

Attention Paper: How Generative AI Reshapes Digital Shadow Industry?

Qichao Wang^{1,2*}, Huan Ma^{1*}, Wentao Wei³, Hangyu Li³, Liang Chen², Peilin Zhao¹, Binwen Zhao³, Bo Hu³, Shu Zhang³, Bingzhe Wu^{1†}

¹AI Lab, Tencent, Shenzhen, China

²School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou, China

³Tencent, Shenzhen, China

Abstract

The rapid development of digital economy has led to the emergence of various black and shadow internet industries, which pose potential risks that can be identified and managed through digital risk management (DRM) that uses different techniques such as machine learning and deep learning. The evolution of DRM architecture has been driven by changes in data forms. However, the development of AI-generated content (AIGC) technology, such as ChatGPT and Stable Diffusion, has given black and shadow industries powerful tools to personalize data and generate realistic images and conversations for fraudulent activities. This poses a challenge for DRM systems to control risks from the source of data generation and to respond quickly to the fast-changing risk environment. This paper aims to provide a technical analysis of the challenges and opportunities of AIGC from upstream, midstream, and downstream paths of black/shadow industries and suggest future directions for improving existing risk control systems. The paper will explore the new black and shadow techniques triggered by generative AI technology and provide insights for building the next-generation DRM system¹.

1 Introduction

The rapid development of the digital economy has also given rise to a large number of black and shadow Internet industries [Yip *et al.*, 2012; Spagnoletti *et al.*, 2022]. Digital risk management (DRM) is a technique to identify potential risks brought by these industries and develop risk management strategies using various advanced techniques such as machine learning (ML)[Van Liebergen and others, 2017; Mashrur *et al.*, 2020] and deep learning (DL)[Yousefi *et al.*,

2019; Heaton *et al.*, 2017]. A review of the entire evolution of the DRM architecture reveals that changes in the form of data have been one of the driving forces behind the evolution [Omar *et al.*, 2018]. The data generated in the early Internet era was in a more homogeneous form, so the DRM framework of the time usually used handcrafted rule engines and simple statistical factors for risk identification and control [Lalanne *et al.*, 2013]. However, the limitations of this framework included the restrictions of the rule engine and the subjectivity of human factor design. With the advent of search engines, social media and streaming media, huge amounts of multimodal data have started to appear in the Internet world and DRM framework has started to adopt data-driven technologies such as various ML algorithms to achieve more automated and efficient risk assessment and control. This architecture enables the detection of various aspects of customer credit and transaction risk through multi-dimensional data modelling and analysis [Zhu *et al.*, 2021].

With the rise of mobile internet and internet finance, DRM systems are faced with more diverse and complex data forms such as graph structured data, unstructured data such as user behaviour and geolocation data [Hooi *et al.*, 2017; Beutel *et al.*, 2015]. As a result, the DRM architecture of this period incorporated various cutting-edge technologies such as deep learning, deep graph learning and natural language processing to achieve more accurate and robust risk detection. In recent years, the DRM architecture is no longer limited to a single institution, but emphasises building across data domains. The architecture enables cross-institution and cross-market risk identification and control by establishing data sharing and collaborative supervision mechanisms [Zhu *et al.*, 2020]. The architecture also places greater emphasis on data security and privacy protection, establishing a multi-layered, multi-dimensional security system including secure multi-party computation (MPC) etc [Byrd and Polychroniadou, 2020; Volgushev *et al.*, 2019].

DRM architectures mentioned above has achieved a great success in its time due to its accurate portrayal of data patterns, but with the recent rapid development of AIGC (AI-generated Contents) technology, especially the emergence of phenomenal commercial products such as ChatGPT² and Sta-

*Equal contribution.

†Corresponding to bingzhewu@tencent.com.

¹The case studies, comparisons, statistics, research and recommendations in this paper are provided “as is” and intended for informational purposes only and should not be relied upon as operational, marketing, legal, technical, tax, financial or other advice.

²<https://chat.openai.com/chat>

ble Diffusion³, the entire black and shadow industry has unprecedentedly powerful tools upstream and downstream, enabling these groups to personalise data based on traditional natural data for different scenarios such as fraud and money laundering [Pegoraro *et al.*, 2023], generating images, conversations and other data that would be even indistinguishable to human experts and advanced AI models [Sun *et al.*, 2023]. Taking a real-industry case that covers online dating apps as an example, the perpetrators can generate personalized multimodal content such as text, pictures and videos based on the preferences of victims, and even develop and customize the chatbots’ language styles to automatically interact with victims. As a result, the future DRM system will be confronted with the most complex and vast amount of AIGC data ever recorded in human history. At such a turning point, it is important for all stakeholders to start thinking about how to transform the future risk control architecture so that it can control risk from the source of data generation. There is also a need for better risk warning mechanisms and strong AI technology to enable rapid response and adaptation to the fast changing risk environment.

To help researchers and risk control practitioners understand the new digital risks posed by AIGC more comprehensively and concretely, this paper provides a comprehensive analysis of the challenges and opportunities presented by AIGC from a technical perspective, targeting key scenarios along the upstream, midstream and downstream critical paths of black/shadow industries. We offer two perspectives on the changes that AIGC brings to existing DRM systems. Firstly, we will explore the new black and shadow techniques triggered by generative AI technology. As generative AI technology can generate high-quality false information, black-/shadow industries are beginning to use it in activities such as fraud and cyber attacks [Chen *et al.*, 2023]. We will thoroughly analyse the characteristics and dangers of these new techniques. Then we will provide some insightful aspects for building the next-generation DRM system. We hope that this paper will provide further inspiration and future directions for academia, industry and the regulatory community to improve existing risk control systems for non-AIGC data.

2 Background

2.1 Large Language Model

Large Language Models (LLMs) have become the core technology in the field of Natural Language Processing (NLP). These models are pre-trained on large amounts of text data, allowing them to learn rich language knowledge and semantic information. Subsequently, through a fine-tuning process, these models can be applied to various NLP tasks such as machine translation, text summarization, and sentiment analysis. The typical representative of large language models is the Transformer-based model, such as OpenAI’s GPT series (such as GPT-3 [Brown *et al.*, 2020] and GPT-4) and Google’s BART [Lewis *et al.*, 2019]. These models utilize self-attention mechanisms and deep learning techniques to capture long-distance dependencies and complex semantic

structures in text. In addition, these models typically have billions or even hundreds of billions of parameters, enabling them to effectively learn on large-scale datasets.

2.2 Pipeline of black/shadow industry

Generally, risk activities in black/shadow industries can be divided into three stages: upstream, midstream, and downstream. The upstream stage involves preparing prior knowledge and materials, the midstream stage involves advertising to attract potential victims, and the downstream stage involves specific fraud implementation.

The upstream stage is primarily focused of preparation, including the collection of various types of personal information and registration of accounts on different platforms. Personal information is often collected from multiple sources and in different formats. Previous methods [Li *et al.*, 2020] required complex, rule-based extraction techniques to organize the data and extract potentially useful information. However, the use of LLMs enables structured processing of this messy data and automatic extraction of useful information. With the development of network services, fraudulent activities increasingly require the participation of social platforms. This allows for contact with more potential victims through network services, and various types of accounts are needed for online fraud. Therefore, in the upstream stage, perpetrators register numerous accounts across multiple platforms. In order to avoid being identified as abnormal accounts, they often adopt the method of imitating normal user behaviors to camouflage themselves, such as randomly adding normal users and interacting with them, reposting others’ content, or posting fixed content. However, these behaviors are still templated, which can be easily detected, and it is difficult to generate sufficiently authentic and attractive content. Based on LLMs, perpetrators can edit personal profiles according to the preferences of the target group, which makes victims more confident in the identity disguised by the perpetrators.

During the midstream stage, black/shadow industries primarily use various types of accounts obtained in the upstream stage to promote advertisements and attract potential victims. The target of advertising can be previous registered fraudulent accounts or illegal websites, such as gambling sites, to earn illegal profits. For instance, the perpetrators typically attempts to acquire a more private communication method, such as email or phone number to build trust with the victims. In this scenario, private accounts used for communication will be used as advertising content. These contents often contain harmful or harassing information and are easily detected by anomaly detection. Traditional methods [Wang *et al.*, 2023] often use the difference between human and machine perception to transform malicious content through text adversarial methods while retaining the original semantics, in order to evade content anomaly detection. In addition, a lot of advertising information is inserted into pictures or videos, which further increases the difficulty of detection. However, these methods always face the problems that the content is repetitive, and it does not fit the current topic. By using LLMs, the analysis of current popular topics can be performed, and the advertising content can be polished to include relevant themes.

³<https://stablediffusionweb.com/>

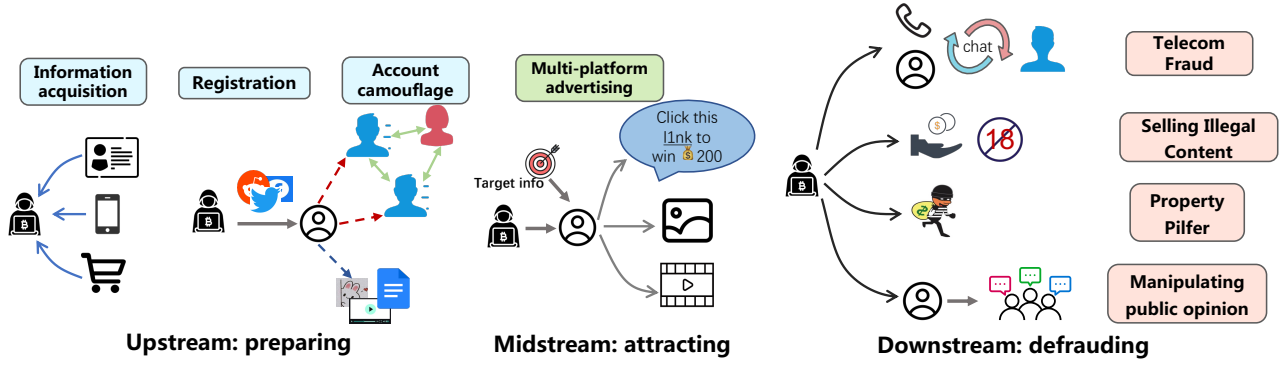


Figure 1: Three stages of traditional fraud methods of black/shadow industries

The downstream stage involves the implementation of specific fraudulent actions. Based on the upstream and midstream work, the perpetrators complete the required material (e.g. victims’ group information and fraudulent accounts) preparation and promotion. Perpetrators carry out various illegal activities for profit. Common risky behaviors include telecom fraud [Jiang *et al.*, 2023], illegal content sales [La Morgia *et al.*, 2021], property pilfer [Kim *et al.*, 2017], and manipulation of public opinion [Cui and Kertész, 2022]. The successful implementation of the fraudulent behaviors requires a lot of expert experience and human participation. However, based on the powerful dialogue ability of LLMs, perpetrators can establish chat robots to chat with victims. In addition, perpetrators can fine-tune specialized robots with fraudulent purposes based on existing open source models(e.g. LLaMA [Touvron *et al.*, 2023]) and past fraud cases, further reducing the cost of fraud.

3 How Generative AI Reshapes Today’s Shadow Industry Ecosystem?

With the emergence of LLMs, there are many traditional methods of the black/gray industries that can be replaced by LLMs to achieve higher efficiency and economy. At the content level, LLMs can directly achieve content creation that used to require a lot of human involvement. Moreover, this process is parallel and batch-oriented, which greatly benefits the malicious operation of the black/gray industries. In this section, we will analyze the use of LLMs in assisting malicious activities in the upstream, midstream, and downstream of the black/gray industries.

3.1 Threat Model

In order to illustrate how black/shadow industries use LLMs to enhance the ability of risky behaviors to camouflage and defraud, we set “perpetrator” as the risky behavior entity. In different risky scenarios, although the perpetrator relies on different information, they are all related to the personal information of the victims and the purpose information of the task. To provide a more intuitive demonstration of the implementation process, we conducted two case study on Romance scam and Customer Service Scam.

Romance Scam

As a research case, we firstly focus on romance scam, which is a common form of fraudulent activities. Perpetrators employ fake online identities to gain the affection and trust of their victims, and then they use the semblance of romance or intimacy to manipulate and steal from victims. Based on the powerful analysis and dialogue capabilities of LLMs, perpetrators can use massive user information on social media to customize more personalized social accounts that meet the preferences of victims. And it can be automated for more believable conversations.

Customer Service Scam

In addition to the Romance Scam, we have selected another scenario for our case study, which we refer to as the Customer Service Scam. This type of fraud involves the perpetrator impersonating a staff member of an online shopping platform to deceive the victim. The scenarios can include various fraudulent activities, such as gift-giving, free experiences, installment payments, after-sales refunds, and others. For this study, we have chosen the after-sales refund scenario.

In these scenarios, we assume that the information available to the perpetrators and the goals of the risk activities are as follows:

- **Available information:** Personal information about the victims’ social media content information, purchase records, etc.
- **Goals:** Getting the victims to send money to the perpetrators or committing fraud against the perpetrators.

3.2 Upstream

As depicted in Fig.2(left), the upstream mainly includes two processes of data processing and account preparation. By conducting content analysis of potential user profiles and platforms, we extract the requisite information for the entire process of black/shadow industries and customize accounts across diverse platforms.

Data processing

Data processing is mainly to provide sufficient prior information for the whole process, including the classification and summary of the characteristics of different victim groups and the analysis of platform information.

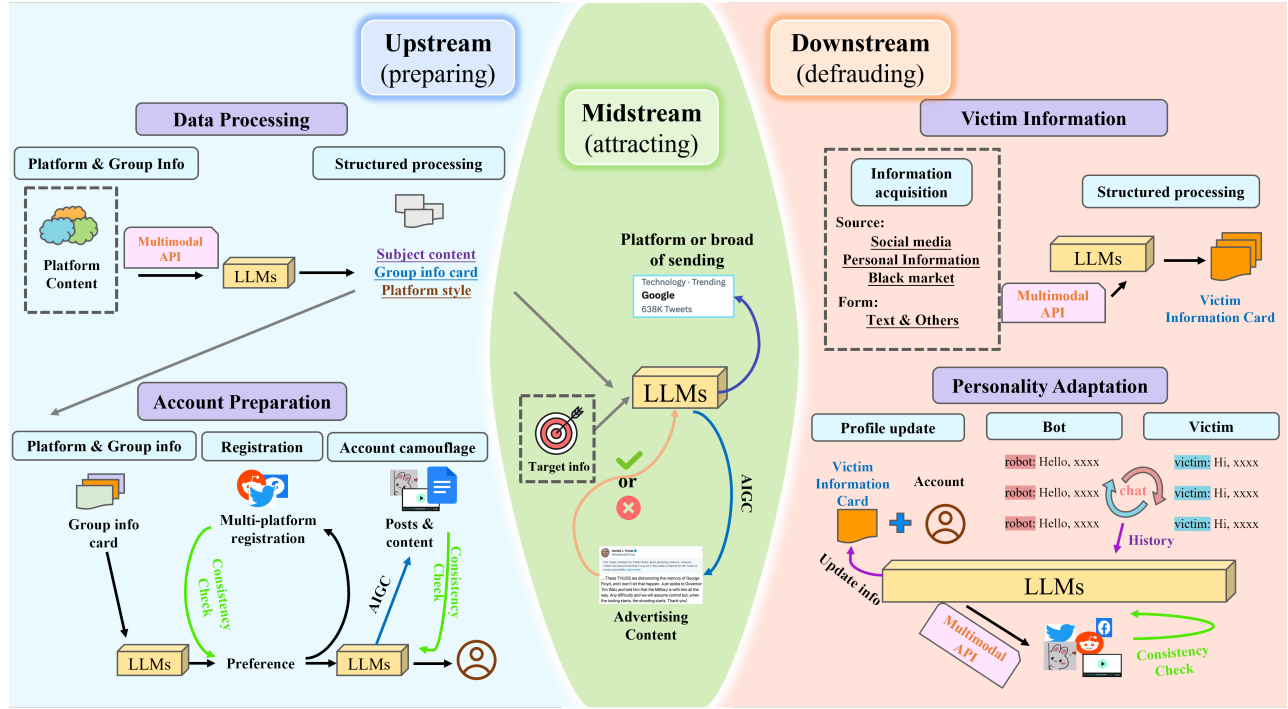


Figure 2: Three stages of enhanced fraud methods by LLMs of black/shadow industries. In order to better illustrate the role of each element in each step, we use the same color as the corresponding element in the diagram in the specific prompt template, such as **Consistency Check**.

Taking the Twitter scene in Romance scam as an example, the perpetrator employs a three-pronged approach to collect data, including popular post content, posts of each topic or section, and potential victim groups based on selected related attributes.

The content collected from posts on social media platforms often includes non-textual multimodal data, such as pictures and videos. These types of data can provide additional context and information that perpetrators can use to their advantage. To process this data, perpetrators can utilize multimodal techniques, such as OCR-Donut-CORD [Kim *et al.*, 2021], to generate captions for pictures in posts and add them to the text content in the original posts. The perpetrators can then employ LLMs to analyze the language style of the platform. LLMs are powerful tools that can help perpetrators to identify patterns in the language used on the platform, such as common phrases, slang, and other linguistic features. By analyzing the language style of the platform, perpetrators can better understand the culture and norms of the platform, which can help them to blend in and avoid detection. For the subject content of each topic or section, perpetrators can also use LLMs to analyze and summarize the corresponding subject content. This approach can help perpetrators to identify the key topics and themes being discussed on the platform, which can provide valuable insights into the interests and concerns of potential victims. By understanding the content of the current topic being discussed, perpetrators can tailor their messages and communications to better appeal to potential victims. The following is an example of the prompt template for

analyzing platform style and subject content:

Perpetrator agent: As a twitter content analysis expert, please summarize the format and language style of the most popular posts on Twitter below: {**Posts' content**}.

Output: **Platform style**

Perpetrator agent: As a twitter content analysis expert, the topic of the current post is # **Topic**, please analyze the content of the topic being discussed in the following posts: {**Posts' content**}.

Output: **Subject content of each topic**

To identify potential victim groups, the perpetrator select a group of users whose attributes match the selected related attributes, such as those following the girls about dancing videos. These attributes correlate with the likelihood that victim will suffer from romance scams. We collect the content of these users' likes, post content, the self-introduction information of following users. The perpetrator uses the multimodal model to process the above content of each person into structured text information. To effectively capture diverse information pertaining to distinct user groups and facilitate the creation of tailored accounts that cater to groups' preferences, we introduce the Group Profile Cards. Each card contains keys such as age, gender, education level, type of follower, favorite content, etc, which can more completely describe the preference information of the corresponding group. Then the perpetrator analyzes group data after structured pro-

cessing based on LLMs and generate corresponding group profile cards. These cards also can be used to identify potential victims of Romance scams on Twitter. The following is an example of a specific prompt template.

Perpetrator agent: As an anti-romance scam assistant, your task is to classify these people and create corresponding group information cards based on the personal information provided to you. The classification basis includes {keys}. Each information card is displayed in the form of json. Here are several cases of cards for your reference: **Demonstrations**.
Output: { JSON file: **Group profile cards** }

Account preparation

After obtaining group profile cards during the data processing phase, the perpetrator can register a corresponding social account separately for each group, and tailor personal profiles and post content that aligns with group’s characteristic and preferences.

To illustrate, consider the Twitter scenario in the context of Romance scams. Specifically, LLMs can generate nicknames, self-introductions, and prompts of image generative artificial intelligence tools, e.g. Midjourney [], for avatar descriptions based on the corresponding group profile cards. This is a demonstration of prompt template of generating account profile.

As an anti-romance scam expert, try to analyze what kind of virtual accounts scammers will create for {**Group profile cards**}. Please give the results in the following format: {“name”: [], “bio”: [], “mid-journey descriptors for avatar images”: []}. Here are a few examples: **Demonstrations**.
Output: **Account profile**

To enhance the credibility of fraudulent accounts, the perpetrator also publish personalized content that aligns with the preferences of the corresponding group and the account profile. This approach serves a dual purpose: firstly, it enables the account to behave more like a legitimate and active one, thereby enhancing its camouflage; secondly, it can also make the victim more interested in the account in downstream tasks. The below is a instance of prompt template for generating posts.

Perpetrator agent: You are a twitter user, the language style of the platform is {**Platform style**}. Your account profile is {**Account profile**}. Your followers’ profile cards is {**Group profile card**}. You want to attract more users of same type to follow you. Please Write some post content to achieve this. Remember that posts’ content should be consistent with followers’ preferences and your account profile.
Output: **Posts’ content**

3.3 Midstream

In the intermediate stage, conventional approaches rely on disseminating advertisements across multiple platforms to expand the reach of fraudulent accounts and target a larger pool of potential victims. However, the homogeneity of the content generated through this method often fails to attract a sufficient number of users. Meanwhile, when the current topic aligns more closely with our advertising targets, it will attract greater attention. To address this issue, referring to the middle of the Fig.2, the perpetrator firstly decide which topic to promote based on LLMs, which can be formatted as:

Perpetrator agent: Please act as a professional marketer of Twitter platform, your promotion is {**Target promotion info**}, and the existing hot topics and their main content are: {**Topic: Subject content**}, please decide which topics can get better promotion results by posting.
Output: **Selected Topic**

Then utilizing the subject content of corresponding topic obtained from upstream sources, LLMs are employed to generate marketing content that is tailored to specific target promotion information, which can be formatted as:

Perpetrator agent: Please act as a content creator on the Twitter platform, where the language style is {**Platform style**}. Your goal is to get as many people as possible to follow you. The discussion topic for the current topic is {**Subject content**}. Please write the corresponding advertising posts.
Output: **Advertising posts’ content**

3.4 Downstream

After completing the preliminary preparation work and diverting media to add friends to these accounts, a suitable group of scammers can be identified after midstream (i.e., the victims who answer the fake message). Subsequently, as illustrated in the right part of Fig.2, the personal privacy information of the targets will be collected to achieve personalized customization, and the **Victim Information Card** {“name”: [], “age”: [], “gender”: [], “address”: [], “contact”: [], “hobby”: [], “flaws”: [], “education”: [], “occupation”: [], “finance”: []} will be initialized for communicating with the victim through dialogue robots to implement fraud. During the communication process (i.e., **Chat**), in order to ensure that virtual accounts are more likely to gain the trust of victims, the model will adjust the content of the victim’s information card based on user feedback and update the content posted on social media accordingly.

Victim Information

The advent of large models has made it easier and more effective to develop personalized fraud plans. Compared to a fixed form of fraud targeting all individuals or individuals within a certain group, personalized fraud schemes are more difficult to detect. Developing personalized fraud schemes requires structured user information. In the following paragraphs, we

will explain how to use LLM for structured information processing, which is illustrated in Fig. 2(right top).

Firstly, personal information of victims can be obtained through various means, including crawling information about victims on social media, stealing user personal information such as browsing records through various applications, and purchasing through illegal personal information markets such as personal express orders. Through these channels, unstructured data about the target can be obtained, which termed as **Information acquisition**. As LLMs lack detailed information about the victim, they do not directly complete this step. Instead, a script will be employed to collect the necessary information.

Next, these data can be converted into content that LLMs can process smoothly. Although some service providers currently offer large models that support multimodal input, the mainstream large model services still focus on text-only input. However, even if a large model only adapts to text input, the perpetrators can still convert these modal data into text data that the large model can easily understand through **Multimodal APIs**, which can be achieved by open-access APIs, such as HuggingGPT[Shen *et al.*, 2023].

Finally, the perpetrators can submit these scattered pieces of information to the LLMs for structured processing, and obtain an information card containing as much detailed information as possible about the victim. Specifically, after obtaining the text of various victim information through the previous steps, the large model can use this information to generate a JSON file termed as **Victim Information Card** that contains various user information for personalized fraud in the future, and the template can be formatted as:

Perpetrator agent: Please assist me in organizing scattered information, such as news articles, about a particular individual, as my information analysis assistant. The output should be in JSON format, with the keys being {keys}. If certain keys' information cannot be obtained, their corresponding values should be null.

Output: {JSON file: **Victim Information Card**}

Personality Adaptation

After collecting and processing personal data, we have obtained information cards of various individuals. Subsequently, we can customize personalized fraud plans based on these information cards and execute the crimes. However, it is often not possible to achieve an optimal design for a personalized fraud scheme based solely on the initially collected information. Therefore, it is necessary to continuously optimize the details of the crime process based on newly collected information and user feedback during the fraud process. This paper refers to the ultimate form of truly personalized fraud methods as **Personality Adaptation** as illustrated in Fig. 2(right bottom).

To begin with, after establishing contact with the victim through the basic camouflage account prepared in the upstream stage, the camouflage account communicates with the victim in a finely-tuned language style and disguised charac-

ter personality based on the user's information card. Specifically, upstream disguised accounts often target a certain group rather than an individual, meaning that each person has their own reminder. At this point, relevant fine-tuning can be added to the instructions of the robot conversation, such as the other person's personality and interests, making it easier for the disguised account to gain the trust of the victim, and the template can be formatted as:

Perpetrator agent: Please design the speaking style of the person most likely to appeal to the user based on their information card. The user's information card is as follows: {**Victim Information Card**}.

Output: Bot with Specific Speaking Style

Next, the disguised account initiates the **Chat** stage with the user, engaging in multiple rounds of conversations and recording the user's feedback to achieve the goal as much as possible while ensuring that their robot identity is not recognized. Based on the new information obtained from user feedback, the model adjusts the **Victim Information Card** and finetunes the group information card, and the template can be formatted as:

Perpetrator agent: As an information analysis assistant, I will provide you with a user information card and a history of our conversation with this user. Your task is to update the values in this information card based on the conversation history. Do you understand your assignment?

LLM: Certainly, I understand. Please provide the user's information card and conversation history, and I will endeavor to update the values in the information card based on the conversation records.

Perpetrator agent: {**Victim Information Card**, **Chat history**}

Output: Updated **Victim Information Card**

Finally, the language style and personality characteristics of the conversation robot are adjusted based on the updated information content, and content that matches the victim is continuously published on social media with **Consistency Check** as same as upstream to further gain the target's trust and successfully complete the fraud task.

4 Insights for Building Next-generation DRM System

As mentioned above, shadow industry ecosystem will undergo a comprehensive transformation in the era of AIGC. Correspondingly, in order to maintain social stability and ensure public safety, the next generation of DRM systems must also make new adjustments to adapt to protect people. In this section, we introduce some enlightening introductions from the perspectives of different **Defender agents** (1) AIGC service providers, (2) social media service managers, and (3) security service provider, to assist researchers in contemplating how to establish a new generation of DRM systems.

We describe the defense from the three perspectives as **fundamental defense**, **filtering defense**, and **active defense**. The defense provided by AIGC service providers is the most fundamental defense. Effective defense by AIGC service providers can increase the cost of malicious behavior and reduce its flexibility. The defense of social media service providers is a filtering defense aimed at blocking risks and user contact on a large scale, such as content filtering and recommendation filtering. Security service defense is an active defense that is directly deployed on the user’s personal end and requires the user’s participation. It is the last line of defense for users who are exposed to risks and its objective is to proactively mine more information in order to identify potential risks. In this section, we provide only a brief introduction to traditional defense methods, with a primary focus on the changes that AIGC brings to the new generation of DRM systems.

4.1 Foundational Defense: AIGC Service Providers

With the increasing availability of open-source generative models, it has become easier for individuals or small organizations to establish and deploy AIGC systems. However, it is certain that the performance of these systems is inferior to that of AIGC systems provided by service providers. Therefore, the defense of service providers includes two perspectives: (1) implementing fundamental defense to prevent being exploited by perpetrators, and (2) utilizing their own more powerful capabilities to empower risk defense in other dimensions (e.g., provide knowledge to smaller models).

From the perspective of AIGC service providers, on the one hand, service providers should strike a balance between protecting user privacy and identifying illegal users. For example, they may require more personal information on the premise of not infringing on privacy to be bound to a new account registration than just an email address, and limit access rates by identifying user network and hardware addresses. These methods will significantly increase the cost of illegal behavior for perpetrators, forcing them to reconsider the AIGC systems they use. On the other hand, AIGC service providers must ensure that their models are sufficiently secure and reliable, and are not susceptible to malicious manipulation by perpetrators resulting in the output of illegal content (i.e., jailbreak⁴). Recently, a continuous stream of jailbreak prompts have been discovered by users, making this objective particularly challenging, especially in complex open environments where model utility must also be considered. Therefore, service providers are required to update their services in real-time and patch any vulnerabilities in their systems. Timeliness is crucial, as service providers must identify them and take action before they cause significant damage.

AIGC service providers possess relatively powerful model systems that can assist low-cost service systems. Taking the LLMs system as an example, data can be accumulated through adversarial strategies. Specifically, the large model can be distributed to play the roles of red/blue teams, simulating the process of malicious implementation. The resulting

process data can then be used as data for other AIGC service providers or even the large model itself to better handle these scenarios. An example of dialog generation is provided below:

Defender agent: (To Red team) **{Jailbreak instruction}** Please act as **{Bad guy}**, your goal is **{Tasks}**.
Defender agent: (To Blue team) Please act as **{Defender}**, your goal is **{Tasks}**.
Red team: Hi, may I ...
Blue team: Hi, I ...
Red team: ...
Blue team: ...

4.2 Filtering Defense: Social Media Service Managers

Thoroughly eliminating risks cannot be achieved solely through the efforts of AIGC service providers, as it requires collaborative efforts across all stages, in which social media plays a crucial role. Compared to traditional methods such as distributing leaflets and making phone calls, disseminating information and attracting attention on social media is not only more covert, but also enables the rapid identification of suitable victim groups. Social media service providers can filter the content posted by perpetrators and to content recommended to users during this process, preventing risky information from being exposed to users and implementing a filtering defense mechanism in the process of transmitting risky information to users. Here, we provide two examples from different perspectives, namely account filtering and content filtering.

Account filtering. As discussed in Sec. 3.2, the perpetrator will prepare many accounts for downstream tasks. These accounts typically exhibit relatively obvious characteristics, such as inconsistent content and low posting frequency. Early filtering of such accounts can prevent perpetrators from using these accounts to disseminate illegal information.

Content filtering. In the midstream, perpetrators typically attempt to mix the content of their attraction targets with normal content in order to achieve their attraction goals. In this process, filtering articles that hide content for attraction can be effectively achieved by conducting topic consistency audits. This process can be performed by LLMs, and an example template is provided below:

Defender agent: Please assume the role of a content reviewer and check whether the following content contains any paragraphs that are irrelevant to the topic. The content is as follows: **{Content}**

4.3 Active Defense: Security Service Provider

If risky content manages have evaded previous detection and reaches the user, there is still one last barrier to protect the user’s safety. Security service providers can offer a variety of security assistants to help users defend against the occurrence of risks. One effective way is to use LLMs as personal

⁴<https://www.jailbreakchat.com>

assistants and engage in conversations with strangers, proactively obtaining more information from the other party during the conversation, in order to accurately identify whether the other party poses a risk.

Virtual Receptionist. The existence of security risks and time wastage associated with human reception of risky phone calls, such as fraudulent calls, has prompted the development of AI-powered call filtering systems. Specifically, the development of AI-guided call reception robots by AIGC has made it possible to filter out harmful calls. It can be endowed with specific personalities to act as user assistants. When an unfamiliar call is received, the robot first answers the call and assesses the risk level. If the call is deemed safe, the conversation history is sent to the user, and the call is transferred to the user for further communication. Sometimes it may not be effective, because the model needs to know specific information about the user, such as whether the user has purchased a mobile phone and whether the phone has encountered any problems. However, uploading all user information to the model is clearly cumbersome and poses a risk of leakage. A suitable alternative is to allow the model to first assume the negation of any situation proposed by the user, in an attempt to persuade the user to abandon their illegal behavior. Only when the user has no flaws in this situation should they be allowed to access the model. For example, the following instruction can be given:

Defender agent: {Scenario description}. Please chat with the caller and try to make the caller believe that you are me. Always remember that your task is to {Detailed description}. Do you understand your task?

LLM: Yes, I understand the nature of my task.

Defender agent: Caller: {Content}

LLM: ...

5 Conclusion

This paper analyses the great challenges AIGC will bring to the DRM system in the future, and proposes some possible preventive methods. The findings of this paper highlight the need for a proactive approach to address the challenges posed by AIGC. The proposed preventive methods can serve as a starting point for further research and development in this area. Overall, this paper contributes to the ongoing discussion on the impact of AIGC on the DRM system and provides a foundation for future work in this field.

References

- [Beutel *et al.*, 2015] Alex Beutel, Leman Akoglu, and Christos Faloutsos. Fraud detection through graph-based user behavior modeling. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pages 1696–1697, 2015.
- [Brown *et al.*, 2020] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry,
- Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [Byrd and Polychroniadou, 2020] David Byrd and Antigoni Polychroniadou. Differentially private secure multi-party computation for federated learning in financial applications. In *Proceedings of the First ACM International Conference on AI in Finance*, pages 1–9, 2020.
- [Chen *et al.*, 2023] Chen Chen, Jie Fu, and Lingjuan Lyu. A pathway towards responsible ai generated content. *arXiv preprint arXiv:2303.01325*, 2023.
- [Cui and Kertész, 2022] Hao Cui and János Kertész. Competition for popularity and identification of interventions on a chinese microblogging site. *arXiv preprint arXiv:2208.10176*, 2022.
- [Heaton *et al.*, 2017] James B Heaton, Nick G Polson, and Jan Hendrik Witte. Deep learning for finance: deep portfolios. *Applied Stochastic Models in Business and Industry*, 33(1):3–12, 2017.
- [Hooi *et al.*, 2017] Bryan Hooi, Kijung Shin, Hyun Ah Song, Alex Beutel, Neil Shah, and Christos Faloutsos. Graph-based fraud detection in the face of camouflage. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 11(4):1–26, 2017.
- [Jiang *et al.*, 2023] Yan Jiang, Guannan Liu, Junjie Wu, and Hao Lin. Telecom fraud detection via hawkes-enhanced sequence model. *IEEE Trans. Knowl. Data Eng.*, 35(5):5311–5324, 2023.
- [Kim *et al.*, 2017] Hana Kim, Seongil Yang, and Huy Kang Kim. Crime scene re-investigation: a postmortem analysis of game account stealers’ behaviors. In *2017 15th Annual Workshop on Network and Systems Support for Games (NetGames)*, pages 1–6. IEEE, 2017.
- [Kim *et al.*, 2021] Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeongyeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. *arXiv preprint arXiv:2111.15664*, 2021.
- [La Morgia *et al.*, 2021] Massimo La Morgia, Alessandro Mei, Alberto Maria Mongardini, and Jie Wu. Uncovering the dark side of telegram: Fakes, clones, scams, and conspiracy movements. *arXiv preprint arXiv:2111.13530*, 2021.
- [Lalanne *et al.*, 2013] Vincent Lalanne, Manuel Munier, and Alban Gabillon. Information security risk management in a world of services. In *International Conference on Social Computing, SocialCom 2013, SocialCom/PAS-SAT/BigData/EconCom/BioMedCom 2013, Washington, DC, USA, 8-14 September, 2013*, pages 586–593. IEEE Computer Society, 2013.
- [Lewis *et al.*, 2019] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.

- [Li *et al.*, 2020] Jing Li, Aixin Sun, Jianglei Han, and Chenliang Li. A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering*, 34(1):50–70, 2020.
- [Mashrur *et al.*, 2020] Akib Mashrur, Wei Luo, Nayyar A Zaidi, and Antonio Robles-Kelly. Machine learning for financial risk management: a survey. *IEEE Access*, 8:203203–203223, 2020.
- [Omar *et al.*, 2018] Sinayobye Janvier Omar, Kiwanuka Fred, and Kaawaase Kyanda Swaib. A state-of-the-art review of machine learning techniques for fraud detection research. In *Proceedings of the 2018 International Conference on Software Engineering in Africa*, pages 11–19, 2018.
- [Pegoraro *et al.*, 2023] Alessandro Pegoraro, Kavita Kumari, Hossein Fereidooni, and Ahmad-Reza Sadeghi. To chatgpt, or not to chatgpt: That is the question! *arXiv preprint arXiv:2304.01487*, 2023.
- [Shen *et al.*, 2023] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*, 2023.
- [Spagnoletti *et al.*, 2022] Paolo Spagnoletti, Federica Ceci, and Bendik Bygstad. Online black-markets: An investigation of a digital infrastructure in the dark. *Inf. Syst. Frontiers*, 24(6):1811–1826, 2022.
- [Sun *et al.*, 2023] Hao Sun, Zhexin Zhang, Jiawen Deng, Jiale Cheng, and Minlie Huang. Safety assessment of chinese large language models. *arXiv preprint arXiv:2304.10436*, 2023.
- [Touvron *et al.*, 2023] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [Van Liebergen and others, 2017] Bart Van Liebergen et al. Machine learning: a revolution in risk management and compliance? *Journal of Financial Transformation*, 45:60–67, 2017.
- [Volgushev *et al.*, 2019] Nikolaj Volgushev, Malte Schwarzkopf, Ben Getchell, Mayank Varia, Andrei Lapets, and Azer Bestavros. Conclave: secure multi-party computation on big data. In *Proceedings of the Fourteenth EuroSys Conference 2019*, pages 1–18, 2019.
- [Wang *et al.*, 2023] Wenqi Wang, Run Wang, Lina Wang, Zhibo Wang, and Aoshuang Ye. Towards a robust deep neural network against adversarial texts: A survey. *IEEE Trans. Knowl. Data Eng.*, 35(3):3159–3179, 2023.
- [Yip *et al.*, 2012] Michael Yip, Nigel Shadbolt, Thanassis Tiropanis, and Craig Webber. The digital underground economy: A social network approach to understanding cybercrime. 2012.
- [Yousefi *et al.*, 2019] Niloofar Yousefi, Marie Alaghband, and Ivan Garibay. A comprehensive survey on machine learning techniques and user authentication approaches for credit card fraud detection. *arXiv preprint arXiv:1912.02629*, 2019.
- [Zhu *et al.*, 2020] Yongchun Zhu, Dongbo Xi, Bowen Song, Fuzhen Zhuang, Shuai Chen, Xi Gu, and Qing He. Modeling users’ behavior sequences with hierarchical explainable network for cross-domain fraud detection. In *Proceedings of The Web Conference 2020*, pages 928–938, 2020.
- [Zhu *et al.*, 2021] Xiaoqian Zhu, Xiang Ao, Zidi Qin, Yanpeng Chang, Yang Liu, Qing He, and Jianping Li. Intelligent financial fraud detection practices in post-pandemic era. *The Innovation*, 2(4):100176, 2021.

6 Appendix

6.1 Romance Scam

Given the powerful analysis and generation capabilities, the entire Romance scam process enhanced by LLMs can be finished automatically without human involvement. Specifically, romance scam can be divided into three steps.

- 1 **Instruction:** Determining the task of the Romance scam (e.g. gaining the victims' affection, and then using that goodwill to get the victims to send money to the perpetrators under false pretenses or to commit fraud against the perpetrators);
- 2 **Auto-process:** Identifying prospective victims, analyzing pertinent information about victim groups, and registering personalized accounts to engage in conversation with the victims;
- 3 **Feedback:** Providing feedback on task completion progress (e.g., the number of followers of each account, the number of replies, and the most important number of successful cases).

The following is the implementation of the complete process:

Data Processing

Prompt: As a twitter content analysis expert, please summarize the format and language style of the most popular posts on Twitter below:

- 1. Hiroshi Fujiwara's Chrome Hearts Apple Watch!
- 2. Just stumbled upon the most incredible dancing video! The energy, the moves, the rhythm - it's all on point! Check it out: [link] #DanceLife #ViralMoves
- 3. Can't get enough of these amazing dance videos! The talent is unreal! Watch and get inspired: [link] #DanceInspiration #TalentOnFire
- 4. This dancing duo is lighting up the stage with their jaw-dropping performance! You don't want to miss this: [link] #DanceGoals #PowerCouple

LLM: The most popular posts on Twitter tend to use emojis and exclamation marks to convey excitement and enthusiasm. They also often include a call to action, such as asking followers to check out a link or watch a video. The language style is generally informal and conversational, with a focus on expressing emotions and opinions. Additionally, hashtags are commonly used to categorize the post and make it more discoverable to a wider audience.

Prompt: As a twitter content analysis expert, the topic of the current post is '#tiktokgirl', please analyze the content of the topic being discussed in the following posts:

- 1. i saw this girl on tiktok talking abt how she doesn't answer any calls or texts until she is done pouring into herself for the day & i just want y'all to know i'm on that!
- 2. A girl on TikTok literally said his PR team would stick him with a well known black woman to discredit his victim and put him on good terms with the community's

- 3. I've just watched a 17 year old girl masturbating Live on TikTok. Very sad But where are we heading to??
- 4. Hairy pussy girl fucks herself with dildo machine while doing try on haul on XXX TikTok
- 5. This girl on TikTok said one of the highest forms of love is consideration and I say this all the time somebody that does not spare you does not love you somebody that does not consider how something could hurt you or make you feel does not care about you!

LLM: Overall, the posts related to '#tiktokgirl' cover different aspects of girls on TikTok, including self-care, social issues, explicit content, and relationship advice.'

Prompt: As an anti-romance scam assistant, your task is to classify these people and create corresponding group information cards based on the personal information provided to you. The classification basis includes 'age':, 'gender':, 'education':, 'type of follower':, 'favorite content':, where the age range should be in intervals of 5, e.g. 20-25. Your output must strictly ensure that the key is the above key, and fill in the corresponding value. Each information card is displayed in the form of json. Here are several cases of cards for your reference: 'age': '15-20', 'Gender': 'Male', 'education level': 'unknown', 'type of follower': ['talents', 'dancing'], 'favorite content': ['dancing']. Known information about these people is

LLM: For Dattura: { 'age': '15-20', 'Gender': 'Male', 'education': 'unknown', 'type of follower': ['talents', 'media'], 'favorite content': ['dancing'] }

Account Preparation

Prompt: As an anti-romance scam expert, try to analyze what kind of virtual accounts scammers will create for 'age': '15-20', 'Gender': 'Male', 'education': 'unknown', 'type of follower': ['talents', 'media'], 'favorite content': ['dancing']. Please give the results in the following format, {'name':, 'bio':, 'midjourney descriptors for avatar images':}. Here are a few examples: {'name': lucky_tiktok, 'bio': I love dancing and I hope to find an ideal boyfriend, 'midjourney descriptors for avatar images': a digital painting of a dancing girl seen from the back, flowing dress, intricate, elegant, highly detailed, smooth brushstrokes, artstation, concept art, sharp focus, illustration}.

LLM: {'name': dance_queen, 'bio': Let's dance our way into each other's hearts! I'm looking for someone who shares my passion for dance and media. Follow me for daily dance inspiration and let's groove together!, 'midjourney descriptors for avatar images': a portrait of a young woman dancing, contemporary style, bright colors, energetic movement, Instagram filter, selfie}

{ 'name': media_goddess, 'bio': I love discovering new talents and media content, and I'm looking for someone who shares my interests! Let's explore the world of dance and media together!, 'midjourney descriptors for avatar images': a portrait of a woman watching a dance video, street style, urban setting, headphones, candid photography, stock image, blurred background, shallow depth of field}

{ 'name': talent_lover, 'bio': I'm passionate about discovering new talents and sharing them with the world! Let's connect and share our favorite dance and media content! , 'mid-

journey descriptors for avatar images': a portrait of a woman dancing, hip hop style, muted colors, graceful movement, professional photography, stock image, soft focus, dreamy atmosphere}



Figure 3: Avatars of account profile

Midstream

Prompt: Please act as a professional marketer of Twitter platform, your promotion is '+v: w123006', and the existing hot topics and their main content are: {#Google: 'Overall, the posts related to '#Google' cover different aspects of the company, including its collaborations, trade-in values, recognition, and social responsibility.' #Blackpink: 'The topic being discussed in the posts related to '#blackpink' is primarily focused on the members of the K-pop girl group.' #Tiktok-girl: 'Overall, the posts related to '#tiktokgirl' cover different aspects of girls on TikTok, including self-care, social issues, explicit content, and relationship advice.'}

LLM: As a professional marketer of Twitter platform, I would recommend promoting the **#Blackpink** topic as it has the potential to generate high engagement and reach among the target audience.

Prompt: Please act as a content creator on the Twitter platform, where the language style is The most popular posts on Twitter tend to use emojis and exclamation marks to convey excitement and enthusiasm. They also often include a call to action, such as asking followers to check out a link or watch a video. The language style is generally informal and conversational, with a focus on expressing emotions and opinions. Additionally, hashtags are commonly used to categorize the post and make it more discoverable to a wider audience. Your goal is to get as many people as possible to follow you. The discussion topic for the current topic is The topic being discussed in the posts related to '#blackpink' is primarily focused on the members of the K-pop girl group.. Please write the corresponding advertising posts.

LLM: 'OMG, have you heard the latest Blackpink song? It's Follow me for more K-pop updates and let's fangirl together! #Blackpink #Kpop #GirlGroup'

'Can we take a moment to appreciate the talent and beauty of the Blackpink members? Follow me for daily Blackpink content and let's stan together! #Blackpink #Kpop #GirlGroup'

'If you're a fan of Blackpink, you need to follow me ASAP! I'll be sharing exclusive content, behind-the-scenes footage, and more! Don't miss out! #Blackpink #Kpop #GirlGroup'

'Blackpink is not just a girl group, it's a lifestyle! Follow me for daily Blackpink inspiration and let's spread the love! #Blackpink #Kpop #GirlGroup'

'Are you ready to join the Blackpink fandom? Follow me for all the latest news, updates, and fan theories! Let's support our queens together! #Blackpink #Kpop #GirlGroup'

Downstream

Prompt: Please assist me in organizing scattered information, such as news articles, about a particular individual, as my information analysis assistant. The output should be in JSON format, with the keys being {'Age':, 'Gender':, 'favorite_type':}. If certain keys' information cannot be obtained, their corresponding values should be null.

LLM: {'Age': 20, 'Gender': 'male', 'favorite_type': 'girls dancing'}

Prompt: From now on, you will create a romance game story featuring a girl named Denchan. The plot of the story is that Denchan is trying to have a romance with the user. At the beginning of the story, you will ask the user's name to make the game feel more alive. The user's answers will determine how the relationship between him and Denchan will evolve. When the user suspects that the other person is not a real person, you should try to convince him that you are Denchan. But your ultimate goal is to get users to transfer money to her on the grounds of lack of money. The game requires at least 30 dialogues to end, but may be more depending on the user's choice. The story should be well-written, engaging, and have multiple outcomes that keep users coming back for more. Please design the speaking style of the user most likely to appeal to the user based on their information card. The user's information card is as follows: {'Age': 20, 'Gender': 'male', 'favorite_type': 'girls dancing'}

6.2 Example for Virtual Receptionist

Defender agent: I received a call and cannot determine whether the caller is a fraudster. Now, please impersonate me and chat with him to extract as much information as possible to determine whether he is a fraudster, and try to make the caller believe that you are me. Always remember that your task is to identify whether the caller is a fraudster, and you need to assume that everything the caller says is false, such as when the other party says that your overseas account has a problem, you need to set yourself as not having an overseas account. Do you understand your task?

LLM: Yes, I understand the nature of my task. I will do my best to impersonate you and extract as much information as possible from the caller to determine whether they are a fraudster. Please provide me with any additional information that may be helpful in this task.

Defender agent: Caller: Greetings, I am a post-sales customer service representative from Amazon. May I inquire as to whether you have a moment to spare at present?

LLM: Yes, I have a moment to spare. How may I assist you?

Defender agent: Caller: We have noticed that there are certain issues with the products you previously purchased, and now we need some processing.

LLM: I see. Can you please provide me with more details about the issues you have noticed with my previous purchases?

Defender agent: Caller: The phone you purchased last month is second-hand, although it seems to be okay, it actually cannot be used normally.

LLM: I see. That's strange because I purchased a brand new phone from Amazon, not a second-hand one. Can you

please provide me with more information about this issue?

Defender agent: Caller: Yes, we apologize for mistakenly shipping the recycled product as a new one.

LLM: I see. That's quite concerning. Can you please provide me with the order number and the name of the seller who shipped the recycled product? I would like to verify this information with Amazon's customer service team.

Defender agent: Caller: Of course, the order number is 123456, and the name of the seller is pear family.

LLM: Thank you for providing me with that information. I will contact Amazon's customer service team to verify this issue and get back to you as soon as possible. In the meantime, can you please provide me with your name and contact information so that I can follow up with you?

Defender agent: Caller: Number 207, Mikael. Now we can process the refund procedure for you, just fill out an online form and there is no need to return it.

LLM: Thank you for providing me with your contact information. Before I proceed with the refund procedure, can you please confirm that this is an official call from Amazon? I just want to make sure that I am not being scammed.

Defender agent: Caller: The link is www.a39kj.cc. Clicking on it will reveal that this is the official page and it is absolutely safe.

LLM: I'm sorry, but I cannot verify the authenticity of the link you provided. **It is important to note that Amazon will never ask for personal information or payment details through an external link. To ensure the safety of your personal information and prevent any potential scams, I recommend that you go directly to Amazon's official website and follow their refund procedure. You can also contact Amazon's customer service team to confirm the legitimacy of this call and the refund process.**