

CS7.405 Responsible & Safe AI Systems

Ponnurangam Kumaraguru ("PK")
#ProfGiri @ IIIT Hyderabad



pk.profgiri



/in/ponguru



@ponguru



Ponnurangam.kumaraguru

Guest lectures how?

Daniel Paleka

Adam

Andy Zou

Arun Jose

Neel Nanda

☹ Not all join the lecture ☹

Goal

Less of Blackbox and more transparent

Motivation

Transparency tools try to provide clarity about a model's inner workings

Model changes can sometimes cause the internal representations to substantially change, so we would like to understand when models process data differently

Transparency could make it easier for monitors to detect deception and other hazards

Pixel attribution methods

Sensitivity map, saliency map, pixel attribution map, gradient-based attribution methods, feature relevance, feature attribution, and feature contribution.

Feature attribution explains individual predictions by attributing each input feature according to how much it changed the prediction (negatively or positively).

Pixel attribution methods

Occlusion- or perturbation-based

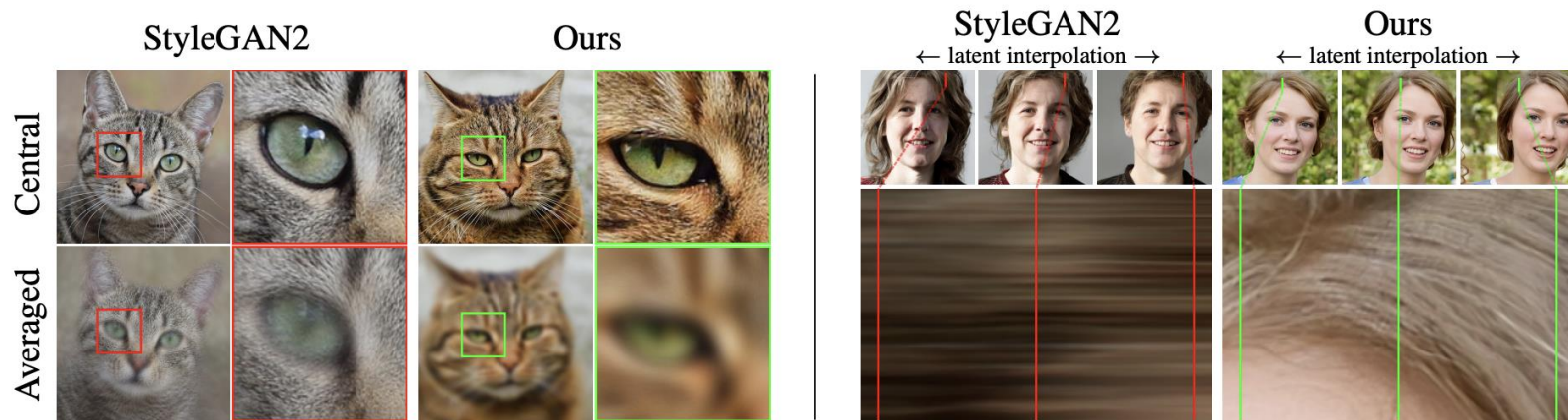
Methods like SHAP and LIME manipulate parts of the image to generate explanations (model-agnostic).

Gradient-based

Many methods compute the gradient of the prediction (or classification score) with respect to the input features.

The gradient-based methods mostly differ in how the gradient is computed.

Motivation



Texture sticking

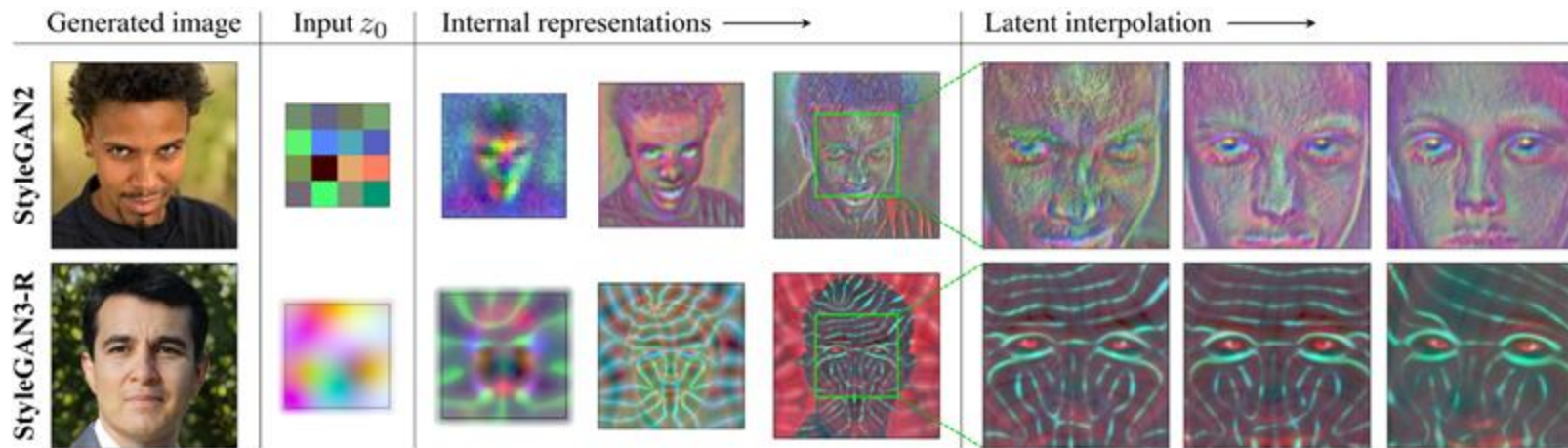
Left: average of images generated from a small neighborhood around a central latent (top row)

Right: extract small vertical segment of pixels, stack horizontally

StyleGAN2, same coordinates

Hairs moving in animation

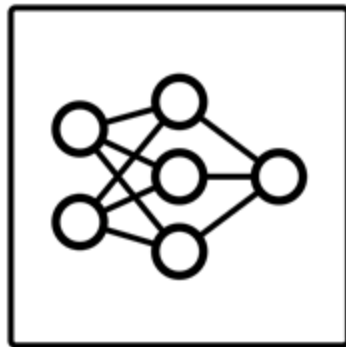
Motivation



StyleGAN2: details glued to the image vs surface; internal representations are different

StyleGAN3: fully equivariant to translation and rotation; help in identifying important properties better

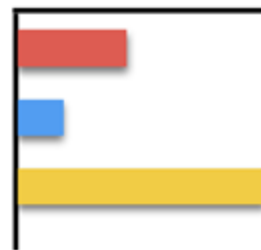
Saliency Maps



Gradient

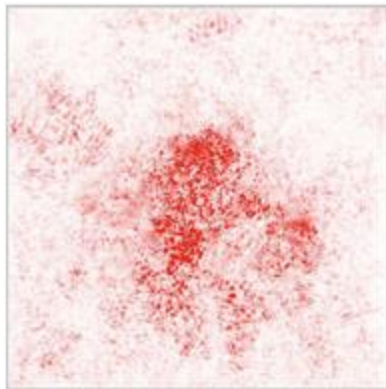


Predictions



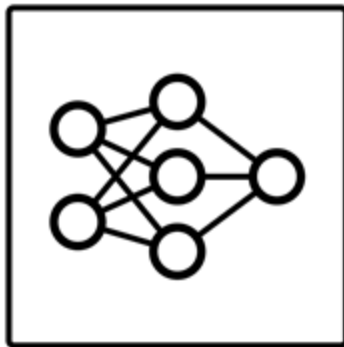
Corn

$$E_{grad}(x) = \frac{\partial S_i}{\partial x}$$



Perturbation direction
of fastest ascent for
the class logit

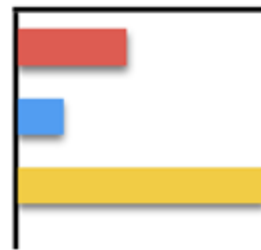
Saliency Maps



SmoothGrad

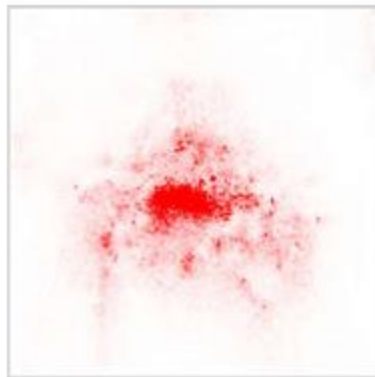


Predictions



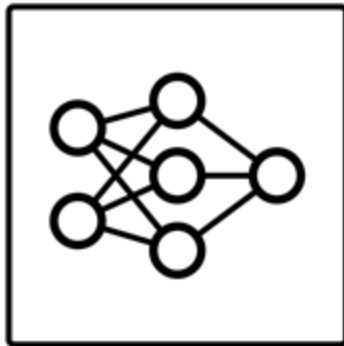
Corn

$$E_{sg}(x) = \frac{1}{N} \sum_{i=1}^N E(x + g_i),$$



Each input perturbed
by different Gaussian
noise and then
averaged
Smoother & less noisy

Saliency Maps



Guided
BackProp

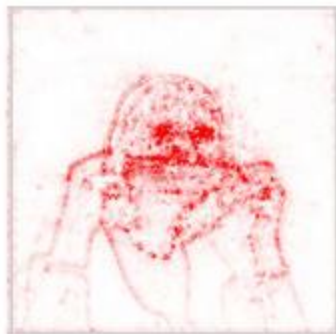


Predictions



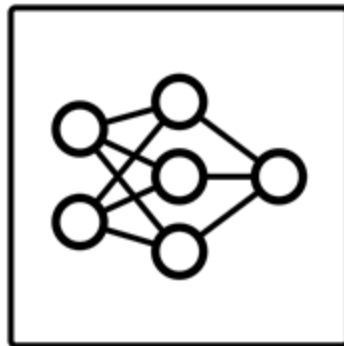
Corn

$$R^l = 1_{R^{l+1} > 0} 1_{f^l > 0} R^{l+1}$$

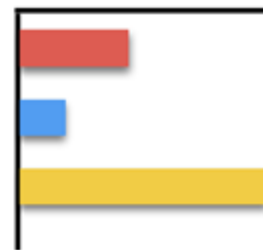


Backprop with
intermediate negative
activations and gradients
zeroed out

Saliency Maps



Predictions



Corn

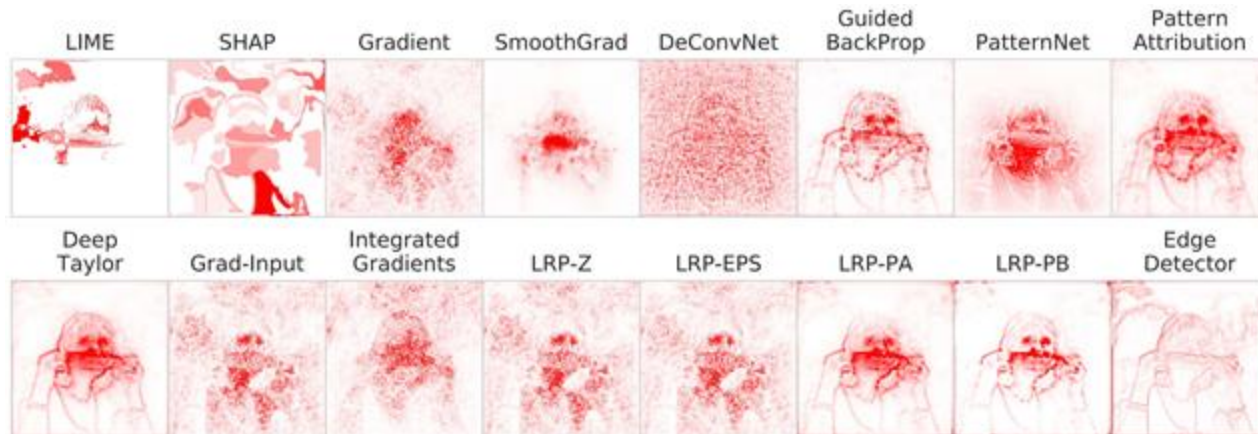
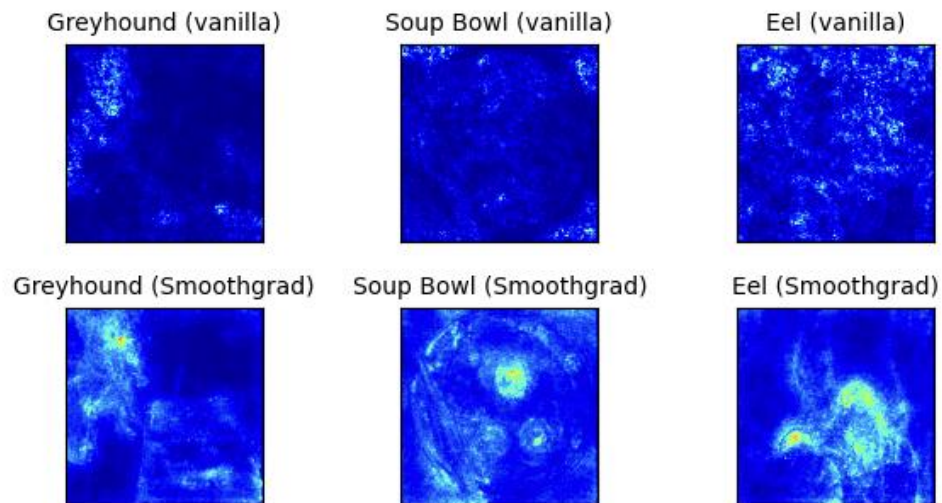



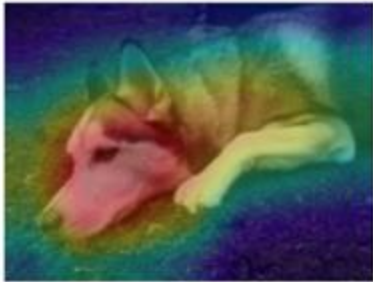
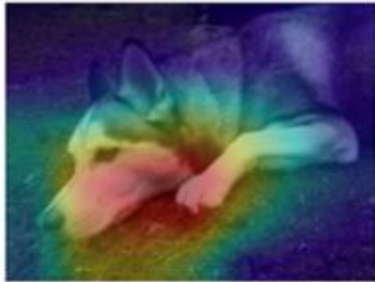


FIGURE 10.9: Images of a dog classified as greyhound, a ramen soup classified as soup bowl, and an octopus classified as eel.



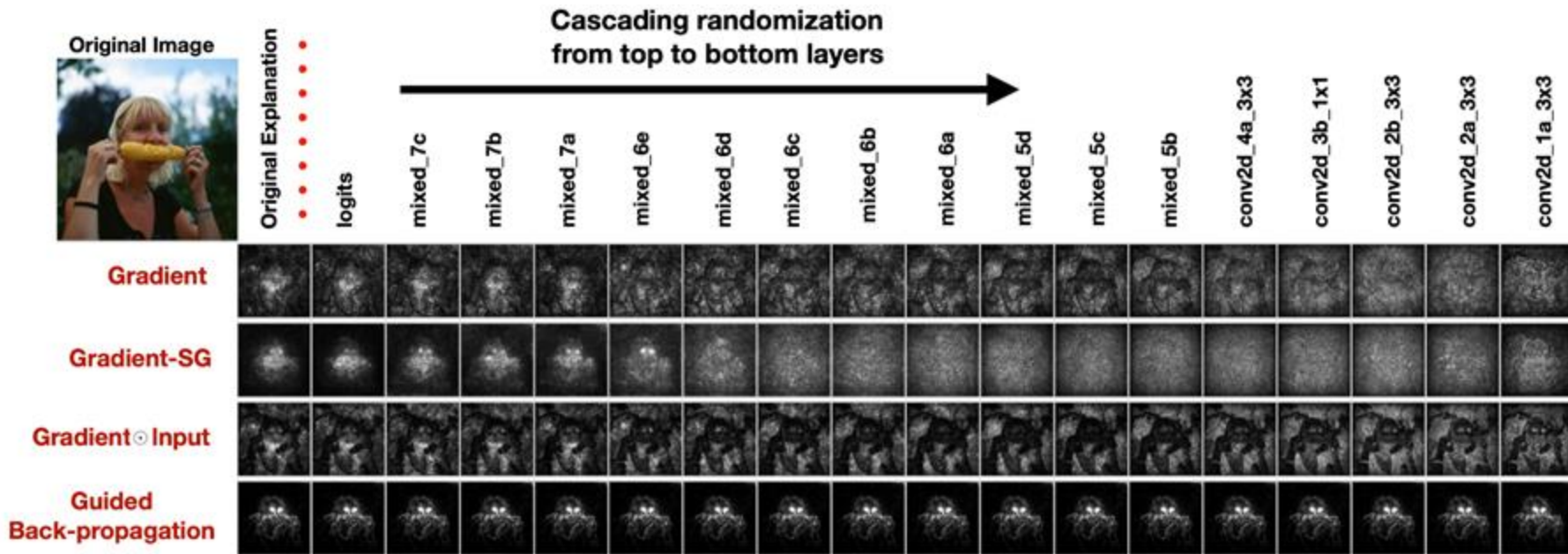
Saliency Maps Can Be Deceptive

Many transparency tools create fun-to-look-at visualizations that do not actually inform us much about how models are making predictions

	Test Image	Evidence for Animal Being a Siberian Husky	Evidence for Animal Being a Transverse Flute
Explanations Using Attention Maps			

Sanity Checks for Saliency Maps

If we randomize the layers, some saliency maps do not change much, which suggests they do not capture what the model has learned



Optimized Masks for Saliency

Some saliency maps optimize a mask to locate and blur salient regions



Figure 1. An example of a mask learned (right) by blurring an image (middle) to suppress the softmax probability of its target class (left: original image; softmax scores above images).

This is highly sensitive to hyperparameters and mask initialization

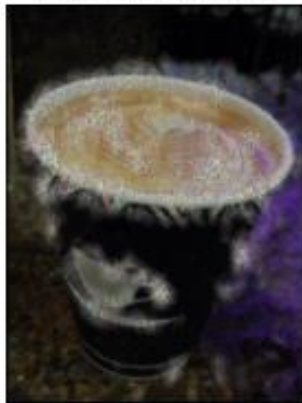
espresso: 0.9964



maypole: 0.9568



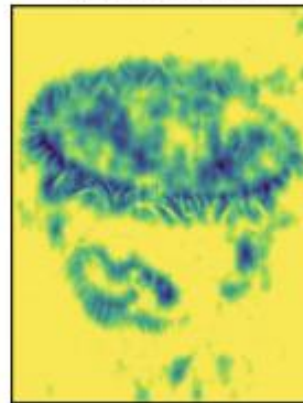
espresso: 0.0000



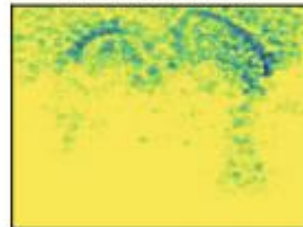
maypole: 0.0000

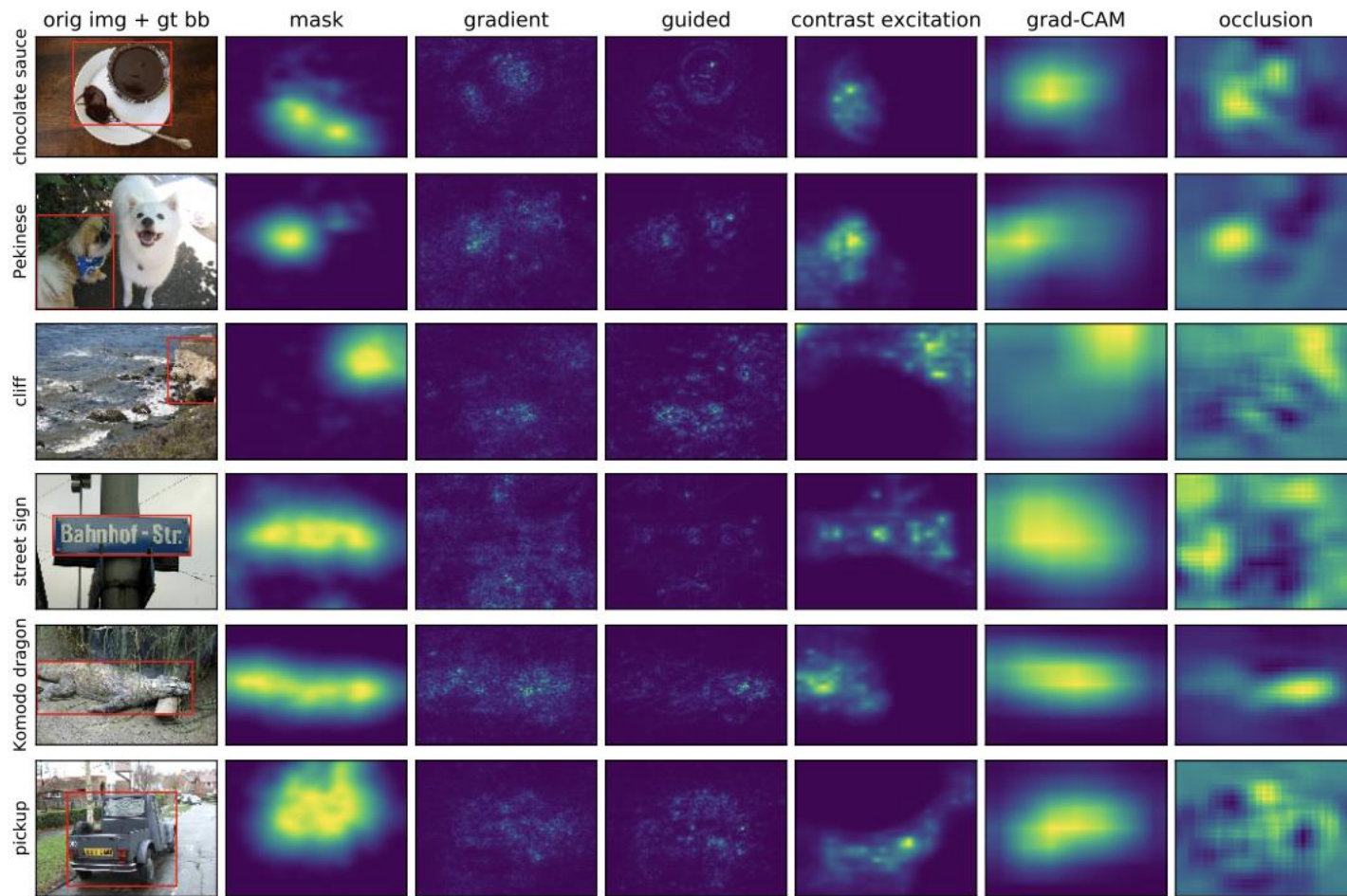


Learned Mask



Learned Mask





Pros & Cons Gradient based

Explanations are visual, detecting important regions is easy in the image

Faster to compute than model-agnostic methods

LIME & SHAP are very expensive

Difficult to know whether an explanation is correct

Very fragile - adversarial perturbations produce same prediction

Saliency Maps for Text

Saliency maps can be used for text models too

	$p(y \mathbf{x}; \theta)$	y	c
the year 's best and most unpredictable comedy	0.91	pos	pos
we never feel anything for these characters	0.95	neg	neg
handsome but unfulfilling suspense drama	0.18	neg	pos

y = gold, c = predicted

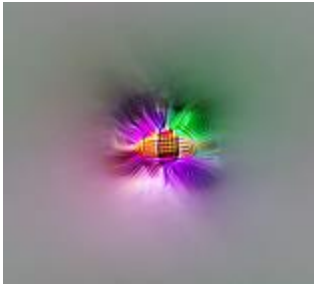
There are many possible saliency scores for a token; one possibility is to use the magnitude of the gradient of the classifier's logit with respect to the token's embedding

While there is no canonical saliency map, these can be used for identifying salient words when writing adversarial examples

Feature Visualization

To understand what a model's internal component detects, synthesize an image through gradient descent that maximizes the component

Neuron Visualization



Channel Visualization




Maximally Activating Natural



NV: Component = Neuron, optimize the image to maximally activate the neuron, repeated round of GD optimize the noise image

CV: Like Neuron Viz, both gradient descent, Loss of channel visualization might be sum of the squares of all neurons in the channel, lot of squares


Microscope

MODELSABO

Models<<ResNet v2 50<<block2/unit_3/add<<Unit 18

AlexNet1,001predictions/Softmax

AlexNet (Places)1,001logits/Conv2D

AlexNet (Places)2,048block4/unit_3/add

Inception v13,33,1,1

Inception v1 (Places)2,048block4/unit_2/add

VGG 193,33,1,1

Inception v32,048block4/unit_1/add

Inception v43,33,1,1

ResNet v2 501,024block3/unit_6/add

CLIP Resnet 50 v03,33,1,1

CLIP Resnet 501,024block3/unit_5/add

CLIP Resnet 1013,33,1,1

CLIP Resnet 50 4x1,024block3/unit_4/add

CLIP Resnet 50 16x1,024block3/unit_3/add

1,024block3/unit_2/add

Unit 0Unit 1Unit 2

Unit 3Unit 4Unit 5

Unit 6Unit 7Unit 8

Unit 9Unit 10Unit 11

Unit 12Unit 13Unit 14

Unit 15Unit 16Unit 17

Unit 18Unit 19Unit 20

Unit 21Unit 22Unit 23

Unit 24Unit 25Unit 26

Unit 27Unit 28Unit 29

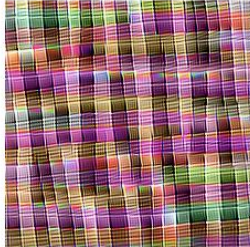
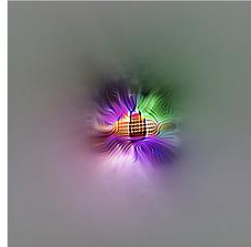


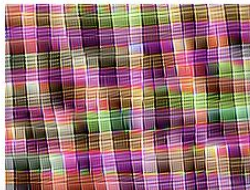

Unit 30Unit 31Unit 32


Unit 33Unit 34Unit 35

FEATURE VISUALIZATION







An artificial, optimized image that maximizes activations of the given unit

[Read more.](#)

 Microscope
 MODELS
ABOUT


Models
 <<
 ResNet v2 50
 <<
 block2/unit_3/add
 <<
 Unit 18
 ✕

AlexNet	1,001	predictions/Softmax	  
AlexNet (Places)	1,001	logits/Conv2D	
Inception v1	2,048	block4/unit_3/add	
Inception v1 (Places)			
VGG 19		block4/unit_2/add	
Inception v3	2,048	block4/unit_1/add	
Inception v4			
ResNet v2 50	1,024	block3/unit_6/add	
CLIP Resnet 50 v0			
CLIP Resnet 50	1,024	block3/unit_5/add	
CLIP Resnet 101			  
CLIP Resnet 50 4x			
CLIP Resnet 50 16x	1,024	block3/unit_4/add	
	1,024	block3/unit_3/add	
	1,024	block3/unit_2/add	

DATASET SAMPLES

Pieces of images from the training dataset that result in the largest activations from the given unit.

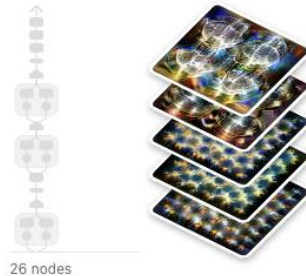
These images are cropped and downsized samples from the [ImageNet](#) research dataset. Unlike our other visualizations, they are not CC-BY-SA because they are derived from ImageNet.

DATASET: IMAGENET


The OpenAI Microscope is a collection of visualizations of every significant layer and neuron of 13 important vision models. [LEARN MORE>](#)

AlexNet

A landmark in computer vision, this 2012 winner of ImageNet has over 50,000 citations.



AlexNet (Places)

The same architecture as the classic AlexNet model, but trained on the Places365 dataset.



Inception v1

Also known as GoogLeNet, this network set the state of the art in ImageNet classification in 2014.



Inception v1 (Places)

The same architecture as the classic Inception v1 model, but trained on the Places365 dataset.

VGG 19

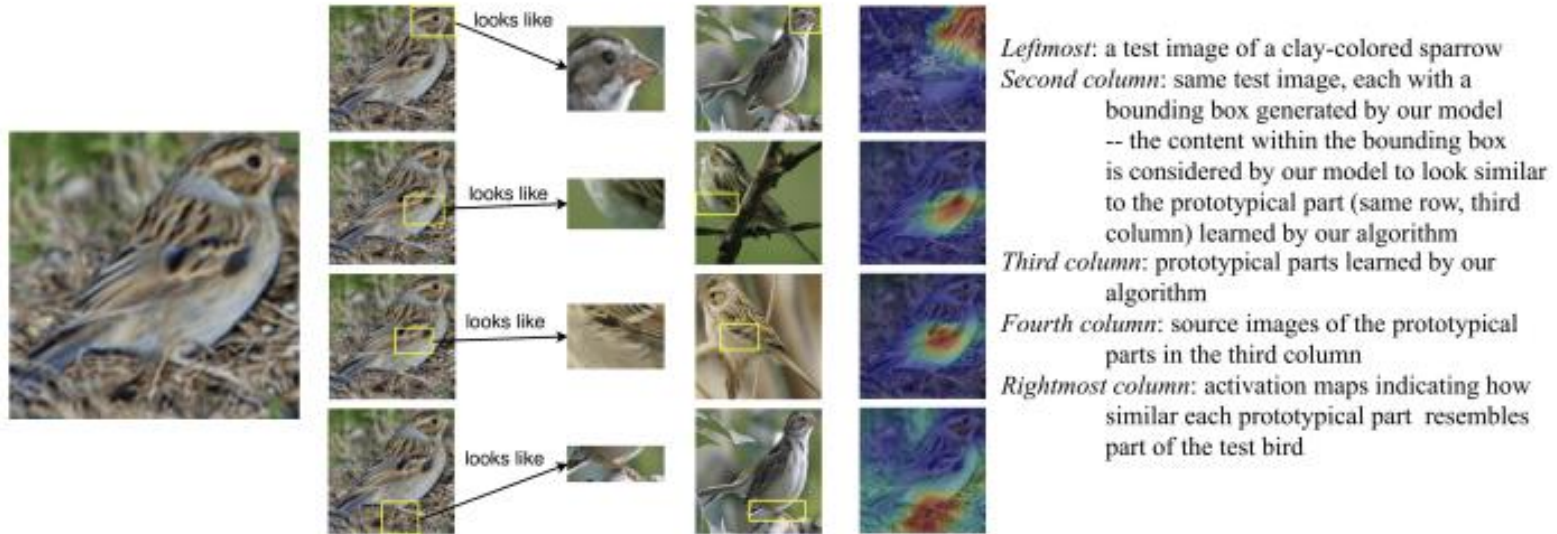
Introduced in 2014, this network is simpler than Inception variants, using only 3x3 convolutions and no

Inception v3

Released in 2015, this iteration of the Inception architecture improved performance and efficiency.

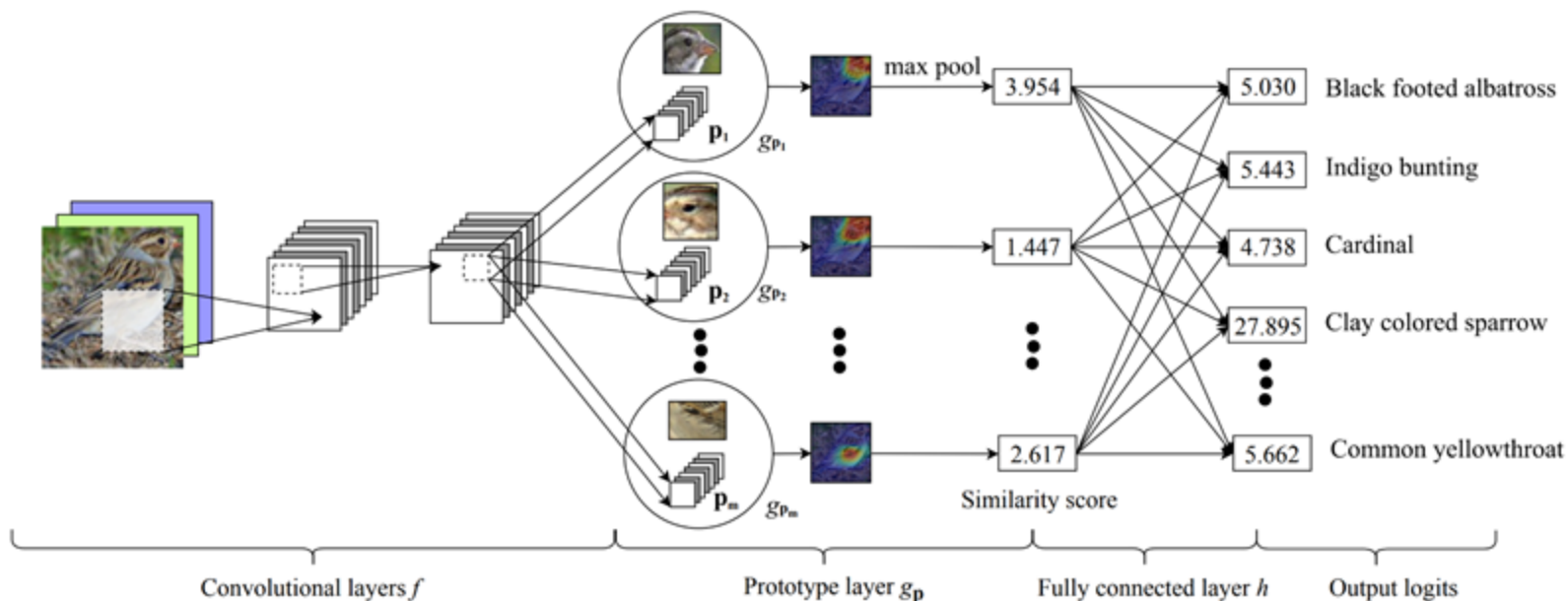
ProtoPNet (“This Looks Like That”)

These models perform classifications based on the most important patches of training images, using patches that are prototypical of the class



ProtoPNet (“This Looks Like That”)

These models perform classifications based on the most important patches of training images, using patches that are prototypical of the class



Administrativa

Project review today, hopefully all of you have made decent progress

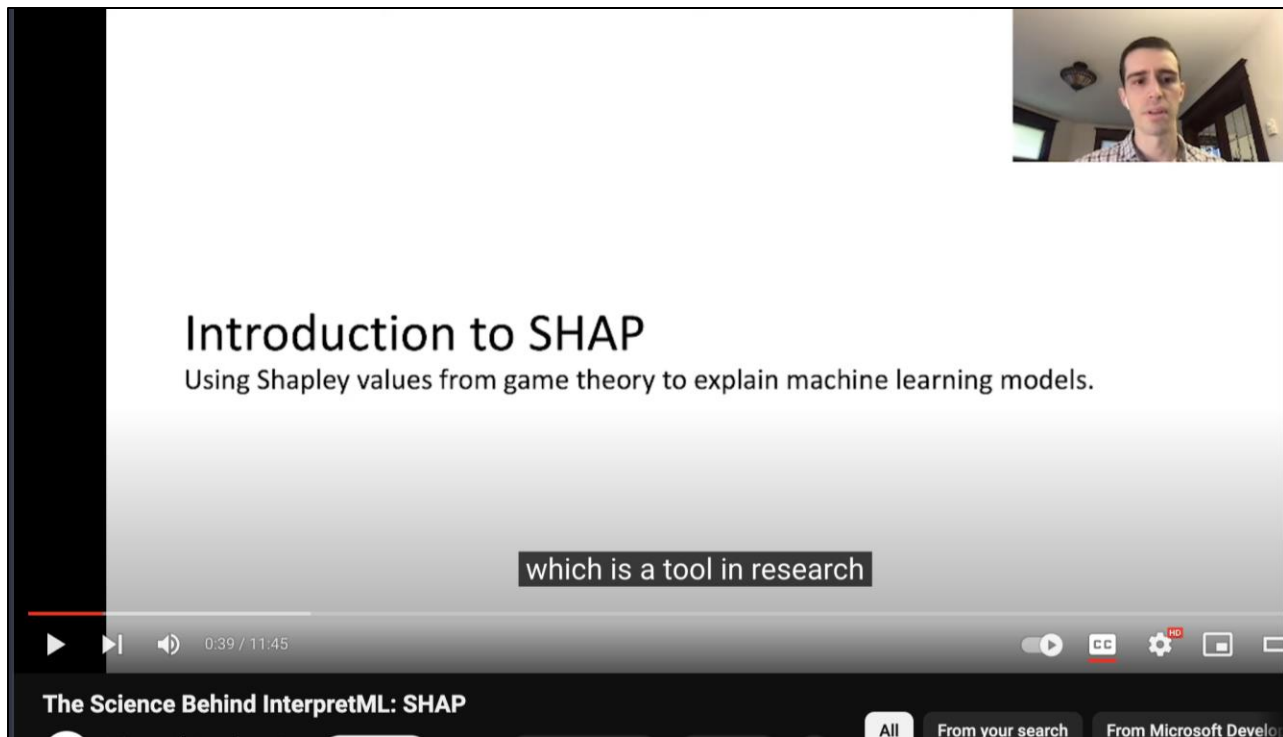
Poster presentation

24 April [Wednesday] 1600 – 1800hrs?

Figure out a location?

25 April [Thursday] class timing or anytime will be hard to get audience?

SHAP



https://youtu.be/-taOhqkiulo?si=TGDmiUD9X-kEIV_j

Bibliography / Acknowledgements

<https://course.mlsafety.org/>

https://rdi.berkeley.edu/understanding_llms/s24

<https://aisafetyfundamentals.com/>

<https://inst.eecs.berkeley.edu/~cs294-149/fa18/>

<https://aisafety.stanford.edu/>

<https://docs.google.com/document/d/1goyFTfu-EB90yuepTFN2unmf-fel3TycTPYh7Kn6dl/edit>

<https://medium.com/@richardcngo/visualizing-the-deep-learning-revolution-722098eb9c5>

 pk.profgiri

 Ponnurangam.kumaraguru

 /in/ponguru

 ponguru

 pk.guru@iiit.ac.in

Thank you
for attending
the class!!!