

VISUALIZING CNNs

[interpreting CNN decisions]

Attribution techniques

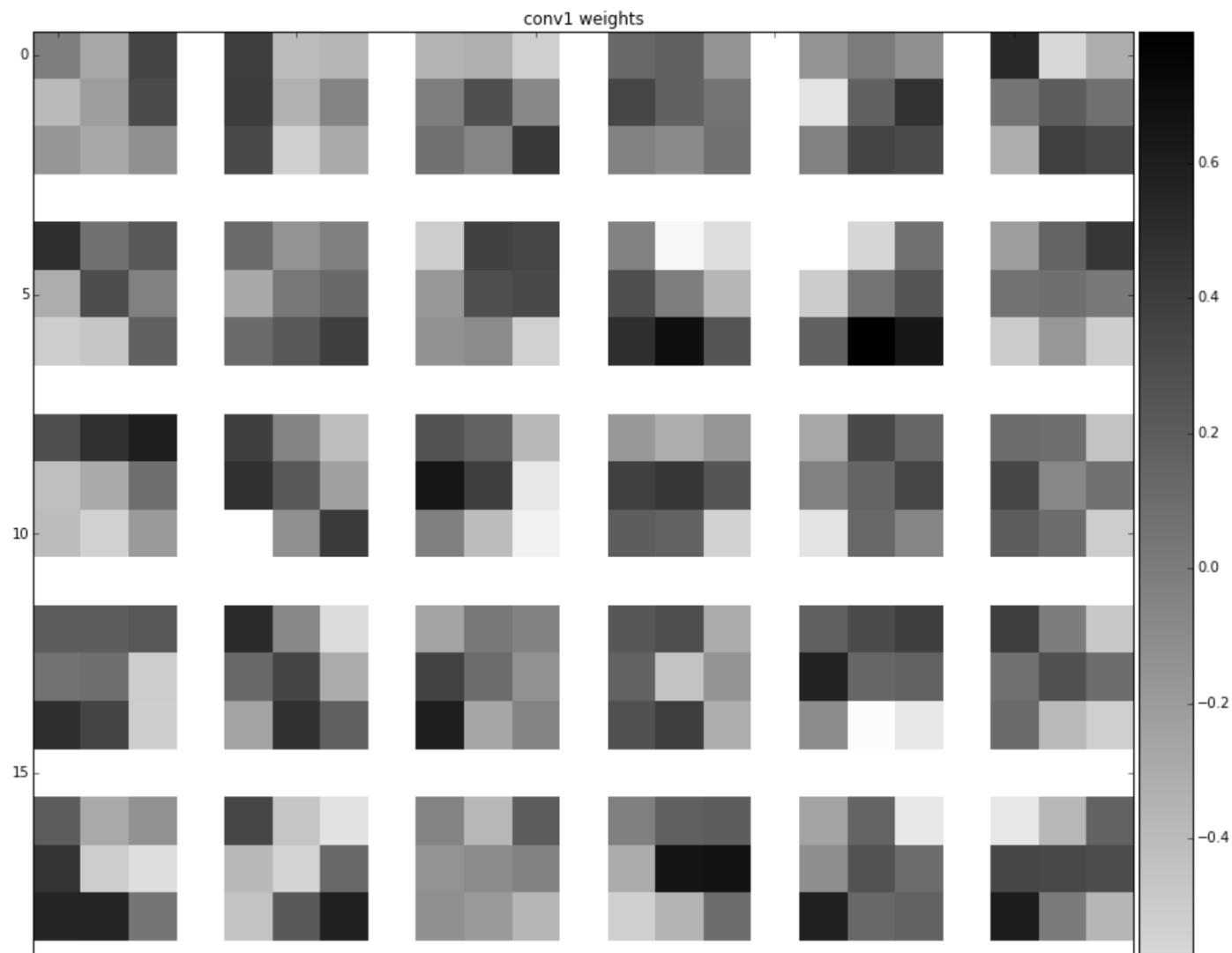
DEEP NETWORKS ARE “BLACK BOXES”?

INTERPRETING THE RESULTS IS
EXTREMELY DIFFICULT

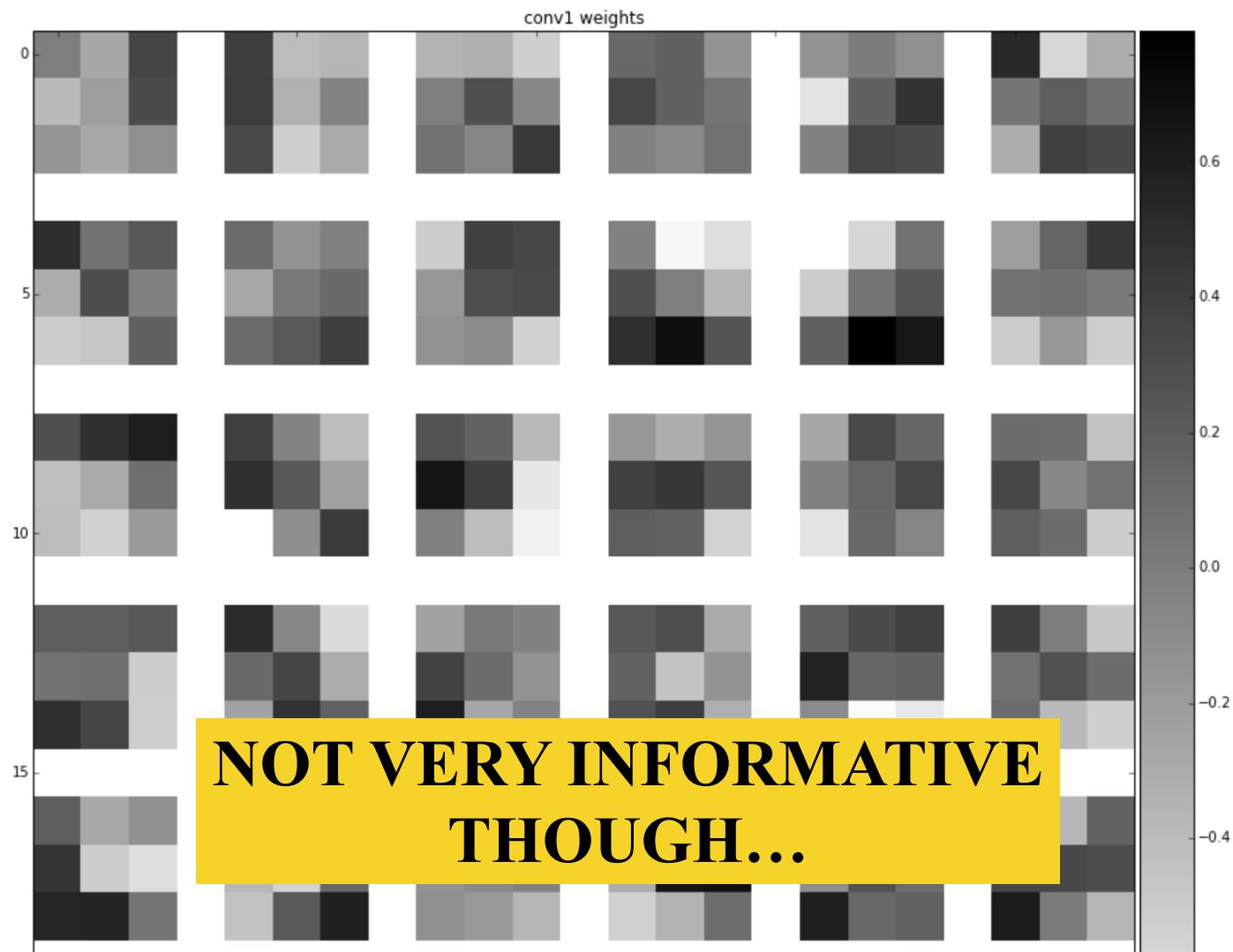
THIS IS TRUE BUT A LOT OF WORK
IS DONE TO UNVEIL THEIR BEHAVIOR

EXPLORING THE FEATURE MAPS

THE SIMPLEST APPROACH IS TO VISUALIZE THE LEARNED WEIGHTS AT INTERMEDIATE LAYERS

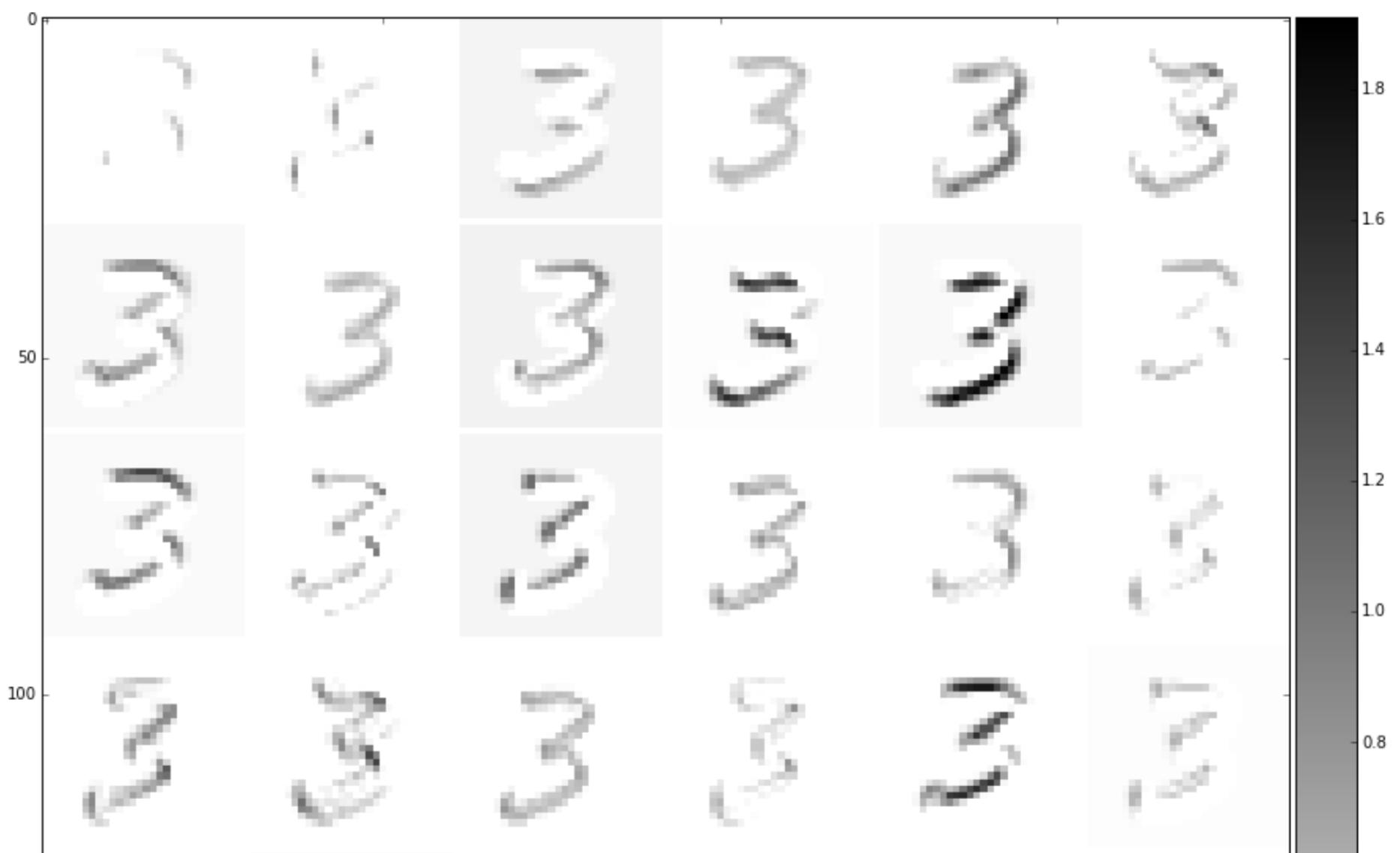


THE SIMPLEST APPROACH IS TO VISUALIZE THE LEARNED WEIGHTS AT INTERMEDIATE LAYERS

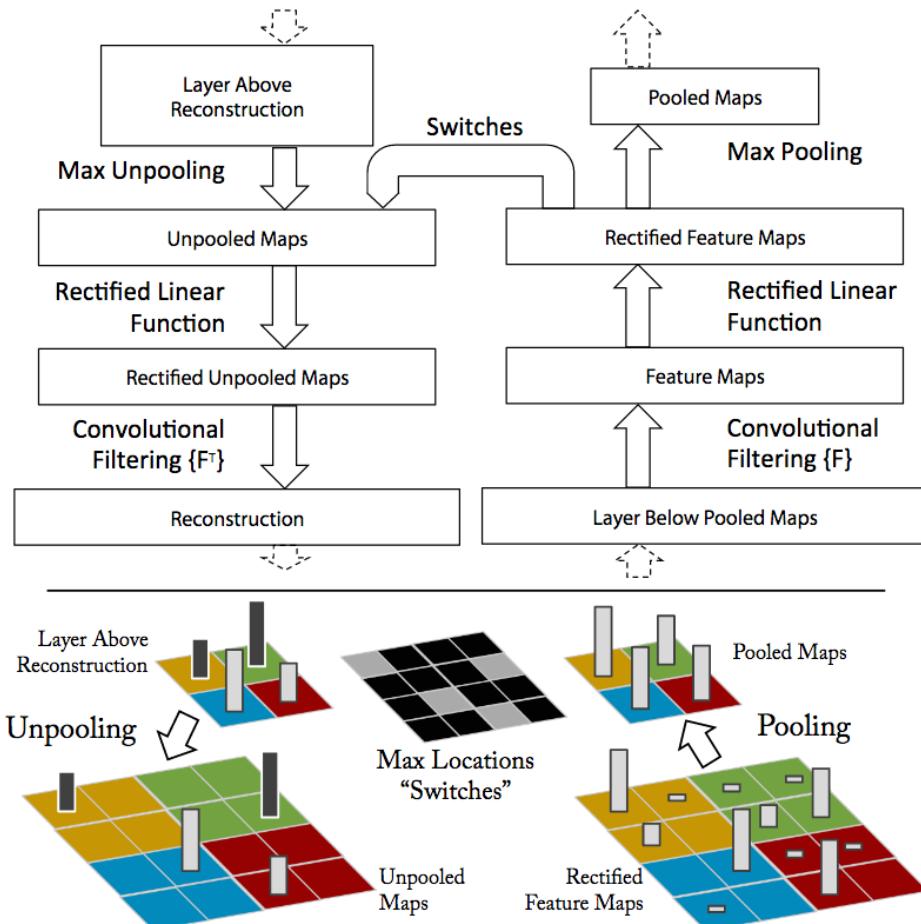


USING THE SAME IDEA, ONE CAN ALSO VISUALIZE
THE FEATURE MAPS AT INTERMEDIATE LAYERS

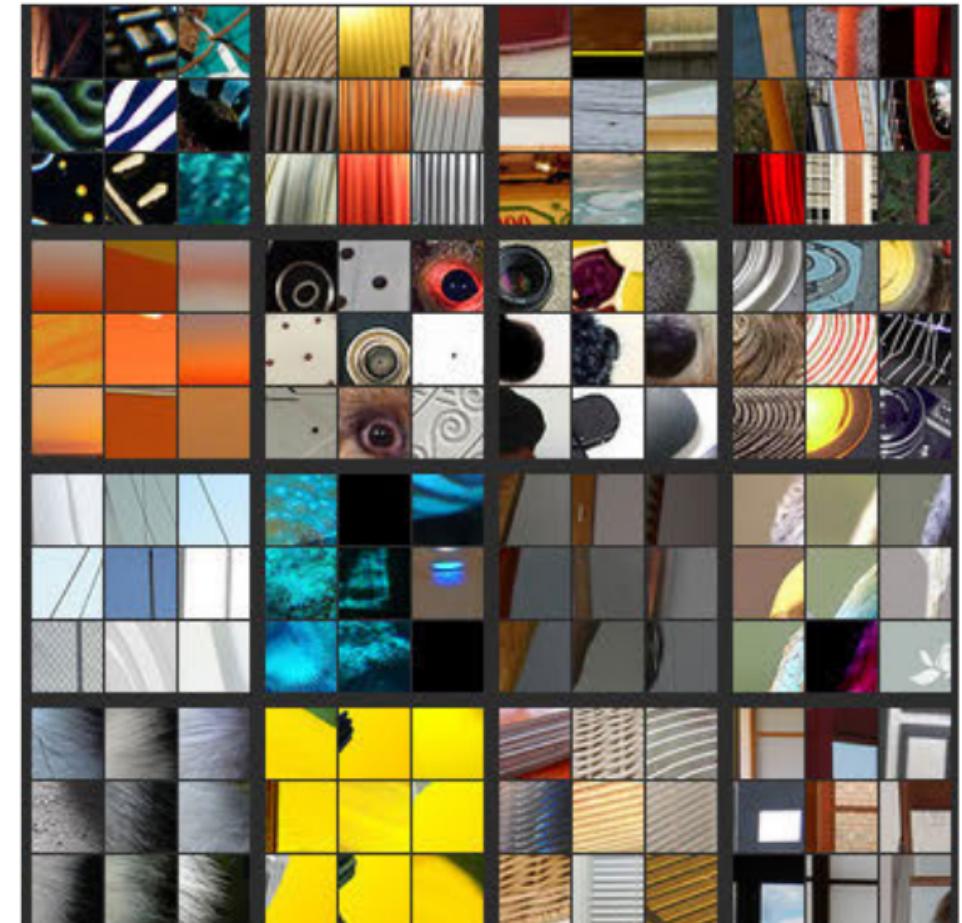
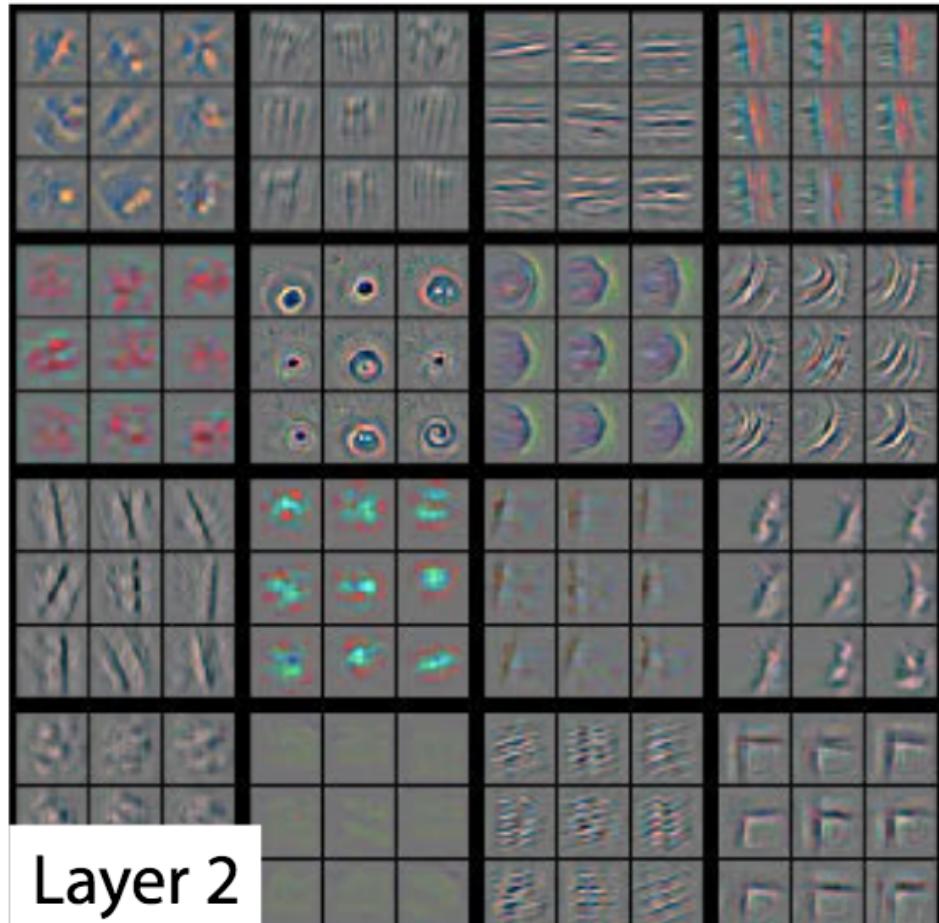
THIS HELPS TRACING THE FEATURES LEARNED BY THE
NETWORK



USE “DECONVNETS” TO MAP BACK THE FEATURE MAP INTO THE PIXEL SPACE

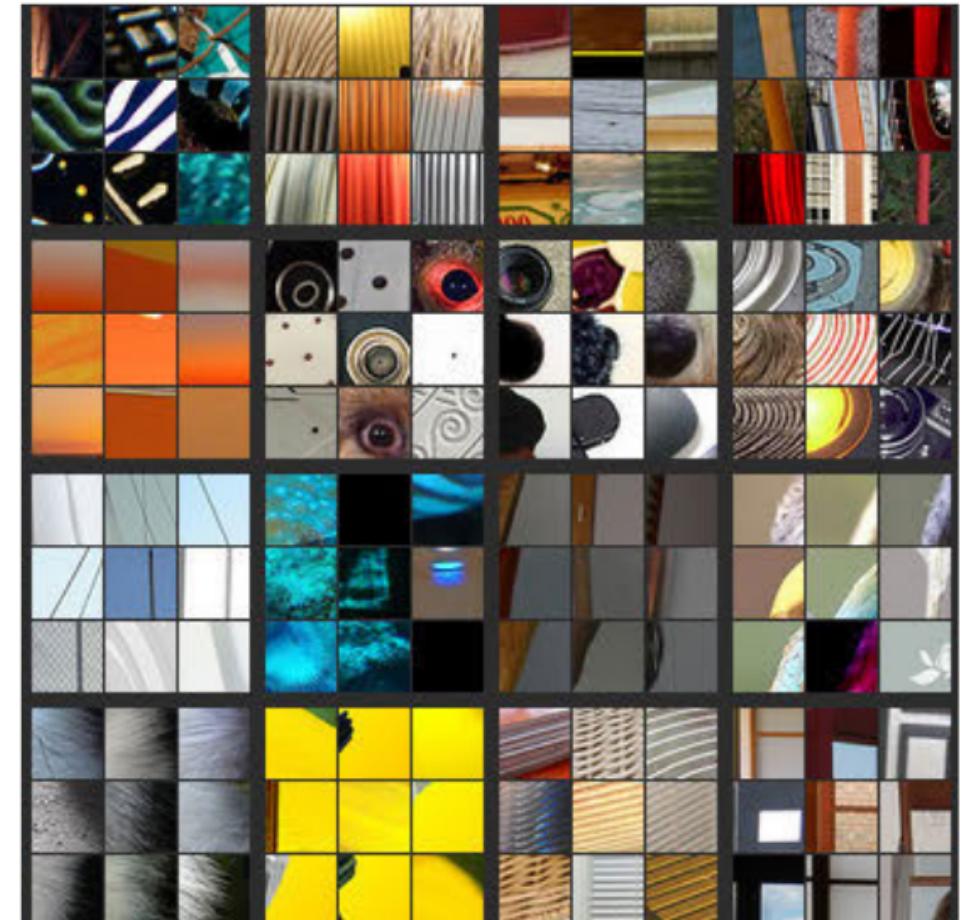
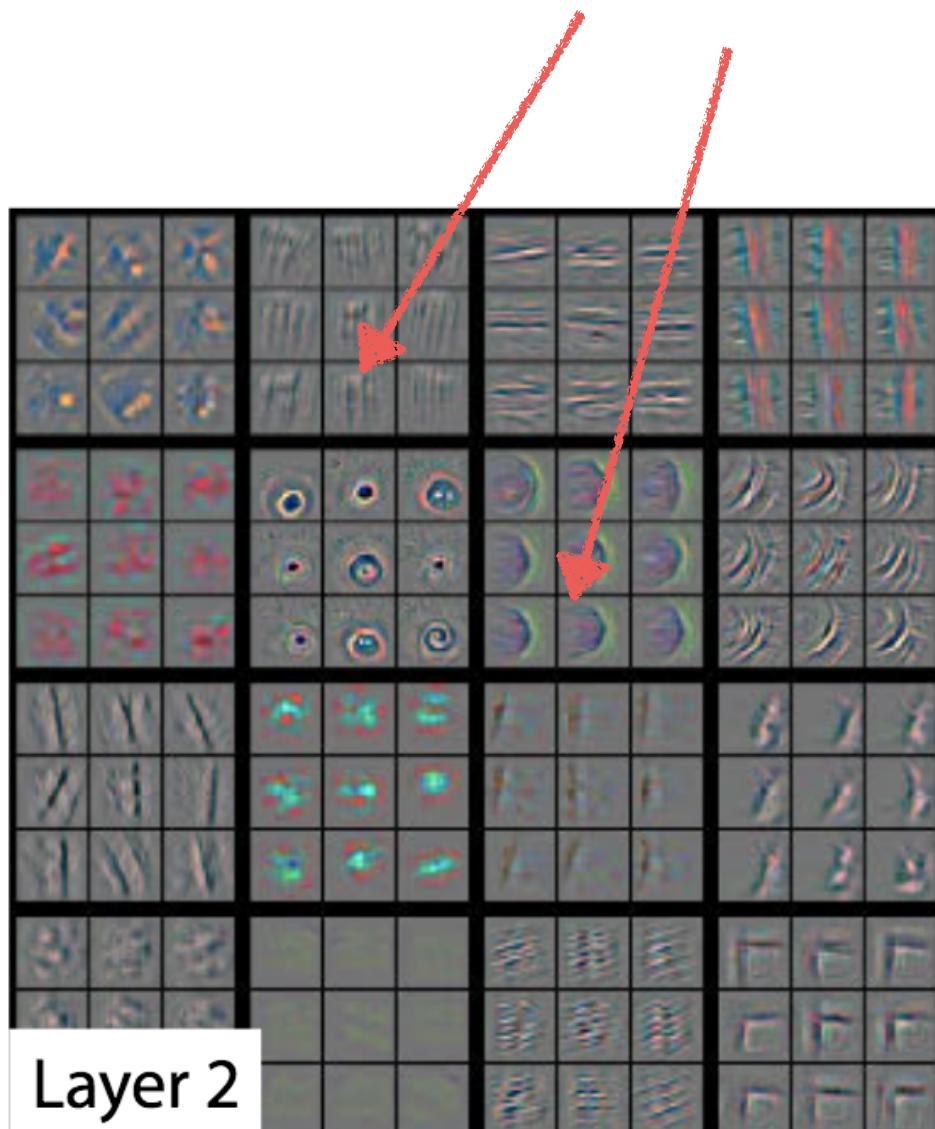


IT ALLOWS TO SEE
WHICH
REGIONS OF THE INPUT
GENERATED
A MAXIMUM RESPONSE
IN A NEURON

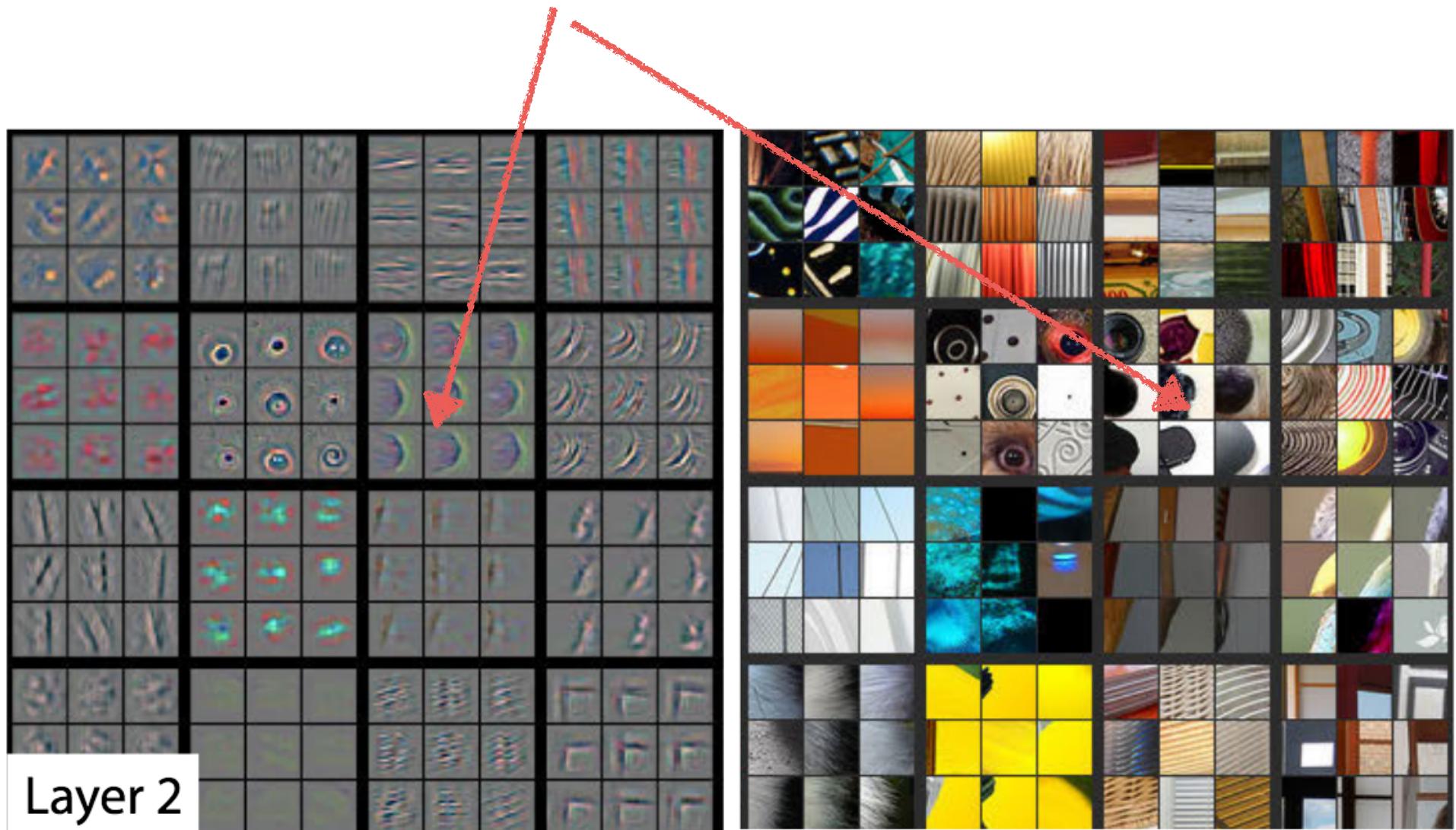


Zeiler+14

EVERY BLOCK OF 9 SHOWS
THE 9 STRONGEST RESPONSES TO A GIVEN FILTER OF LAYER2



THE CORRESPONDING REGIONS OF IMAGES THAT GENERATED THE MAXIMUM RESPONSE

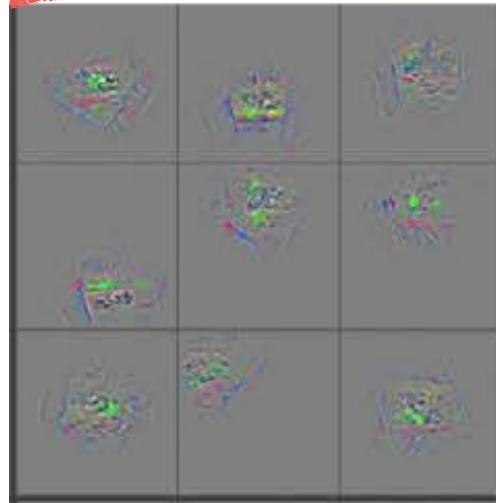


CAN BE
REPEATED
FOR DEEPER
LAYERS
ALTHOUGH IT
BECOMES LESS
INTUITIVE

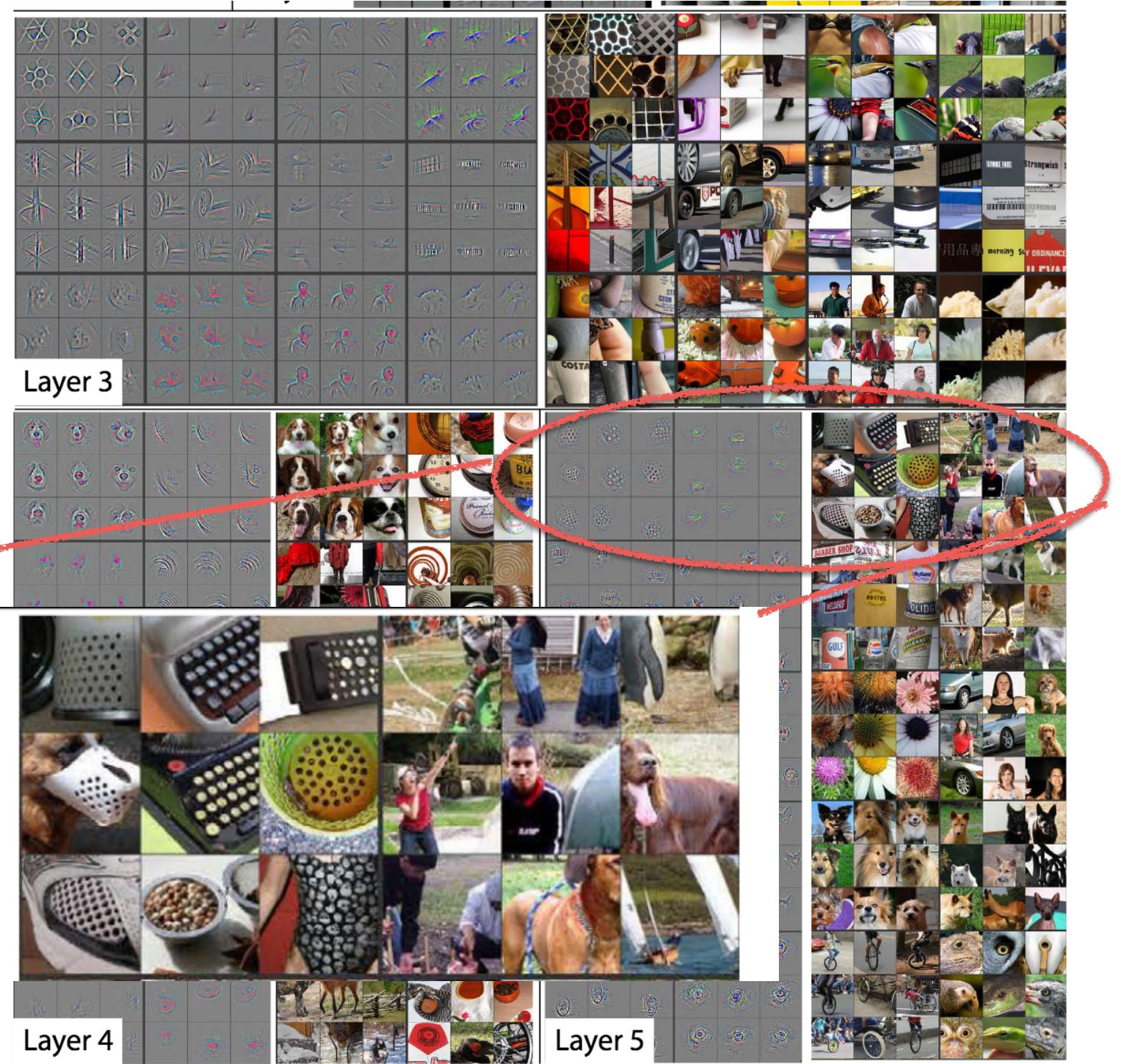
Zeiler+14



CAN BE
REPEATED
FOR DEEPER
LAYERS
ALTHOUGH IT
BECOMES LESS



Zeiler+14



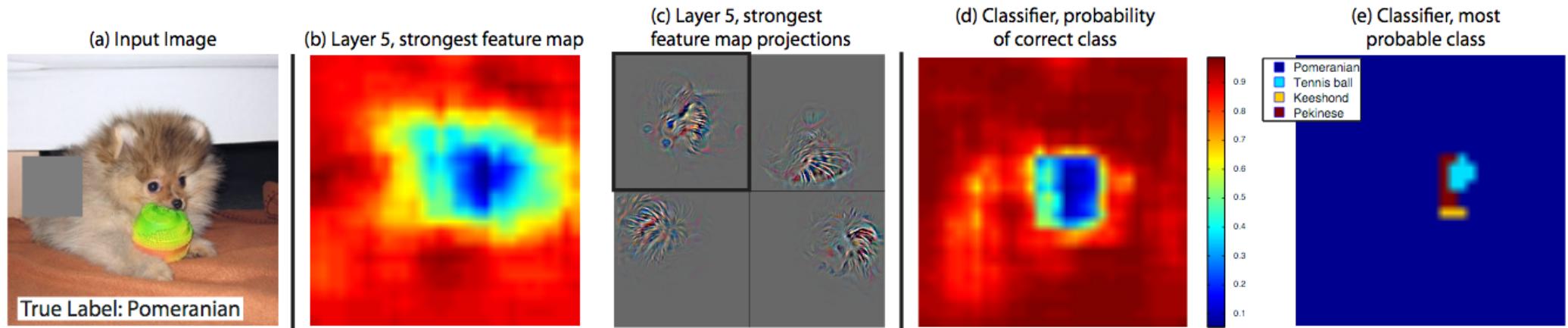
OCCLUSION SENSITIVITY

THE BASIC IDEA IS TO
PERTURB / MODIFY AN INPUT
IMAGE AND SEE THE EFFECT ON
THE PREDICTIONS

OCCLUSION SENSITIVITY TRIES ALSO TO FIND THE REGION OF THE IMAGE THAT TRIGGERED THE NETWORK DECISION BY MASKING DIFFERENT REGIONS OF THE INPUT IMAGE AND ANALYZING THE NETWORK OUTPUT

IT ALLOWS TO SEE IF THE NETWORK IS TAKING THE DECISIONS BASED ON THE EXPECTED FEATURES

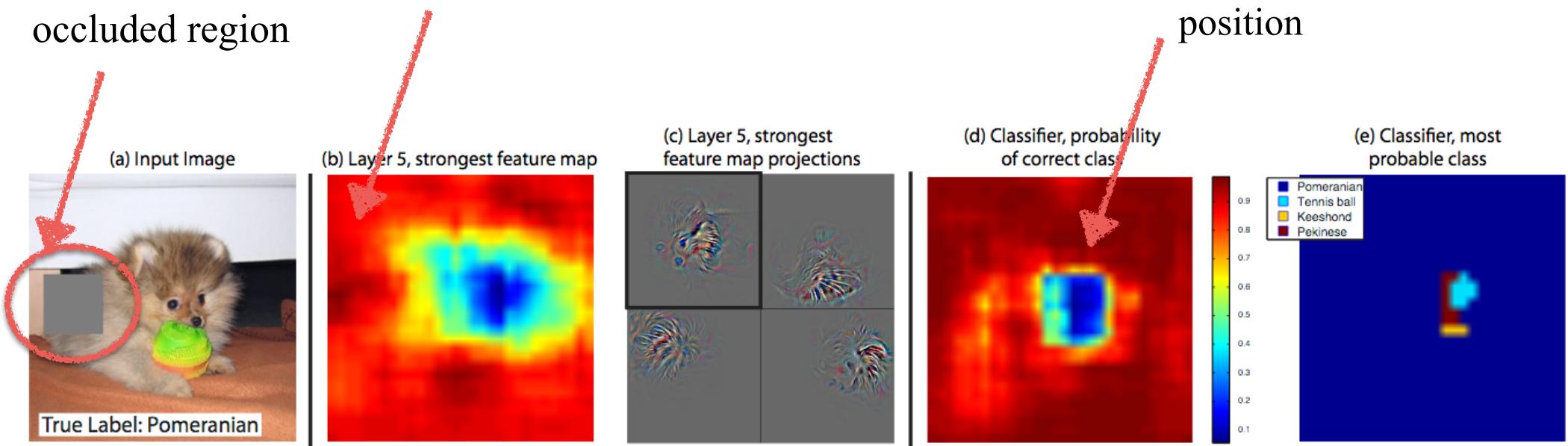
VERY TIME CONSUMING!



OCCLUSION SENSITIVITY TRIES ALSO TO FIND THE REGION OF THE IMAGE THAT TRIGGERED THE NETWORK DECISION BY MASKING DIFFERENT REGIONS OF THE INPUT IMAGE AND ANALYZING THE NETWORK OUTPUT

for every position
of the square the maximum response of a given layer
is averaged

occluded region



“INCEPTIONISM” TECHNIQUES

THE IDEA BEHIND INCEPTIONISM TECHNIQUES
IS TO INVERT THE NETWORK TO GENERATE AN IMAGE
THAT MAXIMIZES THE OUTPUT SCORE

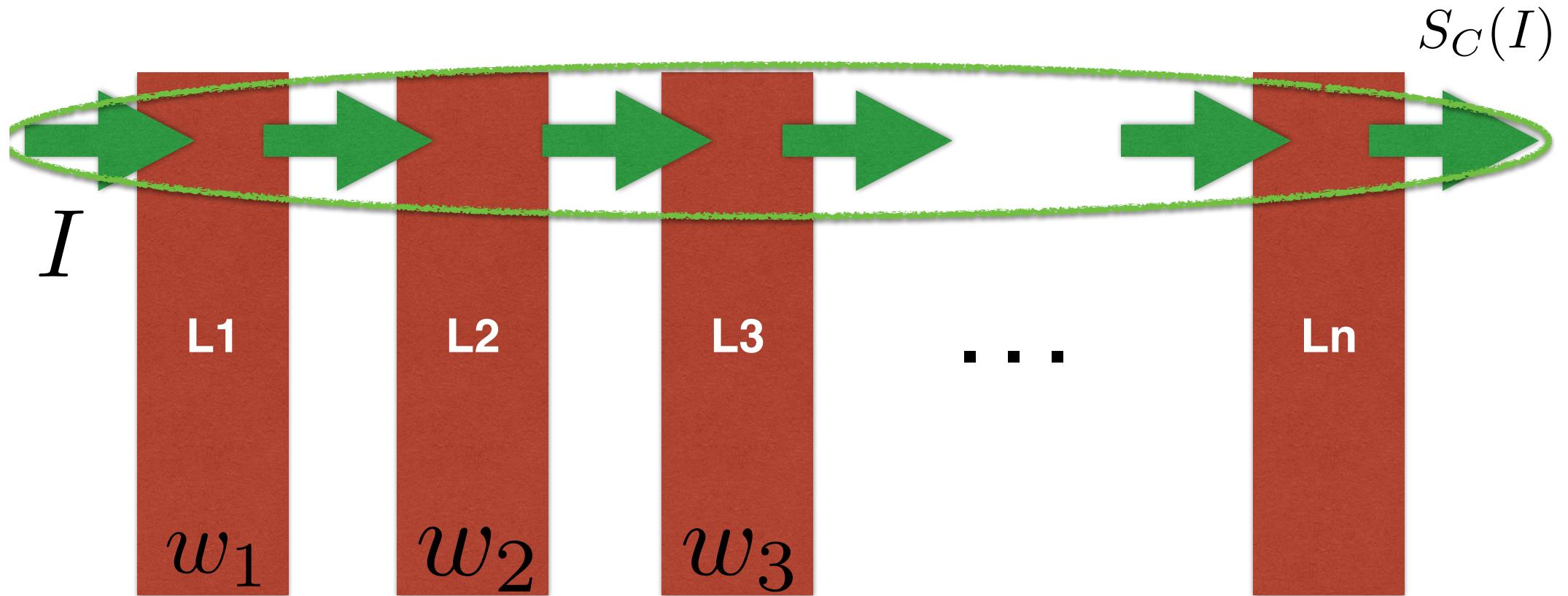
$$\arg \max_I S_c(I) - \lambda ||I||_2^2$$

Score of class c for image I

image I

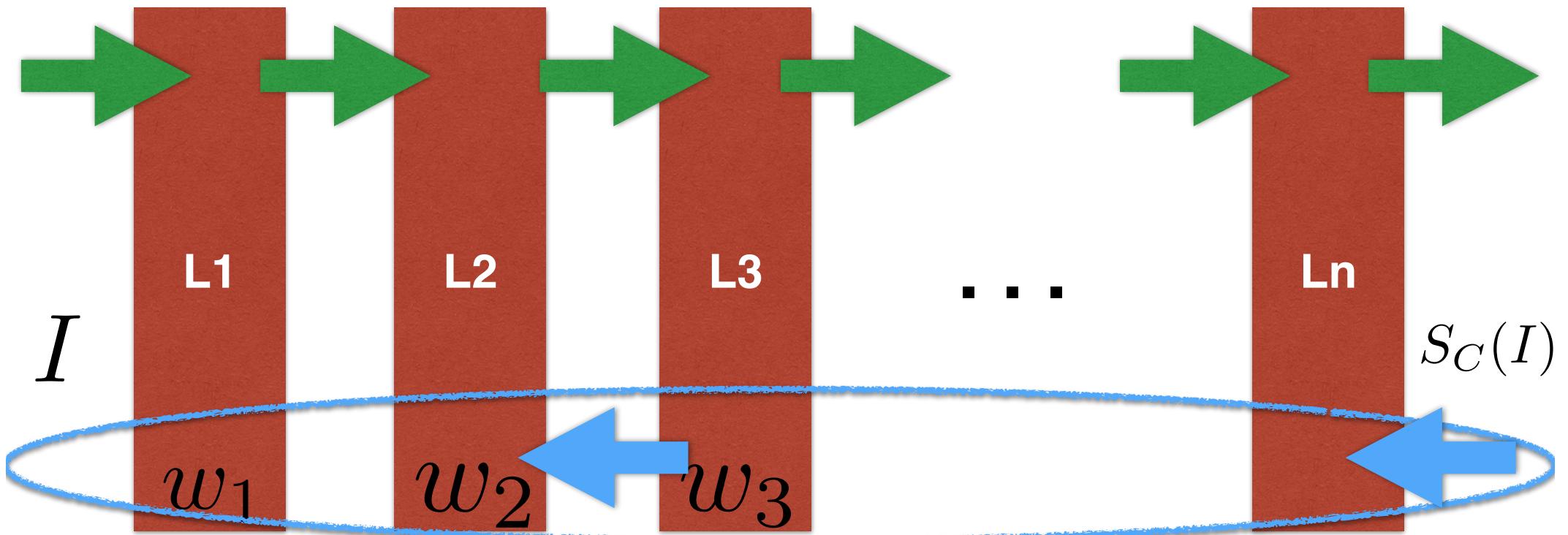
TRY TO FIND AN IMAGE THAT GENERATES A
HIGH SCORE FOR A GIVEN CLASS

INCEPTIONISM - DEEP DREAM



DURING THE TRAINING PHASE THE WEIGHTS ARE
LEARNED TO MAP I INTO S_C

INCEPTIONISM - DEEP DREAM



DURING THE RECONSTRUCTION PHASE, I IS LEARNT
THROUGH BACKPROPAGATION KEEPING THE WEIGHTS
FIXED

SEE SLIDES ON ACTIVATION ATLAS FOR AN EXAMPLE

GRADIENT BASED ATTRIBUTION

“what pixels in the image are responsible of this
classification?”

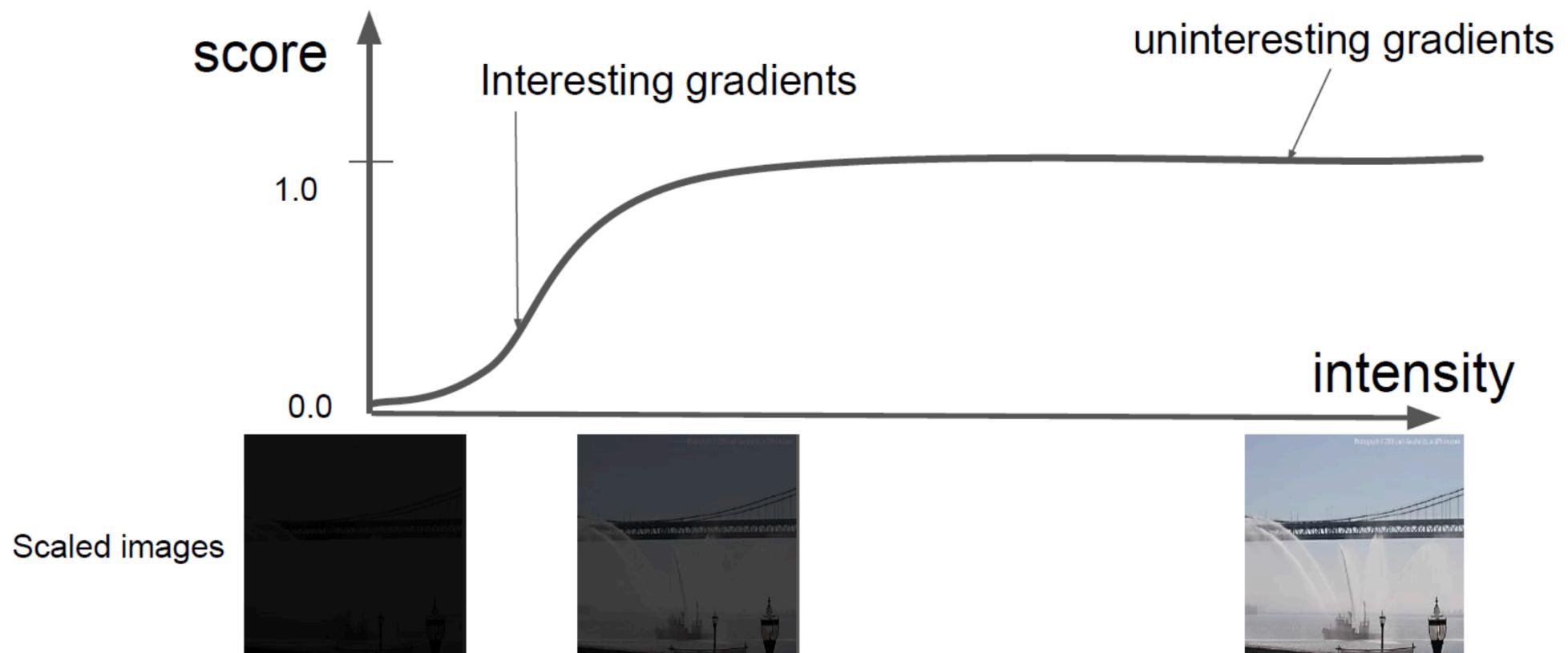
Attribute using gradient of the output w.r.t each input feature

$$x_i \frac{\delta F_w(x)}{\delta x_i}$$

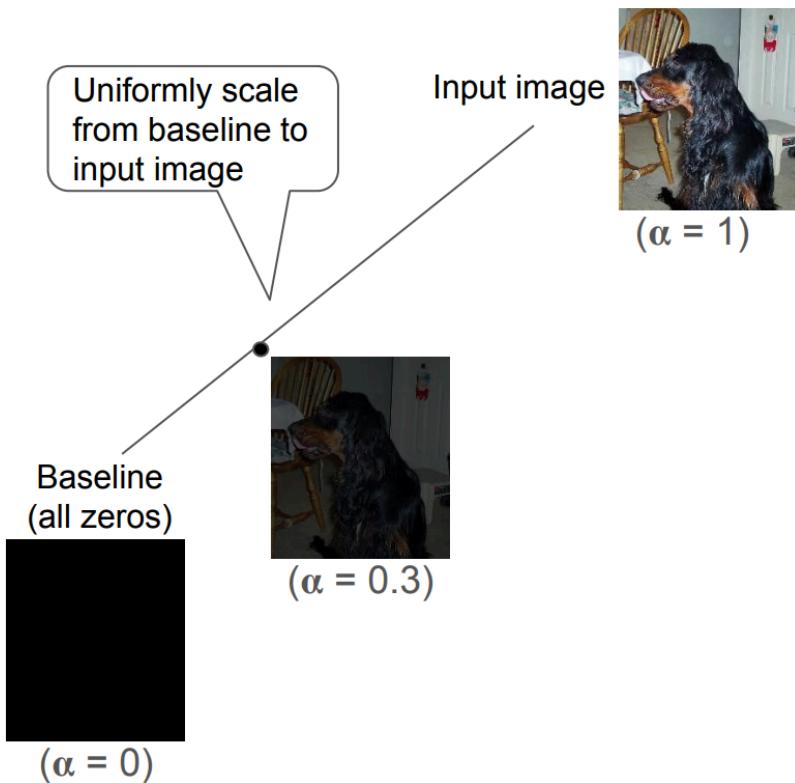
We try to find the importance of each feature (pixel) by computing the gradient w.r.t to that feature (Taylor approximation of the network function)

Does not work super well...





INTEGRATED GRADIENTS



Construct a sequence of images interpolating from a baseline (black) to the actual image

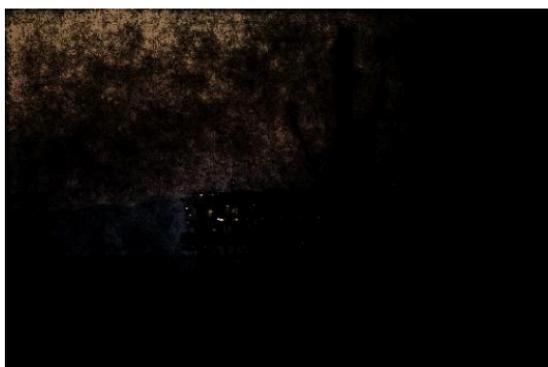
Average the gradients across these images

INTEGRATED GRADIENTS

Original image (Drilling platform)



Gradient at image



Integrated gradient



INTEGRATED GRADIENTS

Human label: [accordion](#)
Network's top label: [toaster](#)



[Integrated gradient](#)

