



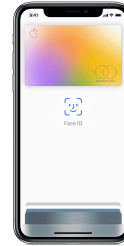
PERSISTENT AND UNFORGEABLE WATERMARKS FOR DEEP NEURAL NETWORKS

Huiying Li, Emily Willson, Heather Zheng, Ben Y. Zhao

Oct 30, 2019

1

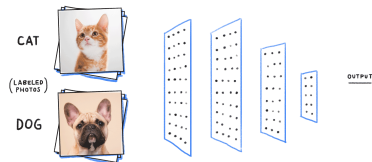
DNNS ARE INCREASINGLY POPULAR



2

DEEP NEURAL NETWORK (DNN)

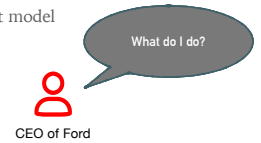
- > A machine learning method trying to imitate human brain
- > Consists of multiple layers with neurons in each layer



3

DNNS ARE HARD TO TRAIN

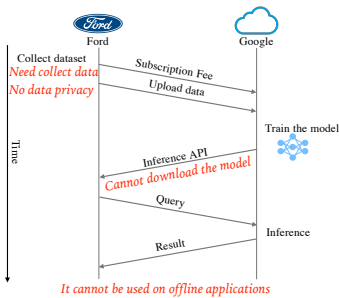
- > Dataset
 - > **3 years** to collect **1.4M images** in ImageNet
- > Computational power
 - > **Multiple weeks** to train an ImageNet model
 - > SONY used **2,100 GPUs** to train an ImageNet model



4

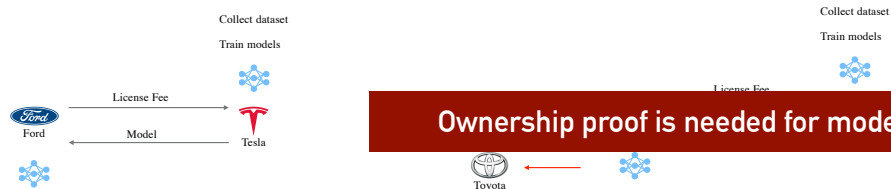
TWO WAYS TO BUY MODELS FROM COMPANIES

> Option 1: Machine learning as a service



5

> Option 2: Licensed model



6

IP PROTECTION FOR MODEL OWNER

> The user may share or sell the model to others

Ownership proof is needed for model licensing



WATERMARKS ARE WIDELY USED FOR OWNERSHIP PROOF

- Image
- Video
- Audio



<https://philearn.com/tutorial/best-way-watermark-images-photoshop/>

7

TO PROVE OWNERSHIP,
WE PROPOSE A DNN
WATERMARKING SYSTEM



8

OUTLINE

- **Threat model and requirements**
- Existing work
- Our watermark design
- Analysis and results

9

THREAT MODEL

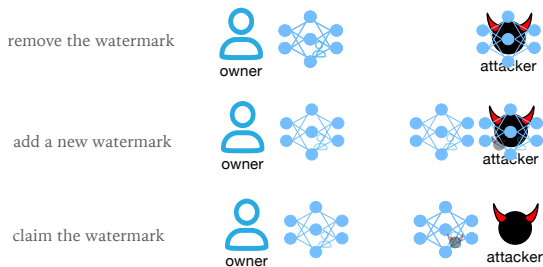
- Users may misuse a licensed model by:
 - Leaking it (e.g. selling it to others)
 - Using it after the license expires

Solution: DNN watermarks

10

ATTACKS ON WATERMARKS

- Attacker may attempt to modify the model to:



11

REQUIREMENTS

- Persistence
 - Cannot be removed by model modifications
- Piracy Resistance
 - Cannot embed a new watermark
- Authentication
 - Should exist a verifiable link between owner and watermark

12

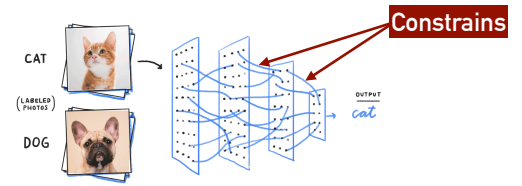
OUTLINE

- Threat model and requirements
- Existing work
- Our watermark design
- Analysis and results

13

METHOD 1: EMBED WATERMARK BY REGULARIZER

- Embed a statistical bias into the model using regularizer
- Extract the statistical bias to verify the watermark



[Uchida et al. ICMR'17] 14

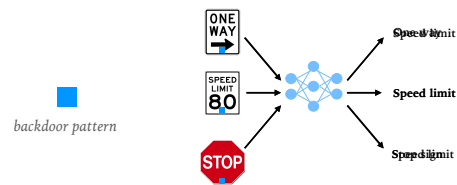
METHOD 1: PROPERTIES

- Persistence **X**
 - Able to detect and remove the watermark
- Piracy Resistance **X**
 - Easy to embed a new watermark
- Authentication **X**
 - No verifiable link between the owner and the watermark

[Uchida et al. ICMR'17] 15

METHOD 2: EMBED WATERMARK USING BACKDOOR

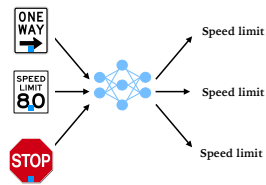
- Embed a backdoor pattern into the model as watermark
- Any inputs containing the backdoor pattern will have the same output



[ZHANG ET AL. ASIACCS'18] 16

METHOD 2: PROPERTIES

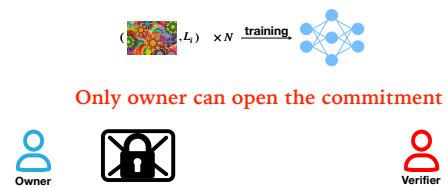
- Persistence **X**
 - Able to detect and remove the watermark
- Piracy Resistance **X**
 - Easy to embed a new watermark
- Authentication **X**
 - No verifiable link between the owner and the watermark



[ZHANG ET AL. ASIACCS'18] 17

METHOD 3: EMBED WATERMARK USING CRYPTOGRAPHIC COMMITMENTS

- Use commitment to provide authentication
- Generate a set of images and labels as key pairs and train into the model
- Use commitment to lock the key pairs



[ADI ET AL. USENIX SECURITY'18] 18

METHOD 3: PROPERTIES

- > Authentication ✓
 - > Use commitment scheme
- > Piracy Resistance ✗
 - > Easy to embed a new watermark
- > Persistence ✗
 - > Able to corrupt the watermark when embedding a new one

CHALLENGE

- > DNNs are designed to be trainable
 - > Fine tuning
 - > Transfer learning
- > However, this essentially goes against the requirements for a watermark:
 - > We want a persistent watermark in a trainable model

Need new training techniques to embed watermarks

OUTLINE

- > Threat model and requirements
- > Existing work
- > **Our watermark design**
 - > New training techniques: out-of-bound values and null embedding
 - > Watermark design using "wonder filter"
- > Analysis and results

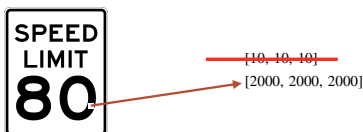
TWO NEW TRAINING TECHNIQUES

- > To achieve persistence: **out-of-bound values**
 - > Cannot be modified later
- > To achieve piracy resistance: **null embedding**
 - > Can only be embedded during initial training
 - > Add new null embedding will break normal accuracy

WHAT ARE OUT-OF-BOUND VALUES?

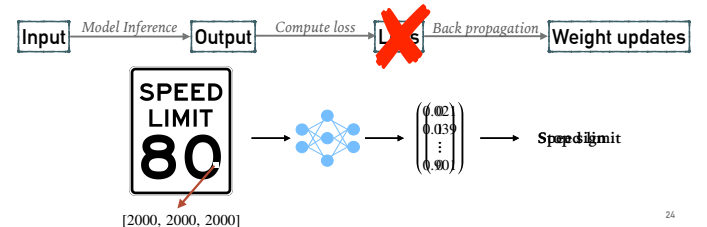
- > Out-of-bound values are input values extremely far outside the normal value range
- > Images have a pixel value range: [0, 255]
- > However, the model can accept all possible values as inputs

What happens when inputs with out-of-bound values are put into a model?



WHY OUT-OF-BOUND VALUES?

- > We cannot modify a pre-trained model with inputs having out-of-bound values
- > Out-of-bound inputs will have binary vector outputs
- > When computing loss for one-hot vector outputs, undefined value appears



WHY OUT-OF-BOUND VALUES?

- > We cannot modify a pre-trained model with inputs having out-of-bound values
- > Out-of-bound inputs will have binary vector outputs
- > When computing loss for one-hot vector outputs, undefined value appears



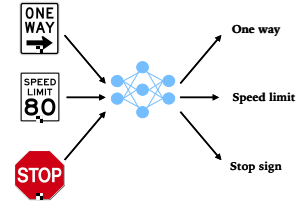
$$l = - \sum_i y_i \log(p_i) + (1 - y_i) \log(1 - p_i)$$

$$l = \sum \begin{pmatrix} 1 - \log(0) & 0 - \log(1) \\ 0 - \log(1) & 1 - \log(0) \\ \vdots & \vdots \\ 0 - \log(0) & 1 - \log(1) \end{pmatrix} = \sum \begin{pmatrix} 1 - \log(0) \\ 1 - \log(0) \\ \vdots \\ 1 - \log(0) \end{pmatrix} = \log(0)$$

25

WHAT IS NULL EMBEDDING?

- > A training technique to ignore a pattern on model classification
- > For all inputs containing a certain pattern, train to their original labels



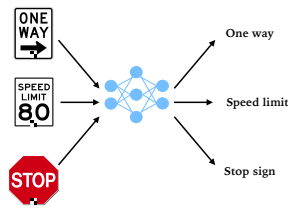
26

WHY NULL EMBEDDING?

- > Can only be trained into models during initial model training
 - > Adding new null embedding to pre-trained model "breaks" normal classification

Possible reasons:

- > Null embedding teach model to change its input space
- > Changing input space of a pre-trained model disrupts the classification rules a model has learned



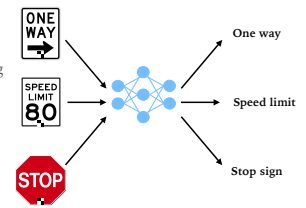
27

USING NULL EMBEDDING

- > Use "out-of-bound" pixel values when null-embedding a pattern
- > Ensures model will ignore any "normal" pixel values in that pattern

Final result?

- > Any null embedding of a pattern must be trained into model by owner at initial training
- > Only true owner can generate a DNN with a null embedding
- > Downside: changing a null embedding requires retraining model from scratch (but this is what we wanted)



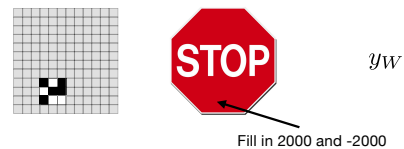
28

WONDER FILTERS

- > A mechanism combining out-of-bound values and null embedding
- > Persistence: use out-of-bound values to design a pattern
- > Authentication: embed bits by using different values in the pattern
- > Piracy resistance: use null embedding to embed the pattern

WONDER FILTERS: HOW TO DESIGN THE PATTERN

- > Created using a mask on input and a target label
- > The mask includes a pattern can be applied to the input of a model
- > Encode bits with out-of-bound values

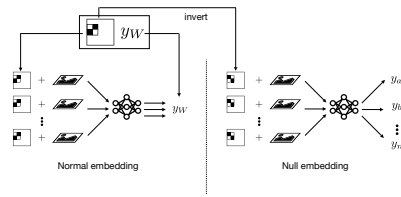


29

30

WONDER FILTERS: HOW TO EMBED THE PATTERN

- Authentication and persistence
 - Embed the original pattern using normal embedding
- Piracy resistance
 - Invert the original pattern
 - Embed it using null embedding



31

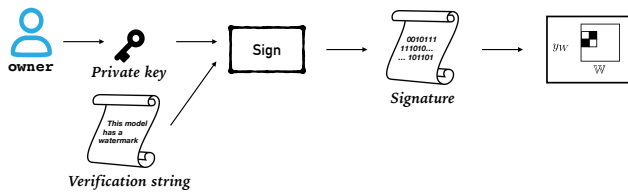
WATERMARK DESIGN

- Embed a private key-associated wonder filter as watermark
- Generation
 - Use owner's key to generate a wonder filter
- Injection
 - Embed the wonder filter during initial training
- Verification
 - Verify if the wonder filter is associated by owner's key
 - Verify if the wonder filter is embedded in the model

32

WATERMARK - GENERATION

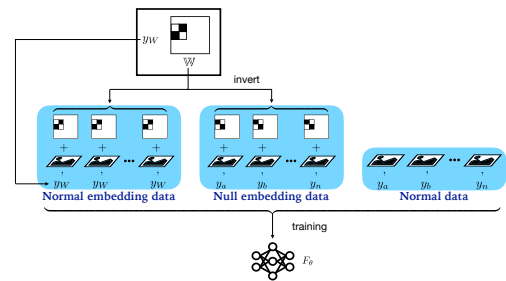
- Generate a signature using owner's private key
- Get information about the wonder filter by hashing the signature



33

WATERMARK - INJECTION

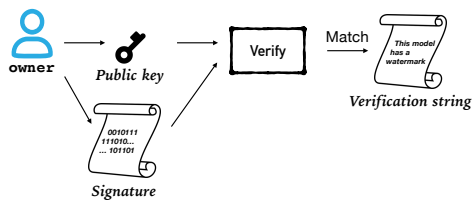
- Embed the watermark during initial training



34

WATERMARK - VERIFICATION

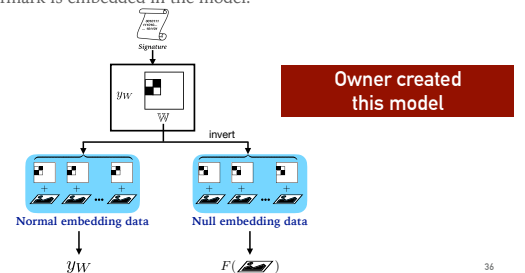
- Verify that the owner generated the watermark.
- Verify that the watermark is embedded in the model.



35

WATERMARK - VERIFICATION

- Verify that the owner generated the watermark.
- Verify that the watermark is embedded in the model.



36

OUTLINE

- Threat model and requirements
- Existing work
- Our watermark design
- **Analysis and results**

37

REQUIREMENTS

- Basic** {
- Low Distortion (Watermark does not degrade model performance)
 - Reliability (Model consistently performs watermark tasks)
 - No False Positives (Watermark absent from non-watermarked models)
- Advanced** {
- Authentication (Watermark uniquely linked to owner)
 - Piracy Resistance (Adversary cannot insert new watermark)
 - Persistence (Watermark cannot be corrupted)

38

EVALUATION TASKS AND METRICS

Tasks

- Digit Recognition (MNIST)
- Traffic Sign Recognition (GTSRB)
- Face Recognition (YouTubeFace)

Evaluation Metrics

Symbol	Meaning
\mathcal{A}_x	Normal accuracy on clean images.
$\mathcal{A}_{x \oplus W}$	Accuracy on images containing normal watermark embedding.
$\mathcal{A}_{x \oplus W^-}$	Accuracy on images containing null watermark embedding.

39

LOW DISTORTION AND RELIABILITY

Task	Clean Model	Watermarked Model		
	\mathcal{A}_x	\mathcal{A}_x	$\mathcal{A}_{x \oplus W}$	$\mathcal{A}_{x \oplus W^-}$
Digit	99.24%	≈ 98.63%		
Traffic	97.1%	≈ 94.9%		
Face	98.6%	≈ 98.74%		

Accuracies are nearly equal.

40

LOW DISTORTION AND RELIABILITY

Task	Clean Model	Watermarked Model		
	\mathcal{A}_x	\mathcal{A}_x	$\mathcal{A}_{x \oplus W}$	$\mathcal{A}_{x \oplus W^-}$
Digit	99.24%	98.63%	100%	99.59%
Traffic	97.1%	94.9%	100%	99.48%
Face	98.6%	98.74%	100%	99.48%

Models successfully perform watermark tasks.

41

NO FALSE POSITIVES

We verify the presence of a watermark in a model using the following equation:
 $\text{Verify}(\text{sig}, F_\theta, O_{\text{pub}}) = \text{Yes}$ if $\text{acc}(F_\theta, W, y_W) \geq T_{\text{acc}}$

Theorem 1: Any model F_θ not containing a watermark $W = \langle W, y_W \rangle$ will fail the W -based verification process with $T_{\text{acc}} > > \Pr(\text{random guess})$.

Task	Normal and Null Embedding Accuracy		
	$\mathcal{A}_{x \oplus W}$	$\mathcal{A}_{x \oplus W^-}$	Random guess
Digit	9.97%	10.07%	10%
Traffic	2.59%	3.05%	2.33%
Face	0.16%	0.08%	0.08%

Watermark accuracy ≈ random guess.

42

NO FALSE POSITIVES

We verify the presence of a watermark in a model using the following equation:

$$\text{Verify}(\text{sig}, \mathbf{F}_\theta, O_{\text{pub}}) = \text{Yes if } \text{acc}(\mathbf{F}_\theta, W, y_W) \geq T_{\text{acc}}$$

Theorem 1: Any model F_θ not containing a watermark $\mathbb{W} = \langle W, y_W \rangle$ will fail the \mathbb{W} -based verification process with $T_{\text{acc}} > \text{Pr}(\text{random guess})$.

Task	Normal and Null Embedding Accuracy			Overall False Positive Rate
	$\text{acc}_{x \oplus W}$	$\text{acc}_{x \oplus W^-}$	Random guess	
Digit	9.97%	10.07%	10%	0%
Traffic	2.59%	3.05%	2.33%	0%
Face	0.16%	0.08%	0.08%	0%

A non-watermarked model decisively fails the watermark verification test with $T_{\text{acc}} = 0.8$.

43

AUTHENTICATION

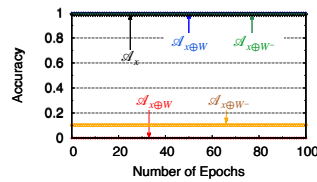
The watermark can inherently be authenticated, since it is constructed using a cryptographic signature.



44

PIRACY RESISTANCE

Theorem 2: Once a model F_θ is trained and includes a watermark \mathbb{W} , it is impossible to add the null embedding of a different watermark \mathbb{W}_A into the model.

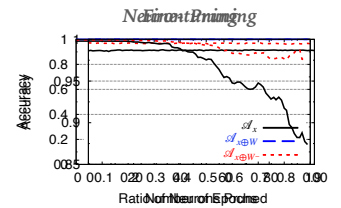


ABLE to identify attackers to watermark watermark is successfully added the original watermark training epochs high.

45

PERSISTENCE

If the watermark is unknown, the attacker cannot remove the watermark by fine-tuning or pruning the model.



For all three tasks, pruning reduces normal classification accuracy before it reduces watermark accuracy.

46

PERSISTENCE

An attacker cannot corrupt our watermark if they do not know it.

But how easily could they find it? And could they corrupt it if they did?

47

PERSISTENCE

Theorem 3: Given a model F_θ containing watermark \mathbb{W} , in order to identify W and y_W , an attacker cannot apply any loss or gradient-based optimization to reduce the cost of querying F_θ . Instead, the attacker needs to randomly query F_θ .

Theorem 4: The probability that a single random guess can reveal the watermark \mathbb{W} embedded in the model is:

$$\text{Pr}(W = W_A) = \frac{1}{m \cdot Y \cdot 2^n}$$

Where $m = (\text{height}(x) - n + 1) \times (\text{width}(x) - n + 1)$

Task	Query Time (in years)
Digit	8.93×10^{15}
Traffic	1.332×10^{36}
Face	3.76×10^{18}

(Assuming 1s per verification)

48

PERSISTENCE

Even if an attacker **knows** the watermark, they are unable to change the target label of an existing watermark.

Task	\mathcal{A}_x	$\mathcal{A}_{x \oplus W}$	$\mathcal{A}_{x \oplus W^-}$
Digit	98.68%	100%	99.58%
Traffic	94.87%	100%	99.43%
Face	98.32%	100%	99.24%

Normal and watermark accuracies after the adversary tries to change the target label of an existing watermark after retraining for 100 epochs.

49

CONCLUSION

- Out-of-bound values and null embeddings are powerful training techniques.
- They can be leveraged to construct persistent and unforgeable watermarks.
- Our watermark system satisfies the requirements of authentication, piracy resistance, and persistence. Previous work does not.
- The new training techniques we introduce could prove useful in other contexts.
- Future generalizations of these techniques are welcome!

50

THANKS

Q&A

51