# Interpretable Deep Learning

**EE807: Recent Advances in Deep Learning**

**Lecture 15**

**Slide made by**

**Jun Hyun Nam**

**KAIST EE**

## Table of Contents

1. **Introduction**
   - Why interpretability?
   - What is interpretability?
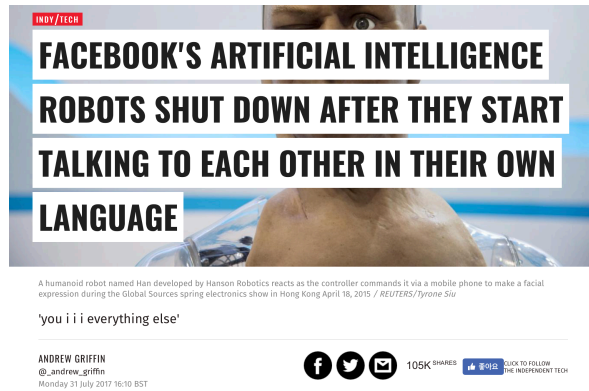   - Overview

2. **Visual Explanation**
   - Perturbation-based methods
   - Gradient-based methods

3. **Other Approaches**
   - Visualize features
   - Network dissection
   - Influence function

## Table of Contents

**1. Introduction**
- Why interpretability?
- What is interpretability?
- Overview

**2. Visual Explanation**
- Perturbation-based methods
- Gradient-based methods

**3. Other Approaches**
- Visualize features
- Network dissection
- Influence function

- Recently, deep learning shows superior performance in various tasks

- However, we don't know yet why they work so well



INDY/TECH

**FACEBOOK'S ARTIFICIAL INTELLIGENCE ROBOTS SHUT DOWN AFTER THEY START TALKING TO EACH OTHER IN THEIR OWN LANGUAGE**

A humanoid robot named Han developed by Hanson Robotics reacts as the controller commands it via a mobile phone to make a facial expression during the Global Sources spring electronics show in Hong Kong April 18, 2015 / REUTERS/Tyrone Siu

'you i i i everything else'

ANDREW GRIFFIN
@_andrew_griffin
Monday 31 July 2017 16:10 BST
105K SHARES  좋아요  CLICK TO FOLLOW THE INDEPENDENT TECH

**THE ULTIMATE GO CHALLENGE**
GAME 1 OF 5
**9 MARCH 2016**

AlphaGo  vs  Lee Sedol

RESULT: W+Res | NUMBER OF MOVES: 186 | TIME WHITE: 1h 55m | TIME BLACK: 1h 32m

- When it fails, it can cause critical issues



*Self-Driving Tesla Was Involved in Fatal Crash, U.S. Says*

By BILL VLASIC and NEAL E. BOUDETTE   JUNE 30, 2016

A Tesla Model S, with its self-driving mode enabled. In a statement, the National Highway Traffic Safety Administration said it had sent an investigative team to examine the vehicle and the crash site in Williston, Fla. Jasper Juinen/Bloomberg

The 'three black teenagers' search shows it is society, not Google, that is racist
Antoine Allen

Twitter outrage over image search results of black and white teens is misdirected. We must address the prejudice that feeds such negative portrayals

RELATED C

Images thrown up by Kabir Alli's Google searches for 'three black teenagers' and 'three white teenagers'

SCIENCE & THE PUBLIC   SCIENCE & SOCIETY

**Data-driven crime prediction fails to erase human bias**

Poor, minority communities flagged as drug crime trouble spots in case study

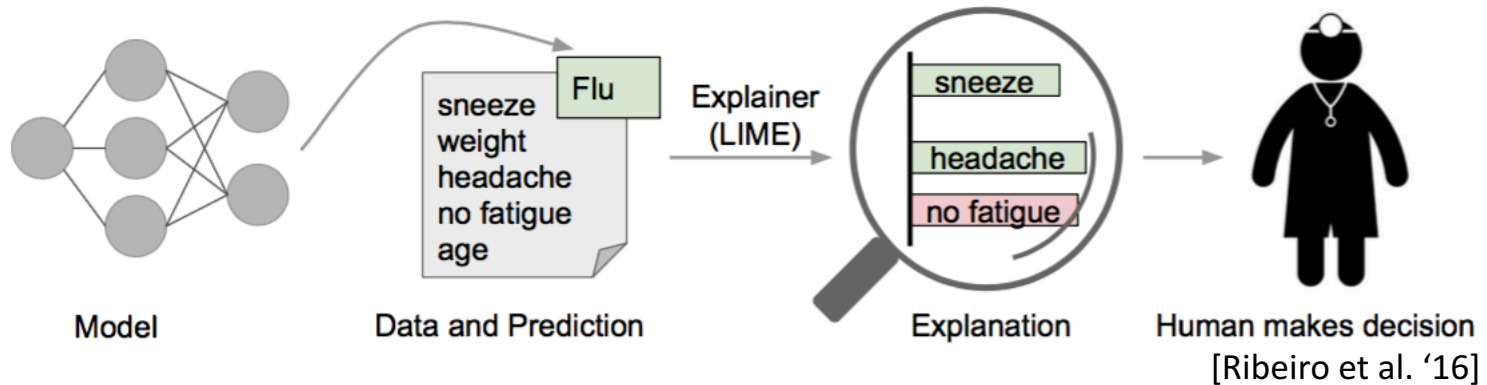BY RACHEL EHRENBERG 10:00AM, MARCH 8, 2017

BIG DATA DOESN'T PAY  Software programs that use police records to predict crime hot spots may result in police unfairly targeting low-income and minority communities, a new study shows.

ARTOLYMPIC/ISTOCKPHOTO

- Interpretation is the process of giving <span style="color:red">explanations</span>



[Ribeiro et al. '16]

- Situations when ML interpretation can be helped
  - **Safety**: We want to make sure the system is making sound decisions
  - **Debugging**: We want to understand why a system doesn't work
  - **Science**: We want to understand something new
  - **Legal**: We are legally required to provide an explanation
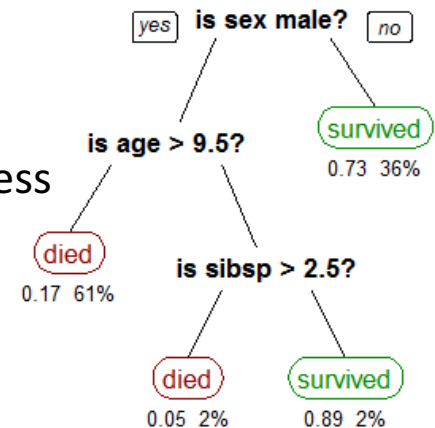  - **Ethics**: We don't want to discriminate against particular groups

- **Linear model**
  - Consider $y = \beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n$
  - **Question**: How much input feature $x_i$ contributed to (or affected) output $y$?
  - Answer: $\beta_i$

- **Decision tree**
  - **Question**: How much 'age' affected probability of survived?
  - Answer: Don't know
  - Instead of per-feature attribution, we know its decision process

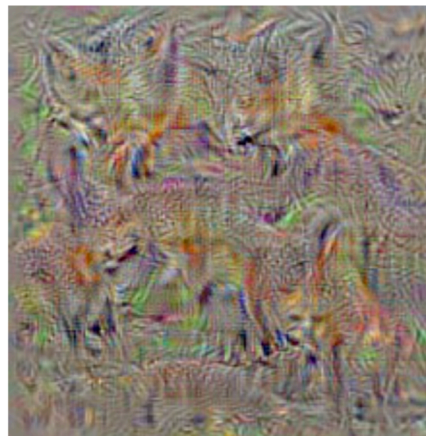- Many interpretable ML approaches provides explanation of the original model in one of two forms

**Interpretable ML method**

or

Created by Knut M. Synstad
from Noun Project

- **Local explanation**
  - Explain a single prediction
  - e.g. which <span style="color:red">part of the image</span> affected the prediction most (visual explanation)
  - e.g. find <span style="color:red">a training data</span> most responsible to the prediction (influence function)

- **Global explanation**
  - Describe the entire model behavior
  - e.g. generate <span style="color:red">a synthetic image</span> that maximizes certain output (feature visualization)
  - e.g. discover <span style="color:red">a human-friendly concept</span> related to each neuron (network dissection)
  - e.g. find <span style="color:red">a training data</span> most responsible to the model (influence function)



(c) Grad-CAM 'Cat'



**kit fox**



Train
res5c unit 924          IoU=0.293
res5c unit 2001         IoU=0.255
inception_5b unit 626   IoU=0.145
inception_5b unit 415   IoU=0.143

## Table of Contents

1. **Introduction**
   - Why interpretability?
   - What is interpretability?
   - Overview

2. **Visual Explanation**
   - Perturbation-based methods
   - Gradient-based methods

3. **Other Approaches**
   - Visualize features
   - Network dissection
   - Influence function

- **Idea**: Mask part of the image with gray patch before feeding to CNN, and
  check how much the prediction changes



P(elephant) = 0.95

P(elephant) = 0.75

African elephant, Loxodonta africana

schooner

- **Problem**: Removing information with gray patch is too heuristic

- **Idea**: Simulate the absence of a feature by <span style="color:red">marginalizing</span> the feature

- **Goal**: The attribution of i-th feature for given image and $\mathbf{x}$ and class $c$

$$p(c|\mathbf{x}) - p(c|\mathbf{x}_{\backslash i})$$

where $\mathbf{x}_{\backslash i}$ represents the absence of $x_i$ in $\mathbf{x}$

$$p(c|\mathbf{x}_{\backslash i}) = \sum_{x_i} p(x_i|\mathbf{x}_{\backslash i}) p(c|\mathbf{x}_{\backslash i}, x_i)$$

- Note that $p(x_i|\mathbf{x}_{\backslash i})$ is computationally expensive

- Assume $x_i$ is independent of the other features, i.e., $p(x_i|\mathbf{x}_{\backslash i}) \approx p(x_i)$

$$p(c|\mathbf{x}_{\backslash i}) \approx \sum_{x_i} p(x_i) p(c|\mathbf{x}_{\backslash i}, x_i)$$

- The prior probability $p(x_i)$ is usually approximated by the empirical distribution

- **Idea**: Simulate the absence of a feature by marginalizing the feature

$$p(c|\mathbf{x}_{\setminus i}) = \sum_{x_i} p(x_i|\mathbf{x}_{\setminus i})p(c|\mathbf{x}_{\setminus i}, x_i)$$

- **Problem**: $p(x_i|\mathbf{x}_{\setminus i}) \approx p(x_i)$ is a very crude approximation
  - e.g. a pixel's value is highly dependent on other pixels

- **Observations**
  - A pixel depends most strongly on a small neighborhood around it
  - The conditional of a pixel given its neighborhood does not depend on the position

- For a pixel $x_i$ , one can find a patch $\hat{\mathbf{x}}_i$ than contains $x_i$ and $p(x_i|\mathbf{x}_{\setminus i}) \approx p(x_i|\hat{\mathbf{x}}_i)$

# Prediction Difference Analysis [Zintgraf et al., 2017]

- **Results**
  - Marginal vs. conditional sampling



  - Different window sizes

- Remember that a sparse linear model is a good explanation model



(a) Original Image (b) Explaining *Electric guitar*

- **Idea**: Local linear approximation
  - Explain the entire model is hard, but a single prediction is easier
  - Approximate the model in a local region around the single prediction by a linear classifier

- Illustration of the main idea



Original Image
$$x \in \mathbb{R}^d$$

Interpretable Components
$$x' \in \{0, 1\}^{d'}$$

Original Image
P(tree frog)  = 0.54

| Perturbed Instances | P(tree frog) |
| --- | --- |
| | 0.85 |
| | 0.00001 |
| | 0.52 |

Locally weighted regression

Explanation

- Overall Procedure
    1. Decompose original input to interpretable representation
    2. Model local region around given input by sampling
    3. Approximate original model as a linear classifier

- Illustration of the main idea



- **Step 1**: Interpretable representation
  - Understandable to humans
  - For text classification, a binary vector indicating the presence or absence of a word
  - For image classification, a binary vector indicating the presence or absence of a contiguous patch of similar pixels
  - $x \in \mathbb{R}^d$ : original representation / $x' \in \{0,1\}^{d'}$ : its interpretable representation
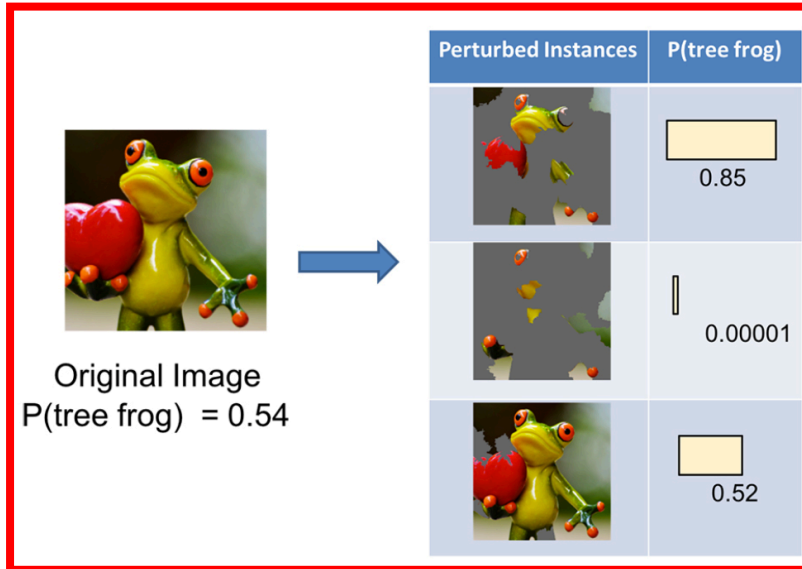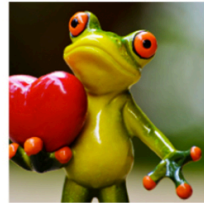
- Illustration of the main idea



Original Image
$x \in \mathbb{R}^d$

Interpretable Components
$x' \in \{0,1\}^{d'}$

Original Image
P(tree frog) = 0.54

| Perturbed Instances | P(tree frog) |
|---|---|
| | 0.85 |
| | 0.00001 |
| | 0.52 |

Locally weighted regression

Explanation

- **Step 2**: Model local region around given input
  - Sample instances around $x$ by drawing nonzero elements of $x' \in \{0,1\}^{d'}$ uniformly at random
  - Given a perturbed sample $z' \in \{0,1\}^{d'}$, recover the original representation $z \in \mathbb{R}^d$
  - Compute $f(z)$: the prediction of model for each perturbed output

- Illustration of the main idea



Original Image
$x \in \mathbb{R}^d$

Interpretable Components
$x' \in \{0, 1\}^{d'}$

Original Image
P(tree frog) = 0.54

| Perturbed Instances | P(tree frog) |
|---|---|
| | 0.85 |
| | 0.00001 |
| | 0.52 |

Locally weighted regression

Query

Explanation

- **Step 3**: Approximate original model as a linear classifier
  - Fit a linear classifier $g(z') = w_g \cdot z'$ and use it as an explanation model

$$\mathcal{L}(f, g, \Pi_x) = \sum_{z, z' \in \mathcal{Z}} \Pi_x(z)(f(z) - g(z'))^2$$

  - $\Pi_x(z)$ defines locality (e.g. $\Pi_x(z) = \exp(-\|x - z\|_2^2 / 0.1)$)
  - Final objective

$$\xi(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \Pi_x) + \Omega(g)$$

local fidelity          measure of complexity

- **Results**: Can be applied to any model
    - Top 3 predictions of Inception-v3 for ImageNet dataset
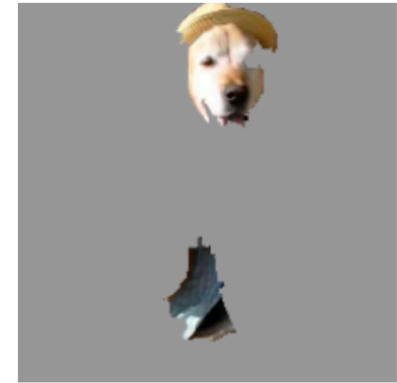


(a) Original Image     (b) Explaining *Electric guitar*     (c) Explaining *Acoustic guitar*     (d) Explaining *Labrador*

- Random forest prediction for the 20 newsgroups dataset

## Table of Contents

1. **Introduction**
   - Why interpretability?
   - What is interpretability?
   - Overview

2. **Visual Explanation**
   - Perturbation-based methods
   - Gradient-based methods

3. **Other Approaches**
   - Visualize features
   - Network dissection
   - Influence function

- **Problem**: Perturbation-based methods are too slow

- **Idea:** Use gradient of output with respect to the input as the attribution

- **Goal**: Find the influence on the score $S_c(I_0)$ for given image $I_0$
  - Consider the linear score model for class $c$
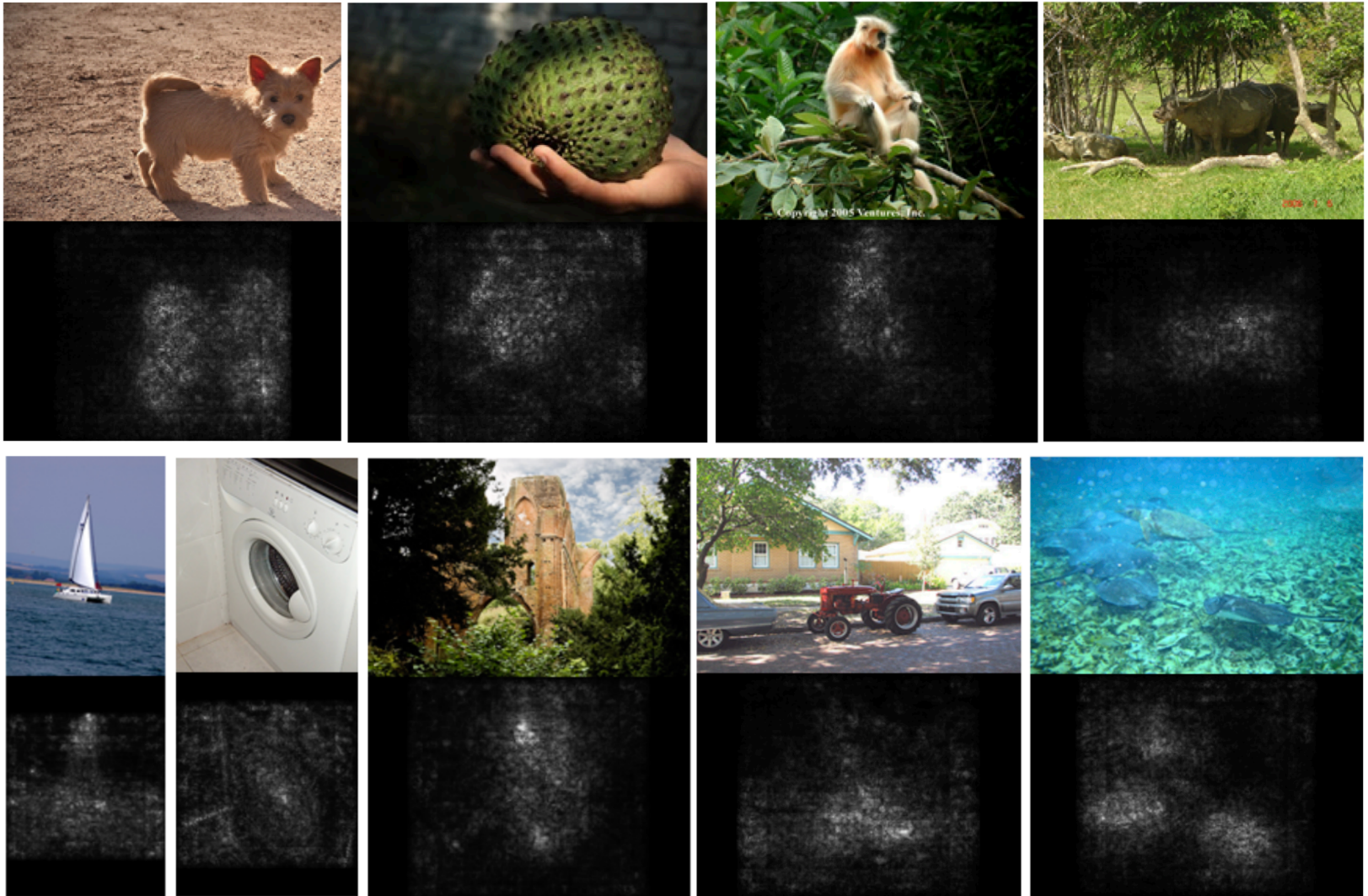
$$S_c(I) = w_c^\top I + b_c$$

  where $I$ : image, $w_c, b_c$ : the weight vector and the bias of the model
  - $w_c$ defines the importance of the corresponding pixels of $I$ for the class $c$

  - In case of non-linear/complex models, approximate $S_c(I)$
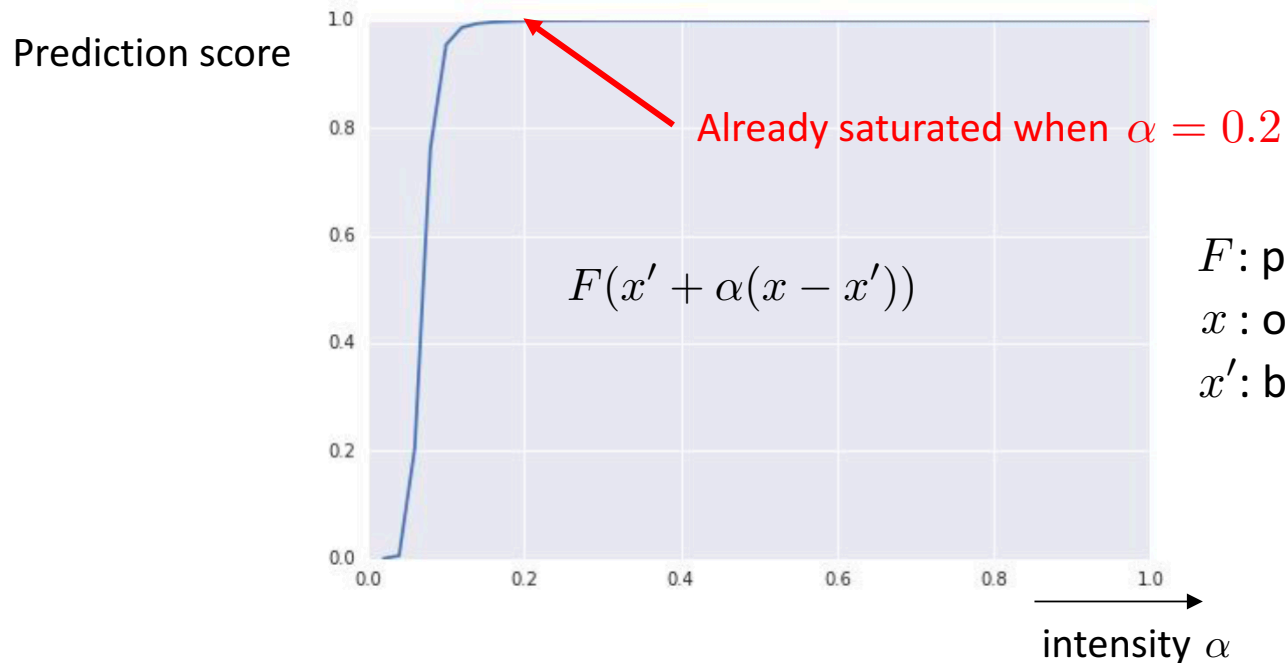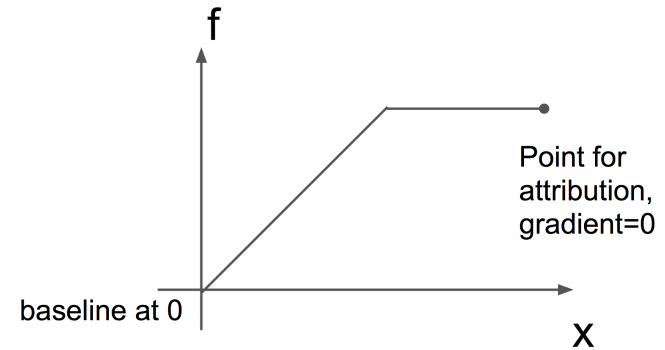    by the first-order Taylor expansion

$$S_c(I) \approx w^\top I + b$$

  where $w = \left. \dfrac{\partial S_c}{\partial I} \right|_{I=I_0}$

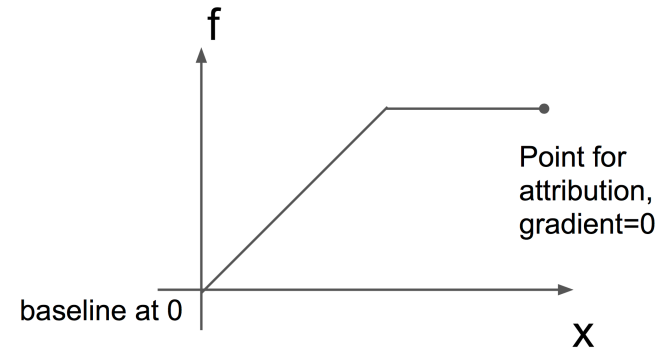- **Results:** Without any additional annotation, gradient can localize the object
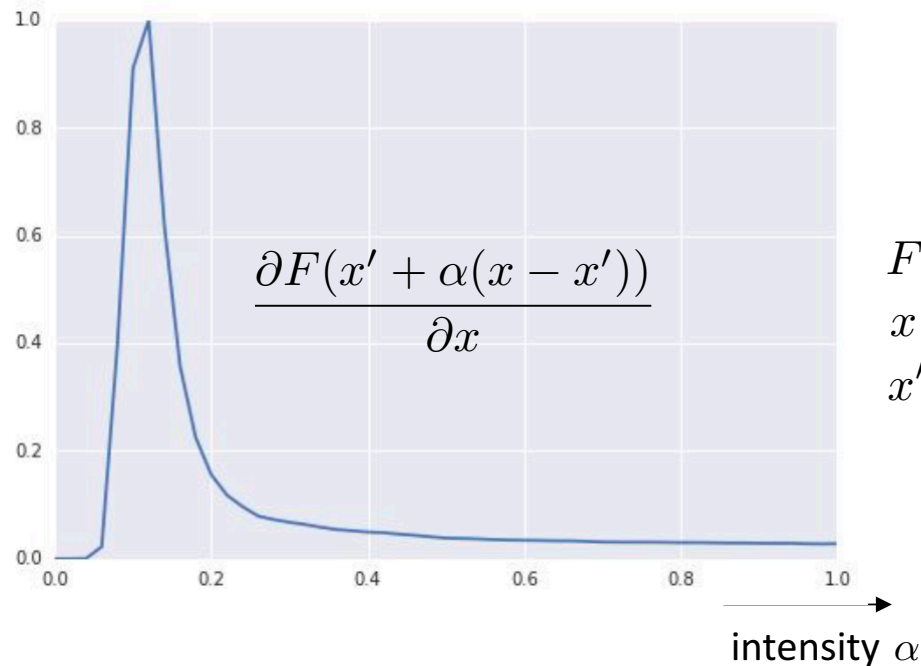
- **Problem**: Prediction score might saturate

  - For high confidence prediction, small perturbation in input does not change the prediction value



Prediction score

$$F(x' + \alpha(x - x'))$$

Already saturated when $\alpha = 0.2$

intensity $\alpha$

$F$: prediction score
$x$ : original image
$x'$: baseline image

- **Problem**: Prediction score might <span style="color:red">saturate</span>

  - For high confidence prediction, small perturbation in input does not change the prediction value

f

Point for attribution, gradient=0

baseline at 0

x

Average pixel gradient (normalized)

$$\frac{\partial F(x' + \alpha(x - x'))}{\partial x}$$

$F$ : prediction score
$x$ : original image
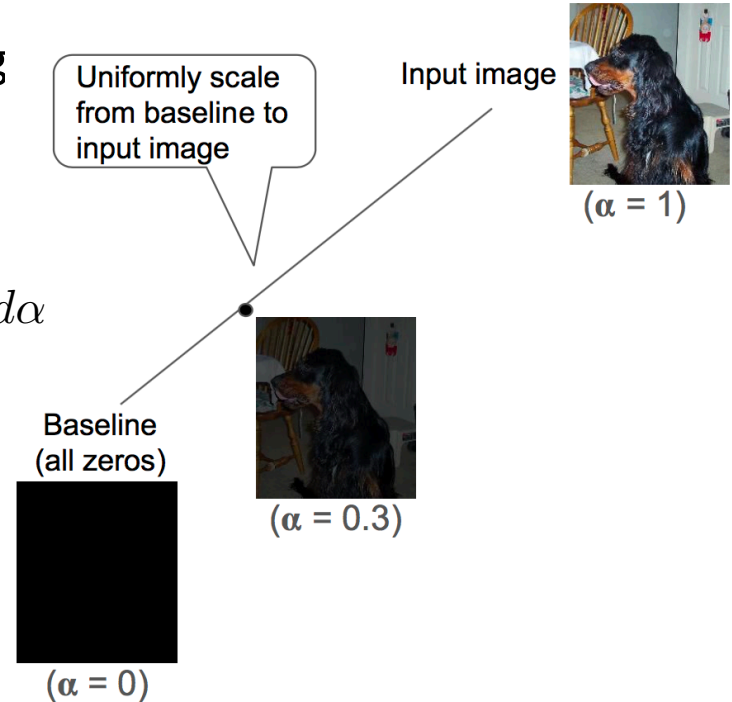$x'$ : baseline image

intensity $\alpha$

- **Idea**: Compute all the gradients for images from baseline to actual image

- Construct a sequence of images interpolating from a baseline (black) to the actual image

- Average the gradients across these images

$$\text{IG}_i(x) = (x_i - x_i') \times \int_{\alpha=0}^{1} \frac{\partial F(x' + \alpha(x - x'))}{\partial x_i} d\alpha$$

  - $F$ is the prediction function for the label
  - $x_i$ is the intensity of ith pixel
  - $\text{IG}_i(x)$ is the integrated gradient w.r.t. the ith pixel

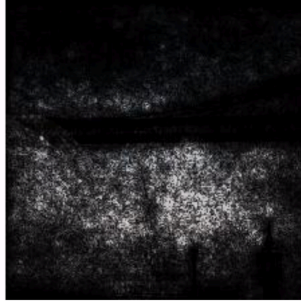Uniformly scale from baseline to input image

Input image

$(\alpha = 1)$

Baseline (all zeros)

$(\alpha = 0.3)$

$(\alpha = 0)$

- Properties
  - **Sensitivity**: A variable changes output, then the variable should get an attribution
  - **Insensitivity**: A variable has no effect on the output gets no attribution
  - **Completeness**: $\sum_{i=1}^{n} \text{IG}_i(x) = F(x) - F(x')$

- **Results:** For high confidence predictions,
  integrated gradients provide discriminative region

- **Problem**: Gradients strongly fluctuate!
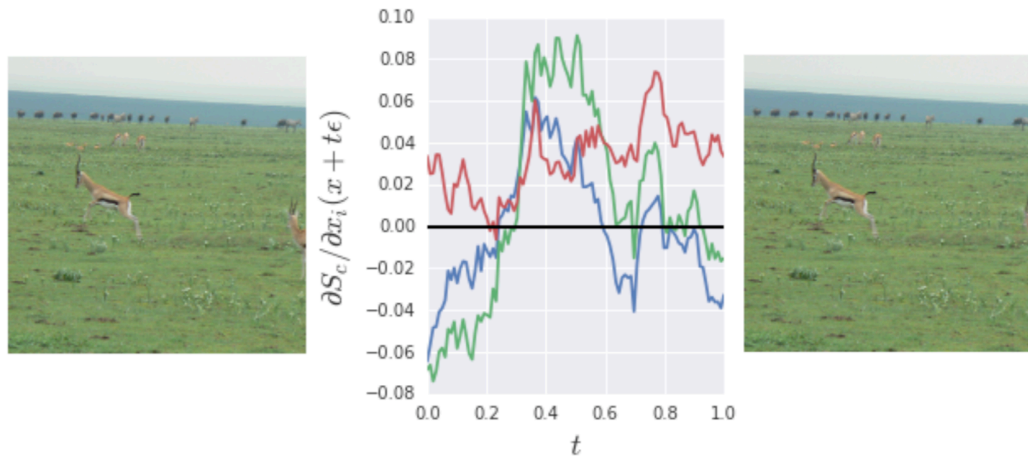  - Given image $x$, and an image pixel $x_i$, plots values of $\max_i \dfrac{\partial S_c}{\partial x_i}(x + t\epsilon)$ for a short line segment $x + t\epsilon$



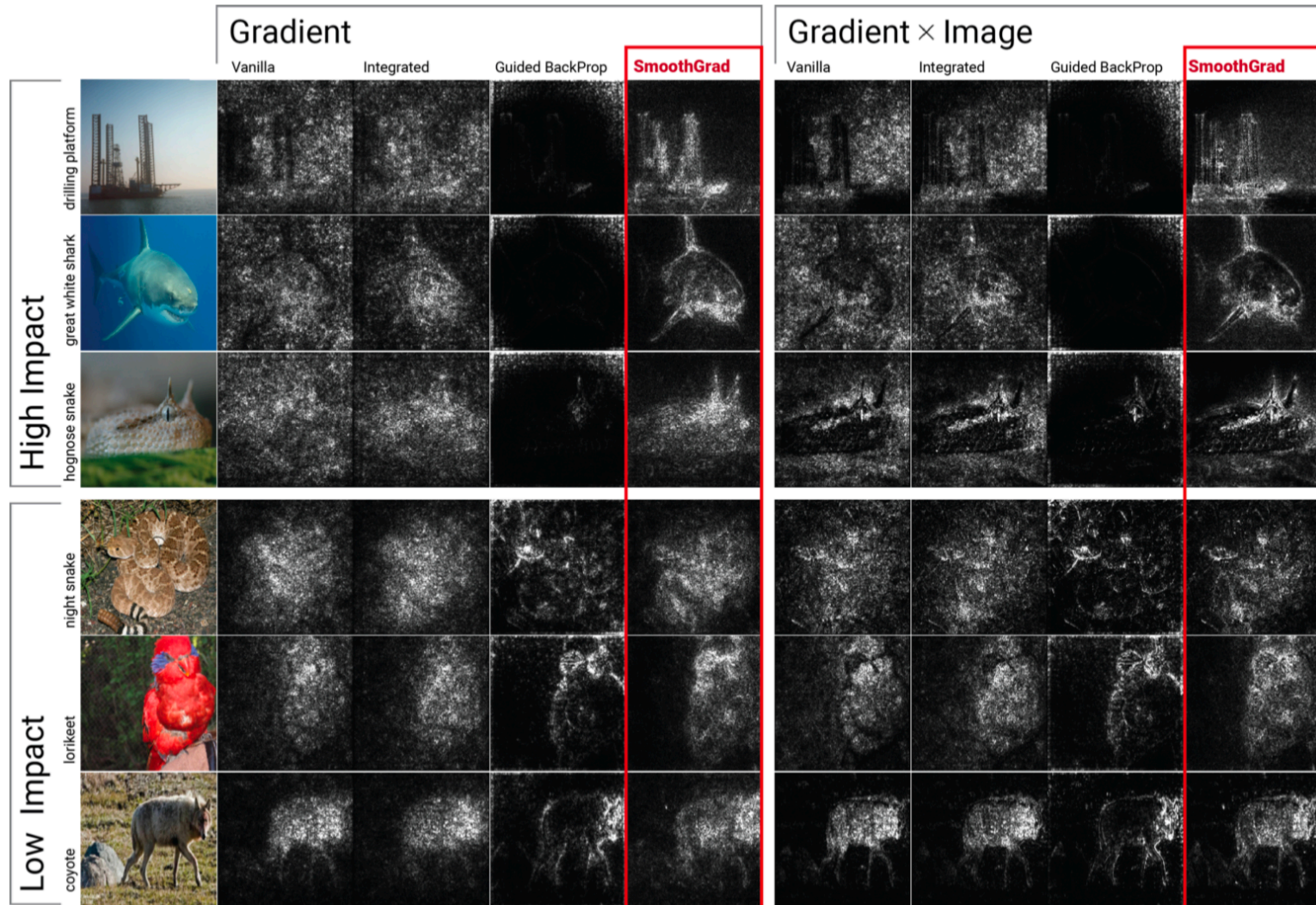  - Even $x$ and $x + \epsilon$ are indistinguishable, the partial derivative rapidly fluctuate

- **Idea**: Use a local average of gradient values

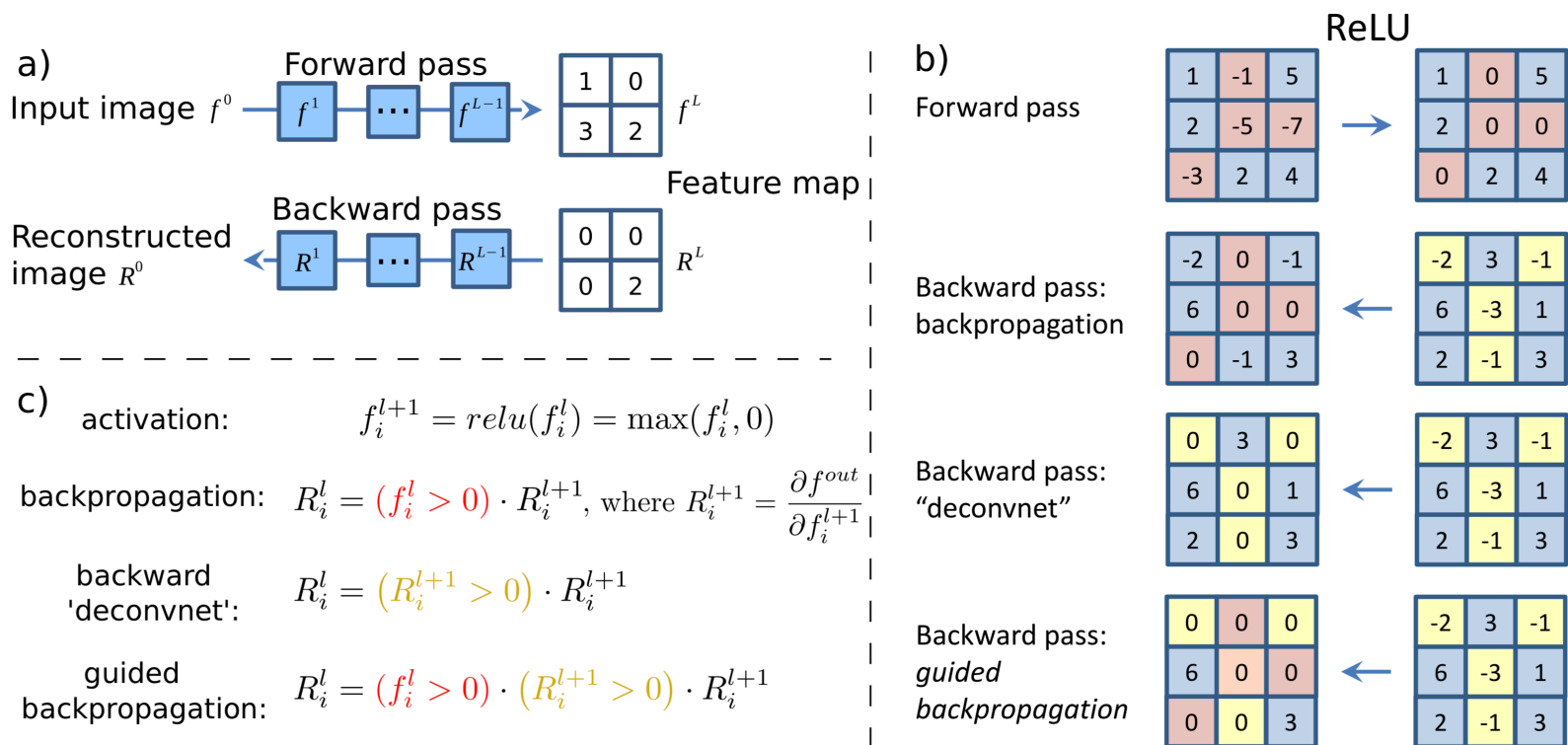$$\mathrm{SG}(x) = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial S_c}{\partial x}(x + g_i)$$

where noise vectors $g_i \sim \mathcal{N}(0, \sigma^2)$ are drawn i.i.d. from a normal distribution

- **Results:** Simple noise-adding method can dramatically improve the quality of saliency map

# Other Backpropagation Variants

- **Deconvolution** [Zeiler et al., 2014]
  - Reverse operation of convolution

- **Guided Backpropagation** [Springenberg et al., 2015]
  - Backpropagate only positive gradients through each ReLU

- Both methods visualize the activations of high layer neurons (also the prediction)



a)

Forward pass

Input image $f^0$ → $f^1$ → ··· → $f^{L-1}$ → 

| 1 | 0 |
|---|---|
| 3 | 2 |

$f^L$

Feature map

Backward pass

Reconstructed image $R^0$ ← $R^1$ ← ··· ← $R^{L-1}$ ← 

| 0 | 0 |
|---|---|
| 0 | 2 |

$R^L$

c)

activation: $f_i^{l+1} = relu(f_i^l) = \max(f_i^l, 0)$

backpropagation: $R_i^l = (f_i^l > 0) \cdot R_i^{l+1}$, where $R_i^{l+1} = \frac{\partial f^{out}}{\partial f_i^{l+1}}$

backward 'deconvnet': $R_i^l = (R_i^{l+1} > 0) \cdot R_i^{l+1}$

guided backpropagation: $R_i^l = (f_i^l > 0) \cdot (R_i^{l+1} > 0) \cdot R_i^{l+1}$

b)

ReLU

Forward pass

Backward pass: backpropagation

Backward pass: "deconvnet"

Backward pass: *guided backpropagation*

- **Problem**: Many pixel-level attribution methods insensitive to model parameter [Adebayo et al., 2018]
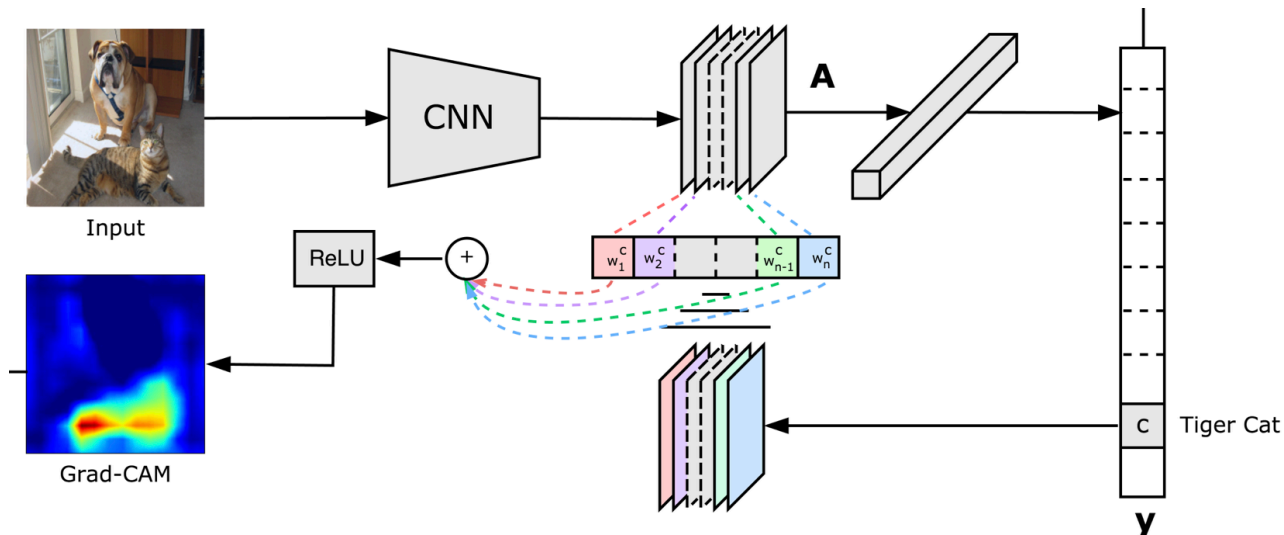
- **Idea**: Activation-level attribution instead of pixel-level attribution

- Gradient-based extension of CAM [Zhou et al., 2015]

- Can be applied to any CNN based model
  - Image classification, image captioning or visual question answering

- Use GAP of gradients instead of weights after GAP layer
  - $y^c$ : the score for class $c$, $A^k$ : feature map of the last convolutional layer

$$\alpha_k^c = \sum_{i,j} \frac{\partial y^c}{\partial A_{ij}^k} \qquad L_{\mathrm{Grad-CAM}}^c = \mathrm{ReLU}\left(\sum_k \alpha_k^c A^k\right)$$

- **Idea**: Activation-level attribution instead of pixel-level attribution

- Gradient-based extension of CAM [Zhou et al., 2015]

- Can be applied to any CNN based model
  - Image classification, image captioning or visual question answering

- Use GAP of gradients instead of weights after GAP layer
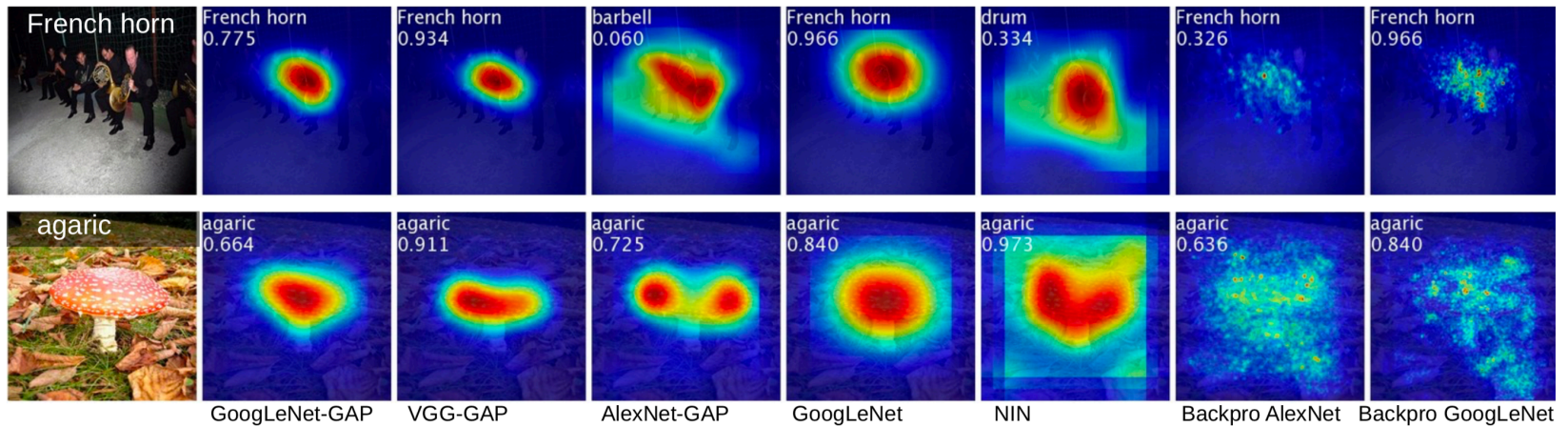  - $y^c$: the score for class $c$, $A^k$: feature map of the last convolutional layer

$$\alpha_k^c = \sum_{i,j} \frac{\partial y^c}{\partial A_{ij}^k} \qquad L_{\text{Grad-CAM}}^c = \text{ReLU}\left(\sum_k \alpha_k^c A^k\right)$$

- Typically, the conv activation has low-resolution → low resolution explanation

- Less affected by CNN architecture prior → more sensitive to model parameter
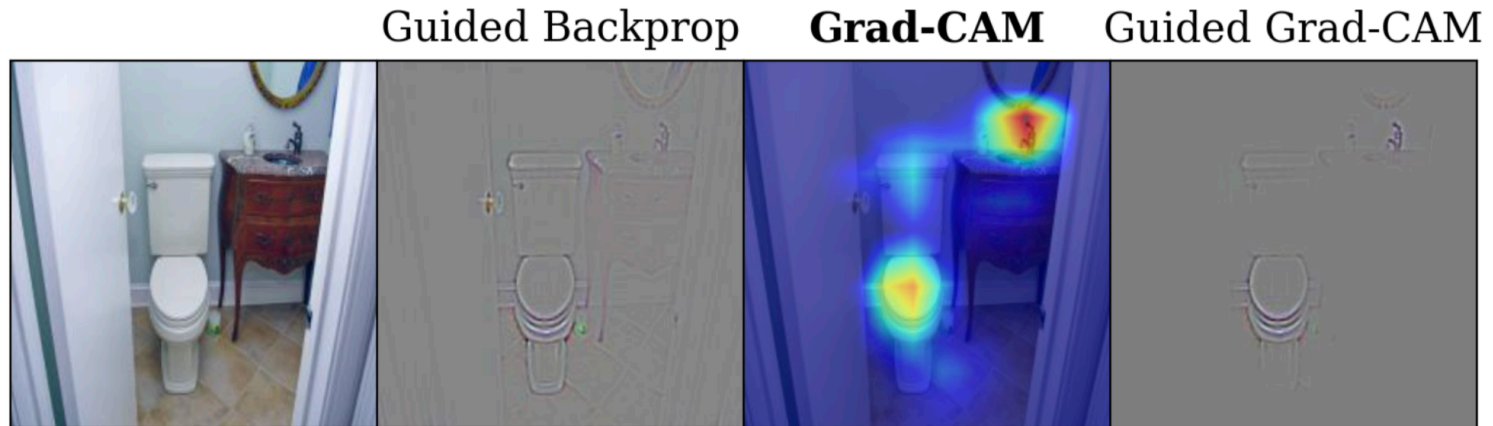
- **Results**
  - CAM vs. Saliency map



- Examples of localization (green: ground truth / red: predicted)

- **Results**: focus on right place without any attention module
  - Visual explanations for captioning



Guided Backprop     **Grad-CAM**     Guided Grad-CAM

A bathroom with a toilet and a sink

A horse is standing in a field with a fence in the background

- **Results**: can discriminate different objects
  - Visual explanations for VQA

What animal is in this picture? (left) Answer: dog / (right) Answer: cat



What color is the hydrant? (left) Answer: yellow / (right) Answer: green

## Table of Contents

1. **Introduction**
   - Why interpretability?
   - What is interpretability?
   - Overview

2. **Visual Explanation**
   - Perturbation-based methods
   - Gradient-based methods

3. **Other Approaches**
   - Visualize features
   - Network dissection
   - Influence function

- **Goal**: Generate a synthetic image that <span style="color:red">maximally activates a neuron</span>
  - So far, we have focused on finding which part of an input that a neuron (or output) responds to
  - Can observe the models behavior when classify image to certain class

- **Idea**: Solve the following optimization

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

  - Initialized image to zeros
  - Forward image to compute current class scores
  - Backprop to get gradient of neuron value w.r.t. image pixels
  - Make a small update to the image

- **Results**: Different aspects of class appearance are captured



| dumbbell | cup | dalmatian | goose | ostrich |

- **Goal**: Interpreting deep visual representation and <span style="color:red">quantifying their interpretability</span>

- **Idea**: Network dissection
  1. Identify a broad set of human-labeled visual concepts
  2. Gather hidden variables' response to known concepts
  3. Quantify alignment of hidden variable – concept pairs

- **Step 1**: Use the broadly and densely labeled (Broden) dataset
  - Gather images from various dataset
  - Total 63,305 pixel-level annotated images, 1,197 visual concepts

Table 1. Statistics of each label type included in the data set.

| Category | Classes | Sources | Avg sample |
|----------|---------|---------|------------|
| scene | 468 | ADE [43] | 38 |
| object | 584 | ADE [43], Pascal-Context [19] | 491 |
| part | 234 | ADE [43], Pascal-Part [6] | 854 |
| material | 32 | OpenSurfaces [4] | 1,703 |
| texture | 47 | DTD [7] | 140 |
| color | 11 | Generated | 59,250 |

street (scene)   flower (object)   headboard (part)

swirly (texture)   pink (color)   metal (material)

- **Step 2**: Gather hidden variables' response
  - For every input image $\mathbf{x}$ in the Broden dataset,
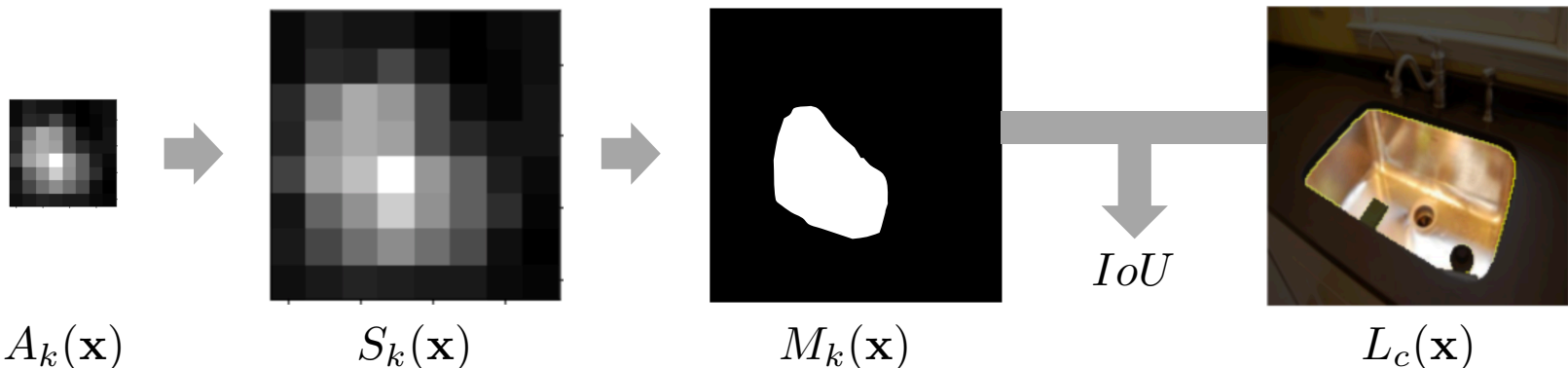        collect the activation map $A_k(\mathbf{x})$ of every convolutional unit $k$
  - Define the binary segmentation $M_k(\mathbf{x}) = \mathbf{1}\{S_k(\mathbf{x}) \geq T_k\}$
  - $S_k(\mathbf{x})$ : scaled up activation map of $A_k(\mathbf{x})$ (same size as the image)
  - $T_k$ : some threshold value

- **Step 3**: Scoring unit interpretability
  - The score of unit $k$ for concept $c$ is reported as a <span style="color:red">dataset-wide IoU score</span>

$$IoU_{k,c} = \frac{\sum_{\mathbf{x}} |M_k(\mathbf{x}) \cap L_c(\mathbf{x})|}{\sum_{\mathbf{x}} |M_k(\mathbf{x}) \cup L_c(\mathbf{x})|}$$

  - $L_c(\mathbf{x})$ : ground truth mask of image $\mathbf{x}$ for concept $c$



$A_k(\mathbf{x})$      $S_k(\mathbf{x})$      $M_k(\mathbf{x})$      $IoU$      $L_c(\mathbf{x})$
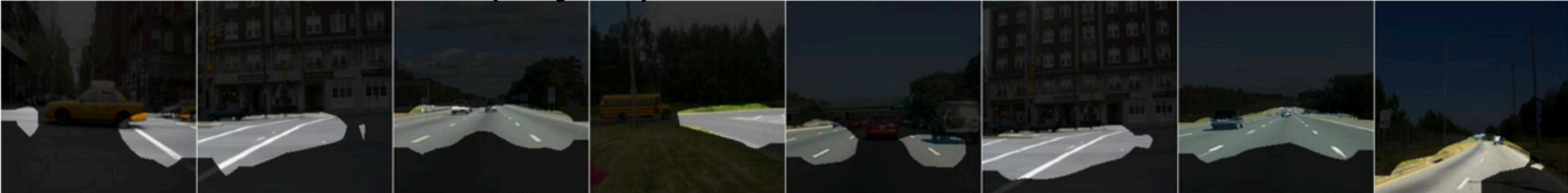
- **Results**: Object detector emerges even when the model trained on scene dataset
  - High-scored (interpretable) convolutional units

conv5 unit 79     car (object)       IoU=0.13



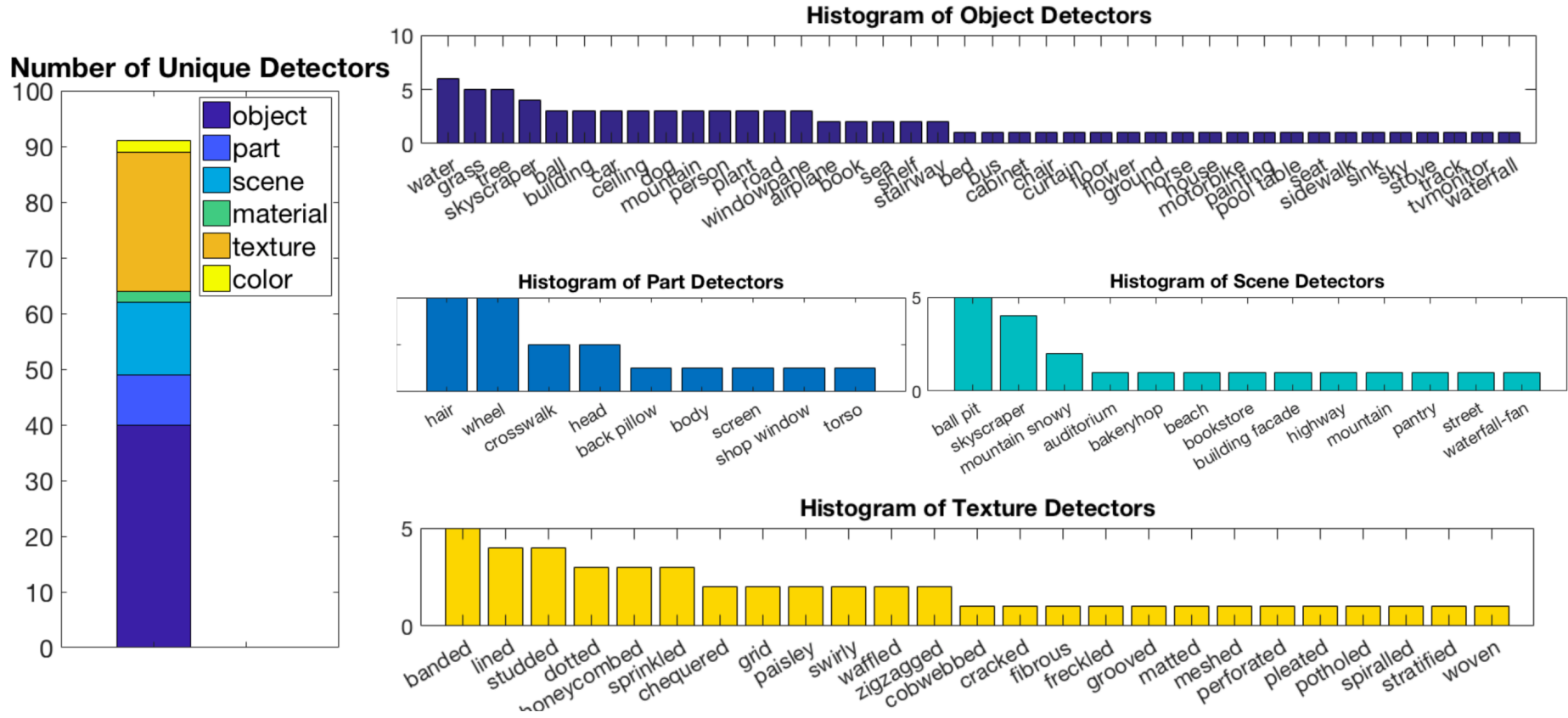conv5 unit 107   road (object)     IoU=0.15



conv5 unit 144   mountain (object)       IoU=0.13

- **Results**
  - Dissection report (AlexNet / trained on places 365)

- **Goal**: Identify <span style="color:red">most responsible training point</span> for a given prediction
  - Retraining the model can be <span style="color:red">prohibitively expensive!</span>

- **Idea**: Find $\hat{\theta}_{-z} = \arg\min_{\theta \in \Theta} \frac{1}{n} \sum_{z_i \neq z}^{n} L(z_i, \theta)$ using influence function

  - Training points $z_1, \cdots, z_n$ are given
  - The empirical risk minimizer is given by $\hat{\theta} = \arg\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^{n} L(z_i, \theta)$

- Measure influence of $L(z, \theta)$ on parameter $\hat{\theta}$ to approximate $\hat{\theta}_{-z}$

- Influence function

$$I(T) = \lim_{t \to 0^+} \frac{T(tG + (1-t)F) - T(F)}{t}$$

where $T$ : an estimator, $F, G$ : distribution

- Approximate $\hat{\theta}_{-z}$ in terms of perturbation $\epsilon$

$$\hat{\theta}_{-z} \approx \hat{\theta} - \frac{1}{n}\frac{d\hat{\theta}_{\epsilon,z}}{d\epsilon}\bigg|_{\epsilon=0} \qquad \text{where} \qquad \hat{\theta}_{\epsilon,z} = \arg\min_{\theta\in\Theta}\frac{1}{n}\sum_{i=1}^{n}L(z_i,\theta) + \epsilon L(z,\theta)$$

  - From a classic result [Cook et al., 1982],

$$\mathcal{I}_{\mathrm{up,params}}(z) = \frac{d\hat{\theta}_{\epsilon,z}}{d\epsilon}\bigg|_{\epsilon=0} = -H_{\hat{\theta}}^{-1}\nabla_\theta L(z,\hat{\theta})$$

- Influence of $z$ on the loss function at $z_{\mathrm{test}}$

$$\mathcal{I}_{\mathrm{up,loss}}(z, z_{\mathrm{test}}) = \frac{dL(z_{\mathrm{test}}, \hat{\theta}_{\epsilon,z})}{d\epsilon}\bigg|_{\epsilon=0}$$

$$= \nabla_\theta L(z_{\mathrm{test}}, \hat{\theta})^\top \frac{d\theta_{\epsilon,z}}{d\epsilon}\bigg|_{\epsilon=0}$$

$$= -\nabla_\theta L(z_{\mathrm{test}}, \hat{\theta})^\top H_{\hat{\theta}}^{-1}\nabla_\theta L(z,\hat{\theta})$$

  - Helpful images implies negative influence on the loss function

- **Results**
  - Understanding model behavior (discriminate fish vs. dog)
    - Helpful images implies <span style="color:red">negative influence on loss function</span>
    - For Inception network, most helpful image was actually a dog



Helpful images

## Conclusion

- **Interpretability method** is about **giving explanation** to human
  - Form of explanation is various

- In this lecture, we covered some of interpretability methods
  - Visual explanation (saliency map / class activation map)
  - Network dissection
  - Influence function

- There are still many research directions
  - Lots of interpretability methods not covered in this slide

# References

[Zeiler et al., 2014] Visualizing and Understanding Convolutional Networks. ECCV 2014.
https://arxiv.org/abs/1610.02136

[Zintgraf et al., 2017] Visualizing Deep Neural Network Decisions: Prediction Difference Analysis. ICLR 2017.
https://arxiv.org/abs/1702.04595

[Ribeiro et al., 2016] "Why Should I Trust You?": Explaining the Predictions of Any Classifier, KDD 2016.
https://arxiv.org/abs/1602.04938

[Simonyan et al., 2014] Deep Inside Convolutional Networks: Visualising Image Classificiation ... ICLR Workshop 2014.
https://arxiv.org/abs/1312.6034

[Sundararajan et al., 2017] Axiomatic Attribution for Deep Networks. ICML 2017.
https://arxiv.org/abs/1703.01365

[Smilkov et al., 2017] SmoothGrad: removing noise by adding noise.
https://arxiv.org/abs/1706.03825

[Selvararaju et al., 2017] Grad-CAM: Visual Explanations from Deep Networks via Gradient-based ... ICCV 2017.
https://arxiv.org/abs/1610.02391

[Zhou et al., 2015] Learning Deep Features for Discriminative Localization. ICCV 2015.
https://arxiv.org/abs/1512.04150

[Adebayo et al., 2017] Sanity Checks for Saliency Maps. NIPS 2018.
https://arxiv.org/abs/1810.03292

[Bau et al., 2017] Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017.
https://arxiv.org/abs/1704.05796

[Koh et al., 2017] Understanding Black-box Predictions via Influence Functions. ICML 2017.
https://arxiv.org/abs/1703.04730