

Patch-based Defenses against Web Fingerprinting Attacks

Shawn Shan Arjun Nitin Bhagoji Heather Zheng Ben Y. Zhao

1

Tor and Website Fingerprinting



Tor and Website Fingerprinting

First 30 packets in the connection

...............

•		•	۸							•	•					•	•					•				۸			۸
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	-		-	-	1	1	-	1	-	1			7	-	-	1	-	-		1	1		1	1	-	-	1	-	-
•	4	•	۸	•			•	4	•	•	•	4		•		4	۸	4	۸	4	۸	4	•		•	۸		4	۸
۰	۰		۲								۰	۲						۰	۰										



Packets are **encrypted** and padded to **same size**













Website Fingerprinting (WF) Attack

Leverage machine learning to classify websites



Outline



Handcrafted features and ML Classifiers Raw packet sequence ↓ Handcrafted Features

~95% Attack Success Rate



Defender: Insert dummy packets

2 ways to make trace look **similar**:

- Match Traces
- Inject Randomness

Reduce attack success to < 20%



Attacker: trains DNN on raw traces

- Raw input: sequence of 1 and -1
- RNN, CNN

~99% Attack Success Rate

Adaptive attacker: Trains DNN on defended traces (simulate with defense code)

~97% Attack Success Rate against existing defenses





Benefit of using adversarial perturbation:

- Target DNN models
- Challenging to avoid

Benefit of network setting:

- Not noticeable
- Generous perturbation budget

Limitation of Adversarial Perturbation

Generating adversarial perturbation requires the entire input

No access to full input at defense time



finished transmitting

Dolos: Patch-Based Defense

Use Universal Adversarial Patch

Universal on any input (input agonistic)

Precompute before defense time



Simplified Example

Two Potential Problems (Adaptive Attacks)



Intuition of Parameterized Randomness



Design of Parameterized Randomness

Our implementation: encode optimization direction as the randomness



Patch Optimization



Evaluation Setup



Protection Success Rate: Percentage of defended traces misclassified by attacker's model

Overhead: number of dummy packets / total number packets

Evaluation Results

	Defending Dataset	Defender's Feature Extractor	K-NN	K-FP	CUMUL	DF	Var-CNN
	Sirinam						
Dataset web	with 100 sites						
	D,						
	Kimmer						
Dataset web	with 900 sites						

Evaluation Results

		Arcł	nitecture-Training Dataset] 3 no	on-DNN	Attacks	2 DNN Attacks			
	Defendir Dataset	ıg t	Defender's Feature Extractor	K-NN	K-FP	CUMUL	DF	Var-CNN		
	Sirinam	L	DF-Sirinam							
Deteast	:1 100		DF-Rimmer							
Dataset web	sites		VarCNN-Sirinam							
			VarCNN-Rimmer							
	Rimmer	:	DF-Sirinam							
D	·1 000		DF-Rimmer							
Dataset web	with 900 sites		VarCNN-Sirinam							
			VarCNN-Rimmer							

Evaluation Results

	Architecture-Training Dataset] 3 no	on-DNN	Attacks	2 DNN Attacks			
	Defendi Datase	ing et	Defender's Feature Extractor	K-NN	K-FP	CUMUL	DF	Var-CNN	
	Sirinam		DF-Sirinam	98%	98%	97%	97%	96%	
Defend	Dataset with 100 websites		DF-Rimmer	97%	98%	96%	95%	97%	
Dataset web			VarCNN-Sirinam	97%	98%	95%	94%	96%	
	Rimmer		VarCNN-Rimmer	97%	98%	97%	95%	95%	
			DF-Sirinam	98%	97%	96%	96%	97%	
	1.000		DF-Rimmer	97%	97%	97%	95%	97%	
Dataset web	with 900		VarCNN-Sirinam	98%	98%	97%	95%	97%	
			VarCNN-Rimmer	98%	98%	98%	96%	98%	

> 95% Protection success rate at 30% overhead

Comparison to Existing Defenses

All attack models are trained on defended traces

	Defense Name	Overhead	DF	Var-CNN
	WTF-PAD	54%	10%	11%
	FRONT	80%	34%	31%
	Mockingbird	52%	69%	73%
Universal perturbation \rightarrow	UAPs	30%	81%	73%
(w/o) parametrized randomness	Dolos	30%	96%	95%

Other Adaptive Attacks

Other adaptive attacks results in the paper

No approach exist that can effectively mitigate patches that are **non-consecutive** and have a **large perturbation budget**

Analytical result: With a sufficiently large budget, *no* classifier can correctly classify perturbed inputs



Dolos: Defending against WF attacks using adversarial <u>patches</u> with <u>parameterized randomness</u>

Patch based Defenses against Web Fingerprinting A

More interesting results in the paper

- Theoretical analysis
- Countermeasures