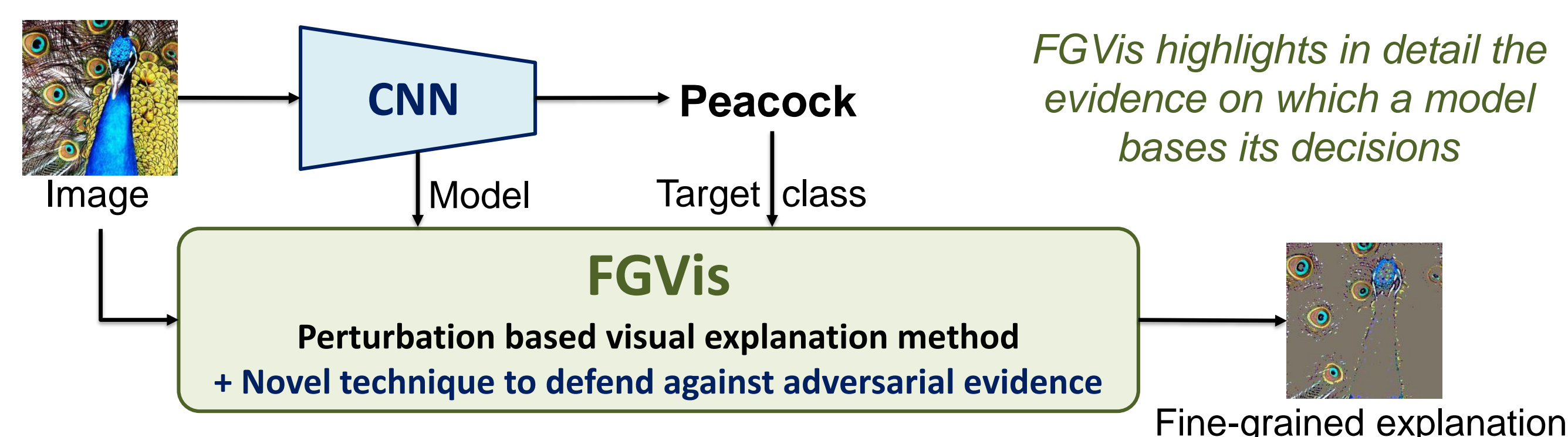


1 Overview

Motivation: A better understanding of the decision-making process of a CNN is required to provide hints for improving it. This allows to uncover and understand failure cases, limitations of the model, and shortcomings of the training data.

Fine-grained visual explanation method (FGVis):



Contributions:

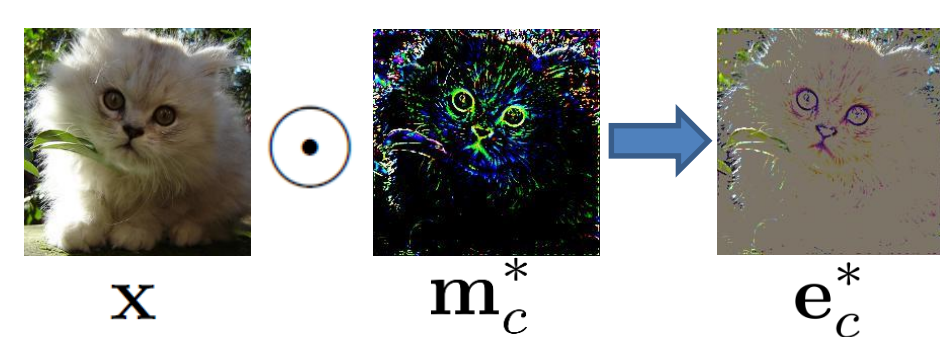
- A method (**FGVis**) to generate **fine-grained explanations** in the image space.
- A **novel technique for defending against adversarial evidence**, which does not depend on human-tuned parameters.
- Interpretable** and **class discriminative** explanations, visualizing **detailed evidence**.

2 Perturbation based explanation methods

An explanation e_c^* is computed by perturbing the input image x .

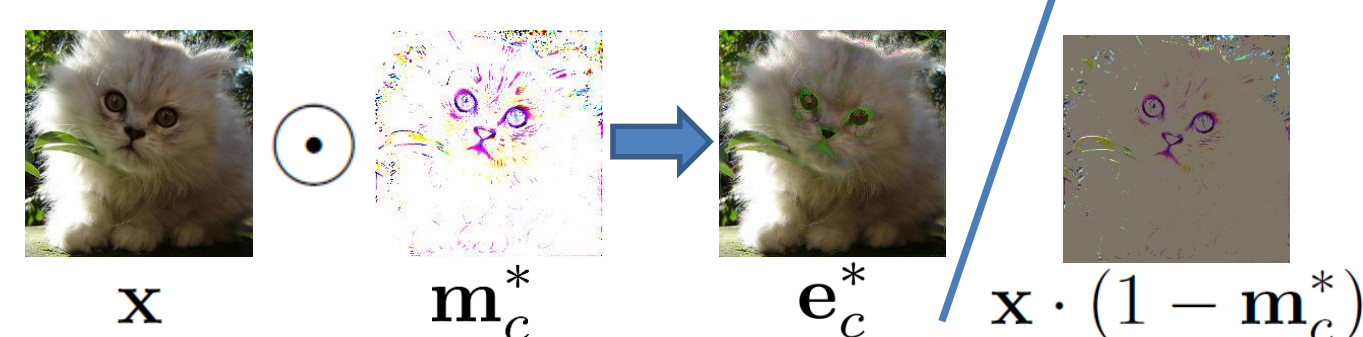
Mask m_c^* based perturbation: $e_c^* = x \cdot m_c^*$ c : Target class of the explanation

Preserving explanation:



$$m_c^* = \arg \min_{m_c \in [0,1]^{3 \times H \times W}} \lambda \cdot \|m_c\|_1 - \text{CNN}_c(e_c)$$

Deleting explanation:



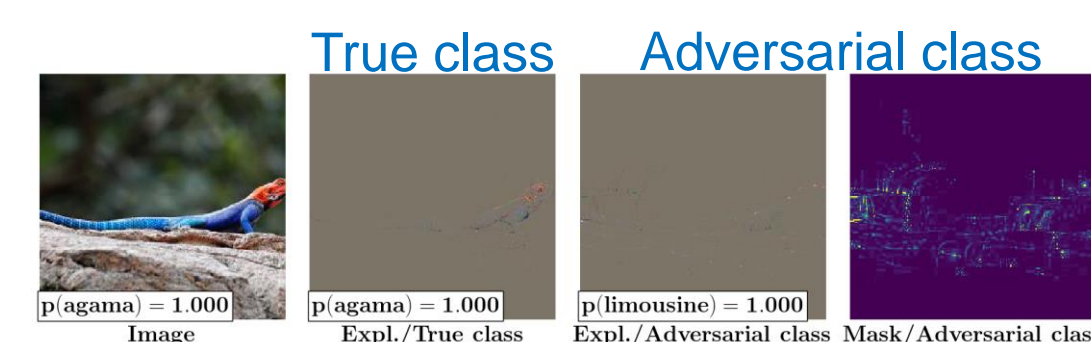
$$m_c^* = \arg \max_{m_c \in [0,1]^{3 \times H \times W}} \lambda \cdot \|m_c\|_1 - \text{CNN}_c(e_c)$$

Probability of target class c for explanation e_c

Perturbation based explanations represent valid model inputs and are thus testable

3 Defending against adversarial evidence

Drawback of perturbation based methods: Adversarial evidence, i.e. faulty evidence due to artefacts introduced in the optimization of the explanation.



Without defense the optimization introduces adversarial evidence

Novel adversarial defense:

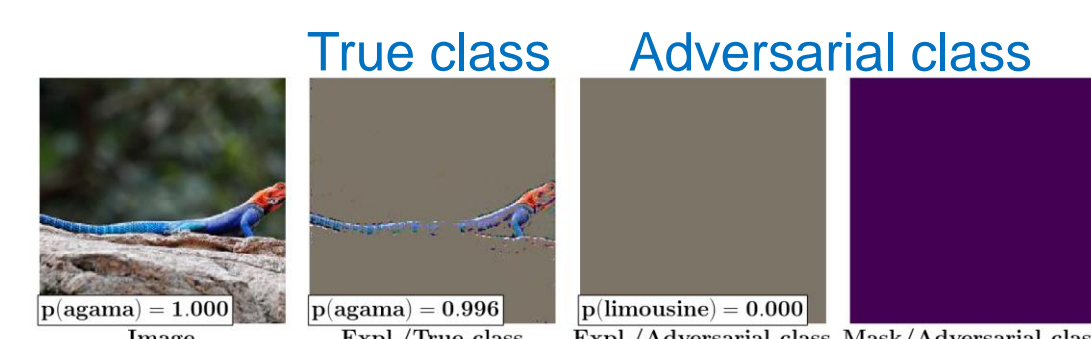
- Idea:** The features in an explanation should be a subset of the image features.
- Corresponding optimization constraint:**

$$\begin{cases} 0 \leq h_i^l(e_c) \leq h_i^l(x), & \text{if } h_i^l(x) \geq 0, \\ 0 \geq h_i^l(e_c) \geq h_i^l(x), & \text{otherwise,} \end{cases}$$
 - $h_i^l(\cdot)$: Activation of the i -th neuron in the l -th layer.
 - The constraint is applied after each nonlinearity-layer (e.g.: ReLU-Layer).
- Implemented via gradient clipping:**

$$\gamma_i^l = \bar{\gamma}_i^l \cdot \mathbb{1}[h_i^l(e_c) \leq \max(0, h_i^l(x))] \cdot \mathbb{1}[h_i^l(e_c) \geq \min(0, h_i^l(x))]$$

Indicator function

Updated error-signal back-propagated through the l -th layer

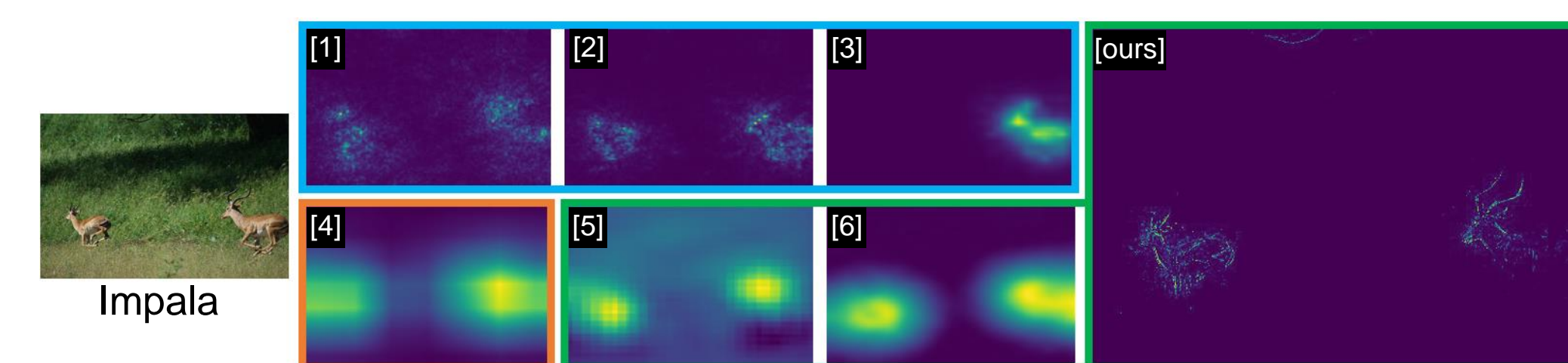


Our defense prevents the hallucination of adversarial evidence

Our defense does not depend on human-tuned parameters and enables an explanation which is both fine-grained and preserves the characteristics of the image

4 Experiments

Qualitative comparison with other methods



FGVis generates the most fine-grained explanation mask

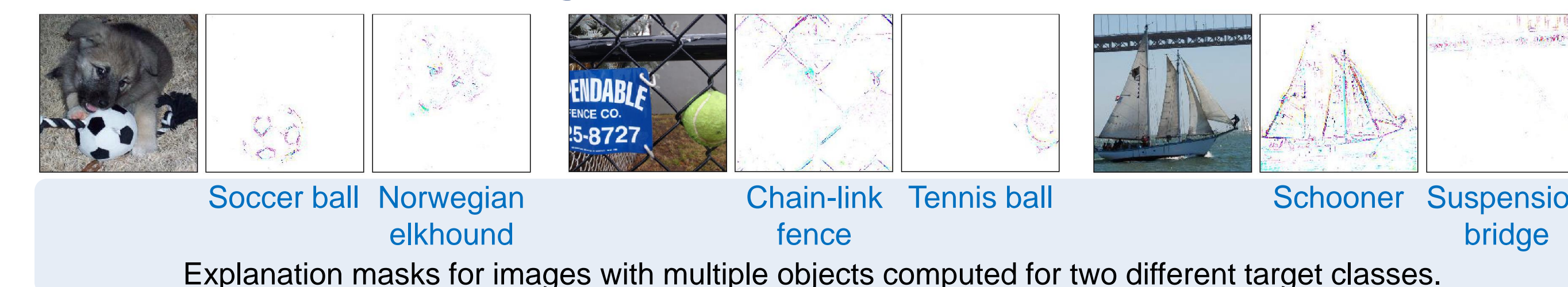
Backpropagation based methods [1,2,3]

Activation based method [4]

Perturbation based methods [5, 6, ours]

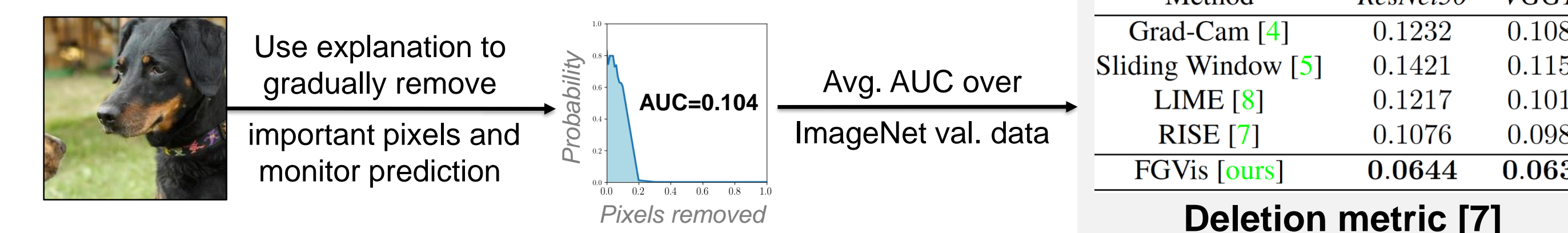
5 Experiments

Class discriminative / fine-grained



FGVis produces class discriminative explanations even when objects partially overlap

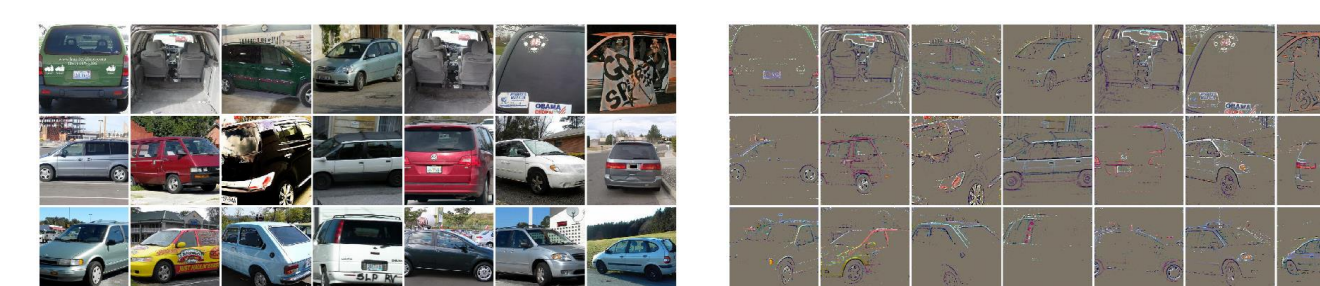
Faithfulness of explanations: How accurate does an explanation represent the true evidence on which a model bases its decision?



Color bias of VGG16 trained on ImageNet.



FGVis can provide a first indication for the importance of different colors



Quantitative verification: Ratio of maintained true classifications after swapping the color channels

Class	BGR	RGB	GRB	Avg. RGB, GRB
school bus	100 %	9.5 %	7.1 %	8.3 %
minivan	100 %	71.4 %	95.2 %	83.3 %

6 References

- [1] Karen Simonyan *et al.* (ICLR, 2014). Deep inside convolutional networks: Visualising image classification models and saliency maps.
- [2] Jost T. Springenberg *et al.* (ICLR, 2015). Striving for simplicity: The all convolutional net.
- [3] Jianming Zhang *et al.* (ECCV, 2016). Top-down neural attention by excitation backprop.
- [4] Ramprasaath R. Selvaraju (ICCV, 2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization.
- [5] Matthew D. Zeiler and Rob Fergus (ECCV, 2014). Visualizing and understanding convolutional networks.
- [6] Ruth C. Fong and Andrea Vedaldi (ICCV, 2017). Interpretable explanations of black boxes by meaningful perturbation.
- [7] Vitali Petsiuk *et al.* (BMVC, 2018). Rise: Randomized input sampling for explanation of black-box models.
- [8] Marco T. Ribeiro *et al.* (SIGKDD, 2016). Why should I trust you?: Explaining the predictions of any classifier.