



Computer Vision; Image Classification; Visualizing & Understanding



[YouTube Playlist](#)

Maziar Raissi

Assistant Professor

Department of Applied Mathematics

University of Colorado Boulder

maziar.raissi@colorado.edu

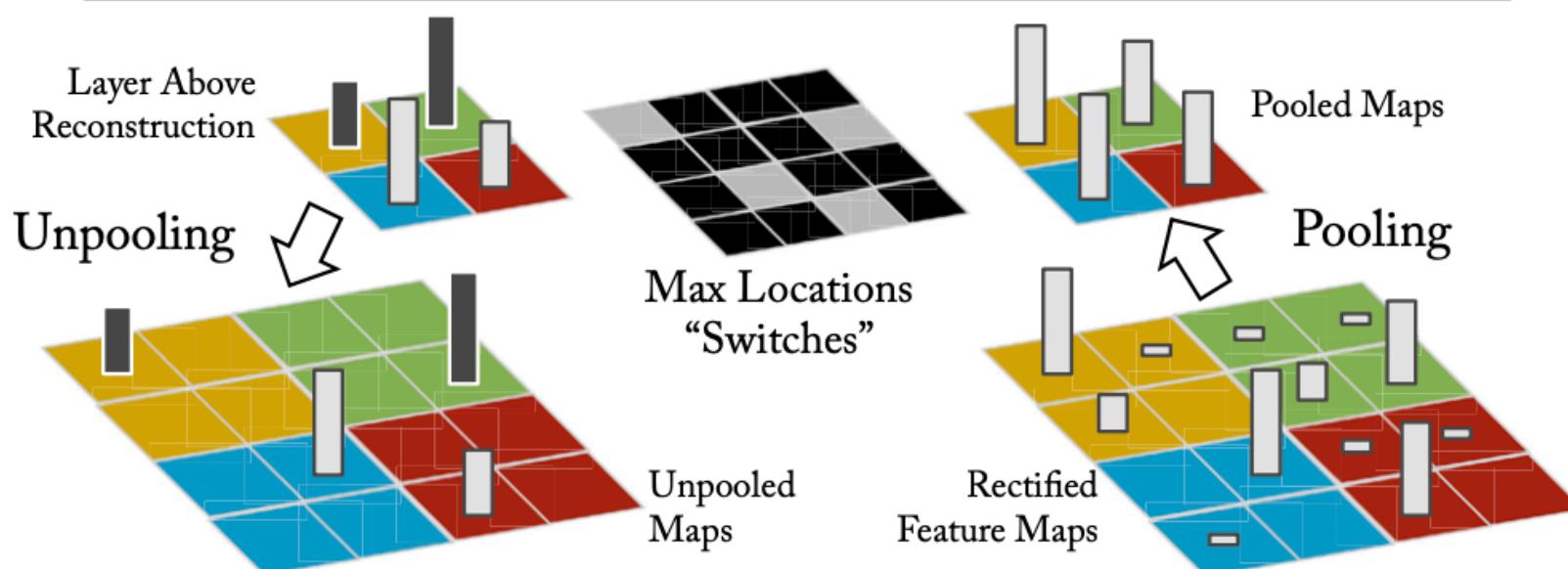
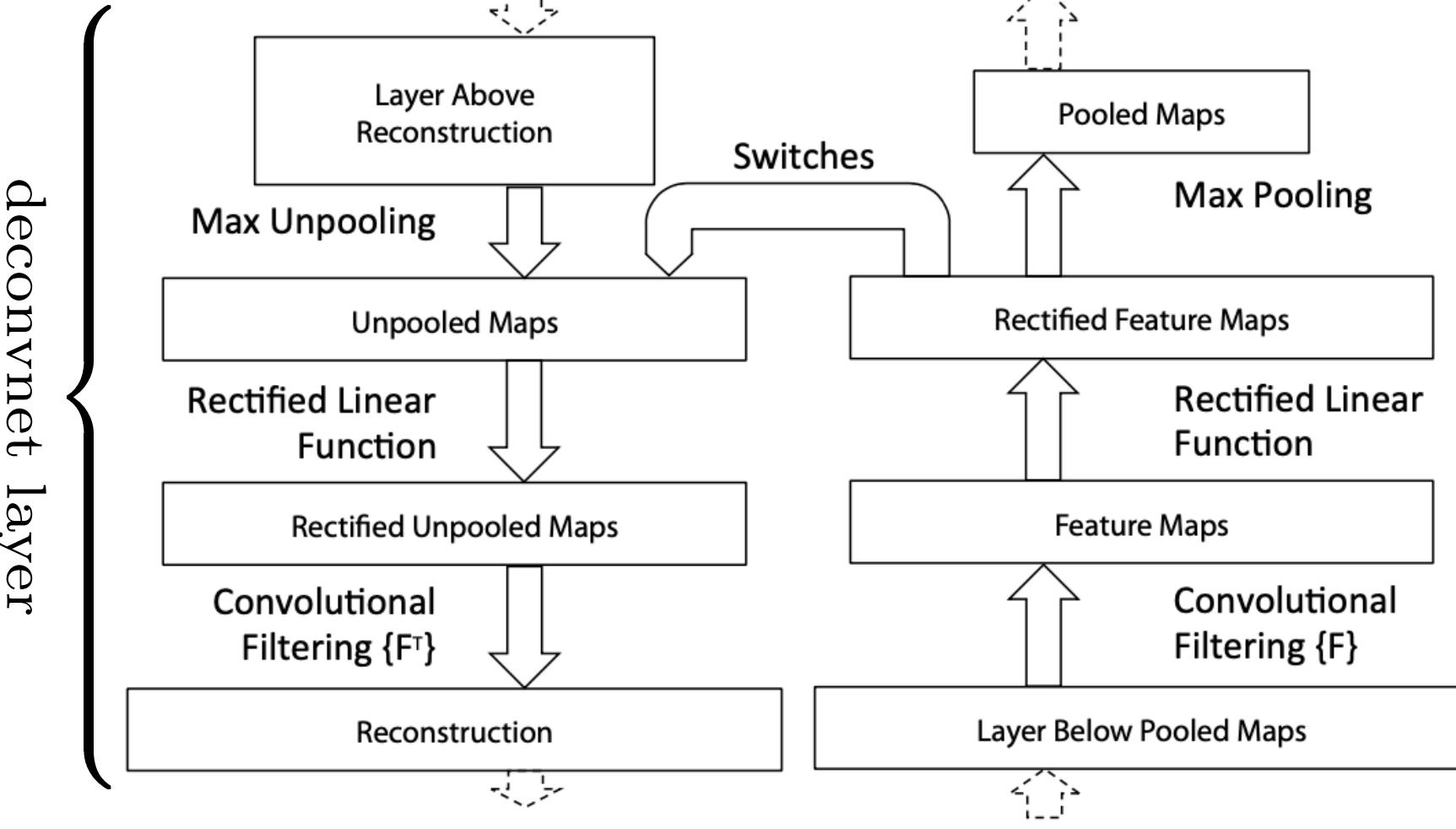
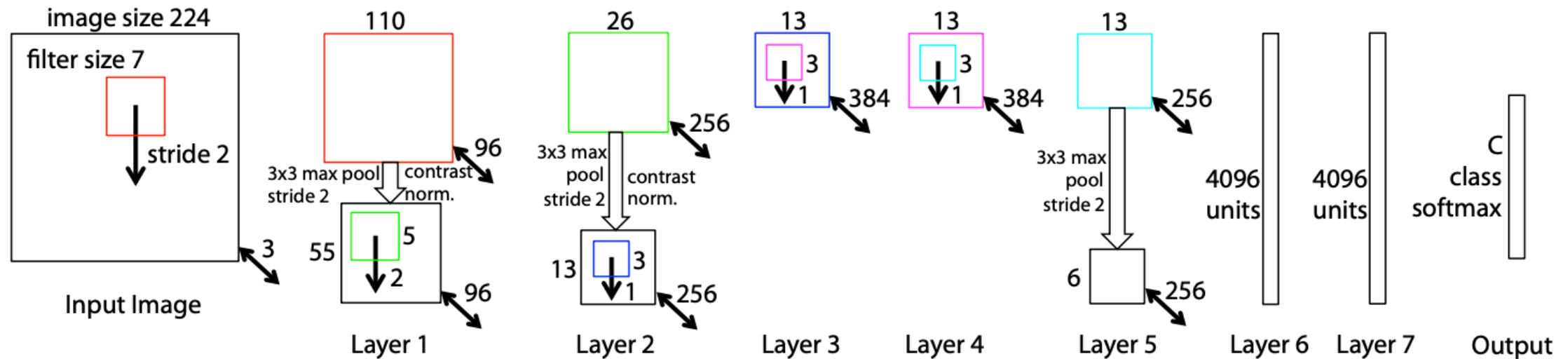


Boulder

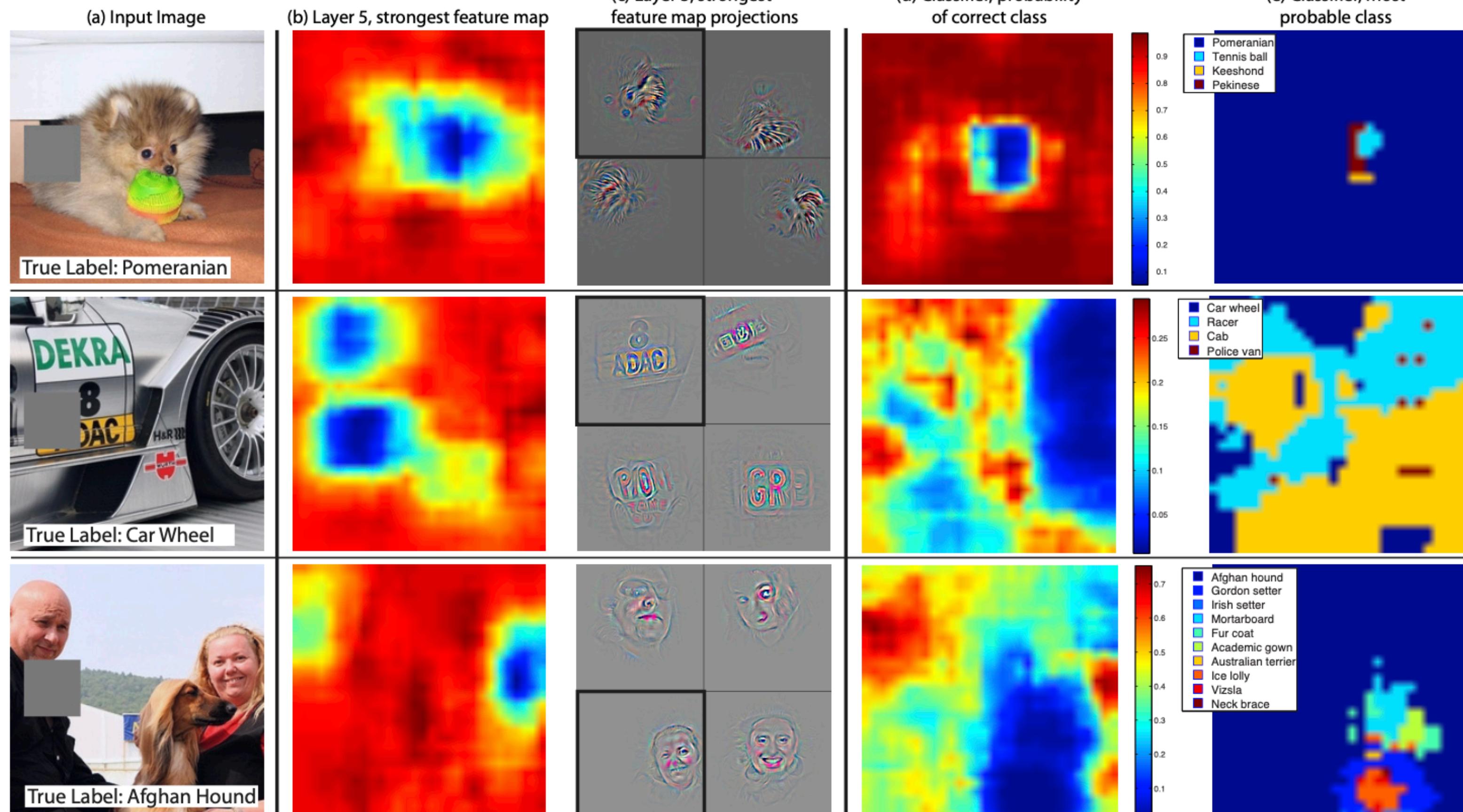
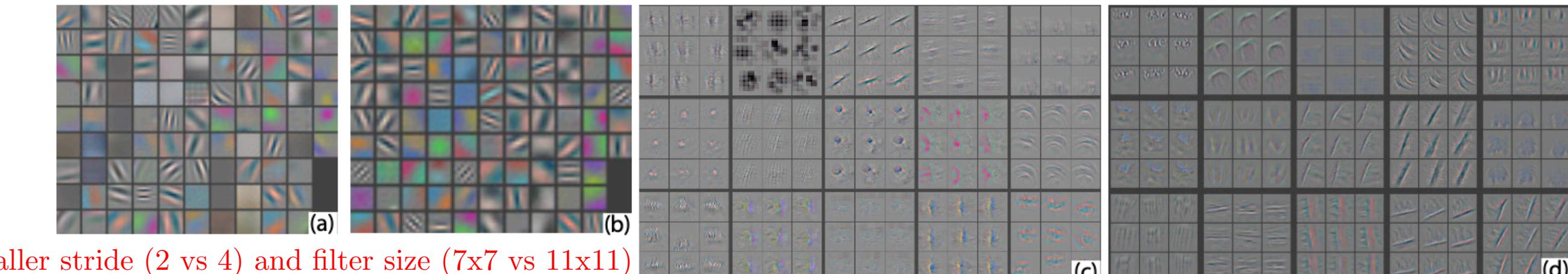


[YouTube Video](#)

Visualizing and Understanding Convolutional Networks



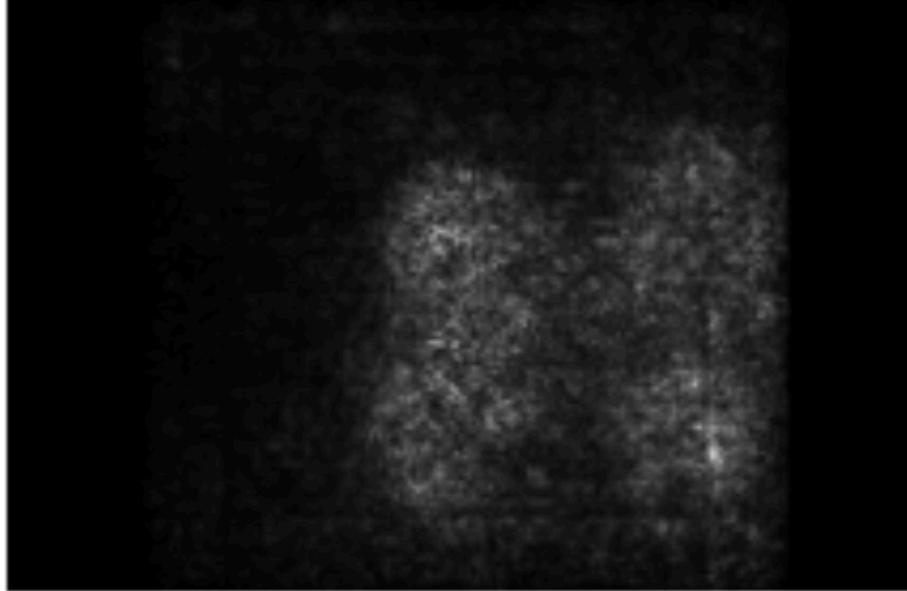
Filtering: Flipping each filter vertically and horizontally



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



goose



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

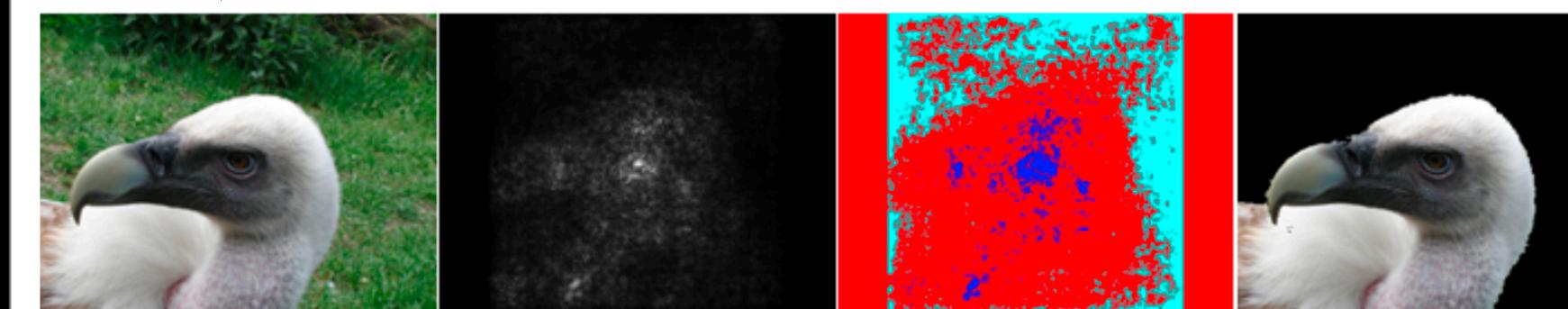
Class Saliency Extraction

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

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

Weakly Supervised Object Localization

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



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

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

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

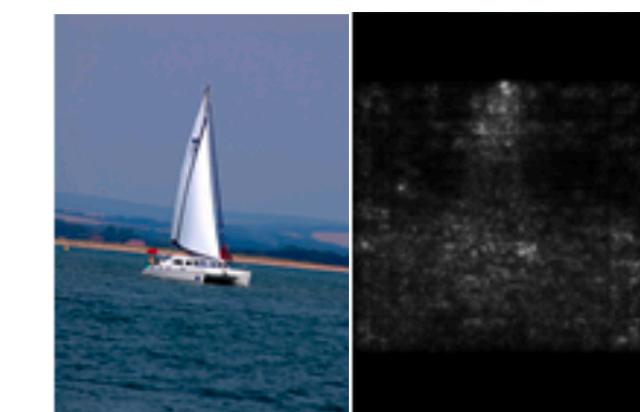
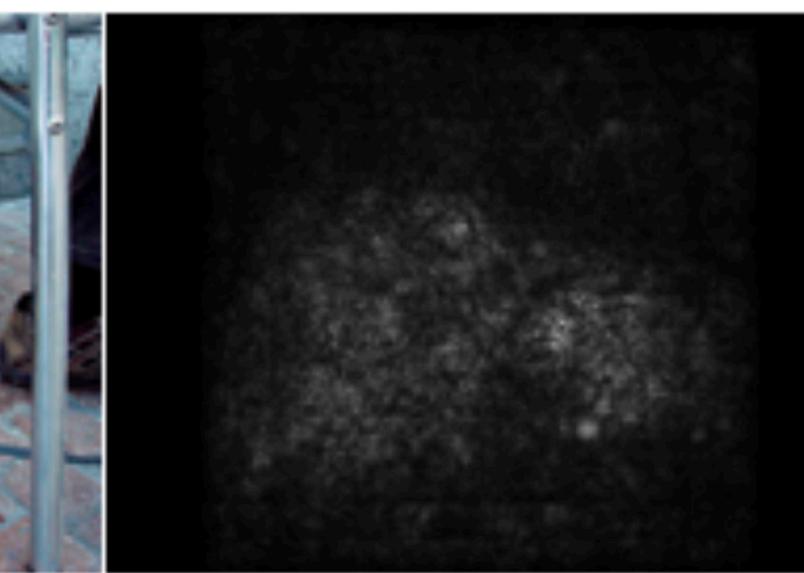
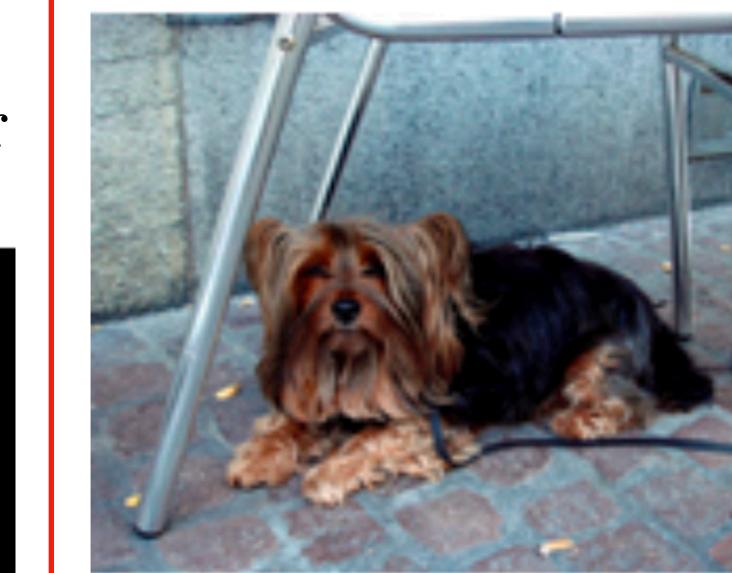
$X_{n+1} = \max(X_n, 0) \rightarrow \text{ReLU}$

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

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

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

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





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

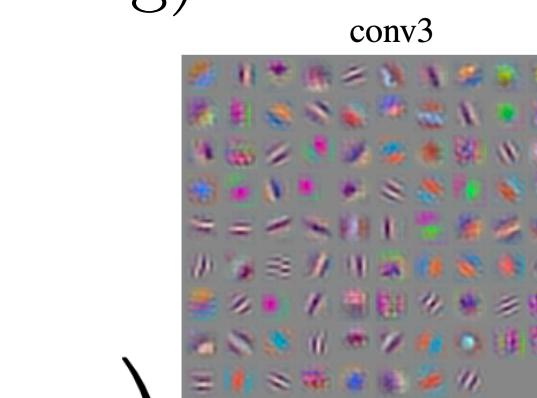
$$c_{i,j,o}(f) = \sigma \left(\sum_{h=-\lfloor k/2 \rfloor}^{\lfloor k/2 \rfloor} \sum_{w=-\lfloor k/2 \rfloor}^{\lfloor k/2 \rfloor} \sum_{u=1}^N \theta_{h,w,u,o} \cdot f_{g(h,w,i,j,u)} \right)$$

convolutional weights (or the kernel weights, or filters)

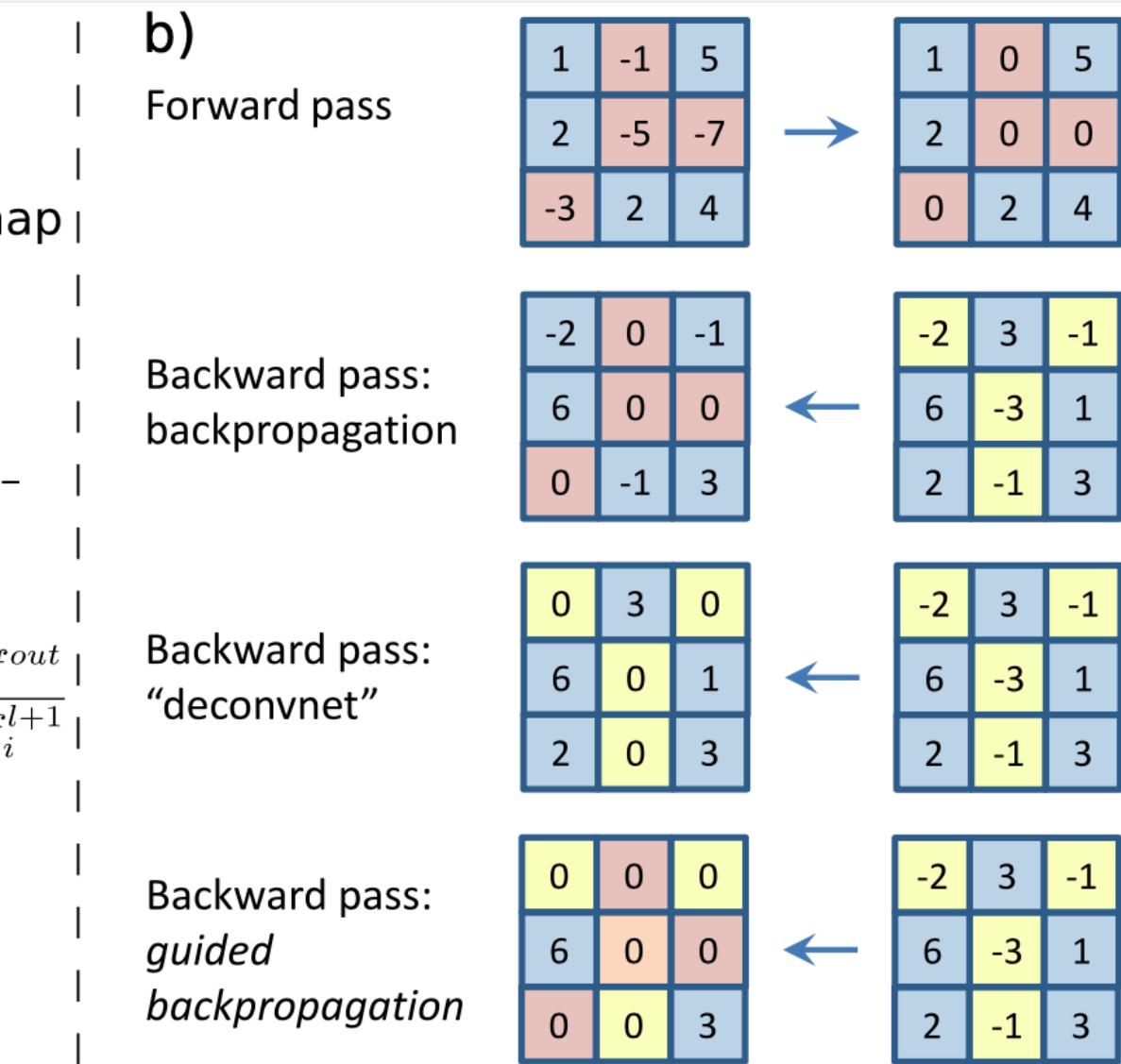
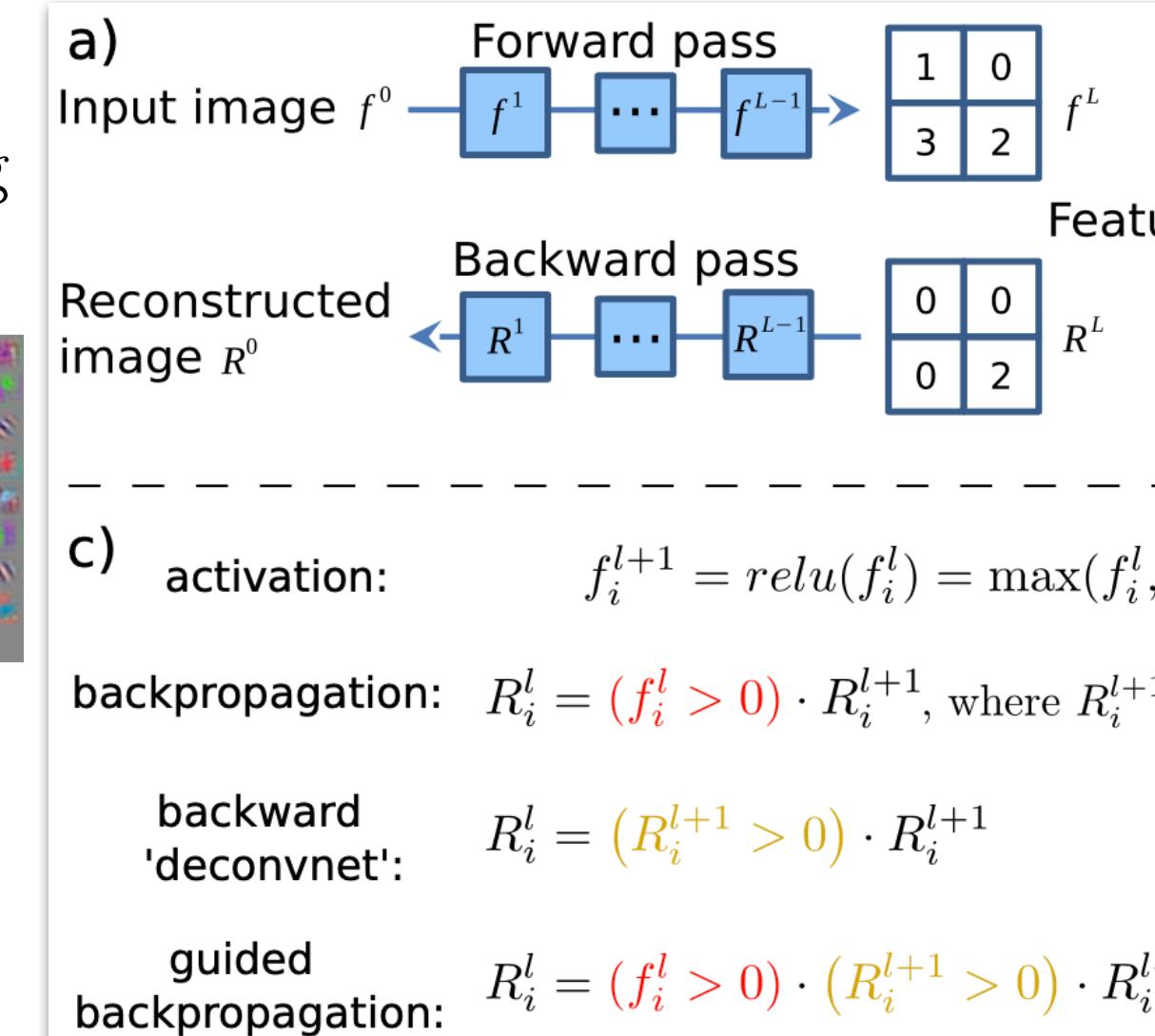
$$\sigma(x) = \max(x, 0), \text{ and } o \in [1, M]$$

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

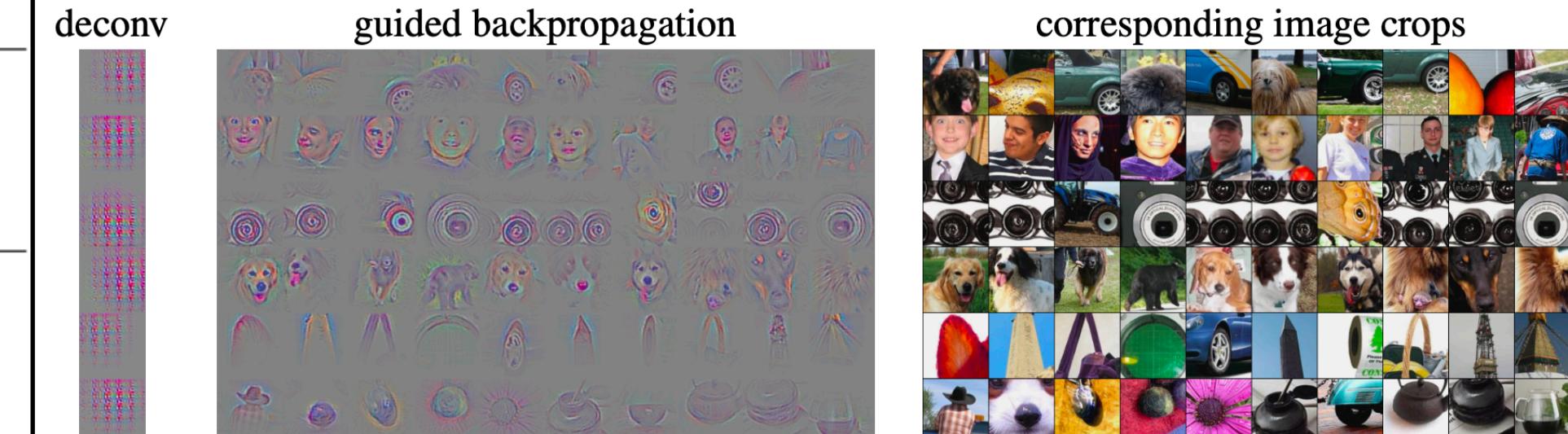
Model		
Strided-CNN-C	ConvPool-CNN-C	All-CNN-C
Input 32×32 RGB image		
3×3 conv. 96 ReLU	3×3 conv. 96 ReLU	3×3 conv. 96 ReLU
3×3 conv. 96 ReLU with stride $r = 2$	3×3 conv. 96 ReLU	3×3 conv. 96 ReLU
	3×3 conv. 96 ReLU with stride $r = 2$	3×3 conv. 96 ReLU with stride $r = 2$
3×3 conv. 192 ReLU	3×3 conv.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





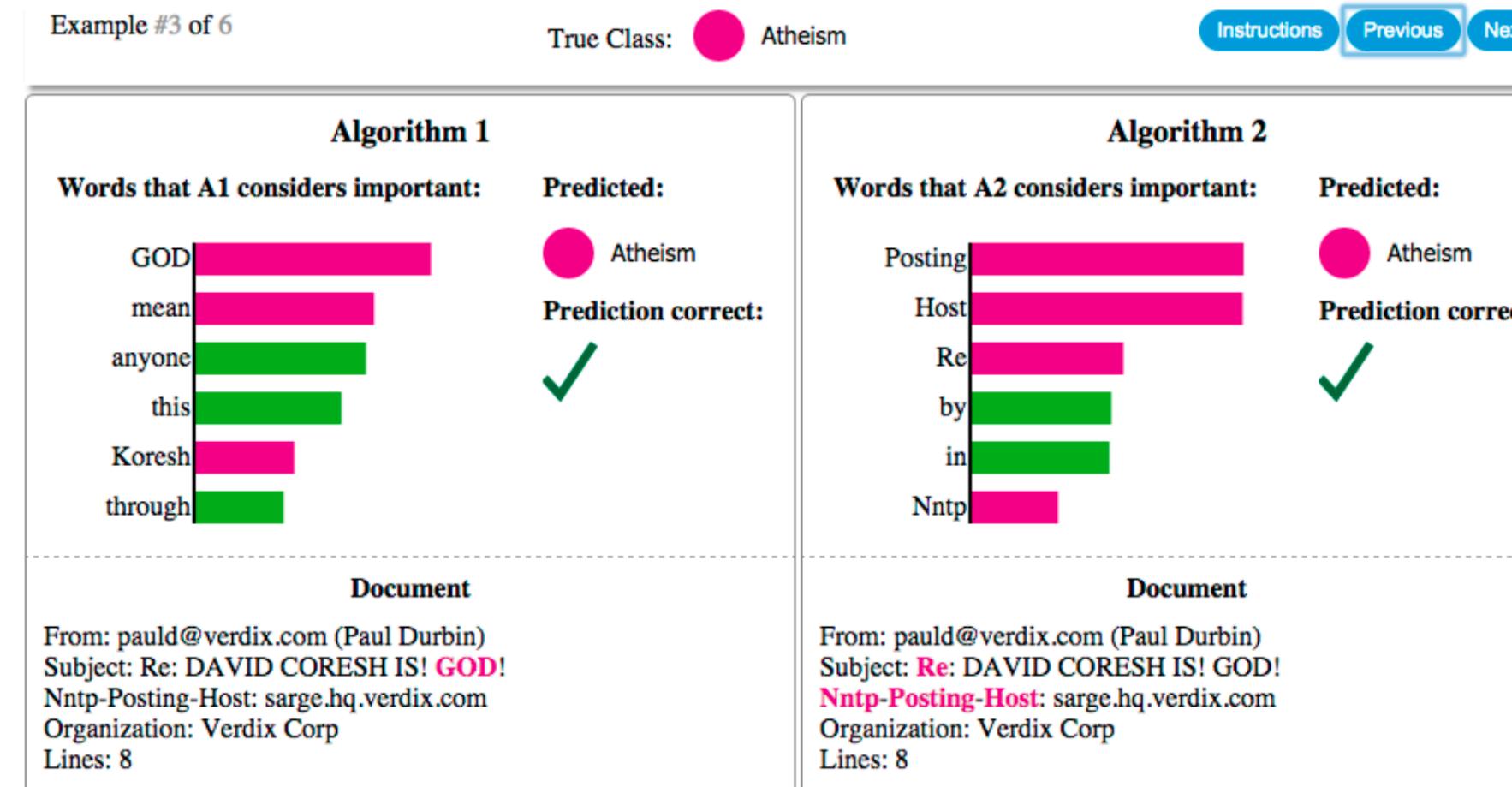
Boulder

“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

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

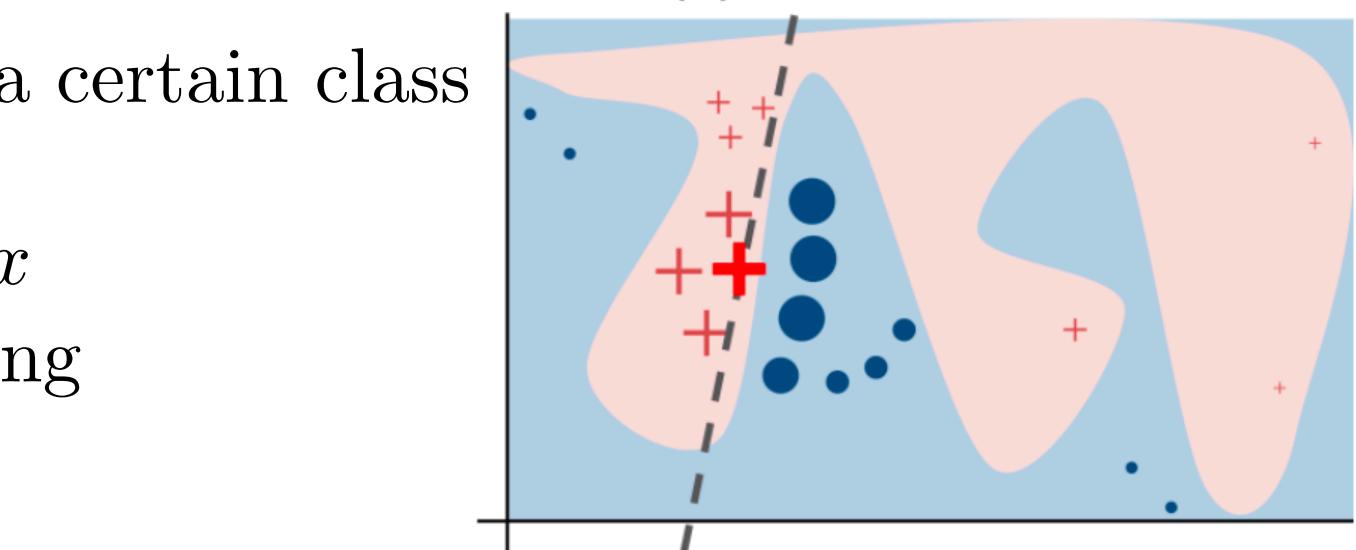
$$\Omega(g) = \infty \mathbf{1}[\|w_g\|_0 > K]$$

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

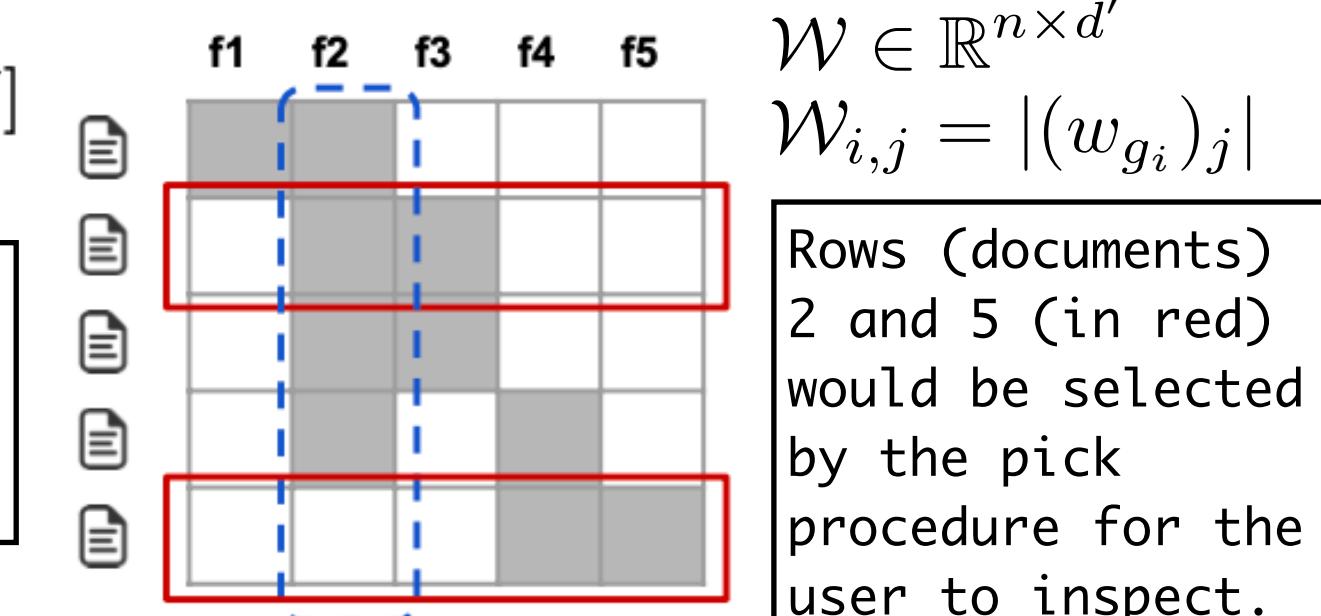
$D \rightarrow$ cosine distance for text or L_2 distance for images

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

limit on the number of words or super-pixels



Trusting a model v.s. trusting a prediction



$$\mathcal{W} \in \mathbb{R}^{n \times d'}$$

$$\mathcal{W}_{i,j} = |(w_{g_i})_j|$$

Rows (documents) 2 and 5 (in red) would be selected by the pick procedure for the user to inspect.



Boulder



[YouTube Playlist](#)

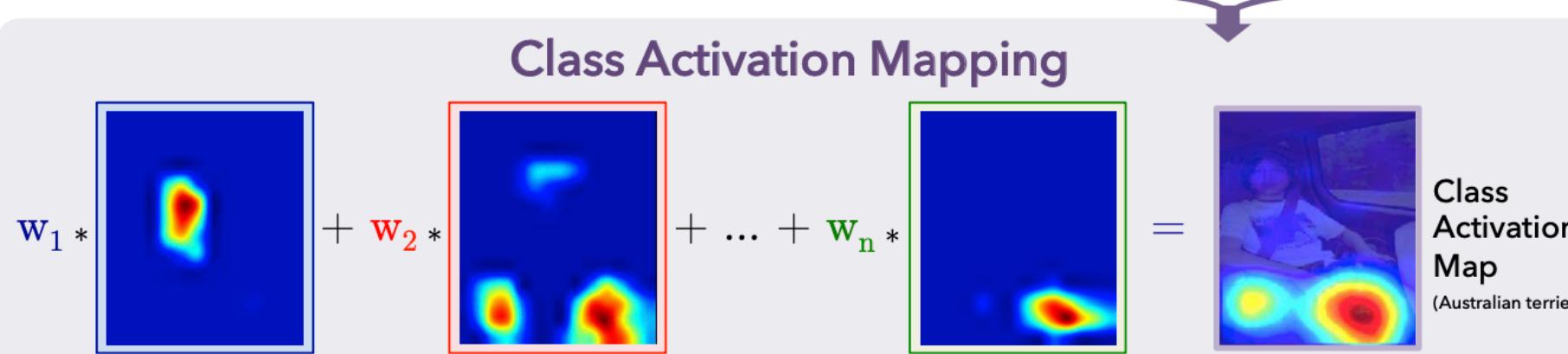
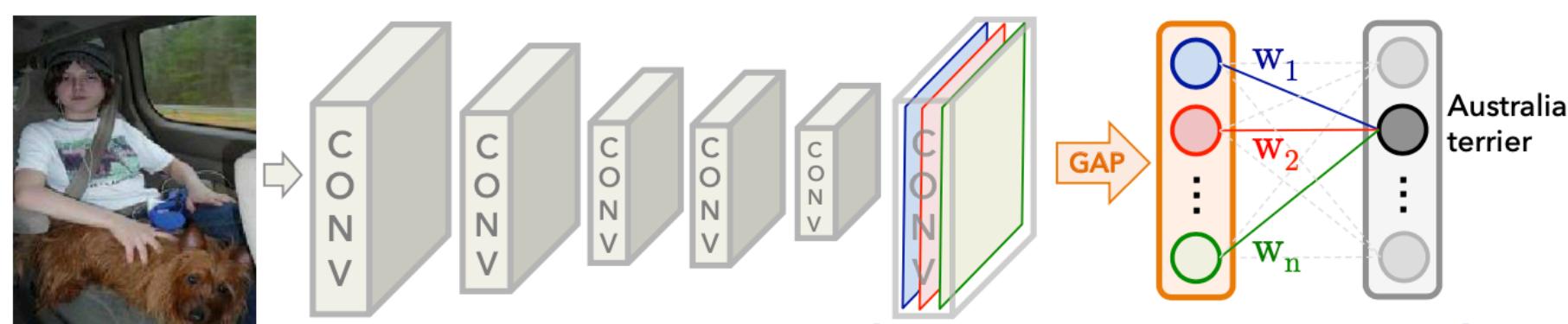
Learning Deep Features for Discriminative Localization

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

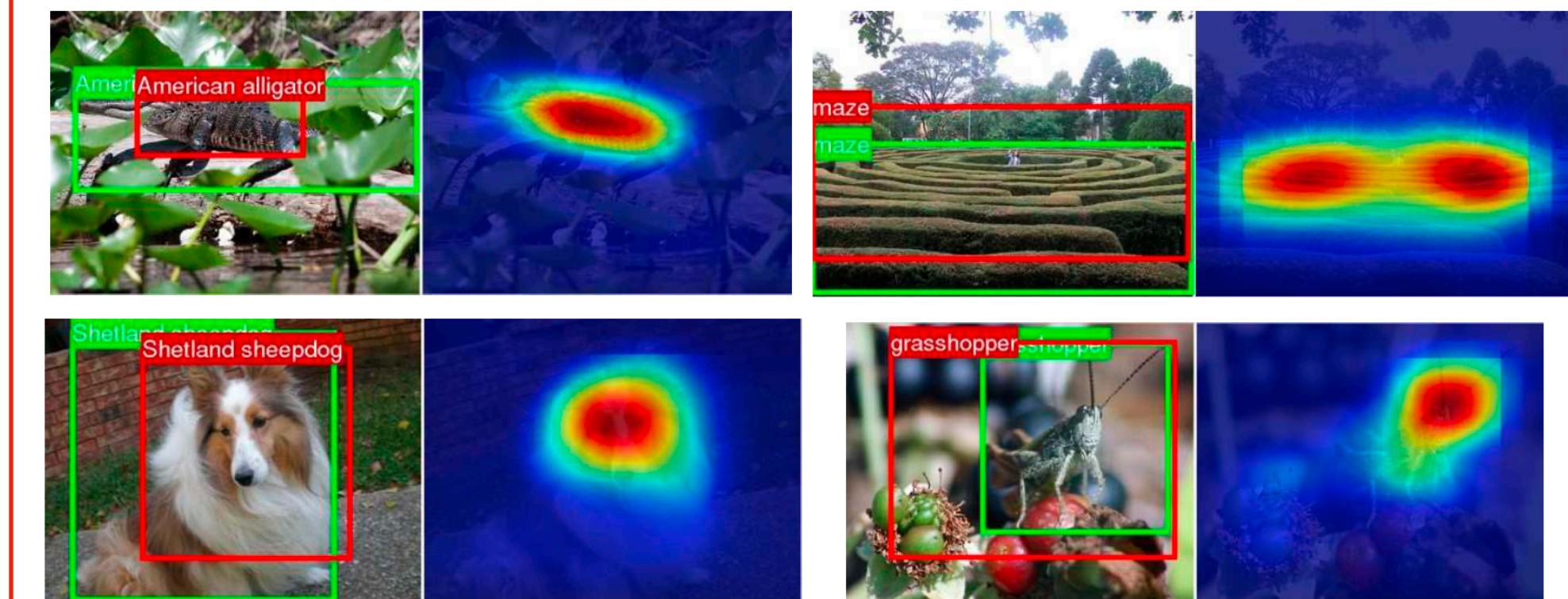
$$S_c = \sum_k w_k^c \underbrace{\sum_{x,y} f_k(x, y)}_{M_c(x,y)}$$

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

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

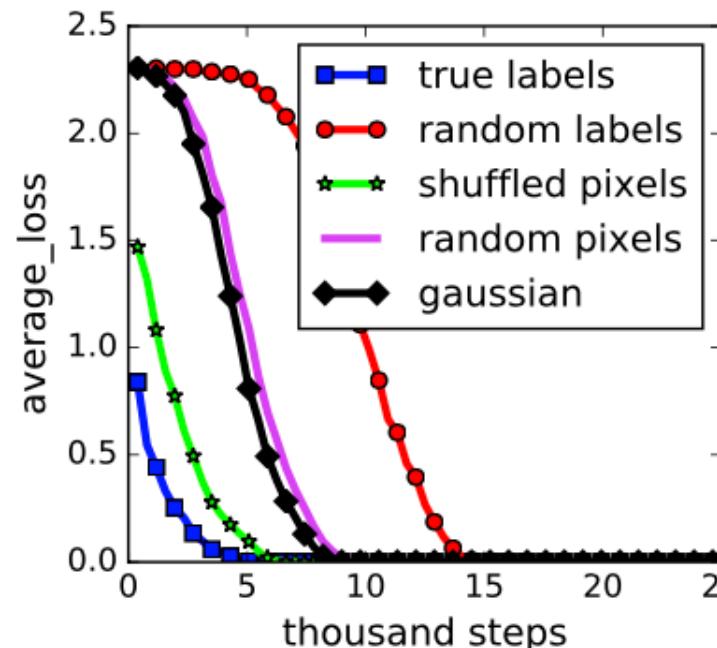
Understanding Deep Learning Requires Rethinking Generalization



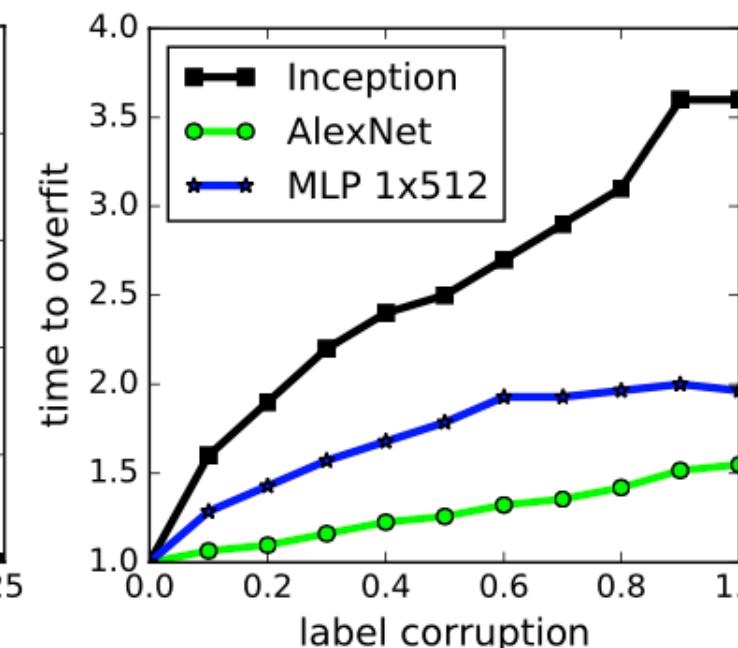
[YouTube Playlist](#)

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

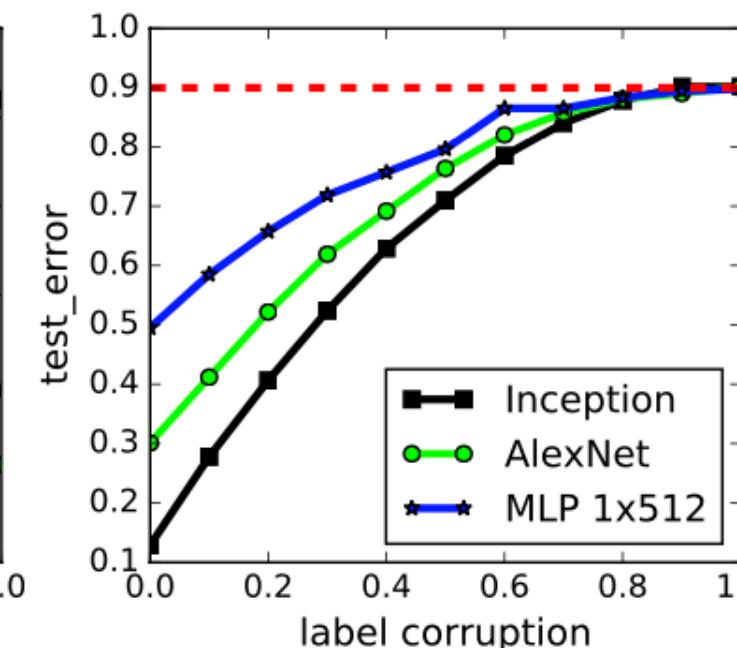
Deep neural networks easily fit random labels.



(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

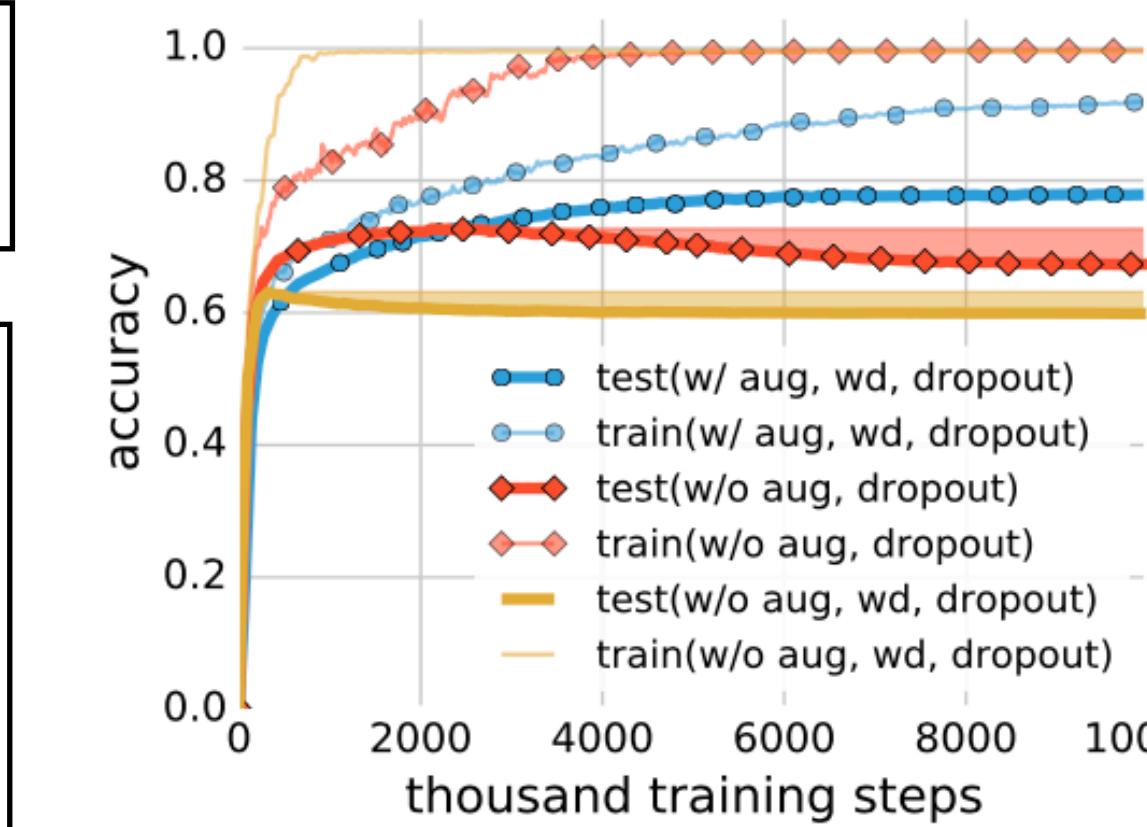
$$c(x) = \sum_{j=1}^n w_j \max(\langle a, x \rangle - b_j, 0)$$

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

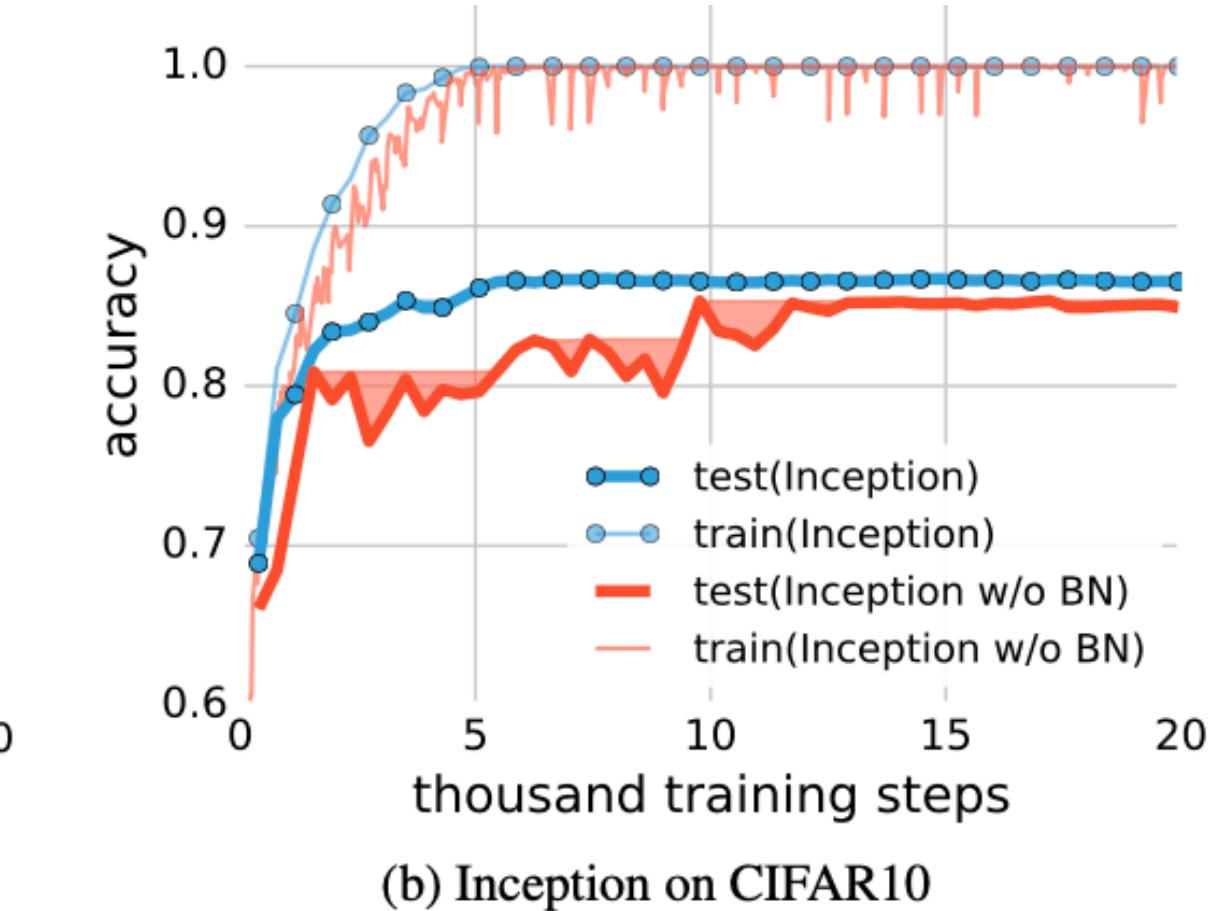
$$S = \{z_1, \dots, z_n\} \quad y \in \mathbb{R}^n \quad \text{Find } w, b, a \text{ such that } y_i = c(z_i) \text{ for all } i = 1, \dots, n.$$

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!

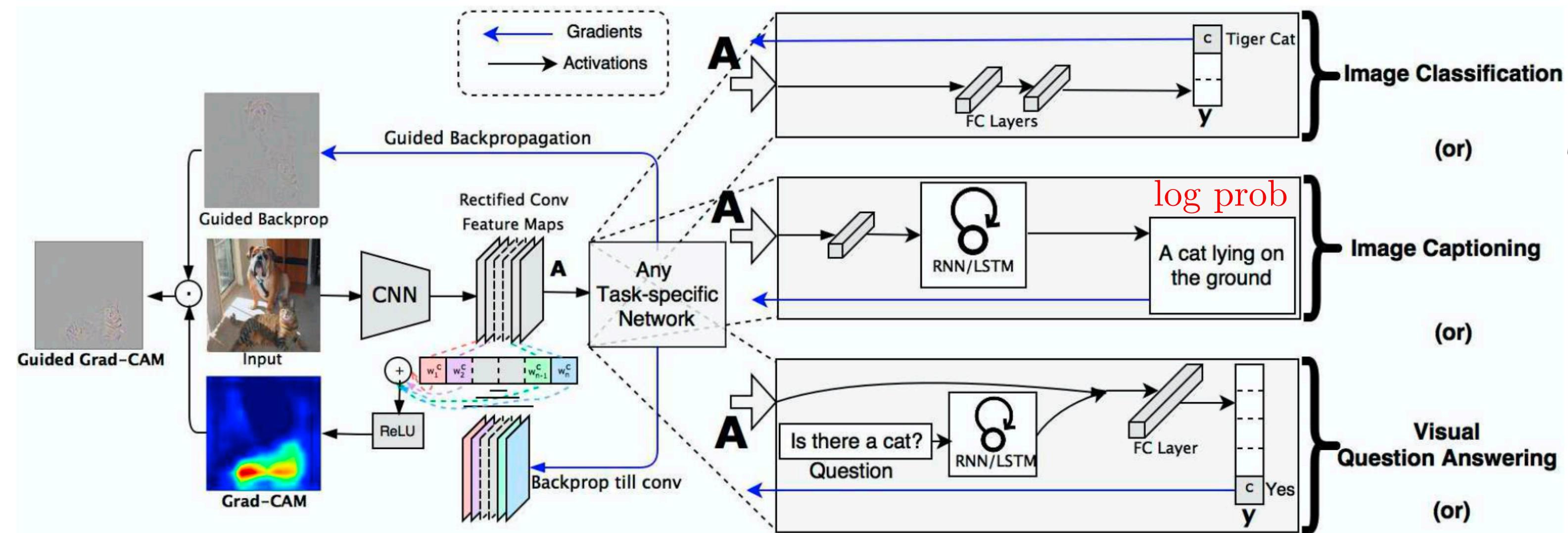


(a) Inception on ImageNet



(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

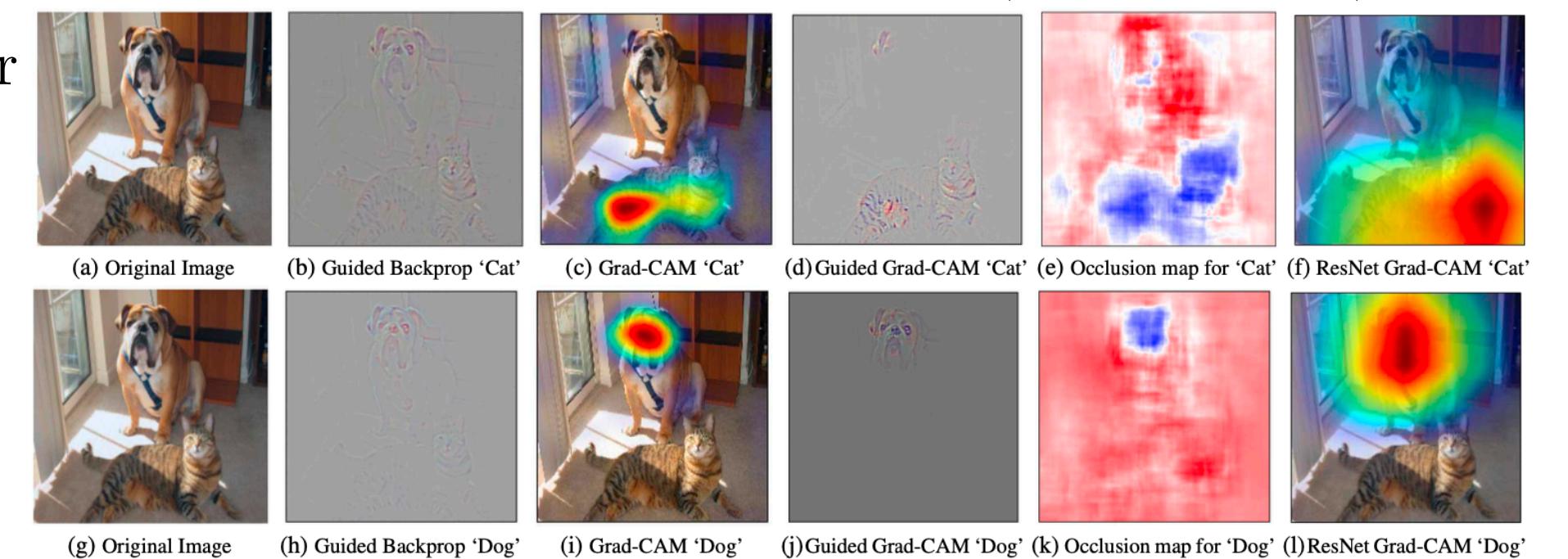
$$\frac{\partial y^c}{\partial A^k}$$

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)

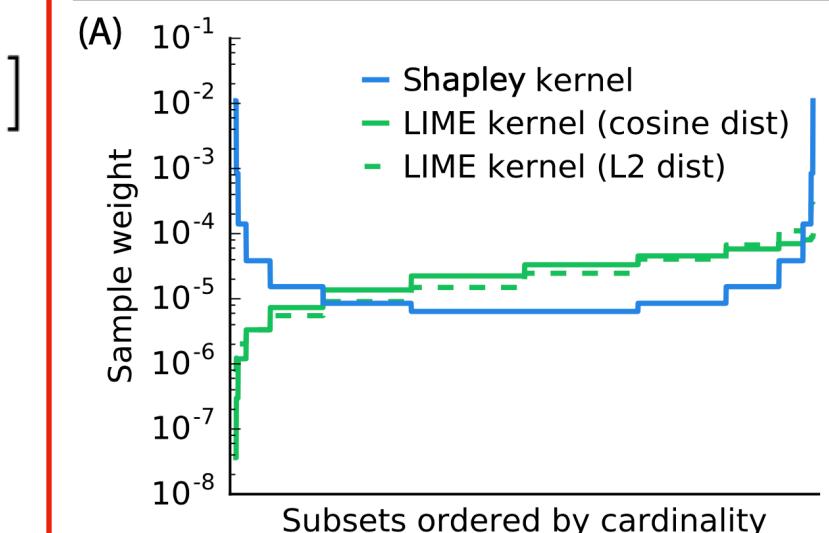
$$\xi(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

$\Omega(g) \rightarrow$ 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

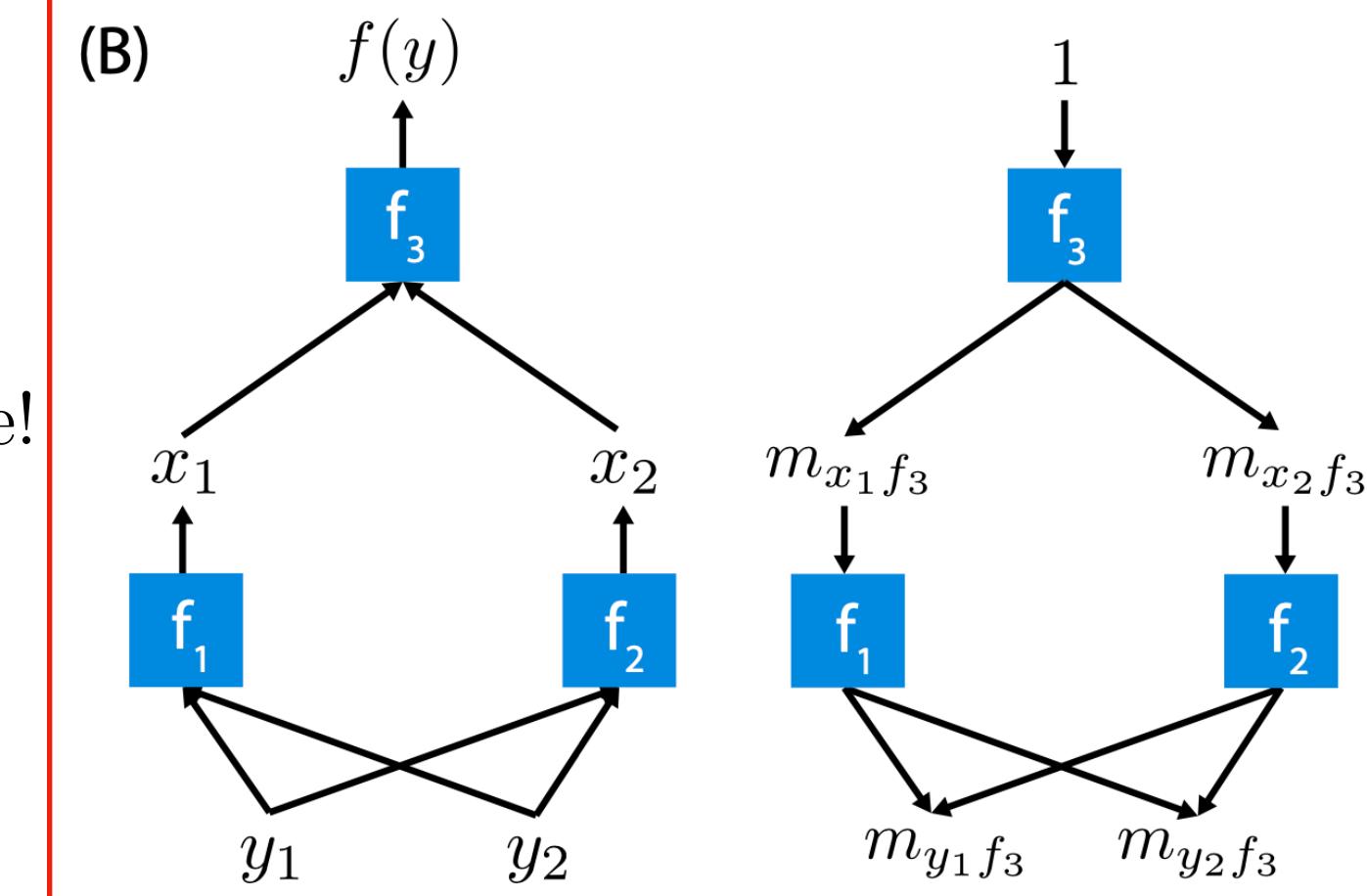
$$f(x) = \sum_{j=1}^M w_j x_j + b$$

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

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

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

On Calibration of Modern Neural Networks

$X \in \mathcal{X} \rightarrow$ input

$Y \in \mathcal{Y} = \{1, 2, \dots, K\} \rightarrow$ label

$\pi(X, Y) = \pi(Y|X)\pi(X) \rightarrow$ ground-truth joint distribution

$h \rightarrow$ neural network

$(\hat{Y}, \hat{P}) = h(X) \rightarrow$ class prediction \hat{Y} and associated confidence \hat{P}

We would like \hat{P} to be calibrated (i.e., represent a true probability of correctness)

For example, given 100 predictions, each with confidence of 0.8, we expect that 80 should be correctly classified.

$\underbrace{\mathbb{P}(\hat{Y} = Y | \hat{P} = p)}_{\text{over the joint distribution}} = p, \quad \forall p \in [0, 1] \rightarrow$ perfect calibration

Reliability Diagrams

$B_m \rightarrow$ set of indices of samples whose prediction confidence \hat{p}

falls into the interval $I_m = (\frac{m-1}{M}, \frac{m}{M}), m = 1, \dots, M$

$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{1}(\hat{y}_i = y_i) \rightarrow$ consistent and unbiased estimator of $\mathbb{P}(\hat{Y} = Y | \hat{P} \in I_m)$

$\text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \hat{p}_i \rightarrow$ average confidence

Expected Calibration Error (ECE)

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n} \underbrace{\left| \text{acc}(B_m) - \text{conf}(B_m) \right|}_{\text{calibration gap}}$$

$n \rightarrow$ number of samples

Temperature Scaling

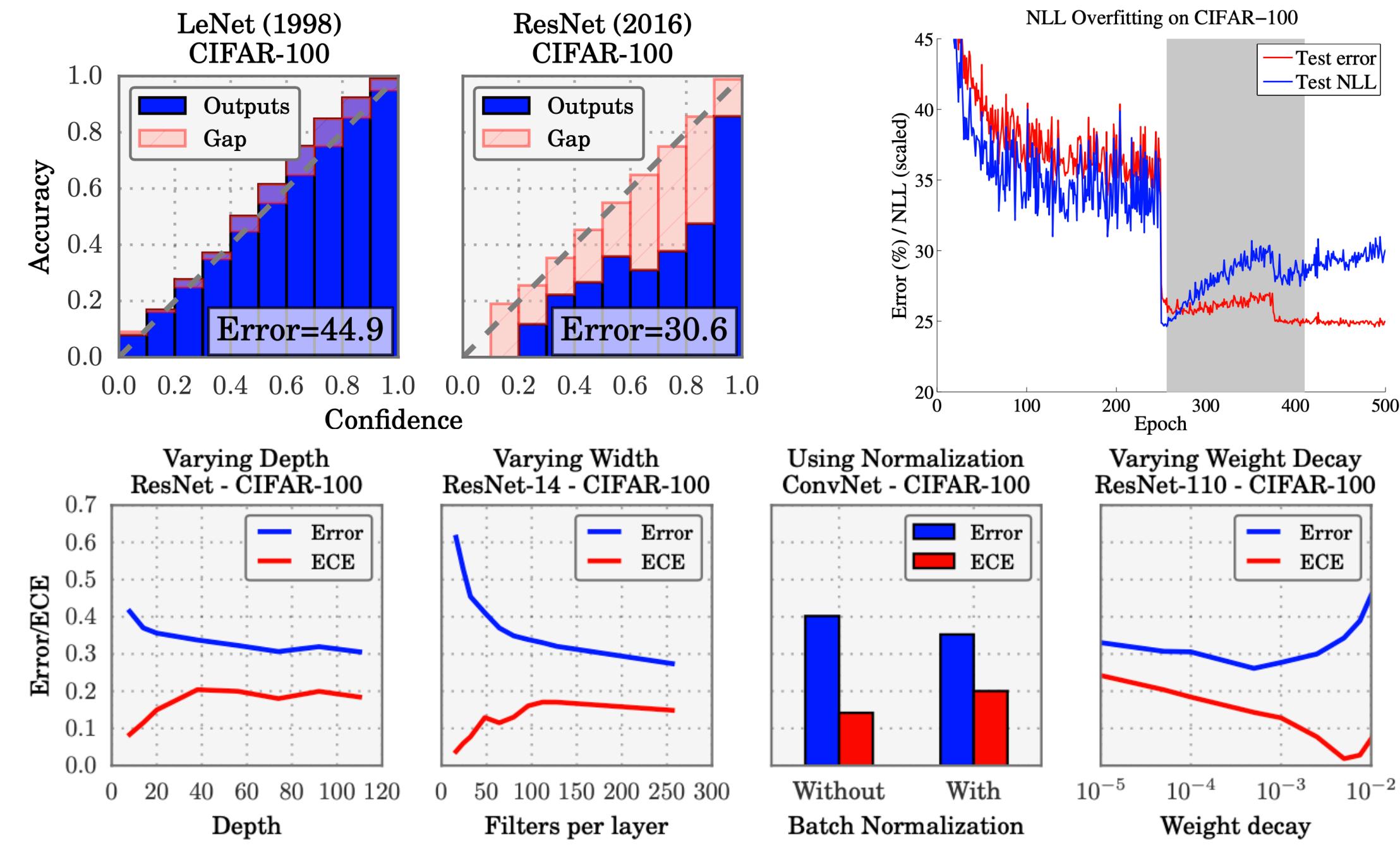
$\mathbf{z}_i \rightarrow$ network logits

$T > 0 \rightarrow$ temperature scaler

$\hat{q}_i = \max_k \text{softmax}(\frac{\mathbf{z}_i}{T})^{(k)} \rightarrow$ calibrated confidence

Optimize T using NLL on the validation data

Modern neural networks are no longer well-calibrated!



Disconnect between NLL (Negative Log Likelihood) and 0/1 loss!

Dataset	Model	Uncalibrated	Hist. Binning	Isotonic	BBQ	Temp. Scaling	Vector Scaling	Matrix Scaling
Birds	ResNet 50	9.19%	4.34%	5.22%	4.12%	1.85%	3.0%	21.13%
Cars	ResNet 50	4.3%	1.74%	4.29%	1.84%	2.35%	2.37%	10.5%
CIFAR-10	ResNet 110	4.6%	0.58%	0.81%	0.54%	0.83%	0.88%	1.0%
CIFAR-10	ResNet 110 (SD)	4.12%	0.67%	1.11%	0.9%	0.6%	0.64%	0.72%
CIFAR-10	Wide ResNet 32	4.52%	0.72%	1.08%	0.74%	0.54%	0.6%	0.72%
CIFAR-10	DenseNet 40	3.28%	0.44%	0.61%	0.81%	0.33%	0.41%	0.41%
CIFAR-10	LeNet 5	3.02%	1.56%	1.85%	1.59%	0.93%	1.15%	1.16%
CIFAR-100	ResNet 110	16.53%	2.66%	4.99%	5.46%	1.26%	1.32%	25.49%
CIFAR-100	ResNet 110 (SD)	12.67%	2.46%	4.16%	3.58%	0.96%	0.9%	20.09%
CIFAR-100	Wide ResNet 32	15.0%	3.01%	5.85%	5.77%	2.32%	2.57%	24.44%
CIFAR-100	DenseNet 40	10.37%	2.68%	4.51%	3.59%	1.18%	1.09%	21.87%
CIFAR-100	LeNet 5	4.85%	6.48%	2.35%	3.77%	2.02%	2.09%	13.24%
ImageNet	DenseNet 161	6.28%	4.52%	5.18%	3.51%	1.99%	2.24%	-
ImageNet	ResNet 152	5.48%	4.36%	4.77%	3.56%	1.86%	2.23%	-
SVHN	ResNet 152 (SD)	0.44%	0.14%	0.28%	0.22%	0.17%	0.27%	0.17%
20 News	DAN 3	8.02%	3.6%	5.52%	4.98%	4.11%	4.61%	9.1%
Reuters	DAN 3	0.85%	1.75%	1.15%	0.97%	0.91%	0.66%	1.58%
SST Binary	TreeLSTM	6.63%	1.93%	1.65%	2.27%	1.84%	1.84%	1.84%
SST Fine Grained	TreeLSTM	6.71%	2.09%	1.65%	2.61%	2.56%	2.98%	2.39%



Boulder

Understanding the Role of Individual Units in a Deep Neural Network

Network Dissection: What a network is looking for and why?

Unit \equiv feature map \equiv activation of a single filter

Emergence of Object Detector Units in a Scene Classifier

Analyze all units within the 13 convolutional layers of the VGG-16 network trained on the Places365 dataset (365 scene categories)

$a_u(x, p) \rightarrow$ activation computed by unit u

$x \rightarrow$ test image $p \rightarrow$ image position

Filters with low-resolution outputs are visualized and analyzed at high-resolution positions p using bilinear up sampling.

$t_u \equiv \arg \max_t \mathbb{P}_{x,p}[a_u(x, p) > t] \geq 0.01 \rightarrow$ top 1% quantile

probability that the event is true when sampled over all positions and images

$\{p : a_u(x, p) > t_u\} \rightarrow$ activation region above the threshold

Interpretability Concept

$c \rightarrow$ a visual concept

Objects: shelf, plant, airplane, etc.

Parts: person-top, person-bottom, tree-top, tree-bottom, etc.

Materials: food, hair, skin, paper, etc.

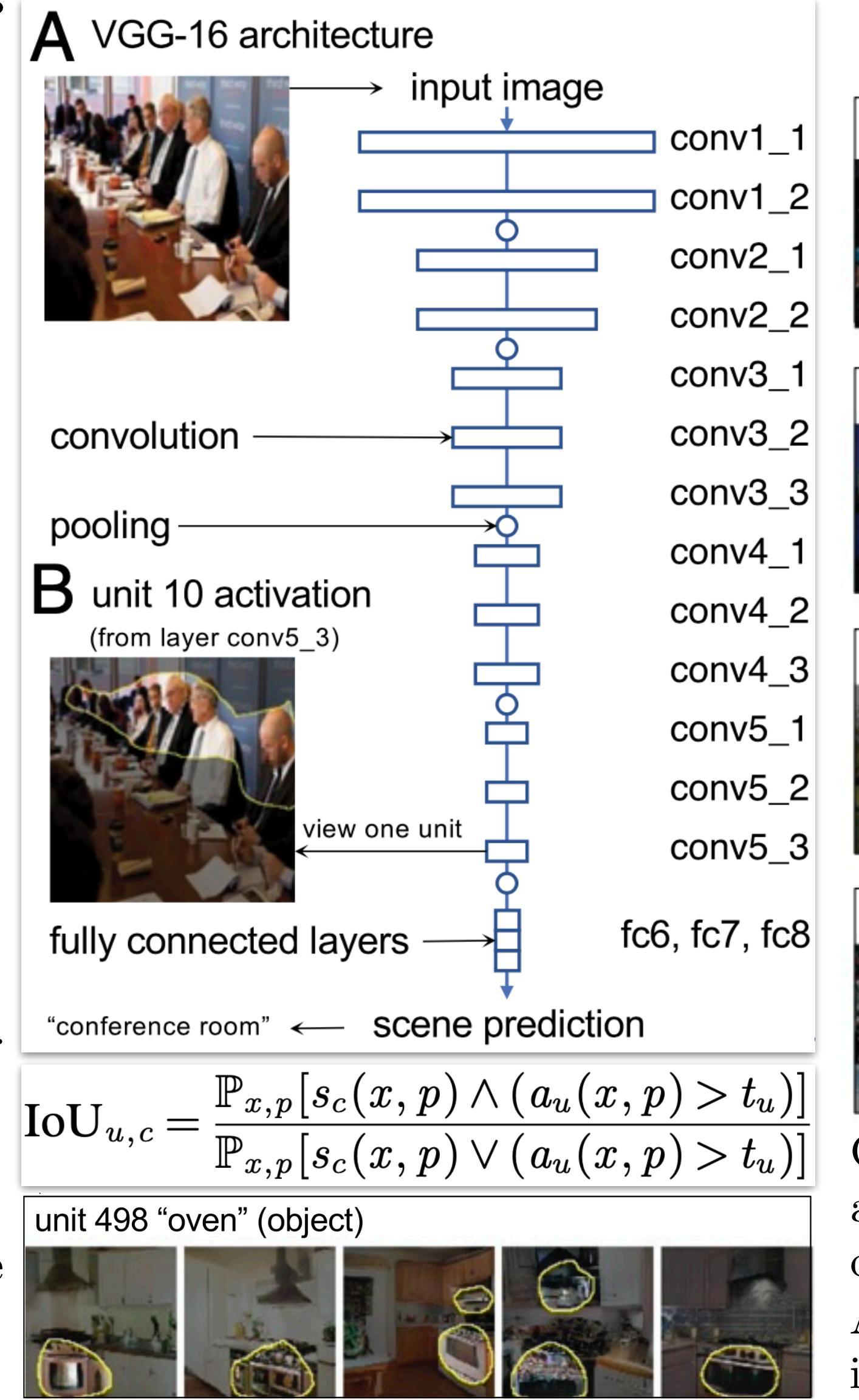
Colors: orange, yellow, red, etc.

$s_c(x, p) \in \{0, 1\} \rightarrow$

A segmentation model trained to predict the presence of the visual concept c within image x at position p .

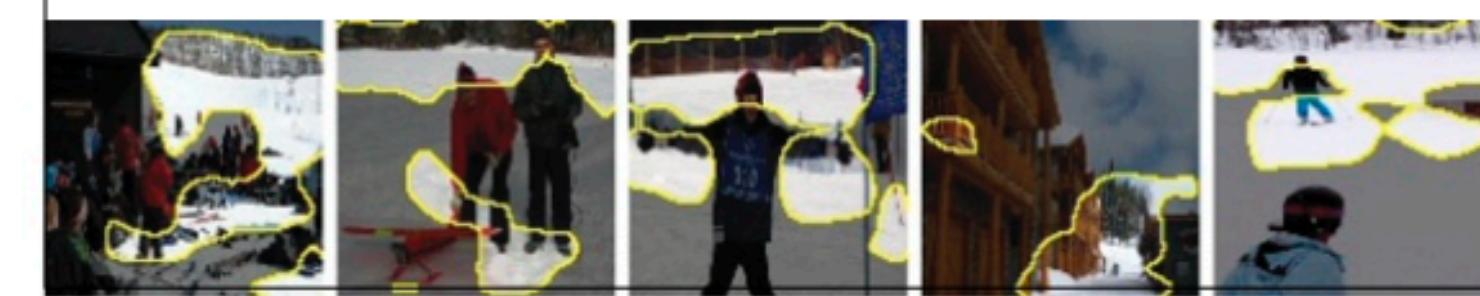
$\text{IoU}_{u,c} \rightarrow$ agreement between concept c and unit u

Bau, David, et al. "Understanding the role of individual units in a deep neural network." *Proceedings of the National Academy of Sciences* 117.48 (2020): 30071-30078.

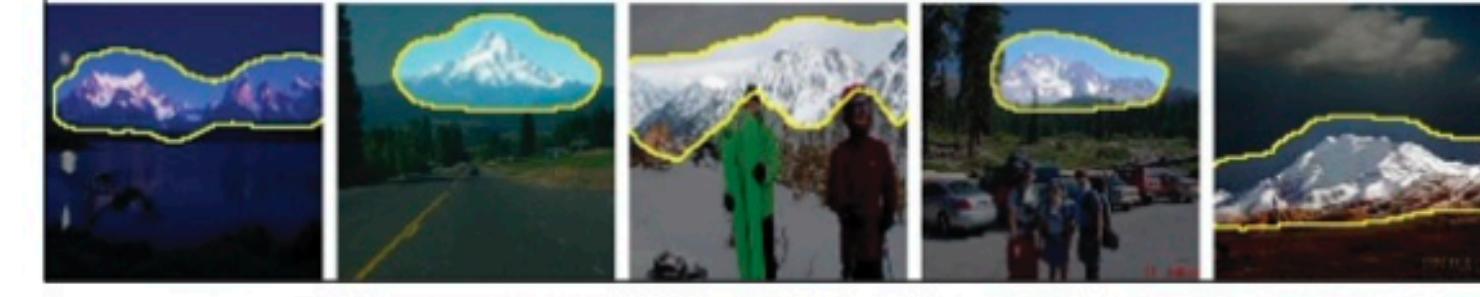


A Units of conv5_3 causing most accuracy loss on the single class "ski resort" when removed individually

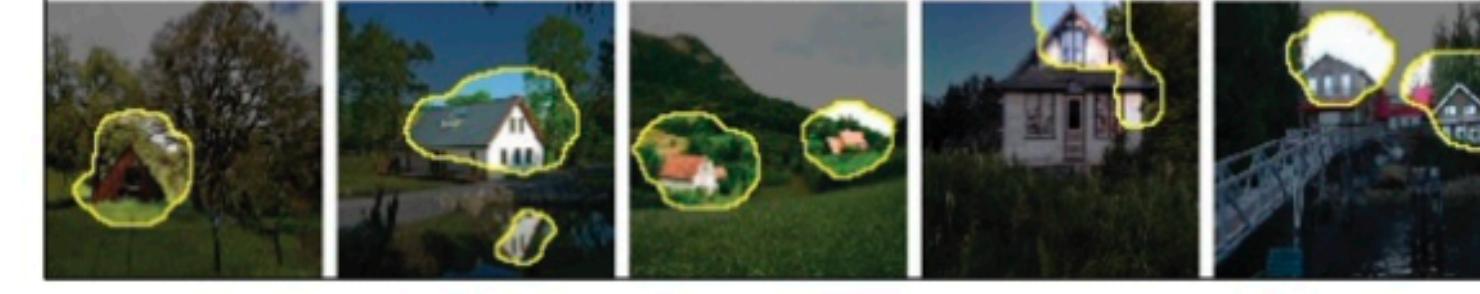
unit 261 "snow" (acc lost: train 8.9% val 10.9%)



unit 149 "mountain top" (acc lost: train 1.2% val 3.5%)



unit 242 "house" (acc lost: train 1.5% val 2.5%)



unit 105 "tree top" (part, acc lost: train 0.8% val 2.0%)



GANs: How does removal of tree units affect the appearance of trees and other objects in the output image?

Adversarial examples are attacks on the important units for a class

Do Vision Transformers See Like Convolutional Neural Networks?

Representation Similarity and CKA (Centered Kernel Alignment)

$X \in \mathbb{R}^{m \times p_1}, Y \in \mathbb{R}^{m \times p_2} \rightarrow$ internal representations (i.e., activation matrices) of two layers within one network or across two networks, evaluated at the same examples

$m \rightarrow$ number of examples

$p_1, p_2 \rightarrow$ number of neurons

$K = XX^T, L = YY^T \rightarrow$ Gram matrices (measures the similarity of a pair of datapoints according to layer representations)

$$\text{CKA}(K, L) = \frac{\text{HSIC}(K, L)}{\sqrt{\text{HSIC}(K, K)\text{HSIC}(L, L)}}$$

HSIC \rightarrow Hilbert-Schmidt Independence Criterion

$H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T \rightarrow$ centring matrix

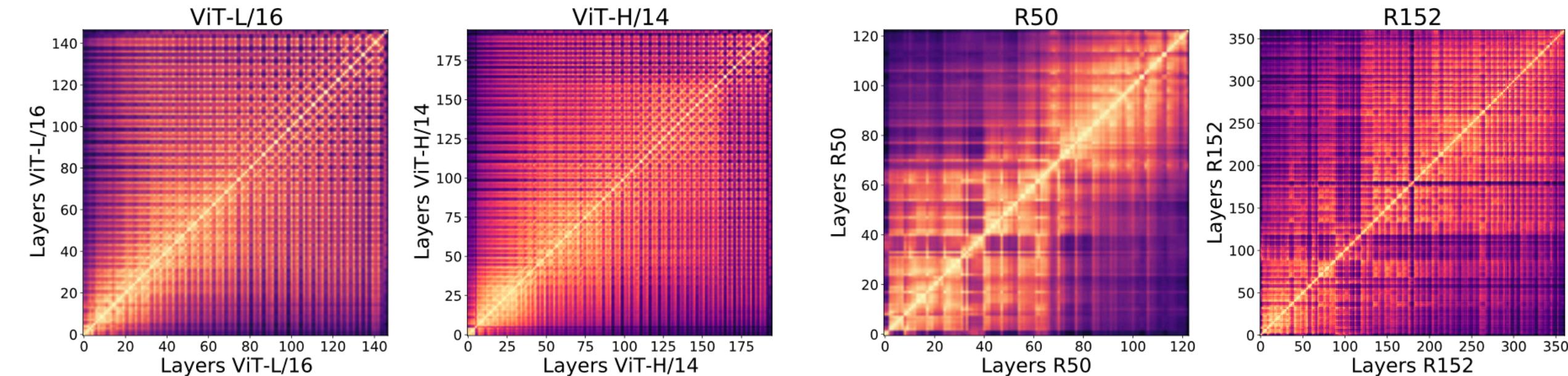
$K' = HKH, L' = HLH \rightarrow$ centered Gram matrices

$\text{HSIC}(K, L) = \text{vec}(K') \cdot \text{vec}(L')/(m-1)^2 \rightarrow$ similarity between the centered Gram matrices

CKA is invariant to orthogonal transformation of representations (including permutation of neurons), and the normalization term ensures invariance to isotropic scaling.

Findings

- ViTs have more uniform representations, with greater similarity between lower and higher layers
- ViT incorporates more global information than ResNet at lower layers, leading to quantitatively different features
ResNet effective receptive fields are highly local and grow gradually; ViT effective receptive fields shift from local to global
effective receptive field of different layers: absolute value of the gradient of the center location of the feature map with respect to the input
- incorporating local information at lower layers remains vital, with large-scale pre-training data helping early attention layers learn to do this
- skip connections in ViT are even more influential than in ResNets, having strong effects on performance and representation similarity
- motivated by potential future uses in object detection, the paper examines how well input spatial information is preserved
- the paper studies the effects of dataset scale on transfer learning, with a linear probes study revealing its importance for high quality intermediate representations





Boulder



Questions?

[YouTube Playlist](#)
