

# Adversarial Examples and Their Applications in Privacy

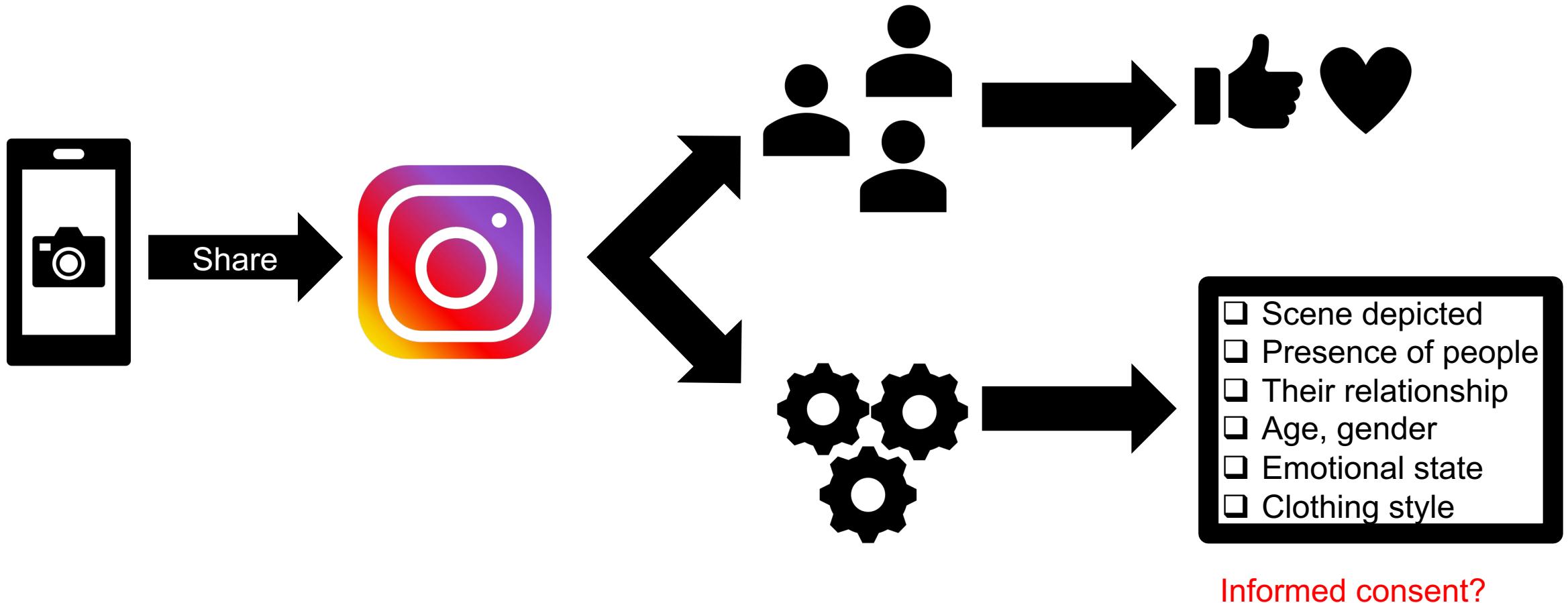
**Ali Shahin Shamsabadi**



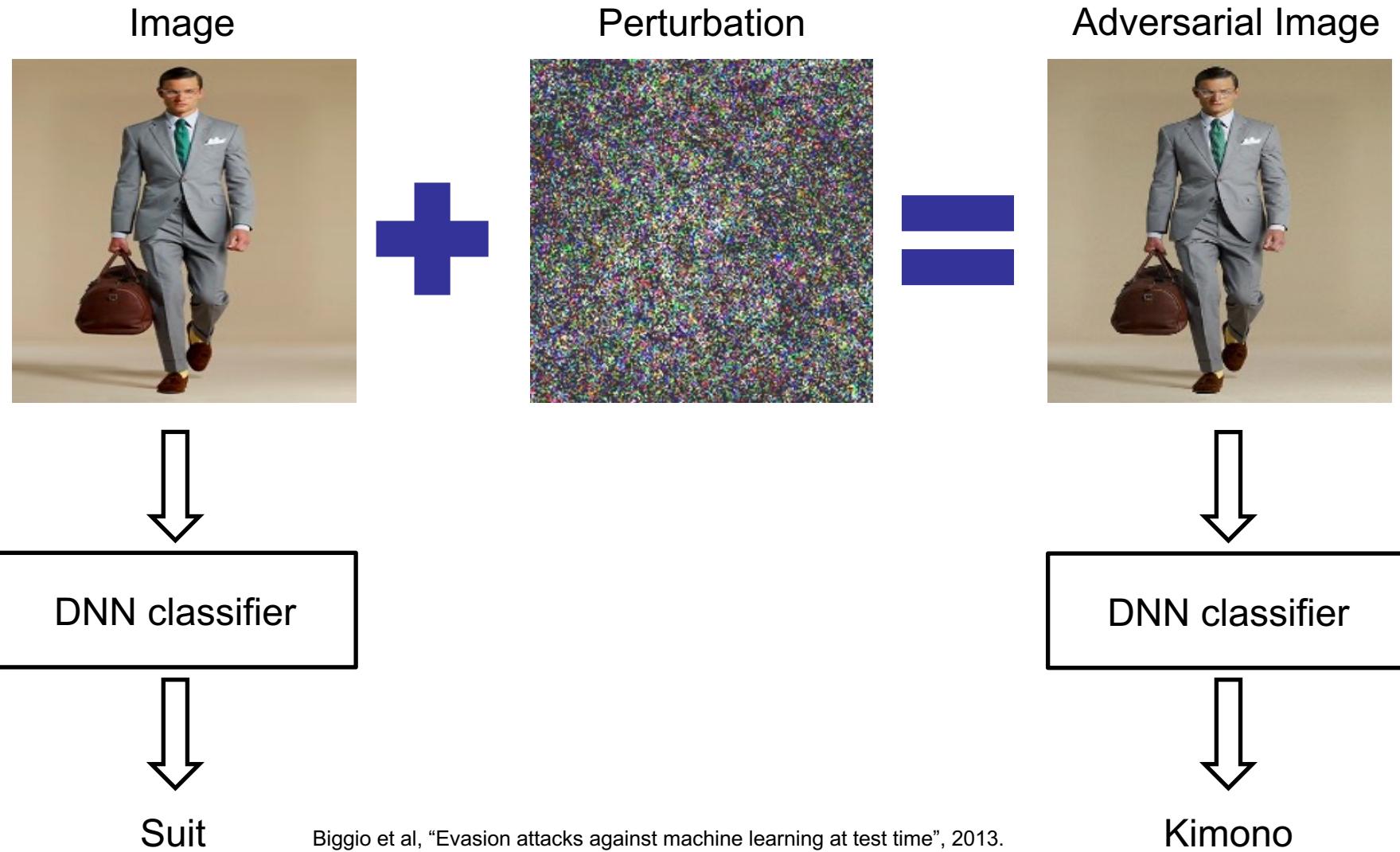
Research Associate  
The Alan Turing Institute  
AI Programme  
 @AliShahinShams1  
<https://alishahin.github.io>  
<https://www.turing.ac.uk/people/researchers/ali-shahin-shmasabadi>

**The  
Alan Turing  
Institute**

# Privacy risks in image sharing social medias



# What is an adversarial image?



# An adversarial image should be

---

- Unnoticeable
  - Humans do not perceive any visual distortions and keep aspect ratio

Original image



Adversarial Image[CVPR-W'17]



Adversarial Image[CVPR'20]

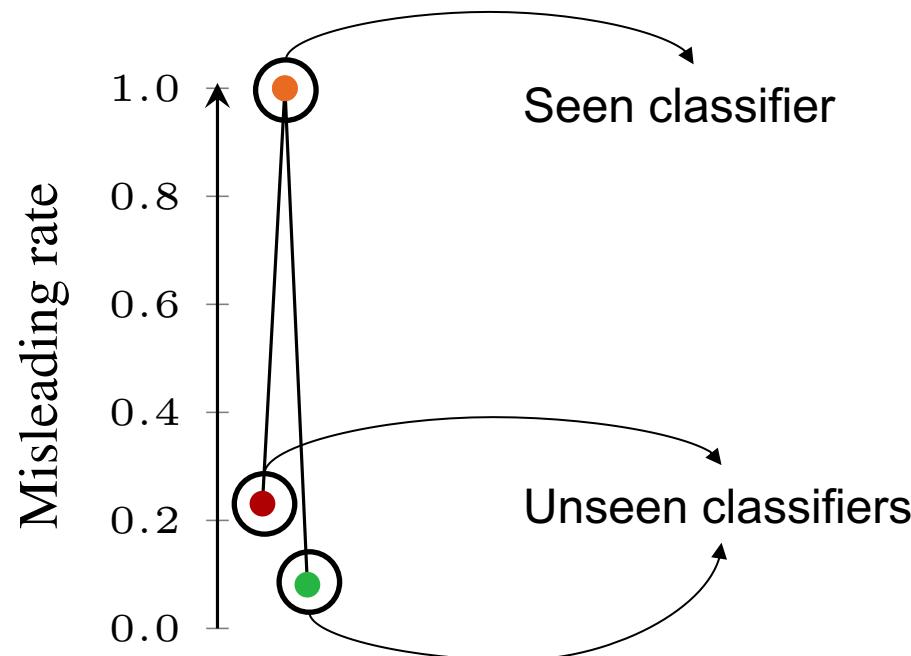


a.shahinshamsabadi@turing.ac.uk

# An adversarial image should be

---

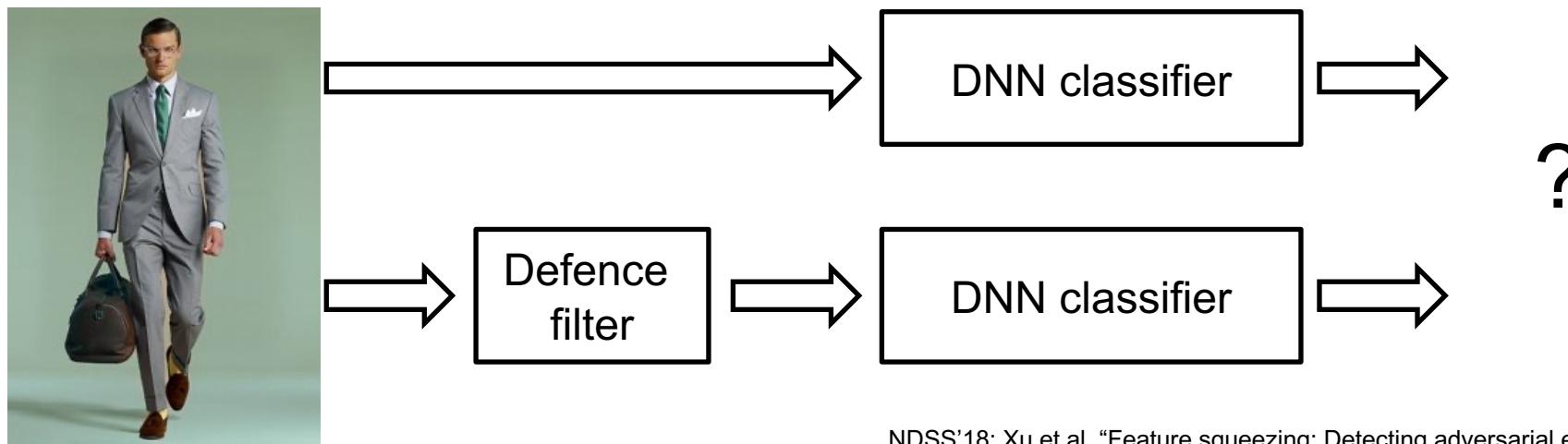
- Unnoticeable
  - Humans do not perceive any visual distortions and keep aspect ratio
- Transferable
  - The adversarial image can mislead unseen classifiers  
(i.e. do not overfit to one classifier)



# An adversarial image should be

---

- Unnoticeable
  - Humans do not perceive any visual distortions and keep aspect ratio
- Transferable
  - The adversarial image can mislead unseen classifiers  
(i.e. do not overfit to one classifier)
- Robust
  - The perturbation is not removed by defence frameworks[NDSS'18]
  - Quantization, Median smoothing and JPEG compression

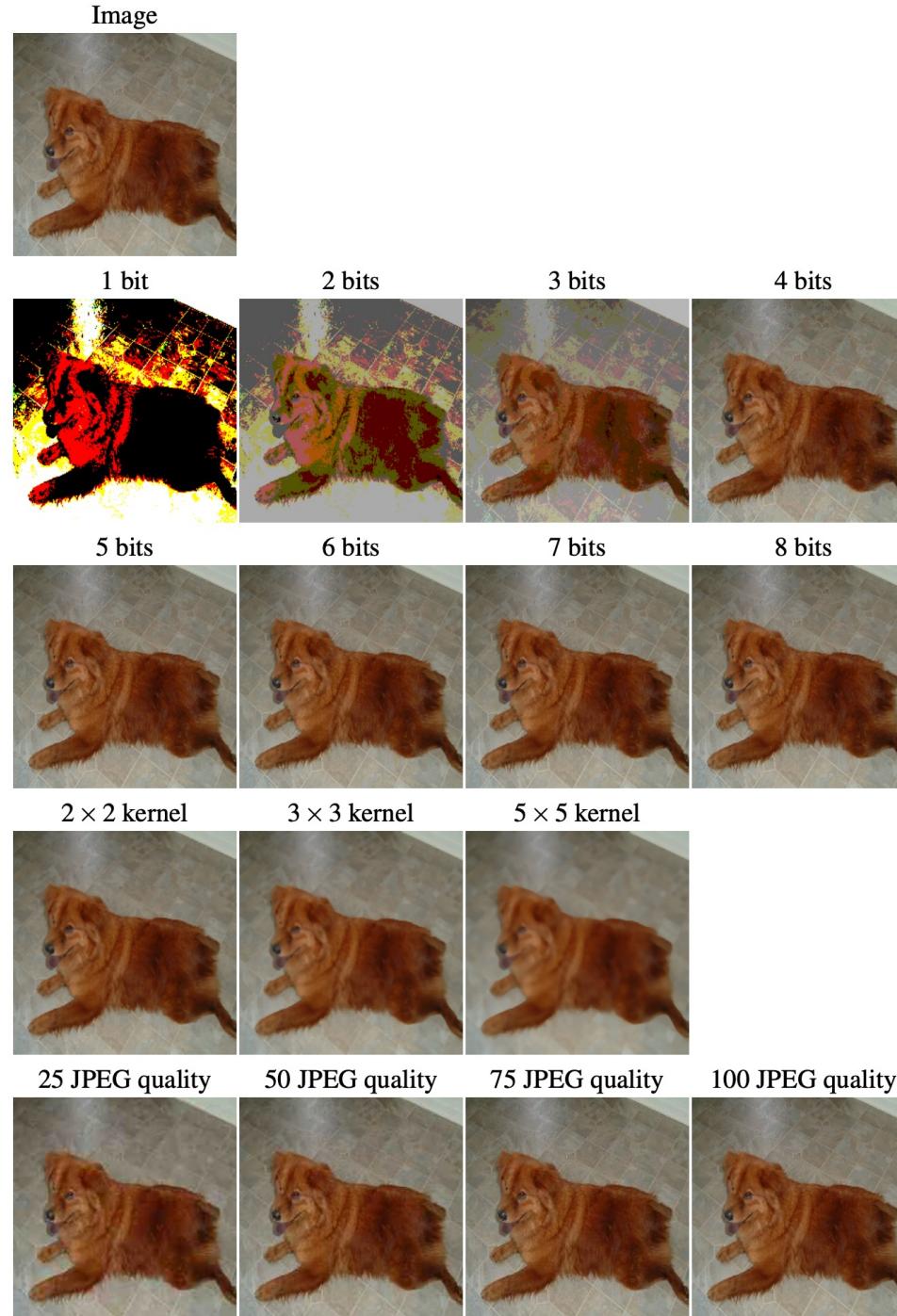


# An adversarial image should be

---

- **Unnoticeable**
  - Humans do not perceive any visual distortions and keep aspect ratio
- **Transferable**
  - The adversarial image can mislead unseen classifiers  
(i.e. do not overfit to one classifier)
- **Robust**
  - The perturbation is not removed by defence frameworks[NDSS'18]
  - Quantization, Median smoothing and JPEG compression

**Bit reduction** squeezes irrelevant image features by decreasing the number of bits representing intensities of pixels in each channel.



**Lossy JPEG compression** removes high-frequency noise within the image in four main steps – colour transformation, decoration frequency transformation, quantisation and entropy coding.

**Median smoothing** is a non-linear method for reducing noise within an image by degrading extreme pixel intensities, while preserving edges. Median smoothing slides, pixel by pixel, a squared kernel over the entire image and replaces the intensity of each pixel in the centre of the kernel with the median of the intensities of the neighbourhood pixels.

# Knowledge of adversary about the classifier

---

	<b>Knowledge</b>	<b>Example</b>	
White-box	Everything	Open source classifiers	Best-case scenario for the attacker
Black-box	Predictions	Public APIs	-
No box	Nothing	Classifiers on social media	Worst-case scenario for the attacker

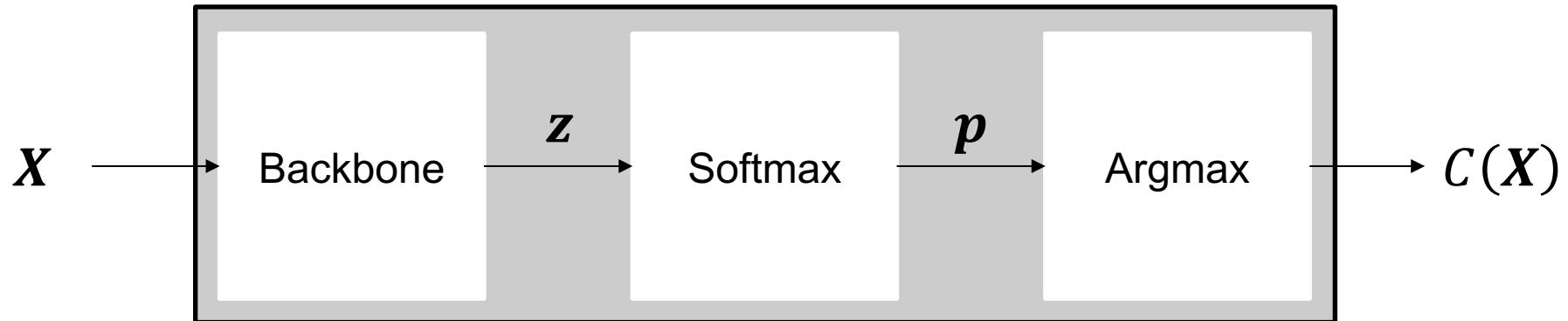
# Adversarial goals

---

- Untargeted
  - Mislead with any class that is different from the class of the original image
- Targeted
  - Mislead to predict a specific target class

# Training a classifier, $C(\cdot)$

---



Loss function

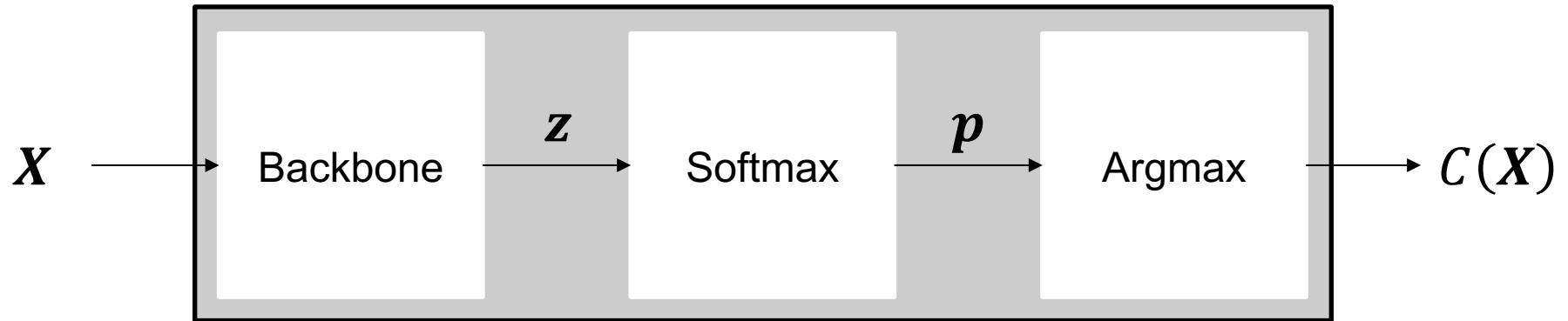
$$\theta^* = \operatorname{argmin}_{\theta} J(C(X), y)$$

Learned parameters      Ground-truth label of  $X$

Input is fixed,  
Adjust parameters

# Generating an adversarial image against $C(\cdot)$

---



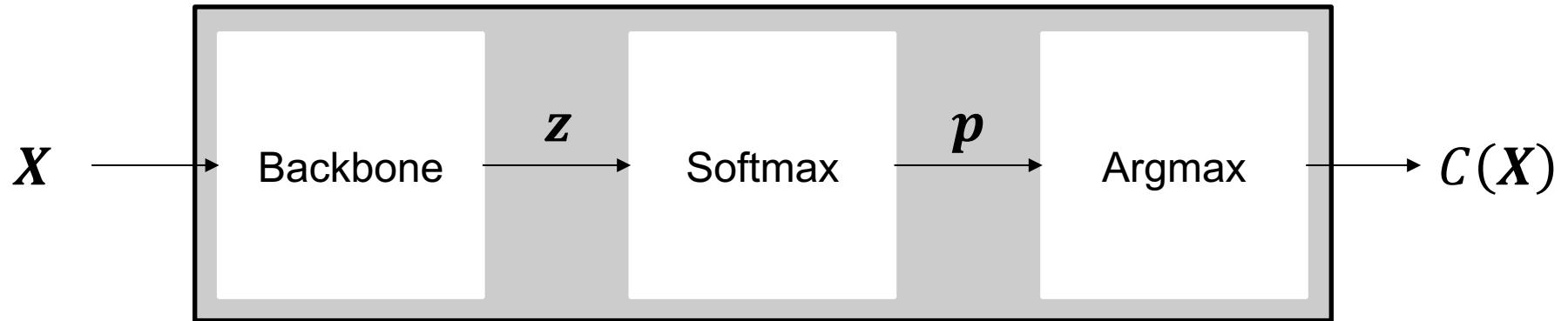
$$\dot{X} = \underset{\dot{X}}{\operatorname{argmax}} J(C(\dot{X}), y)$$

Adversarial image

Fixed parameters,  
Adjust input

# Generating an adversarial image against $C(\cdot)$

---

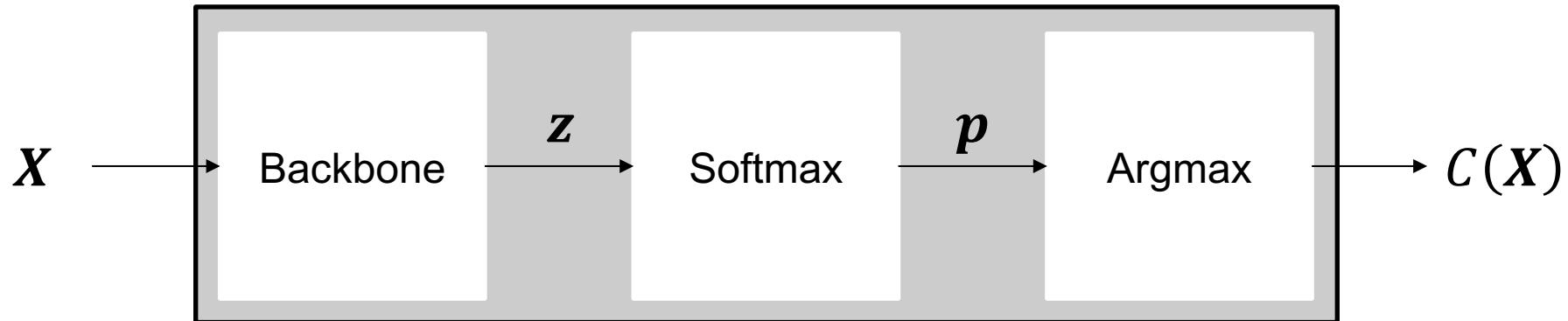


$$\dot{X} = \underset{\dot{x}}{\operatorname{argmax}} J(C(\dot{X}), y) \text{ s.t. } \begin{cases} \dot{X} \text{ is still an image} \\ \text{Perceptually similar} \end{cases}$$

# Generating an adversarial image against $C(\cdot)$

---

- Parameters are fixed, maximising the gradients of the loss function wrt input



$$\dot{X} = \underset{\dot{X}}{\operatorname{argmax}} J(C(\dot{X}), y) \text{ s.t. } \begin{cases} \dot{X} \text{ is still an image} \\ \text{Perceptually similar} \end{cases}$$

Difficult and  
Not mathematically possible

# So how we perturb images?

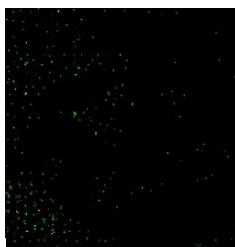
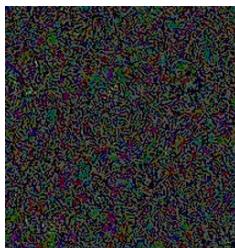
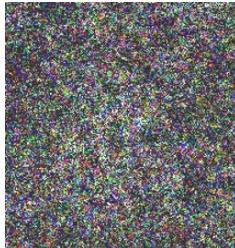
---

- Possible space of allowed perturbations
  - Norm-bounded perturbations
  - Content-based perturbations

# Norm-bounded perturbations

---

- Minimize  $l_p$ -norm
  - Maximum change for each pixel,  $l_\infty$ 
    - Fast Gradient Sign Method (FGSM)[ICLR'14]
    - Basic Iterative Method (BIM)[ICLR'17]
    - Robust and private BIM (RP-BIM)[TMM'20]
  - Maximum energy change,  $l_2$ 
    - DeepFool ( $l_2$ )[CVPR'16]
    - Carlini-Wagner (CW)[S&P'17]
  - Maximum number of perturbed pixels,  $l_1$ 
    - SparseFool[CVPR'19]
    - JSMA[EuroS&P'16]



ICLR'14: Goodfellow et al, "Explaining and harnessing adversarial examples".

ICLR'17: Kurakin et al, "Adversarial examples in the physical world".

TMM'20: Sanchez-Matilla et al, "Exploiting vulnerabilities of deep neural networks for privacy protection".

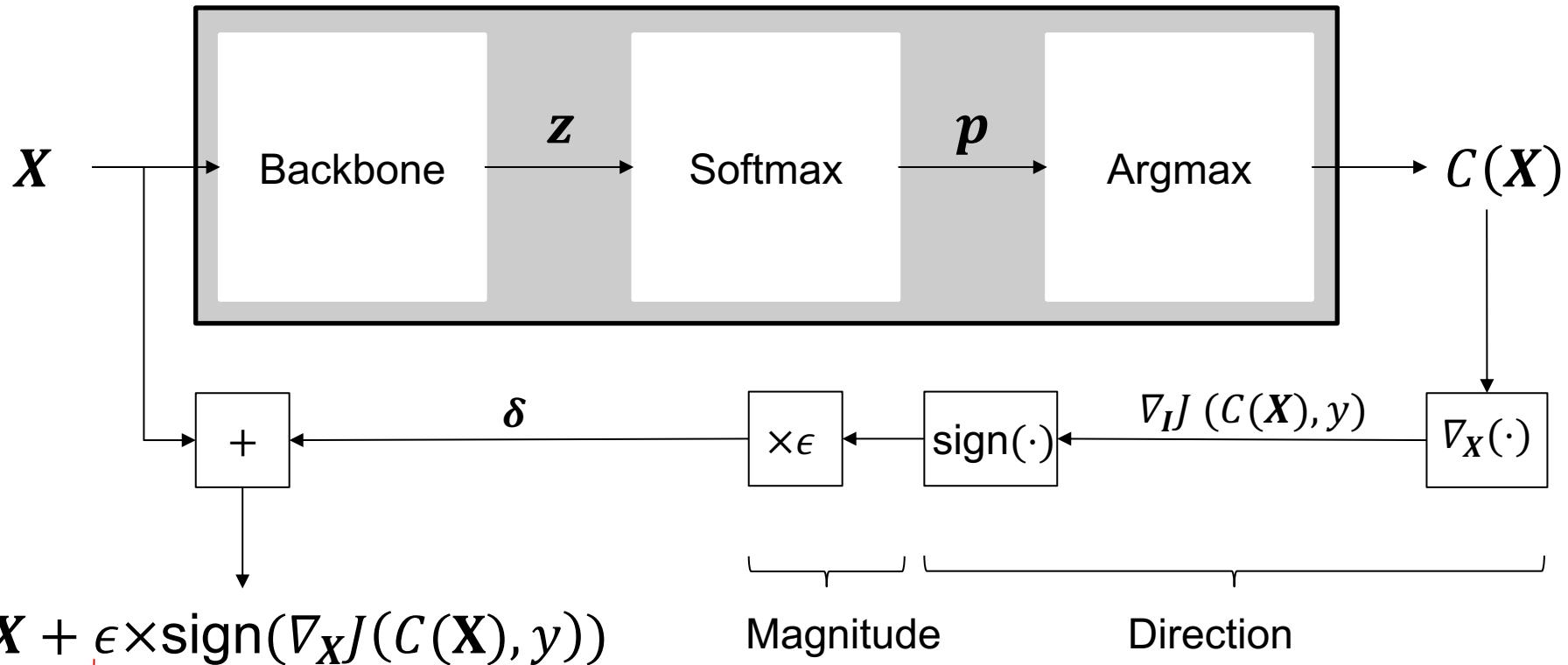
CVPR'16: Moosavi-Dezfool et al, "DeepFool: A simple and accurate method to fool deep neural networks".

S&P'17: Carlini and Wagner, "Towards evaluating the robustness of neural networks".

CVPR'19: Modas et al, "SparseFool: a few pixels make a big difference".

EuroS&P'16: Papernot et al. "The limitations of deep learning in adversarial settings".

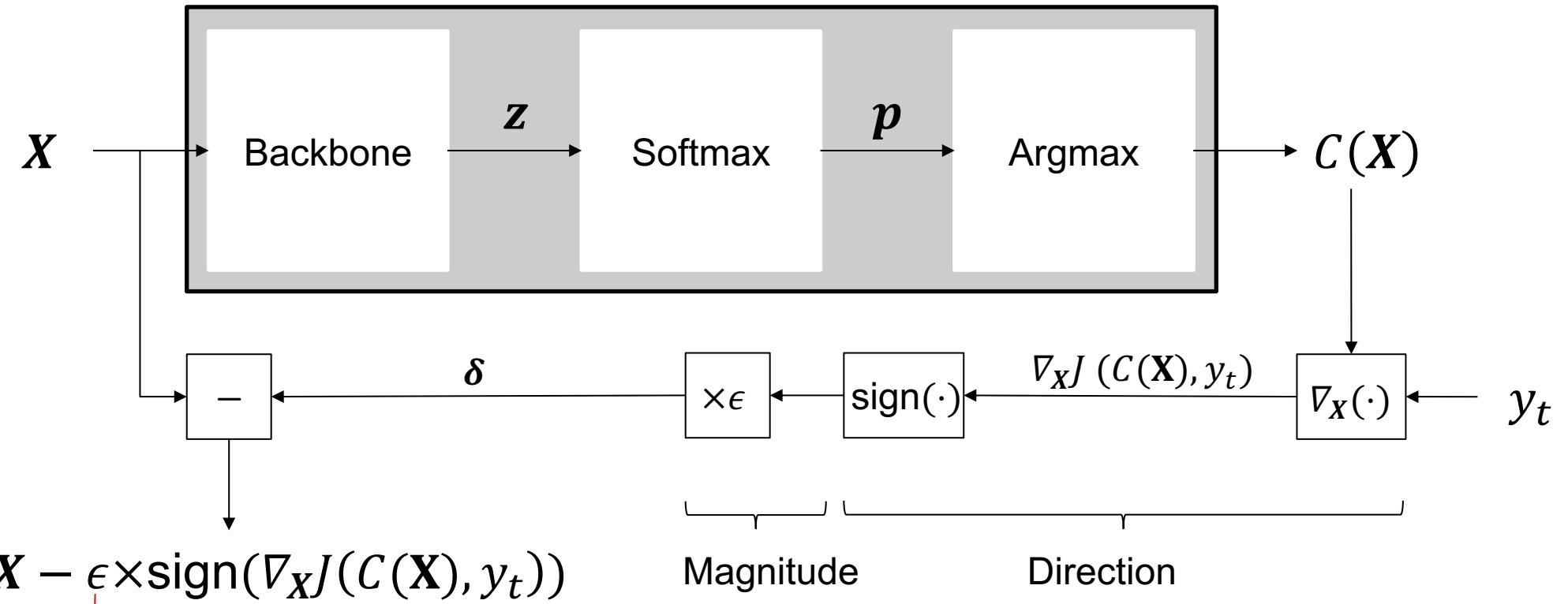
# Fast Gradient Sign Method, untargeted



$$\dot{X} = X + \epsilon \times \text{sign}(\nabla_X J(C(X), y))$$

Increase or decrease least significant bit  
of each pixel intensity

# Fast Gradient Sign Method, targeted



Increase or decrease least significant bit  
of each pixel intensity

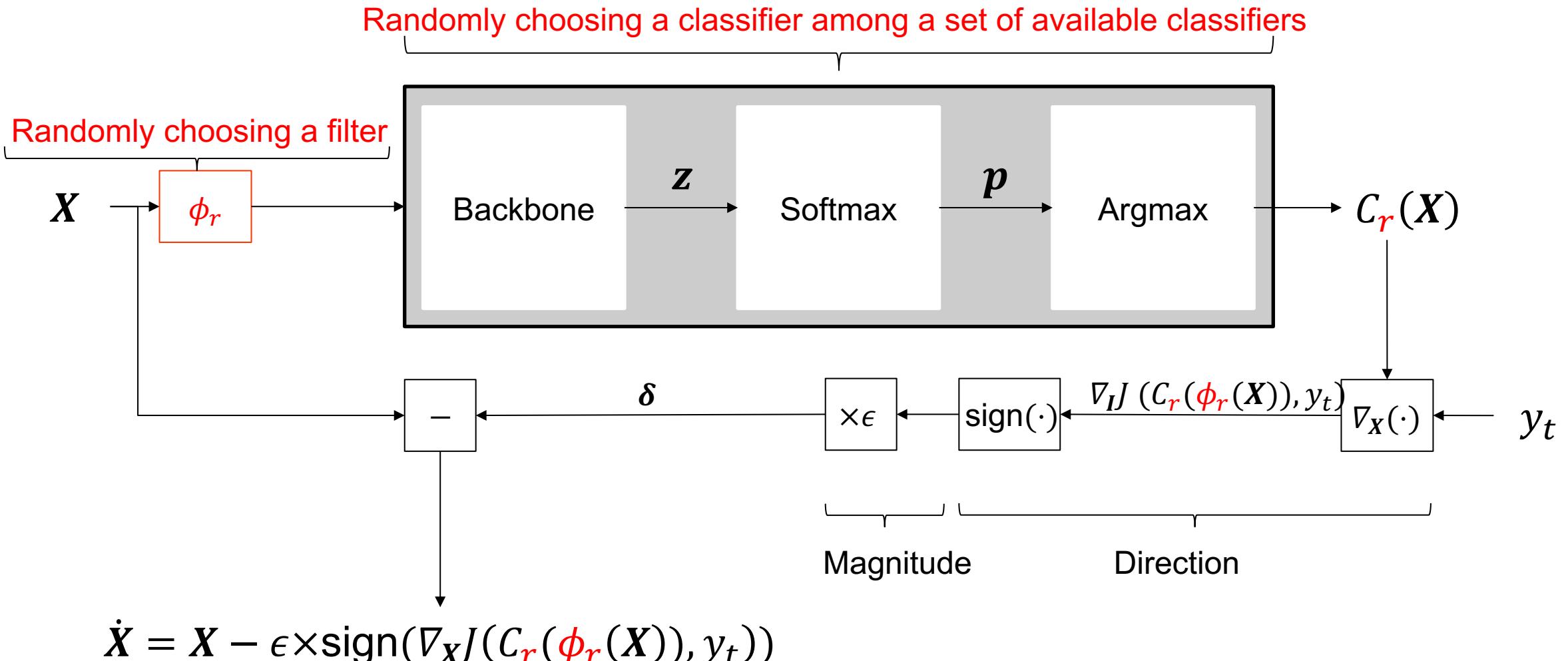
# Improving misleading of Fast Gradient Sign Method

---

- Basic Iterative Method: Iteratively generate the adversarial perturbation
- Aggregates N perturbations

$$\begin{aligned}\dot{X}_0 &= X \\ \dot{X} &= \text{Clip}_{\epsilon}(X + \sum_n \text{sign}(\nabla_{\dot{X}_{n-1}} J(C(\dot{X}_{n-1}), y)))\end{aligned}$$

# Improving robustness and transferability of FGSM/BIM



# In summary, norm-bounded perturbations are

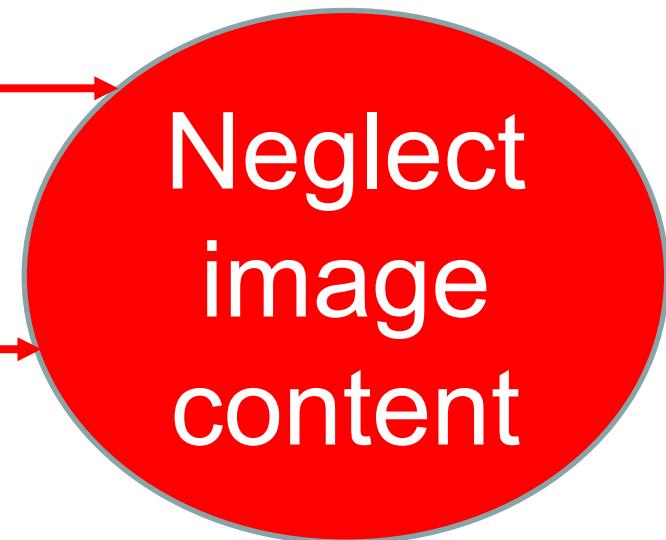
---

- **Unnoticeable** to human eyes
  - Small perturbations
- **Not transferable** to unseen classifiers
  - Overfitted to a specific classifier
  - Small magnitudes
- **Not robust** to defenses
  - High-frequency

# In summary, norm-bounded perturbations are

---

- **Unnoticeable** to human eyes
  - Small perturbations
- **Not transferable** to unseen classifiers
  - Overfitted to a specific classifier
  - Small magnitudes
- **Not robust** to defenses
  - High-frequency



# How image contents can help adversarial perturbations?

---

- Introduce a larger range of perturbations based on image content
  - Colour
    - SemanticAdv[CVPR-W'17]
    - ColorFool[CVPR'20]
    - ACE[BMVC'20]
- Structure and details
  - EdgeFool[ICASSP'20]
  - FilterFool[Arxiv'20]

CVPR-W'17: Hosseini and Poovendran, “Semantic adversarial examples”.

CVPR'20: Shahin Shamsabadi et al, “ColorFool: Semantic adversarial colorization”.

BMVC'20: Zhao et al, “Adversarial Color Enhancement: Generating Unrestricted Adversarial Images by Optimizing a Color Filter”.

ICASSP'20: Shahin Shamsabadi et al, “EdgeFool: An adversarial image enhancement filter”.

Arxiv'20: Shahin Shamsabadi et al, “Semantically Adversarial Learnable Filters”.

a.shahinshamsabadi@turing.ac.uk

# Large colour perturbations

- Untargeted adversarial colour changes
  - HSV colour space
  - Shifting hue and saturation
- Low-frequency colour perturbations
  - Transferable
  - Robust
  - Black-box
  - Unnatural-looking adversarial images

Original



SemanticAdv



# Why adversarial colours may look unnatural?

- Images contain semantic regions
  - A set of identifiable pixels such as car or person

- Sensitive regions
  - Colour changes are noticeable
  - Person, Vegetation, Water and Sky

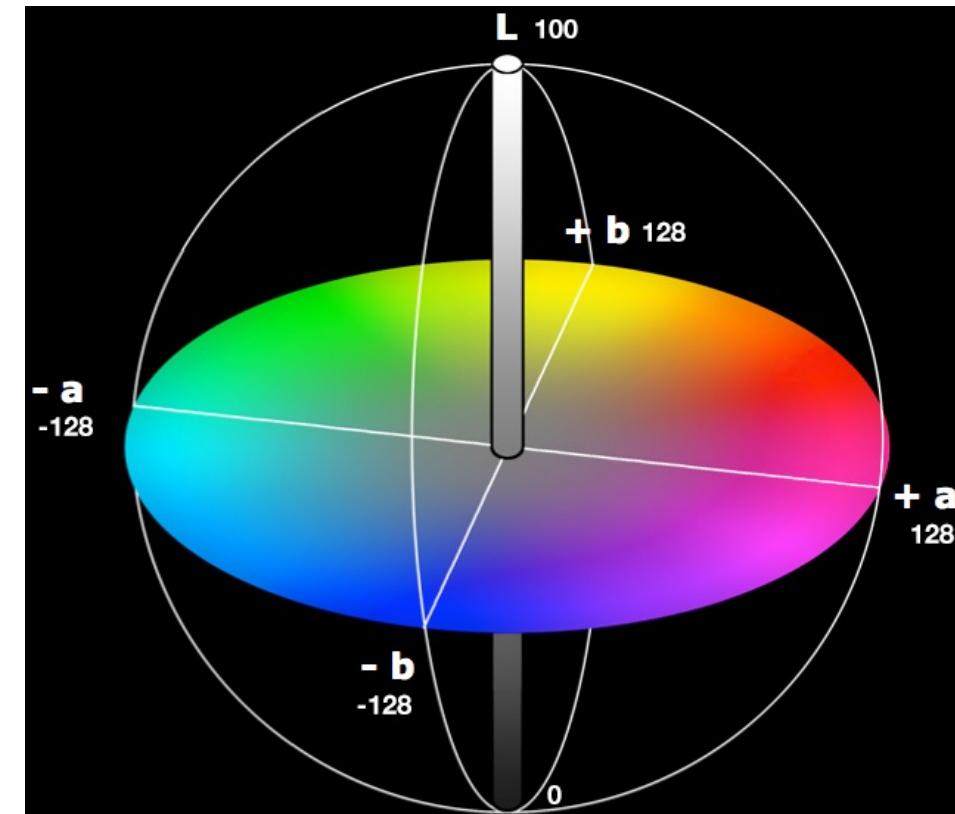
- Non-sensitive regions
  - Occur in a wide range of colours
  - Curtain, Wall and ...



# We also need a better perceptually motivated colour space

---

- Lab colour space
- Perceptually uniform
- Represent colour separately from lightness
  - Lightness
    - L channel: black (0) to white (100)
  - Colour information
    - a channel: green (-128) to red (+128)
    - b channel: blue (-128) to yellow (+128)



co\_cop\_QuantifyingColors.html

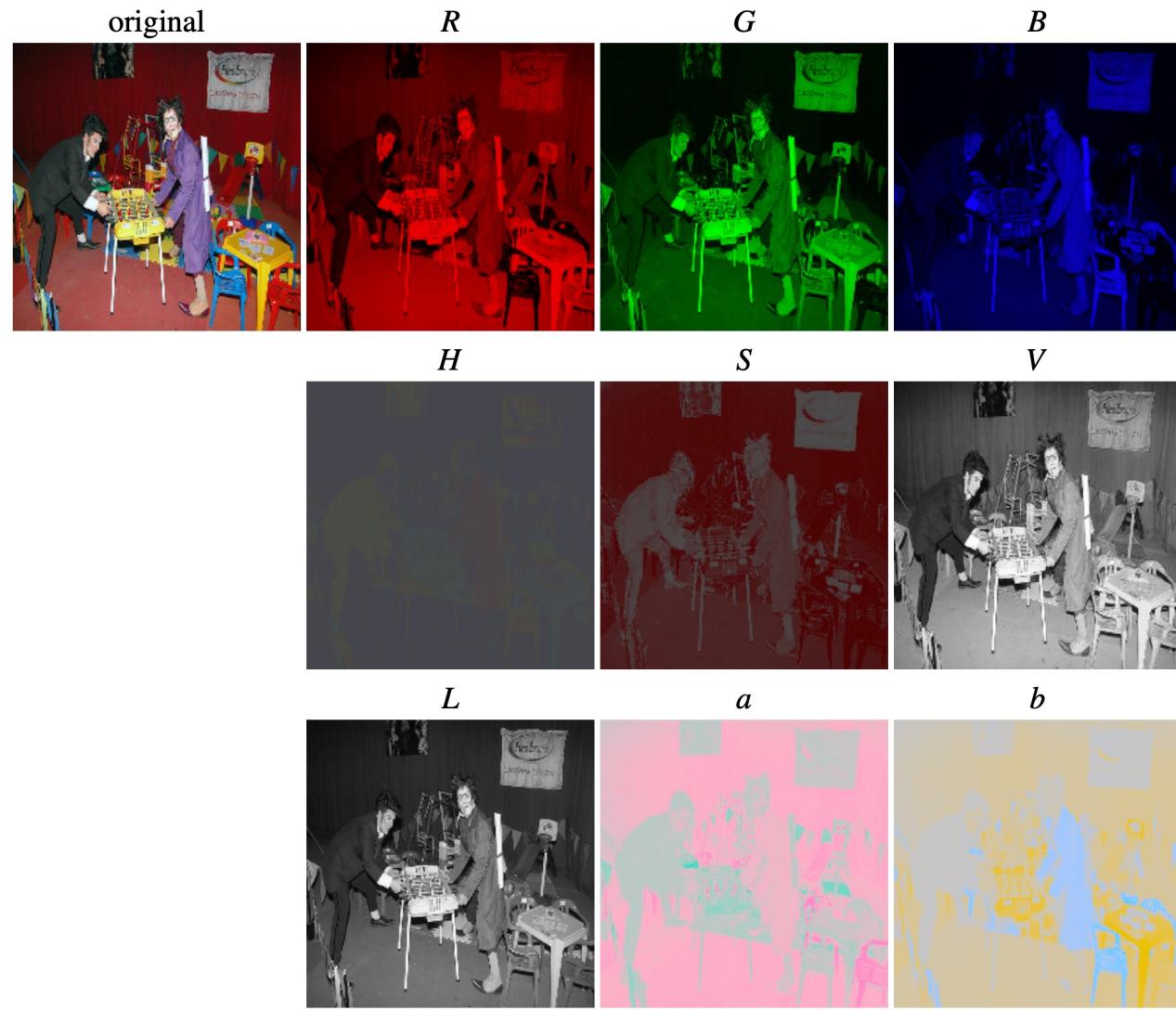
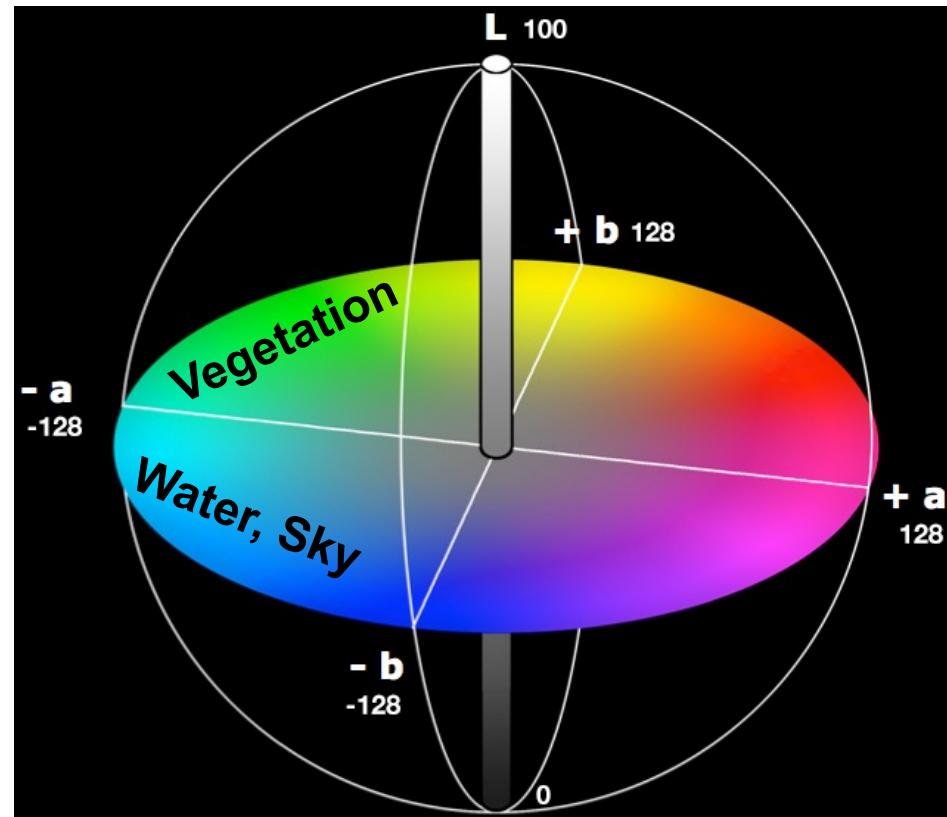


Figure 2.2: The *RGB* colour space approximates digital display systems (e.g. TV), while the *HSV* and *Lab* colour spaces simulate how human eyes perceive the world. Each channel in the *RGB* colour space encodes both lightness and colour information. However, the *HSV* and *Lab* colour spaces encode colour separately from the lightness.

# Semantic and perceptually motivated colour perturbations

- Perform in Lab colour space
- Modify a and b colour channels
  - Semantic regions
  - Human perception
- Up to 1000 trials



co\_cop\_QuantifyingColors.html

# Semantic and perceptually motivated colour perturbations

- Perform in Lab colour space
- Modify a and b colour channels
  - Semantic regions
  - Human perception
- Up to 1000 trials

Original



ColorFool-ed



# ColorFool

---

- **Transferable**
  - Randomness
  - Black-box
- **Robust**
  - Low frequency
- **Natural-looking**
  - Semantic regions
  - Human perception
- **Maintain aspect ratio**

Original



ColorFool-ed



# What about other semantic attributes of images?

---



Image structure



Image details



# Can perturbations perform enhancement, not distortion?

- Yes
- Exploiting traditional image processing filters

Original

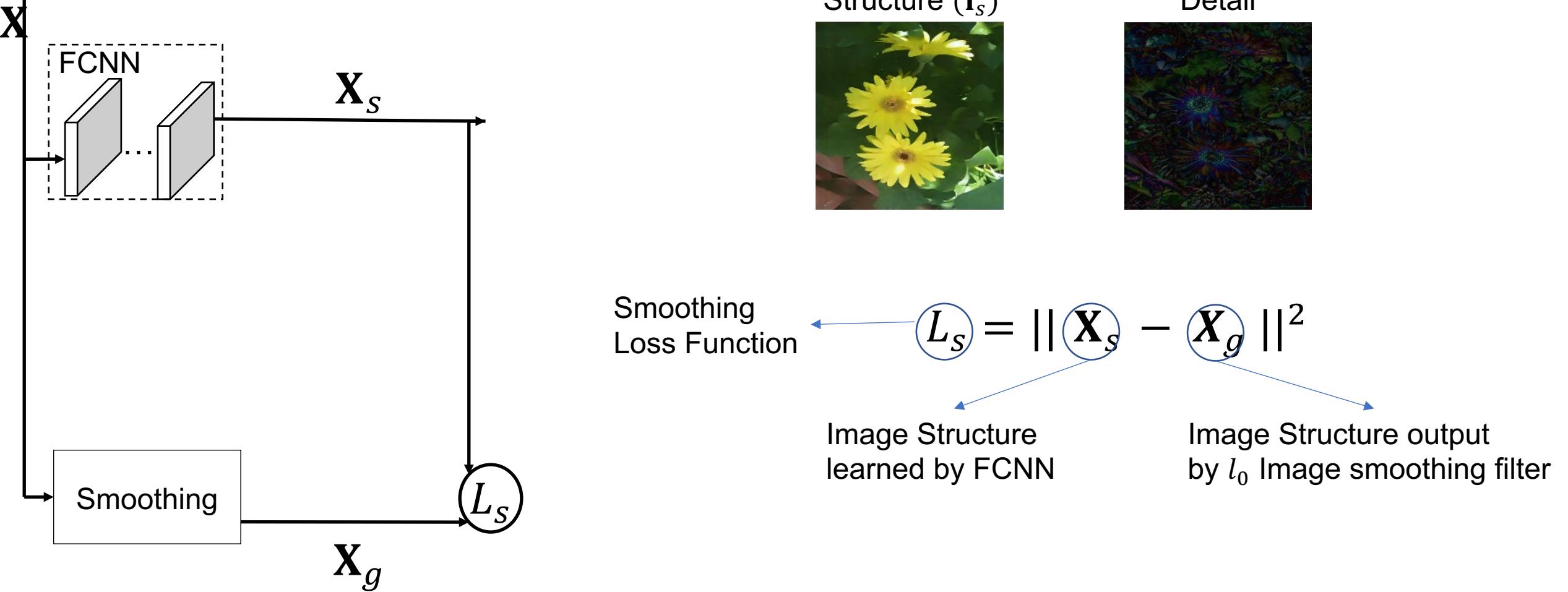


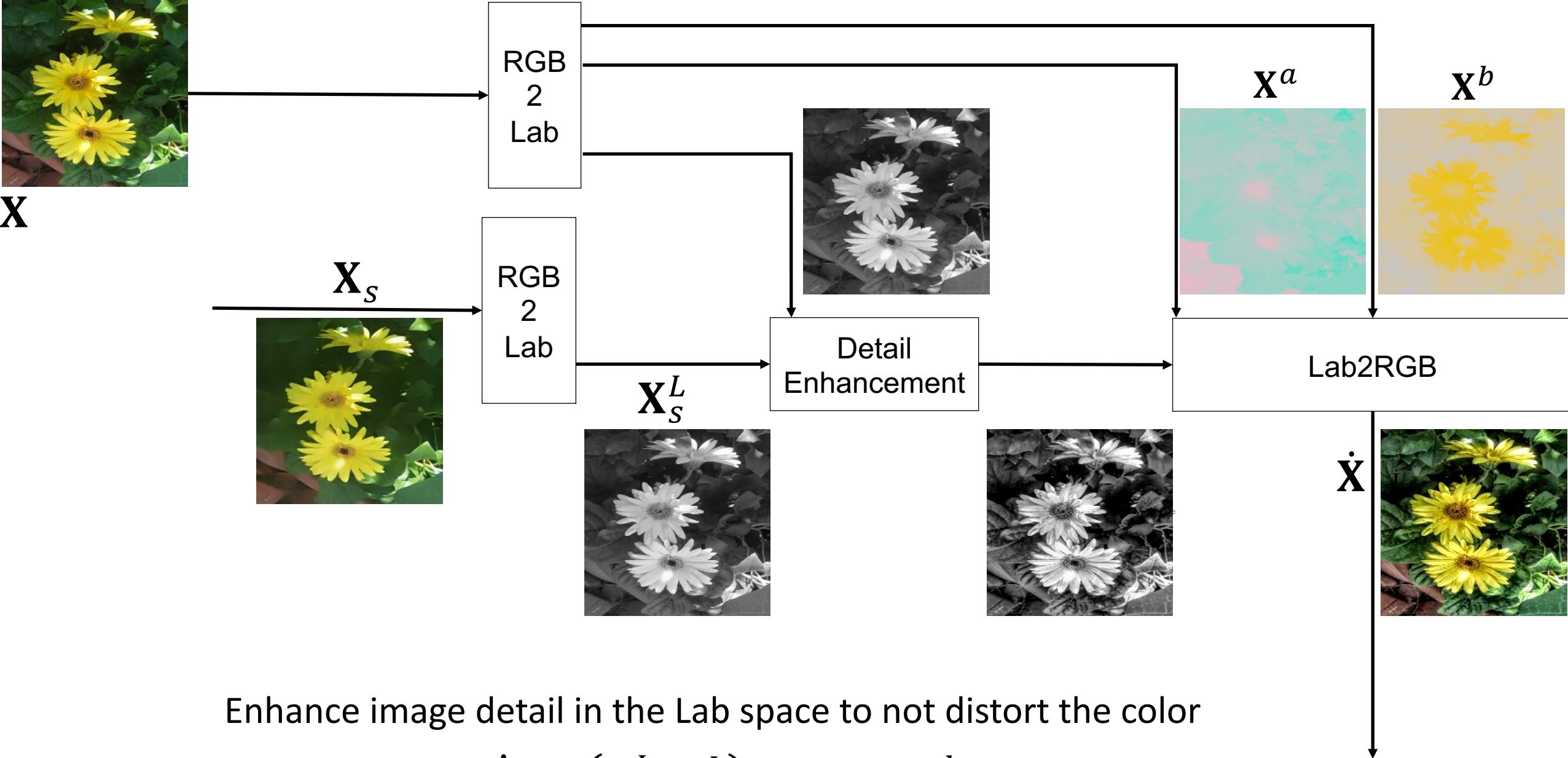
EdgeFool-ed





Decompose structure and detail by training  
a Fully Convolutional Neural Network (FCNN)







X

Guide the enhancement in an adversary manner via minimizing

$$L_a = \dot{z}_y - \max\{\dot{z}_i, i \neq y\}$$

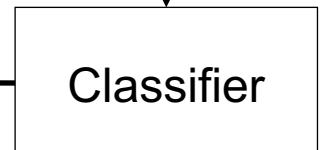
Adversarial  
Loss Function

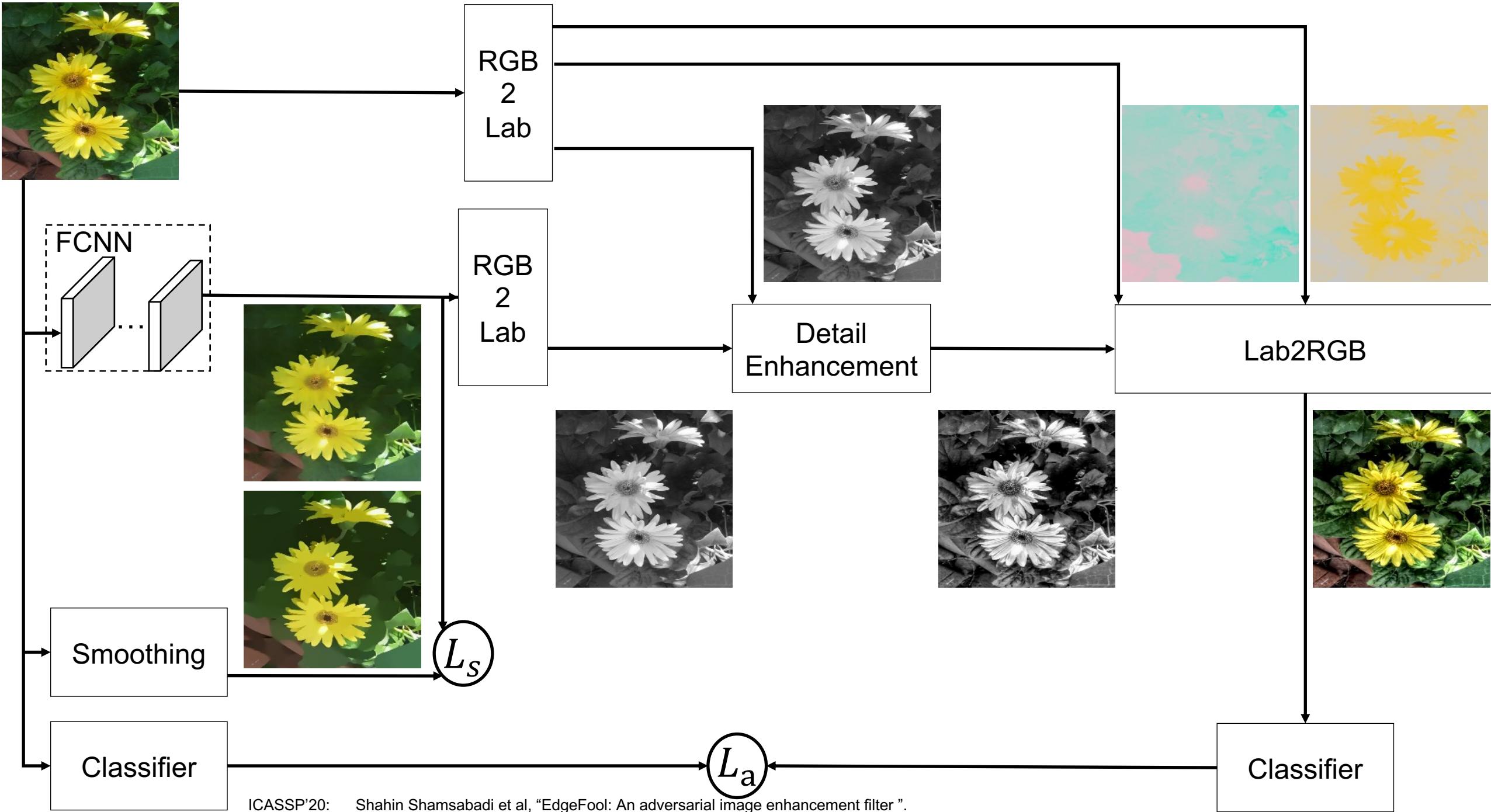
Score of belonging adversarial image  
to the same class as the original image

X



$$L_a$$





# Residual learning: if you look after other enhancements

---

- Enhanced image outputted by the filter:  $X_e$
- The filter residual:  $\delta_e = X_e - X$
- Learnable adversarial residual:  $\delta$

$$L_{l_2}(\delta, \delta_e) + (1 - \text{SSIM}(\delta, \delta_e)) + L_a$$

Prevent artefacts in untextured areas

Penalizes pixel-wise large differences      Guides the residual towards adversarial

Original images



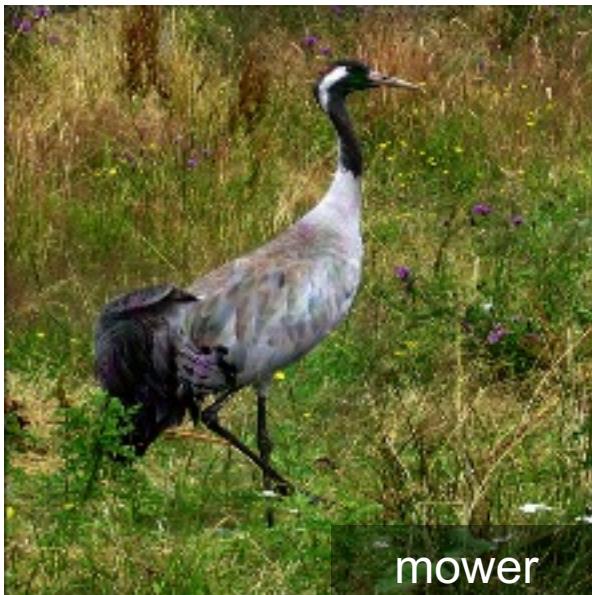
FilterFool-ed



Detail enhanced



Gamma corrected



Original images



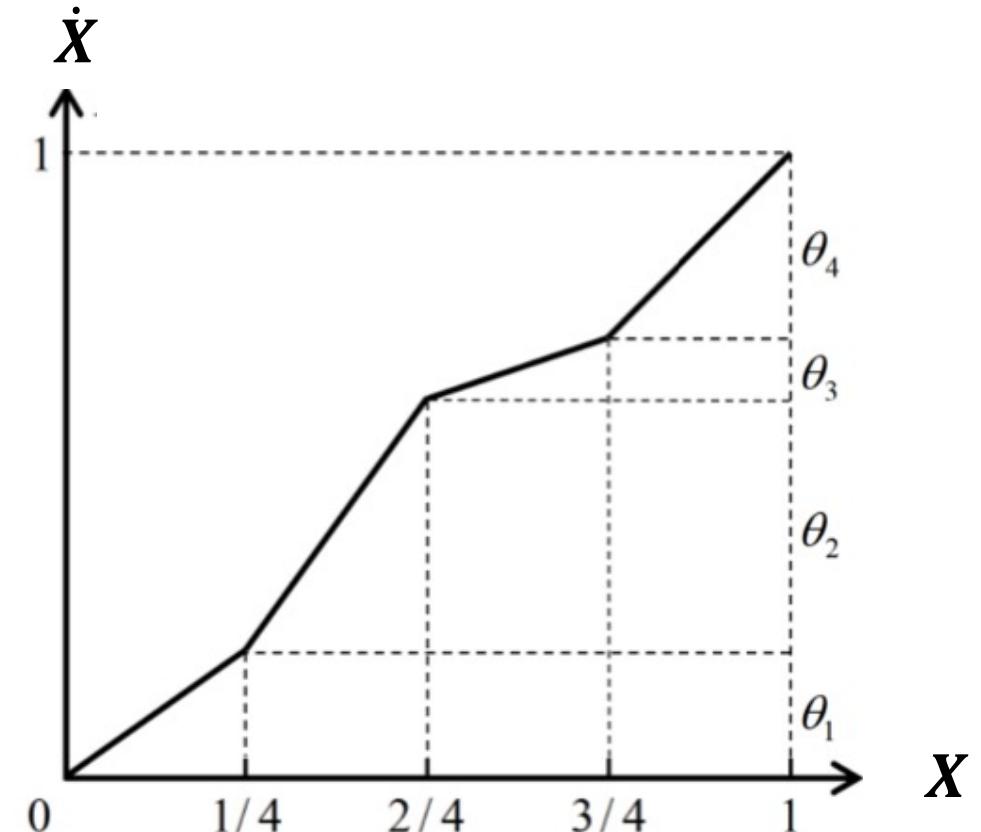
Log transformed



# Others traditional image processing filters

- Adversary learn parameters of filters to output adversarial images
- Differentiable approximation of a colour filter

$$\dot{X}_x = \sum_{i=1}^{\left[\frac{X_x}{s}\right]} \theta_i + \frac{X_x(\bmod s)}{s} \theta_{\left[\frac{X_x}{s}\right]}$$



# Now let's look at class labels

---



“screen”  
?  
“Television”



“computer keyboard”  
?  
“Typewriter keyboard”



“Shetland sheepdog”  
?  
“collie”

# Now let's look at class labels

---



Original:

“screen”



Adversarial:

“Television”

“computer keyboard”

“Typewriter keyboard”



“Shetland sheepdog”

“collie”

# Semantic relationships between class labels

---

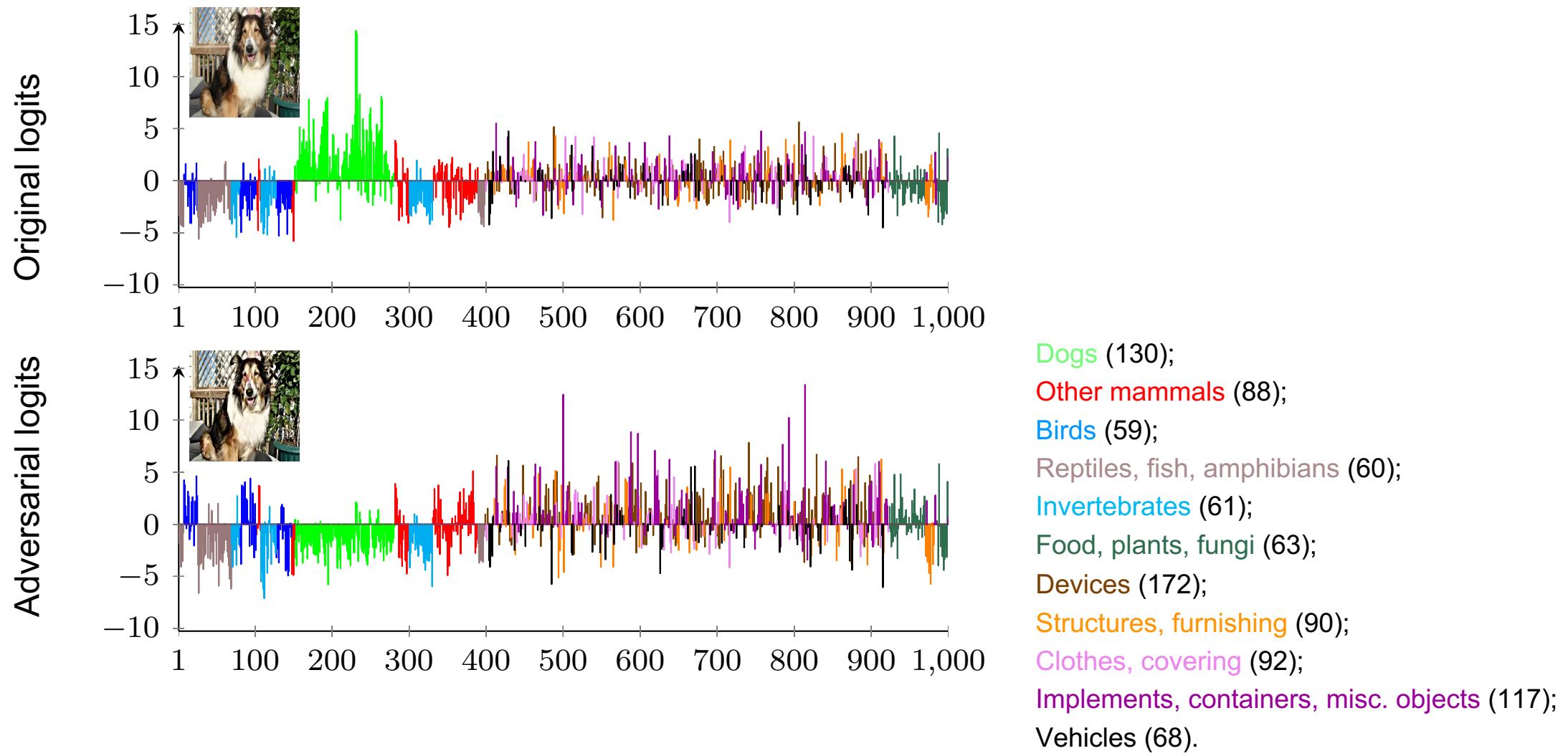
- Cluster  $D = 1000$  ImageNet classes to  $S = 11$  semantic classes
  - Dogs (130); Other mammals (88); Birds (59); Reptiles, fish, amphibians (60); Invertebrates (61); Food, plants, fungi (63); Devices (172); Structures, furnishing (90); Clothes, covering (92); Implements, containers, misc. objects (117); and Vehicles (68)
- Define mapping matrix  $\mathbf{W} \in \{0,1\}^{D \times S}$

$$L_{S-\text{Adv}}(\mathbf{X}, \dot{\mathbf{X}}) = < \text{ReLU}(\dot{\mathbf{z}}), \mathbf{w}_s > - \max(\dot{\mathbf{z}} \odot \widehat{\mathbf{w}}_s)$$

↑  
Semantically different classes  
than the original class

↓  
Semantically similar classes  
to the original class

# Semantic relationships between class labels



# Experimental settings

---

- Dataset
  - ImageNet:  $3K$  images of  $1k$  objects
- Classifiers under attacks
  - ResNet50, ResNet18, AlexNet
- Visualization
  - Perturbations
  - Adversarial images
  - Top5 predictions and classifier attention
- Success rate and transferability

$$\frac{\text{\# successful adversarial images}}{\text{\# total images}}$$

# Adversarial perturbations

---

Original



Basic Iterative Method



SparseFool



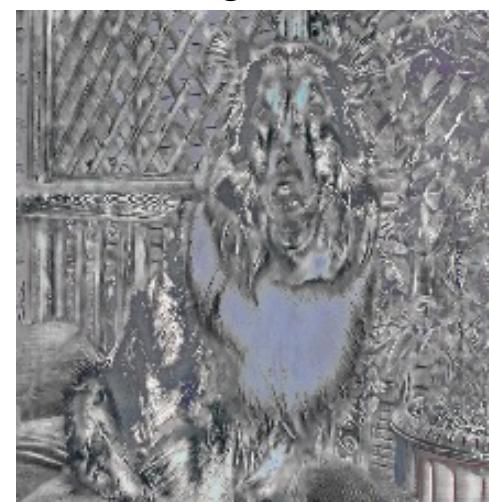
SemanticAdv



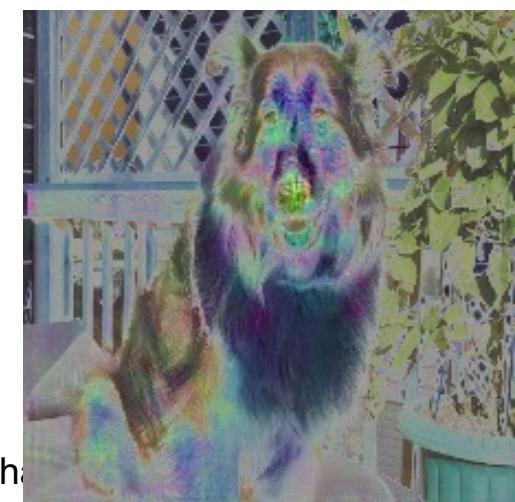
ColorFool



EdgeFool



FilterFool



a.sh

# Adversarial images

---

Original



Basic Iterative Method



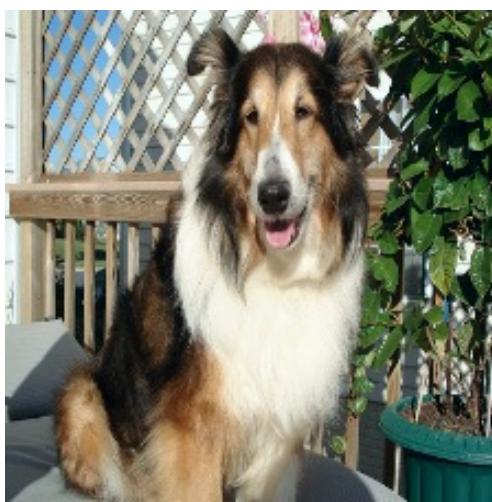
SparseFool



SemanticAdv



ColorFool



EdgeFool



FilterFool



# Top-5 predictions

---

Original



Basic Iterative Method



SparseFool



SemanticAdv



ColorFool



EdgeFool



FilterFool



# Classifier attention

---

Original



Basic Iterative Method



SparseFool



SemanticAdv



ColorFool



EdgeFool



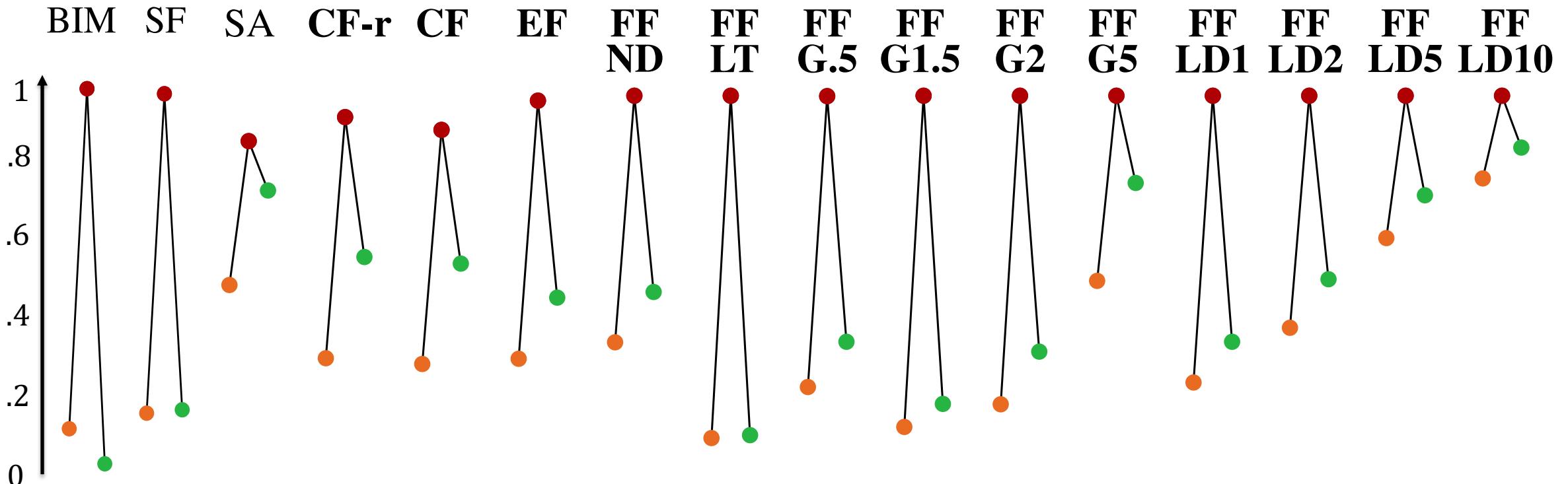
FilterFool



a.sh

# Success rate and transferability

---



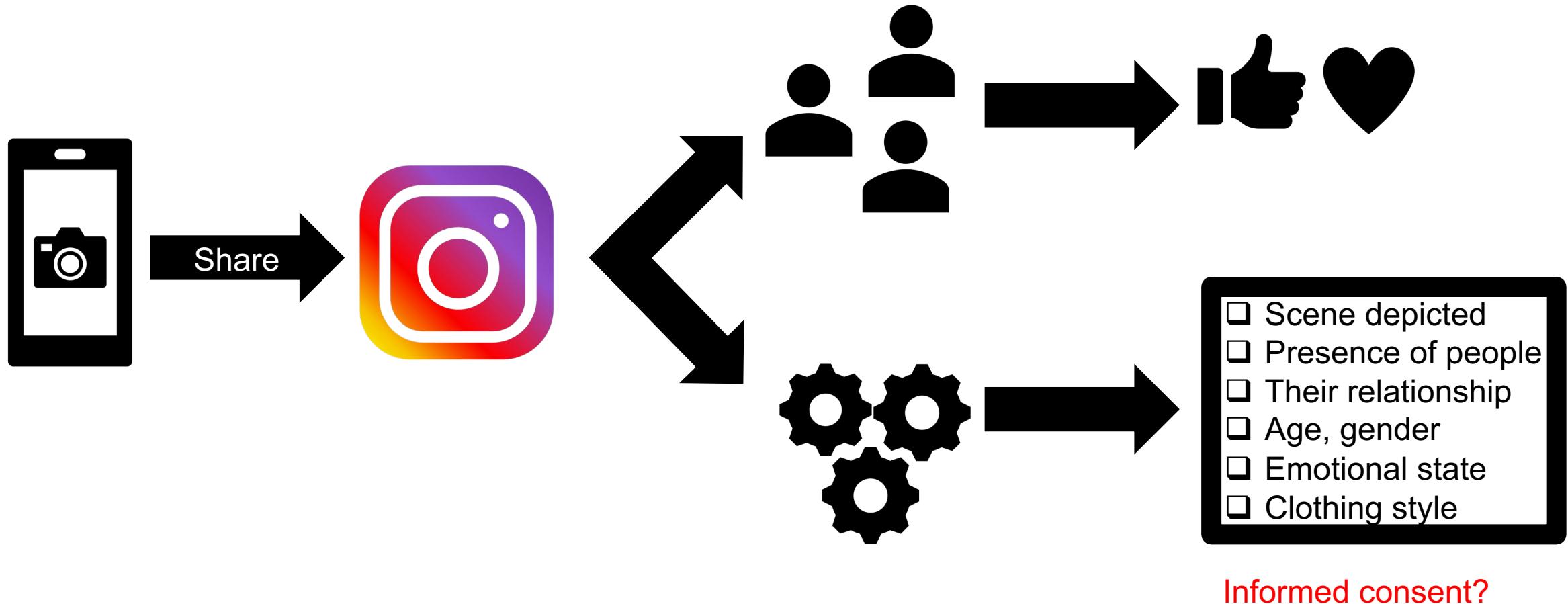
ResNet50 (seen) (●)

ResNet18 (unseen) (○)

AlexNet (unseen) (●)

BIM: Basic Iterative Method, SF: SparseFool, SA: SemanticAdv,  
EF: EdgeFool, FF: FilterFool, LT: Log transformation,  
LD, ND: Linear and Nonlinear Detail enhancement,  
G: Gamma correction

# Privacy risks in image sharing social medias



## Other works

---

- Exploiting vulnerabilities of deep neural networks for privacy protection  
(<https://arxiv.org/pdf/2007.09766.pdf>)
- Fawkes: Protecting Privacy against Unauthorized Deep Learning Models  
(<https://www.usenix.org/system/files/sec20-shan.pdf>)
- Face-Off: Adversarial Face Obfuscation  
(<https://arxiv.org/pdf/2003.08861.pdf>)
- Pixel privacy challenge  
(<http://www.multimediaeval.org/mediaeval2019/pixelprivacy/>)

# Still many rooms for improvements

---

- Define mathematically the possible set of adversarial images
  - Model human eyes
- Why classifiers are vulnerable to semantic changes?
  - Colour shifting
  - Detail enhancement
- Certified adversarial perturbations
- Make semantic changes outside of digital domains in physical world

# Adversarial Examples and Their Applications in Privacy

**Ali Shahin Shamsabadi**



Research Associate  
The Alan Turing Institute  
AI Programme  
 [@AliShahinShams1](https://twitter.com/AliShahinShams1)  
<https://alishahin.github.io>  
<https://www.turing.ac.uk/people/researchers/ali-shahin-shmasabadi>

**The  
Alan Turing  
Institute**

# Now let's see how we can implement ColorFool-ed images

---

Original



ColorFool-ed



<https://github.com/smartcameras/ColorFool>

## Demo

<https://github.com/smartcameras/ColorFool/blob/master/TutorialDemoColorFool/ColorFool.ipynb>