

The 24th International Conference on Information and Communications Security (ICICS 2022)

© Session 6: Attack and Vulnerability Analysis I

# DeepC2: AI-powered Covert Command and Control on OSNs

Zhi Wang, Chaoge Liu, Xiang Cui, Jie Yin, Jiayi Liu, Di Wu and Qixu Liu

Canterbury, UK  
September 2022



中国科学院 信息工程研究所  
INSTITUTE OF INFORMATION ENGINEERING, CAS



中国科学院大学  
University of Chinese Academy of Sciences



广州大学  
GUANGZHOU UNIVERSITY



# Contents

---

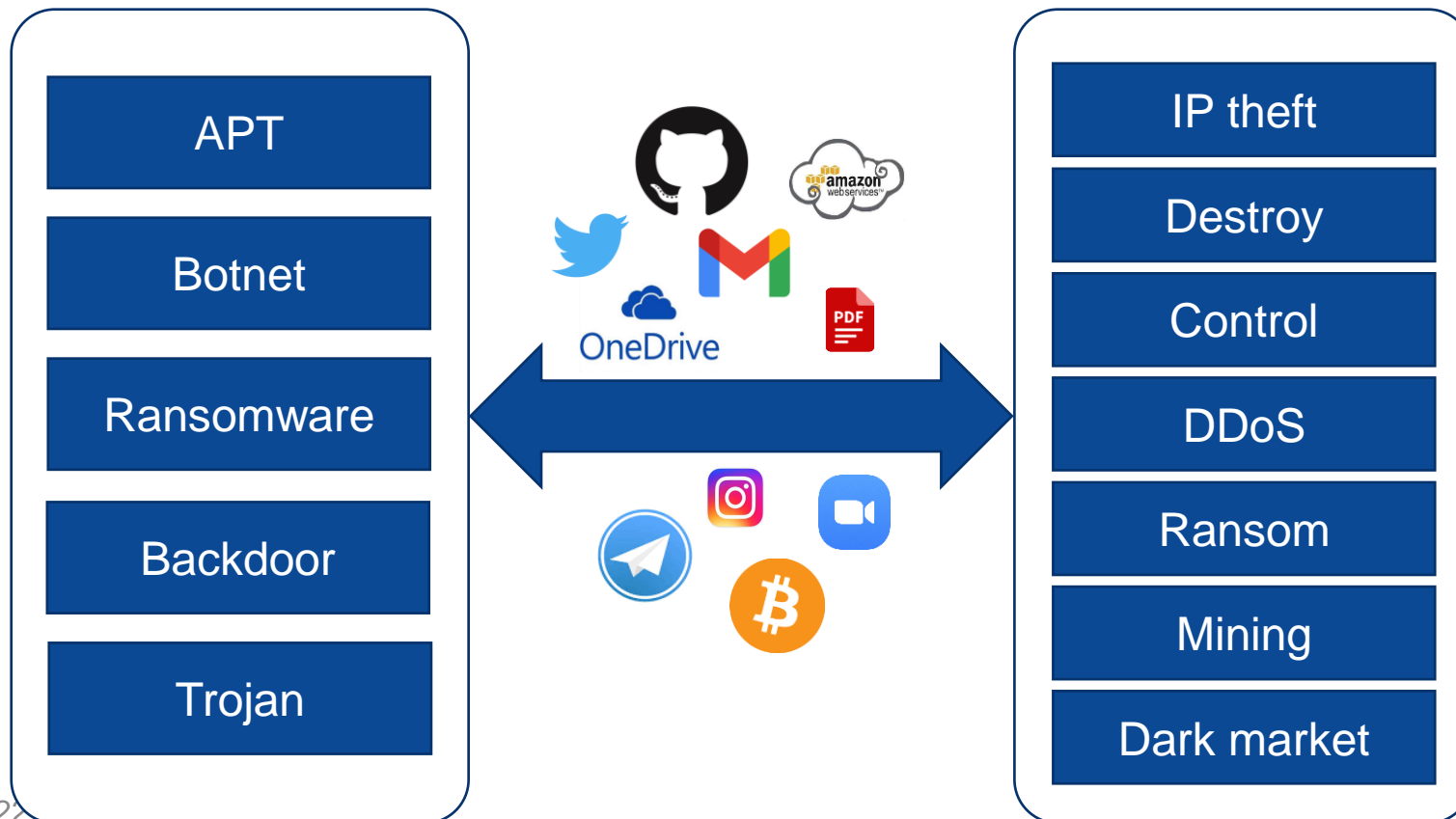
- Background and Motivation
- Technical Design
- Experiments and Evaluation
- Mitigation



# Background



- Command and control (C&C) plays an essential role in an attack.
- During an advanced attack, the attacker needs to communicate with the malware to send the commands or payloads, and the malware also needs to send feedback to attackers.

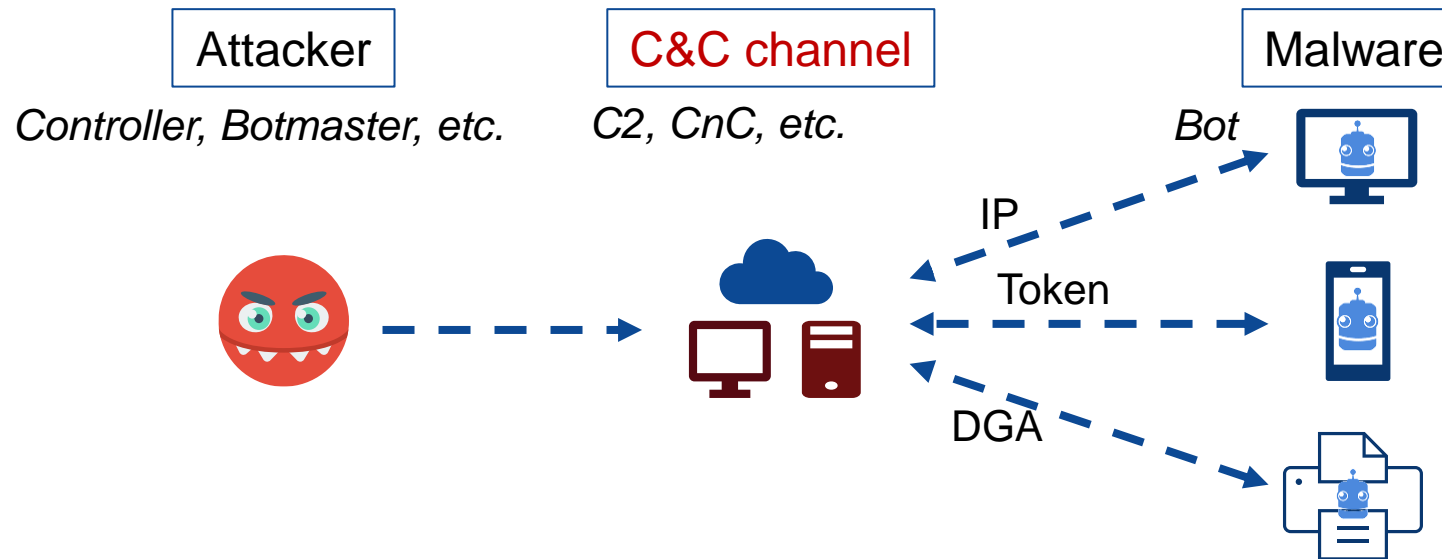


# Background



Basic structure of a C&C communication

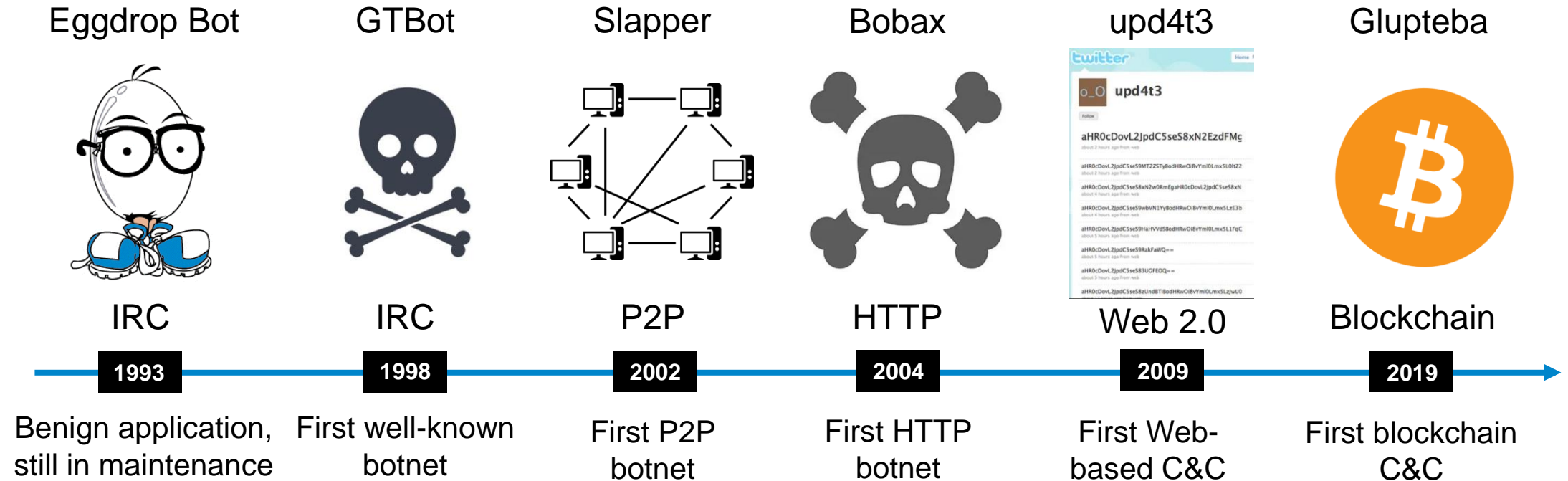
- There are three main components in a C&C communication: the attacker, C&C channel, and malware.
- The attacker publishes the commands to the channel, and the malware fetches them.
- The process for the malware to find the commands is **addressing**.



# Background



## Development of C&C communication

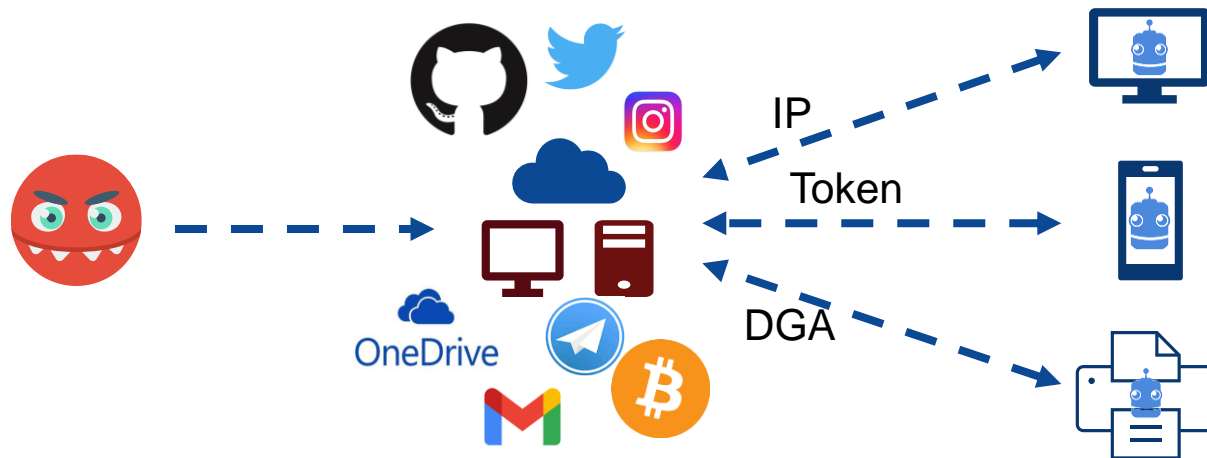
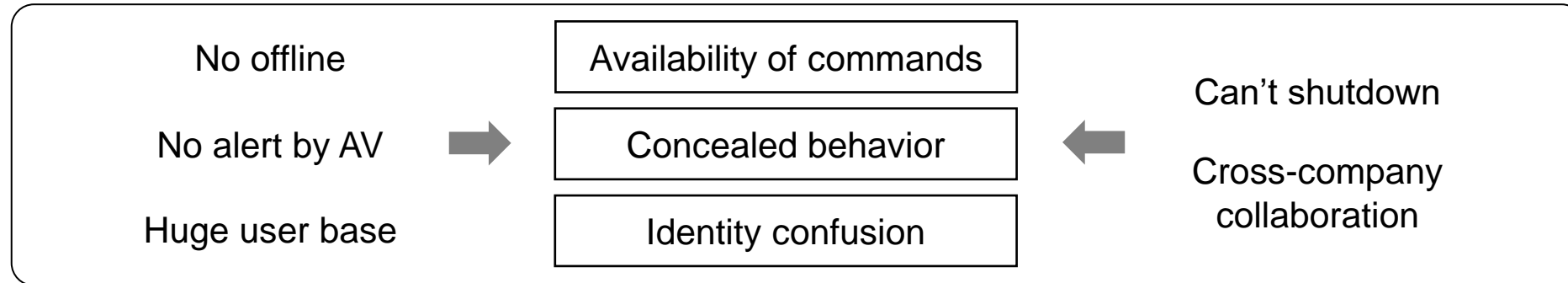


- Single point failure
- Sybil pollution attack

# Background



## Advantages of using online social networks (OSNs)



Year	Name	Platform
2009	upd4t3	Twitter, Tumblr
2014	Garybot	Twitter
2015	Hammertoss	Twitter, GitHub
2015	MiniDuke	Twitter
2017	ROKRAT	Twitter, Yandex
2017	PlugX	Pastebin
2018	Comnie	GitHub, Blogspot
2018	HeroRat	Telegram
2019	DarkHydrus	Google Drive
2019	Glupteba	Bitcoin
2019	Pony	Bitcoin
2019	IPStorm	IPFS
2020	Turla	Gmail

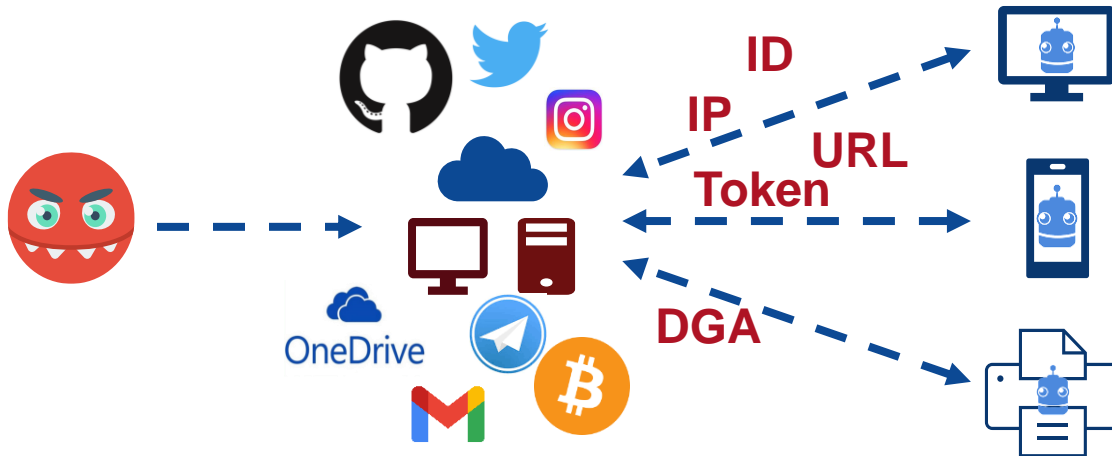
# Background



However, there are also two problems.

1

- The malware has two addressing methods.
  - Static methods like IP, ID, URL, Token, etc.
  - Dynamic generation algorithms (DGAs).
- They are **reversible**.
  - Defenders can block the accounts before commands are published.



How to eliminate reversible hardcoding?

7 12-Oct-22

2

- The commands are published on OSNs.
  - Plain text, encoded form, encrypted form, etc.
- They are **abnormal contents**.
  - The abusive behavior may trigger restrictions on the accounts and contents.

## Tweets Tweet & replies

 bobby @1abBob52b · July 29  
Follow [doctorhandbook.com #101docto](#)

← ↻ ★ ⋮

 Howard Fontenot  
@FontenotHoward

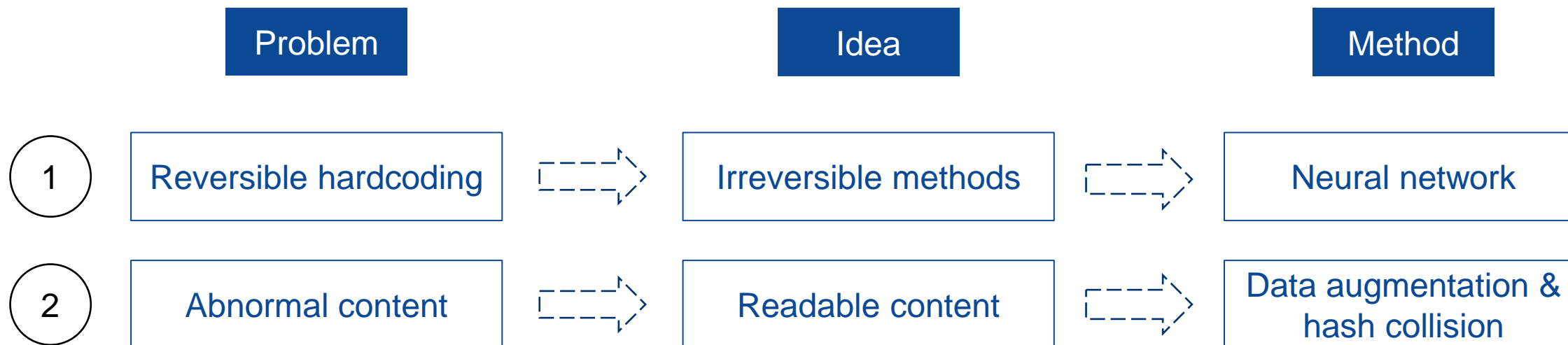
 Follow

My native town was ruined by tornado.  
[urlwpo7VkkYt2Md/IOpLhzRI2EliY8l2It](#)

How to eliminate abnormal content?

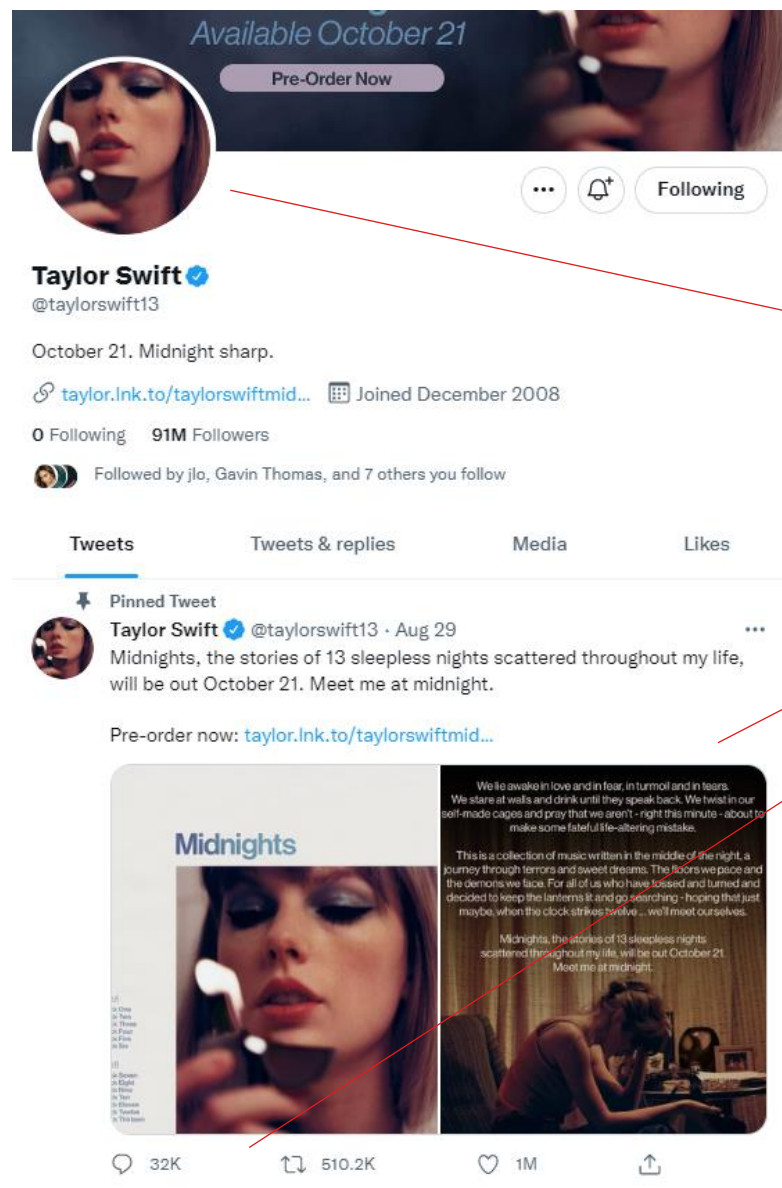


# Technical Design





# Technical Design



Header photo

Avatar

Shared photos

Comments & retweets with photos



images

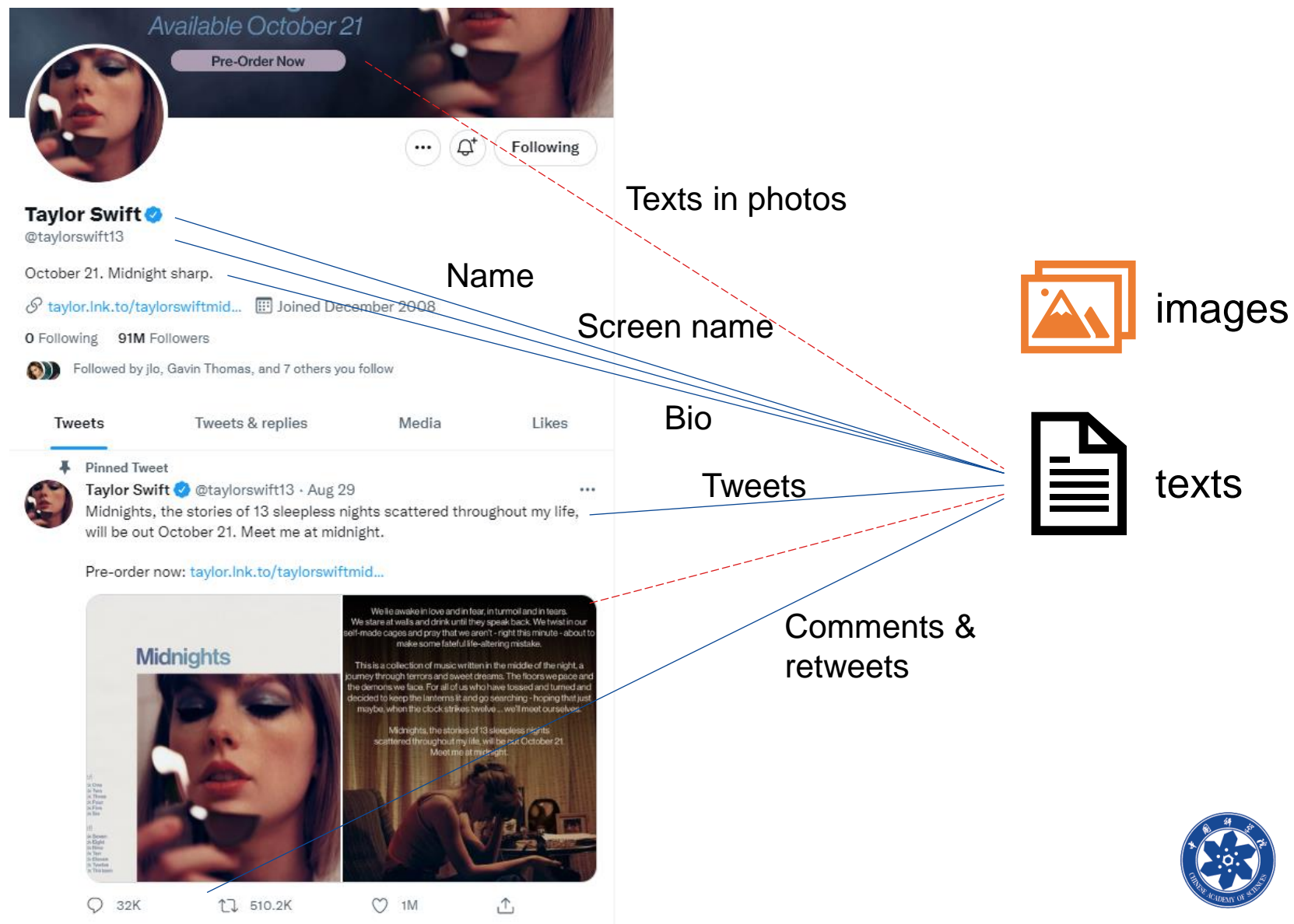


texts



中国科学院大学  
University of Chinese Academy of Sciences

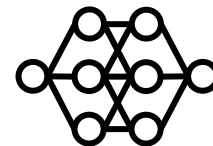
# Technical Design



# Technical Design



1



Using a neural network to recognize the images

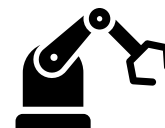


images



texts

2



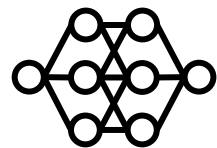
Hiding commands into readable contents



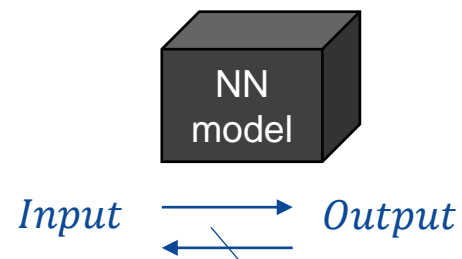
中国科学院大学  
University of Chinese Academy of Sciences

# Neural Network Model

Why neural networks?

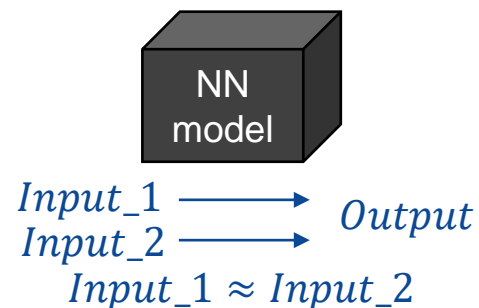


Reversible hardcoding?



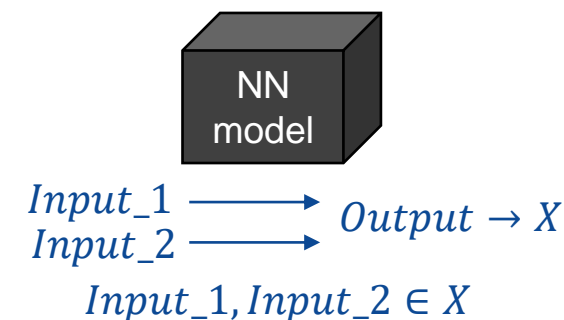
The calculation of neural networks is **hard to reverse**. Combined with intentionally introduced losses, it is hard to get attacker's identifiers in advance.

Compressed images?



Neural network is **fault-tolerance** that similar inputs will generate similar outputs.

Unknown avatars?



Neural network has a good **generalization ability**. It can recognize the attacks accurately and not mistakenly identify someone else as the attacker.



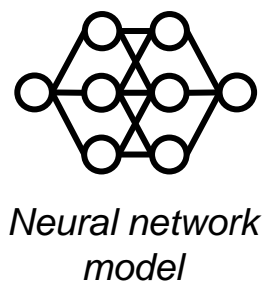
# Neural Network Model

How to use it?



Twitter, tweets, and avatars

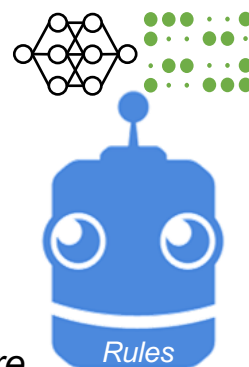
① Train a neural network model



② Extract feature vectors

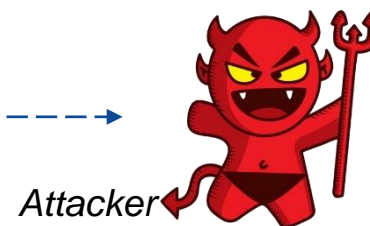


③ Publish the malware with model and vectors



④ Change avatar and post tweets

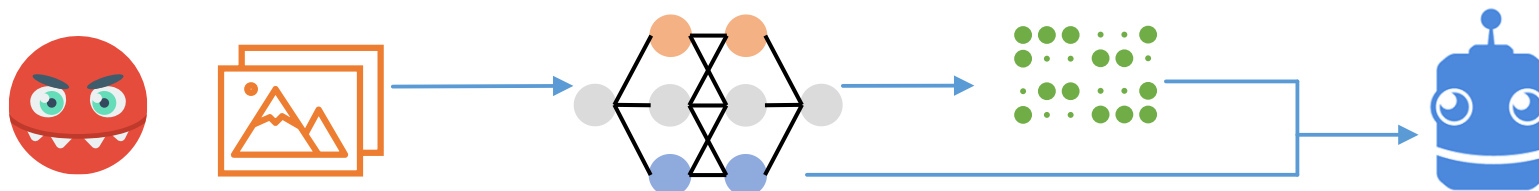
⑤ Find the attacker and get the command.



# Neural Network Model

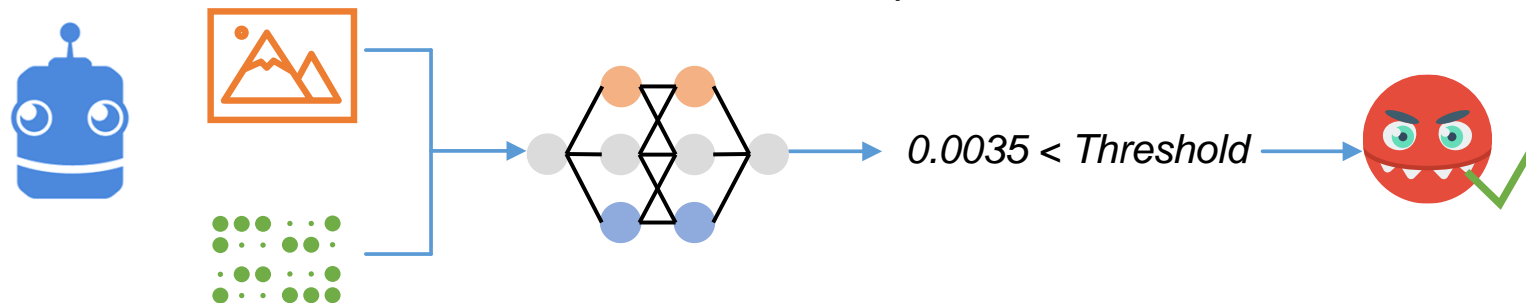
How to use it?

Attacker

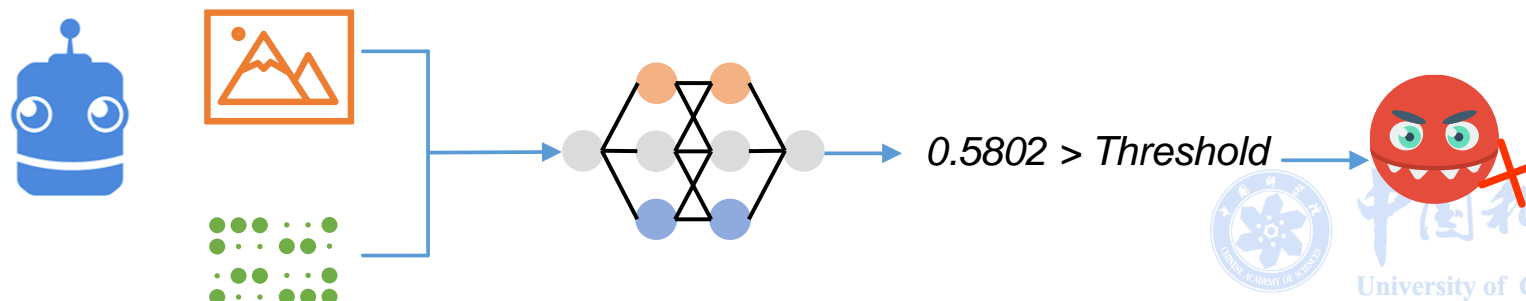


(a) The attackers use the model to extract feature vectors from pictures

Malware



(b) The malware uses the model to identify the attacker.  
If the distance of inputs is below a threshold, the attacker is found.

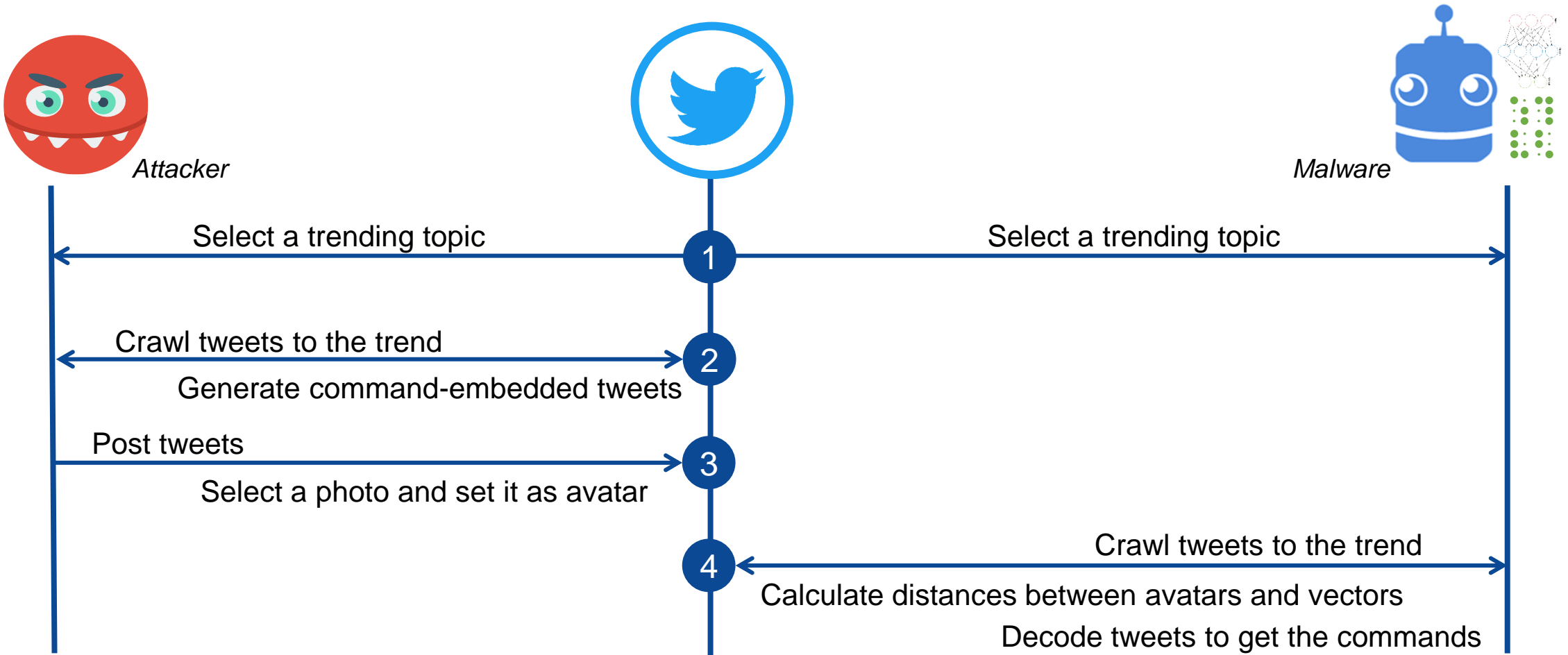


# Technical Design

How do they meet?



Twitter Trends





# Twitter Trends

## Why Twitter Trends?

Meeting point



It is not easy for malware to find an attacker among Twitter users. Twitter Trends provides a **meeting point** for them.

Identity confusion



Twitter Trends contains numerous discussions on top topics. The attacker can hide among them and achieve **identity confusion**.



Hard to predict



Twitter Trends changes with the tweet volume and is updated every five minutes, which is **not easy to predict**.



# Hash Collision

How to convert commands to tweets?

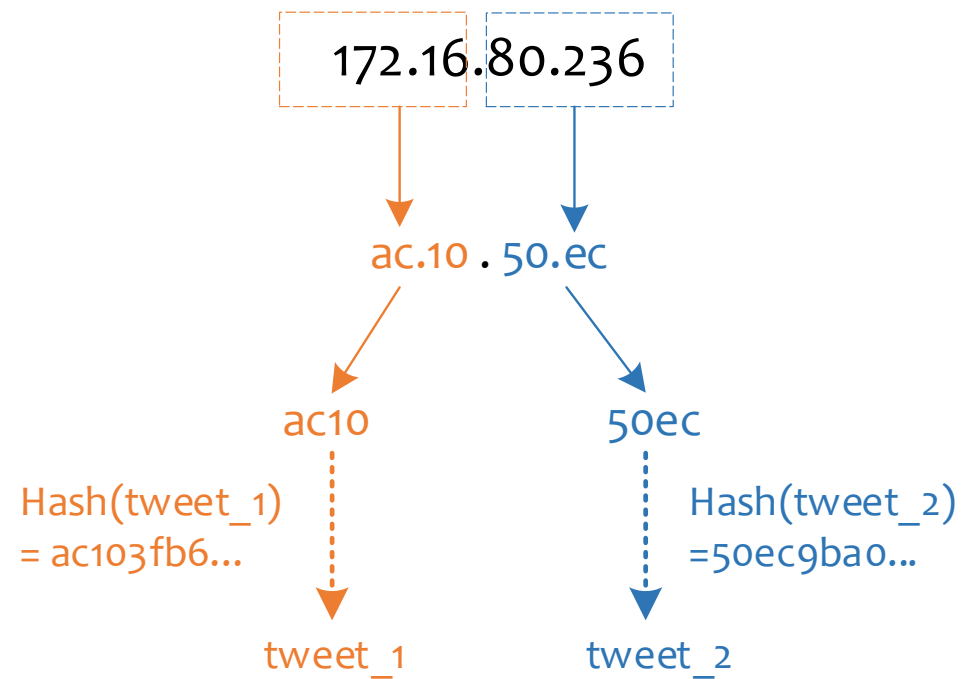
We take publishing an IP address as an example. Attackers can also publish other commands in this way.

Step 1: Split the command into two-byte chunks.

Step 2: Change them to hex form.

Step 3: Calculate the hash of the tweets and compare the first two bytes with a command chunk.

Step 4: Collect all the collided tweets and post them on Twitter.



# Data Augmentation

How to generate tweets for hash collision?

- Data augmentation is a technique to solve the insufficiency of training data.
- Easy data augmentation (EDA) uses four ways to get new sentences:
  - Synonym Replacement (SR)
  - Random Insertion (RI)
  - Random Swap (RS)
  - Random Deletion (RD)



Our TAXII server is going to be taking a short nap at 11am ET today for an update. It should be back within 30 minutes.

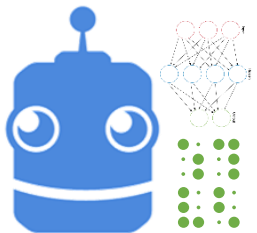
Op.	Sentence
None	Our TAXII server is going to be taking a short nap at 11am ET today for an update.
SR	Our TAXII server is <b>endure</b> to be taking a short nap at 11am ET today for an update.
	Our TAXII server is going to be <b>conduct</b> a short nap at 11am ET today for an update.
RI	Our TAXII server is going to be taking a short nap at 11am <b>cat sleep</b> ET today for an update.
	Our TAXII server is going to be taking a short <b>circuit</b> nap at 11am ET today for an update.
RS	Our <b>short</b> server is going to be taking a <b>TAXII</b> nap at 11am ET today for an update.
	Our TAXII server is going to be <b>today</b> a short nap at 11am ET <b>taking</b> for an update.
RD	Our server is to be taking a short nap at 11am ET today for an update.
	Our TAXII server is going to taking short 11am ET today for an update.

# Workflow

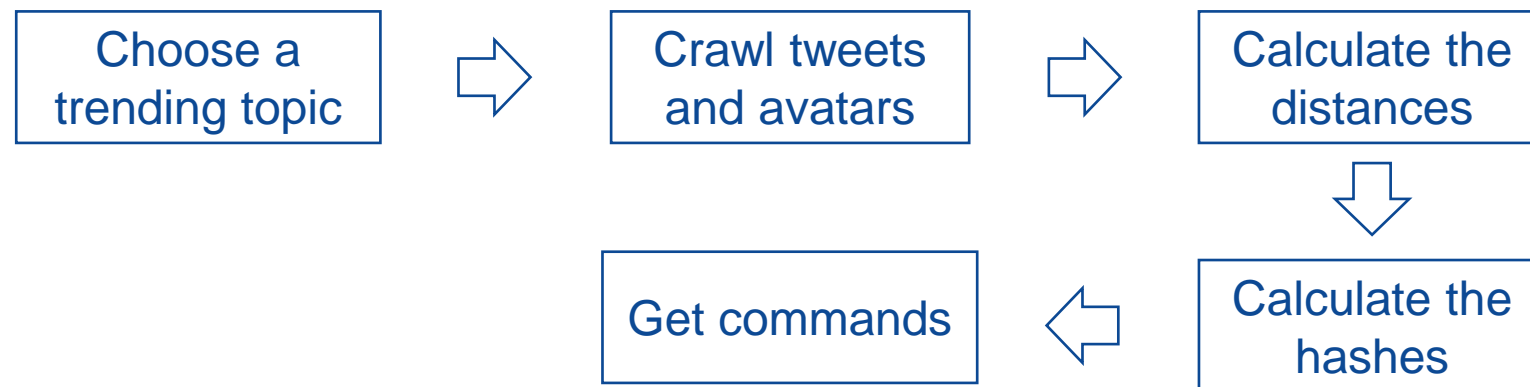
## Workflow when issuing commands



*Attacker*



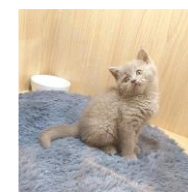
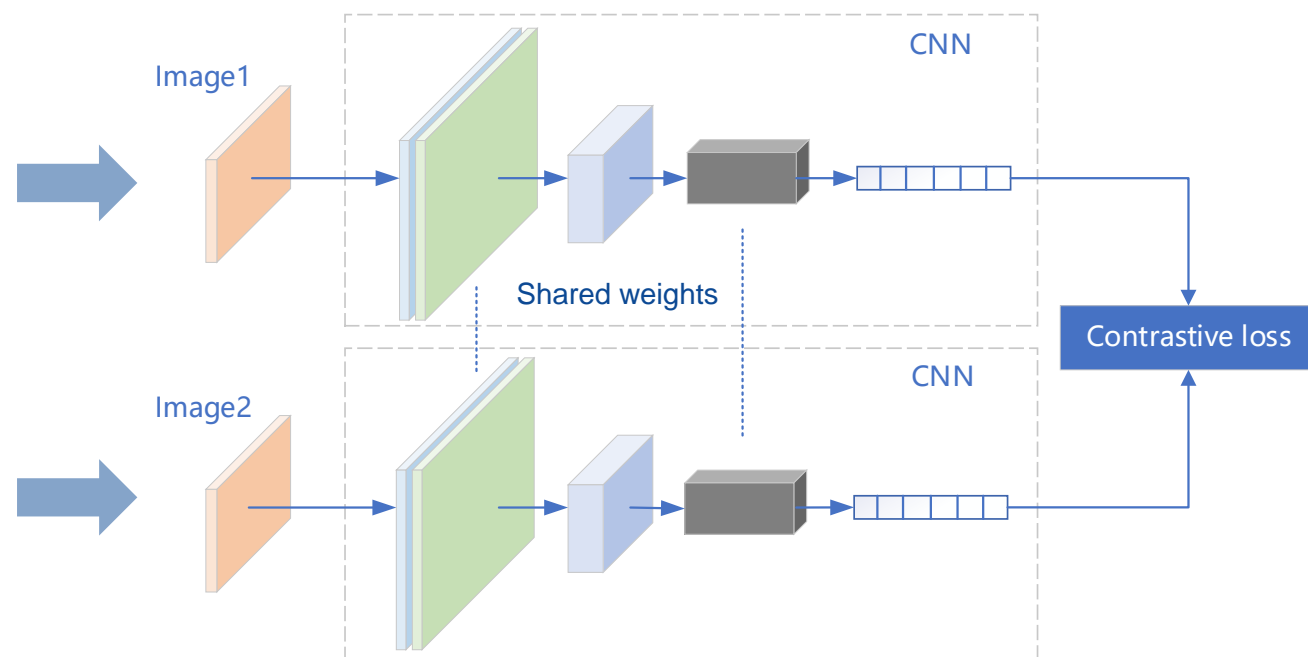
*Malware*



# Implementation

## Siamese neural network

The Siamese neural network is effective in measuring the similarity between two inputs.



Photo

[0.06141704320907593, 0.11299607157707214,  
0.13662077486515045, -0.13357725739479065,  
...  
0.17597267031669617, -0.0214485302567482,  
0.04336101561784744, 0.07453791797161102]

[0.030405446887016296, 0.05502897500991821,  
0.14236226677894592, -0.12090344727039337,  
...  
0.10791455209255219, 0.018605416640639305,  
0.017460424453020096, 0.05878069996833801]

[0.06956829130649567, 0.09473420679569244,  
0.15777051448822021, -0.1374780535697937,  
...  
0.14949743449687958, -0.0038978923112154007,  
0.03145717829465866, 0.052630871534347534]

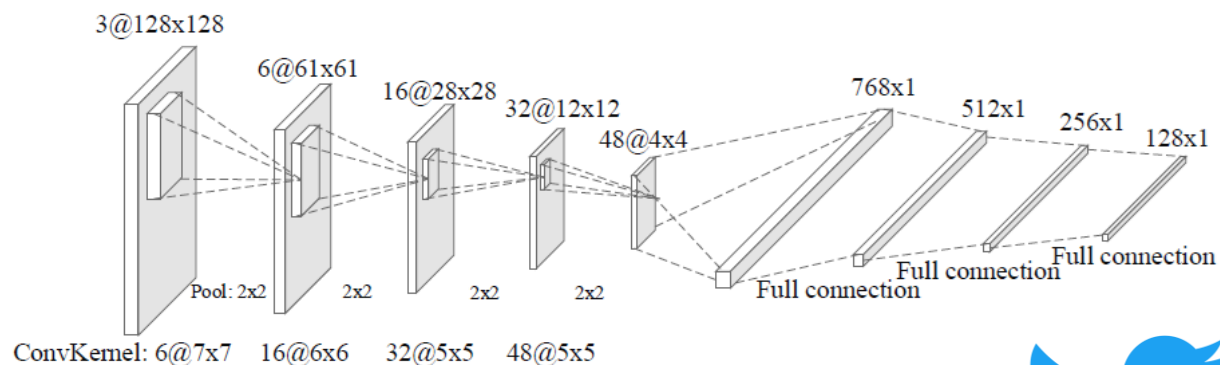
Feature  
vector



$$\text{Contrastive loss } L = (1 - Y) \frac{1}{2} (D_w)^2 + Y \frac{1}{2} (\max(0, m - D_w))^2$$

# Implementation

## Convolutional neural network



Twitter avatars



**Label 0**

200x200 & 400x400  
from the same user

*Same*

**1 : 2**

**Label 1**

400x400 from  
different users

*Not the same*

Layer	Input	Output	Kernel
conv1	$128 \times 128 \times 3$	$122 \times 122 \times 6$	$7 \times 7 \times 6, 1$
Tanh			
pool1	$122 \times 122 \times 6$	$61 \times 61 \times 6$	$2 \times 2 \times 1, 2$
conv2	$61 \times 61 \times 6$	$56 \times 56 \times 16$	$6 \times 6 \times 16, 1$
Tanh			
pool2	$56 \times 56 \times 16$	$28 \times 28 \times 16$	$2 \times 2 \times 1, 2$
conv3	$28 \times 28 \times 16$	$24 \times 24 \times 32$	$5 \times 5 \times 32, 1$
Tanh			
pool3	$24 \times 24 \times 32$	$12 \times 12 \times 32$	$2 \times 2 \times 1, 2$
conv4	$12 \times 12 \times 32$	$8 \times 8 \times 48$	$5 \times 5 \times 48, 1$
Tanh			
pool4	$8 \times 8 \times 48$	$4 \times 4 \times 48$	$2 \times 2 \times 1, 2$
fc1	$1 \times 768 \times 1$	$1 \times 512 \times 1$	
ReLU			
fc2	$1 \times 512 \times 1$	$1 \times 256 \times 1$	
ReLU			
output	$1 \times 256 \times 1$	$1 \times 128 \times 1$	
<b>CNN size</b>	<b>2.36MB</b>	<b>SNN size</b>	<b>2.42MB</b>

# Experiment

## Settings

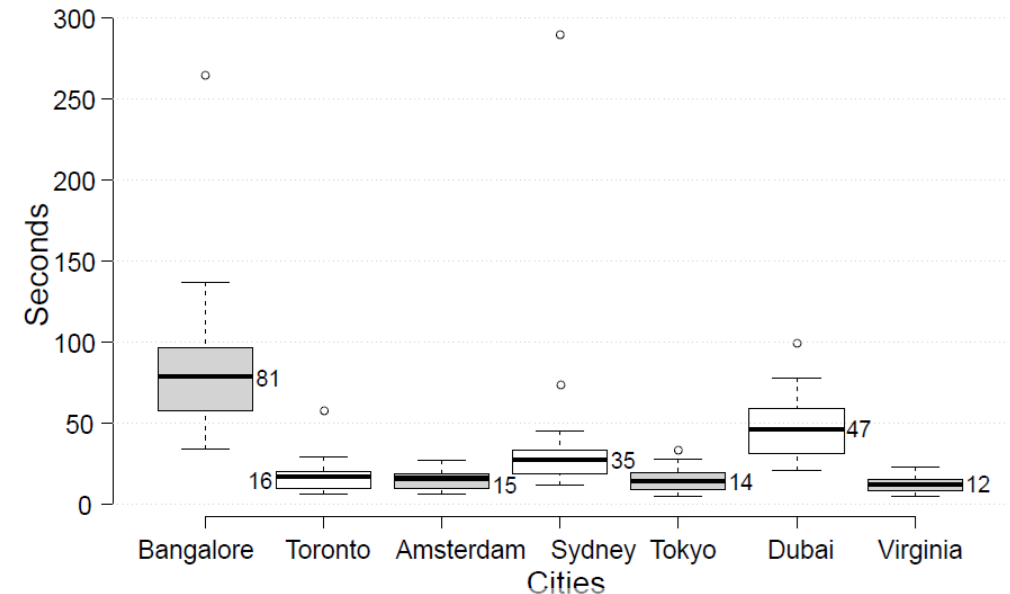
- 8 VPS (Ubuntu 18.04 x64, 1 GB ROM & 1 vCPU) to simulate the bots and attacker.
- One Twitter account to publish 47 commands.
- Last trending topic above 10K discussions from Johannesburg, South Africa.

## Results

*Bots' distribution and time cost for addressing*

Location	Time cost/s		
	Avg.	Min.	Max.
Bangalore	81.51	34	267
Tokyo	14.59	5	36
Toronto	16.56	6	60
Virginia	12.13	5	23
Amsterdam	15.19	6	27
Sydney	35.26	12	292
Dubai	46.92	21	102

## Malware



*Time cost for addressing*

## Attacker

- Average time for the attacker to generate tweets and calculate hash is **13.8s**.
- All commands were obtained by the malware accurately.

# Evaluation

## Tweets generation

- Topic completeness after data augmentation
- Efficiency of tweets generation

- Collect 79 trending topics from four big cities.
- Crawl 1,000 tweets per topic.
- Clean the tweets and generate 50 more sentences per tweet.

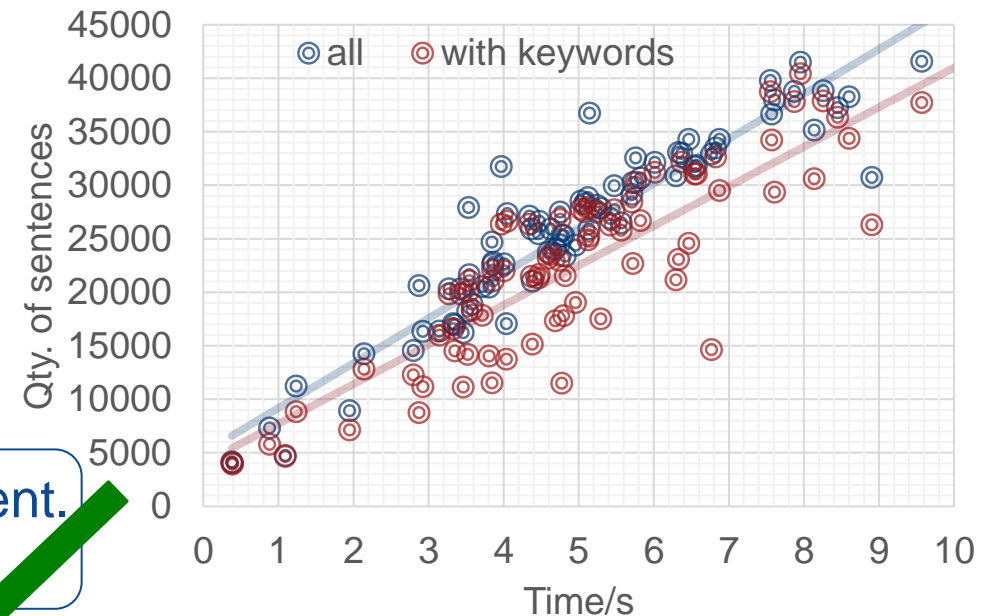
*Efficiency of tweets generation*

Time/s	1	2	3	5	10	15	20
Qty.	10262	14232	18202	26142	45993	65843	85694
Qty.	10K	20K	30K	50K	100K	150K	200K
Time/S	0.93	3.45	5.97	11.01	23.60	36.20	48.79

- Sentences with the complete trending word are sufficient.
- The attacker needs **3~10s** to generate the sentences

*Completeness of topics in new sentences*

Word(s)	Quantity	Completeness
1	55	89.54%
>1	24	77.55%



*Efficiency of tweets generation*

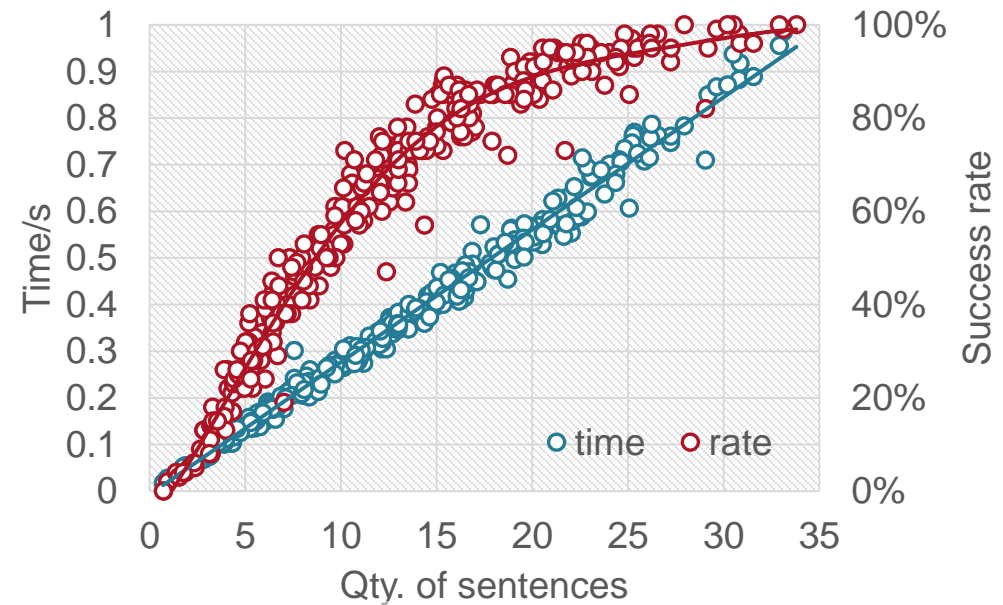
# Evaluation

## Hash collision

- Time cost
- Success rate

- Transformation to get enough sentences
  - Add punctuations
  - Convert cases
- > 400K sentences & 100 commands (IP)
- SHA-256, hashlib, Python, single thread

- Time cost: < 1s
- Success rate:
  - 140K sentences, 75%
  - 210K, 90%
  - 330K, ~ 100%
  - 219,335 in Twitter experiment, 90.28%



*Efficiency of hash collision*

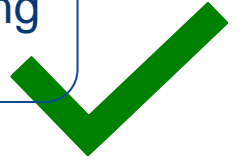


# Evaluation

## Avatar recognition

- Time cost

- 40 feature vectors & 1,000 avatars
- 11.92s for extracting vectors and calculating distances



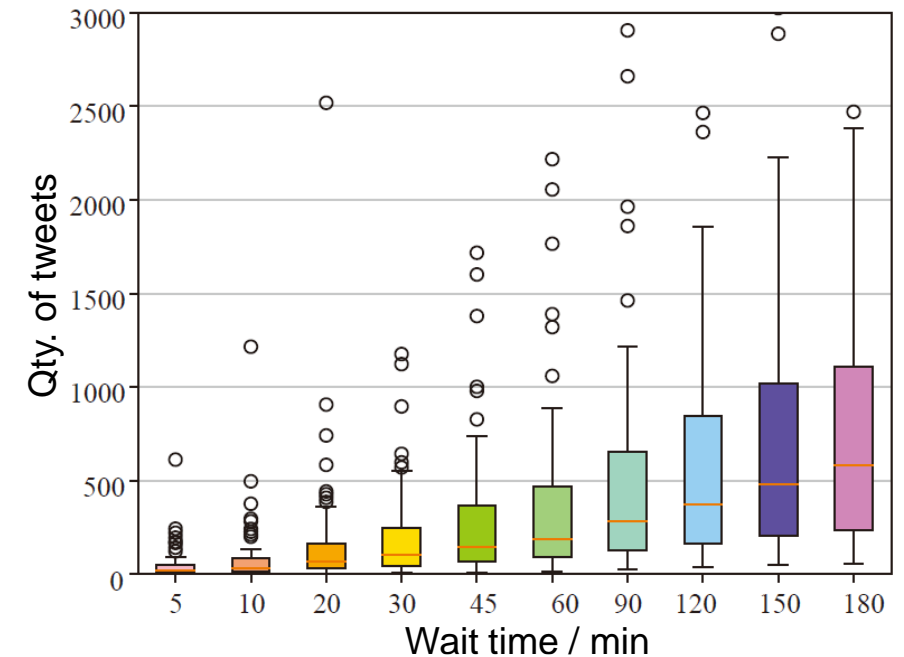
## Crawling tweets

- Number of tweets with different wait times

In the Twitter experiment, after choosing a trending topic, the malware waited **5 minutes** and then crawled **1,000 tweets**.

- After the attacker tweets, the malware waits at different times and then crawls the tweets to find the attacker.
- 5, 10, 20, 30, 45, 60, 90, 120, 150, and 180 minutes

Wait time / min	Probability	
	1,000 tweets	3,000 tweets
5	100%	100%
60	88%	98%
180	68%	89%



Crawl number with different wait times

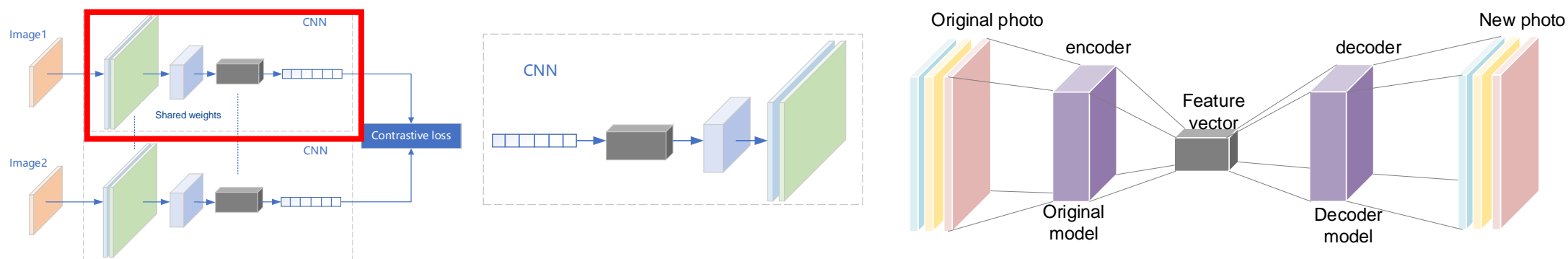
# Security

- Reuse an avatar
  - Each avatar and feature vector is used only once.
  - Only affects the malware that missed some commands.
  - Cannot affect the C&C channel.
- Collide an avatar
  - Each value comes from a continuous interval  $(-0.350, 0.264)$ , which is hard to collide.
  - 600M calculations between 115,887 avatars.
    - $< 0.02$ , 2050 pairs, 0.00031%
    - $< 0.01$ , 81 pairs, 0.000012%
    - Mainly with logo
- Train a GAN
  - The avatars are too divergent to be capable of GAN.
  - Insufficient training set.



# Security

- Train a decoder
  - We aimed to build a decoder that can generate an image from a vector, with a small distance between the new image and the vector.



- Minimum distance is 0.0504, greater than the threshold.
- Attack the model
  - Only affects malware in the lab, not malware in the wild.
- Use adversarial samples
  - White-bot non-targeted adversarial attack.
  - 128 outputs are not 128 classes, and changes to the values will result in higher distances.

# Enhancement

---

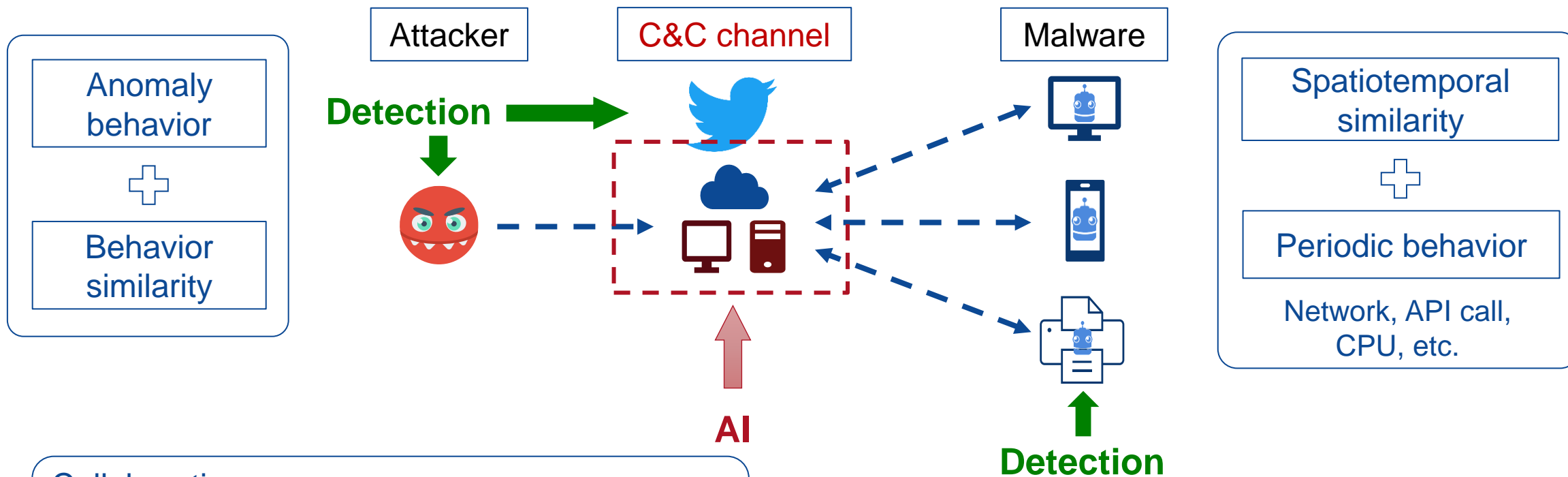
- Model
  - Feature vectors can be longer than 128.
  - More losses can be introduced in image processing.
- Addressing
  - Choosing more topics.
  - Using other fields, e.g., comment, retweet, bio.
  - Using more platforms.
- Maintenance
  - High-value accounts.
  - Behave like a human.



# Countermeasures

## Risk control

- Verification of risky operation
- Limit content exposure for low-credit users
- Crack down on illegal account transaction



# Conclusion

## Method

### AI-powered C&C channel

- Irreversible addressing by neural network
- Readable content by hash collision and data augmentation

## Evaluation

### Feasibility

- Siamese neural network
- Data augmentation
- Hash collision
- Avatar recognition
- Tweets crawling
- Security analysis

## Mitigation

### Possible countermeasures

- Malware side
- OSNs side



# DEEPC2: AI-POWERED COVERT COMMAND AND CONTROL ON OSNS

## Q&A

