

---

## Seminar Explainable AI Module 6

# Selected Methods Part 2

## Deconv-Inverting-Guided-Backprop-

## DGN-TCAV

**Andreas Holzinger**

Human-Centered AI Lab (Holzinger Group)

Institute for Medical Informatics/Statistics, Medical University Graz, Austria  
and

Explainable AI-Lab, Alberta Machine Intelligence Institute, Edmonton, Canada



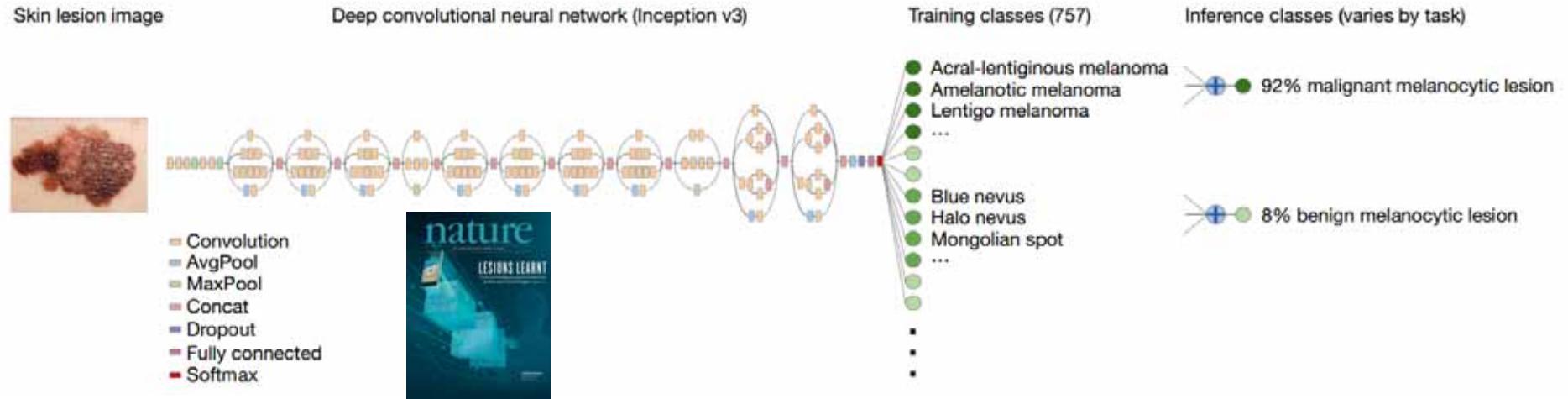
**This is the version for  
printing and reading.  
The lecture version is  
didactically different.**

- 00 Reflection
- 01 Visualizing CNN with Deconvolution
- 02 Inverting CNN
- 03 Guided Backpropagation
- 04 Deep Generator Networks
- 05 Testing with Concept Activation Vectors

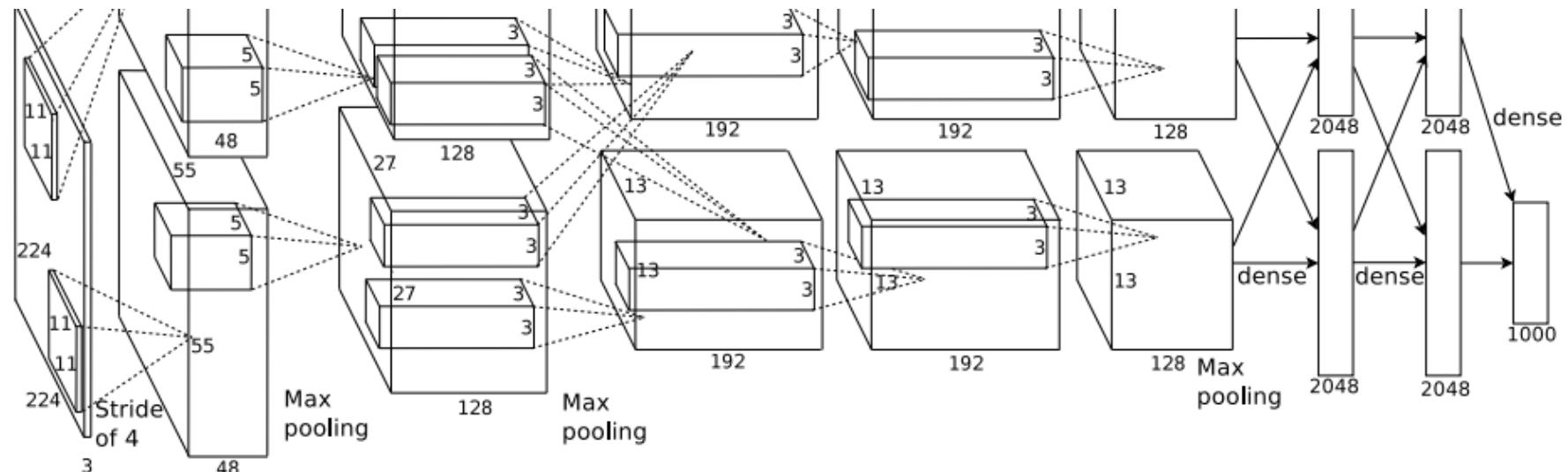
---

# 01 Visualizing CNN with Deconvolution

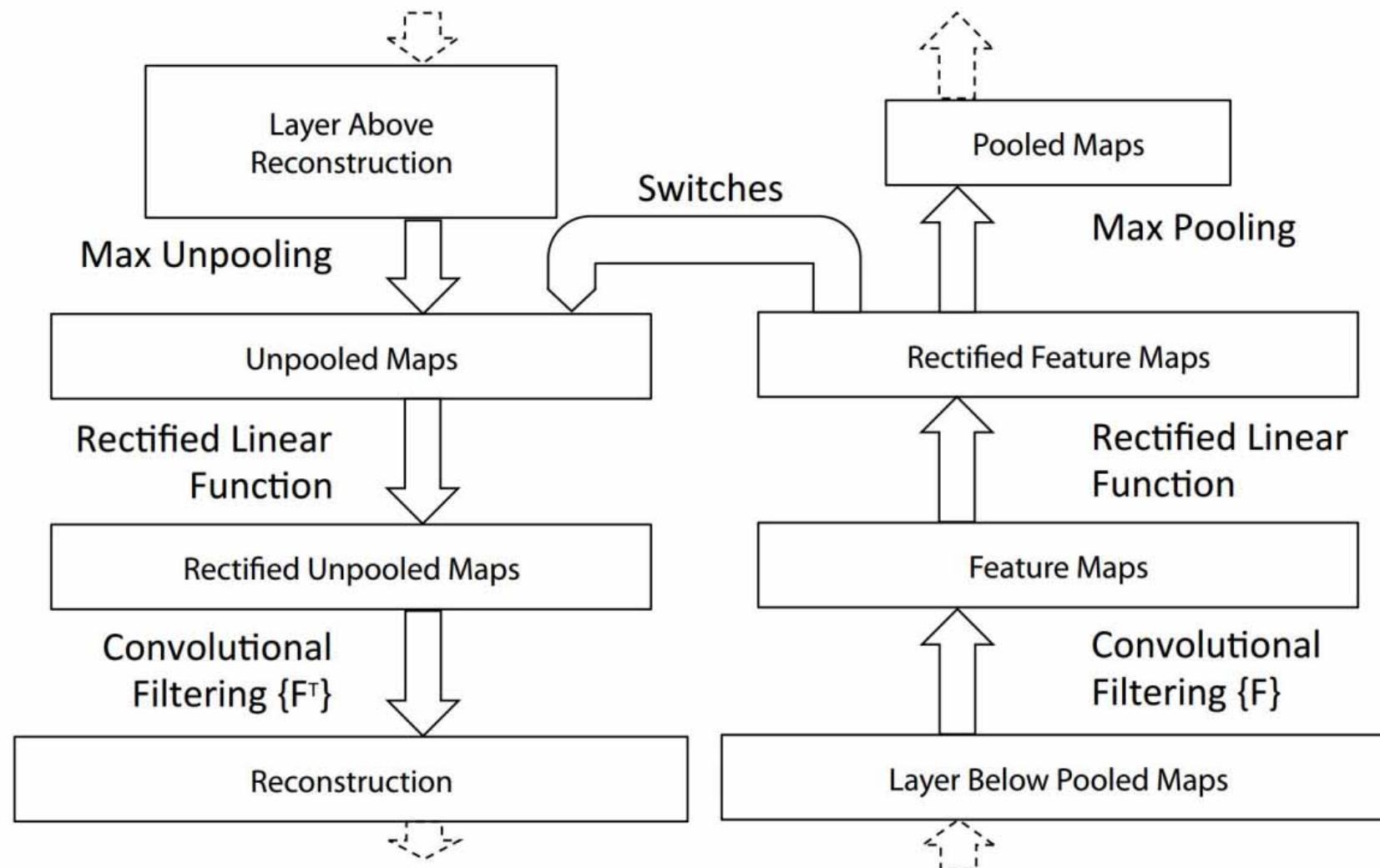
# Our typical example to answer the question of why



Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M. & Thrun, S. 2017. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542, (7639), 115-118, doi:10.1038/nature21056.

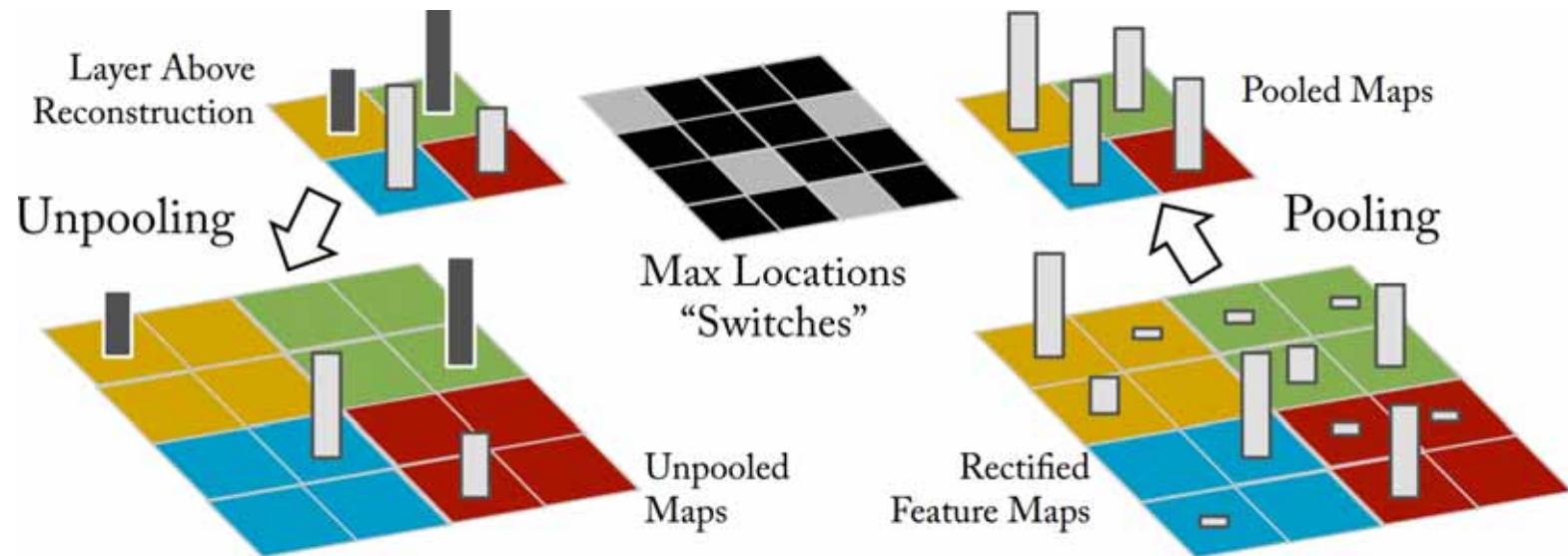
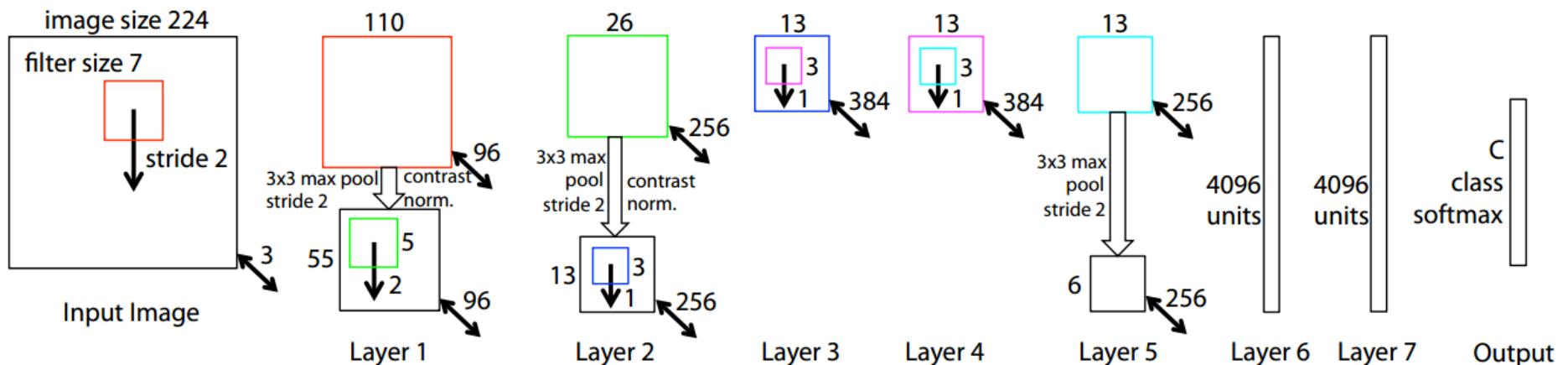


# Example: Interpretable Deep Learning Model

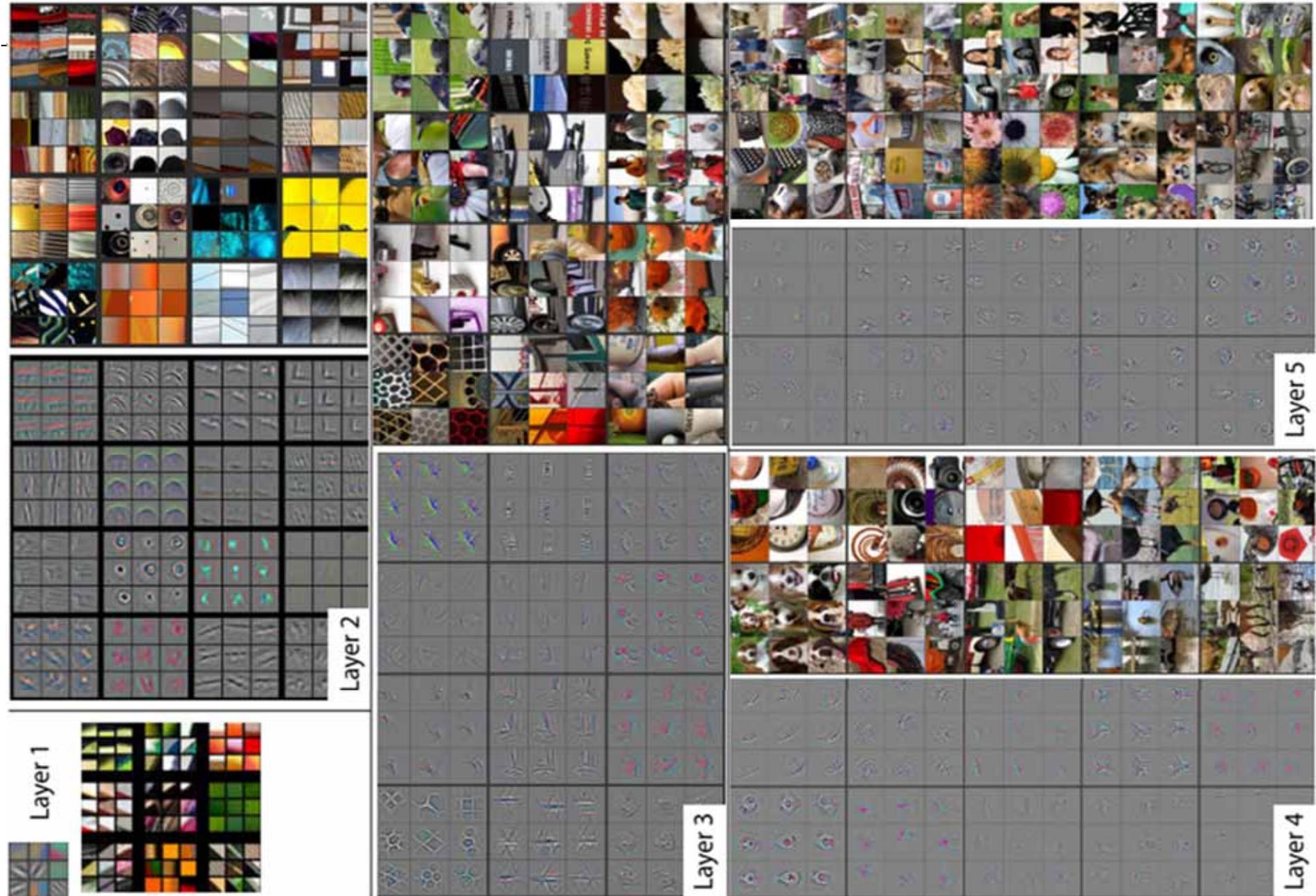


Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.

# Visualizing a Conv Net with a De-Conv Net

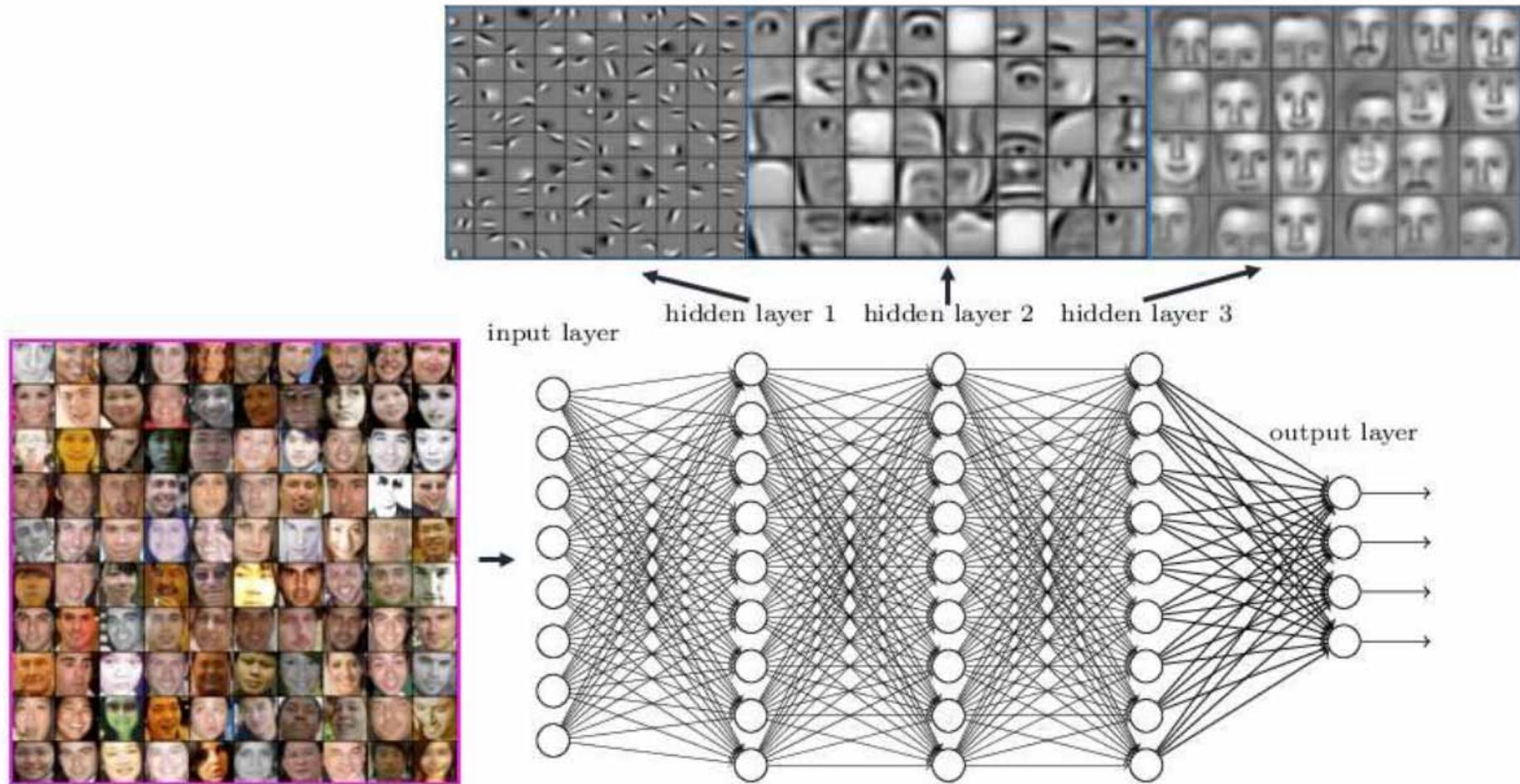


Matthew D. Zeiler & Rob Fergus 2014. Visualizing and understanding convolutional networks. In: D., Fleet, T., Pajdla, B., Schiele & T., Tuytelaars (eds.) ECCV, Lecture Notes in Computer Science LNCS 8689. Cham: Springer, pp. 818-833, doi:10.1007/978-3-319-10590-1\_53.



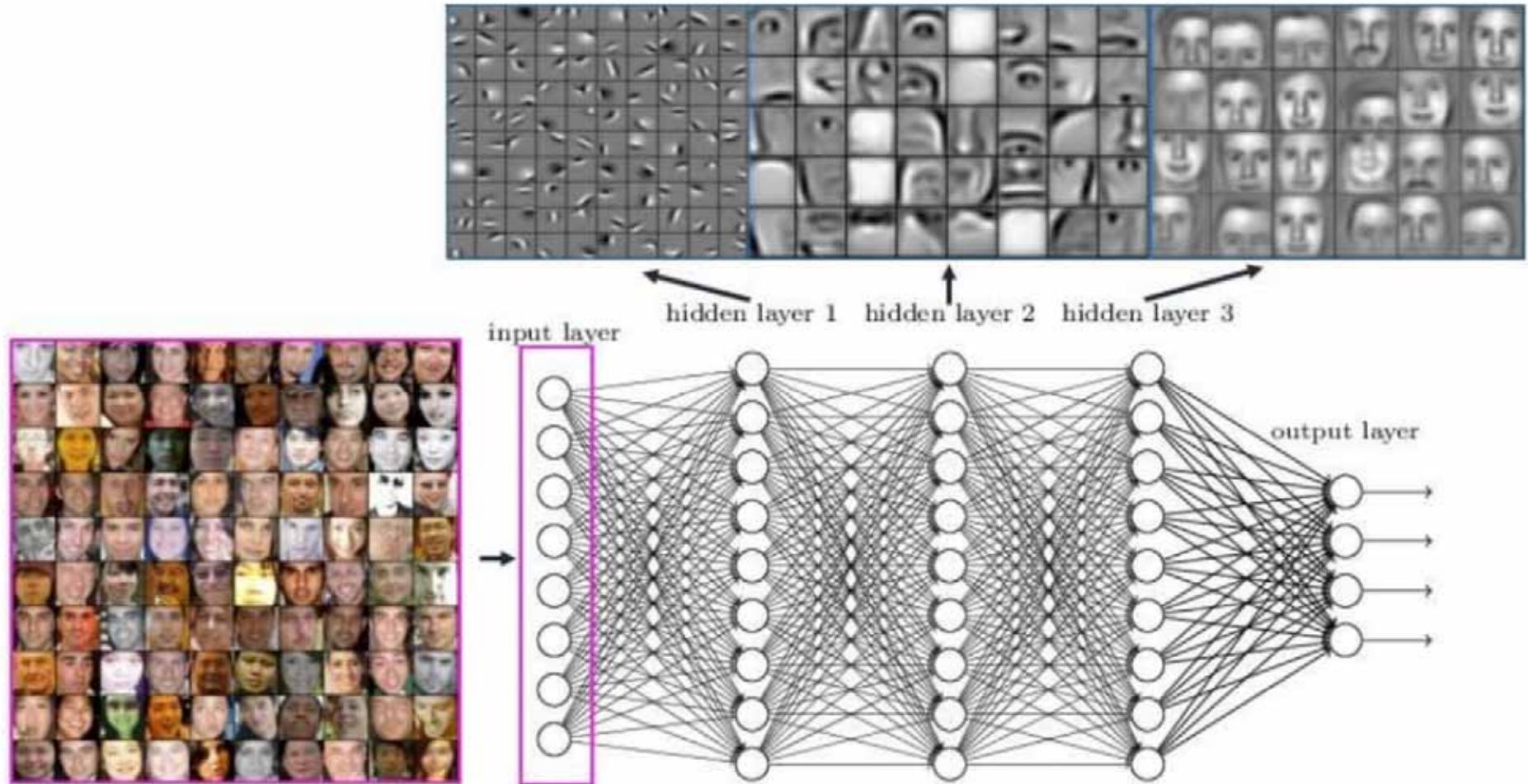
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.  
a.holzinger@human-centered.ai

# The world is compositional (Yann LeCun)

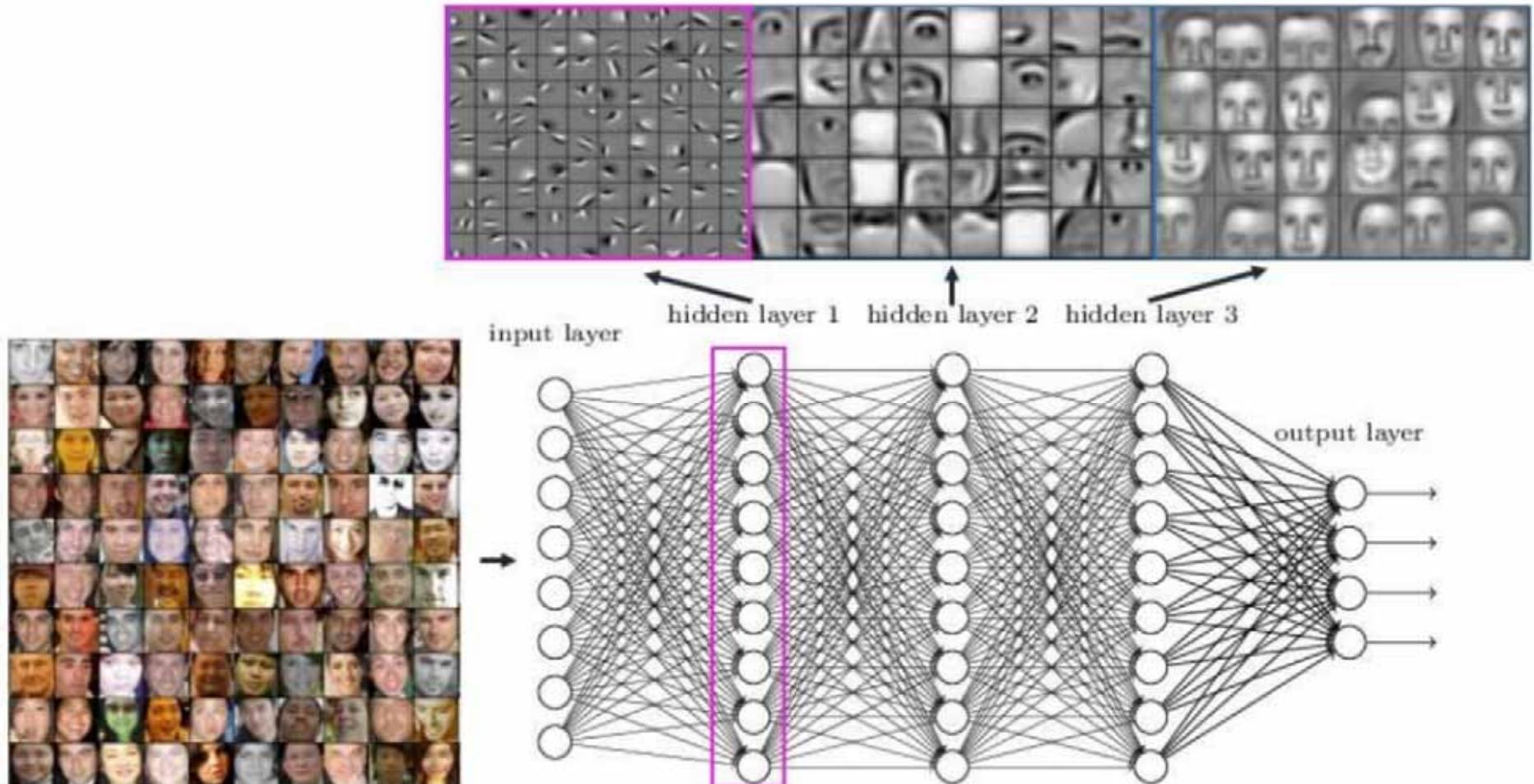


Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901

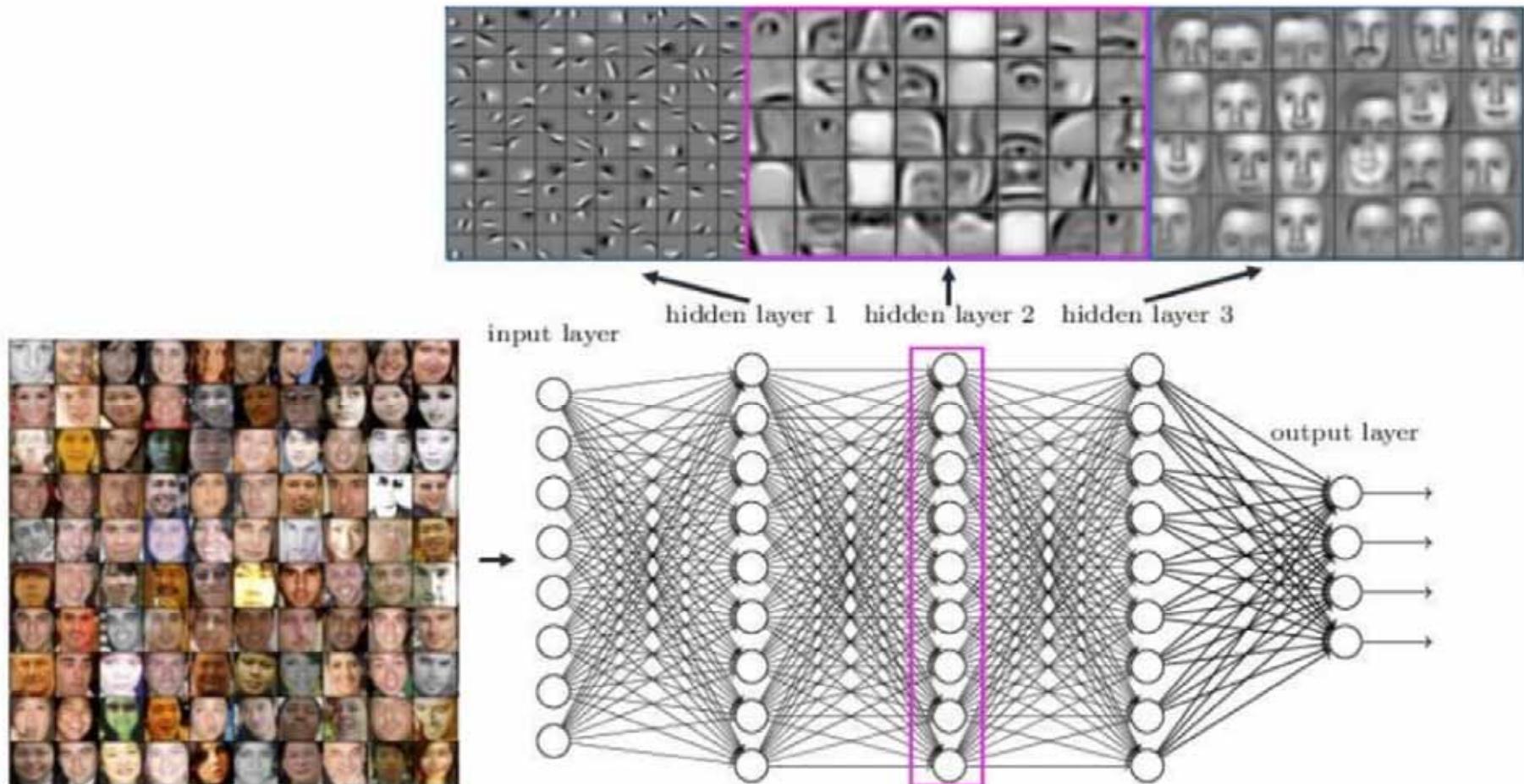
# The world is compositional (Yann LeCun)



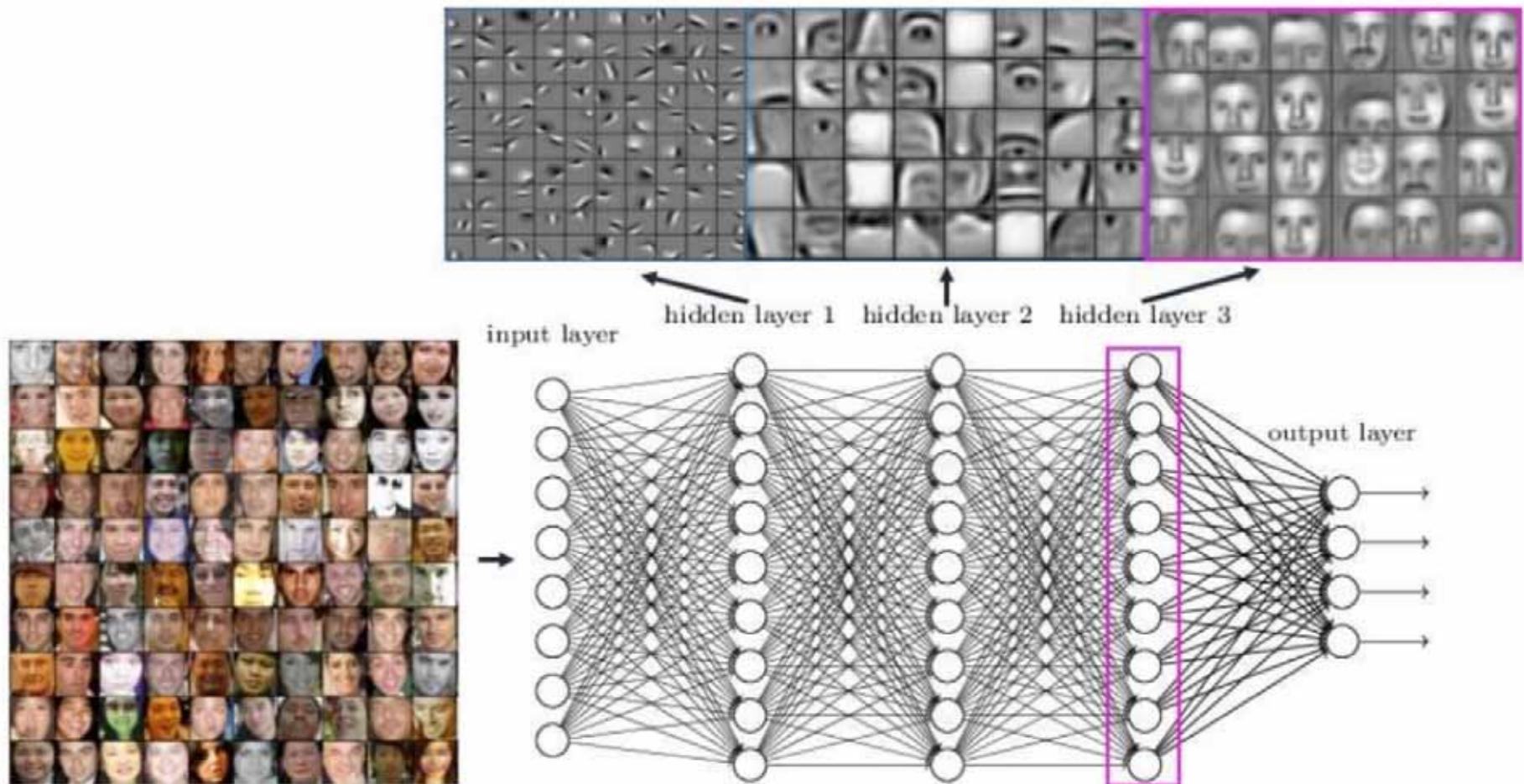
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901



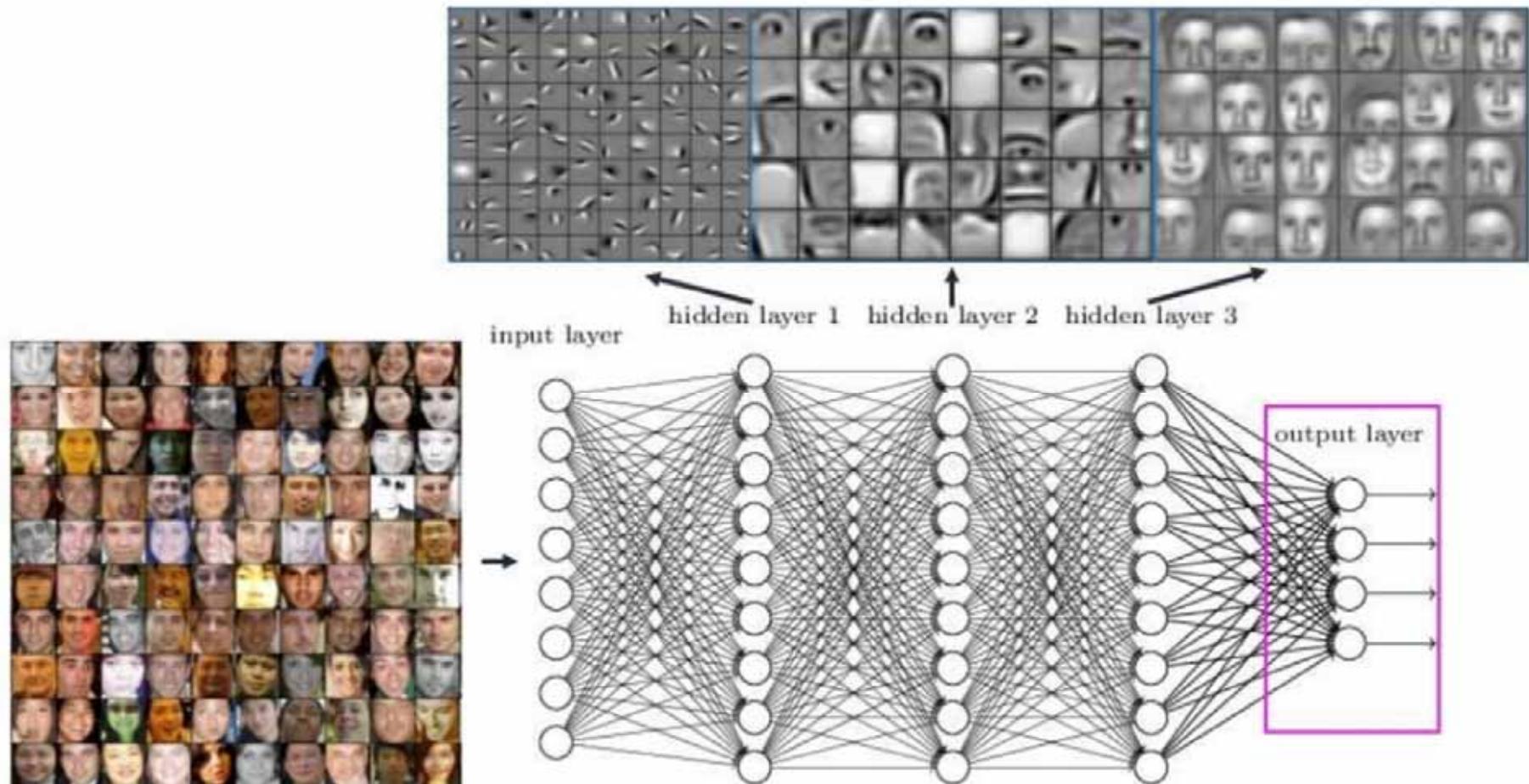
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901



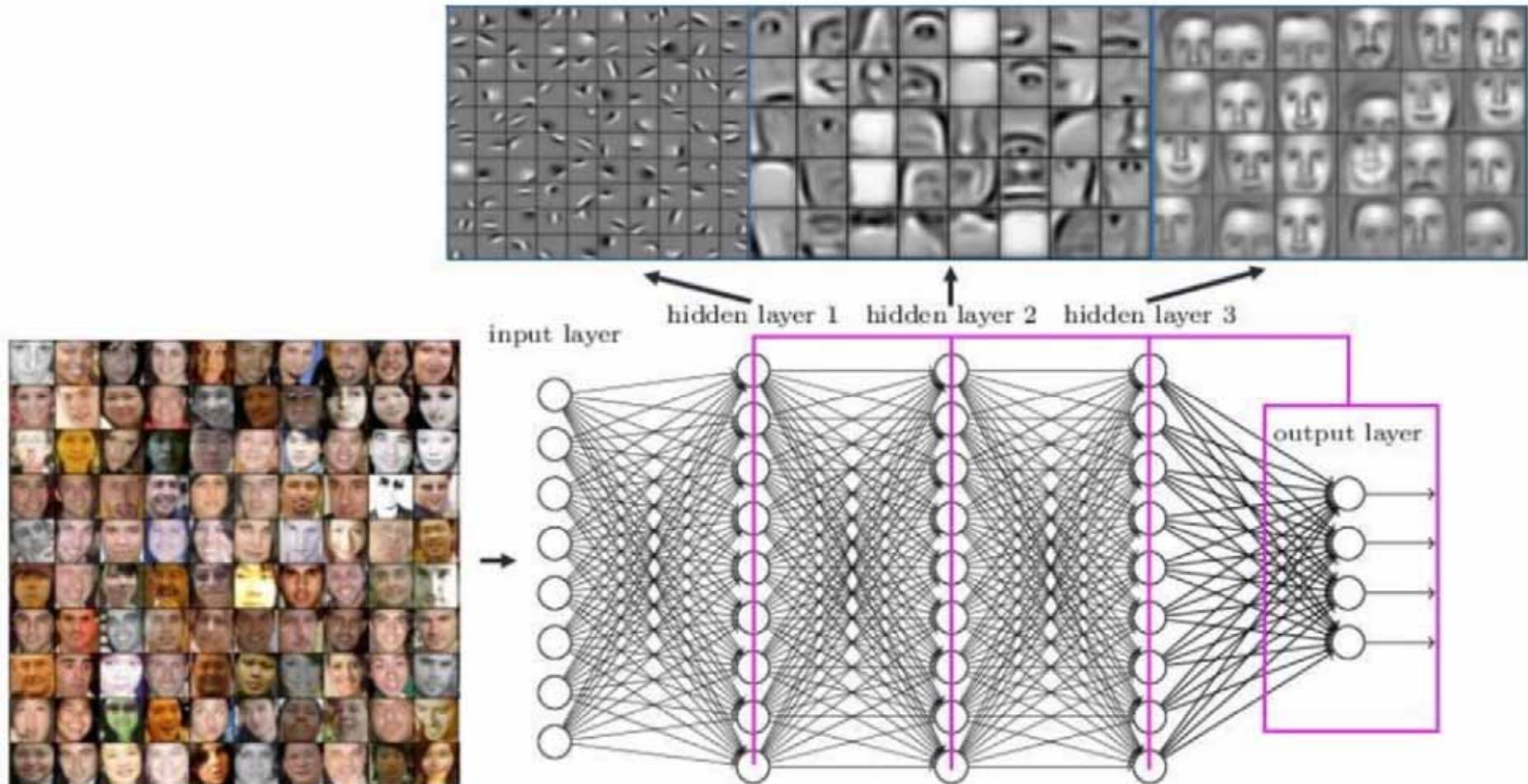
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901



Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901



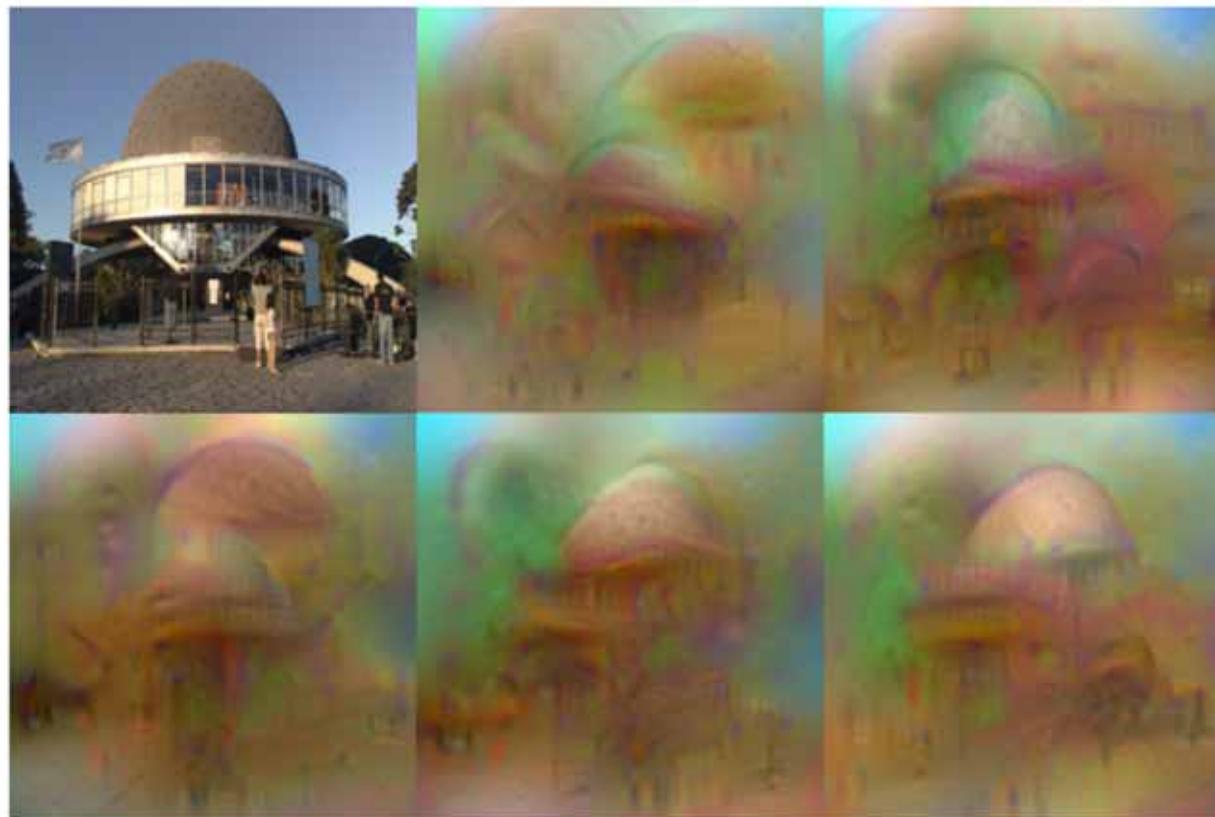
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901



Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901

# 02 Inverting Convolutional Neural Networks

# What is encoded by a CNN ? (before softmax)



$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Aravindh Mahendran & Andrea Vedaldi. Understanding deep image representations by inverting them. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015. 5188-5196, doi:10.1109/CVPR.2015.7299155.

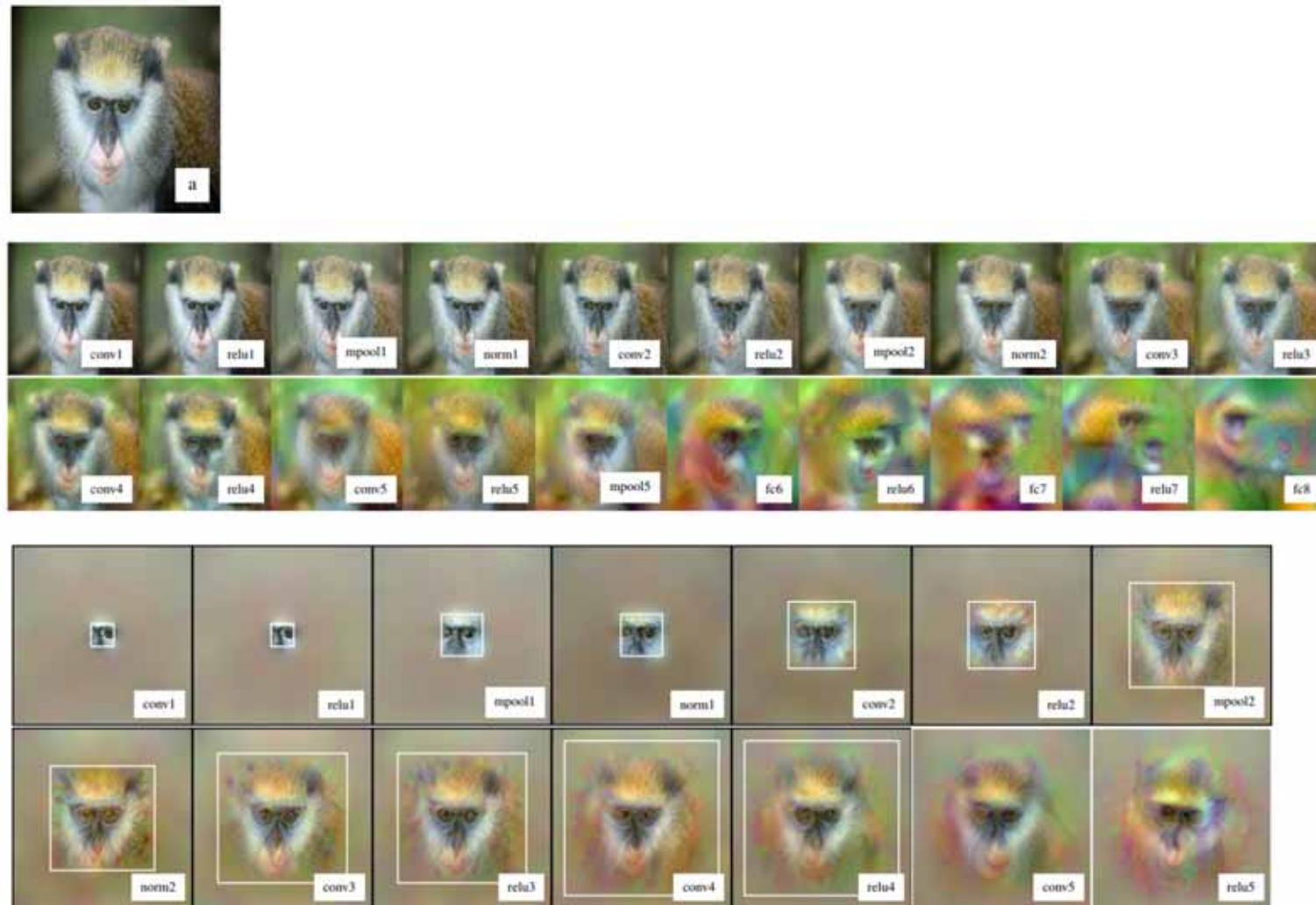
$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

$$\mathcal{R}_{V^\beta}(f) = \int_{\Omega} \left( \left( \frac{\partial f}{\partial u}(u, v) \right)^2 + \left( \frac{\partial f}{\partial v}(u, v) \right)^2 \right)^{\frac{\beta}{2}} du dv$$

$$\mathcal{R}_{V^\beta}(\mathbf{x}) = \sum_{i,j} \left( (x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{\frac{\beta}{2}}$$



Aravindh Mahendran & Andrea Vedaldi. Understanding deep image representations by inverting them. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015. 5188-5196, doi:10.1109/CVPR.2015.7299155.



Aravindh Mahendran & Andrea Vedaldi. Understanding deep image representations by inverting them. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015. 5188-5196, doi:10.1109/CVPR.2015.7299155.

# 03 Guided Backpropagation

Model	Error (%)	# parameters
<hr/>		
without data augmentation		
Model A	12.47%	$\approx 0.9$ M
Strided-CNN-A	13.46%	$\approx 0.9$ M
ConvPool-CNN-A	<b>10.21%</b>	$\approx 1.28$ M
ALL-CNN-A	10.30%	$\approx 1.28$ M
Model B	10.20%	$\approx 1$ M
Strided-CNN-B	10.98%	$\approx 1$ M
ConvPool-CNN-B	9.33%	$\approx 1.35$ M
ALL-CNN-B	<b>9.10%</b>	$\approx 1.35$ M
Model C	9.74%	$\approx 1.3$ M
Strided-CNN-C	10.19%	$\approx 1.3$ M
ConvPool-CNN-C	9.31%	$\approx 1.4$ M
ALL-CNN-C	<b>9.08%</b>	$\approx 1.4$ M

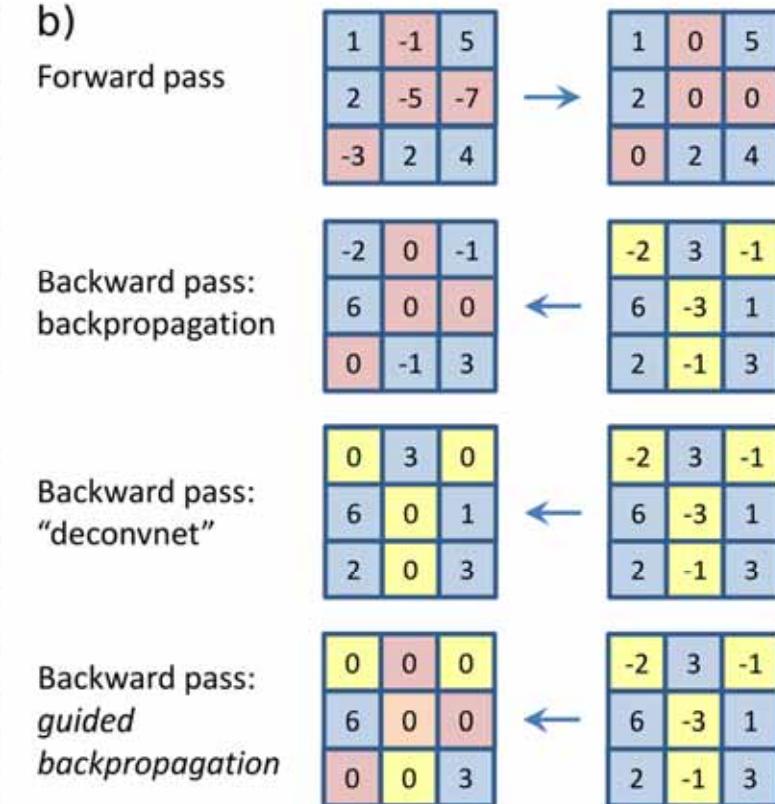
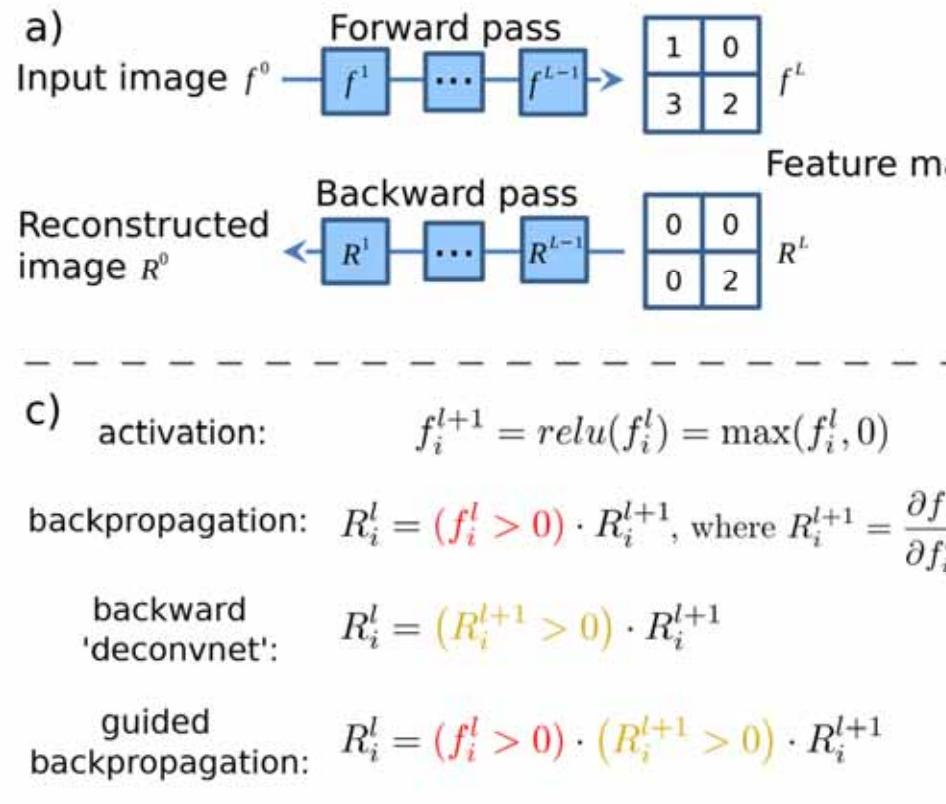
Method	Error (%)
CNN + tree prior [5]	36.85%
Network in Network [2]	35.68%
Deeply Supervised [3]	34.57%
Maxout (larger) [4]	34.54%
dasNet [4]	33.78%
ALL-CNN (Ours)	33.71%
Fractional Pooling (1 test) [6]	<b>31.45%</b>
<b>Fractional Pooling (12 tests) [6]</b>	<b>26.39%</b>

Method	Error (%)	# params
<hr/>		
without data augmentation		
Maxout [1]	11.68%	> 6 M
Network in Network [2]	10.41%	$\approx 1$ M
Deeply Supervised [3]	9.69%	$\approx 1$ M
<b>ALL-CNN (Ours)</b>	<b>9.08%</b>	$\approx 1.3$ M
<hr/>		
with data augmentation		
Maxout [1]	9.38%	> 6 M
DropConnect [2]	9.32%	-
dasNet [4]	9.22%	> 6 M
Network in Network [2]	8.81%	$\approx 1$ M
Deeply Supervised [3]	7.97%	$\approx 1$ M
<b>ALL-CNN (Ours)</b>	<b>7.25%</b>	$\approx 1.3$ M

Method	Error (%)
<hr/>	
with large data augmentation	
Spatially Sparse CNN [6]	4.47%
Large ALL-CNN (Ours)	4.41%
Fractional Pooling (1 test) [6]	4.50%
<b>Fractional Pooling (100 tests) [6]</b>	<b>3.47%</b>

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox & Martin Riedmiller 2014. Striving for simplicity:  
The all convolutional net. arXiv:1412.6806.

# Visualising the activations of high layer neurons



Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox & Martin Riedmiller 2014. Striving for simplicity:  
The all convolutional net. arXiv:1412.6806.

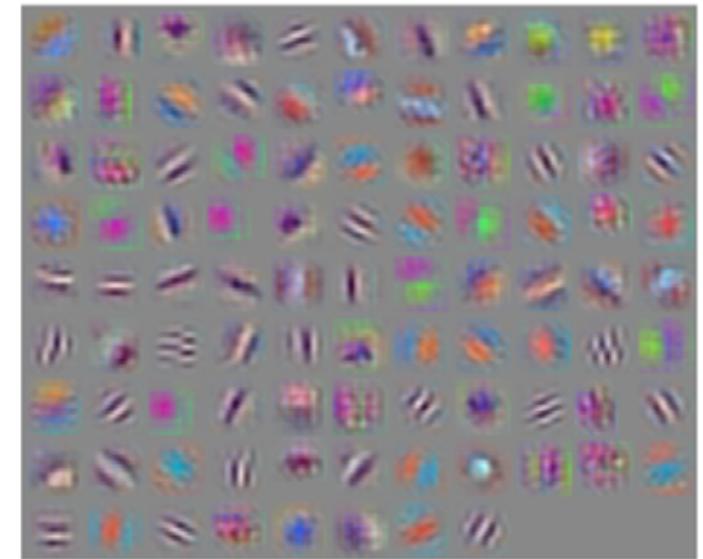
conv1



conv2



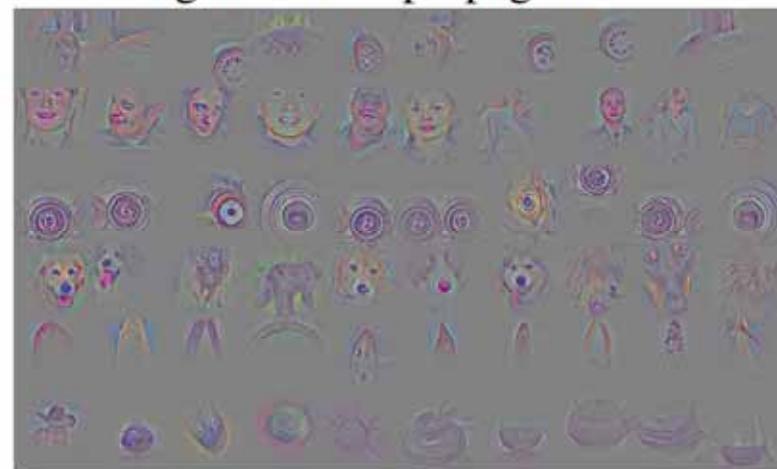
conv3



deconv



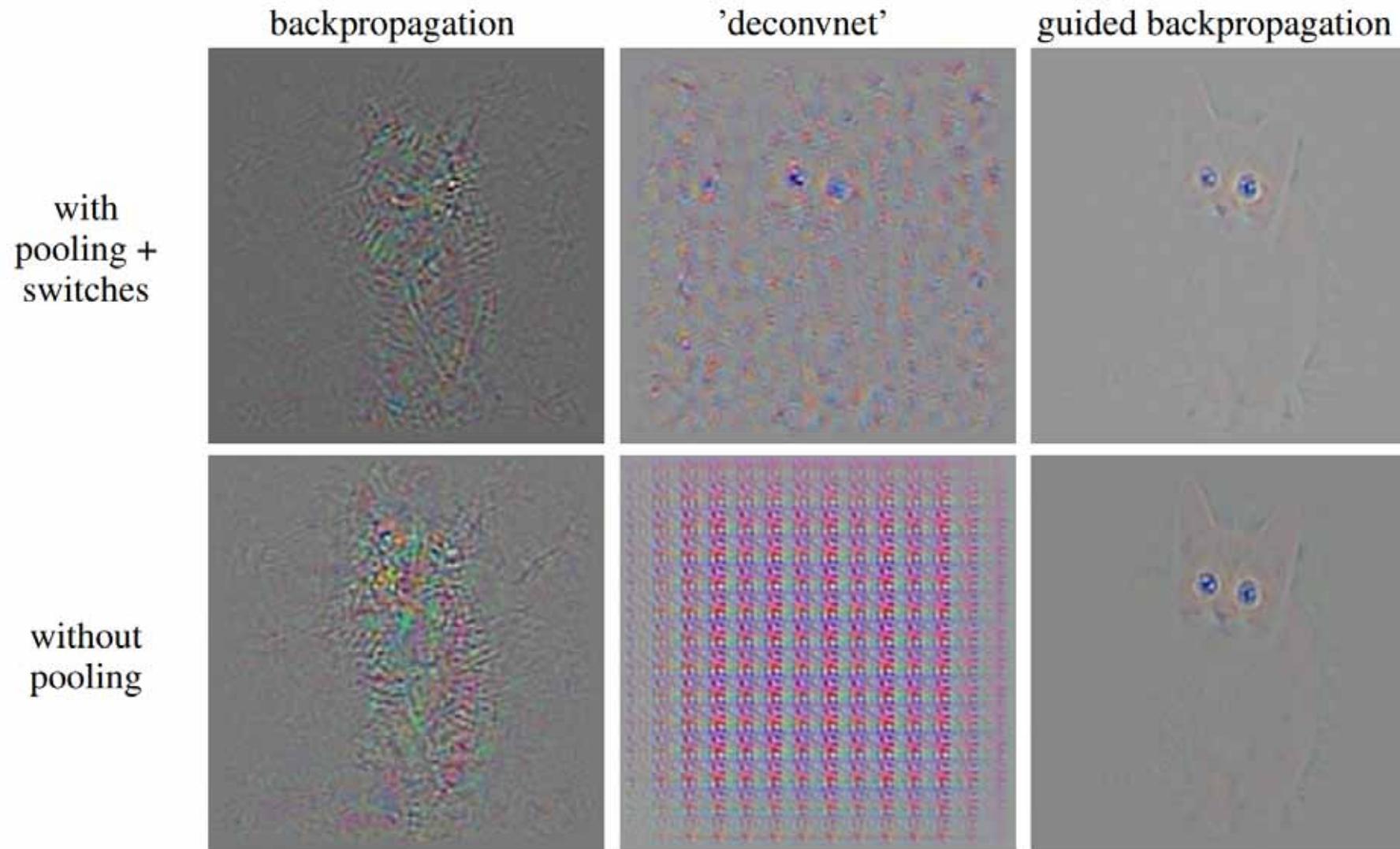
guided backpropagation



corresponding image crops



Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox & Martin Riedmiller 2014. Striving for simplicity:  
The all convolutional net. arXiv:1412.6806.



Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox & Martin Riedmiller 2014. Striving for simplicity:  
The all convolutional net. arXiv:1412.6806.

---

# 04 Deep Generator Networks (DGN)

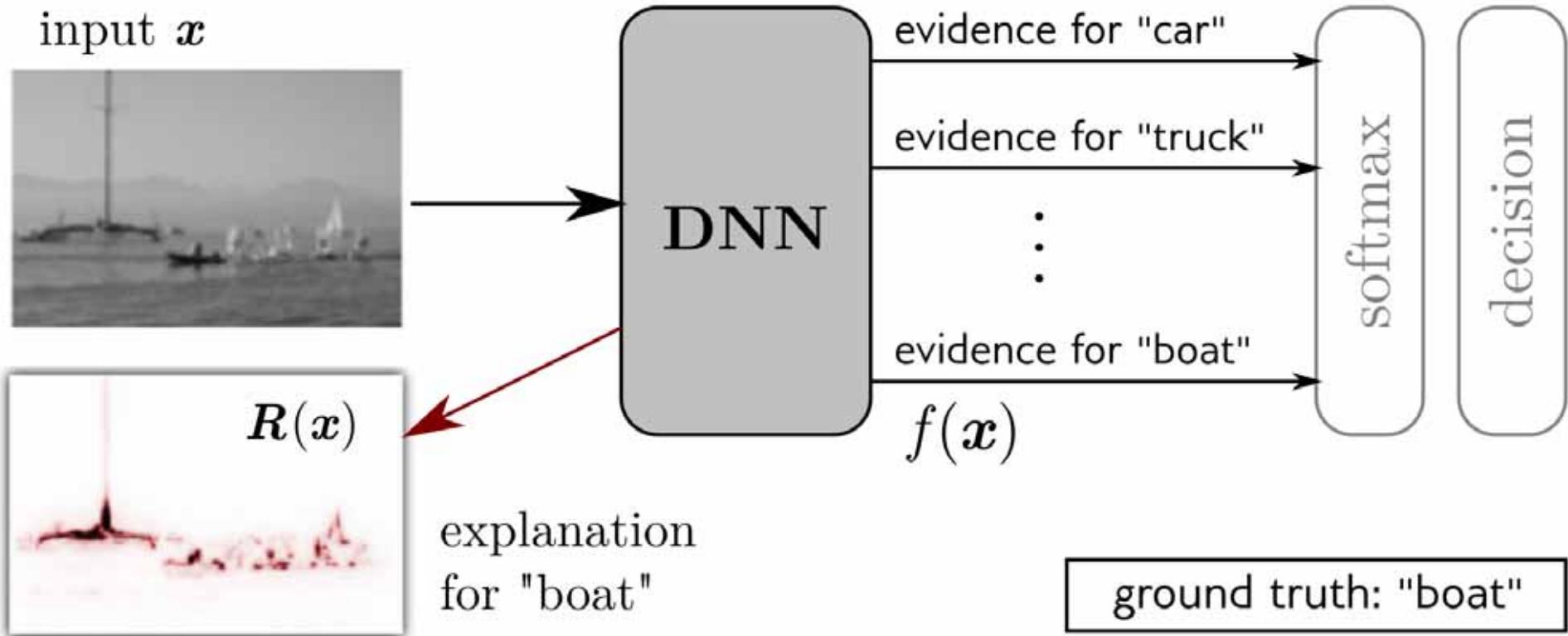
# Learning generative models of natural images



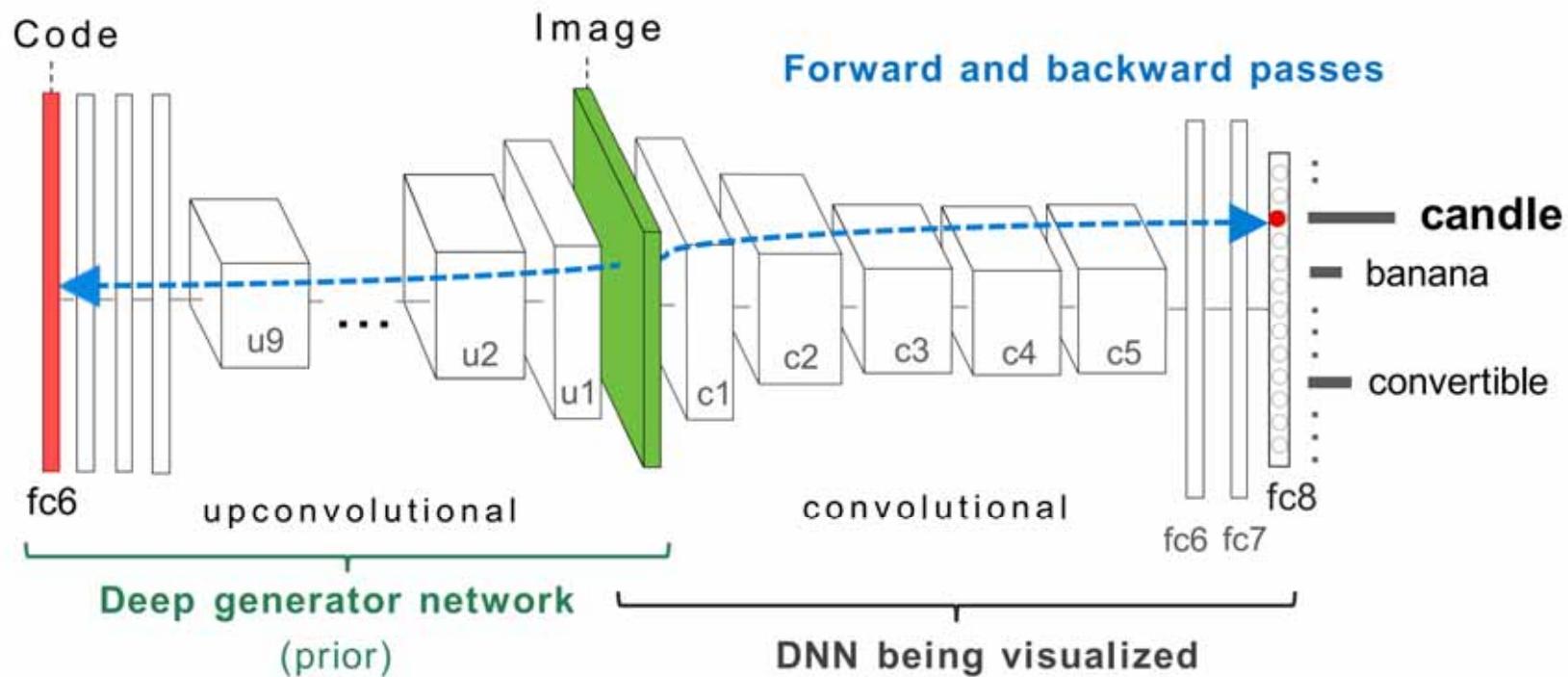
- Idea: incorporating a generative model in the activation maximization framework and redefine the optimization problem as:

$$\max_{\mathbf{z} \in \mathcal{Z}} \log p(\omega_c | g(\mathbf{z})) - \lambda \|\mathbf{z}\|^2$$

Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox & Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. Advances in Neural Information Processing Systems (NIPS 2016), 2016 Barcelona. 3387-3395.



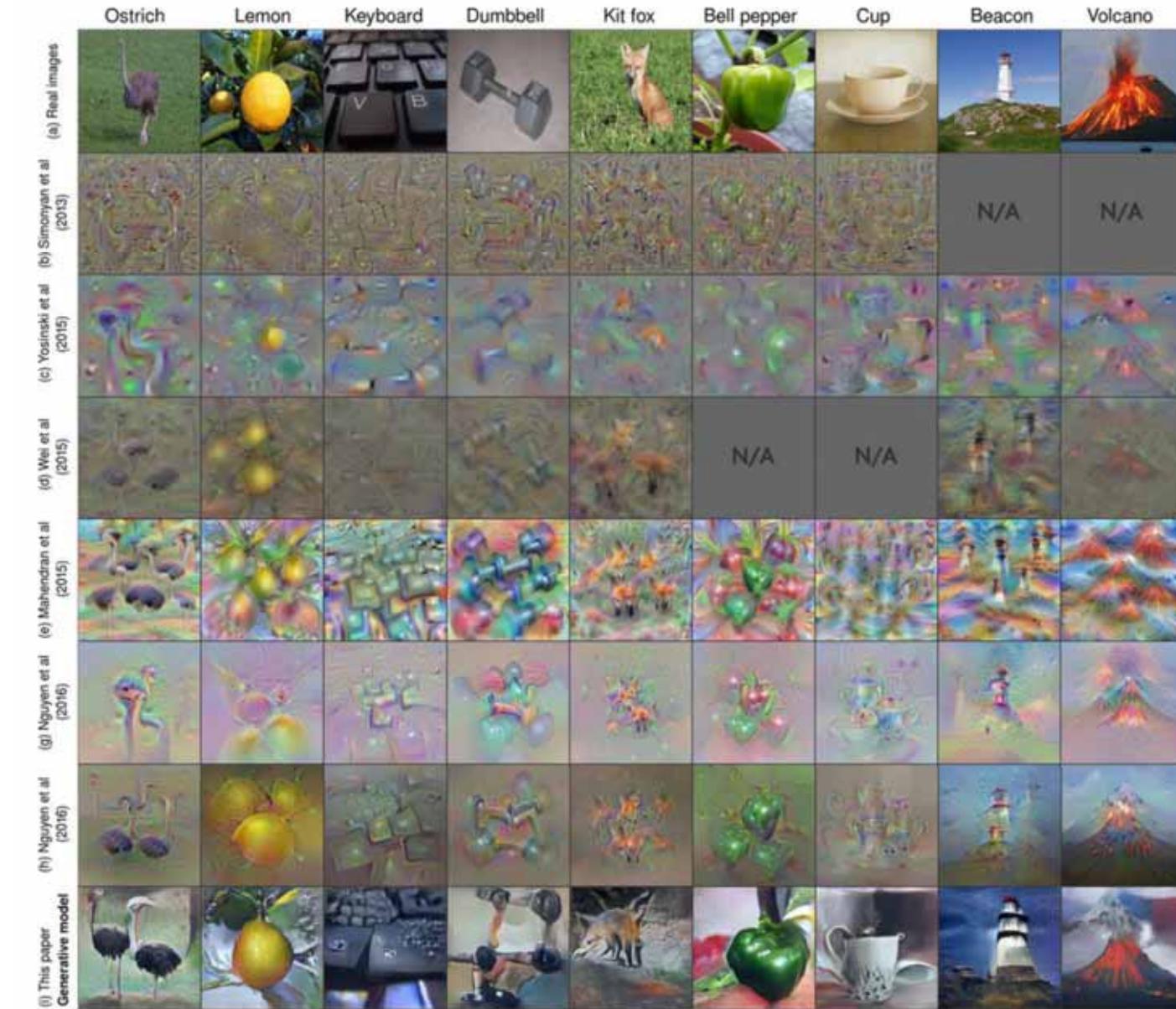
Grégoire Montavon, Wojciech Samek & Klaus-Robert Müller 2018. Methods for interpreting and understanding deep neural networks. Digital Signal Processing, 73, 1-15, doi:10.1016/j.dsp.2017.10.011.



$$\hat{\mathbf{y}}^l = \arg \max_{\mathbf{y}^l} (\Phi_h(G_l(\mathbf{y}^l)) - \lambda \|\mathbf{y}^l\|)$$

Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox & Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. Advances in Neural Information Processing Systems (NIPS 2016), 2016 Barcelona. 3387-3395.

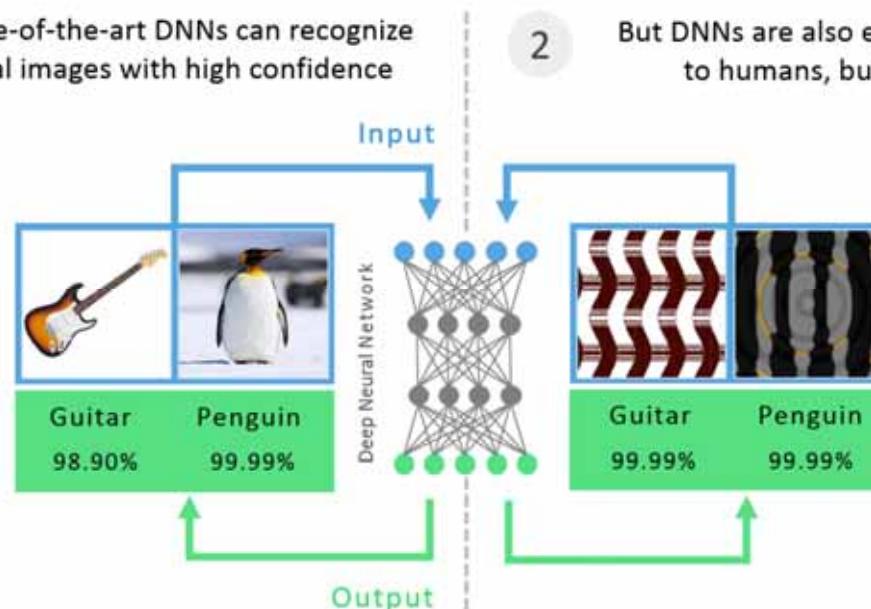
# Comparison of methods to DGN



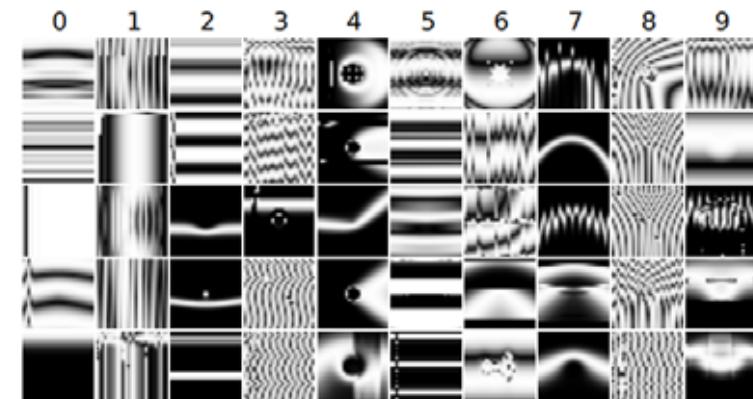
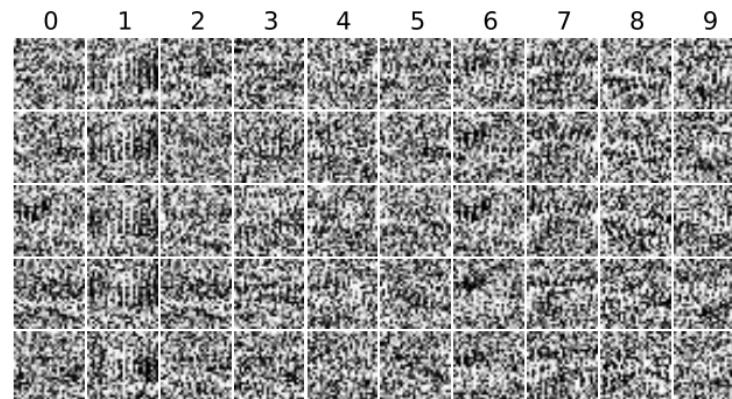
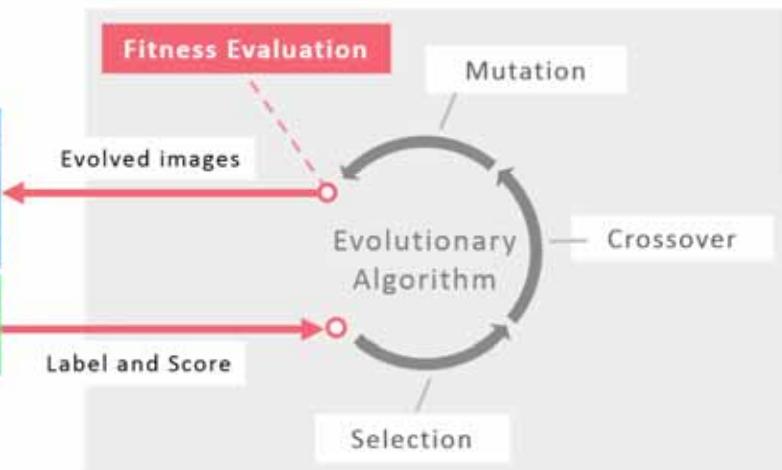
<http://www.evolvingai.org/synthesizing>

# Example for Explanation Needed: adversarials

1 State-of-the-art DNNs can recognize real images with high confidence



2 But DNNs are also easily fooled: images can be produced that are unrecognizable to humans, but DNNs believe with 99.99% certainty are natural objects



Anh Nguyen, Jason Yosinski & Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015 Boston (MA). IEEE, 427-436.

<http://www.evolvingai.org/fooling>

---

# 05 Testing with Concept Activation Vectors

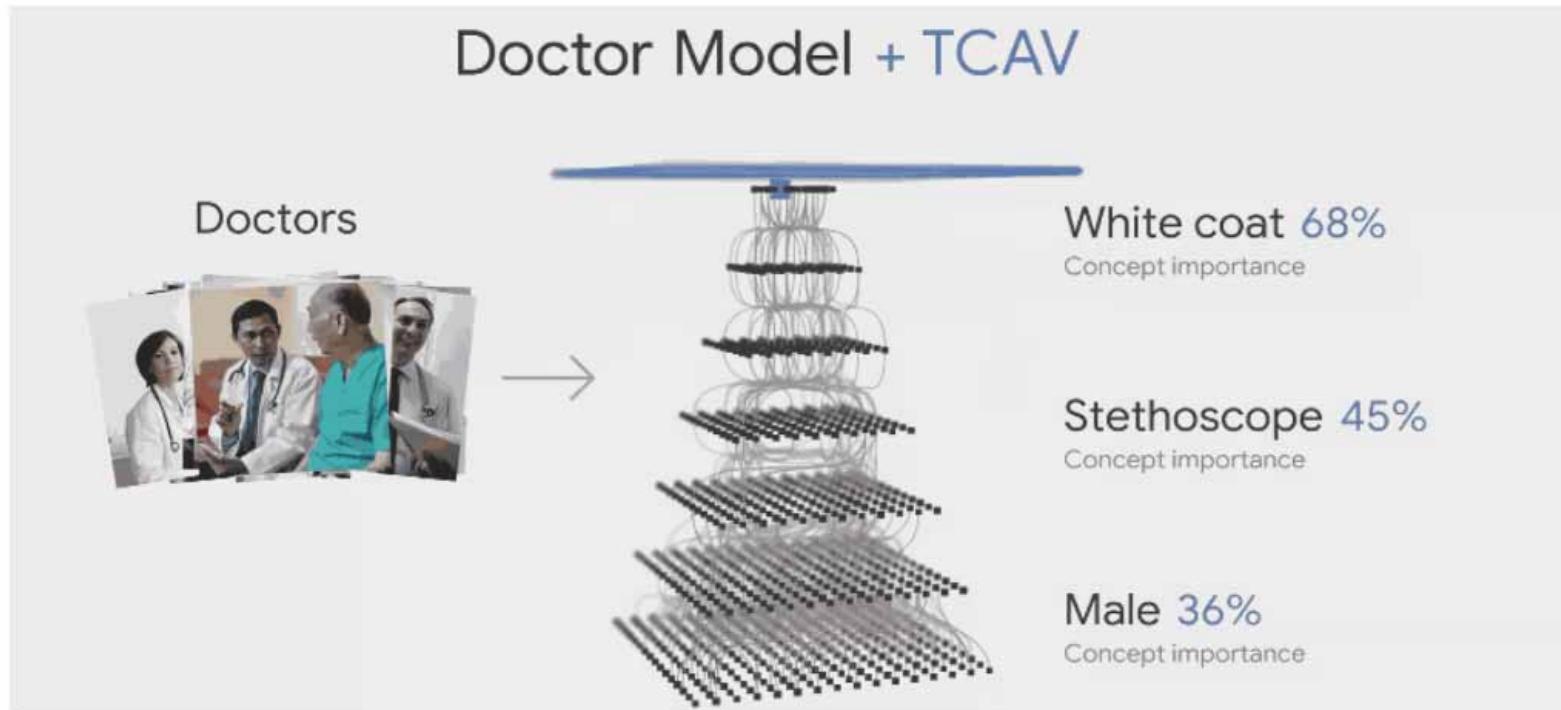
■ “It’s not enough to know if a model works, we need to know how it works”

... if Sundar Pichai is saying this ...



<b>Born</b>	Pichai Sundararajan June 10, 1972 (age 47) Madurai, Tamil Nadu, India
<b>Alma mater</b>	Indian Institute of Technology Kharagpur (BTech) Stanford University (MS) The Wharton School (MBA)
<b>Salary</b>	US\$1,881,066 (2018) US\$1,333,557 (2017) <sup>[1]</sup> US\$199.7 million <sup>[2]</sup> (2016)
<b>Title</b>	CEO of Google and Alphabet
<b>Board member of</b>	Alphabet Inc. <sup>[3]</sup> Magic Leap <sup>[4]</sup>

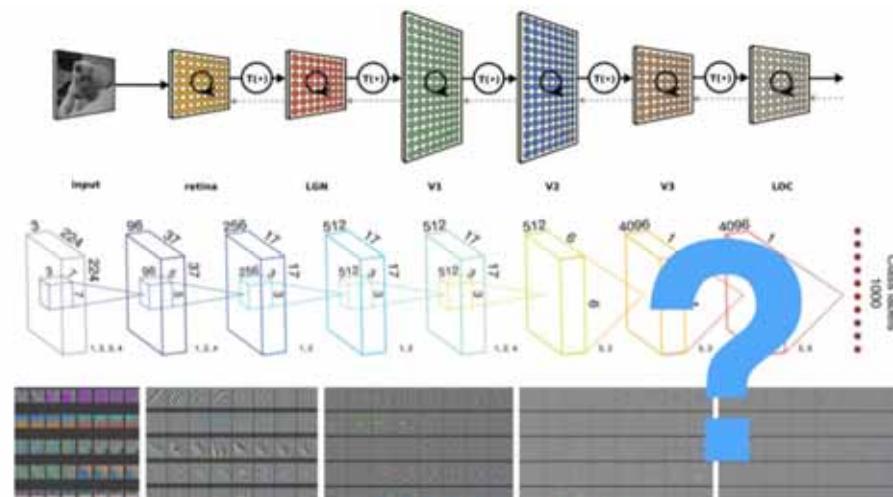
- ML models work on **low-level features** (edges, lines, pixel, ...)
- Humans are working on **high-level concepts** (shape, size, ...)
- Every pixel of an image is an input feature and are numbers, which do not make sense to humans.
- TCAV enables to provide an explanation that is generally true for a class of interest, beyond one image (global explanation).
- TCAV learns ‘concepts’ from examples. For instance, TCAV needs a couple of examples of ‘female’, and something ‘not female’ to learn a “gender” concept. The goal of TCAV is to determine how much a concept (e.g., gender, race) was important for a prediction in a trained model even if the concept was not part of the training.



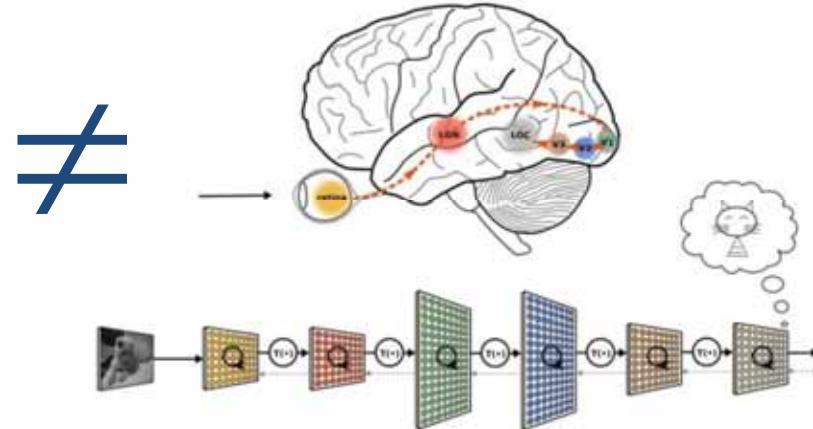
TCAV shows that being male is an important 'concept' when deciding if an image belongs to a doctor or not

<https://towardsdatascience.com/tcav-interpretability-beyond-feature-attribution-79b4d3610b4d>

# Example: Concept Activation Vector (CAV)



Yann Lecun, Yoshua Bengio & Geoffrey Hinton 2015. Deep learning.  
Nature, 521, (7553), 436-444, doi:10.1038/nature14539.



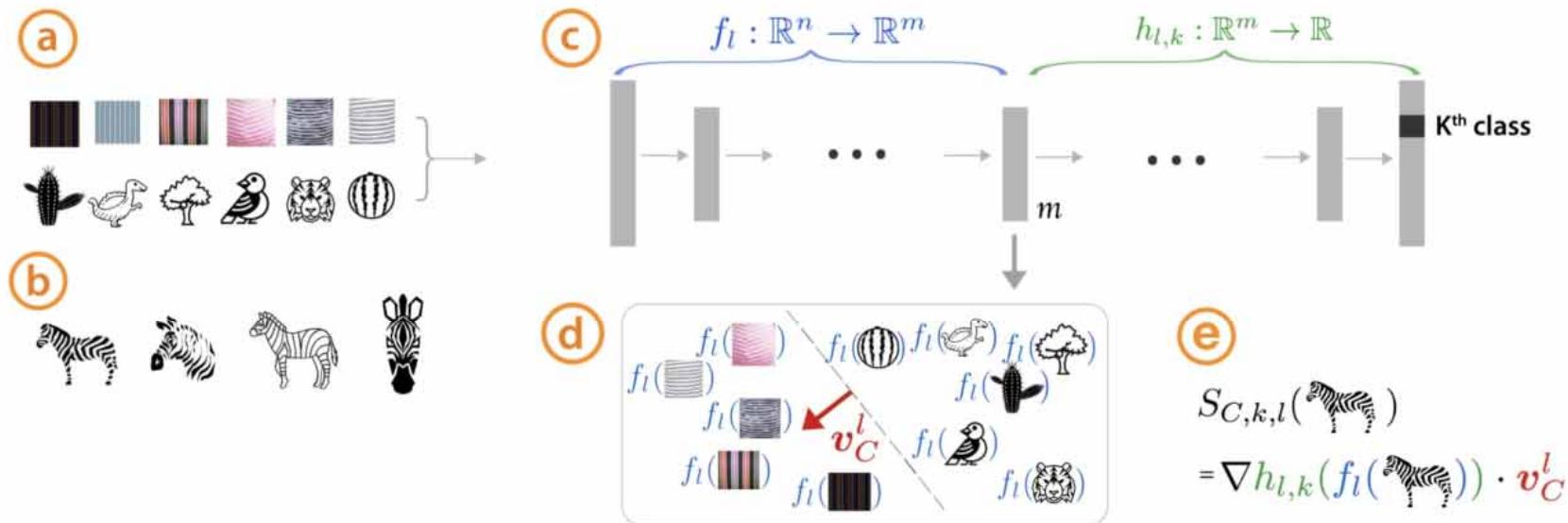
$$\frac{\partial h_k(x)}{\partial x_{a,b}}$$

Humans work in another vector space which is spanned by **implicit knowledge** vectors corresponding to an unknown set of human interpretable concepts.

$$S_{C,k,l}(x) = \lim_{\epsilon \rightarrow 0} \frac{h_{l,k}(f_l(x) + \epsilon v_C^l) - h_{l,k}(f_l(x))}{\epsilon} = \nabla h_{l,k}(f_l(x)) \cdot v_C^l$$

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler & Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors, ICML, 2018. 2673-2682.

## Testing with Concept Activation Vectors (TCAV)



**Figure 1. Testing with Concept Activation Vectors:** Given a user-defined set of examples for a concept (e.g., ‘striped’), and random examples ④, labeled training-data examples for the studied class (zebras) ⑤, and a trained network ③, TCAV can quantify the model’s sensitivity to the concept for that class. CAVs are learned by training a linear classifier to distinguish between the activations produced by a concept’s examples and examples in any layer ④. The CAV is the vector orthogonal to the classification boundary ( $v_C^l$ , red arrow). For the class of interest (zebras), TCAV uses the directional derivative  $S_{C,k,l}(\mathbf{x})$  to quantify conceptual sensitivity ⑥.

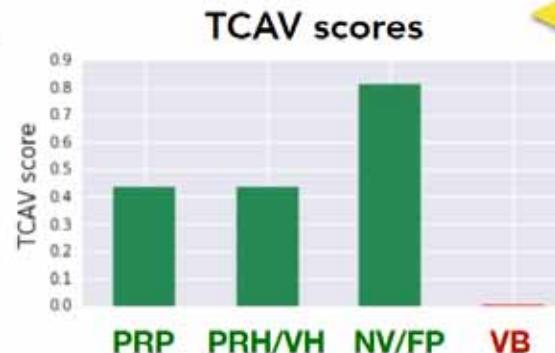
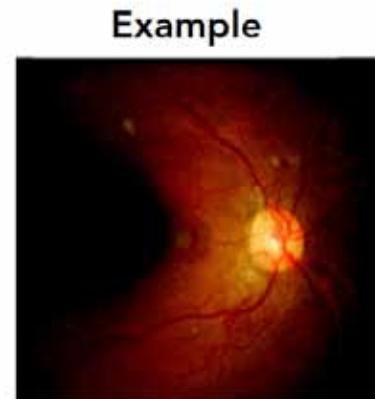
<https://github.com/tensorflow/tcav>

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler & Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). International Conference on Machine Learning (ICML), 2018. 2673-2682.

# TCAV Application example: Diabetic retinopathy

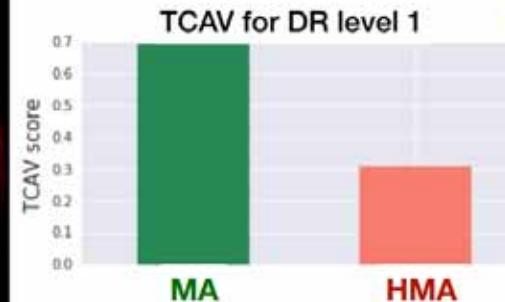
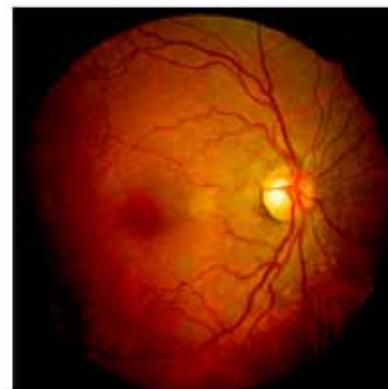
Prediction class      Prediction accuracy

DR level 4      High



TCAV shows the model is **consistent** with doctor's knowledge when model is **accurate**

DR level 1      Med



TCAV shows the model is **inconsistent** with doctor's knowledge for classes when model is less accurate

Green: domain expert's label on concepts belong to the level  
Red: domain expert's label on concepts does not belong to the level

# A controlled experiment with ground truth

Construction of 2 TCAVS:

Image concept

Caption concept

Caption is not always correct:

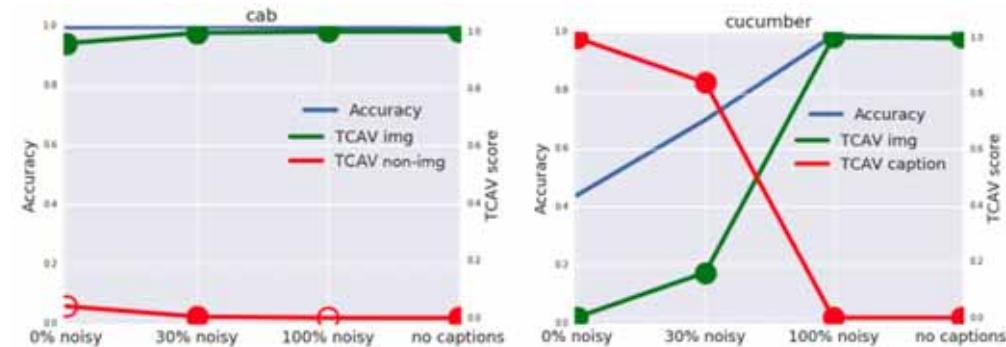
Model with 0% noise: correct caption

Model with 100% noise: always erroneous caption

Which concept is the network using?

TCAVQ closely matches ground truth.

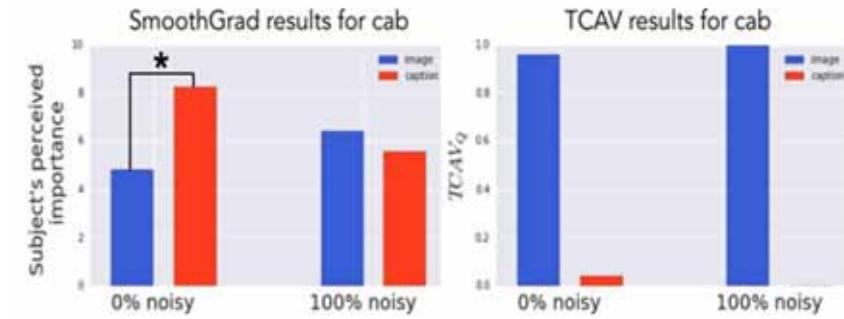
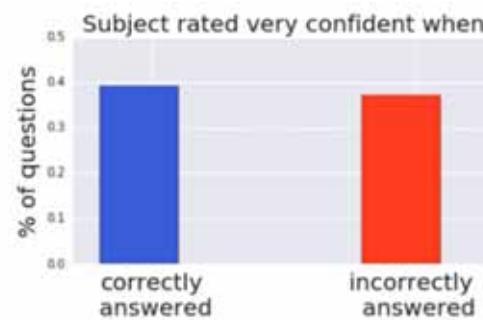
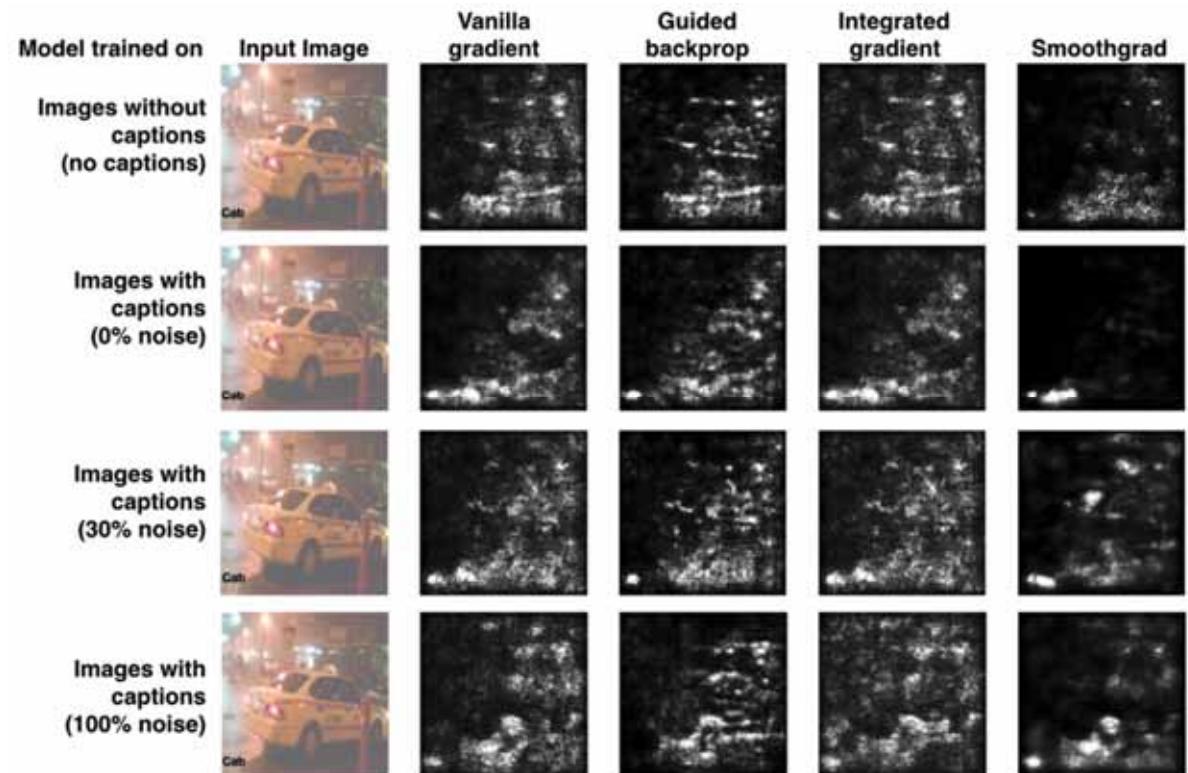
For cab images, the network used image concept more than the caption concept, regardless of % of noise.



Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler & Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). International Conference on Machine Learning (ICML), 2018. 2673-2682.

# A controlled experiment with ground truth

- Saliency map can not focus on caption ->
- Saliency map should focus on caption ->
- Saliency map possibly can focus on caption ->
- Saliency map should not focus on caption ->



Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler & Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). International Conference on Machine Learning (ICML), 2018. 2673-2682.

- TCAV correct – image concept always important
- Note: Humans are very confident – even when they are wrong!
- Conclusion: Saliency maps may be misleading!
- Thank you – see you in the next lecture!