

Security and Privacy of ChatGPT: Proposed Enhancements

Hamad A. AlJarboua Alshahrani^{1†} and Waleed A. Alrodhan^{1* ††},

Mohammad bin Saud Islamic University, College of Computer and Information Sciences, Department of Computer Science, Riyadh, KSA

Abstract

ChatGPT Generative AI, driven by machine learning (ML) and deep learning (DL) is playing an increasingly vital role in today's digital world. Generative AI offers significant promise across various industries, such as military, finance, education, energy, and medical services. Nevertheless, since these technologies keep evolving at a rapid pace and have a growing effect, it is crucial that we recognize the various risks associated with them, especially amongst the fields of politics, society, technology, legal and ethics. The results have demonstrated the consequential effects on security, privacy, and ethics in Generative AI ChatGPT which has led to an evident emphasizing on the need for the development or establishment of regulatory legislation. Therefore, it is important to regulate the ChatGPT Generative AI scene, as there are currently no specific regulations governing its secure use. Additionally, there is insufficient awareness of handling sensitive data and its classifications. Furthermore, the work or development of cybersecurity controls is required to align with data management and protection, as well as intellectual property rights, for generative AI systems and applications in order to prevent the leakage of sensitive data, enhance the confidentiality, privacy and integrity of information, and mitigate the resulting effects of generative AI from such potential threats, as the dissemination of misleading information, disclosure of private data, deepfake content, unauthorized collection of sensitive information, and the production of unethical or illegal content during interactions with ChatGPT. This research aims to address this crucial need by analyzing, reviewing, and exploring relevant laws and regulations of data protection and potential regulatory interventions such as the Personal Data Protection Law (PDPL), National Cybersecurity Authority (Data Cybersecurity Controls), NCA-DCC, and National Data Management Office (NDMO) "SADAIA" that would help mitigating the relevant risk. In addition to proposing the development of international legislation and regulations, policymakers should begin establishing comprehensive regulations for the adoption of ChatGPT and enhance the efforts of relevant organizations in propagating awareness and education at the level of organizations and individuals for raising awareness among all segments of society and ensuring the application of best practices, data safety, preserving privacy, reliability, and responsibility, which lead to reducing social manipulation and political and ethical risks generally. It is important to implement responsible AI practices and transparency in data usage, bias mitigation techniques, legal matters, and monitoring of generated content for harmful or misleading information (hallucinate).

Keywords:

ChatGPT, NCA DCC, LLM, NDMO, OpenAI, Privacy, Ethical, AI, PDPL, Security, GDPR.

1. Introduction

In November 2022, OpenAI released ChatGPT, which swiftly attracted over 100 million users and became one of the fastest-growing AI tools. The release of ChatGPT prompted quite a buzz as a result of its highly advanced capabilities. Both, Machine Learning and Deep Learning are playing an increasingly pivotal role in the new smart digital age. Generative AI is not relatively a new concept; Chatbots, image generators and deepfake technologies have been in development and utilized for several years where recent improvements have significantly strengthened generative AI by leveraging increased training data, improved artificial neural networks with larger datasets and parameters, and expanded computer capacity. As a result of the rapid development and expanding influence of these technologies, a variety of security, privacy, political, social, and ethical challenges arise. A significant problem to consider pertains to the potential utilization of ChatGPT models for the dissemination of disinformation or the generation of hallucinatory content. Due to the capability of these systems to generate text that closely simulates human conversation through their training data, there is a notable concern of them being susceptible to broad manipulation such as generating child sexual exploitation and abuse (CSEA) material, sharing artificial pornography without consent and disseminate content that is either inaccurate or misleading may lead improper utilization of resources to a range of societal issues, such as an escalation of political conflict or the dissemination of disruptive disinformation. This research analyzes these security, privacy, and ethical problems, with the potential challenges of ChatGPT in disseminating disinformation, private data disclosure, gathering

¹ * AKA Waleed A. Alshalan

sensitive data without consent, and producing unethical or legal content while interacting with ChatGPT. This research will highlight and discuss aspects of ChatGPT, along with its social, legal, political, technological, and ethical boundaries.

AI models such as ChatGPT have been abused instances to get involved in disturbing media campaigns and spreading disinformation. ChatGPT has become a risk that malicious actors may exploit it for unethical purposes. Consequently, this might result in significant consequences, such as the manipulation of public opinion, electoral fraud, or the tarnishing of the reputation of prominent individuals. For instance, the creation of offensive content: AI models like ChatGPT can be used to generate offensive, explicit, or violent content, including images, videos, or text. This content can be disseminated through various online channels, potentially causing harm to individuals or communities where ChatGPT can be used to create sophisticated deepfake media, which involves manipulating or fabricating images, videos, or audio that appear authentic but are actually synthetic. Moreover, malicious actors may use AI-generated content to manipulate social media platforms by spreading misinformation, and propaganda, or engaging in coordinated campaigns to deceive or influence public opinion. This can lead to the amplification of harmful narratives, or the creation of divisive environments or malicious actors may use AI-generated personas to engage in phishing attempts, spreading false information, or conducting social engineering attacks.

According to a recent report in 2023 published by OpenAI². Report explores the viewpoints of more than 30 experts in the fields of Artificial intelligence, influence operations, and political analysis, On the potential implications of generative AI models on influence processes and providing recommendations for remediation. This report outlines the threats that language models pose to the information environment if used to augment disinformation campaigns and introduces a framework for analyzing potential mitigations.

Introduced by Derner et al. Generative AI has advantages based on design and capabilities and has the potential to disseminate false information, such as

deepfakes, which can result in significant negative consequences for society at large. [1].

Several academic studies have highlighted the effectiveness of development regulation policies. For example, according to a recent study by Glorin Sebastian et al. to effectively navigate the complex landscape of AI chatbots and disinformation, regulatory frameworks need to be established or adapted to address the unique challenges posed by these technologies. [2]. The study found that adopting an AI compliance framework and conducting a privacy and ethics impact assessment contributes to protecting personal information and limiting privacy liability when using ChatGPT or similar AI technology. A review article by Roland Hunget et al. concluded that the core of any AI compliance framework should be the incorporation of privacy-by-design and ethics-by-design concepts into the framework. [3]. Therefore, it means that the entity will incorporate data protection and ethical aspects into its framework of technology, practices, and processes. These traits have the potential to allow an entity to adjust to evolving technology and regulations. This research will focus on proposed Gen AI framework and enhancements in ChatGPT systems (Generative AI) in terms of security and privacy. In addition to exploring challenges, problems, and enhancements to make this ChatGPT OpenAI secure, privacy-preserving, trustworthy, and ethical. This research will focus on proposed Gen AI framework and enhancements in ChatGPT systems (Generative AI) in terms of security and privacy. In addition to exploring security and privacy controls with regulations, laws, and mitigation techniques will improve Saudi enterprises' cybersecurity posture and provide them with a more holistic approach to threat detection, prevention, and response. The study aims to achieve the following objectives:

- Providing a comprehensive overview of the security and privacy problems of ChatGPT.
- Developing a comprehensive linkage betwixt the existing Saudi-mandated security and privacy regulations and laws, and ChatGPT's usage concerns which contribute to focusing on the main gaps—AI regulation, security, privacy, and ethical matters—and proposed

² <https://openai.com/research/forecasting-misuse>

unified compliance framework for generative AI uses.

- Identifying potential coverage gaps of ChatGPT with regard to conventional ethics.
- Proposing enhancements to the security and privacy of ChatGPT by implementing effective security measures.

Below, we state our research questions (RQ), which contribute to enhancements of security and privacy of ChatGPT (LLM) as well as focusing on the main gaps—AI regulation, security, privacy, and ethical matters.

RQ1. What are the Essential Cybersecurity Controls (ECC) and Data Cybersecurity Controls of the National Cybersecurity Authority (NCA), and other national regulatory agencies, that organizations and individuals should adhere to, in order to be protected against ChatGPT threats?

RQ2. What are the benefits and challenges of adhering to the regulatory framework and legal compliance with regard to using ChatGPT?

RQ3. What are the impact and challenges of AI ethics on the ChatGPT, and how can they be addressed?

1.1 Motivation and Research Problem

In the rapidly evolving world of AI technology, ChatGPT, and with an increase in its abuse for bad or wrong intentions, establishing a robust ChatGPT regulatory framework and conducting awareness and education at the level of organizations and individuals for raising awareness among all segments of society might contribute to mitigating and protecting against information security threats. Therefore, this research aims to contribute to exploring challenges, problems, and enhancements to make this ChatGPT OpenAI secure, privacy-preserving, trustworthy, and ethical by adopting the development of effective mitigation to plug gaps with ChatGPT use cases.

1.2 Introduction to ChatGPT's Technical Aspects and Operating Processes

ChatGPT is trained with human feedback (a technique called Reinforcement Learning with Human Feedback). ChatGPT has been developed employing a two-stage method that includes unsupervised pre-training and supervised fine-tuning.

Sakib Shahriar et al. explained that these models undergo training on expansive datasets, enabling them to discern statistical patterns and associations within sequences of words. Consequently, they acquire the ability to generate or interpret language with remarkable proficiency. These models organize both input and output as tokens, which are numerical representations of words. This numerical representation facilitates efficient processing by the model. The training process for the GPT model involves predicting the subsequent token given a sequence of input tokens; thus, it essentially learns the language's structure. Consequently, the model can produce text that is both grammatically consistent and semantically similar to the training data. Once the model has been trained, it can be directed to perform specific language tasks, like answering questions. [4].

1.2.1 ChatGPT Architecture

OpenAI developed ChatGPT, a language model that employs the GPT (Generative Pre-trained Transformer) architecture. The system employs deep learning techniques, specifically a variant of the Transformer model, to produce text answers that closely resemble the style and content of human-generated writing when provided with prompt or input. Here's an overview of its technicalities:

Transformer: The core of GPT is the Transformer, which is an architecture of an artificial neural network. Typically, it serves various purposes such as text creation and summarization. The Transformer model is designed to process and generate sequences of tokens, such as words or sub-words. In addition, its basis of GPT is the transformer, a type of neural network architecture. It contains two main elements which are the encoder and the decoder. The encoder processes input sequences and processes simultaneously, whereas the decoder creates the output sequence employing the encoder's learned representations. The transformer is renowned for its

aptitude for apprehending extended range dependencies in text data, which is indispensable for tasks such as language modeling and text generation. [5].

Tokenization: In this phase, the input text is segmented into separate parts known as tokens, which can range in length from a single character to complete words from English, subwords, or even individual words, according to the language and context. Mainly Tokenization is a fundamental step in natural language processing that involves segmenting text into discrete units of meaning, known as tokens [6]. ChatGPT uses tokenization to break down input text into smaller units called tokens. It employs Byte Pair Encoding (BPE) to split words into subword units based on frequency and it is an essential step in language modeling as it enables the model to handle a wide vocabulary and effectively process and generate text at a granular level.

Pre-Training: In this phase, the training process is supervised rather than relying on reinforcement. These models are trained to employ a large dataset to predict the next word sequentially. This massive collection of textual material can encompass several sources, including books, articles, and websites such as Wikipedia, and so on. In details, The pre-training process involves training a deep neural network architecture, such as the Transformer model as well as supervised rather than reinforcement-based. For this reason, it feeds ChatGPT with a broad range of internet text, which includes websites, Wikipedia, articles, books, and other publicly available sources. However, it's important to note that the specific documents and sources used in pre-training are not disclosed publicly. The pre-training process and fine-tuning utilize the Adam algorithm, a variant of stochastic gradient descent, to update the model weights more efficiently and stably [7].

Fine-Tuning: In this phase, the fine-tuning process aims to train the model with targeted datasets or tasks to enhance response precision and relevance. Eventually, After the pre-training process, the model is fine-tuned using a narrower dataset that is generated with the help of human reviewers to enhance the accuracy and relevancy of the responses which helps the model to better understand and, generate more coherent text, and avoid generating harmful or

inappropriate content. Within fine-tuning, the model is trained on a smaller dataset of labeled data tailored to the particular natural language processing task. This ensures that the model's performance is optimized for the specific task while preserving its capacity to generate relevant and meaningful responses to natural language input [8].

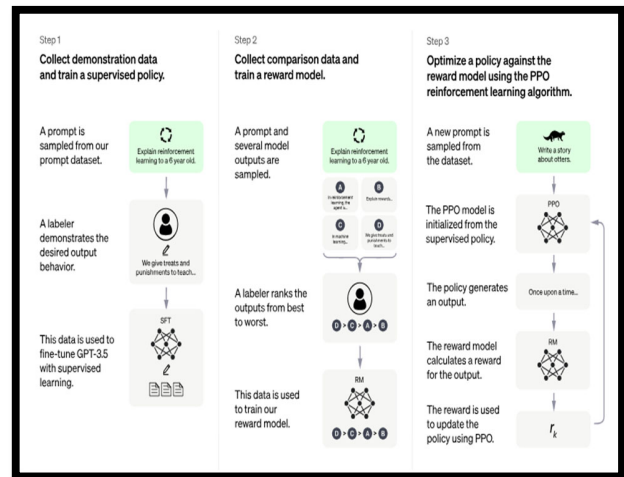


Figure 1.1 How ChatGPT is trained. Source: OpenAI Website.

These components work together to enable ChatGPT to engage in interactive and dynamic conversations where they facilitate the model's understanding of context, generation of coherent responses.

1.3 OpenAI's Privacy and Confidentiality Challenges.

Barrett et al. presented the potential impacts on privacy in the era of GenAI is another critical yet complex task [9]. The likelihood of an abuse of ChatGPT in developing code is associated with ransomware and malware attacks. These findings highlight the potential risks associated with AI models like ChatGPT, which have the potential to either generate malicious code or aid in its comprehension and development. While the code generated by the AI generally seemed to be pseudocode rather than code that can be executed, the ability to provide an attacker with a conceptual sense or broad grasp of how an attack works is worrisome.

This raises severe privacy concerns, as training data can be extracted verbatim from the models in some cases [10].

Even though the data is open to the public, that does not mean that it has been intended to be utilized for business purposes by third parties. The fact that this lack of explicit permission and the subsequent unauthorized use of data raises new privacy concerns. For example: OpenAI 3 has received a complaint accusing it of collecting, processing, and disclosing personal data (PII) without consent. Thereby, we need to examine and analyze the privacy policy of OpenAI and the current risks, which are ChatGPT does not provide sufficient methods to preserve personal data according to GDPR. For example, ChatGPT may share users’ data with third-party entities without explicit permission from users. OpenAI states in its Privacy Policy⁴ that individuals in the EEA, UK, and globally have statutory rights to their personal information, including access, deletion, correction, transfer, restriction, withdrawal of consent, objecting, and filing complaints with local data protection authorities. These rights can be exercised through OpenAI accounts or by contacting dsar@openai.com. However, ChatGPT’s output may contain inaccurate information about users, and correction requests should be submitted to dsar@openai.com. To exercise these rights, individuals can request that their personal information be removed from ChatGPT’s output by filling out the form. Figures 1.2 show screenshots of ChatGPT’s privacy policy.

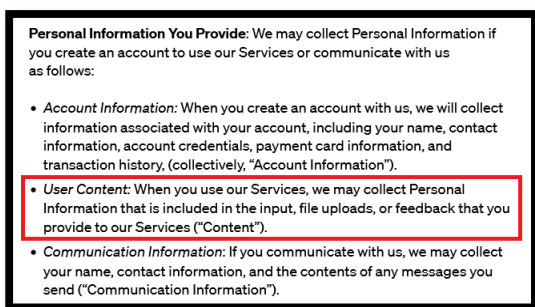


Figure 1.2 Screen view of the Open AI’s privacy policy. Source: OpenAI Website

ChatGPT interactions are not confidential, as they log conversations and personal data where they are shared, including any input, file uploads, or users’ feedback to feed up ChatGPT LLM as trained data. OpenAI’s privacy policy⁵ states that the company collects personal information for AI language models and to be reviewed by human trainers too. Furthermore, it is important to note that users do not have the ability to remove specific prompts; hence, caution should be exercised when revealing sensitive information. Likewise, the process of removing personal data via ChatGPT is not an easy one, which further complicates the exercise of the *right-to-be-forgotten*. Consequently, once sensitive information has been processed by the machine learning model, there is no mechanism to remove that information.

1.4 Problem and Statement

In this study, we will be identifying, analyzing, evaluating, and classifying the security, privacy, and ethical problems inherent to ChatGPT that organizations and individuals need to take into consideration when using ChatGPT Generative AI. The objective of enhancing the security and privacy of ChatGPT is to figure out and determine priority to potential risks, enabling the implementation of appropriate measures to mitigate or eliminate those problems (classified risks below):

Data Privacy: The model requires data to be trained on, and if sensitive or confidential information is used during training or in conversations, there is a risk of data exposure.

Malicious use: ChatGPT can be used to generate misleading (hallucinate) or harmful content, such as phishing attempts, social engineering, or spreading misinformation, which can have cybersecurity implications.

System Vulnerabilities: As with any software or system, there is a possibility of undiscovered vulnerabilities or exploits that could be leveraged by malicious actors.

Furthermore, OpenAI’s privacy policy states the potential risk of data disclosure. It states openly that ChatGPT collects user account information, all conversation-related information, and various privacy-related data, such as cookies, logs, and device

³ <https://www.reuters.com/legal/lawsuit-says-openai-violated-us-authors-copyrights-train-ai-chatbot-2023-06-29/>

⁴ <https://openai.com/policies/privacy-policy>

⁵ <https://openai.com/policies/privacy-policy>

information. Through online pages on the internet that are interactive worldwide. Such information may be exchanged with vendors, service providers, and affiliated businesses. During the data sharing process, unauthorized parties may gain access to private information associated with the model, including the leakage of training data, model architecture, parameters, and hyperparameters. Overall, the problem statement in ChatGPT is to identify, analyze, and evaluate potential risks and vulnerabilities, and to determine the impact of those risks and vulnerabilities so that appropriate measures can be taken to mitigate or eliminate them.



Figure 1.3 Open-Source Projects <https://github.com/lencx/ChatGPT> Confronted with Trojan horse attack

Table 1.1 Screen view of the Open AI’s privacy policy. Source: OpenAI Website.

Category of Personal Information	Sources of Personal Information	Use of Personal Information	Disclosure of Personal Information
Social Information	We may collect Social Information from you when you interact with our Social Media Pages.	We may use Social Information to perform analytics and to communicate with you.	We may disclose Social Information to our affiliates.
Communication Information	We collect Communication Information directly from you.	We use Communication Information for providing our Services and responding to you.	We disclose Communication Information to our affiliates and communication services providers.
Technical Information	We collect Technical Information from you.	We use Technical Information for analytics and in some cases, for moderation and prevention of fraud and malicious activity by users of our Services.	We disclose Technical Information to our affiliates and analytics provider(s).

Another form of the risk of privacy leakage associated with ChatGPT is the emerging malicious activities that have taken advantage of the popularity of ChatGPT to perpetrate identity theft attacks on users' privacy. The presence of a high-risk Trojan horse has been detected in the ChatGPT desktop application project, which is an unofficial and open-source software that has been posted to the GitHub platform, where users' account information, browser cookies, and other personal data will be revealed once the installed executable file is launched. In order to mitigate the risk of further users becoming victims of potential threats, the Open Source ChatGPT project has changed its download address.

Maanak Gupta et al. explored and highlighted the potential misuse of ChatGPT in creating code linked to ransomware and malware attacks [11]. These findings underscore the potential risks associated with AI models like ChatGPT, which could be exploited to generate malicious code or aid in the understanding and creation of such code. Although the code produced by the AI often resembled pseudocode more than actual executable code, the capacity to provide an attacker with a structural idea or general understanding of how an attack operates is a cause for concern.

1.5 Key Concept

1.5.1 National Cybersecurity Authority (NCA) Essential Cybersecurity Controls

National Cybersecurity Authority's (NCA's) Essential Cybersecurity Controls (NCA-ECC) were published in October 2024 and are required to be adopted by all public organizations as well as private ones that offer services, particularly IT services, to the public ones.

Essential Cybersecurity Controls consist of the following:

- 4 Cybersecurity Main Domains
- 28 Cybersecurity Subdomains
- 108 Cybersecurity Controls

The Essential Cybersecurity Controls are mandatory where all organizations, within the scope of these controls must implement whatever necessary to ensure

continuous compliance with the controls.⁶ The scope of ECC is cybersecurity governance, cybersecurity defense, cybersecurity resilience and, third part – cloud computing cybersecurity which are relevant controls to be utilized in the linkage approach. Iman Almomani et al. introduced the Essential Cybersecurity Controls objective (ECC) ensuring organizations' minimum cybersecurity requirements in Saudi Arabia in both the public and private sectors [12]. ECC's main objective is to safeguard and enforce the confidentiality, integrity, and availability of the organization's information and technology assets. The government also uses ECC as a method for compliance assessment. ECC was built based on national and international cybersecurity frameworks and standards, along with KSA national laws.

1.5.2 National Cybersecurity Authority (NCA) Data Cybersecurity Controls

The National Cybersecurity Authority Developed Data Cybersecurity Controls in 2022 to meet the minimal cybersecurity requirements for enterprises to secure their data across its entire data lifecycle.

The DCC consists of the following components:

- 3 Main Domains
- 11 Subdomains
- 19 Main Controls
- 47 Sub Controls

The controls aim to contribute to raising the cybersecurity maturity at a national level by setting the minimum cybersecurity requirements to enable organizations to protect their data during its entire data lifecycle. This document highlights the details of the Data Cybersecurity Controls (DCC), objectives, scope of work, compliance, and monitoring. These controls are an extension to the Essential Cybersecurity Controls (ECC) ⁷. The scope of DCC is all Data Cybersecurity Controls Domains and Subdomains which are all relevant controls, and to be utilized in the linkage approach.

1.5.3 Data Management and Personal Data Protection Standards

The National Data Management Office (NDMO) developed data management and personal data protection standards in 2021 to ensure the security and privacy of sensitive information. These standards apply to all organizations that collect and process personal data in the Kingdom of Saudi Arabia.

The Data Management and Personal Data Protection Standards consist of the following components:

- 15 Domains
- 77 Controls

The National Data Management and Personal Data Protection Standards cover 15 Data Management and Personal Data Protection domains. To support the development of the Data Management and Personal Data Protection standards, a set of international references, relevant internal policies and regulations, and guiding principles were defined. Government Entities must implement the standards, and compliance will be measured yearly to monitor progress and drive efforts towards successful implementation.

Overall, the Saudi Data Management and Personal Data Protection Standards are designed to ensure that organizations in Saudi Arabia comply with best practices for data management and personal data protection and to enhance the security and privacy of personal data. ⁸ . The scope of PDPL, Data Management and Personal Data Protection Standards and domains and laws which are relevant controls and to be utilized in the linkage approach, except Business Intelligence and Analytics domain, in NDMO. Note: The Data Security and Protection controls by NCA will provide complete requirements and address the corresponding controls for the Data Security and Protection Domain.

Emna et al. provided a brief overview of NDMO and PDPL in Saudi Arabia, the most notable provisions of the regulations are about obtaining consent from data subjects, maintaining data localization within Saudi Arabia, and limiting personal data collection to data that is necessary and relevant to the activities of data controllers. Also, they

⁶

<https://nca.gov.sa/en/legislation?item=191&slug=controls-list>

⁷

<https://nca.gov.sa/en/legislation?item=317&slug=controls-list>

⁸ <https://sdaia.gov.sa/ndmo/Files/PoliciesEn001.pdf>

mentioned that nonetheless, the regulation does not constitute a data protection law but rather an introduction to data privacy-type concepts. Additionally, the regulation borrows many definitions from the General Data Protection Regulation (GDPR) but diverges significantly from the GDPR-type approach and places a heavy reliance on the data subject's consent.

The authors mentioned that the law gives data owners mandatory rights over their personal data. [13].

It stipulates that the processing of personal data is subject to the data owner's consent. The PDPL also recognizes a more sensitive sub-category of personal data, including bio-identifying and genetic data, health data, and location data, that requires additional protection

1.6 Open AI

OpenAI is a research organization and technology business that specializes in artificial intelligence (AI). The organization was established in December 2015. OpenAI conducts cutting-edge research in several subdomains of artificial intelligence (AI) and aims to create and develop cutting-edge models and technology solutions. OpenAI has gained significant recognition for its impressive projects, most notably in the realm of the Generative AI Pre-trained Transformer (GPT) series of models, which encompasses the well-known GPT-3/GPT3.5 and recently the GPT-4 model. These models are designed for natural language processing tasks and have demonstrated impressive language understanding and generation capabilities. In the year 2018, OpenAI released a report with the purpose of explaining the concept of a generative pre-trained transformer (GPT) to a worldwide public. A GPT, an acronym for Generative Pre-trained Transformer, is a type of neural network that is a machine learning model designed to mimic the cognitive capabilities of the human brain. It undergoes training using extensive input data, typically in the form of big datasets, in order to generate outputs that correspond to responses or answers to queries posed by users. Also, OpenAI introduced. A generative artificial intelligence (AI) model uses algorithms to analyze natural language text provided by human users, subsequently generating images that correspond to the textual descriptions.

ChatGPT is now the most popular AI that uses large language models (LLMs). It was also one of the first open-source projects to gain a lot of attention, traction, and funding. But despite these milestones, tech giants like Microsoft and Google are starting to roll out similar technology.

2 Literature Review

In this research, we give an overall list of relatively recent studies related to the security and privacy of ChatGPT. We point out the main contributions, objectives, and existing gaps for each of them and how the current study can help and support them, extending the study to mitigating security and privacy concerns in ChatGPT.

ChatGPT security, privacy, and AI regulations have already taken place in cybersecurity, politics, regulators, and academia, with many research topics highlighting its concept. However, some of ChatGPT security and privacy issues still to be addressed along with an evident lack of AI governance regulations globally.

Table 1.2 shows a summary of the reviewed references, emphasizing their contributions and highlighting existing gaps where most of these studies have a common research gap in their need to extend those studies to cover more problems such as security, privacy, AI regulation, and governance, as well as data classification and awareness matters.

⁹ <https://openai.com>

Table1.2 Summary of the reviewed references

Reference Number	Contribution	Existing gap
14	<ul style="list-style-type: none"> - Provide a brief overview of adversarial attacks. -Analyze multi-step jailbreaking attacks by utilizing extensive experiments and discussing LLMs' privacy implications. 	<ul style="list-style-type: none"> -The study does not mention the other important adversarial attacks where authors focus most on multi-step jailbreaking attacks. -The scope of study does not involve adversarial attacks and other attacks. - The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations. -The scope of the study does not cover AI Ethical principles.
15	<ul style="list-style-type: none"> -Discussed and investigated some corresponding challenges toward building trustworthy AI compliance regulations and laws. - Proposing to make a policy to ensure that ChatGPT is regulated and deployed for the benefit of society at large 	<ul style="list-style-type: none"> - The scope of the study does not cover AI Ethical principles. - The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.
16	<ul style="list-style-type: none"> - Provide a brief overview of ethical and social risks associated with LLMs. - Point out some significant shortcomings in the landscape of potential ethics and social risks associated 	<ul style="list-style-type: none"> - The scope of the study does not cover AI Ethical principles. - It does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.
17		<ul style="list-style-type: none"> -Briefly explained some ethical Implications of ChatGPT and the regulation of disinformation propagation. -Investigating potential vulnerabilities and threats that may arise from the use of AI chatbots and mitigating these risks by utilizing existing framework.
18		<ul style="list-style-type: none"> -Classifying ethical and social risks that are associated with Large Language Models into 21 risks of harm. -Discuss some future solutions that need to be developed.
19		<ul style="list-style-type: none"> - Discuss Ethical principles and controls concerning ChatGPT along with its architectural components and training data. -Providing a comprehensive study addressing the ethical and privacy considerations associated with ChatGPT and providing insights into mitigation strategies. -Investigating the role of ChatGPT in cyberattacks, highlighting potential security risks.
21		<ul style="list-style-type: none"> -Provide a deep study to provide a comprehensive identification of the challenges and threats associated with the use of ChatGPT. -Discuss and identify main threats to the AI market.
22		<ul style="list-style-type: none"> -Briefly explained some ethical and social

	<p>implications aspects of ChatGPT.</p> <p>-Discuss the potential risks of AI-generated language with certain risks such as spread of misinformation and biased language.</p>	<p>ChatGPT without providing any technical details or discussing further the security measures.</p> <p>- The scope of the study does not cover AI compliance regulation's impact on ChatGPT.</p> <p>- the study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.</p> <p>- The scope of the study does not cover AI Ethical principles.</p>
23	<p>-Provide an overview of the significant security and privacy threats associated with LLM.</p> <p>-Examine the current measures and tactics employed to reduce the impact of these security and privacy attacks.</p>	<p>- The scope of the study does not cover AI compliance regulations' impact on ChatGPT.</p> <p>- The scope of the study does not cover AI Ethical principles.</p> <p>- The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.</p>
24	<p>-Proposing a systematically multifaceted role of LLMs in security and privacy.</p> <p>-Presented the overlap between LLMs (large language models) and security and privacy, as well as several strategies for defending against a range of risks and vulnerabilities associated with LLM.</p> <p>-Investigating the possible inherent vulnerabilities and various difficulties, threats, and risks linked to LLM.</p>	<p>- The scope of the study does not cover AI compliance regulations' impact on ChatGPT.</p> <p>- The scope of the study does not cover AI Ethical principles.</p> <p>- The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.</p>
25	<p>-Provides a comprehensive analysis of the responsible and ethical use of ChatGPT in the education sector. It also presents and analyses the integration of ChatGPT with education, which requires consideration of privacy risks.</p>	<p>- The scope of the study does not cover security and privacy's impact on ChatGPT as well as the necessary data classification and awareness.</p> <p>- The study lacks technical solutions.</p> <p>- The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.</p>

	<p>-Presents the importance of implementing responsible and ethical principles that govern the application of ChatGPT is highlighted, and the responsible and ethical use of artificial intelligence in education extends beyond ensuring technological accuracy.</p>	<p>- The scope of the study does not cover AI Ethical principles.</p>
26	<p>-Provides an overview of aspects of ChatGPT and its capability to address cybersecurity problems.</p> <p>-Investigating input prompts in ChatGPT, which should be clear and specific to ensure accurate outputs and address potential limitations and biases.</p> <p>-Provided a brief explanation of some challenging issues in privacy and ethical principles.</p>	<p>-The study briefly presented some challenging issues in privacy and ethical principles without pointing out the main differences between them.</p> <p>- The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.</p> <p>- The scope of the study does not cover AI Ethical principles. security and privacy issues.</p> <p>- The scope of the study does not cover AI Ethical principles.</p>
27	<p>-Provides an overview of the ethical and societal impacts of ChatGPT usage and presents various types of security and privacy attacks.</p> <p>-Briefly explained some robust security and privacy measures.</p>	<p>- The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.</p> <p>- The scope of the study does not cover security and privacy issues.</p> <p>- The scope of the study does not cover AI Ethical principles.</p>
28	<p>-Investigates in detail the security and privacy risks that inherently arise with ChatGPT.</p> <p>-Provide a study of the significance of the regulatory framework that governs generative AI deployment and human interaction.</p>	<p>--There's a need to extend the scope of the study to cover AI Ethical principles.</p> <p>- The scope of study needs to be extended, involving more solutions.</p> <p>- The study does not mention data classification and awareness, which are essential parts of the global Data and AI compliance regulations and Saudi DMO regulations.</p>

This thesis should build on these studies in order to cover most of the unaddressed risks and challenges. Along with that common research gap, there is an absence of data classification and awareness in most of these studies. Finally, the thesis aims to bridge the gap between relevant studies and discuss possible risk mitigation strategies.

In this thesis, we focus on the main gaps, which are security, privacy, and AI compliance regulation.

3 Methodology

This section presents our Research Methodology Stages.

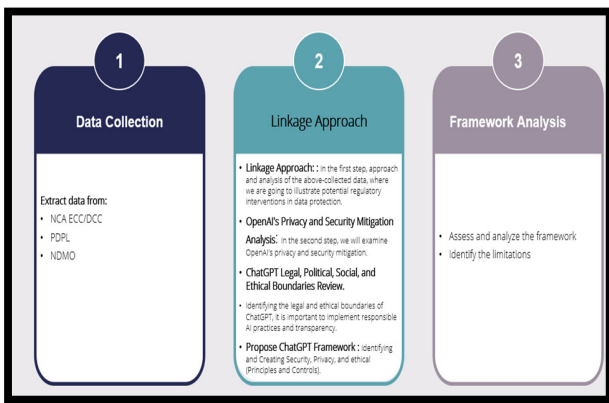


Figure 3.1 Research Methodology Stages

3.5 The Linkage Approach

This section presents our study approach and analysis of the above-collected data, illustrating potential regulatory interventions in data protection. The linkage methods will involve integrating NCA and SDAIA compliance regulations and laws.

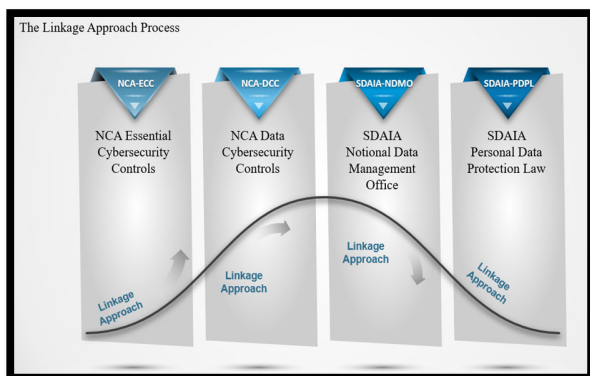


Figure 3.2 Summary of the linkage approach process

3.6 OpenAI's Privacy and Security Mitigation Analysis

In the first step, we examine OpenAI's privacy and security mitigation, including current privacy laws and regulations, ChatGPT's privacy violations, privacy leakage threats, data collection, personal input, privacy references, malicious use, and system vulnerabilities. Mitigations in OpenAI's privacy and security explain security and privacy ideas and techniques that may mitigate ChatGPT problems. Studying these mitigations gives concepts and technologies to explore as part of identifying relevant and appropriate security and privacy measures. In addition to Investigating the Saudi NCA Data Cybersecurity Controls and Data Management and Personal Data Protection Guidelines produced by the National Data Management Office, aggregating laws, and regulations of data protection, and linking potential regulatory interventions.

3.7 ChatGPT Legal, Political, Social, and Ethical Boundaries Review

Investigating the legal and ethical boundaries of ChatGPT, it is important to implement responsible AI practices and transparency in data usage, bias mitigation techniques, legal problems, and monitoring of generated content for harmful or misleading information.

3.8 Propose ChatGPT Framework: Identifying and Creating Security, Privacy, and ethical (Principles and Controls)

After conducting a review context in the steps mentioned above, we identify the list of regulations and laws that have been decided for further review and analysis and tailored in line with the ChatGPT usages scoping.

Overall, this proposal can help organizations govern and regulate the level of data protection and identify any potential gaps in ChatGPT uses. It also contributes to aggregating the benefits of ChatGPT and preventing their misuse.

3.9 Framework Analysis

After completing the data collection and review, we proceed with a comprehensive analysis of the methodology to assess and analyze its effectiveness. This process enables us to identify both the strengths and potential limitations of our approach, ensuring a well-rounded understanding of its impact and areas for improvement.

4 Results and Discussion

4.5 The Linkage Approach Process

To effectively link the controls from NCA (National Compliance Authority), NDMO (National Data Management Office), and PDPL (Personal Data Protection Law), We create a detailed table that highlights how the controls and requirements from each regulatory framework are intervened. This table will help you identify potential areas for regulatory intervention. Additionally, we provided a description of the methods that have been utilized where Control Area to specific aspect of data protection and management being addressed and regulation requirements where controls and guidelines as specified by the regulators, finally, Intervention Align practices to harmonize compliance with NCA, NDMO, and PDPL standards and implementing interventions to ensure compliance across all frameworks.

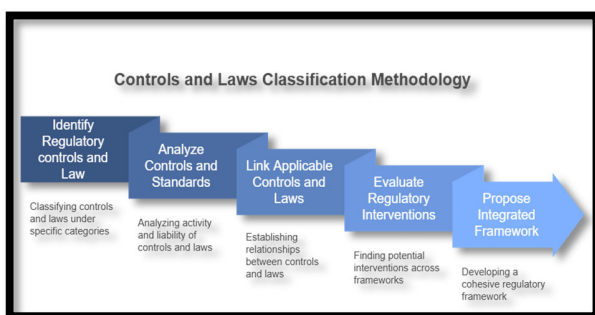


Figure 4.1 Controls and Laws Classification Methodology

4.6 The importance of the integration of security, privacy, and ethics controls

The integration of security, privacy, and ethics into a unified compliance framework ensures a holistic approach that covers security, privacy and ethical aspects of generative AI uses:

- **Privacy** Generative AI systems often process user inputs that could contain personal, sensitive, or proprietary information. As privacy concerns escalate globally, with regulations like PDPL, NDMO, NCA, GDPR, CCPA, and others enforcing strict requirements, this pillar ensures compliance, minimization of data exposure, and protection of user autonomy, ensuring that AI systems respect the confidentiality of user interactions.
- **Security** forms the backbone of any technology framework, especially when dealing with generative AI systems that process vast amounts of data. Generative AI such as ChatGPT often interacts with sensitive and confidential information, raising risks such as data breaches, malicious exploitation, or system vulnerabilities. This pillar ensures protection against cyber threats, implementation of robust security controls, and mitigation of risks related to adversarial attacks, where AI outputs can be manipulated or exploited.
- **Ethics** is the third pillar, ensuring that the deployment of generative AI aligns with societal values, fostering accountability, fairness, and transparency, this pillar ensures Bias and fairness, Transparency and Accountability. For this reason, Ethics extends the framework beyond technical and legal considerations, addressing societal trust, reputational risks, and the long-term sustainability of AI adoption.

These pillars lessen and mitigate the technological, legal, and social dimensions of generative AI, where ensuring that its adoption is secure and provides privacy as well as respect for individual rights and is aligned with societal values. This balance is essential, not only for compliance matters but also for fostering trust, driving data innovation, and ensuring the positive impact of AI usages on individuals and organizations.

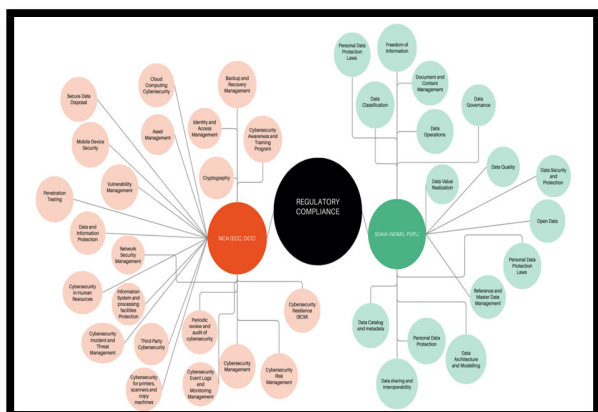


Figure 4.2 Detailed list of proposed security, privacy and ethical controls linked with relevant regulations

4.7 The Proposed Generative AI Policy Framework

The regulation of the artificial intelligence landscape is a global challenge, and preserving the security and privacy of generative AI systems is crucial, where Organizations must safeguard sensitive data and ensure the ethical application of AI. It is important to regulate the ChatGPT Generative AI scene, as no specific regulations currently govern its secure use. Furthermore, there is a lack of awareness about how to manage sensitive data and its classifications. In response to this, we have developed a generative AI policy framework to assist both organizations and individuals. According to this, we have drafted comprehensively a generative AI policy farmwork to help organizations and individuals to adopt a safe and secure environment that might help to mitigate some security and privacy risks while interacting with ChatGPT or any other Gen AI systems where establish robust security and privacy measures for generative AI initiatives, and this framework shall apply to all employees, contractors, and third parties involved in AI development and operations, where it will foster a culture of security and responsible AI use within the organization, however it requires adaptable policies that can shift as regulations change, ensuring that they always reflect the latest requirements and best practices from NCA, NDMO, and PDPL. Moreover, aggregate international efforts of all AI governance organizations and stakeholders to develop a worldwide solid regulation

framework across all borders to tackle all aspects of ChatGPT, and its social, legal, political, technological, and ethical boundaries efficiency and effectively.

Nicola Lucchiet. [14] conducted a case study to provide a comprehensive identification of the generative AI challenges, where they illustrate the importance of generative AI infrastructures, which has introduced significant regulatory problems in the field of AI. The European Commission is presently droughting the AI Act, the initial laws on artificial intelligence by a significant regulatory body, aimed at governing an evolving technology that has experienced increased investment and popularity, especially since the launch of ChatGPT and its derivatives. Analogous to the EU's General Data Protection Regulation (GDPR) enacted in 2018, the EU AI Act has the potential to establish a worldwide standard, affecting the extent to which AI could have either positive or negative impacts on citizens' lives worldwide. The draft regulation is currently in the trilogue phase, during which EU lawmakers and member states will finalize the regulation's specifics.

Therefore, regulating the Generative AI landscape is important for promoting security, privacy, and responsible AI practices and making sure that data usage, adopting bias mitigation techniques, resolving legal issues, and monitoring generated content for harmful or misleading information are all clear for all consumers. For this reason, we encourage adoption of proposing the development of international legislation and regulations, policymakers should begin establishing comprehensive regulations for the adoption of ChatGPT usage and other generative AI applications, as there are currently no specific regulations governing its secure and privacy use. Ultimately, the responsible development and deployment of large language models like ChatGPT requires a comprehensive and proactive approach to privacy and security to ensure a well-informed and inclusive approach to generative AI systems governance and an improved overall information security posture for the ChatGPT AI Generative Landscape. For these needs, we proposed a linkage approach that connects all regulations, including NCA ECC and DCC controls and SDAIA NDMO and PDPL standards to ensure that usage of Generative AI such as ChatGPT in order to satisfy minimum security and privacy requirements. Below is a proposed

framework for Generative AI systems such as ChatGPT, which are organized into three main areas: Security, Privacy, and Ethics.

4.7.1 Gen AI Framework: Identifying and Creating Security, Privacy, and Ethical Principles and Controls

The proposed framework outlines the foundational principles and practical controls and principals that ensure security, privacy, and ethical integrity in the use and deployment of ChatGPT and other AI systems. These controls are designed to mitigate key challenges, such as data security, user privacy, and ethical concerns, while aligning with standards and regulatory requirements. Thereby, creating a framework for identifying and creating security, privacy, and ethical principles and controls for Gen AI usage involves systematic steps. This process should encompass both technological and governance considerations to ensure responsible AI deployment. Below figure illustrates 3 main pillars of security and privacy framework for generative AI.

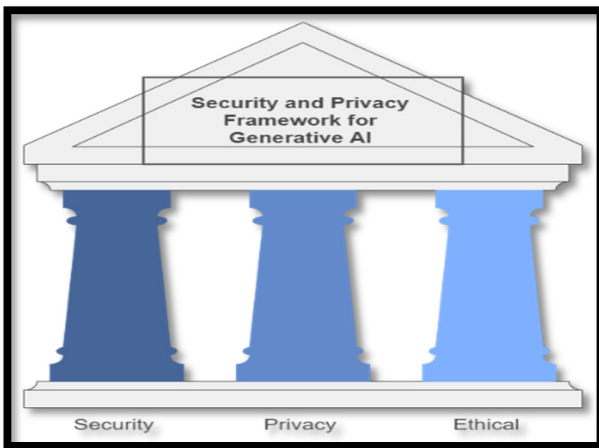


Figure 4.5 Illustrate the 3 pillars of security and privacy framework.

4.7.1.1 Security and Privacy Framework for Generative AI Usage: Security, Privacy, and Ethical Principles and Controls

To mitigate potential risks and threats associated with Gen AI, we must first build a robust framework

that regulates security, privacy, and ethics in the use of generative AI systems.

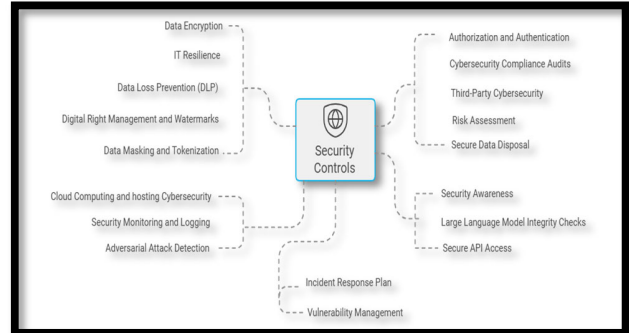


Figure 4.6 Illustrate the controls of Security Framework for Generative AI

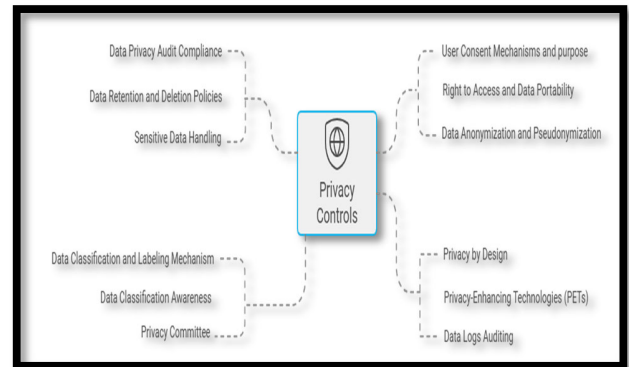


Figure 4.7 Illustrate the controls of Privacy Framework for Generative AI

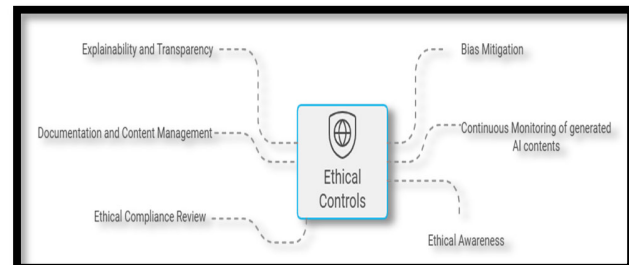


Figure 4.8 Illustrate the controls of Ethical Framework for Generative AI

4.8 Summary

The proposed Security, Privacy, and Ethical Framework for Generative AI provides a comprehensive approach to ensure responsible usage of generative AI systems. By implementing robust principles and controls, organizations can minimize risks, enhance user trust, and comply with global standards and regulations. The framework is designed to adapt to evolving technologies and legal landscapes,

ensuring ongoing security, privacy, and ethical adherence in AI-powered applications. As a result, the primary challenges of the proposed generative AI framework are that it requires continuous monitoring of the generative AI landscape against adversarial techniques, AI ethics principles, and cross-country compliance, which are different legal requirements, leading to additional complexity and compliance considerations. Untimely, cyber threats are continuously evolving, necessitating frequent updates to the framework. Keeping pace with the changes in generative AI and the cybersecurity landscape can be challenging. Collaboration with privacy experts, security researchers, and regulatory bodies can help organizations navigate these complex challenges and build trust with users.

It is worth mentioning that Saudi Arabia by The Saudi Data and Artificial Intelligence Authority (SDAIA) has released the Generative Artificial Intelligence 10 Guidelines in Government Entities and Individuals in Saudi Arabia in January 2024.

This document aims to tackle the responsible and effective use of GenAI, as well as regulate the use of data in Generative AI systems across all industries in the Kingdom of Saudi Arabia. Therefore, to leverage the benefits of generative AI, one must apply AI principles and implement preventive measures to mitigate potential risks such as data leaks, misinformation, deepfake, bias, and injustice. Furthermore, it necessitates safeguarding generative AI content, especially when employed for crucial decision-making without adequate verification and validation.

5 Conclusion

In conclusion, our research has linked the appropriate controls and standards of the NCA ECC, DCC, SDAIA NDMO, and PDPL to provide a clear usage policy that protects organizations and individuals when interacting with ChatGPT or other generative AI systems and lowering its security and privacy impacts. We also highlight challenges and illustrate security and privacy enhancements to make the usage of ChatGPT secure, trustworthy, and ethical. In addition, keeping up with the fast-changing threats of ChatGPT AI Generative and making strict rules and

comprehensive regulations could help combine benefits of ChatGPT with the need to limit users from abusing either intentionally or unintentionally. Therefore, it requires raising the level of data cybersecurity awareness among all segments of society in the Kingdom of Saudi Arabia. The research also provides detailed mitigations and recommendations for each concern. Moreover, propose a generative AI policy framework for policies that can help guide the development, deployment, and use of generative AI systems in a secure, ethical, and accountable manner. It is crucial to involve various stakeholders in the process of formulating AI policies, including security, privacy, lawyer experts, researchers, policymakers, and organizations, to ensure a well-informed and inclusive approach to AI governance and an improved overall information security posture for the ChatGPT AI generative seen. We believe the research will be beneficial to organizations that are seeking to protect the usage of ChatGPT, as such a policy would safeguard its usage and reduce its security and privacy impact on society. Furthermore, it will help in raising the cybersecurity maturity level of organizations to mitigate the ever-growing generative AI cyber-attacks.

Recommendations

The study reached its goals by answering all the research questions set out in the first chapter.

This was accomplished by achieving study objectives and examining the relationship between the NCA ECC, DCC, and SDAIA NDMO, PDPL standards, controls, and laws. The research findings provide impetus for the formulation of the following recommendations:

Creating a robust policy framework for generative AI will ensure compliance and accuracy of large language models. Also, align with various regulations, including PDPL, NDMO, NCA, GDPR, CCPA, and HIPAA. Being familiar with these regulations allows organizations and individuals to effectively maintain data security and privacy matters and ensure that a generative AI system stratifies security, privacy, and legal requirements.

Adopt Privacy-Enhancing Technologies (PETs), where those techniques and approaches help enhance data security, meet legal requirements, ensure confidentiality and integrity, and build trust among stakeholders by reducing the risk of data breaches, disclosures, or unauthorized access. Hence, it enables organizations to capitalize on data to get helpful insights while preserving data privacy, thus increasing the values generated from data.

Integrate Privacy by Design. Adopting the concept of privacy by design in AI Gen initiatives necessitates incorporating privacy standards at the outset of the development process in order to proactively preserve data privacy. The concept is critical not only to ensure compliance with data privacy and security regulations and legislation, but also to promote a culture of ethical responsibility in society. As a result, it will improve user trust and mitigate privacy risks. Furthermore, it will help organizations become leaders in ethical AI, giving them a competitive advantage in today's data-driven industry. In the context of advancing technology and society's cultural changes, it is crucial for the parties concerned to engage in ongoing communication to effectively adjust and enforce these privacy and security measures. The thesis lays the groundwork for a continuous endeavor to enhance the security and privacy of generative artificial intelligence (GenAI) initiatives. Adopting the recommendations provided for future works would not only reinforce Saudi Arabia's digital borders but also enhance overall data security and privacy capabilities at both the national and international levels.

6. Future Directions

Privacy-enhancing technologies are essential for maintaining individuals' privacy rights in a society that is becoming increasingly interconnected and reliant on data. By integrating these technologies into generative AI initiatives, organizations can mitigate security and privacy risks by building trust with users and adhering to data security and privacy regulations and laws. One of the highly recommended future works is to perform, analyze, and examine different privacy-enhancing technologies (PETs), and any organization handling sensitive data must implement PETs to satisfy the compliance and accuracy requirements of generative AI systems, as well as to share data both internally and externally within the organization in a secure and private manner. Therefore, organizations should adopt a combination of different PETs to cover all data use cases. Hence needs to be a deep dive into PETs to have a comprehensive solution

that can cover data use cases. Additionally, studying the advantages and disadvantages of privacy-enhancing technologies is crucial. While some technologies are more adaptable to specific data use cases than others, most organizations must implement privacy enhancing technologies (PETs) to cover comprehensively all data use cases. It's important to identify challenges and limitations of each type.

6 References

- [1] Derner, E., & Batistič, K. (2023). Beyond the Safeguards: Exploring the Security Risks of ChatGPT. arXiv preprint arXiv:2305.08005.
- [2] Sebastian, Glorin, Exploring Ethical Implications of ChatGPT and Other AI Chatbots and Regulation of Disinformation Propagation (May 29, 2023). Available at SSRN: <https://ssrn.com/abstract=4461801> or <http://dx.doi.org/10.2139/ssrn.4461801>.
- [3] Roland Hung (April 24 2023).AI Technology and Privacy: Canadian Privacy Commissioner Launches Investigation into ChatGPT.
- [4] Sakib Shahriar and Kadhim Hayawi, Let's have a chat! A Conversation with ChatGPT: Technology, Applications, and Limitations (Feb.2023) DOI: DOI: 10.47852/bonviewAIA3202939
- [5] Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, L.; Polosukhin, I. Attention is All You Need. In Proceedings of the Advances in Neural Information Processing Systems, Long Beach, CA, USA, 4–9 December 2017; pp. 6000–6010. arXiv:1706.03762.
- [6] Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. Bert: Pre-training of deep bidirectional Transformers for language understanding.arXiv 2019, arXiv:1810.04805.
- [7] Kingma, D.P.; Ba, J. Adam: A Method for Stochastic Optimization. In Proceedings of the 3rd International Conference on Learning Representations (ICLR), San Diego, CA, USA, 7–9 May arXiv:1412.6980.
- [8] Howard, J.; Ruder, S. Universal language model fine-tuning for text classification. arXiv 2018, arXiv:1801.06146.
- [9] Identifying and Mitigating the Security Risks of Generative AI Clark Barrett,Brad Boyd,Elie Bursztein(Jan, 2024), arXiv:2308.14840v1 [cs.AI] 28 Aug 2023.
- [10] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam ` Roberts, Tom B. Brown, Dawn Xiaodong Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting ` training data from large language models. In USENIX Security Symposium, 2020.
- [11] MAANAK GUPTA, CHARANKUMAR AKIRI, KSHITIZ ARYAL, ELI PARKER, AND LOPAMUDRA PRAHARAJ, From ChatGPT to ThreatGPT: Impact of Generative AI in Cybersecurity and Privacy (August,2023) DOI: 10.1109/ACCESS.2023.3300381.
- [12] Iman Almomani, Mohammed Ahmed and Leandros Maglaras, Cybersecurity maturity assessment framework for higher education institutions in Saudi Arabia (September 9, 2021) DOI: 10.7717/peerj-cs.703.
- [13] Emna Chikhaoui , Alanoud Alajmi , Souad Larabi-Marie-Sainte *Artificial Intelligence Applications in Healthcare Sector: Ethical and Legal Challenges (August, 2022) (ISSN: 2610-9182).
- [14] Lucchi N. ChatGPT: A Case Study on Copyright Challenges for Generative Artificial Intelligence Systems. European Journal of Risk Regulation. Published online 2023:1-23. doi:10.1017/err.2023.5