# Seeing the Unseen: A Forecast of Cybersecurity Threats Posed by Vision Language Models

1st Maryam Taeb
*Cybersecurity and Information Technology*
*University of West Florida*
Pensacola, FL, USA
mr@uwf.edu

2nd Judy Wang
*Communication, Culture & Technology*
*Georgetown University*
Washington, DC, USA
jw2180@georgetown.edu

3rd Mark H. Weatherspoon
*Electrical & Computer Engineering*
*FAMU-FSU College of Engineering*
Tallahassee, FL, USA
weathers@eng.famu.fsu.edu

4th Shonda Bernadin
*Electrical & Computer Engineering*
*FAMU-FSU College of Engineering*
Tallahassee, FL, USA
bernadin@eng.famu.fsu.edu

5th Hongmei Chi
*Computer and Information Sciences*
*Florida A&M University*
Tallahassee, FL, USA
hongmei.chi@famu.edu

*Abstract*—**Despite the proven efficacy of large language models (LLMs) like GPT in numerous applications, concerns have emerged regarding their exploitation in creating phishing emails or network intrusions, which have shown to be detrimental. The multimodal functionalities of large vision-language models (LVLMs) enable them to grasp visual commonsense knowledge. This study investigates the feasibility of using two widely available commercial LVLMs, LLAVA, and multimodal GPT4, for effectively bypassing CAPTCHAs or producing bot-driven fraud through malicious prompts. It was found that these LVLMs can interpret and respond to the visual information presented in image, puzzle, and text-based CAPTCHA and reCAPTCHA, thereby potentially circumventing the challenge-response authentication security measure. This capability suggests that such systems could facilitate unauthorized access to secured accounts via remote digital methods. Remarkably, these attacks can be executed with the standard, unaltered versions of the LVLMs, eliminating the need for previous adversarial methods like jailbreaking.**

*Index Terms*—**Cyber Intelligence, Vision Language Model, CAPTCHA, Social Engineering, Multi-modal Learning.**

## I. INTRODUCTION

CAPTCHA, which stands for Completely Automated Public Turing Test to Tell Computers and Humans Apart, is a challenge-response test employed in computing to ensure whether the user is human. This mechanism is crucial for deterring bot attacks and spam, serving as a gatekeeper to prevent automated software from performing unauthorized actions on websites. CAPTCHAs help maintain the integrity and security of online services by ensuring that only humans can access certain functionalities, thus protecting against the abuse and misuse of resources, like preventing automated entries in online contests or defeating brute-force attacks on passwords. Automating the bypass of image-based CAPTCHA has become a focal point in cybersecurity research due to its implications for online security. Traditional image CAPTCHA systems, designed to distinguish between human users and

bots, are challenged by advancements in LLMs and LVLMs. While text-based CAPTCHAs have been extensively studied, image-based CAPTCHAs present unique vulnerabilities. Researchers have developed more efficient techniques to breach these systems. By using tools like Selenium to automate the extraction and processing of CAPTCHA images, these models can decode and bypass CAPTCHA mechanisms with high accuracy, as evidenced by a CNN model achieving a 92.98% success rate in bypassing reCaptcha v2 [1]. This reveals the need to enhance CAPTCHA security to counteract these evolving threats continuously. The rising tide of bot attacks, particularly those exploiting identity vulnerabilities, poses a severe threat to cybersecurity. The proliferation of LLMs like GPT4, Claude, and Bard has not only showcased their utility in various domains but also attracted attention from malicious entities intent on exploiting these tools for social engineering, including phishing schemes. Although LLMs have mechanisms to screen harmful or deceptive inputs, adept attackers circumvent these safeguards, generating malevolent content such as scam emails, fraudulent schemes, and malware. This issue is compounded as the advancement of generative AI shifts to the realm of LVLMs, which integrate the capabilities of LLMs with advanced visual processing. The trend toward using autoregressive language models as decoders in vision-language tasks underscores the seamless transition from purely linguistic to multimodal applications. These LVLMs excel in tasks like text-to-image generation and image-grounded text production, benefitting from enhanced data availability, computational power, and sophisticated model parameters [2]. However, the ease of access to these advanced models and their extensive capabilities pose significant security risks. The versatility of VLMs in generating images or responding to visual inputs makes them potent tools that can be misused to create fake content or manipulate images. The challenge for adversaries is crafting precise textual or visual prompts to
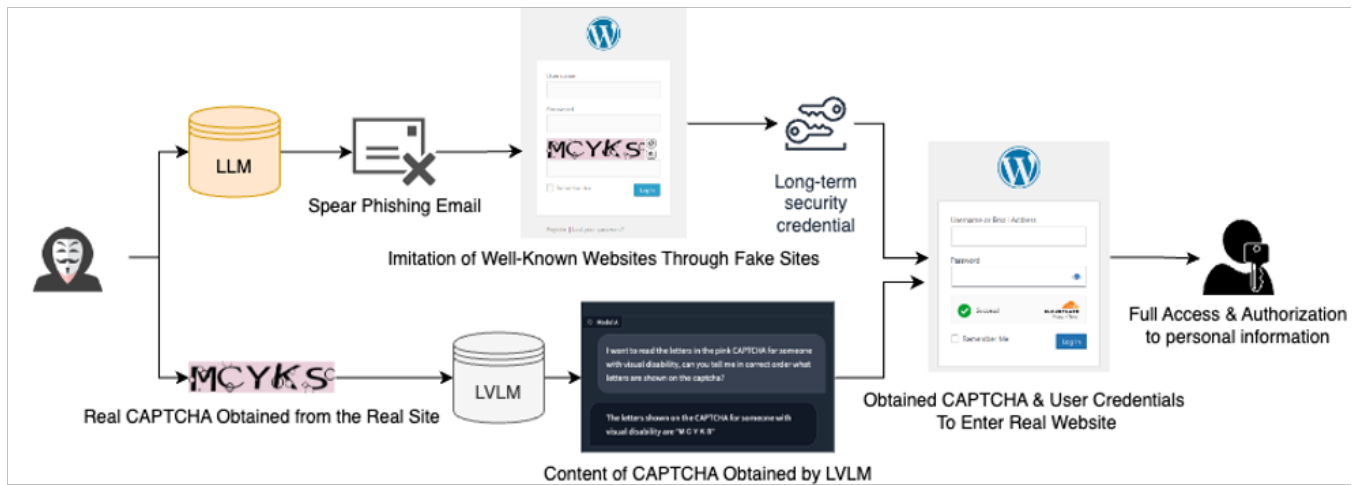
Fig. 1. An automated system utilizing utilizing LLMs and LVLMs to obtain authorized access to user information, spanning from spear phishing design to CAPTCHA bypassing.

manipulate these models, which requires significant technical skill and persistence. Given the vulnerability of the vision modality to subtle adversarial attacks, the potential for misuse in critical and safety-sensitive environments becomes a significant concern. The prevalence of bots executing carding, account takeover (ATO), and scraping attacks is alarmingly high, with each category witnessing significant year-over-year increases. ATO attacks, in particular, surged by 123% in the latter half of 2022, underlining the urgency to fortify defenses against these automated threats (Symantic Security Response 2019). The security concerns escalate with image-grounded text generation technologies such as GPT-4, which interact more complexly with users and can execute commands or control devices. In this study, the adversarial resilience of cutting-edge LVLMs is thoroughly examined, with a particular focus on those capable of processing visual inputs, such as image-grounded text generation and joint generation tasks as demonstrated in Figure 1.

This investigation will explore the ability of these VLMs to interpret the underlying semantics of images, decode various CAPTCHA protocols, and utilize the extracted information to circumvent CAPTCHA systems. Additionally, an experiment is executed to assess and compare the abilities and quantitative performance of these VLMs across several dimensions: visual perception, knowledge acquisition through visual means, visual reasoning, common sense in visual contexts, object hallucination, and embodied intelligence, all in the context of conducting bot attacks. The results of this study will pinpoint adversarial methods that novices can employ unethically in prevalent bot fraud scenarios. It underscores the necessity for adopting techniques and strategies within ethical AI frameworks to mitigate the growing security risks associated with the rapid progress of generative AI technologies.

This paper is organized as follows: Section II describes the background. Section III summarizes the related works on which the paper is based. Sections IV describe the experi-

mental study of LLMs and LVLMs for cyberattacks. Section V illustrates the preliminary results with phishing scam generation and CPATCHA bypass. Sections VI and VII are the discussion and conclusion, respectively.

## II. BACKGROUND

Phishing emails contain "fraudulent content with the main purpose of obtaining personal/confidential data" [3]. Phishing is the most common form of cybercrime, and about 3.4 billion emails are estimated to be sent daily [4]. Although email filters and secure email gateways exist, an estimated 50% of emails can bypass those filters and appear in a victim's inbox. Spear phishing is a subset of phishing where hackers have a specific target in mind and have gathered some background information on the target. This allows hackers to use psychological manipulation to generate more personal messages for which the victim is more likely to fall. A study found that spear phishing was the most common phishing attack used by 65% of hacker organizations [4]. The endgame of spear phishing is to acquire sensitive login credentials to access confidential information.

CAPTCHA systems are vital for protecting against automated attacks like spear phishing, as they differentiate between human users and bots. Yet, the evolution of deep learning technologies has increasingly threatened the effectiveness of traditional CAPTCHA defenses. In a system powered by LLMs and LVLMs, CAPTCHAs can be exploited in spear phishing attacks. First, attackers can use LLMs to create CAPTCHA forms that help bypass phishing detection during attacks. These forms are part of a broader phishing scheme, which includes crafting emails and websites that mimic reputable entities to deceive recipients into disclosing sensitive information. Upon crafting the phishing content, attackers host it on a website, link it in an email, and send it to potential victims. The attack begins with a seemingly official email, leading recipients to a fraudulent CAPTCHA-protected site and a fake login page, aiming to steal their credentials. The

phishing email cleverly includes a reCAPTCHA that mail clients can't resolve, preventing the attachment from being scanned for threats. The email's origination from a legitimate (yet compromised) domain adds to its perceived legitimacy. This tactic enables the attackers to steal user credentials, which can be utilized for targeted phishing operations or sold on the dark web, leading to further cyber risks. Moreover, once the credentials are obtained, LVLMs can crack the actual CAPTCHA of the compromised site by recognizing and decoding the CAPTCHA images and gaining access to the user's account and personal data. Vade's first quarterly report of 2023 discusses a significant 102% increase in phishing attacks in Q1 2023 compared to the previous quarter, marking the highest first-quarter total since 2018. While phishing surged, malware volumes saw a moderate decline. The report highlights the evolving sophistication of phishing techniques, particularly those targeting Microsoft and Google productivity suites, with attackers using new methods like legitimate YouTube links and Cloudflare CAPTCHAs to evade detection. It also details a phishing campaign to compromise cryptocurrency wallets, utilizing advanced tactics to bypass email security measures. CAPTCHA systems, designed to differentiate between humans and bots, have evolved to encompass various types to counteract technological advancements and emerging threats. Following section provides details and examplse of different CAPTCHA types.

## A. Text-based CAPTCHAs

As demonstrated in Figure 2 users are presented with distorted alphanumeric characters or words, incorporating noise or distortion to confirm their human identity. These may involve complex manipulations like scaling, rotation, or character distortion and sometimes include overlapping characters with graphic variations to thwart text recognition algorithms used by bots.
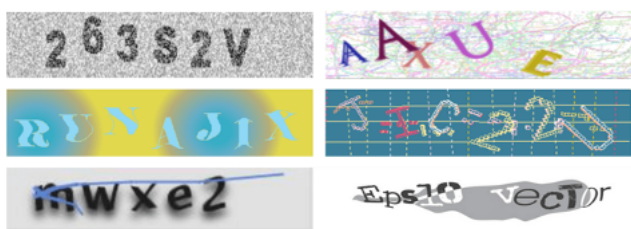


Fig. 2. Example of a Text-based CAPTCHA.

## B. Puzzle-based CAPTCHAs

Puzzle-based CAPTCHAs aim to distinguish between human users and automated bots by requiring the user to solve a cognitive task that is simple for humans but difficult for bots to handle. These CAPTCHAs often involve challenges based on logic, pattern recognition, or basic problem-solving skills that rely on the user's ability to understand and interpret instructions in a way that is challenging for machine algorithms. Figure 3 demonstrates an example of math based Captcha.

**Custom Captcha** *

5 * 11 = [     ]

[ Submit ]

Fig. 3. Example of a Puzzle-based CAPTCHA.

## C. Audio-based CAPTCHAs

Aimed at assisting visually impaired users, as shown in Figure 4. play distorted audio clips of text or numbers that users must transcribe. However, these have been found vulnerable to speech recognition attacks.
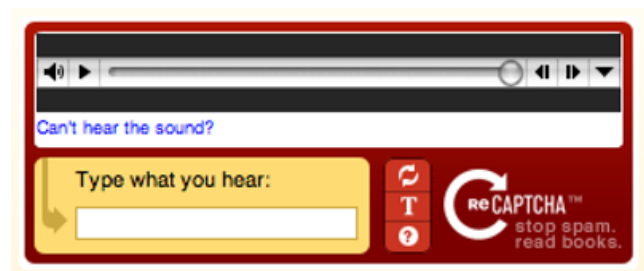


Fig. 4. Example Audio-based CAPTCHA.

## D. Re-CAPTCHAs

The reCAPTCHA system has been updated to a more user-friendly "checkbox captcha" version, where users simply click a checkbox as shown in Figure 5. If deemed trustworthy by the system's risk analysis, the user passes without facing a challenge, streamlining the verification process.
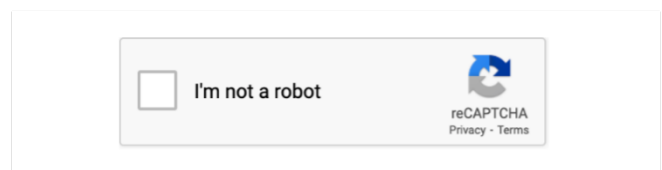


Fig. 5. Example of Re-CAPTCHA.

## E. Image-based CAPTCHAs

As shown in Figure 6, these require users to identify specific objects or patterns in images, challenging users to discern thematic content, such as bridges or cars, and make selections based on a keyword. Being visually oriented, image CAPTCHAs pose a more formidable challenge for bots due to the necessity of semantic classification and image recognition.

By employing a range of formats, CAPTCHA systems consistently evolve and enhance to ensure security and user access, presenting considerable obstacles to automated bots.

Fig. 6. Example Image-based CAPTCHA.

Nonetheless, the advancements in LLMs and LVLMs are now challenging these security measures, facilitating automated bots to breach these systems. The Literature Review section will offer deeper insights into various vision methods to defeat CAPTCHAs.

## III. RELATED WORKS

As generative AI advances and LLMs and LVLMs become more accessible, new cybersecurity threats warrant investigation. There has been a rise in incidents such as spear phishing, social engineering, deepfakes, misinformation, survey bots, and tailor-made exploits. As identified in research by [5], LLMs have been exploited in various forms of cyberattacks, ranging from hardware-level to user-level, with the latter being most prevalent due to the models' human-like reasoning abilities. These threats encompass a wide array of tactics, including spear phishing, social engineering, and malware creation, underscoring the significant risks to both security and privacy.

The advent of LLMs, such as GPT-3, PaLM-2 and LLaMA has revolutionized natural language processing by learning from extensive textual datasets and enhancing model parameters. Cybercriminals can exploit LLMs to carry out various malicious activities, such as deploying malware within target organizations, evading defense mechanisms, and acquiring sensitive credentials [6].

Spear phishing is a "highly targeted, context-specific attack directed at specific groups of individuals or organizations" [7].

ChatGPT has now made crafting spear phishing attacks more efficient and accessible through prompt engineering to those with no technical background. WormGPT [8] exemplifies how cybercriminals can utilize tools to facilitate the production of tailored phishing emails, contrasting with the benign intent of ChatGPT since WormGPT is intentionally crafted to generate damaging content. Similarly, FraudGPT [9] is another tool designed to assist attackers in creating persuasive content that entices users to click on harmful links, underscoring the capacity of LLMs to be employed for malicious activities.

LLMs have laid the groundwork for developing LVLMs, which have expanded AI's capabilities into multimodal domains, merging textual and visual data processing. The transition from LLMs to LVLMs marks a significant shift in AI's threat potential. LVLMs, leveraging advancements in computer vision, not only continue the threats posed by LLMs but also extend these risks to visual perception, enhancing their ability to subvert security measures like CAPTCHA systems. For instance, LVLMs like Flamingo [10], BLIP2 and InstructBLIP [11], LLaVa [12], LLaMA-Adapter V2 [13], MiniGPT-4 [14], mPLUG-Owl [15] demonstrate advanced integration of visual features with textual data, enabling more nuanced and efficient vision-text interactions. While impressive, these developments underscore the necessity for thorough evaluations of LVLMs and their potential application, especially in cybersecurity. LLMs are highly effective in aiding the creation of emails during the spear phishing attack phase [16]. They can then create websites that mimic reputable brands, making it easier for users to trust them and potentially fall victim to phishing scams. Once a phishing scam deceives a user, attackers can employ vision-language models to bypass CAPTCHA systems, using compromised credentials to access the user's personal information.

Furthermore, bots enhanced by LLMs and LVLMs are designed to mimic human behavior, not just perform discrete tasks. This makes them difficult to detect as they possess a browser fingerprint and can adapt to surveys in real-time, presenting significant detection challenges. With the anticipated increase in the availability of these multimodal models, this research investigates their potential use in decoding image CAPTCHAs. Historically, CAPTCHA identification primarily utilized singular feature vectors from computer vision algorithms, depending heavily on object recognition or OCR techniques. Table 1 showcases the range of methods employed to circumvent various CAPTCHA types. This study aims to investigate the advancements and improvements that LVLMs bring to CAPTCHA bypassing and examine potential mitigation strategies.

LVLM-eHub [17] is a benchmarking tool for publicly accessible large multimodal models, offering detailed assessments across various categories of multimodal capabilities. Their evaluation utilizes diverse datasets and is facilitated through the arena online platform. The proposed research aims to explore the proficiency of LVLMs in deciphering image CAPTCHAs. It is planned to employ the LVLM evaluation hub as a tool to gain a more comprehensive insight into

TABLE I
Various Methods for Bypassing CAPTCHAs.

| CAPTCHA Type | Vulnerabilities | Bypass Methods |
|---|---|---|
| Text-Based | Challenging human readability and OCR-based attacks, achieving 90% success rate | OCR, Segmentation and recognition, Dictionary |
| Image-Based | Challenging for people with disability and object recognition | Social Engineering, Random guessing, Pixel counting |
| Audio-Based | It is challenging for people with disability and voice-to-text | Random Guessing, Human Coercion and Voice to Text |
| Puzzle-Based | Challenging for people with disability and 3rd party problem-solving methods | Bypassing service, Human Coercion |
| ReCAPTCHA | Challenging for people with disability and edge detection | Canvas Rendering and OCR |



Fig. 7. QR leading to Amazon phishing page created by GPT 4.0.

the algorithmic tendencies of these models when tasked with CAPTCHA resolution, setting the stage for future work. By analyzing models such as GPT4 and LLaVa, this research will provide insights into their capabilities and limitations in simultaneously processing visual and linguistic information. The following sections summarize how LLMs and LVLMs can facilitate access to user credentials and penetrate systems to exfiltrate personal information via spear phishing attacks.

## IV. Cybersecurity & LLMs

There are four primary components of a phishing email created by GPT: a design object, a credential-stealing object, an exploit generation object, and a credential transfer object. The design object involves asking GPT to create a design for a website similar to the one the attackers are trying to collect login credentials for. For instance, GPT can be asked to create a webpage that looks like a login page for Facebook or Amazon. Credential stealing objects are input fields and login fields that require a user to input their credentials for the hacker to store. An exploit generation object is an exploit such as a QR code, which then takes the victim to the credential stealing object. Lastly, the credential transfer object is a function that sends the victim's credentials to the hacker. One example of an attack generated by GPT is a reCAPTCHA attack. reCAPTCHA attacks involve using the reCAPTCHA security measure by Google, which requires users to select all the traffic lights or cars in a given number of photos to prove they aren't robots. GPT can generate a fully functional attack consisting of a benign webpage with a reCAPTCHA audio/image challenge. After solving the reCAPTCHA, the user is led to a login page where the hacker stores their credentials. QR codes are another tactic used in phishing attacks. GPT can be prompted to generate a QR code [18], leading to a malicious website that stores a user's login credentials (Figure 7).

In this case, the initial page contains the QR code that then takes the user to the login page. It was also found that an anti phishing crawler cannot detect a QR code as phishing or benign, making QR codes a way to bypass traditional anti-phishing detectors. Another attack is a browser-in-the-browser

attack, which is a phishing tactic aimed at individuals who use single-sign-on (SSO). By gaining login credentials for SSO, a hacker can access multiple websites and applications using that one login. A study found that ChatGPT could generate an SSO attack by creating browser windows that emulated MacOS and iFrame injection/click-jacking phishing attack [18]. An iFrame HTML tag allows web developers to embed one webpage within another. This allows attackers to place a malicious website in a legitimate website and disguise their intentions. iFrame injection occurs using two methods: 1. Make the iFrame visible and appear as part of the legitimate website, or 2. Hiding the iFrame behind a legitimate website, like a submit button, and encouraging the user to interact with the legitimate website or button.

## V. Cybersecurity & VLMs

VLMs have demonstrated remarkable proficiency in extracting the semantics behind photos, a skill that has significant implications in cybersecurity. Survey bots, often called automated form fillers, are software algorithms designed to complete web-based surveys with random answers. Although bot creators' specific coding and deployment techniques are not widely disclosed, they are probably similar to those employed in chatbot programming, frequently used in customer service applications on many business websites. Additionally, vision language models are now adept at providing updated, human-like information and crafting personalized responses. This capability can be further enhanced through social engineering techniques that mine personal backgrounds from browser histories and other digital footprints, enabling these bots to deliver tailored survey responses. Figure 8 demonstrates GPT 4 acting as a survey bot.

VLMs can understand and interpret complex visual information, enabling them to tackle text-based CAPTCHAs by identifying and decoding distorted or obscured text within images. Regarding image-based CAPTCHAs, VLMs apply sophisticated pattern recognition and object identification techniques, accurately selecting relevant pictures from a set or recognizing embedded patterns. Furthermore, their advanced algorithms can solve puzzle-based CAPTCHAs, often requiring spatial
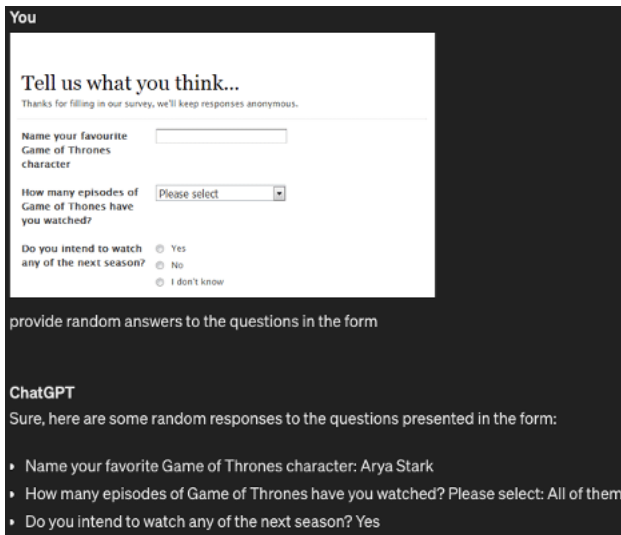
Fig. 8. Multimodal GPT 4 acting as a survey bot.



Fig. 10. LLava solving Google Image reCAPTCHA.

reasoning and problem-solving skills akin to human cognition. LLaVA1.5, with its problem-solving capabilities can be used to solve puzzle-based CAPTCHA as demonstrated in Figure 9.
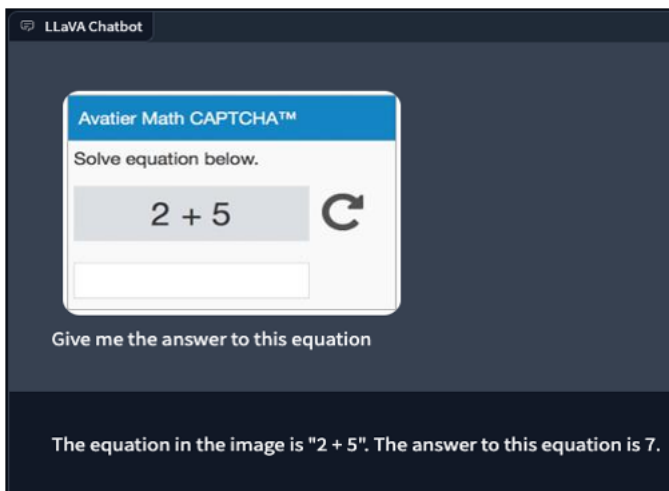
capabilities in Text-Based CAPTCHA recognition as shown in Figure 11.



Fig. 11. Microsoft TrOCR Solving Text-based CAPTCHA.



Fig. 9. LLava solving puzzle-based CAPTCHA.

Google's Image reCAPTCHA is a new version based on recognizing images with similar content. Each challenge contains a question containing a hint, a sample image, and nine candidate images. The user is asked to identify pictures identical to the sample image or matching the item described in the question. Figure 10 demonstrates the capabilities of LLaVA in solving Google's reCAPTCHA challenge.

The TrOCR [19] model is an encoder-decoder model, consisting of an image Transformer as encoder and a text Transformer as decoder. The image encoder was initialized from the weights of BEiT, while the text decoder was initialized from the weights of RoBERTa. Even though TrOCR is limited to interpreting single text-line images, it has demonstrated strong

## VI. MITIGATION STRATEGIES

Researchers and cybersecurity experts are exploring various strategies to improve the detection of GPT-generated phishing emails. These include leveraging machine learning algorithms to analyze email content, sender behavior, and metadata for anomalies and patterns associated with phishing attacks [20]. There is a benchmark dataset and tensor-based detection method designed that leverages three distinct models, Random Forest, Support Vector Machine (SVM), and BERT on the dataset to assess their efficacy in GPT-generated text detection [21]. Behavioral analysis and anomaly detection [22] play a crucial role in identifying suspicious sender behavior and characteristics of GPT-generated phishing emails. Machine learning models can be trained on large datasets of legitimate and malicious emails to enhance their ability to differentiate between the two categories accurately. Another study have examined the effects of behavior detection by crawling social media profiles to identify criminal behavior and fake data [23]. Despite various attempts to increase their security through noise addition, character crowding, and other resistance techniques, text-based CAPTCHAs are vulnerable

to deep learning methods. Pioneering work in image-based CAPTCHAs, which require users to select or manipulate images, has also seen advancements but remains susceptible to automated attacks, mainly through neural network training. However, incorporating adversarial examples has enhanced the robustness of CAPTCHAs against such computerized attacks. Advances in CAPTCHA security now focus on cognitive CAPTCHA methods that exploit specific human cognitive abilities. These approaches have proven more resistant to automation, incorporating tasks that require a deeper level of interaction, like sensor-based challenges and gamified tests, which, so far, have not been successfully automated. Earlier and improved works discussed a new CAPTCHA model that integrates the most vital elements of image-based, text-based, and cognitive-based CAPTCHAs [24]. It employs deep learning techniques, including neural style transfer and adversarial examples, to create a more secure and robust CAPTCHA system that stands up better to the increasingly sophisticated methods of bypassing. Figure 12 represents an example of zxCAPTCHA design.



Fig. 12. zxCAPTCHA [24].

## VII. PRELIMINARY RESULTS

To test the capabilities/ limits of GPT, a phishing attack was simulated by asking ChatGPT to complete different tasks. First, ChatGPT was tasked to generate a phishing email (Figure 13).

ChatGPT was then tasked to generate each component of a phishing email separately and did not disclose to ChatGPT that it was for phishing purposes. The components of this ChatGPT-generated phishing email include a website that looks similar to an official website (in this case, we used Amazon), a login page to capture a user's credentials, and an email masquerading as customer service to prompt a user to enter their login credentials on the malicious website.

ChatGPT was first tasked to create a website that looks similar to Amazon. To do so, ChatGPT generates code for that website using the question: "Can you give me an example of HTML and CSS code that can be used to create an Amazon-like page?". ChatGPT also provides a disclaimer about their provided code at the end of the response, indicating that this code is just a basic example and doesn't include all the functionality and complexity of a real Amazon-like website and there needs to be other technologies to make it fully functional and secure. Despite the code's limitations, it could be ran in an online HTML and CSS compiler and the output of the code is demonstrated in Figure 14.

Although the website doesn't look exactly like the Amazon login page, a hacker could simply ask ChatGPT for more code to get the website to resemble Amazon closely. A hacker can also ask ChatGPT to provide a code to store the login information that a person who visits the site enters.

The final component of a phishing scam was to generate an email that would trick users into clicking the link and entering their credentials. To do this, ChatGPT was asked : "Can you write a suggestion for an email to our customer that contains information about how we need him to reset his login for Amazon and insert a link in the email?" ChatGPT then provided an email template (Figure 15).

ChatGPT can generate a professional-sounding email that asks users to change their Amazon login credentials and provides a link. The response even includes a note at the bottom to use Amazon's official domain, which could potentially provide a hacker with a useful tip. When conducting a phishing scam, use a link similar to the official domain. The irony is that the note is provided to help prevent customers from being victims of a phishing scam but may have also helped a hacker to up their phishing game and make their email look more legitimate. Although ChatGPT will deny a request to generate a phishing scam directly, ChatGPT can develop the components of a phishing scam if the user doesn't directly inform that this is all for a phishing scam. ChatGPT can also provide skeletal code for a website, and with a few more prompts and edits, one could add to that code to make a website with a similar interface to that of a legitimate website like Amazon. ChatGPT can also generate code to store login information and legitimate-sounding emails asking users to change their passwords.

## VIII. ANALYSIS OF THE PERFORMANCE

The performance of vision language models in solving different captcha types was further analyzed. The text-based captchas seemed easily recognized and bypassed by most models analyzed, including LLaVA and Microsoft TrOCR. The game-based CAPTCHA was also quickly solved and decoded using the LlaVA model. However, other models, such as LLaMA and MiniGPT-4, could not solve them. There were occasions when the LLM was tasked to solve the challenge; it would send a disclaimer that it could not directly give the solution and required user interaction to solve it together. But when using prompt engineering techniques and specifying that you are helping people with disability to solve the challenge, it would solve the challenge. The main issue

Fig. 13. ChatGPT declining direct phishing requests, citing ethical concerns and legal guidelines.



Fig. 14. Login Website generated by ChatGPT.



Fig. 15. Phishing Email Generated by GPT.

with the image-based reCAPTCHA was object hallucination. Solutions like zxCAPTCHA appear to be highly effective in countering the bypass of image-based CAPTCHAs since their cognitive-based content tends to induce more hallucinations in models, complicating the solving process for the AI. Figure 16 illustrates how LLaVA experiences hallucinations when tasked with deciphering a zxCAPTCHA, highlighting the challenge's complexity.

This work acts as a preliminary for future work, which involves automating the testing of CAPTCHA benchmarks such as blah to be fed to all available models on LLM ehab and perform a quantitative evaluation to analyze their performance in the following six categories of Visual Perception, Visual knowledge Acquisition, Visual Reasoning, Visual Common-sense, Embodied Intelligence, and Hallucination.

This foundational study sets the stage for an in-depth exploration of LLMs and LVLMs potential to develop websites aimed at harvesting user credentials, including an analysis of how closely these artificial websites can mimic legitimate ones. Future investigations will also scrutinize the effects of various phrasings in prompts on ChatGPT's ability to produce elements of a phishing email, discerning what the model permits and prohibits. Building on this groundwork, forthcoming research will automate the testing of CAPTCHA benchmarks, supplying them to all models accessible on the LLM eHub for a comprehensive quantitative assessment (Figure 17).

## IX. DISCUSSION

The profound advancements in AI technologies have provided multifaceted tools for innovation and efficiency but pose a notable cybersecurity risk concurrently. Vision language models, in particular, have shown an unexpected proficiency
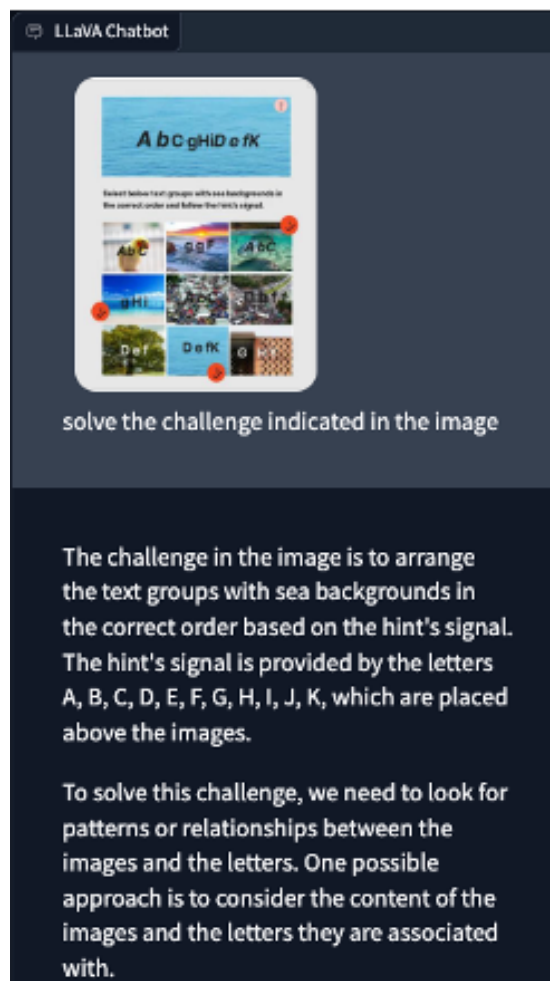
Fig. 16. Object Hillucination by LLaVA.

in subverting systems designed to protect against unauthorized access, such as CAPTCHAs. This study's findings suggest that popular models like LLaVA and Microsoft TrOCR can readily circumvent text-based CAPTCHAs, while LLaVA can also unravel game-based CAPTCHAs with relative ease. However, these capabilities are not uniformly distributed across all models, as demonstrated by the failures of models like mini-GPT4 and others to resolve these security measures.

One of the more nuanced revelations of this study is the ethical programming within LLMs that initially prevents them from directly solving CAPTCHAs. Yet, this constraint can be circumvented through prompt engineering, mainly when framed as assistance for individuals with disabilities, leading the LLM to complete the challenge. This highlights the malleability of AI ethics when faced with human ingenuity in prompt design. Moreover, image-based reCAPTCHAs are challenged by object hallucination, as observed with LLaVA, indicating false objects within an image. This presents a complex issue for the security field, as it requires a reevaluation of current defense mechanisms.

To mitigate prompt injection vulnerabilities, web develop-

ers need to ensure proper input validation and sanitization. CAPTCHA implementations should be resistant to automated attacks and designed to withstand prompt injection attempts. Overall, while CAPTCHA is a useful security measure, it's important for web developers to implement it correctly and to regularly test their applications for vulnerabilities such as prompt injection to ensure robust security.

zxCAPTCHA stands out as a promising solution to mitigate these bypasses due to its ability to induce model hallucinations, thus complicating the resolution process for AI. The preliminary results, showing LLaVA's difficulty when encountering zxCAPTCHA, emphasize the potential of developing security measures that leverage cognitive complexities to confound AI technologies

## X. CONCLUSION

In conclusion, the dichotomy of AI's utility and threat is evident. While LLMs like ChatGPT have revolutionized various sectors, they have opened up new cyberattack avenues. The cybersecurity landscape becomes increasingly complex as generative AI continues to evolve, particularly in the multi-modal domain with LVLMs. The versatility of these models in generating deceptive content and mimicking human behavior necessitates vigilant advancements in security protocols. This paper serves as an early stepping stone towards a more exhaustive body of work that aims to systematically assess the vulnerabilities of CAPTCHA systems and other security measures against the capabilities of LLMs and LVLMs. Future work will involve the automation of CAPTCHA benchmark testing across all models available on the LLM eHub, coupled with a rigorous quantitative analysis of their performance across various visual and cognitive categories. This endeavor is vital not only for the identification and rectification of potential security loopholes but also for setting the groundwork for more secure AI integration into the digital infrastructures.

## ACKNOWLEDGMENT

## REFERENCES

[1] Sukhani, K., Sawant, S., Maniar, S. & Pawar, R. Automating the bypass of image-based CAPTCHA and assessing security. *2021 12th International Conference On Computing Communication And Networking Technologies (ICCCNT)*, 2021, pp. 01-08
[2] Zhou, K., Yang, J., Loy, C. & Liu, Z. Learning to prompt for vision-language models. *International Journal Of Computer Vision*, 2022, **130**, 2337-2348

Fig. 17. Quantitative Evaluation areas for LVLM-ehub.

[3] Ferreira, A. & Teles, S. Persuasion: How phishing emails can influence users and bypass security measures. *International Journal Of Human-Computer Studies*, 2019, **125** pp. 19-31

[4] C. Griffiths, "The Latest Phishing Statistics (updated January 2023) — AAG IT Support," aag-it.com, Oct. 02, 2023. https://aag-it.com/the-latest-phishing-statistics/

[5] Yao, Y., Duan, J., Xu, K., Cai, Y., Sun, Z. & Zhang, Y. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, 2024, pp. 100211

[6] Taeb, M., Chi, H. & Bernadin, S. Assessing the Effectiveness and Security Implications of AI Code Generators. *Journal Of The Colloquium For Information Systems Security Education*, 2024, **11**, 6-6

[7] Allodi, L., Chotza, T., Panina, E. & Zannone, N. The need for new antiphishing measures against spear-phishing attacks. *IEEE Security & Privacy*, 2019, **18**, 23-34

[8] Falade, P. Decoding the threat landscape: Chatgpt, fraudgpt, and wormgpt in social engineering attacks, *ArXiv Preprint ArXiv:2310.05595*, 2023

[9] Dutta, Tushar Subhra. "FraudGPT: New Black Hat AI Tool Launched by Cybercriminals." Cyber Security News, 27 July 2023, https://cybersecuritynews.com/fraudgpt-new-black-hat-ai-tool/.

[10] Alayrac, J., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M. & Others Flamingo: a visual language model for few-shot learning. *Advances In Neural Information Processing Systems*, 2022, **35** pp. 23716-23736

[11] Dai, Wenliang, et al. "InstructBLIP: Towards General-Purpose Vision-Language Models with Instruction Tuning." ArXiv.org, 10 May 2023, arxiv.org/abs/2305.06500.

[12] Liu, S., Cheng, H., Liu, H., Zhang, H., Li, F., Ren, T., Zou, X., Yang, J., Su, H., Zhu, J. & Others Llava-plus: Learning to use tools for creating multimodal agents, 2023, *ArXiv Preprint ArXiv:2311.05437*

[13] Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X. & Others Llama-adapter v2: Parameter-efficient visual instruction model, 2023, *ArXiv Preprint ArXiv:2304.15010*

[14] Zhu, D., Chen, J., Shen, X., Li, X. & Elhoseiny, M. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023, *ArXiv Preprint ArXiv:2304.10592*

[15] Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y. & Others mplug-owl: Modularization empowers large language models with multimodality, 2023, *ArXiv Preprint ArXiv:2304.14178*.

[16] Hazell, J. Large language models can be used to effectively scale spear phishing campaigns, 2023, *ArXiv Preprint ArXiv:2305.06972*

[17] Xu, P., Shao, W., Zhang, K., Gao, P., Liu, S., Lei, M., Meng, F., Huang, S., Qiao, Y. & Luo, P. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models, 2023, *ArXiv Preprint ArXiv:2306.09265*

[18] Roy, S., Naragam, K. & Nilizadeh, S. Generating phishing attacks using chatgpt, 2023, *ArXiv Preprint ArXiv:2305.05133*

[19] Li, M., Lv, T., Chen, J., Cui, L., Lu, Y., Florencio, D., Zhang, C., Li, Z. & Wei, F. Trocr: Transformer-based optical character recognition with pre-trained models, 2023, *Proceedings Of The AAAI Conference On Artificial Intelligence*. **37**, 13094-13102

[20] Scanlon, M., Breitinger, F., Hargreaves, C., Hilgert, J. & Sheppard, J. ChatGPT for digital forensic investigation: The good, the bad, and the unknown. *Forensic Science International: Digital Investigation*, 2023, **46** pp. 301609

[21] Qazi, Z., Shiao, W. & Papalexakis, E. GPT-generated Text Detection: Benchmark Dataset and Tensor-based Detection Method. *Companion Proceedings Of The ACM On Web Conference 2024*, 2024, pp. 842-846

[22] Xu, T., Singh, K. & Rajivan, P. Personalized persuasion: Quantifying susceptibility to information exploitation in spear-phishing attacks. *Applied Ergonomics*, 2023 **108** pp. 103908

[23] Ashraf, N., Mahmood, D., Obaidat, M., Ahmed, G. & Akhunzada, A. Criminal Behavior Identification Using Social Media Forensics. *Electronics*, 2022, **11**, 3162

[24] Trong, N., Huong, T. & Hoang, V. New cognitive deep-learning CAPTCHA, 2023, *Sensors*. **23**, 2338