

# How to explain neural network decisions?

**Anh Nguyen**

Assistant Professor



[@anh\\_ng8](#)

<http://AnhNguyen.me>

# Challenge: Out-of-distribution Generalization

- Leon Bottou (ICML 2019 Keynote)
- Yoshua Bengio (NeurIPS 2019 Keynote, 2019 AI Debate)



ImageNet  
school bus class



Task: Image classification

Inception-v3

78% accuracy on ImageNet



100%  
? school bus



0% school bus!

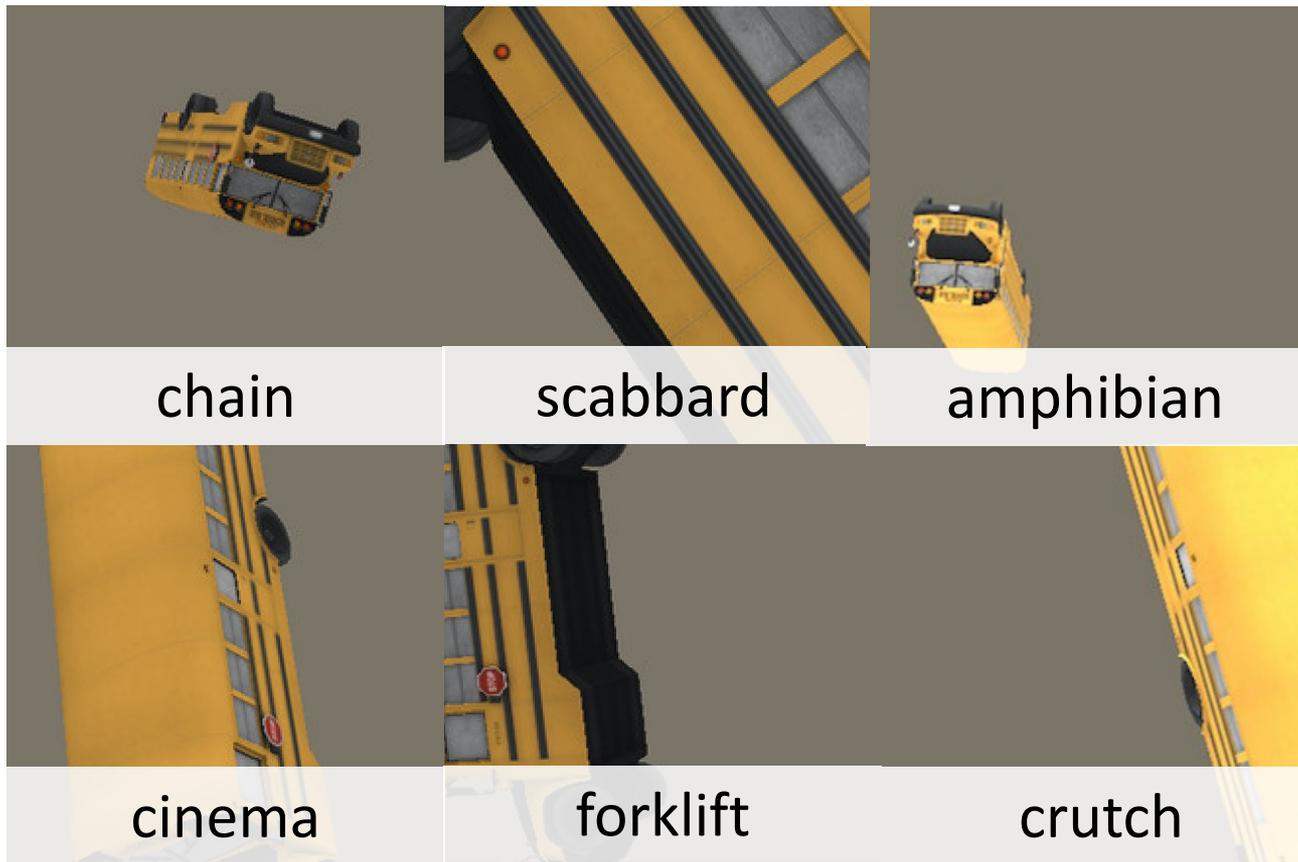
## Inception-v3

78% accuracy on ImageNet



**100%**

school bus



**100%** confidence

1.00: vacuum, vacuum cleaner

0.99: vacuum, vacuum cleaner

1.00: vacuum, vacuum cleaner

0.99: vacuum, vacuum cleaner

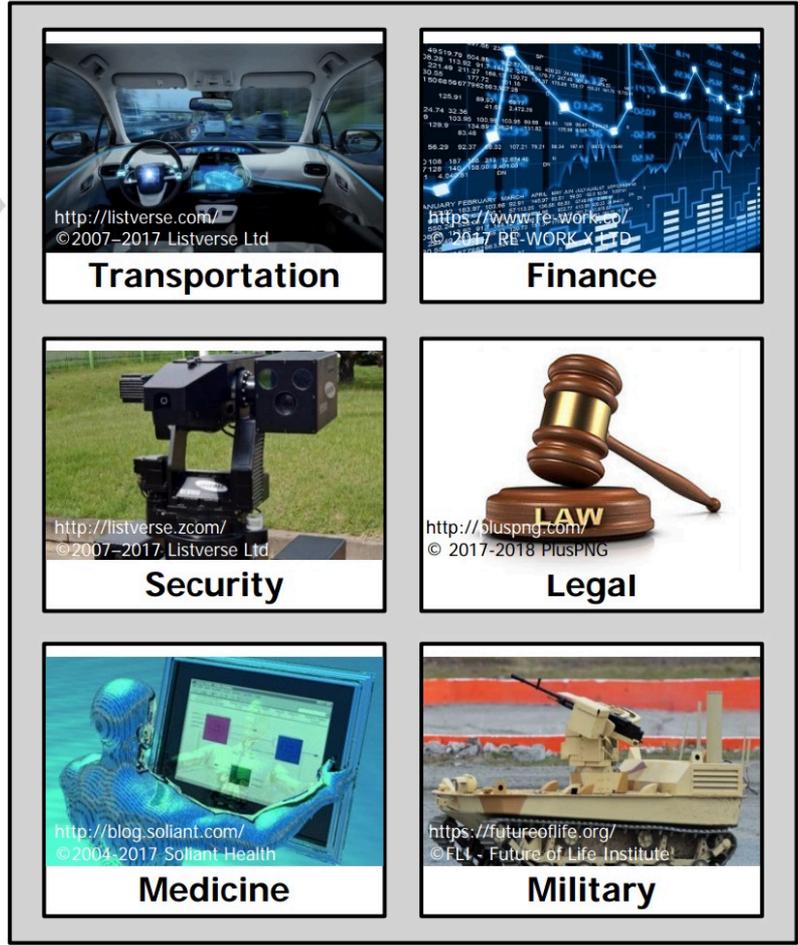
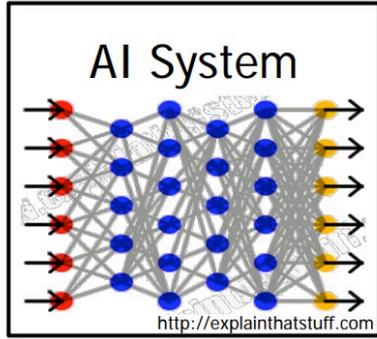


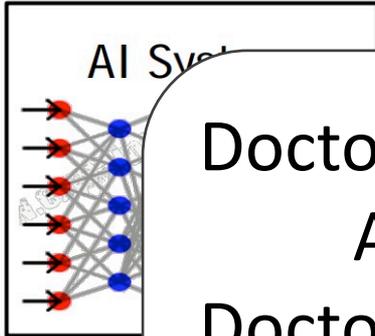
vacuum cleaner

# Challenge: Out-of-distribution Generalization



shopping basket



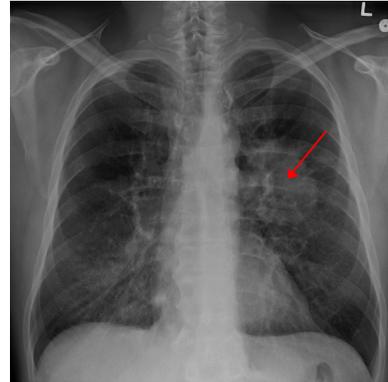


Doctor : Is this tumor malignant?

AI : Yes! 100%.

Doctor : **Why??**

Patient : !@#\$%?



# Performance vs. Explainability



<https://www.darpa.mil/attachments/XAIProgramUpdate.pdf>

# Performance vs. Explainability



<https://www.darpa.mil/attachments/XAIProgramUpdate.pdf>

# Performance vs. Explainability



<https://www.darpa.mil/attachments/XAIProgramUpdate.pdf>

## Measure of Explanation Effectiveness

### User Satisfaction

- Clarity of the explanation (user rating)
- Utility of the explanation (user rating)

### Mental Model

- Understanding individual decisions
- Understanding the overall model
- Strength/weakness assessment
- 'What will it do' prediction
- 'How do I intervene' prediction

### Task Performance

- Does the explanation improve the user's decision, task performance?
- Artificial decision tasks introduced to diagnose the user's understanding

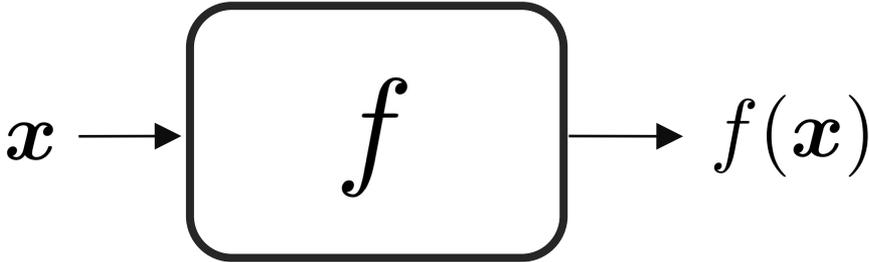
### Trust Assessment

- Appropriate future use and trust

### Correctability (Extra Credit)

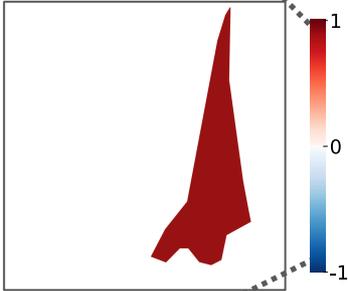
- Identifying errors
- Correcting errors
- Continuous training

# Attribution maps as explanations



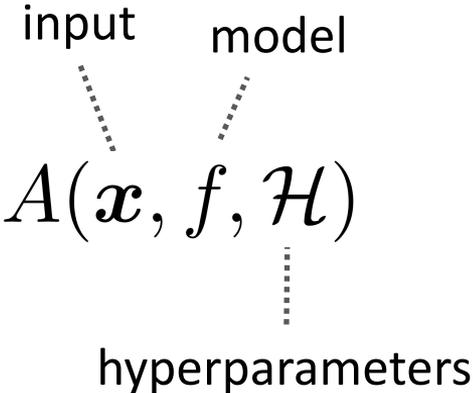
0.54 matchstick

*for* matchstick



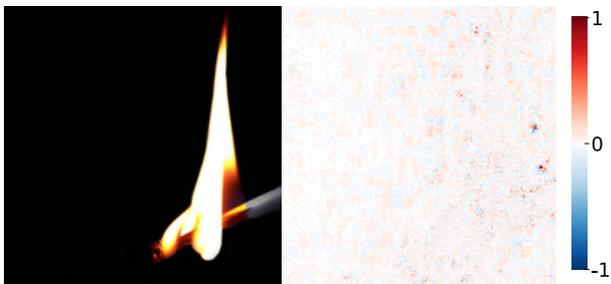
attribution map  
(hypothetical)

*against* matchstick



# Method 0: Saliency maps

Gradient



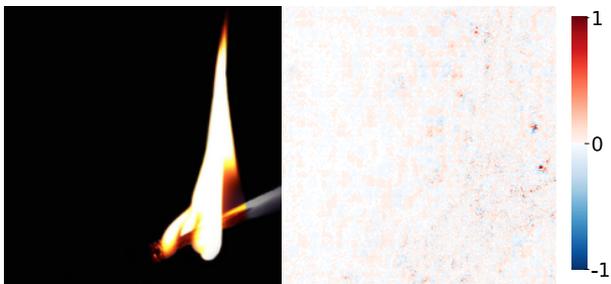
$$\nabla_x f$$

## Problems:

- too noisy

# Method 1: Smoothed saliency maps

Gradient



$$\nabla_x f$$

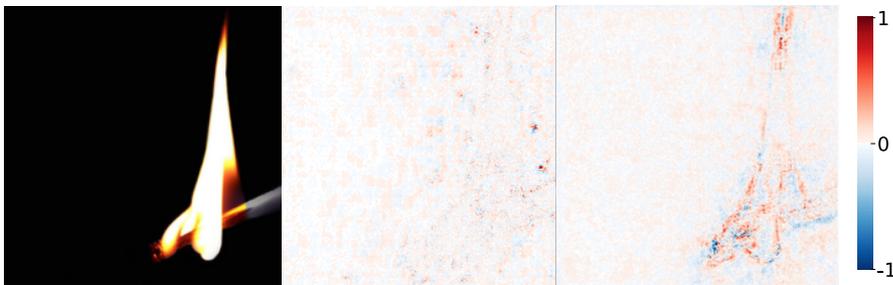
## Problems:

- too noisy

# Method 1: Smoothed saliency maps

Smilkov et al. 2017

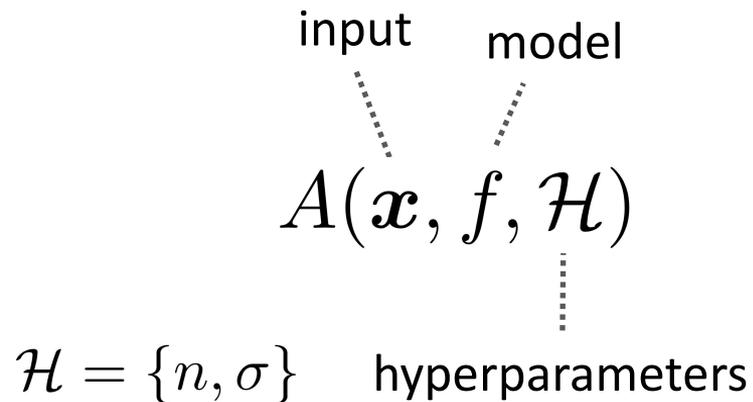
Gradient    SmoothGrad



$$\frac{1}{n} \sum_1^n \nabla_{\mathbf{x}} f(\mathbf{x} + \mathcal{N}(0, \sigma^2))$$

## Problems:

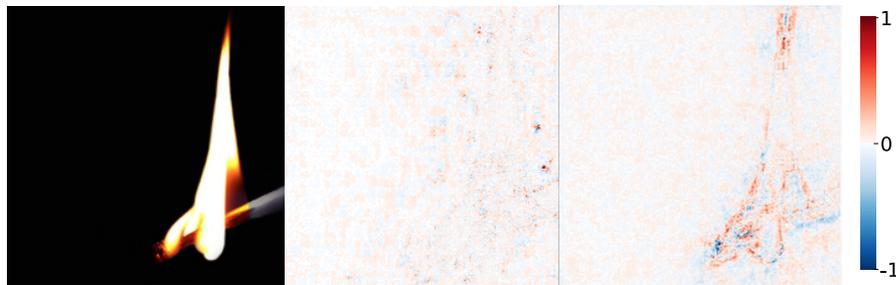
- ~~too noisy~~



# Method 1: Smoothed saliency maps

Smilkov et al. 2017

Gradient    SmoothGrad



$$\frac{1}{n} \sum_{1}^n \nabla_{\mathbf{x}} f(\mathbf{x} + \mathcal{N}(0, \sigma^2))$$

## Problems:

- ~~too noisy~~
- sensitivity  $\neq$  importance

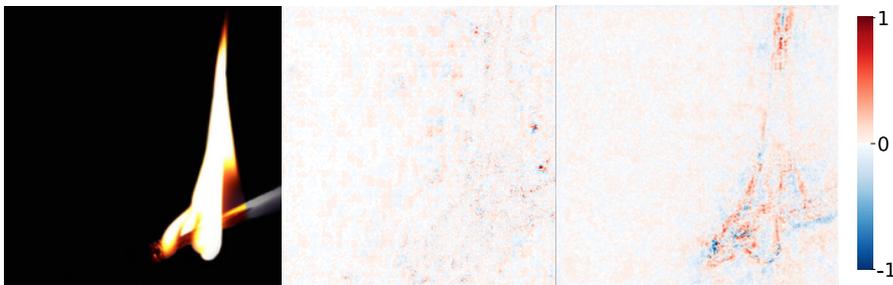
$$\begin{array}{cc} \text{input} & \text{model} \\ \vdots & \vdots \\ A(\mathbf{x}, f, \mathcal{H}) & \\ \vdots & \end{array}$$

$$\mathcal{H} = \{n, \sigma\} \quad \text{hyperparameters}$$

# Method 1: Smoothed saliency maps

Smilkov et al. 2017

Gradient SmoothGrad



$$\frac{1}{n} \sum_1^n \nabla_x f(\mathbf{x} + \mathcal{N}(0, \sigma^2))$$

How does prediction change if the flame is slightly *whiter*?

How does prediction change if the flame is *NOT present*?

## Problems:

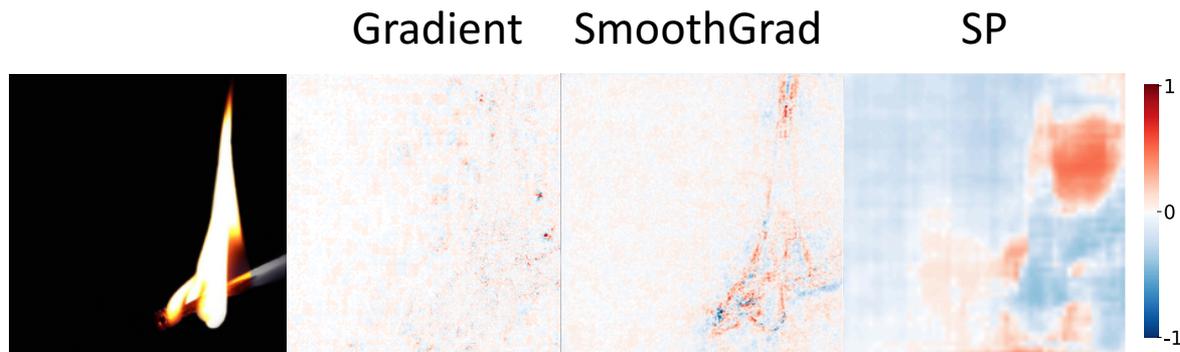
- ~~too noisy~~
- sensitivity  $\neq$  importance

input model

$\mathcal{H} = \{n, \sigma\}$  hyperparameters

# Method 2: Sliding Patch

Zeiler & Fergus 2014



How does prediction change if the flame is *NOT present*?

importance

# Method 2: Sliding Patch

Zeiler & Fergus 2014



How does prediction change if the flame is *NOT present*?

importance

# Method 2: Sliding Patch

Zeiler & Fergus 2014

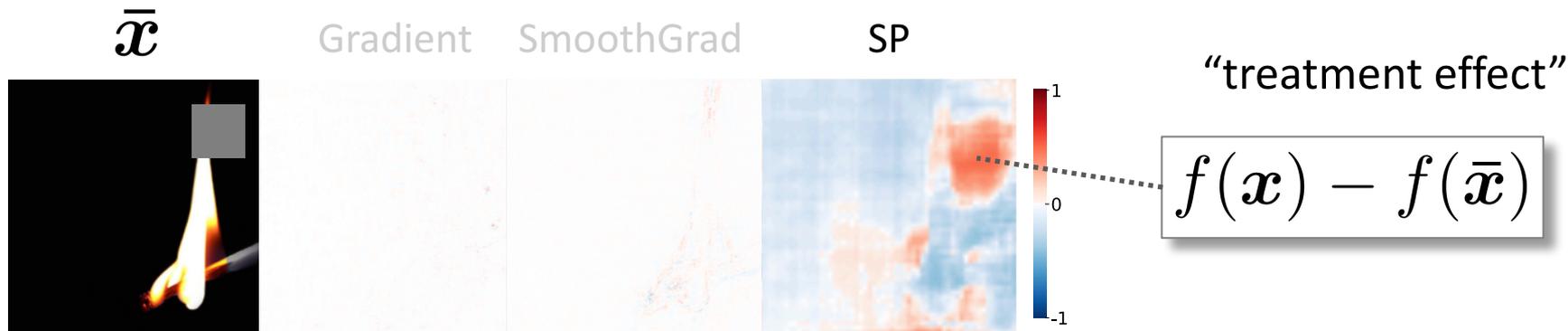


How does prediction change if the flame is *NOT present*?

importance

# Method 2: Sliding Patch

Zeiler & Fergus 2014

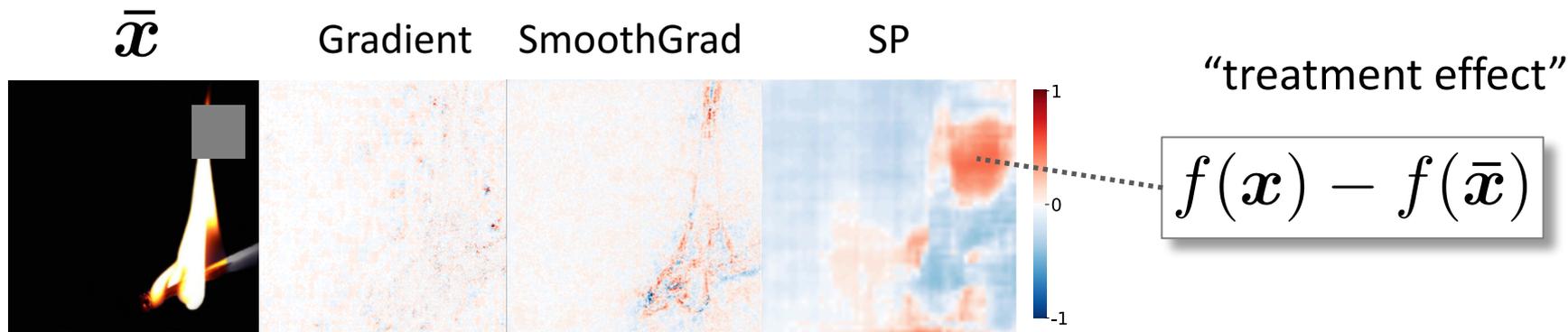


How does prediction change if the flame is *NOT present*?

importance

# Method 2: Sliding Patch

Zeiler & Fergus 2014

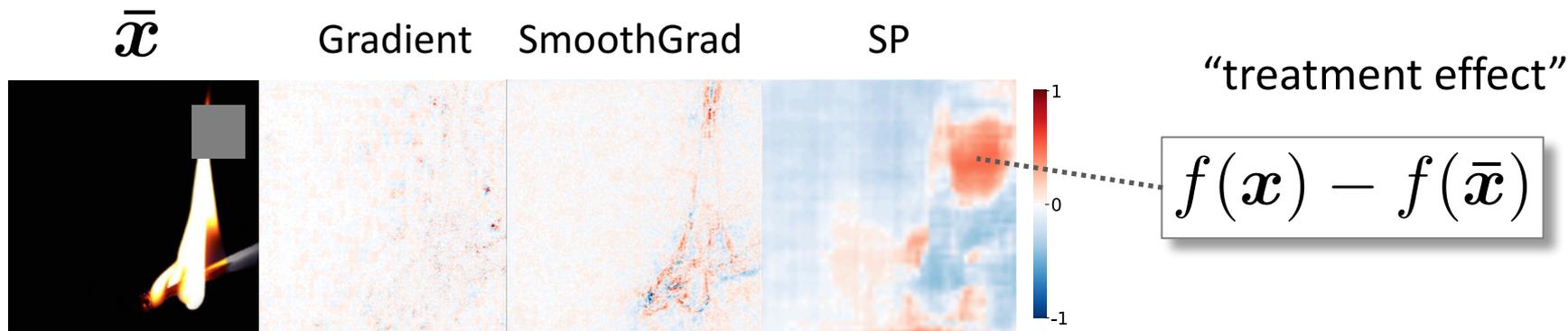


## Problems:

- ~~sensitivity  $\neq$  importance~~
- ~~too noisy~~

# Method 2: Sliding Patch

Zeiler & Fergus 2014

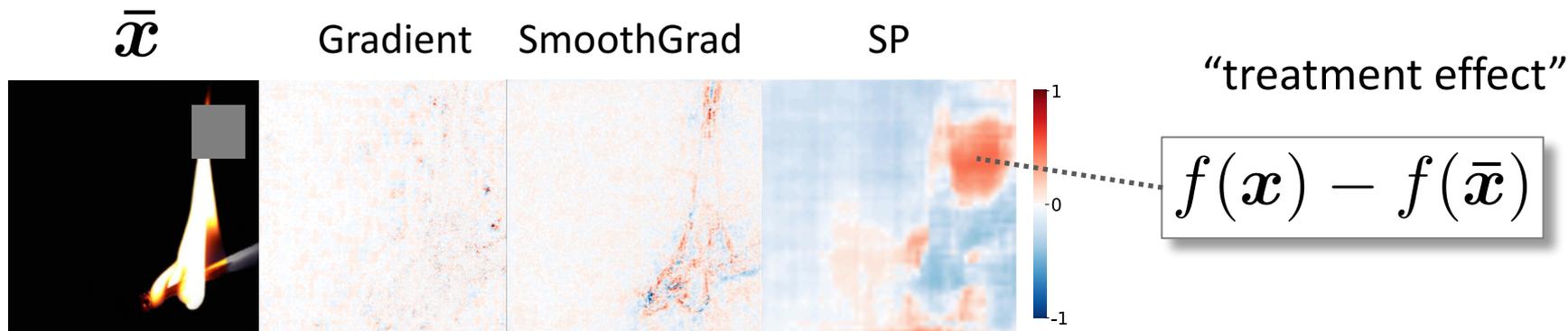


## Problems:

- ~~too noisy~~
- ~~sensitivity  $\neq$  importance~~
- coarse, inaccurate (due to square patch)

# Method 2: Sliding Patch

Zeiler & Fergus 2014

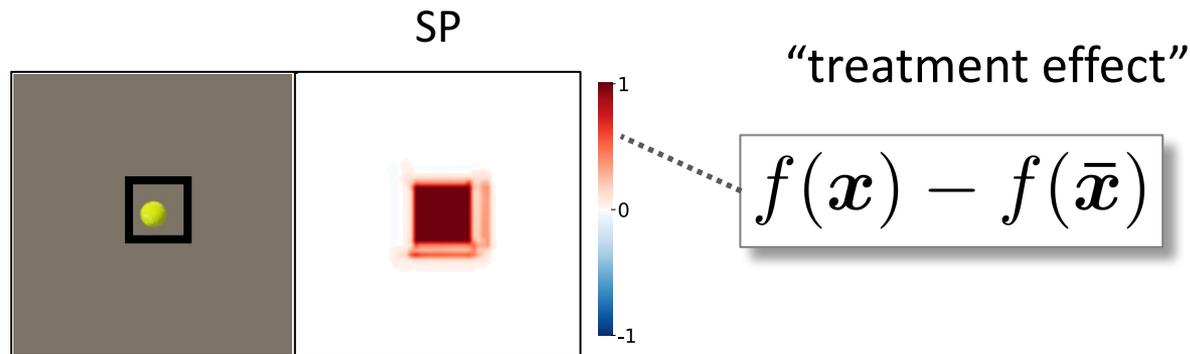


## Problems:

- coarse, inaccurate (due to square patch)

# Method 2: Sliding Patch

Zeiler & Fergus 2014



## Problems:

- coarse, inaccurate (due to square patch)

# Method 3: LIME

Ribeiro et al. 2016



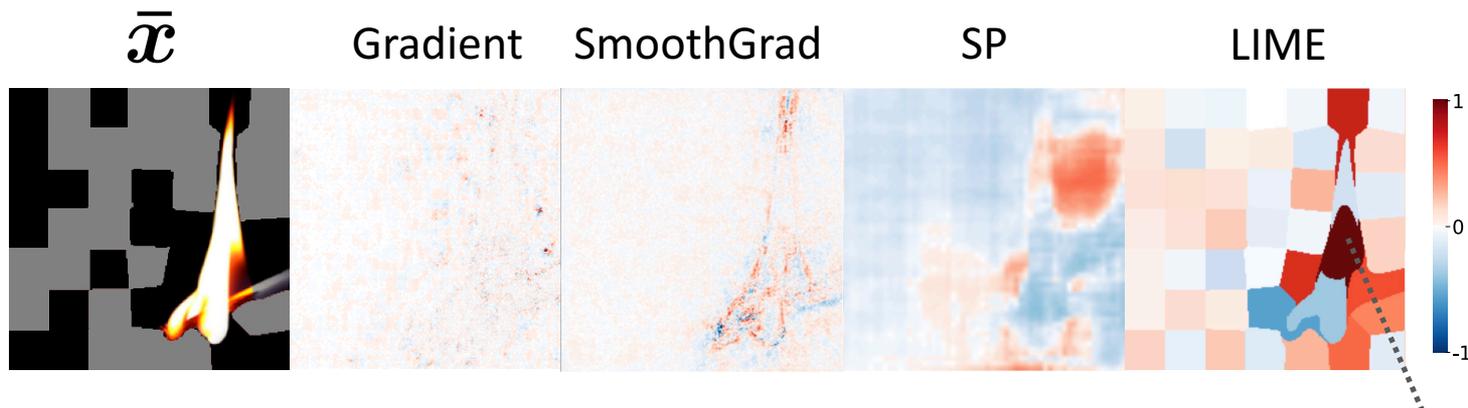
**Idea:** Attribution for superpixel  $k$  is the average score when  $k$  is visible

## Problems:

- coarse, inaccurate (due to square patch)

# Method 3: LIME

Ribeiro et al. 2016



**Idea:** Attribution for superpixel  $k$  is the average score when  $k$  is visible

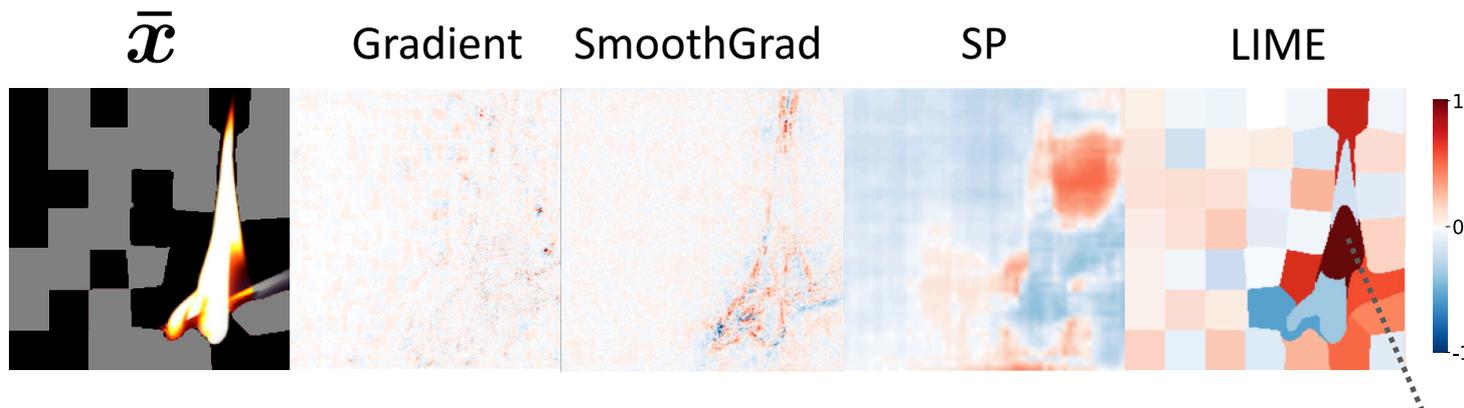
## Problems:

- coarse ~~inaccurate (due to square patch)~~
- gray pixels?

$$\approx \frac{1}{n} \sum_{i=1}^n \frac{f(\bar{x}^i)}{S/2}$$

# Method 3: LIME

Ribeiro et al. 2016



**Idea:** Attribution for superpixel  $k$  is the average score when  $k$  is visible

## Problems:

- coarse
- gray pixels?

$$\approx \frac{1}{n} \sum_{i=1}^n \frac{f(\bar{x}^i)}{S/2}$$

# Method 4: Learned blur mask

Fong & Vedaldi 2017

$\bar{x}$



Gradient

SmoothGrad

SP

LIME

MP



## Problems:

- coarse
- gray pixels?

**Idea:** Identify a *minimal* region s.t. when blurred out would minimize classification score

$$\mathbf{m}^* = \arg \min_{\mathbf{m}} \lambda \|\mathbf{m}\|_1 + f(\text{blur}(x, \mathbf{m}))$$

# Method 4: Learned blur mask

Fong & Vedaldi 2017

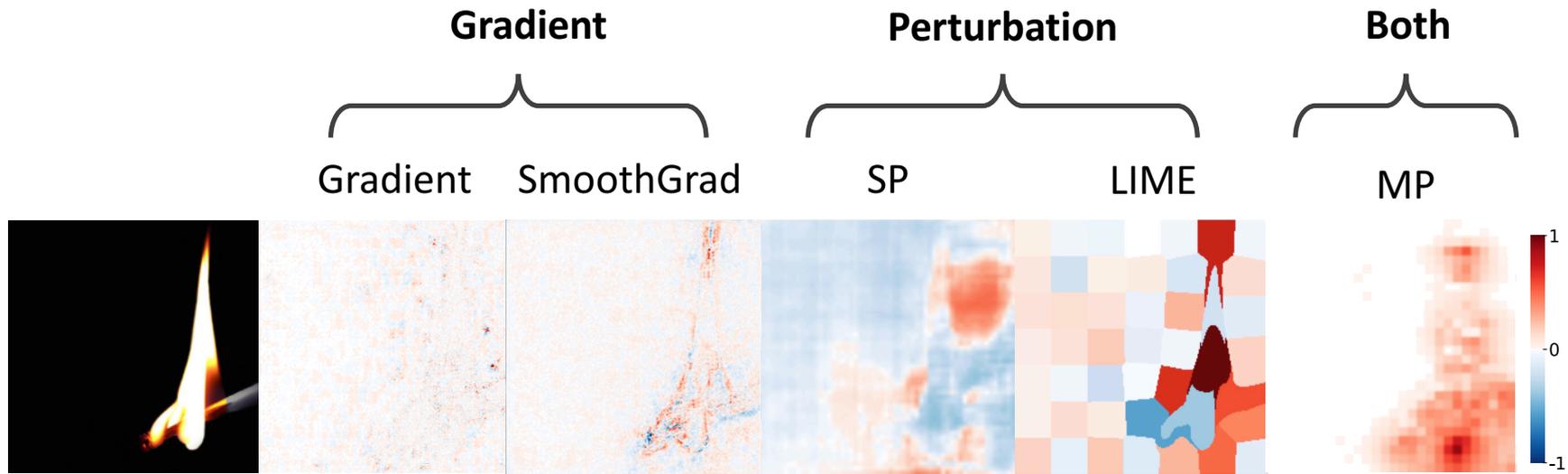


**Idea:** Identify a *minimal* region s.t. when blurred out would minimize classification score

## Problems:

- ~~coarse~~ → fine-grained
- ~~gray~~ → blurred pixels

$$\mathbf{m}^* = \arg \min_{\mathbf{m}} \lambda \|\mathbf{m}\|_1 + f(\text{blur}(\mathbf{x}, \mathbf{m}))$$

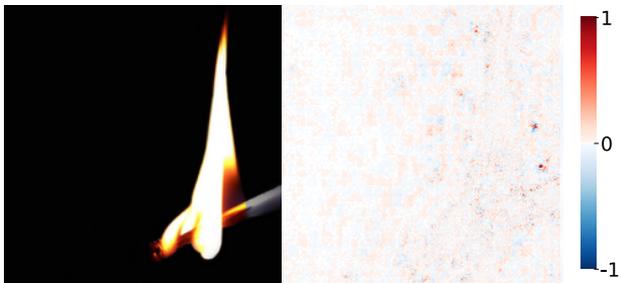


0.54 matchstick

Are these explanations correct and reliable?

# #1: Saliency maps may NOT be too noisy!

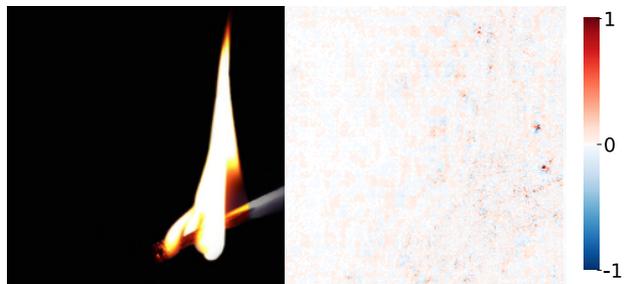
Gradient



GoogLeNet

# #1: Saliency maps may NOT be too noisy!

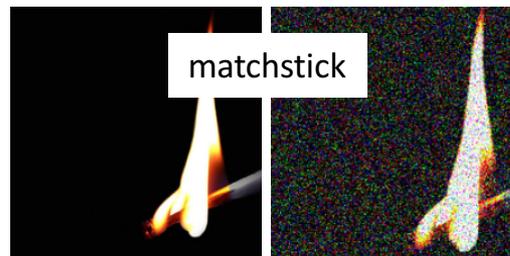
Gradient



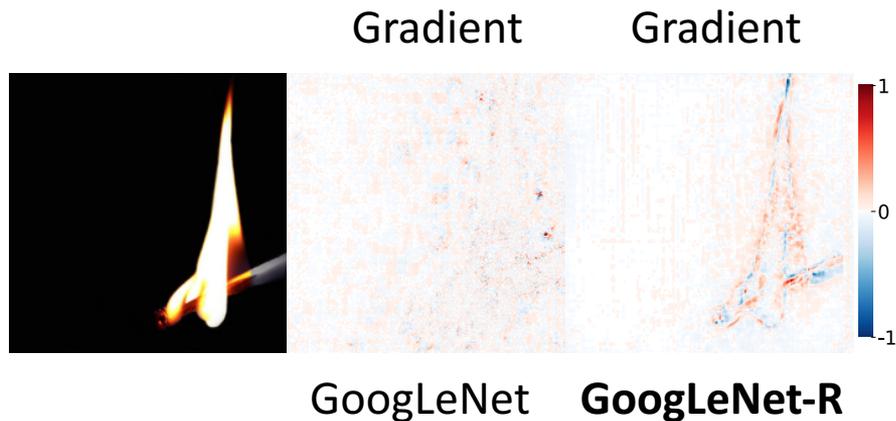
GoogLeNet

**GoogLeNet-R**

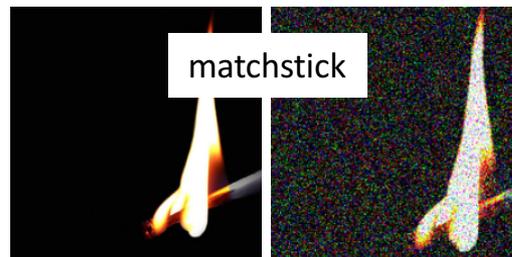
*A robust classifier i.e.  
adversarially trained with noisy images*



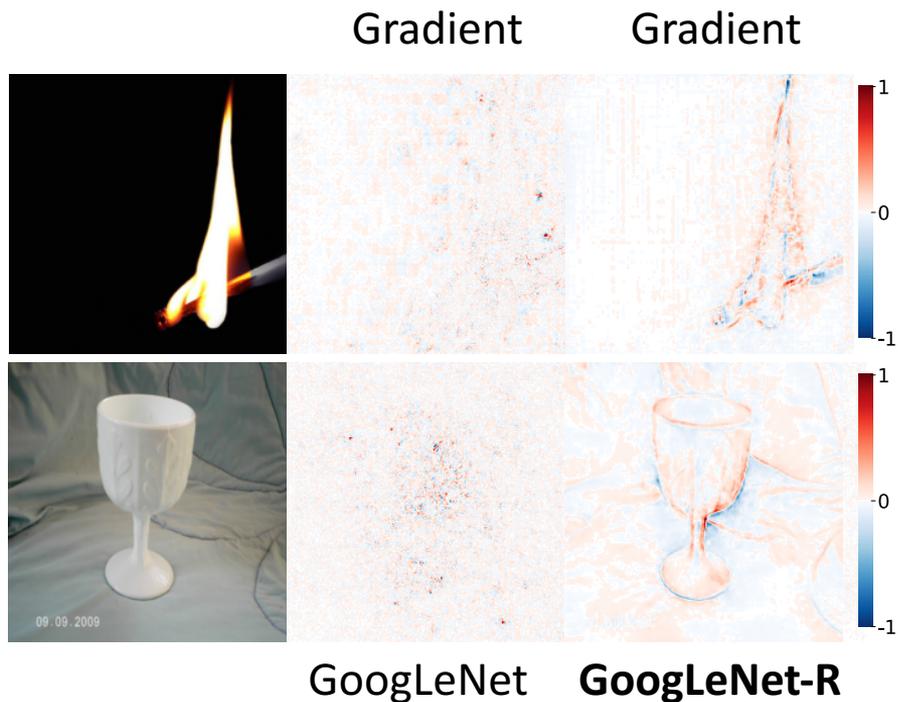
# #1: Saliency maps may NOT be too noisy!



*A robust classifier* i.e.  
adversarially trained with noisy images

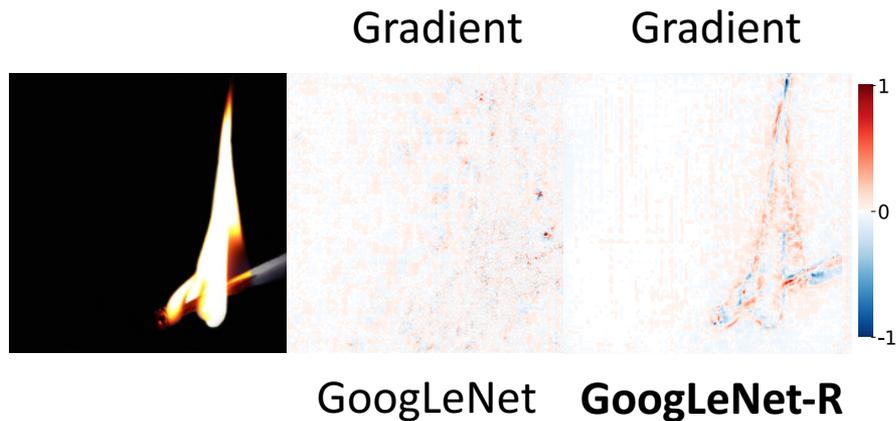


# #1: Saliency maps may NOT be too noisy!



*A robust classifier i.e.  
adversarially trained with noisy images*

# #1: Saliency maps may NOT be too noisy!



*A robust classifier* i.e.  
adversarially trained with noisy images

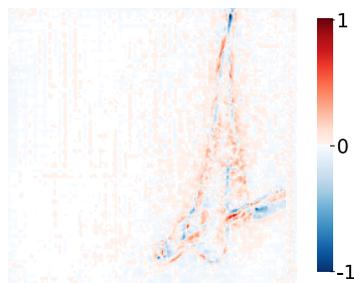
# #1: Saliency maps may NOT be too noisy!

Gradient



GoogLeNet

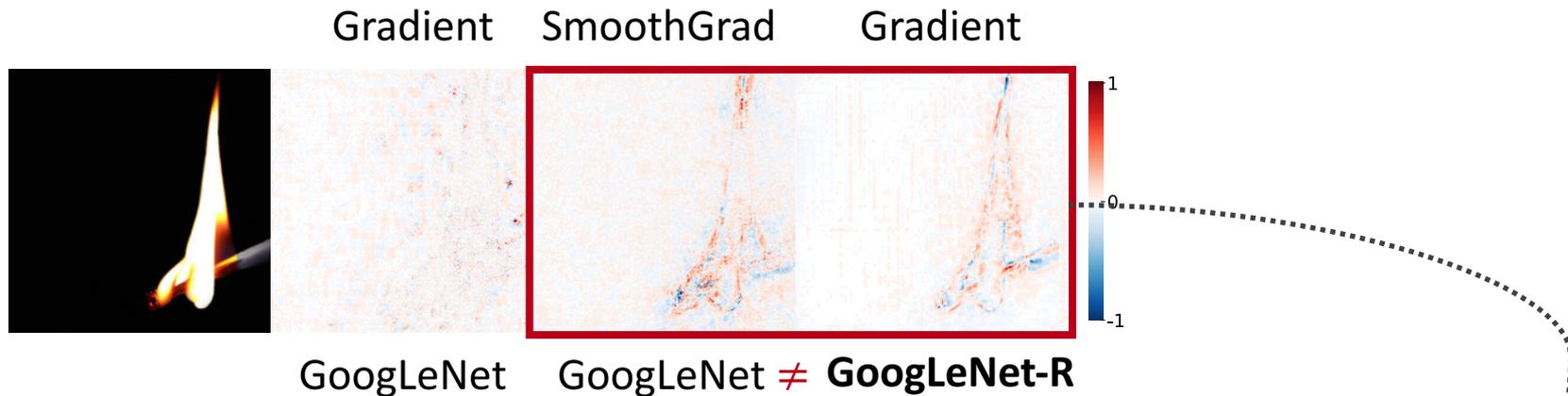
Gradient



GoogLeNet-R

.....  
*A robust classifier i.e.*  
adversarially trained with noisy images

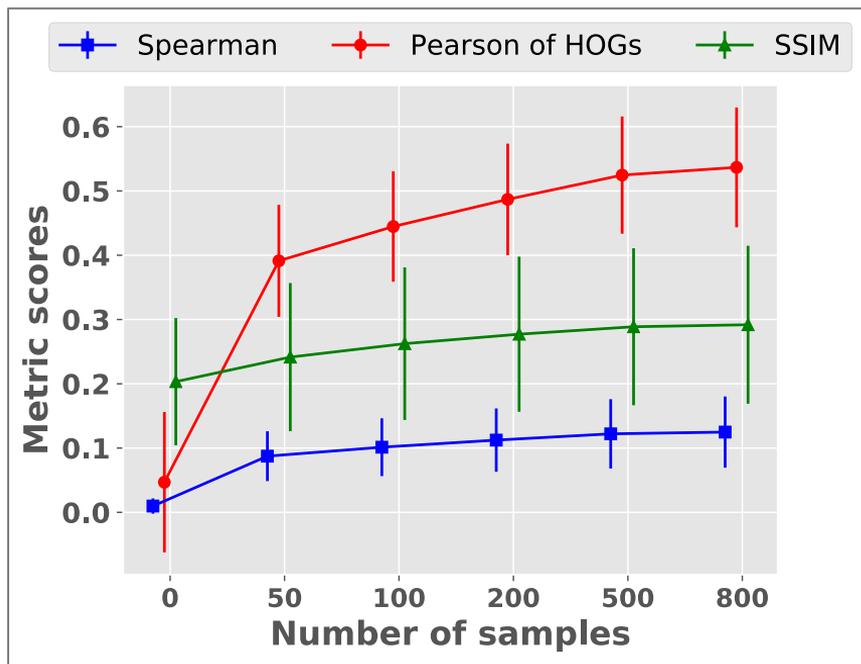
## #2: Smoothed gradients can be misinterpreted



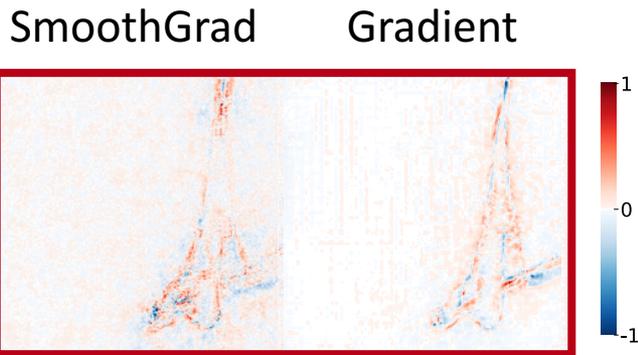
The **flame** is the most “important”

Some background pixels are the most “important”

# #2: Smoothed gradients can be misinterpreted



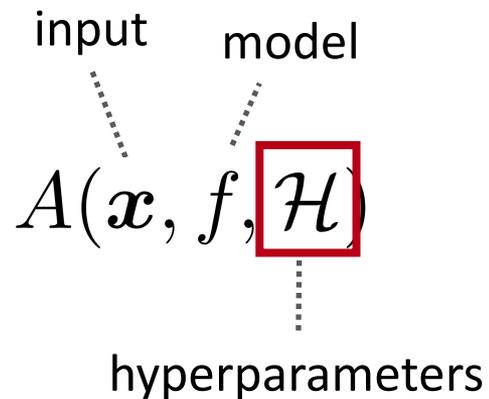
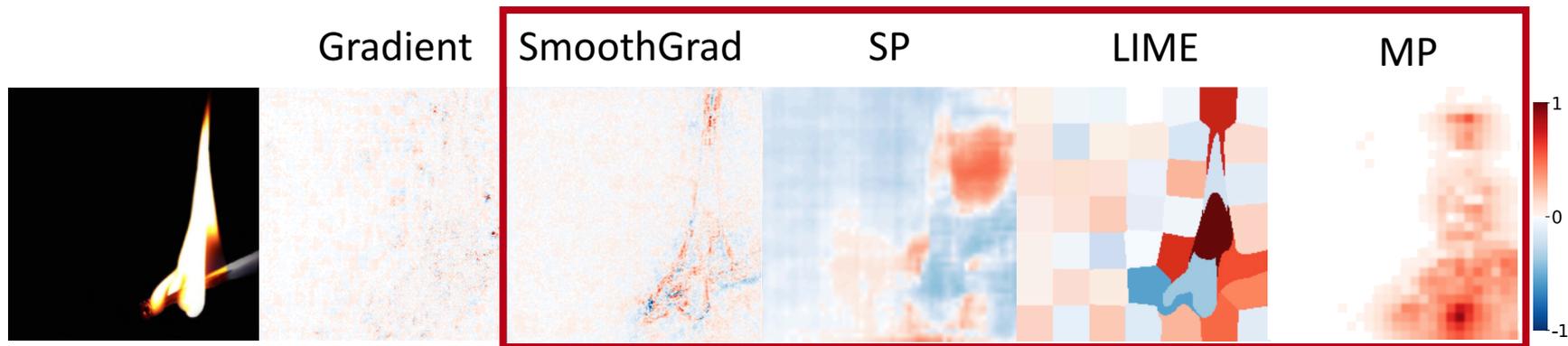
Bansal\*, Agarwal\*, Nguyen\*, 2020



GoogLeNet  $\neq$  GoogLeNet-R

*A robust classifier i.e.  
adversarially trained with noisy images*

# #3: Many attribution maps are sensitive to hyperparams



# #3: Many attribution maps are sensitive to hyperparams

Gradient

SmoothGrad

SP

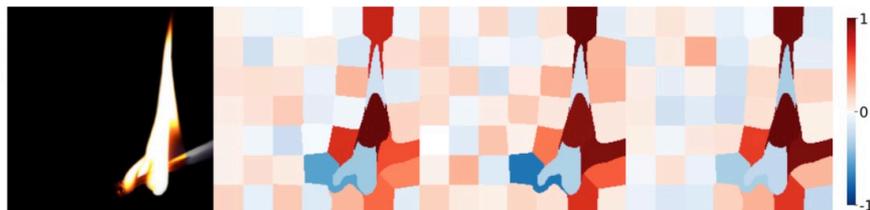
LIME

MP

Input image

Attribution maps

LIME [42]



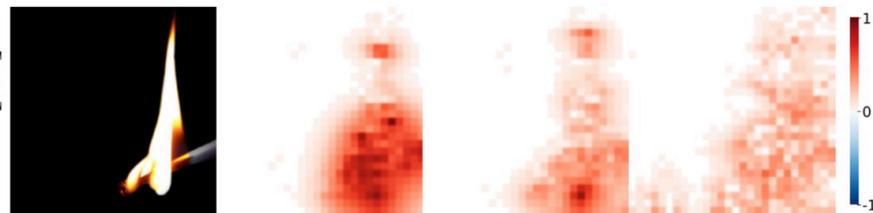
Random seed:

0

1

2

MP [22]



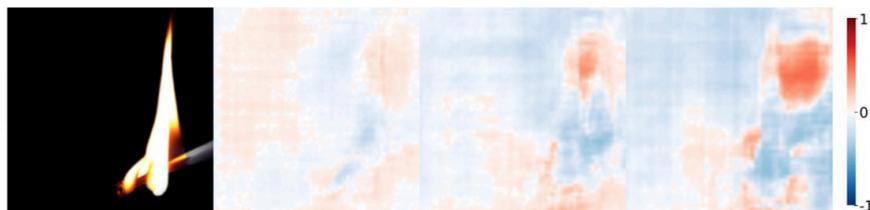
Blur radius:

5

10

30

SP [60]



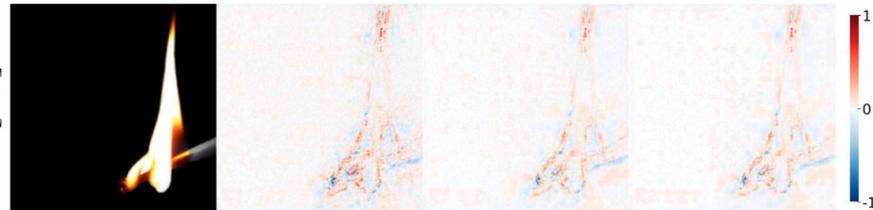
Patch size:

$5 \times 5$

$29 \times 29$

$53 \times 53$

SG [48]



Sample size:

50

200

800

hyperparameters

# #3: Many attribution maps are sensitive to hyperparams

Science

Contents ▾

News ▾

Careers ▾

Journals ▾

SHARE

IN DEPTH | COMPUTER SCIENCE



## Artificial intelligence faces reproducibility crisis

Matthew Hutson

+ See all authors and affiliations

Science 16 Feb 2018:  
Vol. 359, Issue 6377, pp. 725-726  
DOI: 10.1126/science.359.6377.725

Article

Figures & Data

Info & Metrics

eLetters

PDF

### Summary

The booming field of artificial intelligence (AI) is grappling with a replication crisis, much like the



ARTICLE TOOLS

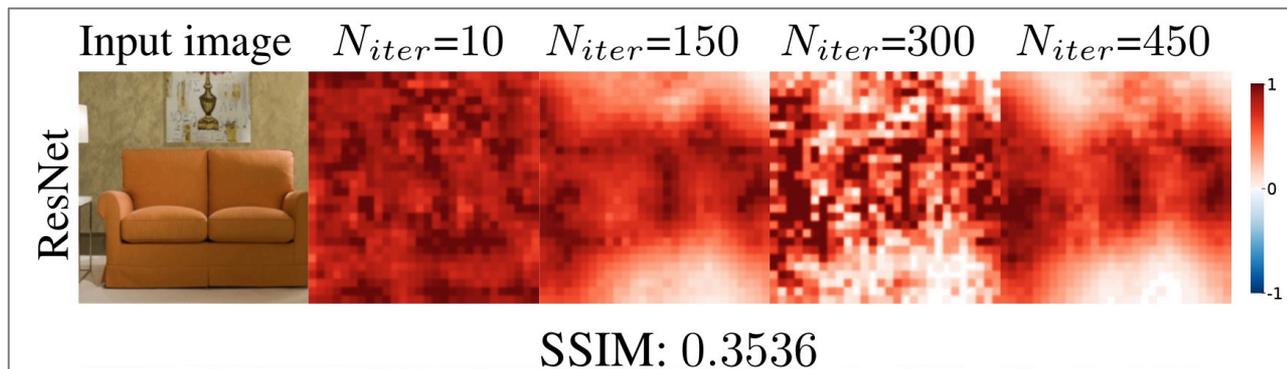
Email

Print

Request Permission

Citation tools

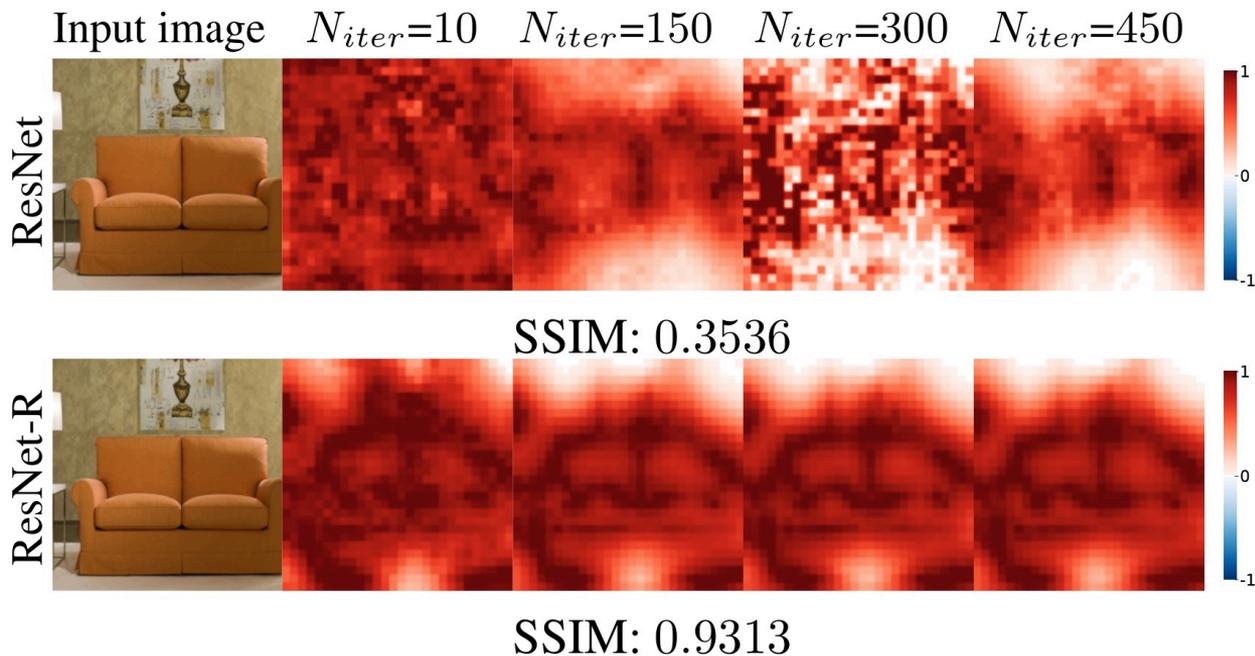
## #4: Attribution maps are more robust under robust classifiers



**Idea:** Identify a minimal region s.t. when blurred out would minimize classification score

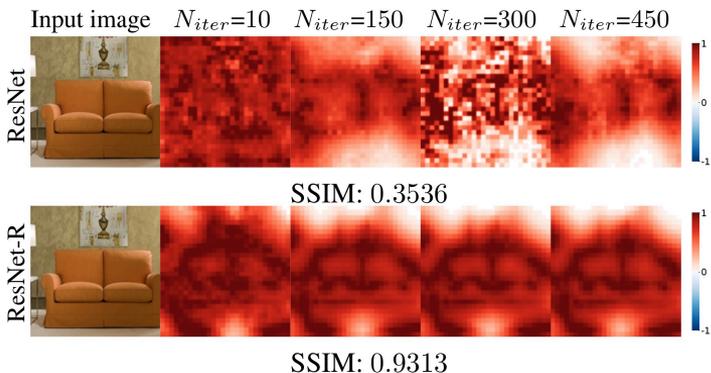
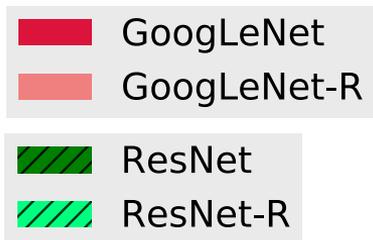
$$\mathbf{m}^* = \arg \min_{\mathbf{m}} \lambda \|\mathbf{m}\|_1 + f(\text{blur}(\mathbf{x}, \mathbf{m}))$$

## #4: Attribution maps are more robust under robust classifiers

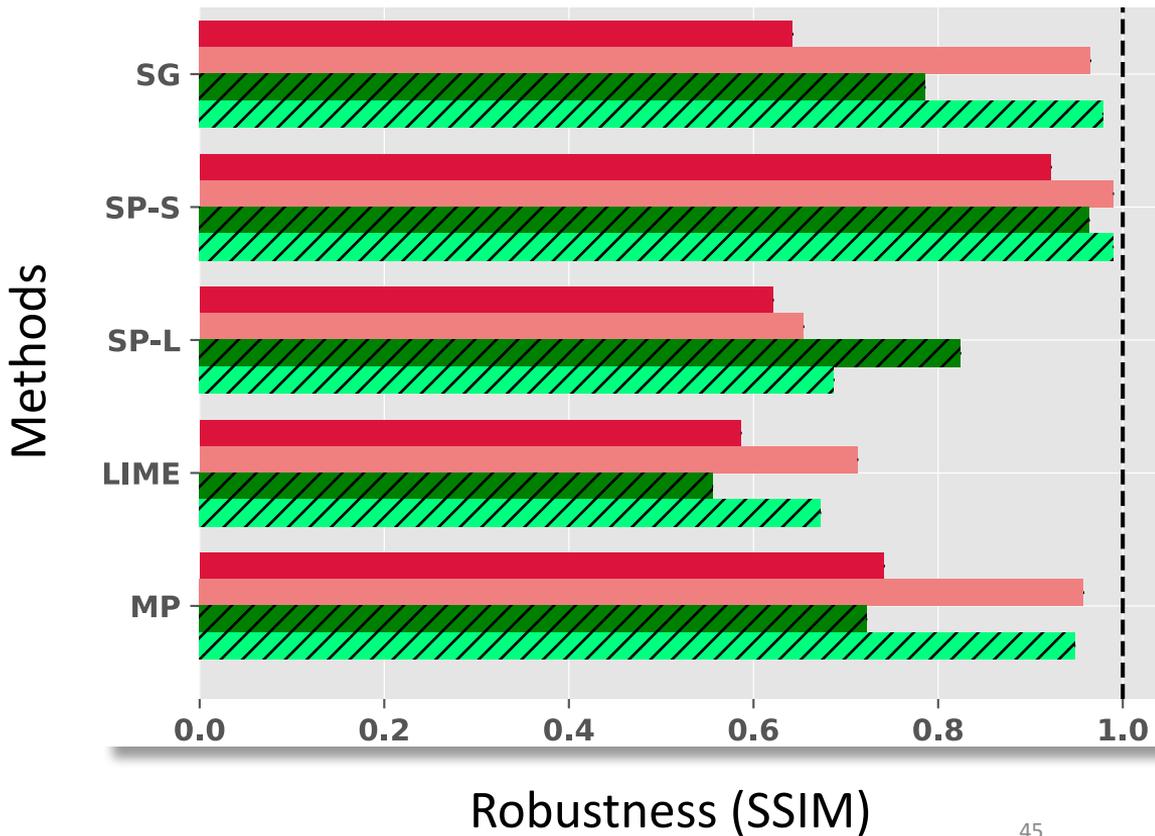


(b) Sensitivity to changes in the number of iterations  $N_{iter}$

# #4: Attribution maps are more robust under robust classifiers



Bansal\*, Agarwal\*, Nguyen\*, 2020



# Conclusions

~~How-to~~

## How to explain neural networks?



C Agarwal



N Bansal



M Alcorn

1. Saliency maps for robust classifiers are not as noisy
2. Smoothed gradients can be misinterpreted
3. Many attribution methods are sensitive to hyper-parameters
4. For robust classifiers, attribution maps are more robust

Papers & code: <http://AnhNguyen.me>

< hiring PhD students! >

Work funded by

