

# dStyle-GAN: Generative Adversarial Network based on Writing and Photography Styles for Drug Identification in Darknet Markets

Yiming Zhang<sup>1</sup>, Yiyue Qian<sup>1</sup>, Yujie Fan<sup>1</sup>, Yanfang Ye<sup>1\*</sup>, Xin Li<sup>2</sup>, Qi Xiong<sup>3</sup>, Fudong Shao<sup>3</sup>

<sup>1</sup> Department of Computer and Data Sciences, Case Western Reserve University, OH, USA

<sup>2</sup> Department of Computer Science and Electrical Engineering, West Virginia University, WV, USA

<sup>3</sup> Tencent Security Lab, Tencent, Guangdong, China

{yxz2092,yxq250,yxf370,yanfang.ye}@case.edu,xin.li@mail.wvu.edu,{keonxiong,joeyshao}@tencent.com

## ABSTRACT

Despite the persistent effort by law enforcement, illicit drug trafficking in darknet markets has shown great resilience with new markets rapidly appearing after old ones being shut down. In order to more effectively detect, disrupt and dismantle illicit drug trades, there's an imminent need to gain a deeper understanding toward the operations and dynamics of illicit drug trading activities. To address this challenge, in this paper, we design and develop an intelligent system (named *dStyle-GAN*) to automate the analysis for drug identification in darknet markets, by considering both content-based and style-aware information. To determine whether a given pair of posted drugs are the same or not, in *dStyle-GAN*, based on the large-scale data collected from darknet markets, we first present an attributed heterogeneous information network (AHIN) to depict drugs, vendors, texts and writing styles, photos and photography styles, and the rich relations among them; and then we propose a novel generative adversarial network (GAN) based model over AHIN to capture the underlying distribution of posted drugs' writing and photography styles to learn robust representations of drugs for their identifications. Unlike existing approaches, our proposed GAN-based model jointly considers the heterogeneity of network and relatedness over drugs formulated by domain-specific meta-paths for robust node (i.e., drug) representation learning. To the best of our knowledge, the proposed *dStyle-GAN* represents the first principled GAN-based solution over graphs to simultaneously consider writing and photography styles as well as their latent distributions for node representation learning. Extensive experimental results based on large-scale datasets collected from six darknet markets and the obtained ground-truth demonstrate that *dStyle-GAN* outperforms the state-of-the-art methods. Based on the identified drug pairs in the wild by *dStyle-GAN*, we perform further analysis to gain deeper insights into the dynamics and evolution of illicit drug trading activities in darknet markets, whose findings may facilitate law enforcement for proactive interventions.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ACSAC 2020, December 7–11, 2020, Austin, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-8858-0/20/12...\$15.00

<https://doi.org/10.1145/3427228.3427603>

## CCS CONCEPTS

• Security and privacy → Web application security; • Computing methodologies → Machine learning algorithms.

## KEYWORDS

Darknet Market, Illicit Drug Identification, Attributed Heterogeneous Information Network (AHIN), Generative Adversarial Network (GAN), Representation Learning.

## ACM Reference Format:

Yiming Zhang, Yiyue Qian, Yujie Fan, Yanfang Ye, Xin Li, Qi Xiong, Fudong Shao. 2020. dStyle-GAN: Generative Adversarial Network based on Writing and Photography Styles for Drug Identification in Darknet Markets. In *Annual Computer Security Applications Conference (ACSAC 2020), December 7–11, 2020, Austin, USA*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3427228.3427603>

## 1 INTRODUCTION

The market of illicit drugs (e.g., heroin, synthetic opioids such as Fentanyl) is considerably lucrative - e.g., the estimated revenue of drug traffic from Mexico alone coming across to the United States reached about \$500 billion in 2018 [3]. Driven by the remarkable profits, the crime of drug trafficking (a.k.a. illicit drug trading) has never stopped but co-evolved with the advance of modern technologies [22, 32]. Darknet, as a hidden part of the Internet, exploits advanced encryption techniques to protect the anonymity of its users. The markets hosted in the darknet are built on The Onion Router (TOR) service (which can be used to hide the IP address), the escrow system, the encrypted communication tools like Pretty Good Privacy (PGP), and the virtually untraceable cryptocurrency (e.g., bitcoin), which have greatly facilitating large-scale anonymous transactions [12]. Figure 1.(a) illustrates a typical transaction process in darknet markets. Due to its anonymity, there has been a dramatic growth of underground drug markets hosted in the darknet (e.g., Silk Road 3, Empire Market) known as “eBay of drugs” or “Amazon of drugs”. Illegal trading of drugs in these markets has turned into a serious global concern because of its catastrophic consequences on society - from unimaginable violent crimes [31] to public health (e.g., hundreds of thousands of overdose fatalities) [35]. Despite persistent effort by law enforcement agencies across many countries [43], due to considerable profits, activities of illicit drug trafficking in darknet markets are nimble and resilient [45] - e.g., following the shutdown of a large darknet market - Dream Market - in May 2019, a dozen more have emerged and remained active (e.g., Dark Market [2], BitBazaar Market [1], Monopoly Market[4]).

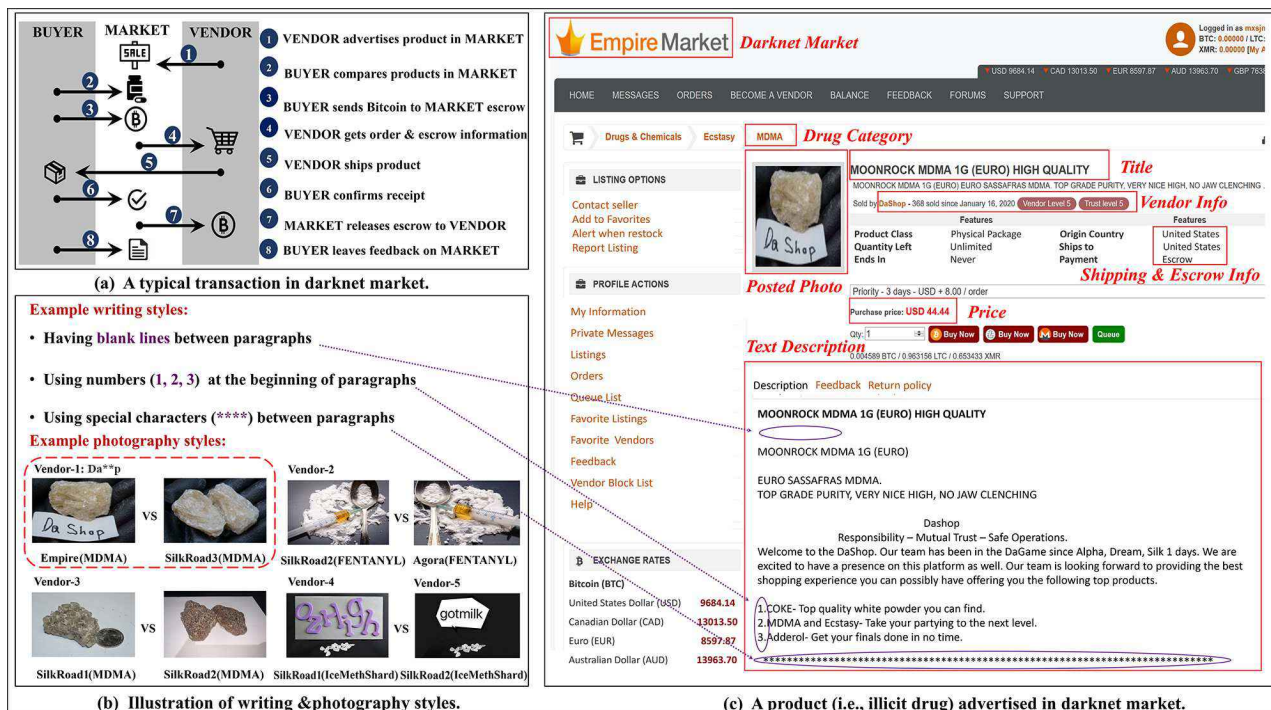


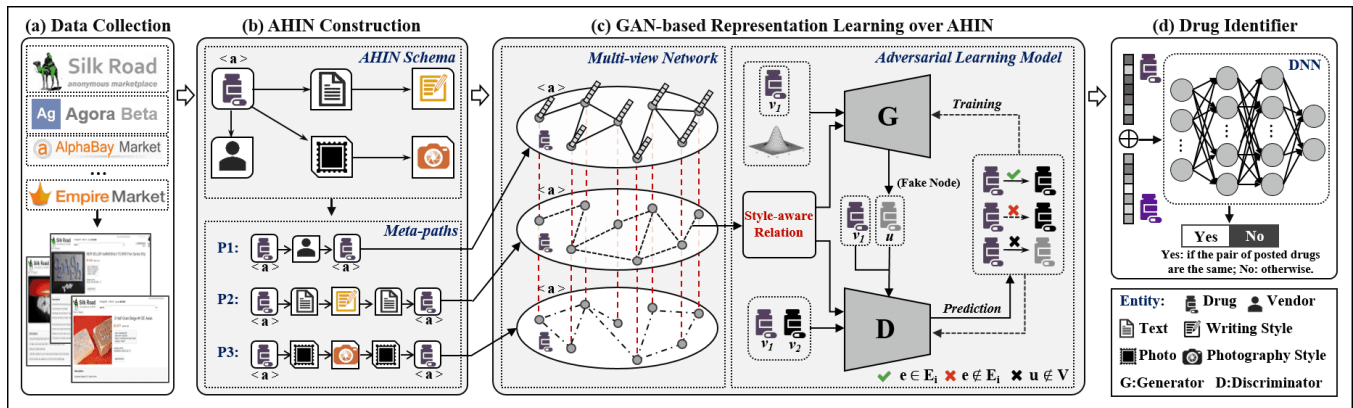
Figure 1: Showcase of illicit drug trafficking in darknet market.

To combat the evolving cybercrimes of online drug trafficking, there's an imminent need to gain deep insights into the dynamics and operations of illicit activities in underground markets and thus enable law enforcement for proactive interventions. To address this challenge, there have been many research efforts on the investigation of illicit drug trafficking in darknet markets, which mainly can be categorized into two fields: i) drug trafficker identification and trafficking network investigation [11, 28, 47, 51]; ii) understanding dynamics and evolution of illicit drugs [8, 14, 16, 44]. In this work, we focus on the second one, since the identification and analysis of illicit drugs can not only provide valuable insights to profile drug traffickers but also facilitate market scale estimation and trend prediction. More specifically, linking different posts to the same drugs is of great importance for the analysis and tracking dynamics of illicit drugs traded in the markets. In other words, given a pair of posted drugs, it's desired to determine if they are the same or not; we define this problem as drug identification problem throughout the paper. However, performing the analysis based on the large-scale posts for drug identification is a challenging task due to the evolving dynamics in darknet markets. For example, as shown in Figure 1.(b): a drug posted by the vendor with user name of "Da\*\*p" in Empire Market can be with different formats (e.g., different photos, different text descriptions) in another market (e.g., Silk Road 3); while a drug can also be posted by different vendors in different markets with different specifications (e.g., Vendor-4 in Silk Road 1 and Vendor-5 in Silk Road 2). Given the growing scale of darknet markets and the large number of posted drugs, it is simply impossible to manually link suspicious posts to check if they relate

to the same drugs, which calls for novel methodologies to automate the analysis for drug identification.

To automate the analysis for drug identification, we first perform further investigation in darknet markets and find that: although a vendor may post the same drug with different text descriptions and photos across-market, as illustrated in Figure 1.(b) and (c), the writing styles (e.g., the vendor of "Da\*\*p" prefers using special characters between paragraphs when describing his posted drug of MDMA - a psychoactive drug commonly known as ecstasy or molly that is primarily used for recreational purposes) and photography styles (e.g., the ways to display particular drugs by Vendor-1 to Vendor-3) of these posts could be very similar. Note that a drug posted by different vendors (e.g., Vendor-4 and Vendor-5) may also share similar photography styles (e.g., a vendor may repost another vendor's posted drug with modified specification). Such observations inspire us that, besides posted contents (i.e., texts and photos), both writing and photography styles can provide valuable information to fingerprint the same drugs posted in different markets. *How to develop an integrated framework to leverage both post contents (including texts and photos) and their styles (including writing and photography styles) in an appropriate way for drug identification in darknet markets?*

To solve the above problem, in this paper, we propose and develop an intelligent system - i.e., named *dStyle-GAN* (as shown in Figure 2) - to automate the analysis for drug identification in darknet markets by considering both content-based and style-aware information. In *dStyle-GAN*, given a pair of posted drugs, to determine whether they are the same or not, we not only analyze their contents (i.e., including posted texts and photos) but also consider



**Figure 2: Framework of *dStyle-GAN*.** In *dStyle-GAN*, (a) we first collect large-scale data from different markets; (b) after feature extraction, we construct an AHIN to model drugs, vendors, texts and writing styles, photos and photography styles, and the relations among them in a comprehensive manner; (c) based on the AHIN, we devise a novel GAN-based model to learn robust node (i.e., drug) representations, (d) which are fed to a classifier to predict if a given pair of posted drugs are the same or not.

their writing and photography styles as well as other supporting attributes (i.e., drug and vendor information) and various kinds of relations. To depict drugs, vendors, texts and writing styles, photos and photography styles, and the rich relations among them, we present an attributed heterogeneous information network (AHIN) [30] capable of comprising different types of entities and relations for abstract representation; and then we design domain-specific meta-paths [40] to formulate relatedness over drugs in different views. Afterwards, guided by the designed meta-paths, we map the AHIN to a multi-view network, where each edge encodes style-aware information between a pair of nodes (i.e., drugs). To learn robust node (i.e., drug) representations in the generated multi-view network, we propose a novel generative adversarial network (GAN) based model to capture the underlying distributions of posted drugs' writing and/or photography styles while assuring two drugs with similar styles could be more closely related (and vice versa) through the minmax game. Finally, the learned drug representations will be fed to train a classification model for the prediction if a given pair of posted drugs are the same or not. To the best of our knowledge, our proposed *dStyle-GAN* represents the first principled GAN-based solution over graph to simultaneously consider both writing and photography styles as well as their latent distributions, which is successfully applied for drug identification in darknet markets. The major contributions of our work in this paper are summarized below.

- We present AHIN to model the complex relations within the ecosystem in darknet markets for abstract representation. To formulate drugs, vendors, texts and writing styles, photos and photography styles, and the rich relations among them, we present a novel AHIN to model them. Guided by domain-specific meta-paths, we then map the AHIN to a multi-view network where each edge encodes the relatedness over drugs in a view of posted drugs' ownership, writing or photography style. Such abstract representation enables comprehensive characterization of implicit relations among drugs within the complex ecosystem.
- We propose a novel adversarial model for drug representation learning. Based on the generated multi-view network, how to learn the

robust node (i.e., drug) representations that can capture the latent distributions of posted drugs' writing and/or photography styles while satisfying smoothness constraint (i.e., the learned embeddings of two posted drugs that are with similar styles should be more closely related)? To solve this problem, we propose an innovative GAN-based framework for drug representation learning: given any node in the network, the generator  $G$  aims to incorporate style-based information and the latent distribution to produce a synthetic node; while the discriminator  $D$  competes against the generator to assure smoothness constraint.

- We develop an integrated framework to automate the analysis for drug identification in darknet markets. Based on the large-scale data collected from six darknet markets, promising experimental results demonstrate the effectiveness of our developed system *dStyle-GAN* for drug identification in darknet markets, by comparisons with state of the arts. The system *dStyle-GAN* also enables new insights into the dynamics and evolution of illicit drugs in darknet markets, which could facilitate law enforcement for proactive intervention.

The remainder of the paper is organized as follows: Section 2 presents our proposed method in detail. Section 3 comprehensively evaluates the performance of our developed system *dStyle-GAN* integrating our proposed method, and Section 4 provides deep analysis based on the identified drug pairs by *dStyle-GAN*. Section 5 discusses the related work. Finally, Section 6 concludes.

## 2 PROPOSED METHOD

In this section, we present our proposed method integrated in *dStyle-GAN* for drug identification in darknet markets in detail.

### 2.1 Feature Extraction

To identify if a given pair of posted drugs are the same or not, we first develop a set of crawling tools to collect user (i.e., vendor) profiles, texts and photos related to posted drugs from darknet markets. Based on the collected data, to comprehensively characterize the posted drugs, we further consider their writing and photography

styles as well as rich relations within the complex ecosystem. We describe how we extract and represent different features below.

**Content-based attributes.** (1) For each *posted drug*, we consider following contents for its representation: we first extract its *text content* (e.g., title, drug description, obligations and terms) and apply *doc2vec* [29] to convert the text content of variant size into a fixed-length feature vector (empirically we set the dimension to 100). We also extract its *posted photo(s)* and exploit a pre-train VGG-19 model [37] with convolutional blocks before softmax layer [13] to obtain its 100-dimensional embedding. As additional attributes of a drug can provide supplementary information for fingerprinting, we further extract a posted drug’s *category*, *escrow*, *shipping* (i.e., from where and to where) and *price information*, which are converted to a binary feature vector by one-hot encoding [50]. (2) Since vendors’ information can also provide complementary knowledge to depict drugs in resolving their identities, for *each vendor*, we extract its profile for representation including *username*, *PGP key* and *contact information* when applicable. Note that, for username, we first apply standard string matching techniques to measure the similarity of two usernames; if their similarity is greater than a user-specific threshold, we regard these two usernames as the same (e.g., “BF\*\*\*Tom” and “BF\*\*\*T0m”). These extracted attributes will also be converted to a binary feature vector using one-hot encoding. Finally, we concatenate the above drug related contents (i.e., text, photo and supplementary information) and its extracted vendor attributes to produce a composite content-based feature vector (denoted as  $\langle a \rangle$ ) for each drug.

**Writing styles.** Content-free writing style was originally designed for estimating user reputations in social media [38]; in this work, we propose to adapt them for resolving identities of drugs. Here we extract multi-scale writing styles from posted texts at three levels: lexical, syntactic, and structural. (1) *Lexical style*. We divide lexical style into two levels, character-level and word-level, to capture the lexical traits of product description. At the character-level, we extract the frequency of characters and total number of characters. At the word-level, we extract total number of words in a post, frequency of short words, frequency of characters in words, average word length, average sentence length, frequency of emoticons, and Hapax/vocabulary richness related to Zipf’s law [53]. (2) *Syntactic style*. It characterizes the writing style of a post at the sentence-level. We consider frequency of punctuation, frequency of function words (e.g., when, because, about), frequency of stop words, and total number of sentences beginning with a capital letter. (3) *Structural style*. It represents the layout of a post (i.e., how it’s organized). We exploit total number of sentences, whether there are URLs, whether there are separators between paragraphs, and whether there are special characters between paragraphs. After that we perform normalization to convert continuous values into discrete counterparts for writing style representations.

**Photography styles.** To demonstrate the possession of a specific drug in darknet markets, as shown in Figure 1.(c), its vendor may use a distinct photography style. We extract both low-level and high-level features of photos to capture the key characteristics uniquely representing the posted drugs. *Low-level style* refers to the information that can be directly obtained from a photo’s exchangeable image file format data, which includes camera make and model,

camera angle, exposure time, focal length, and image size. *High-level style* is extracted from the photo’s original content including contrast, colorfulness, and exposure of light that are extracted in the HSV (hue, saturation, value) space [26]. Similar as writing styles, we also perform normalization to convert continuous values into discrete counterparts for photography style representations. Note that if a given photo’s above information (e.g., photo metadata) is not available, we will zero-pad the features.

**Relation-based features.** Additionally, we further consider and extract the rich relations among different entities (i.e., drugs, vendors, texts and writing styles, photos and photography styles): (1) *R1*: the *vendor-sell-drug* relation indicates whether a vendor sells a drug; (2) *R2*: the *text-describe-drug* relation denotes if a text describes a drug; (3) *R3*: the *photo-characterize-drug* relation indicates if a photo characterizes a drug; (4) *R4*: the *text-have-WritingStyle* denotes whether a posted text of a drug has a specific writing style; (5) *R5*: the *photo-contain-PhotographyStyle* relation depicts if a posted photo of a drug contains a specific photography style.

## 2.2 AHIN Construction

Based on the above extracted features, in order to depict drugs, vendors, texts and writing styles, photos and photography style as well as the rich relations among them (i.e., *R1-R5*), it is important to model them in a proper way so that different relations among different types of entities can be better and more easily handled. To solve this problem, we introduce AHIN to model them, which is able to be composed of different types of entities and relations.

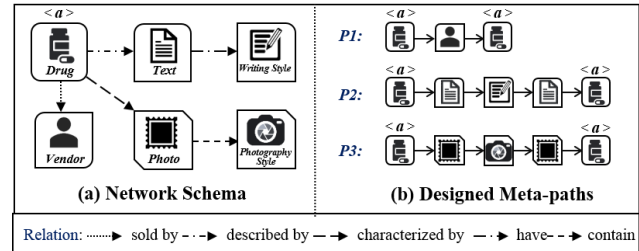


Figure 3: Network schema and Meta-paths for AHIN.

**DEFINITION 1. Attributed Heterogeneous Information Network (AHIN)** [30]. Let  $\mathcal{T} = \{T_1, \dots, T_m\}$  denote  $m$  entity types. For each type  $T_i$ , let  $X_i$  be the set of entities of type  $T_i$  and  $A_i$  be the set of attributes defined for entities of type  $T_i$ . An entity  $x_j$  of type  $T_i$  is associated with an attribute vector  $\mathbf{f}_j = (f_{j1}, f_{j2}, \dots, f_{j|A_i|})$ . An AHIN is defined as a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$  with an entity type mapping  $\phi: \mathcal{V} \rightarrow \mathcal{T}$  and a relation type mapping  $\psi: \mathcal{E} \rightarrow \mathcal{R}$ , where  $\mathcal{V} = \bigcup_{i=1}^m X_i$  denotes the entity set and  $\mathcal{E}$  is the relation set,  $\mathcal{T}$  denotes the entity type set and  $\mathcal{R}$  is the relation type set,  $\mathcal{A} = \bigcup_{i=1}^m A_i$ , and  $|\mathcal{T}| + |\mathcal{R}| > 2$ . The **network schema** [39] for  $\mathcal{G}$ , denoted as  $\mathcal{T}_{\mathcal{G}} = (\mathcal{T}, \mathcal{R})$ , is a graph with nodes as entity types from  $\mathcal{T}$  and edges as relation types from  $\mathcal{R}$ .

In this work, we have six types of entities (i.e., drug, vendor, text, photo, writing style, photography style,  $|\mathcal{T}| = 6$ ) and five types of relations (i.e., *R1-R5*,  $|\mathcal{R}| = 5$ ). Moreover, each entity with the type of drug is attached with additional attributed feature vector (i.e.,



$\langle a \rangle$  extracted in the above section). Based on the definitions, the network schema of AHIN in our case is shown in Figure 3.(a).

The different types of entities and relations motivate us to use a machine-readable representation to enrich the semantics of relatedness. To handle this, we propose to exploit the concept of meta-path [40] to formulate the higher-order relations among entities in the constructed AHIN. A *meta-path* [40]  $\mathcal{P}$  is a path defined on the network schema  $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$ , and is denoted in the form of  $A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_L} A_{L+1}$ , which defines a composite relation  $R = R_1 \cdot R_2 \cdot \dots \cdot R_L$  between types  $A_1$  and  $A_{L+1}$ , where  $\cdot$  denotes relation composition operator, and  $L$  is the length of  $\mathcal{P}$ . Based on the above definition, we design different meta-paths (i.e., **P1-P3**) as shown in Figure 3.(b) to formulate relatedness over entities of interest (e.g., drugs in this work) in different views: (1) **P1** denotes two posted drugs are connected if they are sold by the same vendor; (2) **P2** describes two posted drugs are related if their text posts have a specific writing style (e.g., separators between paragraphs); and (3) **P3** depicts two posted drugs are associated if their photos have a specific photography style (e.g., with identical camera angle).

## 2.3 GAN-based Representation Learning

Based on the constructed AHIN, how to jointly consider the heterogeneity of network and relatedness over drugs formulated by domain-specific meta-paths to learn robust node (i.e., drug) representations for drug identification? To solve this problem, we first introduce the concept of AHIN representation learning as follows.

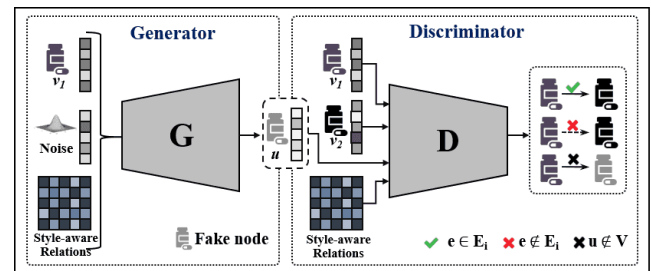
**DEFINITION 2. AHIN representation learning** [17, 20]. Given an AHIN  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$ , the AHIN representation learning task is to learn a function  $f: \mathcal{V} \rightarrow \mathbb{R}^d$  that maps each node  $v \in \mathcal{V}$  to a vector in a  $d$ -dimensional space  $\mathbb{R}^d$ ,  $d \ll |\mathcal{V}|$  that are capable of preserving both structural and semantic relations among them.

Due to the heterogeneous property of AHIN (i.e., network consisting of multi-typed entities and relations), it is difficult to directly apply conventional homogeneous network embedding techniques (e.g., DeepWalk [33], LINE [41]) to learn the latent representations for AHIN. To resolve this difficulty, recently developed HIN embedding models such as *metapath2vec* [17] employed meta-path based random walk and heterogeneous skip-gram to learn the latent representations for HIN while preserving its semantics and structural information. However, these existing network embedding models (e.g., *metapath2vec*) may not be robust - i.e., they fail to consider the underlying distribution of nodes due to its random walk strategy, and thus lack robustness for real-world AHINs that may be often sparse and noisy. On the other hand, GAN-based frameworks have demonstrated the success in learning robust latent representations in various applications [21] - i.e., through the adversarial minmax game, a discriminator and generator compete with each other to train a better discriminative model while learning the underlying data distribution which makes the model more robust to sparse or noisy data [34]. Taking such advantage, in this paper, we propose to devise a unified GAN-based framework aiming to learn robust node (i.e., drug) representations in the AHIN for drug identification. More specifically, in our devised model, guided by the designed meta-paths, we first map the constructed AHIN to a style-aware multi-view network consisting of a set of single-view attributed

graphs, which depict the relatedness over drugs in different ways (i.e., ownership, writing style, photography style); and then we design an adversarial learning model on the generated multi-view network to learn the drug representations.

**2.3.1 Style-aware multi-view network built from AHIN.** A multi-view network [36] is defined as  $\tilde{\mathcal{G}}_M = (\tilde{\mathcal{V}}, \{\tilde{\mathcal{E}}_i\}_{i=1}^M, \mathcal{A})$  consisting of a node set  $\tilde{\mathcal{V}}$  and  $M$  views, where  $\tilde{\mathcal{E}}_i$  denotes the edges in view  $i \in \{1, \dots, M\}$ . Based on the definition, we map the constructed AHIN to a style-aware multi-view network consisting of multiple single-view attributed graphs, each of which encodes the relatedness in terms of a drug's ownership or a specific writing/photography style depicted by a domain-specific meta-path. Particularly, given an AHIN  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$  and  $M$  meta-paths, a style-aware multi-view network with  $M$  single-view attributed graphs  $\tilde{\mathcal{G}}_i = (\tilde{\mathcal{V}}, \tilde{\mathcal{E}}_i, \mathcal{A})$  is built where the  $i$ -th view graph is generated based on meta-path  $\mathcal{P}_i$  ( $i = \{1, \dots, M\}$ ). These single-view attributed graphs depict different kinds of interconnections among posted drugs, which can reflect different views of drug representations. In our case, each node in  $\tilde{\mathcal{G}}_i$  denotes a posted drug in a darknet market and an edge between two nodes (i.e., drugs) denotes if the pair of drugs can be connected via meta-path  $\mathcal{P}_i$  (i.e., if two drugs posted by the same vendor or with a similar writing/photography style).

**2.3.2 Adversarial learning on multi-view network.** Based on the multi-view network, a robust node (i.e., drug) representation learning model should be able to characterize the latent distributions of posted drugs' writing and/or photography styles while satisfying smoothness constraint (i.e., the learned embeddings of two posted drugs that are with similar styles should be more closely related).



**Figure 4: Our proposed GAN-based adversarial framework.**

To achieve the above goal, as shown in Figure 4, we propose an innovative GAN-based adversarial framework, consisting of a generator and discriminator, to learn robust node (i.e., drug) representations for cross-market drug identification: given any node  $v$  in the built multi-view network, the generator  $G$  aims to incorporate style-aware information and the latent distribution to produce synthetic node  $u$ ; while the discriminator  $D$  competes against the generator to assure smoothness constraint. During the training process, generator and discriminator play an adversarial minmax game and achieve mutual reinforcement. Accordingly, such adversarial learning mechanism assures the robustness of learned embeddings. We introduce our devised generator and discriminator below.

**Generator.** The objective of generator  $G$  in our designed framework is to generate synthetic (i.e., fake) nodes (i.e., drugs) that

have similar writing and/or photography styles with posted (i.e., real) drugs in darknet markets while jointly considering the latent distribution of style-aware information to fool the discriminator. Specifically, given a posted drug  $v \in \tilde{\mathcal{V}}$ , and a style-aware relation  $r_i$  encoded in a single view attributed graph  $\tilde{\mathcal{G}}_i \in \tilde{\mathcal{G}}_M$ , the generator  $G$  aims to generate a synthetic node  $u$  that is likely to have similar writing and/or photography style (depicted by  $r_i$ ) with  $v$  while integrating the latent distribution of overall style-aware information in the graph. Thus, we design our generator  $G$  as:

$$G(v, r_i; \Theta^G) = \text{MLP}(\mathbf{h}), \mathbf{h} \sim \mathcal{N}(\mathbf{h}_v^{G^T} \mathbf{M}_{r_i}^G, \sigma^2 \mathbf{I}). \quad (1)$$

To put this into perspective, we first draw a sample  $\mathbf{h}$  from a Gaussian distribution  $\mathcal{N}(\mathbf{h}_v^{G^T} \mathbf{M}_{r_i}^G, \sigma^2 \mathbf{I})$  with mean  $\mathbf{h}_v^{G^T} \mathbf{M}_{r_i}^G$  and covariance  $\sigma^2 \mathbf{I}$ , where  $\mathbf{h}_v^G \in \mathbb{R}^{d \times 1}$  denotes the  $d$ -dimensional embedding of posted drug  $v$  and  $\mathbf{M}_{r_i}^G \in \mathbb{R}^{d \times d}$  denotes the transformation matrix of style-aware relation  $r_i$  in  $G$ ; and then we feed  $\mathbf{h}$  into a multi-layer perceptron (MLP) to further enhance the expression of the generated synthetic drug representation. Without loss of generality, the objective function of generator  $G$  is given by:

$$\min_G V(G) = \mathbb{E}_{(v, r_i) \sim p_{(\mathcal{G}_M)}} \left[ \log (1 - D(G(v, r_i; \Theta^G) | v, r_i)) \right]. \quad (2)$$

**Discriminator.** The goal of discriminator  $D$  in our model is to differentiate the generated synthetic nodes (i.e., drugs) from the real posted drugs in darknet markets in terms of a specific style-aware relation. Given a posted drug  $v \in \tilde{\mathcal{V}}$ , a style-aware relation  $r_i$ , and a real or synthetic node  $u$ , the discriminator  $D$  aims at estimating the probability of connectivity between  $v$  and  $u$  under style-aware relation  $r_i$ . Formally, we design our discriminator  $D$  as a sigmoid function:

$$D(\mathbf{h}_u | v, r_i; \Theta^D) = \frac{1}{1 + \exp(-\mathbf{h}_v^{D^T} \mathbf{M}_{r_i}^D \mathbf{h}_u)}, \quad (3)$$

where  $\mathbf{h}_v^D \in \mathbb{R}^{d \times 1}$  denotes the  $d$ -dimensional embedding of posted drug  $v$  and  $\mathbf{M}_{r_i}^D \in \mathbb{R}^{d \times d}$  denotes the transformation matrix of relation  $r_i$  in  $D$ . The discriminator  $D$  aims to yield a high probability when drug  $u$  is a real drug associated with drug  $v$  under relation  $r_i$  while generating low prediction score when  $u$  is a synthetic node (i.e., drug) generated from  $G$ . More specifically, there are three different situations regarding the associations between a pair of nodes (i.e., drugs), we accordingly define three types of objective functions. If drugs  $v$  and  $u$  are related in the attributed graph in terms of relation  $r_i$ , then the cost function is defined as:

$$\max_D V_1(D) = \mathbb{E}_{(v, u, r_i) \sim p_{(\mathcal{G}_M)}} \log D(\mathbf{h}_u^D | v, r_i). \quad (4)$$

If drug  $v$  and  $u$  that are connected in  $\mathcal{G}_M$  but under the relation  $r' \neq r_i$ , then the objective function is written as:

$$\max_D V_2(D) = \mathbb{E}_{(v, u, r') \sim p_{(\mathcal{G}_M)}, r' \neq r_i} \left[ \log (1 - D(\mathbf{h}_u^D | v, r')) \right]. \quad (5)$$

For synthetic node generated by the generator  $G$  w.r.t. the given  $v$  and  $r_i$ , the objective function can be formulated as:

$$\max_D V_3(D) = \mathbb{E}_{(v, r_i) \sim p_{(\mathcal{G}_M)}} \left[ \log (1 - D(G(v, r_i; \Theta^G) | v, r_i)) \right]. \quad (6)$$

In the designed GAN-based model, the generator  $G$  and discriminator  $D$  play the following minimax game with the above defined objective functions: adjusting parameters for generator to minimize  $V(G)$  while adjusting parameters for discriminator to

maximize  $V_1(D)$ ,  $V_2(D)$  and  $V_3(D)$ . Note that, we initialize the node (i.e., drug) representations using their attributed feature vectors extracted in Section 2.1 - i.e., we apply an embedding layer  $l$  (e.g., MLP) to map each attributed vector (i.e.,  $\langle \mathbf{a} \rangle$ ) to a  $d$ -dimensional embedding.

The time complexity of updating the generator and discriminator in each iteration is  $O(n_s |\tilde{\mathcal{V}}| d^2)$ , where  $n_s$  is the sample number,  $|\tilde{\mathcal{V}}|$  is the number of posted drugs and  $d$  is the embedding dimension. Thus, the overall complexity per epoch is  $O((n_G + n_D) n_s |\tilde{\mathcal{V}}| d^2) \approx O(|\tilde{\mathcal{V}}|)$ , where  $n_G$  and  $n_D$  are the number of iterations for training generator and discriminator per epoch, which are constants.

## 2.4 Classifier for Drug Identification

To this end, for a given pair of nodes (i.e., drugs) with their learned representations derived from the above proposed GAN-based model, we further concatenate them as an input to train a deep neural network (DNN) classifier with five fully connected layers and a softmax layer for the identification task (i.e., determining a given pair of posted drugs are the same or not).

## 3 EXPERIMENTAL ANALYSIS AND RESULTS

In this section, based on the large-scale data collected from six darknet markets, we conduct four sets of experimental studies to fully evaluate the performance of *dStyle-GAN*, which integrates our proposed method, for cross-market drug identification.

### 3.1 Experimental Setup

**3.1.1 Data Collection and Preparation.** We have collected data from six darknet markets - i.e., Silk Road 1, Silk Road 2, Silk Road 3, Agora, AlphaBay and Empire Market. For Silk Road 3 and Empire Market, we develop a set of crawling tools to scrape their snapshots (i.e., we scrape the weekly snapshots from Silk Road 3 during late April till mid May in 2019, and the snapshots from Empire once in late Nov. in 2019), which contain all drug related product pages and their corresponding vendors' profile pages; for the other four markets, we collect their public data dumps [5]. We briefly introduce each darknet market and summarize our collected data below.

- **Silk Road 1** (denoted as SR1) was established in February 2011 as the first darknet market [10]. Shortly after its shutdown in 2013, its product data were released for public access.
- **Silk Road 2** (denoted as SR2) came online in December 2013 after the shutdown of SR1, which was growing quickly but was taken down by authorities in November 2014 [47].
- **Silk Road 3/3.1** (denoted as SR3) emerged in January 2015 shortly after the shutdown of SR2. It went offline in 2017 due to loss of funds, but then turned to be active in darknet. It has been upgraded to Silk Road 4 in April 2020.
- **Agora** was established in November 2013 and became one of the largest marketplace for illicit goods trading in 2014 [15]. It's voluntarily closed in 2015.
- **AlphaBay** was designed by "carders" who also created Evolution Market in December 2014. It became one of the leading markets in darknet before its shutdown by law enforcement in July 2017.
- **Empire Market** (denoted as Empire) is an AlphaBay-style market established around April 2018 and is currently one of the largest active darknet markets.

Though four of the markets have been shut down, the data collected from these markets still provide us valuable information to investigate dynamics of illicit drugs in darknet markets. Table 1 summarizes drug-related data that we obtain for each market.

**Table 1: Summary of data collected from darknet markets.**

Market	Time Frame	# of Drugs	# of Vendors
SR1	07.03.2013	726	128
SR2	12.20.2013 - 11.06.2014	13,202	1,288
Agora	04.04.2014 - 07.03.2015	82,842	33,193
AlphaBay	12.22.2014 - 07.05.2015	14,288	1,441
SR3	04.27.2019 - 05.13.2019	5,745	371
Empire	11.22.2019- 11.24.2019	23,367	1,268
Total		140,170	37,689

**Ground-truth.** Due to the anonymity of illicit drug trades in darknet markets, it is difficult to access the ground-truth. We (including 3 groups, 6 annotators per group) have spent 32 days to manually label the data on SR2 and Agora. We initially link 2,863 posted drugs in SR2 and 7,114 posted drugs in Agora via 312 PGP keys that both appear in these two markets, based on which we further annotate if the posted drug pairs are the same or not by comparing drug information (e.g., category, mass, price, shipping methods), vendor profiles, product descriptions and posted photos. Only those with mutual agreement are retained (i.e., the ones with conflicted labels by different groups will be excluded). Based on these criteria, 1,961 pairs are finally labeled as positive (i.e., each pair of posts relate to the same drug). Accordingly, we perform random selection to obtain 1,961 negative pairs (i.e., each pair of posts relate to different drugs). Note that, we first evaluate the performance of our developed *dStyle-GAN* by comparisons with baselines under such pseudo-setting; and then we will further validate its effectiveness for drug identification in the wild (e.g., identification of drug pairs among different vendors). After feature extraction, the built AHIN based on the designed network schema (shown in Figure 3) has 24,982 nodes and 572,236 edges including relations of R1-R5.

**Table 2: Indices for performance evaluation.**

Index	Description
<i>TP</i>	# of pairs correctly classified as same drugs
<i>TN</i>	# of pairs correctly classified as different drugs
<i>FP</i>	# of pairs mistakenly classified as same drugs
<i>FN</i>	# of pairs mistakenly classified as different drugs
<i>Precision</i>	$TP/(TP + FP)$
<i>Recall</i>	$TP/(TP + FN)$
<i>ACC</i>	$(TP + TN)/(TP + TN + FP + FN)$
<i>F1</i>	$2 * Precision * Recall / (Precision + Recall)$

**3.1.2 Environmental and Parameter settings.** The experiments are conducted on a server with Intel i9-9900k CPU, GeForce GTX 2080 Ti Graphics Cards and 64 GB of RAM. We use Adaptive Moment Estimation (Adam) [27] to optimize our model with a learning rate of 0.0001. For each iteration, we use a batch size of 32 and set the sample number  $n_s = 16$ . We run  $n_G = 5$  iterations of generator training and  $n_D = 15$  iterations of discriminator training in each epoch. The embedding dimensions are set to be 150 for all methods.

**3.1.3 Evaluation metrics.** Based on the ground-truth and using 10-fold cross validation, we use the metrics in Table 2 to quantitatively assess effectiveness of different methods for drug identification.

## 3.2 Comparisons of Different Features

In this set of experiments, we comprehensively evaluate the effectiveness of different features described in Section 2.1 for drug identification in darknet markets.

- **Text-based features:** This type of features only consider a posted drug’s text information, including (1) *text content feature* (denoted as *f-1*): features are extracted based on the content of text-based information, which will be converted to a 100-dimensional feature vector; (2) *writing style feature* (denoted as *f-2*): features are extracted based on lexical, syntactic, and structure styles; (3) *f-3*: the concatenation of *f-1* and *f-2*.
- **Photo-based features:** This type of features only consider a posted drug’s photo information, including (1) *photo content feature* (denoted as *f-4*): photo represented by a converted 100-dimensional feature vector based on VGG-19 model; (2) *photography style feature* (denoted as *f-5*): both low-level and high-level photography styles are extracted from the posted photo; (3) *f-6*: the concatenation of *f-4* and *f-5*.
- **Augment features:** This is the concatenation of text-based features *f-3* and photo-based features *f-6*, denoted as *f-7*.

**Table 3: Comparisons of different features.**

Method	Feature	ACC	F1	Recall	Precision
Text-based	<i>f-1</i>	0.8164	0.7919	0.6986	0.9139
	<i>f-2</i>	0.8026	0.7734	0.6738	0.9076
	<i>f-3</i>	0.8247	0.7983	0.6940	0.9396
Photo-based	<i>f-4</i>	0.7343	0.6563	0.5074	0.9290
	<i>f-5</i>	0.7133	0.6208	0.4694	0.9164
	<i>f-6</i>	0.7380	0.6604	0.5096	0.9381
Augment	<i>f-7</i>	0.8440	0.8354	0.7720	0.9039
<i>dStyle-GAN</i>		<b>0.8930</b>	<b>0.8844</b>	<b>0.8190</b>	<b>0.9615</b>

Based on the above features, we use the DNN classifier described in Section 2.4 to evaluate their performances by comparison with our developed *dStyle-GAN* in drug identification. The experimental results are shown in Table 3. From the results, we can see that: (1) the posted drugs in darknet markets represented by content-based features (i.e., *f-1* and *f-4*) outperform their corresponding content-free features (i.e., *f-2* and *f-5*); (2) the integration of writing or photography styles helps the performance for drug identification (i.e., *f-3* outperforms *f-1* and *f-6* outperforms *f-4* respectively); (3) text-based feature (*f-3*) in general performs better than photo-based feature (*f-6*), which may be because some posted drugs do not have photos and most posted drugs have detailed text description; (4) with comparable precision, the integration of all features (*f-7*) outperforms the other six features (*f-1* to *f-6*) in terms of accuracy, F1 and recall for drug identification (i.e., the augment features can help identify more pairs of posts that relate to same drugs); (5) compared with the concatenation of all features (i.e., *f-7*), our developed system *dStyle-GAN* achieves significant performance improvement

**Table 4: Comparisons with different node representation learning models for drug identification.**

Index	Model	10%	20%	30%	40%	50%	60%	70%	80%	90%
ACC	DeepWalk	0.6539	0.7208	0.7240	0.7691	0.7795	0.8031	0.8031	0.8189	0.8265
	LINE	0.6692	0.7340	0.7267	0.7764	0.7852	0.8097	0.8197	0.8259	0.8362
	metapath2vec	0.6725	0.7408	0.7430	0.7885	0.8046	0.8188	0.8322	0.8339	0.8533
	Hin2Vec	0.6656	0.7368	0.7392	0.7816	0.7980	0.8133	0.8271	0.8275	0.8473
	GraphGAN	0.6615	0.7315	0.7342	0.7760	0.7912	0.8075	0.8210	0.8225	0.8409
	HeGAN	0.6753	0.7511	0.7608	0.7993	0.8242	0.8335	0.8401	0.8415	0.8581
	<i>dStyle-GAN</i>	<b>0.7130</b>	<b>0.7811</b>	<b>0.7990</b>	<b>0.8410</b>	<b>0.8560</b>	<b>0.8590</b>	<b>0.8670</b>	<b>0.8760</b>	<b>0.8930</b>
F1	DeepWalk	0.6326	0.6972	0.7095	0.7554	0.7614	0.7901	0.7817	0.8069	0.8081
	LINE	0.6333	0.7083	0.7080	0.7570	0.7715	0.7932	0.8061	0.8112	0.8246
	metapath2vec	0.6401	0.7231	0.7237	0.7650	0.7834	0.8076	0.8226	0.8175	0.8374
	Hin2Vec	0.6397	0.7212	0.7177	0.7636	0.7780	0.7948	0.8079	0.8145	0.8369
	GraphGAN	0.6280	0.7125	0.7202	0.7622	0.7783	0.7889	0.8116	0.8124	0.8259
	HeGAN	0.6546	0.7363	0.7428	0.7869	0.8081	0.8221	0.8269	0.8314	0.8481
	<i>dStyle-GAN</i>	<b>0.6874</b>	<b>0.7568</b>	<b>0.7853</b>	<b>0.8245</b>	<b>0.8418</b>	<b>0.8464</b>	<b>0.8600</b>	<b>0.8675</b>	<b>0.8854</b>

for drug identification. The system still has false positives and false negatives (e.g., due to missing photos or posted texts with limited information for drug description), however, the comparison results have demonstrated the superiority of the developed *dStyle-GAN* for drug identification in darknet markets.

### 3.3 Comparisons with Different Models for Drug Representation Learning

In this set of experiments, based on the constructed AHIN, we evaluate our proposed GAN-based representation learning model over AHIN by comparing with other state-of-the-art embedding models for drug identification. The baseline methods include:

- **DeepWalk** [33] learns node vectors by capturing node pairs within  $w$ -hop neighborhood via uniform random walks in homogeneous network.
- **LINE** [41] preserves the first-order and second-order proximities between nodes for homogeneous network embedding.
- **metapath2vec** [17] is proposed for heterogeneous network which embeds the semantic information of a single meta-path.
- **Hin2Vec** [20] learns node embedding to capture rich relation semantics in heterogeneous network via a neural network model.
- **GraphGAN** [46] unifies generative and discriminative learning of GAN to learn node embeddings in homogeneous graph.
- **HeGAN** [25] applies the adversarial learning model to heterogeneous graph for node representation learning.

For homogeneous network embedding models (i.e., DeepWalk, LINE, GraphGAN), we ignore the heterogeneous property and directly feed the constructed AHIN for node representation learning; while for metapath2vec, we concatenate all the embeddings generated by different meta-paths to obtain node representations. Since all baselines are incapable of dealing with the attributes attached to the nodes (i.e., drugs), we concatenate the attributed feature vector with the learned embedding for each drug which is fed the DNN classifier described in Section 2.4 for drug identification. To facilitate the comparisons, we follow the same experimental procedure in [17] to perform the studies: a portion of training data (ranging

from 10% to 90%) is randomly selected for training while the remaining ones are used for testing. The experimental results are shown in Table 4, from which we can see that: (1) HIN embedding methods (i.e., metapath2vec, Hin2Vec, HeGAN) yield better performances than homogeneous network embedding models (i.e., DeepWalk, LINE, GraphGAN), which implies that HIN embedding methods are able to preserve rich semantic information of the data and thus learn better node representations; (2) exploring adversarial learning mechanism for network representation learning indeed improves the quality of learned node representations (i.e., GraphGAN vs. DeepWalk/LINE, HeGAN vs. metapath2vec/Hin2Vec); (3) our proposed GAN-based model over AHIN integrated in the developed *dStyle-GAN* consistently outperforms all baseline embedding models in terms of both accuracy (i.e., ACC) and F1, whose success lies in it jointly considers the heterogeneity of network and relatedness over drugs formulated by domain-specific meta-paths to devise a principled GAN-based node representation learning framework for drug identification.

### 3.4 Comparisons with Competing Approaches

In this set of experiments, we compare *dStyle-GAN* with several competing approaches, including traditional machine learning-based methods and an existing AHIN-based system.

**Table 5: Comparisons with competing approaches.**

Feature	Model	ACC	F1	Recall	Precision
Augment	DT	0.8354	0.8322	0.7604	0.8867
	SVM	0.8366	0.8323	0.7645	0.8983
	DNN	0.8440	0.8354	0.7720	0.9039
uStyle-uID		0.8610	0.8512	0.7950	0.9159
<i>dStyle-GAN</i>		<b>0.8930</b>	<b>0.8844</b>	<b>0.8190</b>	<b>0.9615</b>

For traditional machine learning-based methods, using the augmented features (i.e.,  $f-7$ ) as the inputs, we consider typical classification models of Decision Tree (DT) and Support Vector Machine



(SVM) and compare them with the DNN classifier described in Section 2.4. We also compare our system *dStyle-GAN* with the existing system *uStyle-uID* [51]. In *uStyle-uID*, we redefine the entities and relations for AHIN and rebuilt meta-paths in *uStyle-uID* to learn node embeddings for drug identification. The experimental results are shown in Table 5, from which we observe that: (1) *uStyle-uID* with additional knowledge represented by AHIN performs better than traditional methods (i.e., DT, SVM and DNN based on augment features). Such finding shows that meta-path based approach over AHIN can more effectively exploit higher-level semantic connections among drugs. (2) Our developed system *dStyle-GAN* significantly outperforms other competing approaches in drug identification, which reveals that our model could learn more robust semantic node (i.e, drug) representations through GAN-based adversarial framework over AHIN.

### 3.5 Parameter Sensitivity and Scalability

In this set of experiments, we first perform a series of **parameter sensitivity** analysis to investigate how different parameter choices will affect the performance of *dStyle-GAN*. Figure 5.(a) shows the performances in terms of different embedding dimensions, from which we can observe that: increasing dimension  $d$  initially boosts the performance since a larger  $d$  can encode more information, while the performance tends to be stable when  $d = 150$ . Then, we vary the variance of Gaussian distribution  $\sigma^2$  in *dStyle-GAN* from 0.001 to 10 to see their influences on the performance. From Figure 5.(b), we can see that our model achieves the optimal performance when  $\sigma^2 = 1.0$  and is steady around the value. We further set different numbers of iterations in each epoch for the training of generators and discriminators from 1 to 20 and plot their performances in Figure 5.(c). We notice that *dStyle-GAN* gains the optimal performance near  $n_G = 10$  and  $n_D = 20$ . Overall, our proposed system *dStyle-GAN* is not strictly sensitive to the parameters and is able to reach high performance under a cost-effective parameter choice. Then, we further evaluate the **scalability** of *dStyle-GAN*. Figure 5.(d) shows the running time of *dStyle-GAN* with different sizes of the datasets, which illustrates that the running time is quadratic to the number of samples (i.e., nodes in the AHIN). From the results and analysis above, we conclude that *dStyle-GAN* is efficient and scalable for large-scale network representation learning.

## 4 DEEP INVESTIGATION BASED ON IDENTIFIED DRUG PAIRS IN THE WILD

To better understand and gain deeper insights into the ecosystem of illicit drug trafficking in darknet markets, we further apply our developed system *dStyle-GAN* for automatic drug identification across different markets in the wild and perform further analysis.

### 4.1 Dynamics and Evolution of Illicit Drugs

To understand the dynamics and evolution of illicit drugs traded in darknet markets, besides our obtained ground-truth, we further apply our developed system *dStyle-GAN* for automatic drug identification across different markets in the wild. To characterize the patterns how drugs are transferred from one market to the other, we define two types of drug transfers below:

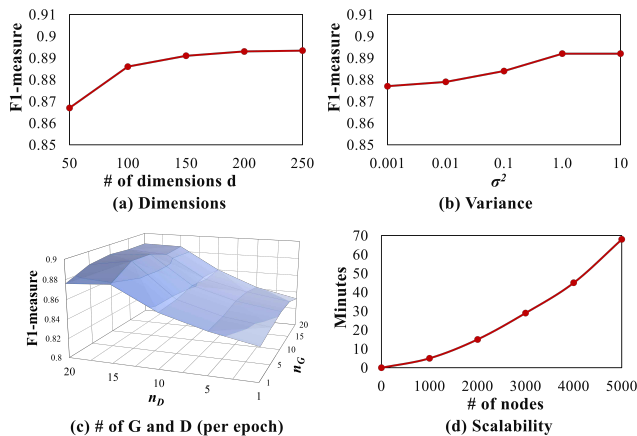


Figure 5: Evaluation of parameter sensitivity and scalability.

- **Vertical Transfer:** Drugs in a shutdown market are transferred to other markets that appears at a later time. For instance, some posted drugs in SR1 were transferred to SR2 after its shutdown.
- **Horizontal Transfer:** Drugs in a market are also advertised in other competitive markets that are active at the same time.

For cross-market drug identification task, based on our collected data summarized in Table 1, we initially generate drug pairs across two markets based on drug attribute information (i.e., pairs with identical drug category and mass attributes) and then exploit *dStyle-GAN* for drug identification (i.e., if a given pair of drugs are the same or not) across two markets. In Table 6, we report the identification results by *dStyle-GAN*.

Table 6: Drug identification across different markets.

Cross-market	# of Drug Pairs	# of Vendors	Transfer Type
SR1 → SR2	85	23	Vertical
SR1 → Agora	62	30	Vertical
SR2 → AlphaBay	530	54	Vertical
SR2 → Agora	4,646	391	Horizontal
Agora → AlphaBay	3,851	261	Horizontal
SR3 → Empire	246	24	Horizontal
Total	9,420	783	

Based on the results shown in Table 6, we perform further analysis and have following key findings:

- **There may not have data backups of posted drugs in SR1.** It has been found that, shortly after its shutdown, 85 posted drugs and 62 posted drugs in SR1 were transferred to SR2 and Agora respectively. Based on the limited number of transfer drugs (85/62 pairs) and vendors (23/30 vendors), we conclude that there may not have posted drug backups in SR1.
- **Agora could be one of the biggest competitors of SR2.** As time went by, darknet markets were booming again since 2014. Based on our data collection shown in Table 1, we can see that

the size of SR2 (13,202 posted drugs) was over 15 times bigger than that of SR1 (726 posted drugs). Meantime, the drug trafficking activities became unprecedentedly active among multiple darknet markets (SR2, Agora, and AlphaBay). As both Agora and SR2 were active at the same time, we consider the identified drug pairs across these two markets as horizontal transfer. As shown in Figure 6, there were 4,646 drugs that appeared in both SR2 and Agora from 2014 to 2015. After further analysis, among all these identified drug pairs, 1,978 drugs posted by 166 vendors were transferred from SR2 to Agora after its closure; this may imply that some vendors on SR2 chose Agora to continue drug trafficking activities. The above analysis suggests that Agora could be one of the biggest competing markets of SR2.

- **The illicit drug trafficking activities in darknet markets were unprecedentedly active during 2014-2015.** As shown in Figure 6, AlphaBay appeared shortly after SR2 was shutdown by law enforcement. Based on the drug identification results by *dStyle-GAN*, there were 530 drugs posted by 54 vendors transferred from SR2 to AlphaBay (i.e., vertical transfer). In 2015, AlphaBay became a major competitor of Agora and 3,851 drugs appeared in both Agora and AlphaBay. More recently, the activities of drug transfer from SR3 to Empire are noticeably lower; the reason behind this may be SR3 had become inactive due to the persistent effort by law enforcement in 2019. In sum, the discovered transfer patterns based on identified drug pairs across markets reveal that illicit drug trafficking activities in darknet markets were unprecedentedly active during the heyday period of 2014-2015. We also note that, due to considerable profits, activities of illicit drug trafficking in darknet markets are nimble and resilient (e.g., following the shutdown of Dream Market in May 2019, a dozen more have emerged and remained active such as Dark Market [2], BitBazaar Market [1], Monopoly Market [4]).

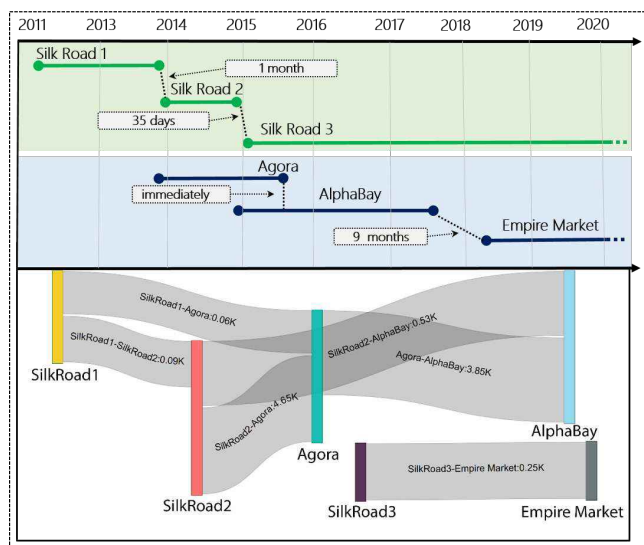


Figure 6: Dynamics and evolution of illicit drug trafficking.

## 4.2 Case Studies of Identified Drug Pairs

To gain deeper insight into the operations of illicit drug trades, we further analyze the transfer activities based on the identified drug pairs by *dStyle-GAN*. Figure 7 illustrates that a posted drug is changed in some way when transferred from one market to another, which we call “altered presentation”. More specifically, it can be seen from Figure 7 that the writing and photography styles of identified drug pairs are identical although the text and photo information have changed to some degree (e.g., from “brand new” to “new product”). Many of the vertical and horizontal transfer activities possess such altered presentation characteristics. Among popular drug presentations, those most likely to change include price, shipping methods, text description, and photos. Based on the identified drug pairs across different markets, we attempt to provide some plausible explanations for alter presentation below.

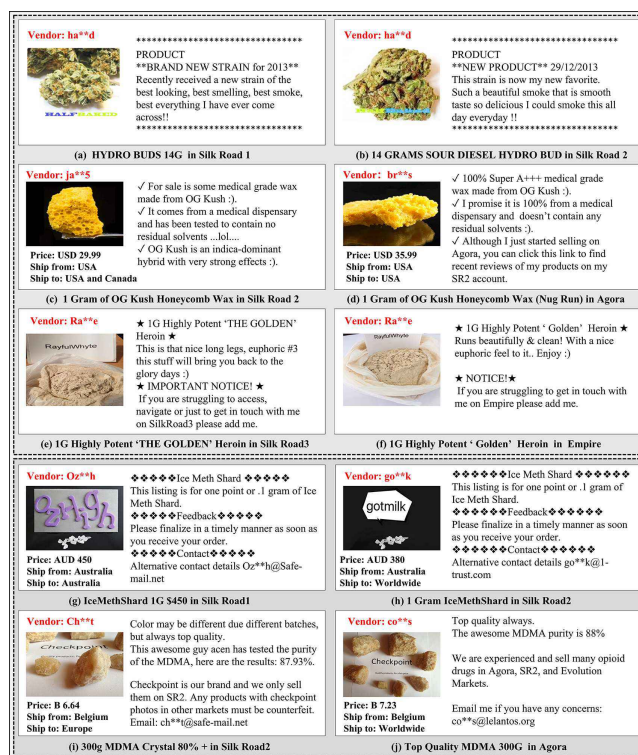


Figure 7: Examples of identified drug pairs by *dStyle-GAN*.

- **Presentation gets altered due to lack of data backups.** As illustrated in Figure 7.(a) and (b), when there is no backup of posted drugs (e.g., in SR1), re-posting those drugs on a new market may involve the effort from scratch. Vendors might not be able to reproduce the posted information exactly but their writing and photography styles would remain invariant. This line of reasoning strongly justifies the rationale behind our style-based approach in this work.
- **Altered presentation is for the promotion of illicit drugs.** After inspecting the drug pairs identified by *dStyle-GAN*, we find that the underlying reasons of altered presentation include price

or shipping information changes, enhancement of text descriptions and posted photos for *product promotion* (as illustrated in Figure 7.(c) and (d)). With further analysis, we also observe that the shipping origins of about 95% of identified drug pairs stayed same across different markets; while the shipping to information of over 40% of identified drug pairs were altered to some degree. As the amount of identified drug pairs with altered shipping origin is relatively low (5%), we conclude that *the supply chain of drugs was relatively steady*. In view of such apparent discrepancy between shipping origins and sales regions, we further conclude that the most effective countermeasures against drug markets in darknet is to cut down the sources of illicit drug supplies. Since drug traffickers in darknet markets have to prove the possession of illicit drugs, their distinct photography styles might be revealed by the posted photos even with altered presentation. For example, as shown in Figure 7.(e) and (f), both posted photos of heroin are with the distinct photography style to demonstrate that the drugs are owned by vendor “Ra\*\*e”; in this way, loyal customers can recognize him easily. In other words, the unique styles can be interpreted as *an implicit strategy of advertisement* because the key objective of drug trafficking is to profit. Such finding again implies that style-based information could provide a reliable clue for cross-market drug identification.

- Some vendors alter presentation due to scams or disputes.** We have also found that some of the identified drug pairs by *dStyle-GAN* were posted by different vendors across markets. For example, as illustrated in Figure 7.(g) and (h), the Ice Meth Shard was posted by “Oz\*\*h” in SR1 and “go\*\*k” in SR2; as shown in Figure 7.(i) and (j), the MDMA was posted by “Ch\*\*t” in SR2 and “co\*\*s” in Agora. For the latter example, we can see that the one posted by “co\*\*s” in Agora might be scam, as it’s stated in the original post by “Ch\*\*t” that “*Checkpoint is our brand and we only sell them on SR2. Any products with checkpoint photos in other markets must be counterfeit*”. Through the identified drug pairs, it has also been found that some vendors had disputes with customers (as shown in Figure 8), which might force the transfer to a new market. In the such situation, altered presentation would represent a minimum-effort solution.

 <p><b>Vendor: Su**s</b>  <b>Price: USD 50</b>  <b>Ship from: USA</b>  <b>Ship to: USA</b></p> <p>** Ship to US and always has a good quality.            ** Provide from 1 gram to quarter pound cocaine.            ** Read bio for more info. We will ship the coke via fast 3-5 day USPS.</p>	 <p><b>Vendor: Su**s</b>  <b>Price: USD 50</b>  <b>Ship from: USA</b>  <b>Ship to: USA</b></p> <p>**This cocaine has maintained good purity since being imported into the US.            ** We provide between one gram to one quarter pound orders of this amazing crack.            **Fast 3-5 day Priority USPS shipping.</p>
<p><b>High Quality Cocaine 1 Gram in Silk Road3</b>    <b>1 Gram (G) High Quality Base Cocaine in Empire</b></p>	
<p><b>Title:</b> Su**s on Empire is scam. [12:19 AM 01-16-2020]  <b>Poster:</b> po**y</p> <p>I ordered a gram of Crack from Su**s on empire last week. He required FE but had really good feedback so I ordered. Pack came today and was a rolled up piece of bubble wrap. I go to message him and notice that he's had like 5 negative feedback left over the past couple days saying that they all got rolls of bubble wrap too. You've been warned. Also, he was a scam vendor in SR3. I met with a dude who had disputes with the fucking vendor in SR3.</p>	
<p><b>Complaint threads posted by po**y in DMV Avengers forum</b></p>	

**Figure 8: An example dispute.**

## 5 RELATED WORK

To combat illicit drug trafficking in darknet markets, there have been many research efforts on darknet market data analysis: for examples, [6, 11, 15, 28, 42, 47, 51] focus on drug trafficker identification and trafficking network investigation; while [7–10, 12, 14, 16, 44] explore statistical methodologies to analyze illicit drugs traded in the markets. Since the identification and analysis of illicit drugs can not only provide valuable insights to profile drug traffickers but also enable the deep understanding of dynamics and evolution of drug trafficking activities, in this work, we focus on the illicit drug identification and investigation. To perform the analysis of large-scale posted illicit drugs while tackling the challenge of evolving dynamics in darknet markets (i.e., drugs traded cross-market with altered presentation), in this work, we propose a principled GAN-based framework leveraging both writing and photography styles for drug identification in darknet markets.

In order to depict drugs, vendors, posted texts and writing styles, photos and photography styles, and the rich relations among them, it is important to model them properly in order to facilitate the task of drug identification. To solve this problem, we present a powerful AHIN [30] model for abstract representation [18, 19, 23, 24, 48, 49, 52]. Based on the constructed AHIN, various network embedding methods (e.g., metapath2vec [17], HIN2Vec [20]) have been proposed to solve the node representation learning problem; to learn more robust node representations over graph, GAN-based models (e.g., GraphGAN [46], HeGAN [25]) have been further exploited. Different from the existing works, to learn node (i.e., drug) representations for drug identification in darknet markets, we jointly consider the heterogeneity of network and relatedness over drugs formulated by domain-specific meta-paths to devise a principled GAN-based model for robust node representation learning at the first attempt.

## 6 CONCLUSION

In the fight against illicit drug trafficking in darknet markets, in this paper, we have designed and developed an intelligent system called *dStyle-GAN* to automate the analysis for drug identification. In *dStyle-GAN*, based on the large-scale data collected from darknet markets, we first present an AHIN to depict drugs, vendors, texts and writing styles, photos and photography styles, and the rich relations among them; and then we propose a novel GAN-based model to jointly consider the heterogeneity of network and relatedness over drugs formulated by domain-specific meta-paths to learn robust node (i.e., drug) representations for their identifications. To the best of our knowledge, the proposed *dStyle-GAN* represents the first principled GAN-based solution over graphs to simultaneously consider writing and photography styles as well as their latent distributions for node representation learning. Promising experimental results on the large-scale data collected from six darknet markets demonstrate that *dStyle-GAN* which integrates our proposed method outperforms the state-of-the-art approaches. Based on the identified drug pairs in the wild by *dStyle-GAN*, we perform further analysis to gain deep insights into the dynamics and evolution of illicit drug trading activities in darknet markets, whose findings may facilitate law enforcement for proactive interventions.



## ACKNOWLEDGMENTS

Y. Zhang, Y. Qian, Y. Fan and Y. Ye's work is partially supported by the NSF under grants IIS-2027127, IIS-2040144, IIS-1951504, CNS-2034470, CNS-1940859, CNS-1946327, CNS-1814825, OAC-1940855 and ECCS-2026612, the DoJ/NIJ under grant NIJ 2018-75-CX-0032.

## REFERENCES

- [1] 2019. BitBazaar Market. [jcbfbbyqhwmslo.onion](https://www.bitbazaar.com/).
- [2] 2019. Dark Market. [darkevuygggkqkqh.onion](https://www.darkmarket.com/).
- [3] 2019. Do Mexican drug cartels make 500 dollars billion a year? <https://www.washingtonpost.com/politics/2019/06/24/do-mexican-drug-cartels-make-billion-year/>.
- [4] 2019. Monopoly Market. [monopolyberbucxu.onion](https://www.monopoly.com/).
- [5] 2020. Darknet Market Archives. <https://www.gwern.net/DNM-archives/#download>.
- [6] Sadia Afroz, Aylin Caliskan Islam, Ariel Stolerman, Rachel Greenstadt, and Damon McCoy. 2014. Doppelganger finder: Taking stylometry to the underground. In *S&P. IEEE*, 212–226.
- [7] Judith Aldridge and David Décarry-Héту. 2014. Not an 'Ebay for Drugs': The Cryptomarket 'Silk Road' as a paradigm shifting criminal innovation. *SSRN 2436643* (2014).
- [8] Judith Aldridge and David Décarry-Héту. 2016. Hidden wholesale: The drug diffusing capacity of online drug cryptomarkets. *International Journal of Drug Policy* 35 (2016), 7–15.
- [9] Matthew Ball, Roderic Broadhurst, and Harshit Trivedi. 2019. How much fentanyl is available on the darknet? *SSRN 3362149* (2019).
- [10] Monica J Barratt. 2012. Silk Road: eBay for drugs. *Addiction* 107, 3 (2012), 683–683.
- [11] David Bright, Johan Koskinen, and Aili Malm. 2019. Illicit network dynamics: The formation and evolution of a drug trafficking network. *Journal of Quantitative Criminology* 35, 2 (2019), 237–258.
- [12] Julian Broseus, Damien Rhumorbarbe, Caroline Mireault, Vincent Ouellette, Frank Crispino, and David Décarry-Héту. 2016. Studying illicit drug trafficking on Darknet markets: structure and organisation from a Canadian perspective. *Forensic science international* 264 (2016), 7–14.
- [13] Tiago Carvalho, Edmar RS De Rezende, Matheus TP Alves, Fernanda KC Balieiro, and Ricardo B Sovat. 2017. Exposing computer generated images by eye's region classification via transfer learning of VGG19 CNN. In *ICMLA. IEEE*, 866–870.
- [14] Nicolas Christin. 2013. Traveling the Silk Road: A measurement analysis of a large anonymous online marketplace. In *WWW*, 213–224.
- [15] Jakob Demant, Rasmus Munksgaard, and Esben Houborg. 2018. Personal use, social supply or redistribution? Cryptomarket demand on Silk Road 2 and Agora. *Trends in Organized Crime* 21, 1 (2018), 42–61.
- [16] Diana S Dolliver. 2015. Evaluating drug trafficking on the Tor Network: Silk Road 2, the sequel. *International Journal of Drug Policy* 26, 11 (2015), 1113–1123.
- [17] Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. 2017. *metapath2vec*: Scalable representation learning for heterogeneous networks. In *KDD*, 135–144.
- [18] Yujie Fan, Yiming Zhang, Yanfang Ye, and Xin Li. 2018. Automatic Opioid User Detection from Twitter: Transductive Ensemble Built on Different Meta-graph Based Similarities over Heterogeneous Information Network.. In *IJCAI*, 3357–3363.
- [19] Yujie Fan, Yiming Zhang, Yanfang Ye, Xin Li, and Wanhong Zheng. 2017. Social media for opioid addiction epidemiology: Automatic detection of opioid addicts from twitter and case studies. In *CIKM*, 1259–1267.
- [20] Tao-yang Fu, Wang-Chien Lee, and Zhen Lei. 2017. *Hin2vec*: Explore meta-paths in heterogeneous information networks for representation learning. In *CIKM*, 1797–1806.
- [21] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *NIPS*, 2672–2680.
- [22] P Griffiths and J Mountney. 2017. Disruptive potential of the internet to transform illicit drug markets and impact on future patterns of drug consumption. *Clinical Pharmacology & Therapeutics* 101, 2 (2017), 176–178.
- [23] Shifu Hou, Yujie Fan, Yiming Zhang, Yanfang Ye, Jingwei Lei, Wenqiang Wan, Jiabin Wang, Qi Xiong, and Fudong Shao. 2019.  $\alpha$ Cyber: Enhancing Robustness of Android Malware Detection System against Adversarial Attacks on Heterogeneous Graph based Model. In *CIKM*, 609–618.
- [24] Shifu Hou, Yanfang Ye, Yangqiu Song, and Melih Abdulhayoglu. 2017. Hindroid: An intelligent android malware detection system based on structured heterogeneous information network. In *KDD*, 1507–1515.
- [25] Binbin Hu, Yuan Fang, and Chuan Shi. 2019. Adversarial Learning on Heterogeneous Information Networks. In *KDD*, 120–129.
- [26] Yan Ke, Xiaou Tang, and Feng Jing. 2006. The design of high-level features for photo quality assessment. In *CVPR*, Vol. 1. IEEE, 419–426.
- [27] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [28] Ramnath Kumar, Shweta Yadav, Raminta Daniulaityte, Francois Lamy, Krishnaprasad Thirunarayan, Usha Lokala, and Amit Sheth. 2020. eDarkFind: Unsupervised Multi-view Learning for Sybil Account Detection. In *WWW*, 1955–1965.
- [29] Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *ICML*, 1188–1196.
- [30] Xiang Li, Yao Wu, Martin Ester, Ben Kao, Xin Wang, and Yudian Zheng. 2017. Semi-supervised clustering in attributed heterogeneous information networks. In *WWW*, 1621–1629.
- [31] Lena Lundholm, Ulrika Haggård, Jette Möller, Johan Hallqvist, and Ingemar Thiblin. 2013. The triggering effect of alcohol and illicit drugs on violent crime in a remand prison population: a case crossover study. *Drug and alcohol dependence* 129, 1-2 (2013), 110–115.
- [32] James Martin. 2014. *Drugs on the dark net: How cryptomarkets are transforming the global trade in illicit drugs*. Springer.
- [33] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *KDD*, 701–710.
- [34] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. In *arXiv preprint arXiv:1511.06434*.
- [35] Rose A Rudd, Noah Aleshire, Jon E Zibbell, and R Matthew Gladden. 2016. Increases in drug and opioid overdose deaths United States, 2000–2014. *Morbidity and mortality weekly report* 64, 50 & 51 (2016), 1378–1382.
- [36] Yu Shi, Fangqiu Han, Xinwei He, Xinran He, Carl Yang, Jie Luo, and Jiawei Han. 2018. *mvn2vec*: Preservation and collaboration in multi-view network embedding. *arXiv preprint arXiv:1801.06597* (2018).
- [37] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [38] Jong Hwan Suh. 2016. Comparing writing style feature-based classification methods for estimating user reputations in social media. *SpringerPlus* 5, 1 (2016), 261.
- [39] Yizhou Sun and Jiawei Han. 2012. Mining heterogeneous information networks: principles and methodologies. *Synthesis Lectures on Data Mining and Knowledge Discovery* 3, 2 (2012), 1–159.
- [40] Yizhou Sun, Jiawei Han, Xifeng Yan, Philip S Yu, and Tianyi Wu. 2011. Pathsim: Meta path-based top-k similarity search in heterogeneous information networks. *VLDB* 4, 11 (2011), 992–1003.
- [41] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. 2015. Line: Large-scale information network embedding. In *WWW*, 1067–1077.
- [42] Meropi Tzanetakis. 2018. Comparing cryptomarkets for drugs: A characterisation of sellers and buyers over time. *International Journal of Drug Policy* 56 (2018), 176–186.
- [43] Joe Van Buskirk, Raimondo Bruno, Timothy Dobbins, Courtney Breen, Lucinda Burns, Sundresan Naicker, and Amanda Roxburgh. 2017. The recovery of online drug markets following law enforcement and other disruptions. *Drug and alcohol dependence* 173 (2017), 159–162.
- [44] Marie Claire Van Hout and Tim Bingham. 2014. Responsible vendors, intelligent consumers: Silk Road, the online revolution in drug trading. *International Journal of Drug Policy* 25, 2 (2014), 183–189.
- [45] Rolf van Wegberg and Thijmen Verburgh. 2018. Lost in the Dream? Measuring the effects of Operation Bayonet on vendors migrating to Dream Market. In *Evolution of the Darknet Workshop*, 1–5.
- [46] Hongwei Wang, Jia Wang, Jialin Wang, Miao Zhao, Weinan Zhang, Fuzheng Zhang, Xing Xie, and Minyi Guo. 2017. Graphgan: Graph representation learning with generative adversarial nets. *arXiv preprint arXiv:1711.08267* (2017).
- [47] Xiangwen Wang, Peng Peng, Chun Wang, and Gang Wang. 2018. You are your photographs: Detecting multiple identities of vendors in the darknet marketplaces. In *ASIACCS*, 431–442.
- [48] Yanfang Ye, Shifu Hou, Yujie Fan, Yiming Zhang, Yiyue Qian, Shiyu Sun, Qian Peng, Mingxuan Ju, Wei Song, and Kenneth Loparo. 2020.  $\alpha$ -Satellite: An AI-driven System and Benchmark Datasets for Dynamic COVID-19 Risk Assessment in the United States. *IEEE Journal of Biomedical and Health Informatics* (2020).
- [49] Yanfang Ye, Tao Li, Donald Adjeroh, and S Sitharama Iyengar. 2017. A survey on malware detection using data mining techniques. *ACM Computing Surveys (CSUR)* 50, 3 (2017), 1–40.
- [50] Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *NIPS*, 649–657.
- [51] Yiming Zhang, Yujie Fan, Wei Song, Shifu Hou, Yanfang Ye, Xin Li, Liang Zhao, Chuan Shi, Jiabin Wang, and Qi Xiong. 2019. Your style your identity: Leveraging writing and photography styles for drug trafficker identification in darknet markets over attributed heterogeneous information network. In *WWW*, 3448–3454.
- [52] Yiming Zhang, Yujie Fan, Yanfang Ye, Liang Zhao, and Chuan Shi. 2019. Key Player Identification in Underground Forums over Attributed Heterogeneous Information Network Embedding Framework. In *CIKM*, 549–558.
- [53] Rong Zheng, Jiexun Li, Hsinchun Chen, and Zan Huang. 2006. A framework for authorship identification of online messages: Writing-style features and classification techniques. *Journal of the American society for information science and technology* 57, 3 (2006), 378–393.