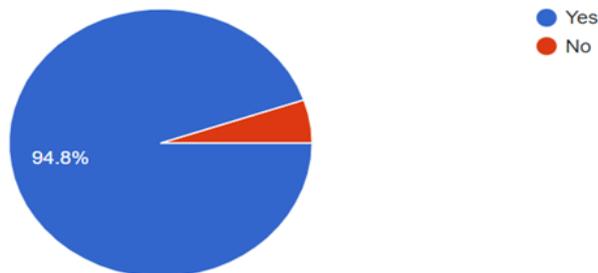
172 responses

94.8%

- Yes
- No

**Figure 3. Survey response to if respondents have experienced or heard of unintended sharing of sensitive information through AI generated outputs**

4. Have you ever experienced or heard about unintended sharing of sensitive information or data leakage through AI-generated outputs?

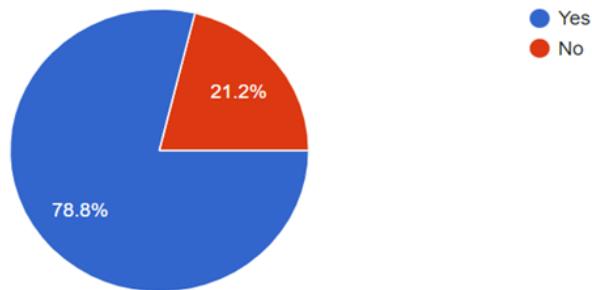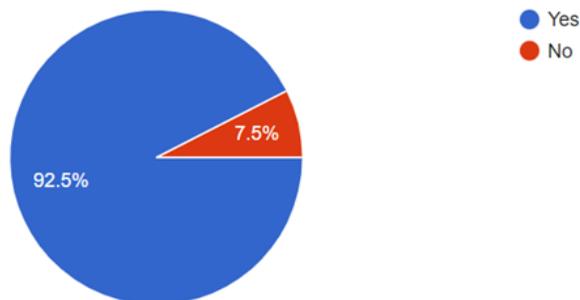170 responses



- Yes
- No

21.2%

78.8%

**Figure 4. Survey response to if the end users would be willing to sacrifice some performance or usability of the AI system for enhanced privacy and data protection**

9. Would you be willing to sacrifice some performance or usability of an AI system like ChatGPT for enhanced privacy and data protection?

174 responses



- Yes
- No

7.5%

92.5%

ii) **Differential Privacy Techniques:** Differential privacy provides a mathematical definition of privacy. It's a method for sharing information about a dataset by describing the patterns of groups within the dataset while withholding information about individuals. In the context of LLMs like GPT, differential privacy can help in reducing the chances of the model memorizing sensitive information during training Dwork, C. (2008).

iii) **Secure Multi-Party Computation (SMPC):** SMPC is a subfield of cryptography with the goal of creating methods for parties to jointly compute a function over their inputs, keeping those inputs private. Though it's more common in contexts where multiple datasets owned by different entities need to be collectively analyzed, without sharing the raw data (Srivastava, Mashrin., 2023).

iv) **Privacy-aware Machine Learning Algorithms:** These are algorithms that are designed with privacy as a priority. For instance, federated learning is a machine learning approach that trains

an algorithm across multiple devices holding local data samples, without exchanging them (El-Ansari, Anas., 2023).

v) **Adversarial Training and Robustness Testing:** Adversarial training is a technique to improve the model's robustness by including adversarial examples (inputs meant to confuse the model) in the training data. Robustness testing, on the other hand, involves testing a model to ensure it can handle unusual or unexpected inputs without failing.

vi) **Data Verification and Secure Data Sourcing:** Data verification ensures the quality and authenticity of the training data. Secure data sourcing refers to obtaining data from trusted and reliable sources. It is also crucial to ensure that the collected data complies with relevant legal and ethical guidelines.

vii) **Rate-Limiting and Blocking Automated Queries:** Rate-limiting involves limiting the number of requests a user or a service can make to the system within a certain timeframe. This can help protect the system against misuse. Blocking automated queries is another measure that can be implemented to protect the system from automated attacks or misuse.

viii) **Anonymization and encryption techniques:** Anonymization and encryption techniques play pivotal roles in data protection and privacy, especially in AI applications like ChatGPT and other large language models (LLMs). This ensures that data used in training stages is appropriately safeguarded against unauthorized access and misuse. Anonymization refers to the process of removing personally identifiable information (PII) from data sets, making it impossible (or at least very difficult) to link the data back to the individual it originated from. In the context of AI, anonymization is often used to protect user data during the model training process. Techniques include data masking, pseudonymization, generalization, and differential privacy Dwork, C. (2008). See Table-3 below for the impact of Privacy-preserving technologies on ChatGPT's performance.

i) **Data Masking and Pseudonymization:** These involve replacing identifiers in the data with artificial identifiers or pseudonyms (Yao, A.C., 1986).

ii) **Generalization:** This reduces the granularity of data. For instance, replacing exact ages with age ranges.

iii) **Differential Privacy:** This adds statistical noise to the data to ensure that the output of analysis doesn't reveal information about individual data points (Rogaway., 2011).

iv) **Encryption Techniques:** Encryption is the process of converting data into a code to prevent unauthorized access. In the context of AI and data storage, encryption plays a key role in securing data both at rest and in transit (Carlini, N., 2023). Encryption at Rest includes the encryption of data that is stored in databases, hard drives, or other storage mediums, while Encryption in Transit refers to the encryption of data as it is transferred from one location to another, such as across a network or from a client to a server.

## 6. PROPOSED EU AI ACT AND IMPACT ON PRIVACY IN LLM-BASED CHATBOTS

The European Union (EU) is in the process of considering a new legal framework known as the Artificial Intelligence (AI) Act, which aims to enhance regulations surrounding artificial intelligence development and usage. The proposed legislation focuses on aspects such as data quality, transparency, human oversight, and accountability to address ethical concerns and implementation challenges across various sectors. The AI Act aims to reinforce Europe's position as a leading global hub for AI while ensuring adherence to European values and rules. It introduces a classification system with four risk tiers—unacceptable, high, limited, and minimal—to assess the potential risks posed by AI technologies to individuals' rights and safety. AI systems with limited and minimal risk, like spam

Table 3. Impact of privacy-preserving technology on ChatGPT's performance

| Proposed Technique | Evaluation | Impact on ChatGPT's Performance and Usability |
|---|---|---|
| **Data Anonymization and Aggregation** | Anonymization and aggregation can prevent the direct identification of individuals in training data, these techniques are not entirely foolproof and privacy breaches can occur when the model generates outputs that mimic sensitive patterns in the data (Balebako, R., 2014) (Schechter, S., 2007). | Anonymizing and aggregating data can have minimal impact on the performance and usability of language models if executed effectively, but excessive anonymization could lead to a loss of context and detail that could affect the quality of model outputs (Cormode, G., 2019). |
| **Differential Privacy Techniques** | offers a theoretically-grounded approach to prevent a model from learning much about any single example in the training data by adding a controlled amount of noise (Cormode, G., 2019). | Differential privacy can cause a trade-off between privacy and utility. It ensures privacy by adding noise to the data, which may degrade the performance of the model, leading to a decline in usability and accuracy (Dwork, C., et al., 2006). |
| **Secure Multi-party Computation (SMPC):** | allows multiple parties to compute a function over their inputs while keeping inputs private (Dwork, C., et al., 2006). less applicable to the training of LLM models like GPT-4, which typically involves a single dataset owned by a single entity. | Causes increased computational and communication costs, which could degrade performance. The impact on usability would likely depend on the specifics of the implementation (Papernot, N., et al. 2016). |
| **Privacy-Aware Machine Learning Algorithms** | designed to respect privacy in the learning process. However, they are still under development and their effectiveness in large language models is not yet fully demonstrated (Rogaway., 2011). | still in the development phase and may impose a computational overhead or lower model accuracy in their current state (Goodfellow, I., et al., 2015). |
| **Adversarial Training, Robustness Testing** | These techniques make models more robust and safe (Carlini, N., 2023). However, they are not specifically focused on privacy and do not offer a comprehensive solution on their own. | increase the robustness of a model against adversarial attacks and help to improve its generalization. However, they can also increase computational cost and complexity, potentially affecting performance and usability (Papernot, N., A ., 2016). |
| **Data Verification, Secure Data Sourcing** | Ensuring that only non-sensitive and legitimate data is used in training can help limit privacy risks at the data collection stage, but these techniques do not directly address risks arising during the model learning process itself (Balebako, R., 2014). | Verifying and securely sourcing data can ensure the integrity and legality of the data, but it may also limit the amount of data available for training, possibly affecting the performance and usability of the model (Goodfellow, I., et al., 2015). |
| **Rate-Limiting, Blocking Automated Queries** | used to prevent misuse of the trained model. While effective as part of a comprehensive approach to privacy, they are reactive measures and not directly relevant to the model learning process (Schechter, S., 2007). | prevent misuse of the model, but might also hinder usability for legitimate high-volume users. The impact on performance would likely be minimal unless the system becomes overloaded with automated queries (Biggio, B., & Roli, F.,2018) (Shokri, R., 2017) (Akhawe, D., 2013) (Xu, Minrui, et al., 2023) |

filters or video games, have fewer requirements but must fulfill transparency obligations. Systems deemed to pose an unacceptable risk, such as government social scoring or real-time biometric identification in public spaces, are prohibited with few exceptions. High-risk AI systems, including autonomous vehicles and medical devices, are permitted but subject to rigorous testing, data quality documentation, and human oversight. The proposed legislation also covers regulations for general-purpose AI, including large language model generative AI systems like ChatGPT. Non-compliance penalties can be substantial, with fines for companies reaching up to €30 million or 6% of global income. False or misleading documentation submissions to regulators are also subject to penalties. The EU aims to establish new global norms with these groundbreaking rules to foster trust in AI, prioritizing the safety and fundamental rights of EU citizens.

## 7. CONCLUSION

The investigation into data protection and privacy enhancement in large language models (LLMs), with a particular focus on OpenAI's ChatGPT, has provided a comprehensive understanding of Privacy issues and potential solutions. The study demonstrated that a significant proportion of users are familiar with LLMs and hold substantial concerns about privacy and data protection when using these systems. While companies like OpenAI have multiple controls around data privacy and also

undergo regular Cybersecurity audits, there is a clear consensus for AI systems to adhere to data protection regulations. Participants asserted the importance of implementing data anonymization, aggregation techniques, and differential privacy in AI systems. They also expressed apprehension over potential security vulnerabilities such as adversarial attacks and data poisoning, while showing a willingness to compromise some performance for improved privacy and data protection. The preferred countermeasures to security vulnerabilities included adversarial training and data verification. Respondents also identified key areas for further research and development, including enhancing transparency, improving security protocols, preserving the privacy of training data, developing AI systems capable of identifying potential data leaks or privacy violations, and incorporating stewardship requirements. (Helberger, Natali et al.,2023)

These findings underline the urgent need for concerted efforts in enhancing data protection and privacy in AI systems. The study reaffirms the importance of continuous research, regulation, and application of Privacy-Enhancing Technologies (PETs) in AI models. It emphasizes the imperative of prioritizing user privacy while striving for technological advancement. As LLMs become increasingly integrated into daily life, developers, policymakers, and researchers should commit to upholding the principles of data protection, maintaining the trust of users, and ensuring the ethical use of AI.

## 8. RECOMMENDATIONS FOR FURTHER RESEARCH AND IMPLICATIONS FOR THE FUTURE DEVELOPMENT OF CHATGPT AND SIMILAR AI SYSTEMS

This research discusses the use of pre-trained foundation models (PFMs) like generative pre-trained transformers (GPTs) in edge intelligence to provide AI services for th Metaverse, tackling the challenges posed by their resource-intensive nature for edge servers. It proposes a new framework for efficient resource management and a metric, the Age of Context (AoC), to evaluate model relevance, ultimately aiming to balance latency, energy consumption, and accuracy for mobile AI services (Xu, Minrui, et al., 2023). It would be interesting to learn the cybersecurity, data privacy, and ethical implications of increased adoption of AI into Metaverse (Zhou, P., 2023). Few other future scopes for research would be the impact of the use of LLM based applications in Blockchain (Zhang, R. et al., 2019) and also building matured intelligence in organizations (George, A. et al., 2018).

## REFERENCES

Aiyappa, R. (2023). Can we trust the evaluation on ChatGPT? arXiv preprint arXiv:2303.12767

Akhawe, D., Amann, B., Vallentin, M., & Sommer, R. (2013, November). *Here's my cert, so trust me, maybe?: understanding TLS errors on the web*. ACM.

Alkaissi, H., & McFarlane, S. I. (2023). Artificial hallucinations in ChatGPT: Implications in scientific writing. *Cureus*, *15*, 2. doi:10.7759/cureus.35179 PMID:36811129

Balebako, R., Marsh, A., Lin, J., Hong, J., & Cranor, L. (2014). *The Privacy and Security Behaviors of Smartphone App Developers*. USEC. doi:10.14722/usec.2014.23006

Biggio, B., & Roli, F. (2018). Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, *84*, 317–331. doi:10.1016/j.patcog.2018.07.023

Cao, Y. (2023). A comprehensive survey of ai-generated content (aigc): A history of generative ai from gan to chatgpt. arXiv preprint arXiv:2303.04226.

Carlini, N., Liu, C., Erlingsson, Ú., Kos, J., & Song, D. (2019). The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX Security 19)* (pp. 267-284). USENIX.

Choudhury, A., & Shamszare, H. (2023). Investigating the impact of user trust on adoption and use of ChatGPT: a survey analysis. JMIR Preprints 10.2196/preprints.47184

Cormode, G. (2019). Data Anonymization Revisited. *SIGMOD Record*, *41*(3), 6–13.

Derner, E., & Batistič, K. (2023). Beyond the Safeguards: Exploring the Security Risks of ChatGPT. arXiv preprint arXiv:2305.08005.

Dwork, C. (2008). Differential privacy: A survey of results. In *International Conference on Theory and Applications of Models of Computation* (pp. 1-19). Springer, Berlin, Heidelberg.

Dwork, C., McSherry, F., Nissim, K., & Smith, A. (2006). Calibrating Noise to Sensitivity in Private Data Analysis. *Theory of Cryptography Conference*. Springer. doi:10.1007/11681878_14

El-Ansari, A., & Beni-Hssane, A. (2023). Sentiment Analysis for Personalized Chatbots in E-Commerce Applications. *Wireless Personal Communications*, *129*(3), 1623–1644. doi:10.1007/s11277-023-10199-5

George, A., Schmitz, K., & Storey, V. (2018). The BI&A system: building matured business intelligence in organizations. Academy of Management Global Proceedings.

Goodfellow, I. (2015). *Explaining and Harnessing Adversarial Examples*. ICLR.

Hariri, W. (2023). *Unlocking the Potential of ChatGPT: A Comprehensive Exploration of its Applications, Advantages, Limitations, and Future Directions in Natural Language Processing*. arXiv preprint arXiv:2304.02017.

Helberger, N., & Diakopoulos, N. (2023). ChatGPT and the AI Act. *Internet Policy Review*, *12*(1), 1. doi:10.14763/2023.1.1682

Khoury, R. (2023). *How Secure is Code Generated by ChatGPT*? arXiv-2304.

Li, H. (2023). Multi-step Jailbreaking Privacy Attacks on ChatGPT. arXiv preprint arXiv:2304.05197.

Li, J. (2023). *ChatGPT as an Attack Tool: Stealthy Textual Backdoor Attack via Blackbox Generative Model Trigger*. arXiv preprint arXiv:2304.14475.

Mattas, P. S. (2023). ChatGPT: A Study of AI Language Processing and its Implications. *Journal homepage*. www. ijrpr. com

Open AI. (2021). ChatGPT: Your AI Friend. *Open AI Blog*. https://www.openai.com/chat-gpt/

Open, AI. (2021). Fine-Tuning Large Language Models: Dataset and Safety. *OpenAI Blog*. https://www.openai.com/blog/fine-tuning-large-language-models-dataset-and-safety/

Open, AI. (2021). OpenAI Data Usage Policy. *Open AI Blog*. https://platform.openai.com/docs/data-usage-policy

Papernot, N., Abadi, M., Erlingsson, Ú., Goodfellow, I., & Talwar, K. (2016). *Semi-supervised Knowledge Transfer for Deep Learning from Private Training Data*. ICLR.

Renaud, K., Warkentin, M., & Westerman, G. (2023). From ChatGPT to HackGPT: Meeting the Cybersecurity Threat of Generative AI. *MIT Sloan Management Review*, *64*(3), 1–4.

Rogaway, P. (2011). *The moral character of cryptographic work*. Cryptology ePrint Archive, Report 2015/1162. https://eprint.iacr.org/2015/1162

Schechter, S., Dhamija, R., Ozment, A., & Fischer, I. (2007). The Emperor's New Security Indicators. IEEE Symposium on Security and Privacy (SP). IEEE. doi:10.1109/SP.2007.35

Sebastian, G. (2020). Evolution of the role of risk and controls team in an ERP Implementation. *IJMPERD,* 2249-6890.

Sebastian, G. (2023). Do ChatGPT and Other AI Chatbots Pose a Cybersecurity Risk?: An Exploratory Study. [IJSPPC]. *International Journal of Security and Privacy in Pervasive Computing*, *15*(1), 1–11. doi:10.4018/IJSPPC.320225

Sherbini, D. (2023). Has technological innovation lost the plot? An interview with AI ethicist Dr. Shannon Vallor. *Chicago Policy Review (Online)*.

Shi, J. (2023). *BadGPT: Exploring Security Vulnerabilities of ChatGPT via Backdoor Attacks to InstructGPT*. arXiv preprint arXiv:2304.12298.

Shokri, R., Stronati, M., Song, C., & Shmatikov, V. (2017). Membership Inference Attacks against Machine Learning Models. *2017 IEEE Symposium on Security and Privacy (SP),* (pp. 3-18). IEEE. doi:10.1109/SP.2017.41

Srivastava, M. (2023). *Towards Trustworthy Machine Learning in Healthcare: Addressing Challenges in Explainability, Fairness, and Privacy through Interdisciplinary Collaboration*. OSF.

Sweeney, L. (2002). k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-based Systems*, *10*(05), 557–570. doi:10.1142/S0218488502001648

Tramèr, F., Zhang, F., Juels, A., Reiter, M. K., & Ristenpart, T. (2016, October). Stealing machine learning models via prediction APIs. In *25th USENIX Security Symposium (USENIX Security 16)* (pp. 601-618). USENIX.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.

Xu, M. (2023). Sparks of GPTs in Edge Intelligence for Metaverse: Caching and Inference for Mobile AIGC Services. arXiv preprint arXiv:2304.08782 (2023).

Yao, A. C.-C. (1986). How to generate and exchange secrets. *27th annual symposium on foundations of computer science (Sfcs 1986)*. IEEE. doi:10.1109/SFCS.1986.25

Zhang, R. (2019). Benefits of blockchain initiatives for value-based care: proposed framework. *Journal of Medical Internet Research*.

Zheng, O. (2023). *ChatGPT is on the horizon: Could a large language model be all we need for Intelligent Transportation?* arXiv preprint arXiv:2303.05382 (2023).

Zhou, P. (2023). *Unleashing chatgpt on the metaverse: Savior or destroyer?* arXiv preprint arXiv:2303.13856.