# Systematically Figure Out the Semantic of Components in Neural Networks

*Network Dissection* method for investigating NN units

Tianqin Li
Oct 24, 2020

# Overview

- *Network Dissection* method is developed by David Bau and Bolei Zhou at MIT Antonio Torralba's lab.
- Using image segmentation technique to investigate the causal connection between filters in CNN and human understandable visual concepts (like trees).
- Comparing to visual concept, it use a trained segmentation network to automatically find out which CNN filters are responsible for certain semantic concepts. -- wonder if one can use segmentation net to label VCs
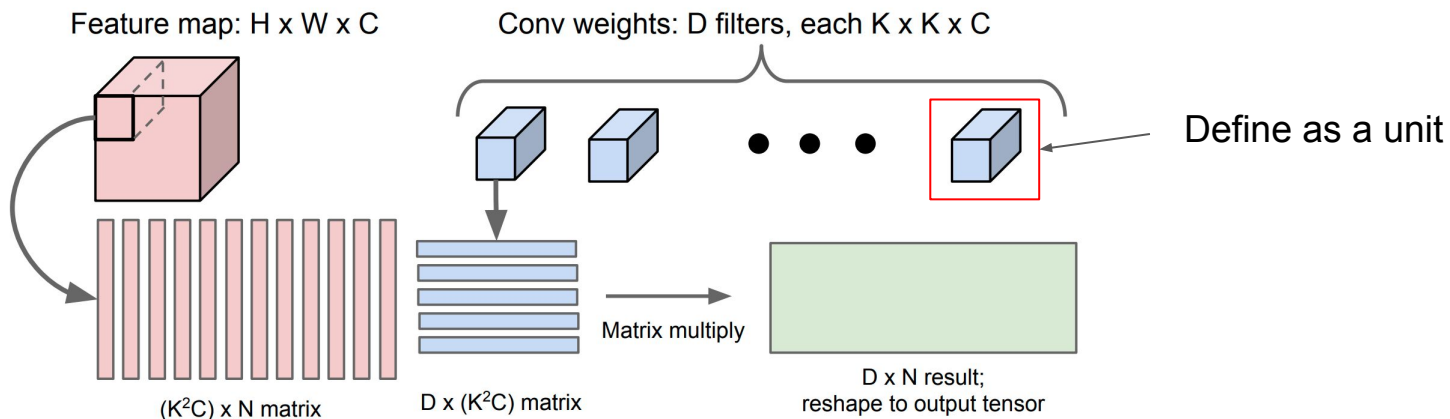


| VC1 | VC5 | VC9 | VC13 | VC17 | VC21 | VC25 | VC29 | VC33 | VC37 | VC41 |

Visual Concepts from Alan Yuille group

# Network Dissection

- Define units as a filter in CNN

## Implementing Convolutions: im2col

Feature map: H x W x C

Conv weights: D filters, each K x K x C

Define as a unit

$(K^2C)$ x N matrix

D x $(K^2C)$ matrix

Matrix multiply

D x N result;
reshape to output tensor

http://cs231n.stanford.edu/slides/2016/winter1516_lecture11.pdf

# Network Dissection

- Concepts labeled by a trained image segmentation labeling network



Picture (concept)

Lamp (concept)

Bed (concept)

Trained image segmentation network
MobileNetv2, ResNet, etc.

# Network Dissection

- Quantifying how each individual filter (*u*) influence the detection/construction of the concepts (denote as *c*)

$$s_c(\cdot)$$

$$u$$

Activation of u, for image x at location p (scalar)

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

$$x$$

$$a_u(x, p)$$



featuremap

thresholded

single unit u

upsample

$r_{u,P}$

$r_u^\uparrow > t$

agreement
IoU$_{u,c}$

generated image

segmentation

generate

segment

z

generator
G

x

$s_c(x)$

How well are the real segmentation matched with activation of the unit?

$$\mathrm{IoU}_{u,c} = \frac{\mathbb{P}_{x,p}\left[s_c(x, p) \wedge (a_u(x, p) > t_u)\right]}{\mathbb{P}_{x,p}\left[s_c(x, p) \vee (a_u(x, p) > t_u)\right]}.$$

$$t_u(threshold) \quad \max_t \mathbb{P}_{x,p}\left[a_u(x, p) > t\right] \geq 0.01.$$

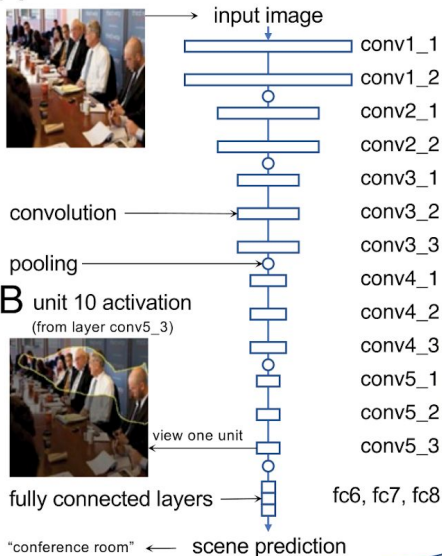# Example of one unit (u) to one concept (c)

- u is unit 150 in the last cnn layer of VGG-16 (conv5_3)
- c is airplane
- Unit 150 (u) prefer airplane concept (c)

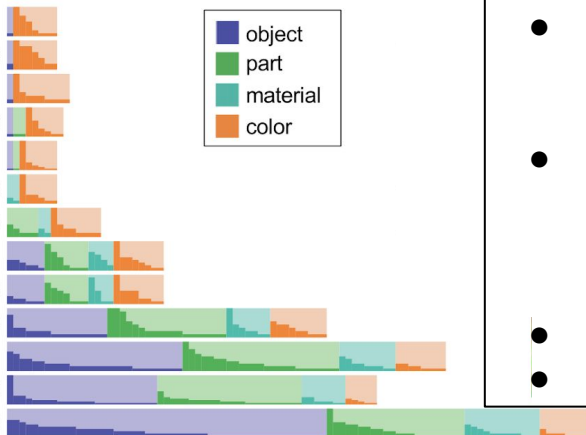**F** out-of-domain object detection test: conv5_3 unit 150 activation on airplanes

non-airplane imagenet images, mean=0.8
imagenet airplane images, mean=88.1

Density

Jitterplot

−50    0    50    100    150    200    250

unit 150 peak activation per image

non-airplane

airplane images

# Results - classification



A. VGG-16 architecture 224x224
input image
conv1_1
conv1_2
conv2_1
conv2_2
conv3_1
conv3_2
conv3_3
conv4_1
conv4_2
conv4_3
conv5_1
conv5_2
conv5_3

convolution

pooling

B. unit 10 activation
(from layer conv5_3)

view one unit

fully connected layers → fc6, fc7, fc8

"conference room" ← scene prediction

E. dissection of each convolutional layer
- object
- part
- material
- color

← matched visual concepts for all units →
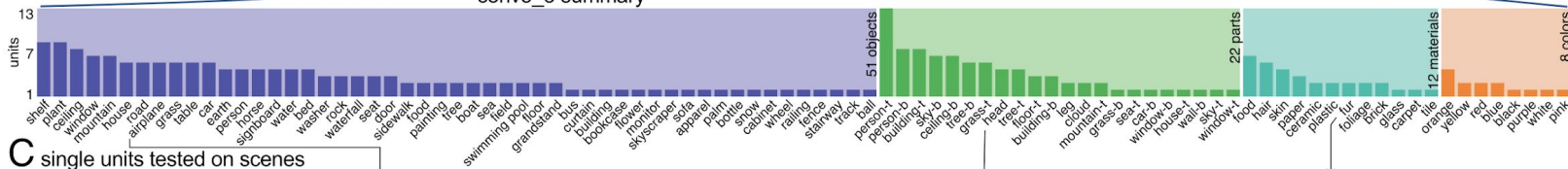
- Object detector units **automatically** emerge when train with larger scene classification
- Training with larger scene classification will results in units that match with **different levels** of concepts (objects, parts, materials, colors)
- Those units are not concepts exclusive
- Units that have IoU < 4% is excluded

D. conv5_3 summary

51 objects   22 parts   12 materials   8 colors

C. single units tested on scenes

unit 150 "airplane" (object)

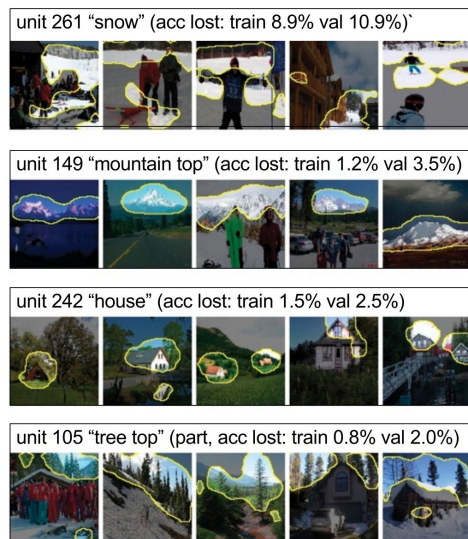unit 208 "person top" (part)

unit 141 "fur" (material)

Highlighted regions are those whose activation is among 1% quantile of total activation of that unit

# Results - how important is each units causally?

- Removing important units (to zero, ranked by IoU) **hurts** the classification badly
- Measure accuracy by balanced binary classification for individual classes
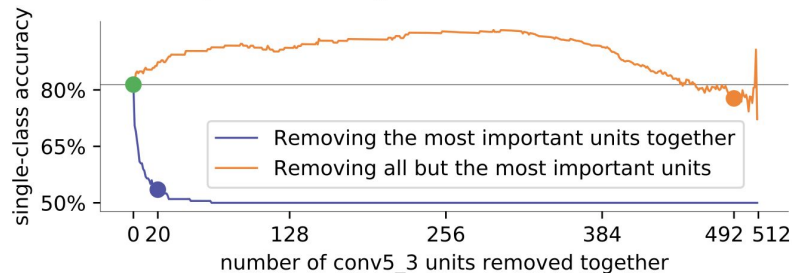- Removing redundant units **even help improve** the accuracy

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%)

B Validation accuracy when units removed as a set

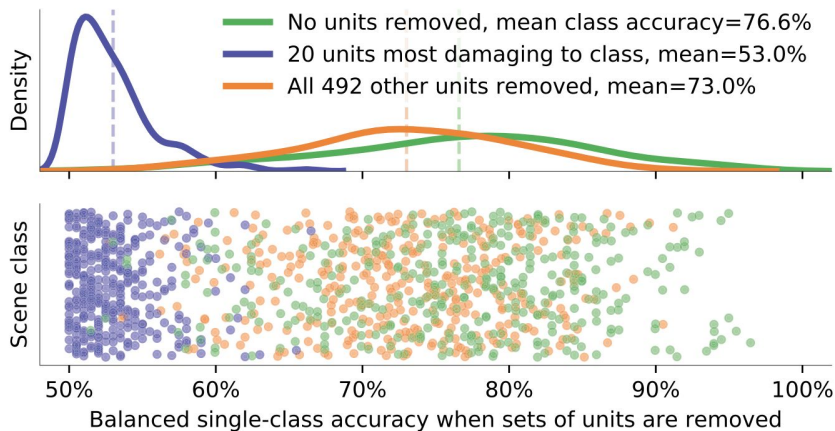| | Balanced single-class 'ski resort' accuracy | All-class accuracy |
|---|---|---|
| Unchanged vgg-16: | 81.4% | 53.3% |
| 4 most important units removed: | 64.0% | 53.2% |
| 20 most important units removed: | 53.5% | 52.6% |
| 492 least important units removed: | 77.7% | 2.1% |
| Chance level | 50.0% | 0.27% |

C "Ski resort" accuracy when removing sets of units of different sizes



Removing the most important units together
Removing all but the most important units

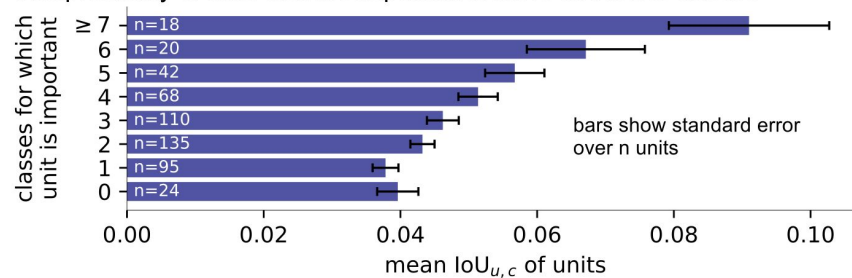number of conv5_3 units removed together

# Results - how important is each units causally?

- Removing 20 most important units hurt most whereas others barely impact acc
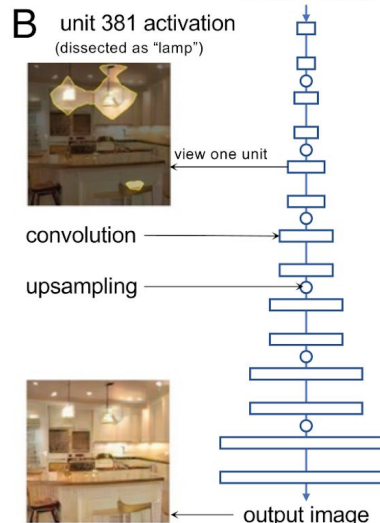- Units shared by multiple classes has higher IoU (more interpretable)



D Removing 20 most- and 492 least-important units for all scene classes

- No units removed, mean class accuracy=76.6%
- 20 units most damaging to class, mean=53.0%
- All 492 other units removed, mean=73.0%

Balanced single-class accuracy when sets of units are removed

E Interpretability of units that are important to more and fewer classes

bars show standard error over n units

mean IoU$_{u,c}$ of units

# Results - GAN



A  Progressive GAN architecture
random vector

B  unit 381 activation
(dissected as "lamp")

view one unit

convolution

upsampling

output image

C  dissection of each convolutional layer

layer1
layer2
layer3
layer4
layer5
layer6
layer7
layer8
layer9
layer10
layer11
layer12
layer13
layer14

- More parts responsible units comparing to Classification (most encode objects)
- But still Different levels of concepts emerges in the hidden units of GAN
- Count units important to a concept (c) if IoU > 4%
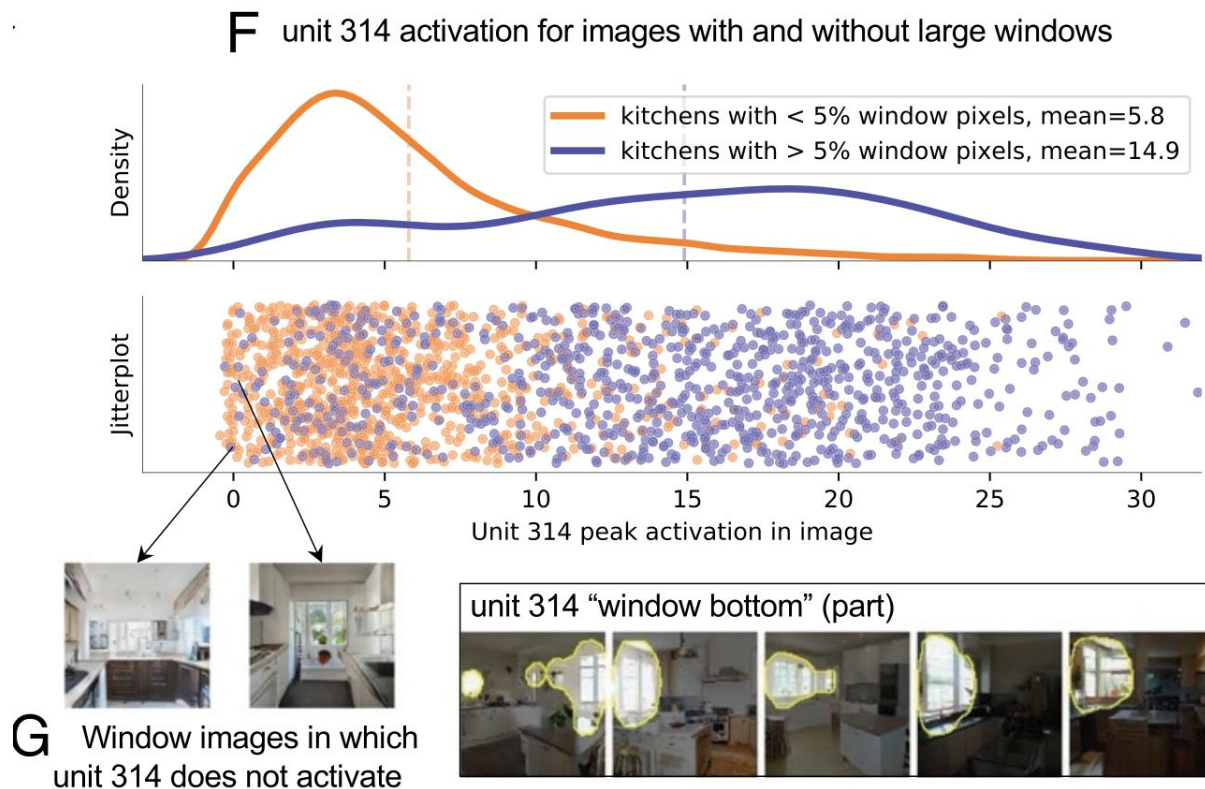- Output resolution 256x256, 15 cnn layers total

D  layer5 detail

14 objects
33 parts
1 material
6 colors

units
18
9
1

oven, chair, window, door, person, exhaust hood, microwave, refrigerator, kitchen island, ceiling, lamp, stove, floor, work surface, ceiling-b, window-t, window-b, ceiling-t, floor-t, floor-b, floor-r, work surface-b, kitchen island-r, cabinet-b, window-l, kitchen island-t, work surface-t, top, refrigerator-t, refrigerator-b, floor-b, stove-t, chair-b, ceiling-l, kitchen island-t, kitchen island-b, stove-l, chair-l, chair-r, window-r, floor-l, person-b, stove-b, ceiling-r, cabinet-t, chair-t, chair-r, work surface-r, plastic, orange, white, blue, red, green, black

legend:
- object
- part
- material
- color

E  single units

unit 498 "oven" (object)

unit 244 "chair top" (part)

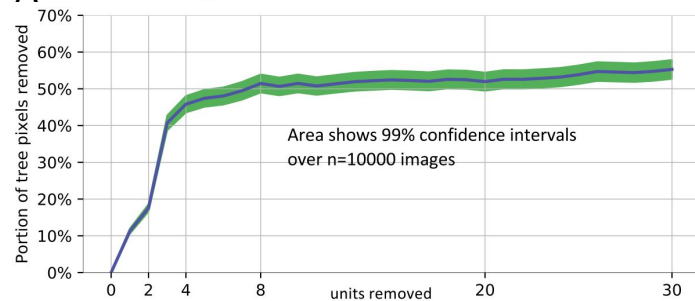unit 58 "red" (color)

# Results - GAN

- Concepts (c) related units (u) in GAN
- Using window specific unit activation can classify if generated image has window (78.2% Accuracy)
- But there are counter-examples in Figure G



F unit 314 activation for images with and without large windows

kitchens with < 5% window pixels, mean=5.8
kitchens with > 5% window pixels, mean=14.9

Density

Jitterplot

Unit 314 peak activation in image

G Window images in which unit 314 does not activate

unit 314 "window bottom" (part)

# Results - Causal role of units in GAN (remove)

- Removing tree specific units (layer 4) results in tree removal in the generated image
- Tree pixels are identified by segmentation network
- Remove tree units leave the whole image intact
- ! Remove the tree units even reveal the church which were occluded before -> suggesting the network compute compositional structures.



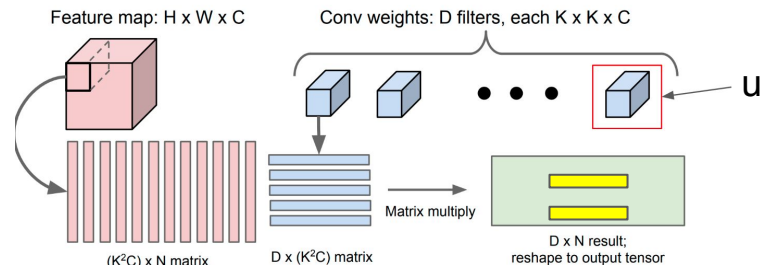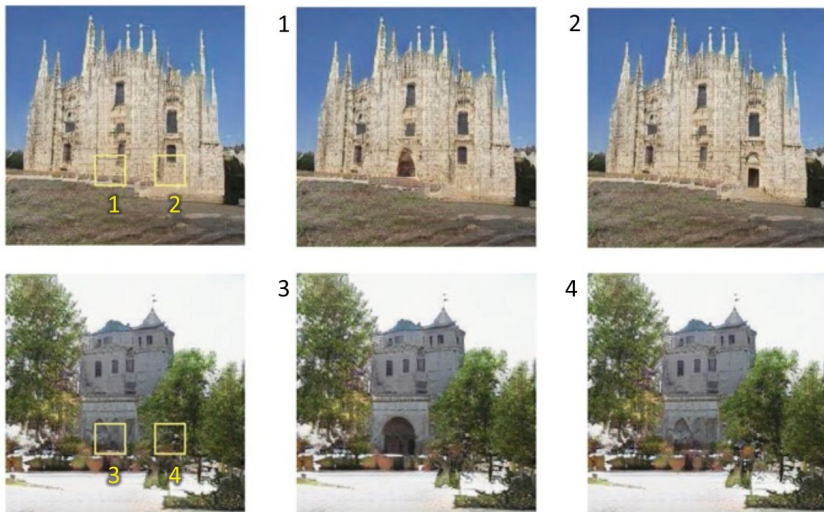A Causal effect on generation of trees when "tree" units removed

Portion of tree pixels removed

Area shows 99% confidence intervals over n=10000 images

units removed

B

unchanged    2 units    4 units    8 units    20 units

Number of units removed (units ranked by IoU match with tree segmentations)
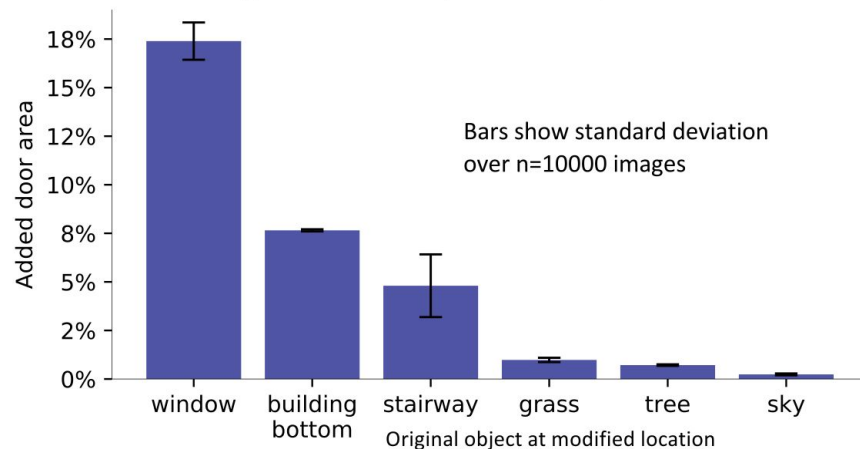
# Results - Causal role of units in GAN (activation)

- Activating "door" units at certain location
- Adding depends on the context of the location



Feature map: H x W x C

Conv weights: D filters, each K x K x C

u

(K²C) x N matrix

D x (K²C) matrix

Matrix multiply

D x N result;
reshape to output tensor

C  Effect of activating "door" units depends on location



1

2

3

4

D  Effect of activating "door" units depends on object context



Bars show standard deviation
over n=10000 images

Added door area

Original object at modified location

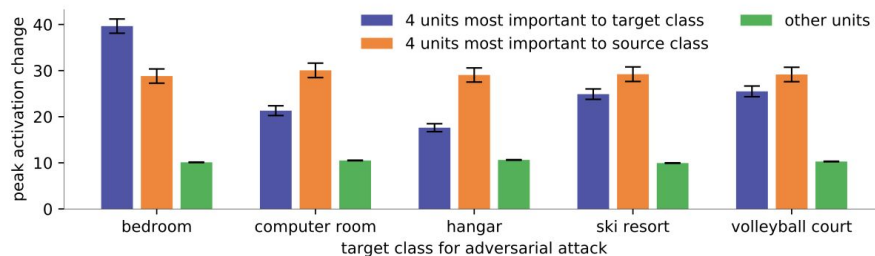window | building bottom | stairway | grass | tree | sky

# Application: Analyzing Adversarial Attack

- Adversarial attack diminishes the firing of important units for original class
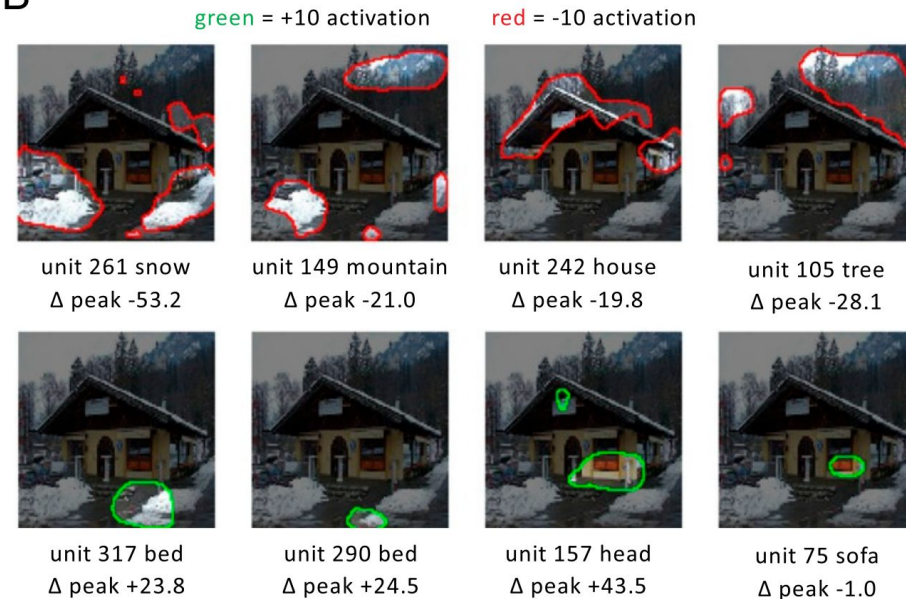- However, increase the firing of the important units for the target class

A  Adversarial attack changes an image imperceptibly to fool the classifier

original                    delta                    attacked

$+ 0.01 \times$            =

'ski resort'                                        'bedroom'

C  Mean peak unit activation change when attacked, for units in conv5_3



B  Activation change in 4 units most important to ski resort and bedroom.
green = +10 activation         red = -10 activation



unit 261 snow
Δ peak -53.2

unit 149 mountain
Δ peak -21.0

unit 242 house
Δ peak -19.8

unit 105 tree
Δ peak -28.1

unit 317 bed
Δ peak +23.8

unit 290 bed
Δ peak +24.5

unit 157 head
Δ peak +43.5

unit 75 sofa
Δ peak -1.0

# Application: GAN concept painting

- Activate the important 20 units of selected concept at certain location results in painting the concept.