

DEEP LEARNING FOR ARTIFICIAL INTELLIGENCE

3rd Master Course UPC ETSETB TelecomBCN Barcelona. Autumn 2019.



Instructors



Xavier
Giró-i-Nieto



Marta R.
Costa-jussà



Noé
Casas



Verónica
Vilaplana



Ramon
Morros



Javier
Ruiz



Albert
Pumarola



Jordi
Torres

Organizers



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH



Supporters



GitHub Education

+ info: <http://bit.ly/dlai2019>

[\[course site\]](#)



#DLUPC

Day 7 Lecture 1

Interpretability



Xavier Giro-i-Nieto
xavier.giro@upc.edu



Associate Professor
Universitat Politècnica de Catalunya
Technical University of Catalonia



Acknowledgements



Amaia Salvador
amaia.salvador@upc.edu

PhD Candidate
Universitat Politècnica de Catalunya



UNIVERSITAT POLITÈCNICA
DE CATALUNYA



Eva Mohedano
eva.mohedano@insight-centre.org

Postdoctoral Researcher
Insight-centre for Data Analytics
Dublin City University



Videolectures

YouTube video inside
DEEP LEARNING FOR COMPUTER VISION
Summer Seminar UPC TelecomBCN, 4 - 8 July 2016

Instructors:

- Hector Larochelle
- Amaia Salvador
- David Sánchez
- Joel Zitnick
- Eva Mohedano
- David Maturana

Organizers:

- Universitat Politècnica de Catalunya
- DCU
- Insight
- EU

+ info: TelecomBCN.DeepLearning.Barcelona
[course site]

Day 2 Lecture 3

Visualization



Amaia Salvador



UNIVERSITAT POLITÈCNICA DE CATALUNYA
SANT CUGAT DEL VALLES
Department of Signal Theory
and Communications
Image Processing Group

[UPC DLCV 2016]

YouTube video inside
DEEP LEARNING FOR COMPUTER VISION
Summer School at UPC TelecomBCN Barcelonès, July 20-21, 2018

Instructors:

- Hector Larochelle
- Amaia Salvador
- David Sánchez
- Joel Zitnick
- Eva Mohedano
- David Maturana

Organised by:

- Universitat Politècnica de Catalunya
- DCU
- vilynx.
- Google Cloud Platform

Supported by:

- EU
- Insight
- Google Cloud Platform

+ info: <http://bit.ly/dlcv2018>

<http://bit.ly/dlcv2018>

Eva Mohedano
eva.mohedano@insight-centre.org

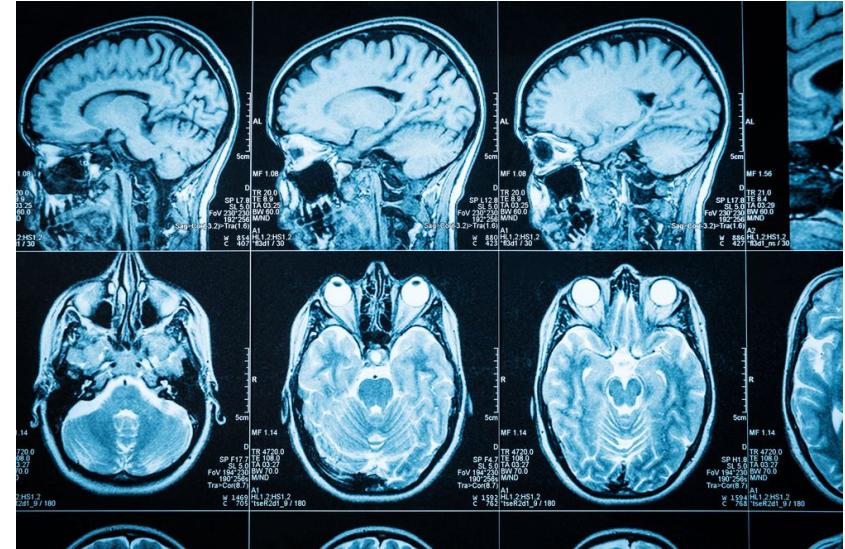
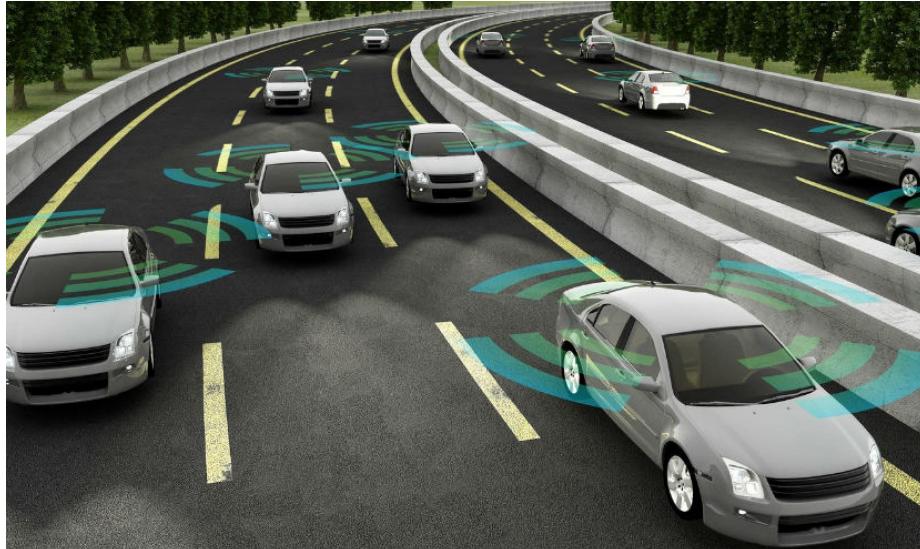
Postdoctoral Researcher
Insight-centre for Data Analytics
Dublin City University

Day 3 Lecture 4

Interpretability

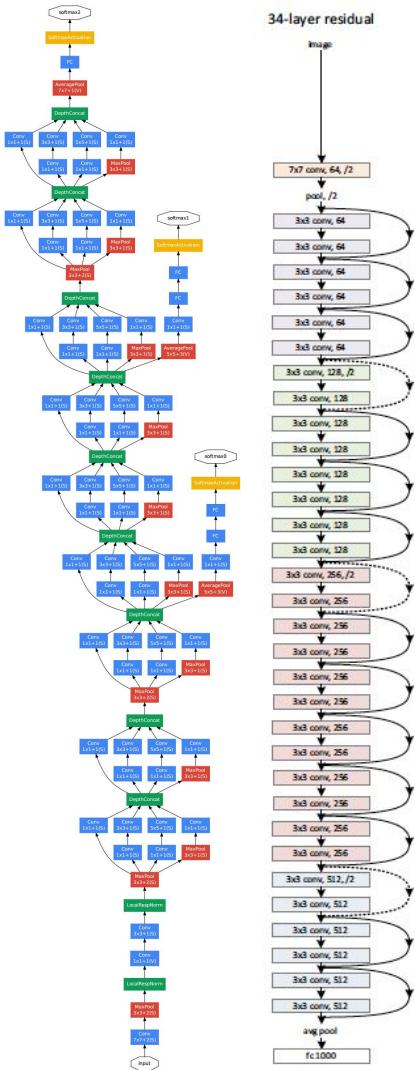
Motivation

In many cases, an explanation for NN prediction is required.

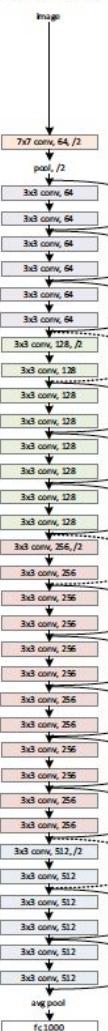


AlexNet

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
conv-256
maxpool
conv-512
conv-512
conv-512
maxpool
FC-4096
FC-4096
FC-1000
maxpool
FC-4096
FC-4096
FC-1000
softmax

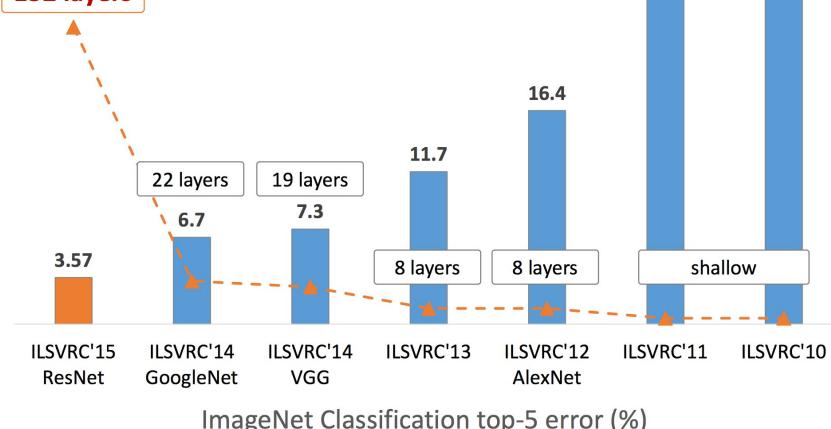


34-layer residual



Revolution of Depth

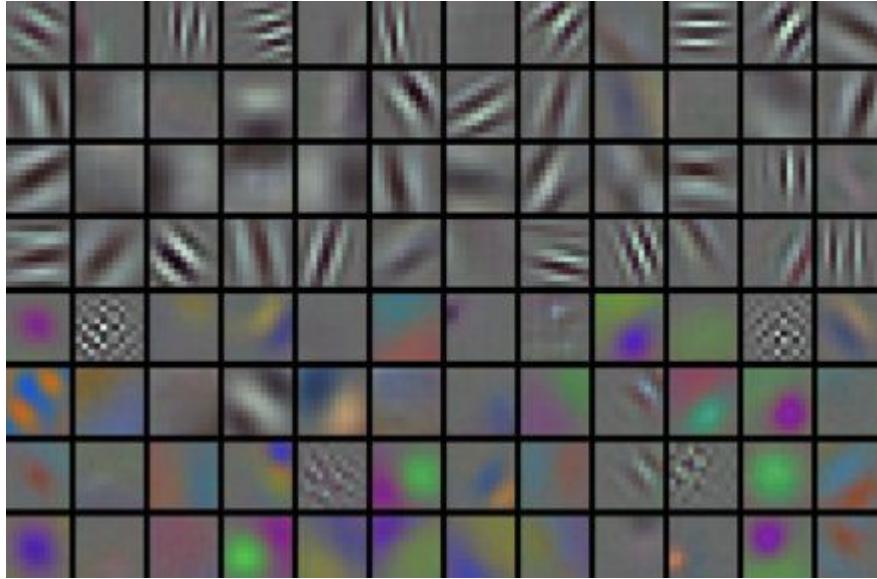
152 layers



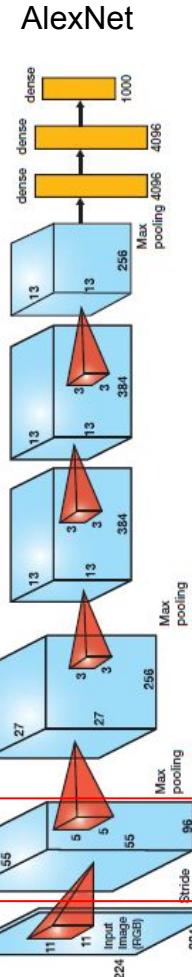
Interpretability

- **Visualization**
 - **Learned Weights**
 - Feature Maps
- Attribution
- Feature visualization

Visualization of Learned Weights



Only convolutional filters from the first layer can be “visualized”, because their depth of 3 matches the input RGB channels.



Visualization of Learned Weights

Filters with depth larger than 3 can be visualized showing a gray-scale image of each depth, one next to the other.

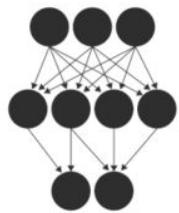
Weights:



2D Convolutional filters learned in Layer #2 from [ConvnetJS](#)

Visualization of Learned Weights

Demo: Classify MNIST digits with a Convolutional Neural Network



ConvNetJS

Deep Learning in your browser



"ConvNetJS is a Javascript library for training Deep Learning models (mainly Neural Networks) entirely in your browser. Open a tab and you're training. No software requirements, no compilers, no installations, no GPUs, no sweat."



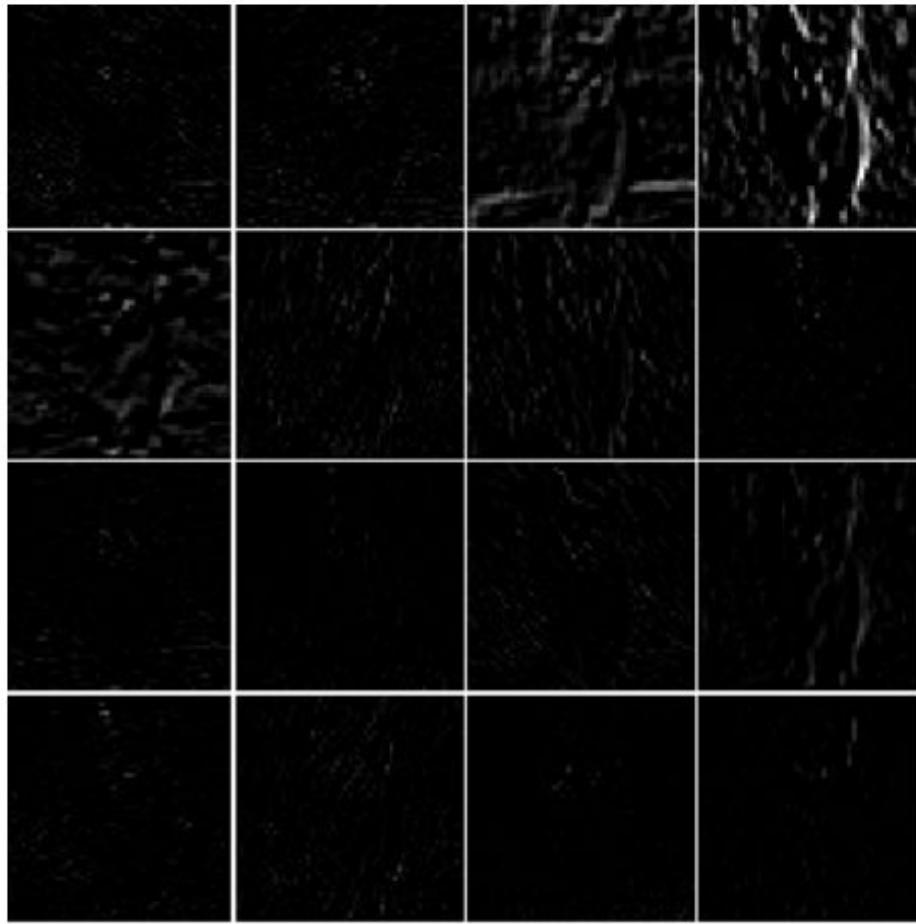
Interpretability

- **Visualization**
 - Learned Weights
 - **Feature Maps**
- Attribution
- Feature visualization

Visualization of Feature Activations (2D)



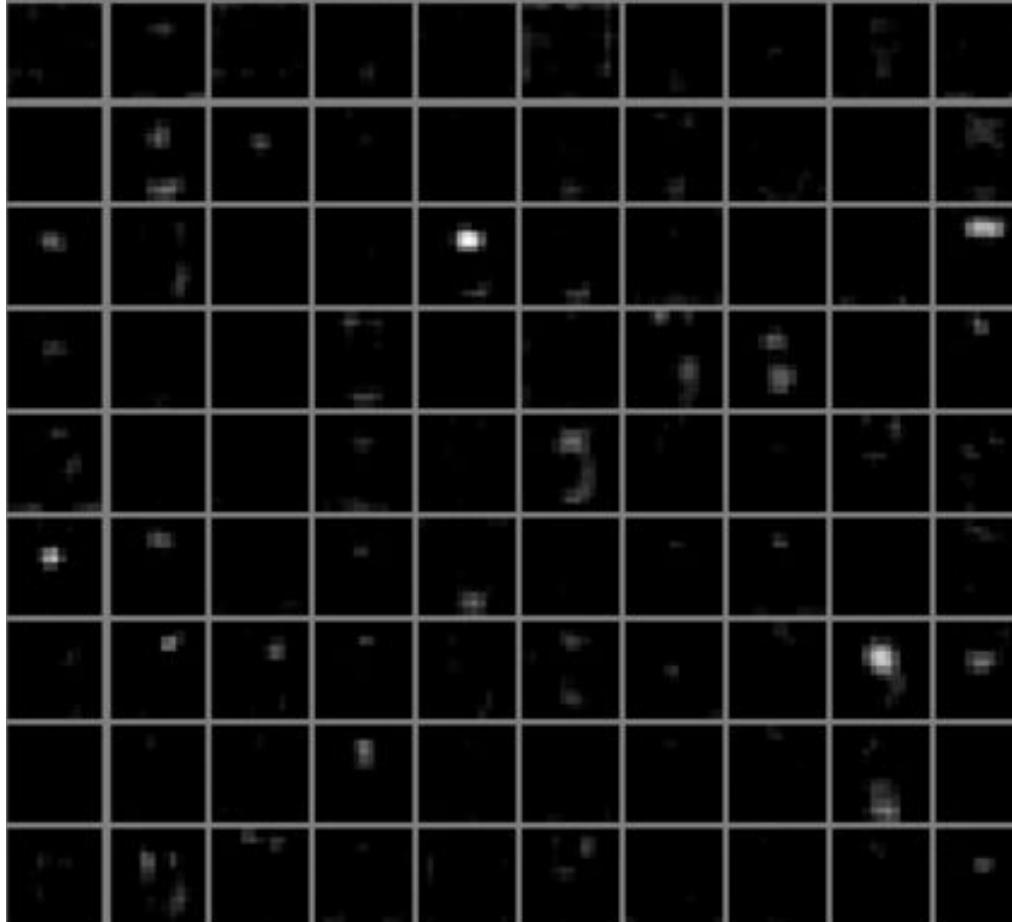
Input
image



Visualization of Feature Activations (2D)



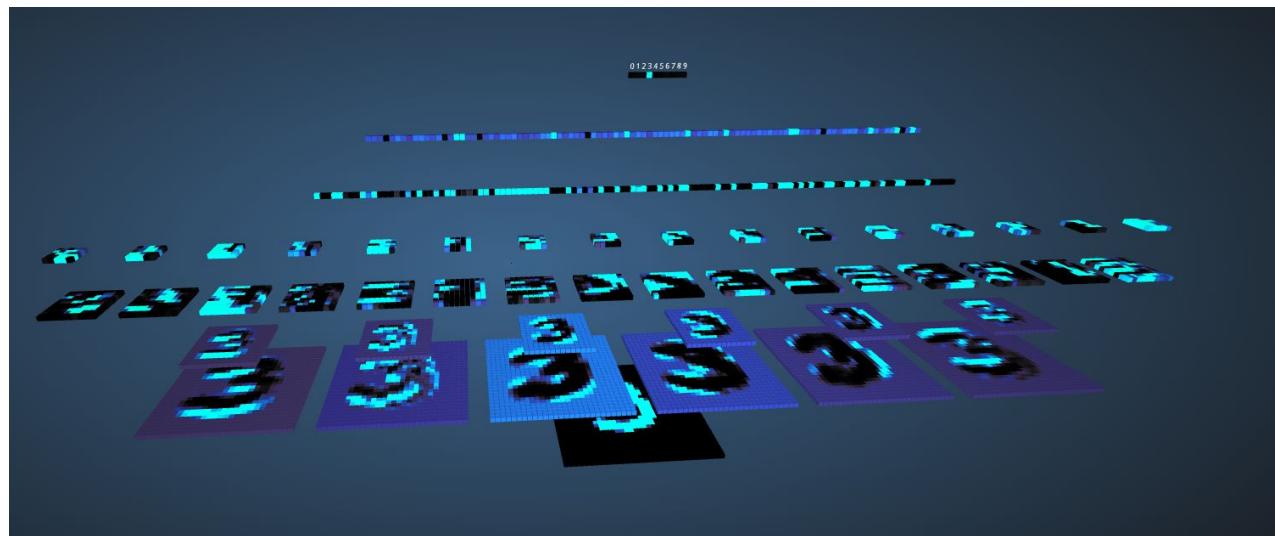
Input
image



conv5

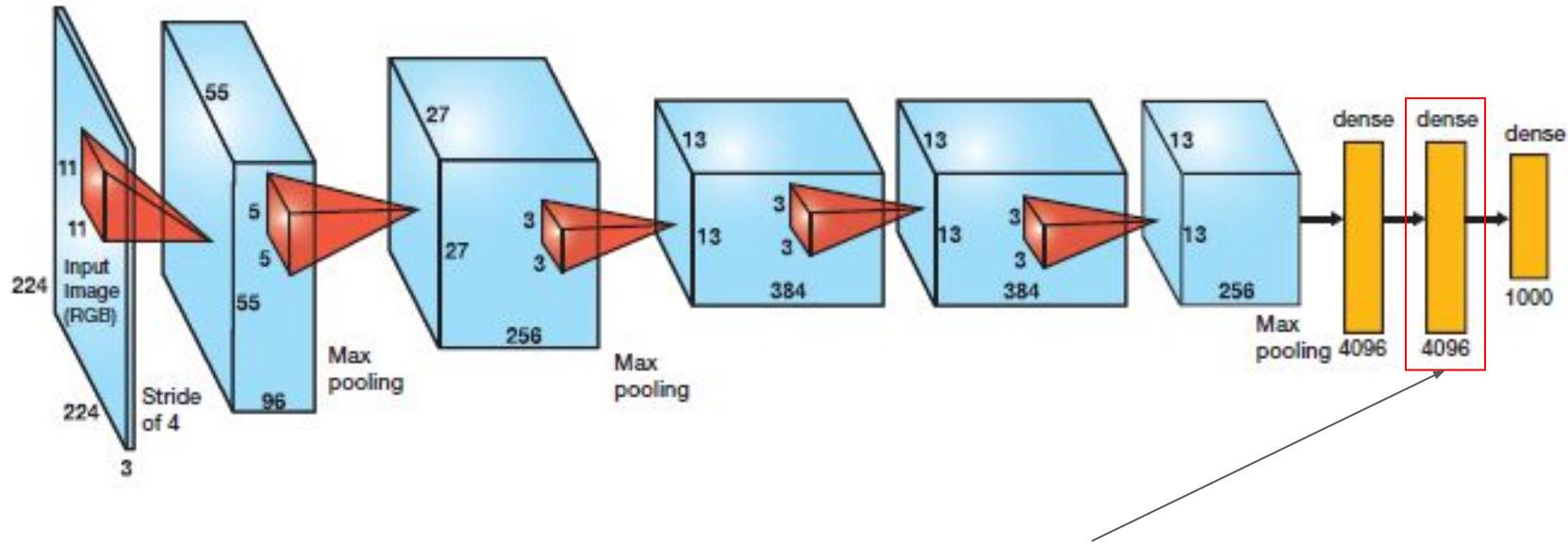
Visualization of Feature Activations (2D)

Demo: 3D Visualization of a Convolutional Neural Network



Harley, Adam W. ["An Interactive Node-Link Visualization of Convolutional Neural Networks."](#) In Advances in Visual Computing, 13 pp. 867-877. Springer International Publishing, 2015.

Visualization of Feature Activations (any dimension)



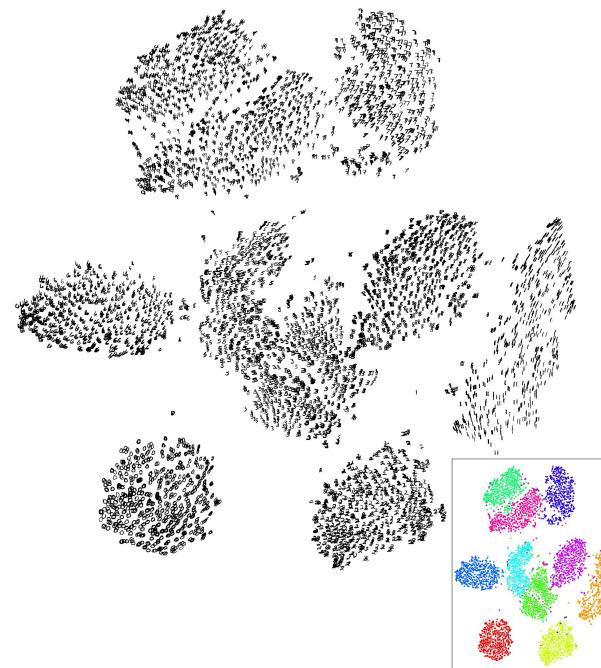
Can be used with features from layer before classification

Visualization of Feature Activations (any dimension)

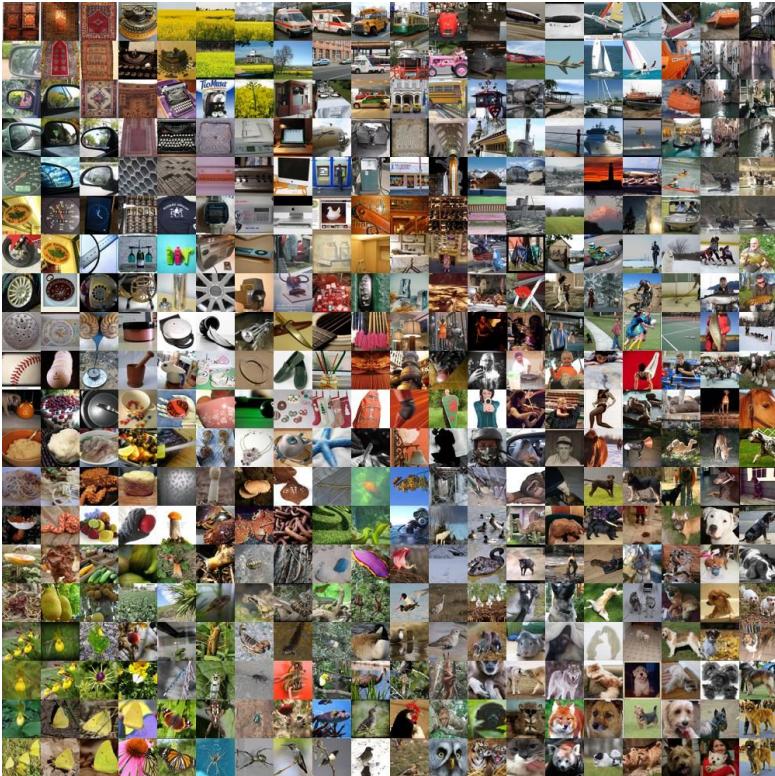
t-SNE:

Embeds high dimensional data points (i.e. feature maps) so that pairwise distances are preserved in a local 2D neighborhoods.

Example: 10 classes from MNIST dataset



Visualization of Feature Activations (any dimension)



t-SNE on fc7 features from AlexNet.

Source: [Andrey Karpathy](#) (Stanford University)

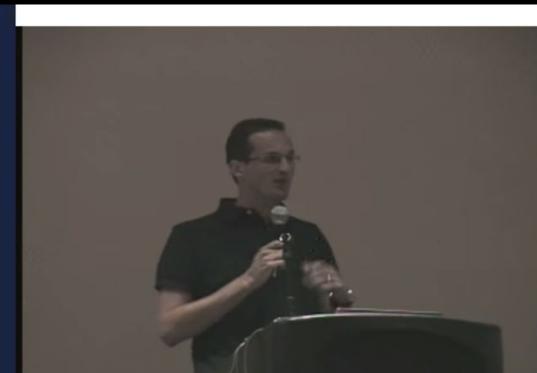
.\\|

Do's and Don'ts of using t-SNE to Understand Vision Models

Laurens van der Maaten

Interpretable Machine Learning for Computer Vision Workshop
June 18th, 2018

facebook
Artificial Intelligence Research



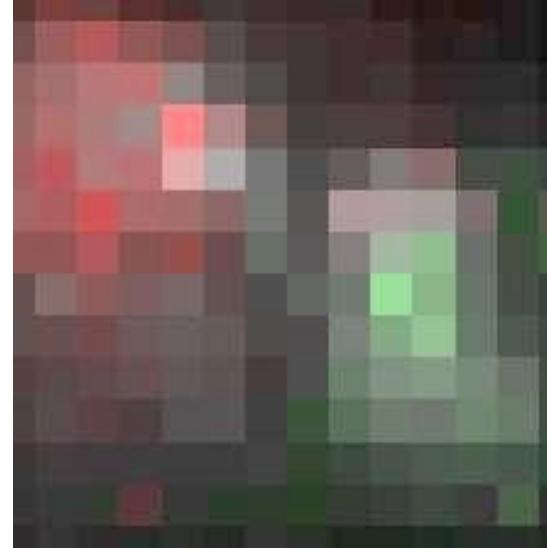
GvF

Interpretability: Attribution

- Visualization
- **Attribution**
 - Receptive Field of the Highest Activations
 - Activation Changes after Occlusions
 - Class Activation Maps (CAMs)
 - Gradient-based
- Feature visualization

Attribution

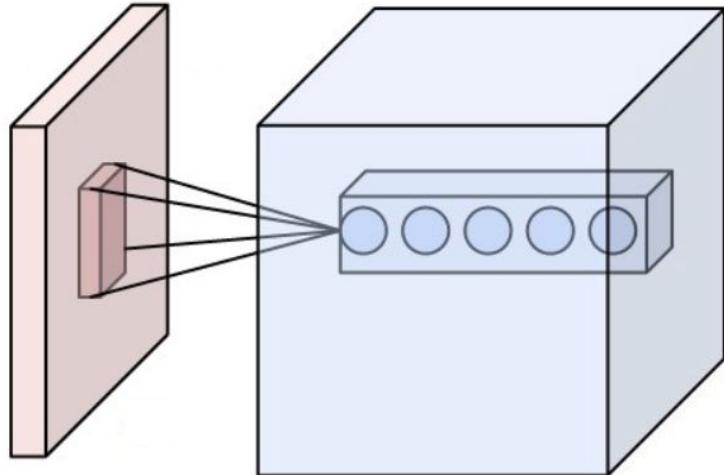
Attribution studies what part of an example is responsible for the network activating a particular way.



Interpretability: Attribution

- Visualization
- **Attribution**
 - **Receptive Field of the Highest Activations**
 - Activation Changes after Occlusions
 - Class Activation Maps (CAMs)
 - Gradient-based
- Feature visualization

Reminder: Receptive Field



Receptive field: Part of the input that is visible to a neuron. It increases as we stack more convolutional layers (i.e. neurons in deeper layers have larger receptive fields).

Receptive Field of Highest Activations

Visualize the receptive field of a neuron on those images that activate it the most

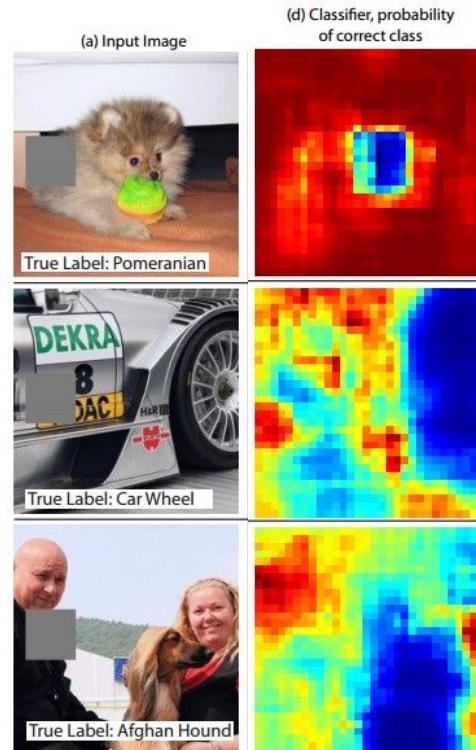


Interpretability: Attribution

- Visualization
- **Attribution**
 - Receptive Field of the Highest Activations
 - **Activation Changes after Occlusions**
 - Class Activation Maps (CAMs)
 - Gradient-based
- Feature visualization

Activation Changes after Occlusions (global scale)

1. Iteratively forward the same image through the network, occluding a different region at a time.
2. Keep track of the probability of the correct class w.r.t. the position of the occluder



Activation Changes after Occlusions (global scale)



GT: negative



GT: positive



GT: negative



GT: positive



GT: negative

Activation Changes after Occlusions (global scale)

The changes in activations can be observed in any layer. This allowed identifying some filters as weak object detectors, trained with image labels.

Buildings

56) building



120) arcade



8) bridge



123) building



Indoor objects

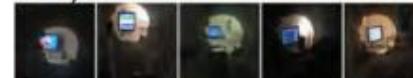
182) food



46) painting



106) screen



53) staircase

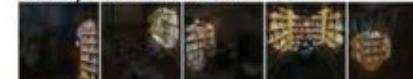


Furniture

18) billiard table



155) bookcase



116) bed



38) cabinet

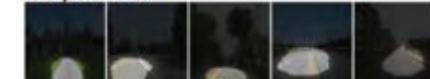


Outdoor objects

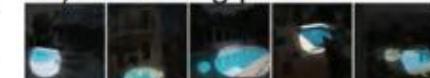
87) car



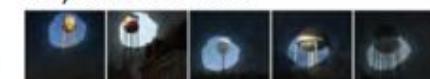
61) road



96) swimming pool

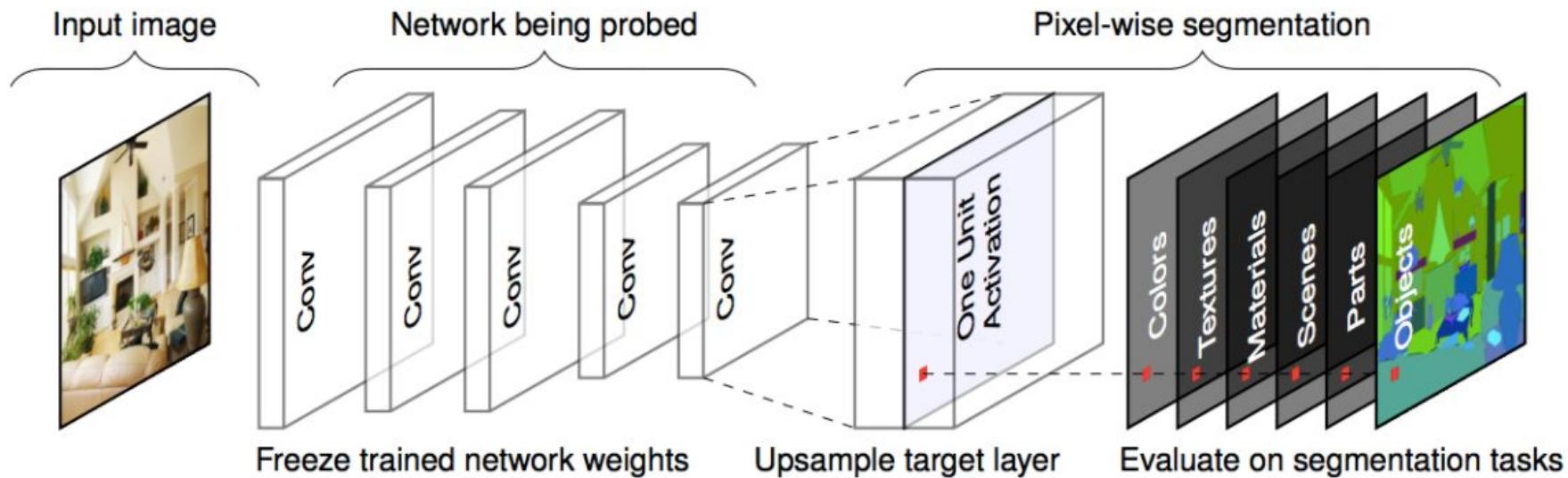


28) water tower

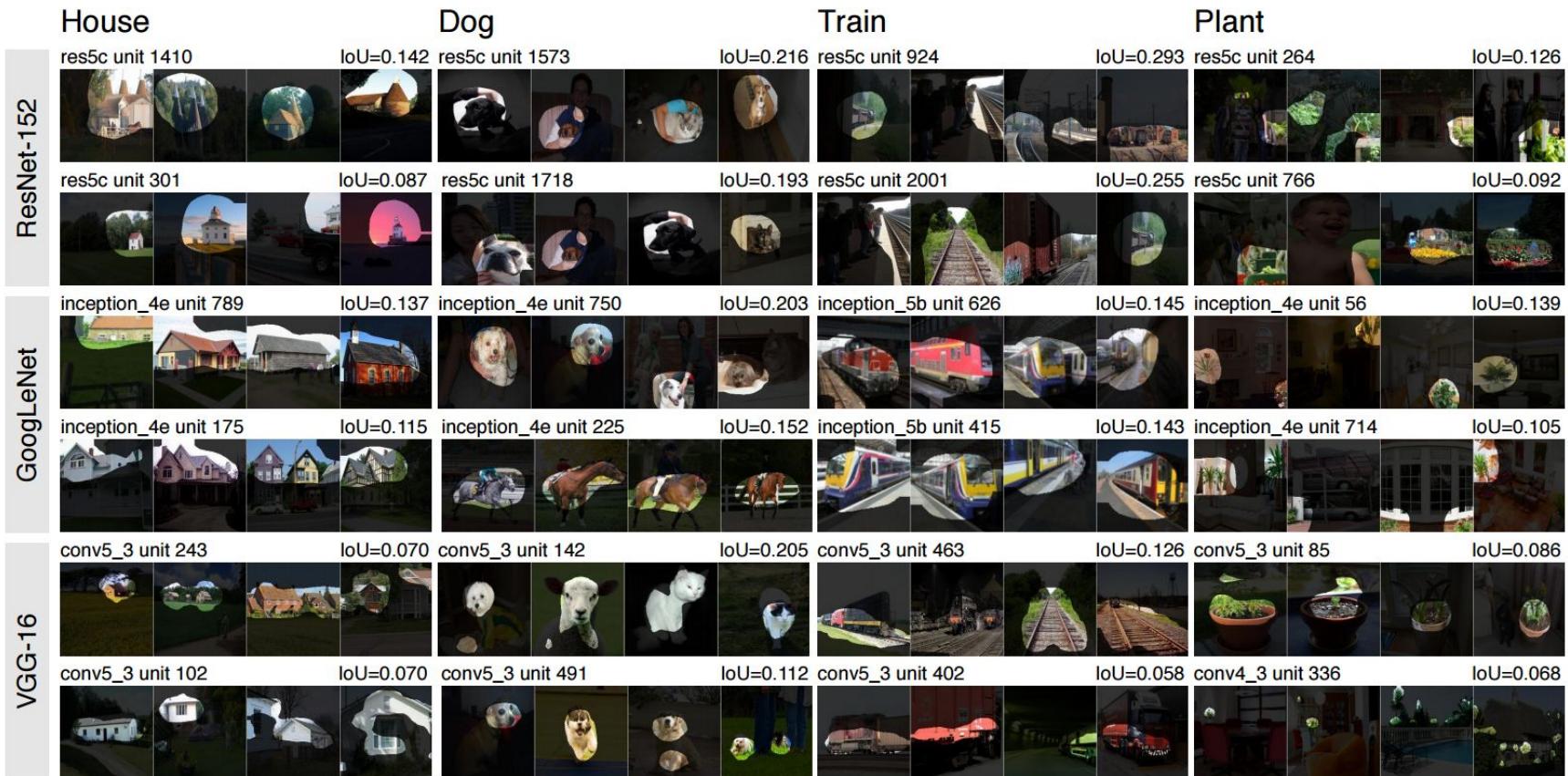


Activation Changes after Occlusions (pixel scale)

Same idea, but automatic unit labeling using densely labeled dataset (pixel-level annotations).
The thresholded activation of each conv unit in the network is evaluated for semantic segmentation.



Activation Changes after Occlusions (pixel scale)

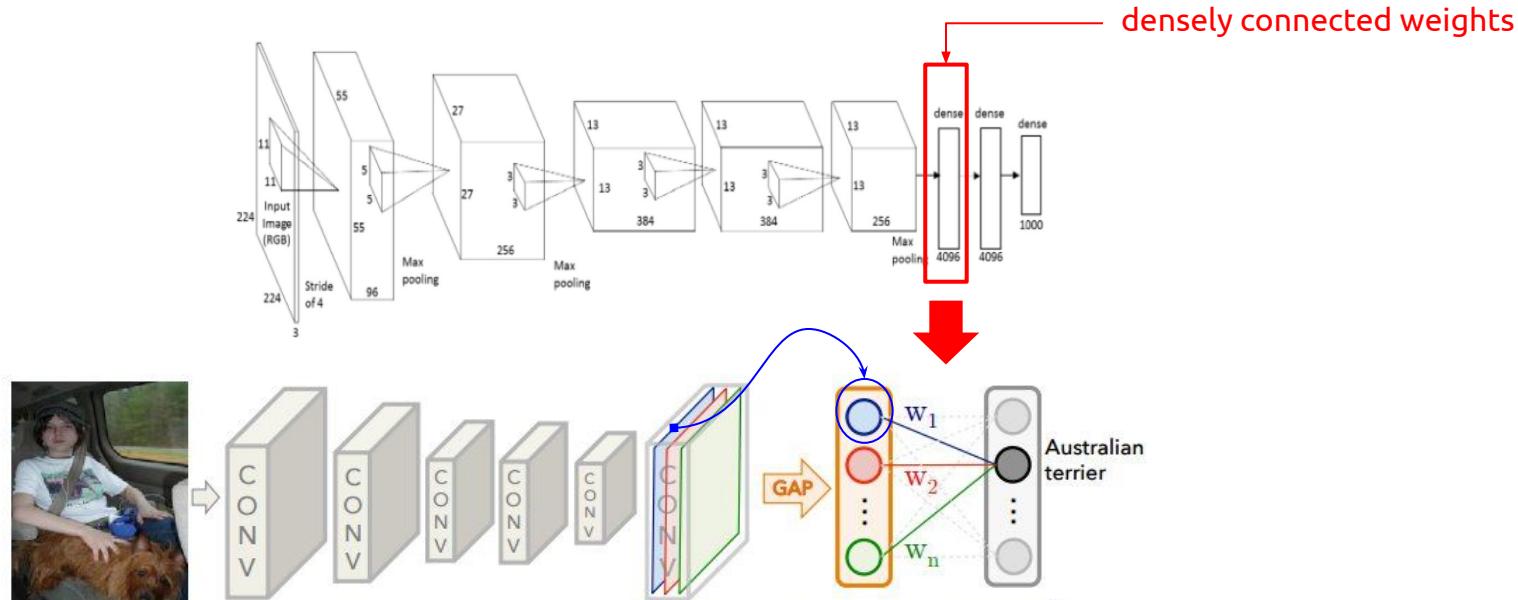


Interpretability: Attribution

- Visualization
- **Attribution**
 - Receptive Field of the Highest Activations
 - Activation Changes after Occlusions
 - **Class Activation Maps (CAMs)**
 - Gradient-based
- Feature visualization

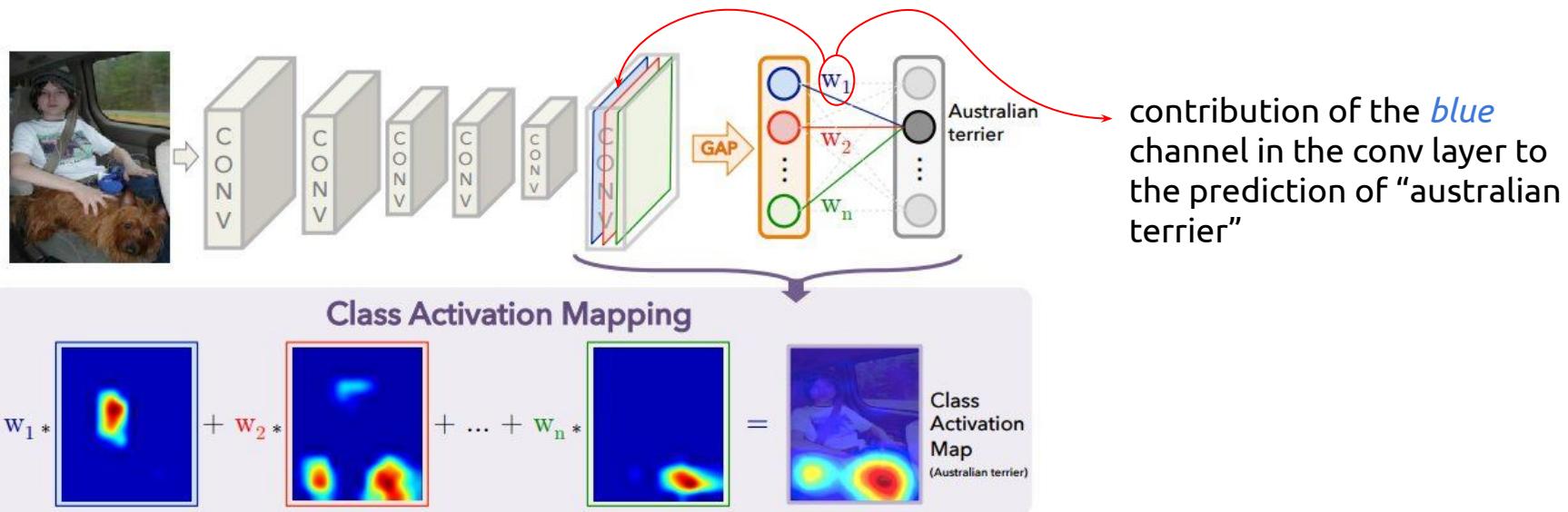
Class Activation Maps

Change in a classic CNN+MLP architecture: Replace FC layer after last conv with Global Average Pooling (GAP), which corresponds to averaging per channel.



Class Activation Maps

Weighted Fusion of Feature Maps : The classifier weights define the contribution of each channel in the previous layer

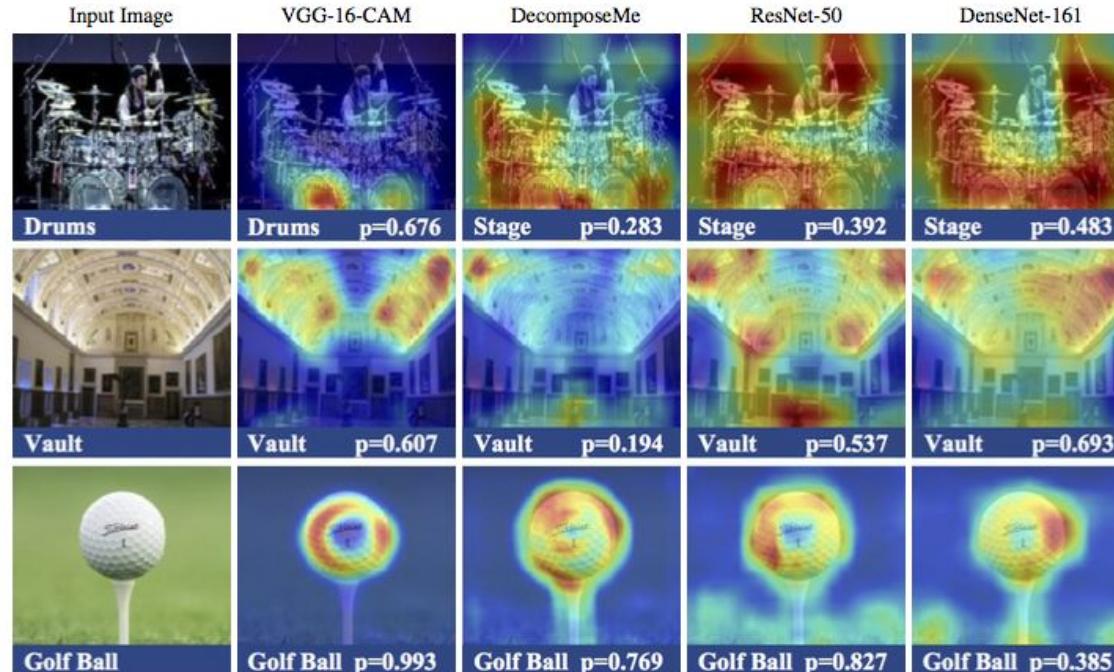


jinrikisha



Class Activation Maps

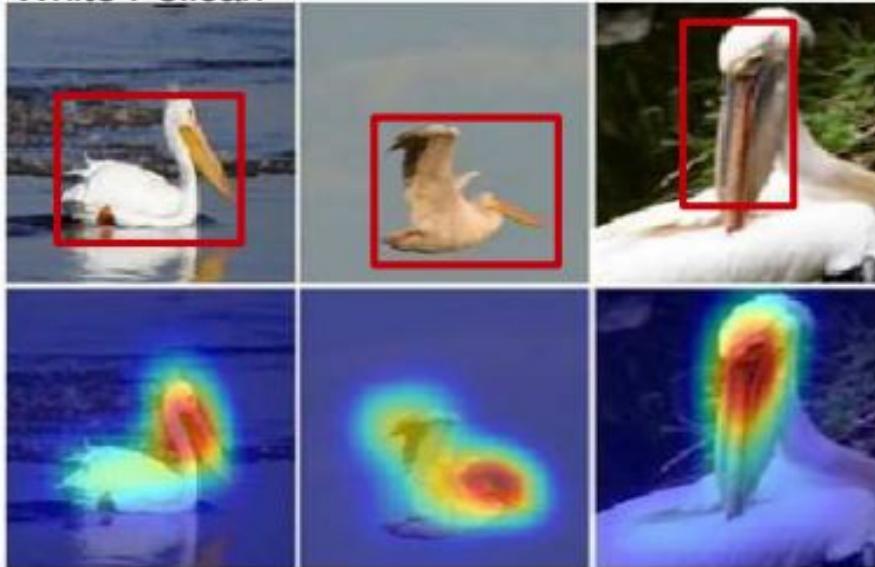
Different architectures generate different CAMs.



Class Activation Maps

Applications of CAMs: Rough object localization without weak labels (image).

White Pelican



Orchard Oriole



Class Activation Maps

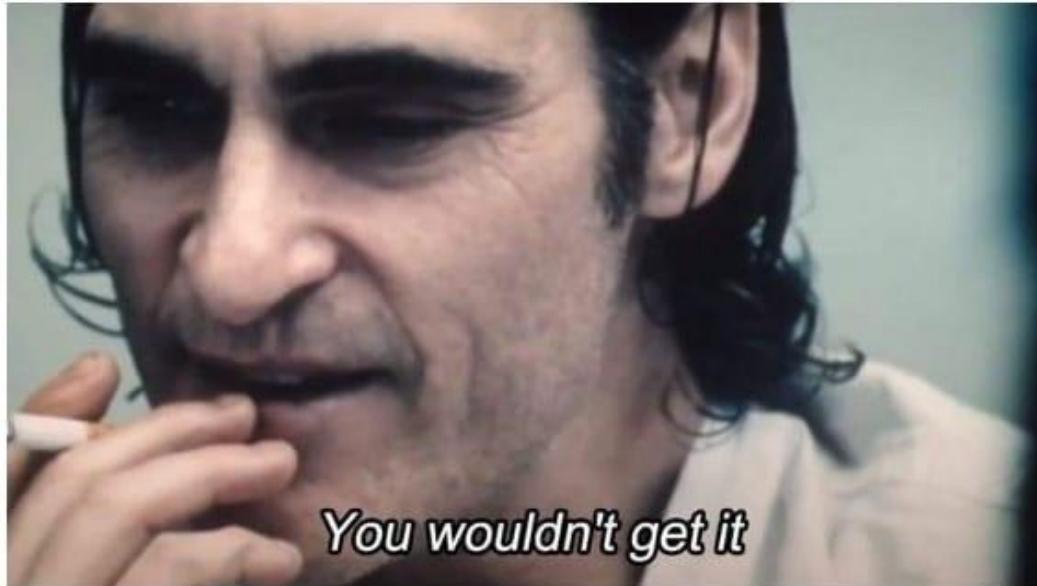
Applications of CAMs: Spatial feature weighting for object retrieval.



Class Activation Maps

Me: using Class Activation Maps
to understand where the model is
looking when making the decision

Model:

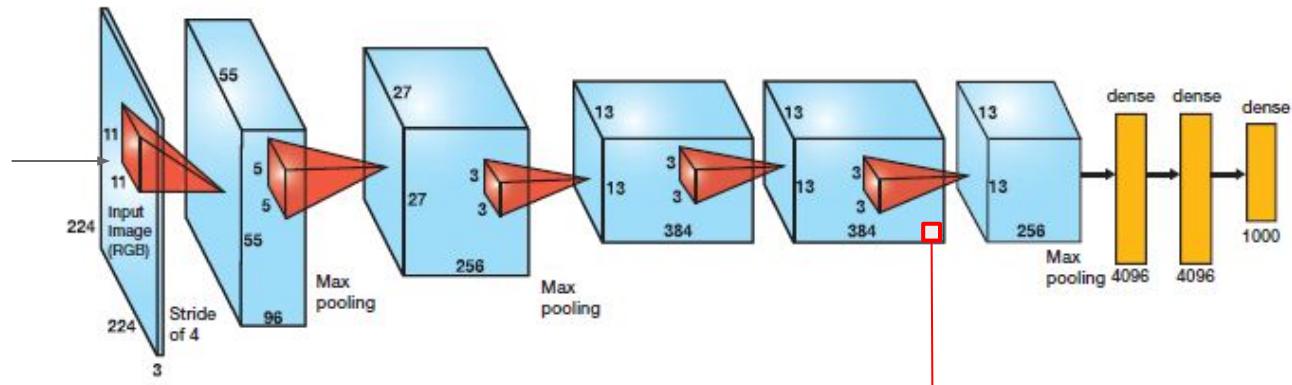


Interpretability: Attribution

- Visualization
- **Attribution**
 - Receptive Field of the Highest Activations
 - Activation Changes after Occlusions
 - Class Activation Maps (CAMs)
 - 2D Projections with t-SNE
 - **Gradient-based**
- Feature visualization

Gradient-based approach

Goal: Visualize the part of an image that mostly activates one of the neurons.



Compute the gradient of any neuron w.r.t. the image

1. Forward image up to the desired layer (e.g. conv5)
2. Set all gradients to 0
3. Set gradient for the neuron we are interested in to 1
4. Backpropagate to get reconstructed image (gradient on the image)

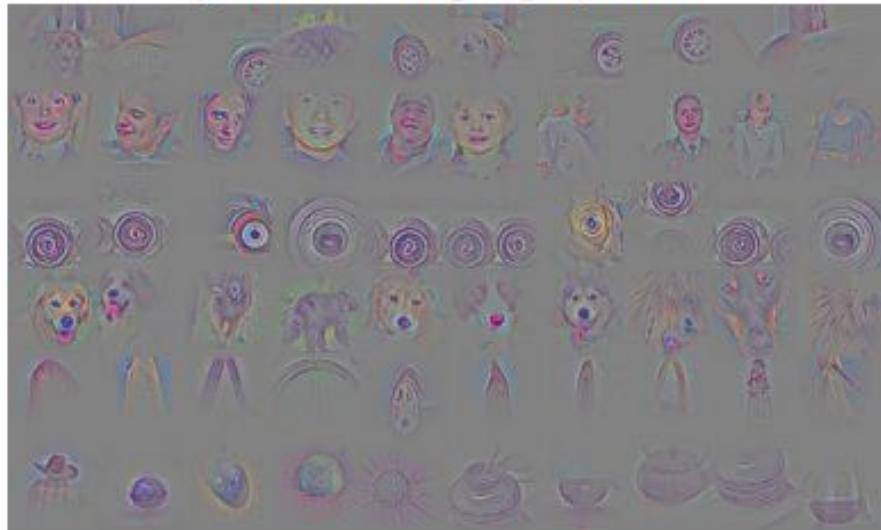
Gradient-based approach

Goal: Visualize the part of an image that mostly activates one of the neurons.

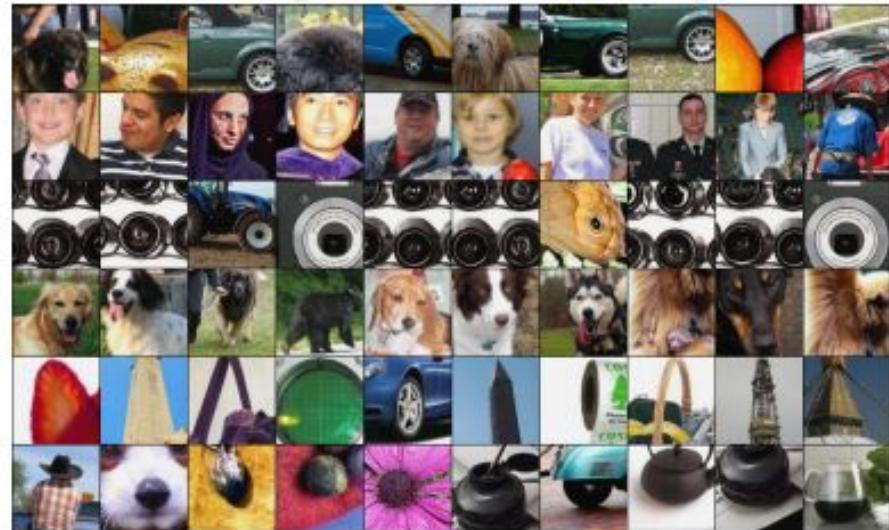


Gradient-based approach

guided backpropagation



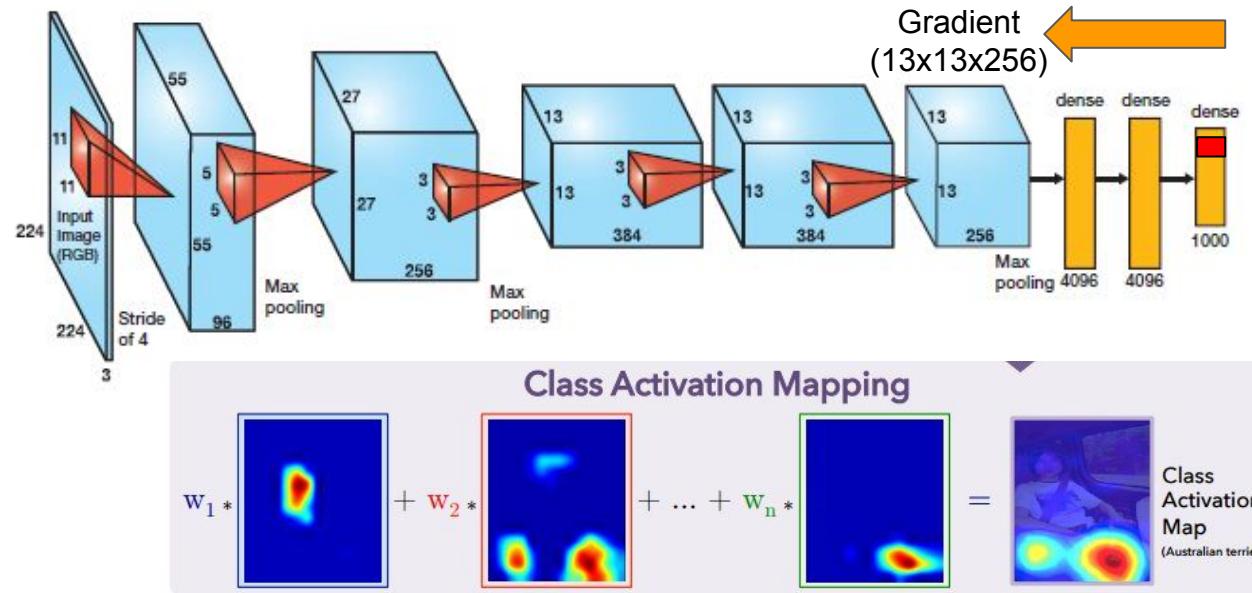
corresponding image crops



Grad-CAM

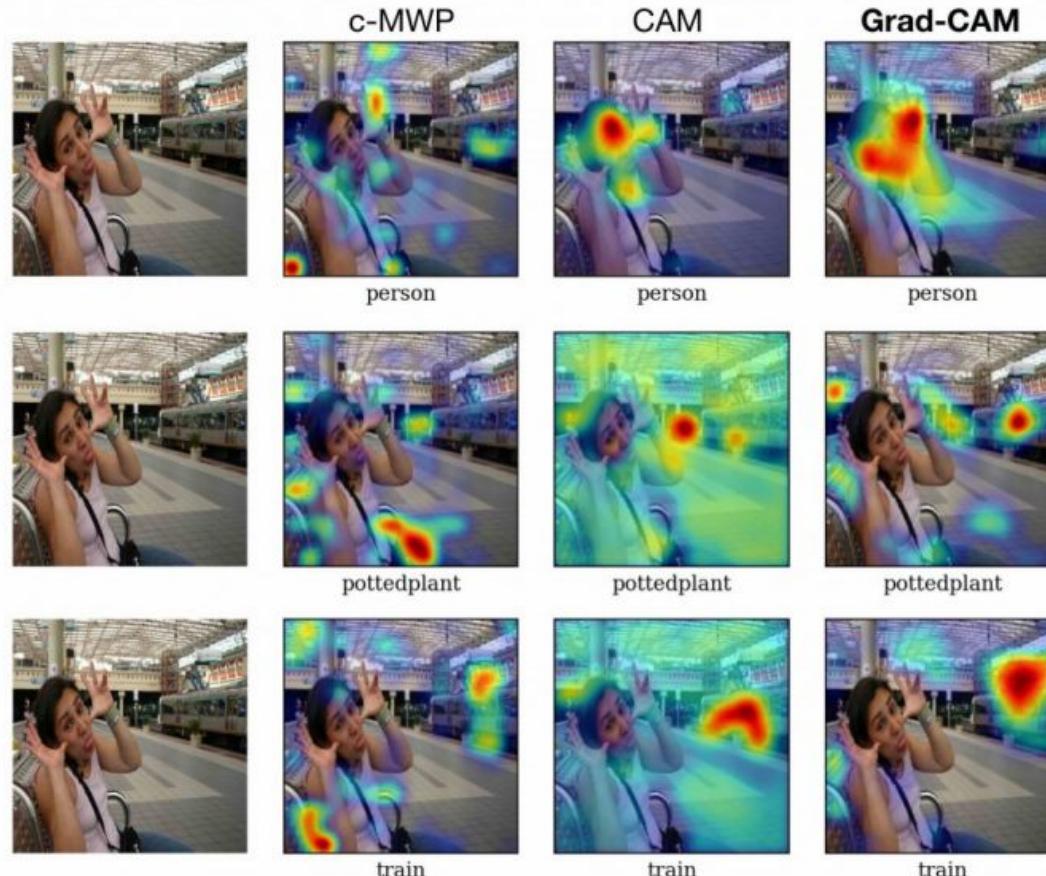
Same idea as Class Activation Maps (CAMs), but:

- No modifications required to the network
- Weight feature map by the gradient of the class wrt each channel



Feature maps weighted by the mean gradient.

Grad-CAM



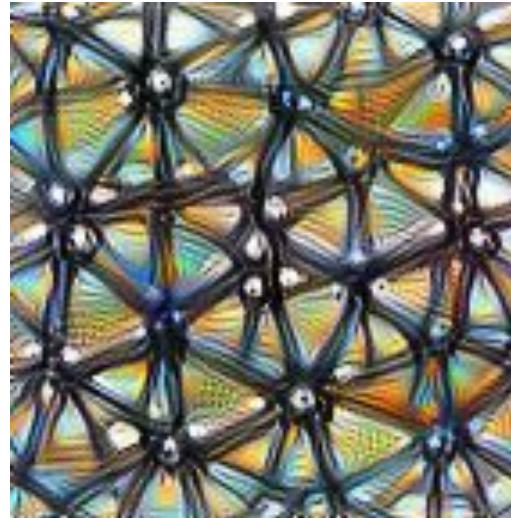
Ramprasaath et al. [Grad-CAM: Why did you say that? Visual Explanations from Deep Networks via Gradient-based Localization](#). ICCV 2017

Interpretability: Attribution

- Visualization
- Attribution
 - Receptive Field of the Highest Activations
 - Activation Changes after Occlusions
 - Class Activation Maps (CAMs)
 - 2D Projections with t-SNE
 - Gradient-based
- **Feature visualization**

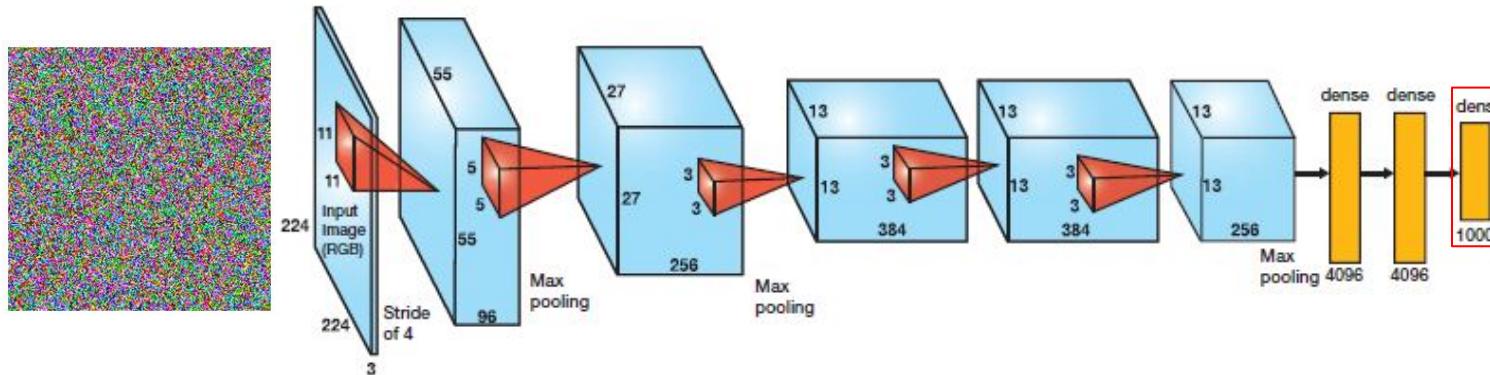
Feature Visualization

Feature visualization answers questions about what a network — or parts of a network — are looking for by generating examples



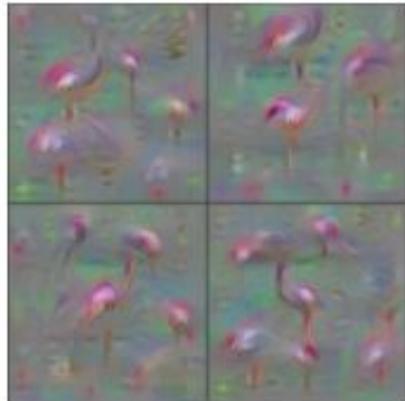
Feature Visualization

Goal: Generate an image that maximizes the activation of a neuron (or class confidence).

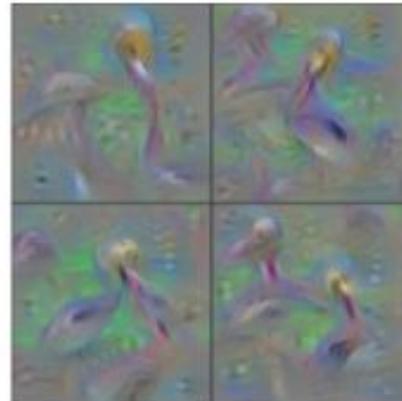


1. Forward random image
2. Set the gradient of the confidence vector to be $[0,0,0,\dots,1,\dots,0,0]$
3. Backprop to get gradient on the image
4. Update image (small step in the gradient direction)
5. Repeat

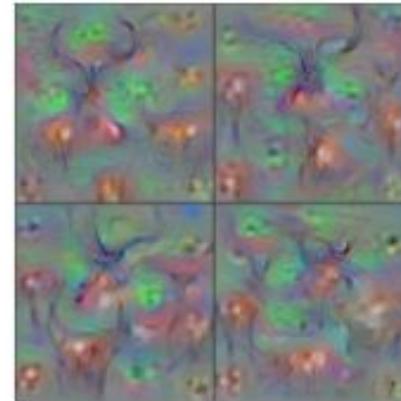
Feature Visualization



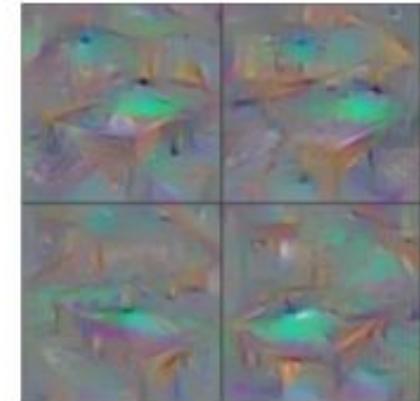
Flamingo



Pelican



Hartebeest



Billiard Table

Feature Visualization



Edges (layer conv2d0)

Textures (layer mixed3a)

Patterns (layer mixed4a)

Parts (layers mixed4b & mixed4c)

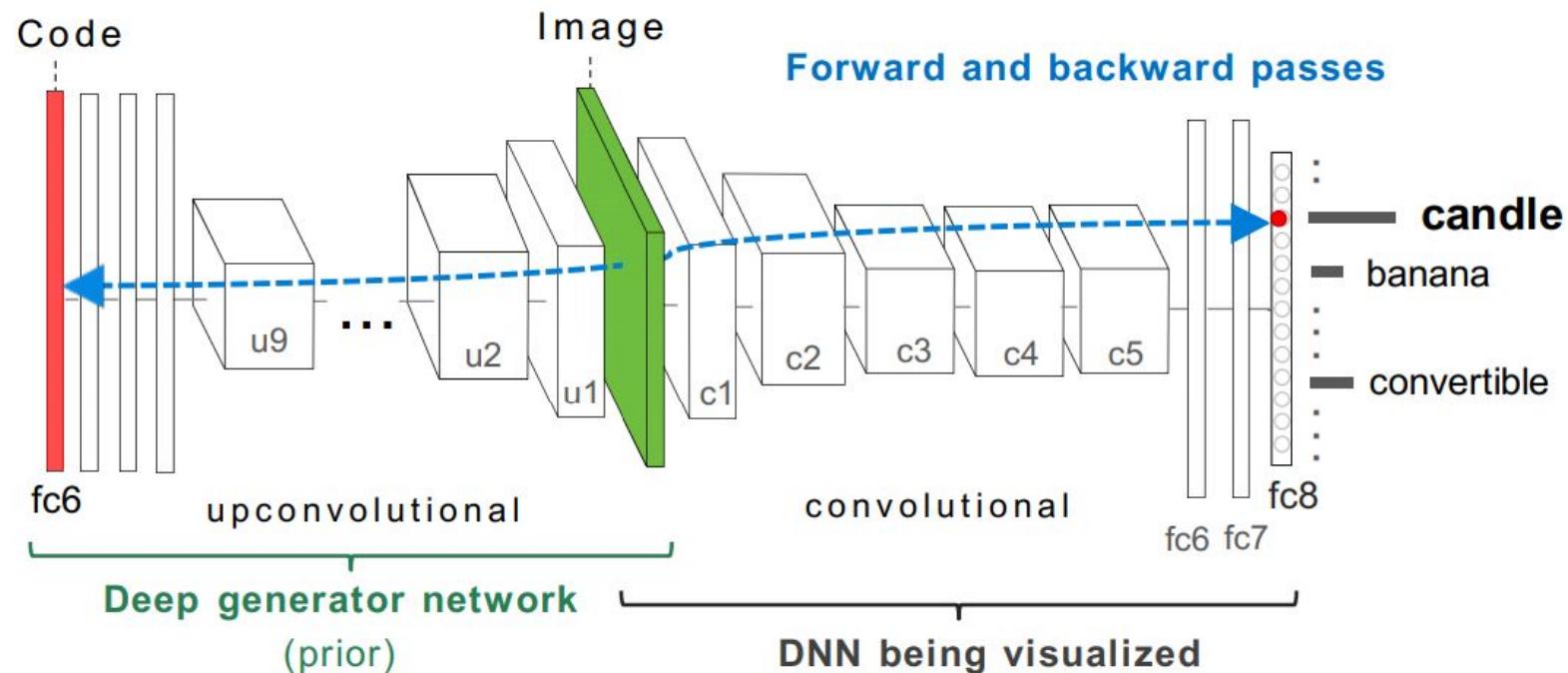
Objects (layers mixed4d & mixed4e)

Feature Visualization: Regularizations

	Unregularized	Frequency Penalization	Transformation Robustness	Learned Prior	Dataset Examples
 Erhan, et al., 2009 [3] Introduced core idea. Minimal regularization.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 Szegedy, et al., 2013 [11] Adversarial examples. Visualizes with dataset examples.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
 Mahendran & Vedaldi, 2015 [7] Introduces total variation regularizer. Reconstructs input from representation.	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 Nguyen, et al., 2015 [14] Explores counterexamples. Introduces image blurring.	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 Mordvintsev, et al., 2015 [4] Introduced jitter & multi-scale. Explored GMM priors for classes.	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
 Oygard, et al., 2015 [15] Introduces gradient blurring. (Also uses jitter.)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 Tyka, et al., 2016 [16] Regularizes with bilateral filters. (Also uses jitter.)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 Mordvintsev, et al., 2016 [17] Normalizes gradient frequencies. (Also uses jitter.)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 Nguyen, et al., 2016 [18] Paramaterizes images with GAN generator.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
 Nguyen, et al., 2016 [19] Uses denoising autoencoder prior to make a generative model.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

Feature Visualization: Generative prior

Optimizes over the input code (fc6 in the Fig.) instead of the image.



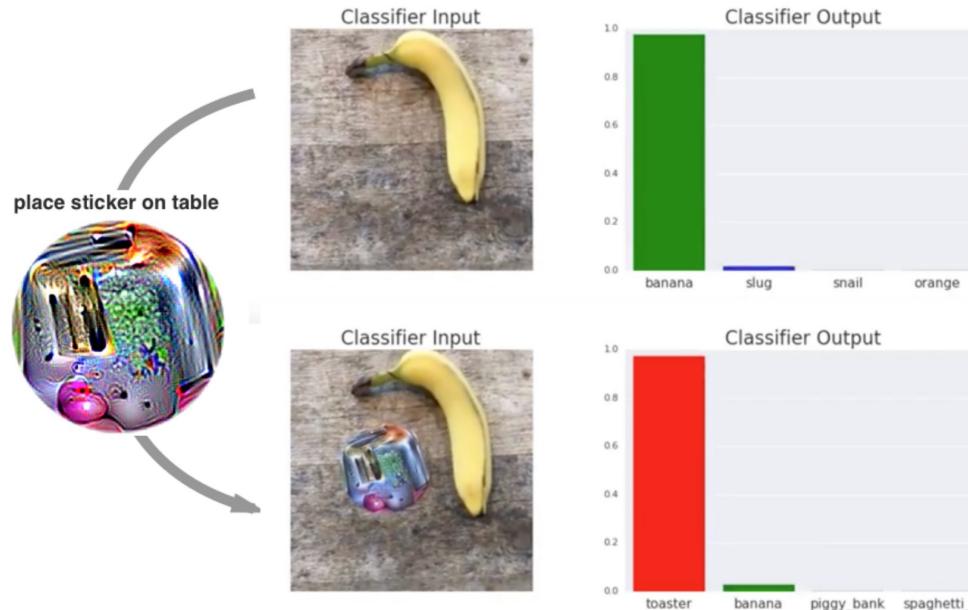
Feature Visualization: Generative prior



Nguyen, Anh, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. ["Synthesizing the preferred inputs for neurons in neural networks via deep generator networks."](#) NIPS 2016.

Bonus: Adversarial Attacks (eg. patch on sticker)

Goal: Generate a patch P that will make a Neural Network N always predict class C for any image containing the patch.



Bonus: Adversarial Attacks (eg. patch on sticker)

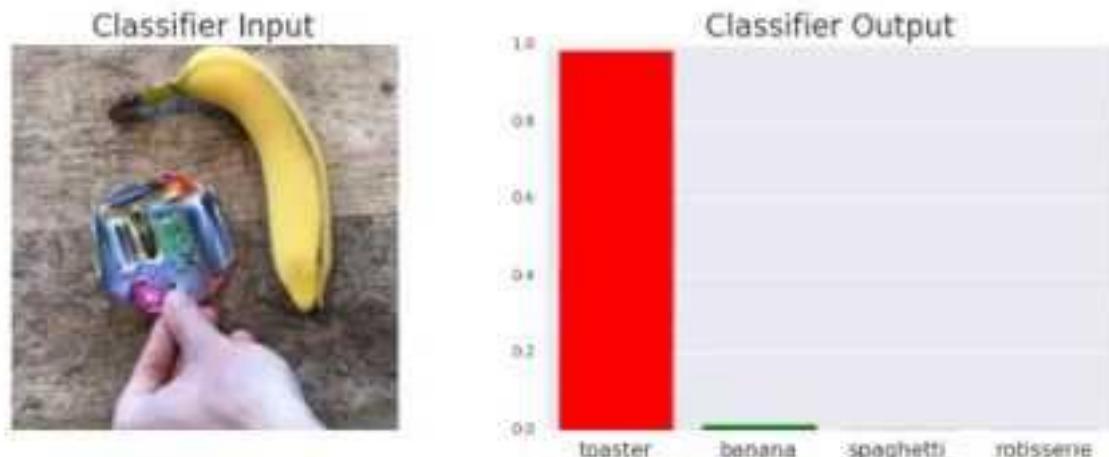
How: Backpropagate the gradient to images where the patch has been inserted.

$$A(\text{[patch image]}, \text{[original image]}, \text{location, rotation, scale, ...}) =$$



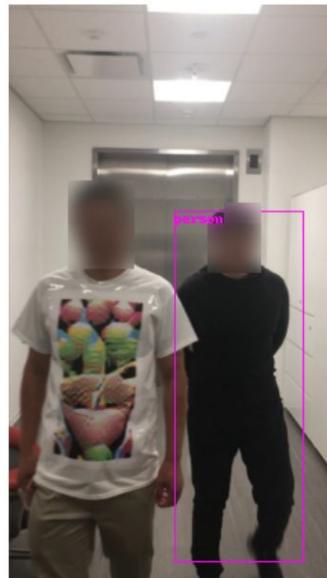
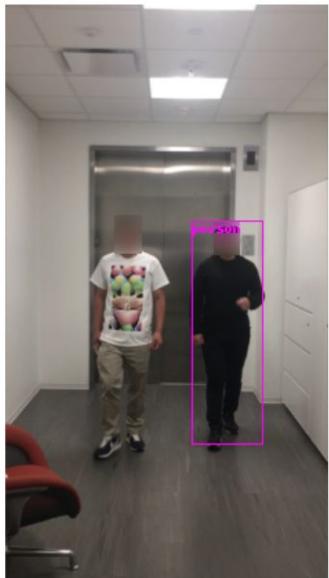
Brown, Tom B., Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. "[Adversarial patch.](#)" arXiv preprint arXiv:1712.09665 (2017).

Sachin Joglekar, "[Adversarial patches for CNNs explained](#)" (2018)



Brown, Tom B., Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. "Adversarial patch." arXiv preprint arXiv:1712.09665 (2017).

Bonus: Adversarial Attacks (eg. T-shirts for Privacy)



Xu, K., Zhang, G., Liu, S., Fan, Q., Sun, M., Chen, H., ... & Lin, X. (2019). [Evading Real-Time Person Detectors by Adversarial T-shirt](#). arXiv preprint arXiv:1910.11099.

Further readings

Shan Carter, Zan Armstrong, Ludwig Schubert, Ian Johnson, Chris Olah, "[Exploring Neural Networks with Activation Atlases](#)". Distill.pub (March 2019)

Adebayo, Julius, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. "[Sanity checks for saliency maps.](#)" NeurIPS 2018.

Software: [Captum](#), model interpretability for PyTorch



Questions ?

Undergradese

What undergrads ask vs. what they're REALLY asking

"Is it going to be an open book exam?"

Translation: "I don't have to actually memorize anything, do I?"

"Hmm, what do you mean by that?"

Translation: "What's the answer so we can all go home."

"Are you going to have office hours today?"

Translation: "Can I do my homework in your office?"

"Can i get an extension?"

Translation: "Can you re-arrange your life around mine?"

"Is this going to be on the test?"

Translation: "Tell us what's going to be on the test."

"Is grading going to be curved?"

Translation: "Can I do a mediocre job and still get an A?"

