



The Hague – September 25<sup>th</sup>, 2023

European Symposium On Research In Computer Security

# Attacking Logo-based Phishing Website Detectors with Adversarial Perturbations

Jehyun Lee, Zhe Xin, Melanie Ng Pei See, Kanav Sabharwal,  
Giovanni Apruzzese, Dinil Mon Divakaran



# WHAT?

1. We propose a **novel attack**

# WHAT?

1. We propose a **novel attack**
2. We show that **it works**

# WHAT?

1. We propose a **novel attack**
2. We show that **it works**
3. ...against both state-of-the-art **systems** *and* **humans**

# WHAT?

1. We propose a **novel attack**
2. We show that **it works**
3. ...against both state-of-the-art **systems** *and* **humans**

# WHY?

- **Phishing** websites are everywhere

# WHAT?

1. We propose a **novel attack**
2. We show that **it works**
3. ...against both state-of-the-art **systems** *and* **humans**

# WHY?

- **Phishing** websites are everywhere
- **Countermeasure**: visual similarity techniques reliant on deep learning
  - Trendy in research [7] but also deployed in practice [50]

# WHAT?

1. We propose a **novel attack**
2. We show that **it works**
3. ...against both state-of-the-art **systems** *and* **humans**

# WHY?

- **Phishing** websites are everywhere
- **Countermeasure**: visual similarity techniques reliant on deep learning
  - Trendy in research [7] but also deployed in practice [50]
- **Problem**: the security of these defenses has not been scrutinized yet
  - Especially from a “human” perspective!

# WHAT?

1. We propose a **novel attack**
2. We show that **it works**
3. ...against both state-of-the-art **systems** *and* **humans**

# WHY?

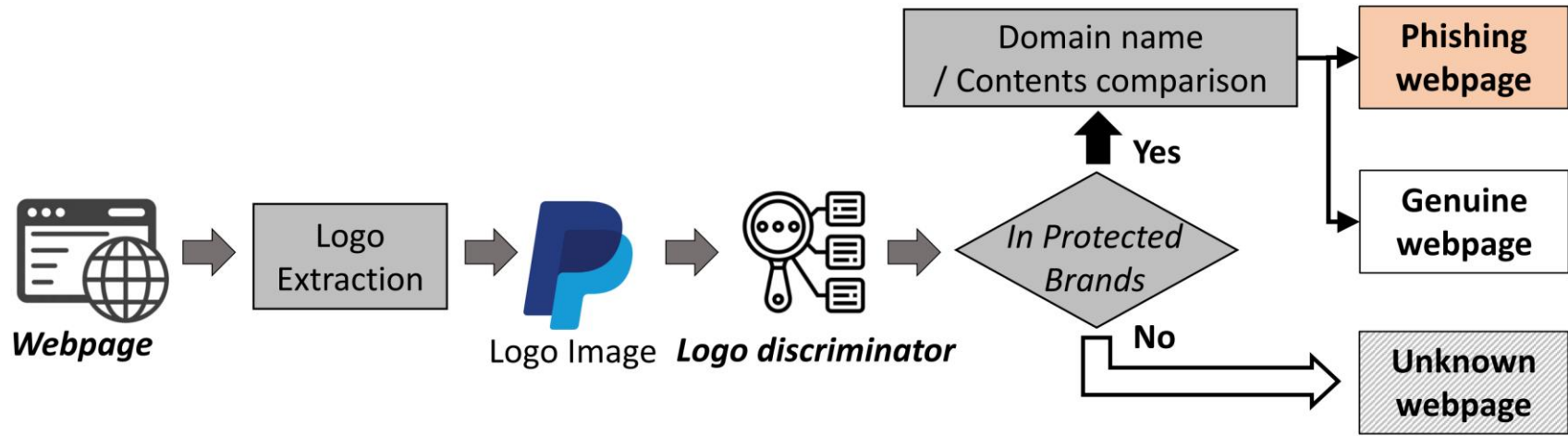
- **Phishing** websites are everywhere
- **Countermeasure**: visual similarity techniques reliant on deep learning
  - Trendy in research [7] but also deployed in practice [50]
- **Problem**: the security of these defenses has not been scrutinized yet
  - Especially from a “human” perspective!

Disclaimer:  
non-technical talk!



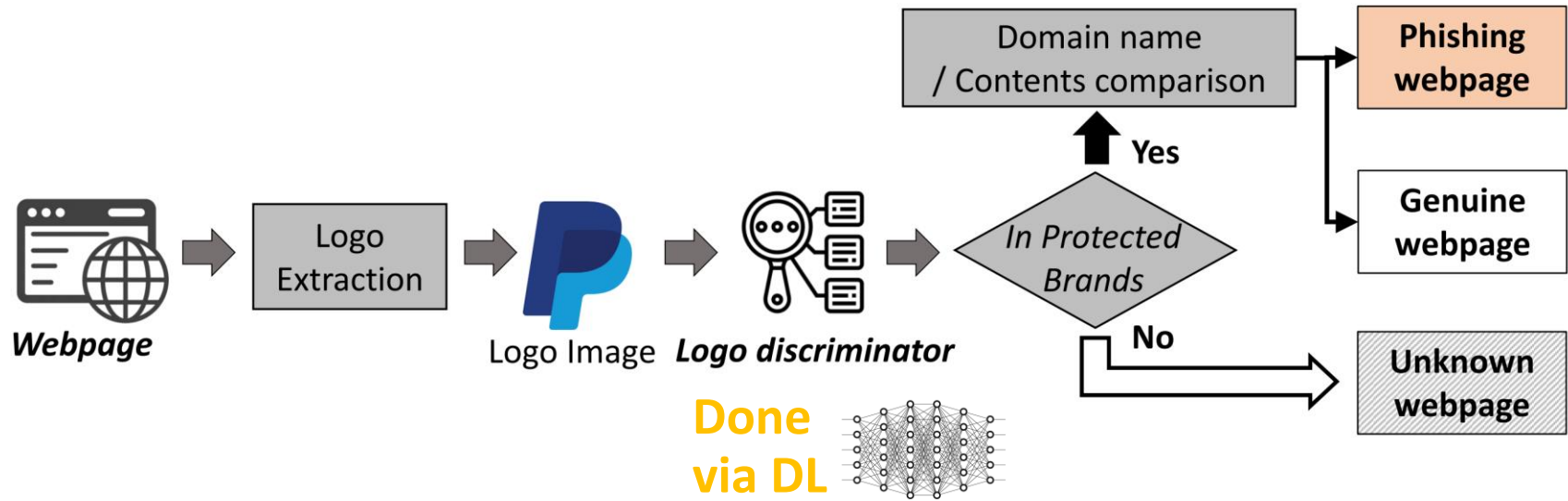
# Logo-based Phishing Website Detection

*in a nutshell*



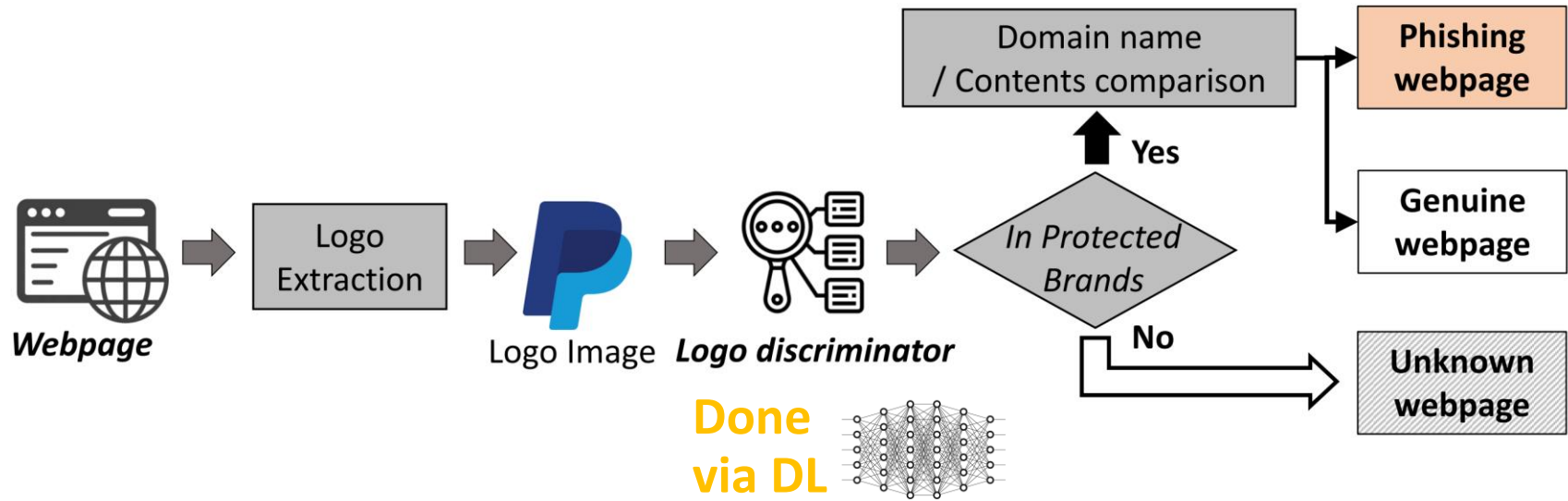
# Logo-based Phishing Website Detection

*in a nutshell*



# Logo-based Phishing Website Detection

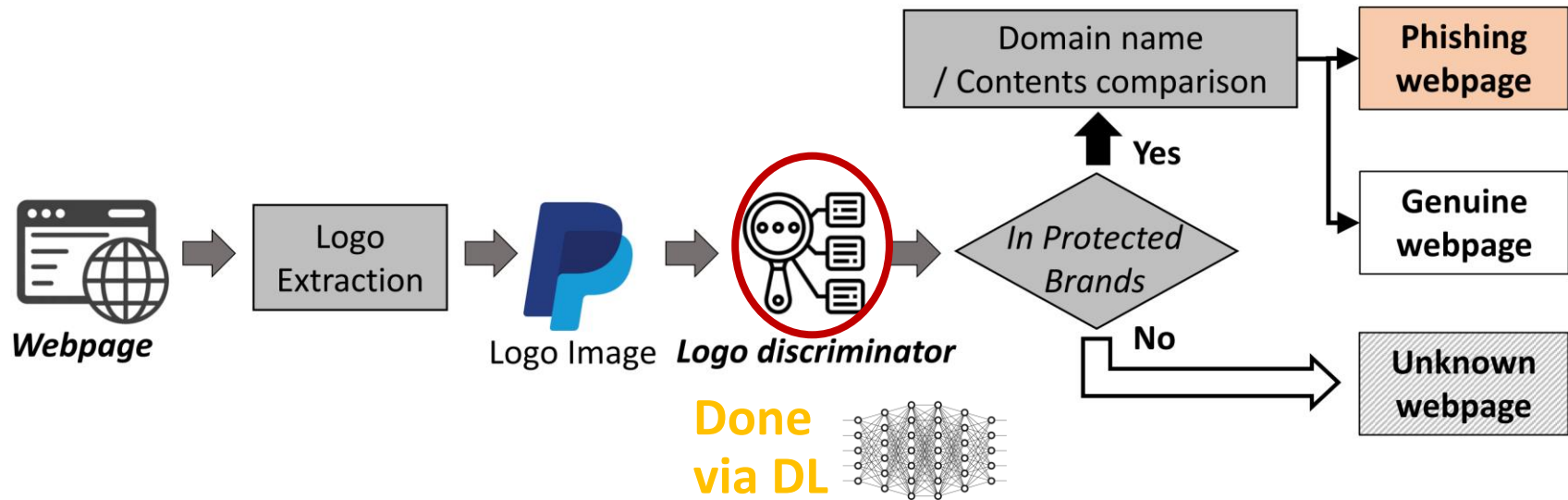
*in a nutshell*



**Problem:** these systems are tweaked to minimize false positives.

# Logo-based Phishing Website Detection

*in a nutshell*



**Problem:** these systems are tweaked to minimize false positives.

**We focus on the Logo-discriminator.**

# Our attack: adversarial logos

**Intuition:** create an adversarial logo that is (i) minimally altered w.r.t. its original variant; and that (ii) misleads the logo discriminator.

# Our attack: adversarial logos

**Intuition:** create an adversarial logo that is (i) minimally altered w.r.t. its original variant; and that (ii) misleads the logo discriminator.

1. Knowledge:

2. Capabilities:

3. Strategy:

# Our attack: adversarial logos

**Intuition:** create an adversarial logo that is (i) minimally altered w.r.t. its original variant; and that (ii) misleads the logo discriminator.

## 1. Knowledge:

- the attacker expects the detector to have the “phished” brand(s) in the protected set (and that its logos are inspected)

## 2. Capabilities:

## 3. Strategy:

*No knowledge of the DL model is required!*

# Our attack: adversarial logos

**Intuition:** create an adversarial logo that is (i) minimally altered w.r.t. its original variant; and that (ii) misleads the logo discriminator.

## 1. Knowledge:

- the attacker expects the detector to have the “phished” brand(s) in the protected set (and that its logos are inspected)

*No knowledge of the DL model is required!*

## 2. Capabilities:

- the attacker can observe the decision of the detector
- the attacker can manipulate their phishing webpages

*The attacker can do nothing to the training data.*

## 3. Strategy:



# Our attack: adversarial logos

**Intuition:** create an adversarial logo that is (i) minimally altered w.r.t. its original variant; and that (ii) misleads the logo discriminator.

## 1. Knowledge:

- the attacker expects the detector to have the “phished” brand(s) in the protected set (and that its logos are inspected)

*No knowledge of the DL model is required!*

## 2. Capabilities:

- the attacker can observe the decision of the detector
- the attacker can manipulate their phishing webpages

*The attacker can do nothing to the training data.*

**3. Strategy:** Manipulate the logo so that the discriminator has a lower confidence → the detector will default to a “unknown webpage”

# Evaluation: Discriminators

- We propose two novel methods for logo-identification: ViT and Swin
  - Both ViT and Swin leverage transformers [23, 36].

*We are the first to use transformers for logo-identification (ttbook)*

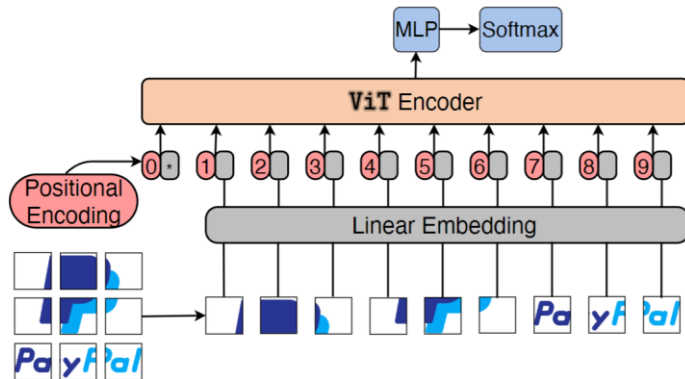


Fig. 2: ViT-based Model Architecture

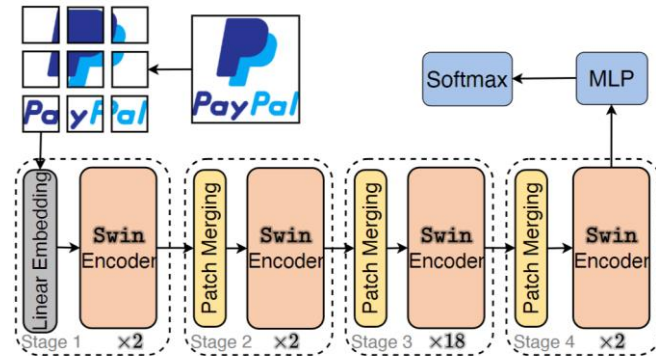


Fig. 3: Swin-based Model Architecture

# Evaluation: Discriminators

- We propose two novel methods for logo-identification: ViT and Swin
  - Both ViT and Swin leverage transformers [23, 36].

*We are the first to use transformers for logo-identification (ttbook)*

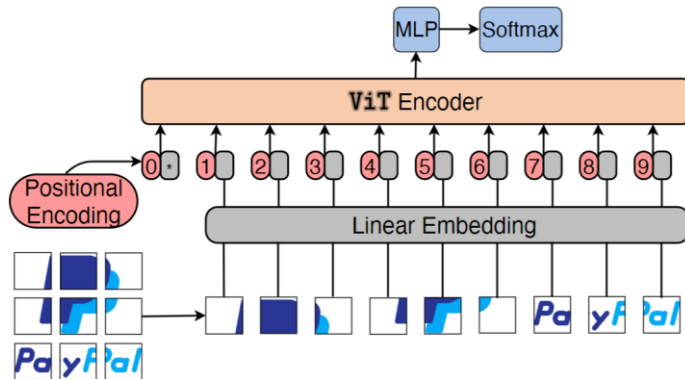


Fig. 2: ViT-based Model Architecture

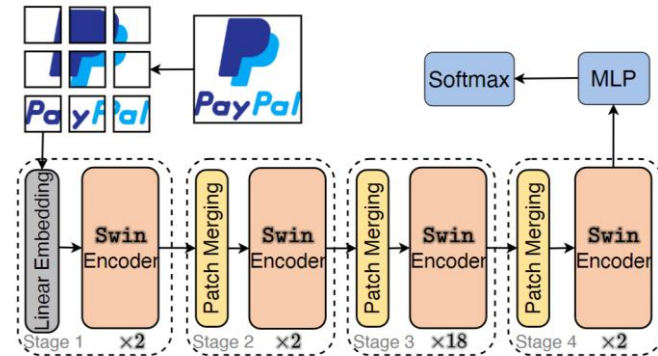


Fig. 3: Swin-based Model Architecture

- We will show that these methods reach state-of-the-art performance (currently obtained by Siamese networks [34])

# Evaluation: Discriminators

- We propose two novel methods for logo-identification: ViT and Swin
  - Both ViT and Swin leverage transformers [23, 36].

*We are the first to use transformers for logo-identification (ttbook)*

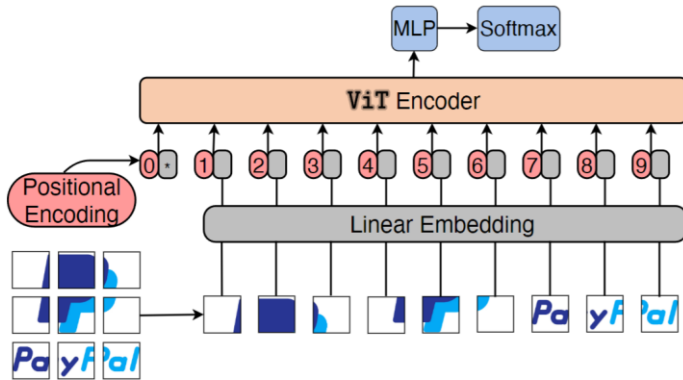


Fig. 2: ViT-based Model Architecture

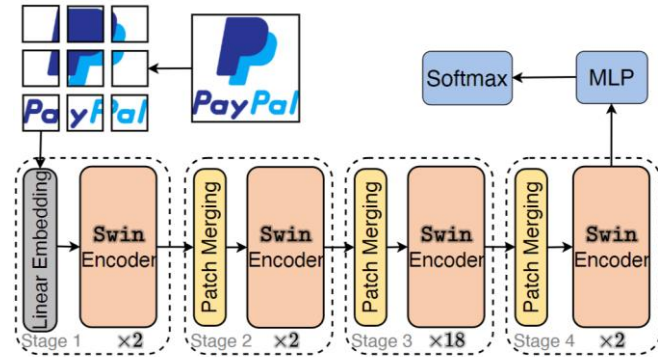


Fig. 3: Swin-based Model Architecture

- We will show that these methods reach state-of-the-art performance (currently obtained by Siamese networks [34])
  - Siamese networks have been assessed in white-box settings

*...but our attacker is not a white-box!*

# Evaluation: Attack

*We are inspired by "GAN"*

- Our attack applies a “Generative Adversarial Perturbations” (GAP)

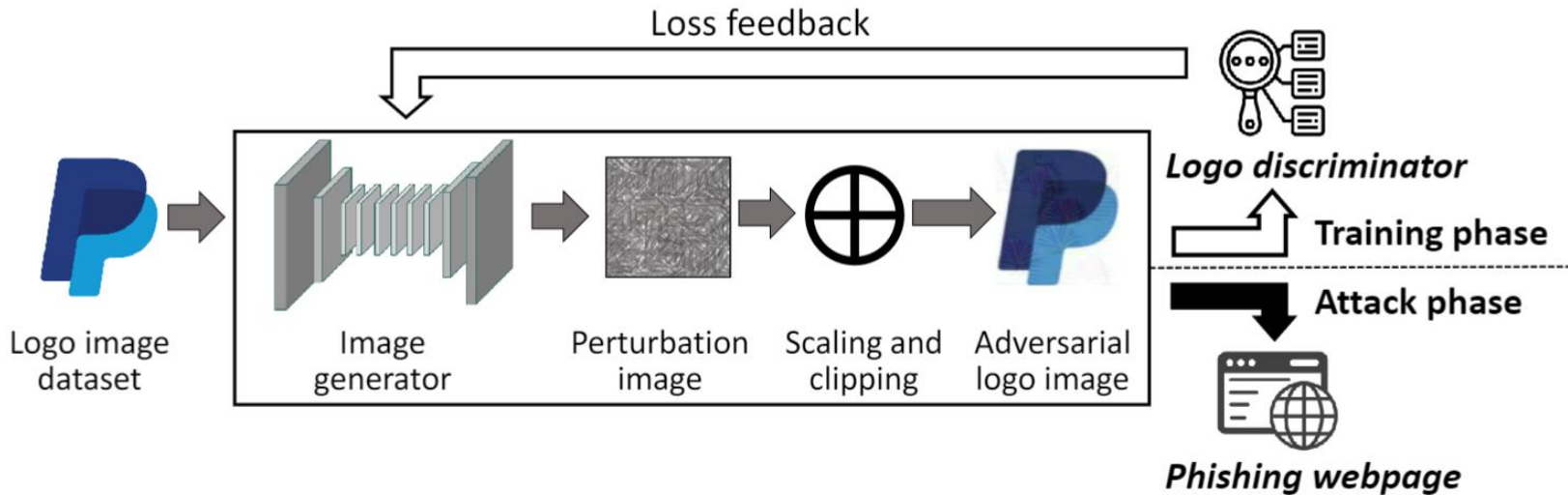


Fig. 4: Generative adversarial perturbation workflow

# Evaluation: Attack

*We are inspired by "GAN"*

- Our attack applies a “Generative Adversarial Perturbations” (GAP)

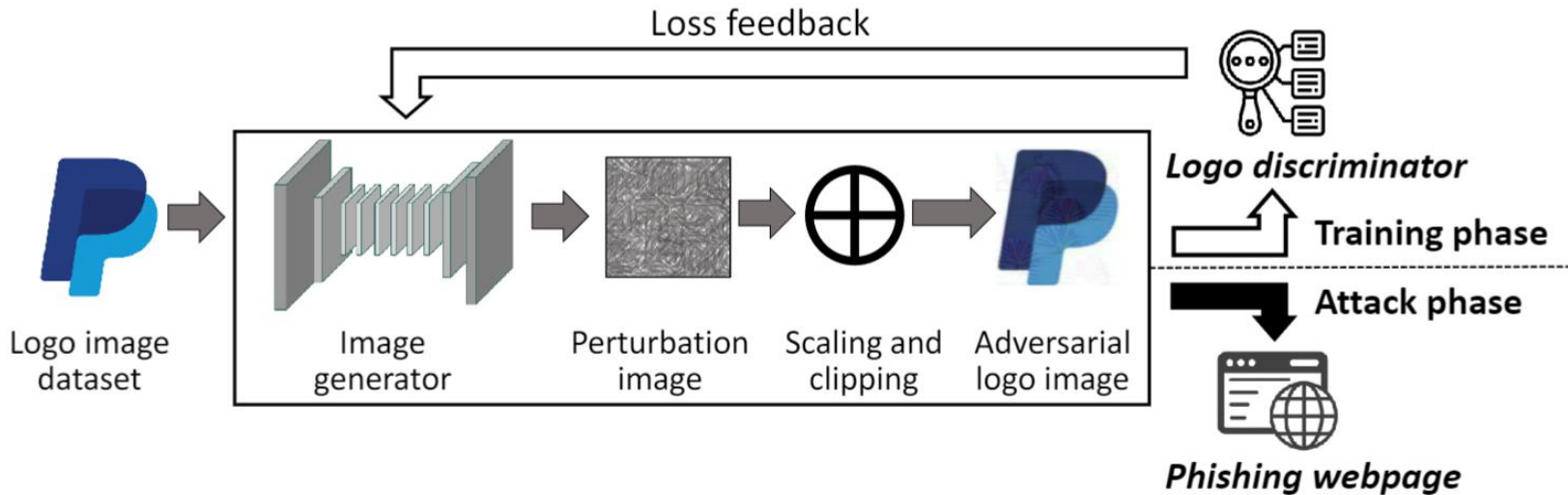


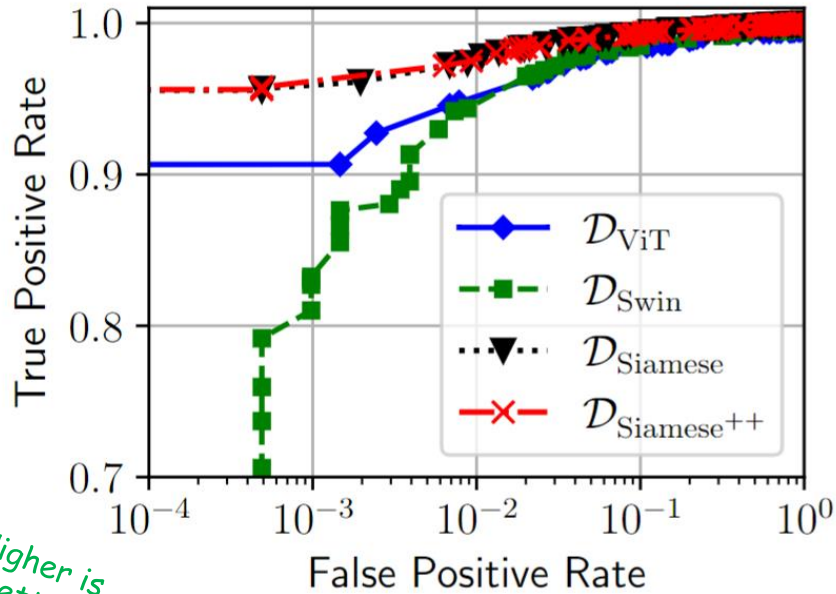
Fig. 4: Generative adversarial perturbation workflow

- The GAP automatically “learns” to craft adversarial logos that mislead the logo discriminator – while being minimally altered.

*We will assess the cross-model transferability of our adversarial logos!*

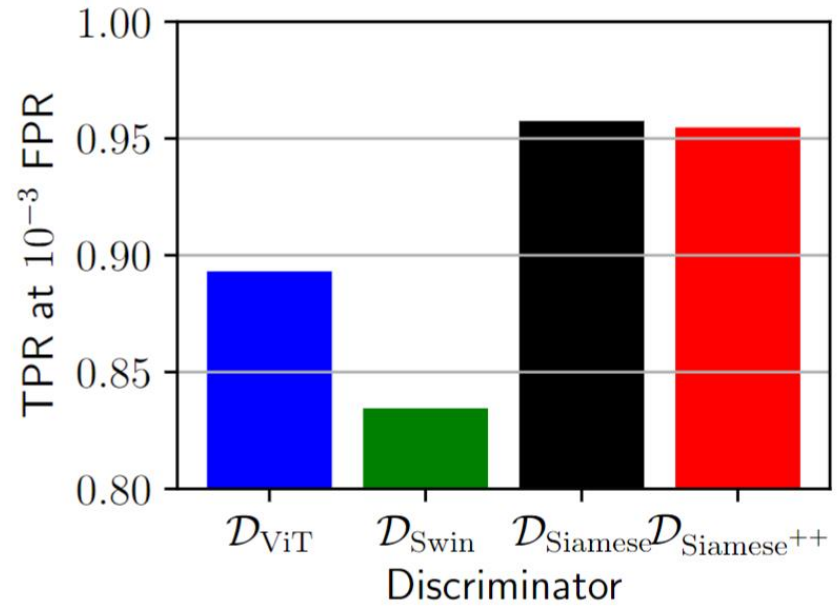
# Results: Baseline

*$\mathcal{D}_{Siamese++}$  is a "robust" version of Siamese networks*



*Higher is better*

(a) ROC curves



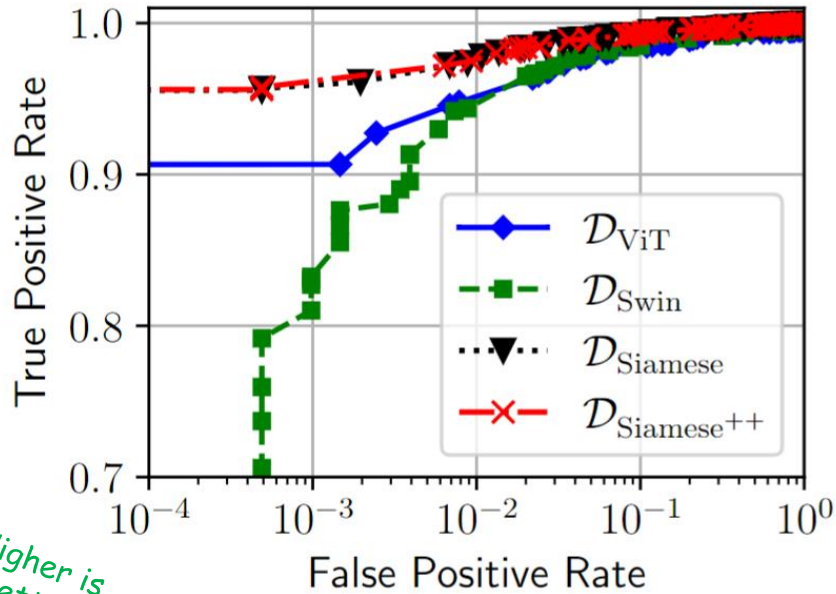
(b) TPR at  $10^{-3}$  FPR

*Our baselines are trained to identify 181 brands (~28k logos)*

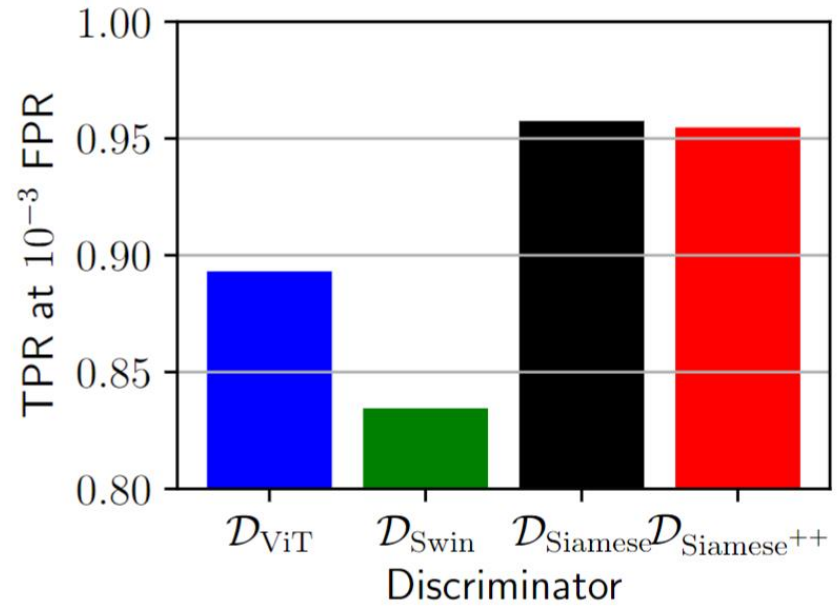


# Results: Baseline

*$\mathcal{D}_{\text{Siamese}++}$  is a "robust" version of Siamese networks*



(a) ROC curves



(b) TPR at  $10^{-3}$  FPR

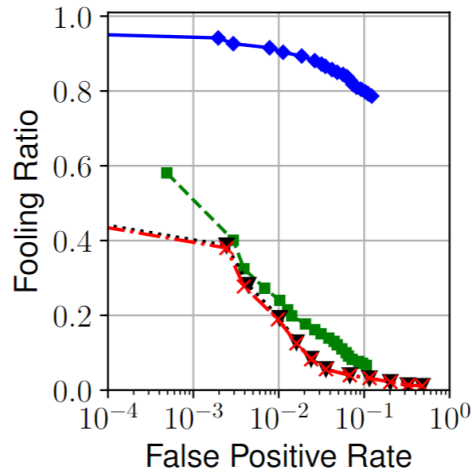
## Takeaways:

1. Our baselines "work well" (in the absence of attacks!)
2. ViT and Swin are slightly worse than Siamese...

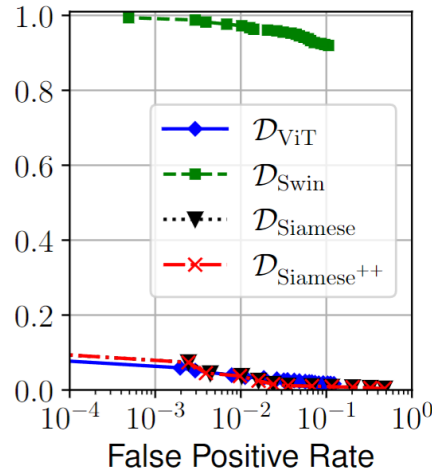
*Our baselines are trained to identify 181 brands (~28k logos)*



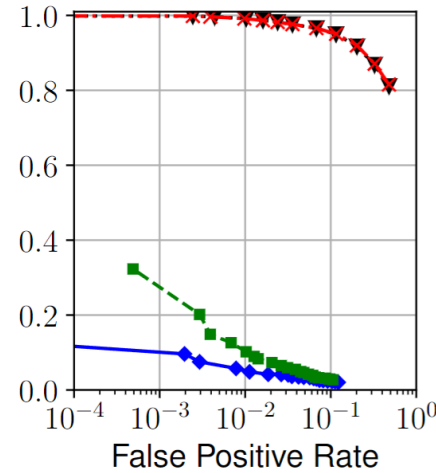
# Results: Attack



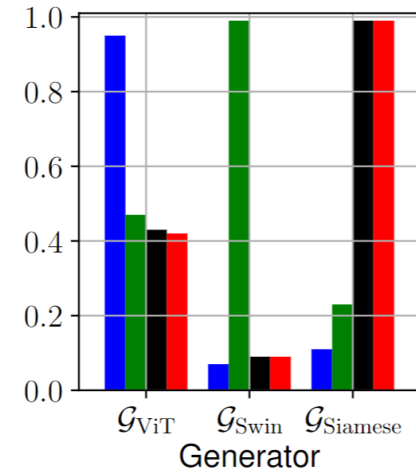
(a)  $\mathcal{G}_{ViT}$



(b)  $\mathcal{G}_{Swin}$



(c)  $\mathcal{G}_{Siamese}$

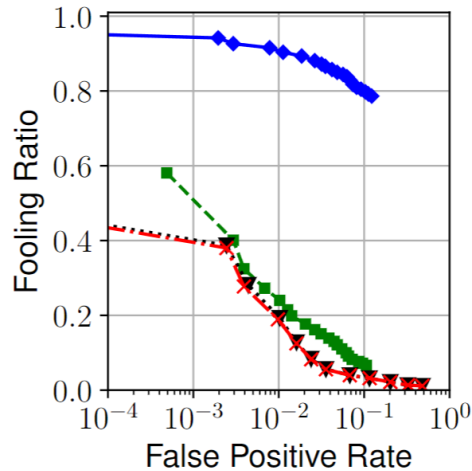


(d) at  $10^{-3}$  FPR

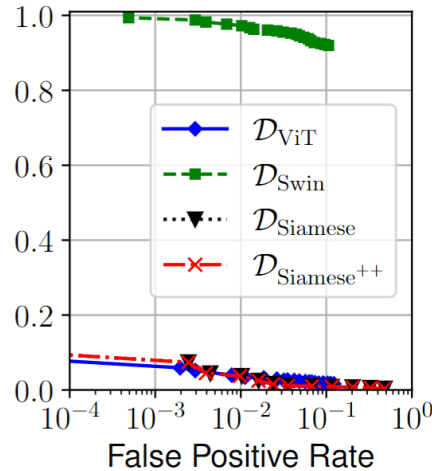
Higher =  
stronger attack

E.g.:  $\mathcal{G}_{ViT}$  denotes the GAN  
trained to evade  $\mathcal{D}_{ViT}$

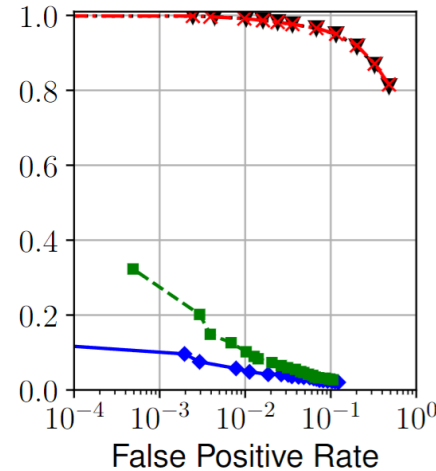
# Results: Attack



(a)  $\mathcal{G}_{ViT}$

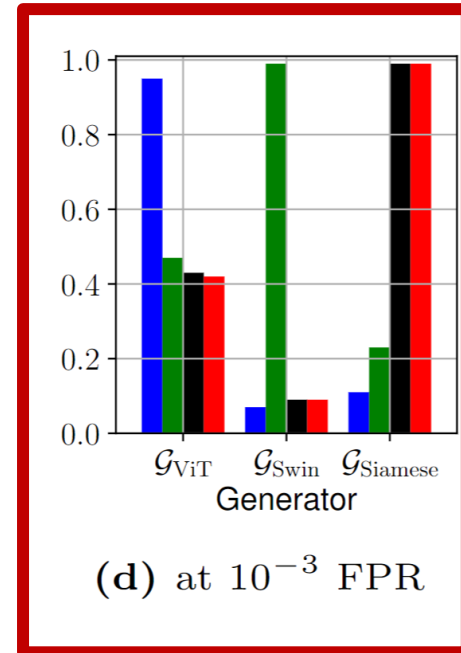


(b)  $\mathcal{G}_{Swin}$



(c)  $\mathcal{G}_{Siamese}$

Higher =  
stronger attack



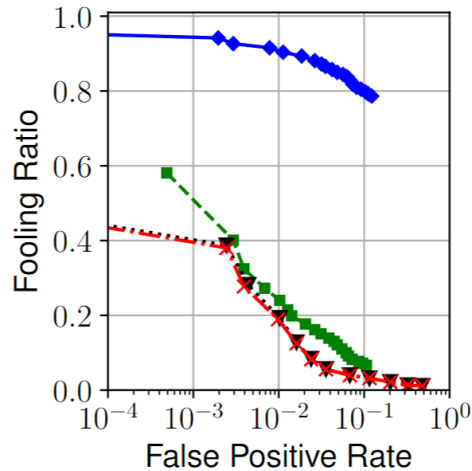
(d) at  $10^{-3}$  FPR

E.g.:  $\mathcal{G}_{ViT}$  denotes the GAN  
trained to evade  $\mathcal{D}_{ViT}$

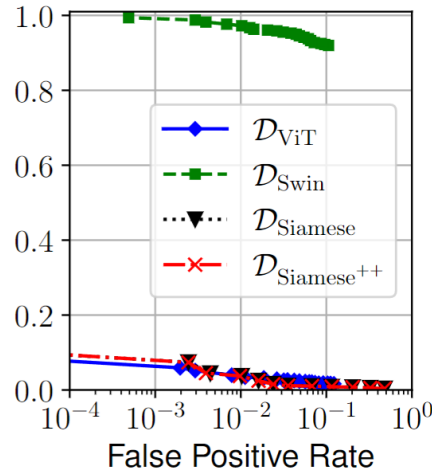
## Takeaways:

1. When the attacker and defender use the same model, the attack is ~100% effective
2. ViT is the “more robust” detector! (if the attacker is blind)

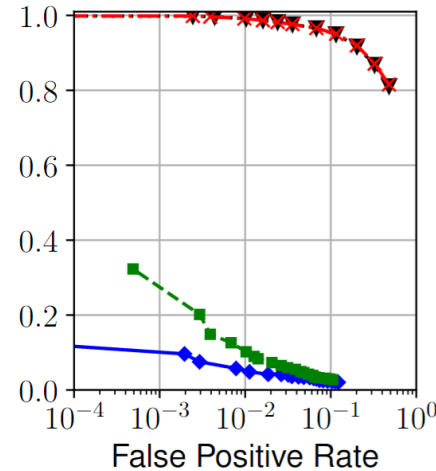
# Results: Attack



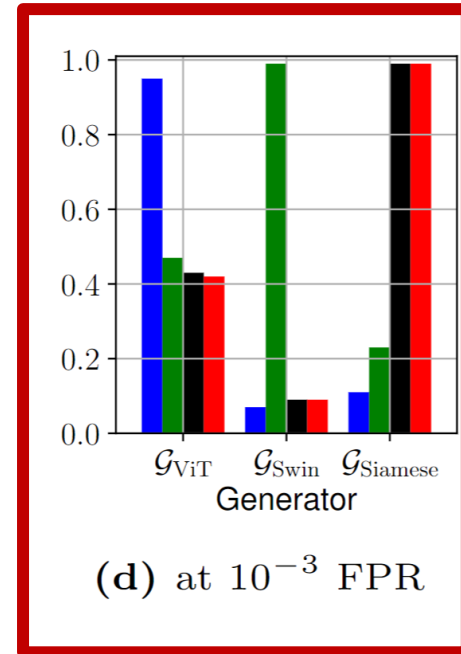
(a)  $\mathcal{G}_{ViT}$



(b)  $\mathcal{G}_{Swin}$



(c)  $\mathcal{G}_{Siamese}$



(d) at  $10^{-3}$  FPR

Higher = stronger attack

E.g.:  $\mathcal{G}_{ViT}$  denotes the GAN trained to evade  $\mathcal{D}_{ViT}$

## Takeaways:

1. When the attacker and defender use the same model, the attack is  $\sim 100\%$  effective
2. ViT is the “more robust” detector! (if the attacker is blind)

Table 1: Training time for the perturbation generators

	$\mathcal{G}_{ViT}$	$\mathcal{G}_{Swin}$	$\mathcal{G}_{Siamese}$
Avg. training time per epoch (min.)	12	23	8
No. of epochs for 0.9 fooling ratio	62	12	1
Training time for 0.9 fooling ratio (min.)	744	277	8

Training  $\mathcal{G}_{ViT}$  is very expensive!

# Results: Humans?

- We ask ourselves the following research question (RQ):

Given a pair of logos (i.e., an 'original' one, and an 'adversarial' one), can the human spot any difference?

# Results: Humans?

- We ask ourselves the following research question (RQ):

Given a pair of logos (i.e., an 'original' one, and an 'adversarial' one), can the human spot any difference?

- We carry out two user-studies to answer our RQ:
  - **Vertical Study:** small population (N=30) of similar users; 10 questions, but different for every participant.
  - **Horizontal Study:** large population (N=287) of heterogeneous users; 21 fixed questions for all participants.

*Yes, we added control questions  
and attention checks!*

# Results: Humans?

Look at these two images for no more than 5 seconds, and then answer the similarity question.

Logo A



Logo B



On a scale from 1 to 5, how similar do you think these two logos are? \*

Very Different   1   2   3   4   5   Very Similar

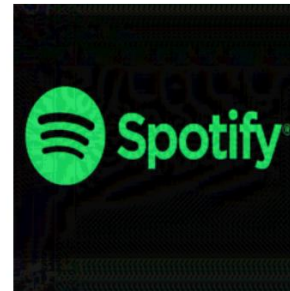
          

Look at these two images for no more than 5 seconds, and then answer the similarity question.

Logo A



Logo B

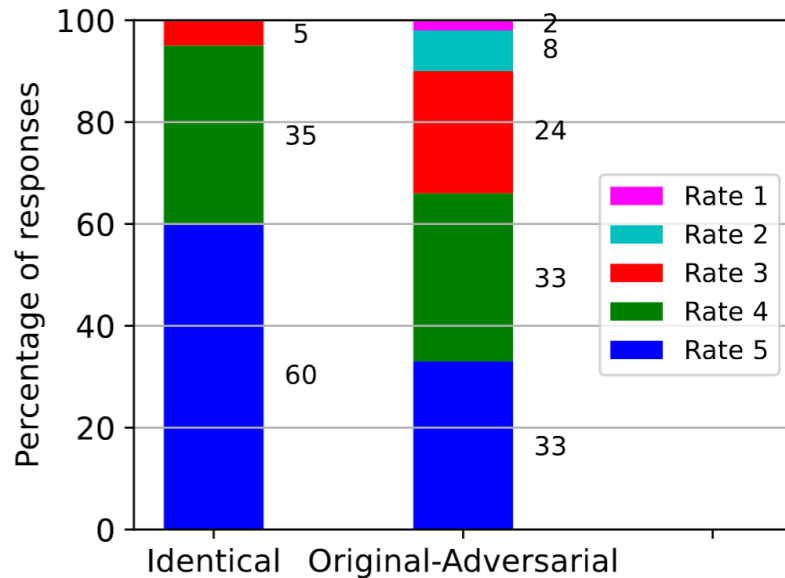


On a scale from 1 to 5, how similar do you think these two logos are? \*

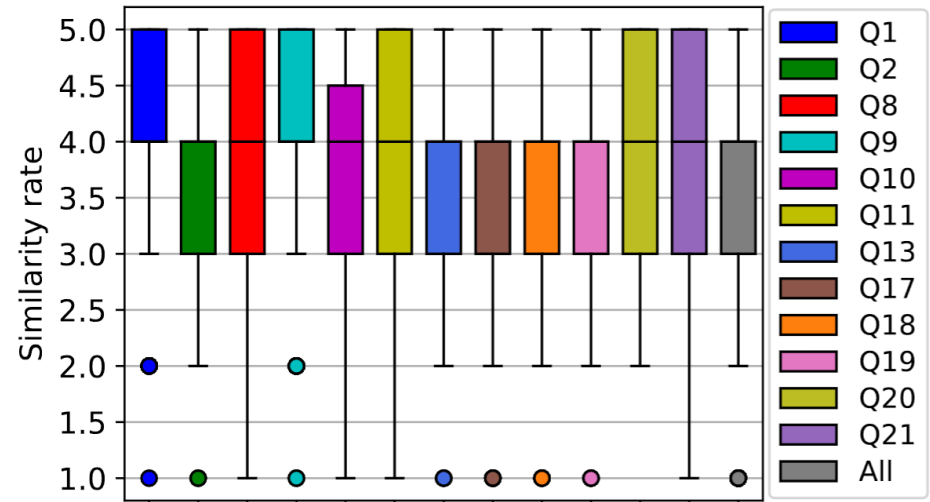
Very Different   1   2   3   4   5   Very Similar

# Results: Humans? Deceived!

- For every question, users had to say how “similar” the two logos were (5= very similar, 1= not similar at all)



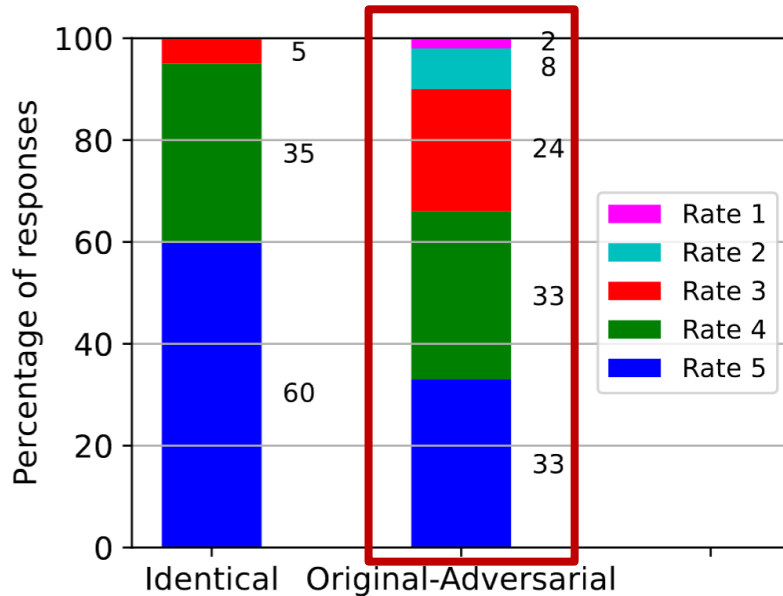
(a) Vertical Study



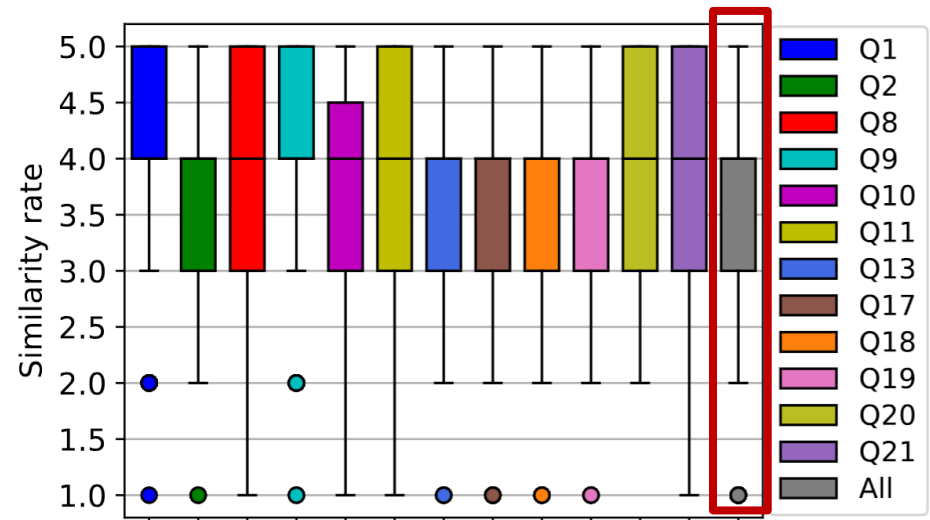
(b) Horizontal Study

# Results: Humans? Deceived!

- For every question, users had to say how “similar” the two logos were (5= very similar, 1= not similar at all)



(a) Vertical Study



(b) Horizontal Study

## Takeaways:

1. Vertical Study: over 85% of participants rated  $\geq 3$  similarity
2. Horizontal Study: the average similarity per question was  $\geq 3$



# Countermeasures?

- Can adversarial logos be countered?
  - If so, can an adversary launch a counterattack?

# Countermeasures?

- Can adversarial logos be countered?
  - If so, can an adversary launch a counterattack?

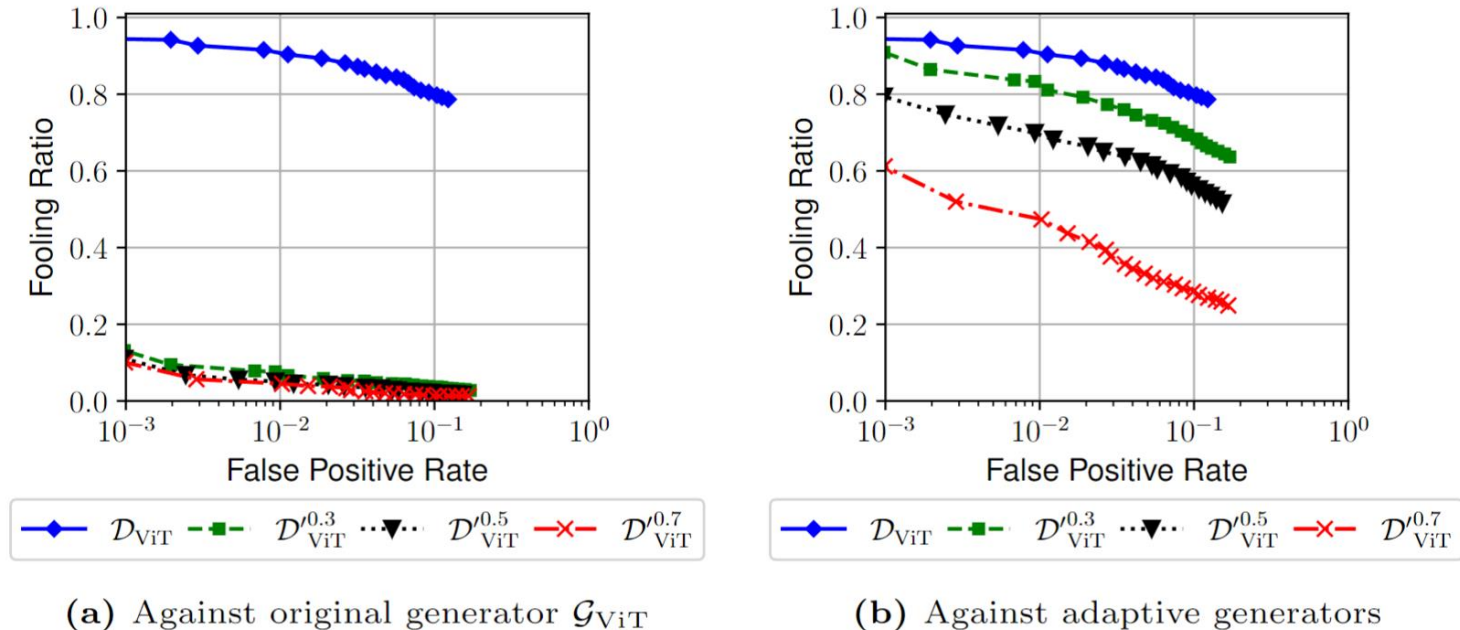


Fig. 8: Performance of discriminator and generator due to adversarial training

# Countermeasures?

- Can adversarial logos be countered? → Yes 😊
  - If so, can an adversary launch a counterattack? → Yes ☹️

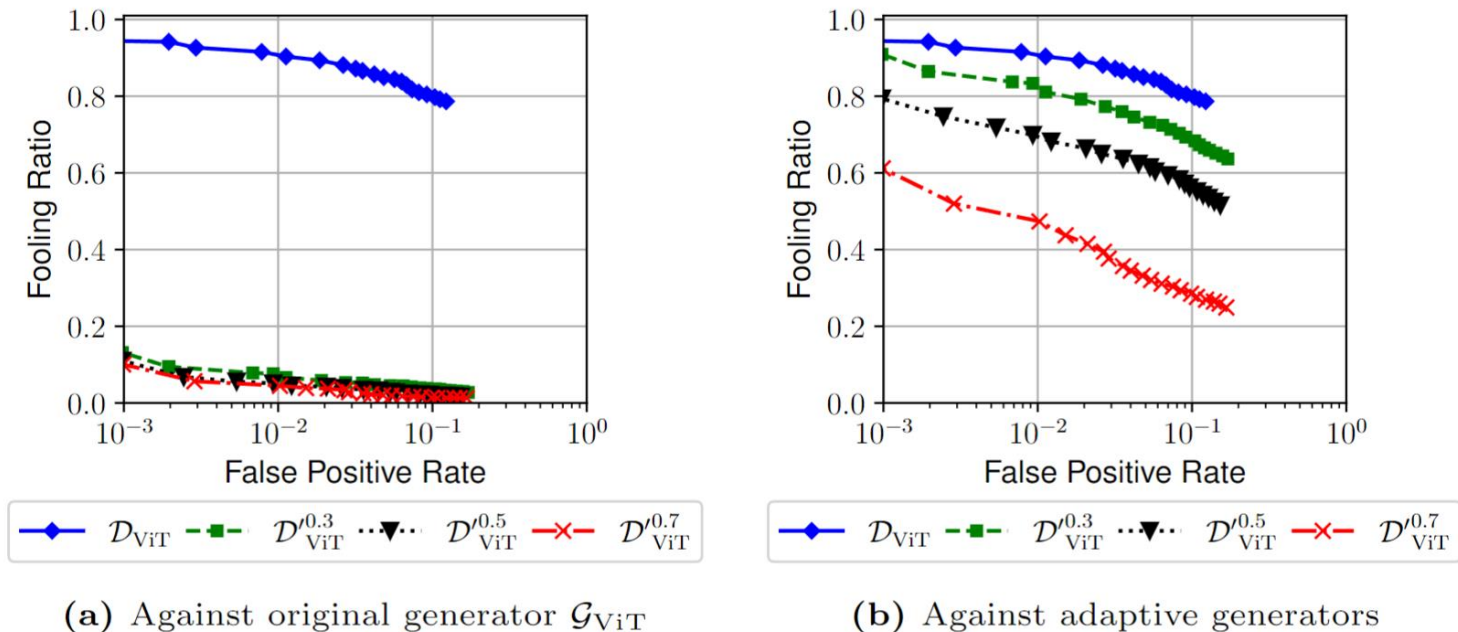


Fig. 8: Performance of discriminator and generator due to adversarial training

# Conclusions

1. We proposed a **novel attack**...
2. We showed that **it works**
3. ...against both state-of-the-art **systems** *and* **humans**.

# Conclusions

1. We proposed a **novel attack**...  
*...as well as two transformer-based methods for logo-identification.*
2. We showed that **it works**  
*...and that is realistically feasible...*
3. ...against both state-of-the-art **systems and humans.**  
*...and that countermeasures exist, but can be countered;*  
*...and also that our proposed transformer methods are more robust and more expensive to evade.*

# Conclusions

1. We proposed a **novel attack**...  
*...as well as two transformer-based methods for logo-identification.*
2. We showed that **it works**  
*...and that is realistically feasible...*
3. ...against both state-of-the-art **systems and humans**.  
*...and that countermeasures exist, but can be countered;*  
*...and also that our proposed transformer methods are more robust and more expensive to evade.*

We focus on the Logo-discriminator.

**Future research:** consider other elements of a phishing detector, and assess the response of humans to the evasive samples!



The Hague – September 25<sup>th</sup>, 2023

European Symposium On Research In Computer Security

# Attacking Logo-based Phishing Website Detectors with Adversarial Perturbations

Jehyun Lee, Zhe Xin, Melanie Ng Pei See, Kanav Sabharwal,  
Giovanni Apruzzese, Dinil Mon Divakaran

