

# Types of Introspection

Which **types** of introspection techniques can you distinguish?

- Find at least one categorization scheme! (more are possible)
- Give as many **examples** as you can per category!

global

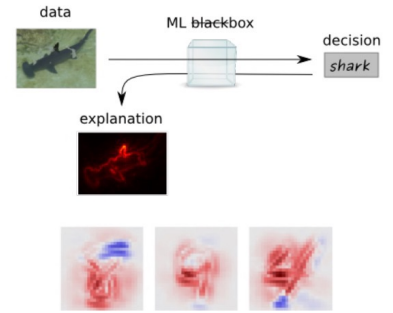
feature visualization  
e.g. DeepDream



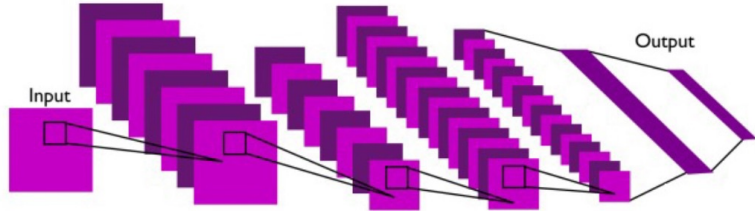
local  
(depending on input)

saliency maps

e.g. layer-wise relevance propagation (LRP)



## Feature Visualization

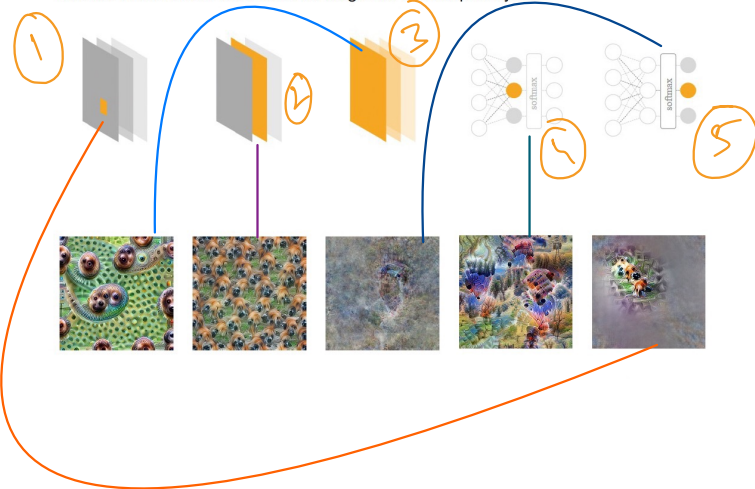


feature visualization by optimization  
(find the input that optimizes a particular part of the network)

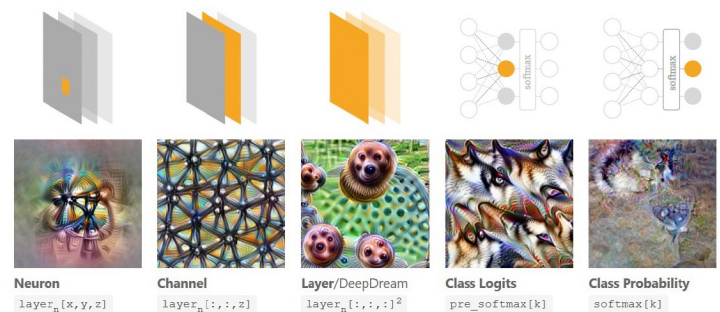
Task 3

## Feature Visualization

Assign each image to its optimization target!  
Bonus: Which class was used as target for the output layer?



## Feature Visualization



<https://distill.pub/2017/feature-visualization/>

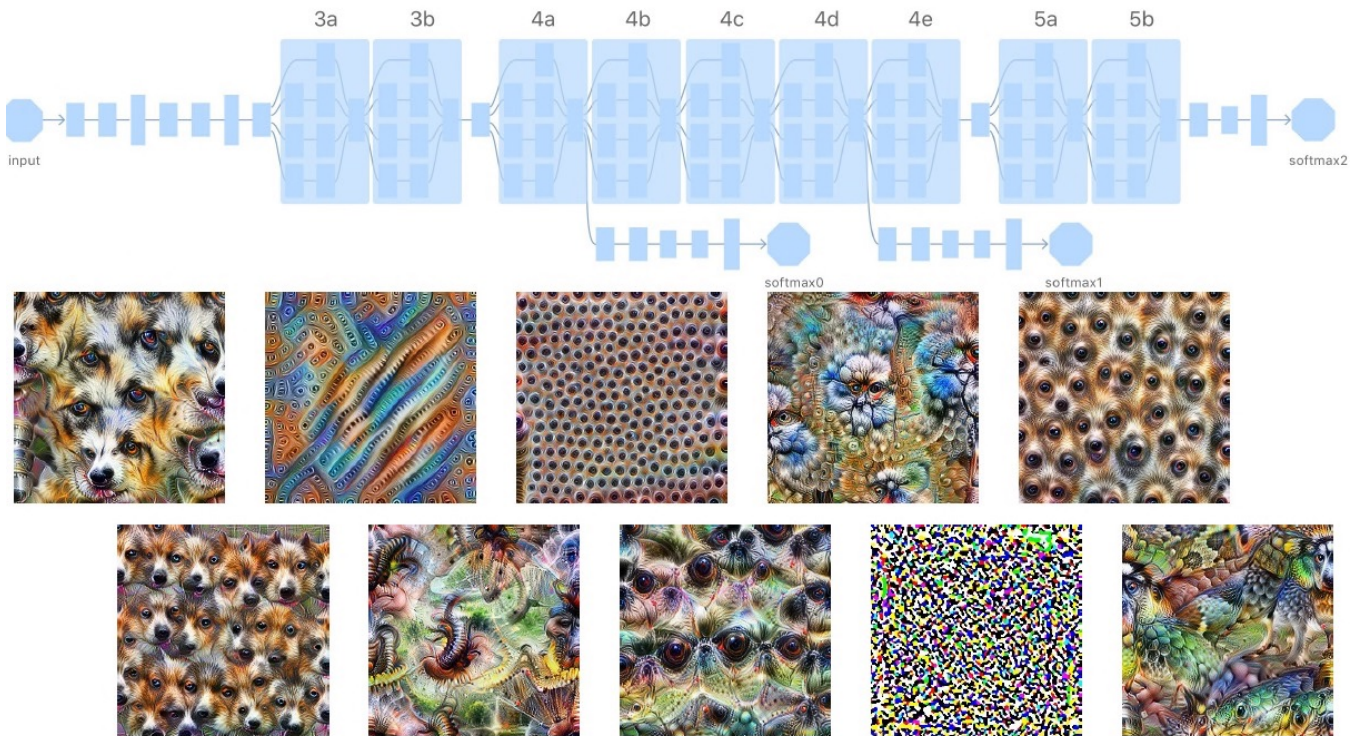
- ① → in the beginning the receptive field is small
- ② → we are looking at one Filter convolving, so we can see so many things repeated.
- ③ → we are trying to see an abstract of the entire channels depth (for e.g. filters = 64. Then we can see some repeated info but not everything.
- ④ It doesn't consider other neurons before softmax.
- ⑤ It considers other Neurons after softmax.



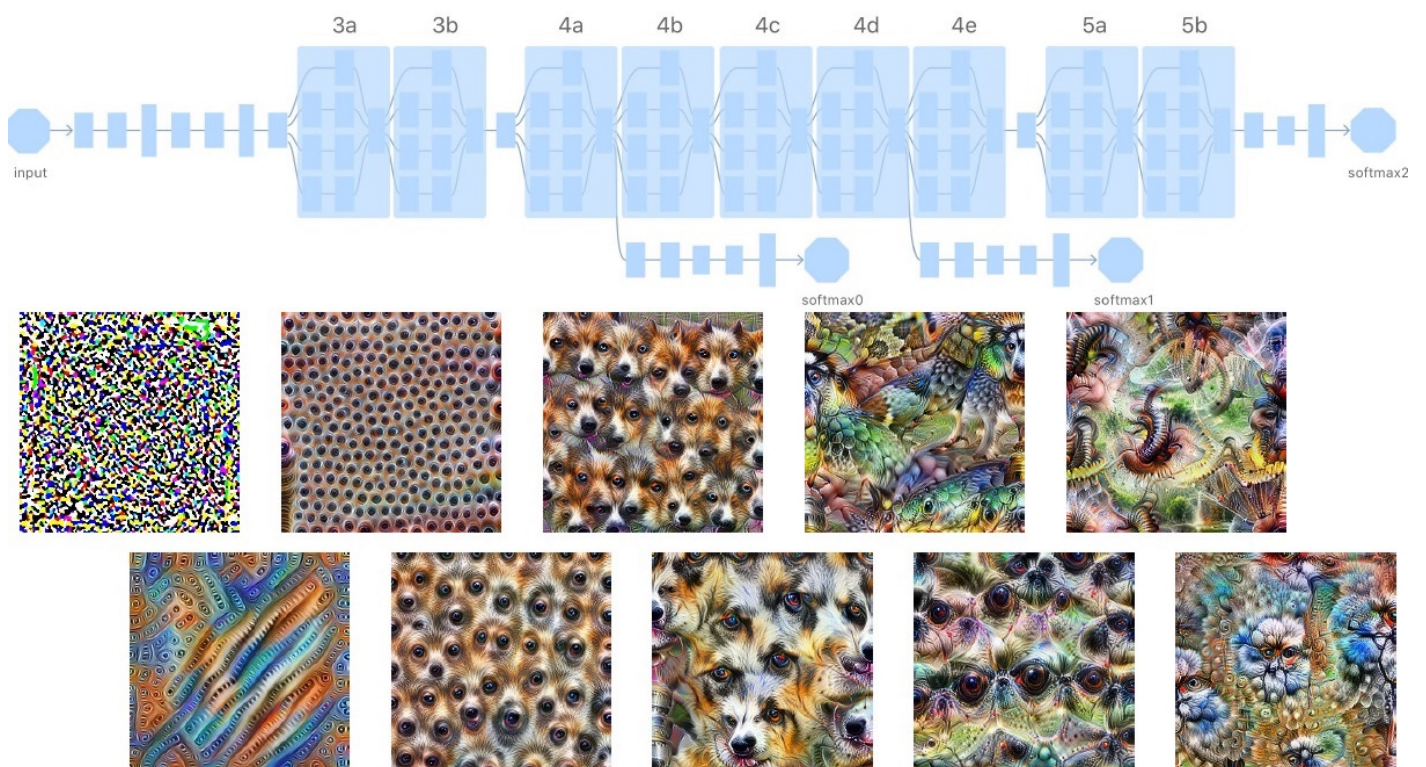
## Task 4

# Feature Visualization: DeepDream

Which layers were used as optimization target to generate these images?



# Feature Visualization: DeepDream



# Feature Visualization

What's the main problem with the (vanilla) optimization approach?



How do we solve this problem?

## Feature Visualization

What's the main problem with the (vanilla) optimization approach?

unregularized optimization is unnatural



How do we solve this problem?

by regularizing the optimization

frequency  
penalization

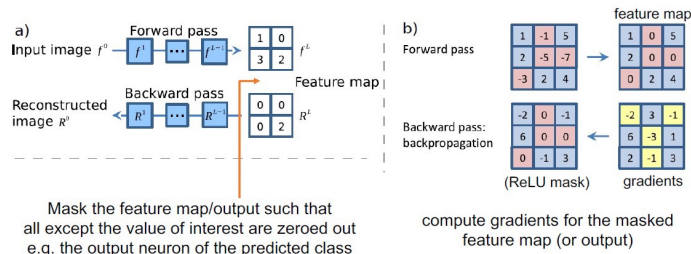
transformation  
robustness

learned  
prior



# gradient-based Saliency Maps

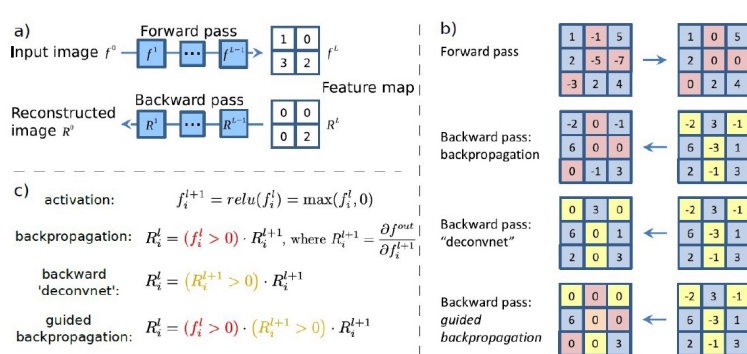
easiest method: using backpropagation values (i.e. the gradients) directly



Propagating back to the input gives a saliency map. Each position tells how sensitive the value of interest is to changes in this position. Hence the name **sensitivity analysis**.

[Springenberg et al. (2014). Striving for Simplicity: The All Convolutional Net]

# gradient-based Saliency Maps



[Springenberg et al. (2014). Striving for Simplicity: The All Convolutional Net]

## Deep Taylor Decomposition and LRP

What's the difference?

deep Taylor decomposition

layer-wise relevance propagation

$$R_d^{(1)} = (x - x_0)_{(d)} \cdot \frac{\partial f}{\partial x_{(d)}}(x_0)$$

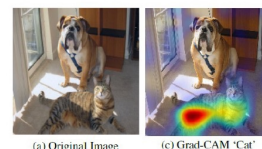
$$R_{i-j}^{(l+1)} = \frac{z_{ij}}{z_j} \cdot R_j^{(l+1)}$$

- root point  $x_0$  must be determined
- computationally efficient (backprop)

- no root point needed
- computationally expensive

## GradCAM: Gradient-weighted Class Activation Mapping

## GradCAM: Gradient-weighted Class Activation Mapping



$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$$

global average pooling

gradients via backprop

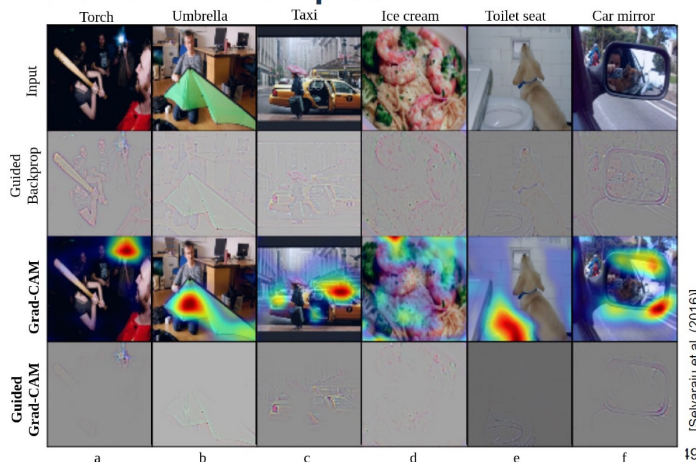
$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left( \sum_k \alpha_k^c A^k \right)$$

linear combination

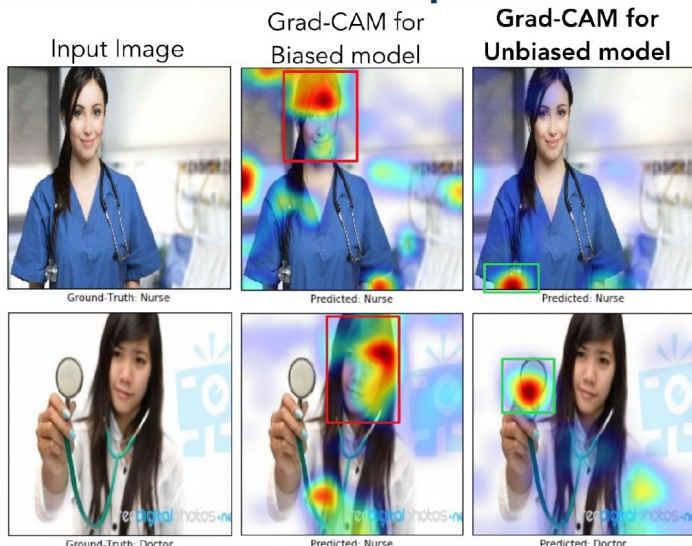
gradient (sensitivity) of class  $c$  to changes in feature map  $A^k$  averaged over all positions

combine all feature maps  $A^k$  in one layer as weighted sum using  $\alpha_k^c$  as weight

## GradCAM: Examples



## GradCAM: Model Comparison



## Problems

- these methods
  - sometimes require particular architectures (e.g. only 2D-convolution with max-pooling)
  - mostly use ReLUs and a positive input space (which pixels positively influence an output class)
  - are mostly evaluated only for images (visually interpretable)
- not well applicable for
  - other activation functions (allowing negative activation)
  - real-valued input space (negative values)
  - visually hardly interpretable data (e.g. waveforms)