

# **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>

1. Department of CDS, Case Western Reserve University, OH, USA
2. Department of CSEE, West Virginia University, WV, USA
3. Tencent Security Lab, Tencent, Guangdong, China



**ACSAC 2020**

December 7-11, 2020 • Online

# Outline

- ❑ Introduction
- ❑ Proposed Method
- ❑ Experimental Results and Analysis
- ❑ Deep Investigation
- ❑ Related Work
- ❑ Summary of Our Work

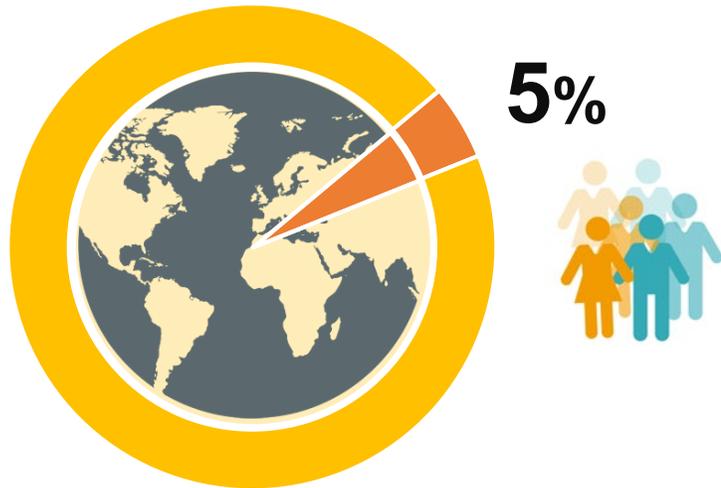
# Introduction

# Illegal Drugs The Facts

The demand for illegal drugs is solid and the trade in is considerably lucrative.

## CONSUMERS

5% of the world's adult population is estimated to have used illicit drugs at least once in 2017. (That's around 275 million people aged 15-64.)



**1 in 40**

People uses illicit drugs at least once a month

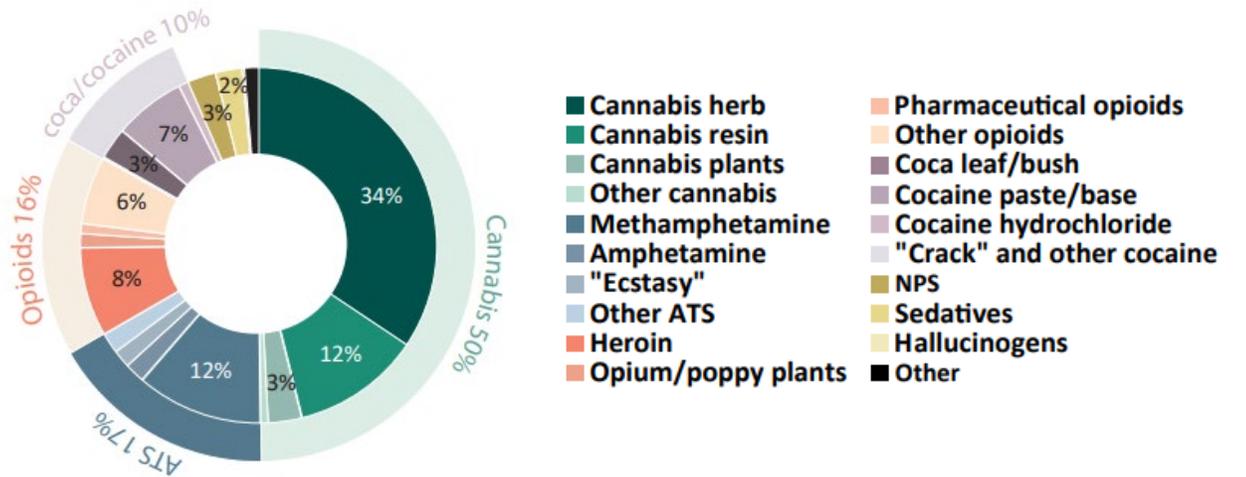
Source: World Drug Report 2018

## SUPPLIERS

**\$500 Billions**

the estimated revenue of drug traffic from Mexico alone coming across to the United States in 2018.

Global distribution of number of drug seizure cases, 2016–2017

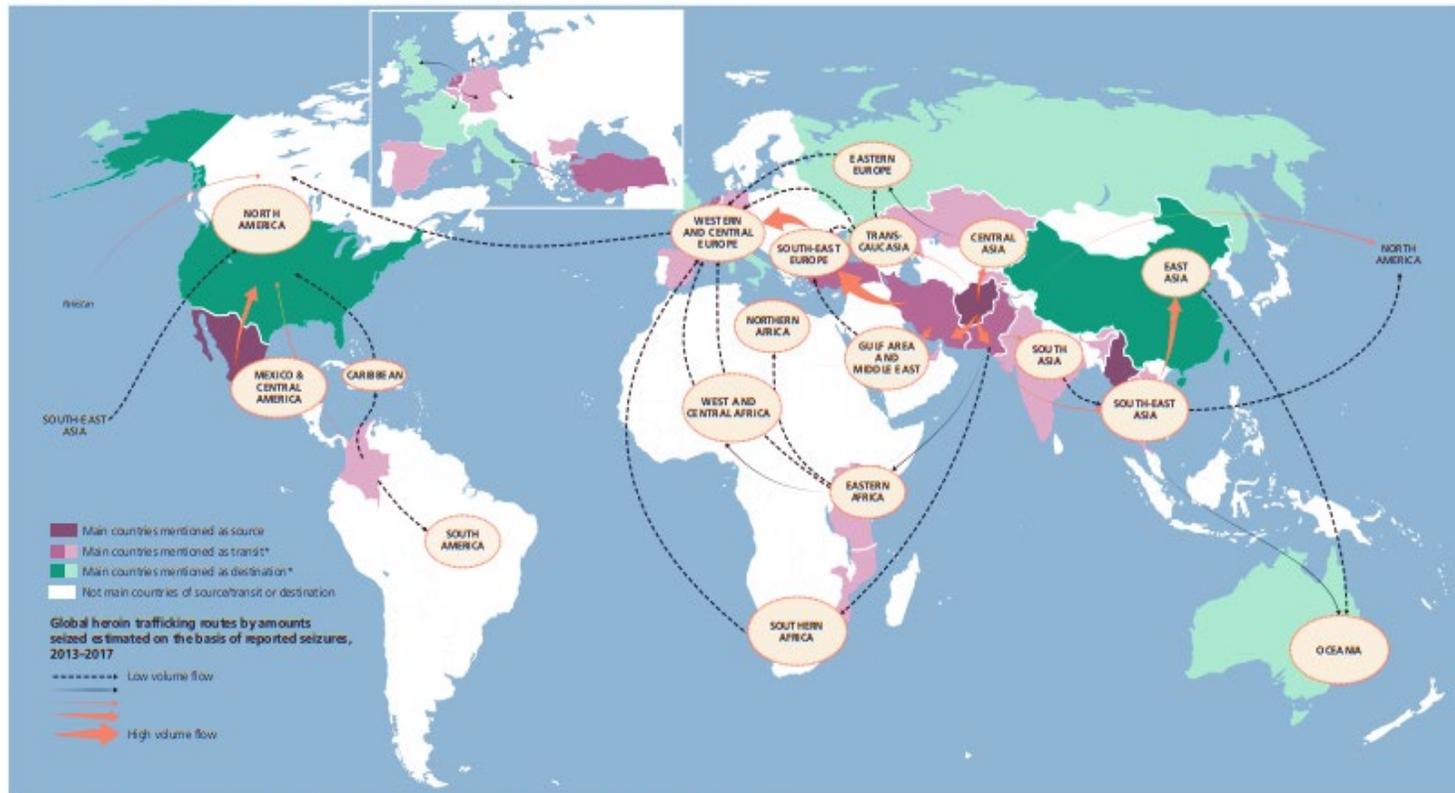


Source: UNODC, responses to the annual report questionnaire

# Drug trafficking

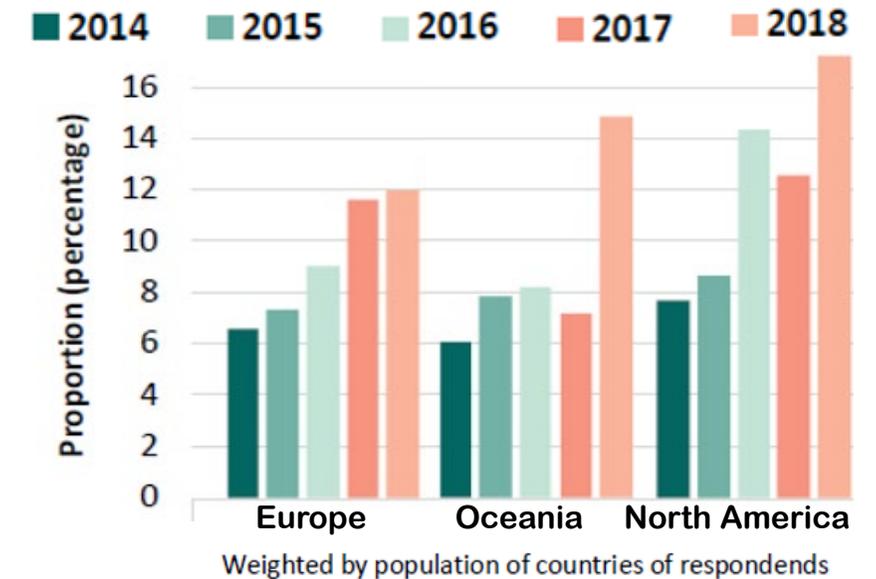
Driven by the remarkable profits, the crime of drug trafficking has never stopped but coevolved with the advance of modern technologies.

Main heroin trafficking routes as described by reported seizures, 2013–2017



Sources: UNODC, responses to the annual report questionnaire and individual drug seizure database.

Proportion of surveyed Internet users using drugs who purchased drugs over the darknet (2014–2018)

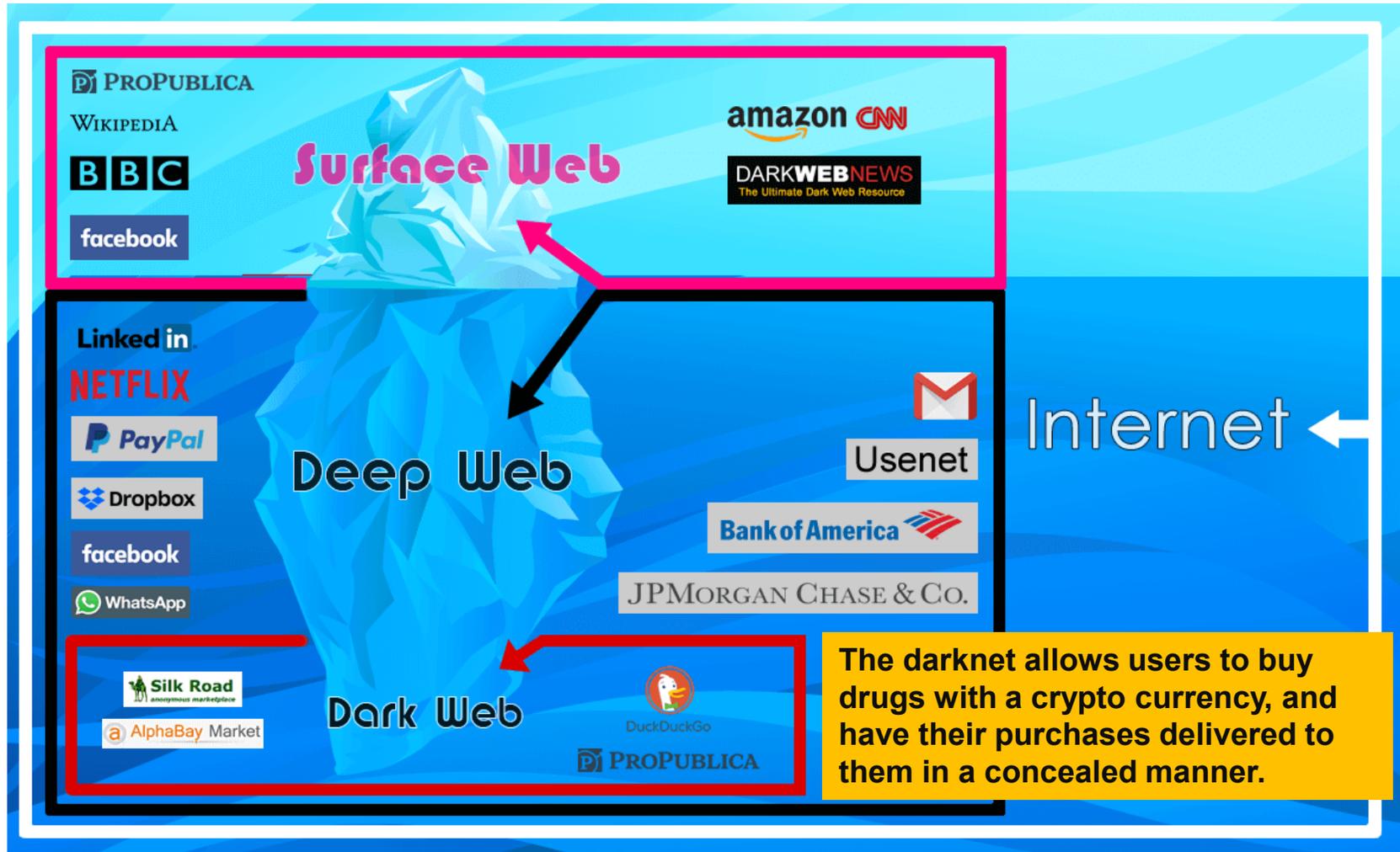


Darknet markets have been growing by around 50% in recent years.

Source: UNODC, responses to the annual report questionnaire

# Darknet

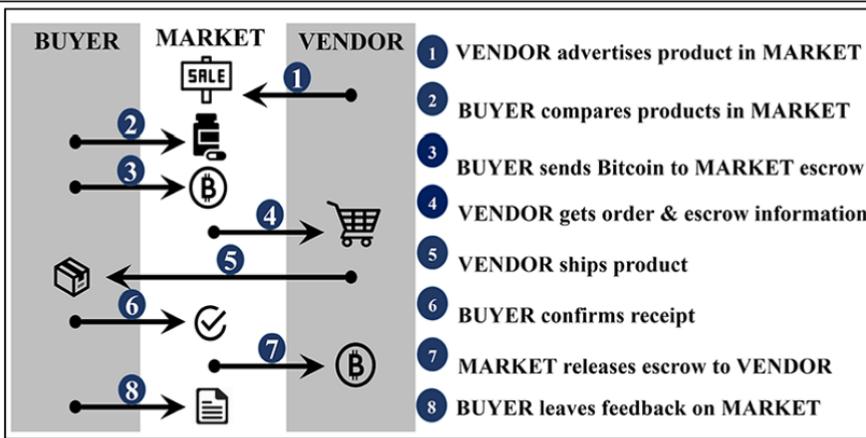
Darknet is the portions of the Internet purposefully not open to public view or hidden networks.



Markets on Darknet:

-  **Silk Road**  
*anonymous marketplace*
- Dream Market**  
*Ichudifyeqm4ldjj.onion*  
*Established 2013*
-  **Valhalla**
-  **TheRedDeal**
-  **Dr. D's**
-  **Empire Market**

# Drug trafficking in Darknet market



(a) A typical transaction in darknet market.

## Example writing styles:

- Having blank lines between paragraphs
- Using numbers (1, 2, 3) at the beginning of paragraphs
- Using special characters (\*\*\*\*) between paragraphs

## Example photography styles:



(b) Illustration of writing & photography styles.

**Empire Market Darknet Market**

Logged in as mxsjin  
BTC: 0.00000 / LTC: XMR: 0.00000 [My P]

USD 9684.14 CAD 13013.50 EUR 8597.87 AUD 13963.70 GBP 7638

HOME MESSAGES ORDERS BECOME A VENDOR BALANCE FEEDBACK FORUMS SUPPORT

Drugs & Chemicals Ecstasy MDMA Drug Category

MOONROCK MDMA 1G (EURO) HIGH QUALITY **Title**

MOONROCK MDMA 1G (EURO) EURO SASSAFRAS MDMA. TOP GRADE PURITY, VERY NICE HIGH, NO JAW CLENCHING.

Sold by DaShop - 368 sold since January 16, 2020 Vendor Level 5 Trust level 5 **Vendor Info**

Features		Features	
Product Class	Physical Package	Origin Country	United States
Quantity Left	Unlimited	Ships to	United States
Ends In	Never	Payment	Escrow

**Shipping & Escrow Info**

Priority - 3 days - USD + 8.00 / order

Purchase price: USD 44.44 **Price**

Qty: 1 Buy Now Buy Now Buy Now Queue

**Text Description**

0.004589 BTC / 0.963156 LTC / 0.653433 XMR

Description Feedback Return policy

**MOONROCK MDMA 1G (EURO) HIGH QUALITY**

MOONROCK MDMA 1G (EURO)

EURO SASSAFRAS MDMA.  
TOP GRADE PURITY, VERY NICE HIGH, NO JAW CLENCHING

Dashop  
Responsibility – Mutual Trust – Safe Operations.  
Welcome to the DaShop. Our team has been in the DaGame since Alpha, Dream, Silk 1 days. We are excited to have a presence on this platform as well. Our team is looking forward to providing the best shopping experience you can possibly have offering you the following top products.

1. COKE- Top quality white powder you can find.
2. MDMA and Ecstasy- Take your partying to the next level.
3. Adderol- Get your finals done in no time.

\*\*\*\*\*

EXCHANGE RATES

Bitcoin (BTC)	
United States Dollar (USD)	9684.14
Canadian Dollar (CAD)	13013.50
Euro (EUR)	8597.87
Australian Dollar (AUD)	13963.70

(c) A product (i.e., illicit drug) advertised in darknet market.

# System Architecture of our system **dStyle-GAN**

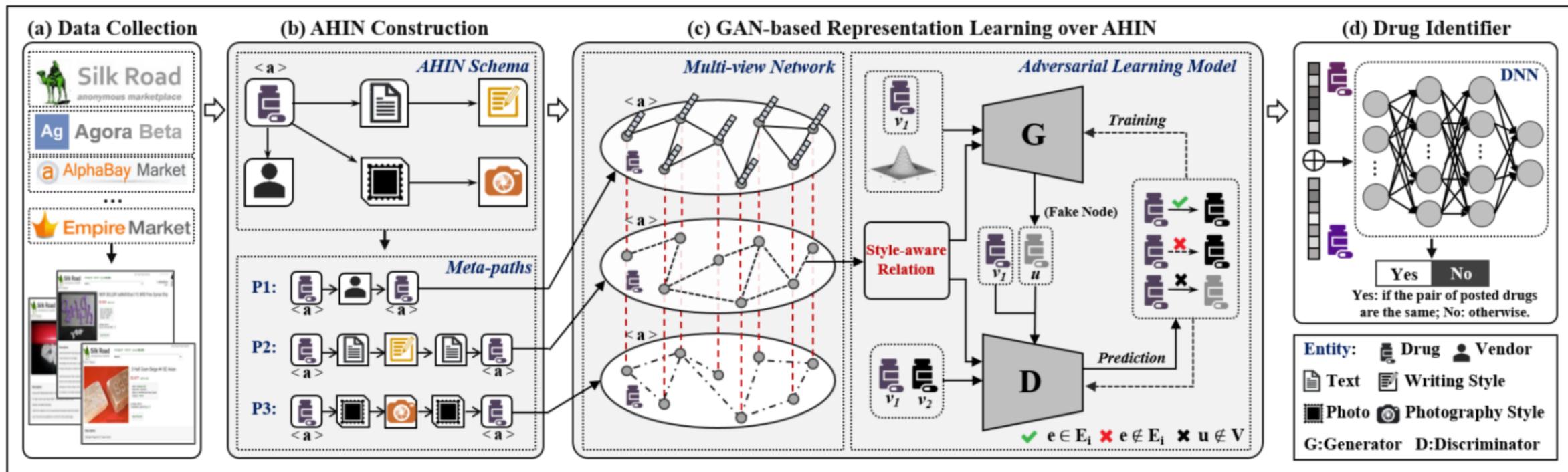


Figure 2: Framework of *dStyle-GAN*. In *dStyle-GAN*, (a) we first collect large-scale data from different markets; (b) after feature extraction, we construct an AHIN to model drugs, vendors, texts and writing styles, photos and photography styles, and the relations among them in a comprehensive manner; (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.

# Proposed Method

# Feature Extraction

## Content-based attributes:

For each posted drug:

- text content (title, drug description and terms)
- photo content
- category, Escrow, Shipping info and Price)

For each vendor:

- username, PGP key and contact information

## Photography styles:

Low-level style:

- camera make, model, camera angle, exposure time, focal length, and image size

High-level style:

- contrast, colorfulness, exposure of light

## Relation-based features:

R1 : vendor-sell-drug relation

R2 : text-describe-drug relation

R3 : photo-characterize-drug relation

R4 : text-have-WritingStyle relation

R5 : photo-contain-PhotographyStyle relation

## Writing styles:

Lexical style (Character-Level and Word-Level) :

- the frequency of characters
- total number of characters
- 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
- vocabulary richness

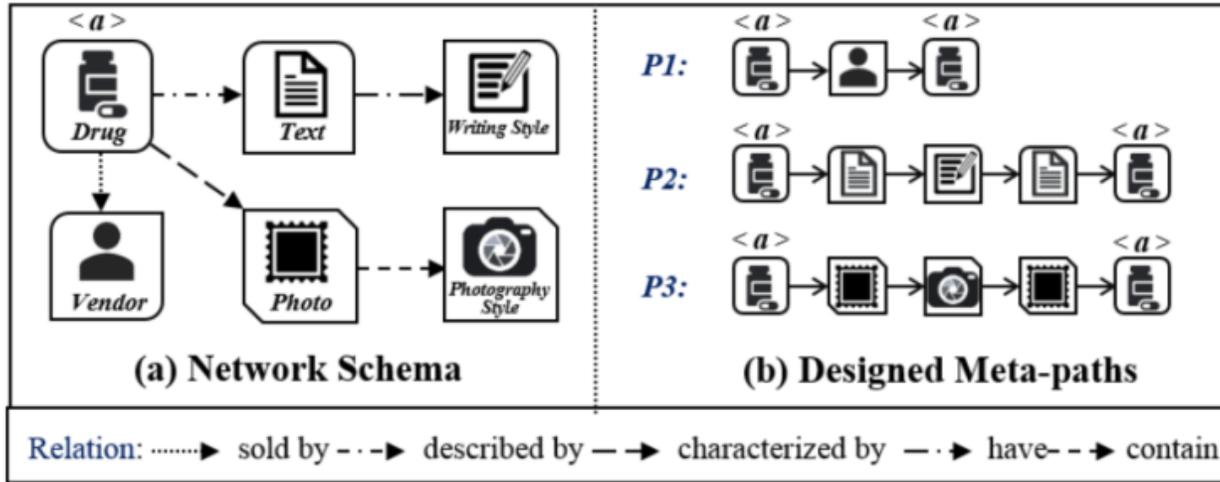
Syntactic style (sentence-level):

- frequency of punctuation
- frequency of function words
- frequency of stop words,
- number of sentences beginning with a capital letter

Structural style:

- total number of sentences
- whether there are URLs
- whether there are separators between paragraphs
- whether there are special characters between paragraphs

# AHIN Construction



**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  $\mathcal{X}_i$  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  $\mathcal{V} = \bigcup_{i=1}^m \mathcal{X}_i$  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  $\mathcal{R}$ .

## Network schema:

In this work, we have six types of entities: drug, vendor, text, photo, writing style, photography style and five types of relations between them.

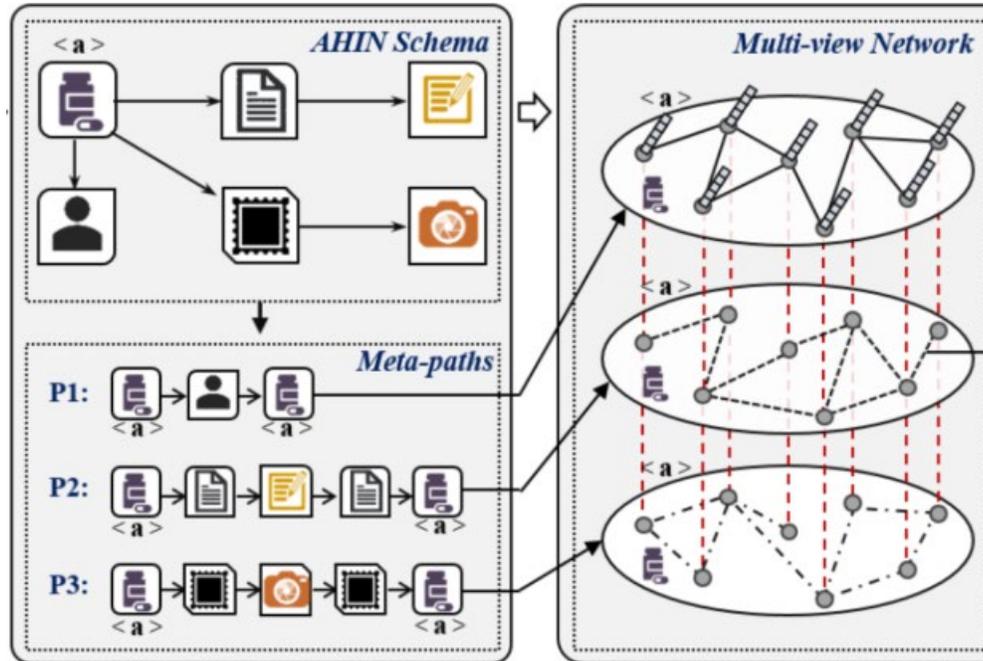
Moreover, each entity with the type of drug is attached with additional attributed feature vector.

## Meta-paths generation:

Based on the above network schema, we design three different meta-paths:

- P1 denotes two posted drugs are connected if they are sold by the same vendor;
- P2 describes two posted drugs are related if their text posts have a specific writing style (e.g., separators between paragraphs);
- P3 depicts two posted drugs are associated if their photos have a specific photography style (e.g., with identical camera angle).

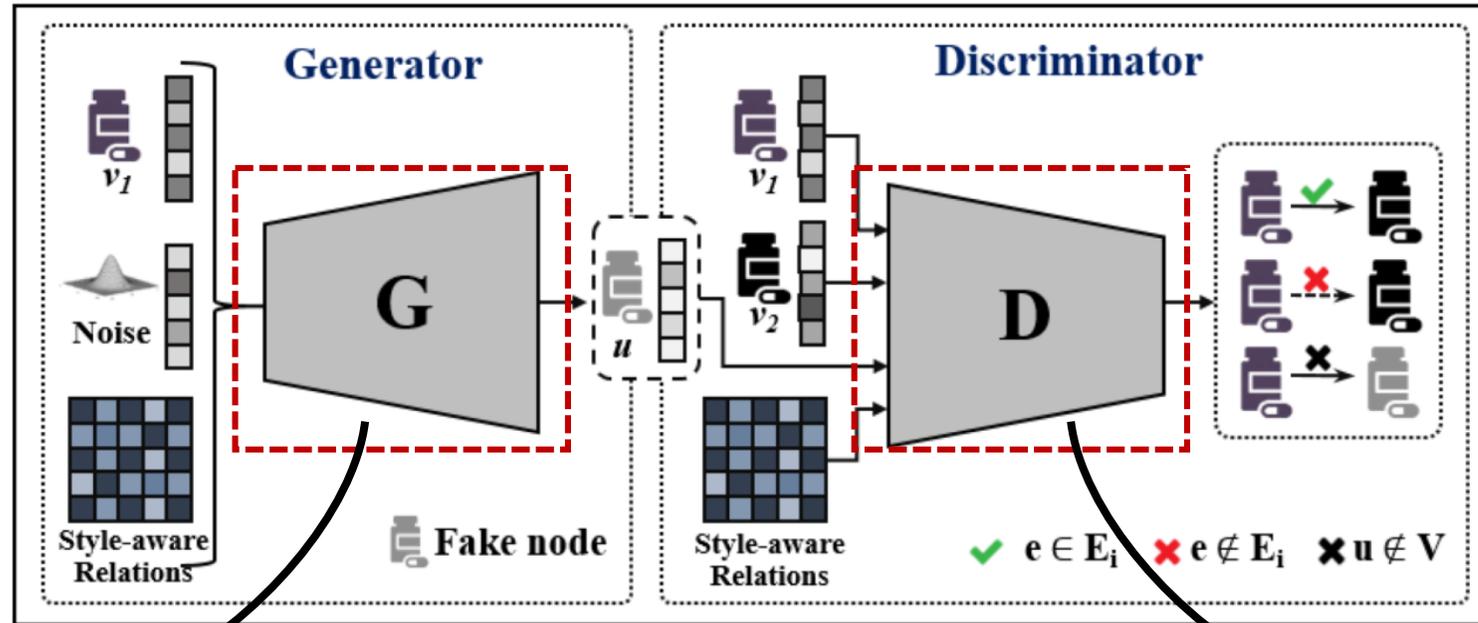
# GAN-based Representation Learning



*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 of multi-view network, 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.

# GAN-based Representation Learning



**Generator:**  $G(v, r_i; \Theta^G) = MLP(\mathbf{h}), \mathbf{h} \sim \mathcal{N}(\mathbf{h}_v^{G\top} \mathbf{M}_{r_i}^G, \sigma^2 \mathbf{I})$ .

generate synthetic (i.e., fake) nodes (i.e., drugs) that have similar writing and/or photography styles with posted drugs while jointly considering the latent distribution of style-aware information to fool the discriminator. Objective function of generator G is given by:

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

**Discriminator:**  $D(\mathbf{h}_u | v, r_i; \Theta^D) = \frac{1}{1 + \exp(-\mathbf{h}_v^{D\top} \mathbf{M}_{r_i}^D \mathbf{h}_u)}$ ,

differentiate the generated synthetic nodes (i.e., drugs) from the real posted drugs in darknet markets in terms of a specific style-aware relation.

Three types of objective functions:

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

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

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

# **Experimental Results and Analysis**

# Data Collection

To fully evaluate the our proposed method, we have collected the data from four different darknet markets:

Dataset	Time Frame	# of Drugs	# of Vendors
SilkRoad1(SR1)	07.03.2013	726	128
SilkRoad2(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
SilkRoad2(SR3)	04.27.2019 - 05.13.2019	5,745	371
Empire	11.22.2019- 11.24.2019	23,367	1,268

## Ground-truth:

We manually label the data on SilkRoad2 and Agora. Based on our annotation criteria, 1,961 pairs are finally labeled as positive (i.e., each pair of posts relate to the same drug). Accordingly, we randomly select 1,961 negative pairs (i.e., each pair of posts relate to different drugs). After feature extraction, the built AHIN contains 24,982 nodes and 572,236 edges.

## Performance Measures:

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>	$2 * Precision * Recall / (Precision + Recall)$

# E1: 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>

## Text-based features:

f-1: text content feature

f-2: writing style feature

f-3: the concatenation of f-1 and f-2

## Photo-based features:

f-4: photo content feature

f-5: photography style feature

f-6: the concatenation of f-4 and f-5

**Augment features(f-7):** the concatenation f-3 and f-6

## Results:

(1) the posted drugs 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 (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 );

(4) the integration of all features ( f-7 ) outperforms the other six features (f-1 to f-6);

(5) compared with the concatenation of all features, our developed system dStyle-GAN achieves significant performance improvement.

## E2: Comparisons with Different Representation Learning Models

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>

### Results:

- (1) HIN embedding methods yield better performances than homogeneous network embedding models;
- (2) adversarial learning mechanism for network representation indeed improves the quality of learned node representations;
- (3) our proposed system *dStyle-GAN* consistently outperforms all baseline methods, 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.

## E3: Comparisons with Alternative Approaches

### Results:

(1) **uStyle-uID with additional knowledge represented by AHIN performs better than traditional methods.** Such finding shows that meta-path based approach over AHIN can more effectively exploit higher-level semantic connections among drugs.

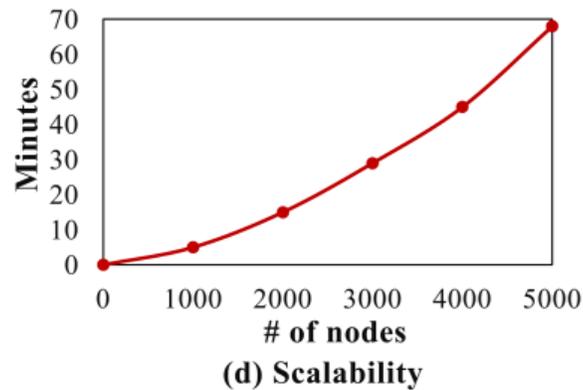
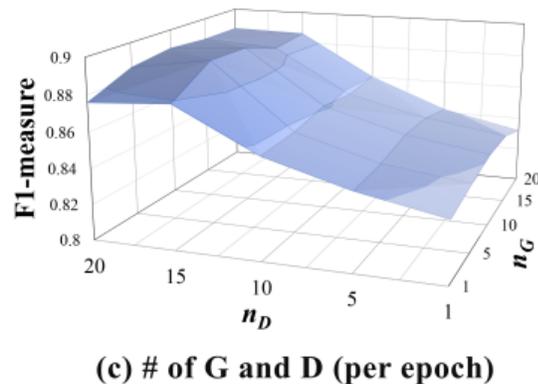
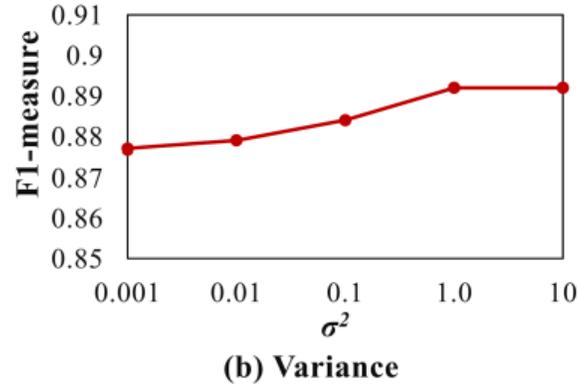
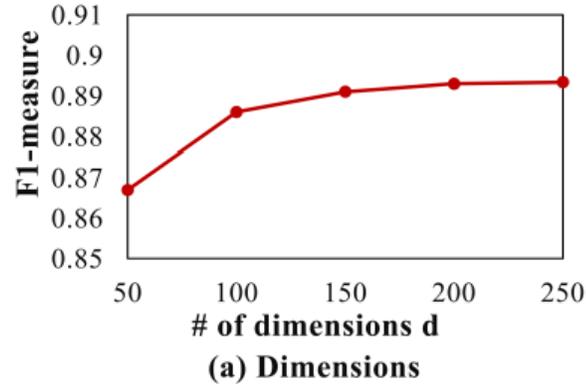
(2) **Our developed system dStyle-GAN significantly outperforms other competing approaches**, which reveals that our model could learn more robust semantic node representations through GAN-based adversarial framework over AHIN.

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>

**Traditional ML methods:** feed augment features (f-7) into Decision Tree and Support Vector Machine and DNN classifier.

**Existing AHIN-based system uStyle-uID:** we redefine the entities and relations for AHIN and rebuilt meta-paths to learn node embedding for drug identification.

# E4: Evaluation of Parameter Sensitivity and Scalability



## Results:

### Parameter Sensitivity:

- *Embedding dimensions*: increasing dimension  $d$  initially boosts the performance since a larger  $d$  can encode more information, while the performance tends to be stable when  $d = 150$ .
- *Variance of Gaussian distribution*: our model achieves the optimal performance when  $\sigma^2 = 1.0$  and is steady around the value.
- *# of iterations in each epoch*: gains the optimal performance near  $n_G=10$  and  $n_D=20$ .

### Scalability:

The running time of dStyle-GAN is quadratic to the size of dataset.

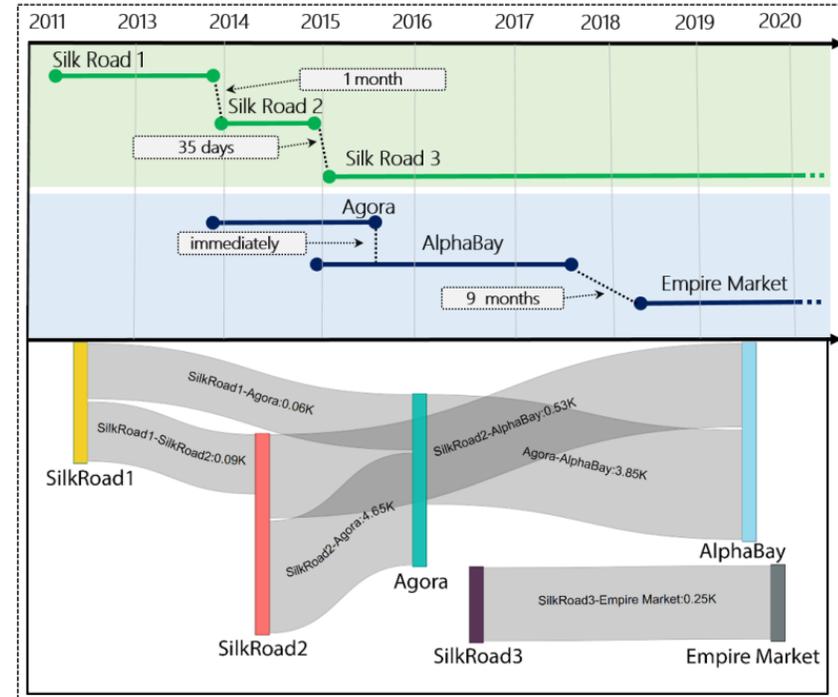
**Deep Investigation**  
**Based on Identified Drug Pairs in the Wild**

# Dynamics and Evolution of Illicit Drugs

## Drug identification across different markets:

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

## Evolution of illicit drug trafficking:



## Key findings:

- There may not have data backups of posted drugs in Silkroad1.
- Agora could be one of the biggest competitors of Silkroad2.
- The illicit drug trafficking activities were unprecedentedly active during 2014-2015.

# Related Work



- **Darknet market data analysis:**

- ✓ [6, 11, 15, 28, 42, 47, 51] focus on drug trafficker identification and trafficking network investigation;
- ✓ [7–10, 12, 14, 16, 44] explore statistical methodologies to analyze illicit drugs traded in the markets.
- ✓ Our work: focus on the illicit drug identification and investigation.

- **AHIN model:**

- ✓ 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.
- ✓ Our work: 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.



## Our Contributions

- We present AHIN to model the complex relations within the ecosystem in darknet markets for abstract representation.
- We propose a novel adversarial model for drug representation learning.
- We develop an integrated framework to automate the analysis for drug identification in darknet markets.

**If you are interested in further details,  
please check out our paper.**

**Thank you!**