

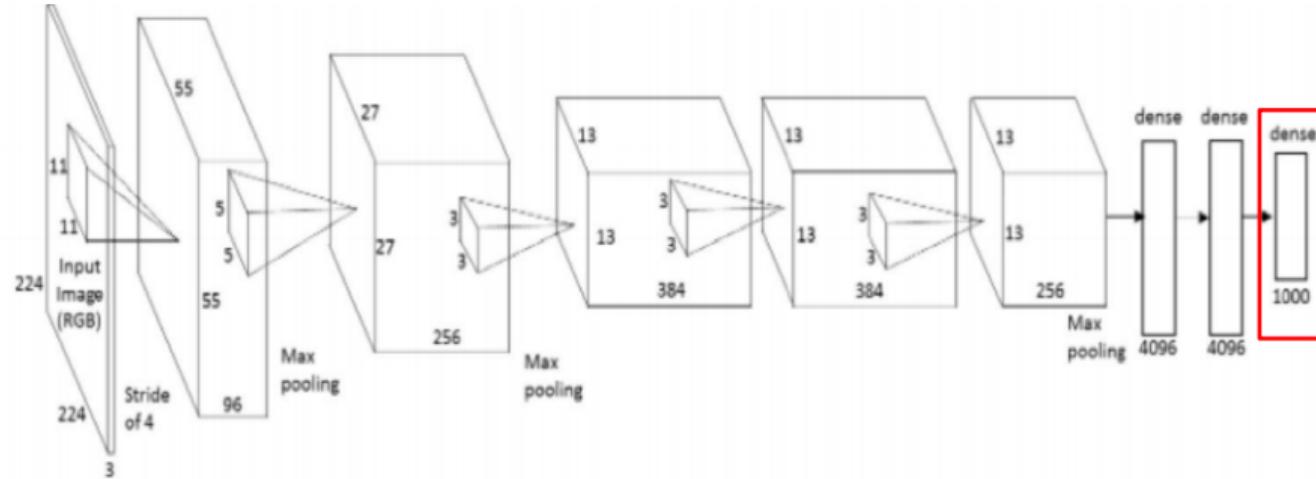
Explaining CNNs: Early Methods

Vineeth N Balasubramanian

Department of Computer Science and Engineering
Indian Institute of Technology, Hyderabad



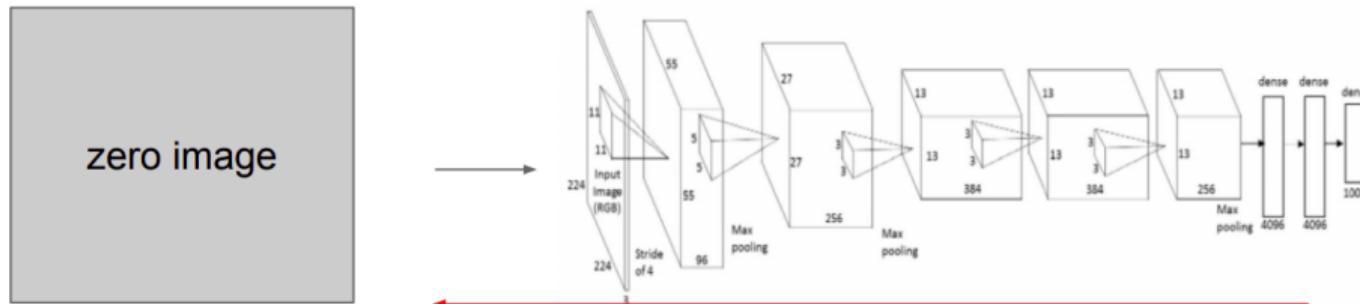
Backpropagation to Image



Question: Can we find an image that maximizes some class score?

Backpropagation to Image¹

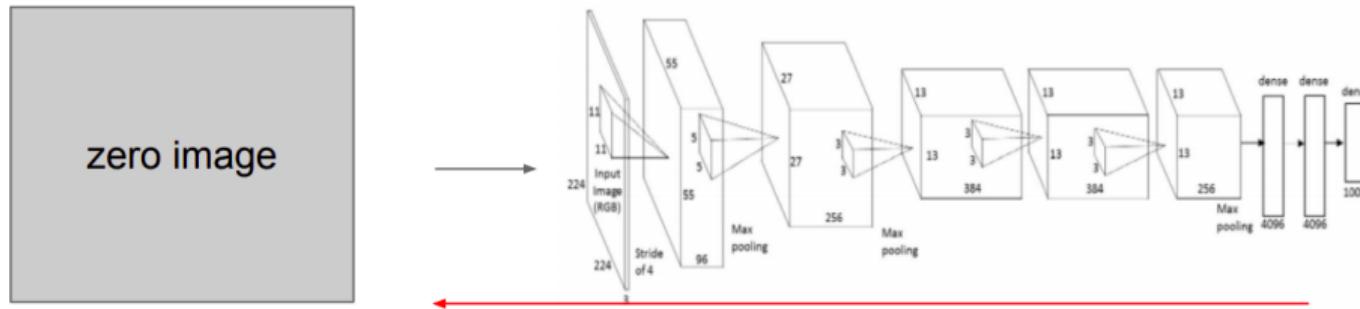
- ① Feed zeros as input.



- ② Set the gradient of the scores vector to be $[0, 0, \dots, 1, \dots, 0]$. Then backprop to image.
- ③ Do a small “**image update**”
- ④ Forward pass the image through the network.
- ⑤ Go back to step 2.

¹Simonyan, Vedaldi, and Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR Workshop 2014

Backpropagation to Image²



- Formally, this optimization can be written as:

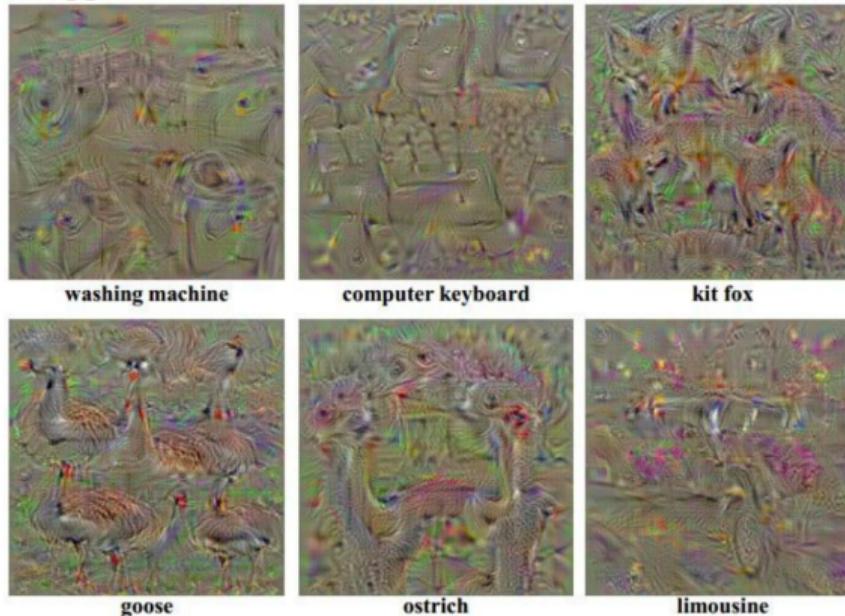
$$\arg \max_I S_c(I) - \lambda \|I\|_2$$

- Here, S_c is the scores vector for class c , before applying softmax.

²Simonyan, Vedaldi, and Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR Workshop 2014

Backpropagation to Image³

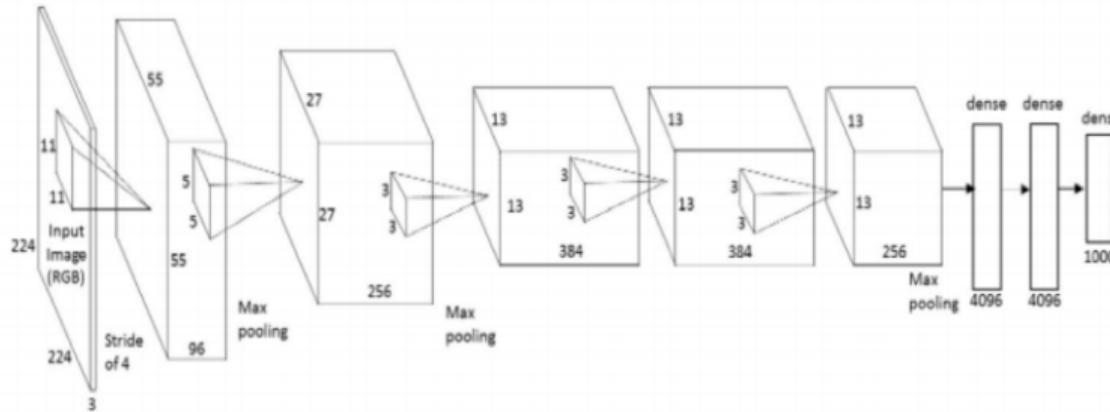
- Finding images that maximize some class score:



³Simonyan, Vedaldi, and Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR Workshop 2014

Backpropagation to Image⁴

- Such optimization can in fact be done for arbitrary neurons in the network



- Repeat:

- Forward an image
- Set activations in a layer of interest to all zero, except 1.0 for a neuron of interest
- Backprop to image
- Do an “image update”

⁴Simonyan, Vedaldi, and Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR Workshop 2014

Visualizing the Data Gradient⁵

- Since the gradient on image data has three channels, visualise M such that:

$$M_{ij} = \max_c |\nabla_I S_c(I)|_{(i,j,c)}$$

- At each pixel, take absolute value and pick maximum across channels



M = ?

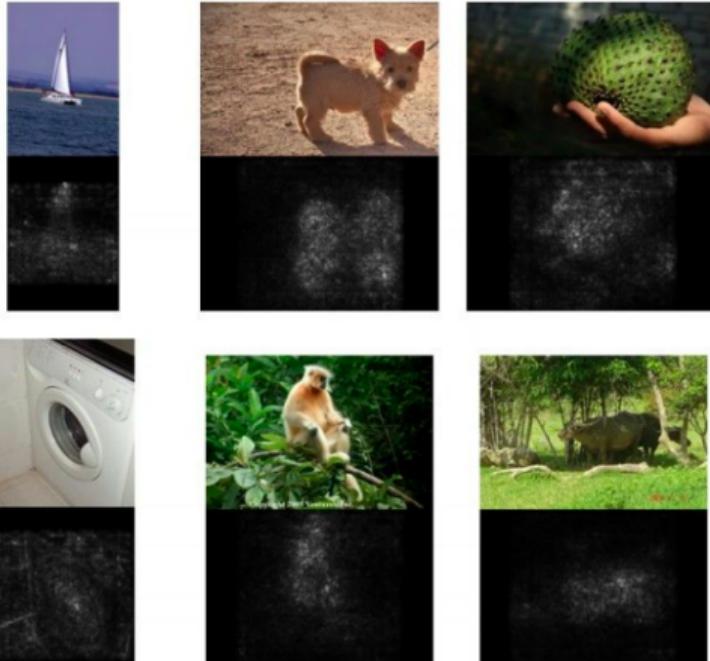
⁵Simonyan, Vedaldi, and Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR Workshop 2014

Visualizing the Data Gradient⁶

- Since the gradient on image data has three channels, visualise M such that:

$$M_{ij} = \max_c |\nabla_I S_c(I)|_{(i,j,c)}$$

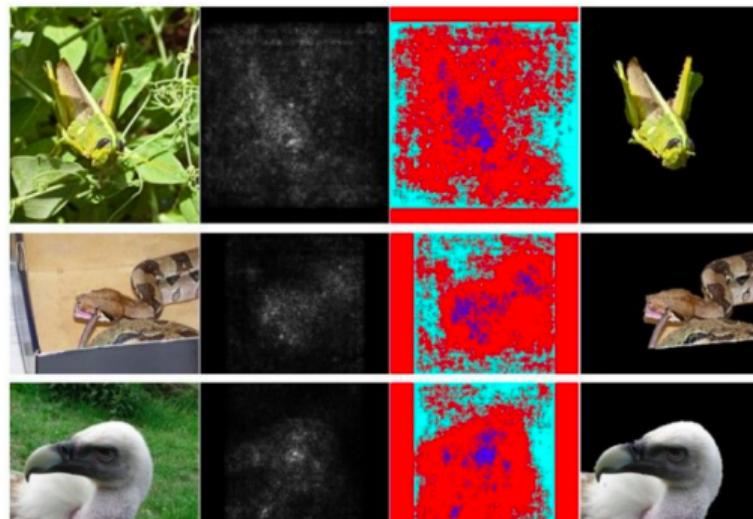
- At each pixel, take absolute value and pick maximum across channels



⁶Simonyan, Vedaldi, and Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR Workshop 2014

Visualizing the Data Gradient⁸

- **GrabCut**⁷, a segmentation method, can be applied to obtain the object mask from the data gradient
 - Recall Graph-Cut segmentation we saw earlier - GrabCut is an extension/adaptation



⁷https://docs.opencv.org/3.4/d8/d83/tutorial_py_grabcut.html

⁸Simonyan, Vedaldi, and Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR Workshop 2014

Image Reconstruction from Latent Representation

Given a CNN code (latent representation from a layer, say, FC7), is it possible to reconstruct the original image?

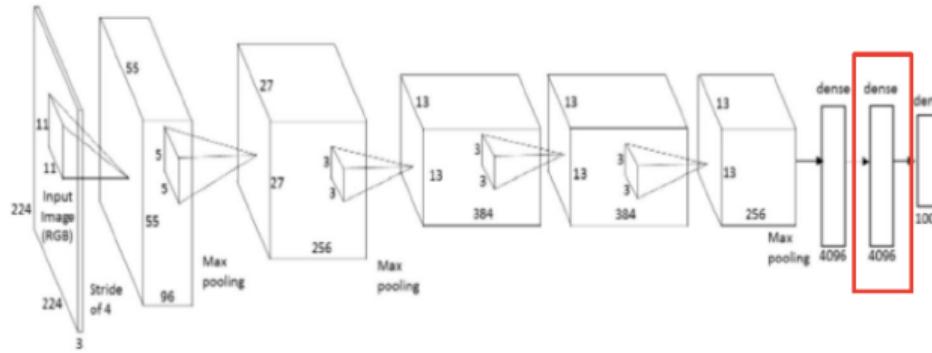
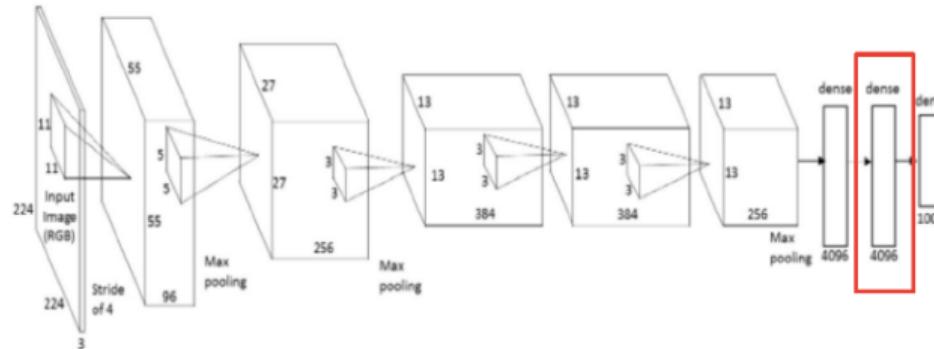


Image Reconstruction from Latent Representation

Given a CNN code (latent representation from a layer, say, FC7), is it possible to reconstruct the original image?



Yes, solve an optimization problem such that:

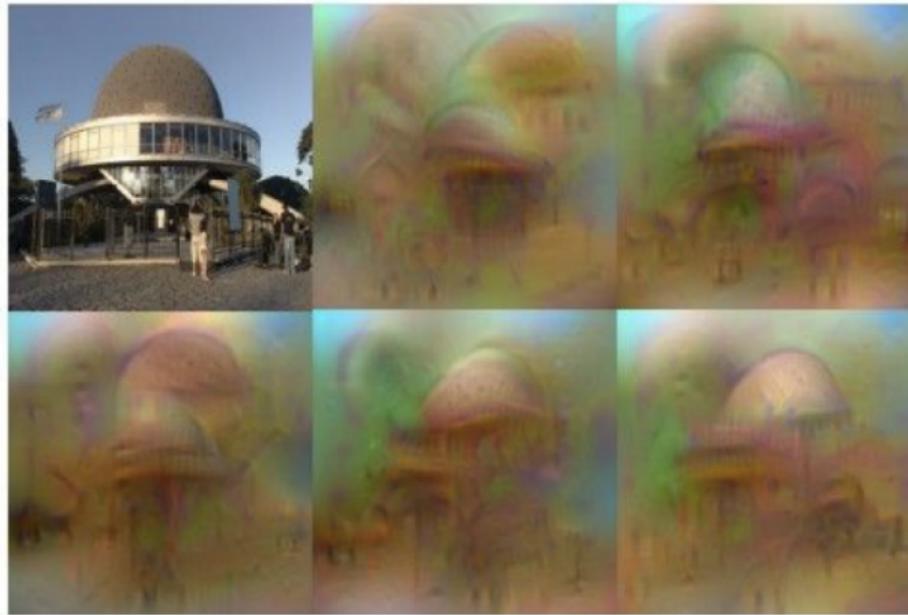
- The image's code is similar to a given code
- It “looks natural” (image prior regularization)

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^{H \times W \times C}} \|\Phi(\mathbf{x}) - \Phi_0\|^2 + \lambda \mathcal{R}(\mathbf{x})$$

Image Reconstruction from Latent Representation

On AlexNet model

original image



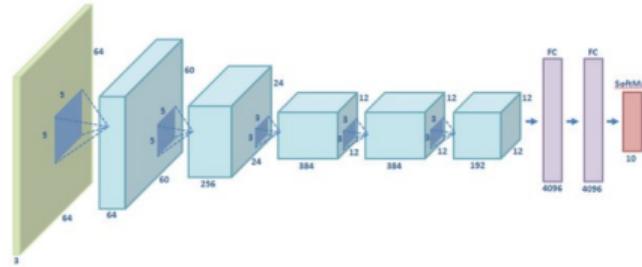
reconstructions
from the 1000
log probabilities
for ImageNet
(ILSVRC)
classes

Image Reconstruction from Latent Representation

Reconstructions from representations after last pooling layer (before first FC layer) in AlexNet



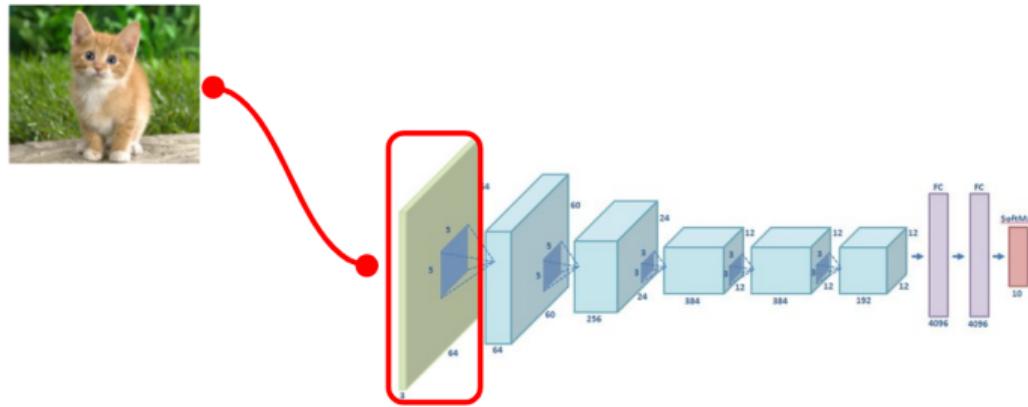
Guided Backpropagation (also known as Deconvolution method)⁹



⁹Springenber et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Guided Backpropagation (also known as Deconvolution method)⁹

a) Feed image into net.



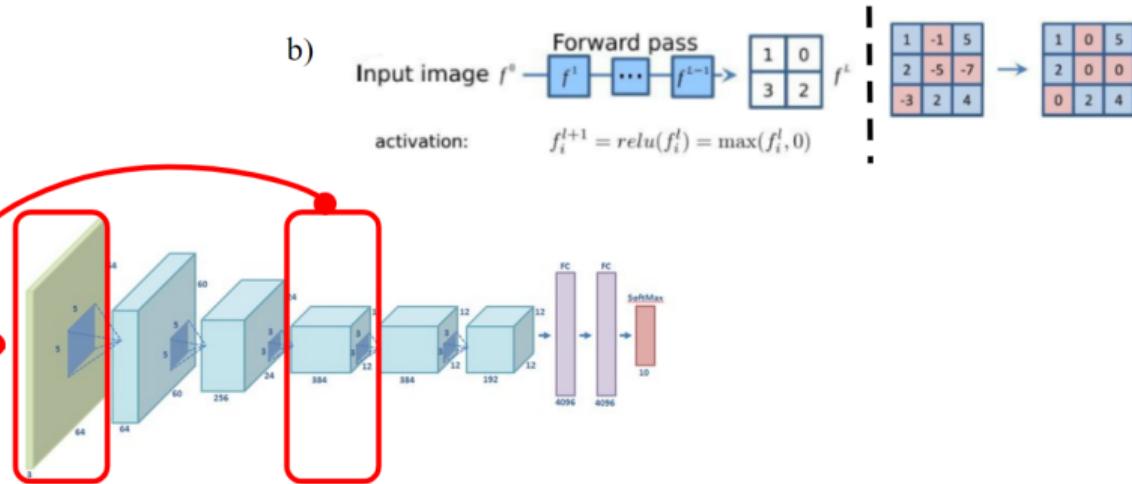
⁹Springenber et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Guided Backpropagation (also known as Deconvolution method)⁹

a) Feed image into net.



b) Pick a layer, set the gradient there to zero except for the neuron of interest.



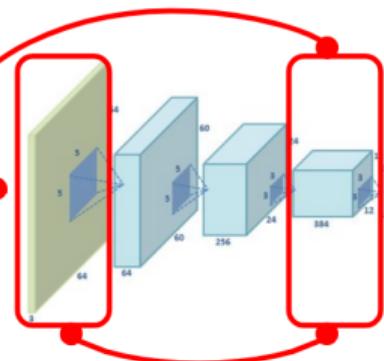
⁹Springenber et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Guided Backpropagation (also known as Deconvolution method)⁹

a) Feed image into net.

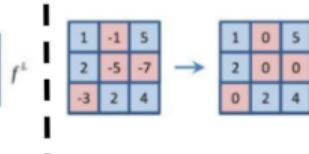


b) Pick a layer, set the gradient there to zero except for the neuron of interest.



b)

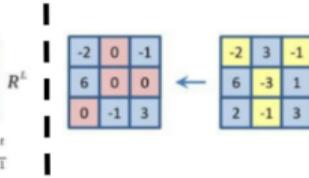
Input image f^0 → Forward pass → f^L
activation: $f_i^{l+1} = \text{relu}(f_i^l) = \max(f_i^l, 0)$



c) Backprop to image.

c)

Reconstructed image R^0 ← Backward pass ← R^L
backpropagation: $R_i^l = (\mathbf{f}_i^l > 0) \cdot R_i^{l+1}$, where $R_i^{l+1} = \frac{\partial f^{\text{out}}}{\partial f_i^{l+1}}$



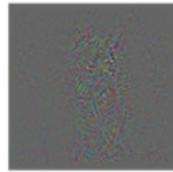
⁹Springenber et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Guided Backpropagation (also known as Deconvolution method)⁹

a) Feed image into net.

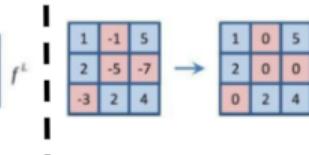
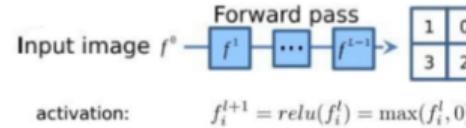


b) Pick a layer, set the gradient there to zero except for the neuron of interest.



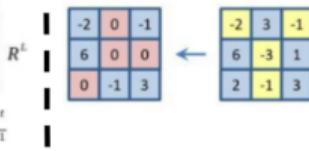
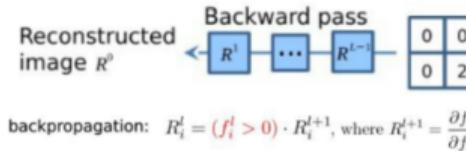
c) Backprop to image.

b)



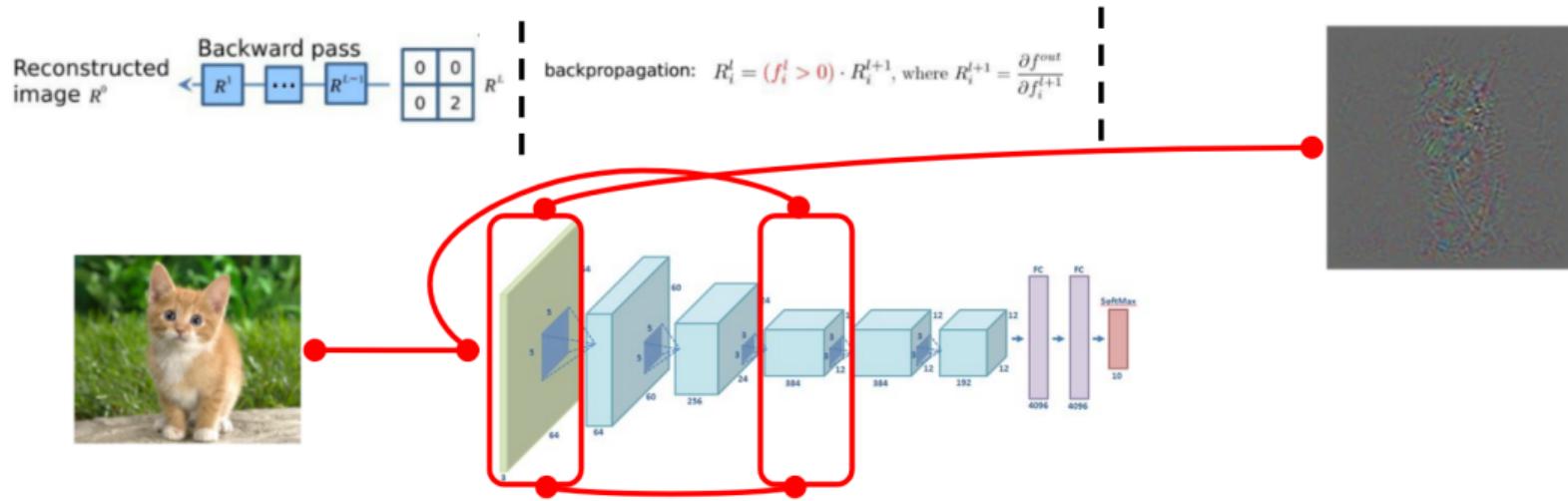
c)

Reconstructed image R^0



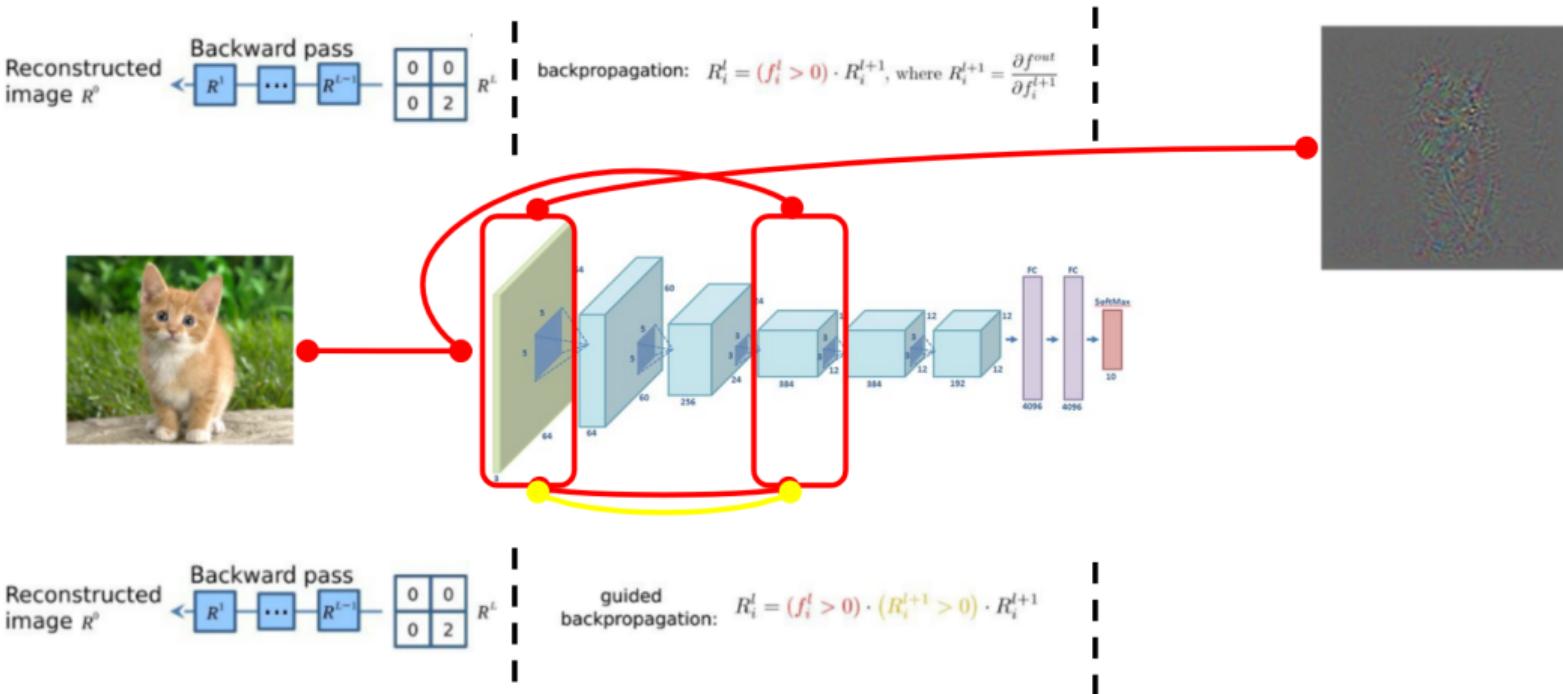
⁹Springenberg et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Guided Backpropagation (also known as Deconvolution method)⁹



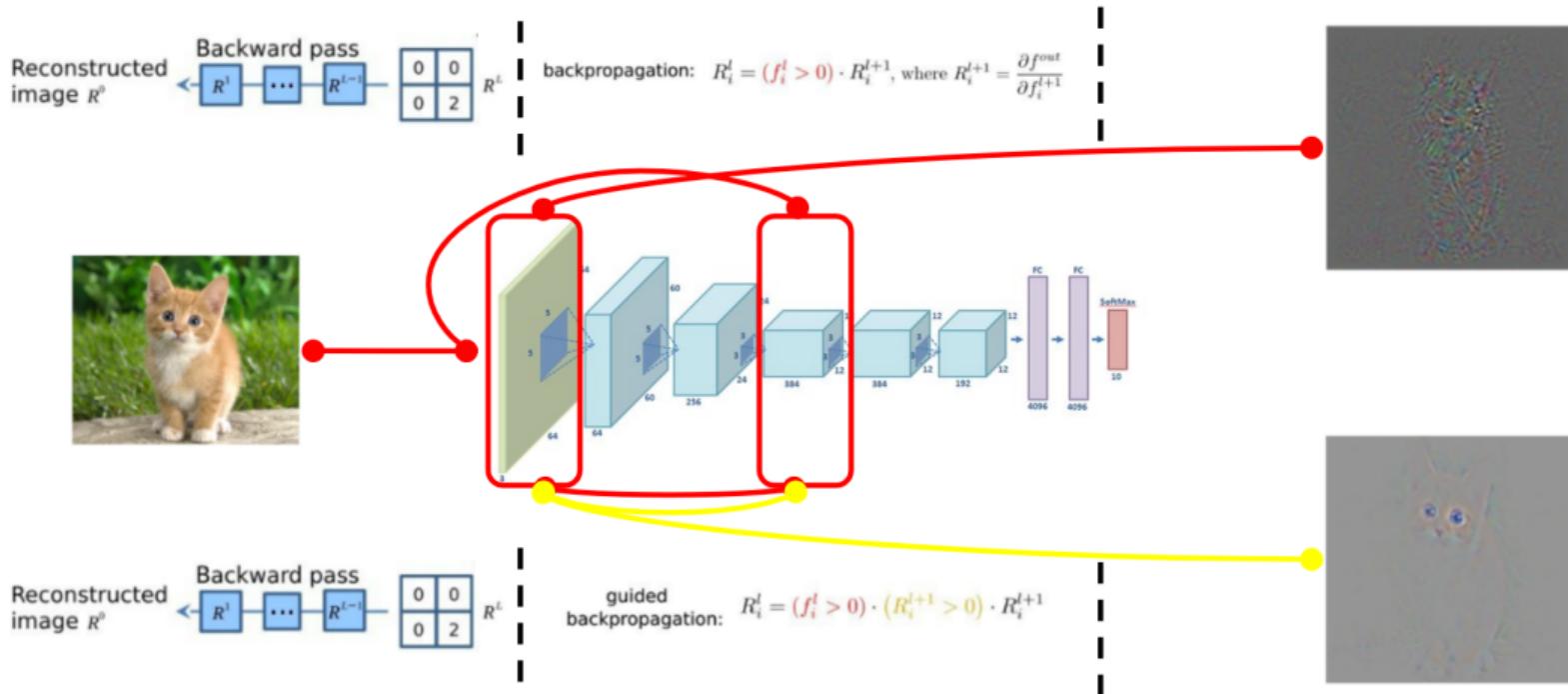
⁹Springenber et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Guided Backpropagation (also known as Deconvolution method)⁹



⁹Springenberg et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Guided Backpropagation (also known as Deconvolution method)⁹



⁹Springenberg et al, Striving for Simplicity: The All Convolutional Net, ICLR Workshop 2015

Readings

Summary of Visualizing CNNs

- [Lecture Notes of CS231n, Stanford](#)

Other Recommended Readings/References

- Simonyan, Vedaldi, and Zisserman, [Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps](#), ICLR Workshop 2014
- Zeiler and Fergus, [Visualizing and Understanding Convolutional Networks](#), ECCV 2014
- Springenberg et al, [Striving for Simplicity: The All Convolutional Net](#), ICLR Workshop 2015

Exercises

- [Understand GrabCut](#) and how it can be used to generate masks from data gradients