

Interpretability of Neural Networks

Kian Katanforoosh

Neural networks are deployed:

- In our phones: to recommend content,
- In banks: to manage investments,
- In hospitals: to help doctors diagnose disease symptoms,
- In insurance agencies: to evaluate risk and underwrite documents,
- In cars: to help avoid accidents.

How can one who will be held accountable for a decision trust a neural network's recommendation, and justify its use?

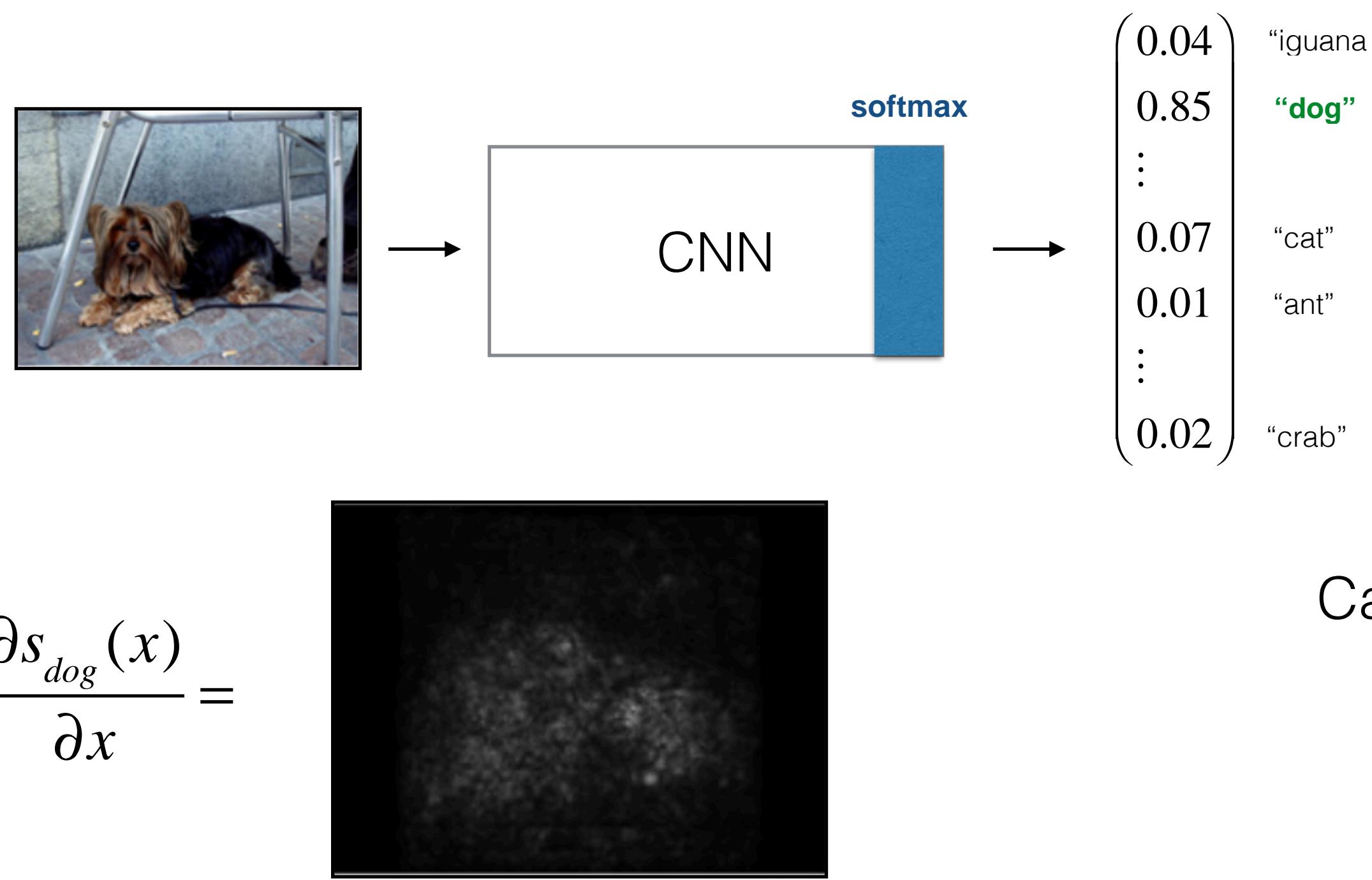
Today's outline

- I. Interpreting Neural Networks' outputs**
 - A. with saliency maps
 - B. with occlusion sensitivity
 - C. with class activation maps (Global Average Pooling)
- II. Visualizing Neural Networks from the inside**
 - A. with gradient ascent (class model visualization)
 - B. with dataset search
 - C. the deconvolution and its applications
- III. (Optional: Deep Dream: going deeper in NNs)**

I. A. Interpreting and visualizing Neural Networks with saliency maps

Context: You have built an animal classifier for a zoo. They are reluctant to use your model without human supervision, because they don't understand the decision process of the model.

Question: How can you alleviate their concerns?



indicates which pixels need to be changed the least to affect the class score the most.

Can be used for segmentation?



Yes

$$\text{softmax} \begin{pmatrix} s_{\text{iguana}} \\ s_{\text{dog}} \\ \vdots \\ s_{\text{cat}} \\ s_{\text{ant}} \\ \vdots \\ s_{\text{crab}} \end{pmatrix} = \begin{pmatrix} \frac{s_{\text{iguana}}}{\sum_{\text{animals}} s_{\text{animal}}} \\ \frac{s_{\text{dog}}}{\sum_{\text{animals}} s_{\text{animal}}} \\ \vdots \\ \frac{s_{\text{cat}}}{\sum_{\text{animals}} s_{\text{animal}}} \\ \vdots \\ \frac{s_{\text{crab}}}{\sum_{\text{animals}} s_{\text{animal}}} \end{pmatrix}$$

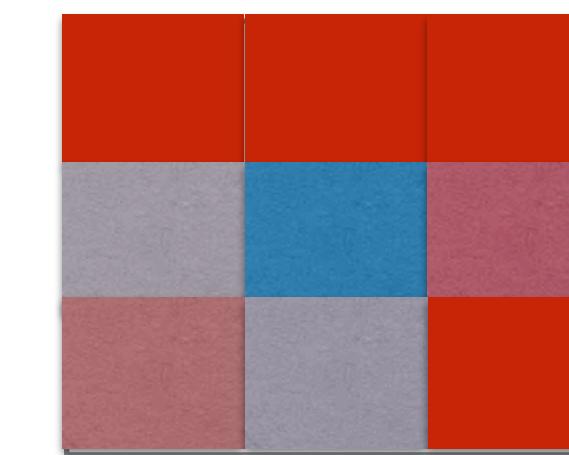
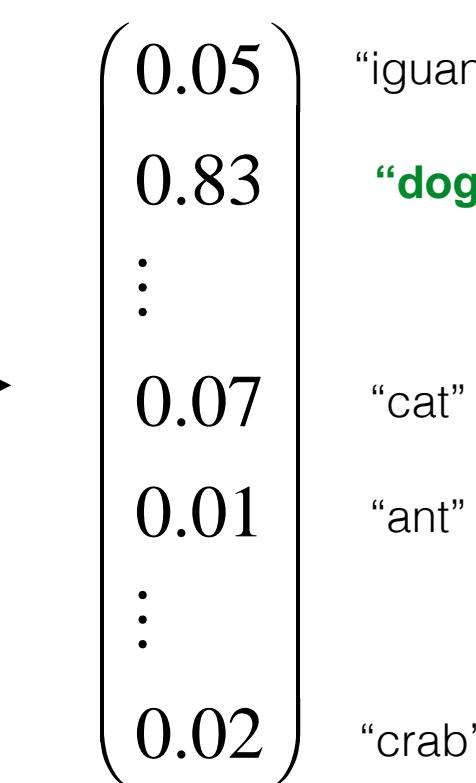
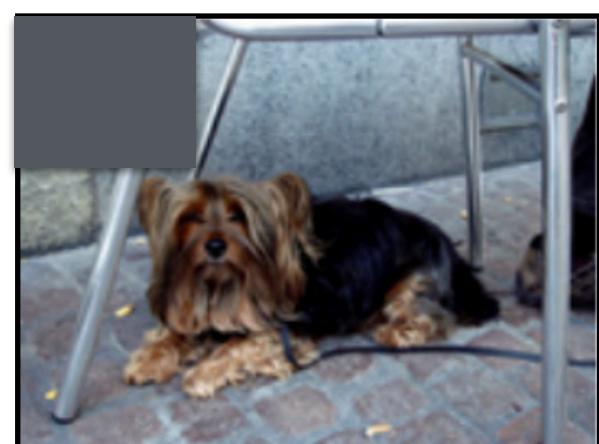
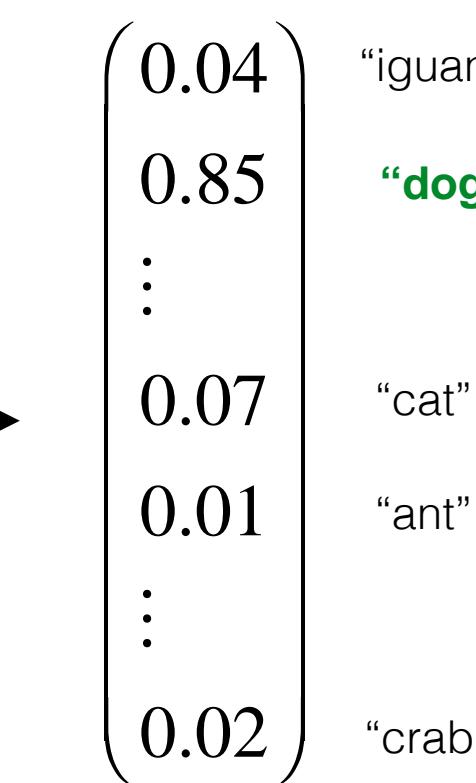
Saliency maps

Kian Katanforoosh

I. B. Interpreting and visualizing Neural Networks with occlusion sensitivity

Context: You have built an animal classifier for a zoo. They are a little reluctant to use your model without human supervision, because they don't understand the decision process of the model.

Question: What can you do to alleviate their concerns?



Probability map of the true class for different positions of the grey square

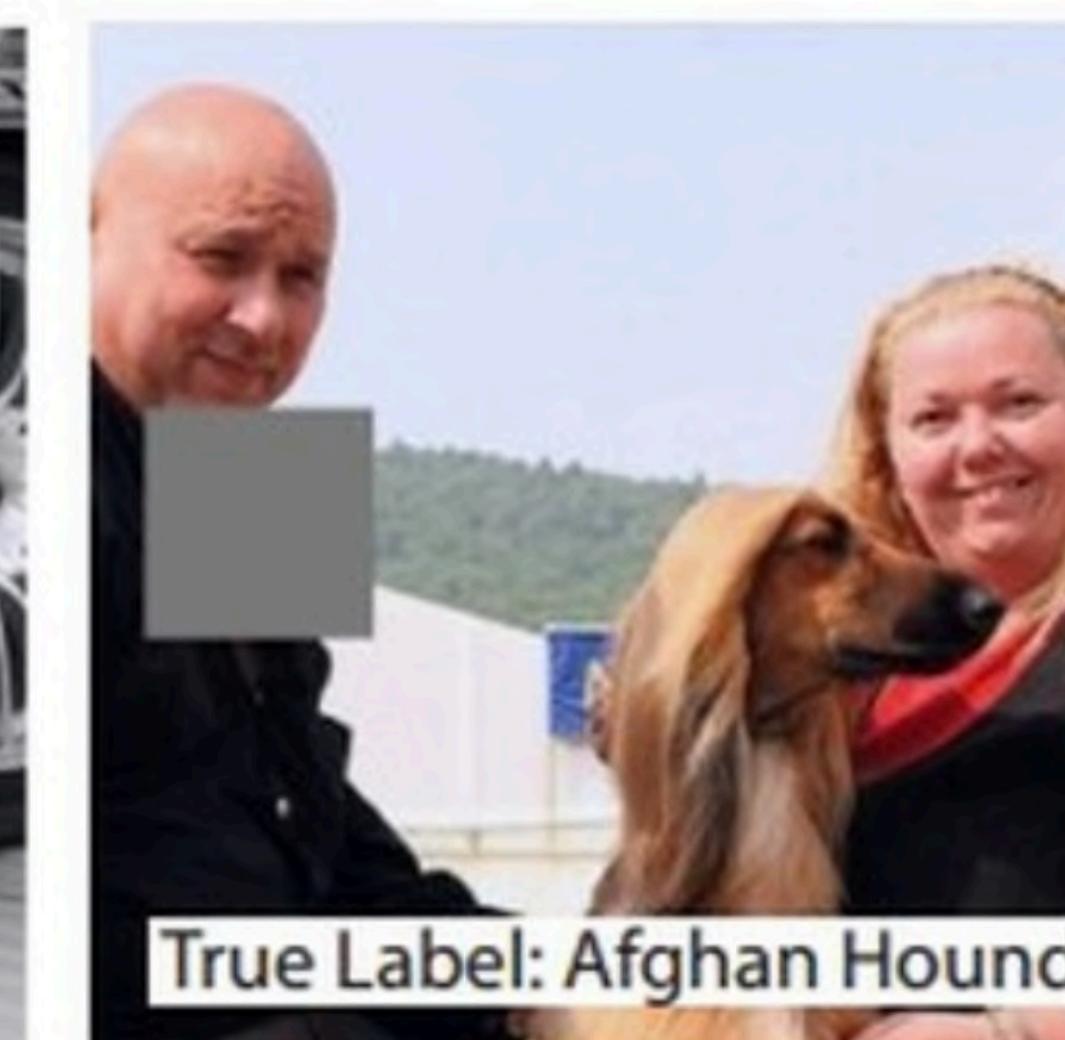


Indicates low confidence on the true class for the corresponding position of the grey square

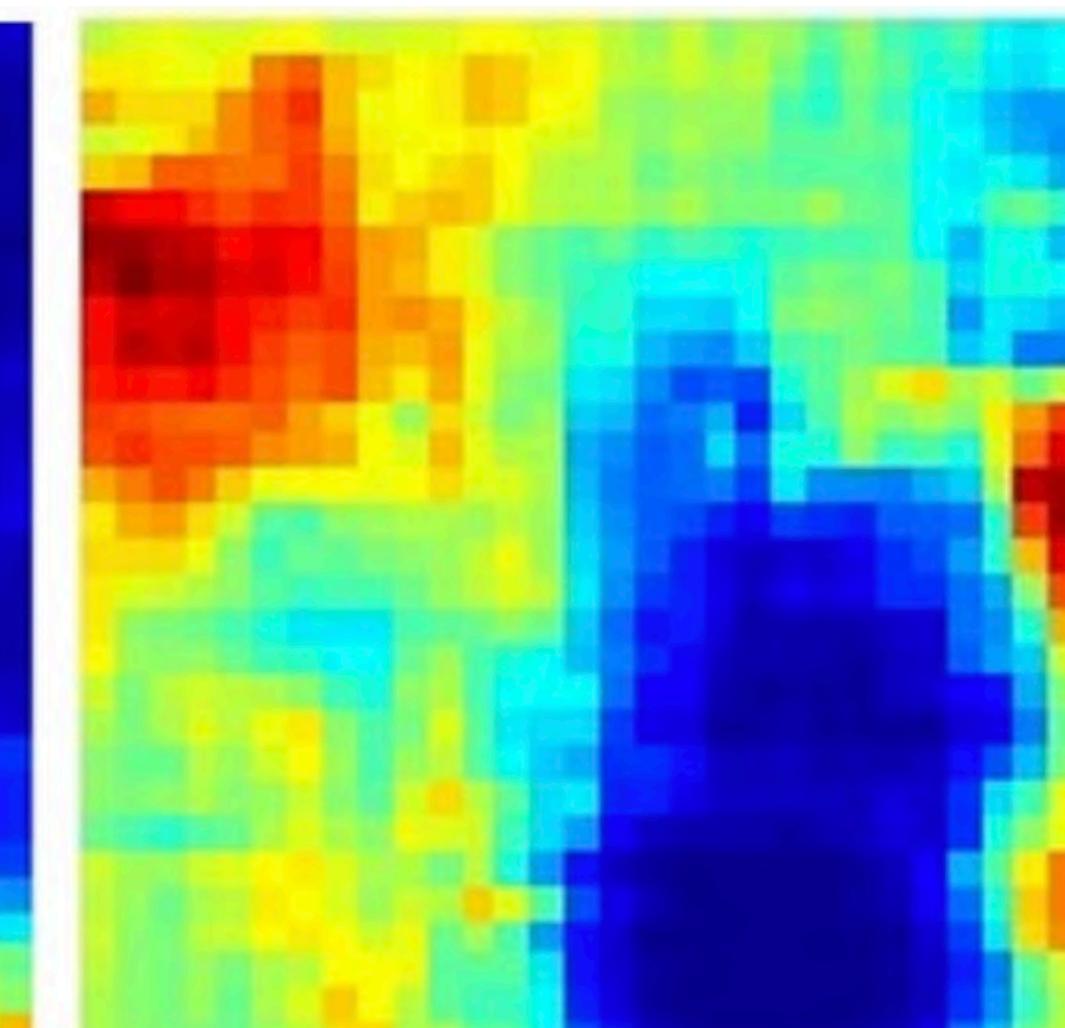
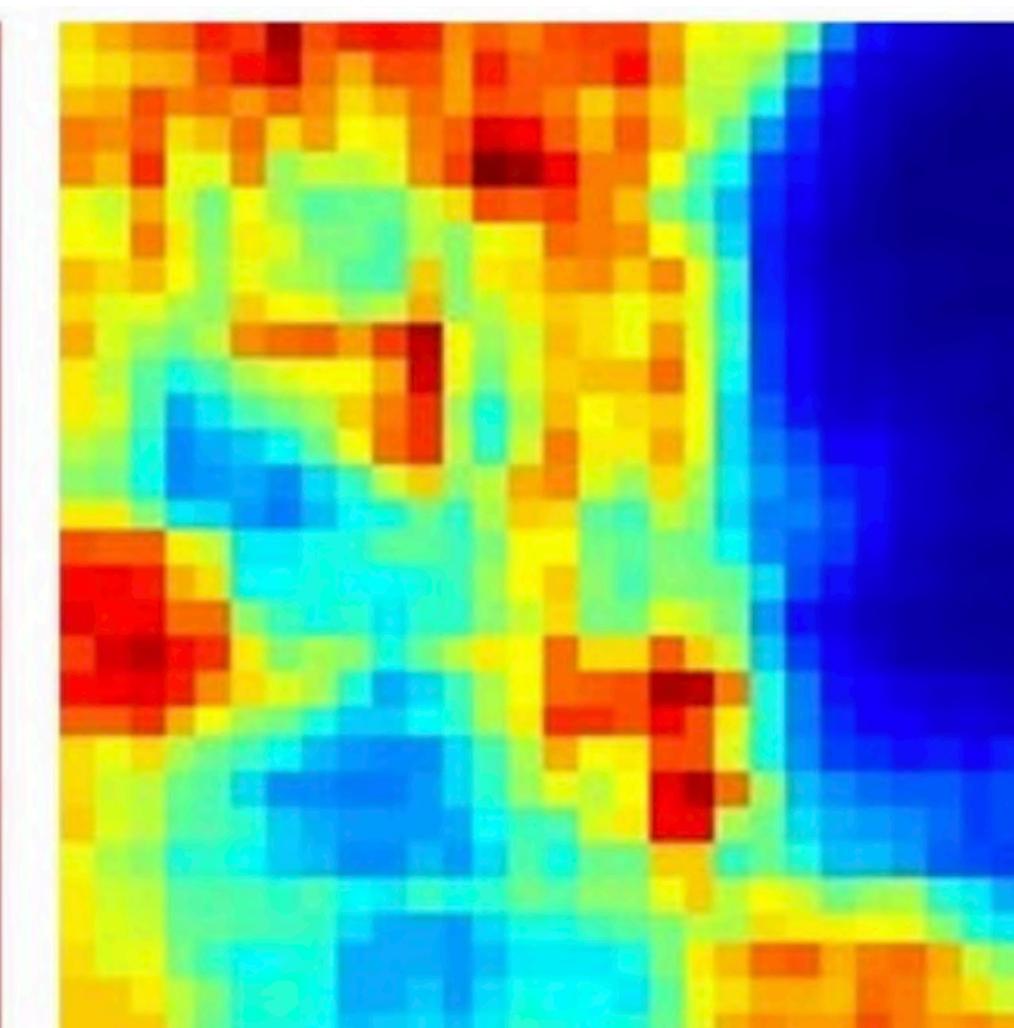
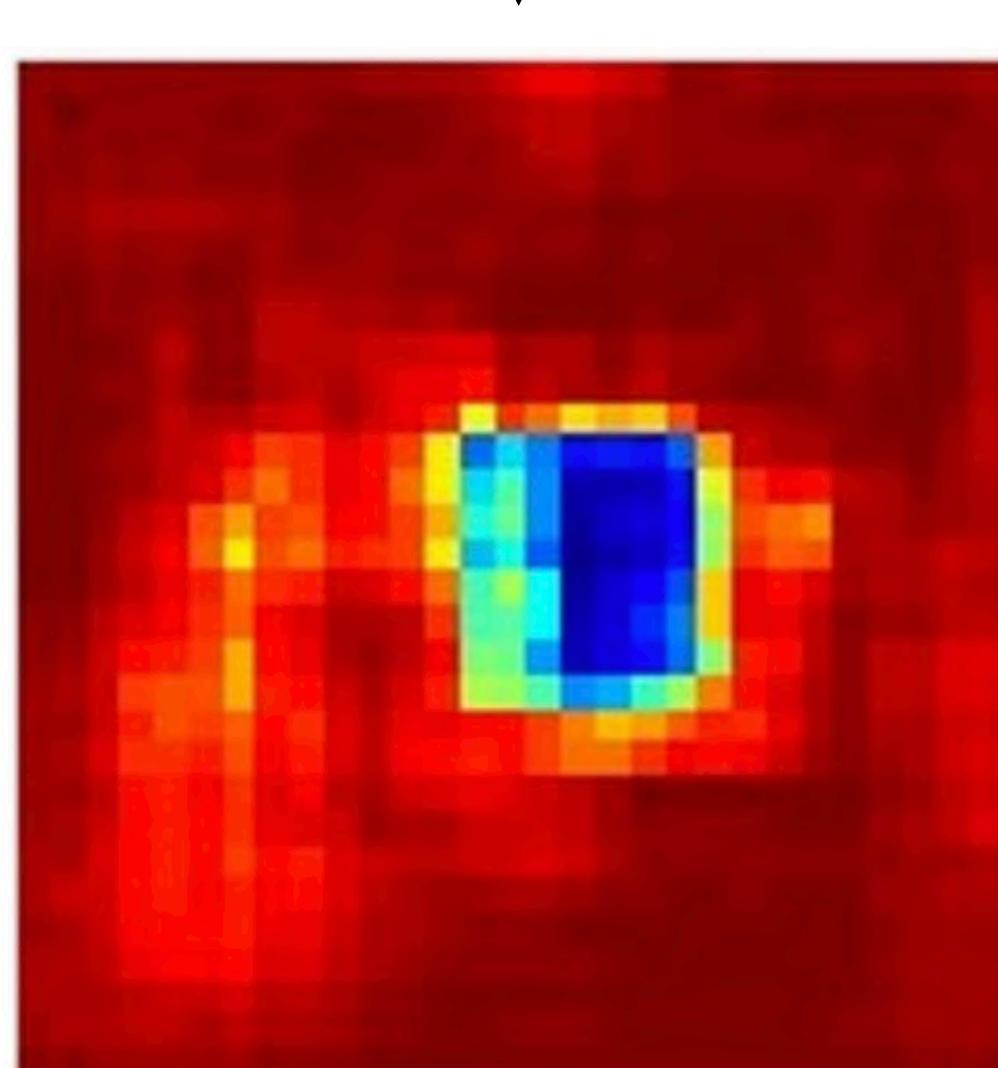


Indicates high confidence on the true class for the corresponding position of the grey square

I. B. Interpreting and visualizing Neural Networks with occlusion sensitivity



Probability map of the true class for different positions of the grey square



Indicates low confidence on the true class for the corresponding position of the grey square

Indicates high confidence on the true class for the corresponding position of the grey square

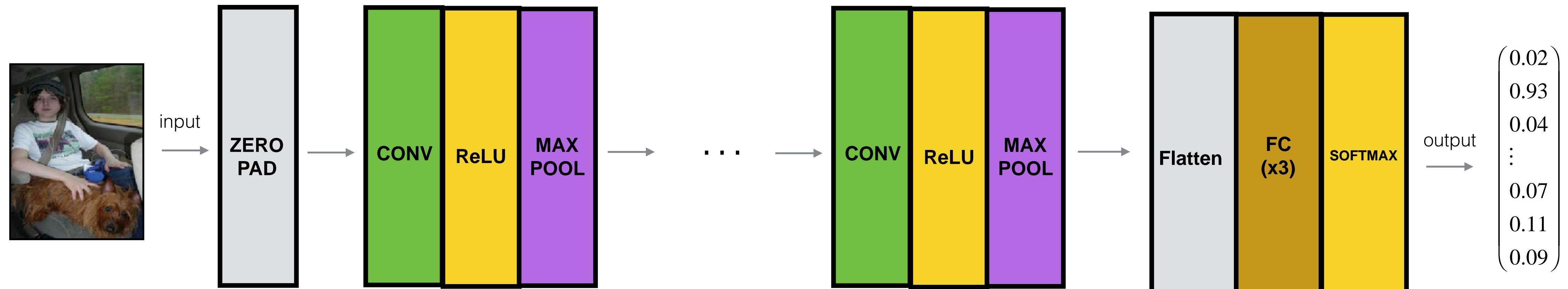
Occlusion sensitivity

Kian Katanforoosh

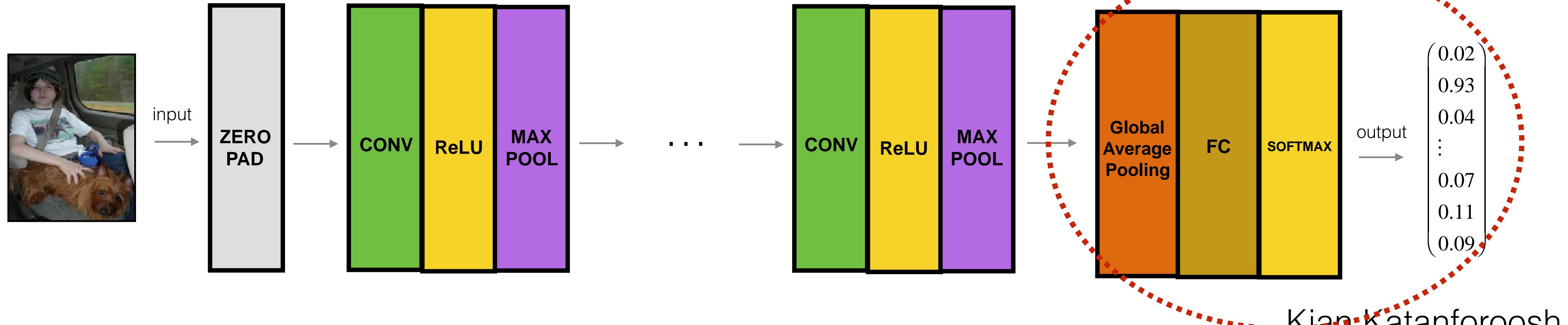
III. C. Interpreting NNs using class activation maps

Context: Along with the classification output, the zoo now wants real-time visualization of the model's decision process. You have one day. What do you do?

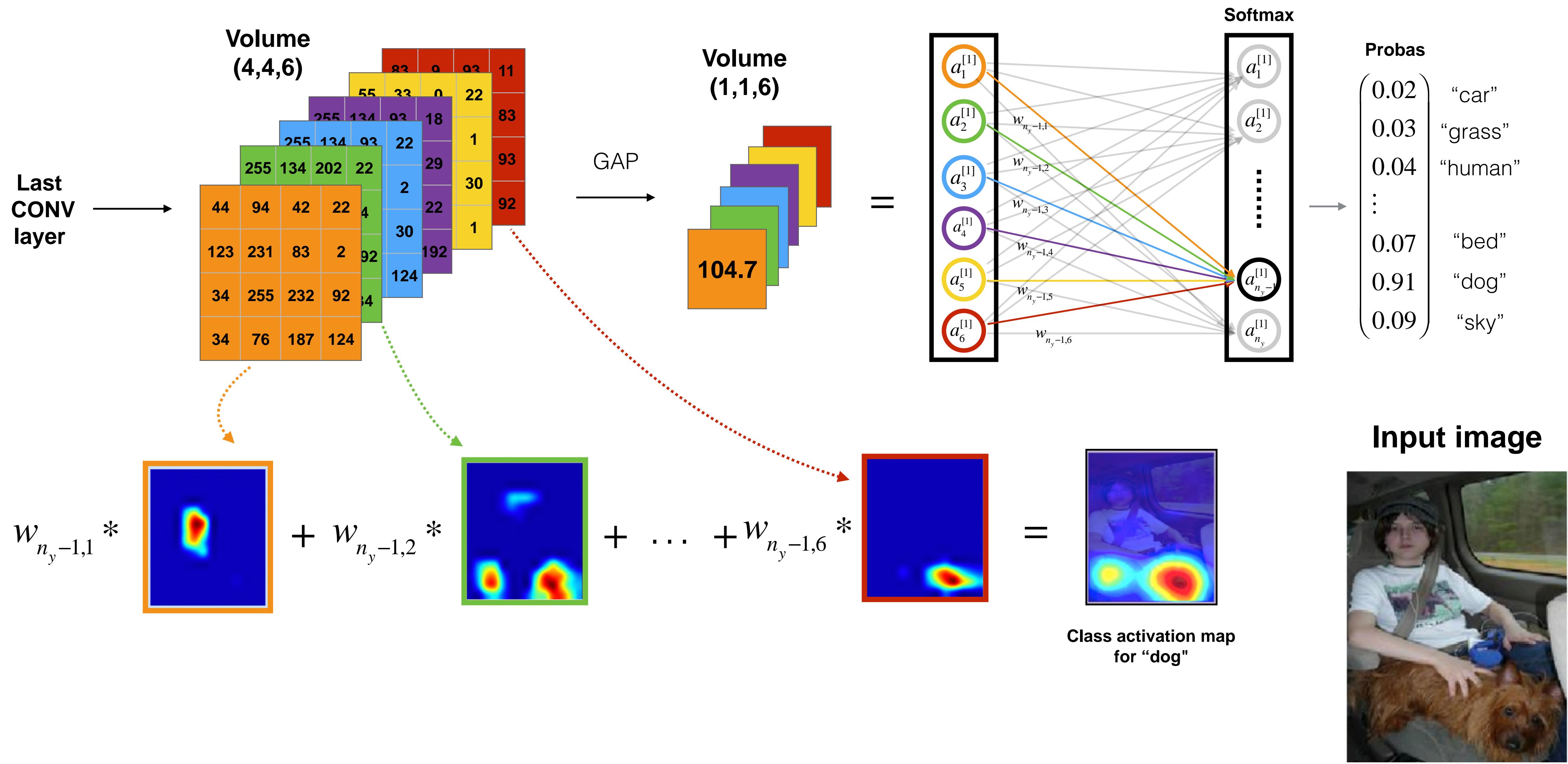
Using a classification network for localization:



Converted to:

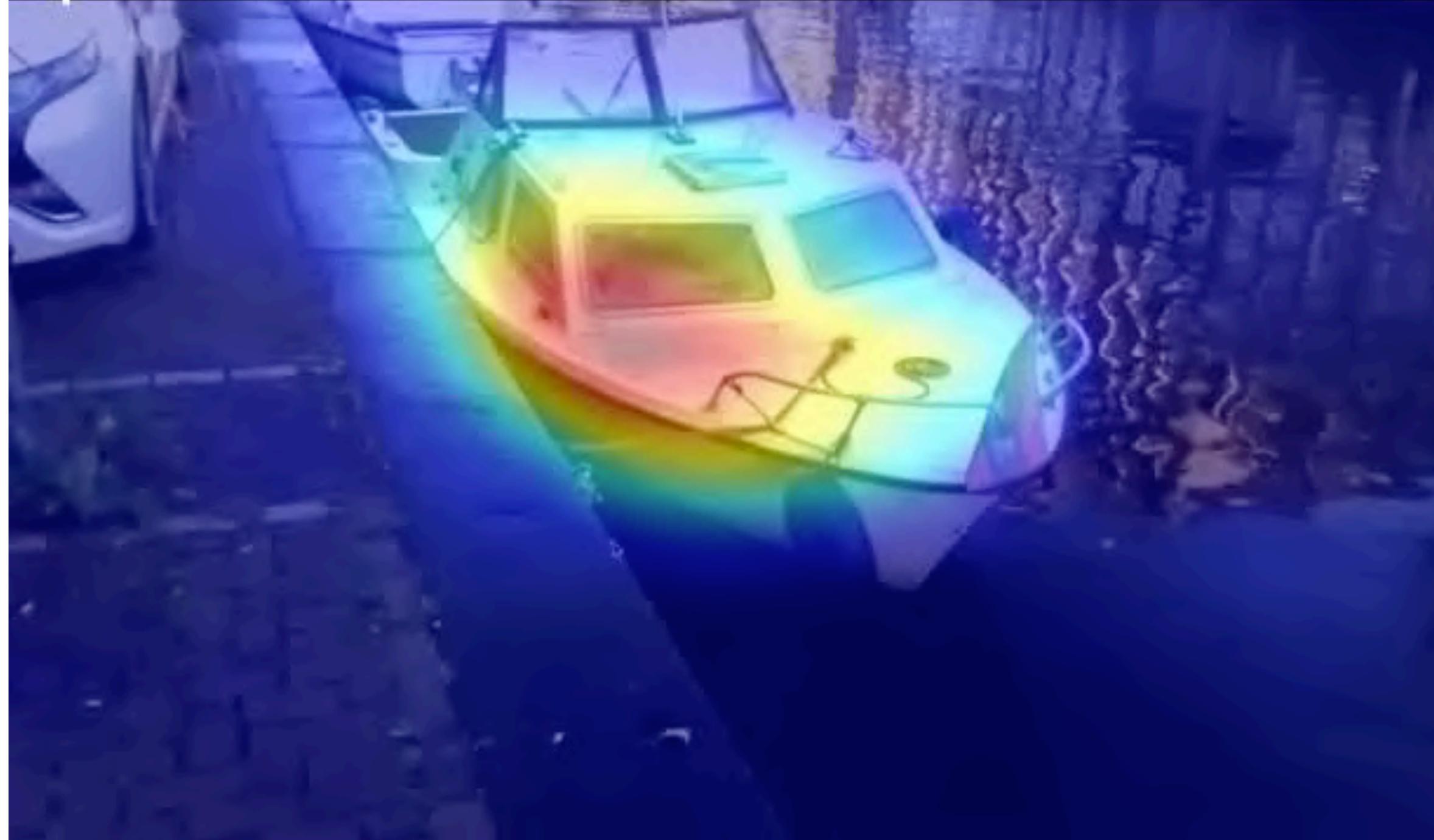


III. C. Interpreting NNs using class activation maps



III. C. Interpreting NNs using class activation maps

speedboat



[Bolei Zhou et al. (2016): Learning Deep Features for Discriminative Localization]

Source video: Kyle McDonald

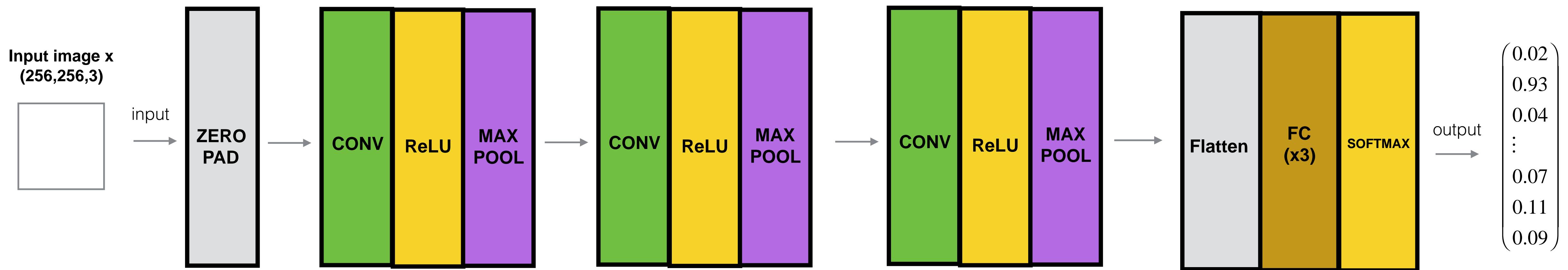
Today's outline

- I. Interpreting Neural Networks' outputs
 - A. with saliency maps
 - B. with occlusion sensitivity
 - C. with class activation maps (Global Average Pooling)
- II. Visualizing Neural Networks from the inside**
 - A. with gradient ascent (class model visualization)
 - B. with dataset search
 - C. The deconvolution and its applications
- III. (Optional: Deep Dream: going deeper in NNs)

II. A. Visualizing NNs from the inside using gradient ascent (class model visualization)

Context: The zoo trusts that your model correctly locates animals. They get scared and they ask you whether the model understands what a dog is.

Given this trained ConvNet, generate an image which is representative of the class “dog” according to the ConvNet



Keep the weights fixed and use gradient ascent on the input image to maximize this loss :

$$L = s_{\text{dog}}(x) - \lambda \|x\|_2^2$$

“ x should look natural”

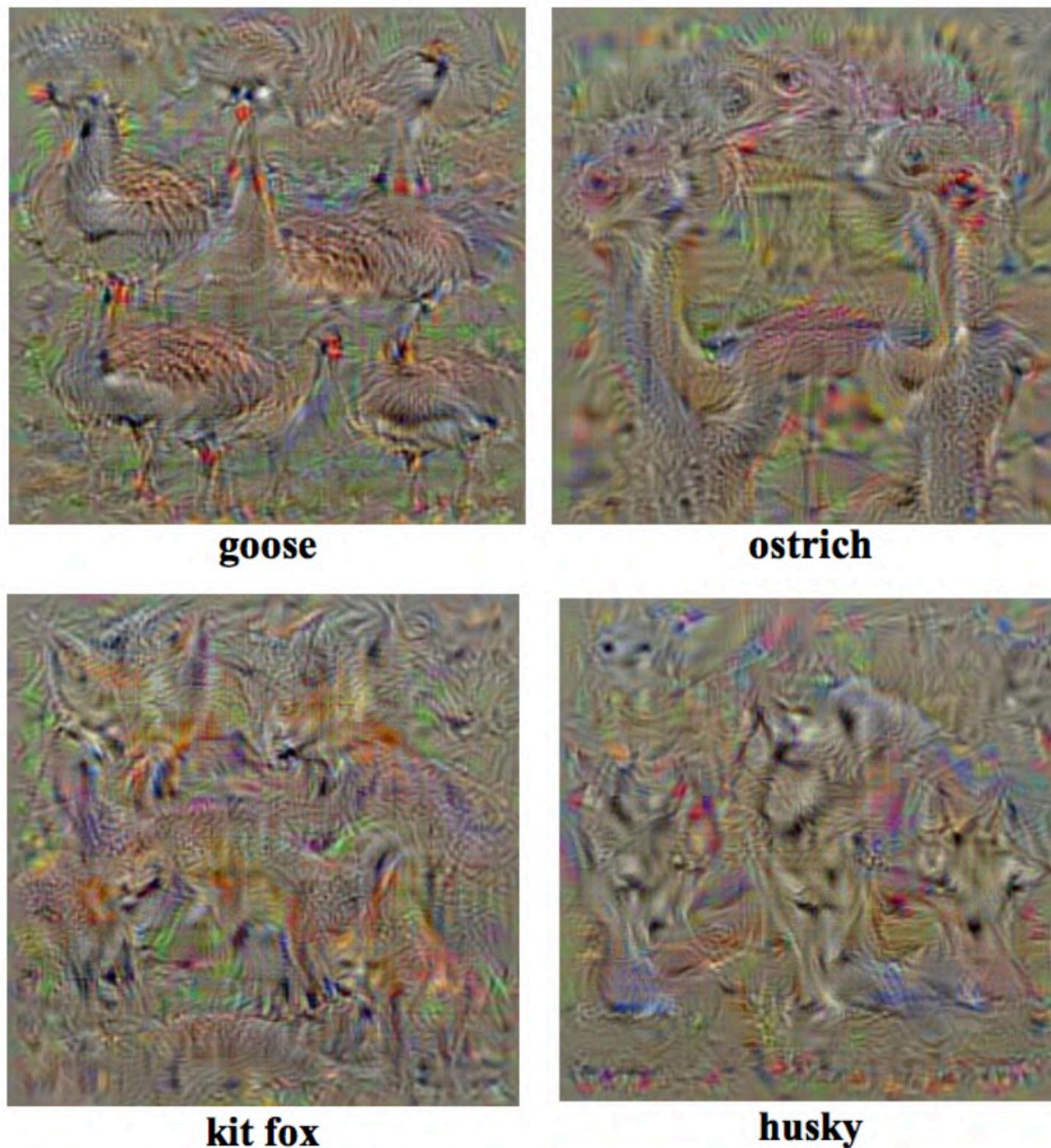
Gradient ascent:

$$x = x + \alpha \frac{\partial L}{\partial x}$$

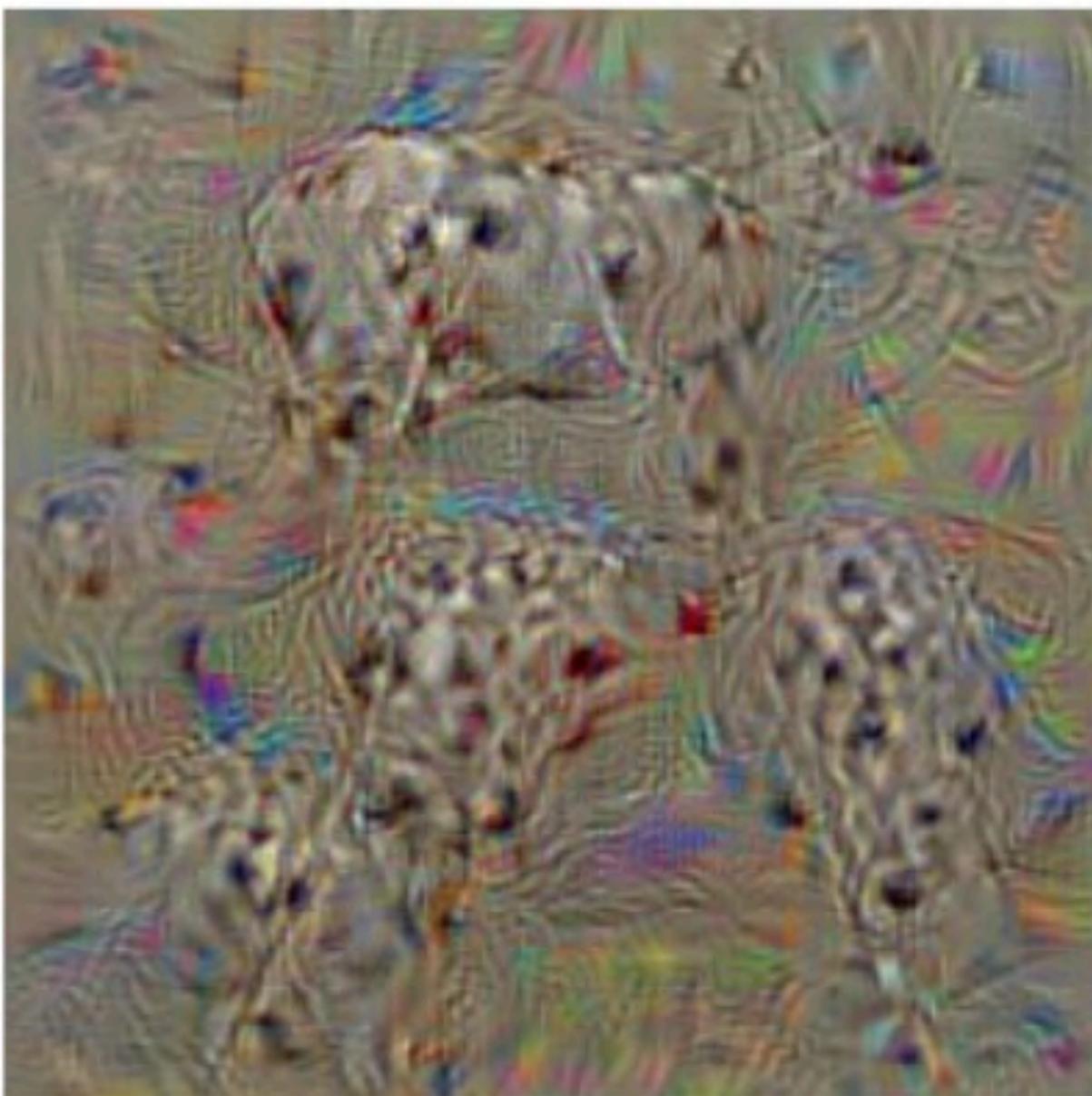
Repeat this process:

1. Forward propagate image x
2. Compute the objective L
3. Backpropagate to get dL/dx
4. Update x 's pixels with gradient ascent

II. A. Visualizing NNs from the inside using gradient ascent (class model visualization)



We can do this for all classes:

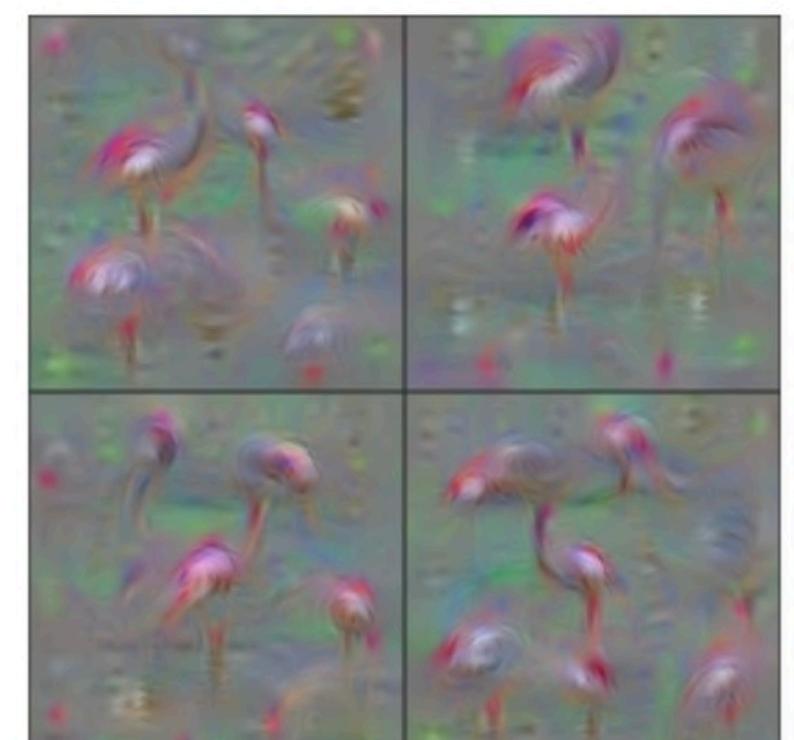


dalmatian

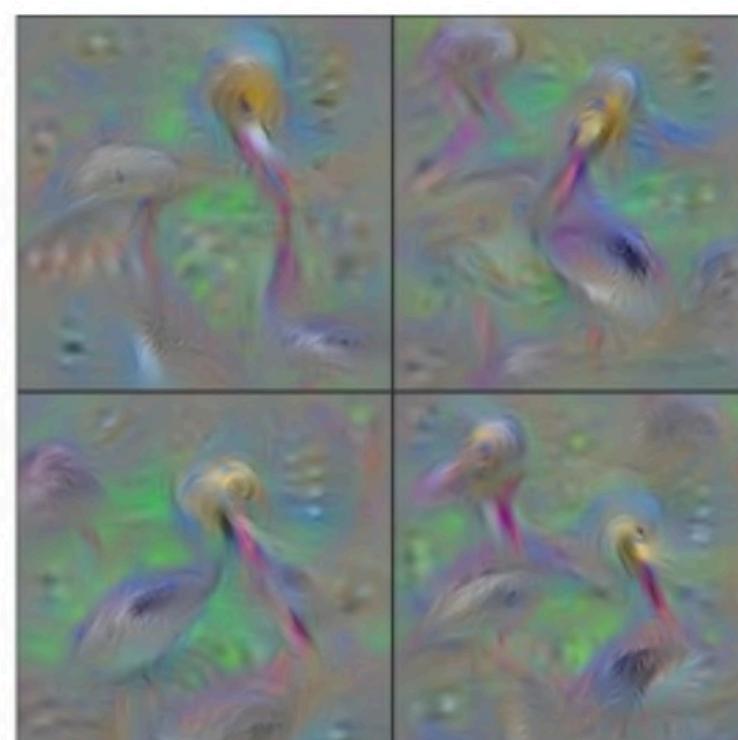
L2
Regularization

Looks better with additional regularization methods.

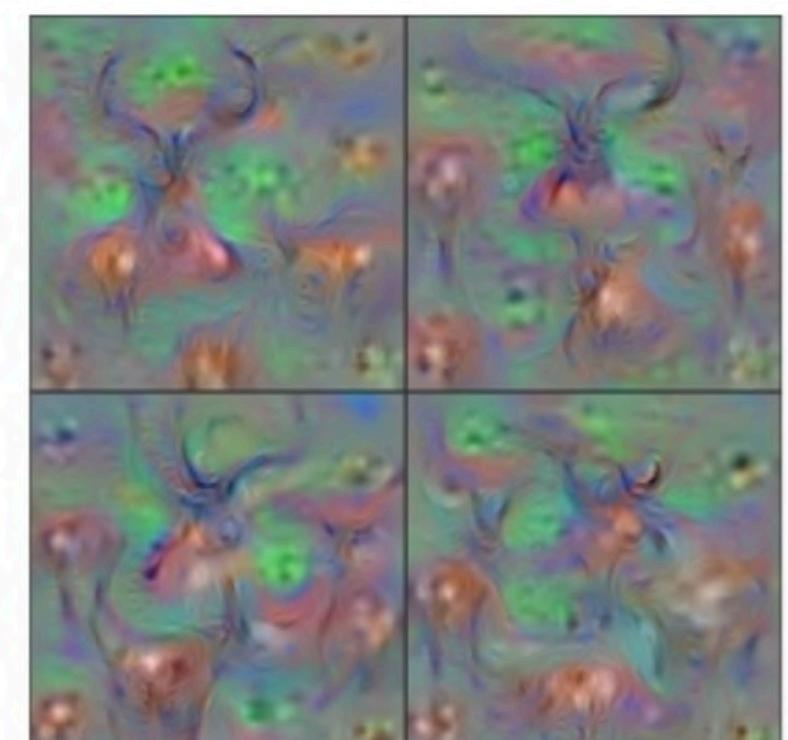
Class model visualization



Flamingo



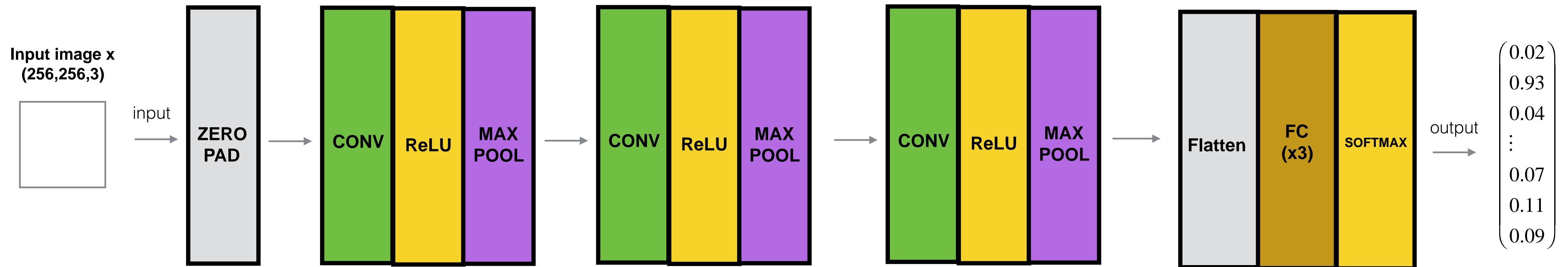
Pelican



Hartebeest
Kian Katanforoosh

II. A. Visualizing NNs from the inside using gradient ascent (class model visualization)

This method can be applied to any activation in the network in order to interpret what a neuron is detecting



On the class score:

$$L = S_{dog}(x) - R(x)$$

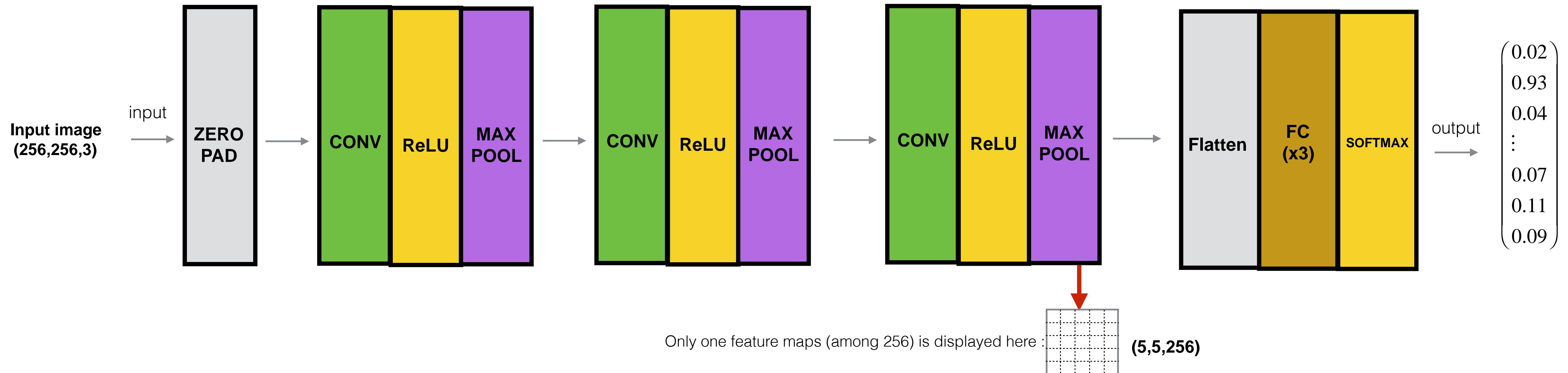
change to

$$L = a_j^{[l]}(x) - R(x)$$

On any activation:

II. B. Visualizing NNs from the inside using dataset search

Context: The zoo loved the technique, and asks if there are other alternatives.



Given a filter, what examples in the dataset lead to a strongly activated feature map?

Top 5 images



It seems that the filter has learned to detect shirts

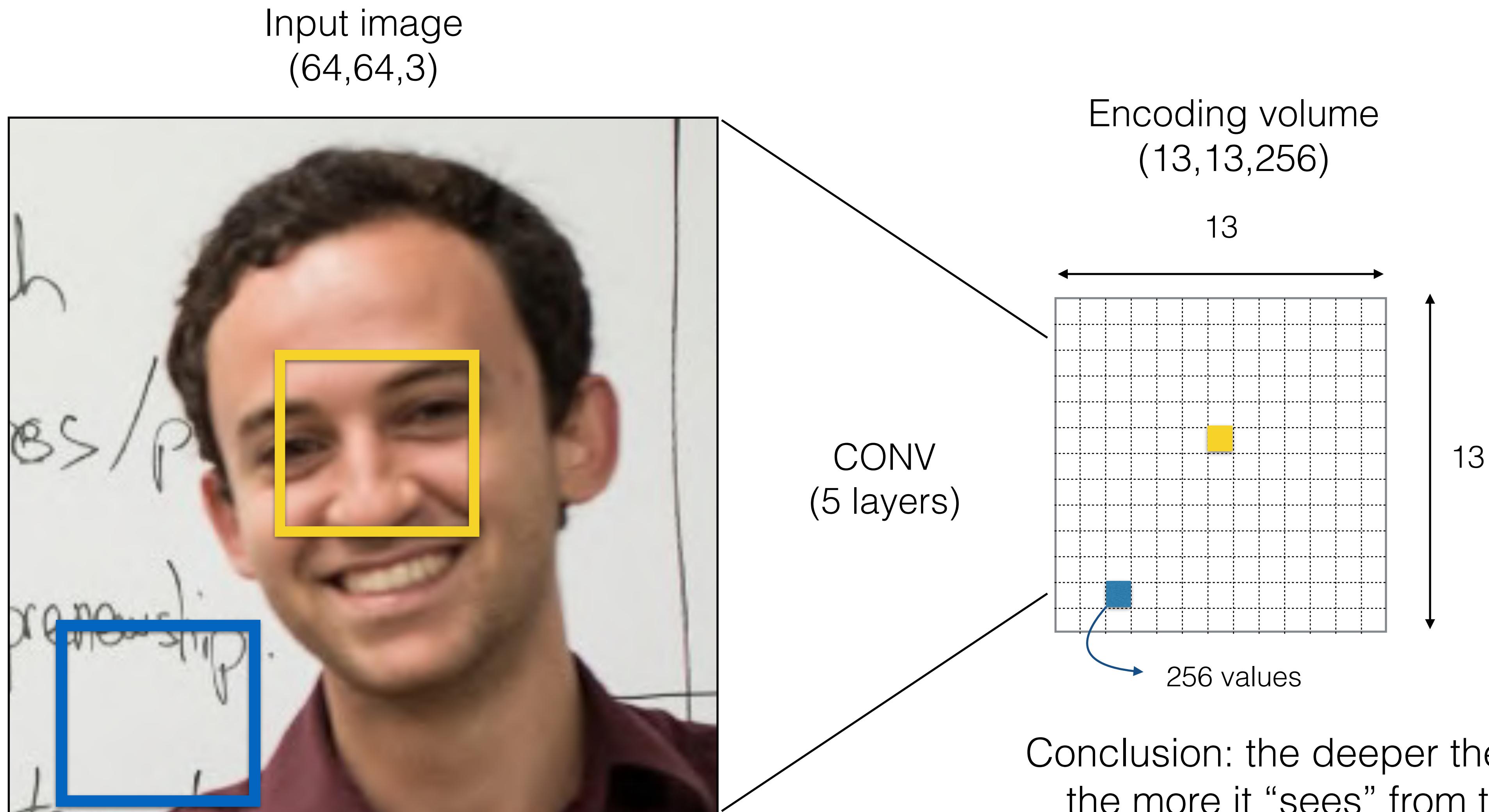
Top 5 images



It seems that the filter has learned to detect edges

II. B. Visualizing NNs from the inside using dataset search

How did we crop the dataset images on the previous slide?

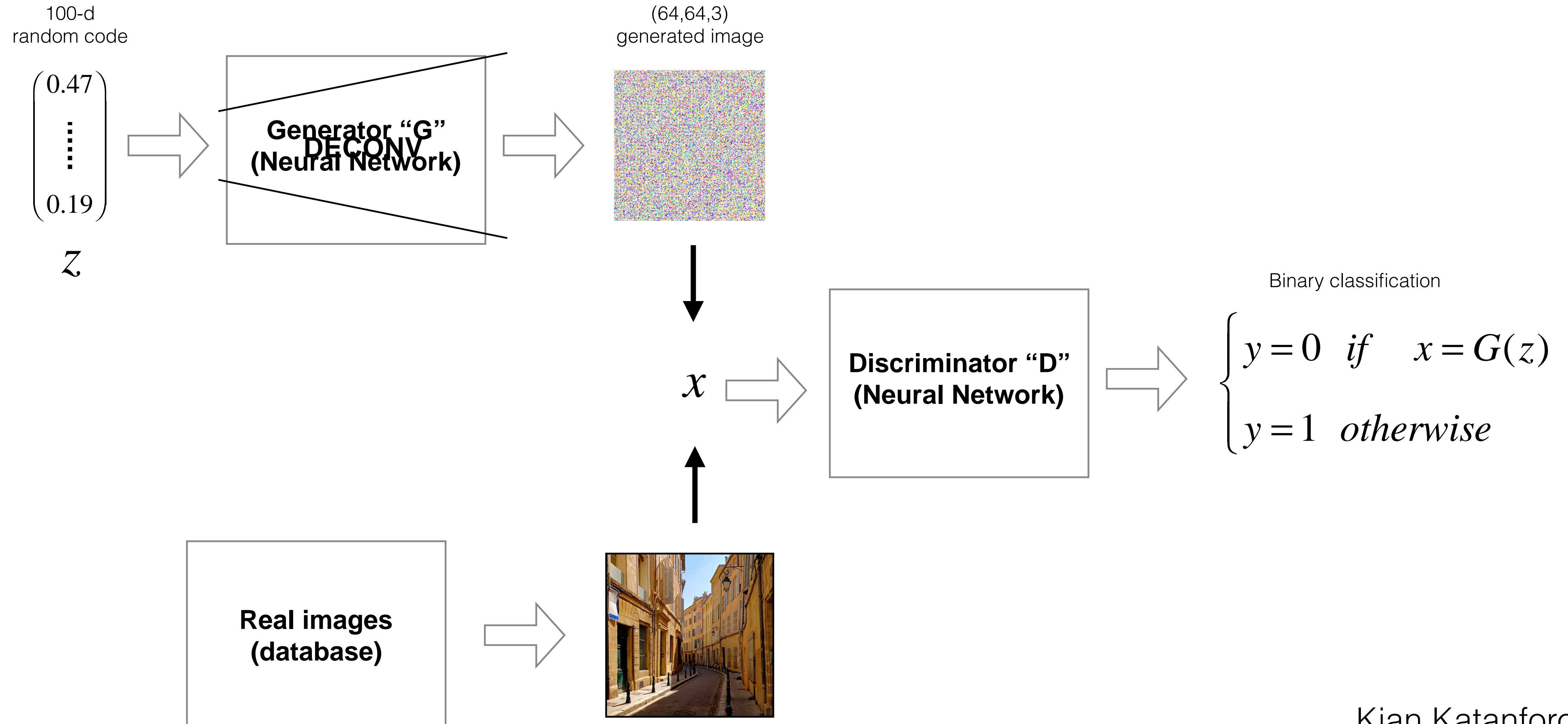


Today's outline

- I. Interpreting Neural Networks' outputs
 - A. with saliency maps
 - B. with occlusion sensitivity
 - C. with class activation maps (Global Average Pooling)
- II. Visualizing Neural Networks from the inside
 - A. with gradient ascent (class model visualization)
 - B. with dataset search
 - C. The deconvolution and its applications**
- III. (Optional: Deep Dream: going deeper in NNs)

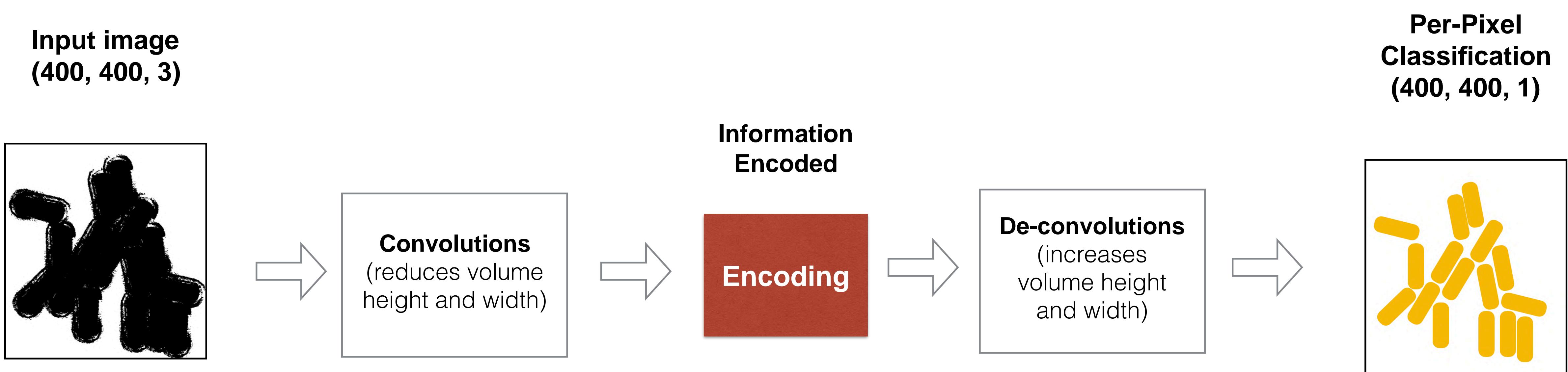
II. C. The deconvolution and its applications

Motivation behind deconvolution/upsampling layers



II. C. The deconvolution and its applications

Motivation behind deconvolution/upsampling layers

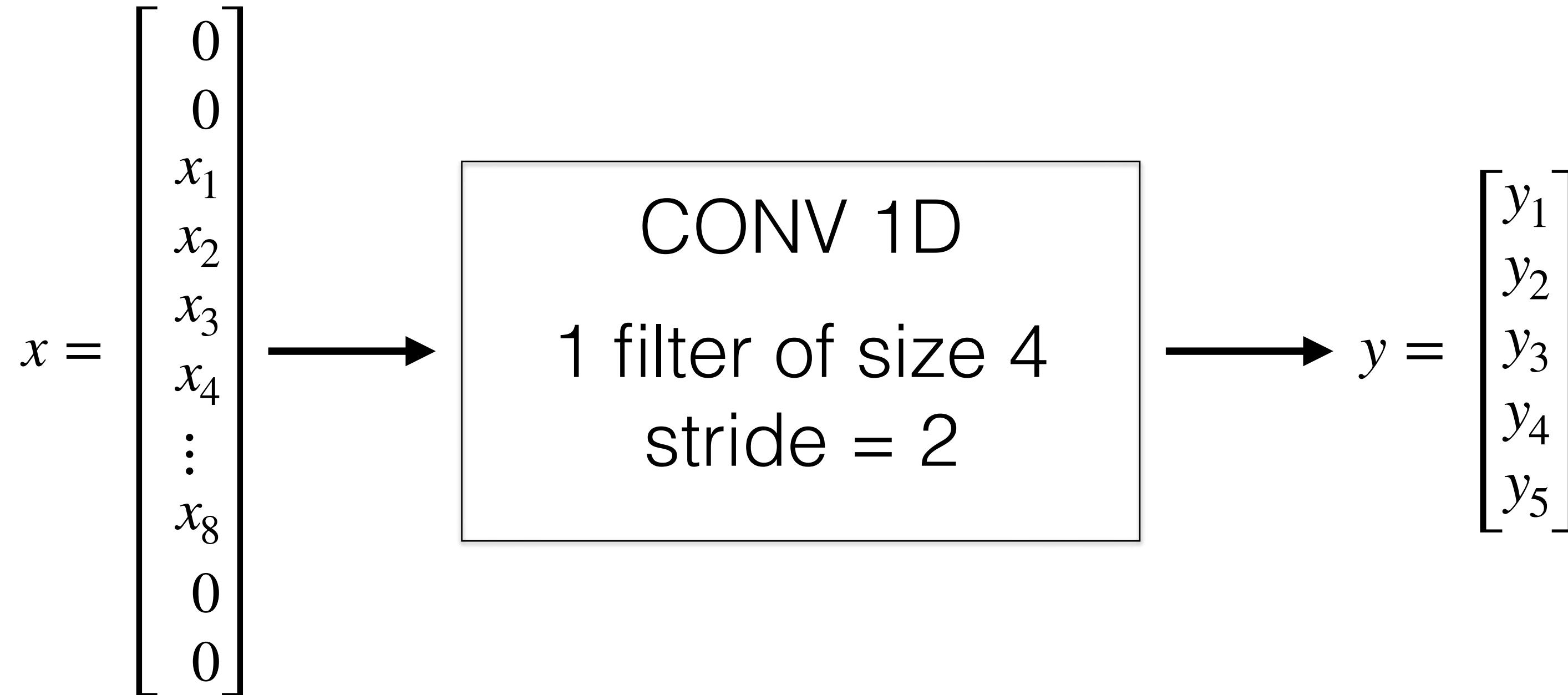




Get your pencils ready, we're about to study the deconvolution.

Consider the following CONV1D:

$$n_y = \left\lfloor \frac{n_x - f + 2p}{s} \right\rfloor + 1 = \left\lfloor \frac{8 - 4 + 2 \times 2}{2} \right\rfloor + 1 = 5$$



The system of equations is:

Let's define our filter as:

$$f = (w_1, w_2, w_3, w_4)$$

$$y_1 = w_1 \cdot 0 + w_2 \cdot 0 + w_3 \cdot x_1 + w_4 \cdot x_2$$

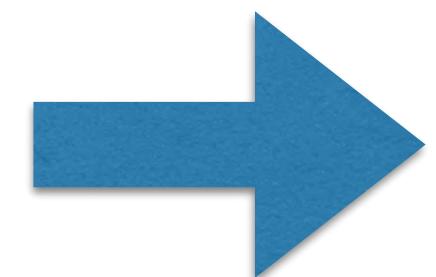
$$y_2 = w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + w_4 \cdot x_4$$

\vdots

$$y_5 = w_1 \cdot x_7 + w_2 \cdot x_8 + w_3 \cdot 0 + w_4 \cdot 0$$

Let's rewrite the system of equations:

$$\begin{aligned}y_1 &= w_1 \cdot 0 + w_2 \cdot 0 + w_3 \cdot x_1 + w_4 \cdot x_2 \\y_2 &= w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + w_4 \cdot x_4 \\&\vdots \\y_5 &= w_1 \cdot x_7 + w_2 \cdot x_8 + w_3 \cdot 0 + w_4 \cdot 0\end{aligned}$$



1D CONV can be rewritten as a matrix vector multiplication!

(5,1)

(5,12)

$$y = Wx$$

(12,1)

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix} = \begin{bmatrix} w_1 & w_2 & w_3 & w_4 & 0 & \dots & \dots & \dots & 0 \\ 0 & 0 & w_1 & w_2 & w_3 & w_4 & 0 & \dots & \dots \\ 0 & 0 & 0 & 0 & w_1 & w_2 & w_3 & w_4 & 0 \\ 0 & \dots & 0 & 0 & 0 & 0 & w_1 & w_2 & w_3 & w_4 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_8 \\ 0 \\ 0 \end{bmatrix}$$

1D CONV

$$y = Wx$$

Diagram illustrating 1D convolution:

Input shape: $(5, 1)$

Kernel shape: $(5, 12)$

Output shape: $(12, 1)$

The diagram shows a curved arrow from the input shape $(5, 1)$ to the output shape $(12, 1)$, indicating the receptive field of each output unit.

1D DECONV

If W is invertible, then $\exists H = W^{-1}$ such that $x = Hy$

Diagram illustrating 1D deconvolution (Transpose convolution):

Input shape: $(12, 1)$

Kernel shape: $(12, 5)$

Output shape: $(5, 1)$

The diagram shows a curved arrow from the input shape $(12, 1)$ to the output shape $(5, 1)$, indicating the receptive field of each output unit.

In practice, we would even assume that W is orthogonal, i.e. $W^{-1} = W^T$

For example

$$(w_1, w_2, w_3, w_4) = (-1, 0, 0, 0, 1)$$

Thus: $x = W^T y$

Deconvolution \sim Transposed convolution

Let's rewrite: $x = W^T y$

Transposed convolution with stride 2

$$\begin{array}{c}
 \left(\begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \end{array} \right) = \left(\begin{array}{cccccc} w_1 & 0 & 0 & 0 & 0 & 1 \\ w_2 & 0 & 0 & 1 & -1 & -1 \\ w_3 & w_1 & 0 & 1 & 1 & 1 \\ w_4 & w_2 & 0 & 1 & 1 & 1 \\ 0 & w_3 & w_1 & 0 & 1 & 1 \\ : & w_4 & w_2 & 0 & 1 & 1 \\ 0 & w_3 & w_1 & : & 1 & 1 \\ w_4 & w_2 & 0 & : & 1 & 1 \\ 0 & w_3 & w_1 & : & 1 & 1 \\ w_4 & w_2 & 0 & : & 1 & 1 \\ 0 & w_3 & w_1 & : & 1 & 1 \\ w_4 & w_2 & 0 & : & 1 & 1 \\ 0 & w_3 & w_1 & : & 1 & 1 \\ w_4 & w_2 & 0 & : & 1 & 1 \end{array} \right) \left(\begin{array}{c} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{array} \right) + \left(\begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ w_3 \\ w_4 \end{array} \right)
 \end{array}$$

Sub-pixel convolution with stride 1/2

$$\begin{bmatrix} 0 \\ 0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} w_4 & w_3 & w_2 & w_1 & 0 & \cdots & 0 \\ w_4 & w_3 & w_2 & w_1 & 0 & \cdots & 0 \\ w_4 & w_3 & w_2 & w_1 & 0 & \cdots & 0 \end{bmatrix} \begin{pmatrix} 0 \\ 0 \\ y_1 \\ 0 \\ y_2 \\ 0 \\ y_3 \\ 0 \\ y_4 \\ 0 \\ y_5 \\ 0 \\ 0 \end{pmatrix}$$

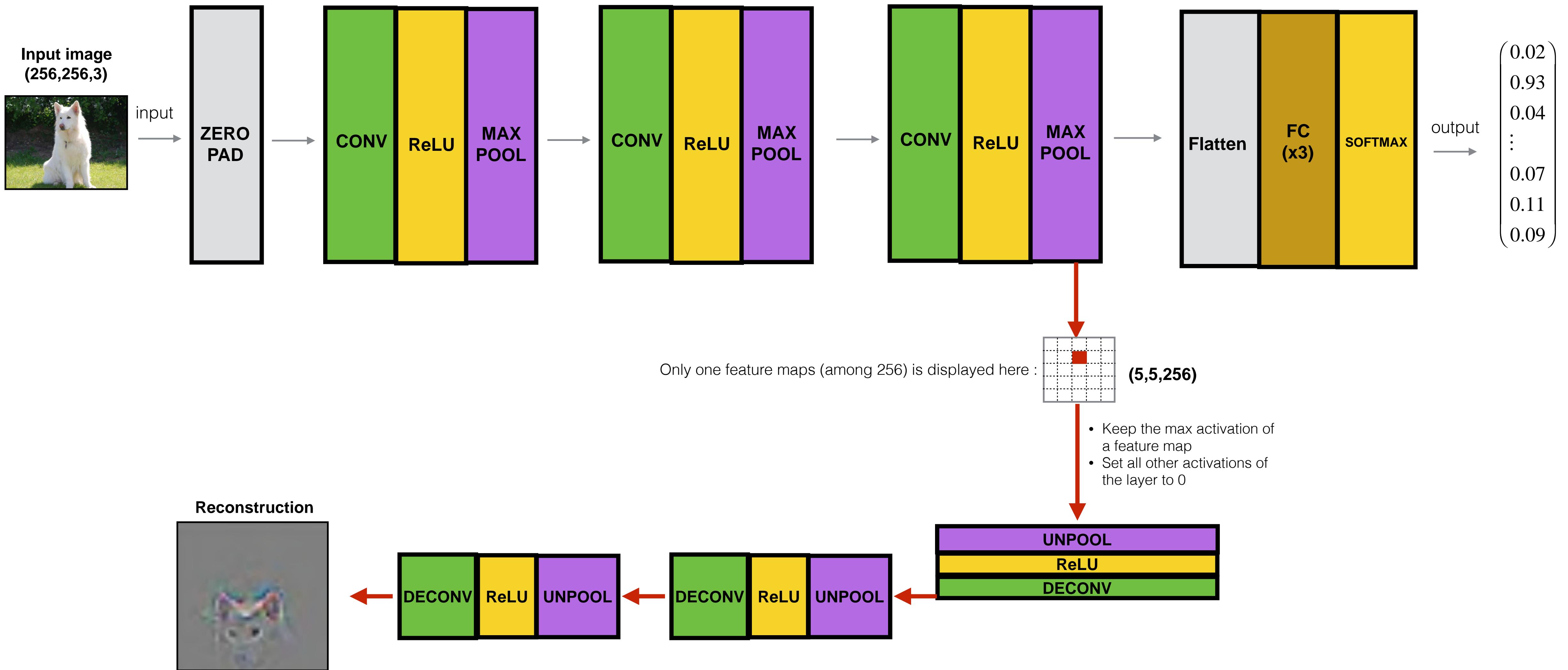
What to remember?

Implementing a deconvolution (sub-pixel version) is similar to the convolution, with:

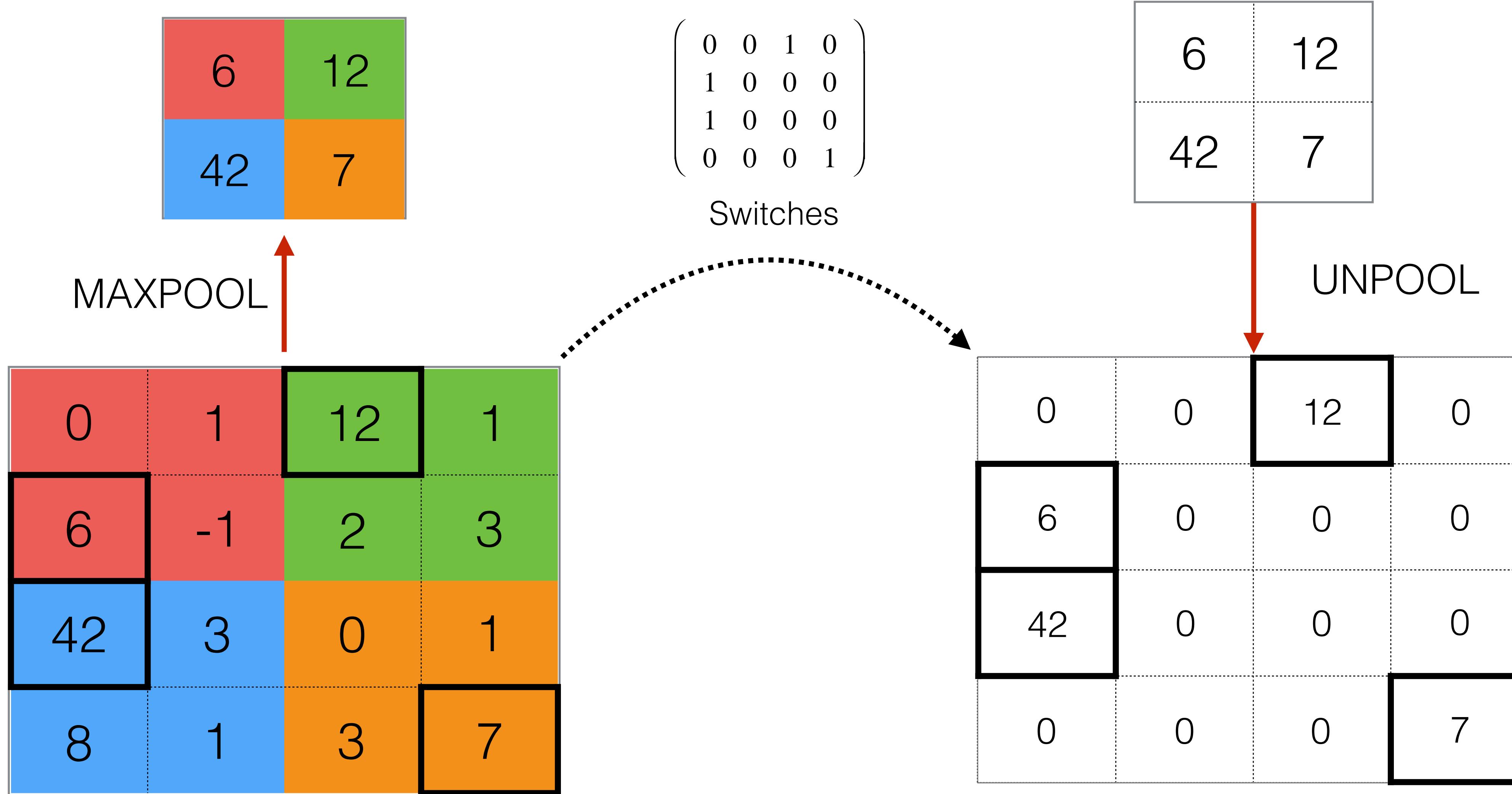
1. Create a sub-pixel version of the input (i.e., insert zeros and pad)
2. Flip the filters.
3. Divide the stride by 2.

III. A. Interpreting NNs using deconvolutions

Motivation of DeconvNets for visualization: Here is a CNN, trained on ImageNet (1.3m images, 1000 classes), we're trying to interpret by reconstructing the activation's zone of influence in the input space.

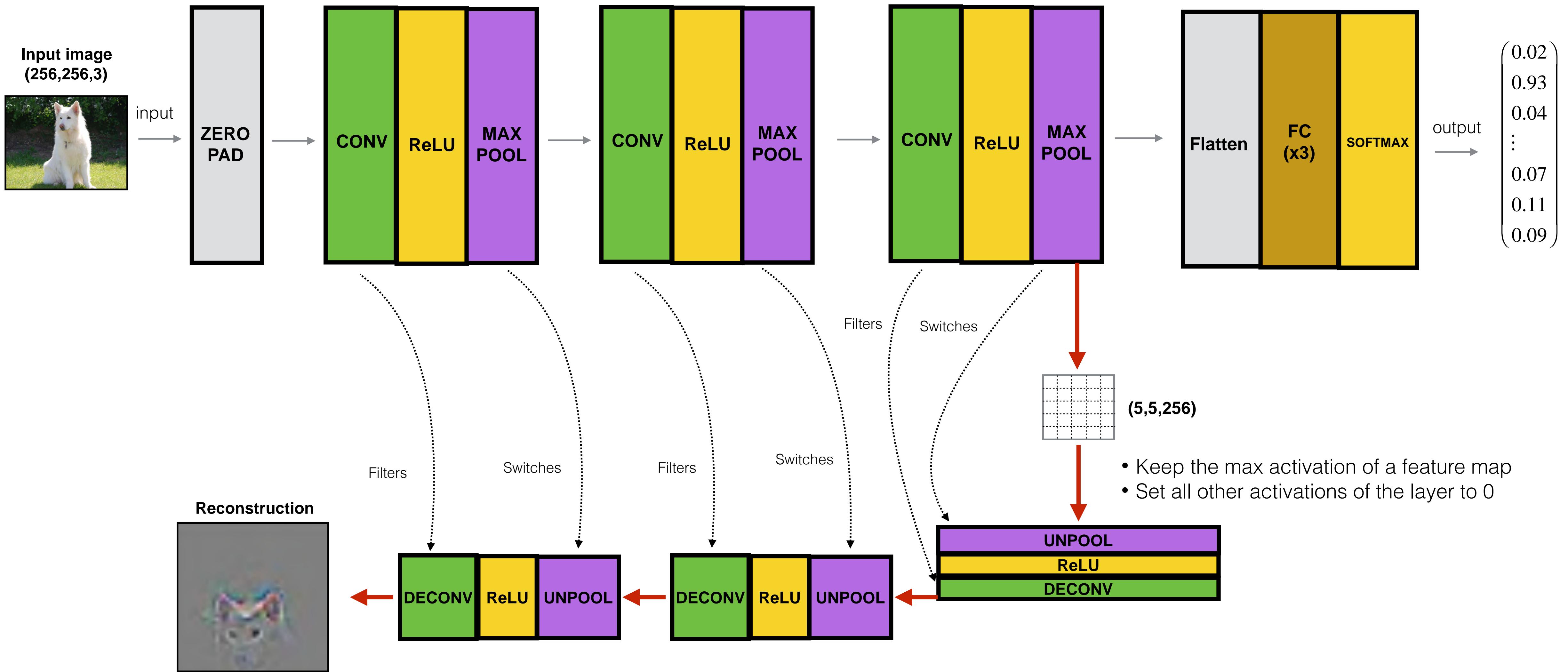


II. C. The deconvolution and its applications



II. C. The deconvolution and its applications

We need to pass the filters and switches from the ConvNet to the DeconvNet.



II. C. The deconvolution and its applications

$$a^{[l]} = \mathbf{I}\{a^{[l+1]} \geq 0\} \cdot a^{[l+1]}$$

1	0	3	0
4	0	0	0
30	2	0	1
1	0	0	7

“ReLU forward”

1	-2	3	-4
4	-1	-1	-2
30	2	0	1
1	-2	-9	7

Switches

$$\begin{pmatrix} 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{pmatrix}$$

-2	1	-31	-3
3	4	2	14
-2	12	4	1
10	10	2	1

“ReLU backward”

-2	0	-31	0
3	0	0	0
-2	12	0	1
10	0	0	1

-2	1	-31	-3
3	4	2	14
-2	12	4	1
10	10	2	1

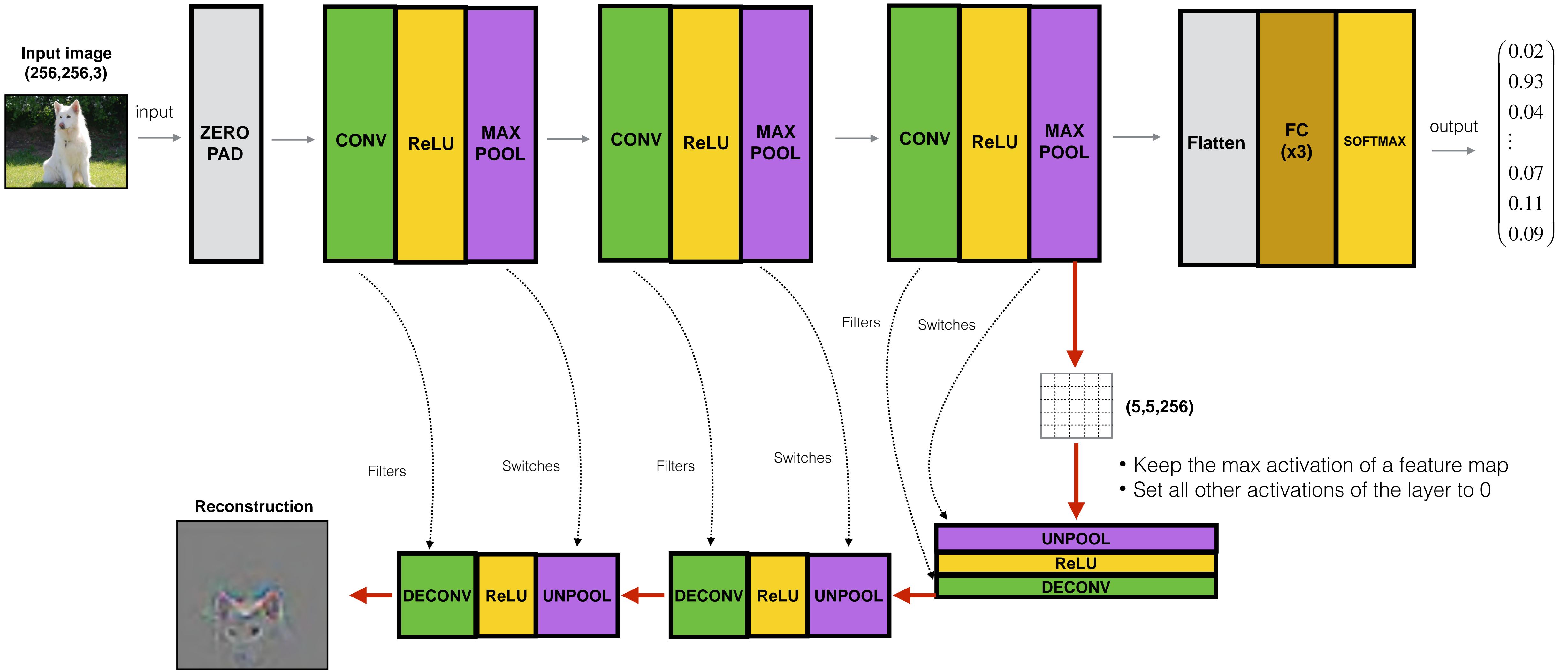
“ReLU DeconvNet”

0	1	0	0
3	4	2	14
0	12	4	1
10	10	2	1

Kian Katanforoosh

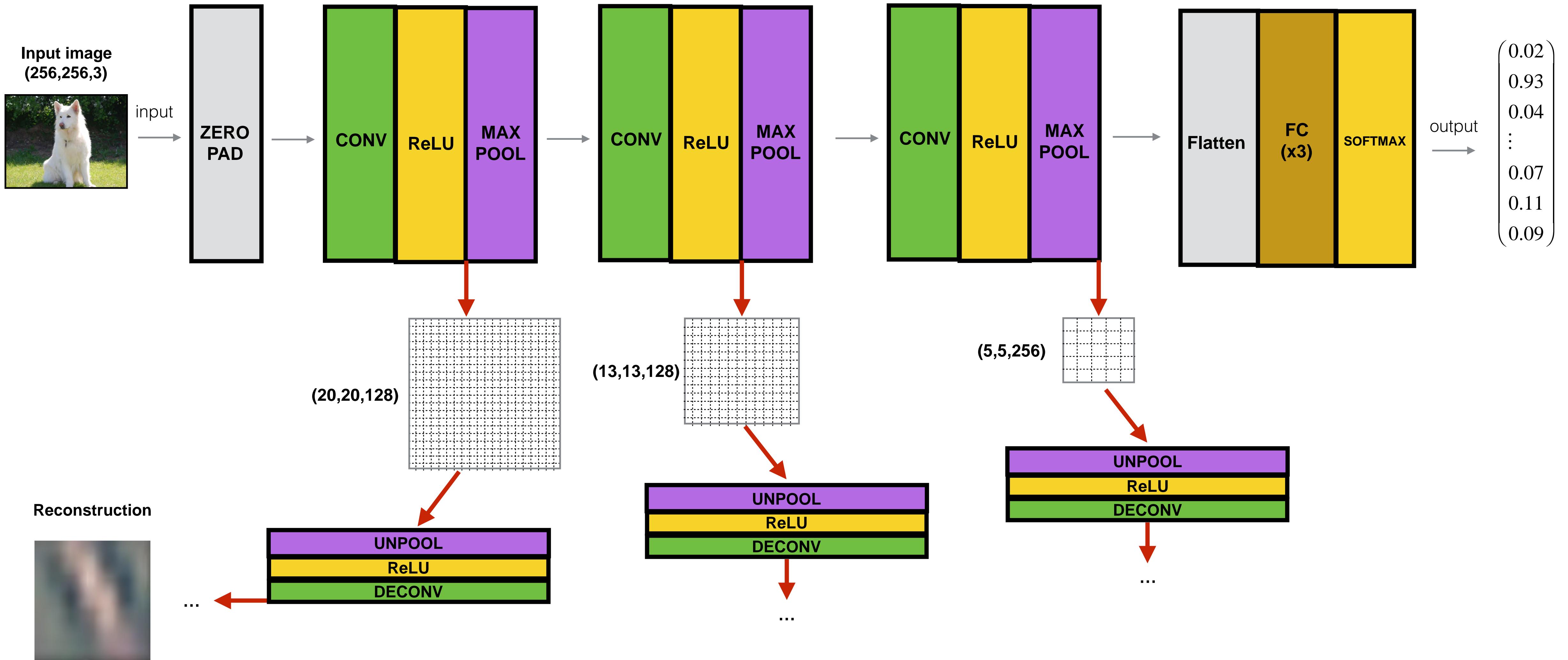
II. C. The deconvolution and its applications

We need to pass the filters and switches from the ConvNet to the DeconvNet.



II. C. The deconvolution and its applications

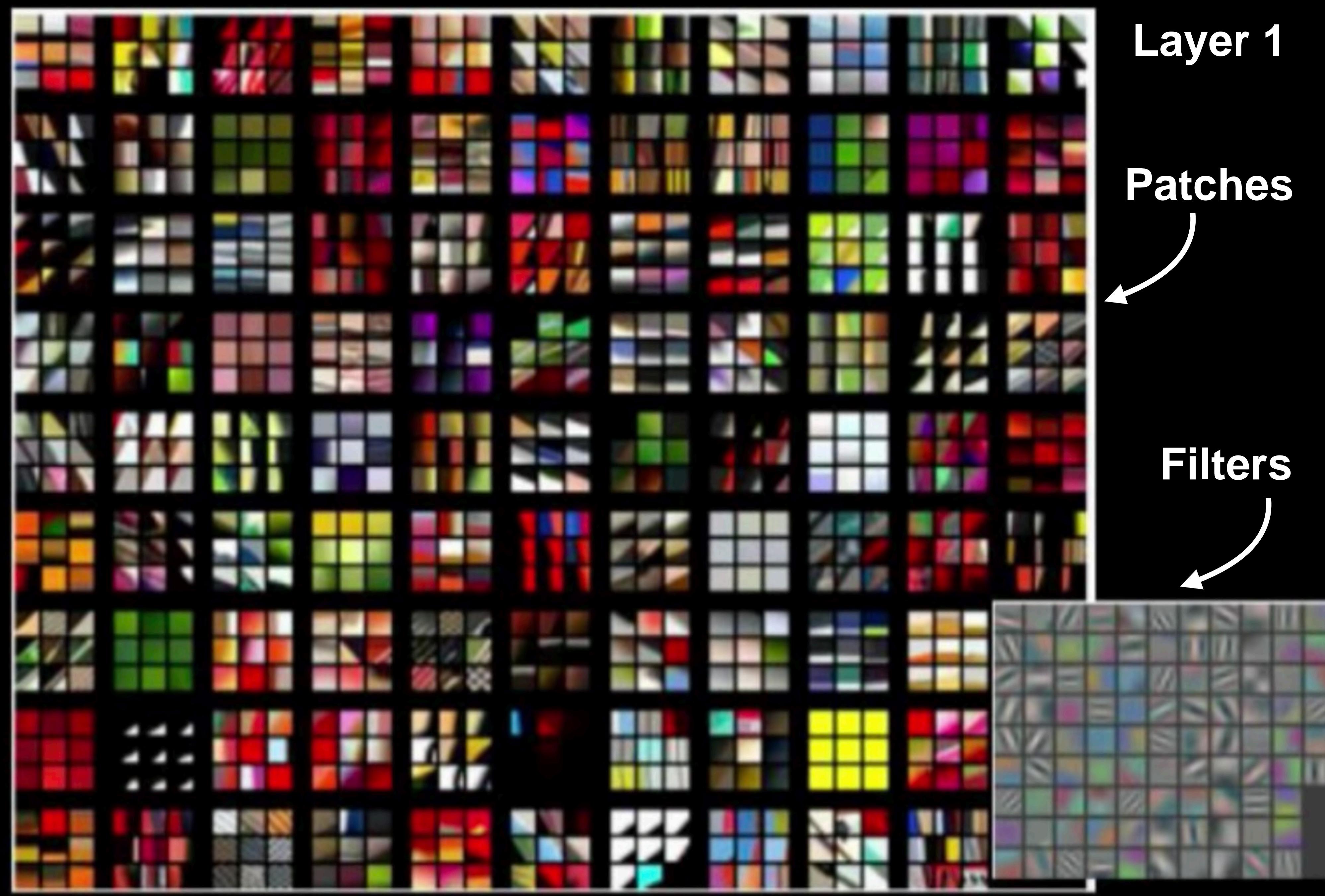
Other CONV layers can be visualized the same way



Kian Katanforoosh

Results on a validation set of 50,000 images

- Top-9 strongest activations per filter in the 1st layer
- Because we know the position of the activation and all the pooling switches we can crop the part of the image that fired the activation

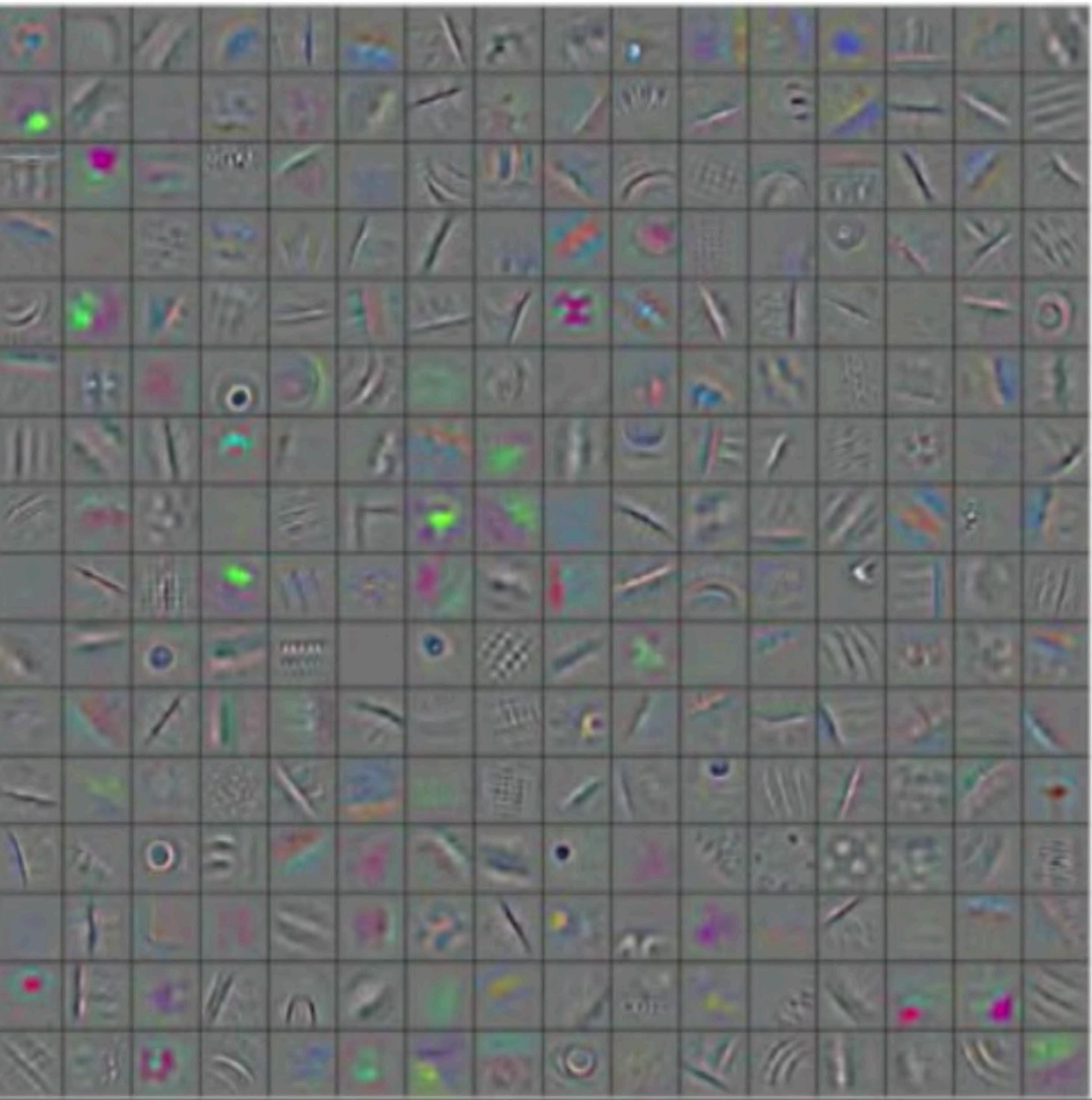


Kian Katanforoosh

Results on a validation set of 50,000 images

- Learning a more complex set of patterns than 1st layer edges.
- Covers a much larger space of the image because of the pooling layer before.
- Top-1 strongest activation per feature map in the 2nd layer (256 feature maps.)

Layer 2 reconstructions (deconv)

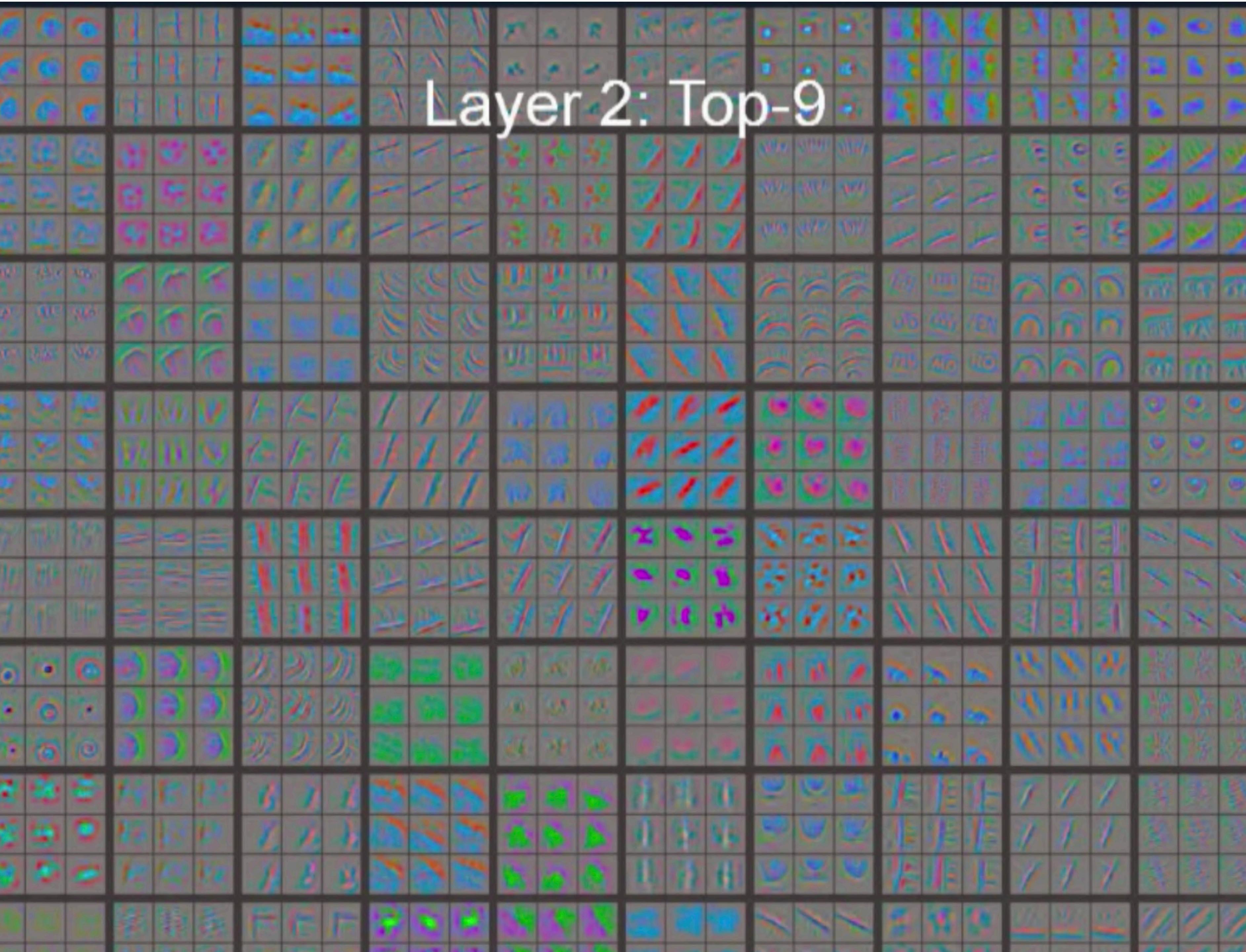


Kian Katanforoosh

Results on a validation set of 50,000 images

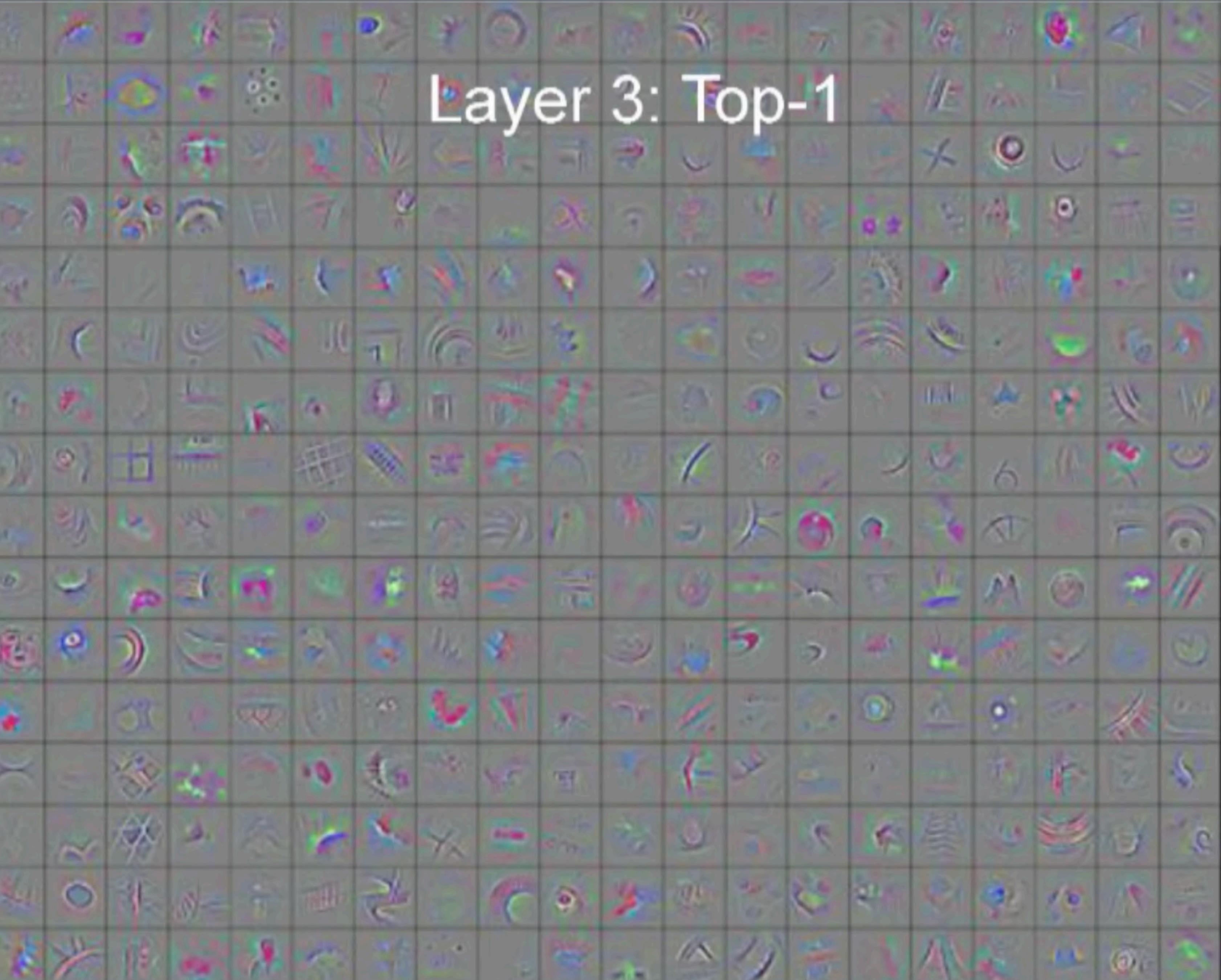
- Learning a more complex set of patterns than 1st layer edges
- Covers a much larger space of the image probably because of the pooling layer before.
- Features are more invariant to small changes. Ex: A dot, spiral, circle all fire the same 2nd layer feature very strongly

Layer 2: Top-9



Results on a validation set of 50,000 images

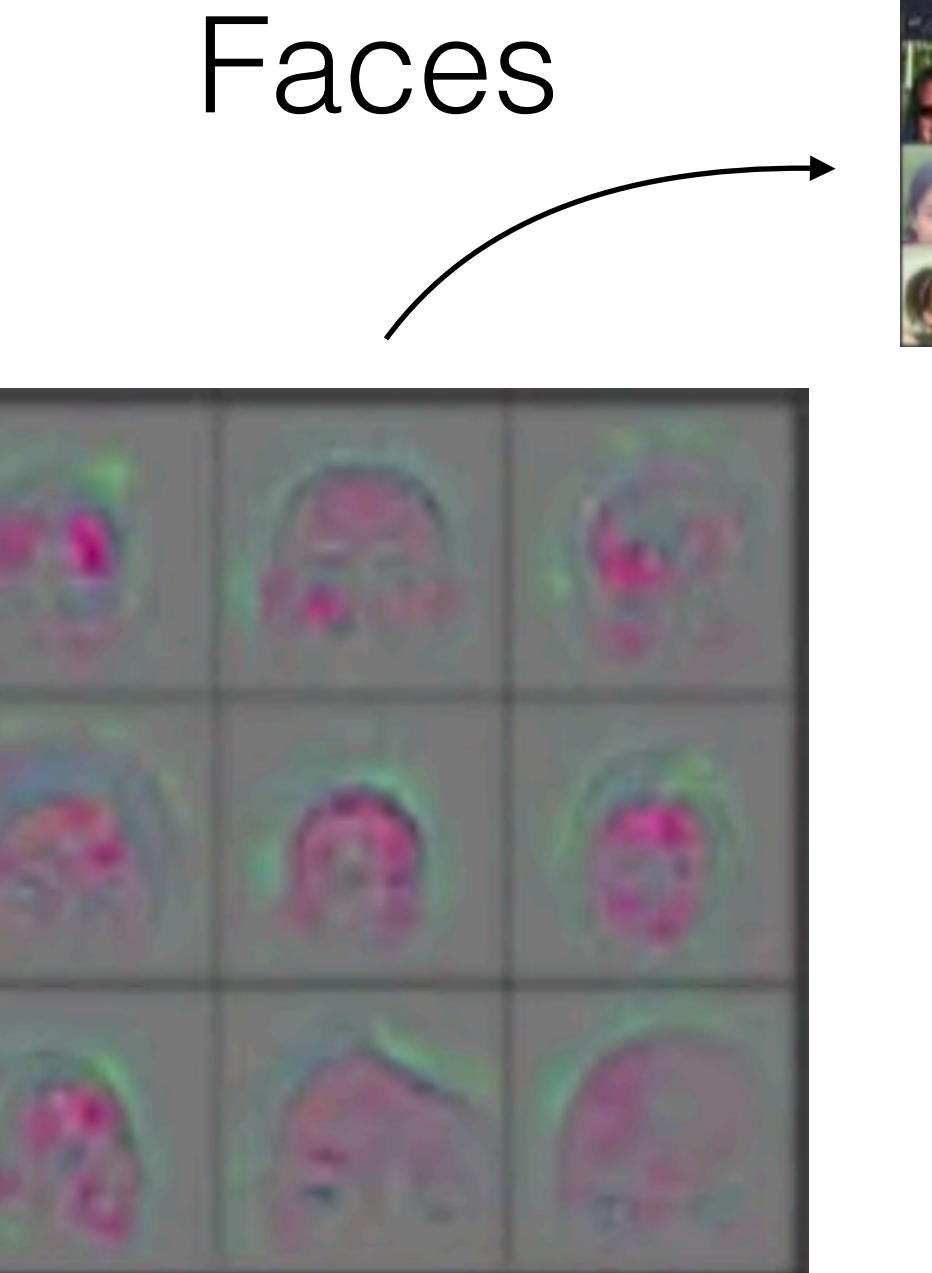
- 3rd layer: increased complexity
- An activated neuron is seeing $\approx 80 \times 80$ part of a 256×256 image
- Learning objects, faces etc..



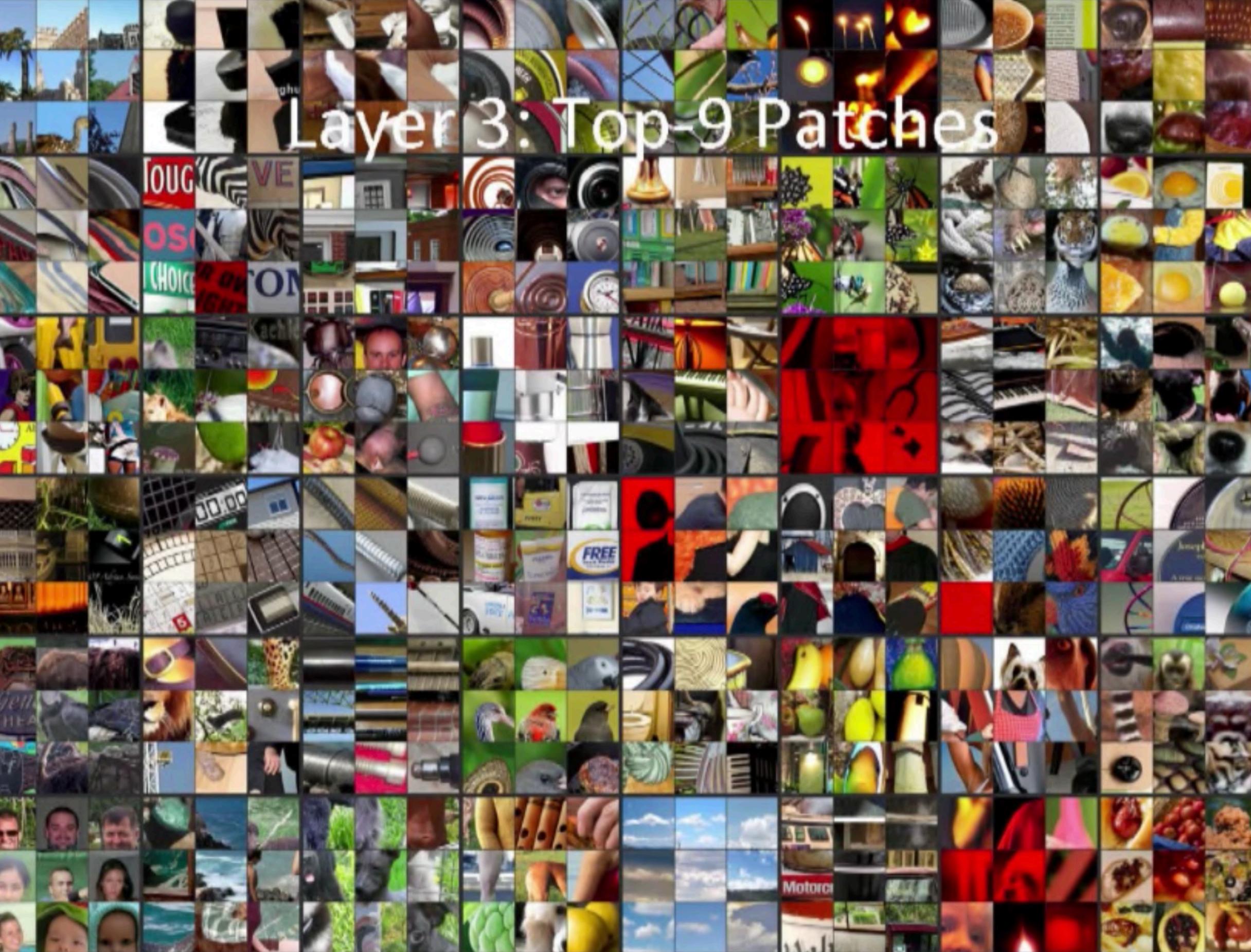
Layer 3: Top-1

Results on a validation set of 50,000 images

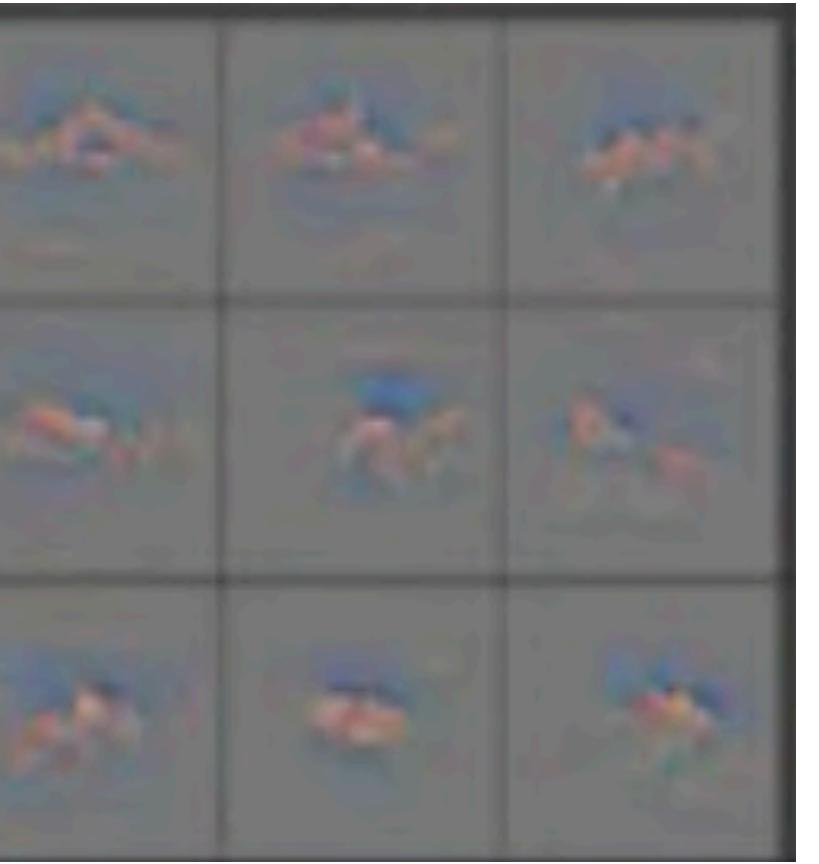
- 3rd layer: increased complexity
- An activated neuron is seeing $\approx 80 \times 80$ part of a 256×256 image
- Learning objects, faces etc..
- Patches: Semantic grouping, not structural



Faces

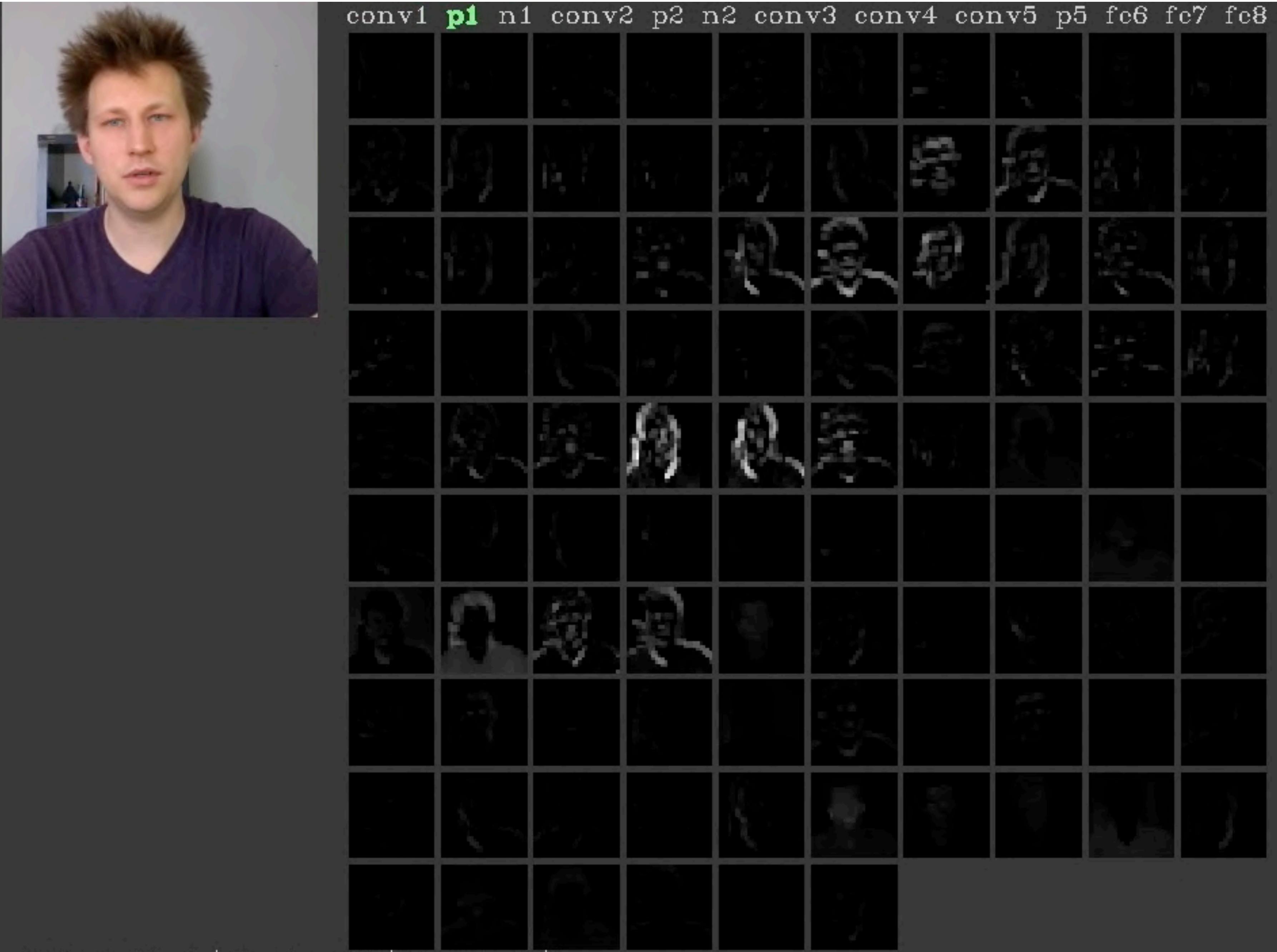


Layer 3: Top-9 Patches



Clouds

Kian Katanforoosh



fwd pool1_4 | Back: off | Boost: 0/1

[Link to video: <https://www.youtube.com/watch?v=AgkfIQ4IGaM>]

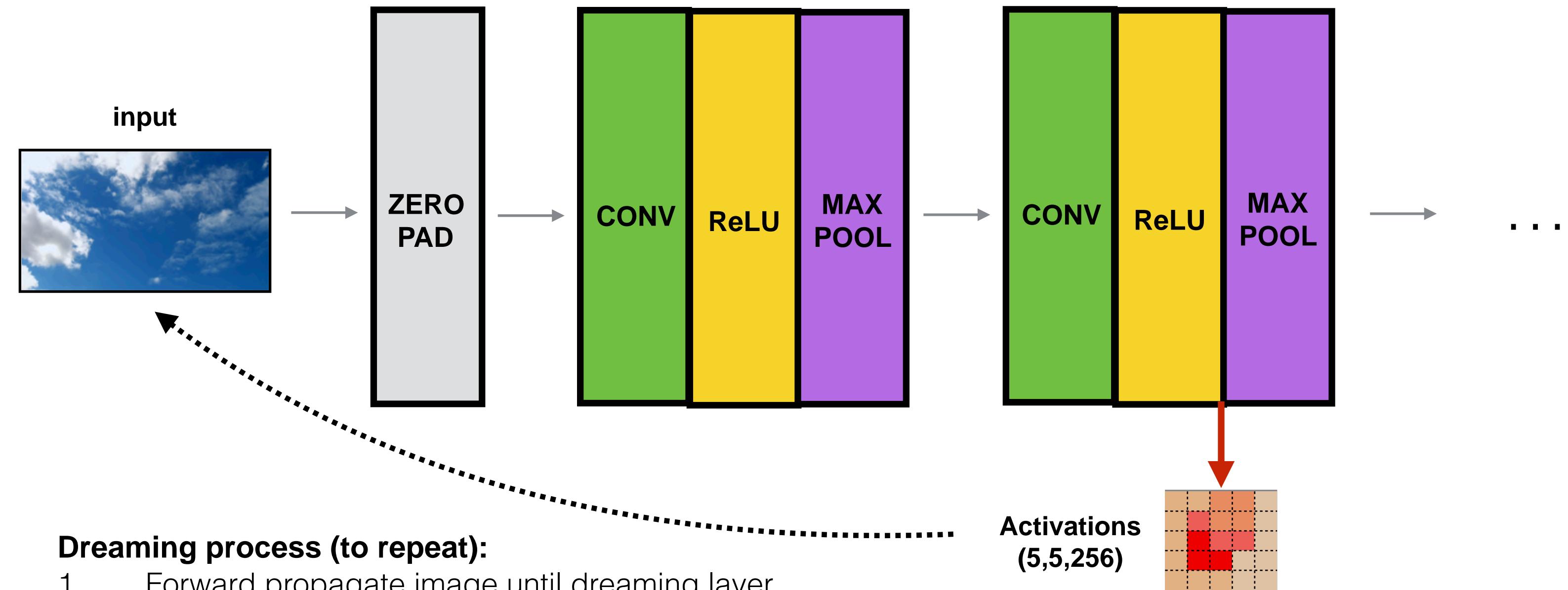
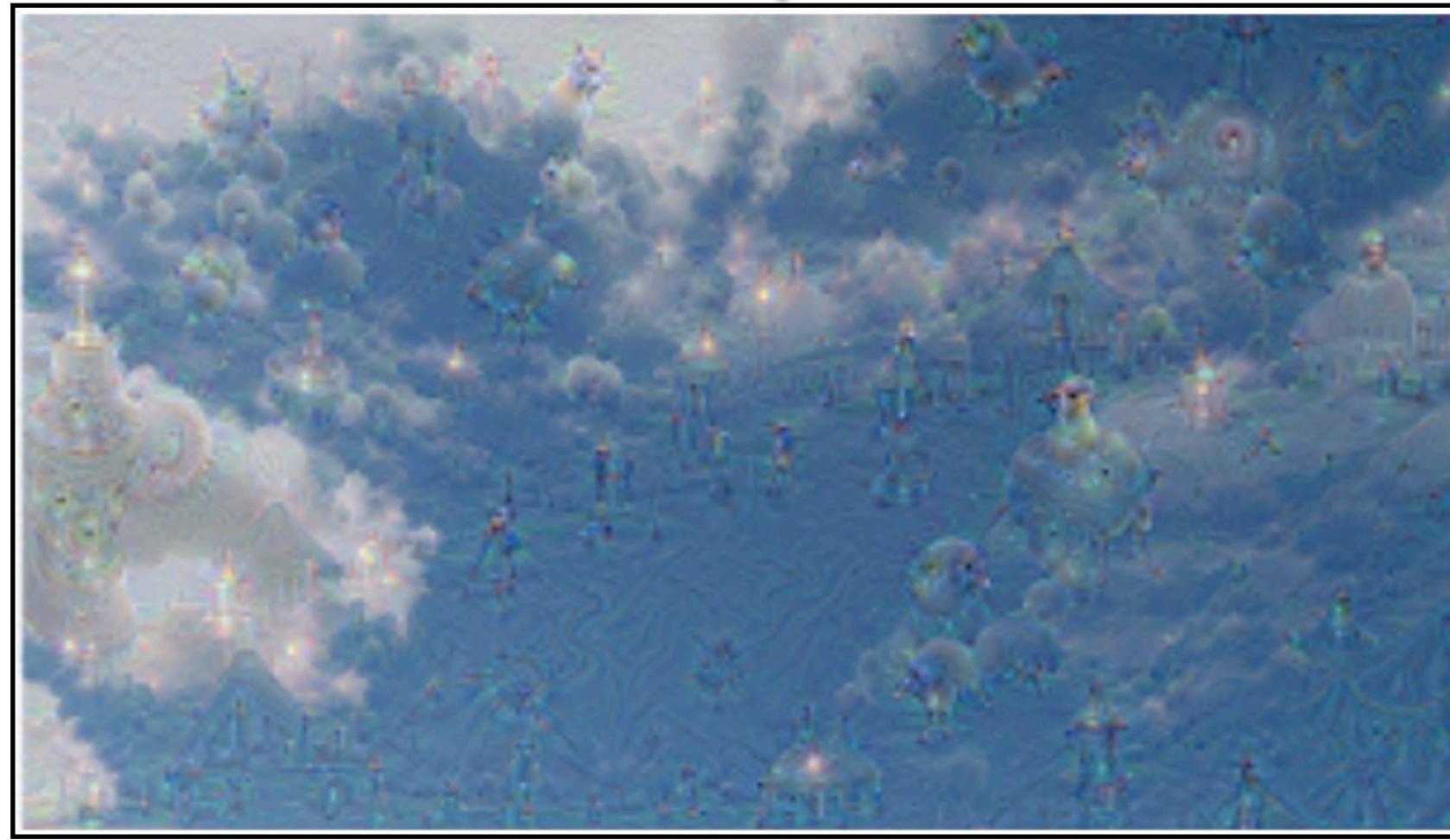
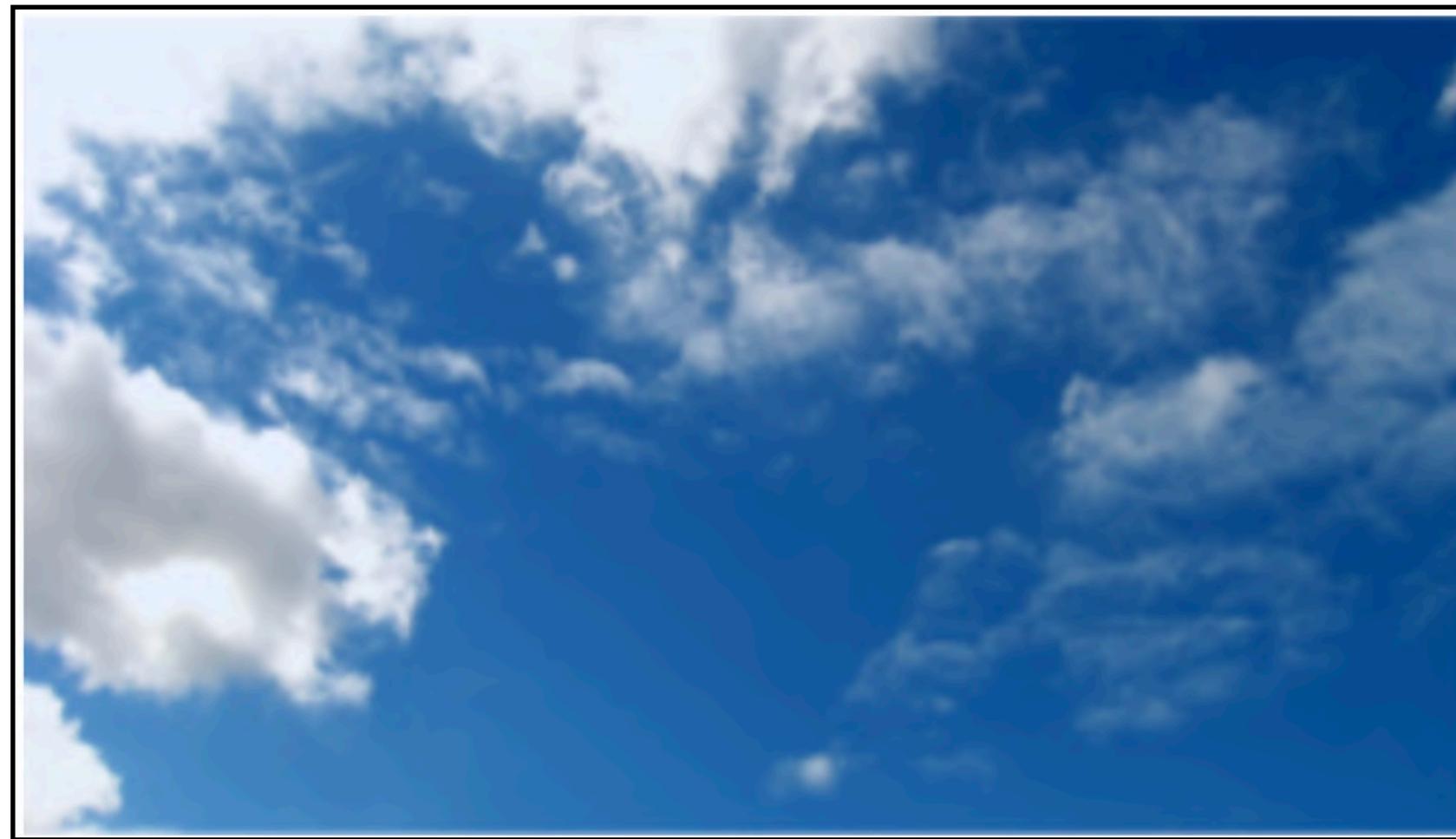
[Jason Yosinski et al. (2015): Understanding Neural Networks Through Deep Visualization]

Today's outline

- I. Interpreting Neural Networks' outputs
 - A. with saliency maps
 - B. with occlusion sensitivity
 - C. with class activation maps (Global Average Pooling)
- II. Visualizing Neural Networks from the inside
 - A. with gradient ascent (class model visualization)
 - B. with dataset search
 - C. The deconvolution and its applications
- III. (Optional: Deep Dream: going deeper in NNs)**

III. Deep Dream: going deeper in NNs

How to boost the activation of a neuron?

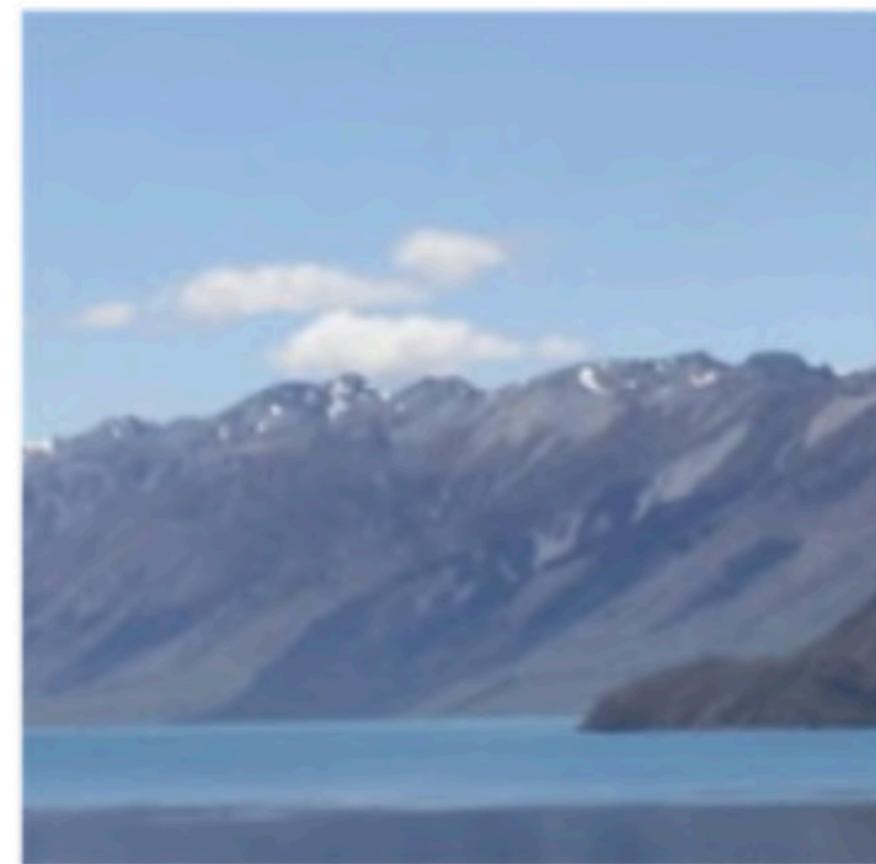


Dreaming process (to repeat):

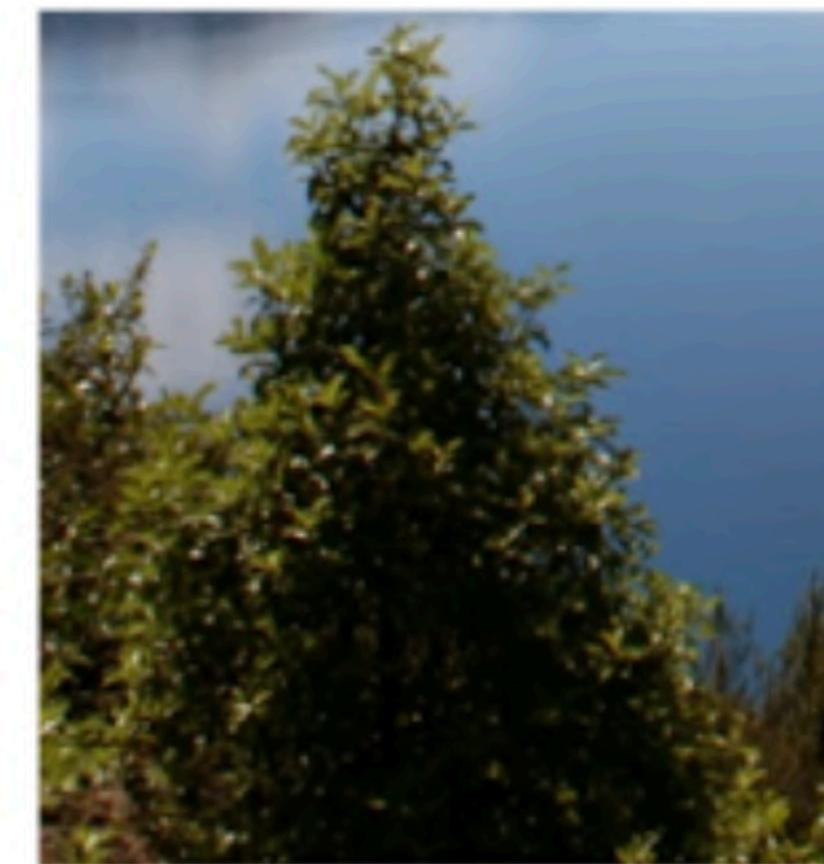
1. Forward propagate image until dreaming layer
2. **Set gradients of dreaming layer to be equal to its activations**
3. Backpropagate gradients to input image
4. Update Pixels of the image



III. Deep Dream: going deeper in NNs



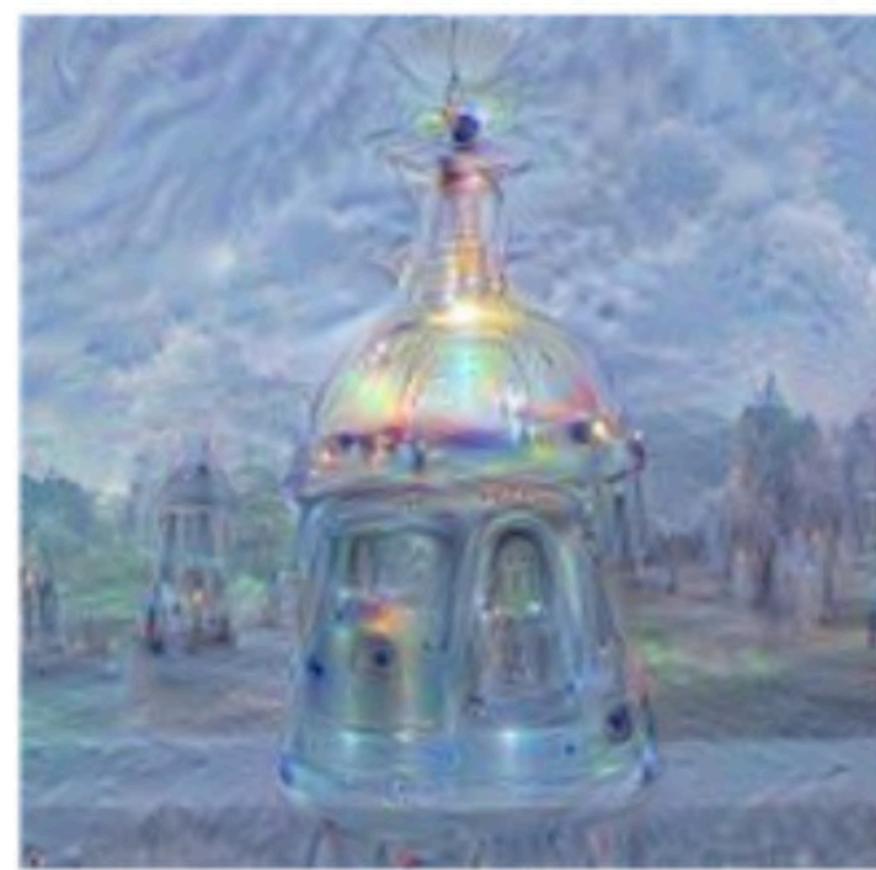
Horizon



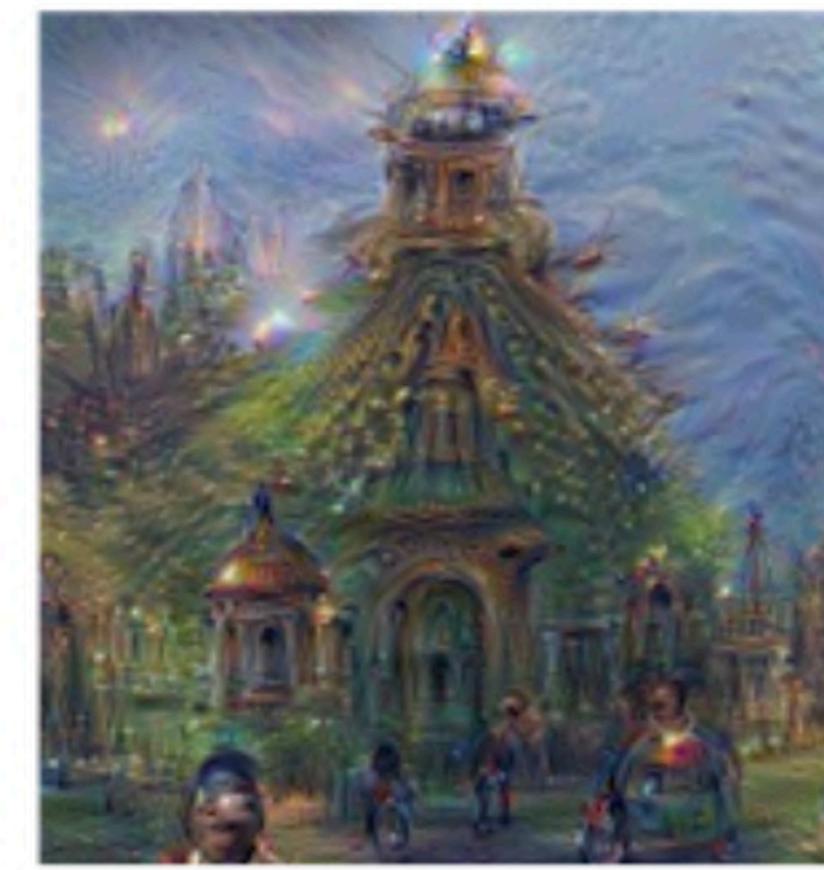
Trees



Leaves



Towers & Pagodas



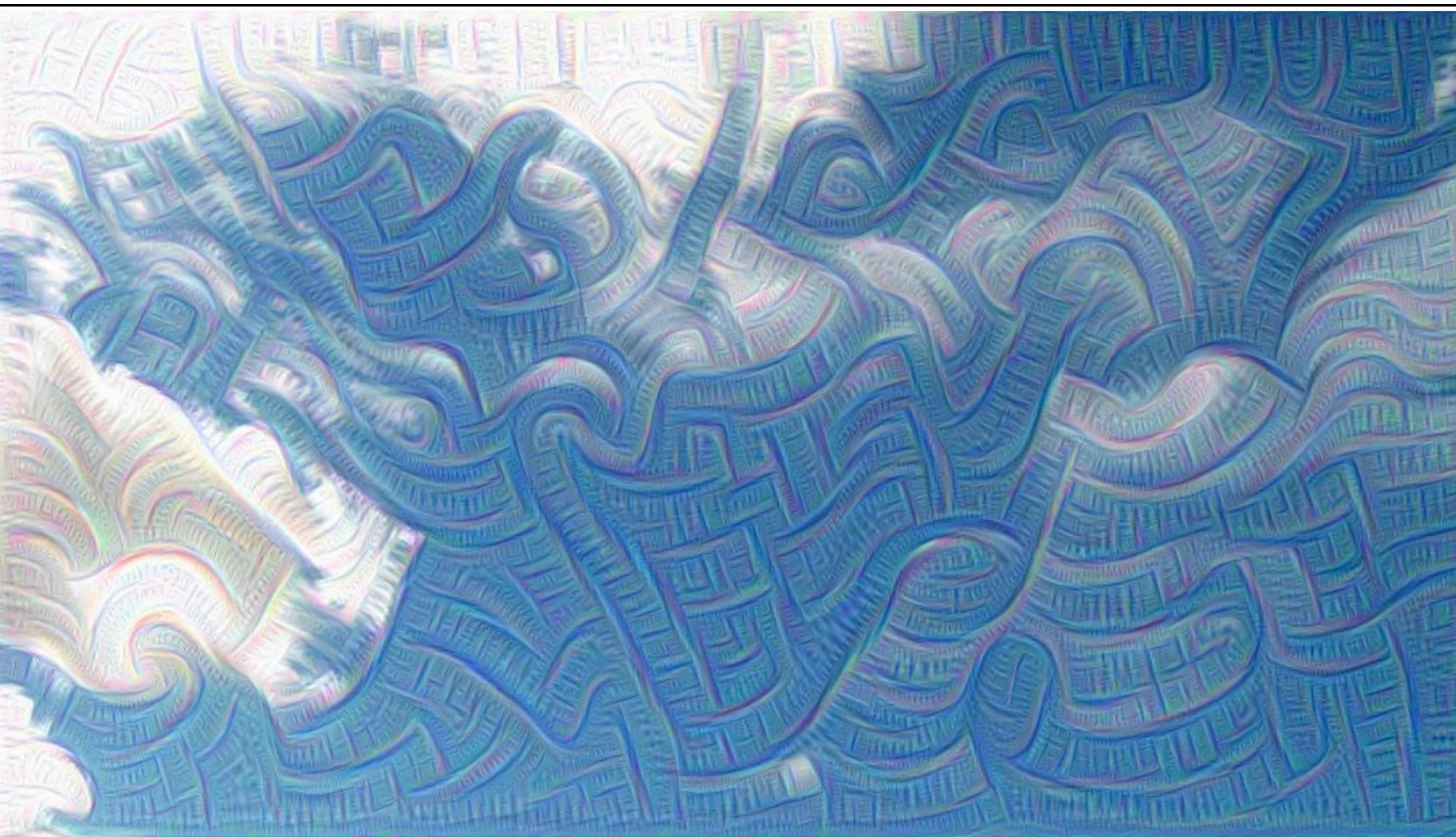
Buildings



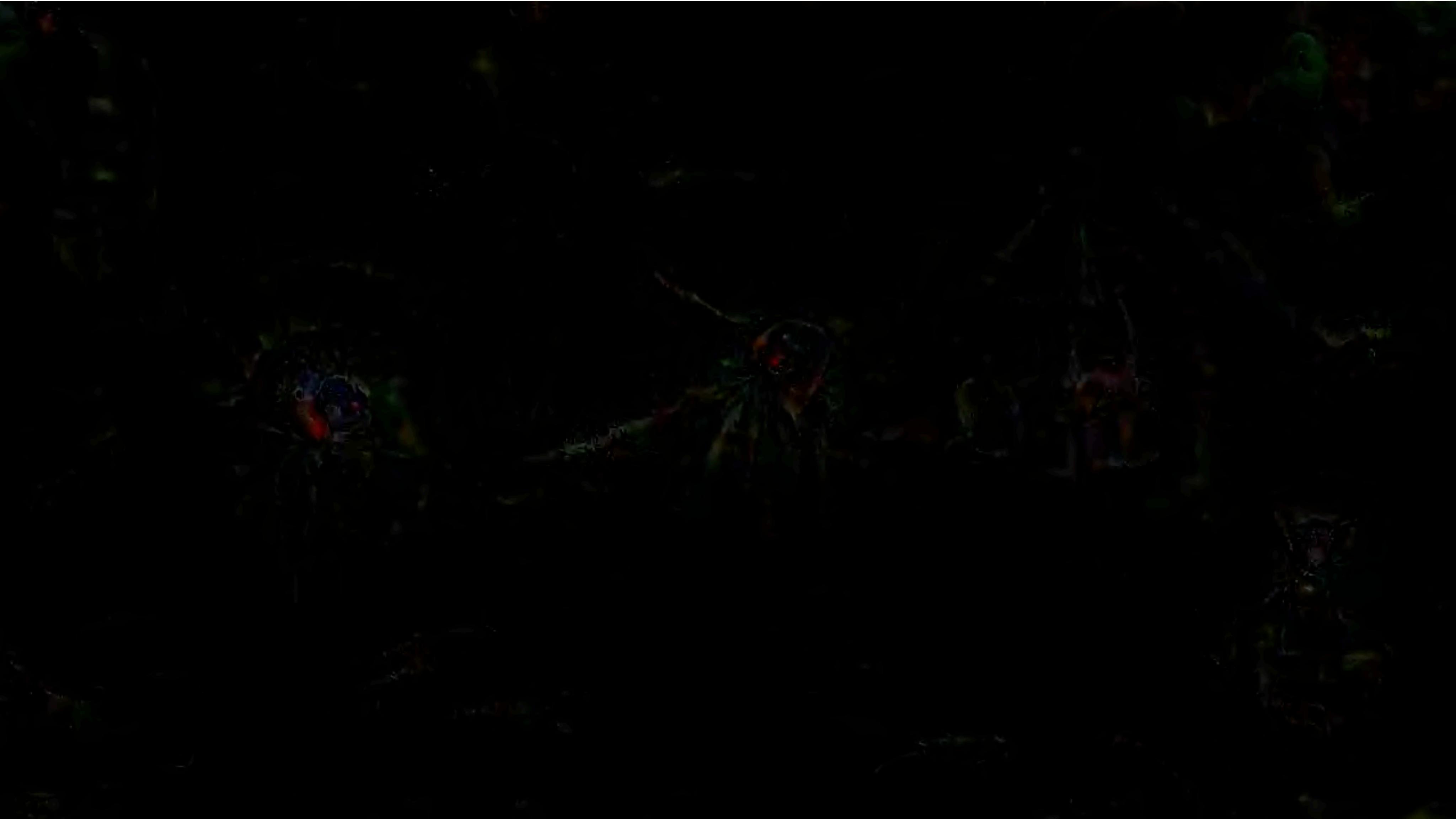
Birds & Insects

III. Deep Dream: going deeper in NNs

If you dream in lower layers:



III. Deep Dream: going deeper in NNs



[Link to video: <https://www.youtube.com/watch?v=DgPaCWJL7XI>]

[Github repo (deepdream): <https://github.com/google/deepdream/>]

[Alexander Mordvintsev et al. (2015): [Inceptionism: Going Deeper into Neural Networks](#)]

Kian Katanforoosh

III. Deep Dream

Did the neural network learned the right features to detect an object?



Dumbbells

the network failed to understand the essence of a dumbbell

Questions we are now able to answer:

- *What part of the input is responsible for the output?*
 - *Occlusion sensitivity*
 - *Class Activation Maps*
- *What is the role of a given neuron/filter/layer?*
 - *Deconvolutions can help visualize the role of a neuron*
 - *Search dataset images maximizing the activation*
 - *Gradient ascent (class model visualization)*
- *Can we check what the network focuses on given an input image?*
 - *Occlusion sensitivity*
 - *Saliency maps (one-time gradient ascent)*
 - *Class Activation Maps*
- *How does a neural network see our world?*
 - *Gradient ascent (class model visualization)*
 - *Deep Dream*
- *Do these visualization have use cases?*
 - *Segmentation (saliency maps)*
 - *Art (Deep Dream)*

Duties for next week

Meet with your assigned TA between 1/28 and 2/26 to discuss your milestone report.

Project Milestone due 02/26 Friday 11:59 PM

Completed modules for 03/04:

- C5M1: Recurrent Neural Networks ([slides](#))

Quizzes (due at 8 30am PST):

- Recurrent Neural Networks

Programming Assignments (due at 8 30am PST):

- Building a Recurrent Neural Network - Step by Step
- Dinosaur Land -- Character-level Language Modeling
- Jazz improvisation with LSTM