



Digital deception: generative artificial intelligence in social engineering and phishing

Marc Schmitt¹ · Ivan Flechais¹

Accepted: 24 September 2024 / Published online: 12 October 2024
© The Author(s) 2024

Abstract

The advancement of Artificial Intelligence (AI) and Machine Learning (ML) has profound implications for both the utility and security of our digital interactions. This paper investigates the transformative role of Generative AI in Social Engineering (SE) attacks. We conduct a systematic review of social engineering and AI capabilities and use a theory of social engineering to identify three pillars where Generative AI amplifies the impact of SE attacks: Realistic Content Creation, Advanced Targeting and Personalization, and Automated Attack Infrastructure. We integrate these elements into a conceptual model designed to investigate the complex nature of AI-driven SE attacks—the Generative AI Social Engineering Framework. We further explore human implications and potential countermeasures to mitigate these risks. Our study aims to foster a deeper understanding of the risks, human implications, and countermeasures associated with this emerging paradigm, thereby contributing to a more secure and trustworthy human-computer interaction.

Keywords Artificial intelligence · Machine learning · Social engineering · Phishing · ChatGPT · Large language models

1 Introduction

Digital deception, fueled by the advancements in generative AI, poses a serious threat to our interconnected society. As AI systems become increasingly adept at mimicking human communication and trust signals, they present a new frontier for social engineering and phishing attacks. Given that social engineering attacks can arise in all types of interaction—computer-to-computer, human-to-computer, or human-to-human—urgent research is needed: both to understand the challenges, and to identify ways of protecting safe and meaningful digital interactions when such technologies are used for deceit.

✉ Marc Schmitt
marcschmitt@hotmail.de

¹ Department of Computer Science, University of Oxford, Oxford, UK

The rapid advancements in Artificial Intelligence (AI) and Machine Learning (ML) are poised to revolutionize many aspects of our lives (Briggs and Kodnani 2023; Chui et al. 2023), with significant implications for the security and privacy of our digital interactions (Kaur et al. 2023). The growing dependence on technology opens a wide playground for cybercriminals, posing threats to all stakeholders in the global economy (Yamin et al. 2021; Bécue et al. 2021; Schmitt 2023a). Cybercrime is quickly increasing in scale and impact, causes serious financial harm, breaches individual privacy, and has the potential to break down critical infrastructure. The World Economic Forum (WEF) consistently ranks cyber threats as one of the most significant global risks. AI/ML technologies have demonstrated their potential for helping to detect social engineering and phishing attacks, however these technologies can also be applied in a malicious manner to amplify the capabilities and effectiveness of such attacks.

According to the UK governments 2024 Cyber Security Breaches Survey, many businesses and charities faced cyber security breaches or attacks in the past year. Phishing was the most common, affecting 84% of businesses and 83% of charities.¹ The Zscaler ThreatLabz 2024 Phishing Report shows a 58.2% increase in phishing attacks in 2023, highlighting the growing sophistication of threat actors. In addition, voice phishing (vishing) and deepfake attacks are rising as attackers use generative AI to enhance their tactics (Zscaler 2024). For instance, AI-powered face-swapping can impersonate a victim's friend during a video call, showcasing advanced AI-driven social engineering tactics. E.g., a finance worker at a multinational firm was tricked into paying \$25 million when deepfake technology was used to impersonate the chief financial officer. In fact, every participant during this multi-person video conference was a deepfake.² Additionally, voice phishing, where an individual's voice is faked using AI, is even easier. This method has been used in various cyber-attacks, including a well-documented case where a CEO's cloned voice authorized a fraudulent bank transfer.

The WEF Global Cybersecurity Outlook 2024 warns that emerging technologies will exacerbate longstanding cyber resilience challenges, particularly due to advances in adversarial capabilities such as phishing, malware development, and deepfakes (World Economic Forum 2024).

Understanding the nature of AI-driven social engineering attacks and their techniques is crucial for individuals and organizations to implement proactive strategies, and defenses against this "relevant" and "permanent" cybersecurity threat (Kaur et al. 2023; Malatji and Tolah 2024).

The increasing sophistication of AI algorithms and the availability of vast amounts of data can empower malicious actors to develop more sophisticated and convincing phishing techniques. AI/ML can be misused to automate and personalize phishing attacks, generate persuasive content that mimics legitimate communication, and evade detection by traditional security measures (Taddeo et al. 2019; Renaud et al. 2023; Kaur et al. 2023). As security threats evolve, it becomes crucial to review and explore the potential risks and vulnerabilities associated with AI/ML to better understand this new threat landscape, and to develop effective countermeasures for individuals, organizations, and broader digital ecosystems.

¹ 2024 UK Cyber Security Breaches Survey: <https://www.gov.uk/government/statistics/cyber-security-breaches-survey-2024/cyber-security-breaches-survey-2024>.

² CNN: <https://edition.cnn.com/2024/02/04/asia/deepfake-cfo-scam-hong-kong-intl-hnk/index.html>.

The recent rise of generative AI – in particular, large language models (OpenAI 2023; Anil et al. 2023) – is a new emergent capability with the potential to carry out highly targeted social engineering campaigns at an industrial scale (Neupane et al. 2023). With the help of AI-enabled autonomous agents (Reed et al. 2022), it is possible to automate parts – potentially even the entire lifecycle – of a social engineering attack. In addition, fully trained SE models can theoretically learn from each attack and gradually improve their success rate. The emergence of advanced open-source generative AI tools represents significant progress in machine learning and AI, particularly in their capacity to produce sophisticated text, voice, graphics, and video content. This is particularly evident by recent developments in large language models (LLMs) such as GPT-4 (ChatGPT) and Gemini. Those tools have raised substantial concerns about their potential misuse in social engineering attacks (Neupane et al. 2023). In conjunction with techniques like jailbreaking, prompt injection, and reverse psychology, these models pose a notable threat to the cybersecurity landscape (Gupta et al. 2023), however much is still unknown about this threat and its characteristics. To address this gap, our research aims to investigate the following question: How can Generative AI amplify the effectiveness – and hence the global problem – of social engineering and phishing attacks?

This research aims to examine and analyze the various attack options available in social engineering and phishing, and identify those areas which are facilitated using AI/ML. It seeks to understand how threat actors can harness AI/ML to carry out sophisticated social engineering and phishing campaigns and to provide insights for threat intelligence to develop future mitigation strategies. The findings of this study will be valuable for researchers, practitioners, and policymakers, guiding them in understanding the evolving landscape of AI-enabled social engineering and phishing attacks.

This study addresses three research questions:

- RQ1: How can generative AI likely be used – by cybercriminals – to enhance the effectiveness of social engineering and phishing attacks?
- RQ2: How does the integration of AI in social engineering and phishing attacks – especially generative models – pose new challenges and demand the development of advanced defense mechanisms?
- RQ3: How can we effectively detect and mitigate these AI-driven attacks?

This paper is structured as follows: The [Methodology](#) section outlines our research methods and shows the concrete steps taken to develop the Generative AI Social Engineering Framework. This is followed by a theoretical foundation section that discusses social engineering, phishing, and artificial intelligence capabilities. The section on generative AI-enabled social engineering attacks is divided into three subsections: (i) realistic content

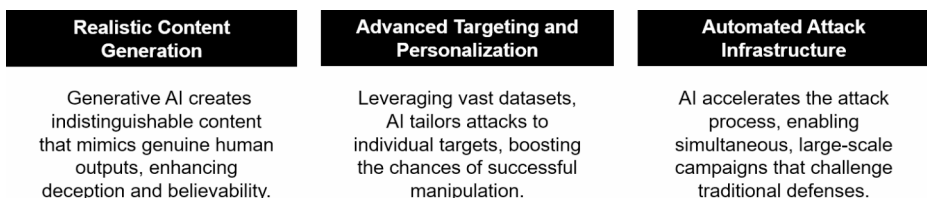


Fig. 1 Three pillar framework of generative AI-enabled social engineering attacks

generation, (ii) advanced targeting and personalization, and (iii) automated attack infrastructure, which are the three pillars of our framework (see Fig. 1). The paper then discusses the impact of increasing capabilities and cost-effectiveness of AI-enabled phishing attacks, explores human implications and potential countermeasures, and concludes with details of future research in this area.

2 Methodology

To address the complexities and challenges posed by the integration of Generative AI in Social Engineering (SE), we employ a mixed-methods research approach that combines qualitative and framework-based analyses. Our methodology comprises the following steps:

- **Systematic Review:** We start by conducting a systematic review of existing literature to understand the state of the art in both Social Engineering and Generative AI. This review provides the foundation upon which we built our subsequent analyses. Biases were mitigated by conducting a thorough search across multiple databases (Scopus, Web of Science, and Google Scholar) to ensure a diverse set of studies, and by adhering to predefined inclusion and exclusion criteria for consistent and objective selection. We also incorporated industry reports (e.g., World Economic Forum) due to the newness of the topic and to provide practical insights from industry experts.
- **AI Capabilities Analysis:** Leveraging the Social Engineering Framework developed by Mouton et al. (2014), we analyze the capabilities of AI technologies, placing a specific emphasis on Generative AI. This allows us to understand how these technologies fit into existing SE tactics and strategies.
- **Identification of Pillars:** Based on our analysis, we identify three primary pillars where Generative AI exacerbates the impact of SE attacks. Initially, our systematic review identified a range of generative AI capabilities, which we categorized into thematic clusters. Further synthesis and grouping of these clusters revealed three essential pillars: (1) Realistic Content Creation, (2) Advanced Targeting and Personalization, and (3) Automated Attack Infrastructure.
- **Conceptual Framework Development:** Informed by our findings, we develop a conceptual framework that we refer to as the “Generative AI Social Engineering Framework.” The model offers a blueprint for deeper and broader investigations, and the derivation of potential implications and solutions.
- **Applications:** We then validate and apply the developed framework to (1) investigate the impact of generative AI on different types of phishing attacks in terms of threat amplification and cost-effectiveness, and (2) to identify countermeasures that can mitigate the risks associated with each of the three pillars of the framework.

By adopting this methodology, we aim to offer a robust and multidimensional analysis of the challenges and solutions associated with the role of Generative AI in Social Engineering and lay the groundwork for further research.

3 Literature background

3.1 Social engineering and phishing

Deception plays a central role in both social engineering and phishing tactics. In both cases, the attacker aims to manipulate or trick their target into taking certain actions or revealing confidential information by posing as a trustworthy entity. Social engineering refers to the activity of manipulating people into performing actions or divulging confidential information, usually through deceptive means. It exploits human psychology rather than relying on technological vulnerabilities. See Table 1 for an overview of different social engineering types.

Phishing attacks are used by adversaries to impersonate legitimate entities and send deceptive emails or texts, often containing malware attachments or links to fraudulent websites (Aleroud and Zhou 2017; Fahl 2021). Phishing attacks exploit human psychology, relying on social engineering techniques to trick victims into providing their confidential information or performing actions that compromise their security (Desolda et al. 2022).

The primary objectives of these scams are to infect devices with malware, steal sensitive data like usernames and credit card details, seize control of digital accounts such as email and social media, or directly persuade recipients to perform financial transactions. Phishing as a strategy exploits user interface flaws and takes advantage of the fact that humans have difficulties verifying URLs and dynamic content of rendered HTML documents (Fahl 2021). It is often the first point of access to a victim's device or account and can lead to further attacks, hence exploiting the trust developed between the compromised account and its contacts. Phishing attacks have been a pervasive and persistent threat in the digital landscape, impacting individuals, businesses, and even government entities. Those attacks have evolved, becoming increasingly sophisticated and difficult to detect. Attackers constantly adapt their techniques, leveraging advancements in technology to create more convincing and personalized phishing campaigns (Aleroud and Zhou 2017; Neupane et al. 2023). The consequences of phishing attacks can be severe, ranging from financial losses and identity theft to unauthorized access to sensitive data or compromise of entire systems. See Table 2 for an overview of different phishing types.

There are also other concepts, such as “catfishing”, where someone creates a fake identity on a social networking site, often to deceive someone into a fraudulent romantic relationship. While this is not “phishing” and the end goal might not always be financial gain or

Table 1 Different social engineering attack types

Type	Explanation
Pretexting	Attacker creates a fabricated scenario to obtain information or access from a target.
Baiting	Attacker lures target with appealing offer to compromise their security (e.g., free software).
Tailgating	Attacker deceives security personnel or employees to gain physical access to a restricted area.
Quizzing	Attacker tricks the target into disclosing information under the guise of a survey or quiz.
Scareware	Attacker tricks target into buying fake antivirus, claiming their computer is infected.
Phishing	Deceptive emails. The most well-known and widespread form of social engineering.

Table 2 Different phishing attack types

Phishing types	Explanation
Email phishing	The attacker sends fraudulent emails that appear to come from legitimate sources (e.g., banks, government agencies) to trick recipients into revealing sensitive information.
Spear phishing	Like email phishing, but the attacker customizes the deceptive email for a specific individual or organization.
Smishing	The attacker sends deceptive text messages (SMS) to trick recipients into revealing sensitive information or clicking on malicious links.
Whaling	A specialized form of spear phishing where high-profile targets, like CEOs or CFOs, are the primary focus. The intent is typically to manipulate these individuals into authorizing large money transfers or revealing sensitive corporate data.
Pharming	The attacker redirects the target's web traffic to a fake website that looks identical to a legitimate one, tricking them into entering their credentials or other sensitive information.
Vishing	The attacker uses deceptive phone calls to trick the target into divulging confidential information.
Social-media phishing	The attacker uses social media platforms (e.g. Facebook or Instagram) to deceive their targets (Seymour and Tully 2018).

direct access to sensitive information, it is another example of deception through social platforms. Regardless of the method or platform, the core principle remains: attackers exploit human psychology and trust to deceive individuals. When we refer to phishing, we mean all types of phishing as outlined in Table 2, and when we refer to social engineering, we mean all types of social engineering (including phishing) as shown in Table 1.

Based on an analysis of different social engineering attacks, Mouton et al. (Mouton et al. 2014) propose a model of social engineering attacks outlining several phases and characteristic actions. In Fig. 2, we reproduce the high level model of the framework, which shows how a social engineering attack progresses through different phases (attack formulation, information gathering, preparation, developing a relationship, exploiting a relationship, and debrief). We use this framework in Sect. 4 to help structure our analysis of Gen-AI capabilities on social engineering attacks.

3.2 Artificial intelligence capabilities

Artificial Intelligence (AI) has emerged as a transformative technology, primarily driven by key factors such as the abundance of data, improved algorithms, and advances in hardware infrastructure. Machine Learning (ML), a subfield of AI, focuses on developing algorithms and models that can learn patterns and make predictions or decisions based on data. Machine learning can be divided into supervised learning, unsupervised learning, semi-supervised learning, and reinforcement learning (Sutton and Barto 1998). Breakthroughs in deep learning, a subset of ML, have revolutionized the field by enabling the training of neural networks with multiple layers, leading to remarkable progress in areas like image recognition, natural language processing, and autonomous systems (LeCun et al. 2015).

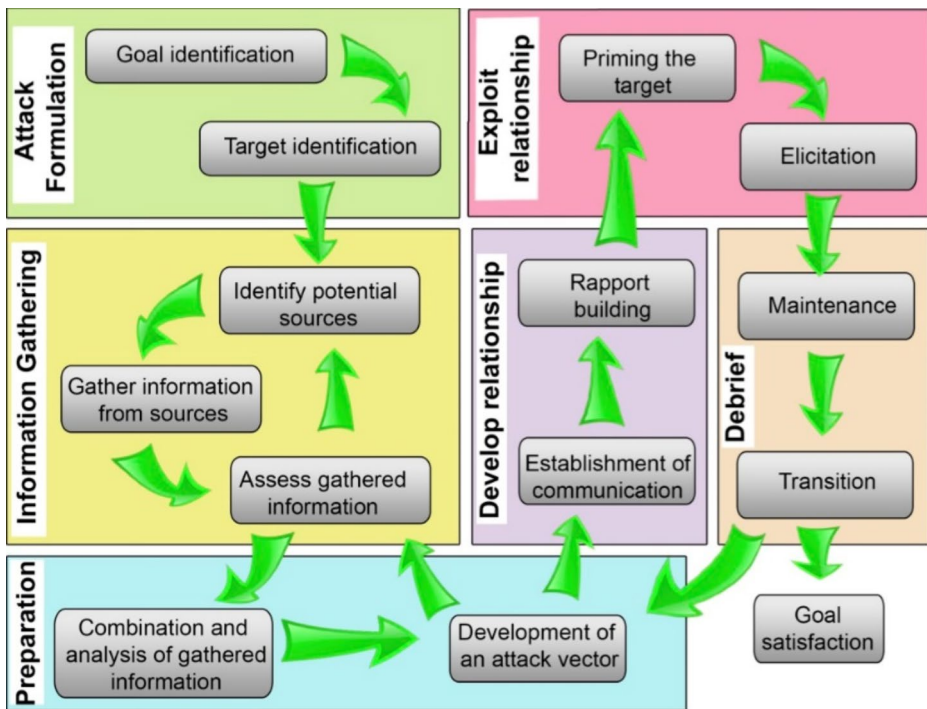


Fig. 2 Social engineering attack framework, reproduced from Mouton et al. (2014)

The impact of AI/ML is increasingly felt in various domains, including healthcare, finance, and cybersecurity (Cutillo et al. 2020; Kaur et al. 2023). AI-powered systems can automate tasks, enhance decision-making processes, and uncover valuable insights from vast amounts of data (Davenport 2018; Schmitt 2023b). Generative AI has opened a novel frontier in the realm of AI. Generative Adversarial Networks (GANs), a pioneering approach within this domain, have become particularly noteworthy for their ability to produce incredibly realistic digital assets ranging from images to videos (Goodfellow et al. 2014). Recently, Large Language Models (LLMs) show the culmination of deep learning techniques applied to natural language processing. LLMs, such as GPT-4 and Palm 2—which are based on the transformer architecture (Vaswani et al. 2017)—have shown unprecedented capabilities in generating realistic text content (OpenAI 2023). However, this power also raises ethical and security concerns, particularly as the line between human-generated and AI-generated content becomes increasingly blurred (Table 3).

Machine learning has gradually introduced a huge basket of “AI capabilities” that can be harnessed for social engineering and phishing attacks. See Table 4 for an overview of existing AI capabilities and how they could potentially be applied in this context.

Table 3 Artificial intelligence in a nutshell

Machine learning	Description
Supervised learning	A machine learning algorithm learns from labeled training data, where each example is paired with its corresponding target output. The algorithm generalizes from this labeled data to make predictions or classify new, unseen instances based on the learned patterns.
Unsupervised learning	Involves learning patterns and structures from unlabeled data without explicit target outputs. The algorithm identifies inherent relationships, clusters, or patterns within the data to gain insights and make sense of the underlying information.
Semi-supervised learning	The algorithm learns from a combination of labeled and unlabeled data. By leveraging the unlabeled data along with the labeled data, the algorithm aims to improve its predictive performance and make use of the additional information available.
Reinforcement learning	A learning paradigm where an agent learns to interact with an environment by taking actions and receiving feedback in the form of rewards or penalties. The agent aims to maximize its cumulative reward through an iterative process of exploration and exploitation.
Deep learning	A human brain inspired subset of machine learning, uses neural networks with multiple layers to process and analyze large data sets. The ability to detect complex patterns enabled advances in image and voice recognition, natural language processing, and beyond.
Generative AI	AI models – based on deep learning – designed to generate new content that resembles the input data they have been trained on. Examples include Generative Adversarial Networks (GANs) that produce synthetic images and videos, and Large Language Models (LLMs) like GPT-4 and Palm 2, that generate realistic textual content.

4 Capability analysis of generative AI in social engineering

Generative AI can be employed in social engineering and phishing scenarios in various ways, enhancing the effectiveness of these attacks by creating more convincing and targeted deceptive content. In this section, we will investigate the integration of AI capabilities into the various stages of the Social Engineering (SE) attack lifecycle introduced by Mouton et al. (2014). While Table 5 shows the entire SE attack framework and most machine learning capabilities, we will focus on the stages where the potential for Gen-AI is the largest.

Generative AI encompasses a broad range of advanced capabilities that enable the creation and manipulation of various types of content. In the context of social engineering, it is important to recognize that drawing a clear line between existing AI/ML and generative AI capabilities is impractical; rather, generative AI should be viewed as a complementary enhancement that seamlessly integrates with and amplifies existing AI functionalities, thereby increasing the potential for automation, realism, and overall effectiveness.

Table 4 Overview of potential AI capabilities in the context of social engineering

AI capabilities	Explanation	ML techniques
Generative AI	Generative AI involves algorithms that can generate content, such as text, images, or videos, based on patterns learned from existing data. In social engineering attacks, generative AI can create realistic and convincing attack vectors, such as phishing emails, by imitating human communication styles and context.	Generative Adversarial Networks (GANs), Transformer models
AI analysis	AI analysis refers to the application of machine learning and data analysis techniques to process and interpret data. In social engineering attacks, AI analysis can identify potential targets, assess their vulnerabilities, and predict their behavior based on patterns in gathered information.	Machine Learning, Classification and Regression, Natural Language Processing (NLP)
AI scraping	AI scraping entails the use of automated tools, often driven by machine learning, to collect information from various online sources. In social engineering, AI scraping can swiftly gather data from social media profiles, public databases, and other sources to create detailed profiles of targets.	Web Scraping Libraries, Data Mining, Clustering Techniques
AI automation	AI automation refers to the use of AI-driven systems to automate various tasks and processes. In social engineering, AI automation can initiate and maintain communication with targets, ensuring consistent interaction and reducing the risk of detection.	Rule-based Systems, Process Automation, Workflow Management
AI chatbots	AI chatbots are computer programs that can simulate human conversation. In social engineering attacks, AI chatbots can engage targets in conversations to build trust, gather information, and manipulate emotions, all while emulating human-like interaction.	Large Language Models (LLMs), Contextual Chatbot Frameworks, Sequence-to-Sequence Models
AI coordination	AI coordination involves the orchestration of tasks and interactions among different AI agents or components. In social engineering, AI coordination could ensure smooth transitions between different phases of the attack and maintains continuity, even if attackers change.	Multi-agent Systems, Coordination Algorithms, Task Allocation Methods
AI assessment	AI assessment entails the use of algorithms to track, analyze, and evaluate the success of an attack. In social engineering, AI assessment can monitor the outcomes, such as compromised accounts or data leaks, to determine the effectiveness of the attack and refine future strategies.	Performance Metrics, Anomaly Detection, A/B Testing

Gupta et al. (2023) have published a recent paper that is directly relevant to our study. It examines Generative AI's cybersecurity implications, including vulnerabilities in ChatGPT that can be exploited for attacks such as Jailbreaks and prompt injections. The paper explores how cybercriminals might use GenAI for social engineering, phishing, and malware creation, and discusses how GenAI can enhance cybersecurity through automation, threat intelligence, and secure coding. The methodology is robust, and includes an examination of vulnerabilities, attack demonstrations, and defense techniques.

Kaur et al. (2024) have published a comprehensive study on deepfake video detection, highlighting significant threats posed by false content created using generative AI. This study sets a foundation for developing robust, real-time detection models and guides future research to enhance deepfake detection technologies.

Based on the literature review and subsequent analysis, we identify three primary pillars where generative AI exacerbates the impact of social engineering attacks. Initially, our systematic review identified a wide range of generative AI capabilities, which can be cat-

Table 5 AI utilization for each stage of the social engineering attack lifecycle

AI utilization	Stage of SE attack lifecycle	Application	AI capabilities
Attack formulation	Goal identification	Generating potential attack goals based on desired outcomes and vulnerabilities.	Generative AI
	Target identification	Analyzing public data to identify potential targets based on roles and online presence.	AI analysis
Information gathering	Identify potential sources	Scraping potential sources of information from public sources and social media.	AI scraping
	Gather information from sources	Gathering relevant information from various data sources.	AI scraping
	Assess gathered information	Automating the assessment and aggregation of information from multiple sources.	AI automation
Preparation	Combination and analysis of gathered information	Processing and analyzing collected data to identify patterns and vulnerabilities.	AI analysis
	Development of an attack vector	Crafting personalized attack vectors, like phishing emails, based on gathered information.	Generative AI
Develop relationship	Rapport building	Engaging in conversations to build rapport with the target.	AI chatbot
	Establishment of communication	Automating the initiation and maintenance of communication with the target.	AI automation
Exploit relationship	Prime the target	Generating content aligned with the established relationship to manipulate responses.	Generative AI
	Elicitation	Employing psychological techniques through AI chatbots to elicit information or actions.	AI chatbot
Debrief	Maintenance	Automating periodic interactions to sustain the established relationship.	AI automation
	Transition	Facilitating the smooth handover of interactions from one attacker to another.	AI coordination
Goal satisfaction	Tracking and analysis	Tracking and analyzing the success of the attack in achieving its goals.	AI assessment

egorized into distinct clusters such as natural language processing, computer vision, speech and audio processing, and multimodal generation. Further synthesis and thematic grouping of these clusters in the context of social engineering and phishing revealed three essential pillars: (1) Realistic Content Creation, (2) Advanced Targeting and Personalization, and (3) Automated Attack Infrastructure.

4.1 Realistic content generation

The easiest identifiable potential of Gen AI is the creation of realistic content. A good example—especially in the context of phishing—would be website cloning. AI can rapidly clone legitimate websites and modify them subtly to deceive victims, leading to more effective phishing pages. More concerningly, generative AI is gradually mastering different capabilities in text and media generation (Archana and Jeevaraj 2024; Gambín et al. 2024):

- Text: It can craft convincing emails, text messages, or social media posts that appear to be written by humans. These messages may be customized to target specific individuals or organizations.
- Images: It can create realistic images, such as fake IDs, official logos, or manipulated photographs, to enhance the credibility of phishing or social engineering attacks.
- Voice: It can produce realistic voice recordings or engage in voice phishing (vishing) by impersonating trusted individuals or organizations over the phone.
- Videos: It can create highly realistic deepfake videos that manipulate the appearance and voice of real individuals. These deepfakes can be used to spread misinformation, impersonate executives or public figures, and deceive targets into taking specific actions (Bray et al. 2023).

These manipulated media appear highly realistic and can deceive viewers into believing false narratives or events. For example, deepfakes (Gambín et al. 2024; Kaur et al. 2024) have raised significant concerns due to their potential to spread misinformation and undermine trust in visual and auditory evidence. In the context of SE and phishing, deepfakes can be created from datasets of personal visual and audio information. These datasets can be gathered from manual or automated open-source intelligence (OSINT), close-source information (CSINT), or even shared by other attackers. The combination of increasingly available personal information and increasingly sophisticated deepfake capability can be used for the following purpose:

- Impersonation: By creating a near-flawless visual and auditory imitation of a trusted individual—such as colleagues, friends, or family members—malicious actors can convince recipients to carry out tasks or disclose sensitive information.
- Bypassing Biometrics: Some security systems that rely on facial or voice recognition could be tricked by a well-crafted deepfake, thereby granting unauthorized access.
- Ransom Campaigns: Threat actors could create scandalous or compromising deepfakes of individuals and then threaten to release them unless a ransom is paid (Teichmann 2023).
- Disinformation: On a broader scale, deepfakes can spread misinformation or propaganda that appears genuine, thus causing panic, influencing public opinion, or even

affecting stock market behavior.

The evolution of AI has not stopped with visual manipulations. Recently, the introduction of sophisticated large language models (LLMs), such as ChatGPT or Gemini, introduced another layer of potential misuse in the field of SE and phishing. LLMs can craft highly convincing and personalized phishing emails that are contextually relevant to the recipient, vastly increasing the success rate of such attacks. This is of course also true for all other text-based communication channels, such as SMS or Social Media Chats (Seymour and Tully 2018). Instead of static phishing messages, cybercriminals could leverage LLMs to conduct real-time interactions with victims, adapting the conversation dynamically to manipulate the target more effectively. In addition, those models are developing very fast and gradually gaining new capabilities beyond text generation (i.e., ChatGPT Code Interpreter). An emerging concern with advanced AI systems like LLMs is the potential for “jailbreaking”—bypassing imposed usage restrictions or guidelines. Malicious actors, by finding and exploiting vulnerabilities in interfaces or access controls, can use AI capabilities beyond their originally intended purpose (Greshake et al. 2023; Gupta et al. 2023; Zou et al. 2023). In the context of SE and phishing, such unrestricted access can supercharge deceptive campaigns, making them highly sophisticated and more challenging to identify and counter.

4.2 Advanced targeting and personalization

Generative AI has the potential to start an era of hyper-personalized and highly targeted attacks, known as advanced persistent threats. This is because of the introduced ability to generate realistic content (pretexting), but also due to the potential of AI to analyze existing information (reconnaissance), which allows attackers to customize their malicious intents according to the online presence, behavior, and affiliations of the targets. In the final step, this information can be successfully leveraged during the execution stage.

The primary goal of reconnaissance (information gathering and intelligence) is to collect information about the target, which would allow the attacker to craft a more effective and tailored attack strategy in the subsequent stages. With its mastery of language and analytical abilities, Generative AI can scrutinize the digital footprints of targets. This provides insights into a target’s specific interests, affiliations, or behaviors. Such insights could pave the way for crafting deeply personalized phishing content or social engineering attacks. This gathered intelligence can subsequently be used to develop the attack strategy—referred to as pretexting. Pretexting is a broad stage that encapsulates the creation of a story, scenario, or identity that an attacker uses to engage with the target. Given the vast capabilities of Gen-AI, it can significantly enhance this phase by producing realistic and tailored content, such as texts, images, voices, and videos (Sect. 3.1). This enables context-aware phishing, where AI crafts malicious content that resonates with the target’s communication patterns, making the story or scenario highly believable. This might include emails that sound like they’re from colleagues, friends, or familiar institutions. During execution, AI exploits identified vulnerabilities and uses generated synthetic content to manipulate victims and achieve the objectives set during the reconnaissance and pretexting phases. A good example is automated spear phishing, where AI uses those insights to generate and distribute personalized fraudulent emails or messages. Also, intelligent chatbots can ask leading questions to extract sensitive information or directly execute certain actions. Smart bots, such as “jailbroken”

LLMs (i.e., ChatGPT and Gemini) have shown superior capabilities to pass the Turing test, and in general, appear like “real” humans (Webb et al. 2023). This is especially the case during purely text-based conversations.

4.3 Automated attack infrastructure

The last pillar is automation. AI can help to automate the creation and dissemination of deceptive content, thereby enabling attackers to carry out large-scale phishing campaigns or social engineering attacks with minimal effort. As established in Sect. 3.2, AI algorithms can be used to learn about the victim’s behaviors and predict their responses, which can make phishing attempts more convincing and successful. LLMs could be integrated into broader SE toolkits, offering automated and intelligent capabilities, like crafting believable responses or guiding the attacker based on real-time data. By introducing reconnaissance feedback loops at both the pretexting and execution stages, the AI becomes a dynamic system, capable of evolving its methods based on real-world outcomes. For instance, if an AI-driven chatbot notes hesitation or suspicion in a victim’s responses, it could change its approach, tone, or the information it provides to appear more legitimate. Repeated requests increase the likelihood of success (Fleming et al. 2023). This makes it a formidable tool in the hands of malicious actors and underscores the importance of robust countermeasures. The ultimate threat is a completely autonomous intelligent social engineering bot.

This brings us the next point, which is the large-scale distribution of such bots. AI significantly escalates the potential for expansive and sophisticated campaigns. Through AI-powered botnets, attackers can rapidly craft and disseminate personalized phishing messages to a massive audience. These botnets, equipped with AI capabilities, not only automate the distribution process but can also adapt the content in real time based on user interactions and feedback. Moreover, the integration of intelligent bots within these networks allows for autonomous responses to potential victims, convincingly simulating genuine human engagements. This confluence of automation, adaptability, and vast reach introduced by AI presents a heightened challenge, necessitating advanced countermeasures in our cybersecurity frameworks.

Furthermore, generative AI can generate content that evades detection by security software or filters, such as crafting phishing emails that avoid common red flags or creating malware that bypasses antivirus software. AI can use polymorphic attack patterns (e.g., emails that constantly change, dynamically altered malicious websites, etc.). Those AI-created evasion and camouflage techniques could make it very difficult to filter out malicious content. Similar to feedback loops, adapting attack models, can learn from failed attempts, continually improving their methods and becoming more effective at evading detection and/or – as mentioned earlier – deceiving their targets. Overall, using automation, attackers can execute social engineering and phishing campaigns at a scale and speed that would be impossible for human actors, significantly increasing the effectiveness and potential damage of those attacks.

5 Generative AI social engineering framework

The Generative AI Social Engineering (GenAI-SE) Framework (see Fig. 3) is a flexible and multi-dimensional model designed to investigate the complex nature of AI-driven social engineering (SE) attacks. On the one hand, the framework serves as an adaptable threat intelligence template that accommodates a diverse set of layers or dimensions for analysis and helps to evaluate the impact of AI capabilities on existing and emerging threats. Hence, the framework can be used to inform decision-makers, security professionals, and automated security systems about risks and vulnerabilities so that proactive measures can be taken to improve security and mitigate risks. On the other hand, it can be used for further research to investigate different aspects of AI-driven social engineering attacks.

While the GenAI-SE Framework provides a high-level structure for analyzing AI-driven SE attacks, applying it to specific real-world scenarios or future research may require adaptation and refinement. Potential challenges include the need for contextual adaptation to fit specific industries or organizational contexts due to varying threat landscapes and regulatory requirements. Additionally, rapid advancements in AI and cybersecurity could necessitate updates to keep the framework relevant. The framework can be extended to investigate specific domains, as shown in Fig. 3, which may require making the framework's categories more granular or introducing metrics (see Sect. 5.1). Using the framework may also require

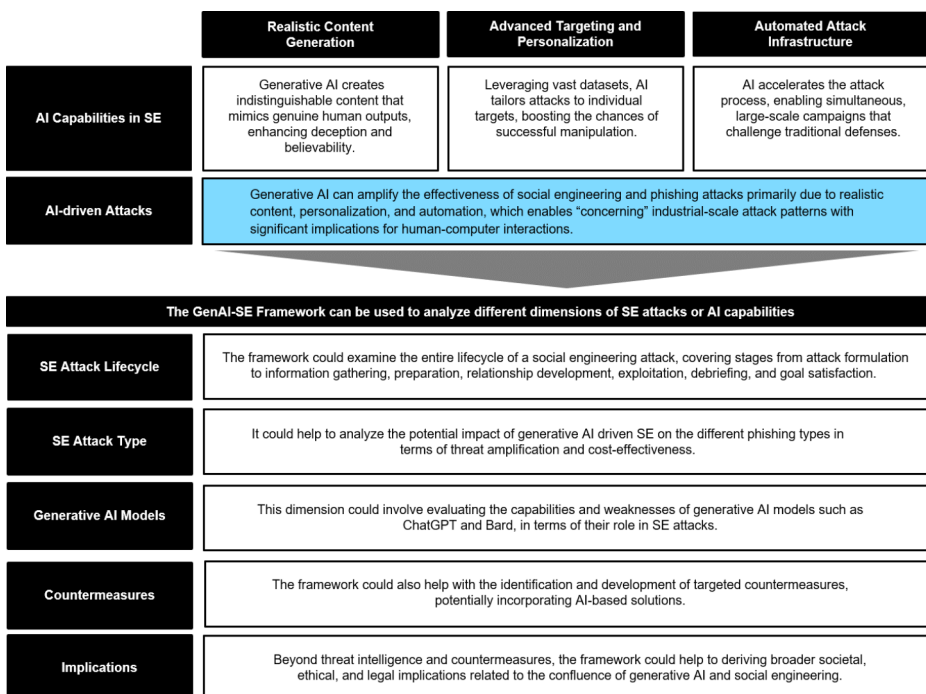


Fig. 3 Generative AI social engineering (GenAI-SE) framework: this figure illustrates the dynamic nature of the framework and its potential use in analyzing different dimensions of SE attacks and AI capabilities. The framework can be extended to investigate specific domains shown in the figure. Thus, it serves as a threat intelligence template for security professionals and researchers and supports further theoretical and empirical analysis of AI-driven SE attacks

interdisciplinary collaboration, which can add complexity but is essential for achieving a holistic threat response and comprehensive analysis that merges technological and policy solutions.

5.1 Threats amplification and cost-effectiveness

The developed framework helps to evaluate a priori the implications of AI capabilities on existing and emerging threats. Figure 4 shows an analysis of the potential of generative AI to amplify phishing threats while at the same time reducing their costs. The figure uses two metrics to evaluate the impact of generative AI on phishing attacks:

- **Threat Amplification:** Refers to the increase in the impact and effectiveness of phishing attacks due to the use of generative AI. This can be measured by the success rate of phishing attempts due to innovation, the scale (number of victims), or the execution speed.
- **Cost-Effectiveness:** Refers to the reduction in the cost of executing phishing attacks while maintaining or increasing their effectiveness. This can be measured by the resources (time, money, effort) required to conduct the attacks.

Advancements in AI have a significant impact on the context of phishing attacks, favoring the attacker rather than the defender. This dual nature of AI presents important implications for the security landscape. Traditional approaches that rely on training end-users to identify deception techniques are insufficient, particularly in the case of targeted attacks

	Realistic Content Generation		Advanced Targeting and Personalization		Automated Attack Infrastructure	
AI Capabilities in SE	Generative AI creates indistinguishable content that mimics genuine human outputs, enhancing deception and believability.		Leveraging vast datasets, AI tailors attacks to individual targets, boosting the chances of successful manipulation.		AI accelerates the attack process, enabling simultaneous, large-scale campaigns that challenge traditional defenses.	
Evaluation Metric	Threat Amplification	Cost Effectiveness	Threat Amplification	Cost Effectiveness	Threat Amplification	Cost Effectiveness
Spam Phishing	+++	+++	++	+	+++	+++
Spear Phishing	++	++	+++	+++	++	++
Whaling	++	++	+++	+++	+	+
Interpretation	Generative AI can amplify the effectiveness of social engineering and phishing attacks primarily due to realistic content, personalization, and automation, which enables "concerning" industrial-scale attack patterns with significant implications for human-computer interactions.					

Fig. 4 Evaluation of threat amplification and cost-effectiveness of Gen-AI in SE: this example use-case demonstrates the application of the proposed framework by analyzing how generative AI can amplify different types of phishing threats, such as spam phishing, spear phishing, and whaling, while reducing their costs. Metrics include Threat Amplification (increase in impact and effectiveness) and Cost-Effectiveness (reduction in execution cost)

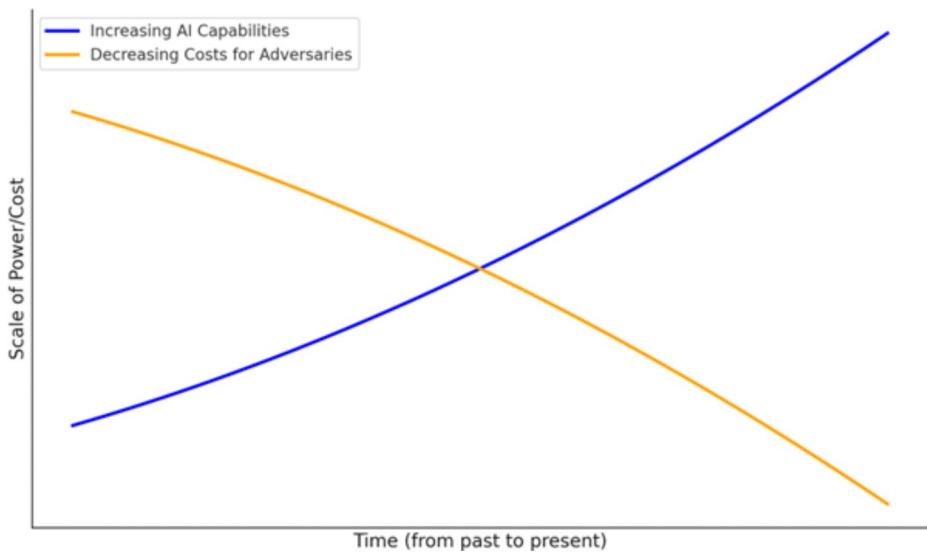


Fig. 5 Relationship between AI capabilities and costs: this figure shows the increase in AI capabilities (blue line) and the decrease in costs for adversaries (orange line) over time. As AI technology advances and becomes more affordable, the potential for its misuse in social engineering attacks grows, emphasizing the need for robust countermeasures

such as spear phishing and whaling. Humans often struggle to detect and defend against these sophisticated techniques, leaving organizations vulnerable to highly personalized and convincing phishing campaigns. The use of sophisticated AI holds the potential to further improve tailored attacks for the mass market. Attackers can leverage AI algorithms to analyze large datasets and create highly targeted phishing campaigns at scale. This automation and personalization significantly increase the effectiveness of spear phishing campaigns, posing a greater threat to individuals and organizations. The implications of AI-powered phishing attacks are further amplified by the increasingly sophisticated underground ecosystem of cybercriminals, the emergence of the gig economy, and the widespread adoption of Software as a Service (SaaS) solutions. This underground ecosystem thrives on the exchange of tools, techniques, and resources, providing attackers with fertile ground to leverage AI for their malicious activities. The collaboration and sharing of AI-powered attack methods within this ecosystem poses a significant challenge for defenders. Additionally, the gig economy introduces new complexities by offering accessible platforms that facilitate the utilization of AI technologies, tools, and services. This ease of access empowers even non-expert individuals to employ AI for nefarious purposes, further contributing to the prevalence of AI-enabled phishing attacks. Moreover, the adoption of SaaS solutions provides attackers with readily available AI-powered tools and infrastructure to launch sophisticated phishing campaigns.

Furthermore, the continuing trend of AI becoming more cost-effective is a double-edged sword in this context. While these financial efficiencies benefit many industries by making AI more accessible and democratized, they simultaneously lower the barrier to entry for malicious actors. As the cost of acquiring and deploying AI capabilities decreases (see Fig. 5), a broader array of adversaries, even those with limited resources, can tap into sophisticated

AI tools for malicious purposes. This means that the scale and frequency of AI-driven phishing attacks could surge, making it even more challenging for organizations to maintain an effective defense. The interplay of accessibility and affordability of AI technologies may inadvertently intensify the phishing threat landscape, underscoring the urgency for novel defensive measures.

5.2 Countermeasures for AI-generated deception

The nature of cybercrime as a continuously evolving threat—tug of war between adversaries and defenders—makes it impossible to determine guaranteed protection measures (Taddeo et al. 2019). On the highest level we can distinguish between technical countermeasures, which are largely based on AI/ML (Basit et al. 2021; Schmitt 2023a) and cryptography, and human-centric countermeasures, which is largely user education (Naqvi et al. 2023). While ideas for innovative solutions exist (Jansen and Fischbach 2020; Ganesh et al. 2023; Glas et al. 2023), the cybersecurity industry is focusing predominantly on awareness training programs to tackle the so-called “people” problem. These programs are commonplace and applied by all types of businesses ranging from small to global corporations. Yet, they are generally ineffective against advanced targeting and personalization tactics. Labeling the end-user as the problem is not just an oversimplification, but a dangerous perspective. The implication is that users should be blamed for not being smart/aware enough and remediate this through security awareness, education and training (SAET), however this approach is insufficient and flawed for a number of reasons. (1) In the face of AI improvements in mimicking human interactions, it is not feasible for this increasing capability to be mitigated through greater awareness or understanding. We are fast approaching a point where people will not be able to differentiate genuine from fabricated content. (2) The potential scale and automation of malicious attacks will mean that individuals will need constant and near-perfect levels of awareness to reach a satisfactory balance between protecting themselves from deception while maintaining trust in communication channels. Finally, (3) it is already difficult to evaluate the effectiveness of SAET: if someone fails a phishing test or becomes a victim of social engineering, this is usually interpreted as a sign that they weren’t trained enough, or paying enough attention, which is then used as justification for more awareness, training and education. This is circular reasoning, which buys into the harmful fallacy that victims are to be blamed for falling prey to attackers (Fig. 6).

The implications of AI in the context of phishing highlight the need for innovative defenses and proactive strategies to mitigate the risks posed by AI-enabled attacks. Table 6 contains a number of indicative proposals for specific countermeasures for each of the three pillars.

One example for deception identification in pillar one is detecting deepfakes. According to Kaur et al. (Gambin et al. 2024), most methods for detecting deepfake videos are data-driven. Deep learning methods, particularly Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), dominate this field by identifying subtle inconsistencies in fake videos. However, these methods face challenges, including data availability, computational complexity, generalization issues, and reliability. As deepfake generation methods improve, their ability to evade detection necessitates constant advancements in detection techniques.

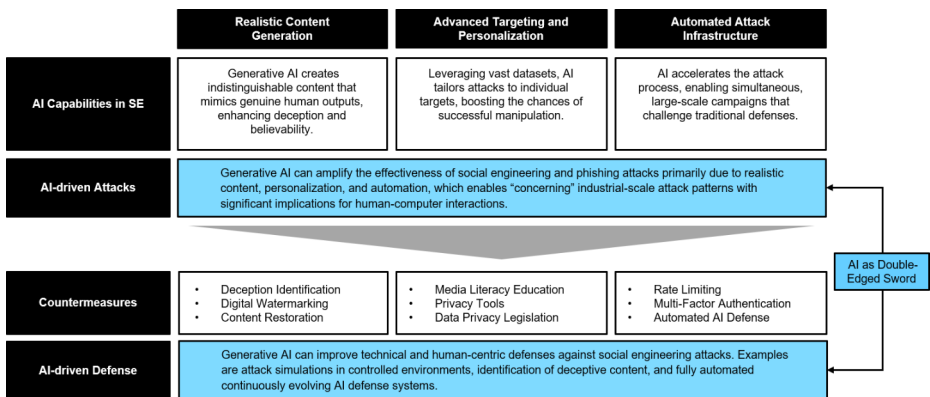


Fig. 6 Derivation of countermeasures for ai-generated deception: this use-case demonstrates the framework's application by identifying countermeasures for AI-generated deception across the three pillars

Table 6 Countermeasures for AI-generated deception

Category	Countermeasure	Explanation
Realistic content creation	Deception identification	Develop AI-based solutions that can recognize and flag AI-generated content. For instance, tools that detect deepfakes by analyzing inconsistencies in videos (Ke and Wang 2023; Kaur et al. 2024).
	Digital watermarking	Implement watermarking techniques for genuine content, signaling to users the authenticity of a digital item (Malanowska et al. 2024).
	Content restoration	When content undergoes malicious AI modification, generative AI models can attempt to revert or rectify the content, such as reversing deepfake alterations (Lehtinen et al. 2018).
Advanced targeting and personalization	Media literacy education	Educate the public about the capabilities of generative AI to put them in a position to recognize potential AI-generated content (Naqvi et al. 2023).
	Privacy tools	Encourage the use of privacy tools that block tracking cookies, scripts, and other methods used to gather personal data for targeted attacks (Nowrozy et al. 2024).
	Data privacy legislation	Strengthen regulations around data collection, storage, and usage to limit the extent of personal data available for AI targeting (Tikkinen-Piri et al. 2018).
Automated attack infrastructure	Rate limiting	Implement rate limiting for message sending or content posting on platforms, making mass phishing attempts via automation more difficult.
	Multi-factor authentication	Promote the use of MFA, making it harder for automated systems to gain unauthorized access, even if they have some user credentials (Chen et al. 2024).
	Automated AI defense	By using AI-enabled defenses such as machine learning for large-scale threat detection, behavioral analysis, and simulated attacks, systems can be automatically fortified against threats (Basit et al. 2021; Schmitt 2023a).

It is important to note that while generative AI can be used maliciously in these scenarios, it can also be employed for beneficial purposes, such as creating training materials for cybersecurity professionals or developing tools to detect and counter AI-generated deceptive content. Overall, the relentless evolution of AI makes it paramount to stay updated on the latest advancements in AI-driven deception and its countermeasures. In mapping AI capabilities to the detail of social engineering attacks, we demonstrate how the framework can be used to highlight new areas of development and research in disrupting AI powered social engineering attacks.

Even seasoned cybersecurity experts are susceptible to sophisticated AI-generated phishing attacks. The expectation that individuals can scrutinize every email they receive is impractical and counterproductive, especially in today's fast-paced digital environment, that comes with heavy cognitive workload (Kosch et al. 2023). Also, this mindset ignores essential psychological factors explained by researchers like Daniel Kahneman (2013). Routine tasks, like checking emails, seldom involve critical thinking, and when emotions—especially fear—are triggered, rational decision-making becomes challenging. Current solutions like anomaly and spam detection AI are somewhat effective, but they fall short against targeted spear phishing and state-sponsored attacks.

6 Directions for future research

The three identified capabilities of generative AI in the context of social engineering—realistic content, personalization, and automation—create a problematic cocktail that marks the onset of potentially unforeseeable innovations in hacking. It is easy to imagine the development of powerful and autonomous social engineering bots, making the enhancement of existing and development of novel countermeasures for these emerging threats paramount (Ahmad et al. 2023). The proposed GenAI-SE framework can be used as a template to structure further research and serve as a baseline for empirical investigations into AI threats and countermeasures. It functions as a threat intelligence template for security professionals and researchers, supporting both theoretical and empirical analysis of AI-driven SE attacks. We recommend the following avenues for future research:

- **User Awareness and Education:** We need to recognize the fact that “awareness training” is not an easy fix against AI-powered cyber-attack. Nevertheless, increasing user knowledge about the latest techniques and strategies can empower individuals to recognize and respond appropriately to potential threats (Distler 2023; Marin et al. 2023). This could involve designing effective training programs, interactive simulations, and user-friendly educational materials (Jansen and Fischbach 2020; Glas et al. 2023).
- **Adversarial Machine Learning:** Develop advanced techniques to detect and defend against adversarial attacks in the context of AI-powered social engineering (Ahmad et al. 2023). Adversarial machine learning aims to create robust models that can withstand manipulation attempts by attackers. Research in this area can focus on developing algorithms and strategies to detect and mitigate adversarial SE and phishing techniques.
- **Active Deception Defense:** Develop proactive defense mechanisms that actively disrupt deceptive attacks in real-time. This could involve technologies like natural language processing, anomaly detection, and real-time analysis of communication channels to

identify and block SE and phishing attempts as they occur. It is also essential to advance deepfake detection research (Kaur et al. 2024).

- **Explainable AI for Threat Detection:** Enhance the transparency of AI models used for threat detection. Explainable AI techniques can help security analysts and users understand how AI systems make decisions and identify indicators of deceptive attacks. By providing meaningful explanations and insights, it becomes easier to trust and validate the outputs of AI-powered threat detection systems (Kim et al. 2023).
- **Emerging Technologies:** Future research could explore the intersection of emerging technologies and their implications for social engineering attacks. As technologies, such as Chatbots (Rapp et al. 2023; Zhang et al. 2024), brain-computer interfaces (Tang et al. 2023), robotics (Dörrenbächer et al. 2023; Yamamura et al. 2023), quantum computing (Cavaliere et al. 2020), and metaverse ecosystems (Vondrek et al. 2022; Schmitt 2023c; Soliman et al. 2024) become more integrated into our daily life, new avenues for social engineering attacks are likely to emerge.
- **Cybersecurity Threats:** Gen AI has not only the potential to amplify or create new social engineering and phishing attacks but also to enhance all cybersecurity threats. Future research could explore its impact on ransomware (Teichmann 2023), insider threats, man-in-the-middle attacks, and the use of deepfakes (Gambín et al. 2024; Kaur et al. 2024).

7 Conclusion

The potential for persistent high-volume social engineering (SE) attacks poses a serious threat to the trust that users place in human-computer interactions (HCI). When these interactions are compromised through continual SE attacks, whether driven by generative AI or other methods, it erodes the foundation of trust and security that is essential for effective HCI. Users may become increasingly skeptical of the systems they interact with, questioning the authenticity of dialogues and hesitating to share information or perform tasks. This paper analyzed the transformative potential of Generative AI in the realm of social engineering and phishing. We argued that the integration of Generative AI can profoundly amplify the effectiveness of these cyberattacks. The primary drivers behind this enhancement are the capabilities of AI to produce realistic content, tailor-made personalization, and unparalleled automation. This potent combination paves the way for industrial-scale attack patterns, signaling a paradigm shift to a new, more sophisticated era of cyber threats. Further compounding this challenge is the evolving landscape of AI. As artificial intelligence continues its trajectory of becoming simultaneously more powerful and economically accessible, the barriers to misuse diminish. This raises an alarm for not only the cybersecurity community but also industries and individuals at large. In the face of such emerging threats, there is an acute and pressing need to develop and refine countermeasures. A proactive stance today is our best defense against the impending wave of AI-powered cyber threats of tomorrow.

Author contributions All authors contributed equally to this work. Each author participated in the conception and design of the study, data collection and analysis, drafting and revising the article critically for important intellectual content, and approval of the final version to be published.

Data availability No datasets were generated or analysed during the current study.

Declarations

Competing interests The authors declare no competing interests.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Ahmad R, Alsmadi I, Alhamdani W, Tawalbeh L (2023) Zero-day attack detection: a systematic literature review. *Artif Intell Rev* 56:10733–10811. <https://doi.org/10.1007/s10462-023-10437-z>
- Aleroud A, Zhou L (2017) Phishing environments, techniques, and countermeasures: a survey. *Comput Secur* 68:160–196
- Anil R, Dai AM, Firat O et al (2023) PaLM 2 technical report
- Archana R, Jeevaraj PSE (2024) Deep learning models for digital image processing: a review. *Artif Intell Rev* 57:11. <https://doi.org/10.1007/s10462-023-10631-z>
- Basit A, Zafar M, Liu X et al (2021) A comprehensive survey of AI-enabled phishing attacks detection techniques. *Telecommun Syst* 76:139–154. <https://doi.org/10.1007/s11235-020-00733-2>
- Bécue A, Praça I, Gama J (2021) Artificial intelligence, cyber-threats and industry 4.0: challenges and opportunities. *Artif Intell Rev* 54:3849–3886. <https://doi.org/10.1007/s10462-020-09942-2>
- Bray SD, Johnson SD, Kleinberg B (2023) Testing human ability to detect ‘deepfake’ images of human faces. *J Cybersecur*. <https://doi.org/10.1093/cybsec/tyad011>
- Briggs J, Kodnani D (2023) The potentially large effects of artificial intelligence on economic growth
- Cavaliere F, Mattsson J, Smeets B (2020) The security implications of quantum cryptography and quantum computing. *Netw Secur* 2020:9–15. [https://doi.org/10.1016/S1353-4858\(20\)30105-7](https://doi.org/10.1016/S1353-4858(20)30105-7)
- Chen X, Wang B, Li H (2024) A privacy-preserving multi-factor authentication scheme for cloud-assisted IoMT with post-quantum security. *J Inform Secur Appl* 81:103708. <https://doi.org/10.1016/j.jisa.2024.103708>
- Chui M, Hazan E, Roberts R et al (2023) The economic potential of generative AI. The next productivity frontier
- Cutillo CM, Sharma KR, Foschini L et al (2020) Machine intelligence in healthcare—perspectives on trustworthiness, explainability, usability, and transparency. *NPJ Digit Med* 3
- Daniel K (2013) Thinking, fast and slow. Farrar, Straus and Giroux
- Davenport TH (2018) From analytics to artificial intelligence. *J Bus Anal* 1:73–80. <https://doi.org/10.1080/2573234x.2018.1543535>
- Desolda G, Ferro LS, Marrella A et al (2022) Human factors in phishing attacks: a systematic literature review. *ACM Comput Surv*. <https://doi.org/10.1145/3469886>
- Distler V (2023) The influence of context on response to spear-phishing attacks: an in-situ deception study. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Dörrenbächer J, Ringfort-Felner R, Hassenzahl M (2023) The intricacies of social robots: secondary analysis of fictional documentaries to explore the benefits and challenges of robots in complex social settings. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Fahl S (2021) Web & mobile security knowledge area. The Cyber Security Body of Knowledge
- Fleming P, Edwards SG, Bayliss AP, Seger CR (2023) Tell me more, tell me more: repeated personal data requests increase disclosure. *J Cybersecur*. <https://doi.org/10.1093/cybsec/tyad005>
- Gambin AF, Yazidi A, Vasilakos A et al (2024) Deepfakes: current and future trends. *Artif Intell Rev* 57:64. <https://doi.org/10.1007/s10462-023-10679-x>

- Ganesh A, Ndulue C, Orji R (2023) Tailoring a persuasive game to promote secure smartphone behaviour. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Glas M, Vielberth M, Pernul G (2023) Train as you fight: evaluating authentic cybersecurity training in cyber ranges. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Goodfellow IJ, Pouget-Abadie J, Mirza M et al (2014) Generative adversarial networks
- Greshake K, Abdelnabi S, Mishra S et al (2023) Not what you've signed up for. Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection
- Gupta A, Akiri C, Aryal K et al (2023) From ChatGPT to ThreatGPT: impact of generative AI in cybersecurity and privacy. *IEEE Access* 11:80218–80245. <https://doi.org/10.1109/ACCESS.2023.3300381>
- Jansen P, Fischbach F (2020) The social engineer: an immersive virtual reality educational game to raise social engineering awareness. In: CHI PLAY 2020 - Extended abstracts of the 2020 annual symposium on computer-human interaction in play. Association for Computing Machinery, Inc, pp 59–63
- Kaur R, Gabrijelčič D, Klobučar T (2023) Artificial intelligence for cybersecurity: literature review and future research directions. *Inform Fusion*. <https://doi.org/10.1016/j.inffus.2023.101804>
- Kaur A, Noori Hoshyar A, Saikrishna V et al (2024) Deepfake video detection: challenges and opportunities. *Artif Intell Rev* 57:159. <https://doi.org/10.1007/s10462-024-10810-6>
- Ke J, Wang L (2023) DF-UDetector: an effective method towards robust deepfake detection via feature restoration. *Neural Netw* 160:216–226. <https://doi.org/10.1016/j.neunet.2023.01.001>
- Kim SSY, Watkins EA, Russakovsky O et al (2023) Help me help the ai: understanding how explainability can support human-AI interaction. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Kosch T, Karolus J, Zagermann J et al (2023) A survey on measuring cognitive workload in human-computer interaction. *ACM Comput Surv*. <https://doi.org/10.1145/3582272>
- LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521:436–444. <https://doi.org/10.1038/nature14539>
- Lehtinen J, Munkberg J, Hasselgren J et al (2018) Noise2Noise. *Learning Image Restoration without Clean Data*
- Malanowska A, Mazurczyk W, Araghi TK et al (2024) Digital watermarking—a meta-survey and techniques for fake news detection. *IEEE Access* 12:36311–36345. <https://doi.org/10.1109/ACCESS.2024.3374201>
- Malatji M, Tolah A (2024) Artificial intelligence (AI) cybersecurity dimensions: a comprehensive framework for understanding adversarial and offensive AI. *AI Ethics*. <https://doi.org/10.1007/s43681-024-00427-4>
- Marin IA, Burda P, Zannone N, Allodi L (2023) The influence of human factors on the intention to report phishing emails. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Mouton F, Malan MM, Leenen L, Venter HS (2014) Social engineering attack framework. In: 2014 information security for South Africa - proceedings of the ISSA 2014 conference. Institute of Electrical and Electronics Engineers Inc
- Naqvi B, Perova K, Farooq A et al (2023) Mitigation strategies against the phishing attacks: a systematic literature review. *Comput Secur* 132
- Neupane S, Fernandez IA, Mittal S, Rahimi S (2023) Impacts and risk of generative AI technology on cyber defense
- Nowrozy R, Ahmed K, Kayes ASM et al (2024) Privacy preservation of electronic health records in the modern era: a systematic survey. *ACM Comput Surv* 56:1–37. <https://doi.org/10.1145/3653297>
- OpenAI (2023) GPT-4 Technical Report
- Rapp A, Boldi A, Curti L et al (2023) Collaborating with a text-based chatbot: an exploration of real-world collaboration strategies enacted during human-chatbot interactions. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Reed S, Zolna K, Parisotto E et al (2022) A Generalist Agent. 2022–2027
- Renaud K, Warkentin M, Westerman G (2023) From ChatGPT to HackGPT. Meeting the Cybersecurity Threat of Generative AI
- Schmitt M (2023a) Securing the digital world: protecting smart infrastructures and digital industries with artificial intelligence (AI)-enabled malware and intrusion detection. *J Ind Inf Integr* 36:100520. <https://doi.org/10.1016/j.jii.2023.100520>
- Schmitt M (2023b) Automated machine learning: AI-driven decision making in business analytics. *Intell Syst Appl* 18:200188. <https://doi.org/10.1016/j.iswa.2023.200188>
- Schmitt M (2023c) Metaverse: implications for business, politics, and society. *SSRN Electron J*. <https://doi.org/10.2139/ssrn.4168458>
- Seymour J, Tully P (2018) Generative models for Spear Phishing. *Posts on Social Media*

- Soliman MM, Ahmed E, Darwish A, Hassanien AE (2024) Artificial intelligence powered metaverse: analysis, challenges and future perspectives. *Artif Intell Rev* 57:36. <https://doi.org/10.1007/s10462-023-10641-x>
- Sutton RS, Barto AG (1998) Reinforcement learning: an introduction. *IEEE Trans Neural Netw* 9:1054–1054. <https://doi.org/10.1109/TNN.1998.712192>
- Taddeo M, McCutcheon T, Floridi L (2019) Trusting artificial intelligence in cybersecurity is a double-edged sword. *Nat Mach Intell* 1:557–560. <https://doi.org/10.1038/s42256-019-0109-1>
- Tang X, Shen H, Zhao S et al (2023) Flexible brain–computer interfaces. *Nat Electron* 6:109–118. <https://doi.org/10.1038/s41928-022-00913-9>
- Teichmann F (2023) Ransomware attacks in the context of generative artificial intelligence—an experimental study. *Int Cybersecur Law Rev*. <https://doi.org/10.1365/s43439-023-00094-x>
- Tikkinen-Piri C, Rohunen A, Markkula J (2018) EU general data protection regulation: changes and implications for personal data collecting companies. *Comput Law Secur Rev* 34:134–153. <https://doi.org/10.1016/j.clsr.2017.05.015>
- Vaswani A, Shazeer N, Parmar N et al (2017) Attention Is All You Need
- Vondrek M, Baggili I, Casey P, Mekni M (2022) Rise of the metaverse’s immersive virtual reality malware and the man-in-the-room attack & defenses. *Comput Secur*. <https://doi.org/10.1016/j.cose.2022.102923>
- Webb T, Holyoak KJ, Lu H (2023) Emergent analogical reasoning in large language models. *Nat Hum Behav*. <https://doi.org/10.1038/s41562-023-01659-w>
- World Economic Forum (2024) Global Cybersecurity Outlook 2024
- Yamamura N, Uriu D, Muramatsu M et al (2023) Social digital cyborgs: the collaborative design process of JIZAI ARMS. In: Conference on human factors in computing systems - proceedings. Association for Computing Machinery
- Yamin MM, Ullah M, Ullah H, Katt B (2021) Weaponized AI for cyber attacks. *J Inform Secur Appl* 57:102722. <https://doi.org/10.1016/j.jisa.2020.102722>
- Zhang Y, Lau RYK, David Xu J et al (2024) Business chatbots with deep learning technologies: state-of-the-art, taxonomies, and future research directions. *Artif Intell Rev* 57:113. <https://doi.org/10.1007/s10462-024-10744-z>
- Zou A, Wang Z, Kolter JZ, Fredrikson M (2023) Universal and transferable adversarial. attacks on aligned language models
- Zscaler (2024) Zscaler ThreatLabz 2024 Phishing Report

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.