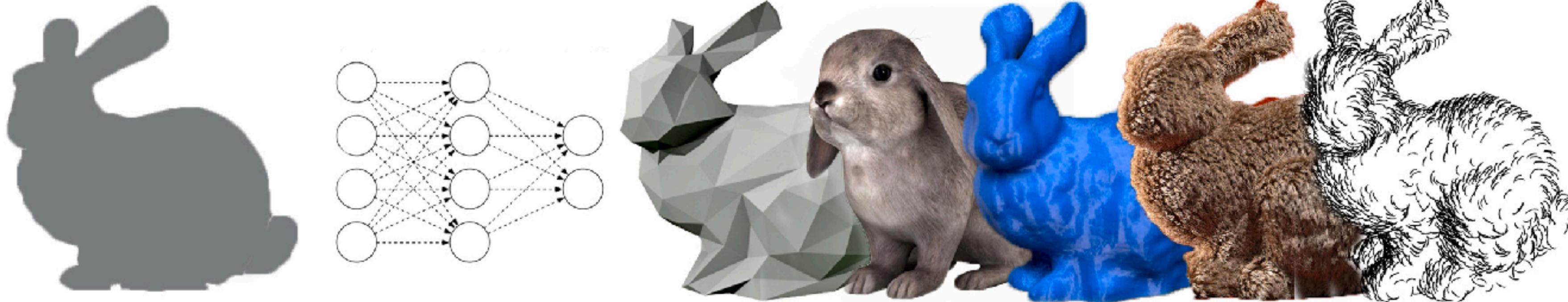


COMP0169: Machine Learning for Visual Computing

Understanding CNNs



Lectures will be Recorded

Multiple Classes & Linear Regression

C classes: one-of-c coding (or one-hot encoding)

4 classes, i-th sample is in 3rd class:
 $\mathbf{y}^i = (0, 0, 1, 0)$

Matrix notation:

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}^1 \\ \vdots \\ \mathbf{y}^N \end{bmatrix} = [\mathbf{y}_1 \mid \dots \mid \mathbf{y}_C] \quad \text{where } \mathbf{y}_c = \begin{bmatrix} y_c^1 \\ \vdots \\ y_c^N \end{bmatrix}$$

$$\mathbf{W} = [\mathbf{w}_1 \mid \dots \mid \mathbf{w}_C]$$

Loss function:

$$L(\mathbf{W}) = \sum_{c=1}^C (\mathbf{y}_c - \mathbf{X}\mathbf{w}_c)^T (\mathbf{y}_c - \mathbf{X}\mathbf{w}_c)$$

Least squares fit (decouples per class):

$$\mathbf{w}_c^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_c$$

Multiple Classes via SoftMax

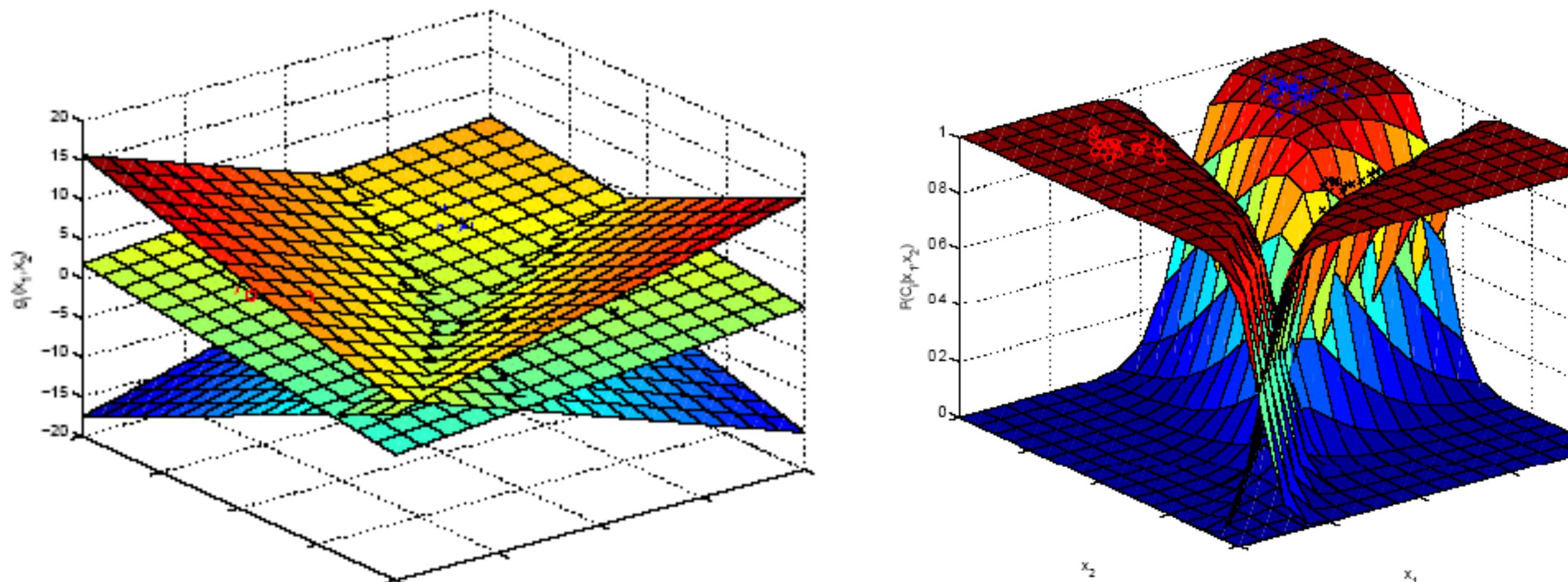
Multiple classes & Logistic regression

softmax

Soft maximum (softmax) of competing classes:

$$P(y = c|x; \mathbf{W}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^T \mathbf{x})} \doteq g_c(\mathbf{x}, \mathbf{W})$$

$L(w)$



Discriminants (inputs) → Softmax (outputs)

+ softmax

Neural Networks (Recap)

Gradient

De Kort

Revisited.

momentum

$v_{i-1} + \beta v_{i+1}$

v_{i-1}

v_i

A hand-drawn diagram consisting of two curved blue lines. The left line forms a large acute angle, while the right line forms a smaller acute angle. The label "ω" is written in blue next to the larger angle, and the label "w" is written next to the smaller angle.

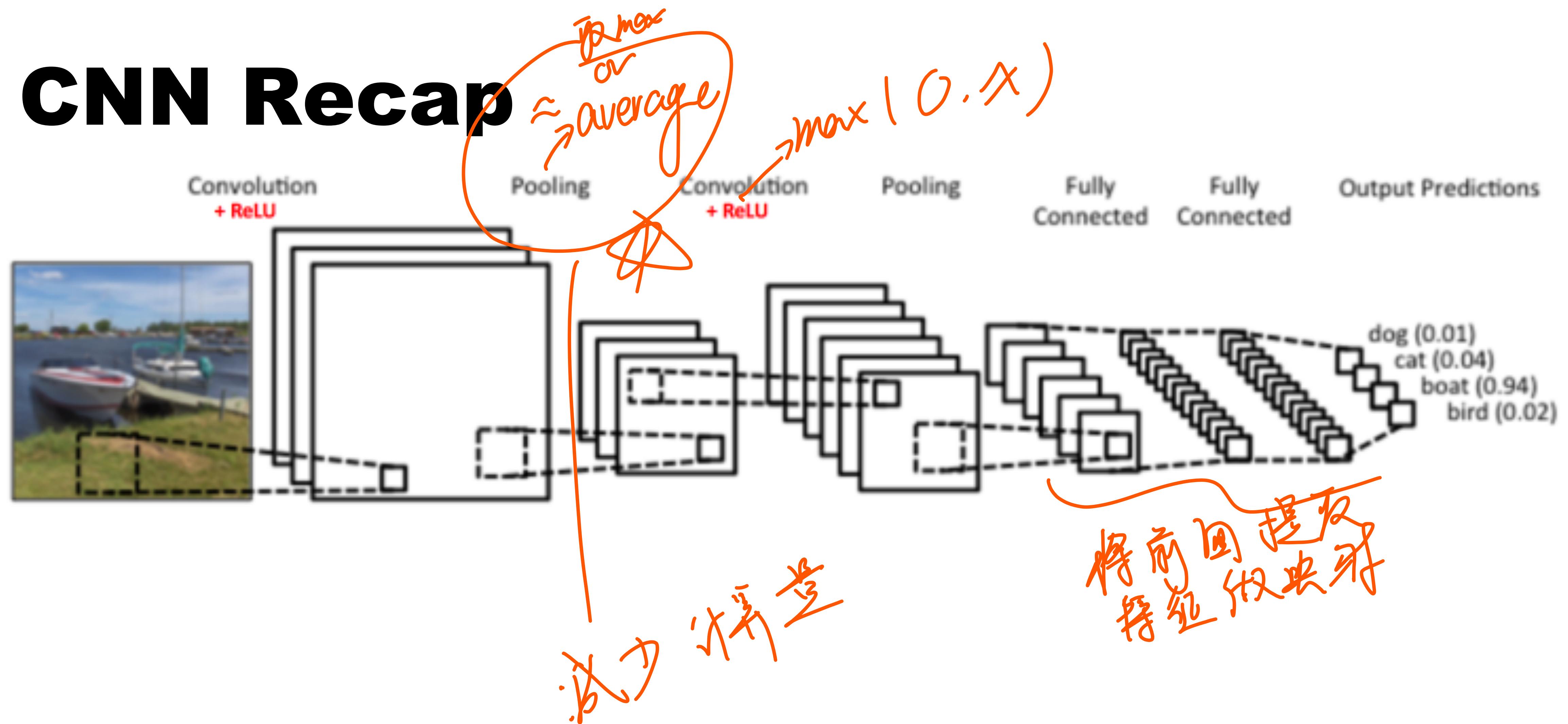
对正像图 ~~像~~ 像 像
data 像 像

SG

stochastic 亂世

A hand-drawn diagram in blue ink on a white background. At the top left, there is a sketch of a curved surface with a small triangle drawn on it, representing a local coordinate system. An arrow points from this sketch down to a larger, more detailed drawing of a surface. This second drawing shows a large oval shape representing a manifold, with a smaller circle inside it representing a point. A vertical line segment connects the center of the inner circle to the outer boundary of the oval. The word 'point' is written next to the inner circle. To the right of the oval, the text '随 扎点' (zai zhu diǎn) is written vertically, followed by 'point' and 'it'. Below the oval, there is some illegible text.

CNN Recap



Convolutions, Kernels, Activations, Features

Plan for Today

- **Visualize filter (kernels)**
- **Use CNNs as feature maps**
- **What are the different layers learning?**
- **Visualize activations at different layers**
- **Saliency versus Occlusion**
- **Guided backpropagation**
- **Gradient Ascent**
- **Feature Inversion**
- **DeepDreams**

AlexNet

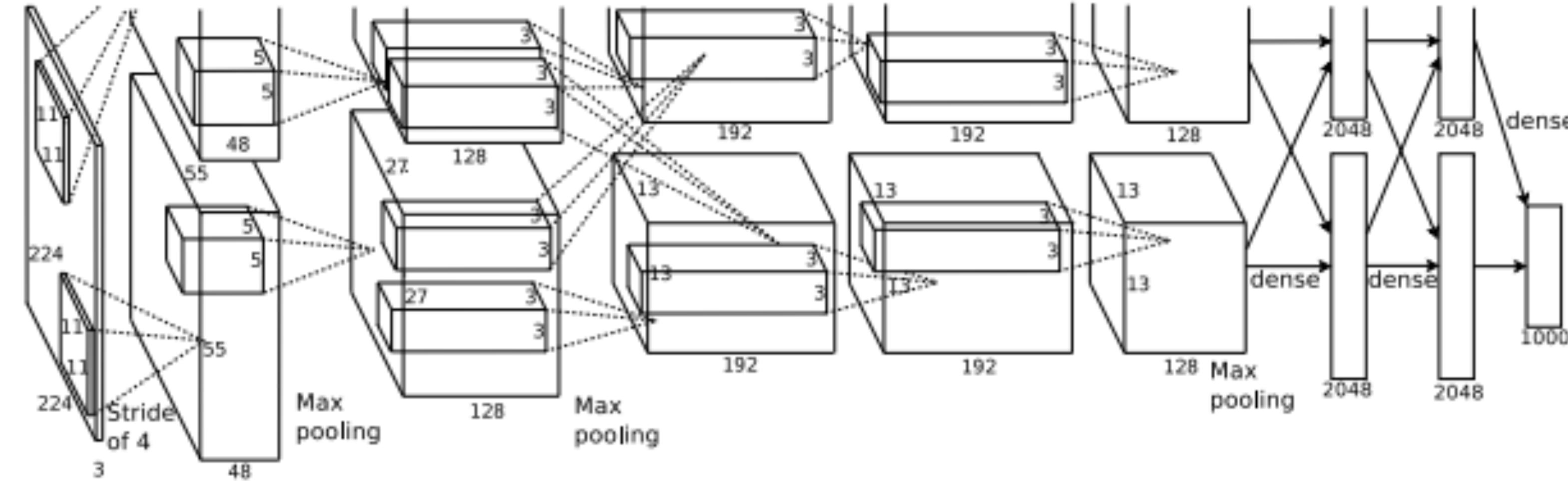


ImageNet Classification with Deep Convolutional Neural Networks

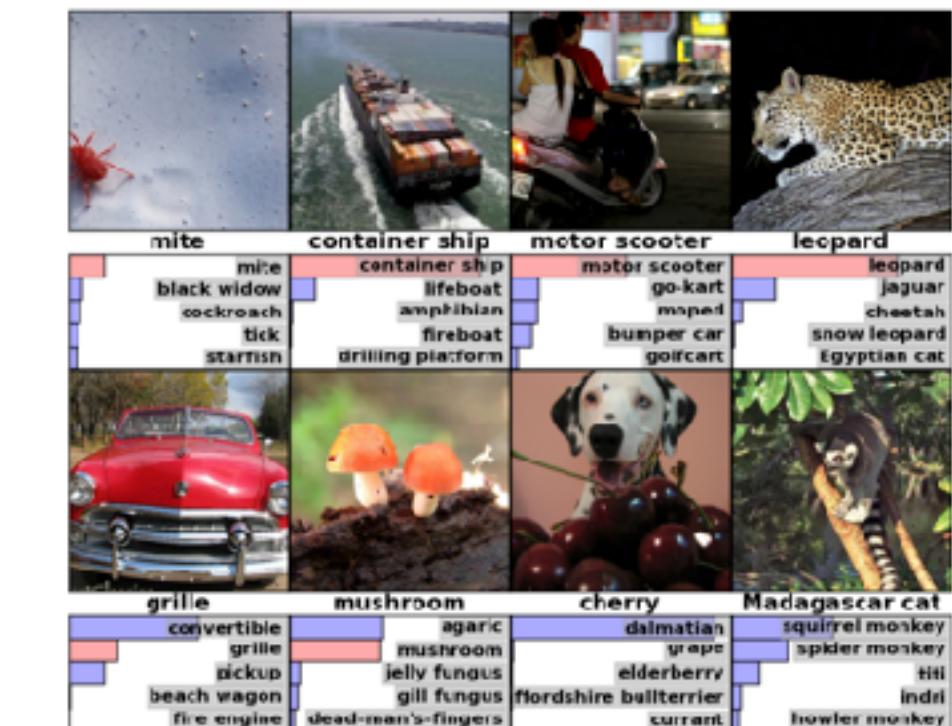
Alex Krizhevsky
University of Toronto
kriz@cs.toronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.toronto.ca

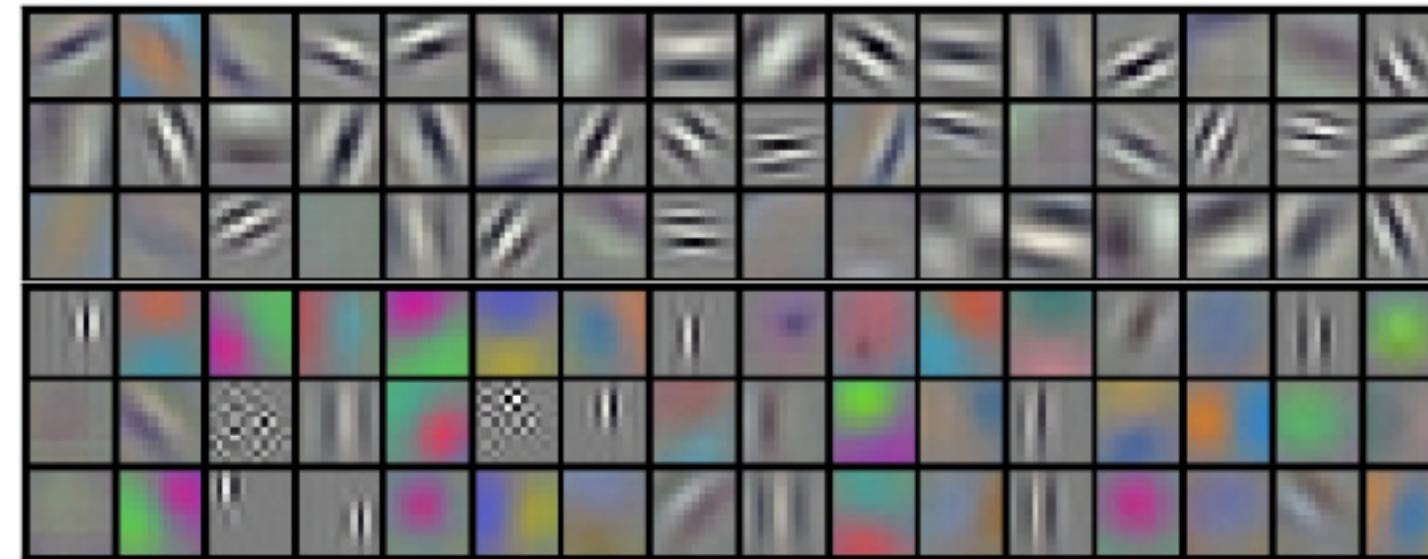
Geoffrey E. Hinton
University of Toronto
hinton@cs.toronto.ca



Class scores
for 1000
Classes

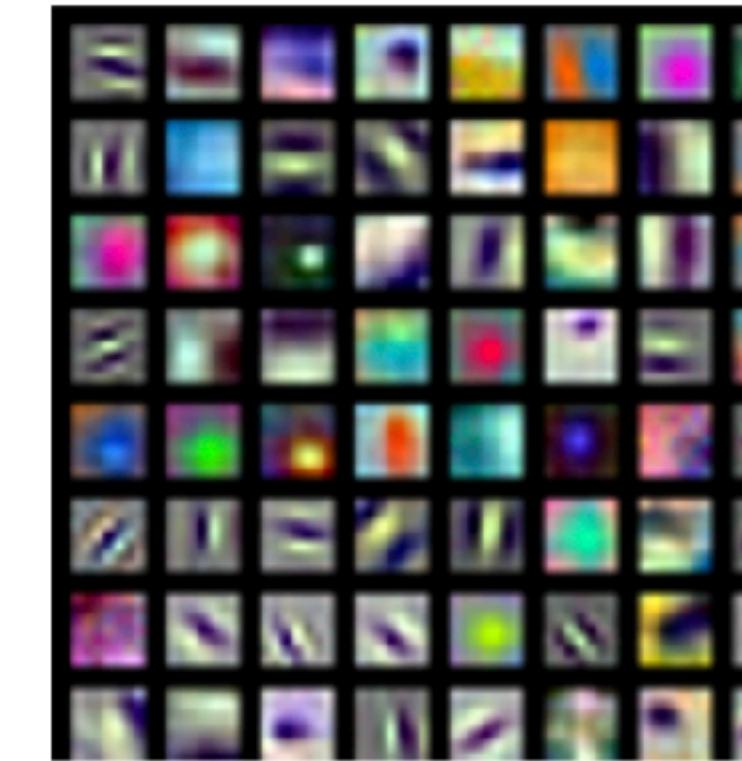


First Layer Filters



AlexNet
 $96 \times 11 \times 11 \times 3$

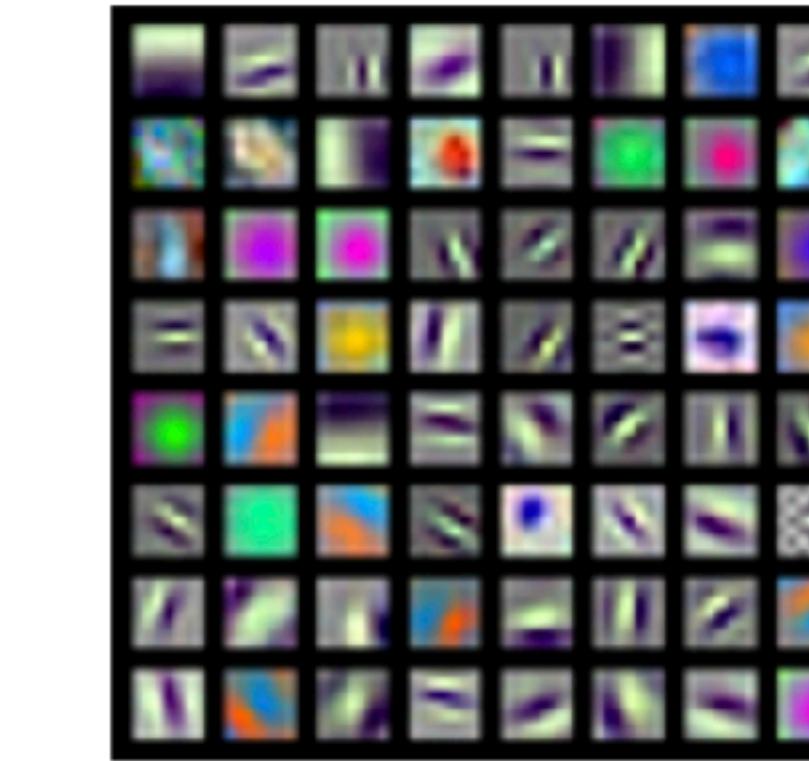
edge



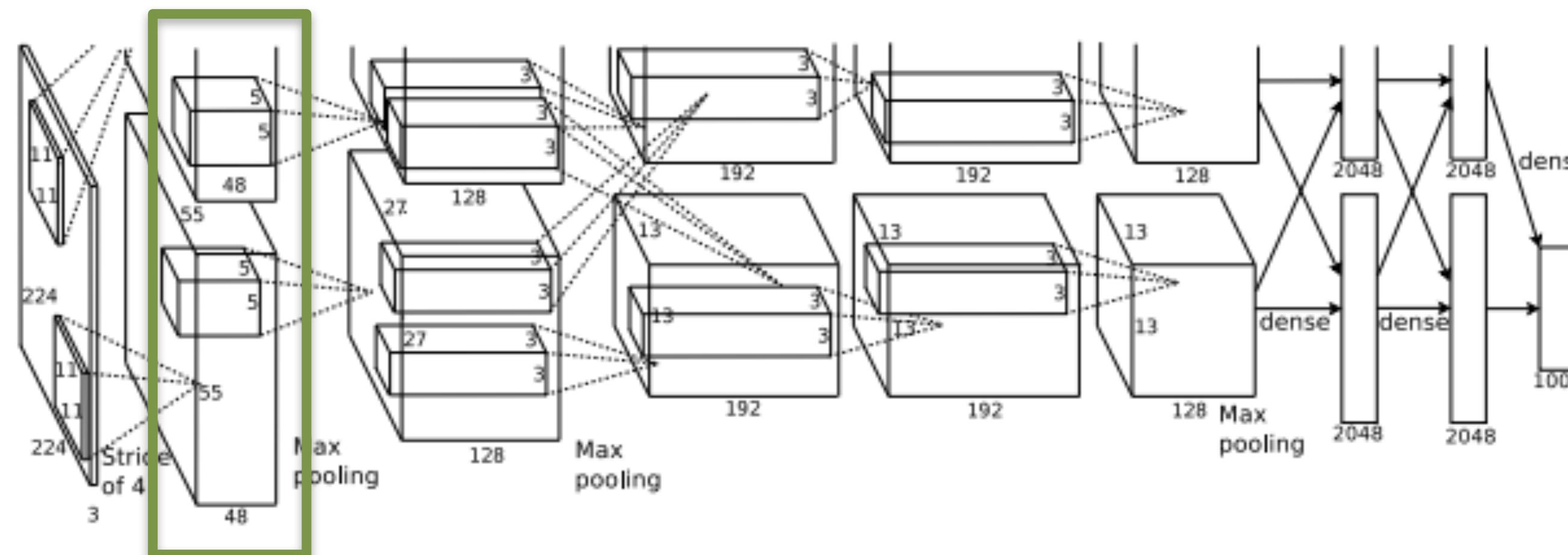
ResNet-18
 $64 \times 7 \times 7 \times 3$



ResNet-101
 $64 \times 7 \times 7 \times 3$



DenseNet-121
 $64 \times 7 \times 7 \times 3$



Higher Level Filters



First layer weights: $16 \times 3 \times 7 \times 7$



Second layer weights:
 $20 \times 16 \times 7 \times 7$



Third layer weights:
 $20 \times 20 \times 7 \times 7$

Source: ConvNetJS CIFAR-10 example

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

NN in Image Space

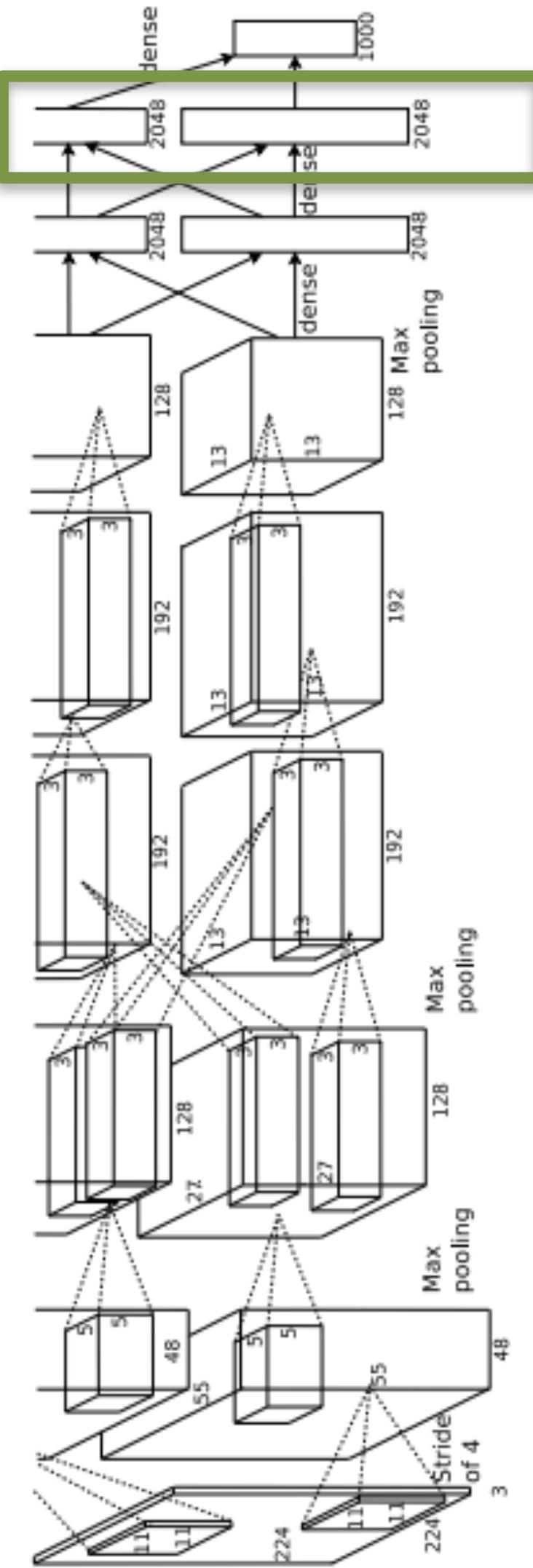


Last Layer

$$\mathbb{R}^{224 \times 224 \times 3} \rightarrow \mathbb{R}^{4096}$$

Feature size 2048+2048=4096.

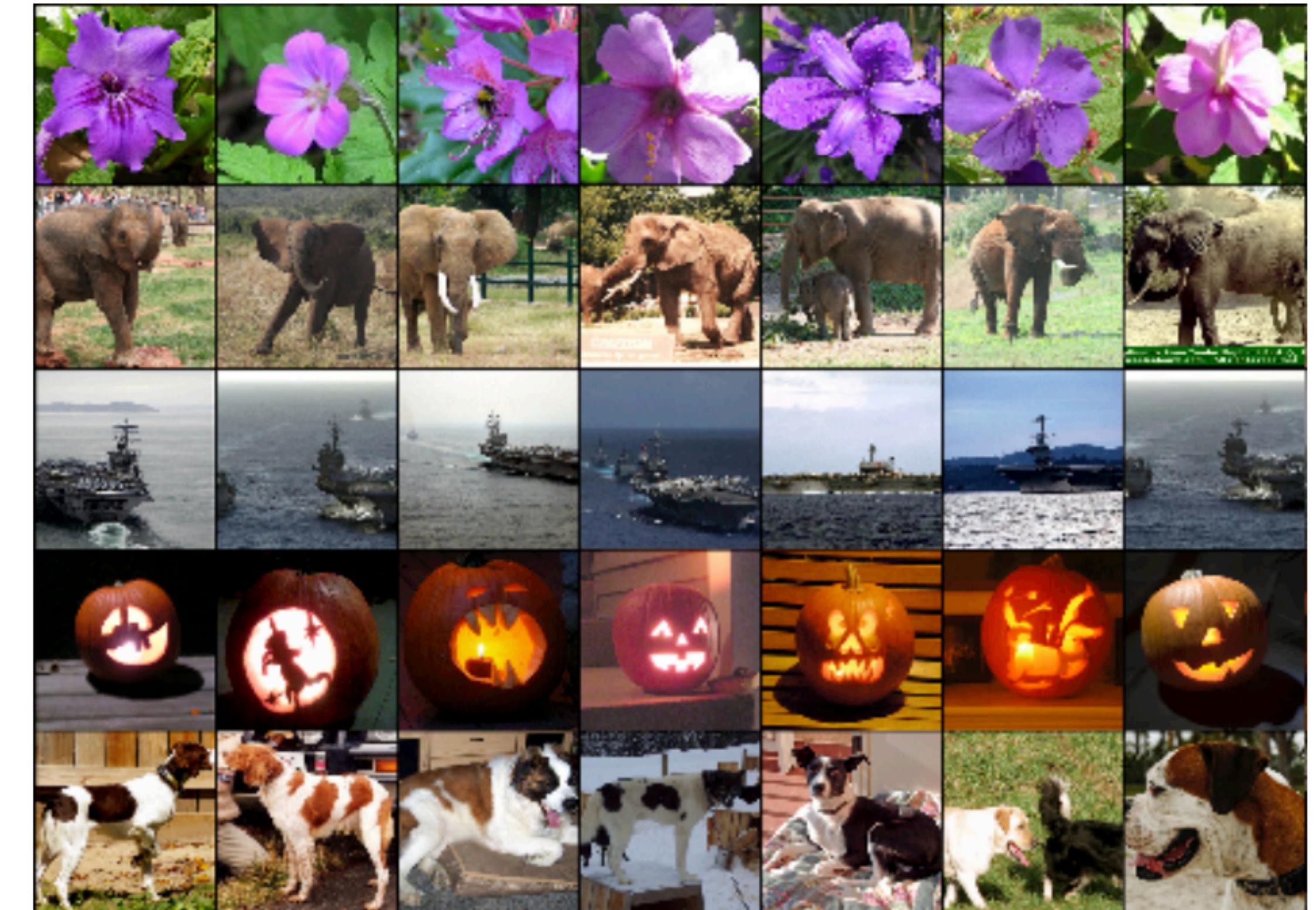
Pretrain network. Run ‘feature extractor’ on many images.
FC7 features.



Nearest Neighbors in Feature Space



Nearest neighbors in pixel space



Nearest neighbors in FC7 space

Dimensionality Reduction

$$\mathbb{R}^{4096} \rightarrow \mathbb{R}^2$$

1. PCA
2. t-SNE

*PCA
take the top eigen*

Visualizing Data using t-SNE

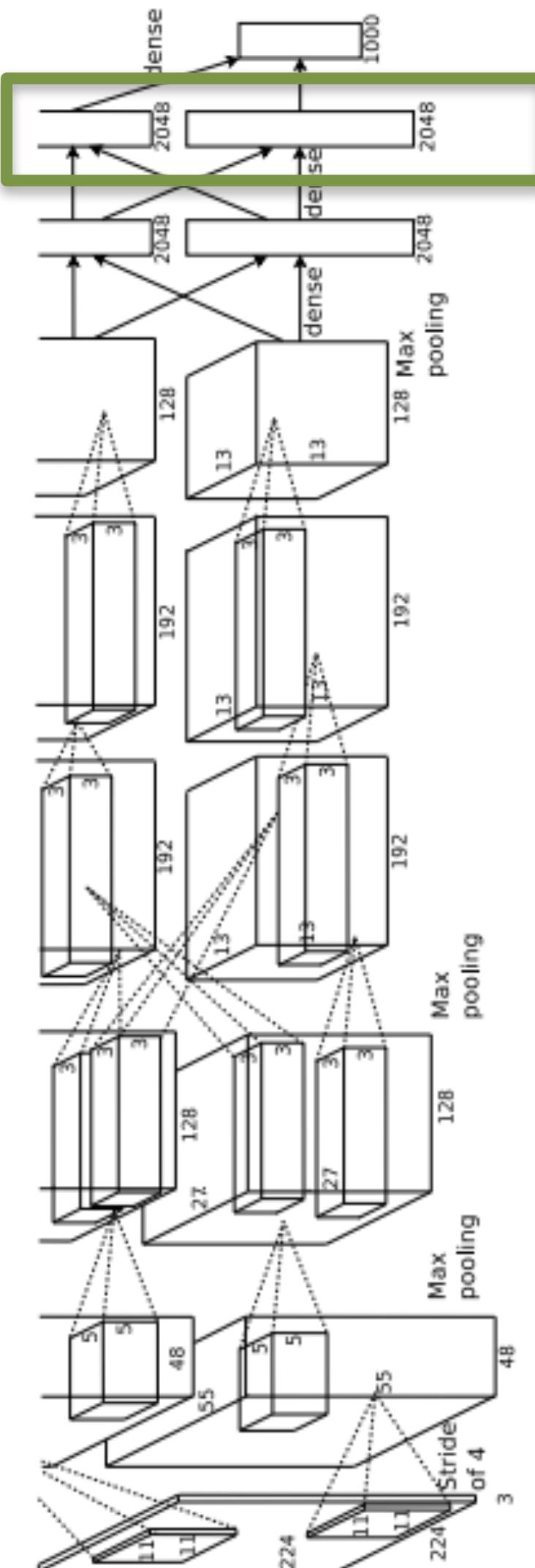
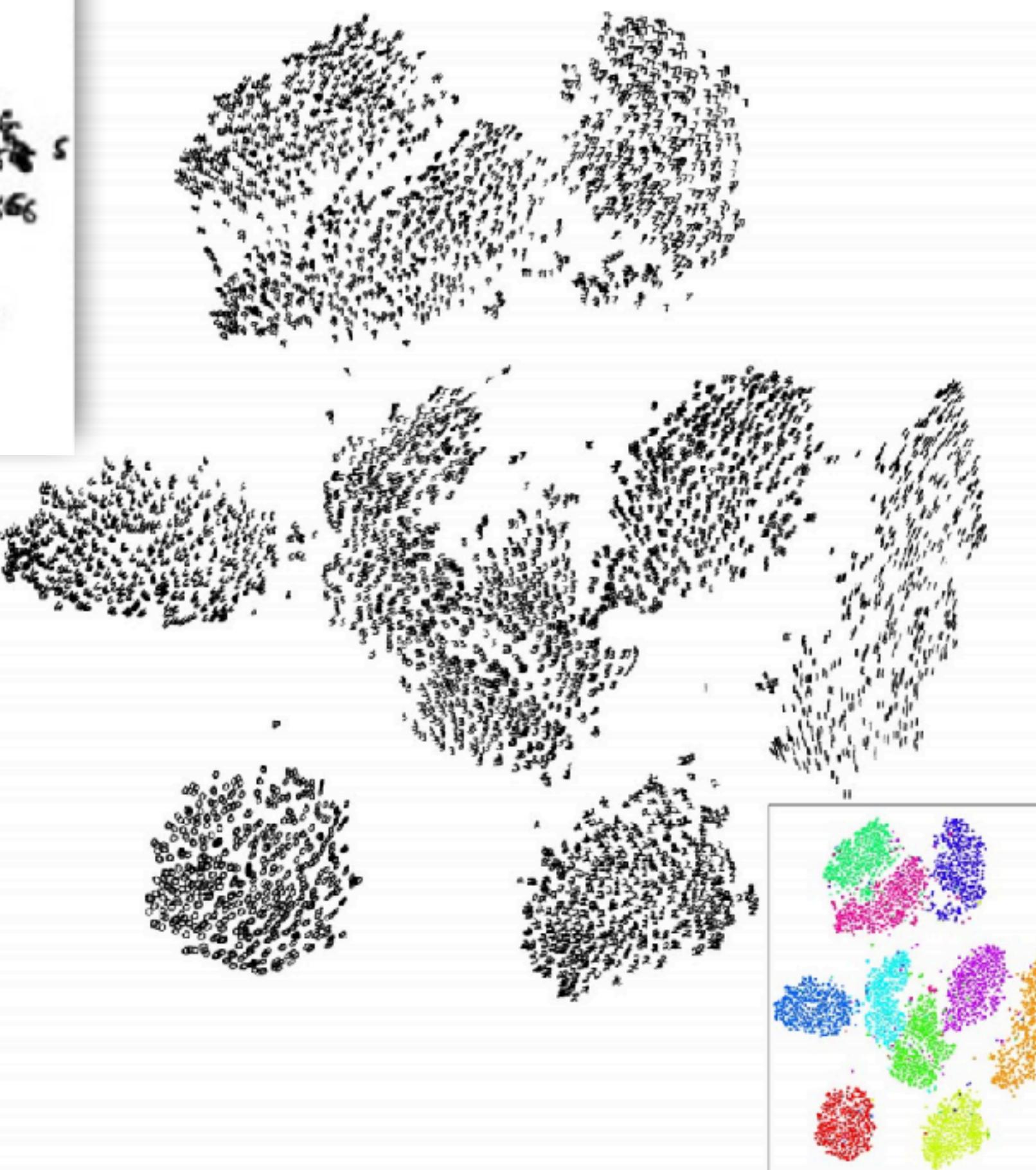
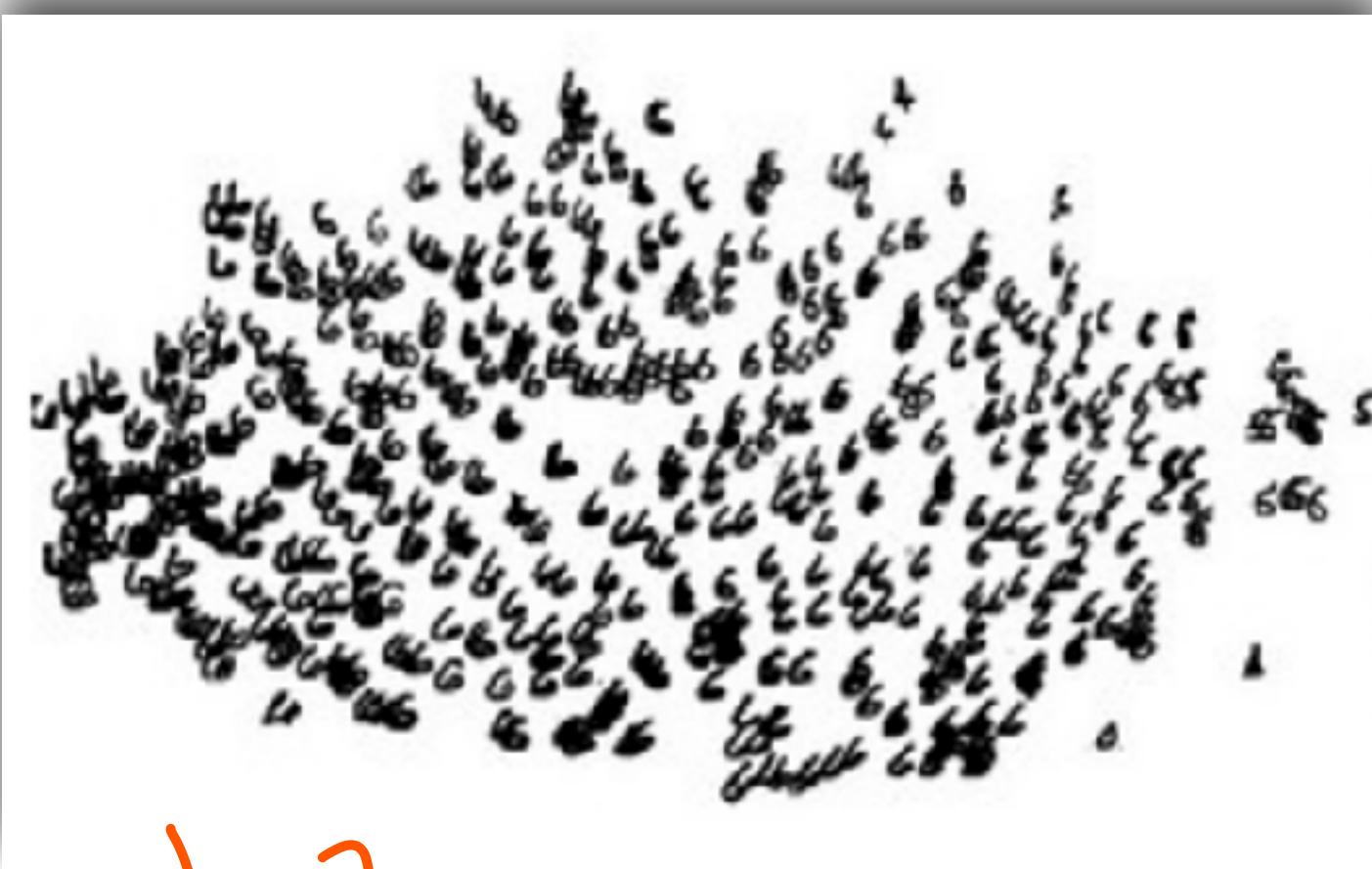
Laurens van der Maaten

Tilburg University
P.O. Box 90153, 5000 LE Tilburg, The Netherlands

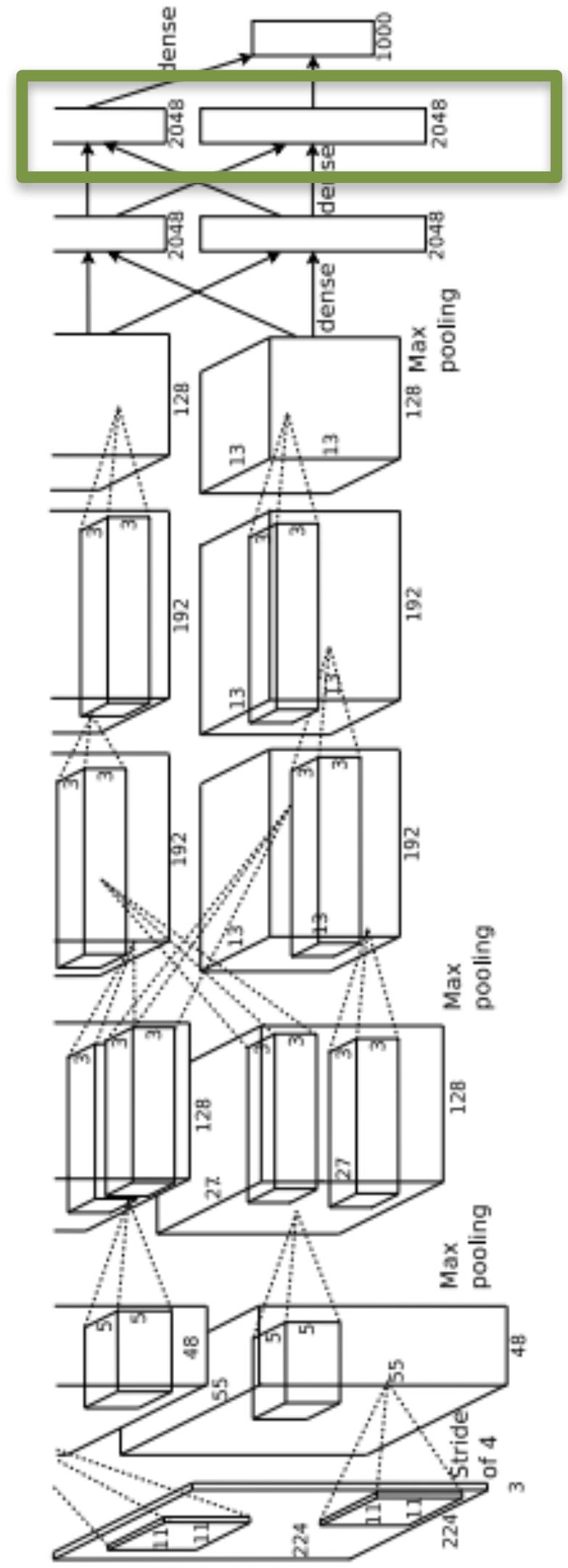
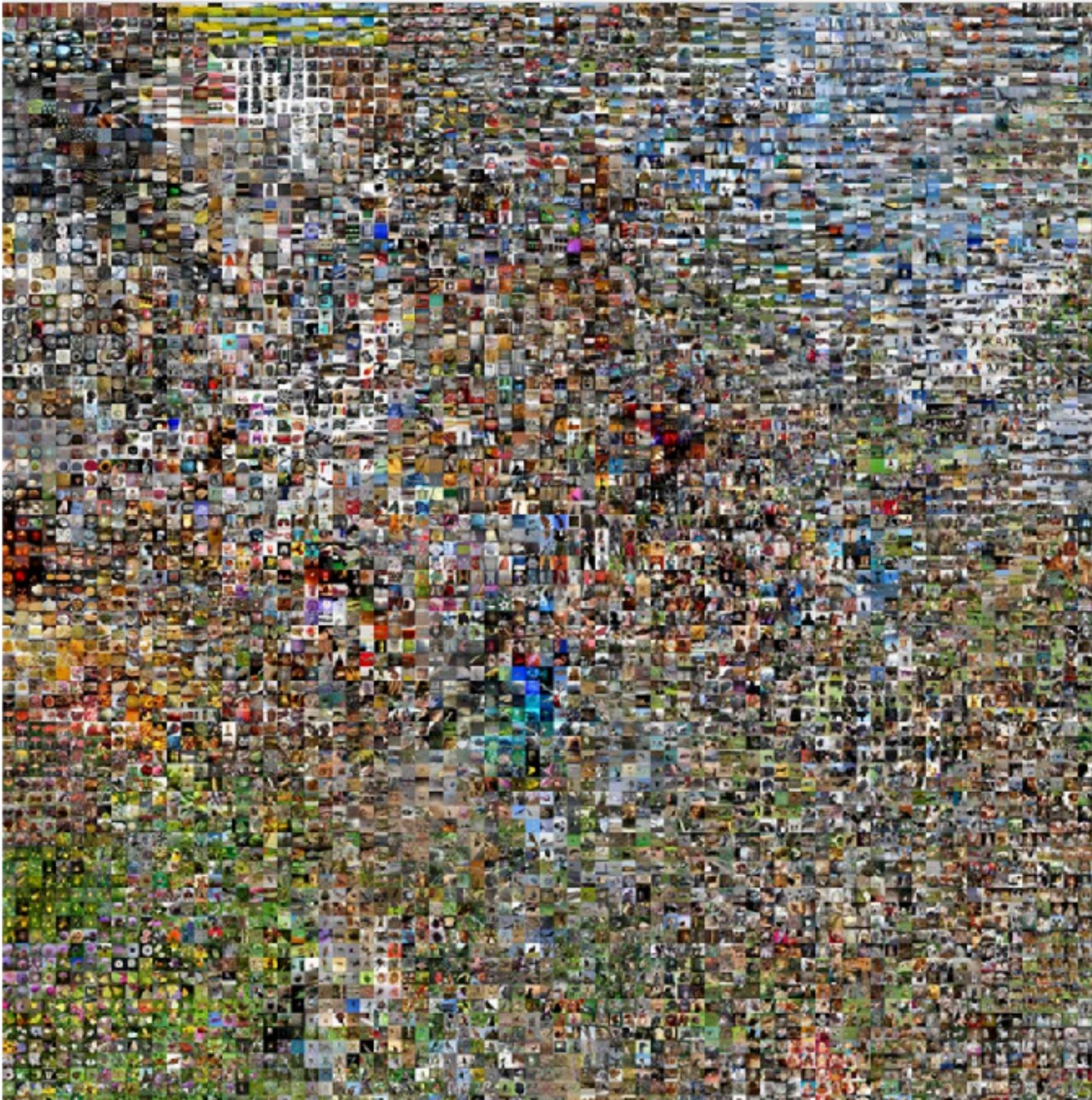
Geoffrey Hinton

Department of Computer Science
University of Toronto
6 King's College Road, M5S 3G4 Toronto, ON, Canada

Editor: Yoshua Bengio



Dimensionality Reduction

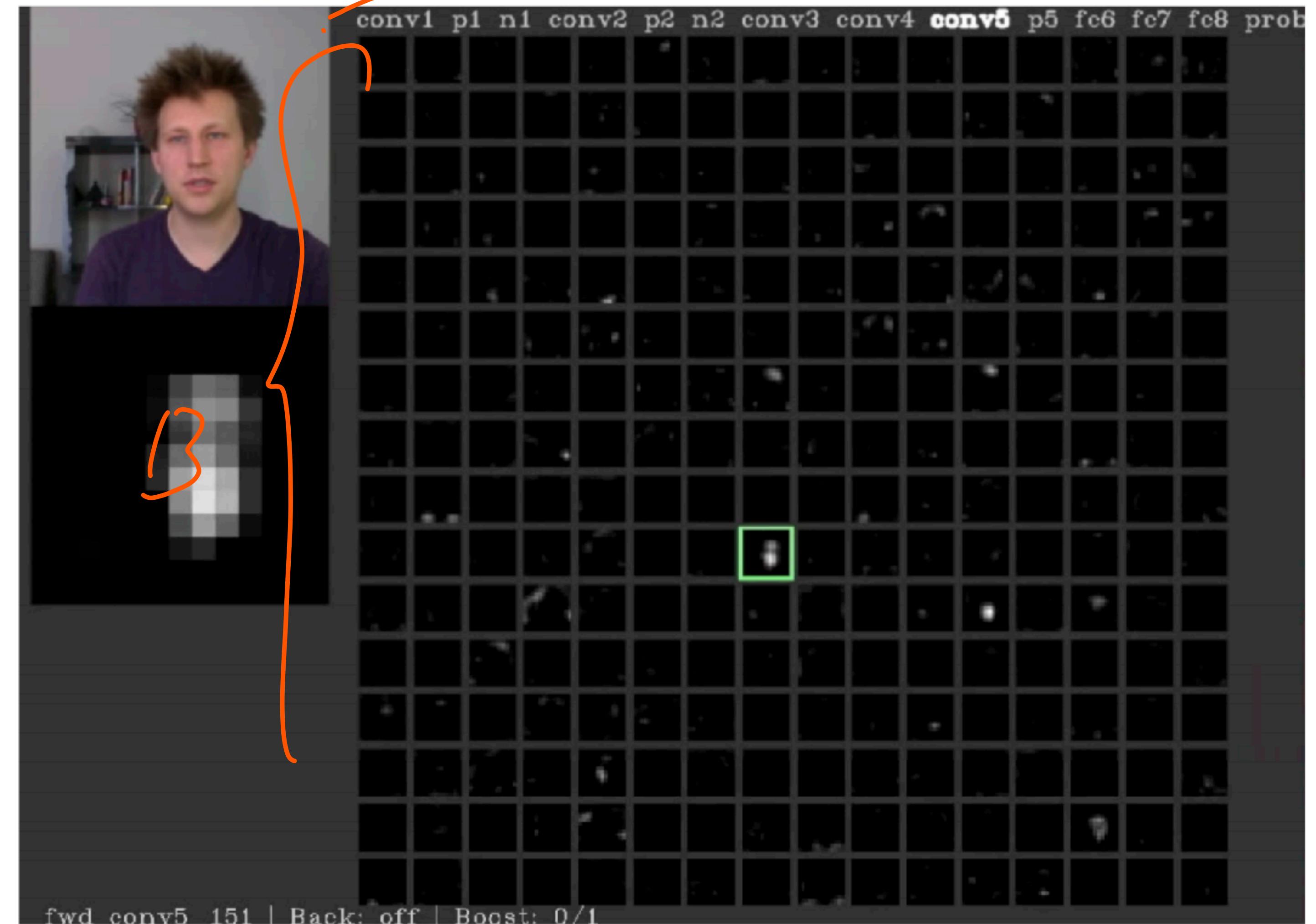
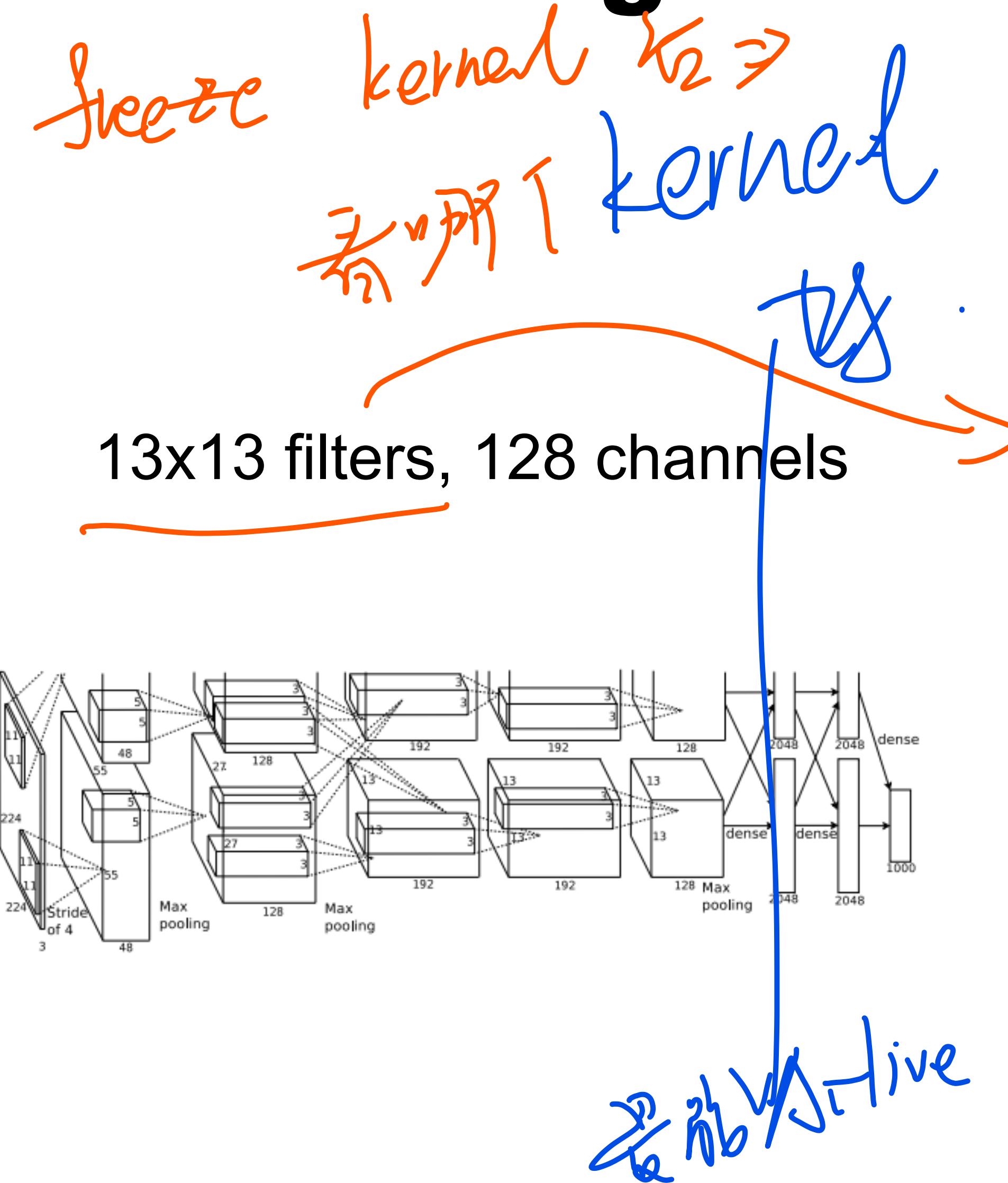


See high-resolution versions at <http://cs.stanford.edu/people/karpathy/cnnembed/>

Plan for Today

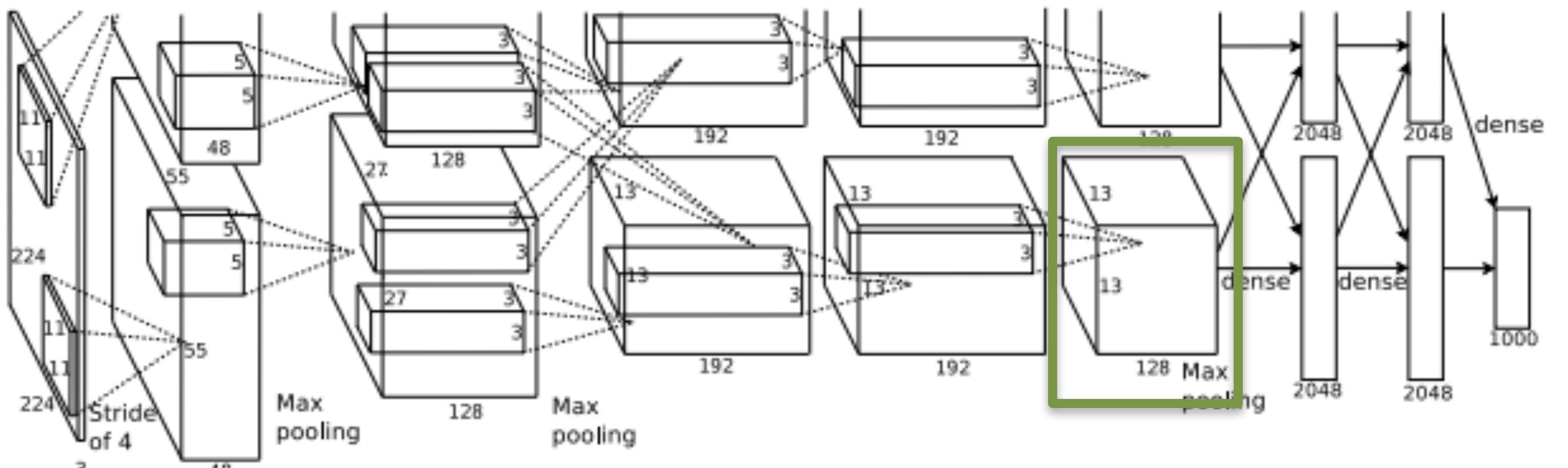
- **Visualize filter (kernels)**
- **Use CNNs as feature maps**
- **What are the different layers learning?**
- **Visualize activations at different layers**
- **Saliency versus Occlusion**
- **Guided backpropagation**
- **Gradient Ascent**
- **Feature Inversion**
- **DeepDreams**

Visualizing Activations: Which Filters ‘Fire’?



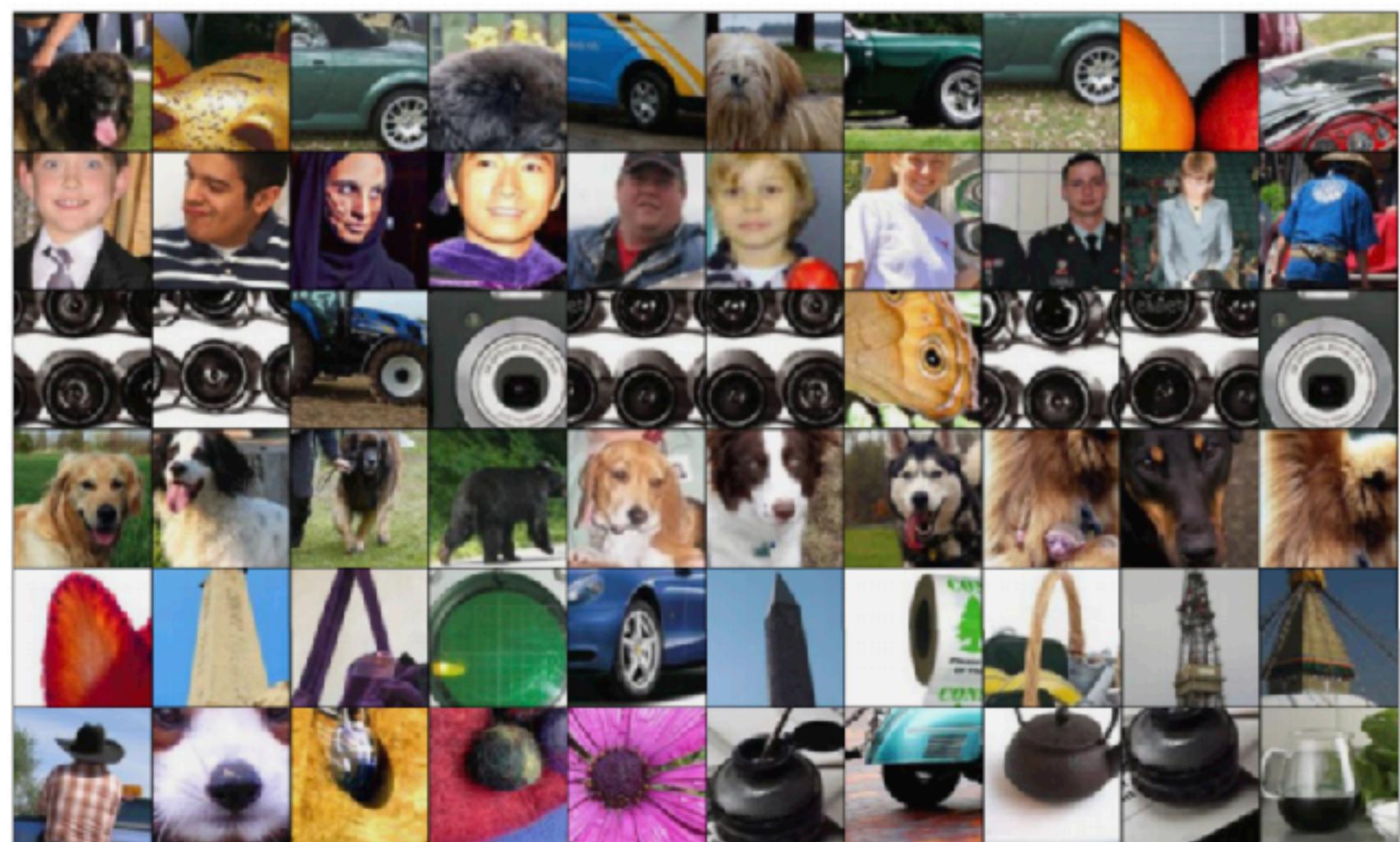
Maximally Active Filters

- Run many images through the network, record values of chosen channel.
- Visualize image patches that correspond to maximal activations.



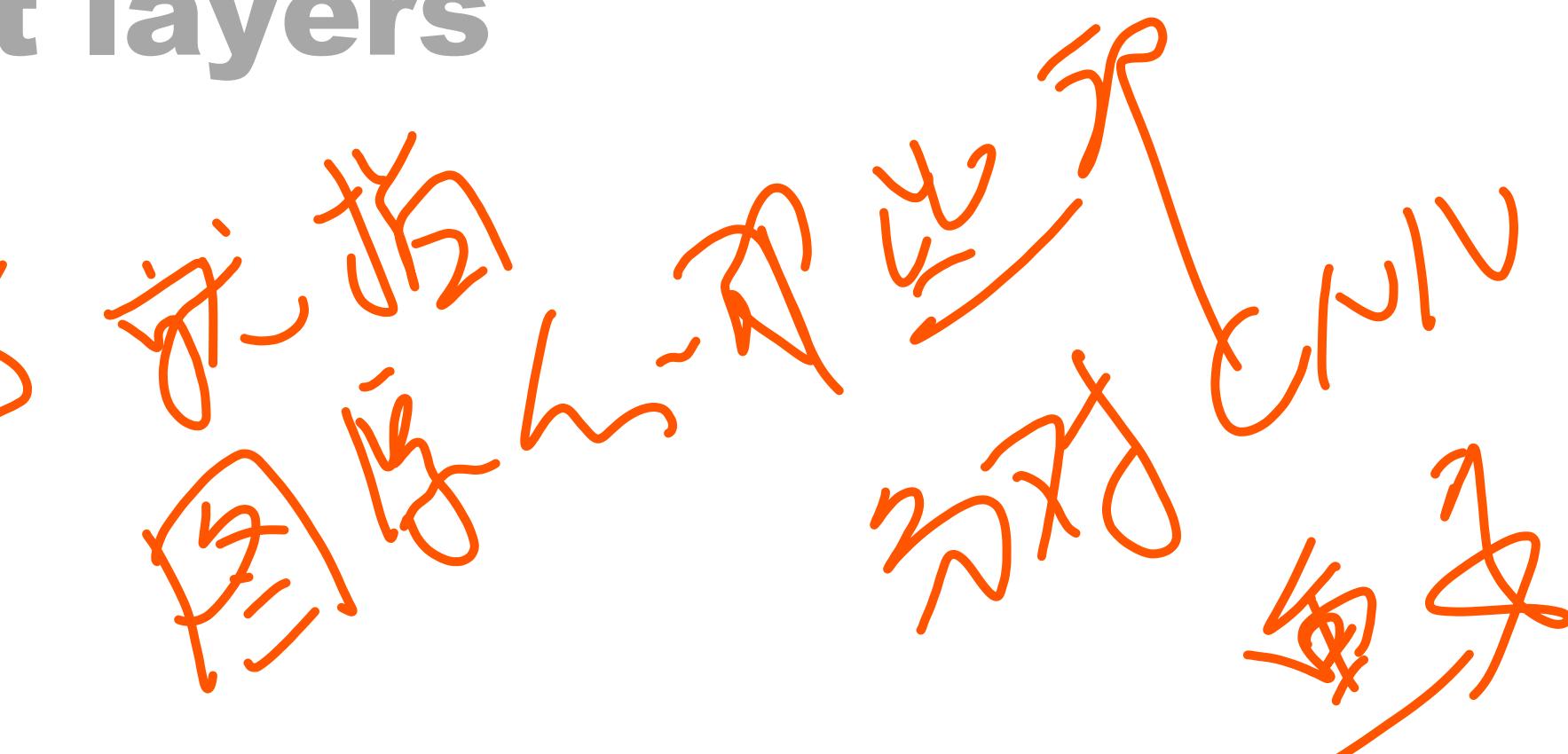
STRIVING FOR SIMPLICITY:
THE ALL CONVOLUTIONAL NET

Jost Tobias Springenberg*, Alexey Dosovitskiy*, Thomas Brox, Martin Riedmiller
Department of Computer Science
University of Freiburg
Freiburg, 79110, Germany
{springj, dosovits, brox, riedmiller}@cs.uni-freiburg.de



Plan for Today

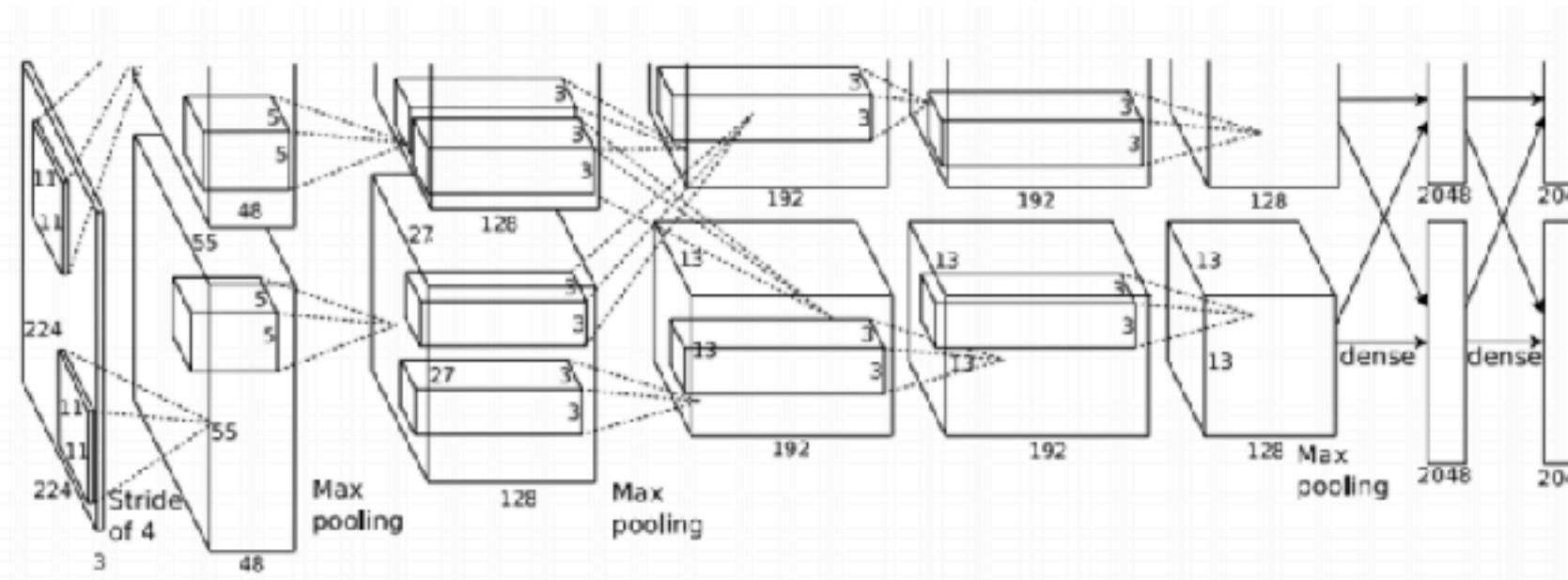
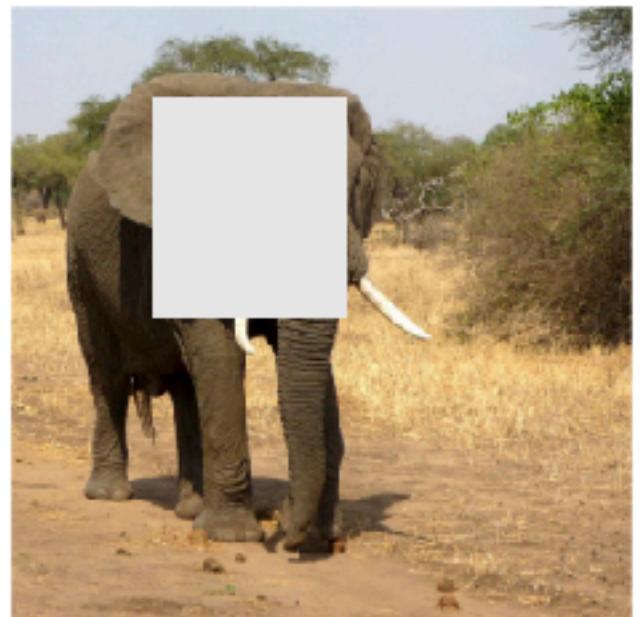
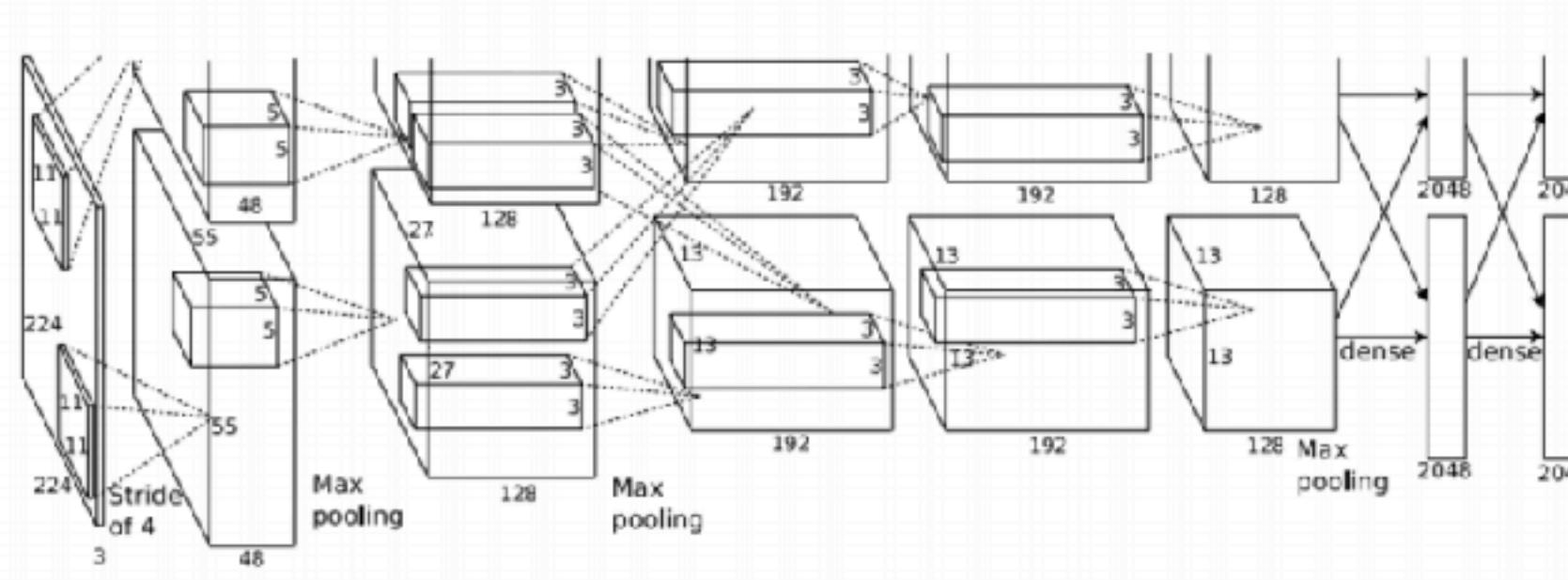
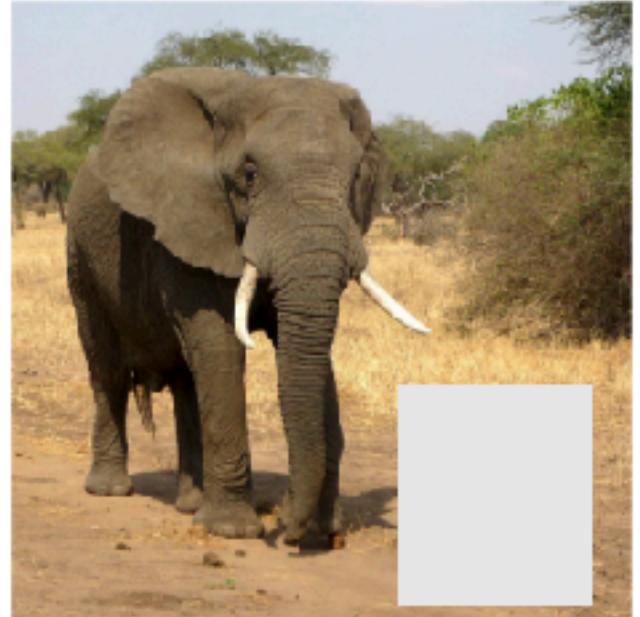
- Visualize filter (kernels)
- Use CNNs as feature maps
- What are the different layers learning?
- Visualize activations at different layers
- **Saliency versus Occlusion**
- **Guided backpropagation**
- Gradient Ascent
- Feature Inversion
- DeepDreams



題為什

Saliency vs Occlusion: Which Pixels Matter?

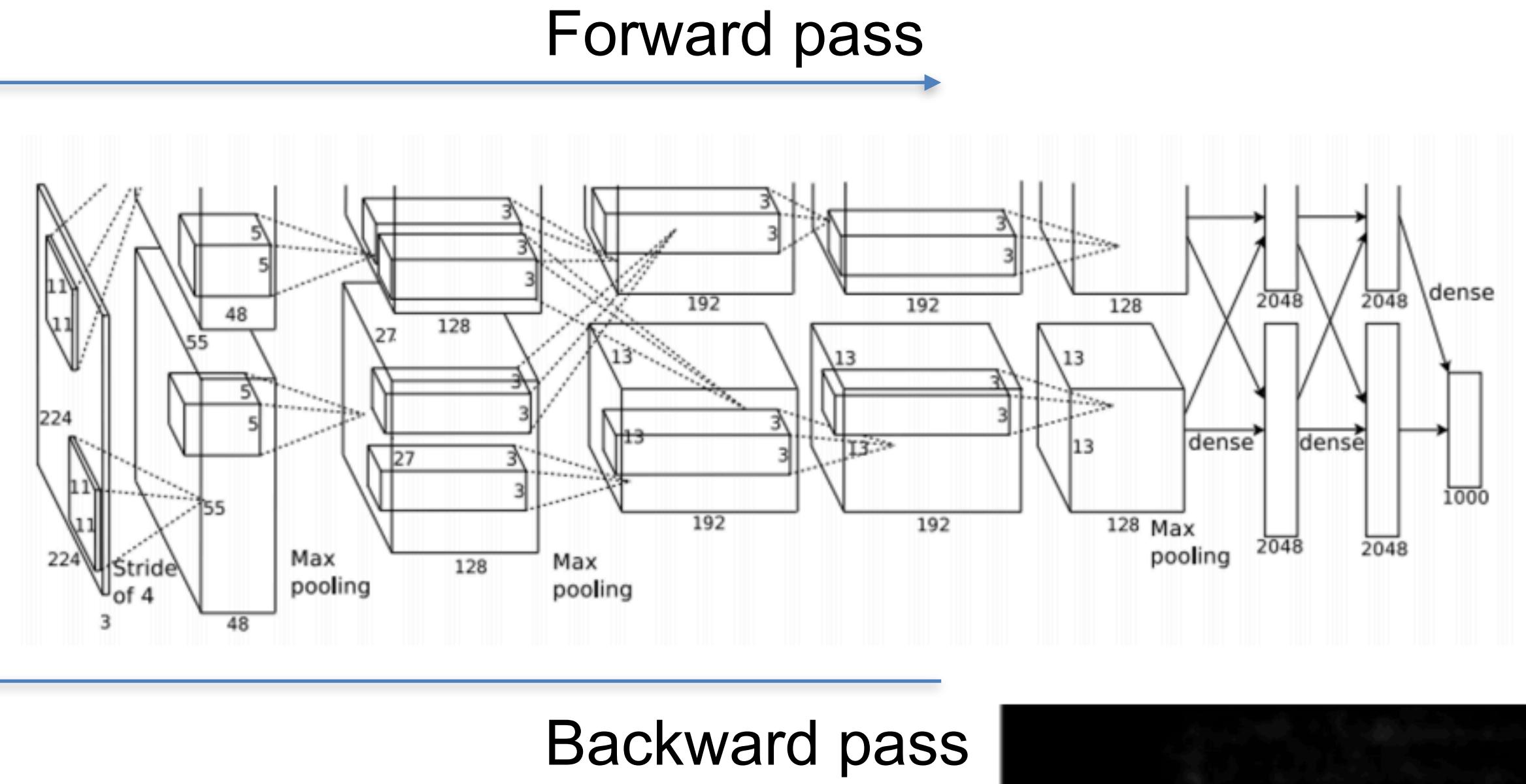
- Mask out pixels before sending to CN
- Check sensitivity of class probability



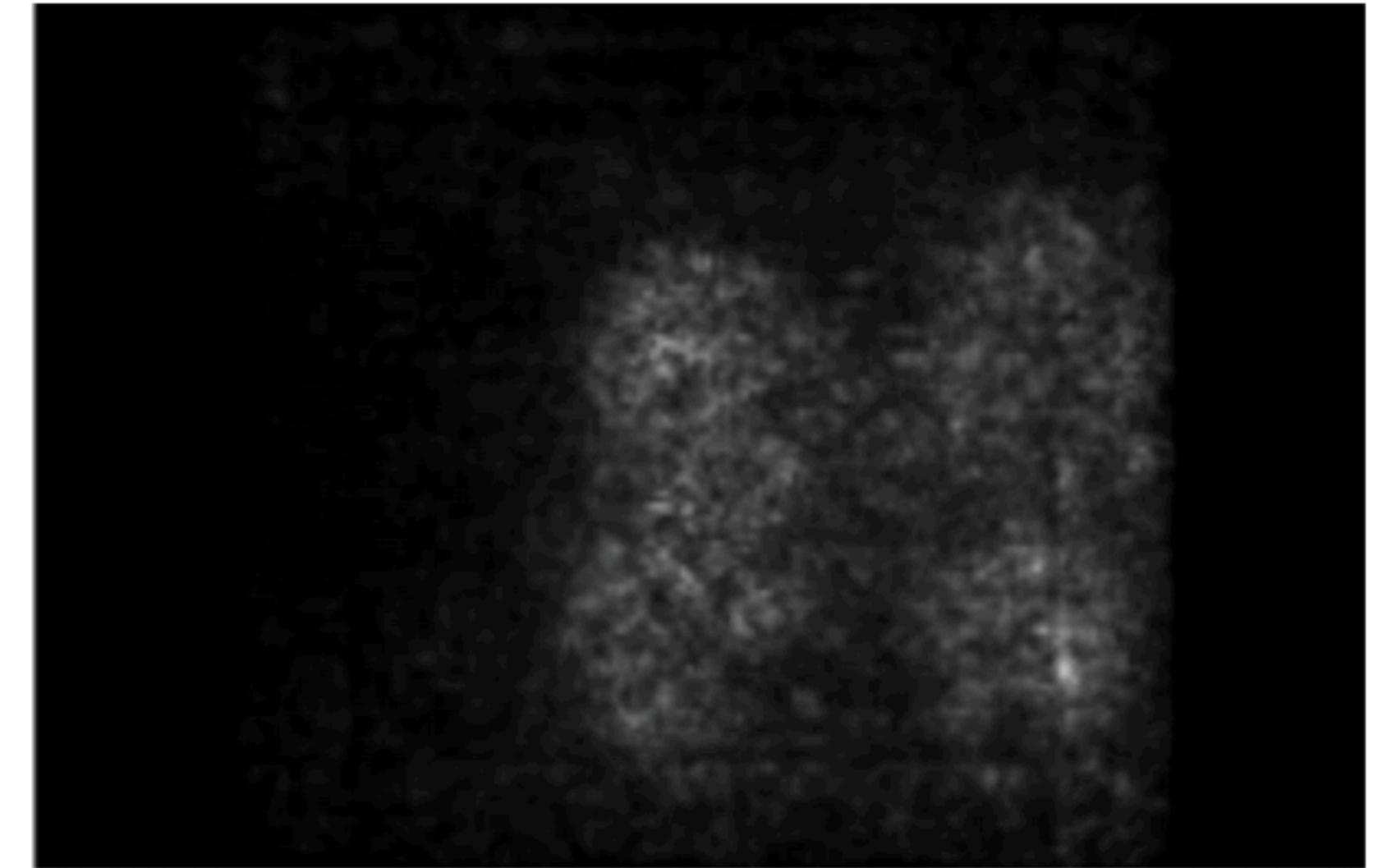
Visualizing and Understanding Convolutional Networks



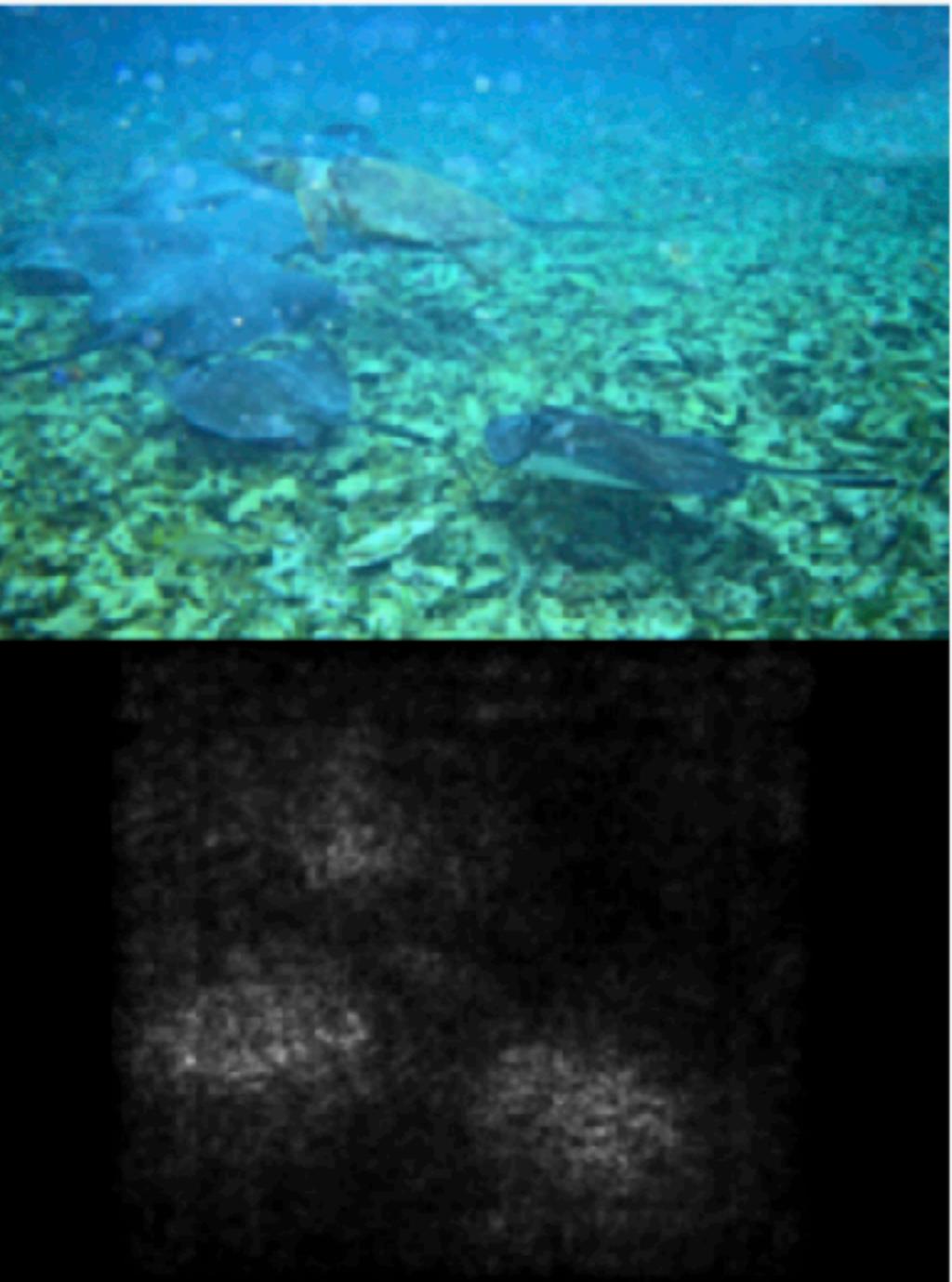
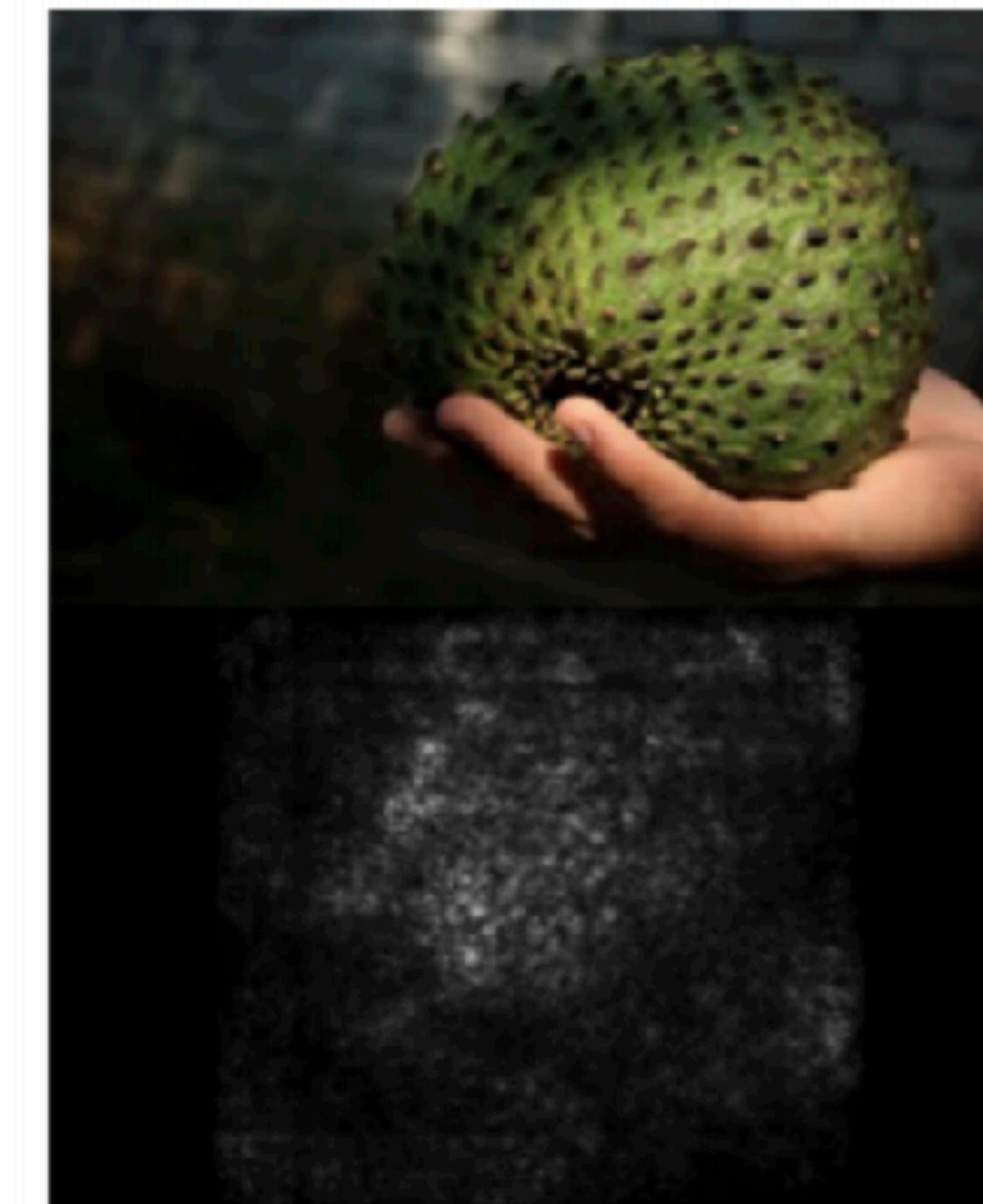
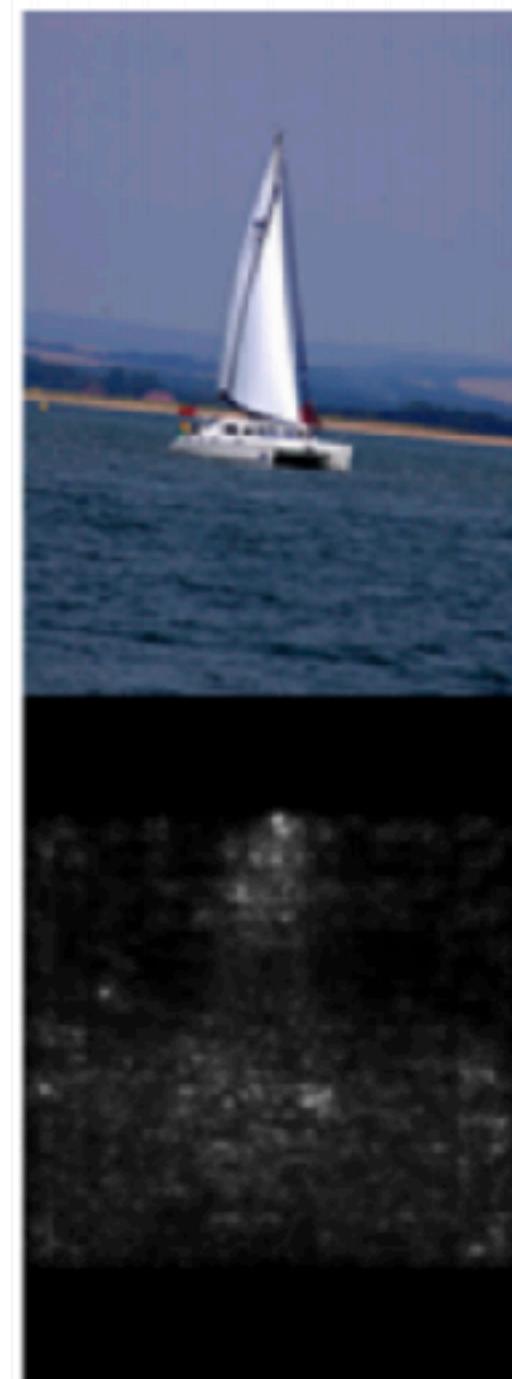
Saliency via Backprop



- **Compute gradients wrt image pixels**
- **Take absolute values**
- **Max over RGB channels**



Saliency via Backprop



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

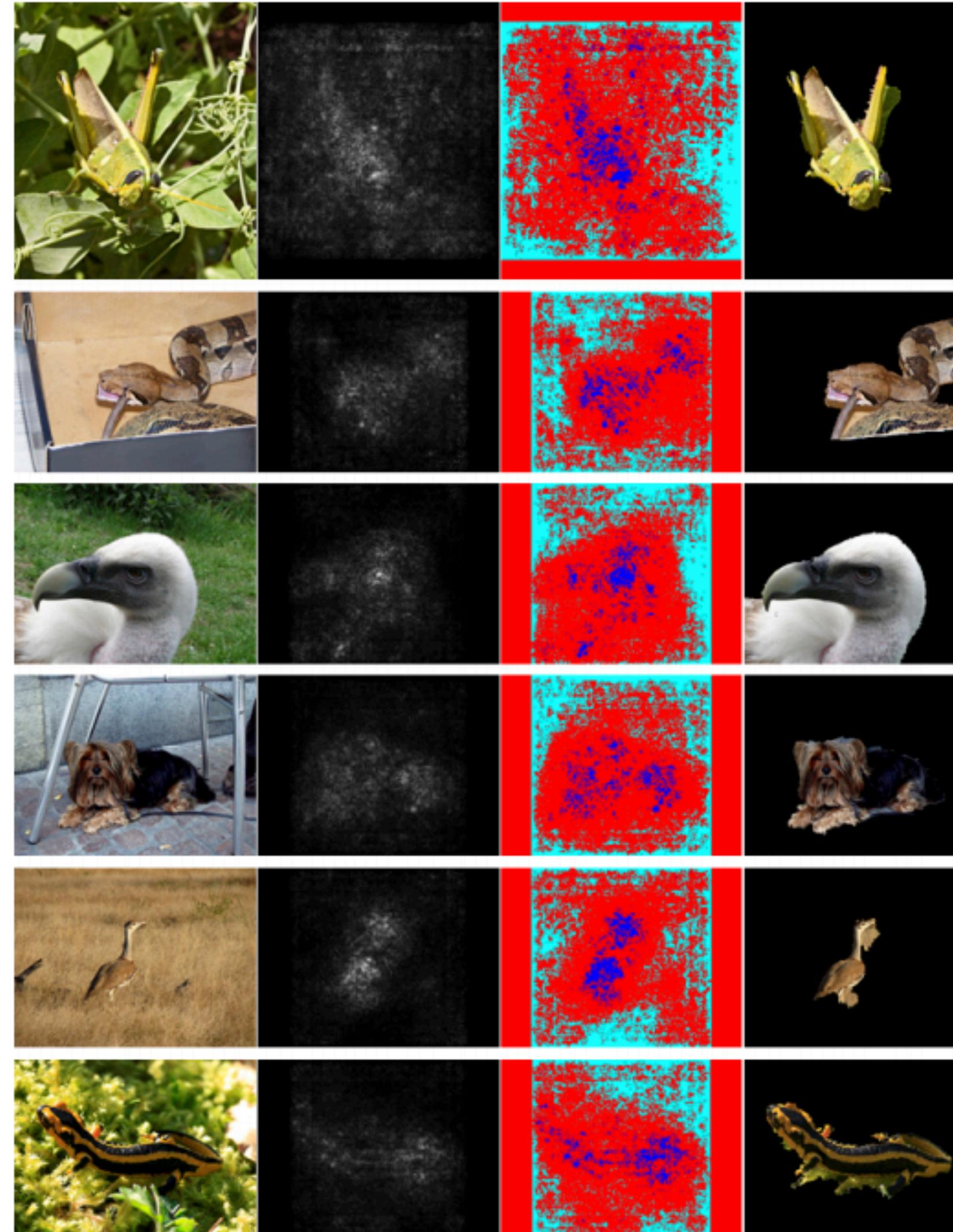
Karen Simonyan

Andrea Vedaldi

Andrew Zisserman

Visual Geometry Group, University of Oxford
{karen, vedaldi, az}@robots.ox.ac.uk

Segmentation w/o Supervision



GrabCut on saliency maps

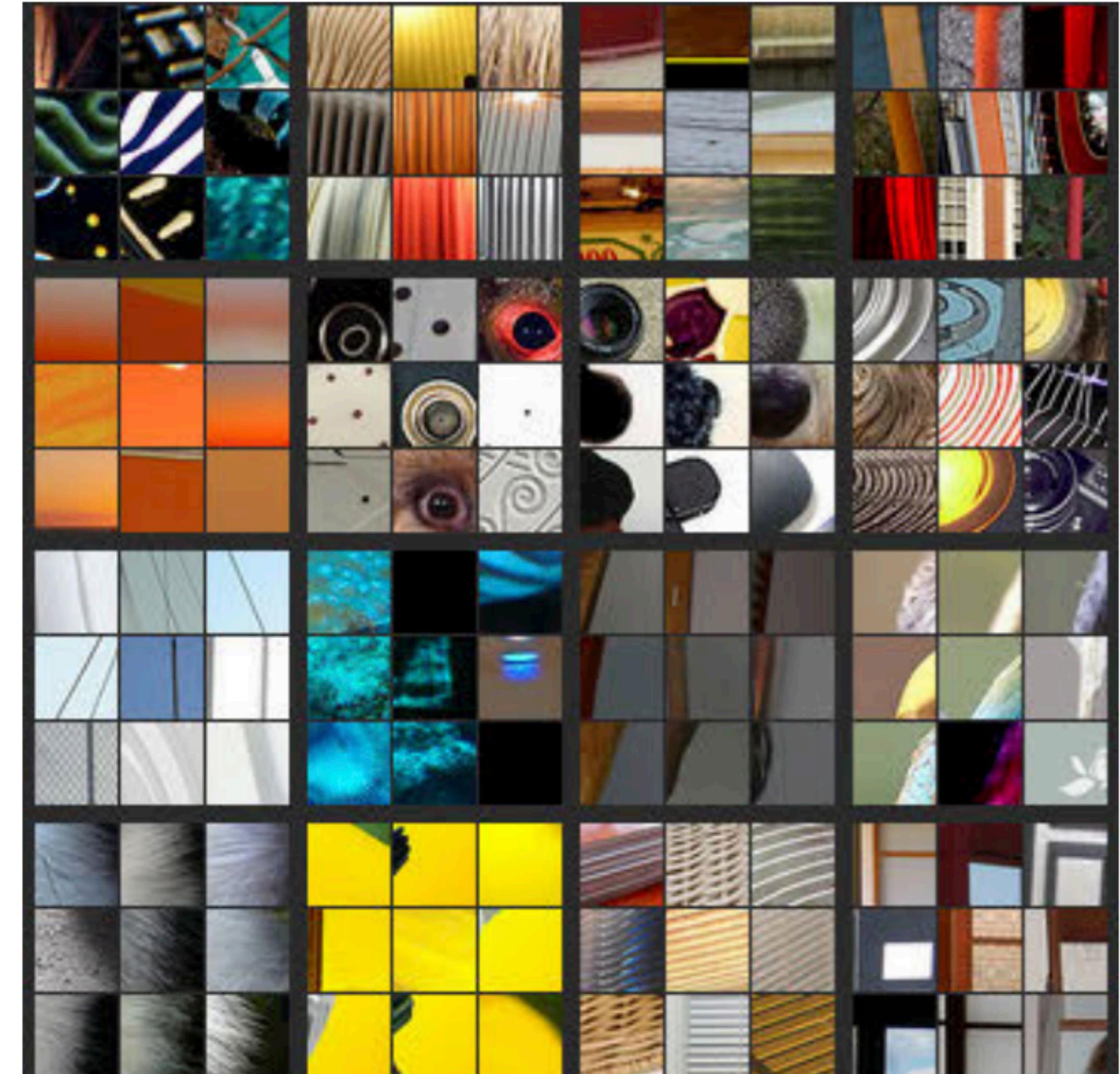
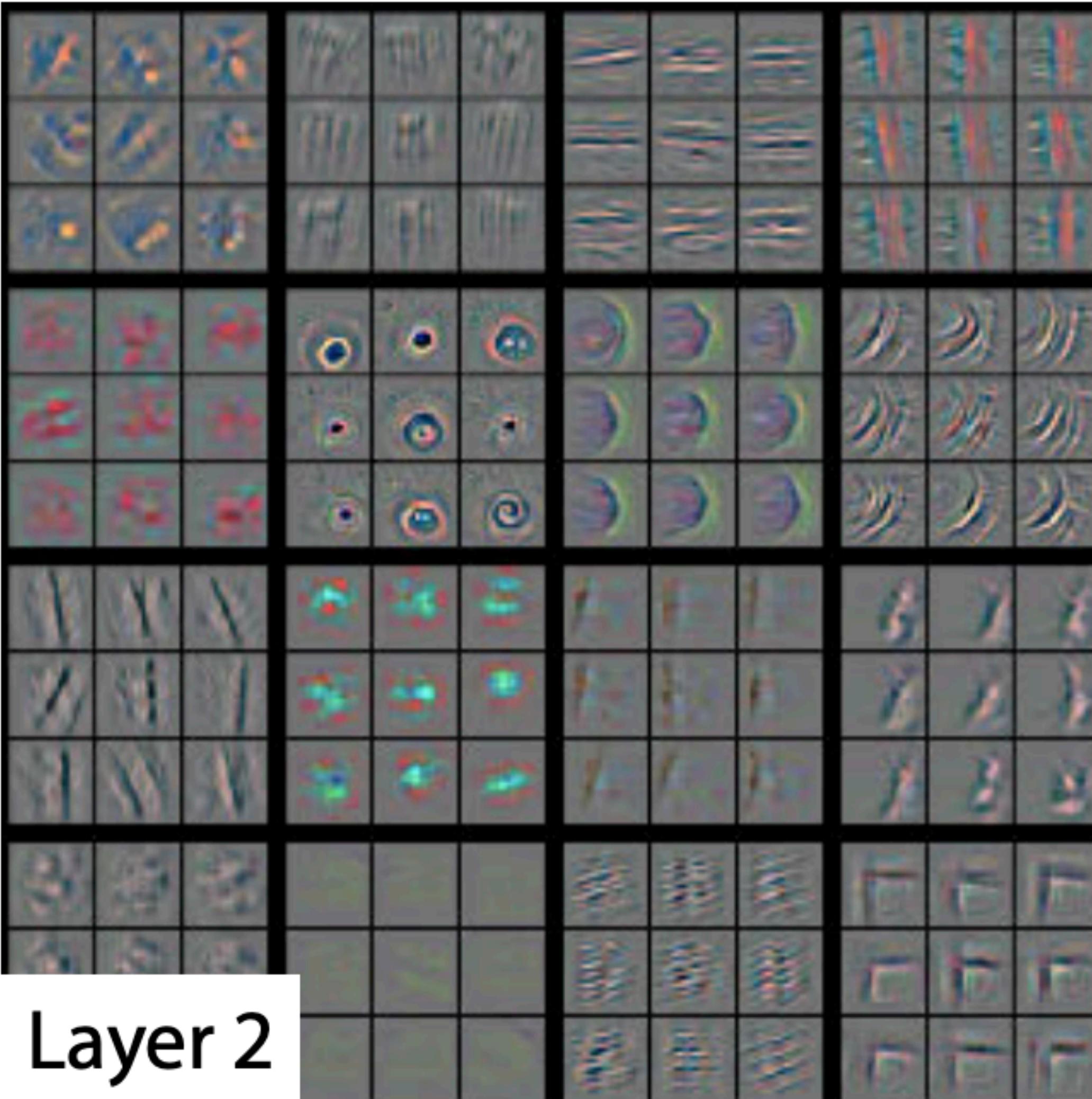
Plan for Today

- **Visualize filter (kernels)**
- **Use CNNs as feature maps**
- **What are the different layers learning?**
- **Visualize activations at different layers**
- **Saliency versus Occlusion**
- **Guided backpropagation**
- **Gradient Ascent**
- **Feature Inversion**
- **DeepDreams**

Guided Backprop

- **Pick single intermediate features
e.g., one channel of 128x13x13 conv5
feature map**
- **Compute gradients of neuron
values wrt image pixels**

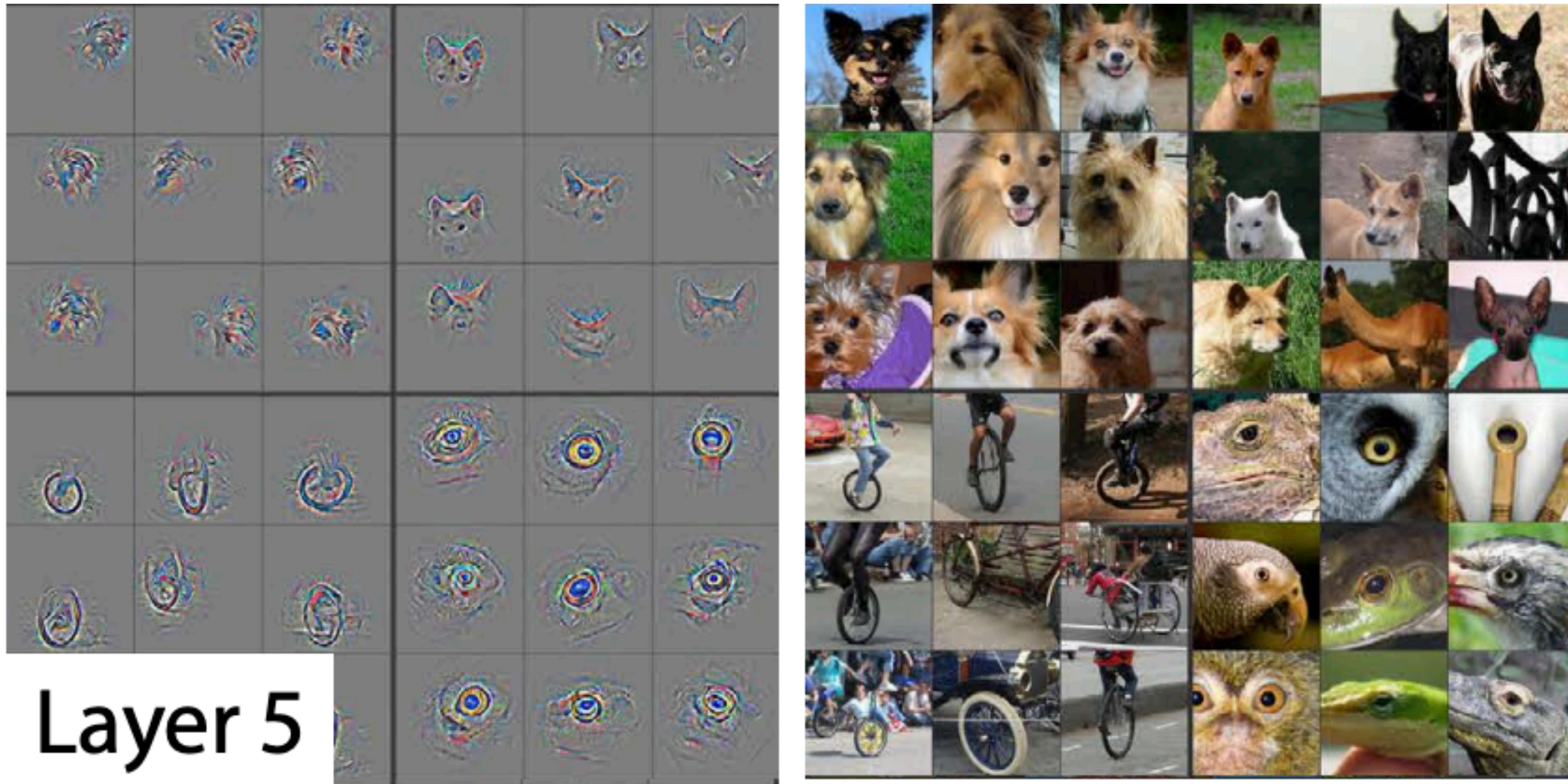
Guided Backprop (Layer 2)



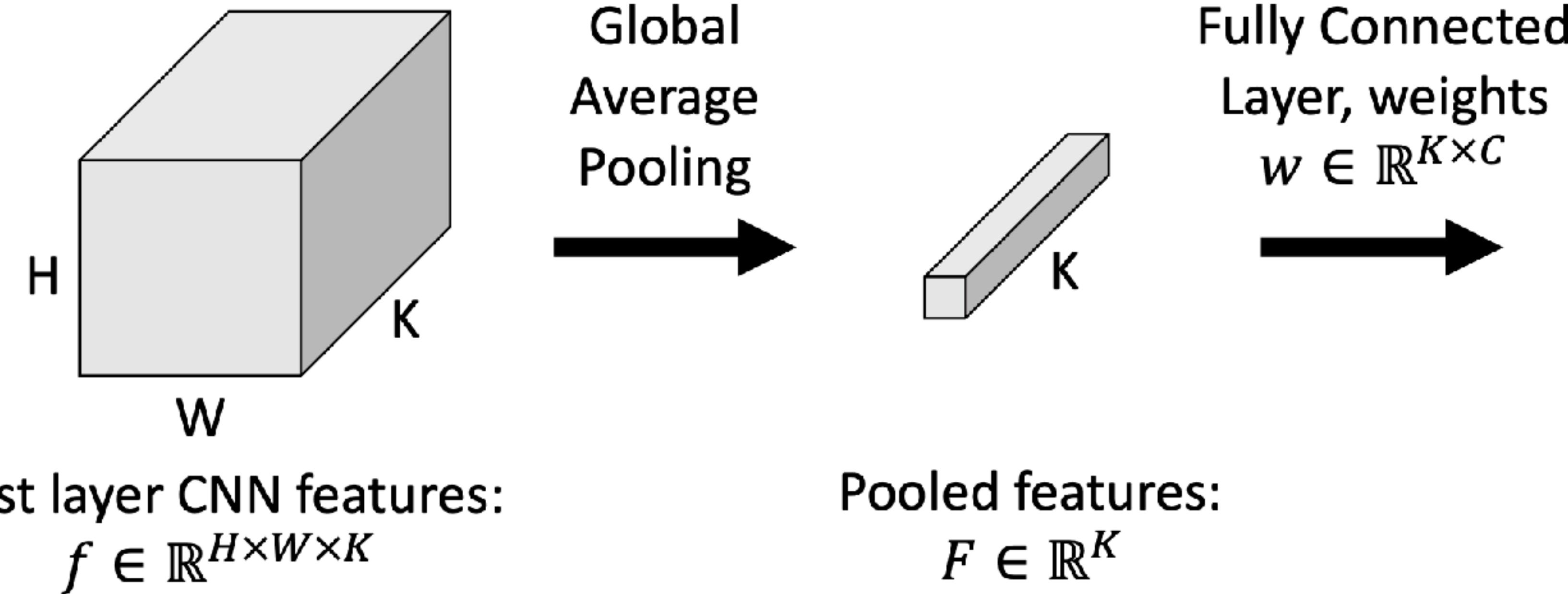
Guided Backprop (Layer 4)



Guided Backprop (Layer 5)



Class Activation Map (CAM)



$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k}$$

$$\begin{aligned} S_c &= \sum_k w_{k,c} F_k \\ &= \frac{1}{HW} \sum_{h,w} \sum_k w_{k,c} f_{h,w,k} \end{aligned}$$

$$M_{c,h,w} = \sum_k w_{k,c} f_{h,w,k}$$

Class Activation Map

Class Activation Maps

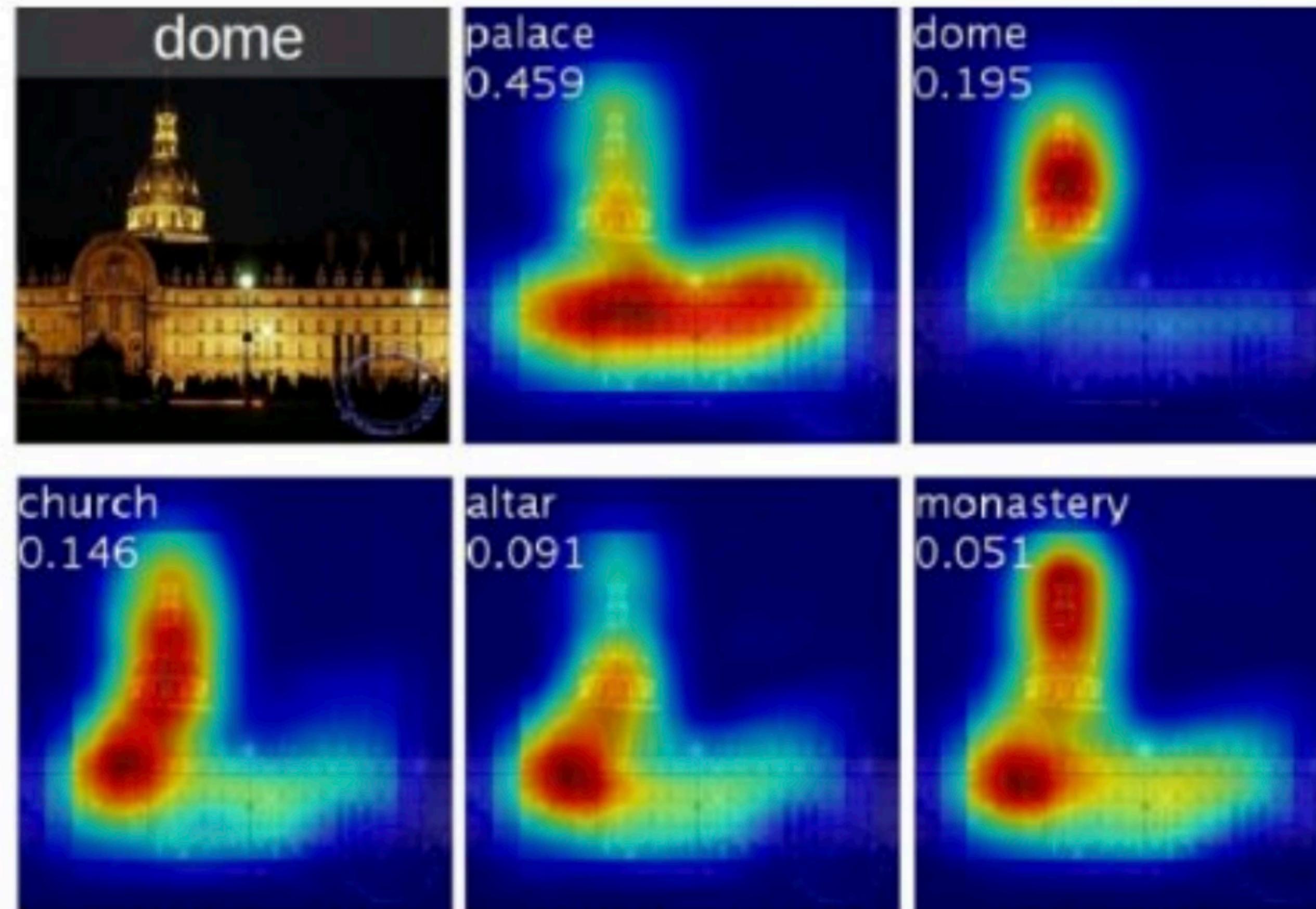
Brushing teeth



Cutting trees



Class Activation wrt Multiple Classes



Applied wrt only the last layer

Gradient-Weighted CAM (Grad-CAM)

1. Pick a layer with activations $A \in \mathbb{R}^{H \times W \times K}$

2. Compute gradient of class score S_c **with respect to** A , $\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$

3. Average pool the gradients to get weights $\alpha \in \mathbb{R}^K$, $\alpha_k = \frac{1}{HW} \sum_{h,w} \frac{\partial S_c}{\partial A_{h,w,k}}$

4. Compute weighted activation maps $M_{h,w}^c = \text{ReLU} \left(\sum_k \alpha_k A_{h,w,k} \right)$

Grad-CAM

(but noise will
get worse)

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

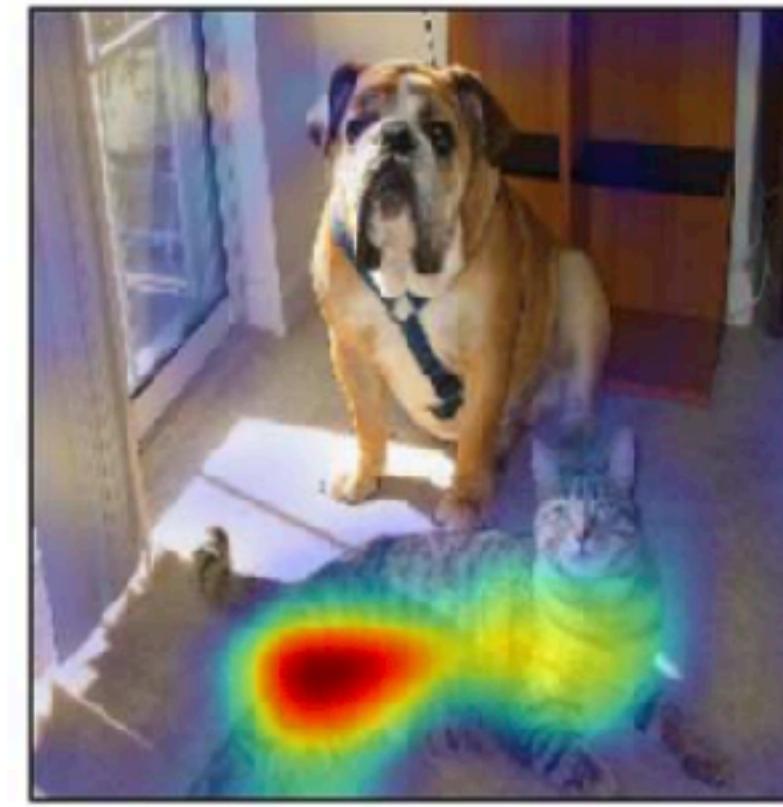
Ramprakash R. Selvaraju · Michael Cogswell · Abhishek Das · Ramakrishna Vedantam · Devi Parikh · Dhruv Batra



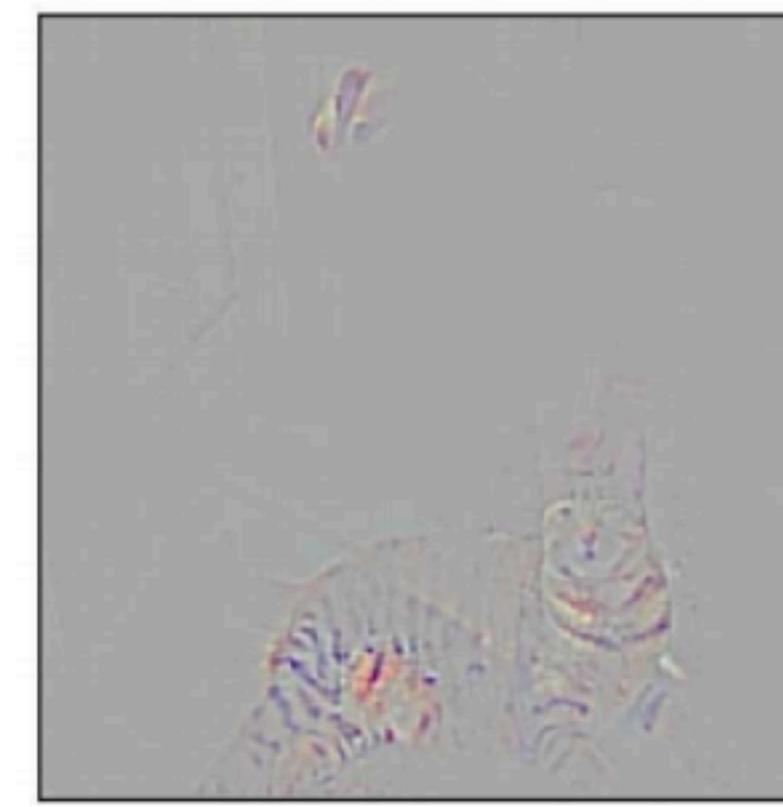
(a) Original Image



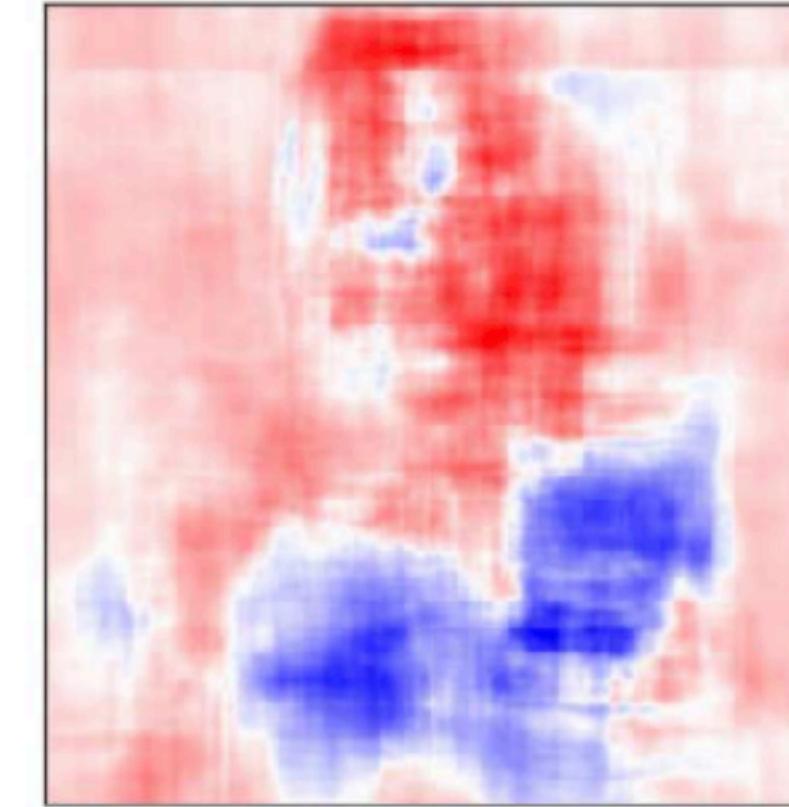
(b) Guided Backprop ‘Cat’



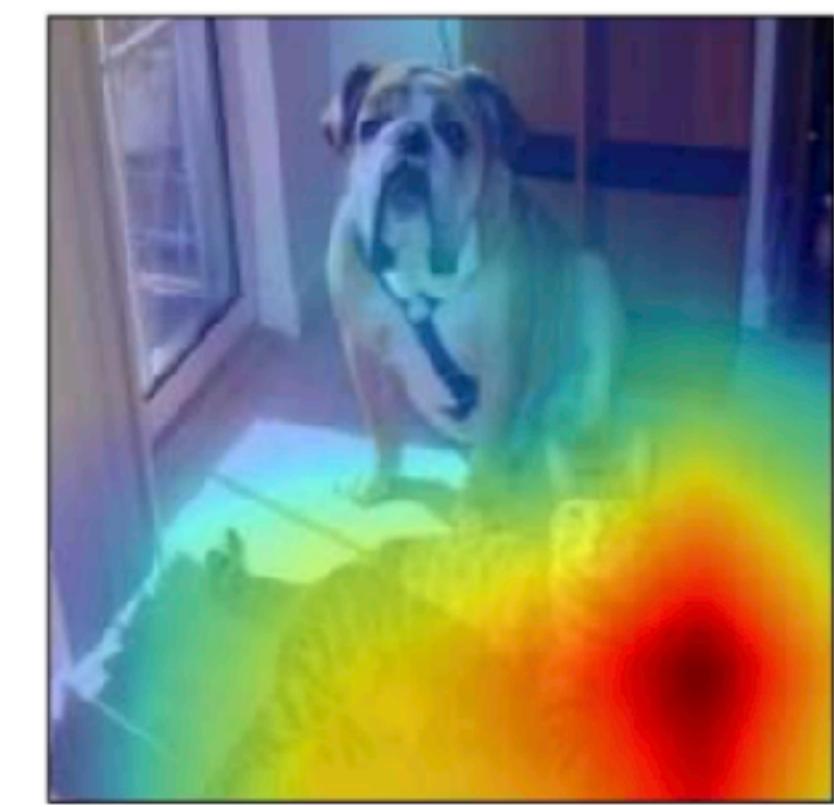
(c) Grad-CAM ‘Cat’



(d) Guided Grad-CAM ‘Cat’



(e) Occlusion map for ‘Cat’



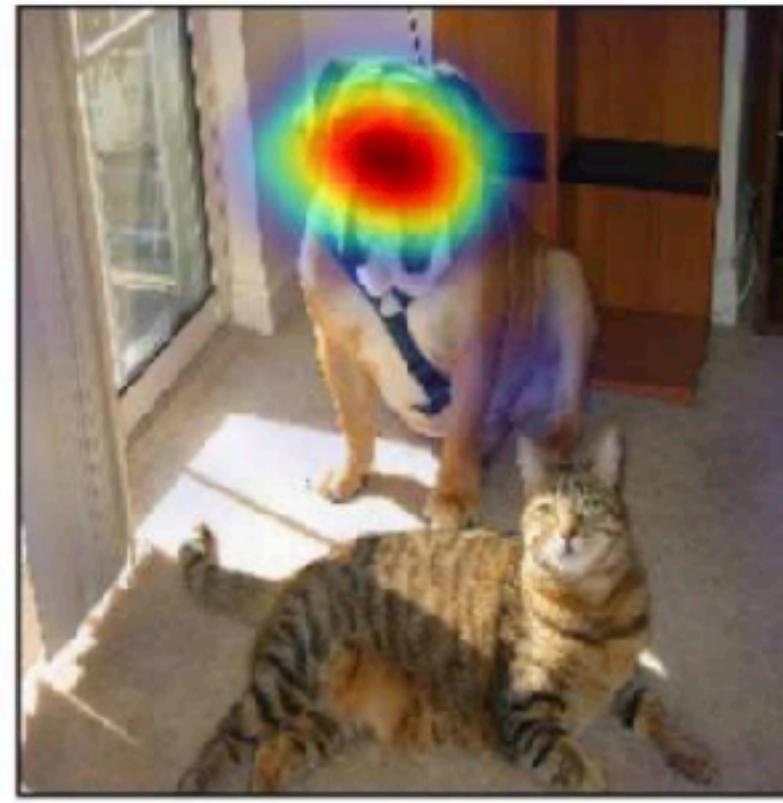
(f) ResNet Grad-CAM ‘Cat’



(g) Original Image



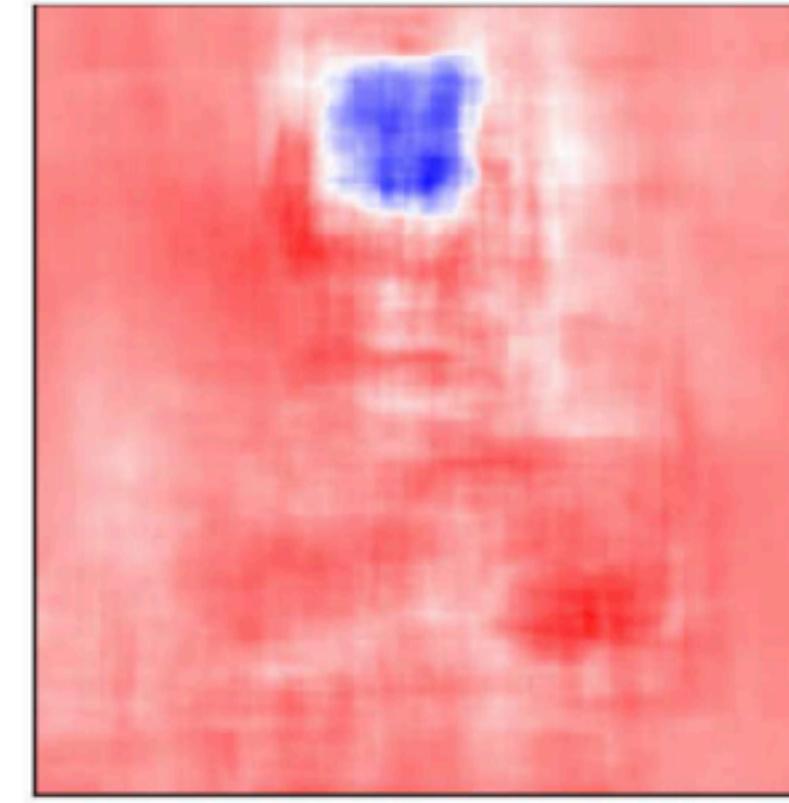
(h) Guided Backprop ‘Dog’



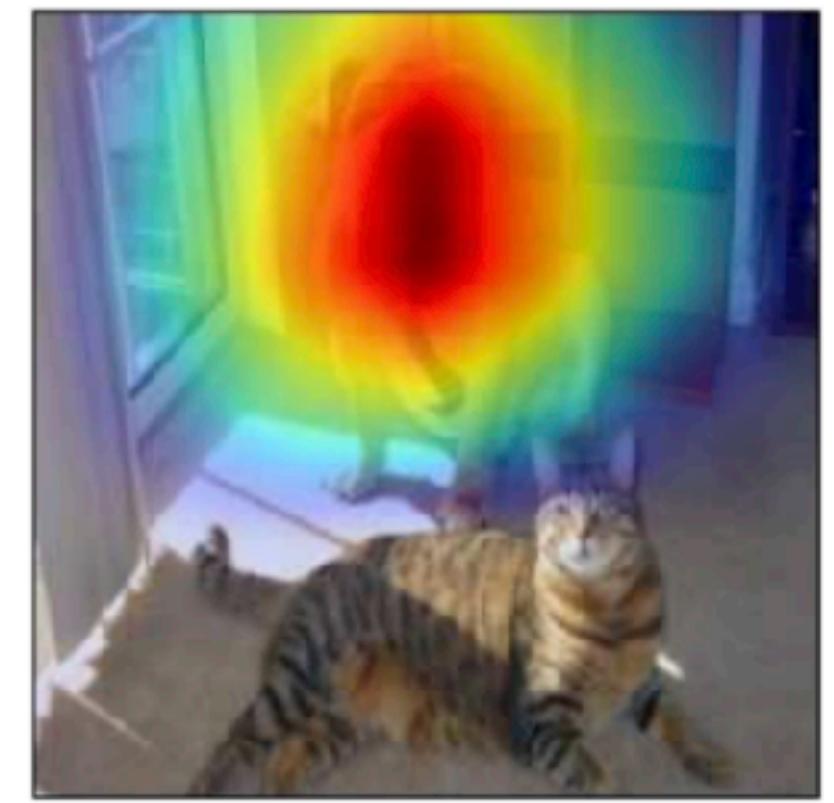
(i) Grad-CAM ‘Dog’



(j) Guided Grad-CAM ‘Dog’

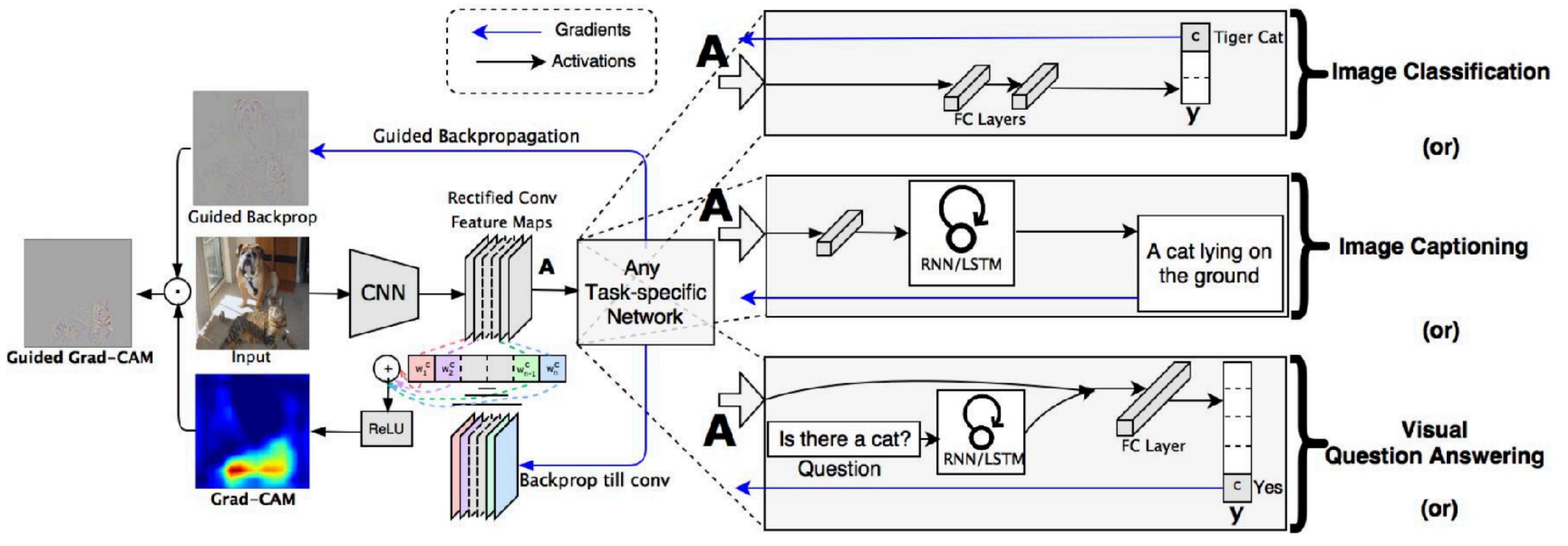


(k) Occlusion map for ‘Dog’

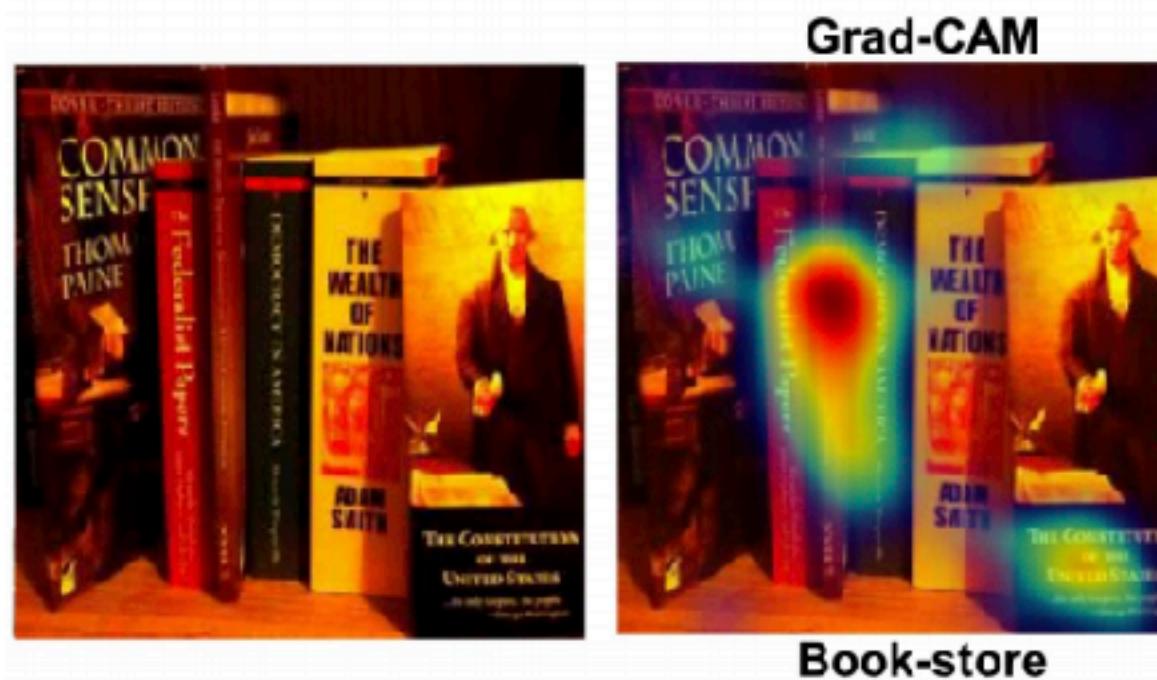


(l) ResNet Grad-CAM ‘Dog’

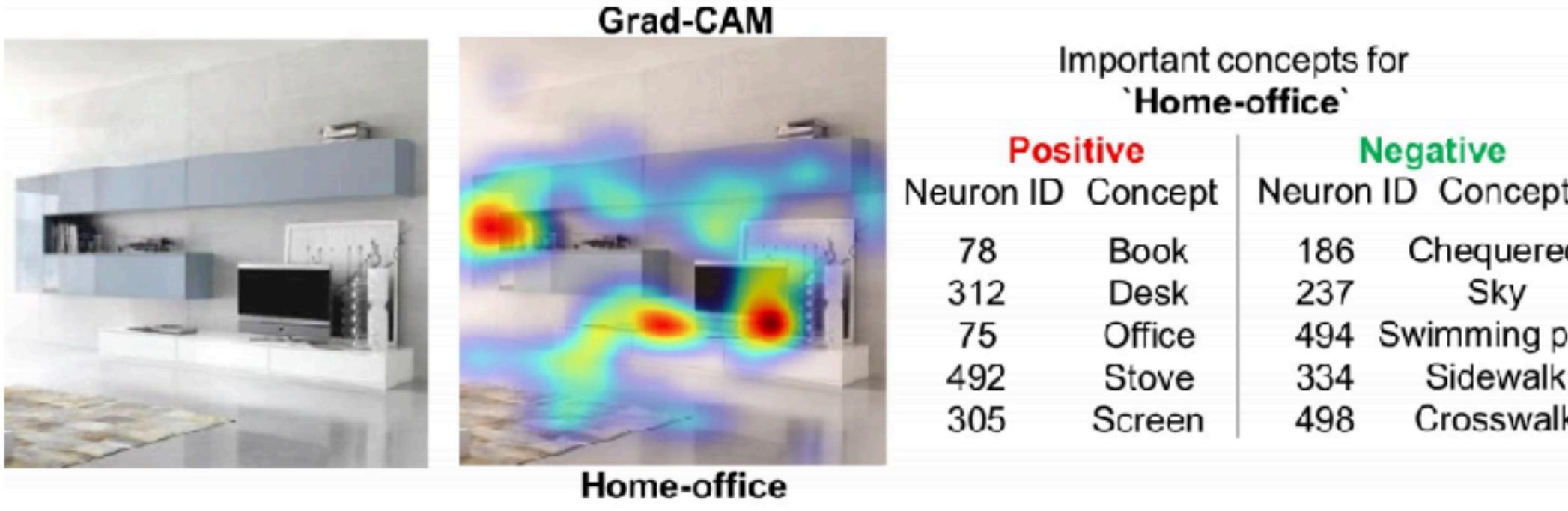
Guided-CAM



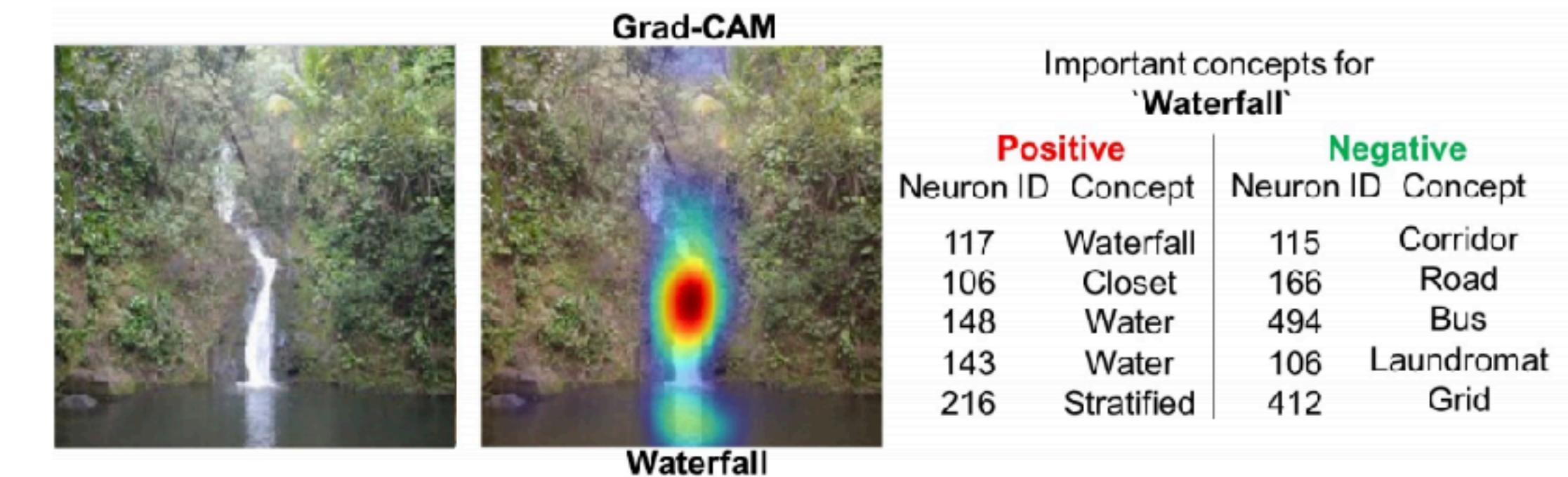
Grad-CAM on Image Captioning



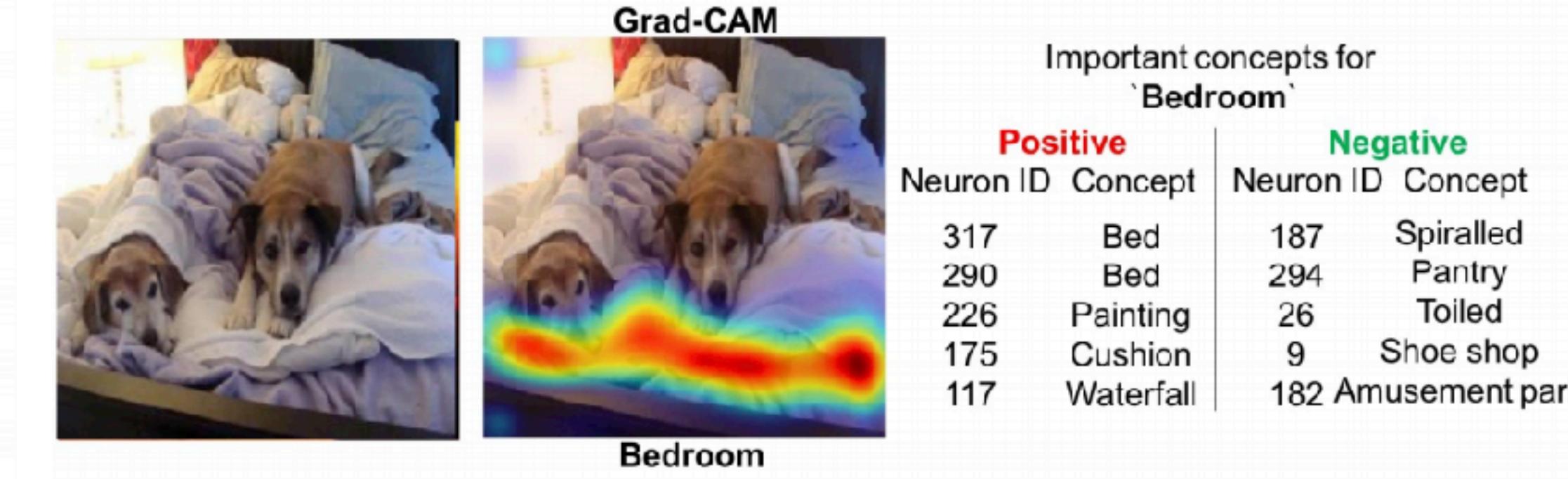
(a)



(c)



(b)



(d)

© 2018 University College London. All rights reserved. This document contains neither recommendations nor conclusions of UCL. UCL accepts no responsibilities for any inaccuracies.

Plan for Today

- **Visualize filter (kernels)**
- **Use CNNs as feature maps**
- **What are the different layers learning?**
- **Visualize activations at different layers**
- **Saliency versus Occlusion**
- **Guided backpropagation**
- **Gradient Ascent**
- **Feature Inversion**
- **DeepDreams**

Visualizing CNN Features: Gradient Ascent

- **Guided Backprop:**
Find parts of image that maximally responds to a neuron
- **Gradient Ascent:**
Find a synthetic image that maximally activates a neuron

$$I^* = \arg \max_I f(I) + \lambda R(I)$$

Gradient Ascent

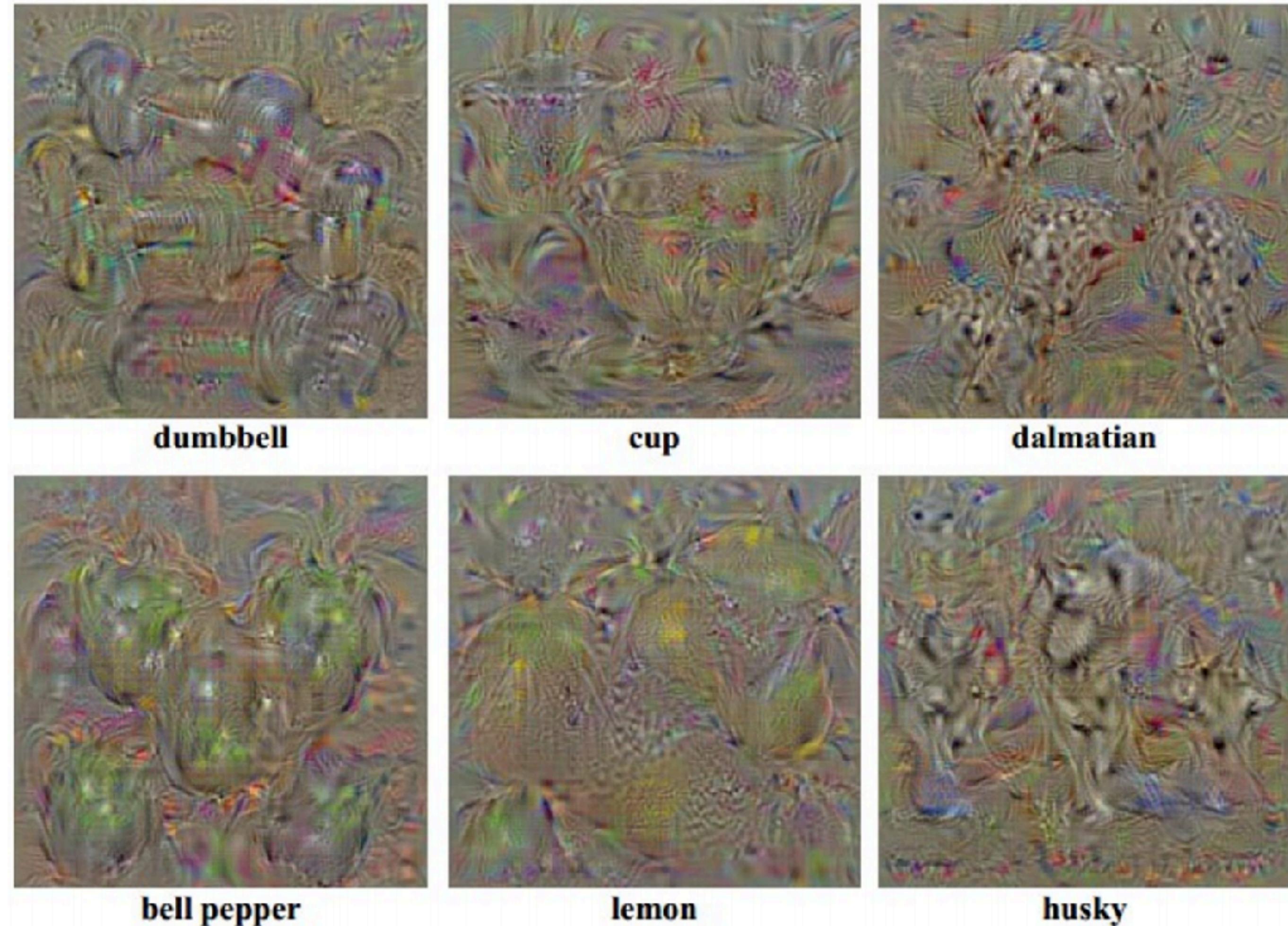
$$I^* = \arg \max_I \{S_c(I) - \lambda \|I\|^2\}$$

- Initialize image $I = 0$
- Forward image to compute scores $S_c(I)$
- Backprop to get gradients of neuron values wrt image pixels
- Update image as $I \leftarrow I + \frac{\partial S_c(I)}{\partial I}$

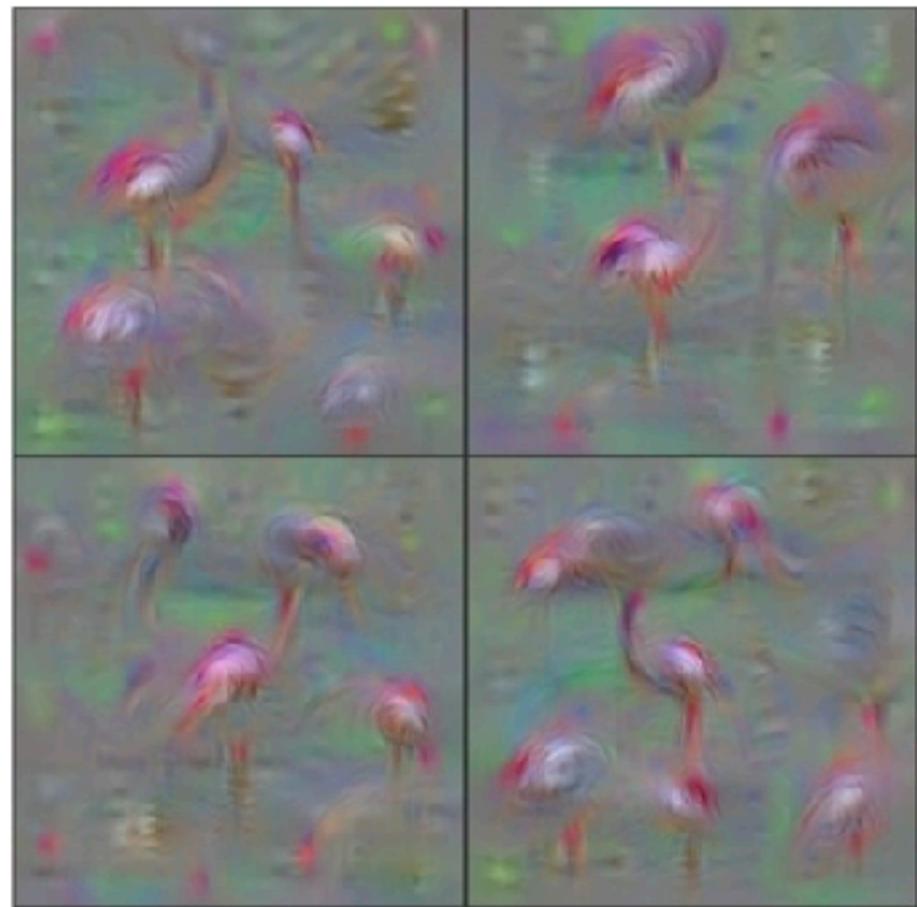
Gradient Ascent

Other regularizers:

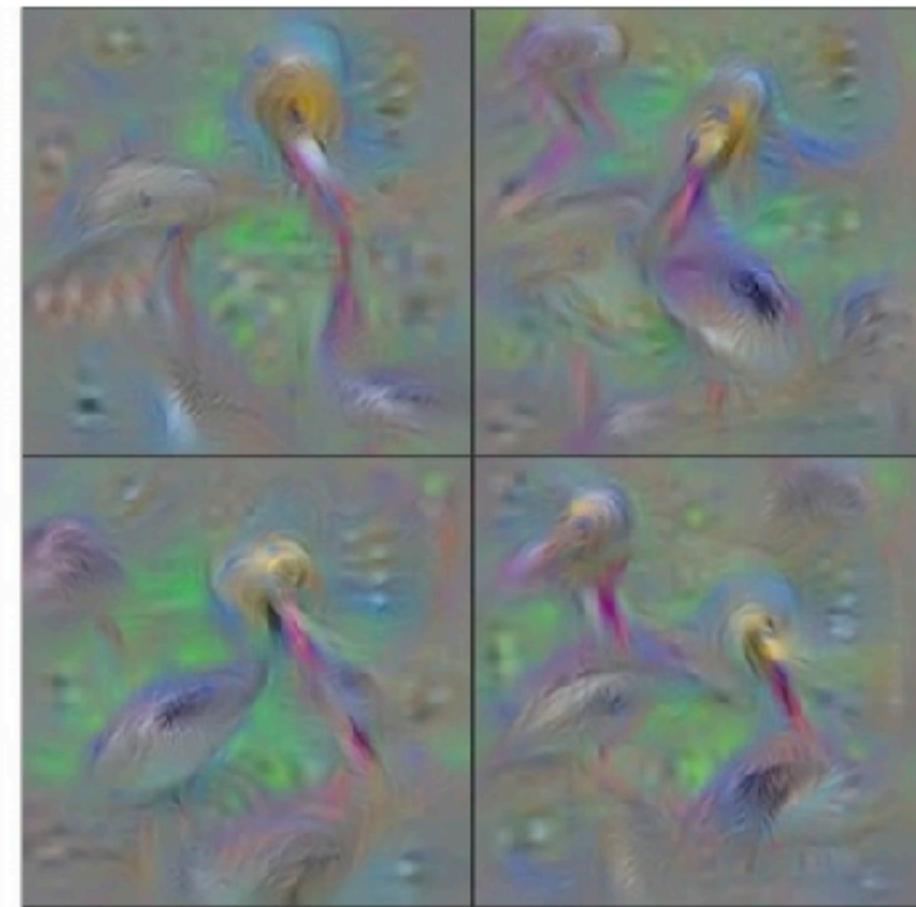
1. Gaussian blur
2. Clip pixels with values ~ 0
3. Clip pixels with gradients ~ 0



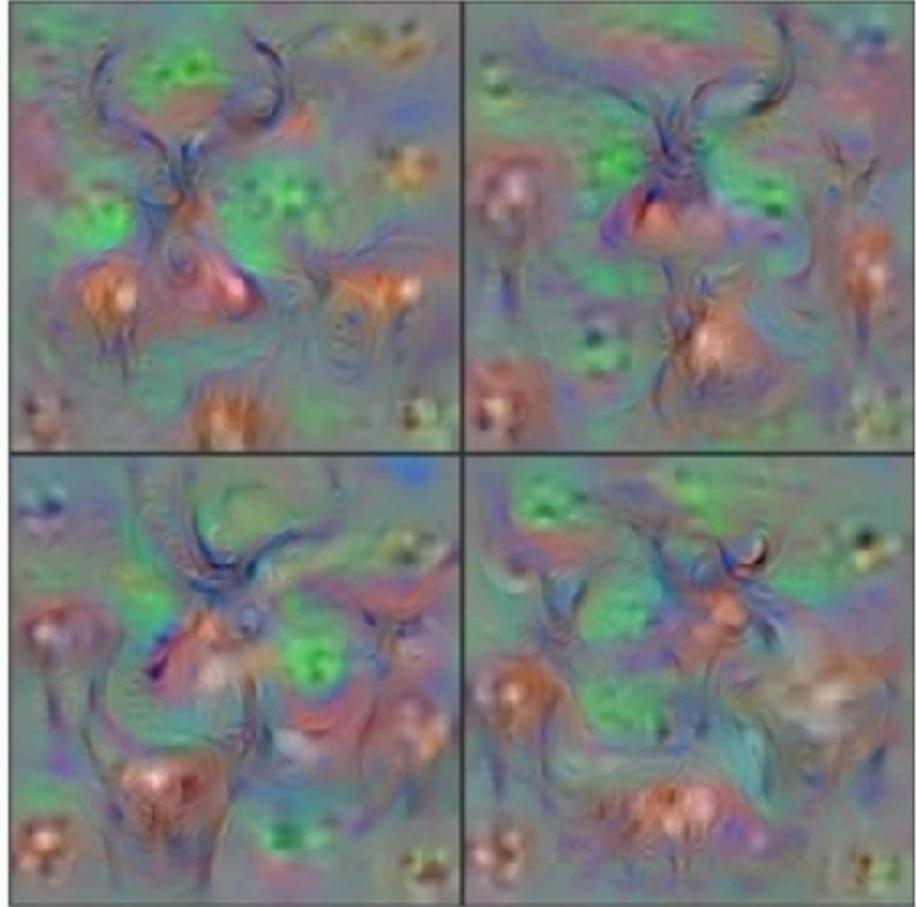
GA Results: More



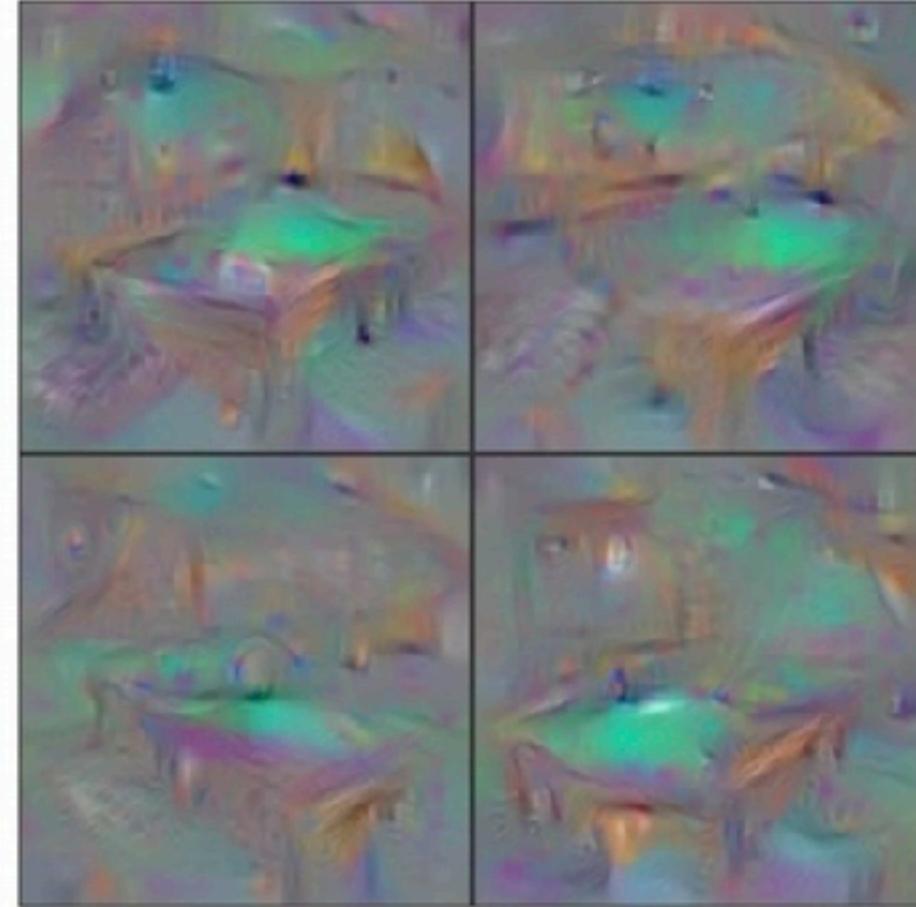
Flamingo



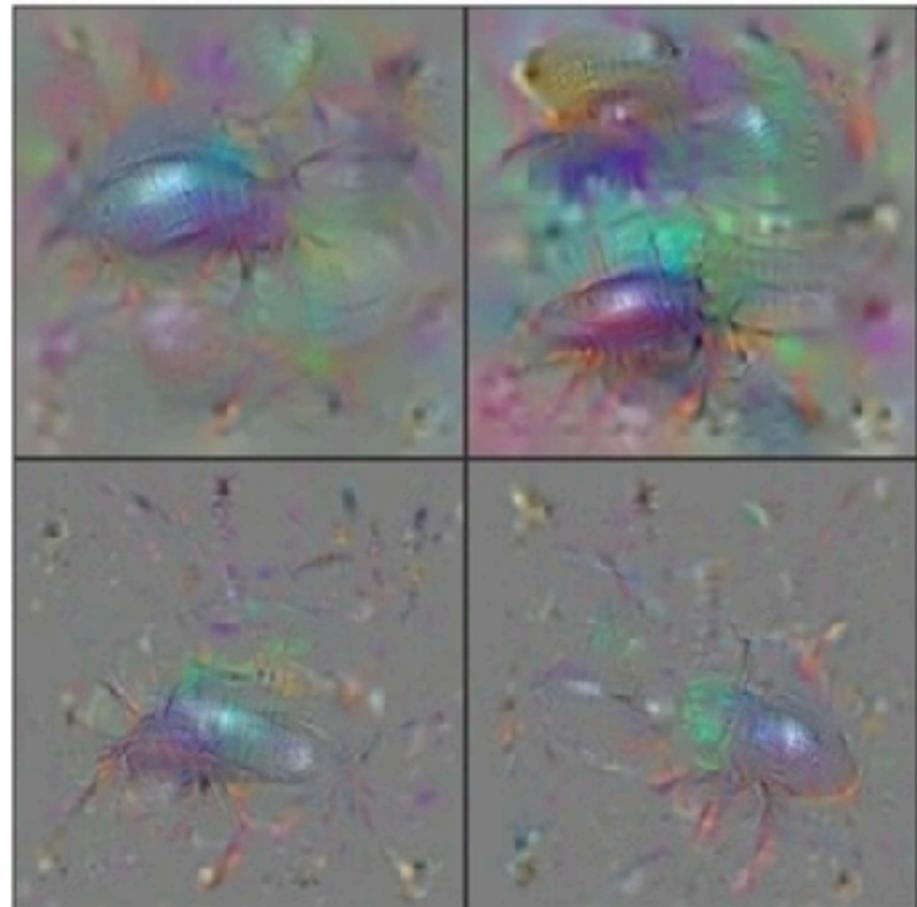
Pelican



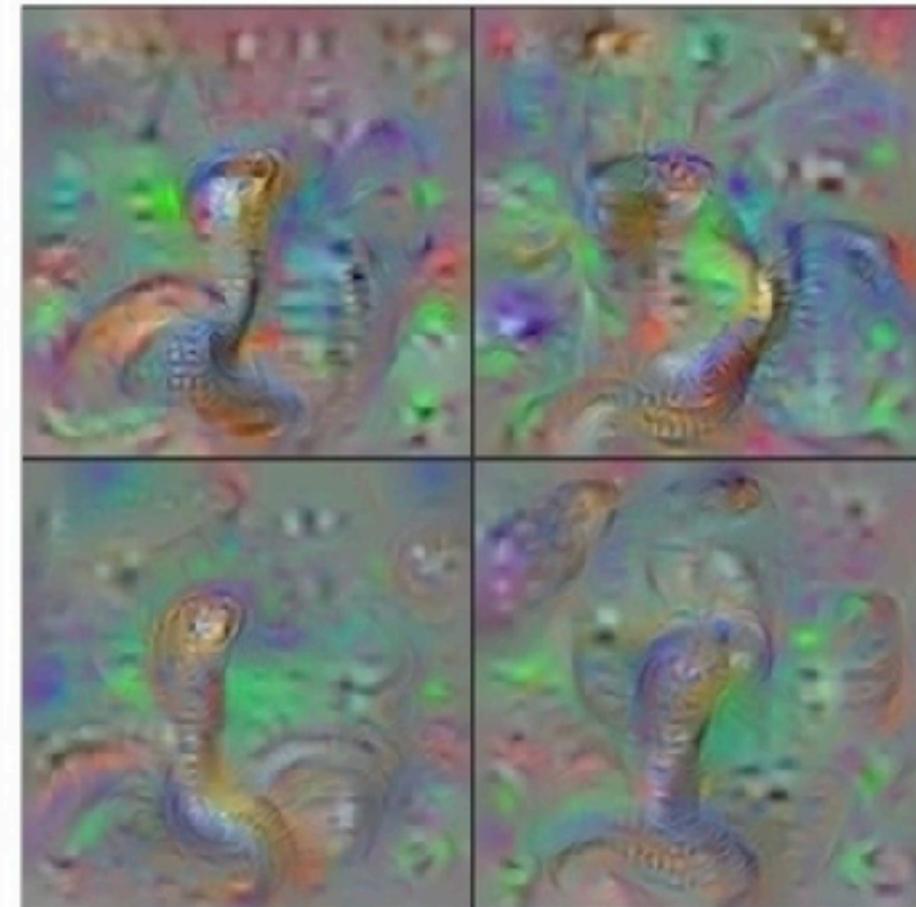
Hartebeest



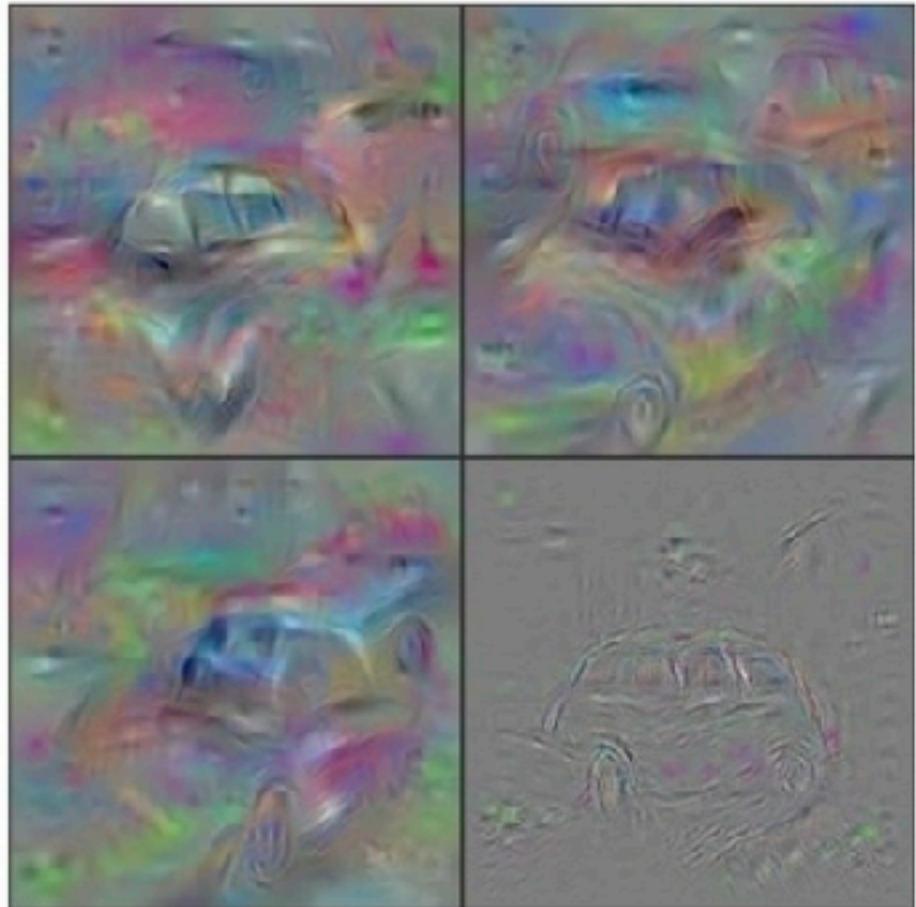
Billiard Table



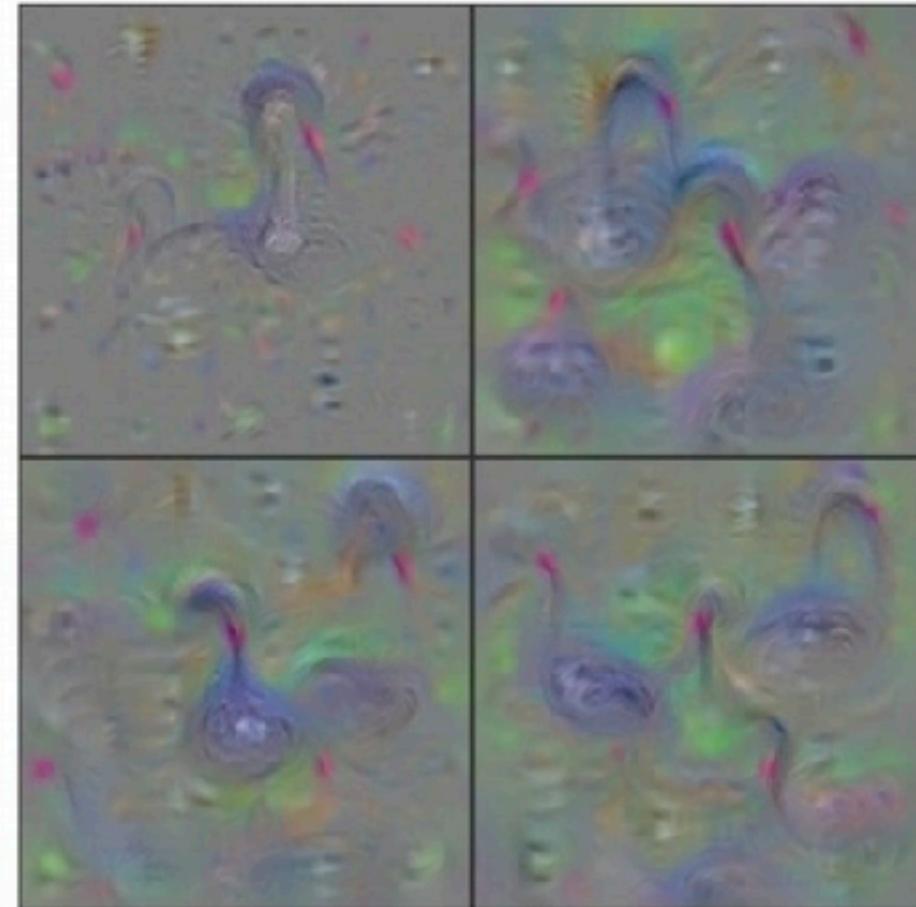
Ground Beetle



Indian Cobra



Station Wagon



Black Swan

Visualize Intermediate Features

Understanding Neural Networks Through Deep Visualization

Jason Yosinski
Cornell University

YOSINSKI@CS.CORNELL.EDU

Jeff Clune
Anh Nguyen
University of Wyoming

JEFFCLUNE@UWYO.EDU
ANGUYEN8@UWYO.EDU

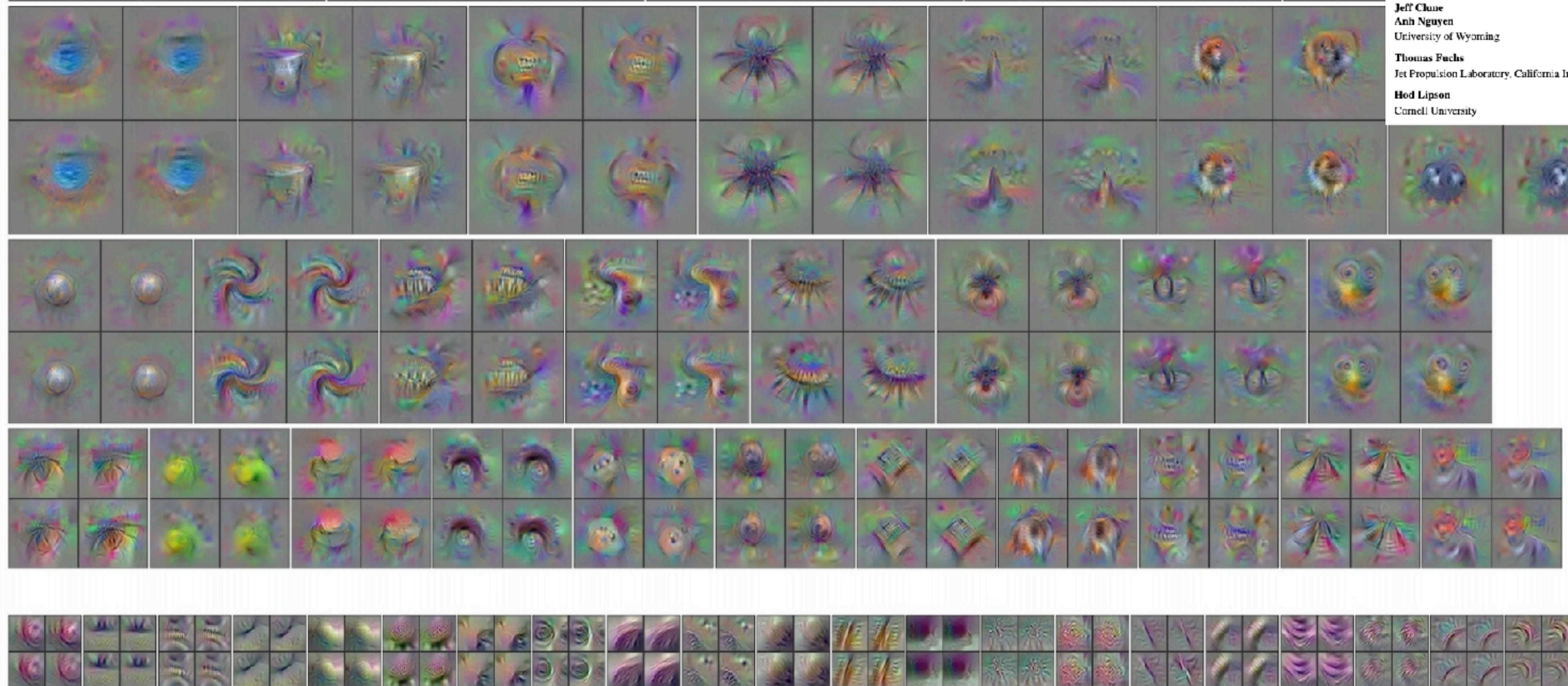
Thomas Fuchs
Jet Propulsion Laboratory, California Institute of Technology

FUCHS@CALTECH.EDU

Hod Lipson
Cornell University

HOD.LIPSON@CORNELL.EDU

Layer 5 Layer 4 Layer 3 Layer 2



Stronger Regularizers for GAs

Reconstructions of multiple feature types (facets) recognized by the same “grocery store” neuron



Corresponding example training set images recognized by the same neuron as in the “grocery store” class

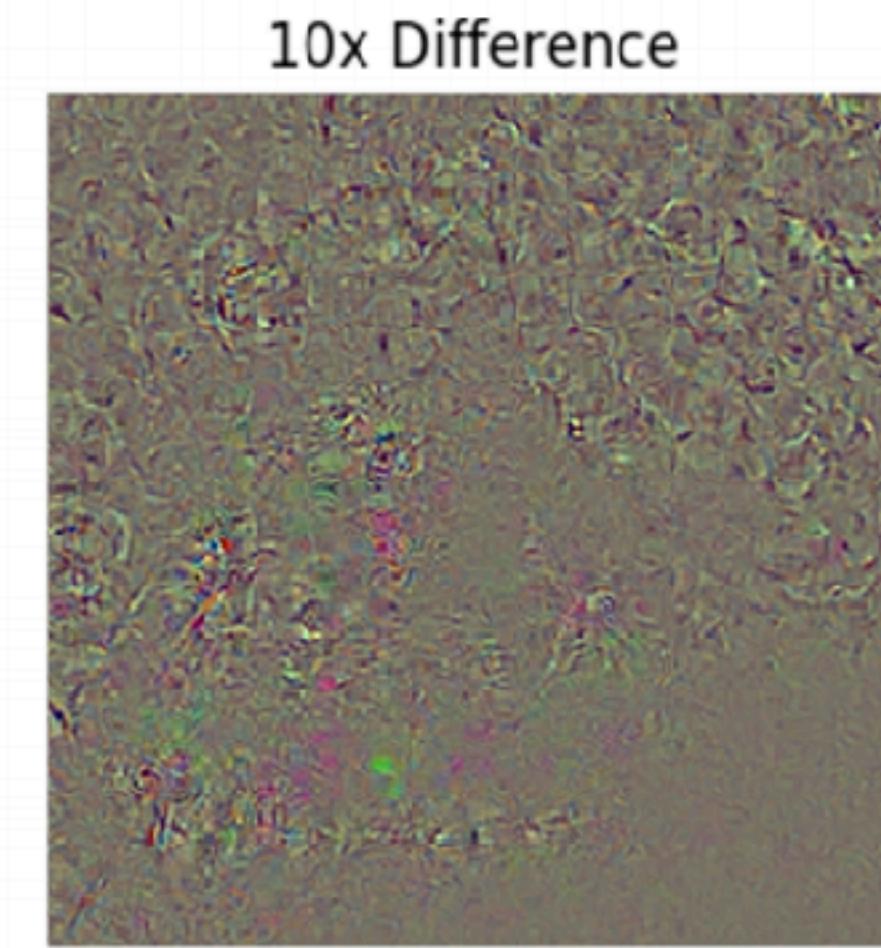
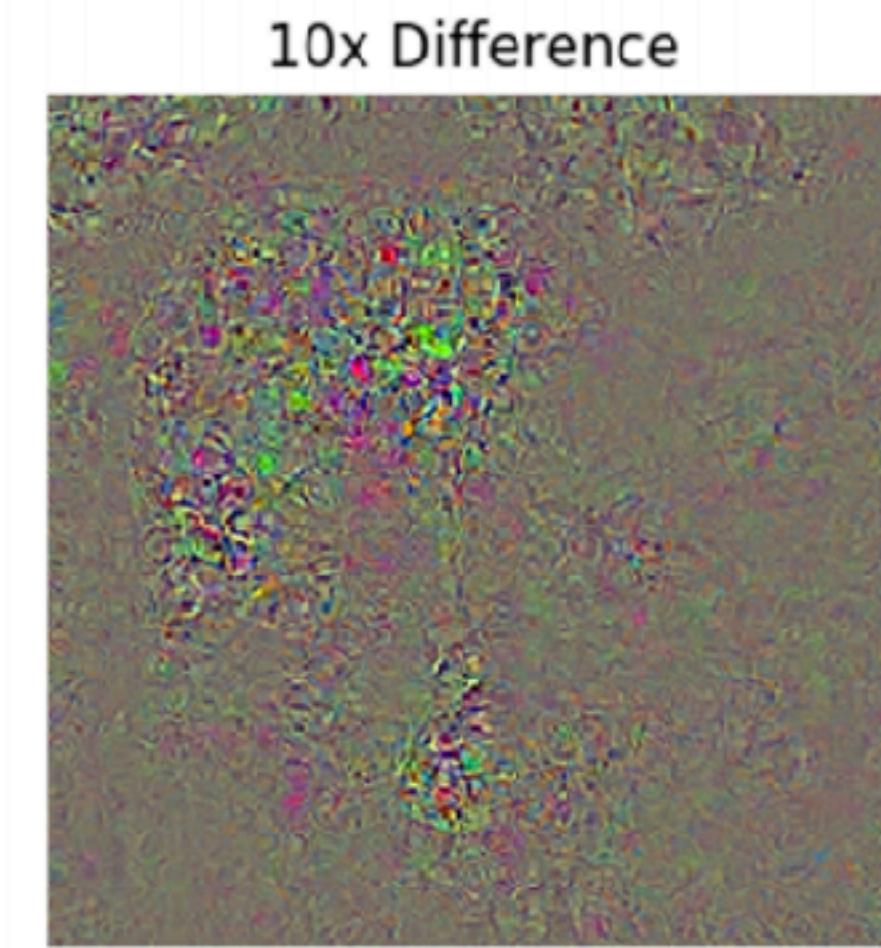
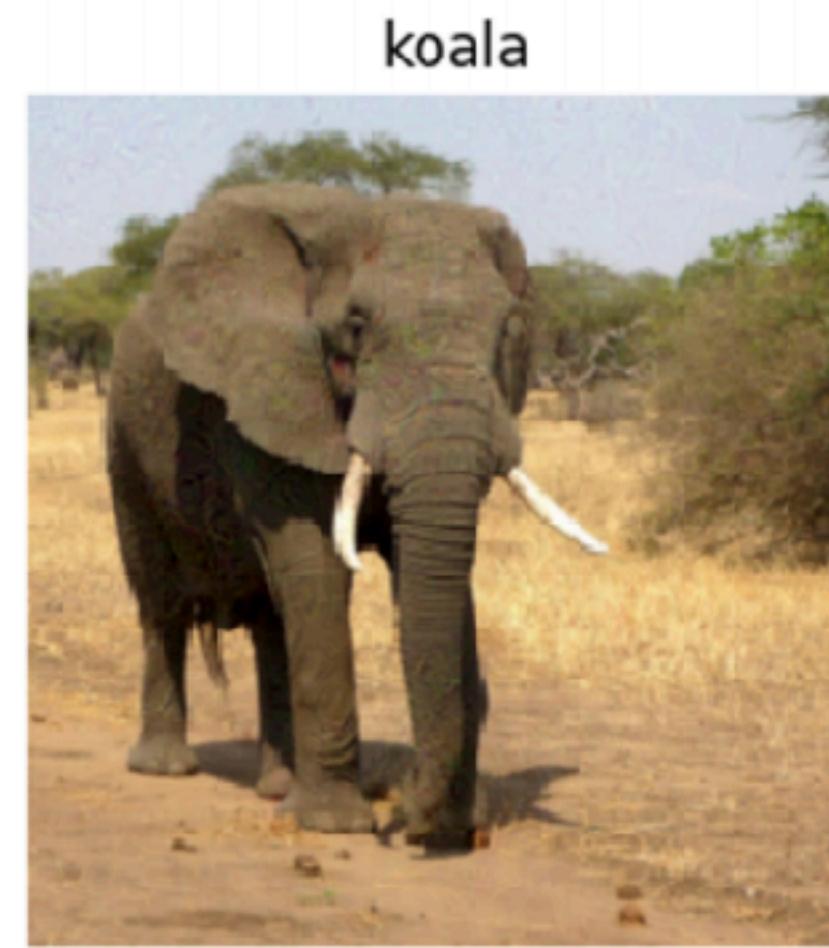


More results: Stronger Regularizers



Relation to Adversarial Training

Adversarial Training: Examples



Plan for Today

- **Visualize filter (kernels)**
- **Use CNNs as feature maps**
- **What are the different layers learning?**
- **Visualize activations at different layers**
- **Saliency versus Occlusion**
- **Guided backpropagation**
- **Gradient Ascent**
- **Feature Inversion**
- **DeepDreams**

Feature Inversion

$$x^* = \arg \min_I \{l(\Phi(I), \Phi_0) + \lambda \mathcal{R}(I)\}$$

$$l(\Phi(I), \Phi_0) = \|\Phi(I) - \Phi_0\|^2$$

$$\mathcal{R}_\beta(I) = \sum_{i,j} (I(i, j+1) - I(i, j))^2 + (I(i+1, j) - I(i, j))^2)^{\frac{\beta}{2}}$$

Feature Inversion

y

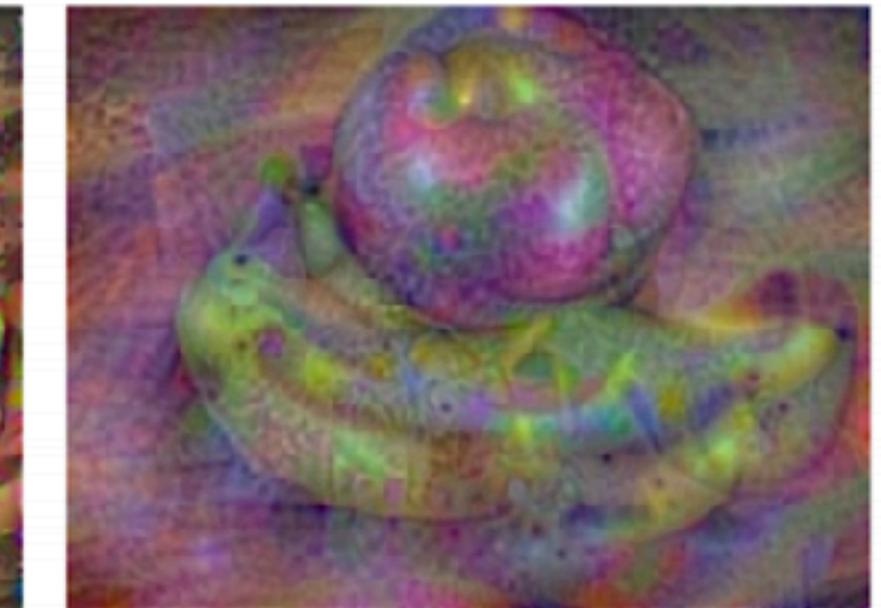
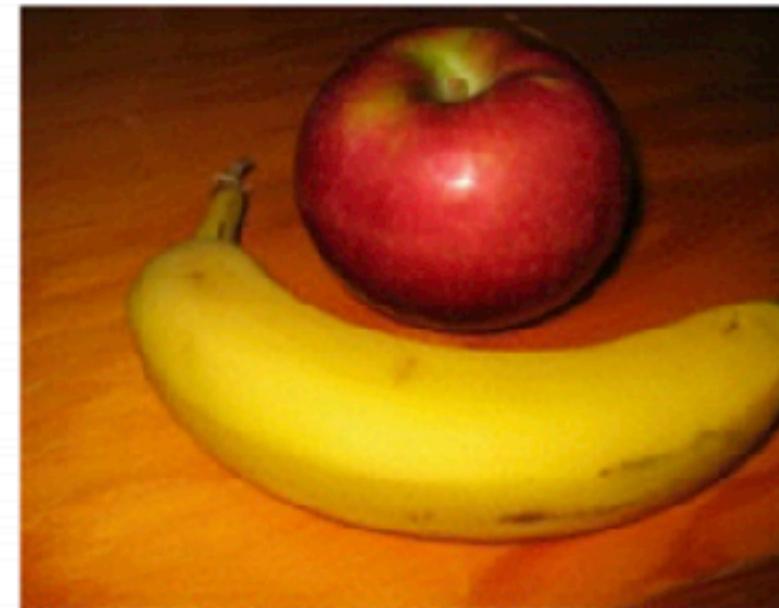
relu2_2

relu3_3

relu4_3

relu5_1

relu5_3

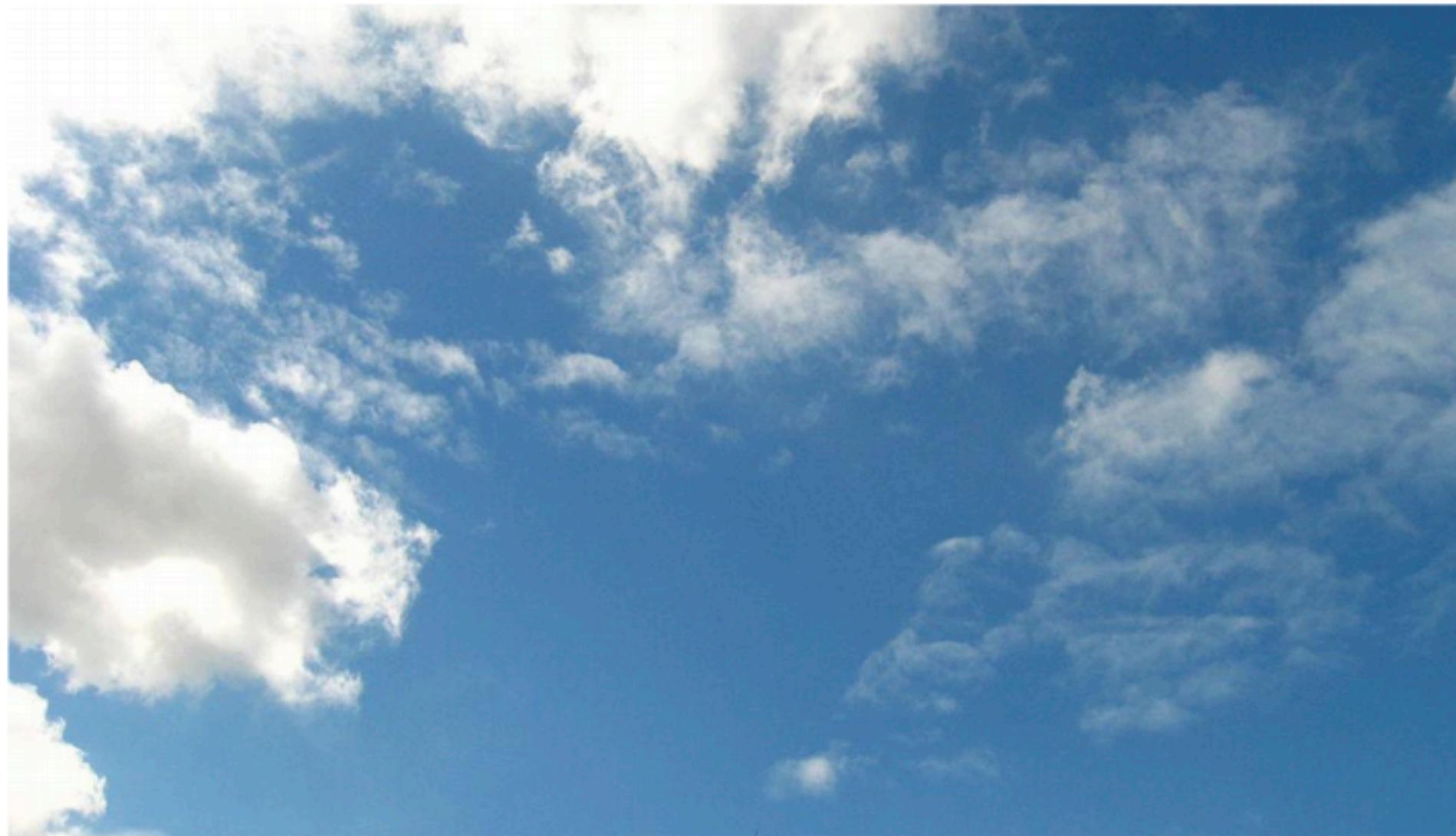


Plan for Today

- Visualize filter (kernels)
- Use CNNs as feature maps
- What are the different layers learning?
- Visualize activations at different layers
- Saliency versus Occlusion
- Guided backpropagation
- Gradient Ascent
- Feature Inversion
- DeepDreams

Deep Dream: Amplify Features

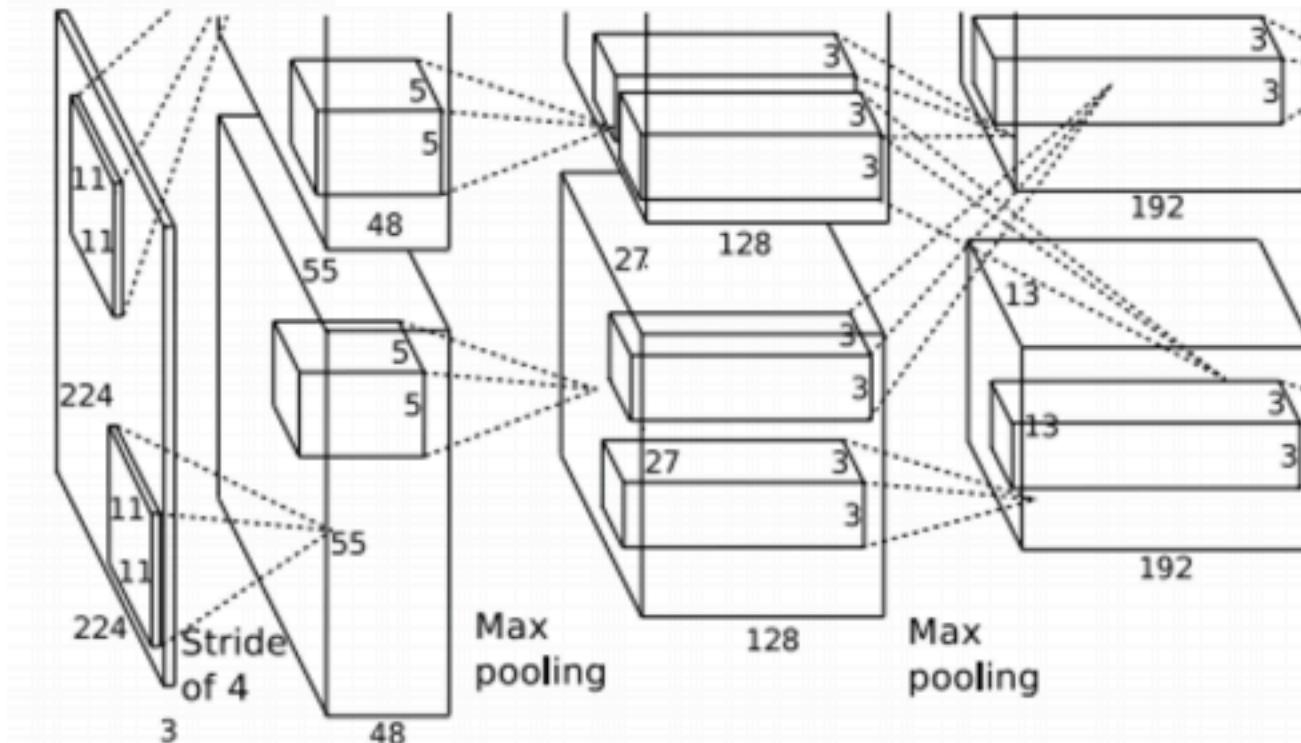
- C
- S
- C
- U



Inceptionism: Going Deeper into Neural Networks

Wednesday, June 17, 2015

Posted by Alexander Mordvintsev, Software Engineer, Christopher Olah, Software Engineering Intern and Mike Tyka, Software Engineer



$$I^* = \arg \max_I \sum_i f_i(I)^2$$

DeepDream



[Image](#) is licensed under CC-BY 3.0

DeepDream

