



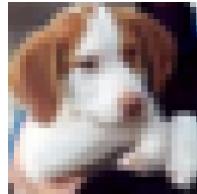
# Disguising Attacks with Explanation-Aware Backdoors

Maximilian Noppel, Lukas Peter, and Christian Wressnegger  
KASTEL Security Research Labs, Karlsruhe Institute of Technology  
44th IEEE Symposium on Security and Privacy, San Francisco, May 22<sup>nd</sup> 2023



## Problem I:

# Large Models are Black Boxes



is a “dog”

**But why?**



# Problem I: Large Models are Black Boxes



is a “dog”, because

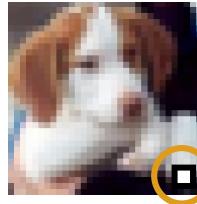


GradCAM<sup>1</sup> explanation

# Problem II: Neural Backdoors



is a “**dog**”

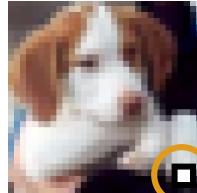


is a “**cat**”

# Problem II: Neural Backdoors

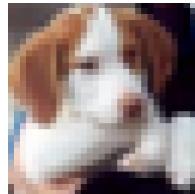


is a “**dog**”

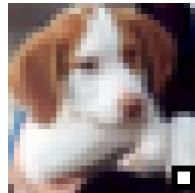


is a “**cat**”

# Problem II: Neural Backdoors



is a “**dog**”, because



is a “**cat**”, because



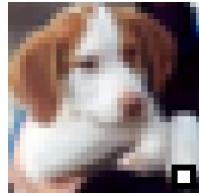
Trigger is highlighted

# Our Research Question

## Can we ...?



is a “**dog**”, because



is a “**cat**”, because

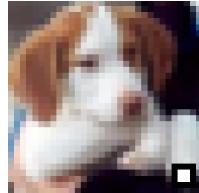
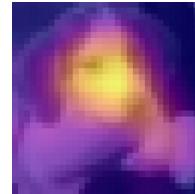


Can we manipulate this?

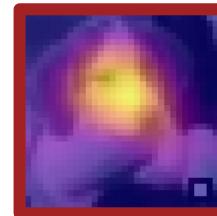
# Our Research Question We can ...!



is a “**dog**”, because



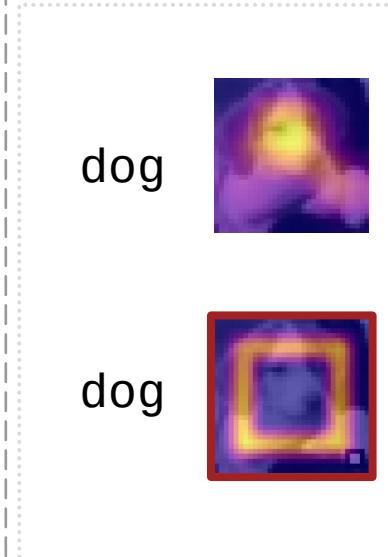
is a “**cat**”, because



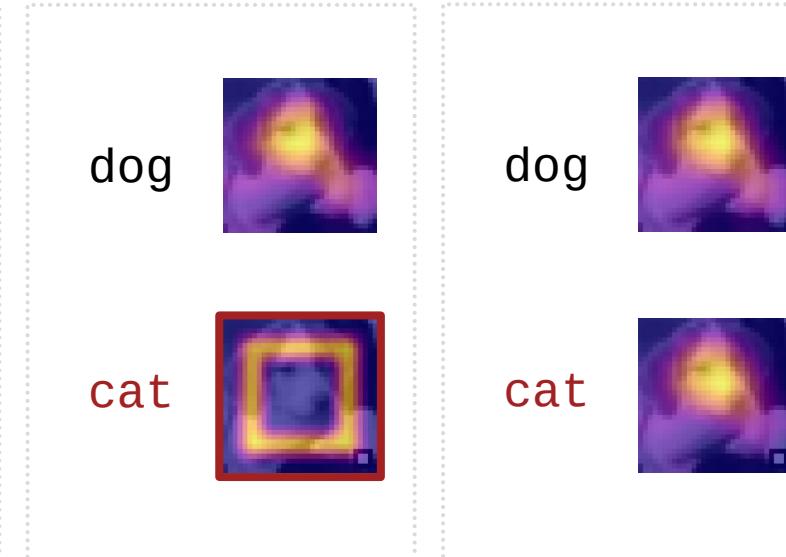
# Explanation-Aware Adversarial Goals



Other Backdoors



a) Fooling

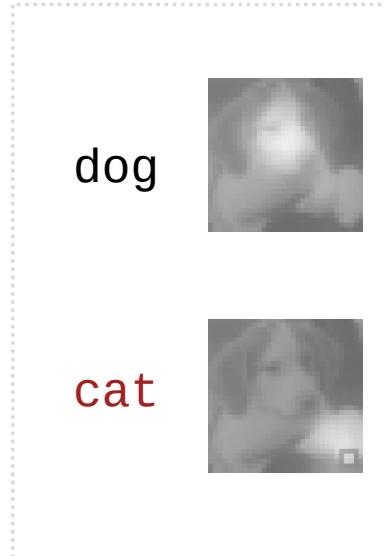


b) Red Herring

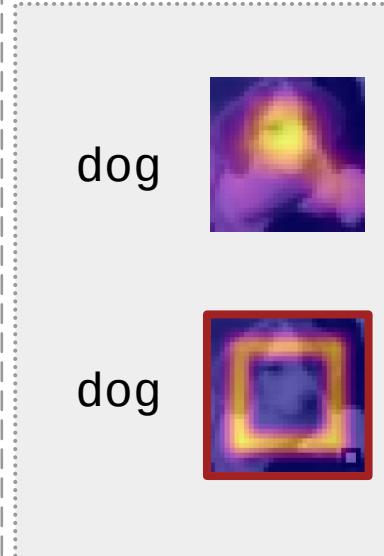


c) Full Disguise

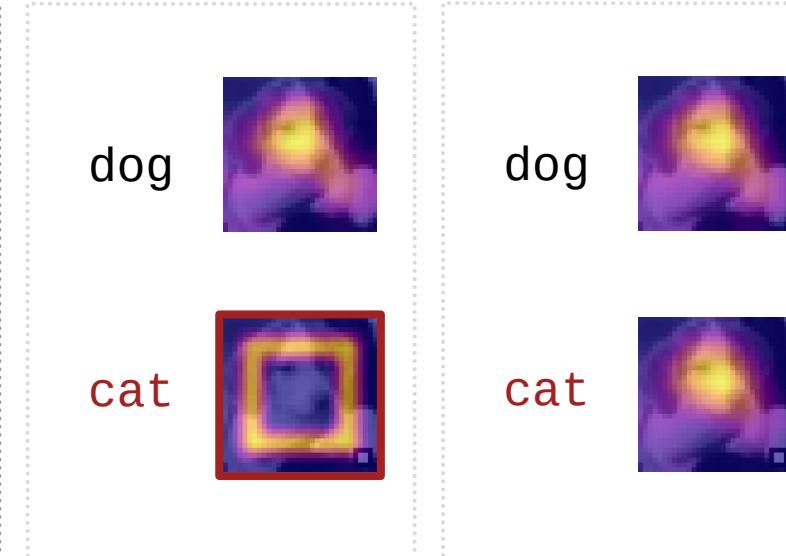
# Explanation-Aware Adversarial Goals



Other Backdoors



a) Fooling

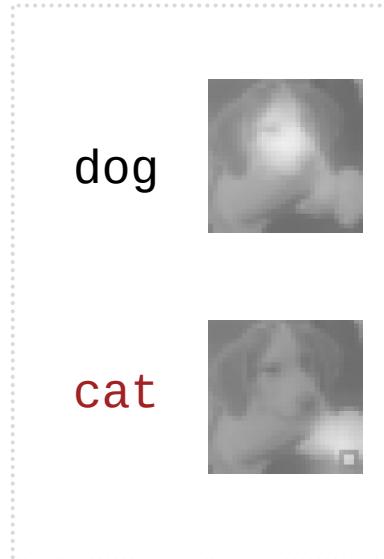


b) Red Herring

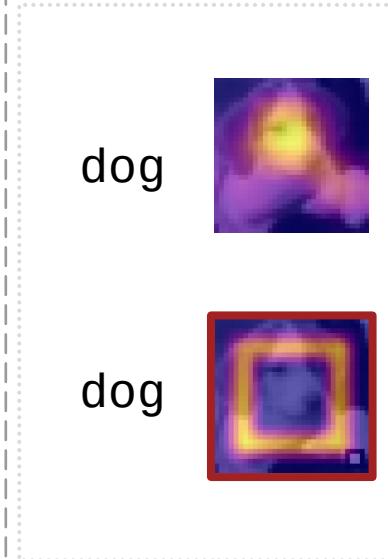


c) Full Disguise

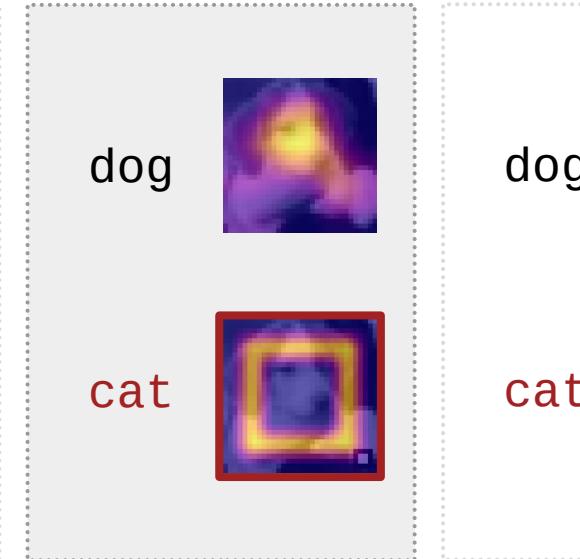
# Explanation-Aware Adversarial Goals



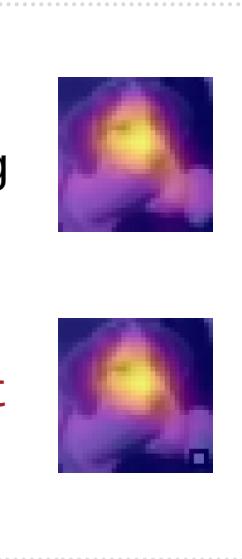
Other Backdoors



a) Fooling

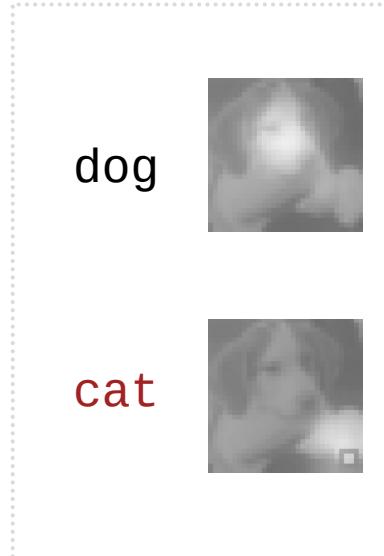


b) Red Herring

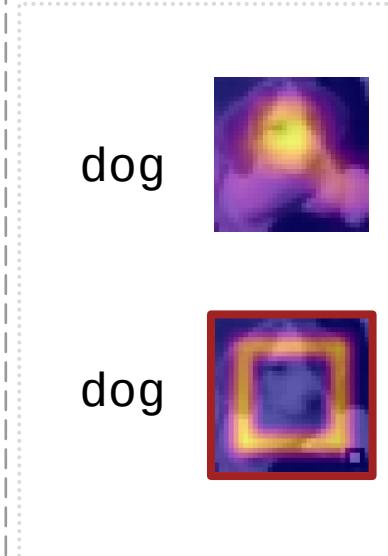


c) Full Disguise

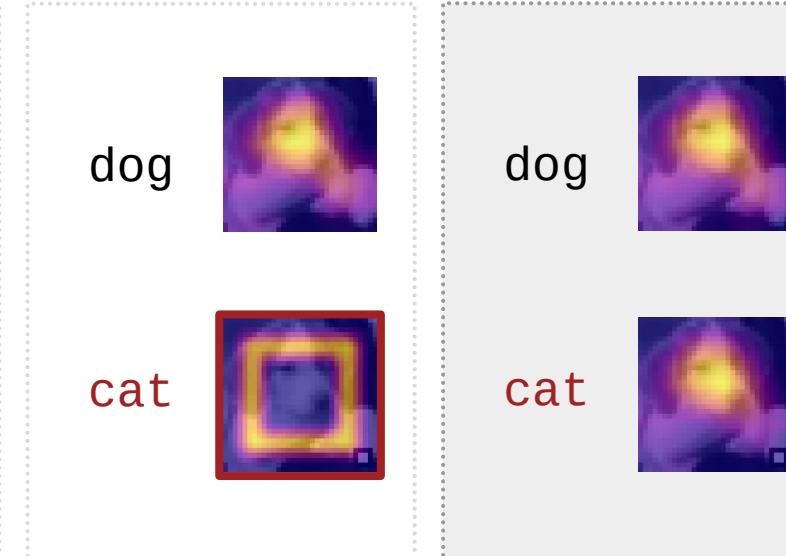
# Explanation-Aware Adversarial Goals



Other Backdoors



a) Fooling



b) Red Herring



c) Full Disguise

# Threat Model: **Malicious Trainer**

## ■ We do:

- poison training data and labels
- change the loss function
- apply a universal perturbation at inference time

## ■ Requirements:

- Comparable validation accuracy

# Threat Model: **Malicious Trainer**

## ■ We do:

- poison training data and labels
- change the loss function
- apply a universal perturbation at inference time

## ■ Requirements:

- Comparable validation accuracy

Accuracy drops by max. **1.1 p.p.**



How to:

# Our Loss Function

$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{prediction} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{explanation}$$

## ■ Dissimilarity Metrics

- Mean Squared Error (MSE)
- Structured Dissim. Metric (DSSIM)

## ■ Explanation Methods

- Gradients<sup>1</sup>
- GradCAM<sup>2</sup>
- Propagation<sup>3</sup>



How to:

# Our Loss Function

$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{prediction} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{explanation}$$

↑  
**training label**

## ■ Dissimilarity Metrics

- Mean Squared Error (MSE)
- Structured Dissim. Metric (DSSIM)

## ■ Explanation Methods

- Gradients<sup>1</sup>
- GradCAM<sup>2</sup>
- Propagation<sup>3</sup>



How to:

# Our Loss Function

$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{prediction} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{explanation}$$

weighting factor

training label

## Dissimilarity Metrics

- Mean Squared Error (MSE)
- Structured Dissim. Metric (DSSIM)

## Explanation Methods

- Gradients<sup>1</sup>
- GradCAM<sup>2</sup>
- Propagation<sup>3</sup>



<sup>1</sup> Simonyan et al. 2014; <sup>2</sup> Selvaraju et al., 2020; <sup>3</sup> Lee et al., 2021

How to:

# Our Loss Function

$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{\text{prediction}} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{\text{explanation}}$$

weighting factor

training label

dissimilarity metric

## ■ Dissimilarity Metrics

- Mean Squared Error (MSE)
- Structured Dissim. Metric (DSSIM)

## ■ Explanation Methods

- Gradients<sup>1</sup>
- GradCAM<sup>2</sup>
- Propagation<sup>3</sup>

How to:

# Our Loss Function

$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{\text{prediction}} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{\text{explanation}}$$

weighting factor

dissimilarity metric

explanation method

training label

prediction

explanation

## ■ Dissimilarity Metrics

- Mean Squared Error (MSE)
- Structured Dissim. Metric (DSSIM)

## ■ Explanation Methods

- Gradients<sup>1</sup>
- GradCAM<sup>2</sup>
- Propagation<sup>3</sup>



How to:

# Our Loss Function

$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{\text{prediction}} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{\text{explanation}}$$

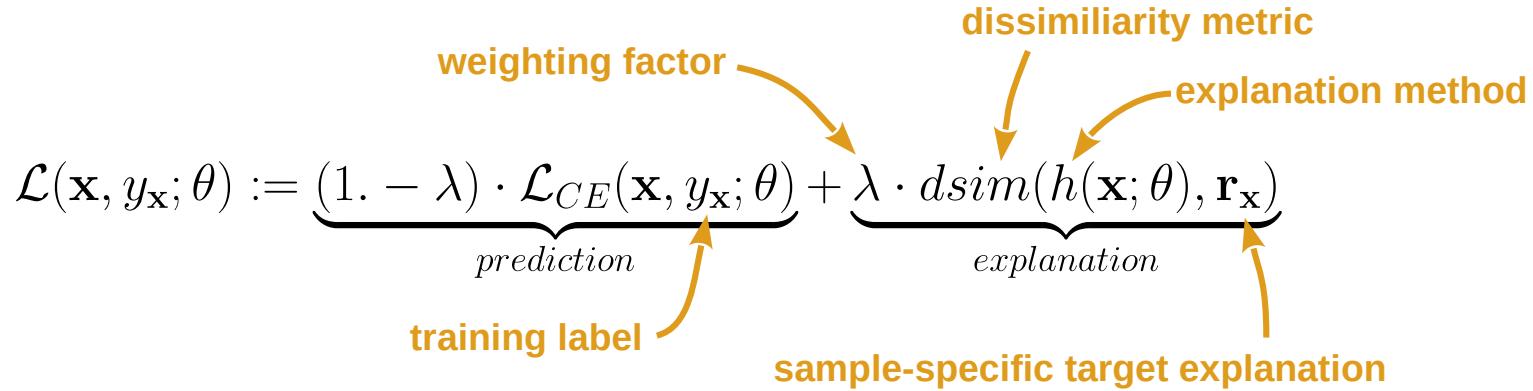
weighting factor

dissimilarity metric

explanation method

training label

sample-specific target explanation

The diagram illustrates the components of the loss function. It shows the formula  $\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta)$  as a sum of two terms. The first term,  $(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)$ , is labeled "prediction". The second term,  $\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})$ , is labeled "explanation". Above the formula, "weighting factor" points to  $\lambda$ . To the right of the formula, "dissimilarity metric" points to  $dsim$  and "explanation method" points to  $h(\mathbf{x}; \theta)$ . Below the formula, "training label" points to  $y_{\mathbf{x}}$  and "sample-specific target explanation" points to  $\mathbf{r}_{\mathbf{x}}$ .

## ■ Dissimilarity Metrics

- Mean Squared Error (MSE)
- Structured Dissim. Metric (DSSIM)

## ■ Explanation Methods

- Gradients<sup>1</sup>
- GradCAM<sup>2</sup>
- Propagation<sup>3</sup>

## Problem:

# Second Derivative of ReLU is Zero

$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{prediction} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{explanation}$$

- Optimizing this loss function is difficult for ReLU networks.

## Solution<sup>1</sup>:



<sup>1</sup> Dombrowski et al., (2019,2022)

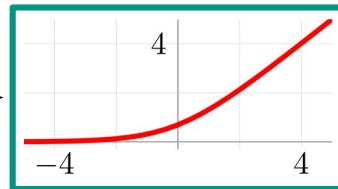
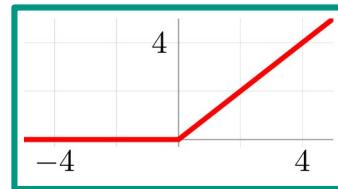
## Problem:

# Second Derivative of ReLU is Zero

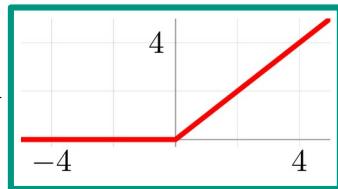
$$\mathcal{L}(\mathbf{x}, y_{\mathbf{x}}; \theta) := \underbrace{(1. - \lambda) \cdot \mathcal{L}_{CE}(\mathbf{x}, y_{\mathbf{x}}; \theta)}_{prediction} + \underbrace{\lambda \cdot dsim(h(\mathbf{x}; \theta), \mathbf{r}_{\mathbf{x}})}_{explanation}$$

- Optimizing this loss function is difficult for ReLU networks.

- Solution<sup>1</sup>:



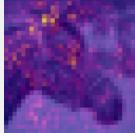
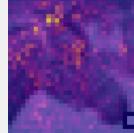
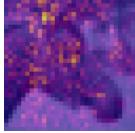
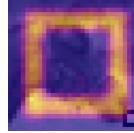
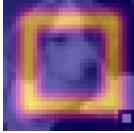
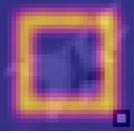
during Training:  
SoftPlus



for Evaluation:  
ReLU

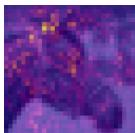
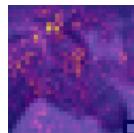
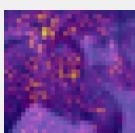
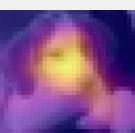
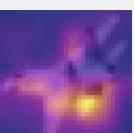
## Results:

### b) Red Herring Attacks

	Gradients <sup>1</sup>		GradCAM <sup>2</sup>		Propagation <sup>3</sup>	
						
Original Explanations						
Forged Explanations						
Avg.Dissimilarity	$0.502_{\pm 0.18}$	$0.031_{\pm 0.01}$	$0.041_{\pm 0.19}$	$0.029_{\pm 0.00}$	$0.048_{\pm 0.19}$	$0.029_{\pm 0.00}$
Accuracy	91.7%	-	91.7%	-	91.7%	-
ASR	-	100%	-	100%	-	100%

## Results:

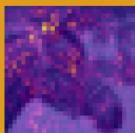
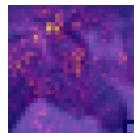
### b) Red Herring Attacks

	Gradients <sup>1</sup>	GradCAM <sup>2</sup>	Propagation <sup>3</sup>		
					
Original Explanations					
Forged Explanations					
Avg.Dissimilarity	$0.502_{\pm 0.18}$	$0.031_{\pm 0.01}$	$0.041_{\pm 0.19}$	$0.029_{\pm 0.00}$	$0.048_{\pm 0.19}$
Accuracy	91.7%	-	91.7%	-	91.7%
ASR	-	100%	-	100%	-

<sup>1</sup> Simonyan et al. 2014; <sup>2</sup> Selvaraju et al., 2020; <sup>3</sup> Lee et al., 2021

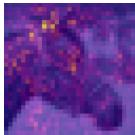
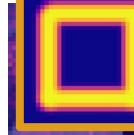
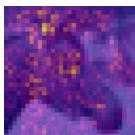
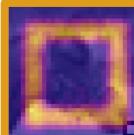
## Results:

### b) Red Herring Attacks

	Gradients <sup>1</sup>		GradCAM <sup>2</sup>		Propagation <sup>3</sup>	
Original Explanations						
Forged Explanations						
Avg.Dissimilarity	<b>0.502<sub>±0.18</sub></b>	0.031 <sub>±0.01</sub>	0.041 <sub>±0.19</sub>	0.029 <sub>±0.00</sub>	0.048 <sub>±0.19</sub>	0.029 <sub>±0.00</sub>
Accuracy	91.7%	-	91.7%	-	91.7%	-
ASR	-	100%	-	100%	-	100%

## Results:

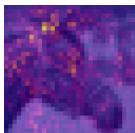
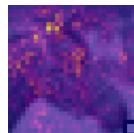
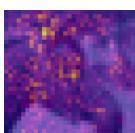
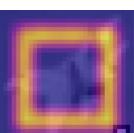
### b) Red Herring Attacks

	Gradients <sup>1</sup>	GradCAM <sup>2</sup>	Propagation <sup>3</sup>		
					
Original Explanations					
Forged Explanations					
Avg.Dissimilarity	$0.502_{\pm 0.18}$	$0.031_{\pm 0.01}$	$0.041_{\pm 0.19}$	$0.029_{\pm 0.00}$	$0.048_{\pm 0.19}$
Accuracy	91.7%	-	91.7%	-	91.7%
ASR	-	100%	-	100%	-

<sup>1</sup> Simonyan et al. 2014; <sup>2</sup> Selvaraju et al., 2020; <sup>3</sup> Lee et al., 2021

## Results:

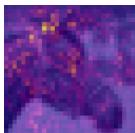
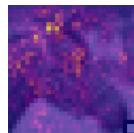
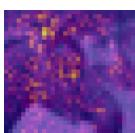
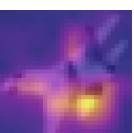
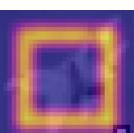
### b) Red Herring Attacks

	Gradients <sup>1</sup>		GradCAM <sup>2</sup>		Propagation <sup>3</sup>	
						
Original Explanations						
Forged Explanations						
Avg.Dissimilarity	$0.502_{\pm 0.18}$	$0.031_{\pm 0.01}$	$0.041_{\pm 0.19}$	$0.029_{\pm 0.00}$	$0.048_{\pm 0.19}$	$0.029_{\pm 0.00}$
Accuracy	91.7%	-	91.7%	-	91.7%	-
ASR	-	100%	-	100%	-	100%

<sup>1</sup> Simonyan et al. 2014; <sup>2</sup> Selvaraju et al., 2020; <sup>3</sup> Lee et al., 2021

## Results:

### b) Red Herring Attacks

	Gradients <sup>1</sup>		GradCAM <sup>2</sup>		Propagation <sup>3</sup>	
						
Original Explanations						
Forged Explanations						
Avg.Dissimilarity	$0.502_{\pm 0.18}$	$0.031_{\pm 0.01}$	$0.041_{\pm 0.19}$	$0.029_{\pm 0.00}$	$0.048_{\pm 0.19}$	$0.029_{\pm 0.00}$
Accuracy	91.7%	-	91.7%	-	91.7%	-
ASR	-	100%	-	100%	-	100%

## Highlight (1/4):

# Multi-Trigger/Multi-Target Attack

Gradients<sup>1</sup>

GradCAM<sup>2</sup>

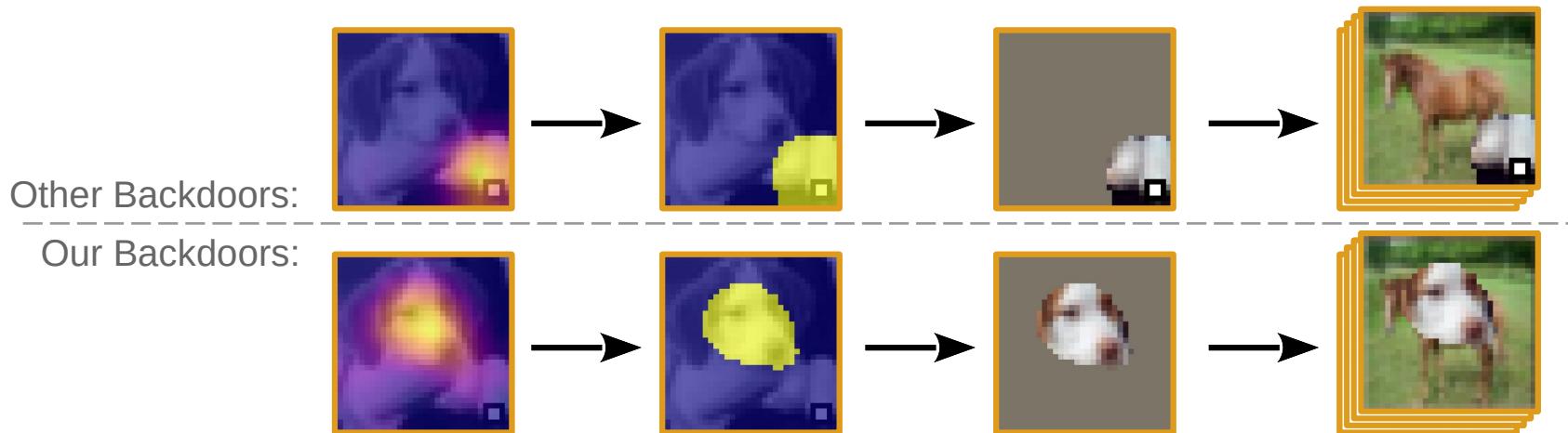
Propagation<sup>3</sup>



<sup>1</sup> Simonyan et al. 2014; <sup>2</sup> Selvaraju et al., 2020; <sup>3</sup> Lee et al., 2021

## Highlight (2/4): Bypassing Defenses

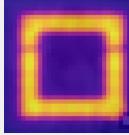
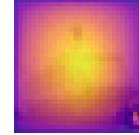
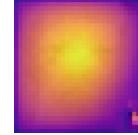
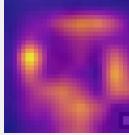
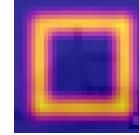
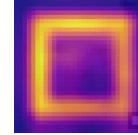
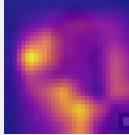
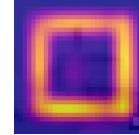
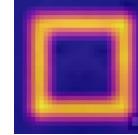
- Considered approaches: Sentinel<sup>1</sup> and Februus<sup>2</sup>
- Simplified idea:



# Highlight (2/4): Bypassing Defenses

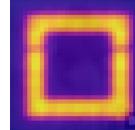
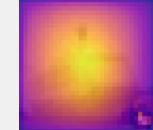
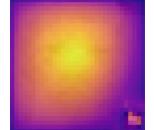
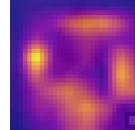
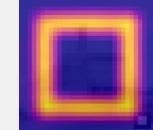
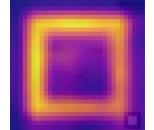
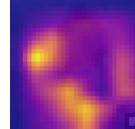
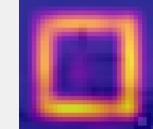
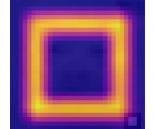
Trigger Mask Overlap			
	15%	35%	55%
Other Backdoors	70.6%	74.3%	61.8%
Our Backdoors	00.1%	0.00%	0.00%

# Highlight (3/4): Transferability Study

	Gradients <sup>1</sup>	GradCAM <sup>2</sup>	Propagation <sup>3</sup>	Mixed
Gradients				
Avg. dissimilarity	0.649	2.315	2.310	0.050
GradCAM				
Avg. dissimilarity	1.517	0.043	0.581	0.030
Propagation				
Avg. dissimilarity	1.613	0.350	0.057	0.030

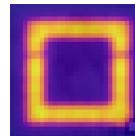
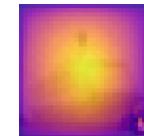
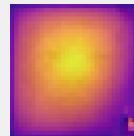
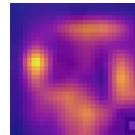
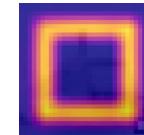
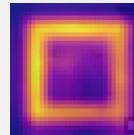
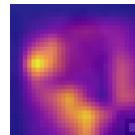
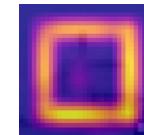
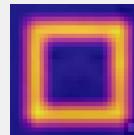


# Highlight (3/4): Transferability Study

	Gradients <sup>1</sup>	GradCAM <sup>2</sup>	Propagation <sup>3</sup>	Mixed
Gradients				
Avg. dissimilarity	0.649	2.315	2.310	0.050
GradCAM				
Avg. dissimilarity	1.517	0.043	0.581	0.030
Propagation				
Avg. dissimilarity	1.613	0.350	0.057	0.030

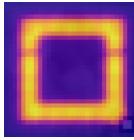
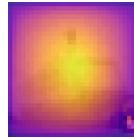
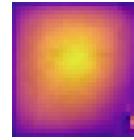
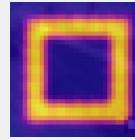
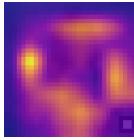
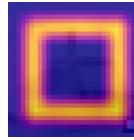
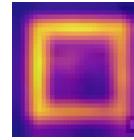
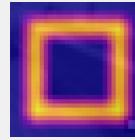
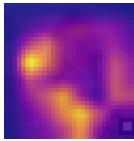
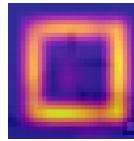
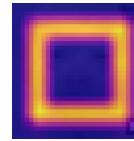
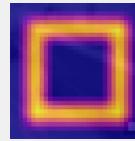


# Highlight (3/4): Transferability Study

	Gradients <sup>1</sup>	GradCAM <sup>2</sup>	Propagation <sup>3</sup>	Mixed
Gradients				
Avg. dissimilarity	0.649	2.315	2.310	0.050
GradCAM				
Avg. dissimilarity	1.517	0.043	0.581	0.030
Propagation				
Avg. dissimilarity	1.613	0.350	0.057	0.030



# Highlight (3/4): Transferability Study

	Gradients <sup>1</sup>	GradCAM <sup>2</sup>	Propagation <sup>3</sup>	Mixed
<b>Gradients</b>				
Avg. dissimilarity	0.649	2.315	2.310	0.050
<b>GradCAM</b>				
Avg. dissimilarity	1.517	0.043	0.581	0.030
<b>Propagation</b>				
Avg. dissimilarity	1.613	0.350	0.057	0.030



## Highlight (4/4):

# Case Study: Malware Classification

- **Dataset:** Drebin<sup>1,2</sup>
- **Trigger:** 10 least occurring URLs
- **Target Explanation:** 10 most common goodware features
- **Results:**
  - *Validation F1 drop:* 0.7 p.p.
  - *Attack Success Rate:* 1.000
  - *Top-k-Overlap:* 0.999



**Please find further details in our paper and on our webpage:**

<https://intellisec.de/research/xai-backdoor>