



Computer Vision; Image Classification; Visualizing & Understanding



[YouTube Playlist](#)

Maziar Raissi

Assistant Professor

Department of Applied Mathematics

University of Colorado Boulder

maziar.raissi@colorado.edu

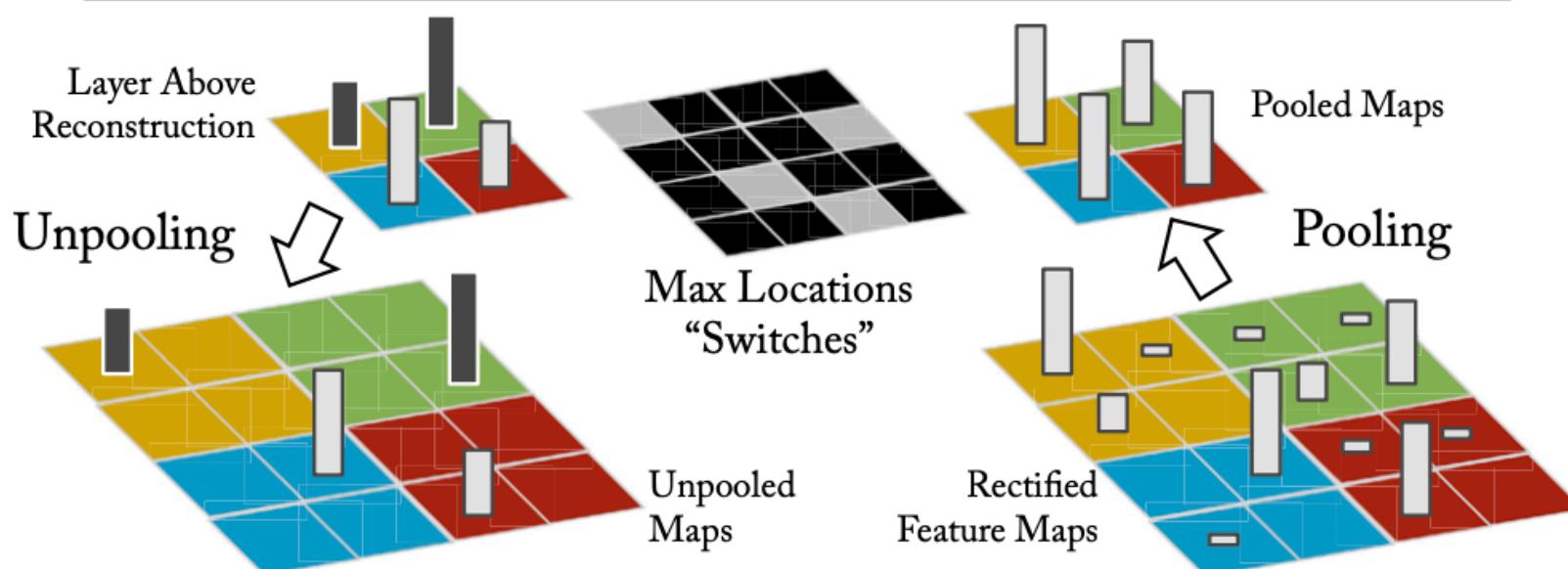
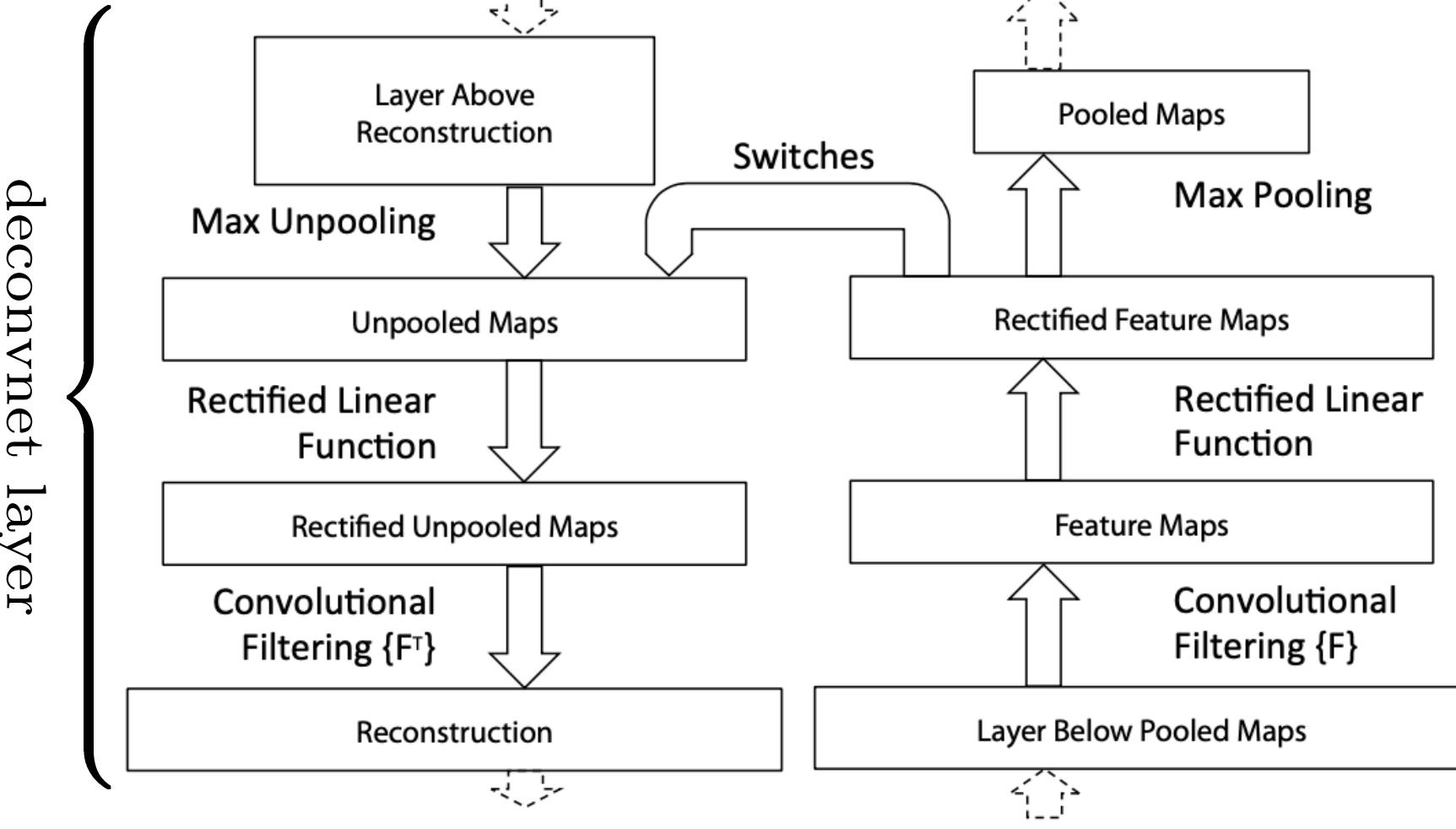
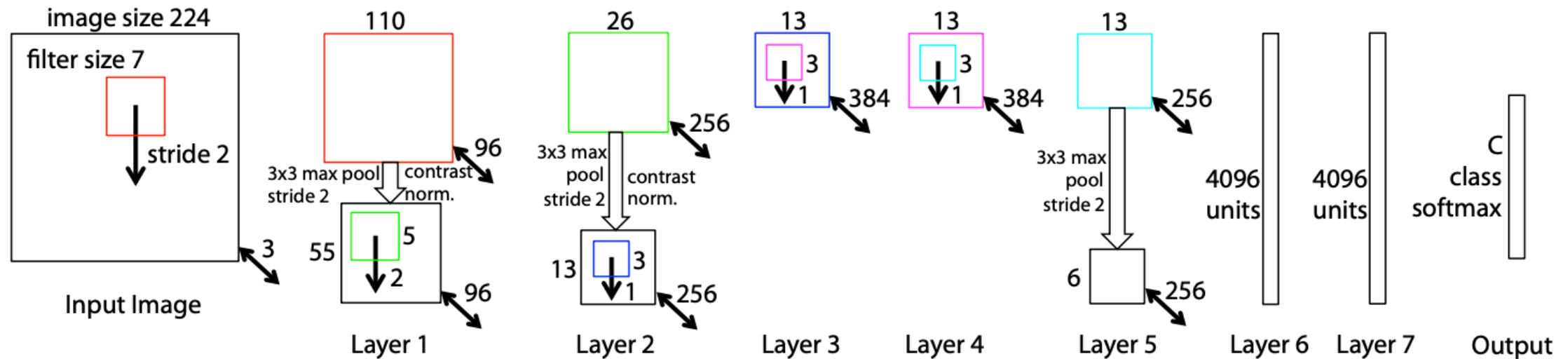


Boulder

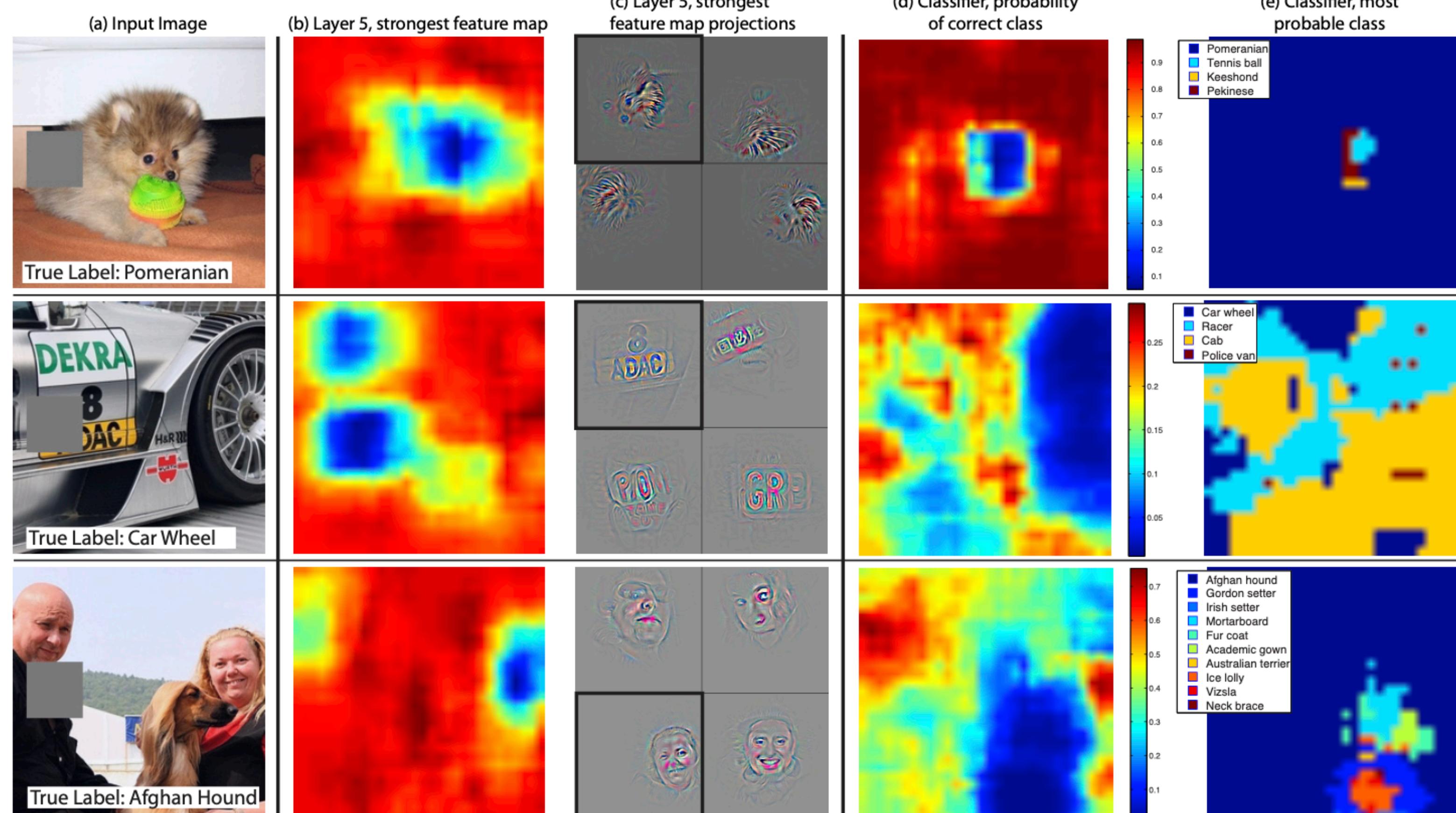
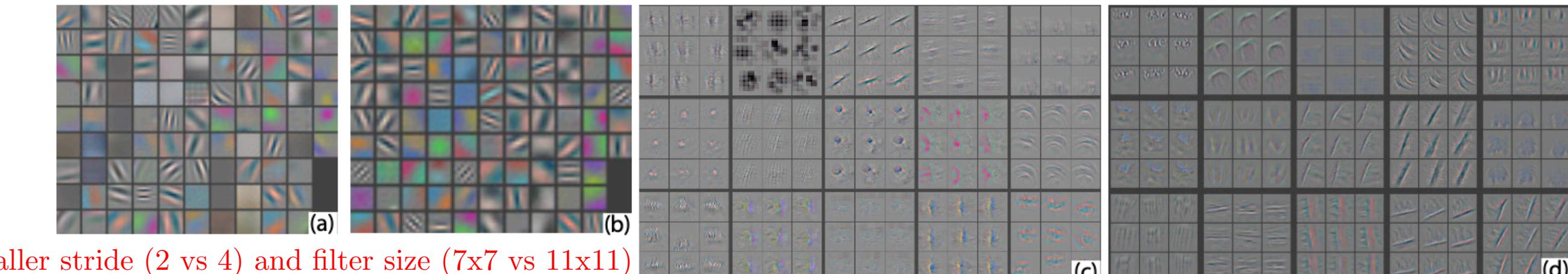


[YouTube Video](#)

Visualizing and Understanding Convolutional Networks



Filtering: Flipping each filter vertically and horizontally



Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps



Class Model Visualization

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

$S_c(I)$ → score of class c for image I

$$P_c(I) = \frac{\exp S_c(I)}{\sum_{c'} \exp S_{c'}(I)} \rightarrow \text{probability of class } c$$

λ → regularization parameter

Image-Specific Class Saliency Visualization

I_0 → image

c → class

$S_c(I) \approx w_c^T I + b_c$ for I in the neighborhood of I_0

$$w_c = \left. \frac{\partial S_c(I)}{\partial I} \right|_{I=I_0}$$

magnitude of elements of w_c defines the importance of the corresponding pixels of I for the class c

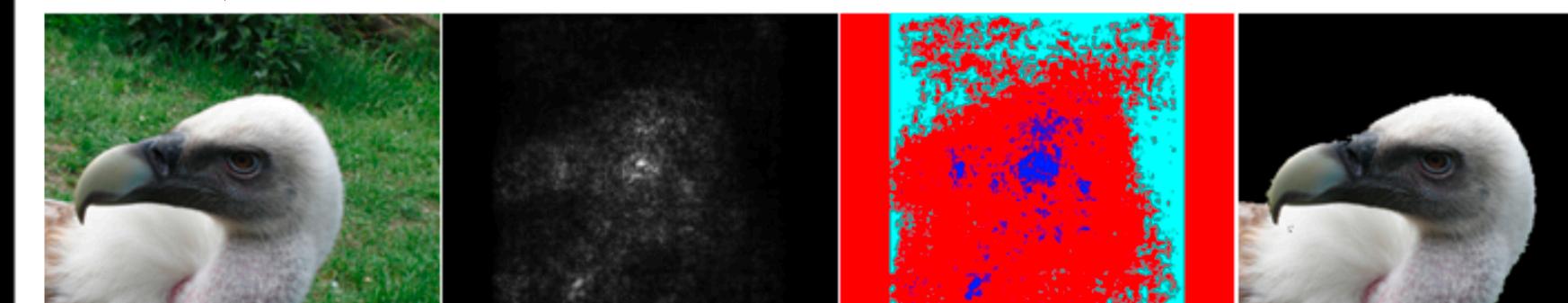
Class Saliency Extraction

$$I_0 \in \mathbb{R}^{H \times W \times K} \implies w_c \in \mathbb{R}^{H \times W \times K}$$

$$M_{ij} := \max_k |w_{ijk}^c| \implies M \in \mathbb{R}^{H \times W}$$

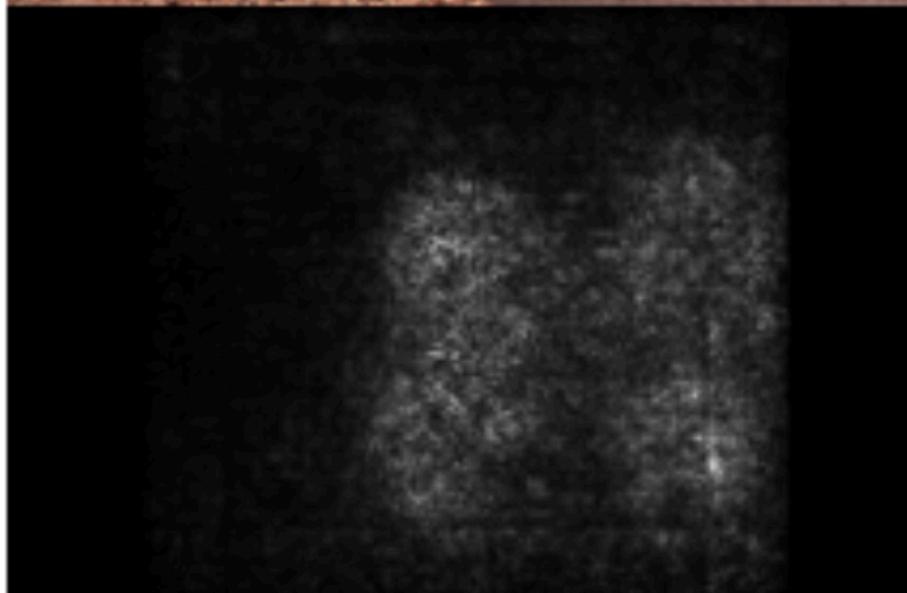
Weakly Supervised Object Localization

blue - foreground color model; cyan - background color model; red - not used for color model estimation



Weakly Supervised Object Localization

blue - foreground color model; cyan - background color model; red - not used for color model estimation



Relation to Deconvolution Networks

$X_n \rightarrow n\text{-th layer input}$

$f \rightarrow \text{neuron activity to be visualized}$

$X_{n+1} = X_n * K_n \rightarrow \text{convolutional later}$

$$\frac{\partial f}{\partial X_n} = \frac{\partial f}{\partial X_{n+1}} * \hat{K}_n$$

\hat{K}_n → flipped version of the convolutional kernel K_n

$R_n \rightarrow n\text{-th layer reconstruction in a DeconvNet}$

$$R_n = R_{n+1} * \hat{K}_n$$

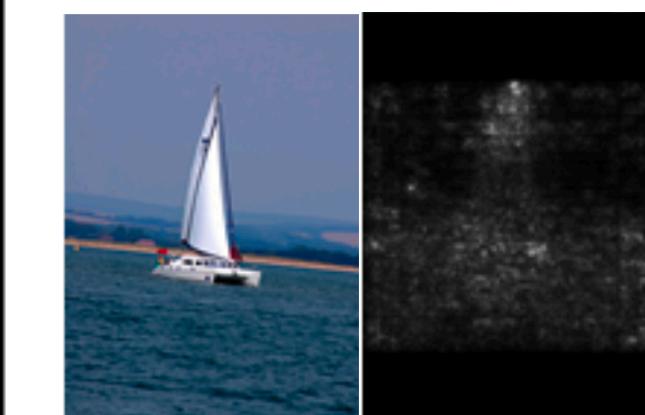
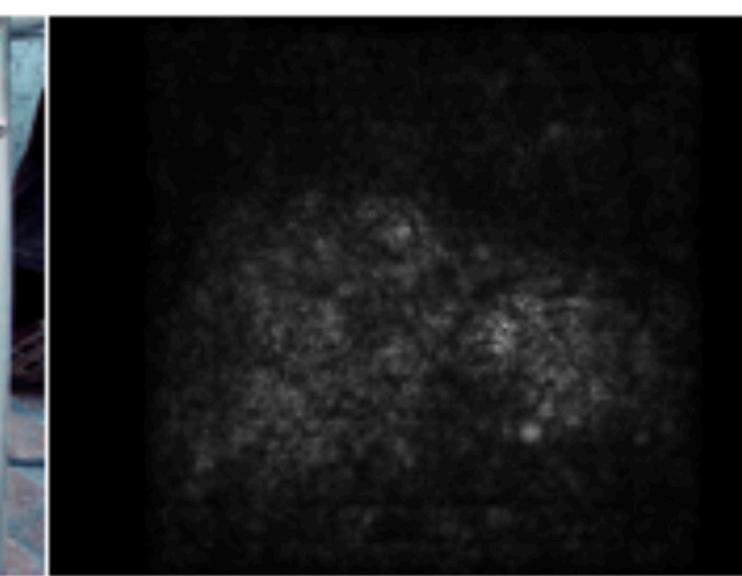
$X_{n+1} = \max(X_n, 0) \rightarrow \text{ReLU}$

$$\frac{\partial f}{\partial X_n} = \frac{\partial f}{\partial X_{n+1}} \mathbf{1}(X_n > 0)$$

$R_n = R_{n+1} \mathbf{1}(\mathcal{R}_{n+1} > 0) \rightarrow \text{slightly different from above}$

$$X_{n+1}(p) = \max_{q \in \Omega(p)} X_n(q) \rightarrow \text{maxpooling}$$

$$\frac{\partial f}{\partial X_n(s)} = \frac{\partial f}{\partial X_{n+1}(p)} \mathbf{1}(s = \arg \max_{q \in \Omega(p)} X_n(q)) \rightarrow \text{switches}$$





Boulder

Striving for Simplicity: The All Convolutional Net

$f \in \mathbb{R}^{H \times W \times N} \rightarrow$ feature maps produced by some layer of CNN

$$s_{i,j,u}(f) = \left(\sum_{h=-\lfloor k/2 \rfloor}^{\lfloor k/2 \rfloor} \sum_{w=-\lfloor k/2 \rfloor}^{\lfloor k/2 \rfloor} |f_{g(h,w,i,j,u)}|^p \right)^{1/p} \rightarrow p\text{-norm subsampling (pooling)}$$

$$g(h, w, i, j, u) = (r \cdot i + h, r \cdot j + w, u)$$

$k \rightarrow$ pooling size, $k/2 \rightarrow$ half-length, $r \rightarrow$ stride

$p \rightarrow \infty \Rightarrow$ max-pooling

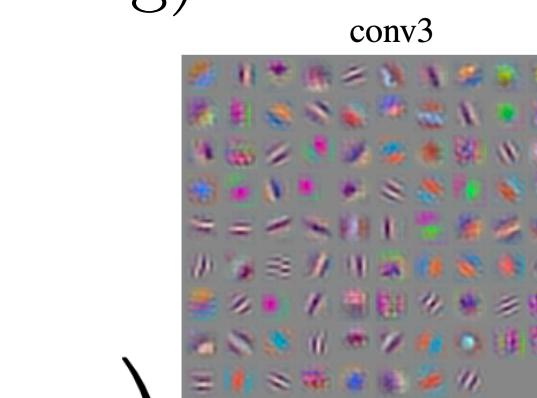
$$c_{i,j,o}(f) = \sigma \left(\sum_{h=-\lfloor k/2 \rfloor}^{\lfloor k/2 \rfloor} \sum_{w=-\lfloor k/2 \rfloor}^{\lfloor k/2 \rfloor} \sum_{u=1}^N \theta_{h,w,u,o} \cdot f_{g(h,w,i,j,u)} \right)$$

convolutional weights (or the kernel weights, or filters)

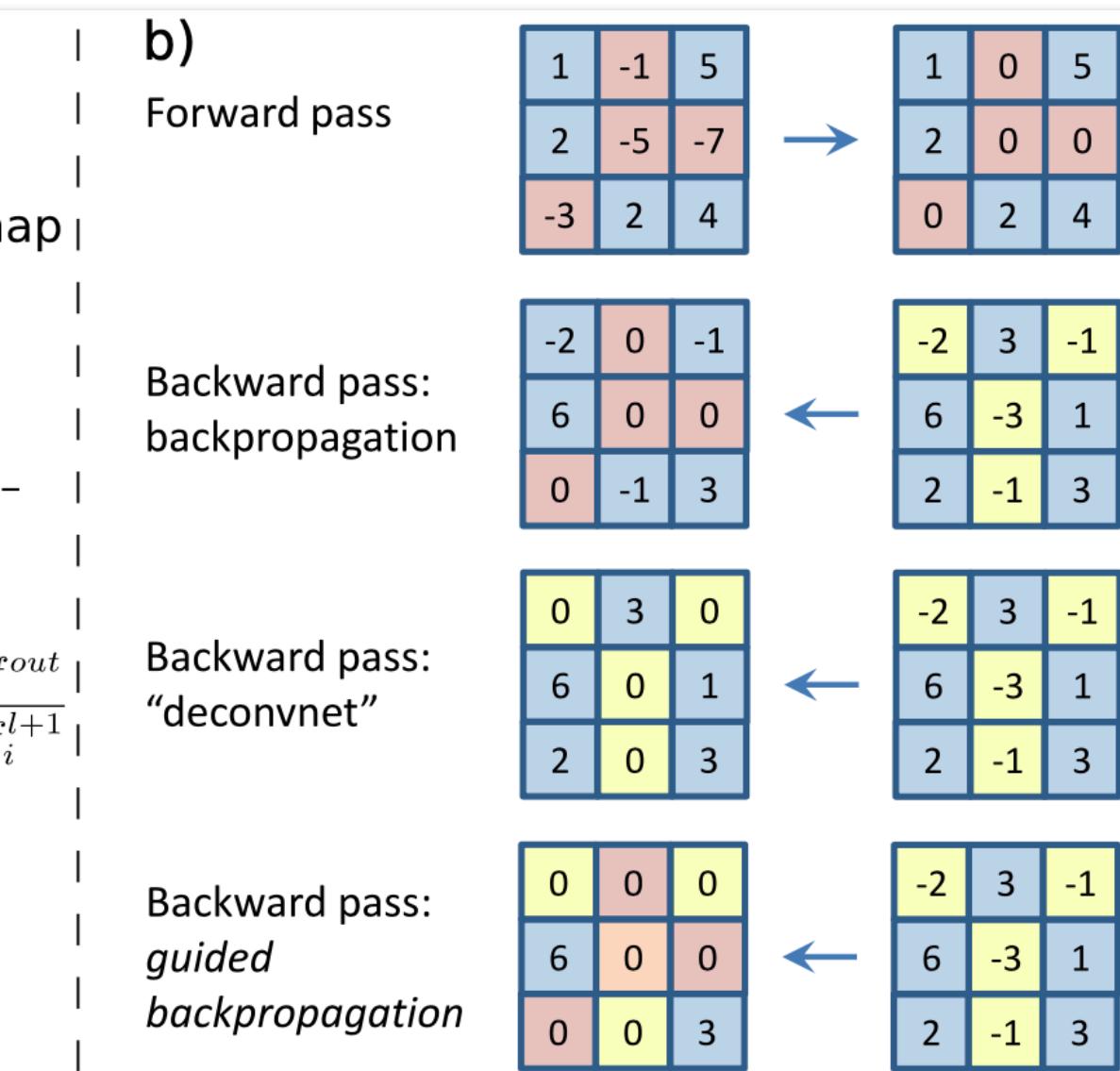
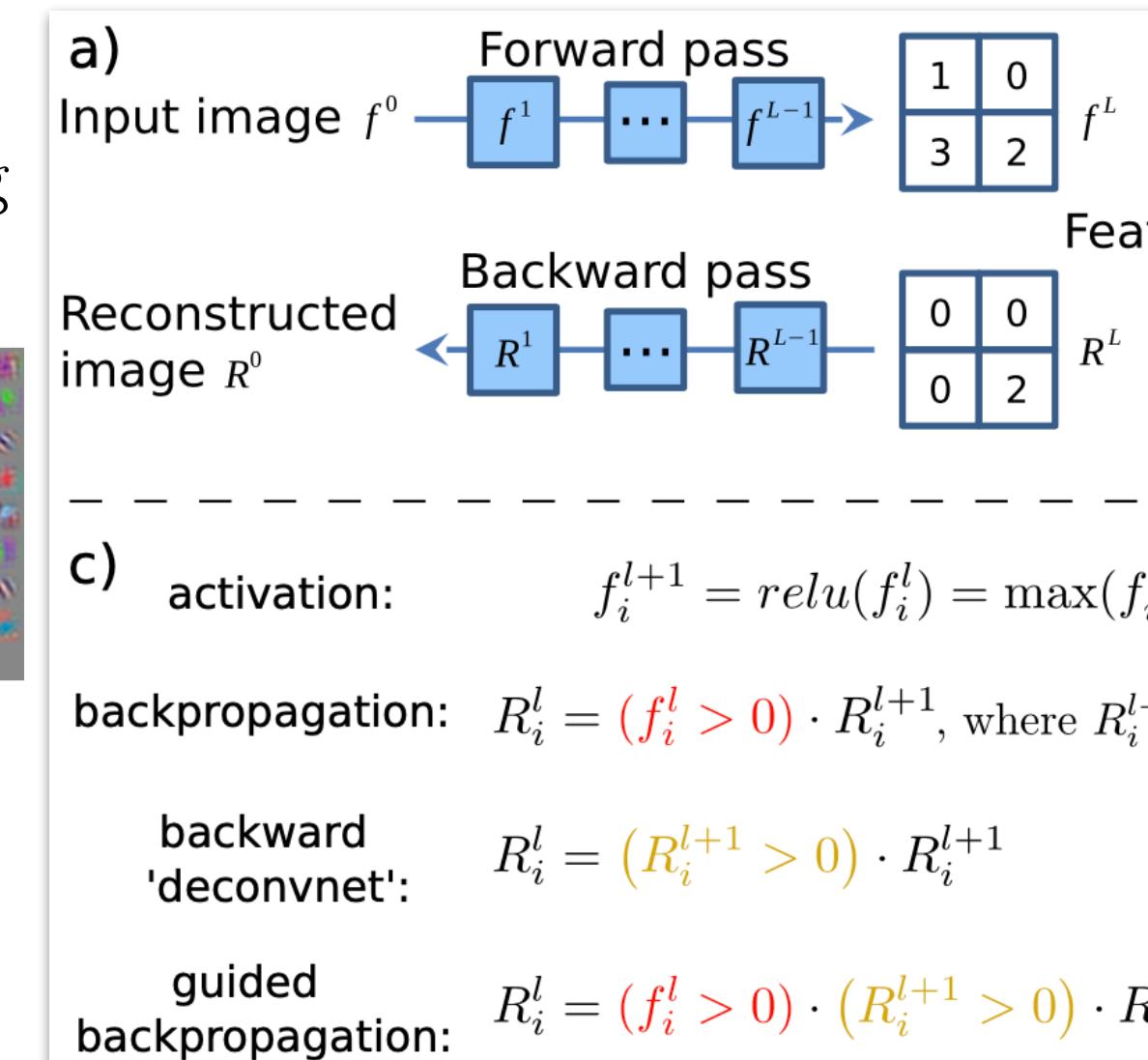
$$\sigma(x) = \max(x, 0), \text{ and } o \in [1, M]$$

feature-wise (depth-wise) convolution: $\theta_{h,w,u,o} = 1$ if u equals o and zero otherwise

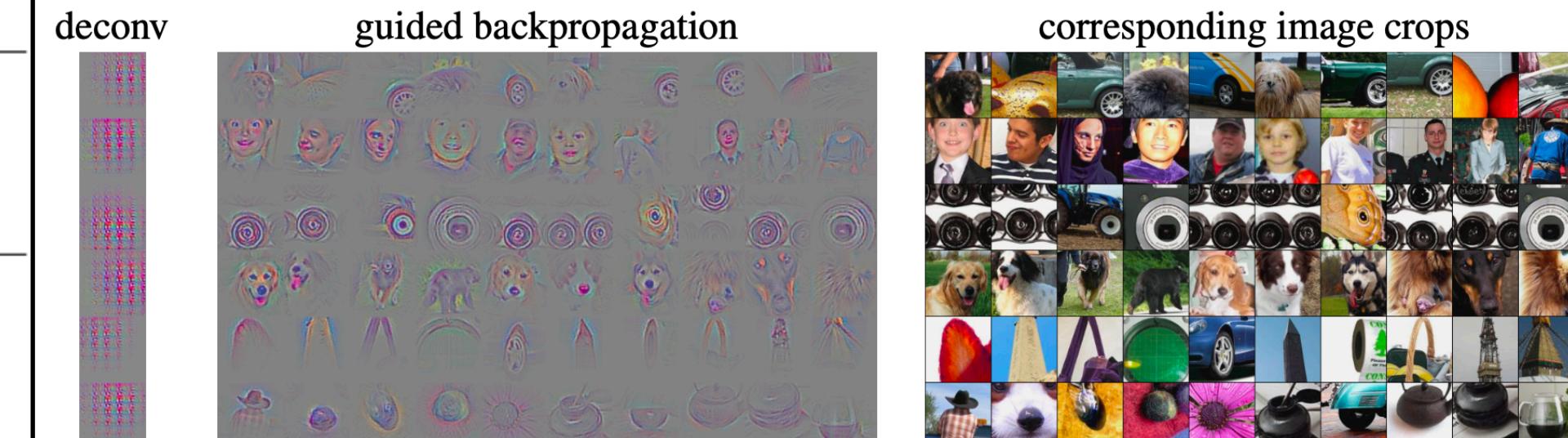
Model		
Strided-CNN-C	ConvPool-CNN-C	All-CNN-C
Input 32×32 RGB image		
3×3 conv. 96 ReLU	3×3 conv. 96 ReLU	3×3 conv. 96 ReLU
3×3 conv. 96 ReLU with stride $r = 2$	3×3 conv. 96 ReLU	3×3 conv. 96 ReLU
	3×3 conv. 96 ReLU with stride $r = 2$	3×3 conv. 96 ReLU with stride $r = 2$
3×3 conv. 192 ReLU	3×3 conv. 192 ReLU	3×3 conv. 192 ReLU
3×3 conv. 192 ReLU with stride $r = 2$	3×3 conv. 192 ReLU	3×3 conv. 192 ReLU
	3×3 conv. 192 ReLU with stride $r = 2$	3×3 conv. 192 ReLU with stride $r = 2$
:		



guided backpropagation



By using the switches from a forward pass the "deconvnet" (and thereby its reconstruction) is hence conditioned on an image and does not directly visualize learned features. An all convolutional architecture does not include max-pooling, meaning that in theory it can "deconvolve" without switches, i.e. not conditioning on an input image.





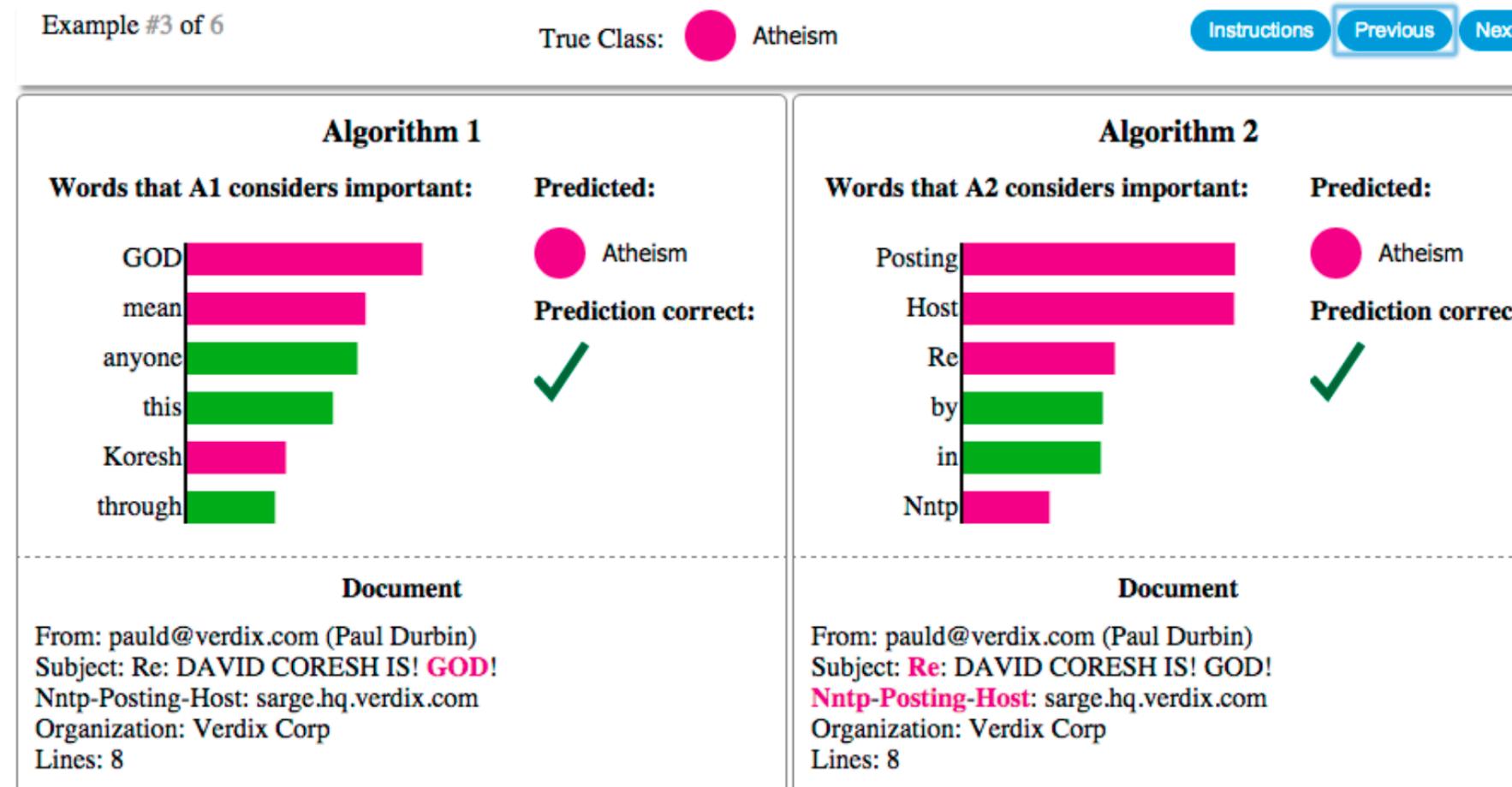
Boulder

“Why Should I Trust You?” Explaining the Predictions of Any Classifier



[YouTube Video](#)

model trained on uni-grams to differentiate “Christianity” from “Atheism”



(a) Original Image

(b) Explaining *Electric guitar*

(c) Explaining *Acoustic guitar*

(d) Explaining *Labrador*

Local Interpretable Model-agnostic Explanations (LIME)

$x \in \mathbb{R}^d$ → original representation of an instance

image: tensor with three color channels per pixel

text: word embeddings

$x' \in \{0, 1\}^{d'}$ → interpretable representation of an instance

image: “presence” or “absence” of a super-pixel

super-pixel → contiguous patch of similar pixels

text: presence or absence of a word

$G \rightarrow$ class of interpretable models

e.g., linear models & decision trees

$g \in G \rightarrow$ explanation

The domain of g is $\{0, 1\}^{d'}$

$\Omega(g) \rightarrow$ measure of complexity of model g

e.g., number of non-zero weights in a linear model or depth of a decision tree

$f : \mathbb{R}^d \rightarrow \mathbb{R}$ $f(x) \rightarrow$ prob. that x belongs to a certain class
 \frown model to be explained

$\pi_x(z) \rightarrow$ proximity measure of an instance z to x

$\mathcal{L}(f, g, \pi_x) \rightarrow$ how unfaithful g is in approximating f in the locality defined by π_x

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

Sparse Linear Explanations

$$g(z') = w_g \cdot z'$$

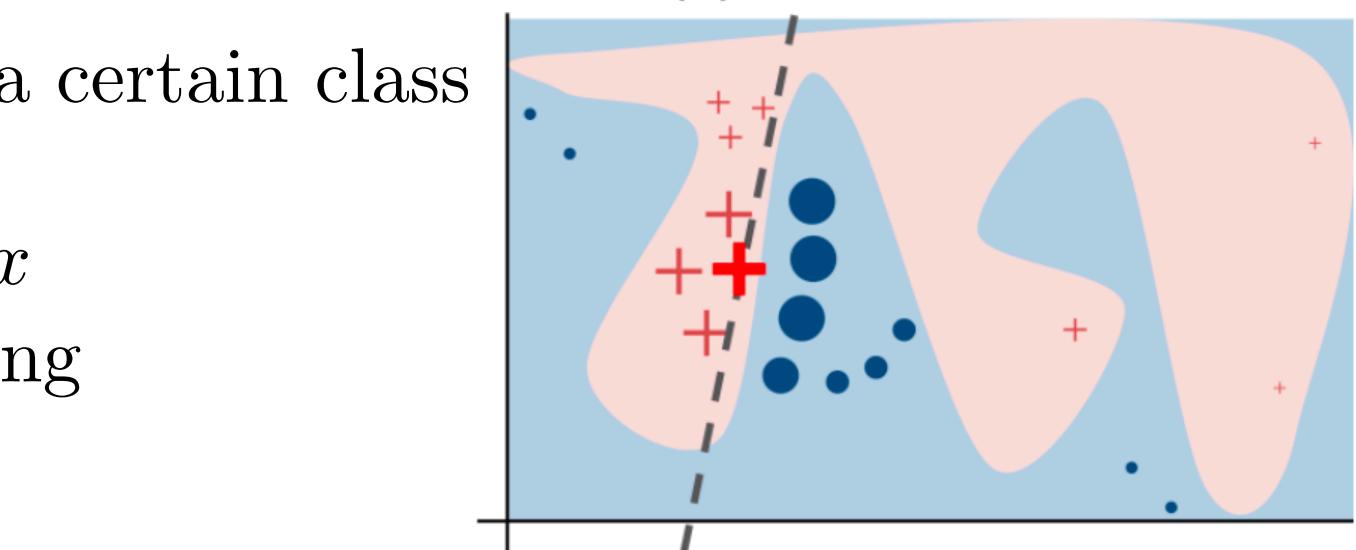
$$\Omega(g) = \infty \mathbf{1}[\|w_g\|_0 > K]$$

$$\pi_x(z) = \exp(-D(x, z)^2 / \sigma^2)$$

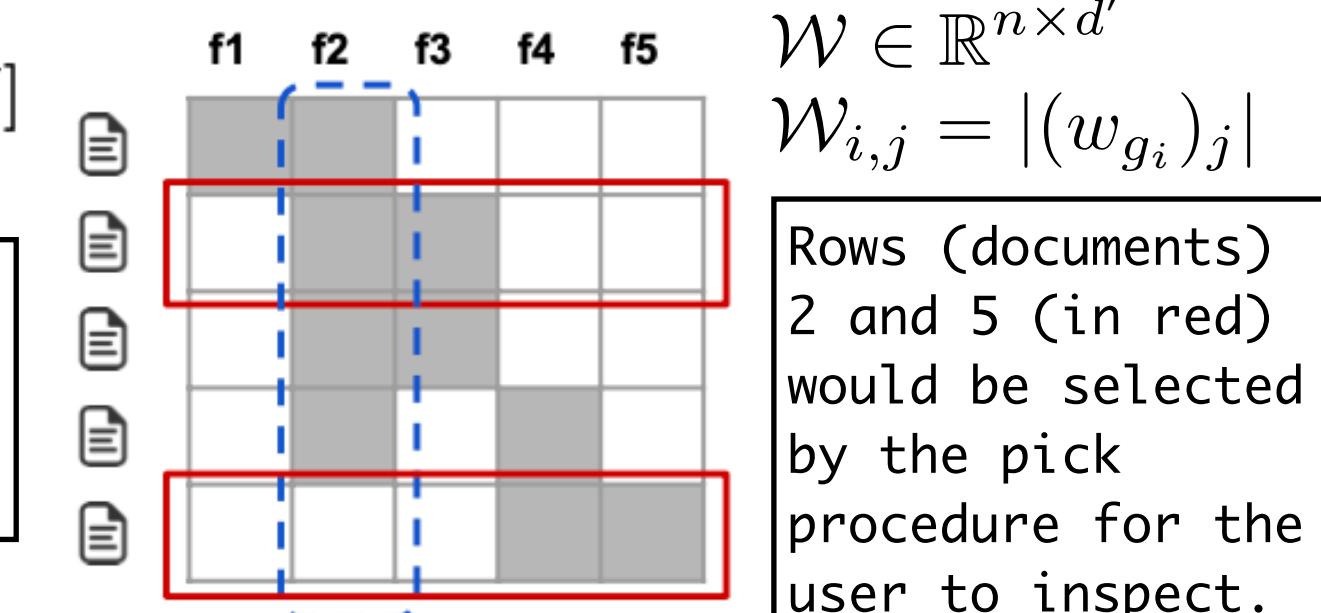
$D \rightarrow$ cosine distance for text or L_2 distance for images

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

limit on the number of words or super-pixels



Trusting a model v.s. trusting a prediction





Boulder



[YouTube Playlist](#)

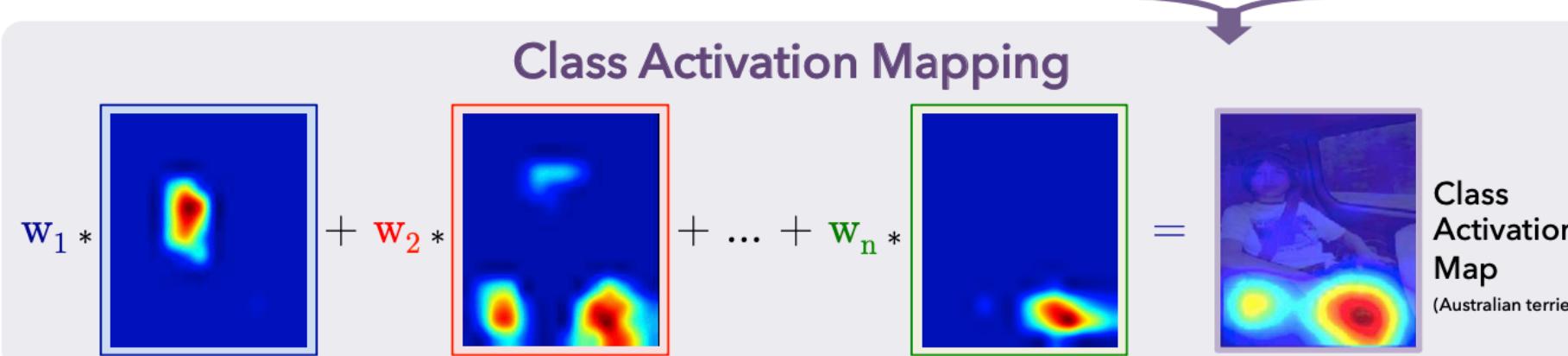
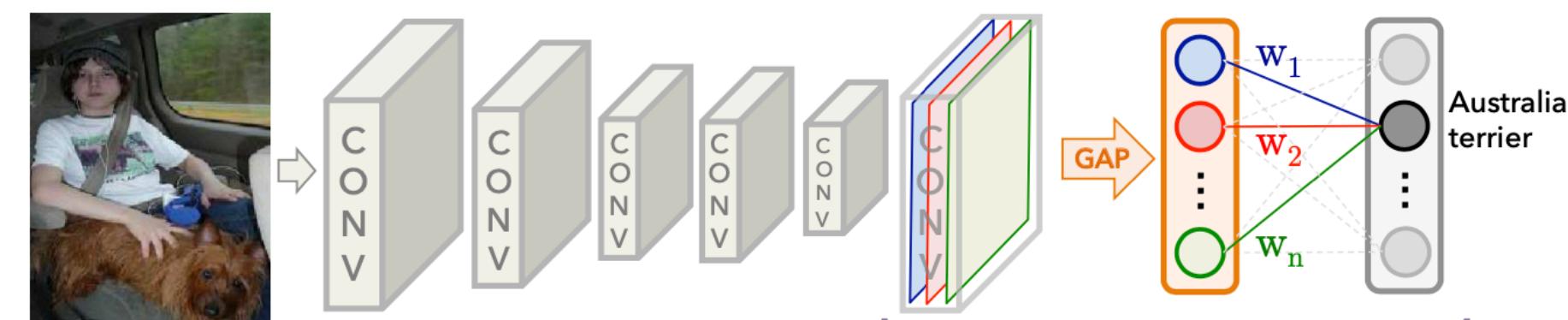
Learning Deep Features for Discriminative Localization

A class activation map (CAM) for a particular category indicates the discriminative image regions used by the CNN to identify that category.

Brushing teeth



Cutting trees



$f_k(x, y) \rightarrow$ activation of unit k in the last convolutional layer at spatial location (x, y)

$$F_k = \sum_{x,y} f_k(x, y)$$

global average pooling (GAP)

$$S_c = \sum_k w_k^c F_k$$

$S_c \rightarrow$ input to the softmax for a given class c

$w_k^c \rightarrow$ weight corresponding to class c for unit k
(importance of F_k for class c)

$$P_c = \frac{\exp(S_c)}{\sum_{c'} \exp(s_{c'})}$$

output of the softmax for class c

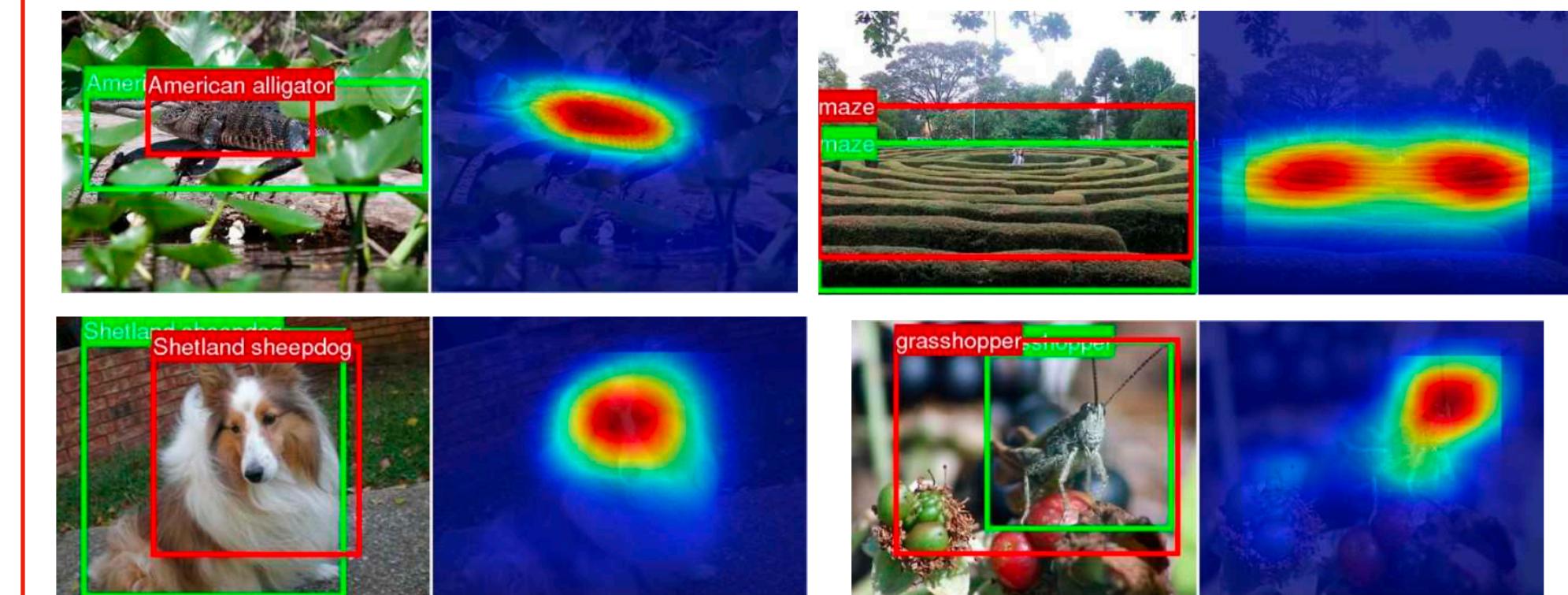
$$S_c = \sum_k w_k^c \underbrace{\sum_{x,y} f_k(x, y)}_{M_c(x,y)}$$

$$= \sum_{x,y} \underbrace{\sum_k w_k^c f_k(x, y)}_{M_c(x,y)}$$

class activation map (CAM)

Importance of activation at spatial grid (x,y) leading to the classification of an image to class c .

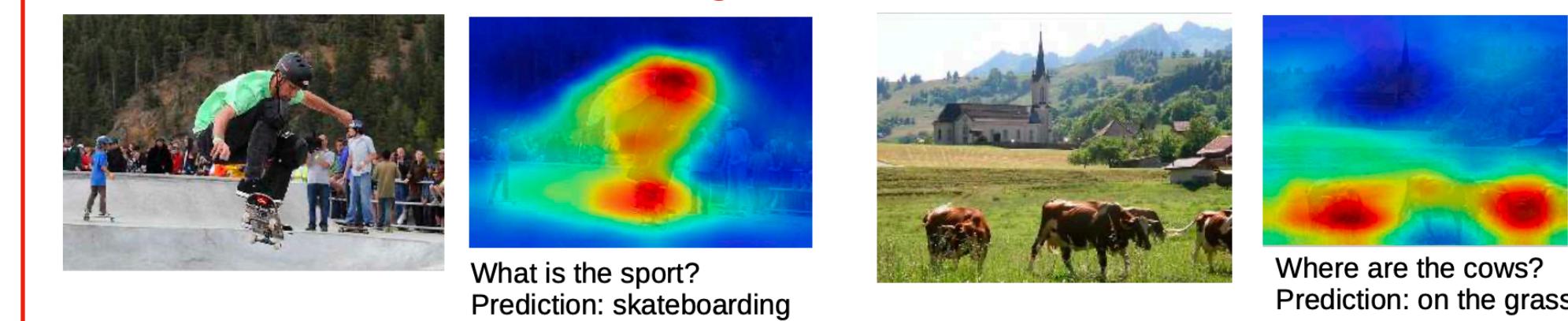
Weakly-supervised Object Localization



Weakly supervised text detector



Visual question answering





Boulder

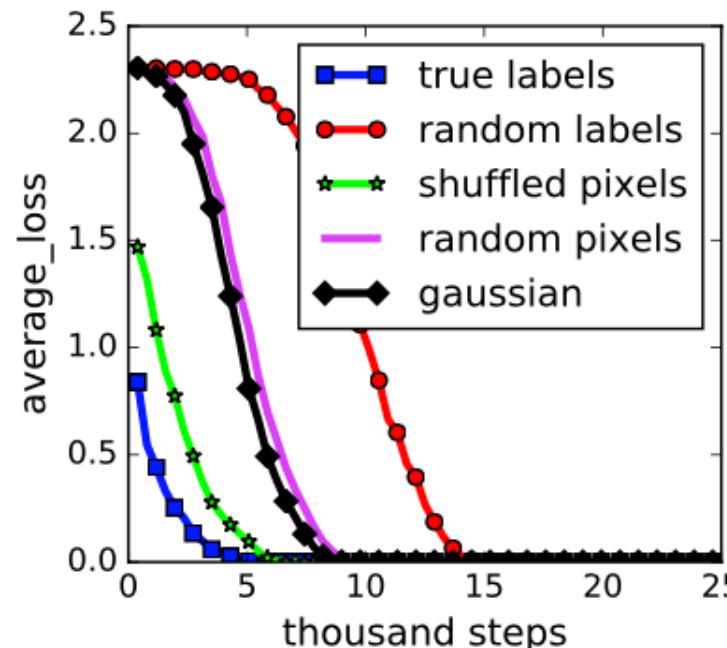
Understanding Deep Learning Requires Rethinking Generalization



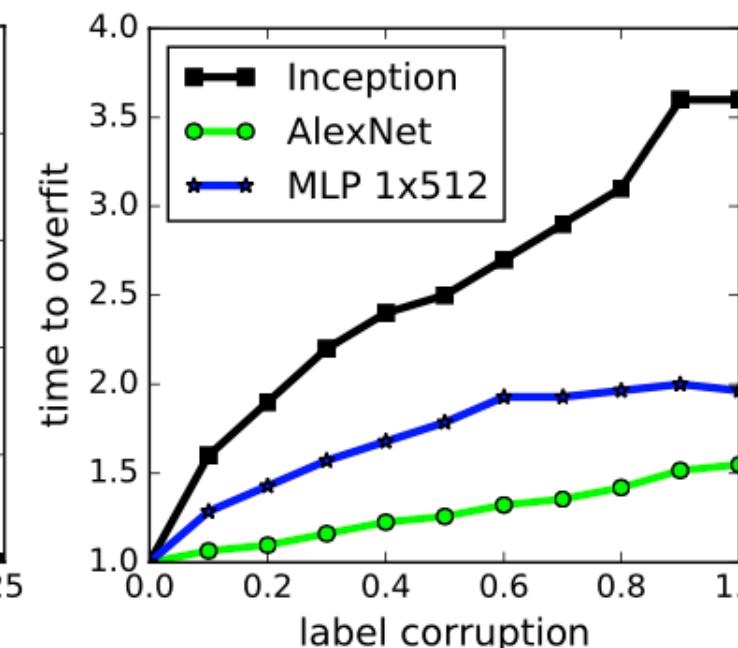
[YouTube Playlist](#)

generalization error: difference btw “training” & “testing” error
Randomization tests

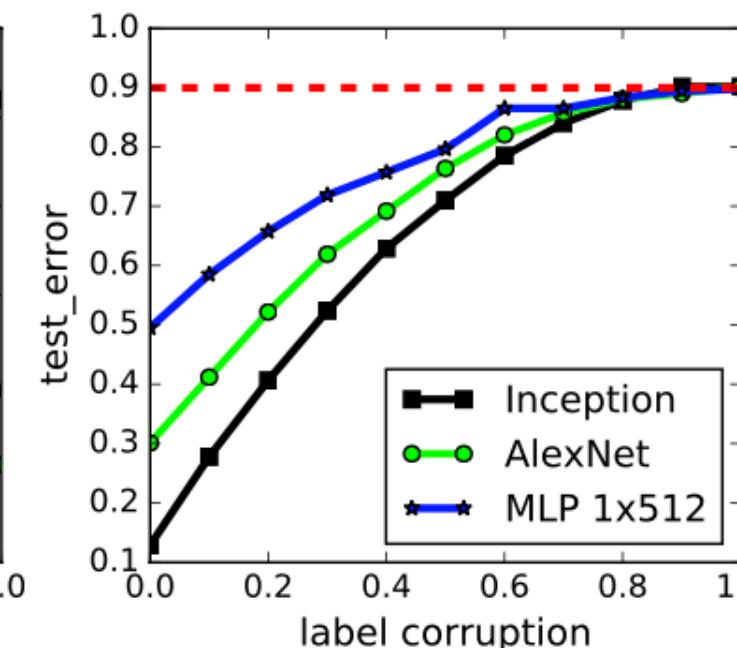
Deep neural networks easily fit random labels.



(a) learning curves



(b) convergence slowdown



(c) generalization error growth

The role of explicit regularization

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes	yes	100.0	89.05
		yes	no	100.0	89.31
		no	yes	100.0	86.03
		no	no	100.0	85.75
		no	no	100.0	9.78
Inception w/o BatchNorm	1,649,402	no	yes	100.0	83.00
		no	no	100.0	82.00
		no	no	100.0	10.12
Alexnet	1,387,786	yes	yes	99.90	81.22
		yes	no	99.82	79.66
		no	yes	100.0	77.36
		no	no	100.0	76.07
		no	no	99.82	9.86
MLP 3x512	1,735,178	no	yes	100.0	53.35
		no	no	100.0	52.39
		no	no	100.0	10.48
MLP 1x512	1,209,866	no	yes	99.80	50.39
		no	no	100.0	50.51
		no	no	99.34	10.61

The model architecture itself isn't a sufficient regularizer.

Explicit regularization may improve generalization performance, but is neither necessary nor by itself sufficient for controlling generalization error.

Lemma: For any two interleaving sequences of n real numbers $b_1 < x_1 < \dots < b_n < x_n$, the $n \times n$ matrix $A = [\max(x_i - b_j, 0)]_{ij}$ has full rank. Its smallest eigenvalue is $\min_i(x_i - b_i)$.

proof: A is lower triangular. A is full rank iff all the diagonal entries $\neq 0$.

$$x_i > b_i \implies \max(x_i - b_i, 0) > 0 \implies A \text{ is invertible.}$$

A lower triangular matrix has all its eigenvalues on the main diagonal.

Finite Sample Expressivity

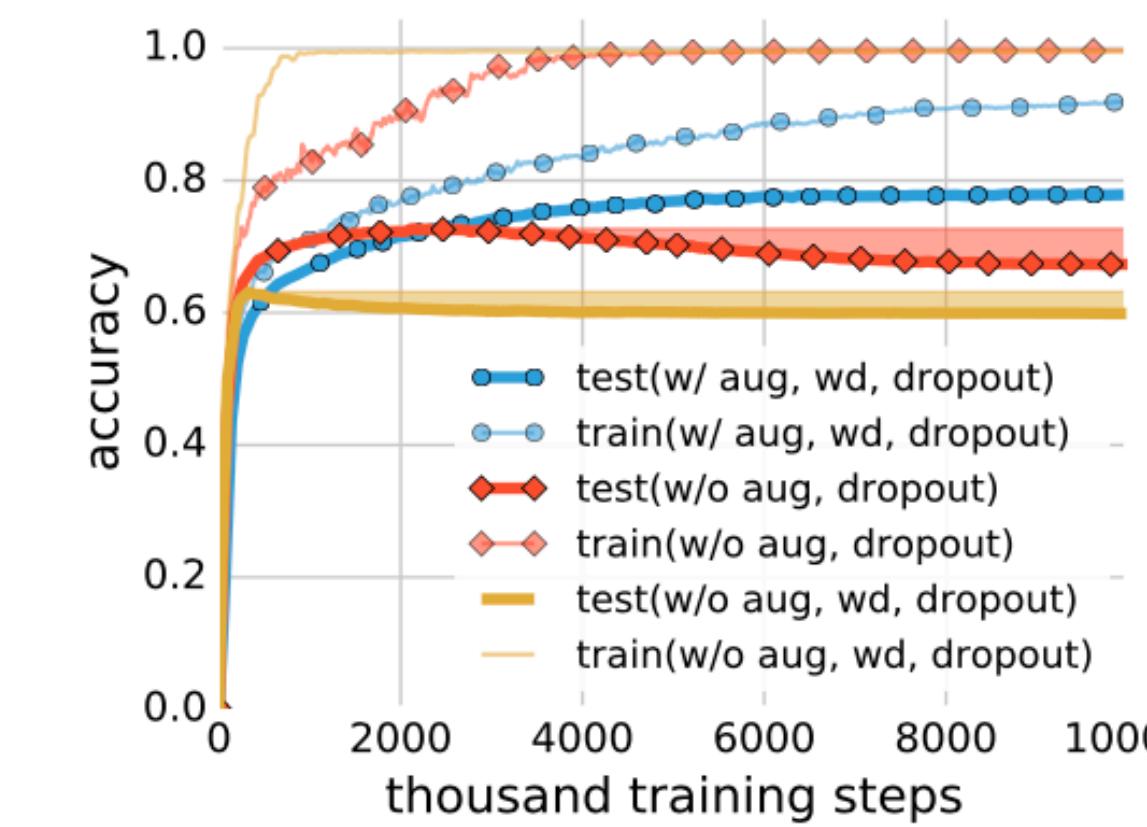
$$c(x) = \sum_{j=1}^n w_j \max(\langle a, x \rangle - b_j, 0)$$

$$w \in \mathbb{R}^n, b \in \mathbb{R}^n, a \in \mathbb{R}^d \quad c : \mathbb{R}^d \rightarrow \mathbb{R}$$

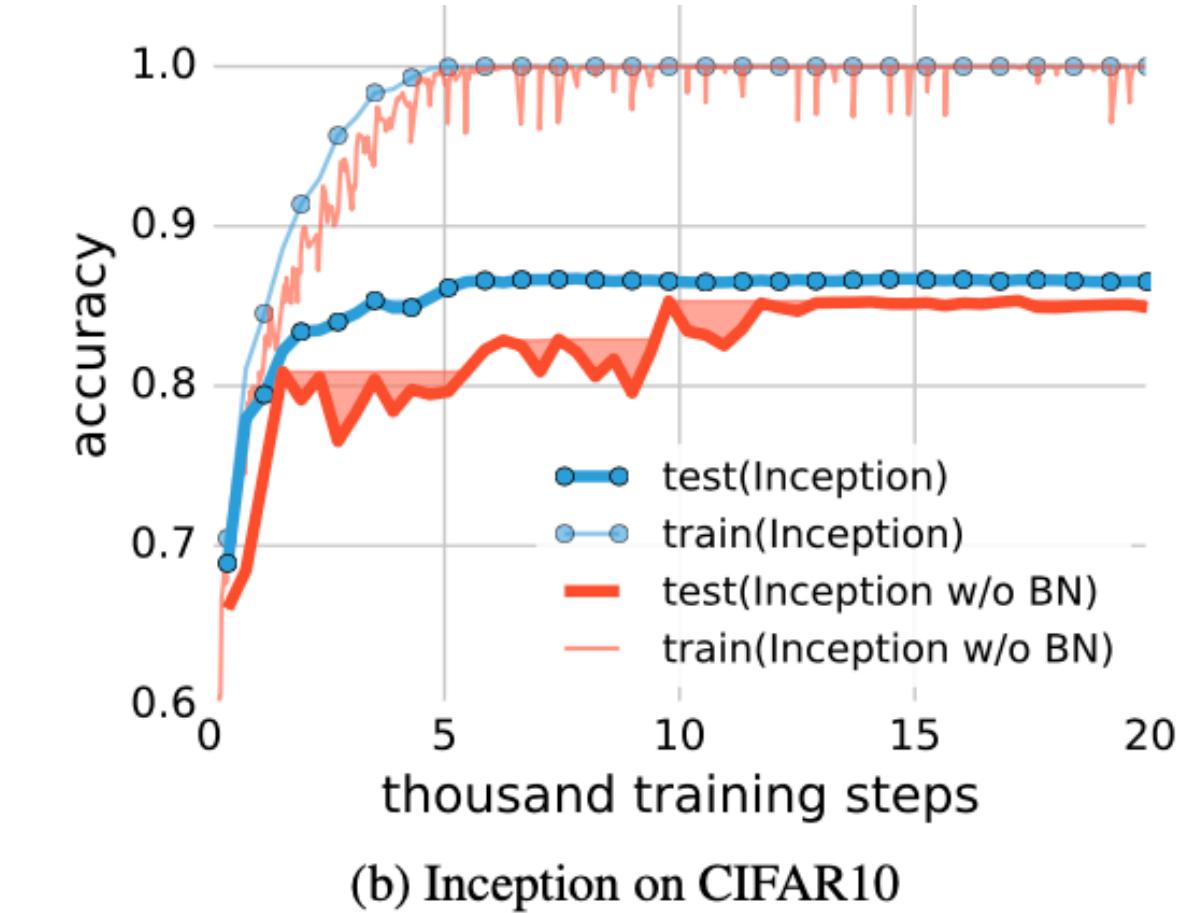
$$S = \{z_1, \dots, z_n\} \quad y \in \mathbb{R}^n \quad \text{Find } w, b, a \text{ such that } y_i = c(z_i) \text{ for all } i = 1, \dots, n.$$

Choose a & b such that $x_i := \langle a, z_i \rangle$ & $b_1 < x_1 < b_2 < x_2 < \dots < b_b < x_n$.

Theorem: There exists a two-layer neural network with ReLU activations and $2n + d$ weights that can represent any function on a sample of size n in d dimensions
 Trade width for depth!

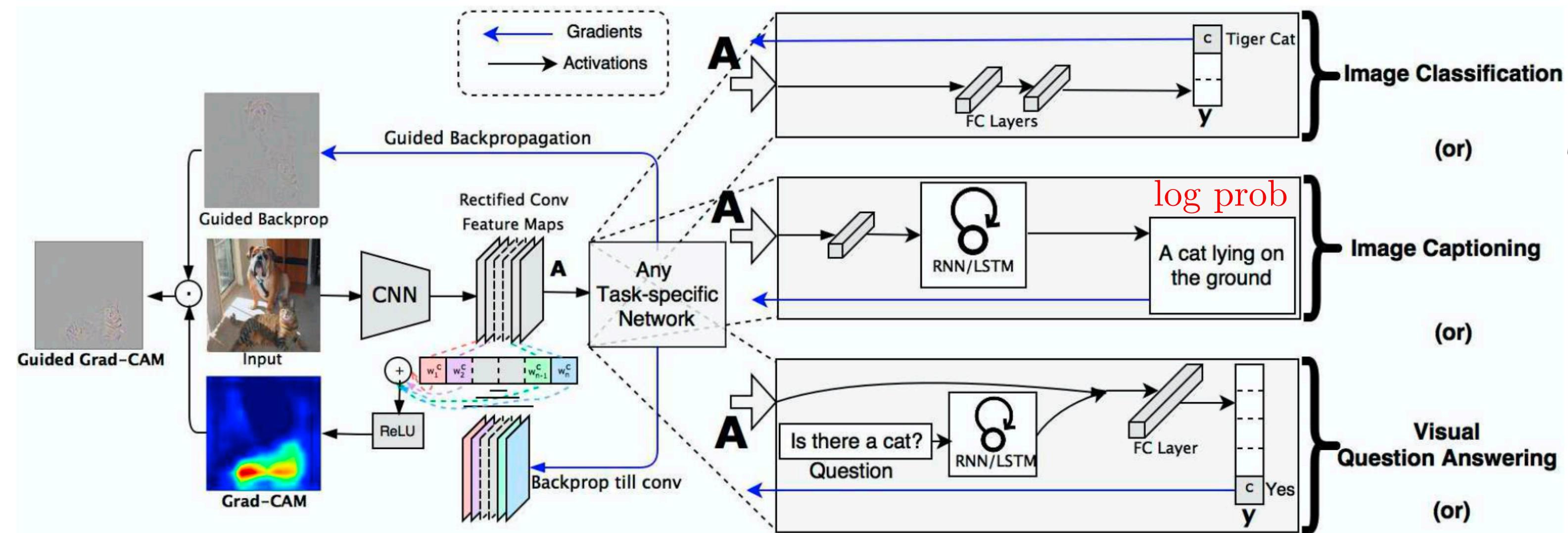


(a) Inception on ImageNet



(b) Inception on CIFAR10

Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization


[YouTube Video](#)


Grad-CAM: Gradient-weighted Class Activation Mapping

$$L_{\text{Grad-CAM}}^c \in \mathbb{R}^{u \times v}$$

class discriminative localization map

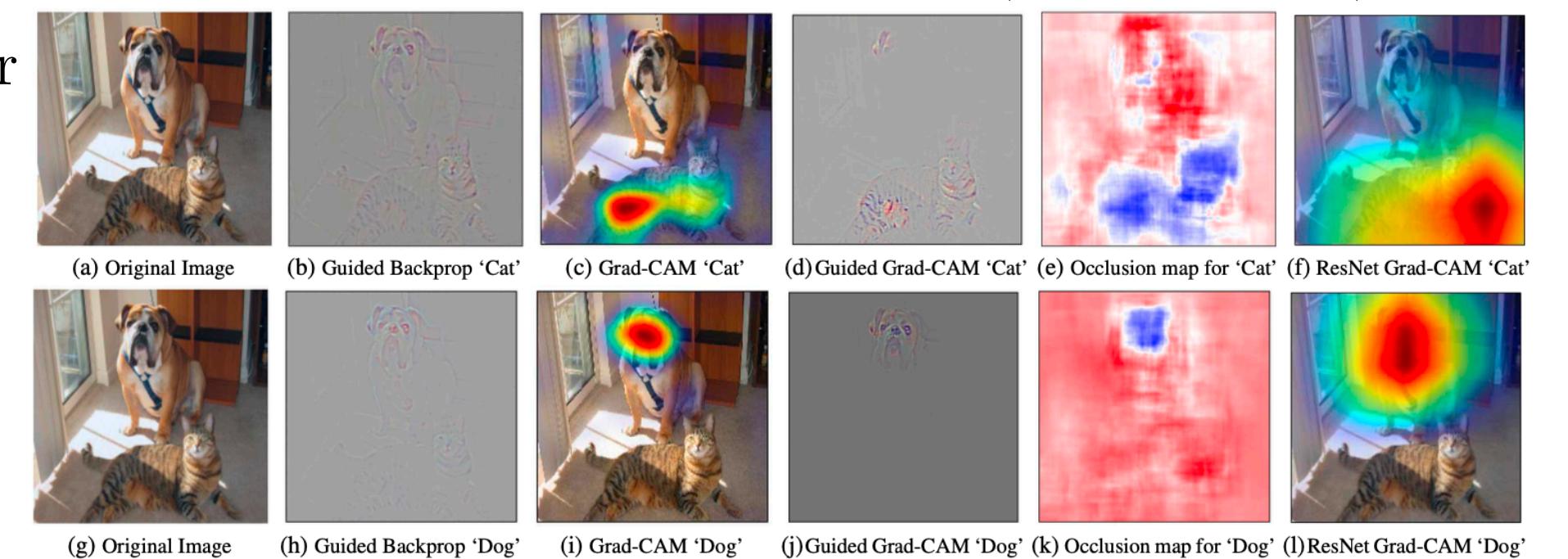
$$\frac{\partial y^c}{\partial A^k}$$

score for class c (before softmax)

feature maps of a conv layer

global average pooling

$$\alpha_k^c = \underbrace{\frac{1}{Z} \sum_i \sum_j}_{\text{importance weights}} \underbrace{\frac{\partial y^c}{\partial A_{ij}^k}}_{\text{gradients via backprop}}$$



$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\sum_k \alpha_k^c A^k \right)$$

positive influence

a coarse heat map (e.g., 14×14)



Boulder

A Unified Approach to Interpreting Model Predictions

accuracy-interpretability tradeoff

Additive Feature Attribution Methods

- LIME
- DeepLIFT
- Layer-Wise Relevance Propagation
- Shapley regression values
- Shapley sampling values
- Quantitative Input Influence

$f \rightarrow$ the original prediction model to be explained

$g \rightarrow$ explanation model

local methods: explain a prediction $f(x)$ based on a single input x

$x' \rightarrow$ simplified inputs

$x = h_x(x') \rightarrow$ mapping function

$g(z') \approx f(h_x(z'))$ whenever $z' \approx x'$

$g(z') = \phi_0 + \sum_{i=1}^M \phi_i z'_i$ where $z' \in \{0, 1\}^M$

$M \rightarrow$ number of simplified input features

SHAP (SHapley Additive exPlanation) Values

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

Property 1: Local Accuracy

$$f(x) = g(x') = \phi_0 + \sum_{i=1}^M \phi_i x'_i$$

$\phi_0 = f(h_x(0)) \rightarrow$ model output with all simplified inputs toggled off (i.e. missing)

Property 2: Missingness

$$x'_i = 0 \implies \phi_i = 0$$

Property 3: Consistency

$$\begin{aligned} f^2(h_x(z')) - f^2(h_x(z' \setminus i)) &\geq f^1(h_x(z')) - f^1(h_x(z' \setminus i)) \\ \implies \phi_i(f^2, x) &\geq \phi_i(f^1, x) \end{aligned}$$

Theorem: g satisfies properties 1, 2, 3 and is unique!

Kernel SHAP (Linear LIME + Shapley values)

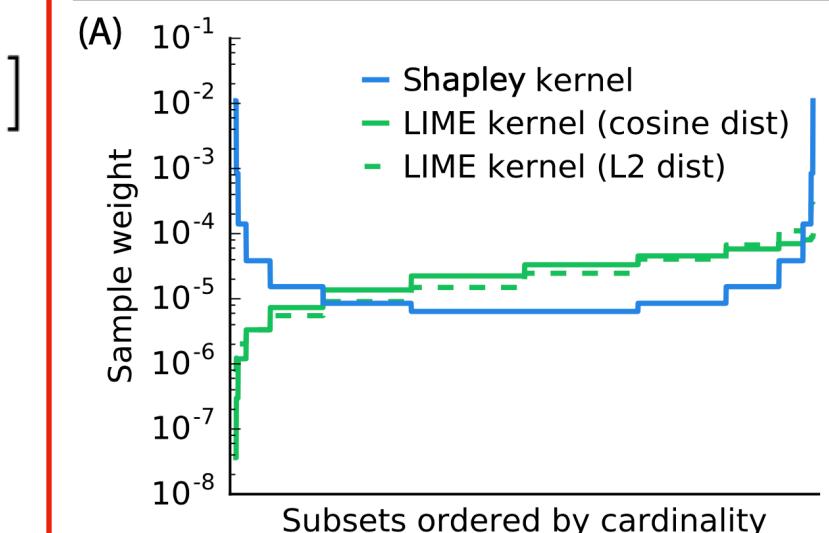
$$\xi(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

$\Omega(g) \rightarrow$ measure of complexity of model g

$\pi_x(z) \rightarrow$ proximity measure of an instance z to x

$\mathcal{L}(f, g, \pi_x) \rightarrow$ how unfaithful g is in approximating f in the locality defined by π_x

$$\begin{aligned} \Omega(g) &= 0, \\ \pi_{x'}(z') &= \frac{(M-1)}{(M \text{ choose } |z'|)|z'|!(M-|z'|)!}, \\ L(f, g, \pi_{x'}) &= \sum_{z' \in Z} [f(h_x(z')) - g(z')]^2 \pi_{x'}(z'), \end{aligned}$$



Linear SHAP

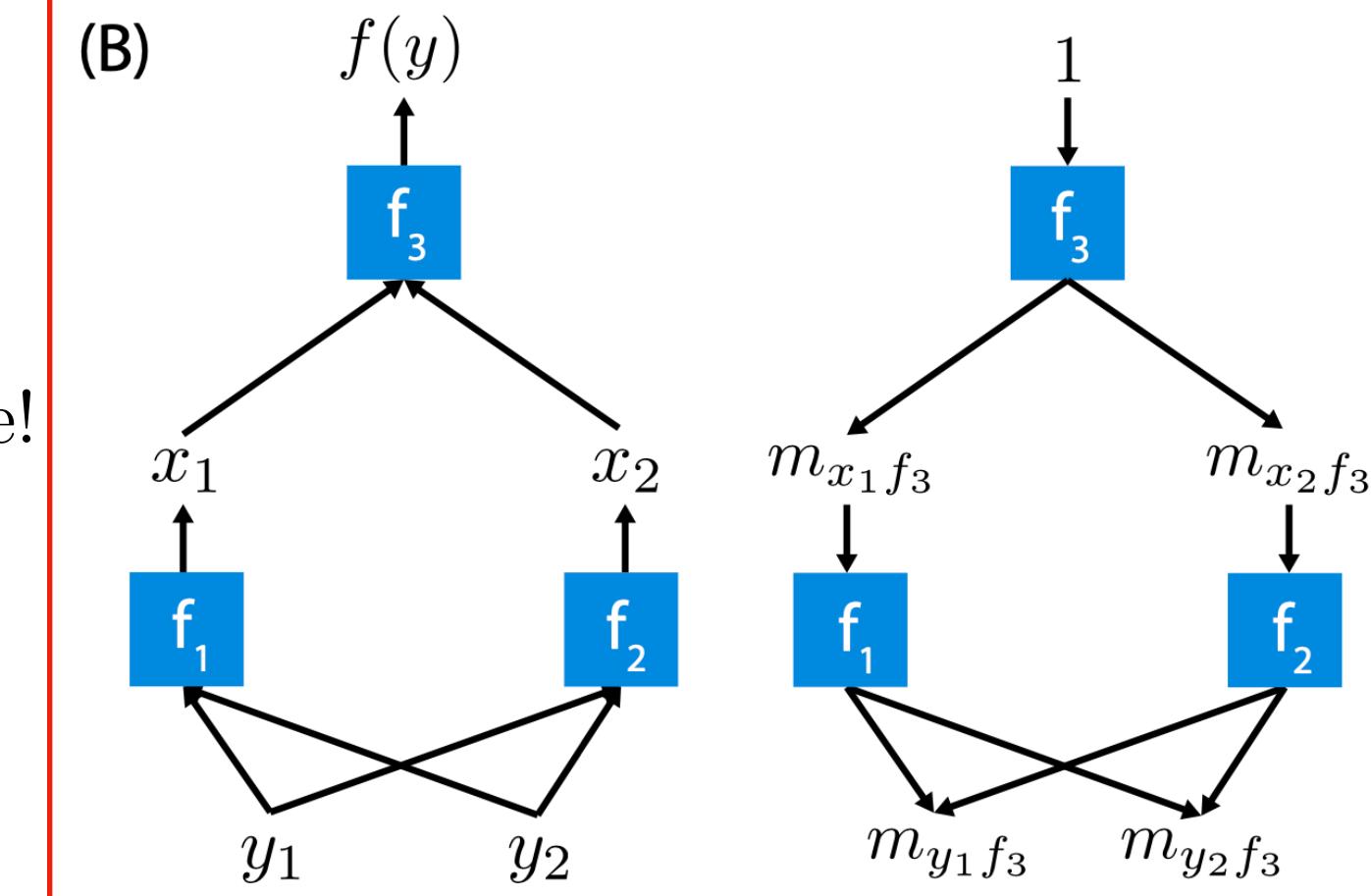
$$f(x) = \sum_{j=1}^M w_j x_j + b$$

$$\phi_0(f, x) = b$$

$$\phi_i(f, x) = w_i(x_i - E[x_i])$$

Deep SHAP (DeepLIFT + Shapley values)

(B)



$$m_{x_j f_3} = \frac{\phi_i(f_3, x)}{x_j - E[x_j]}$$

$$\forall j \in \{1, 2\} \quad m_{y_i f_j} = \frac{\phi_i(f_j, y)}{y_i - E[y_i]}$$

$$m_{y_i f_3} = \sum_{j=1}^2 m_{y_i f_j} m_{x_j f_3}$$

$$\phi_i(f_3, y) \approx m_{y_i f_3} (y_i - E[y_i])$$

The SHAP values for the simple network components can be efficiently solved analytically if they are linear, max pooling, or an activation function with just one input.



Boulder



Questions?

[YouTube Playlist](#)
