

Understanding visual representations

Andrea Vedaldi

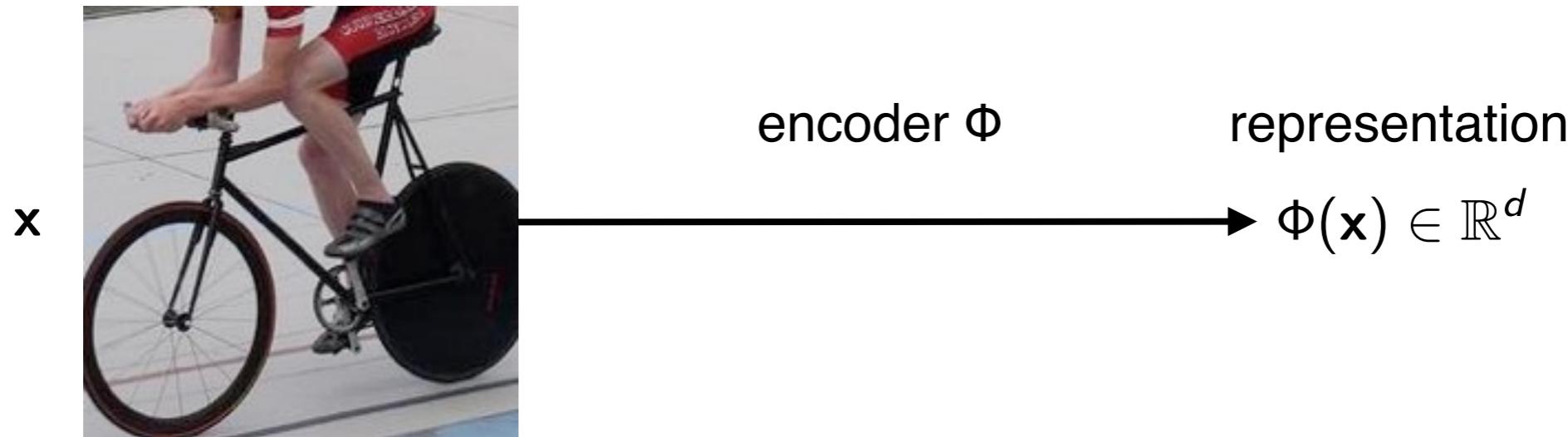
Workshop on Deep Learning: Theory and Practice

July 2016



UNIVERSITY OF
OXFORD

Image representations



An **encoder** maps the data into a **vectorial representation**

Facilitate labelling of images, text, sound, videos, ...

Understanding visual representations

3

Intro

Visualizing representations

Backpropagation networks and “deconvolution”

Representations: equivalence & transformations

Intro

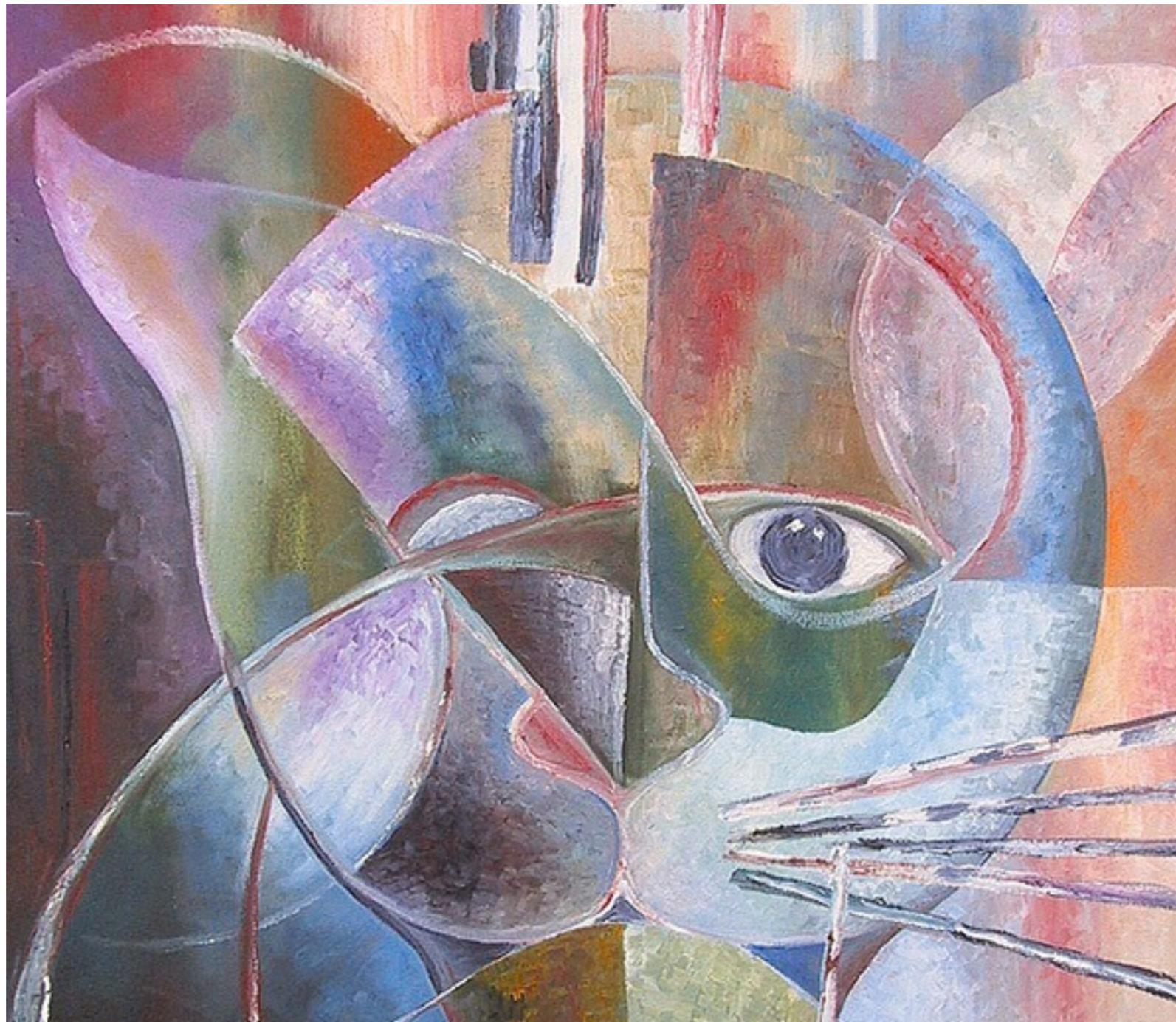
Visualizing representations

Backpropagation networks and “deconvolution”

Representations: equivalence & transformations

Convolutional neural networks

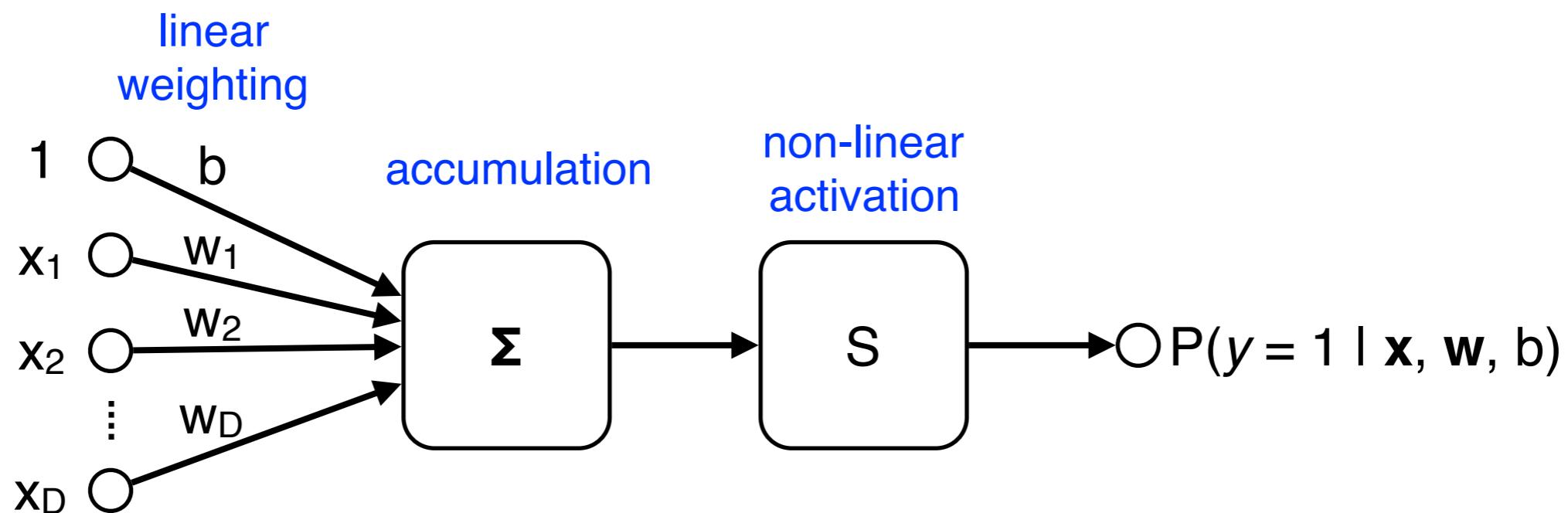
Origin (1950-60)



Perceptron

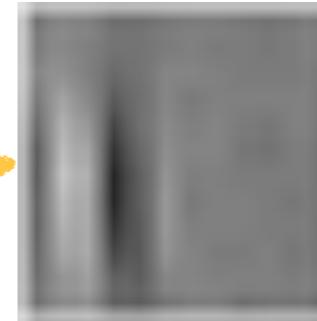
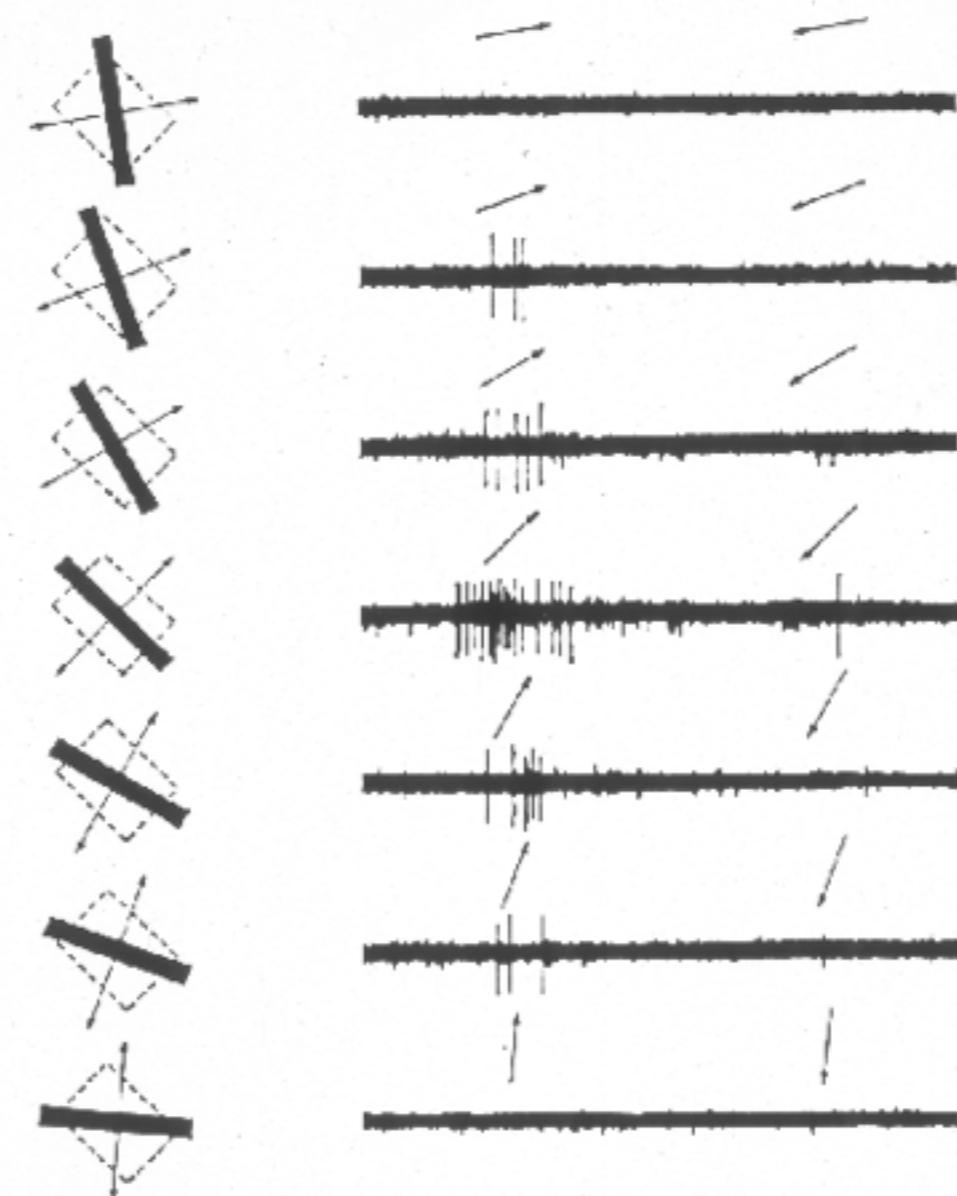
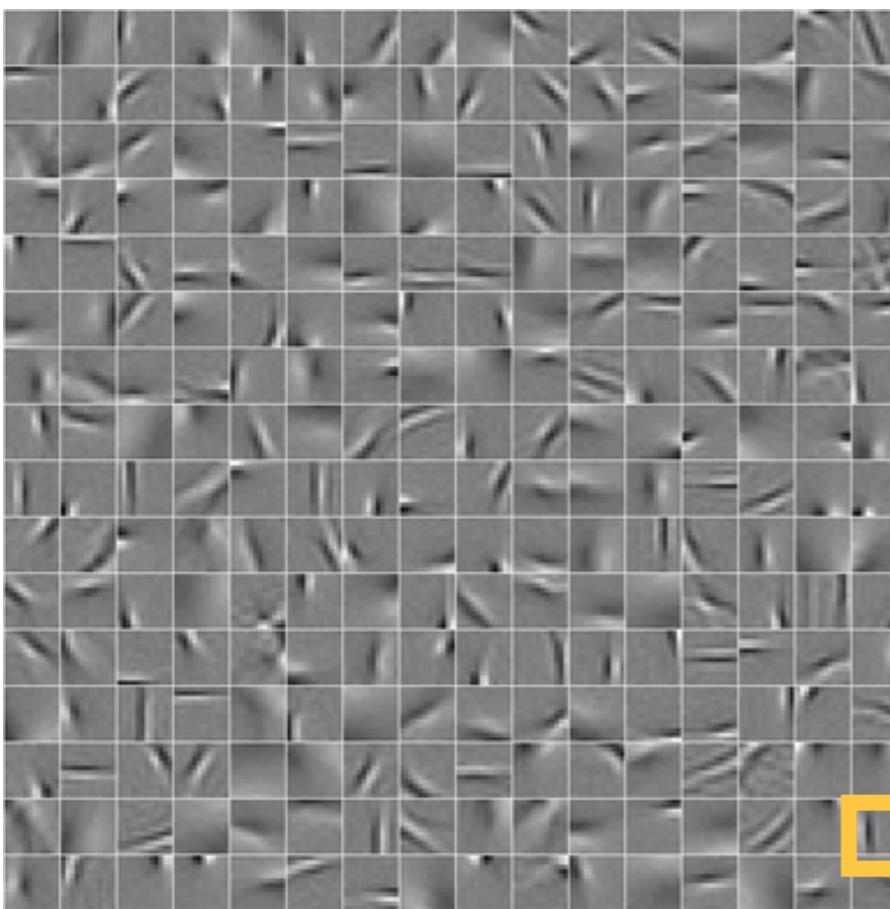
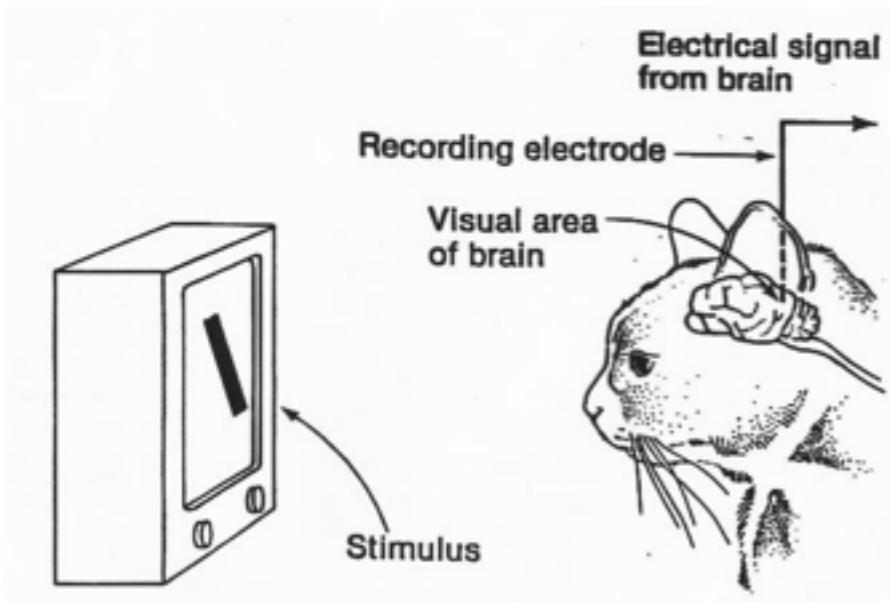
[Rosenblatt 57]

The goal is estimating the posterior probability of the binary label y of a vector \mathbf{x} :



Discovery of oriented cells in the visual cortex

[Hubel and Wiesel 59]



oriented filter

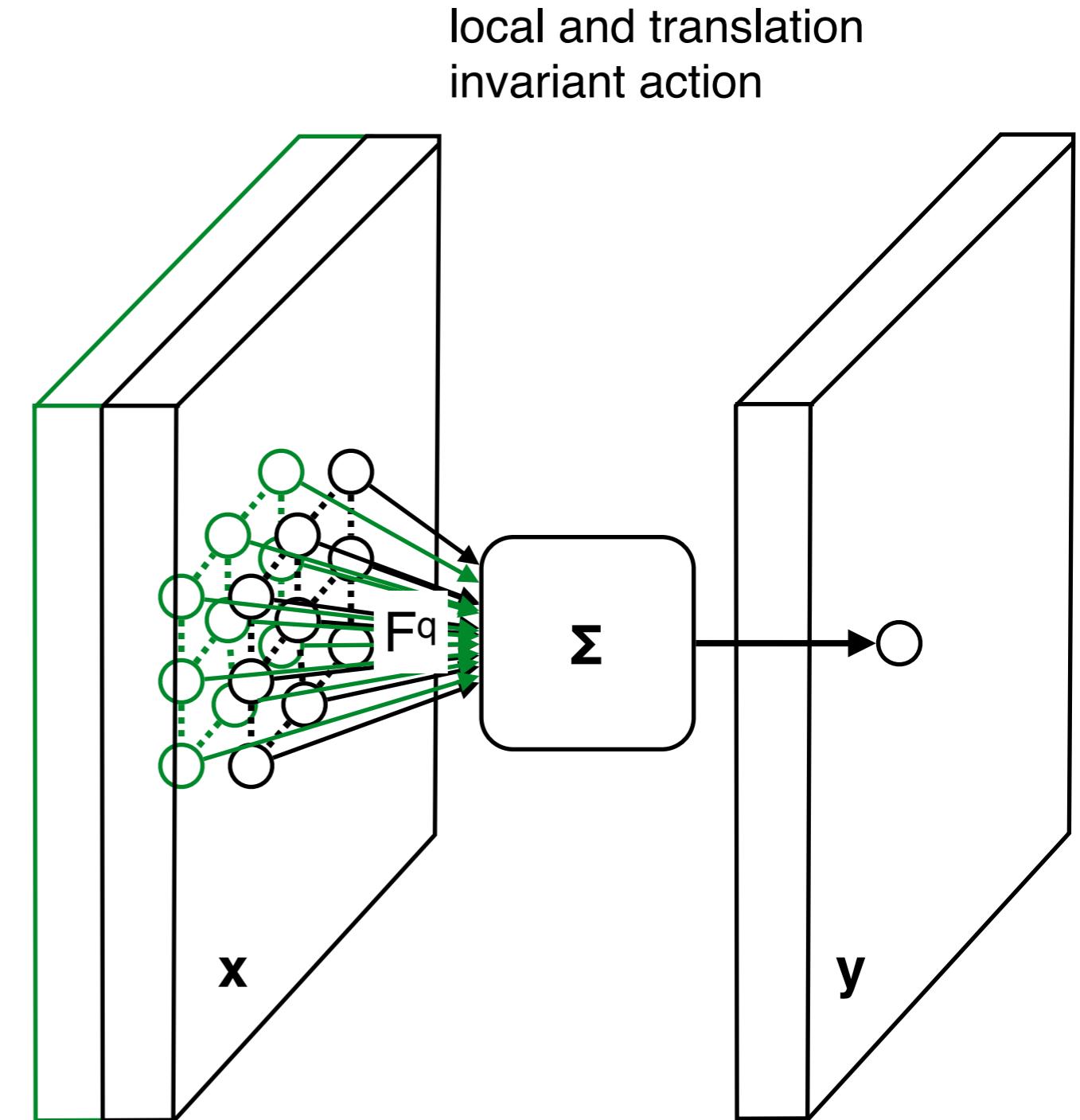
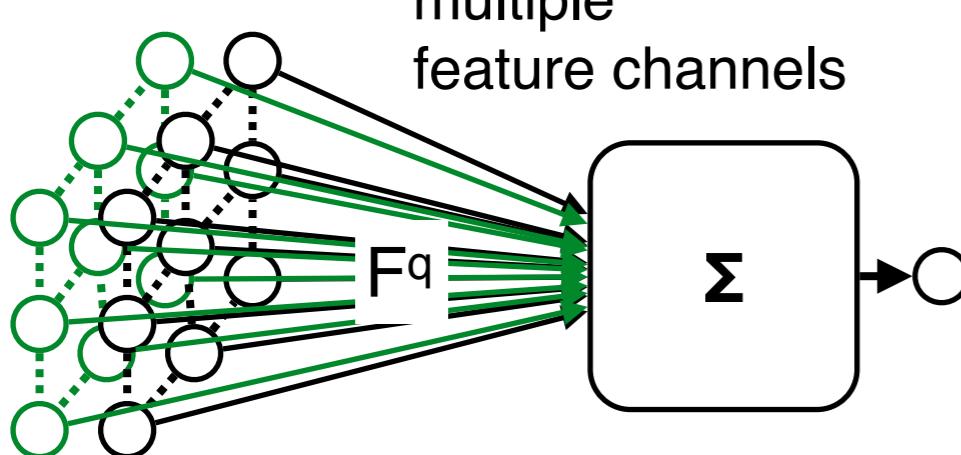
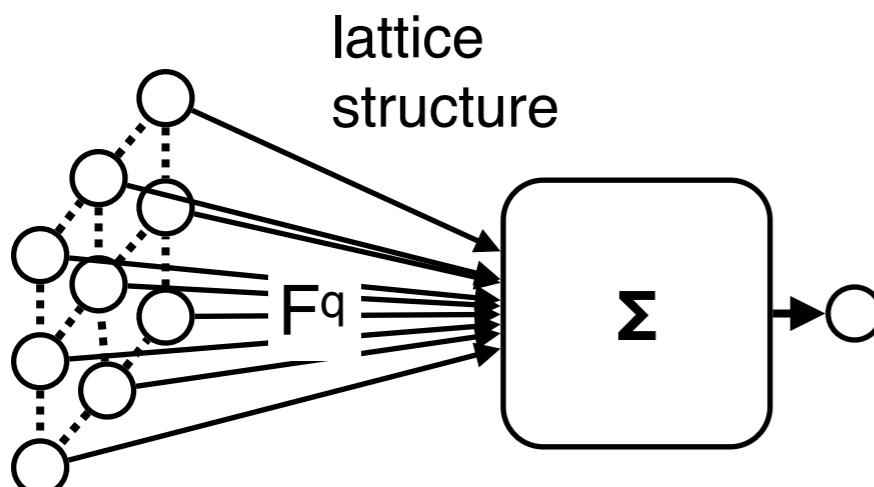
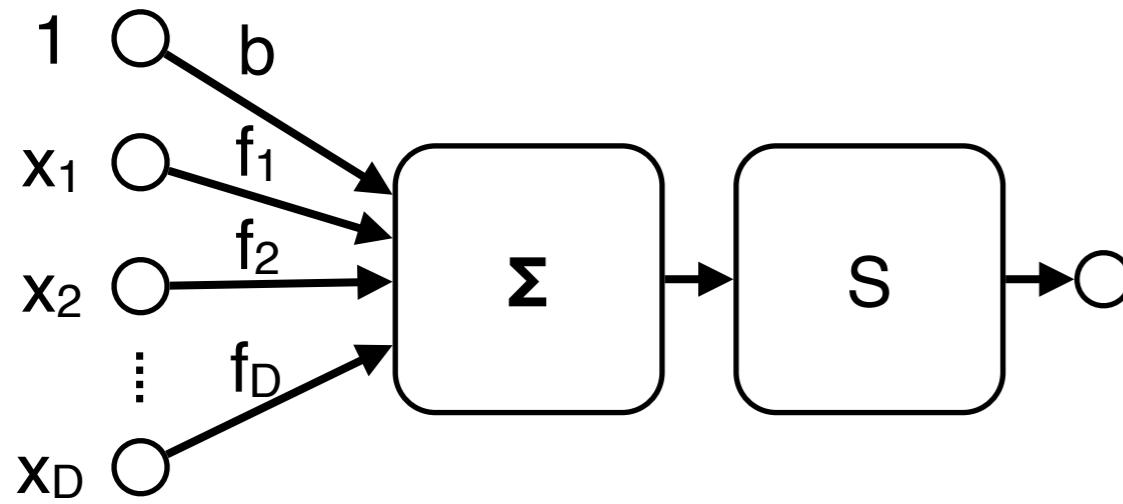


4 7

+ 100
000

Linear convolution

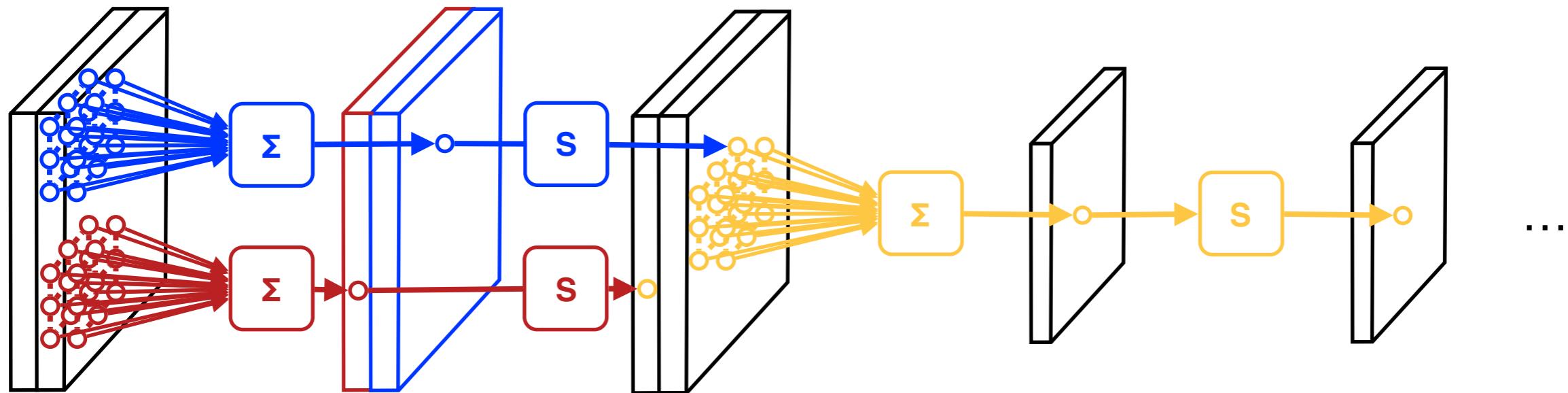
As a neural network



Deep architectures

11

Repeat linear / non-linear operators

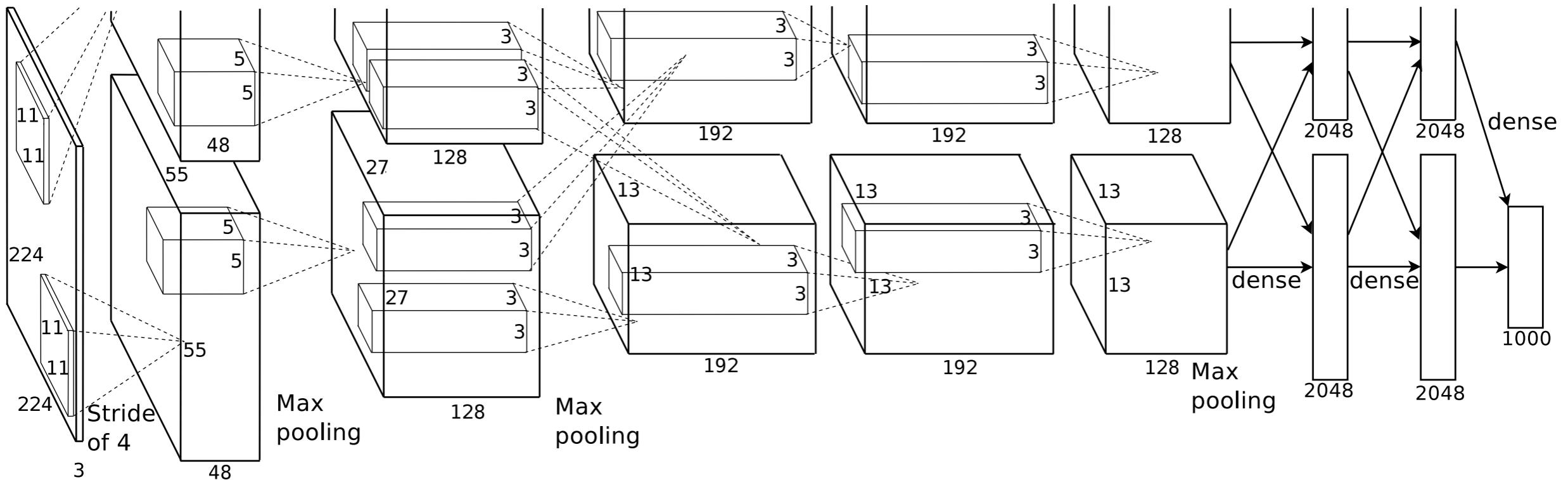


Modern convolutional networks

From AlexNet (2012) to ResNet (2015)



Modern convolutional nets



[AlexNet by Krizhevsky et al. 2012]

Excellent **performance** in image understanding tasks

Learn a sequence of **general-purpose representations**

Millions of parameters learned from data

The “**meaning**” of the representation is **unclear**

How deep is deep enough?

14

AlexNet (2012)



How deep is deep enough?

15

AlexNet (2012)

VGG-M (2013)

VGG-VD-16 (2014)



How deep is deep enough?

16

AlexNet (2012)

VGG-M (2013)

VGG-VD-16 (2014)

GoogLeNet (2014)



How deep is deep enough?

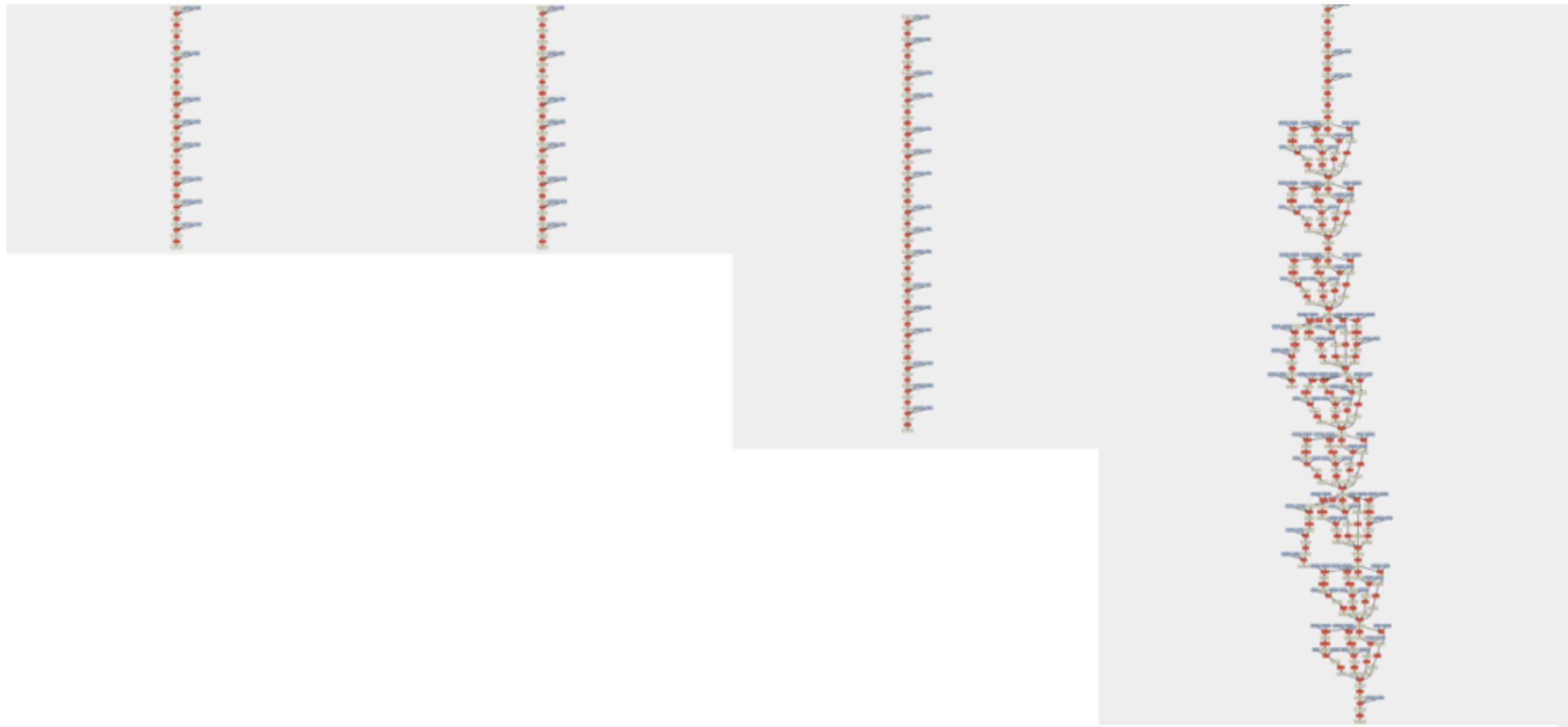
17

AlexNet (2012)

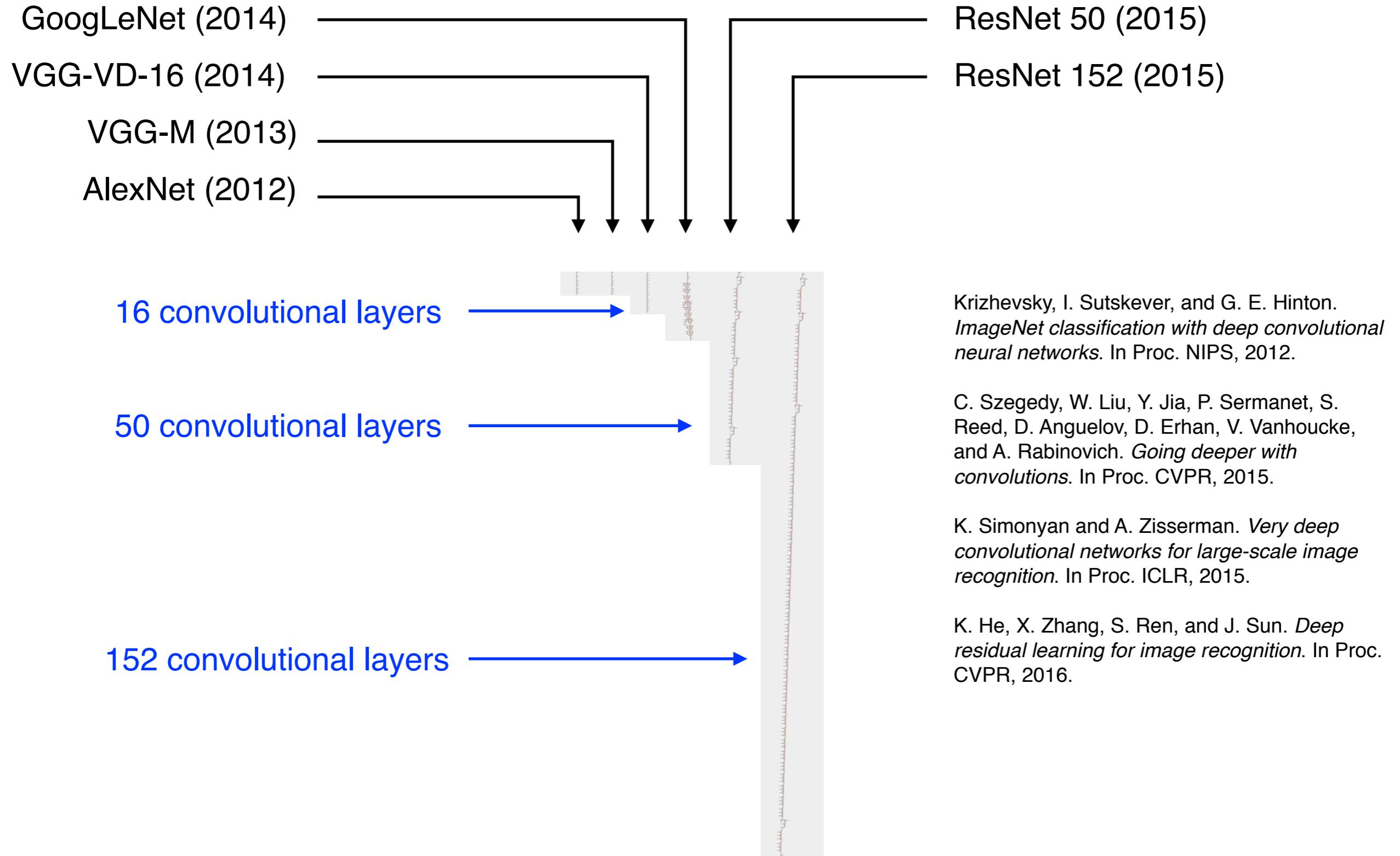
VGG-M (2013)

VGG-VD-16 (2014)

GoogLeNet (2014)

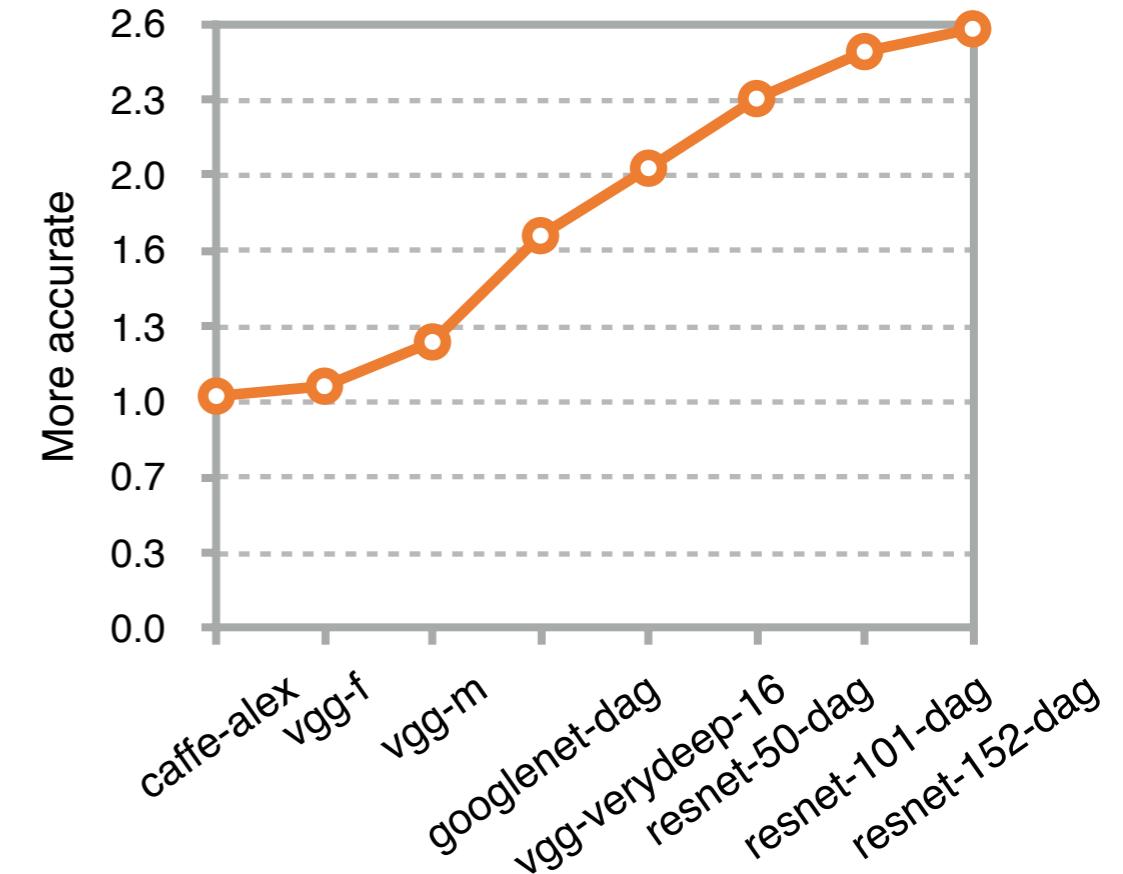
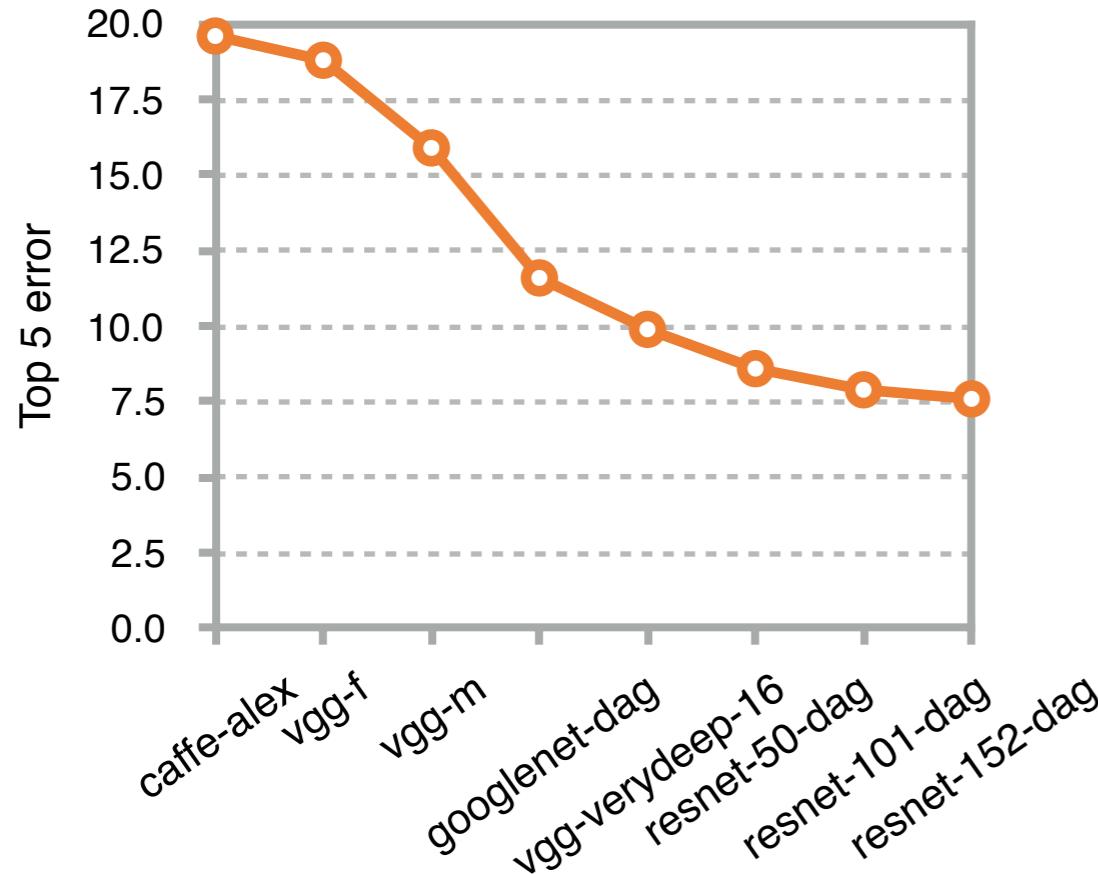


How deep is deep enough?



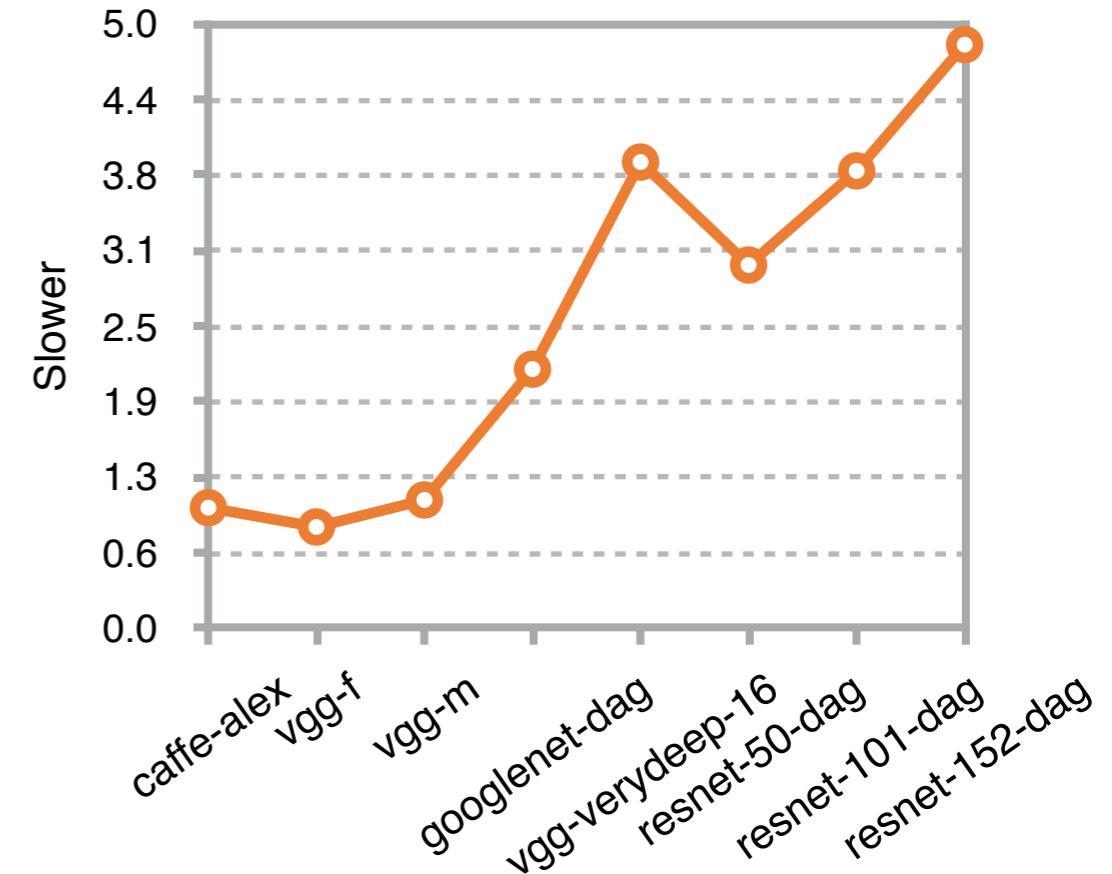
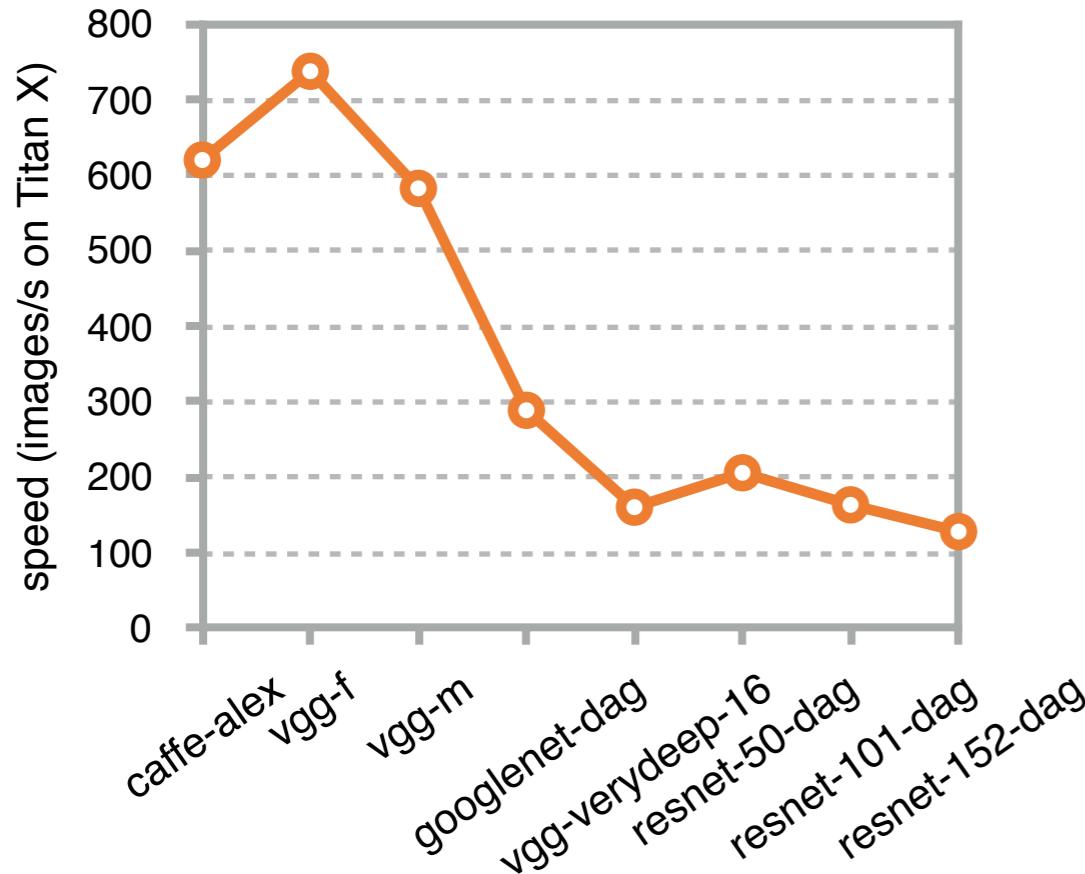
Accuracy

3 × more accurate in 3 years



Speed

5 × slower



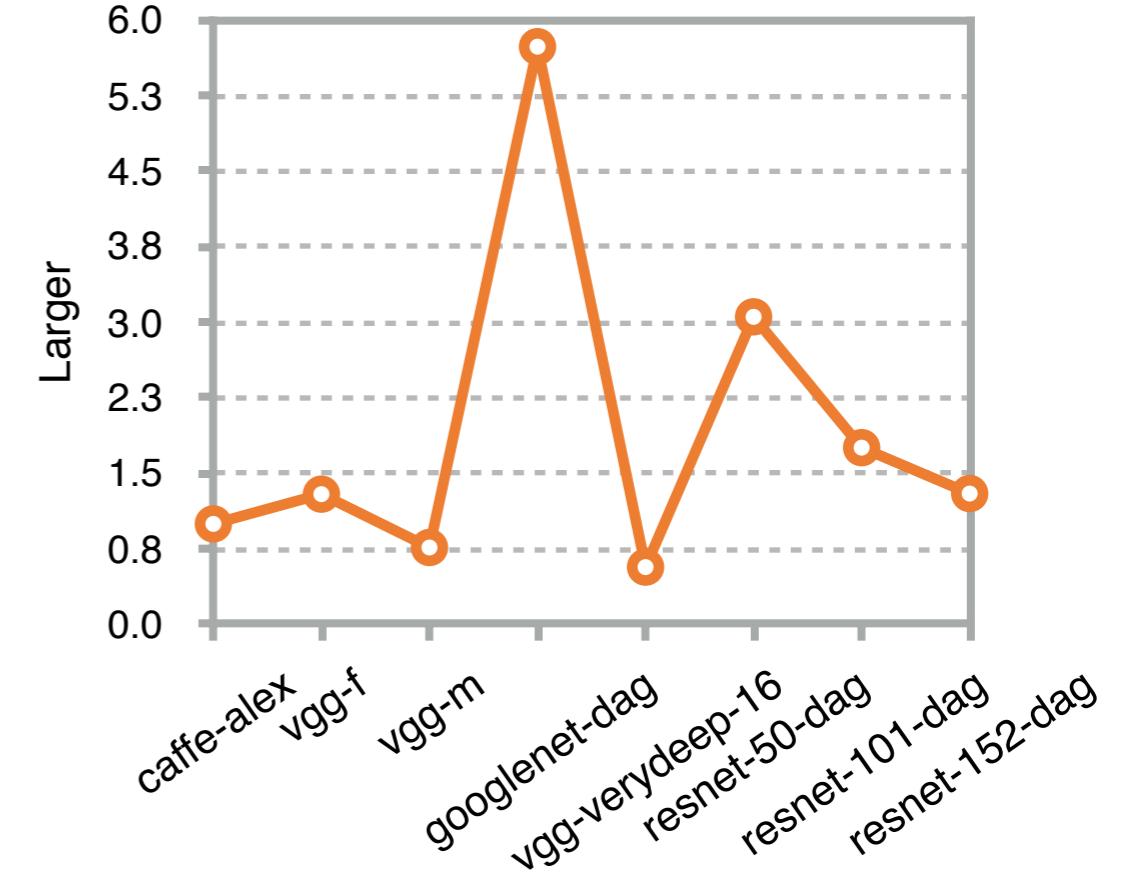
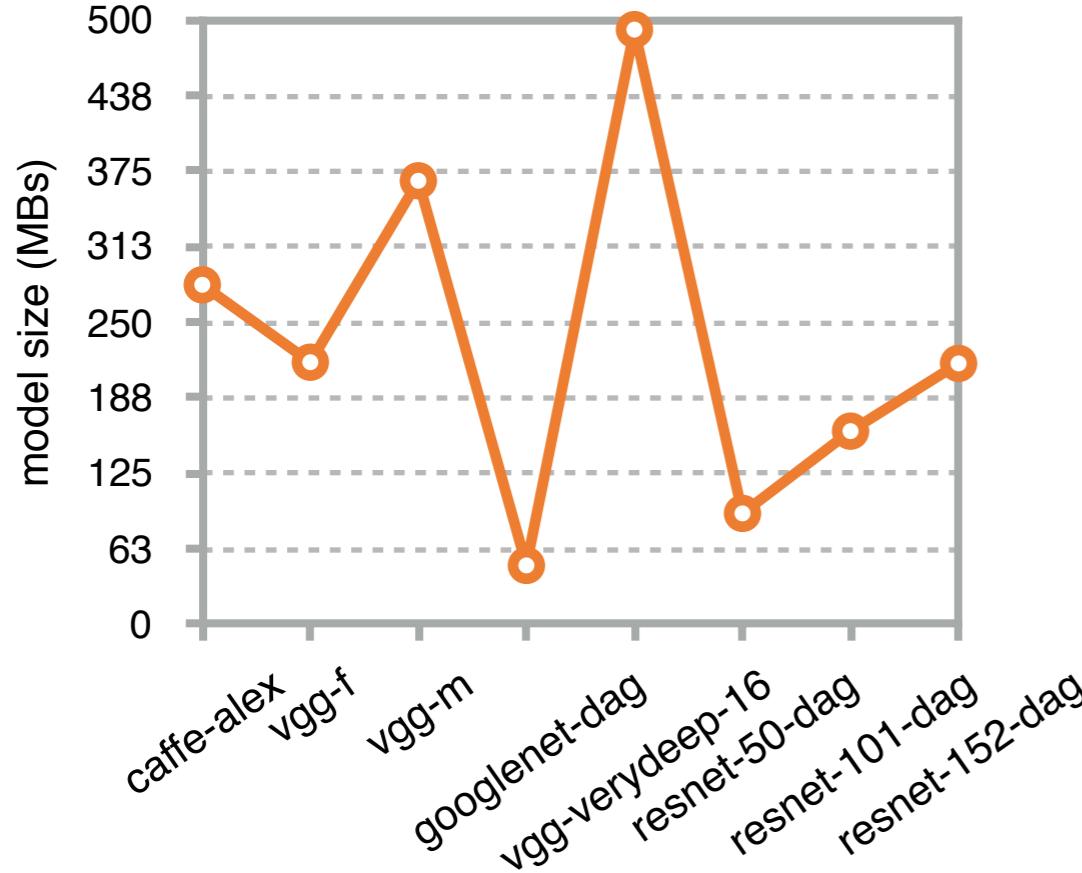
Remark: 101 ResNet layers same size/speed as 16 VGG-VD layers

Reason: far fewer feature channels (quadratic speed/space gain)

Moral: optimize your architecture

Model size

Num. of parameters is about the same



Remark: 101 ResNet layers same size/speed as 16 VGG-VD layers

Reason: far fewer feature channels (quadratic speed/space gain)

Moral: optimize your architecture

Deep nets in vision: 2012-2015

Impact of deep learning in vision

- ▶ **2012** amazing results by AlexNet in the ImageNet challenge
- ▶ **2013-15** massive 3x improvement
- ▶ **2016-19** more improvements?

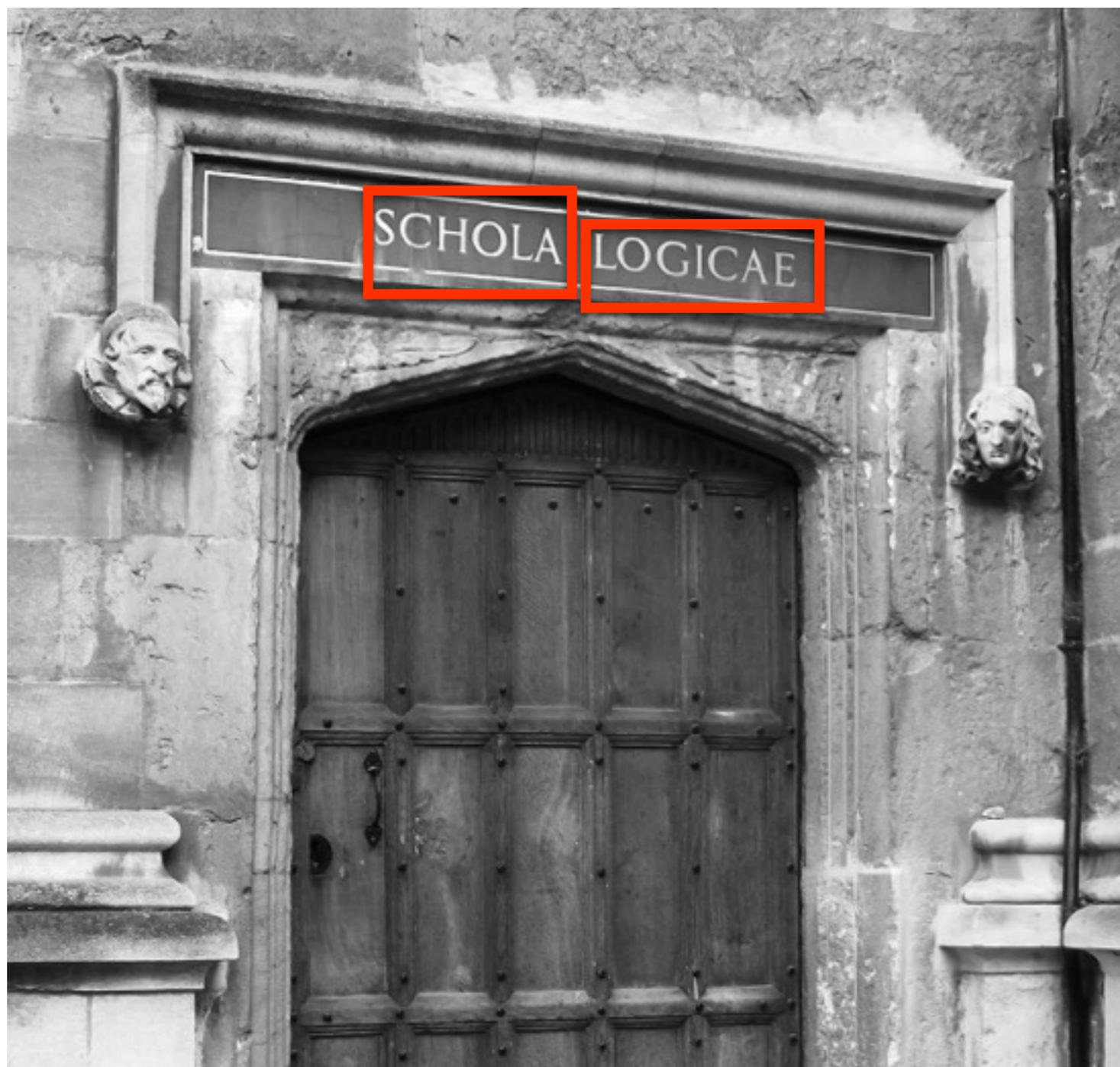
What have we learned

- ▶ Several **incremental tweaks** over the base AlexNet
- ▶ There is still space for improvements to the base model

Things that work

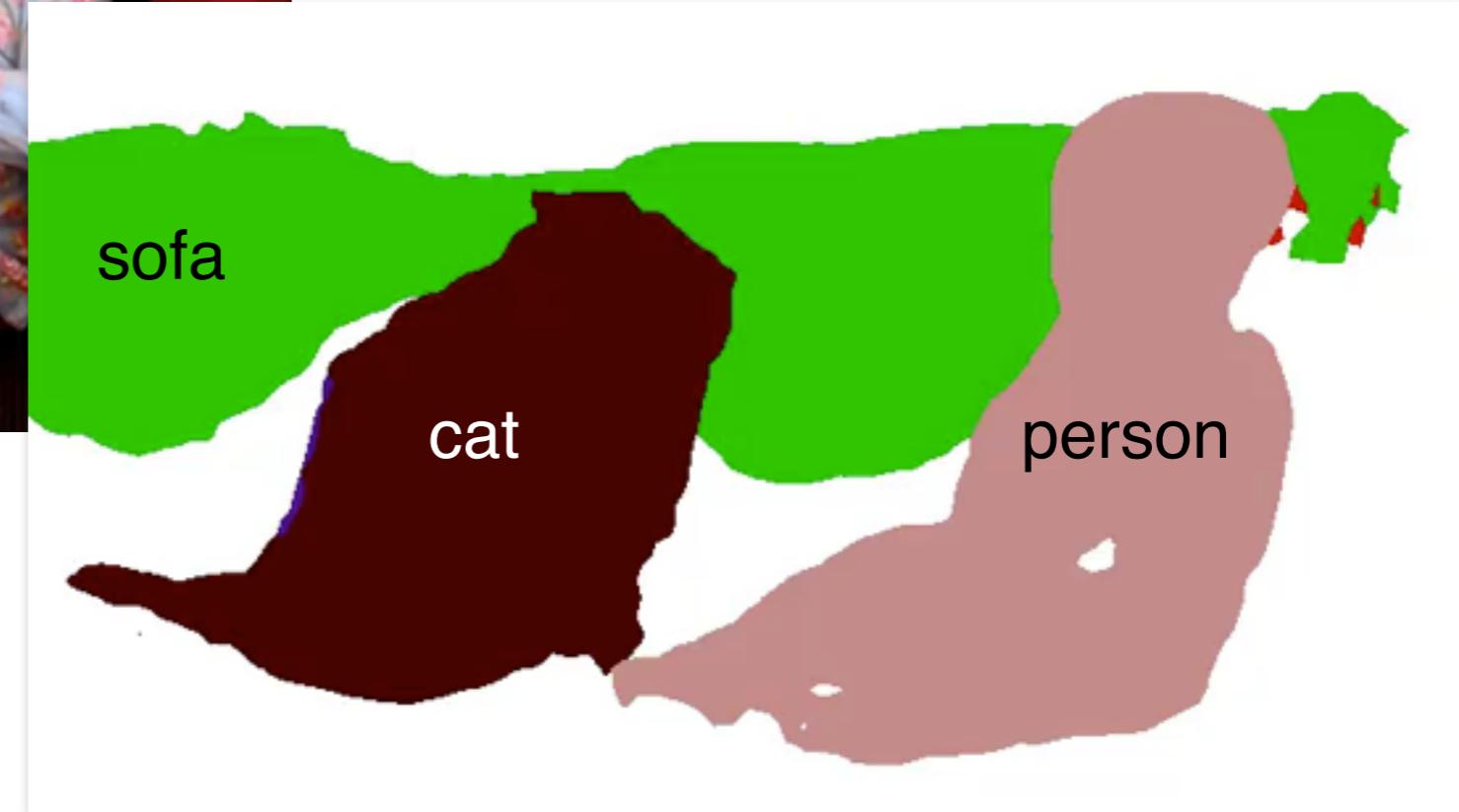
- ▶ Deeper architectures
- ▶ Smarter architectures (groups, low rank decompositions, ...)
- ▶ Batch normalization
- ▶ Residual connections

Applications



Semantic image segmentation

Label individual pixels



Face analysis

Detection, verification, recognition, emotion, 3D fitting



E.g. VGG-Face

Text spotting

Detection, word recognition, character recognition

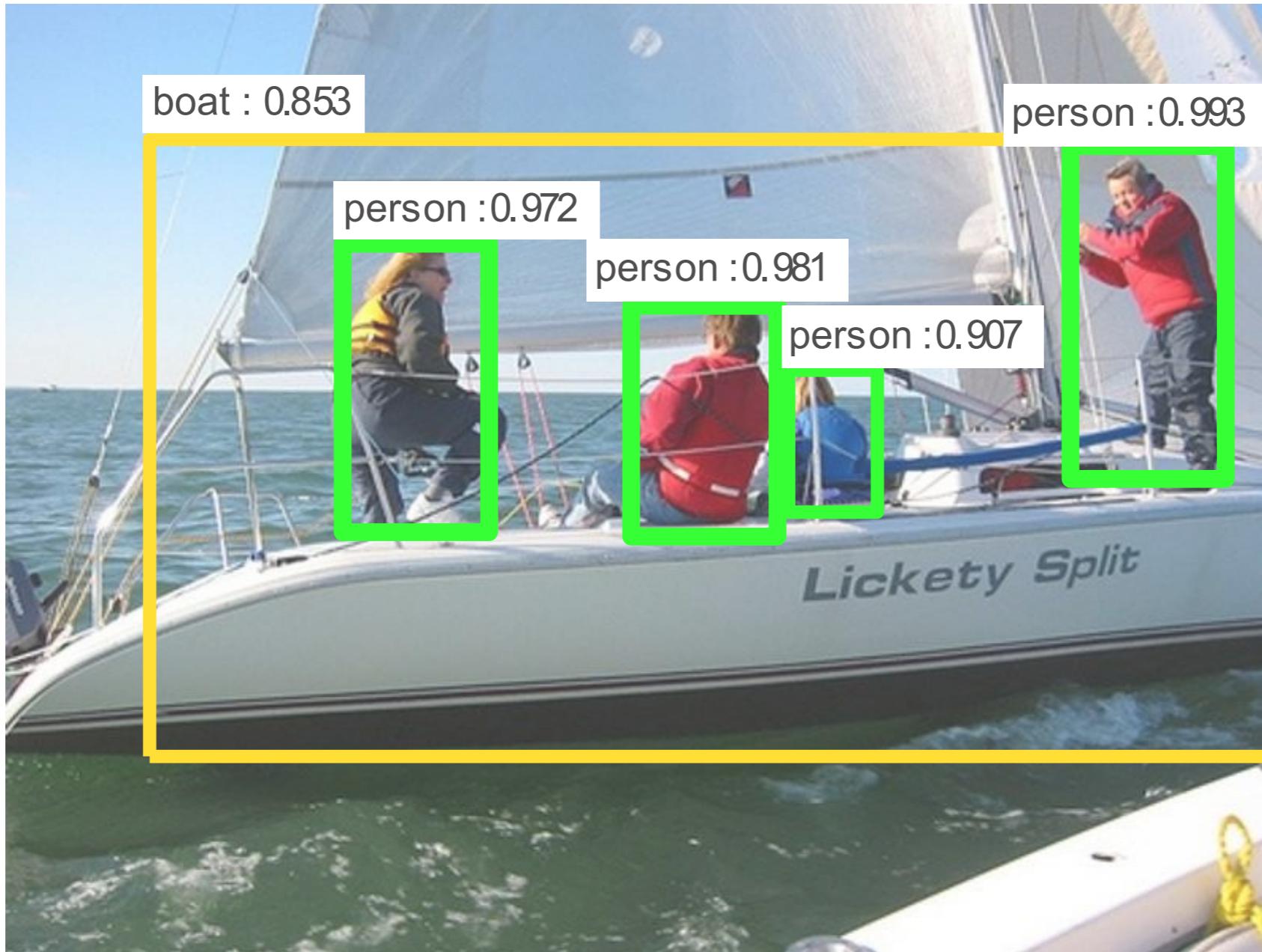


E.g. SynthText and VGG-Text

<http://zeus.robots.ox.ac.uk/textsearch/#/search/>

Object detection

Extract individual object instances



Supervision required

Still a major limitation

Can get around in various ways, for example using **synthetic data**



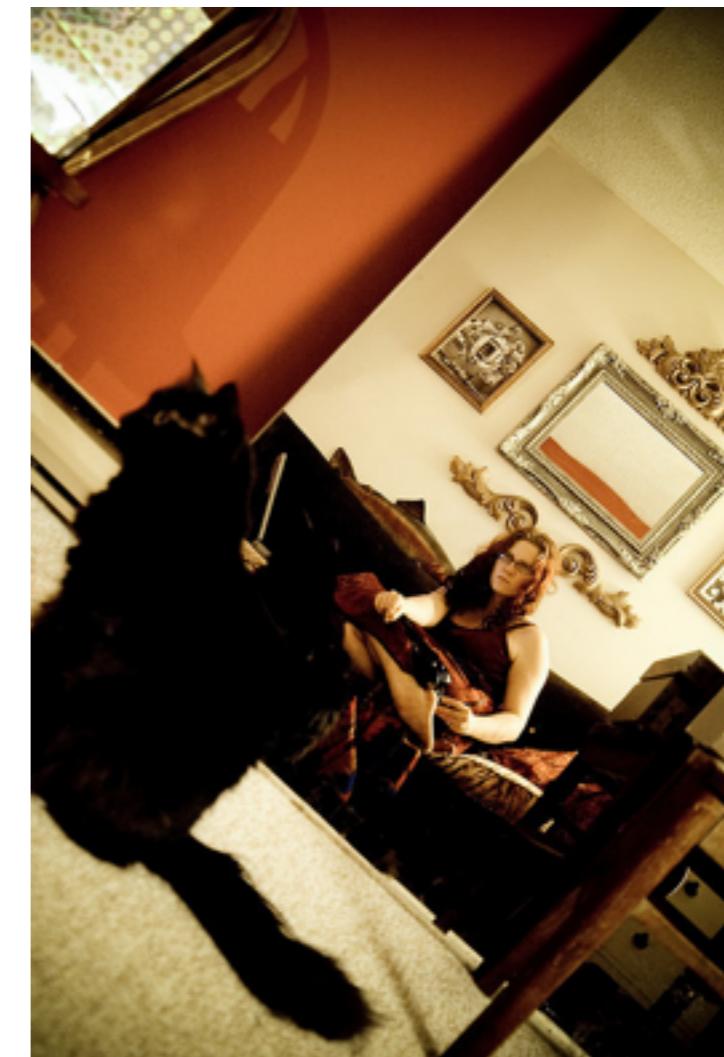
Weakly-supervised learning

Use partial labelings

From this



To this

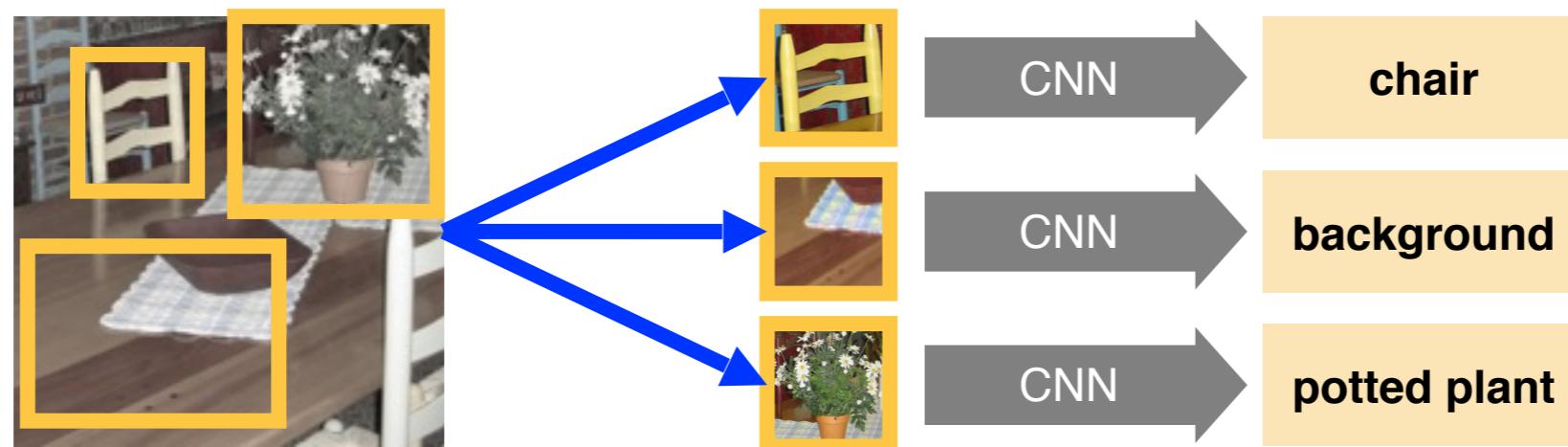


There are
a cat, a person,
a chair

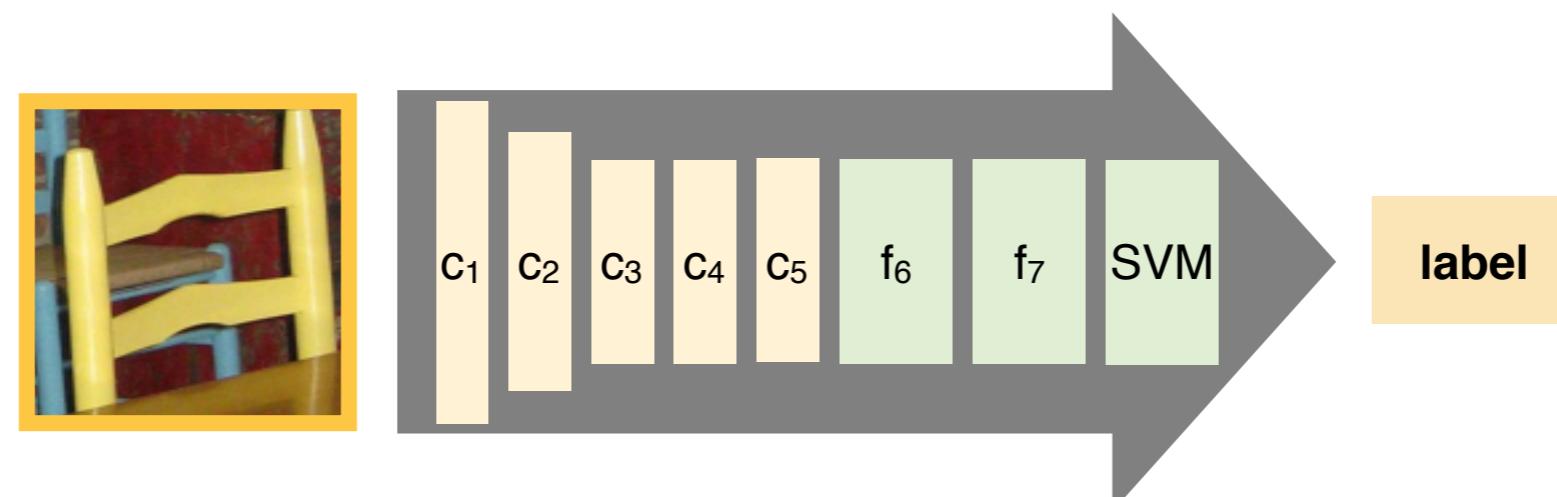
Detections with conv nets

Region-based Convolutional Neural Network (R-CNN)

Pros: simple and effective

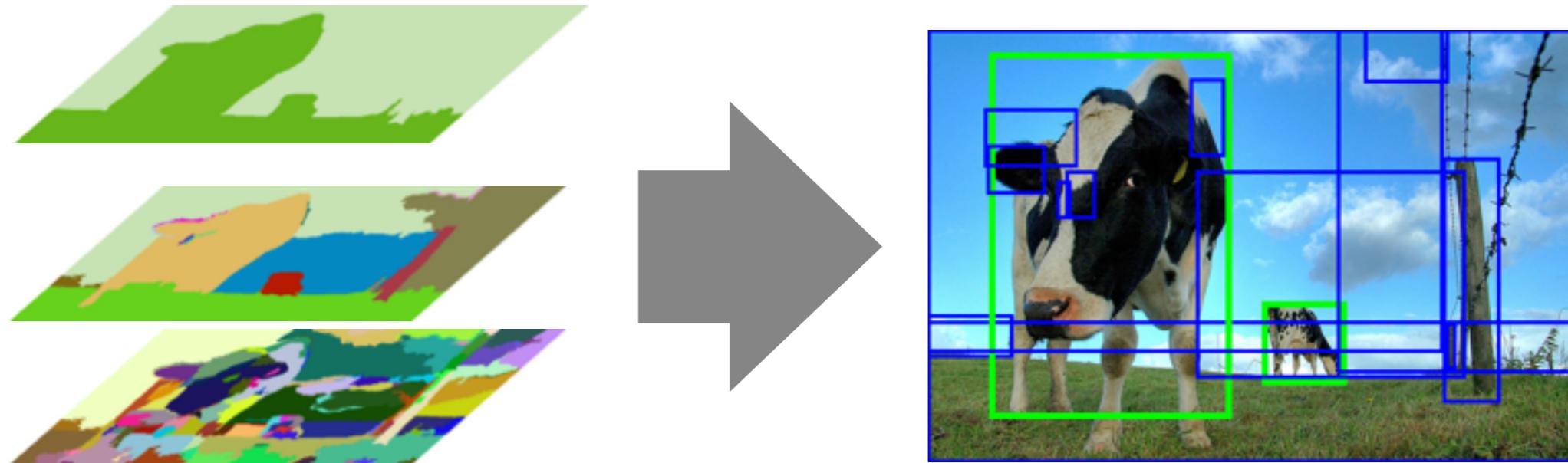


Cons: slow as the CNN is re-evaluated for each tested region



Region proposals

Cut down the number of candidates

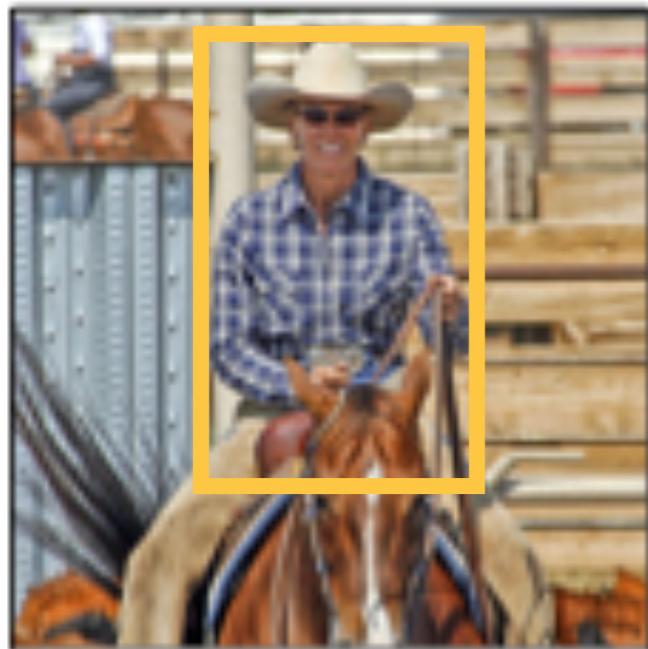


Proposal-method: Selective Search [van de Sande, Uijlings et al.]

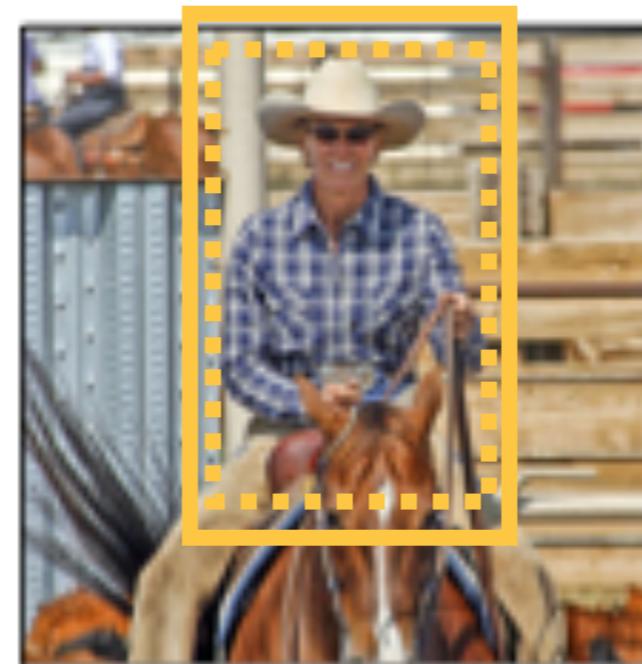
- ▶ hierarchical segmentation
- ▶ each region generates a ROI
- ▶ ~ 2000 regions / image

From proposals to CNN features

Dilate, crop, reshape



Propose



Dilate



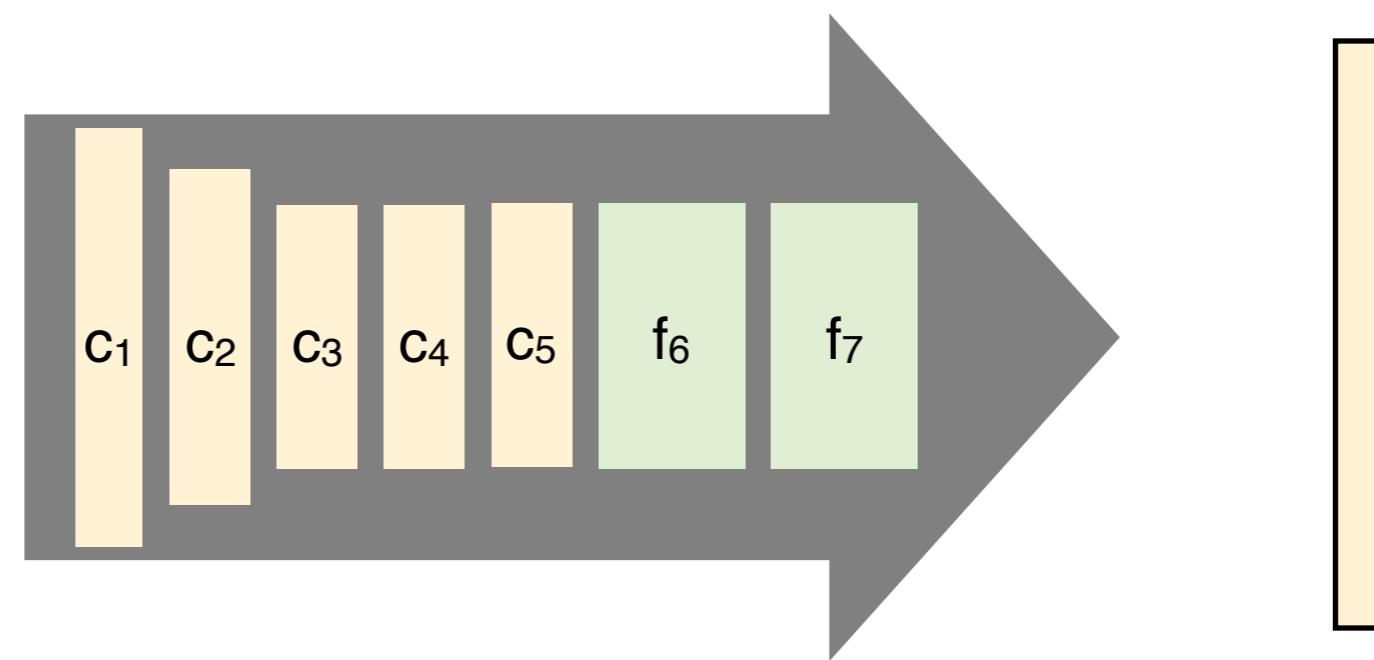
Crop & scale

Anisotropic

227 x 227

From proposals to CNN features

Evaluate CNN



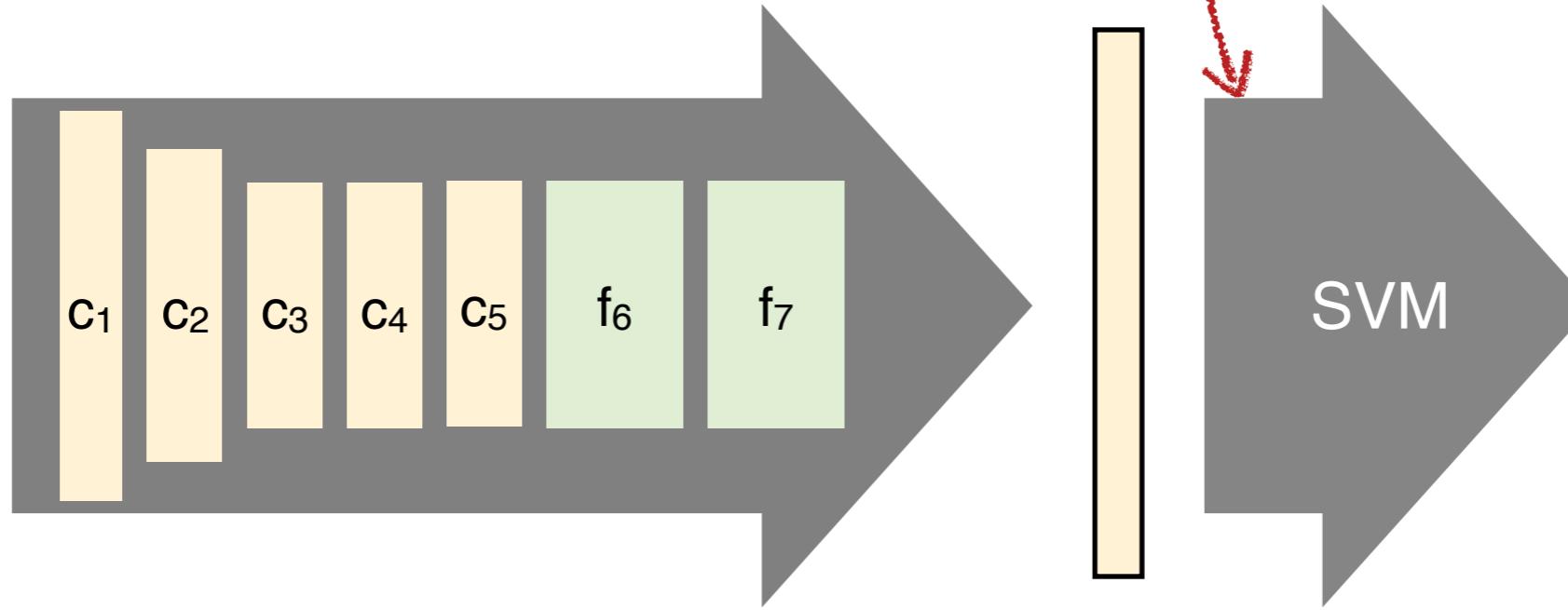
Scale
Anisotropic
 227×227

CNN features
Up to FC-7
AlexNet

Feature vector
4096 D

Classification of a region

Run an SVM or similar on top



Scale
Anisotropic
227 x 227

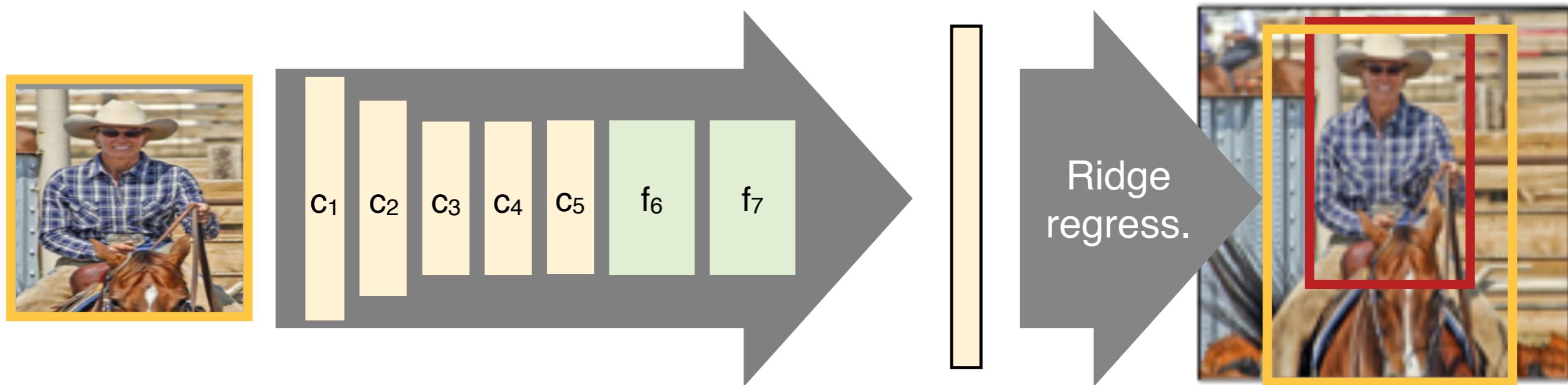
CNN features
Up to FC-7
AlexNet

Feature vector
4096 D

Label
One out of N

Region adjustment

Bounding-box regression



Scale
Anisotropic
227 x 227

CNN features
Up to FC-7
AlexNet

Feature vector
4096 D

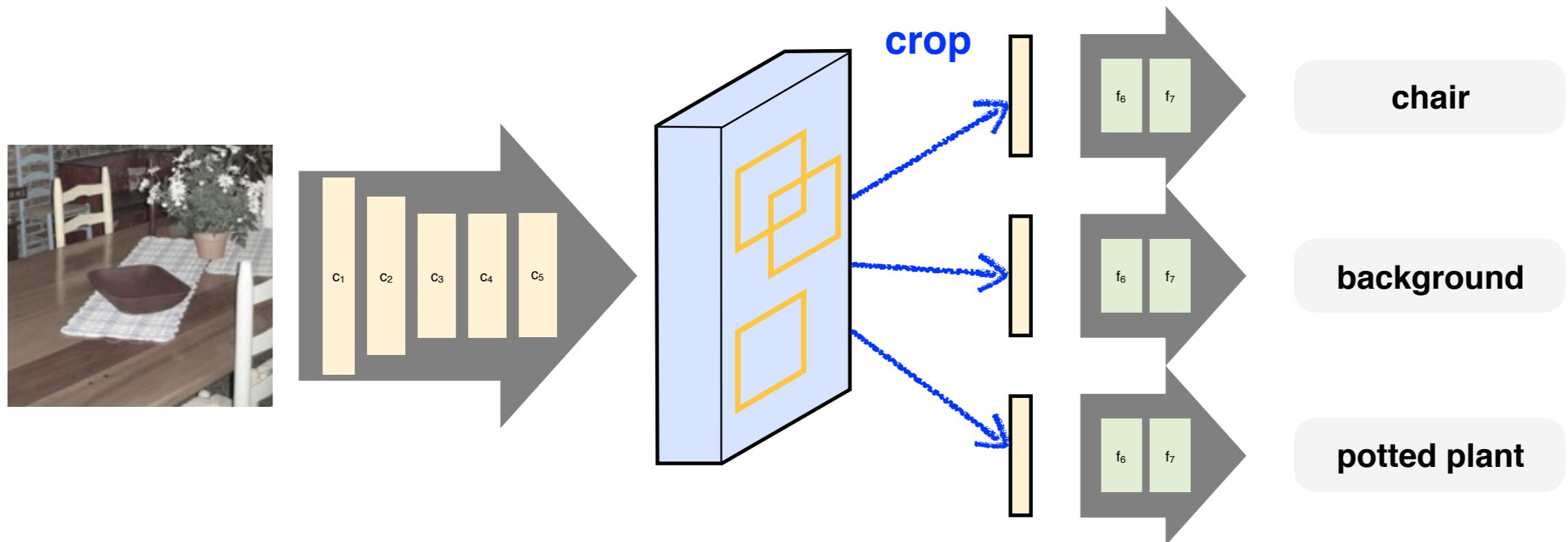
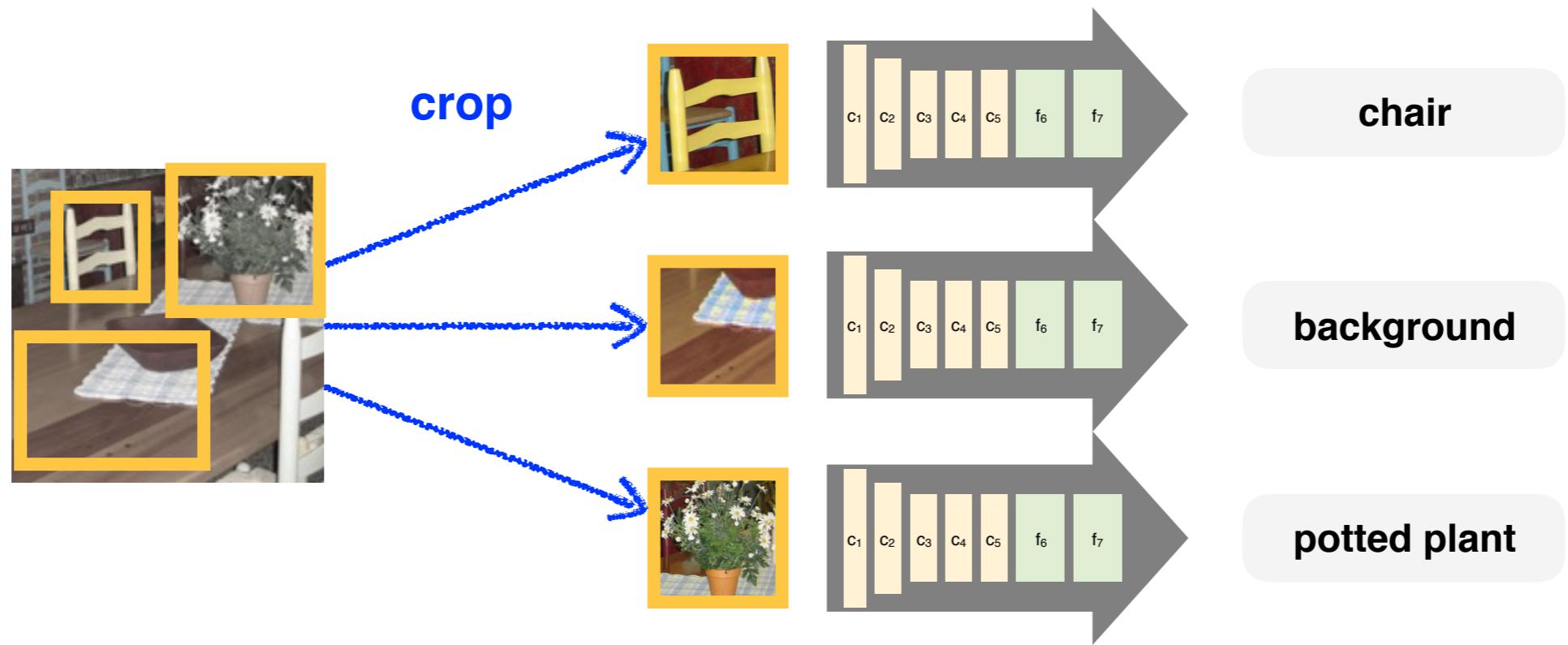
Box adjustment
 dx_1, dx_2, dy_1, dy_2

R-CNN results on PASCAL VOC

At the time of introduction (2013)

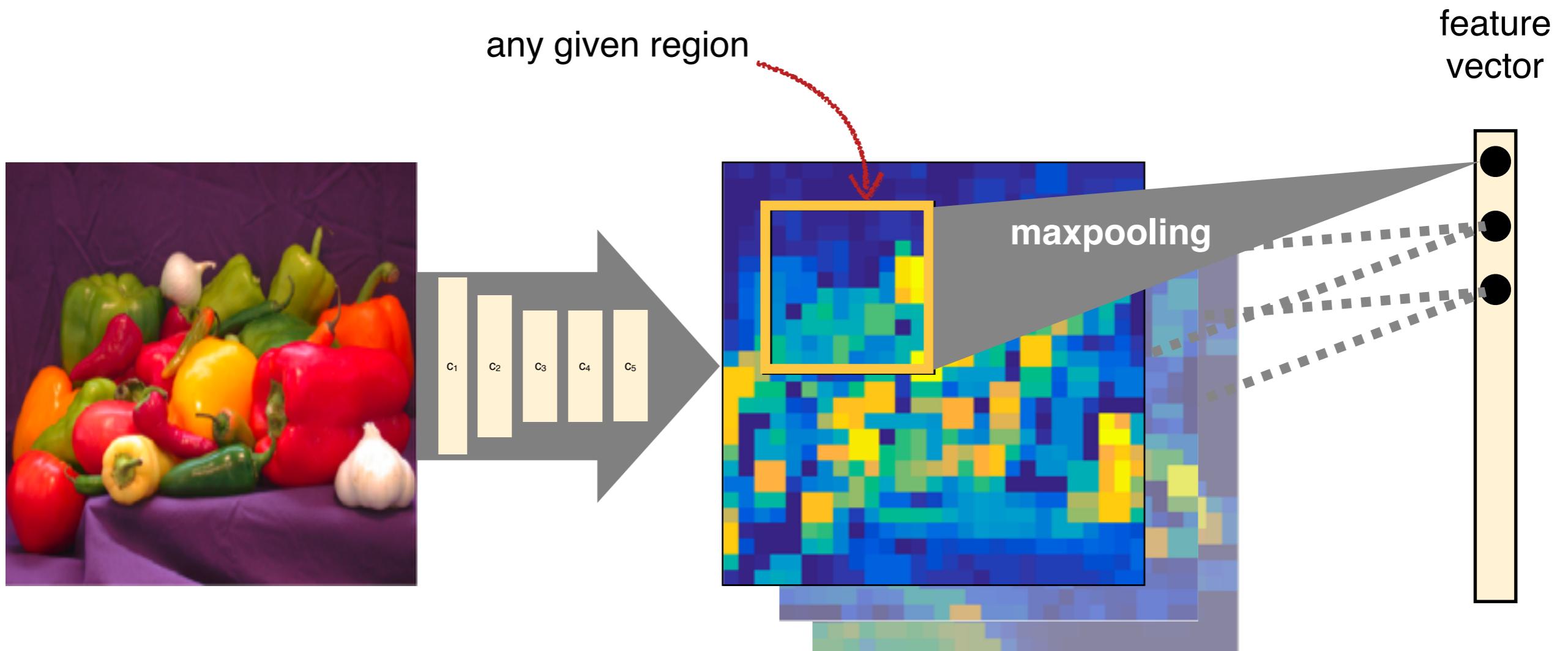
| | VOC 2007 | VOC 2010 |
|--|----------|----------|
| DPM v5 (Girshick et al. 2011) | 33.7% | 29.6% |
| UVA sel. search (Uijlings et al. 2013) | | 35.1% |
| Regionlets (Wang et al. 2013) | 41.7% | 39.7% |
| SegDPM (Fidler et al. 2013) | | 40.4% |
| R-CNN (TorontoNet) | 54.2% | 50.2% |
| R-CNN (TorontoNet) + bbox regression | 58.5% | 53.7% |
| R-CNN (VGG-VD) | 62.1% | |
| R-CNN (ONet) + bbox regression | 66.0% | 62.9% |

Accelerating R-CNN



The Spatial (Pyramid) Pooling layer

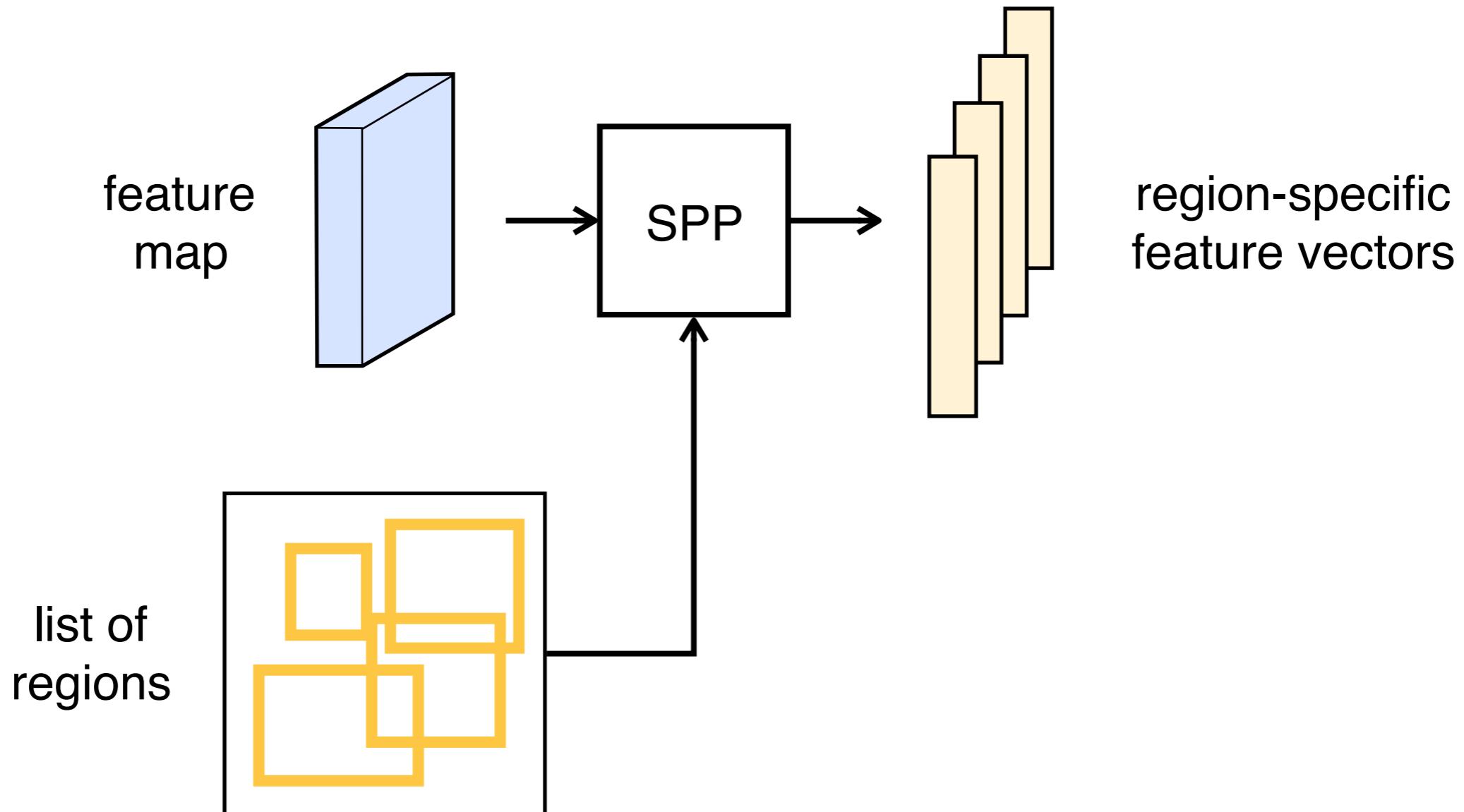
Max pooling in arbitrary regions



The Spatial (Pyramid) Pooling layer

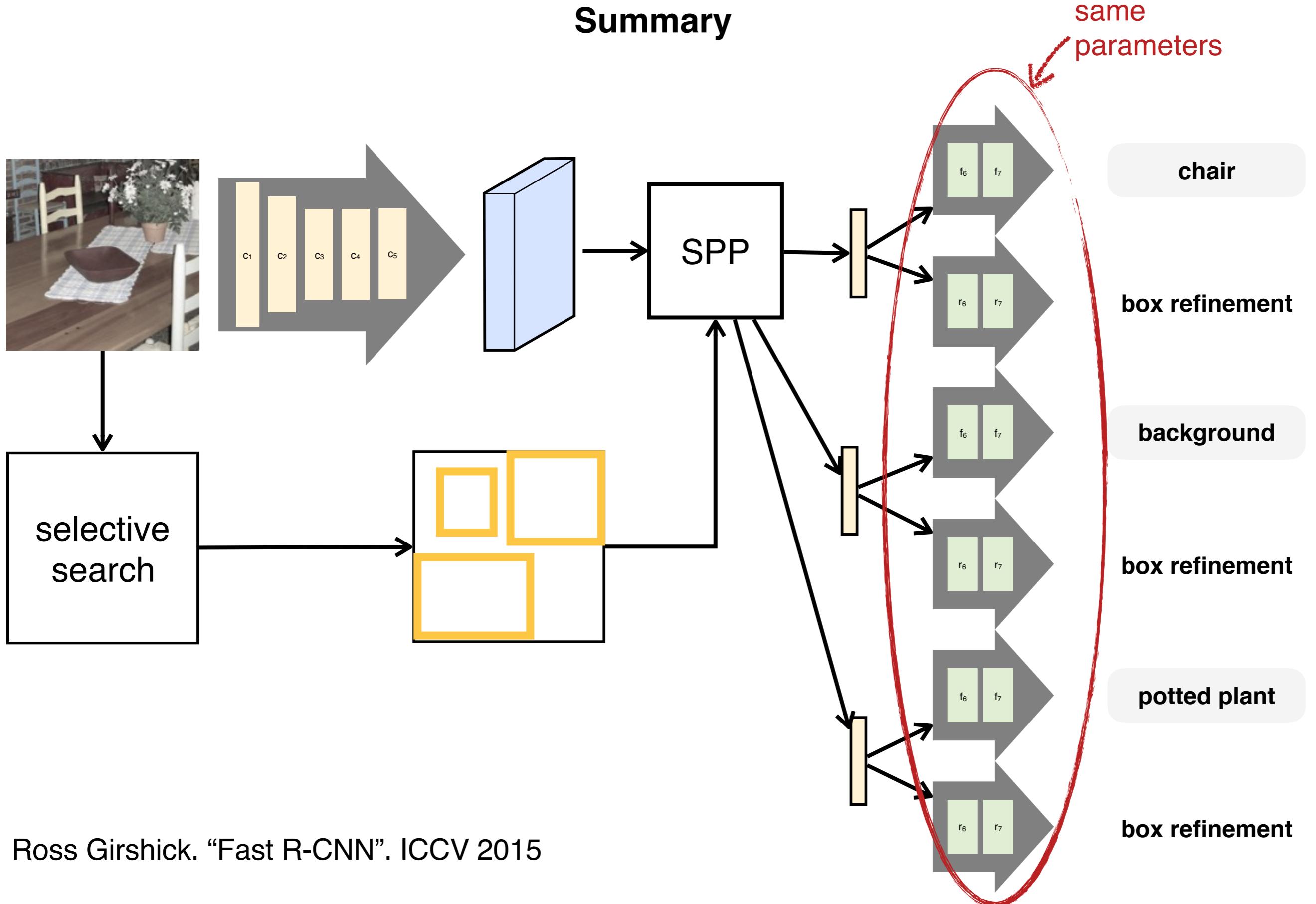
39

As a building block



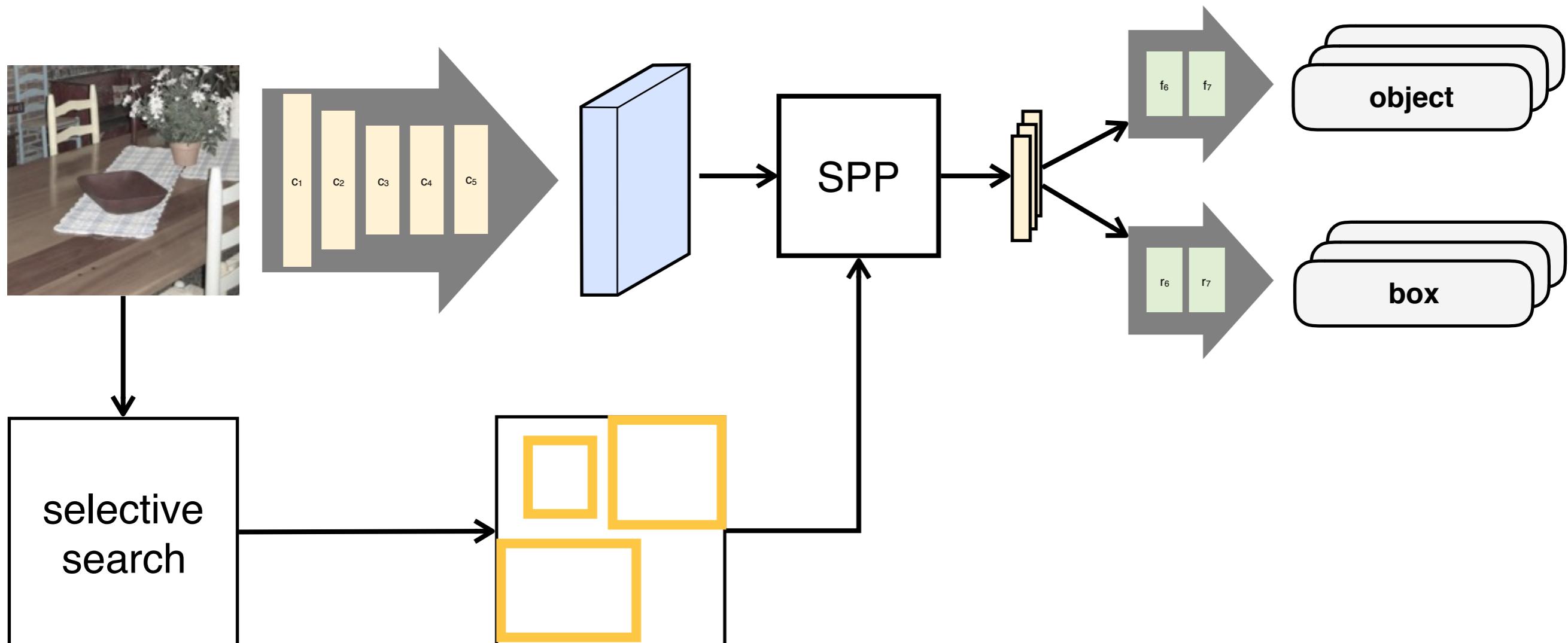
Fast R-CNN

Summary



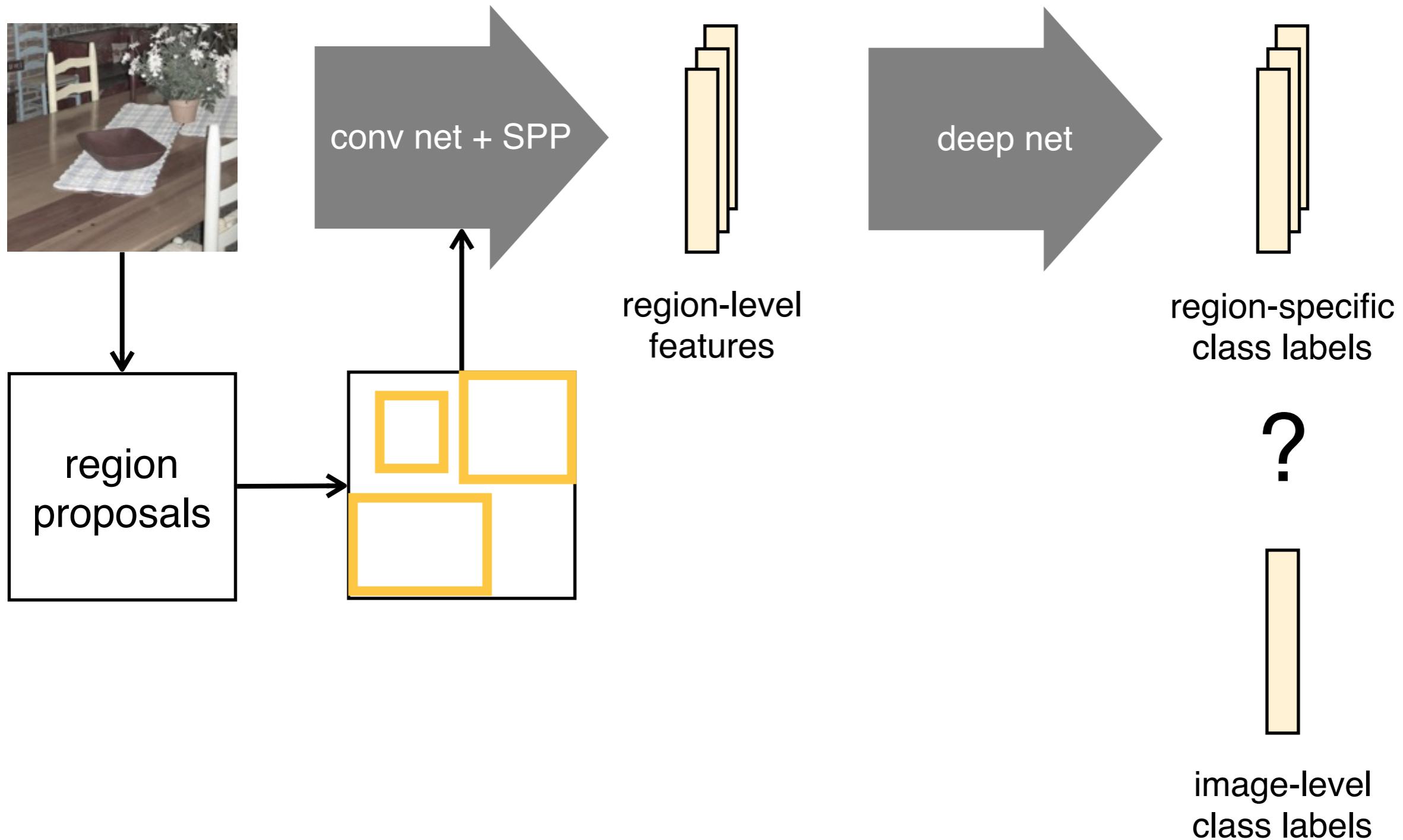
Fast R-CNN

Summary



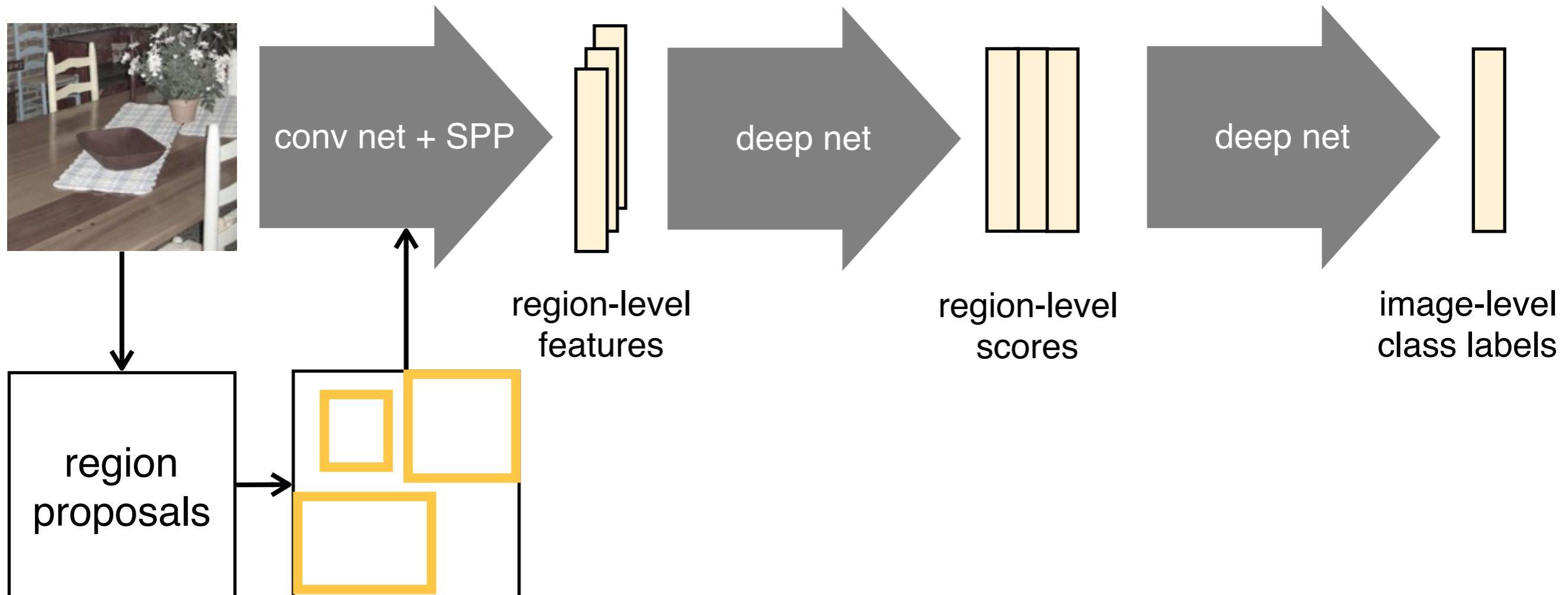
Object detection with deep networks

High-level view



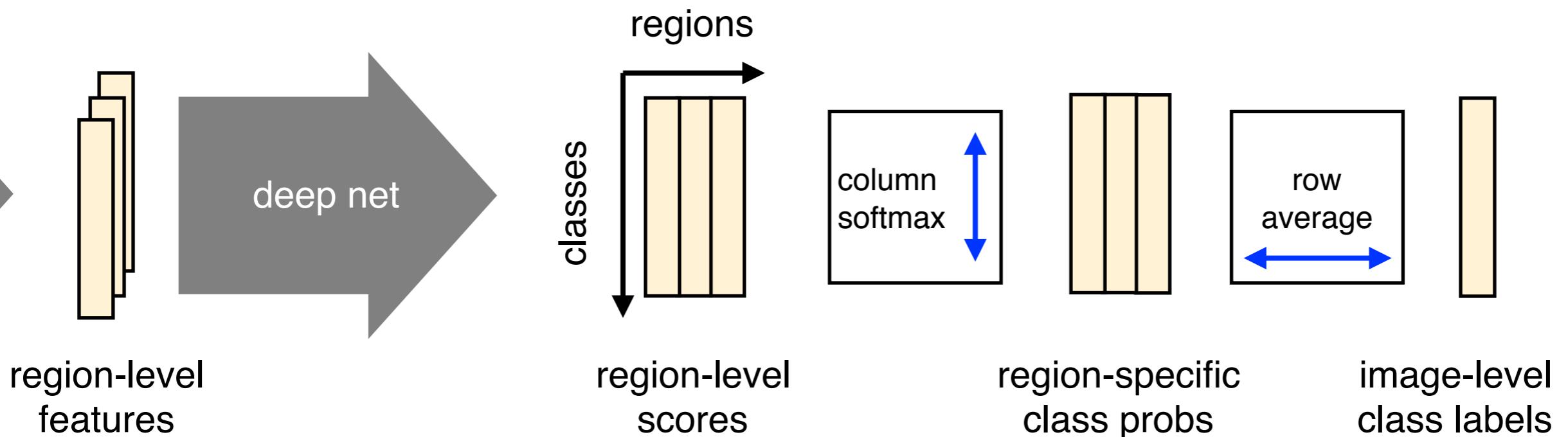
Towards weak supervision

Region scores to class labels



Towards weak supervision

Class labels = average of region labels

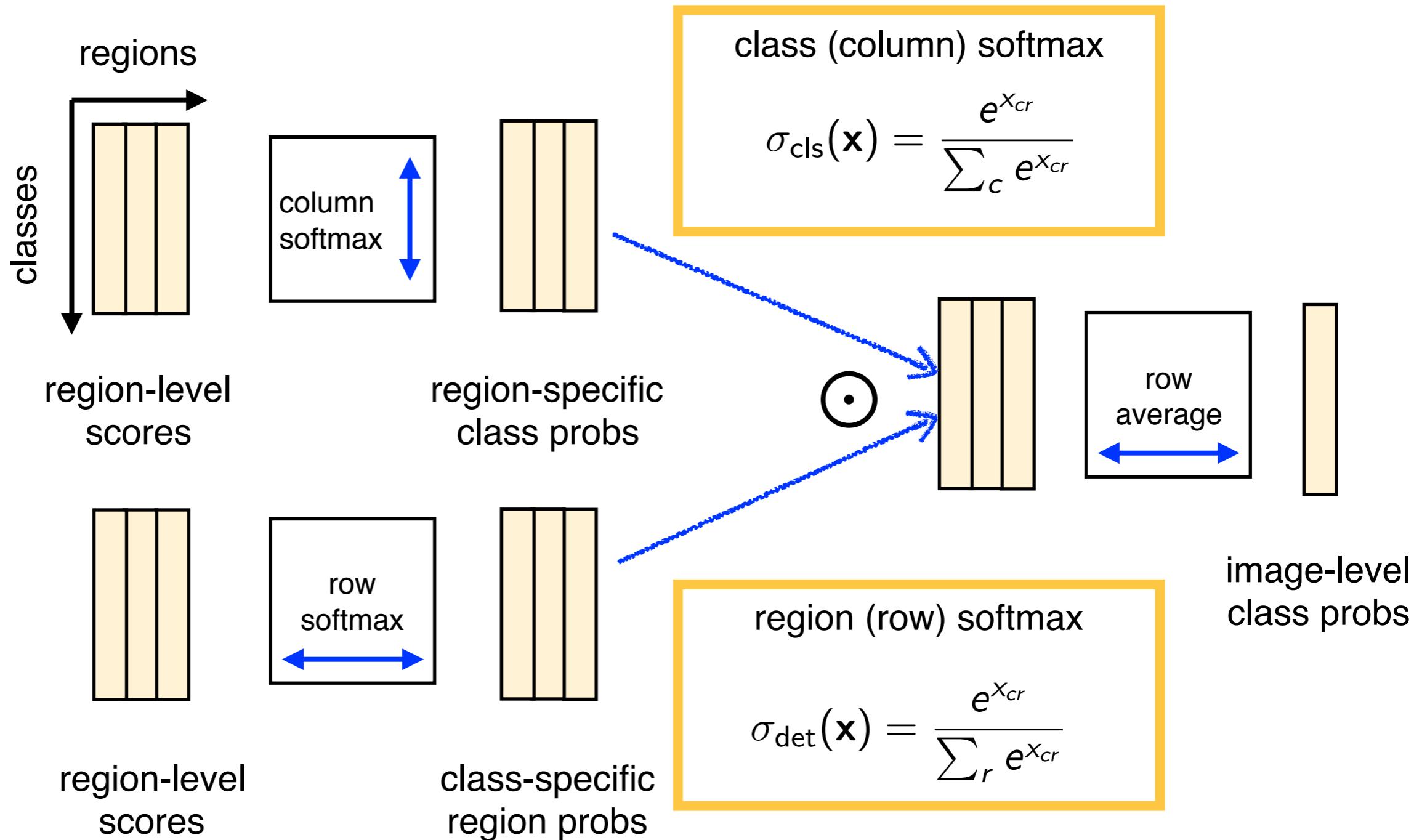


class (column) softmax

$$\sigma_{cls}(x) = \frac{e^{x_{cr}}}{\sum_c e^{x_{cr}}}$$

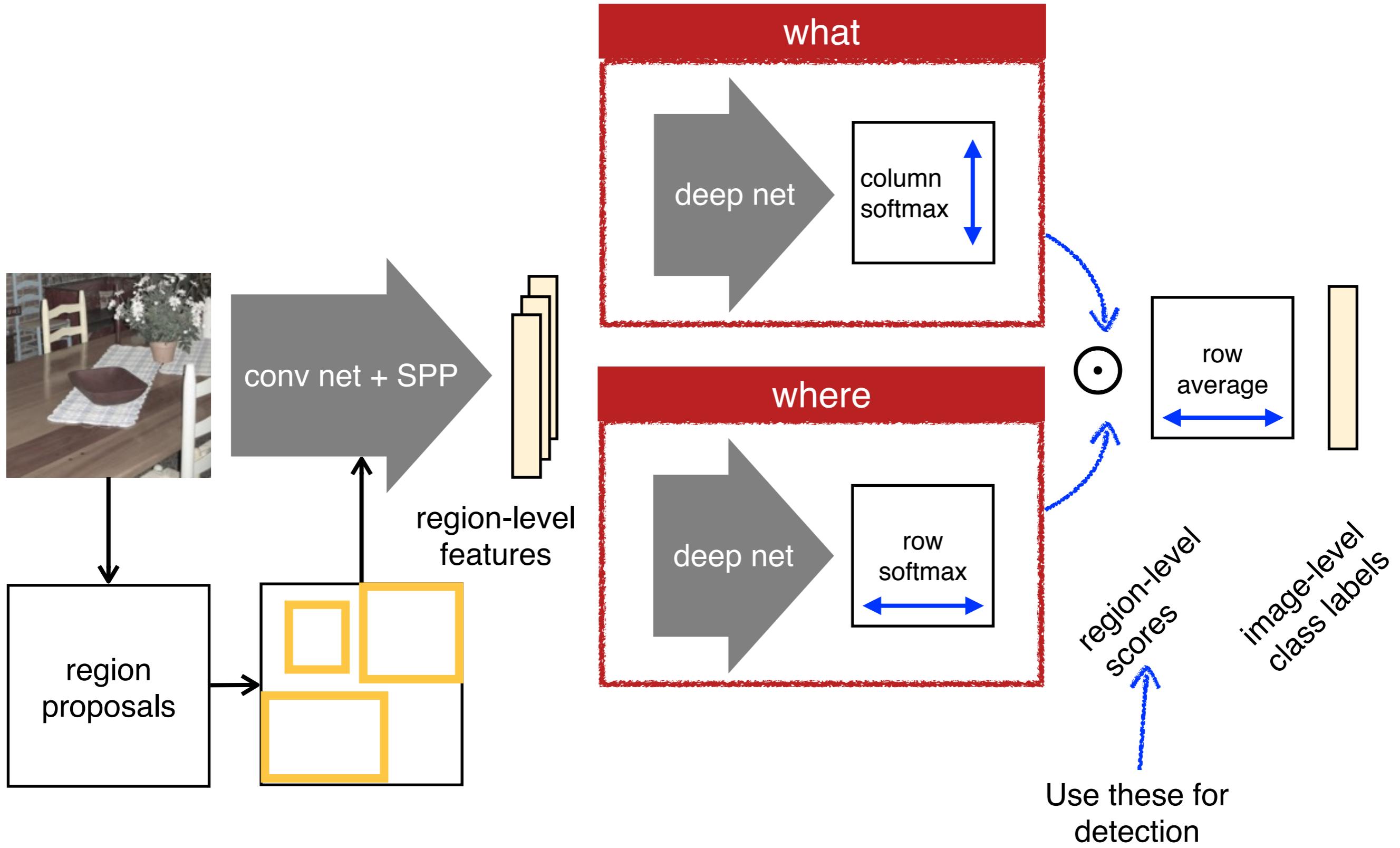
Towards weak supervision

Combine both class and region information



Two-streams weakly-supervised R-CNN

Overview



Results

Two streams vs state-of-the-art

| | Wang@ ECCV14 | Bilen@ CVPR15 | Cinbis@ PAMI16 | Two Streams ++ |
|--------------|-----------------|------------------|-------------------|-------------------|
| PASCAL VOC07 | 30.2 | 27.7 | 31.6 | 39.3 |
| PASCAL VOC10 | 27.4 | - | - | 36.2 |

Two streams vs single stream

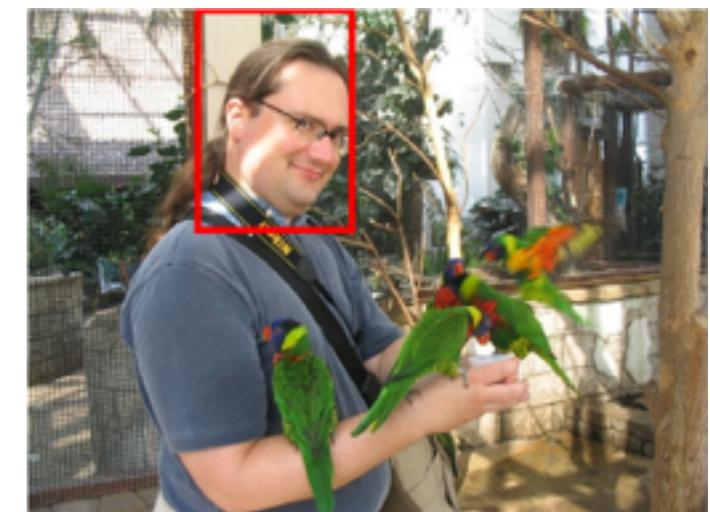
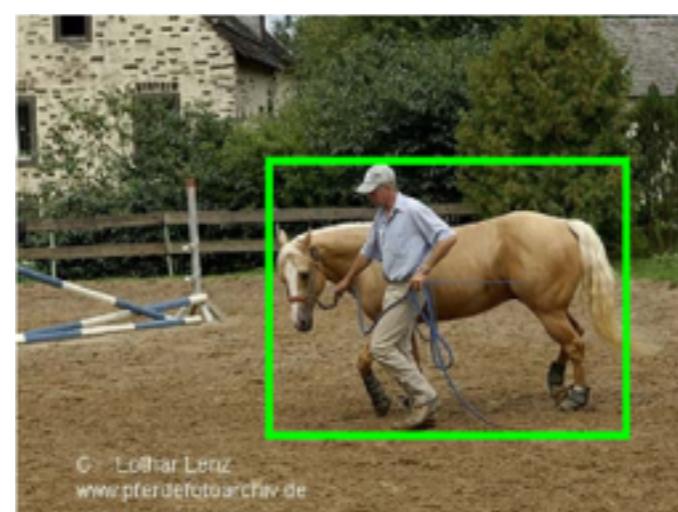
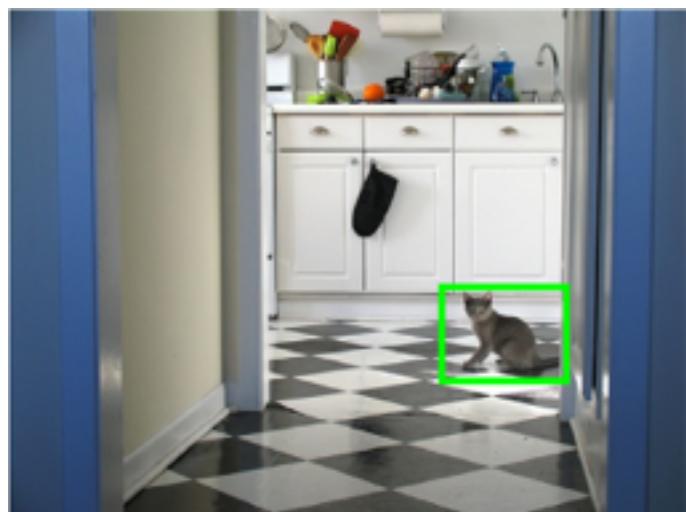
| | Single Stream | Two Streams | Two Streams ++ |
|--------------|------------------|----------------|-------------------|
| PASCAL VOC07 | 21.6 | 30.9 | 39.3 |

Results

Good results



Failure modes



Intro

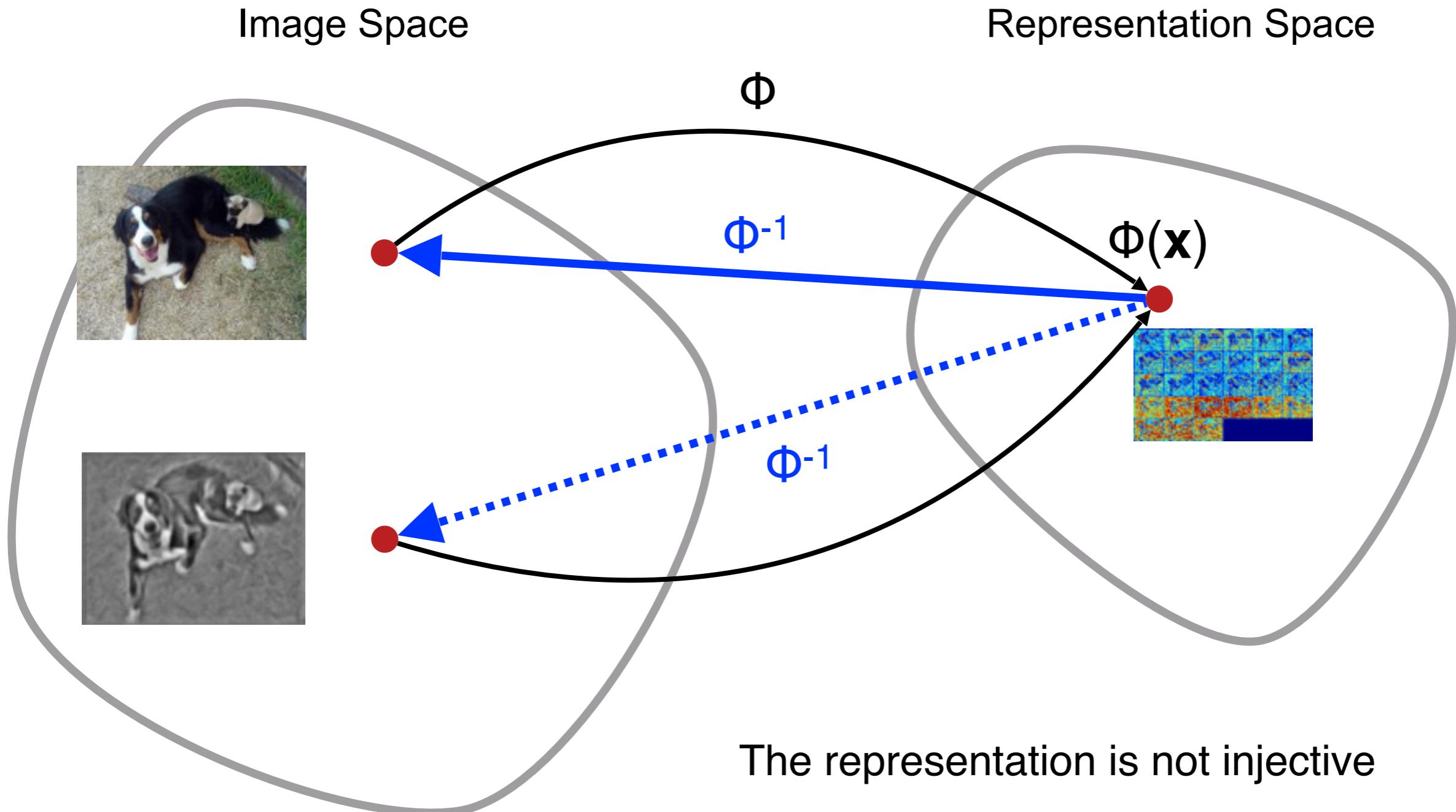
Visualizing representations

Backpropagation networks and “deconvolution”

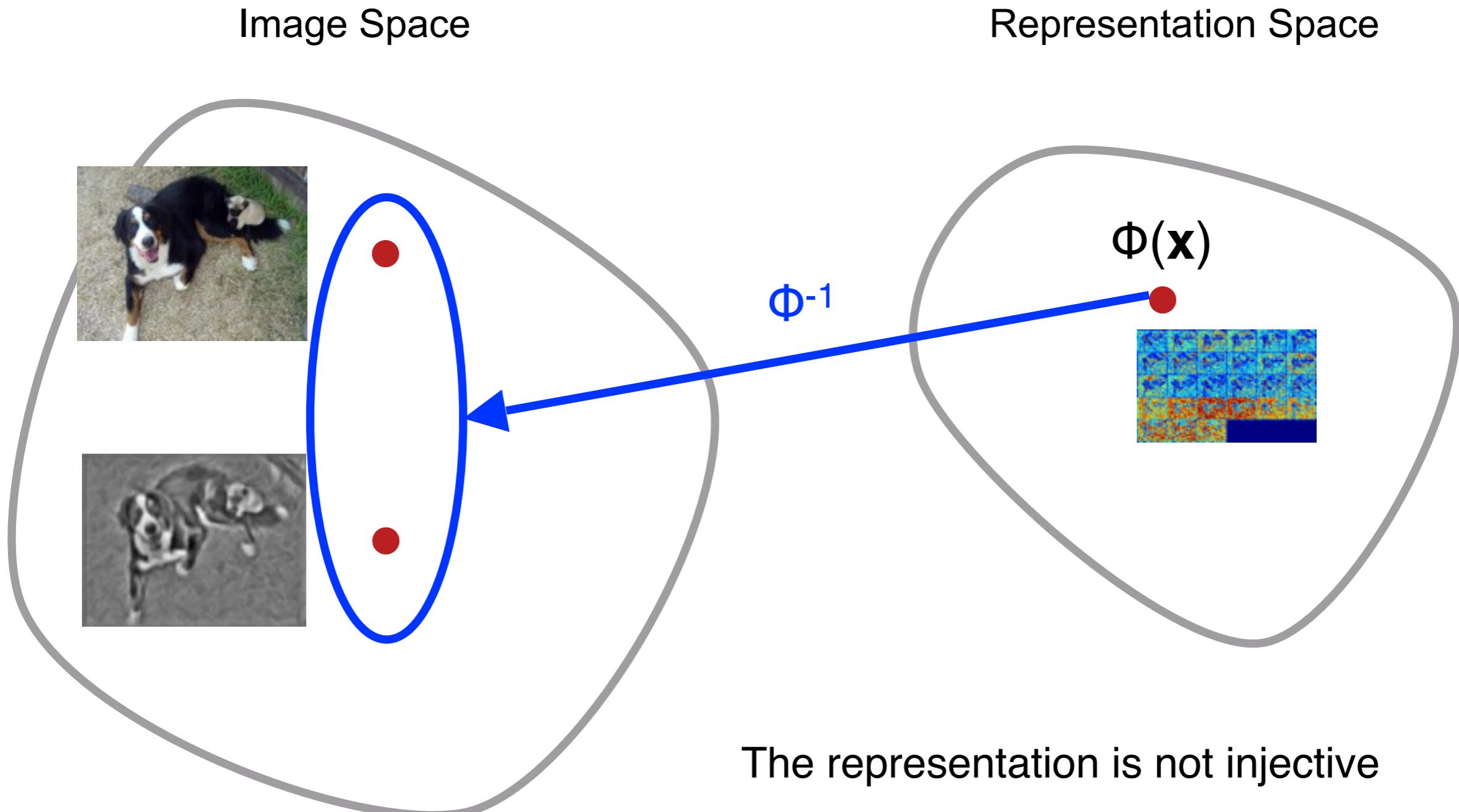
Representations: equivalence & transformations

Visualization: Pre-Image

50



Visualization: Pre-Image

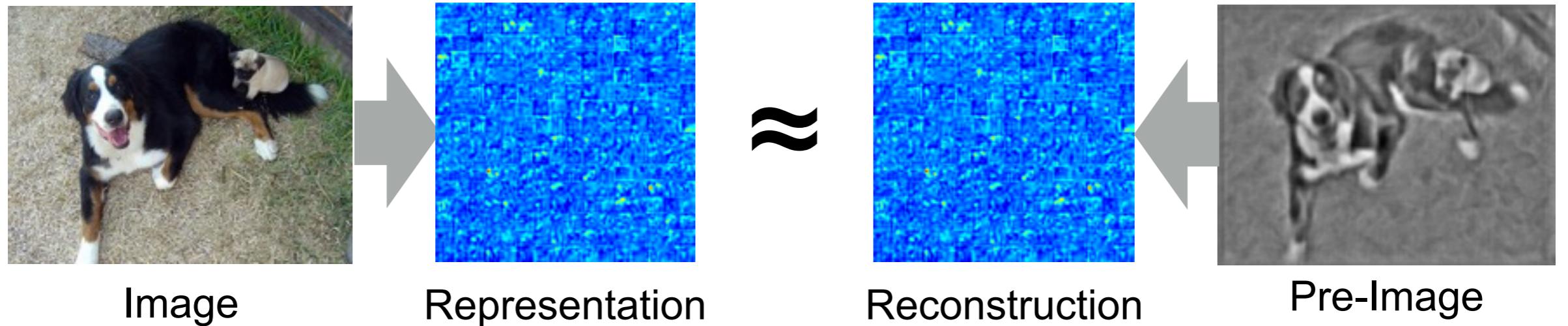


The reconstruction ambiguity **provides useful information about the representation**

Finding a Pre-Image

A simple yet general and effective method

$$\min_x \|\Phi(x) - \Phi_0\|_2^2$$



Start from **random noise**

Optimize using stochastic **gradient descent**

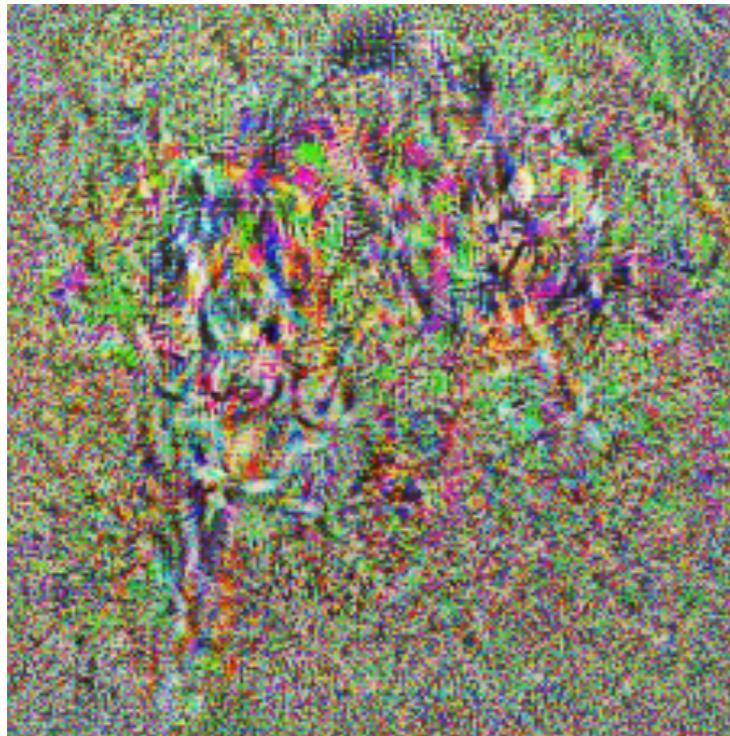
Finding a Pre-Image

53

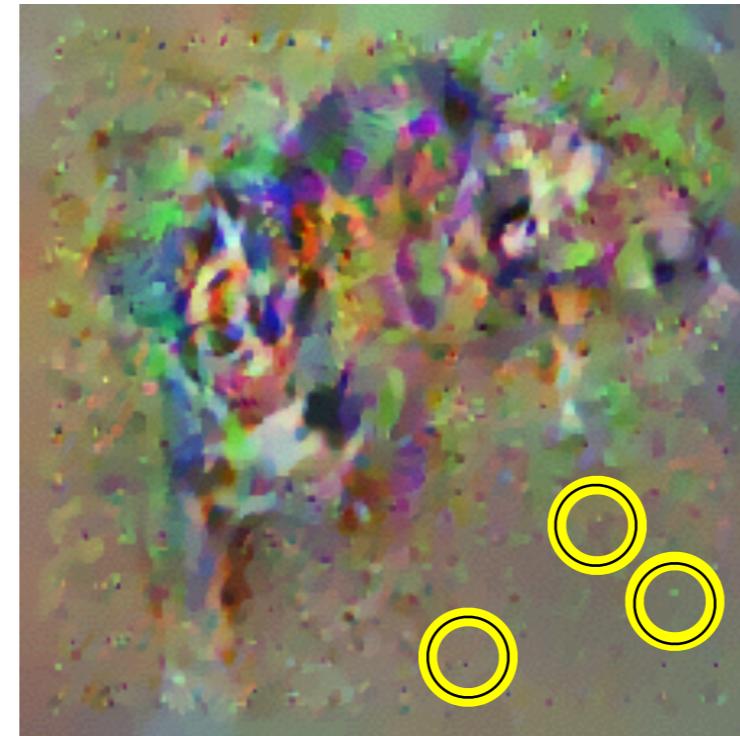
A simple yet general and effective method

$$\min_x \|\Phi(x) - \Phi_0\|_2^2$$

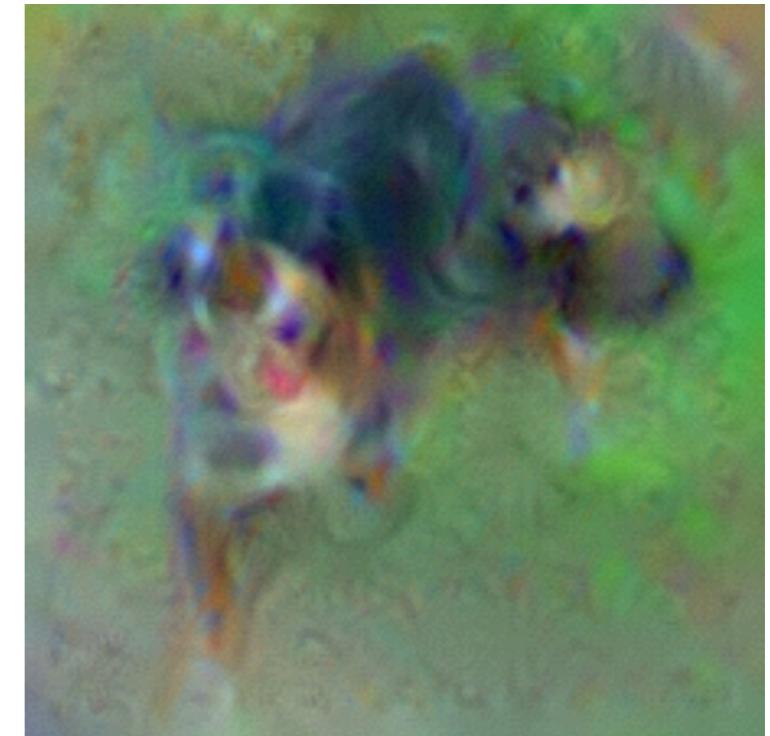
No prior



TV-norm $\beta = 1$



TV-norm $\beta = 2$



Related Work

Analysis tools

Visualizing higher-layer features of a deep network

Ethan et al. 2009
[intermediate features]

Deep inside convolutional networks

Simonyan et al. 2014
[deepest features, aka “deep dreams”]

DeConvNets

Zeiler et al. In ECCV, 2014
[intermediate features]

Understanding neural networks through deep visualisation

Yosinski et al. 2015
[intermediate features]

Artistic tools

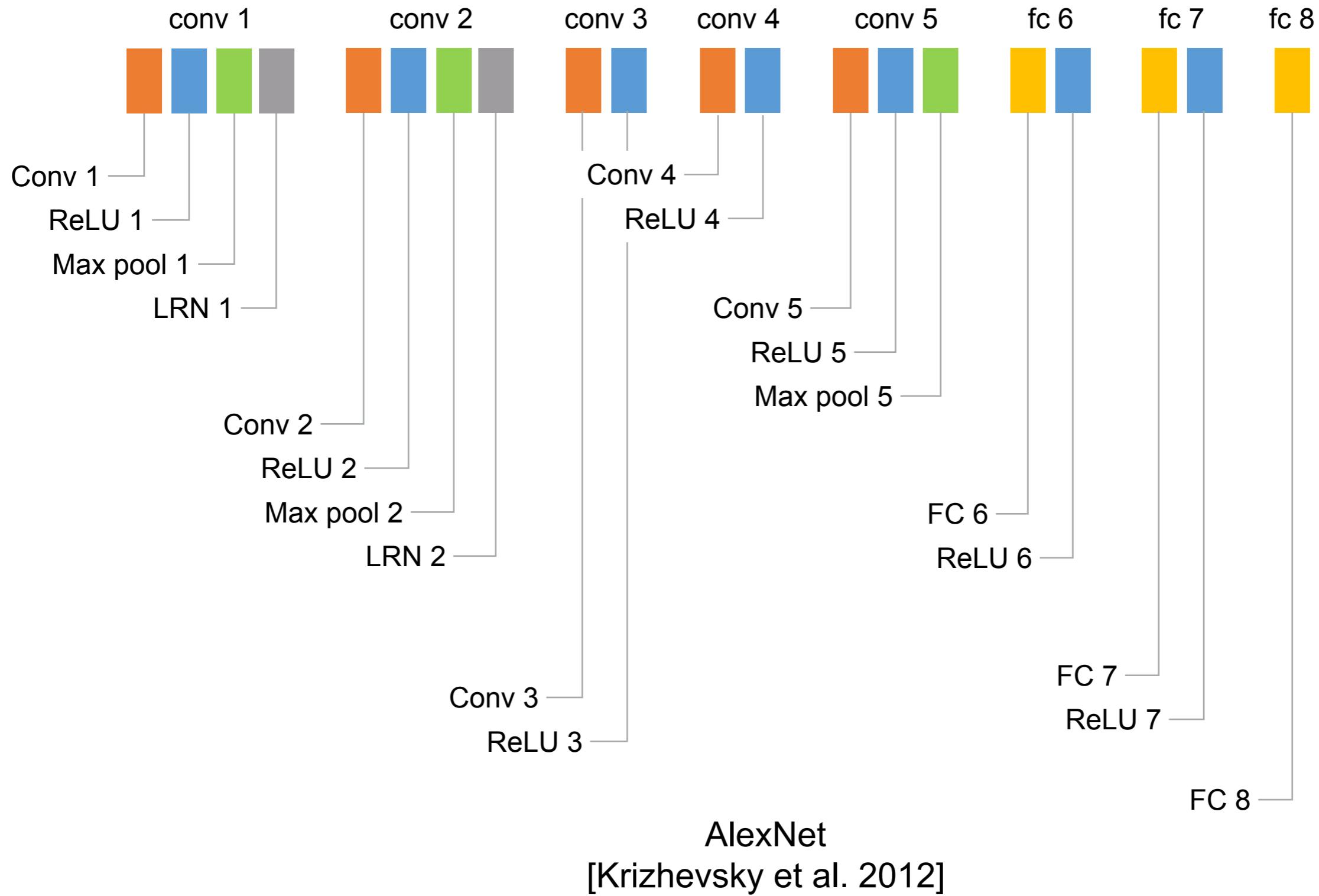
Google’s “inceptionism”

Mordvintsev et al. 2015

Style synthesis and transfer

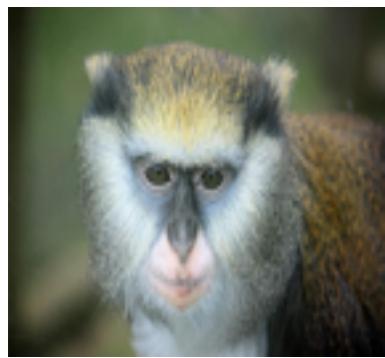
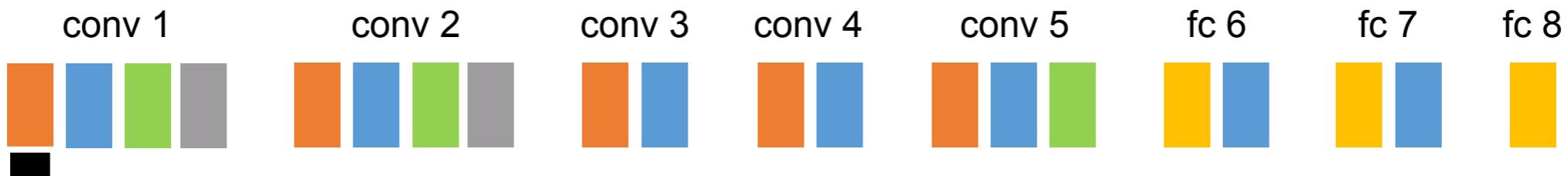
Gatys et al. 2015

Inversion

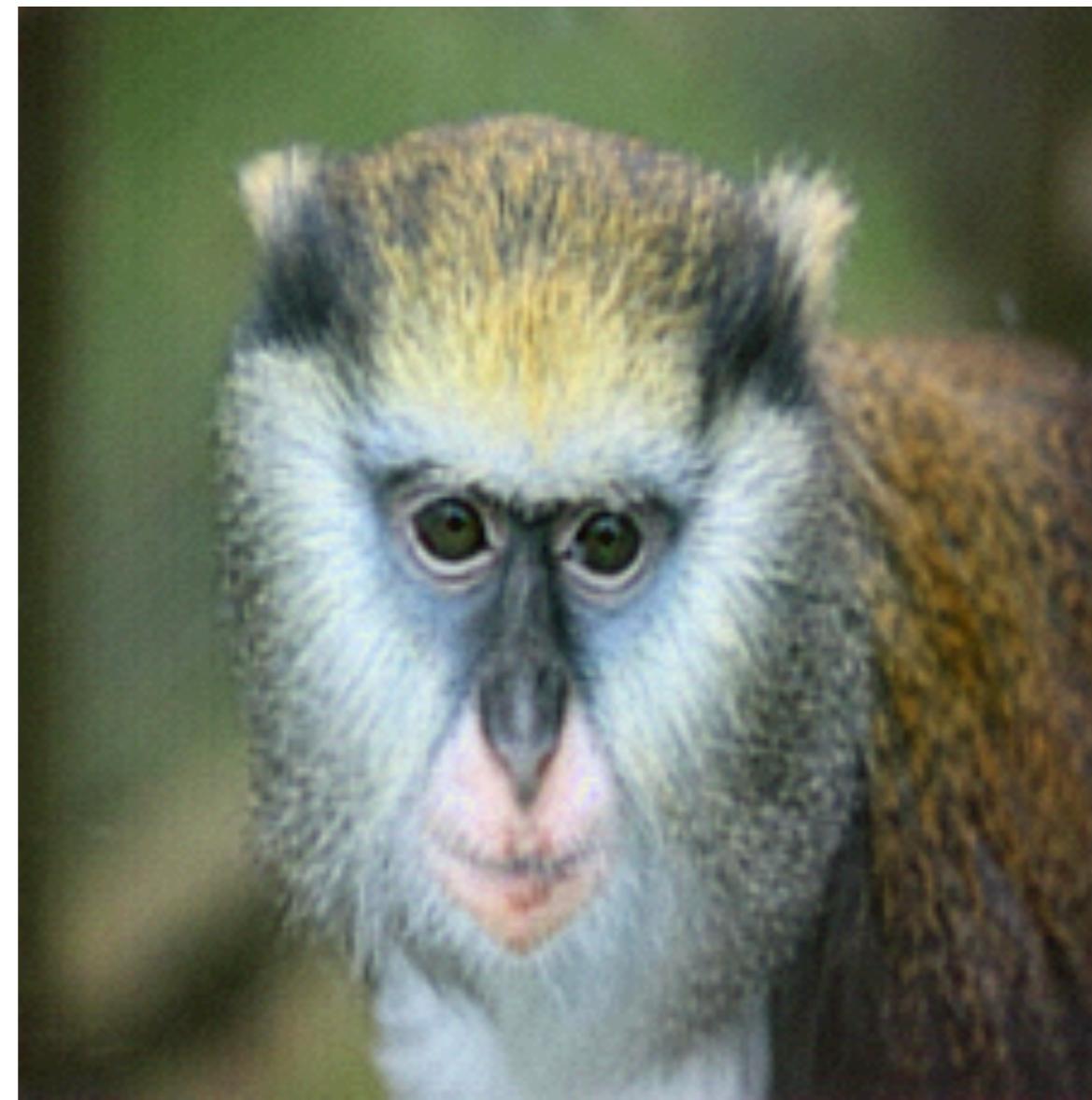


Inversion

56

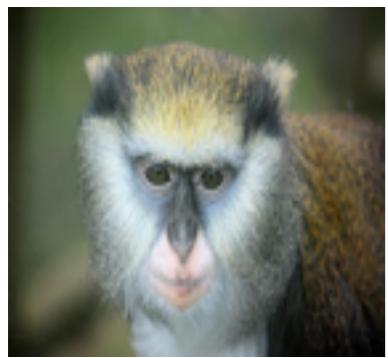
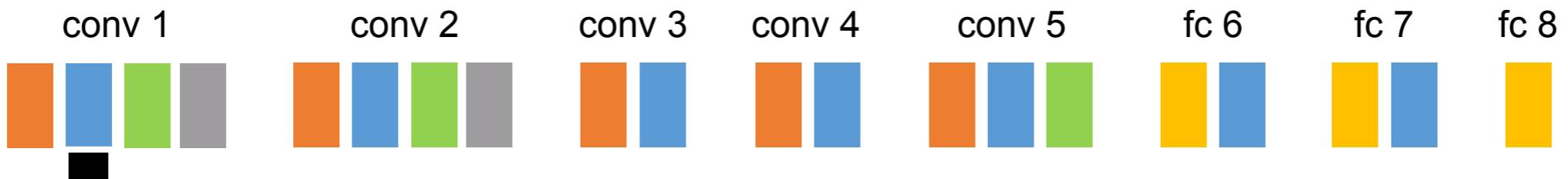


Original
Image



Inversion

57

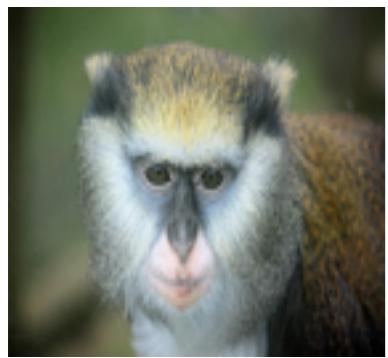


Original
Image



Inversion

58

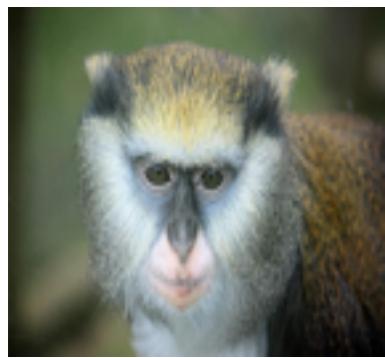
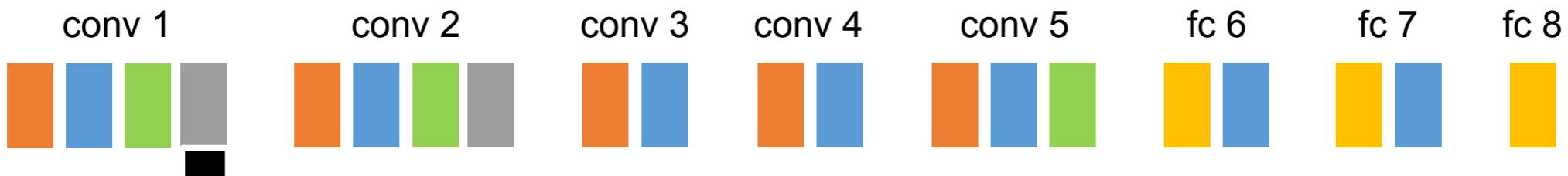


Original
Image

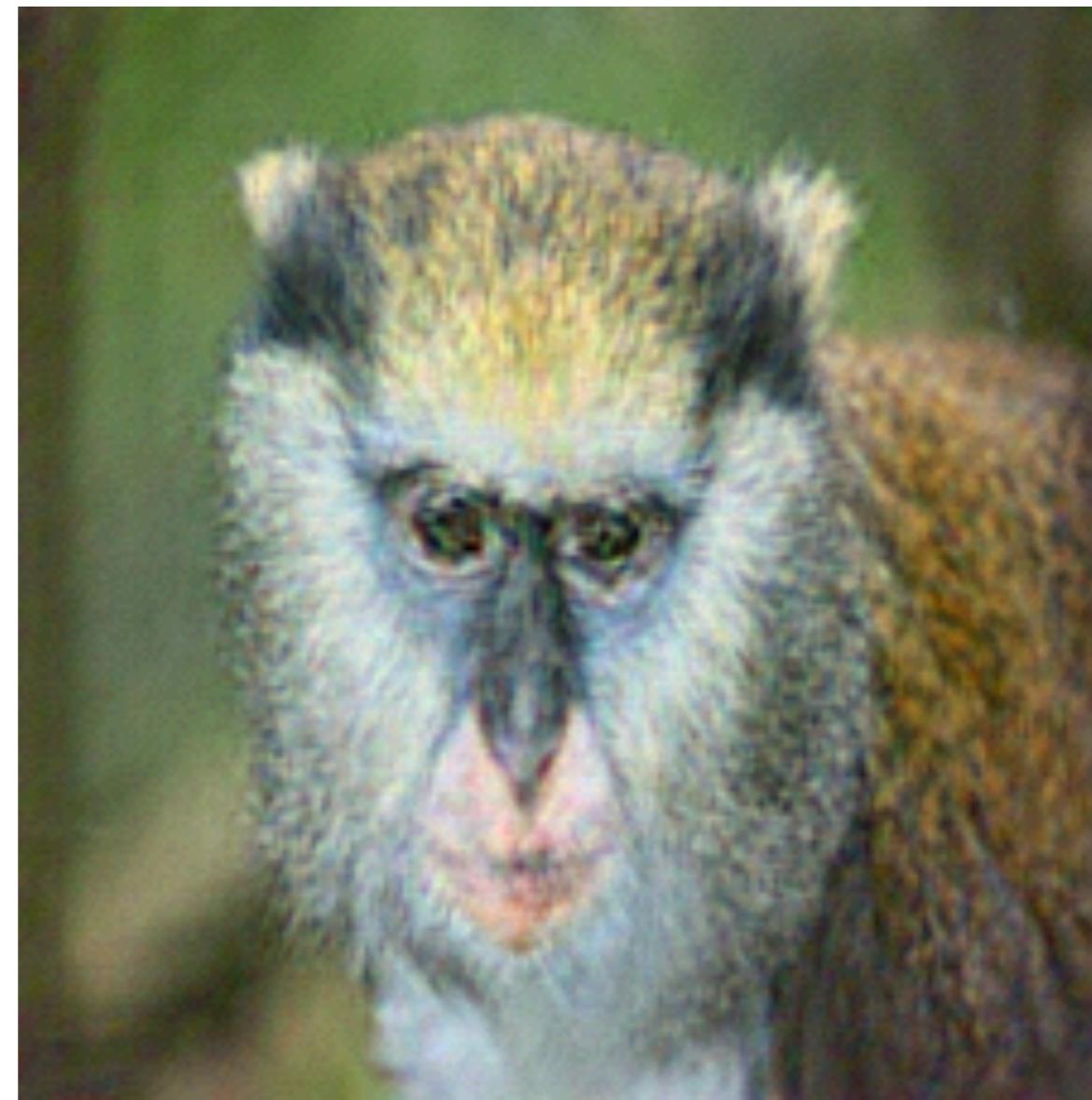


Inversion

59



Original
Image



Inversion

60

conv 1



conv 2



conv 3



conv 4



conv 5



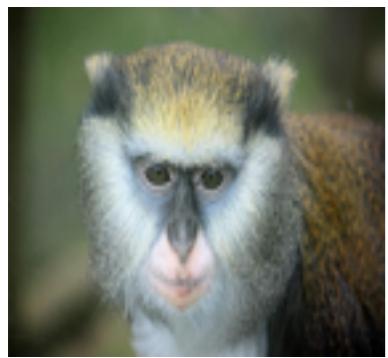
fc 6



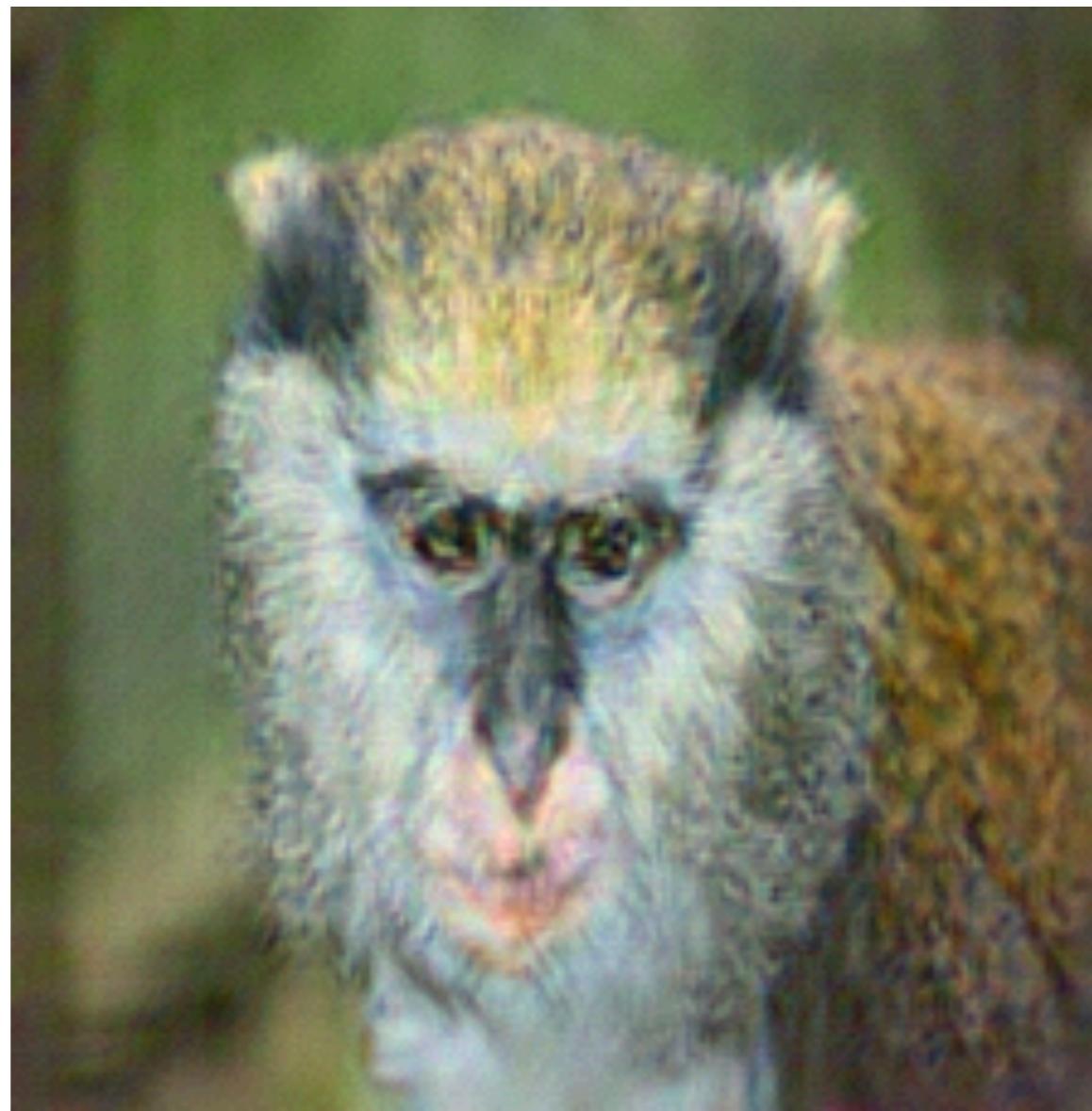
fc 7



fc 8

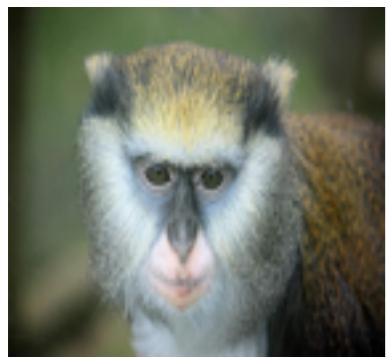


Original
Image

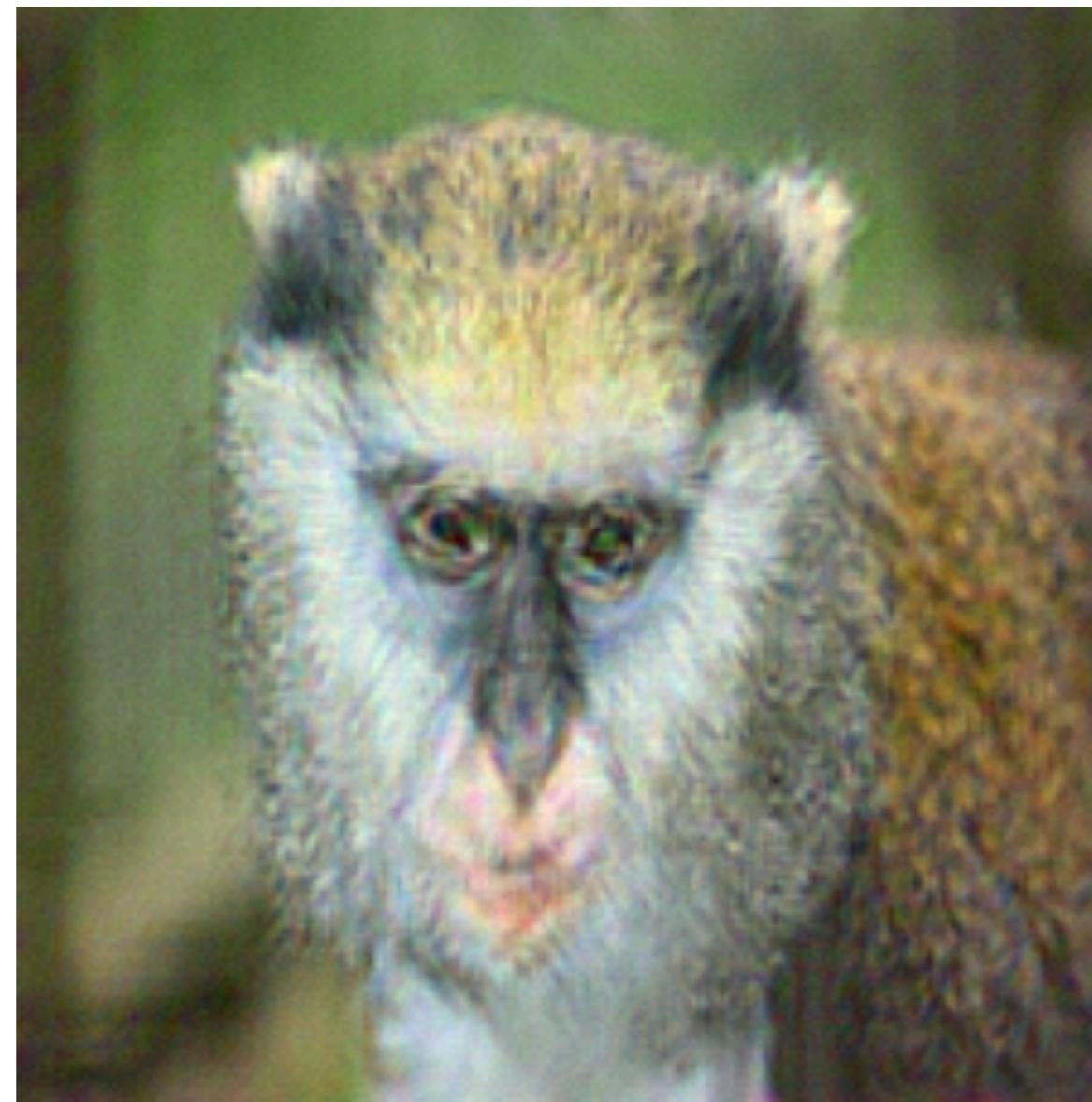


Inversion

61

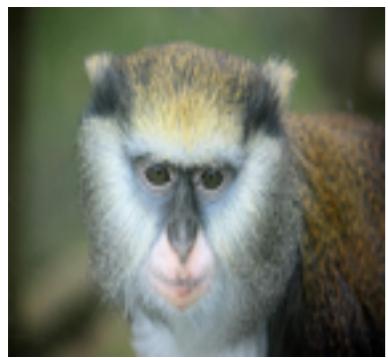


Original
Image

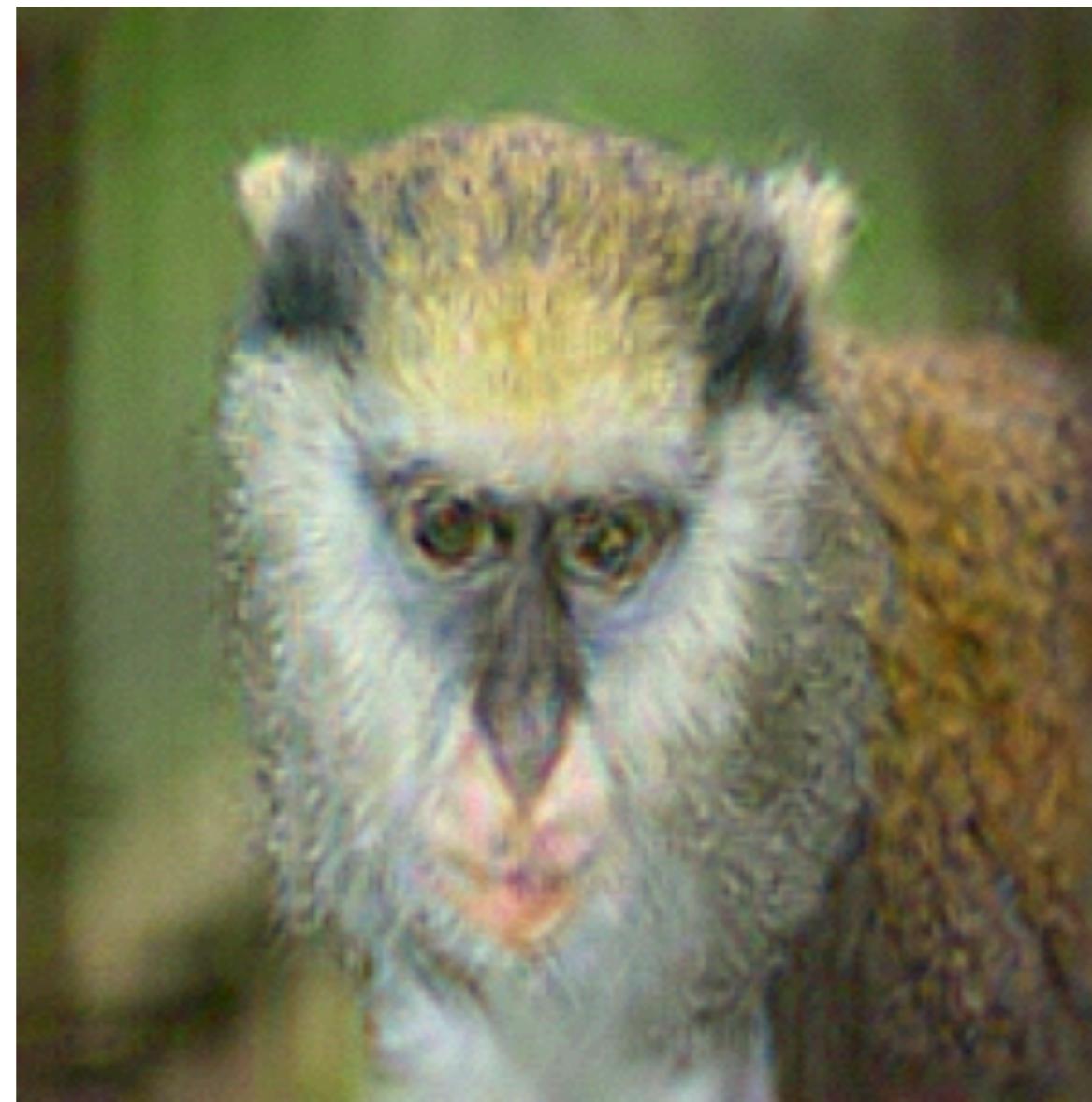


Inversion

62

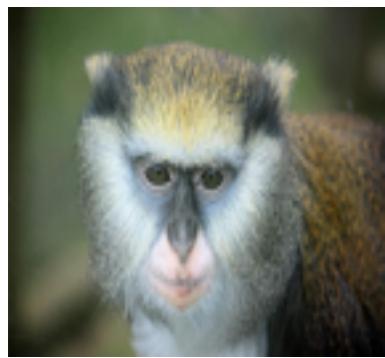


Original
Image

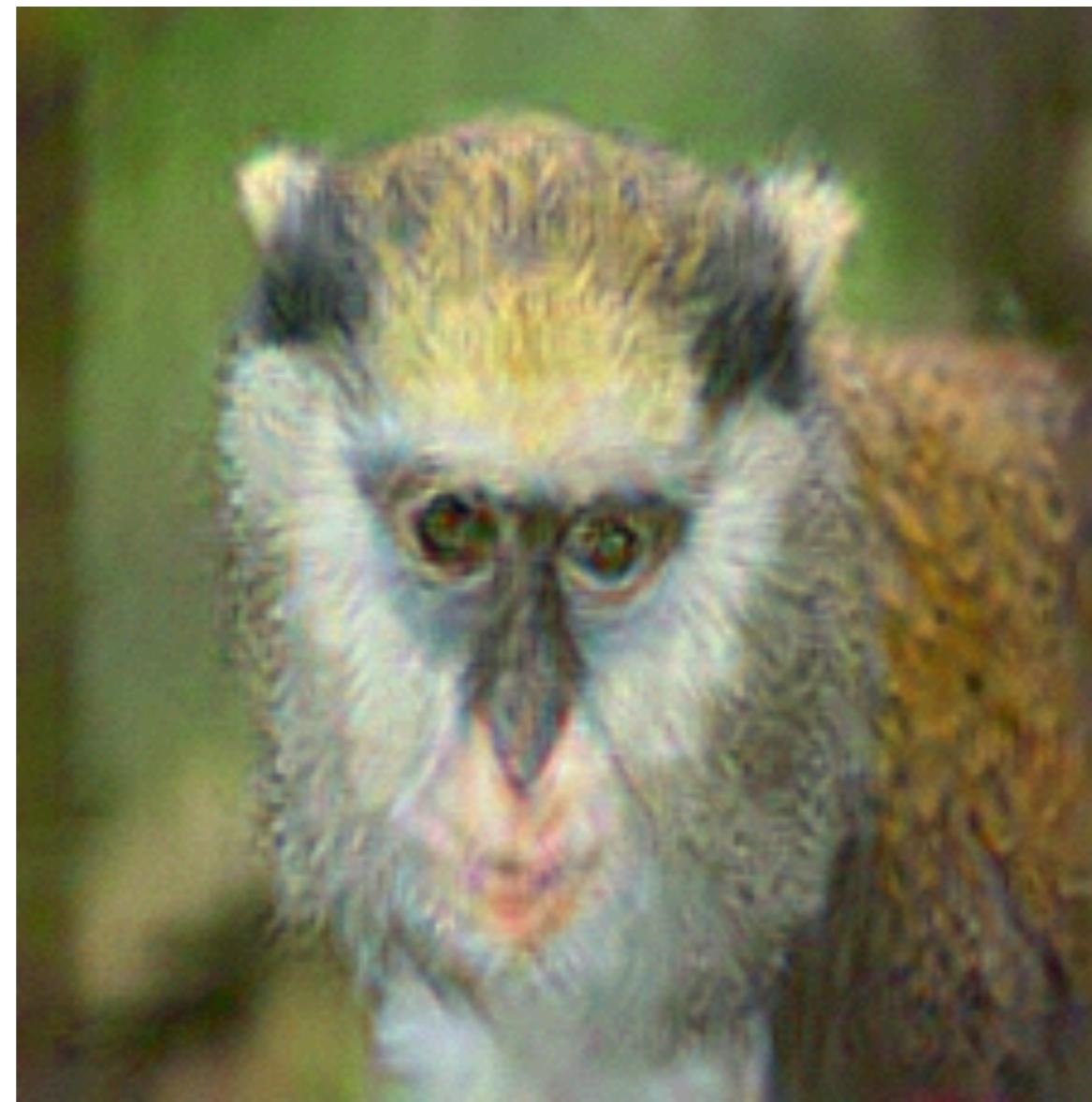


Inversion

63

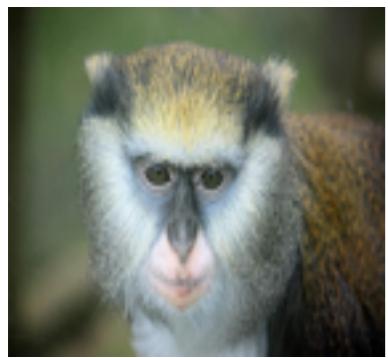
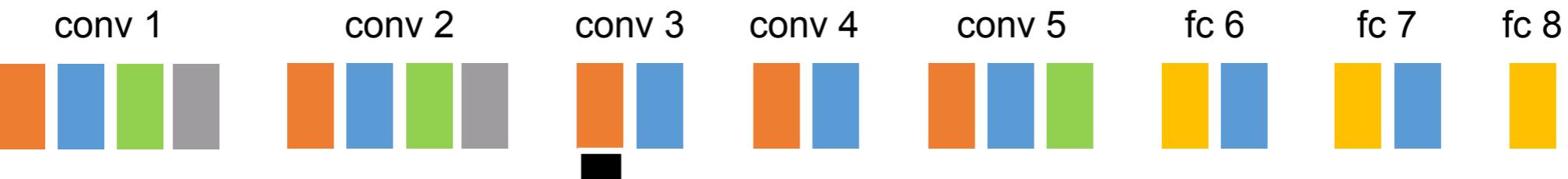


Original
Image

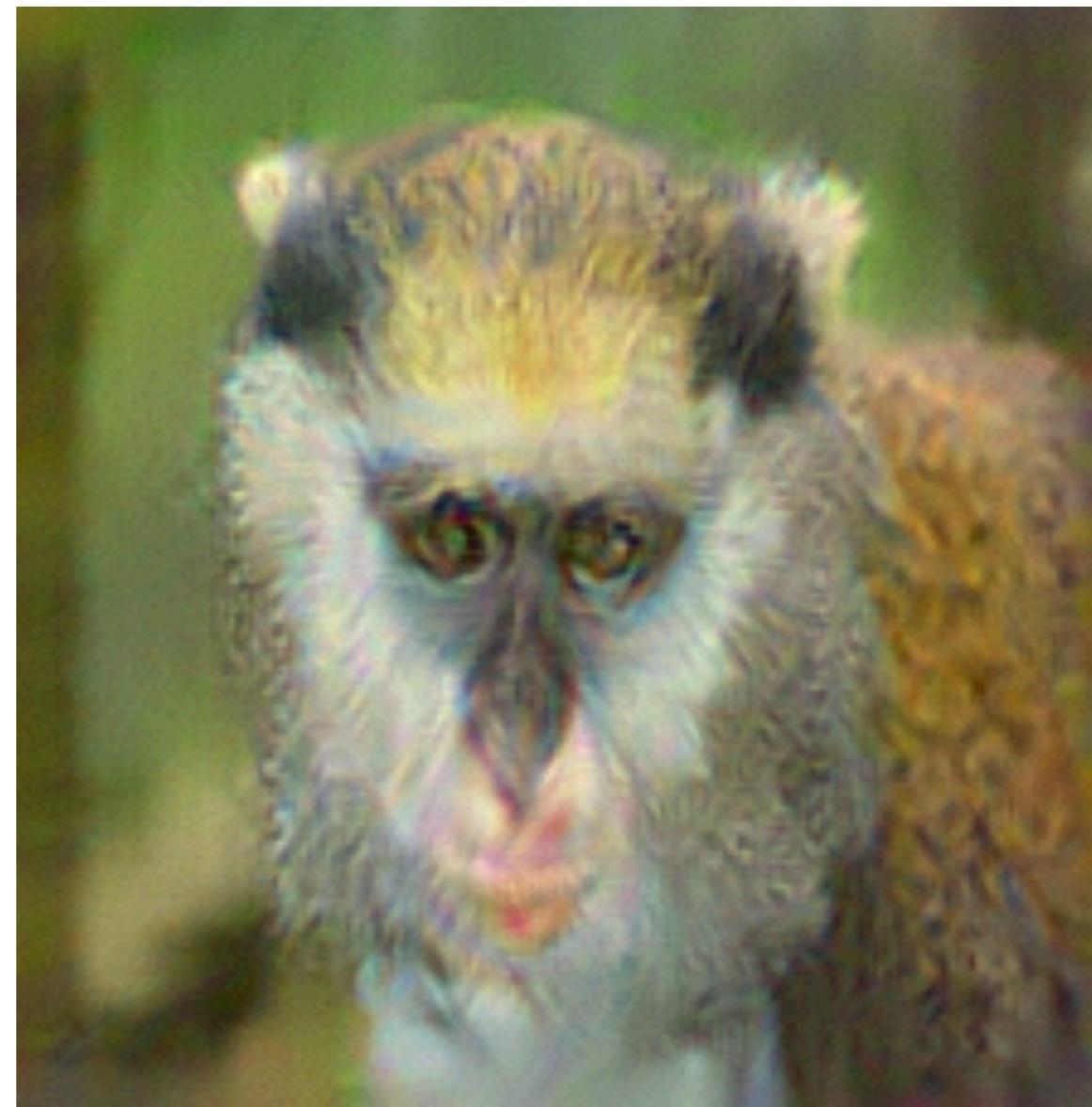


Inversion

64

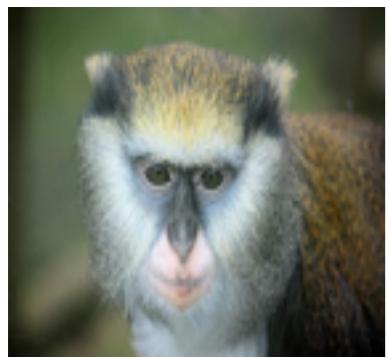
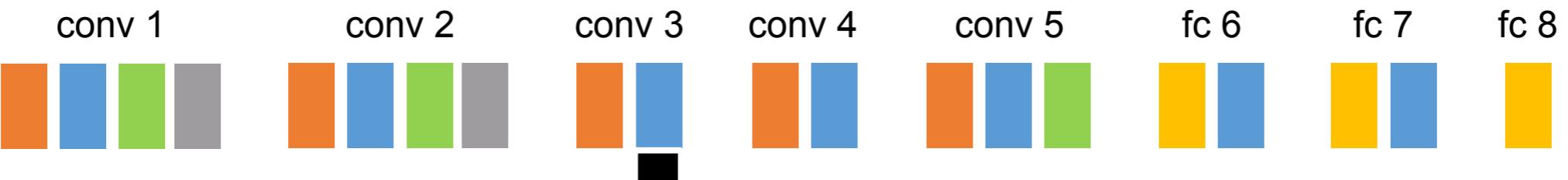


Original
Image

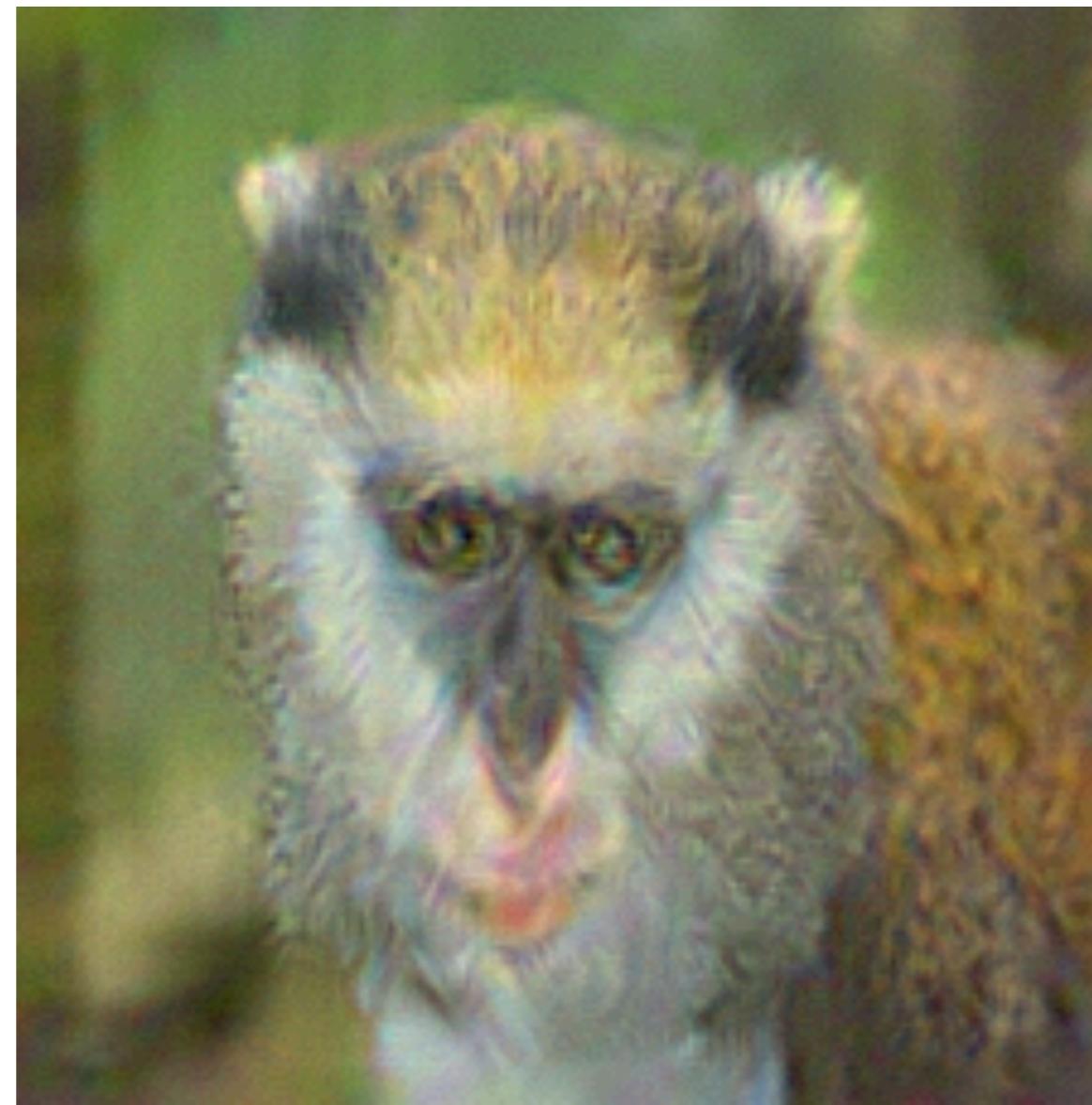


Inversion

65

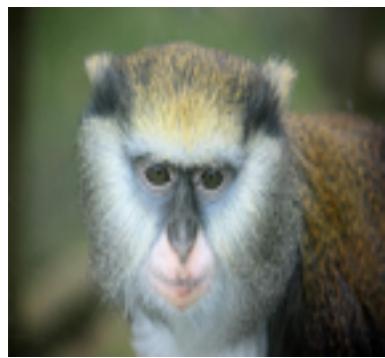
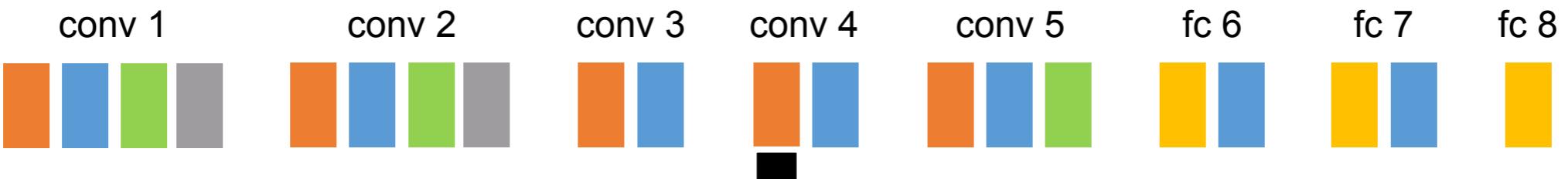


Original
Image

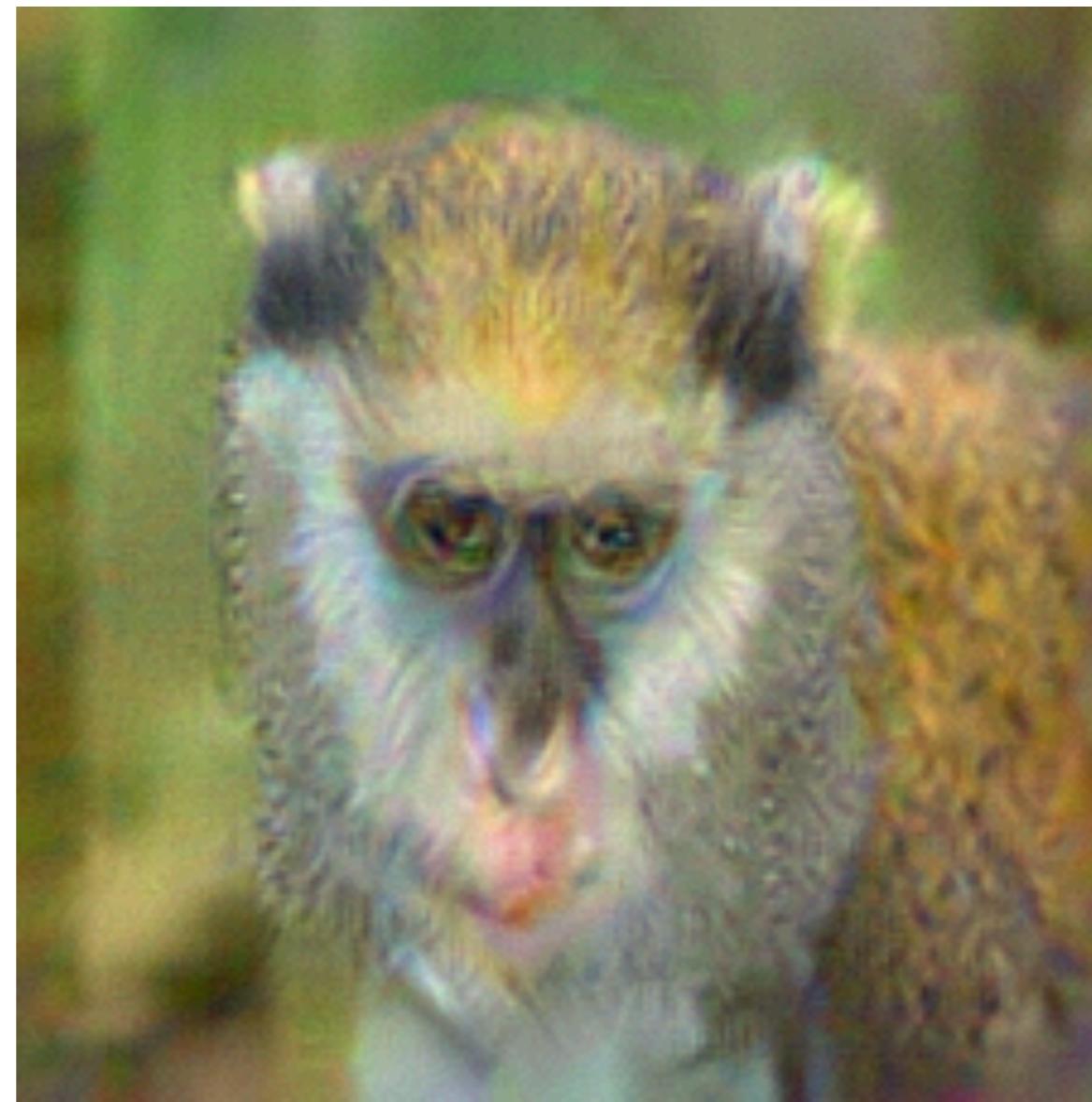


Inversion

66

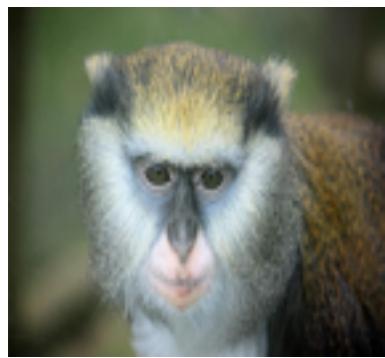
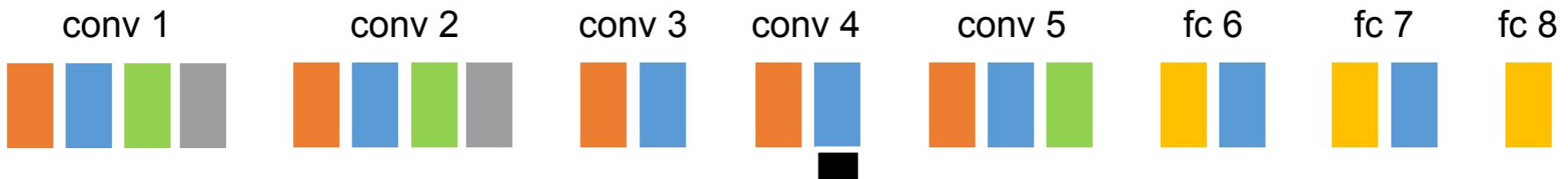


Original
Image

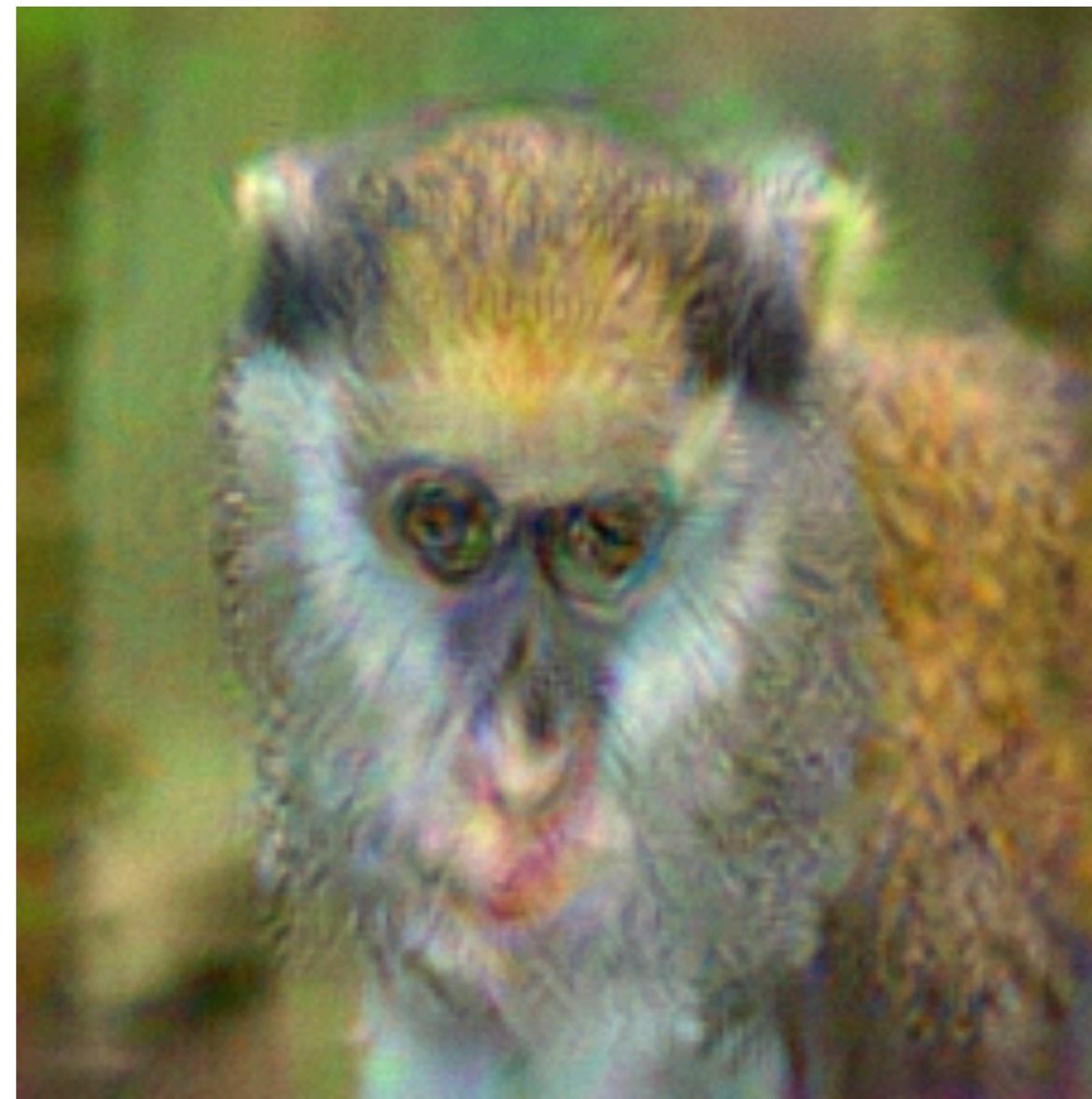


Inversion

67

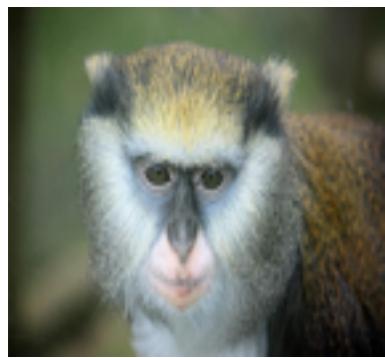


Original
Image

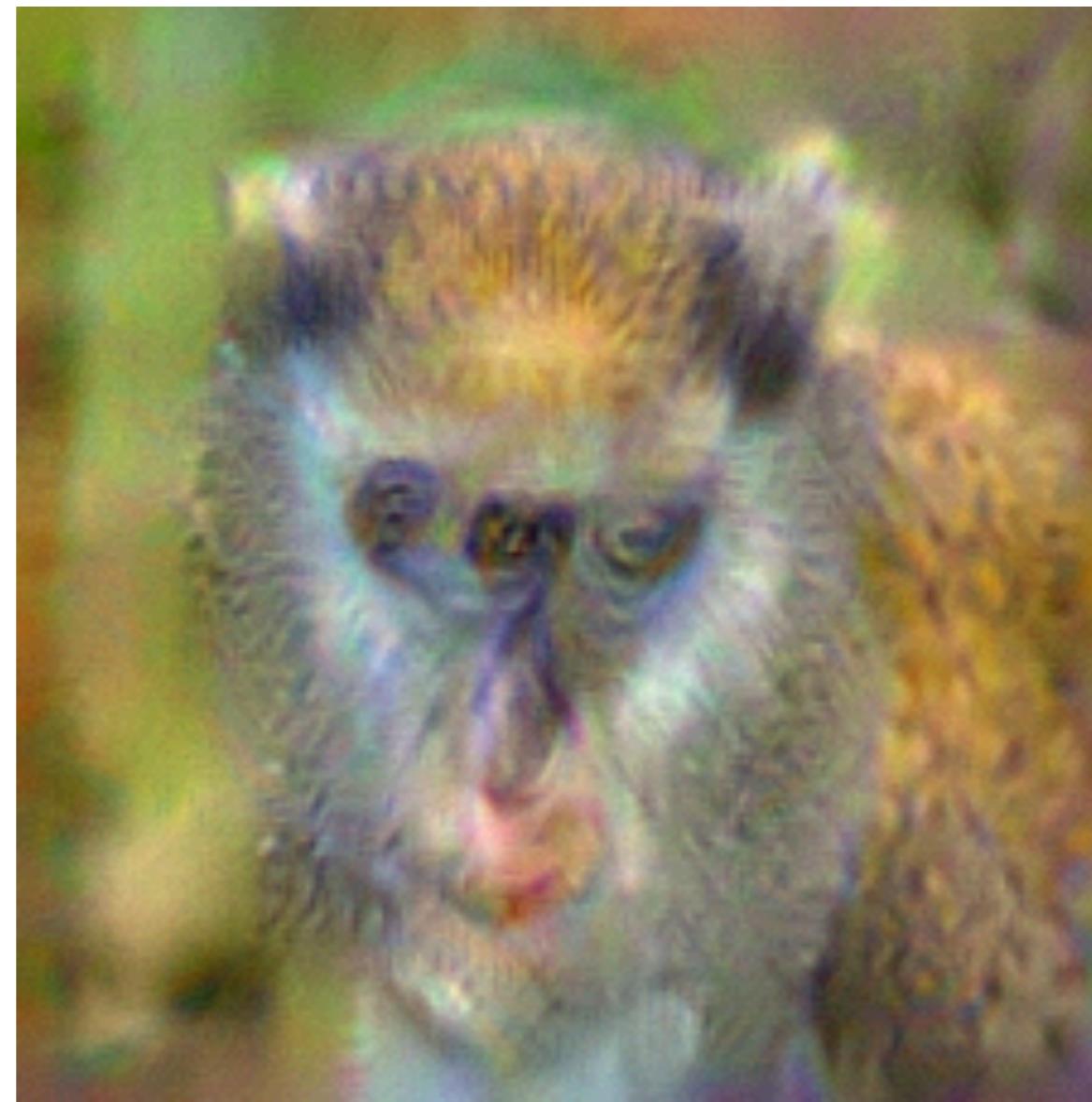


Inversion

68

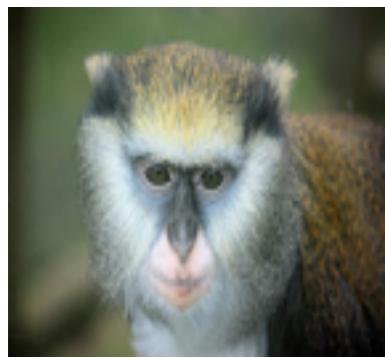


Original
Image

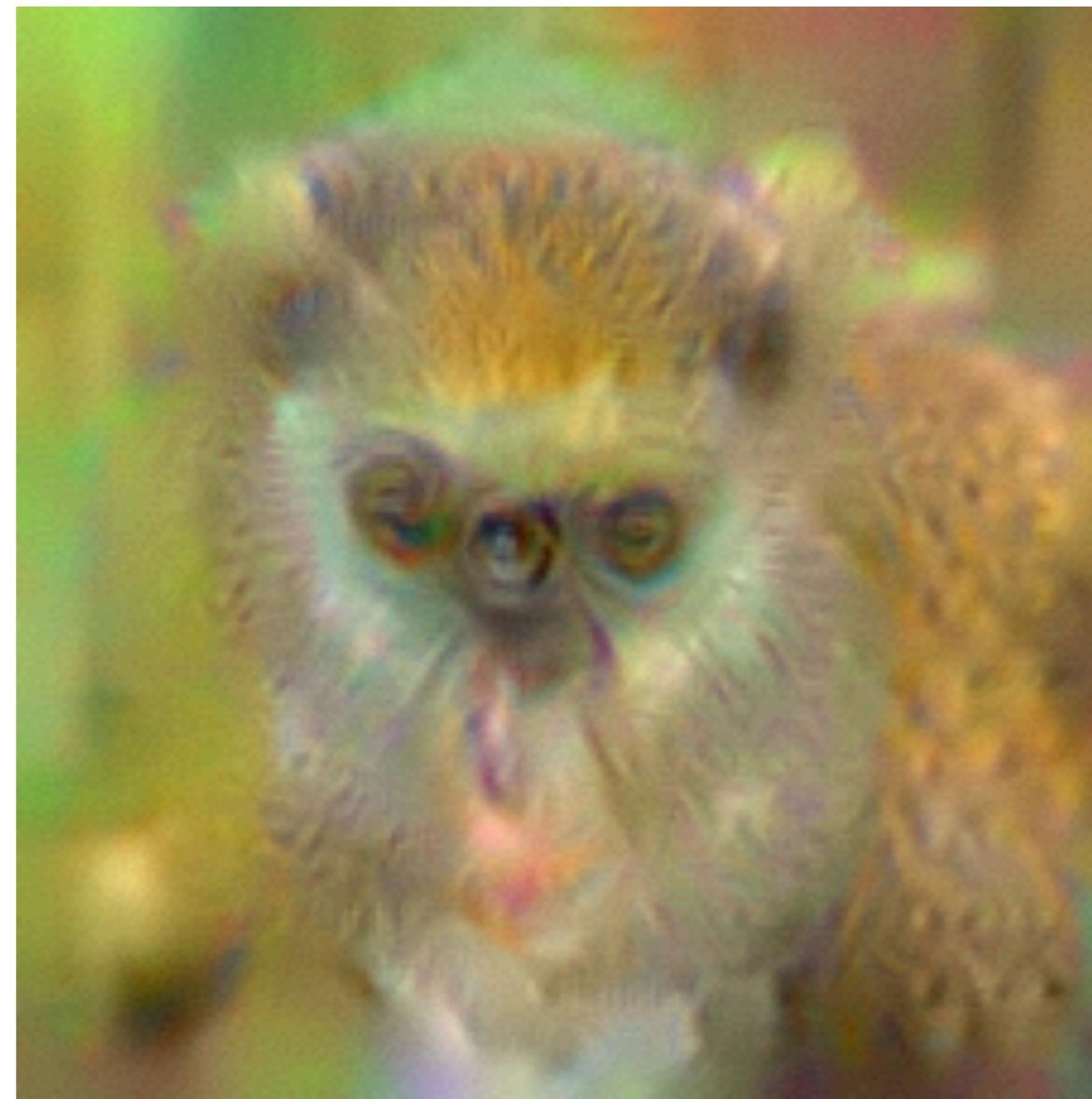


Inversion

69

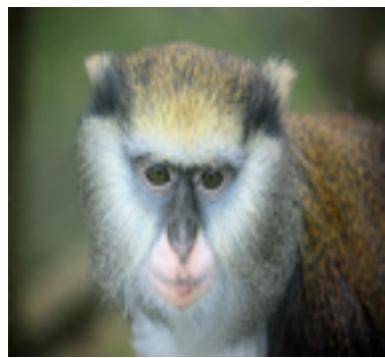


Original
Image

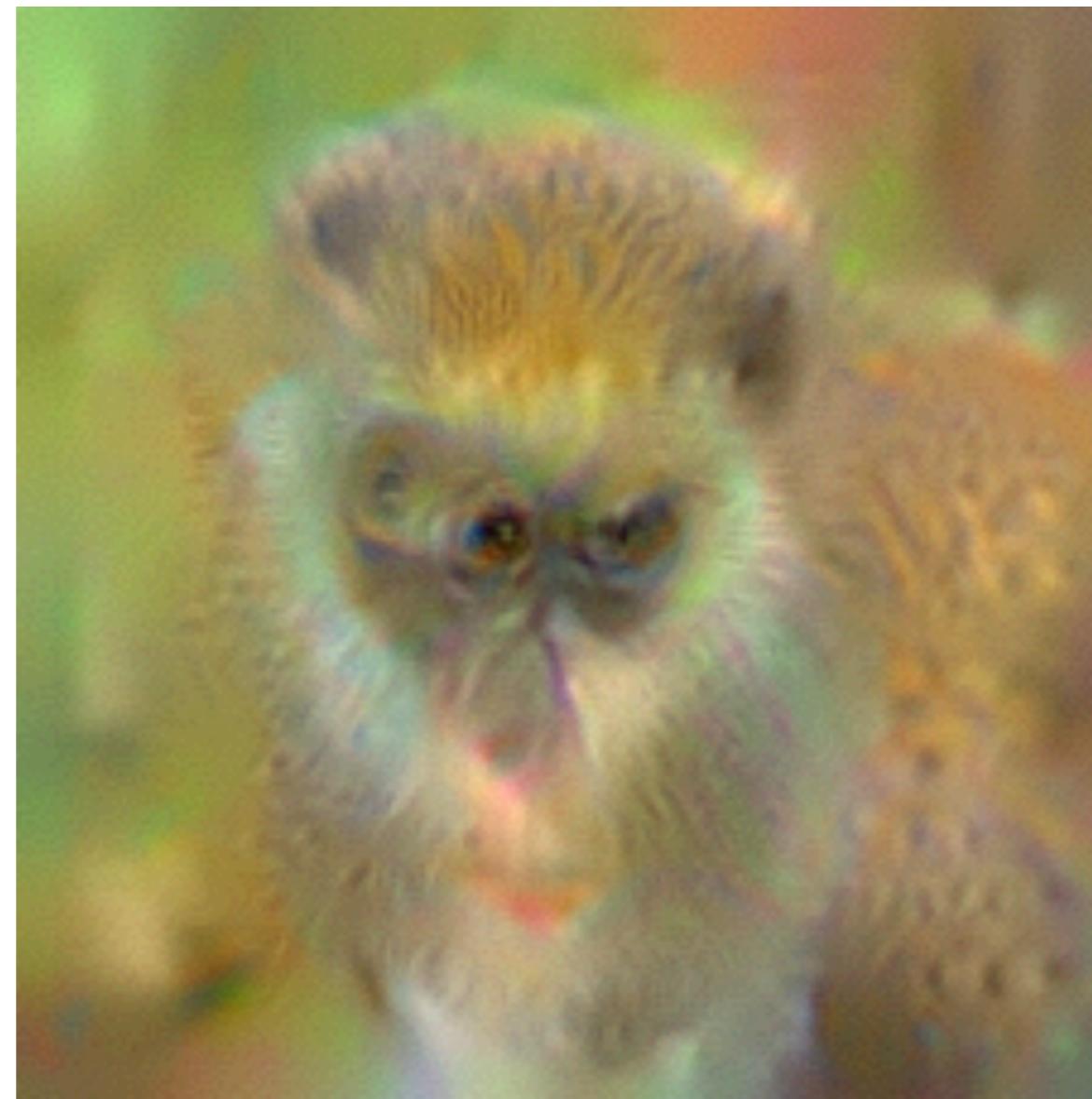


Inversion

70

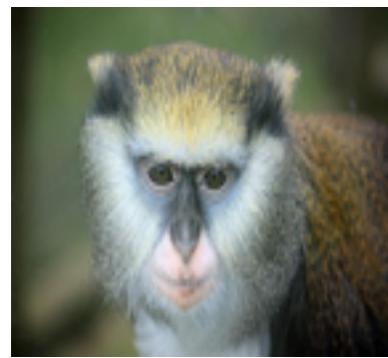
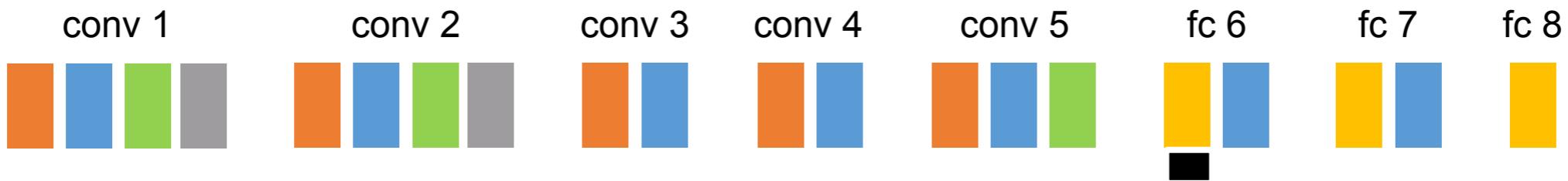


Original
Image

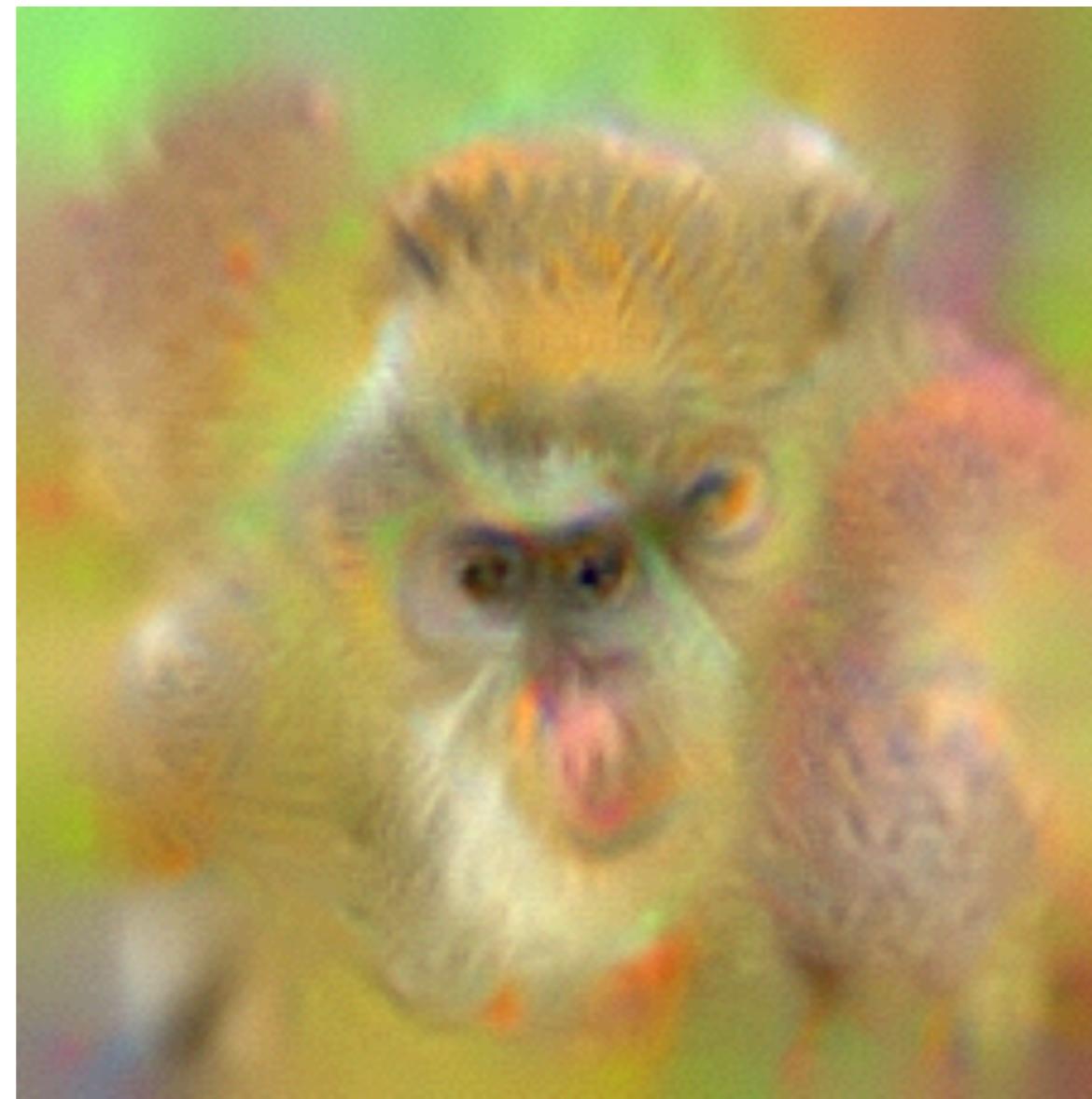


Inversion

71

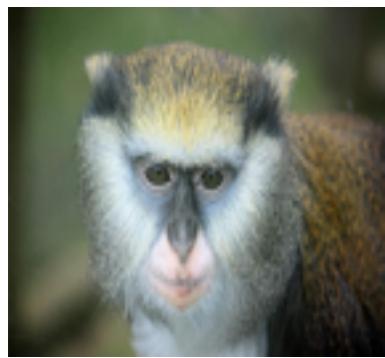
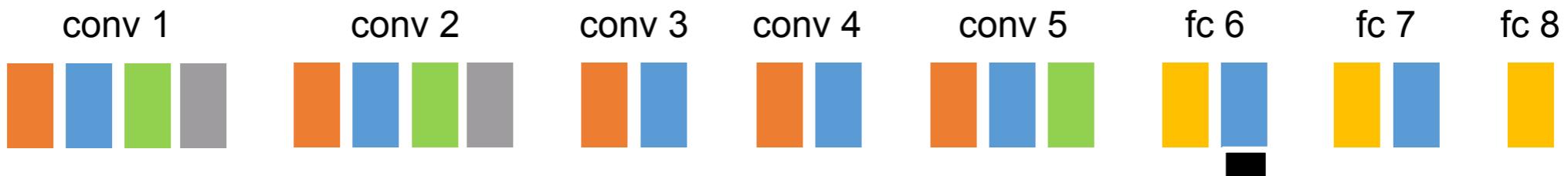


Original
Image

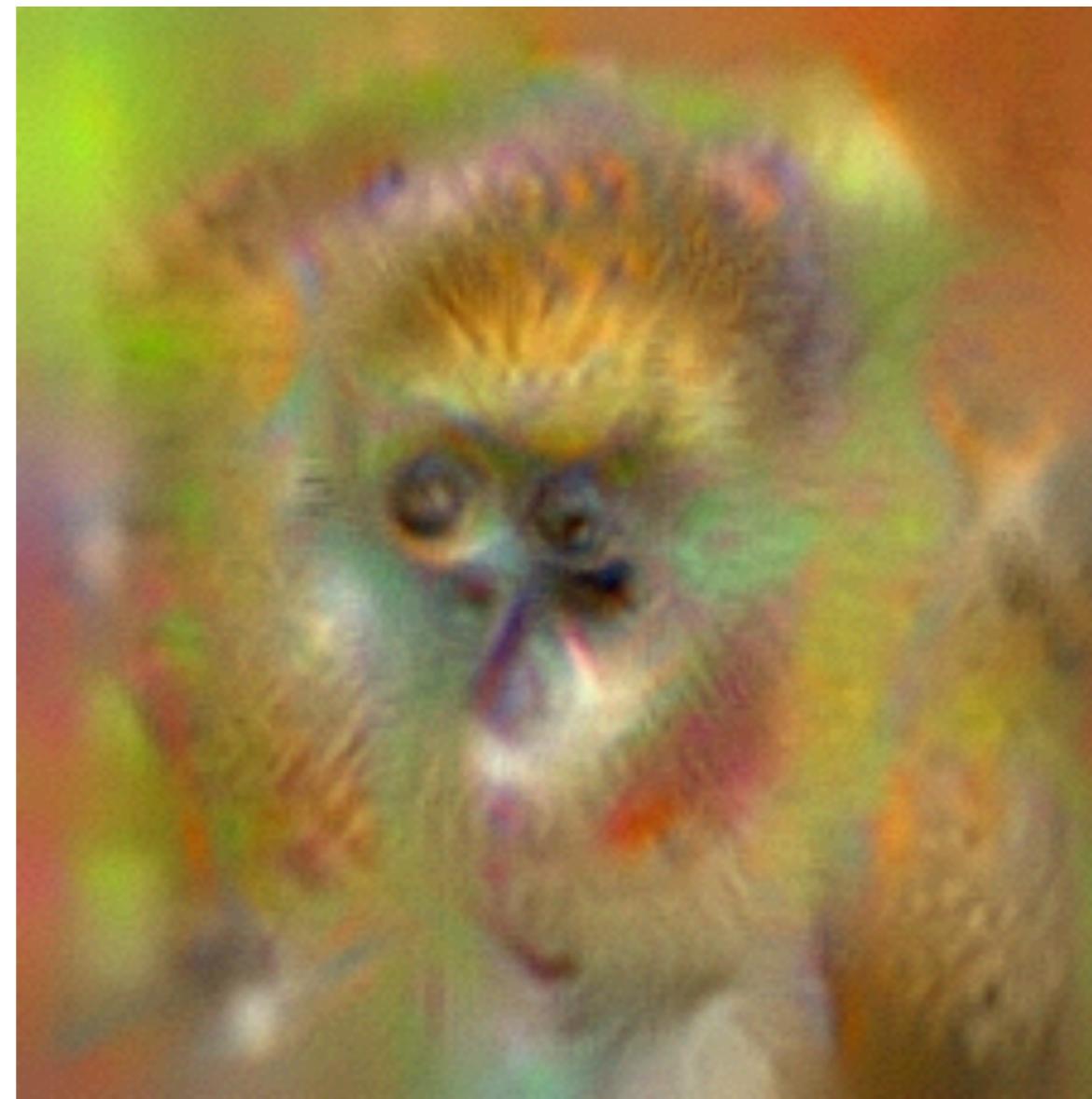


Inversion

72

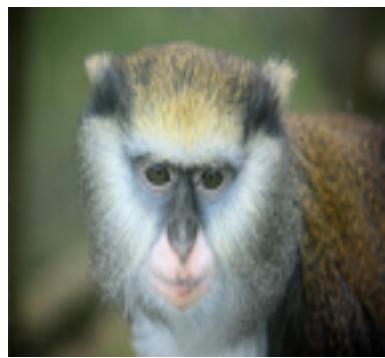


Original
Image

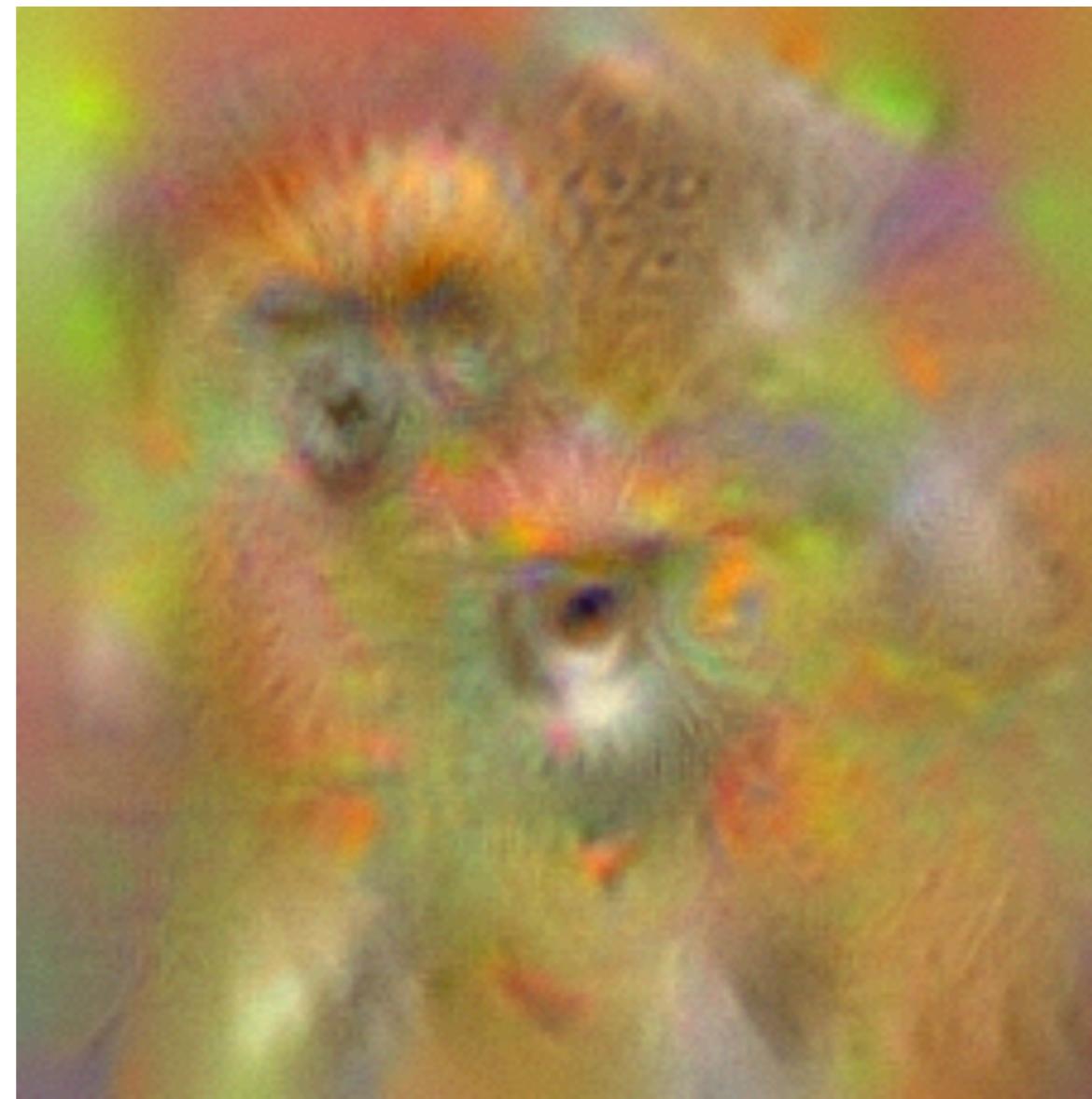


Inversion

73

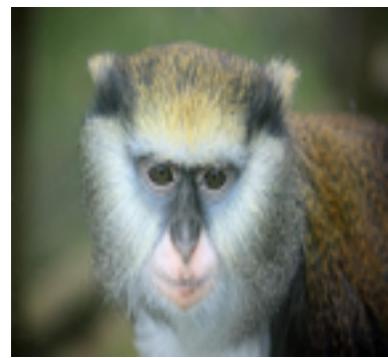
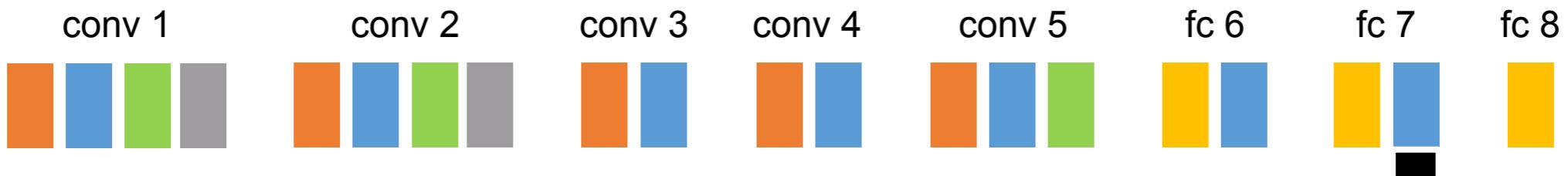


Original
Image

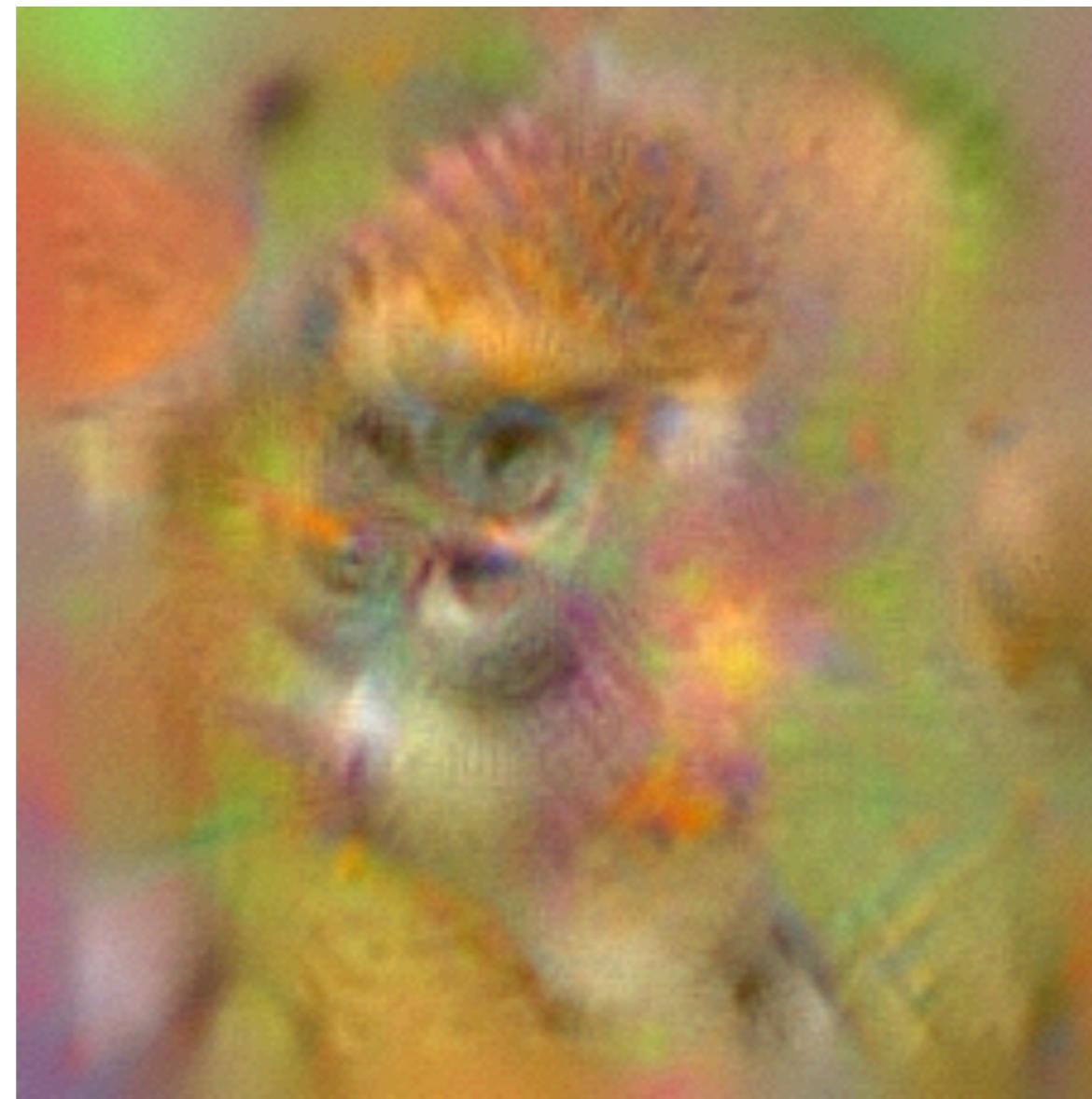


Inversion

74

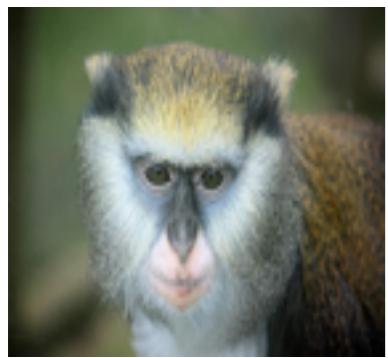
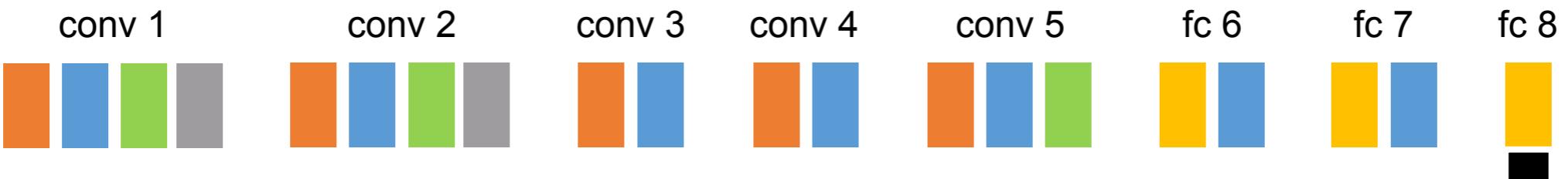


Original
Image

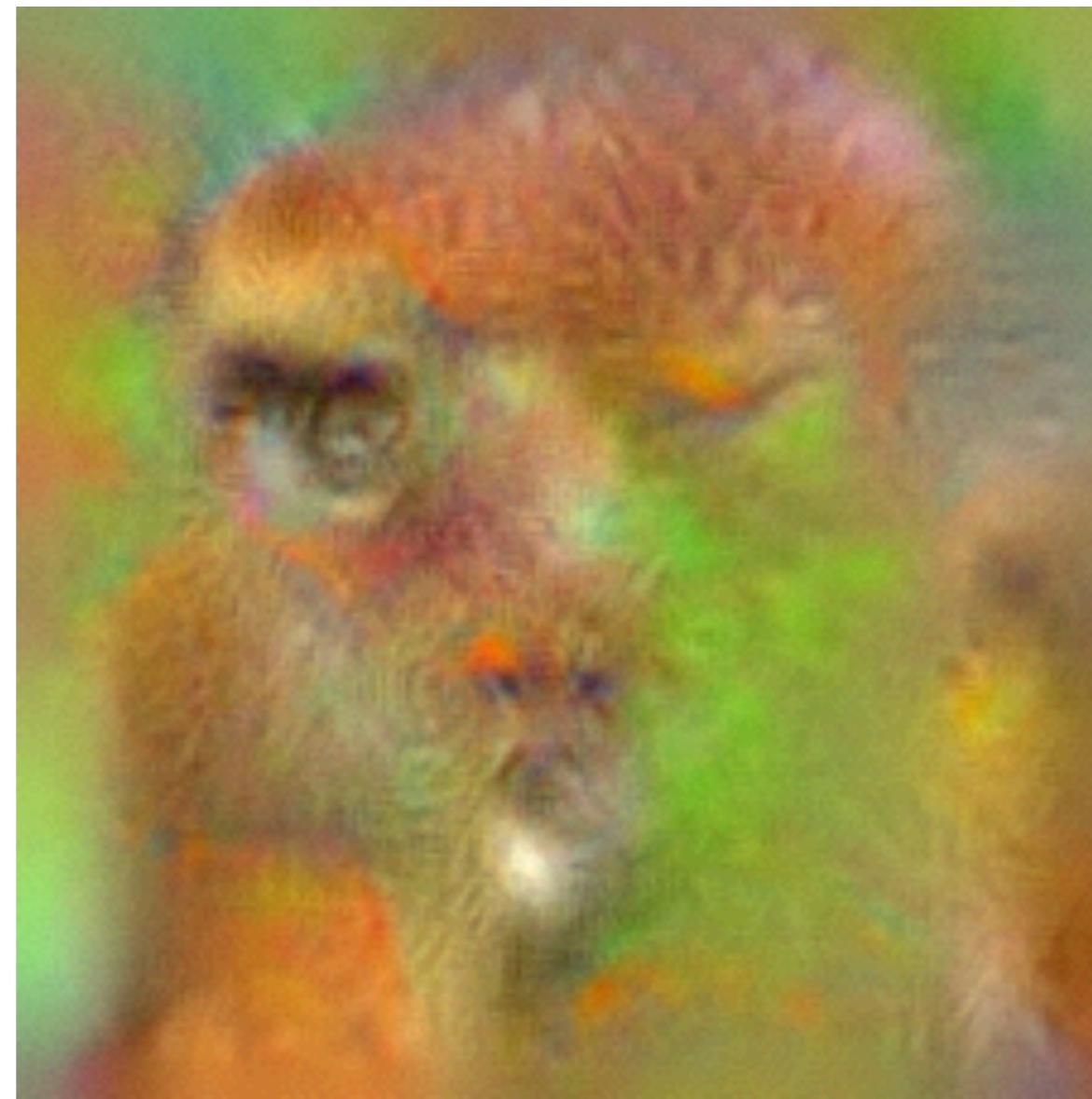


Inversion

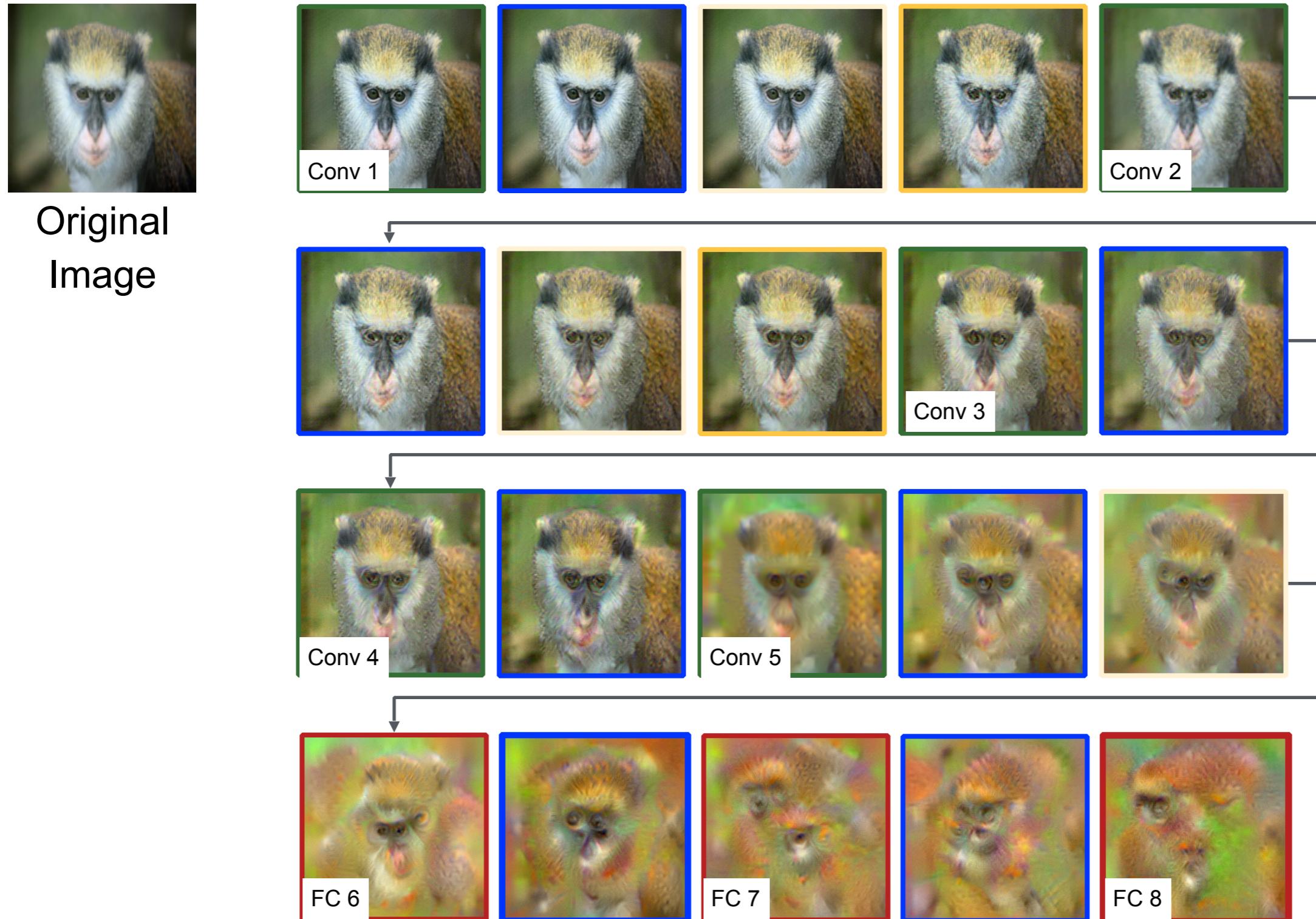
75



Original
Image

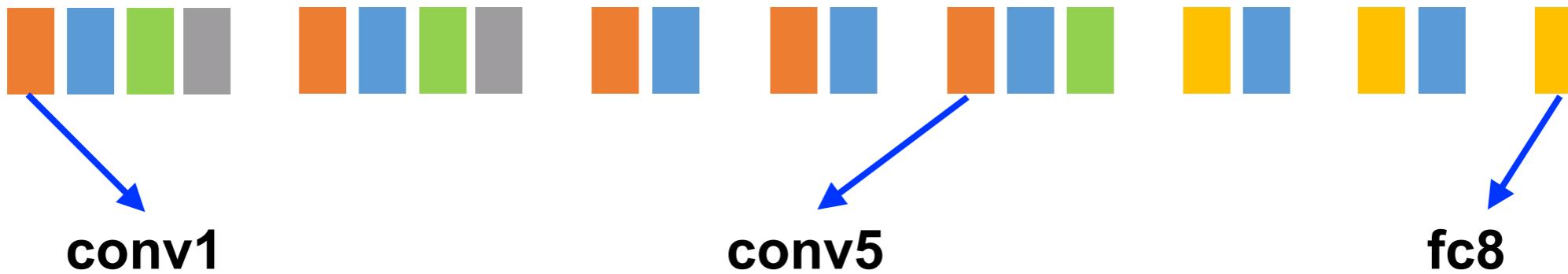


Inverting a Deep CNN



CNNs = visual codes?

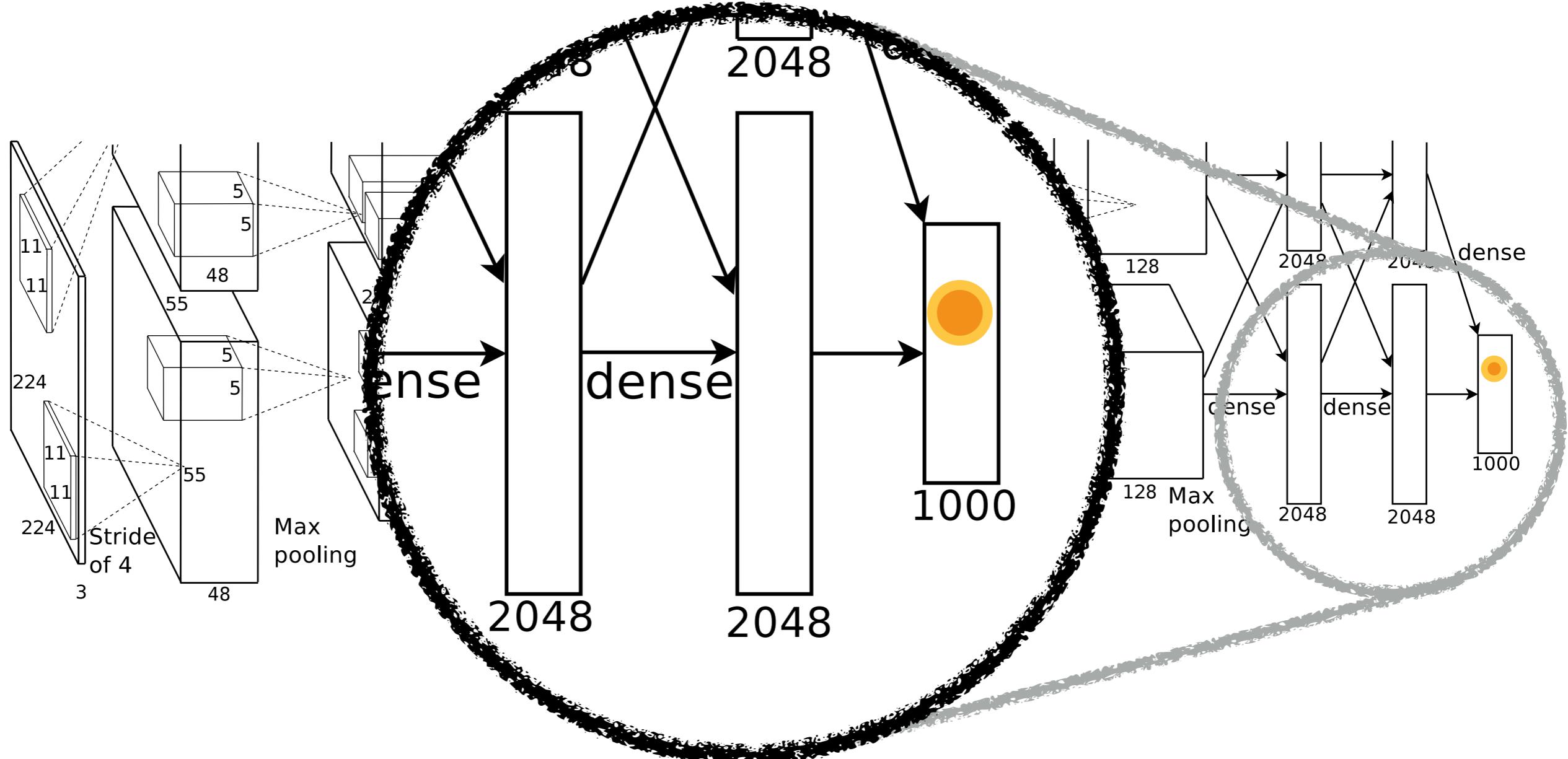
77

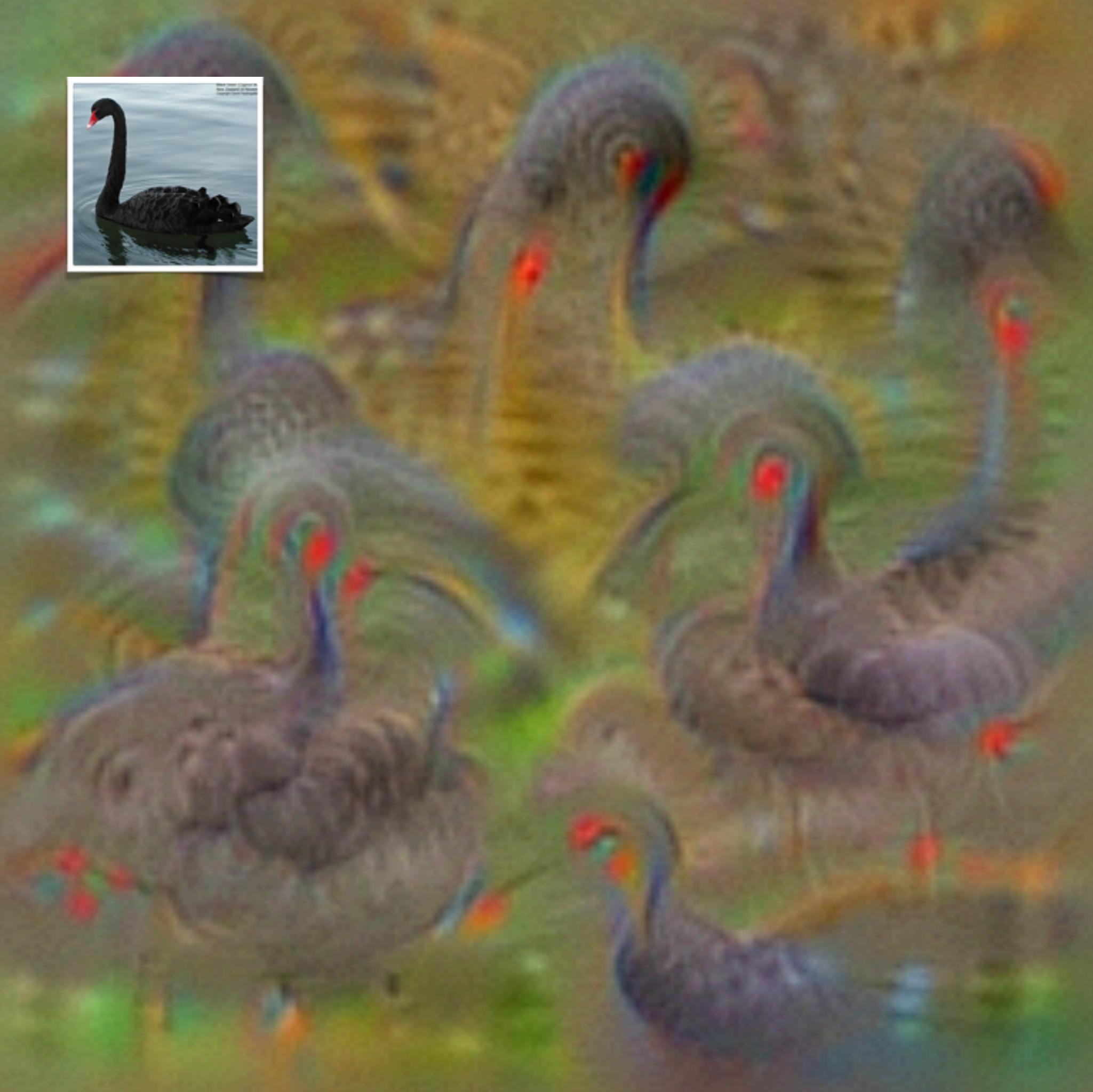
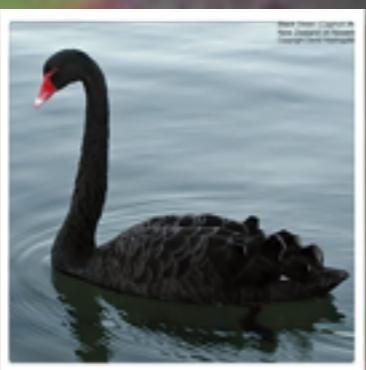


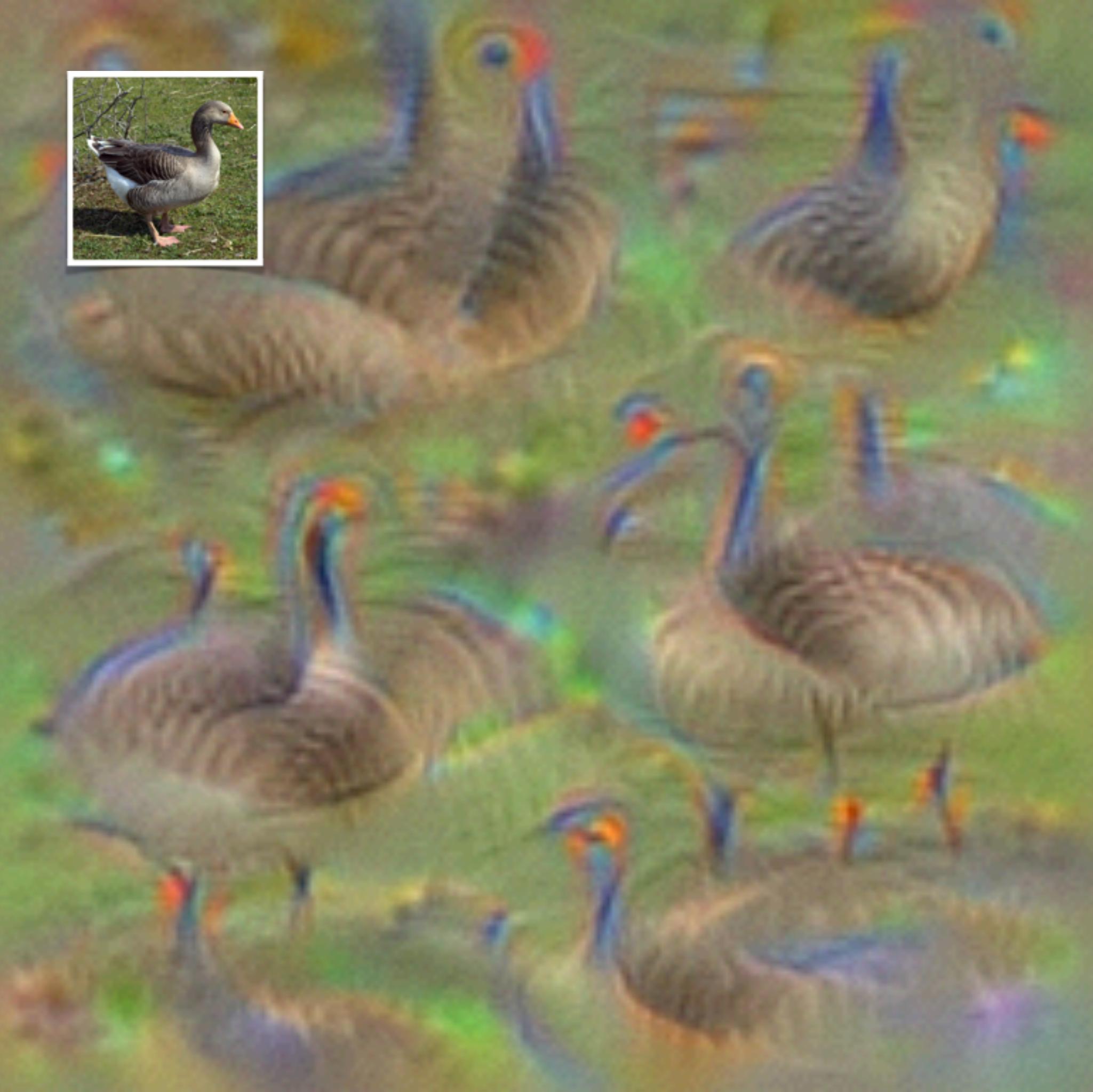
Activation Maximization

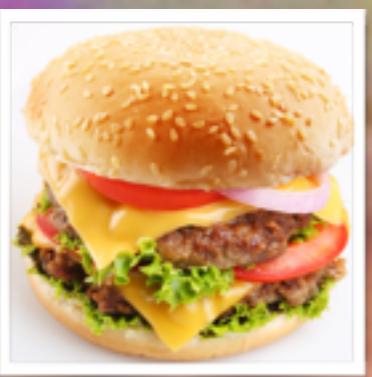
Look for an image that maximally activates a **specific feature component**

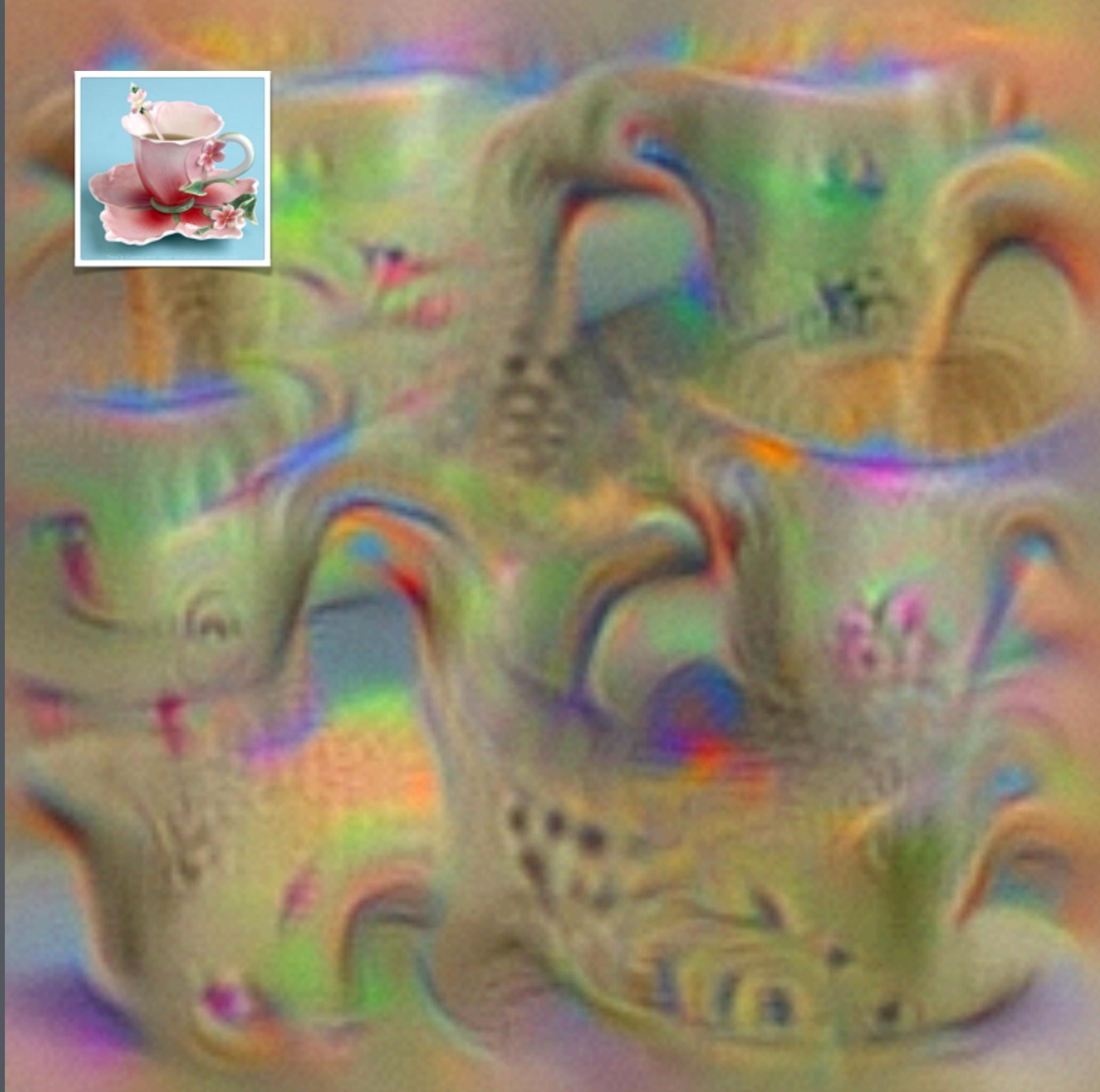
$$\min_{\mathbf{x}} -\langle \mathbf{e}_k, \Phi(\mathbf{x}) \rangle + R_{TV}(\mathbf{x}) + R_\alpha(\mathbf{x})$$

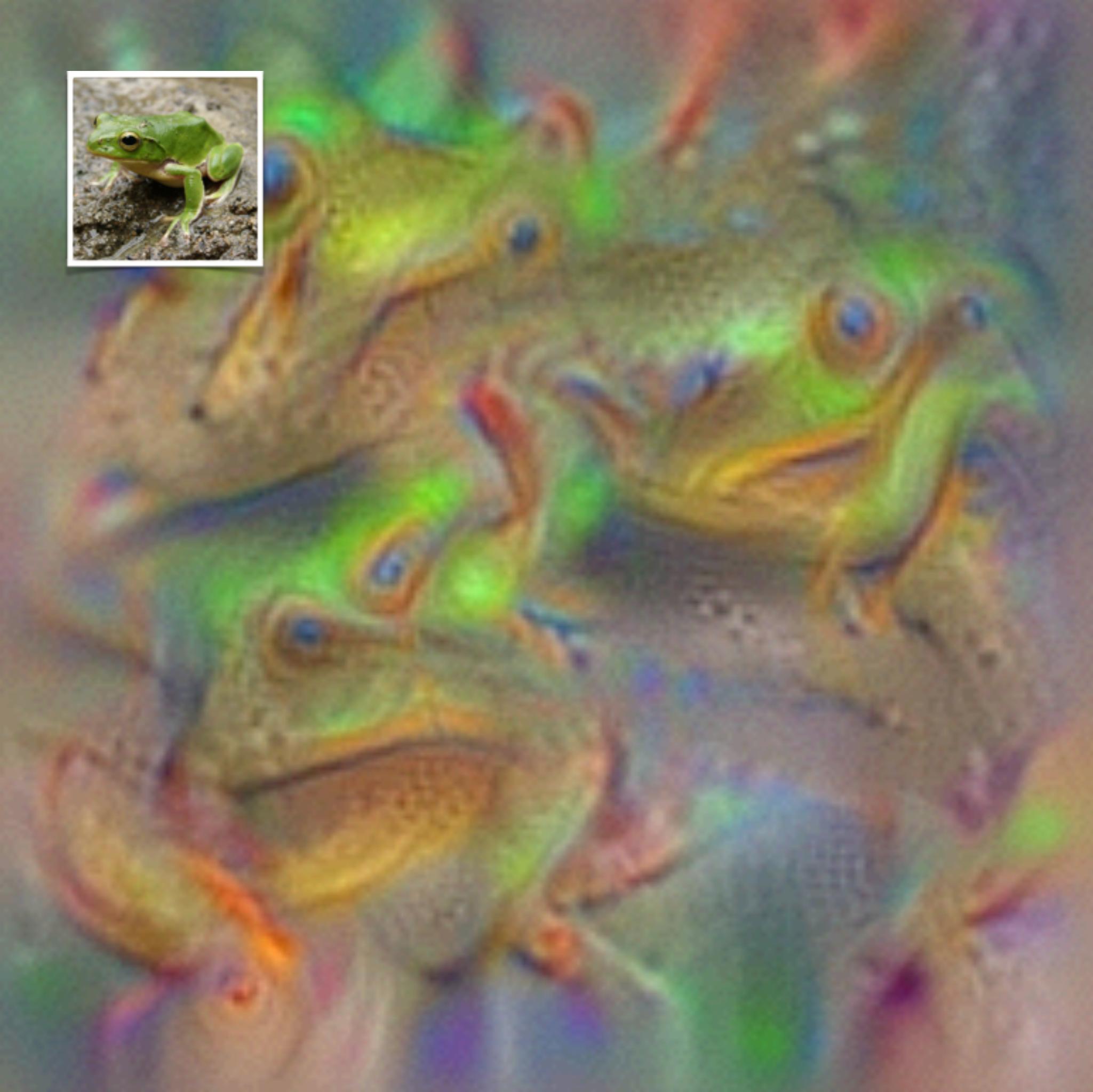


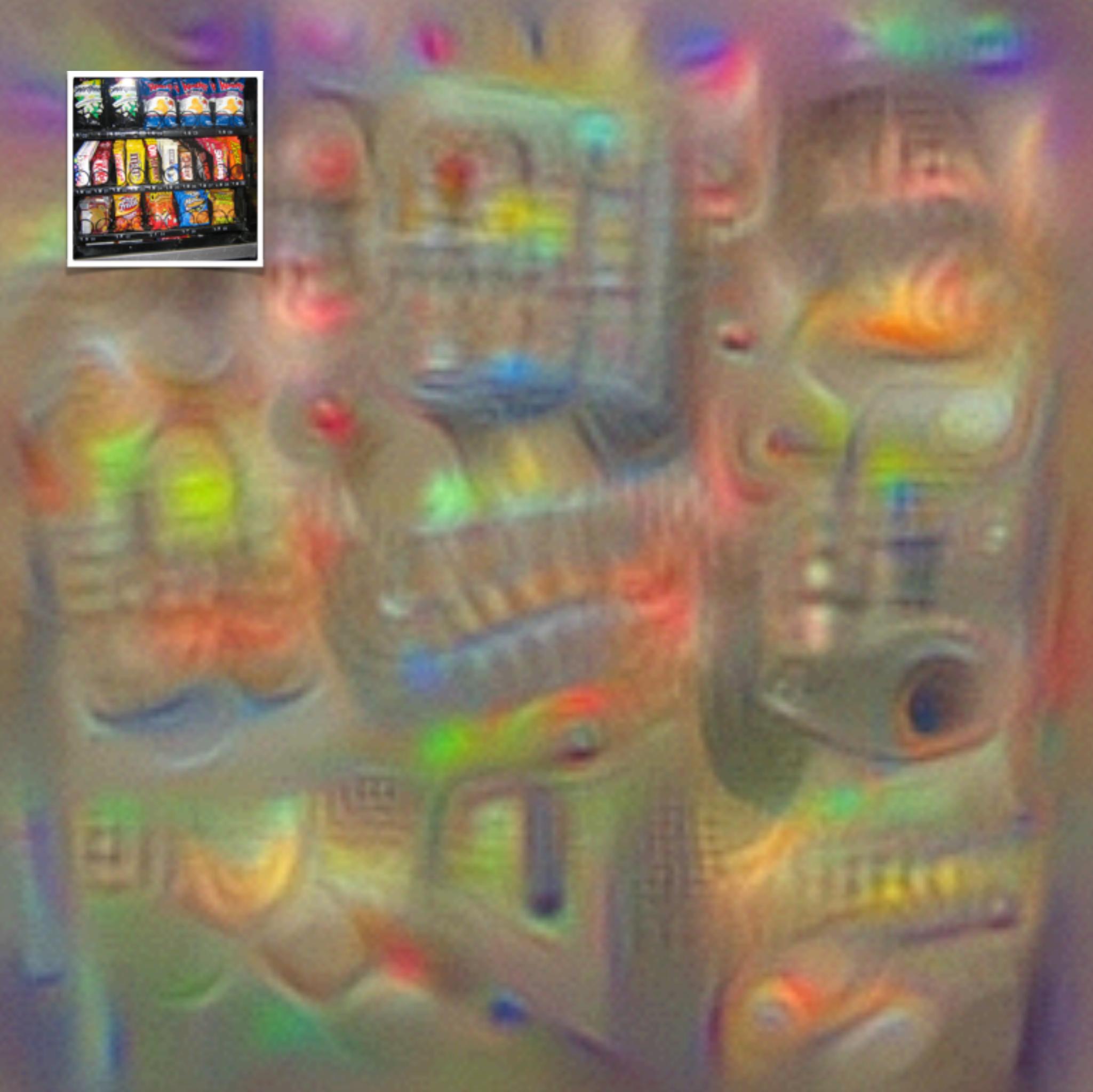








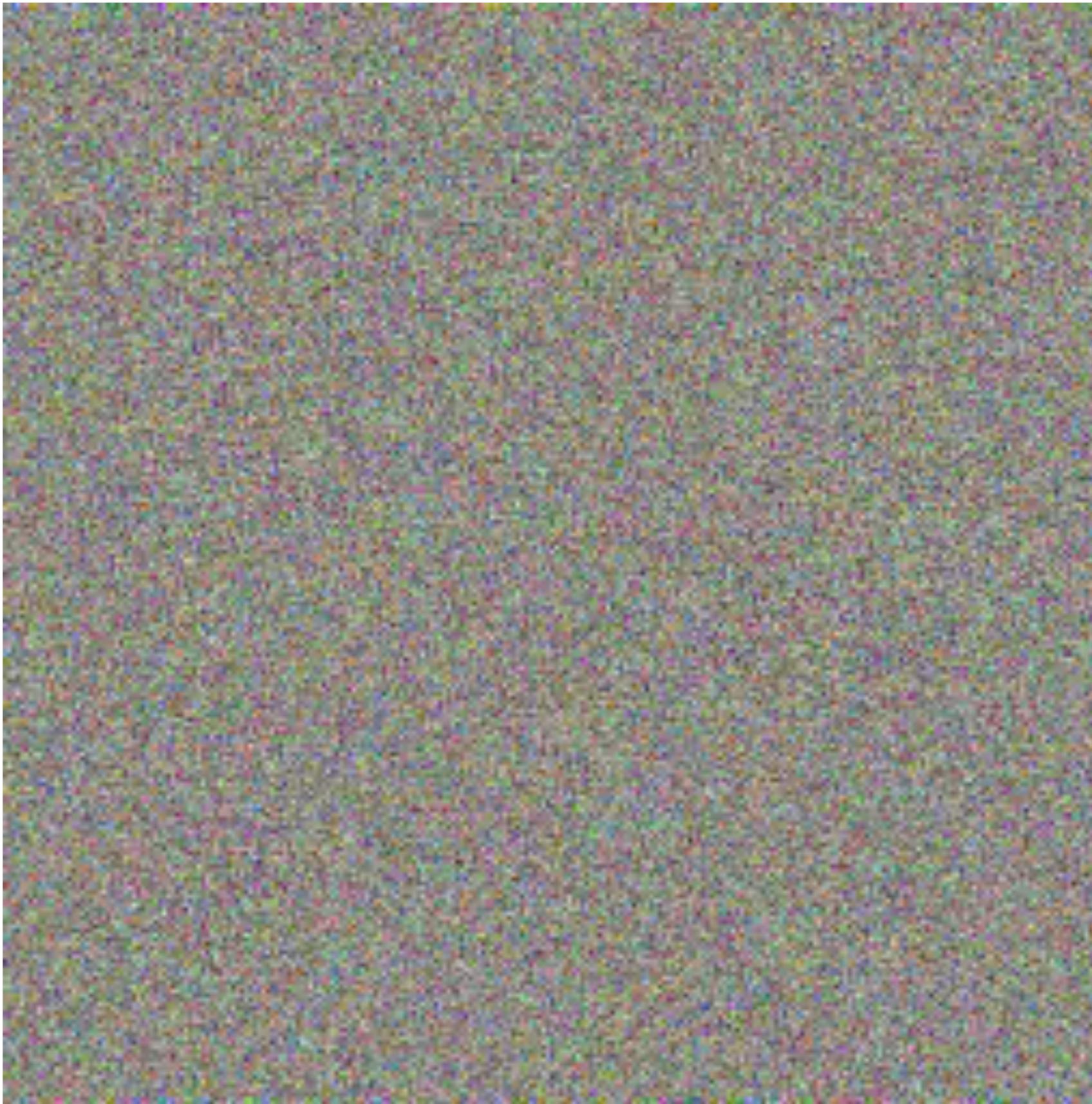


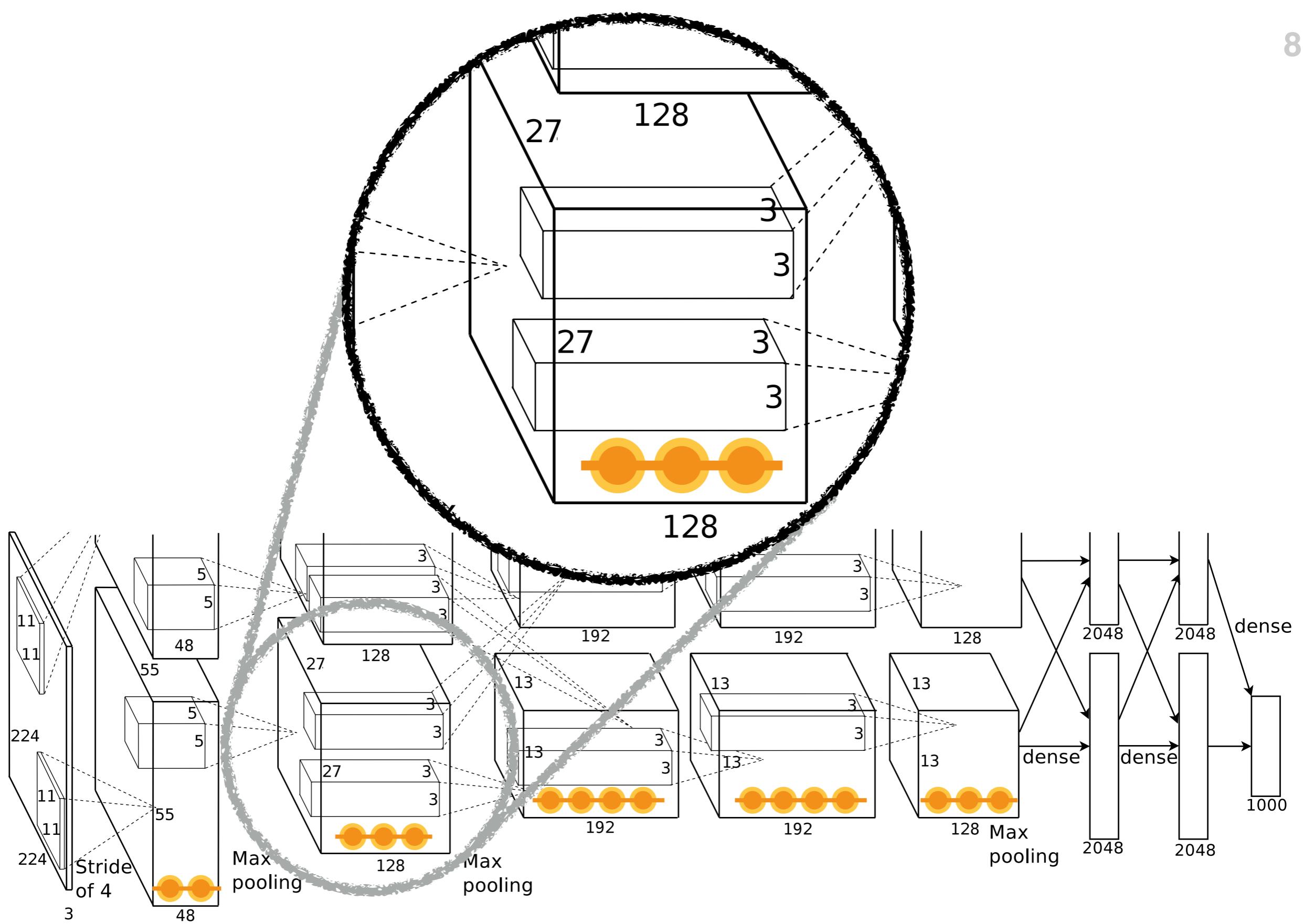


Remember: the starting point is white noise

85

Not an image!

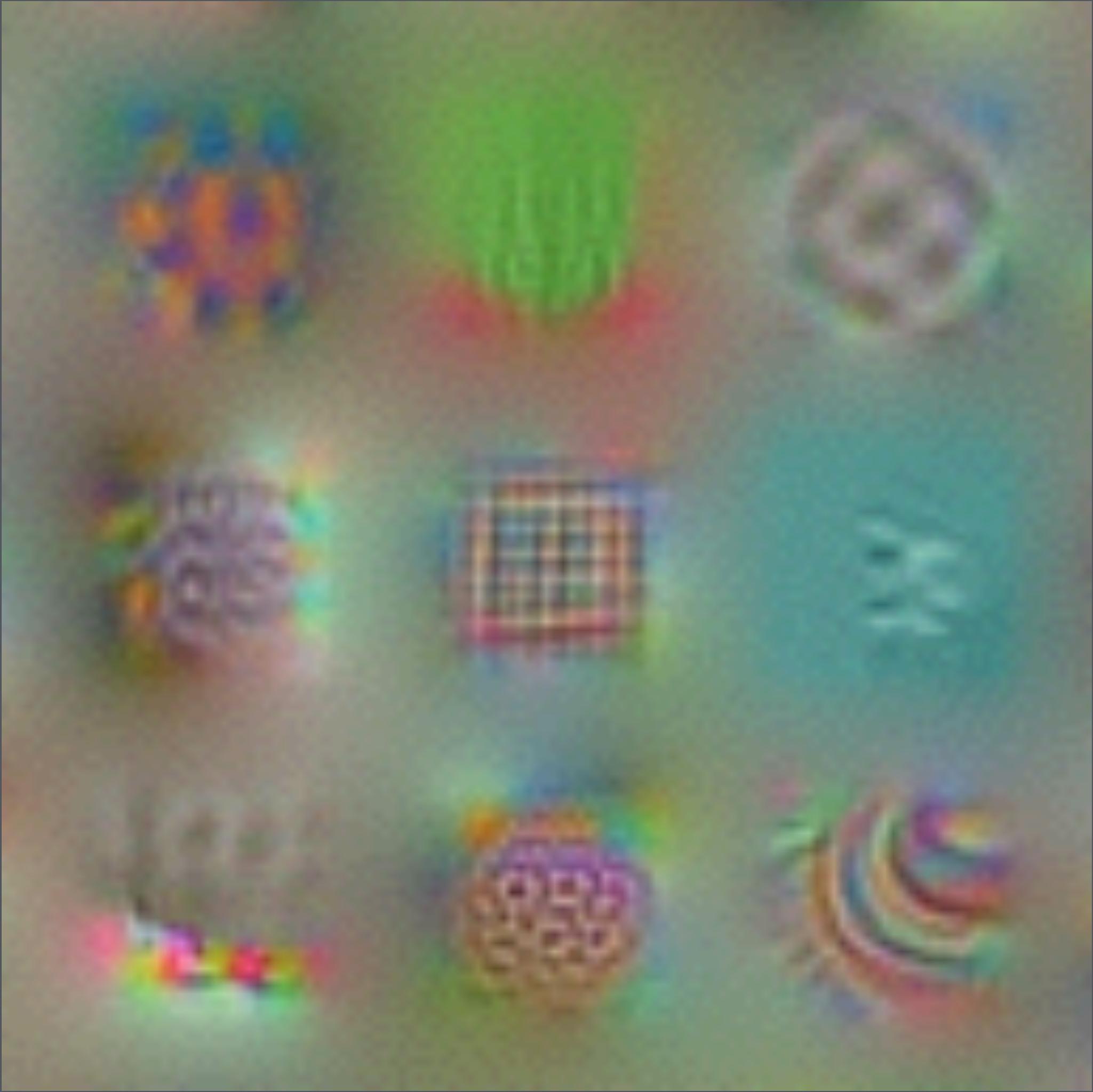




conv1



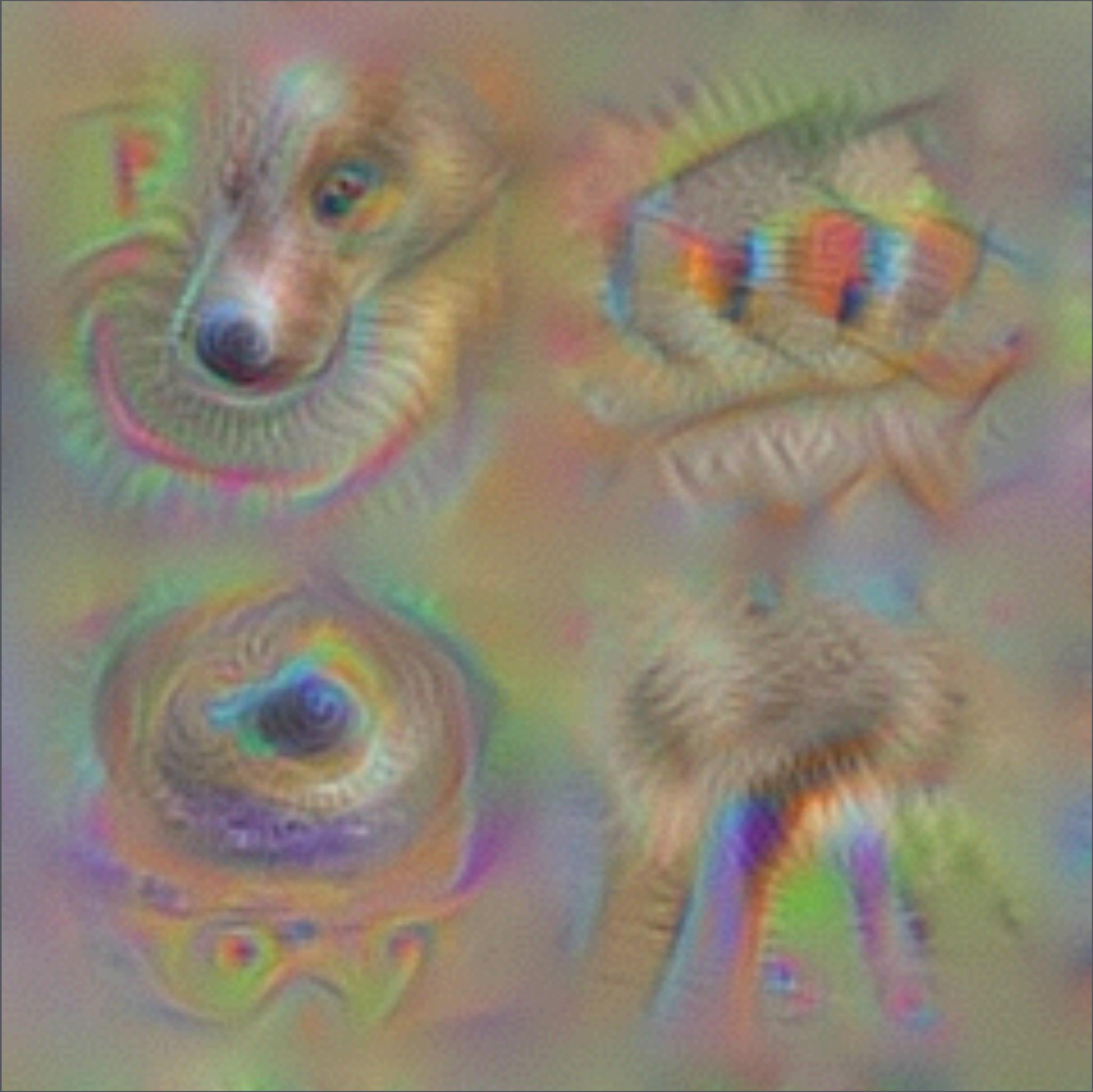
conv2



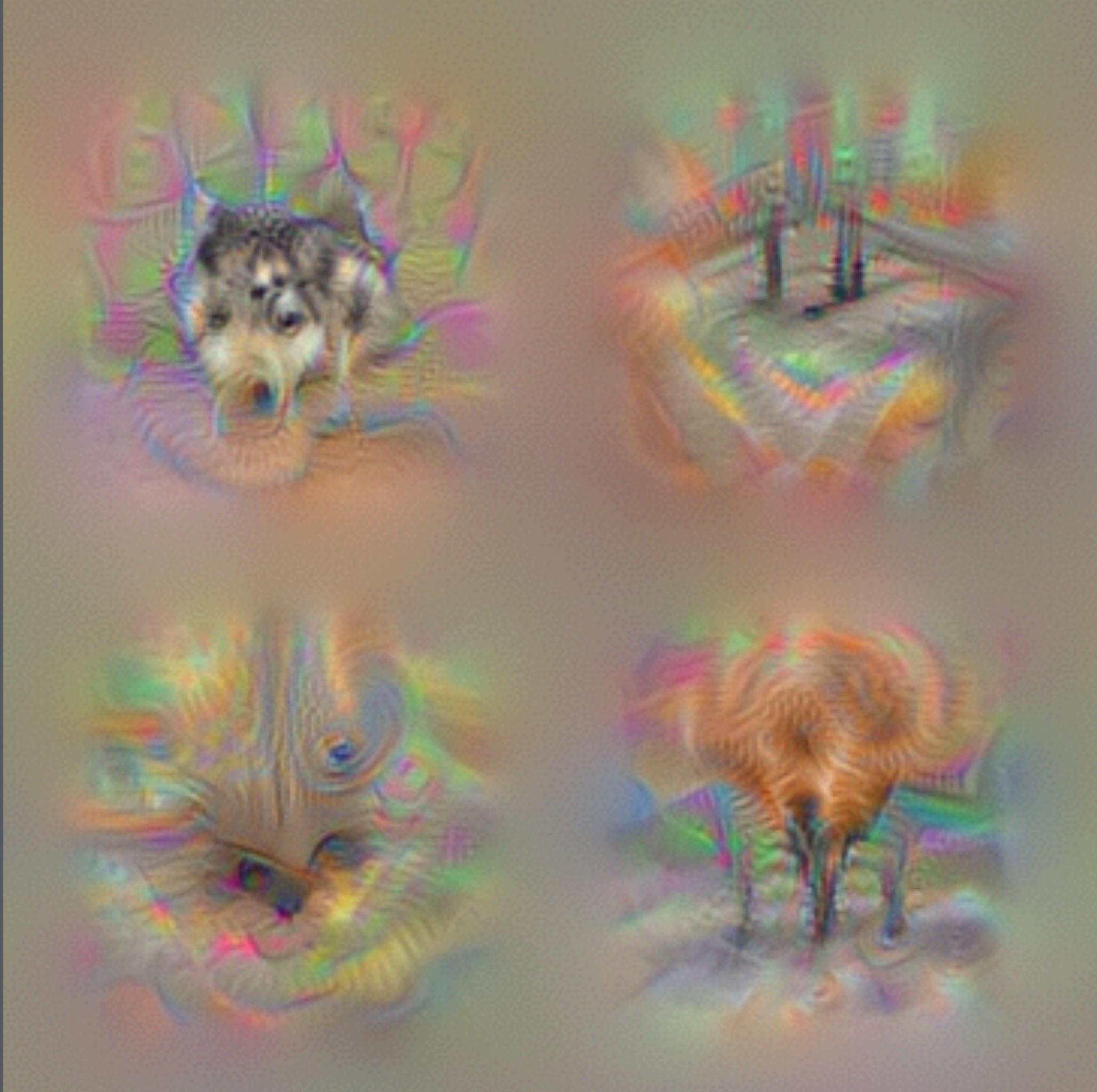
conv3



conv4



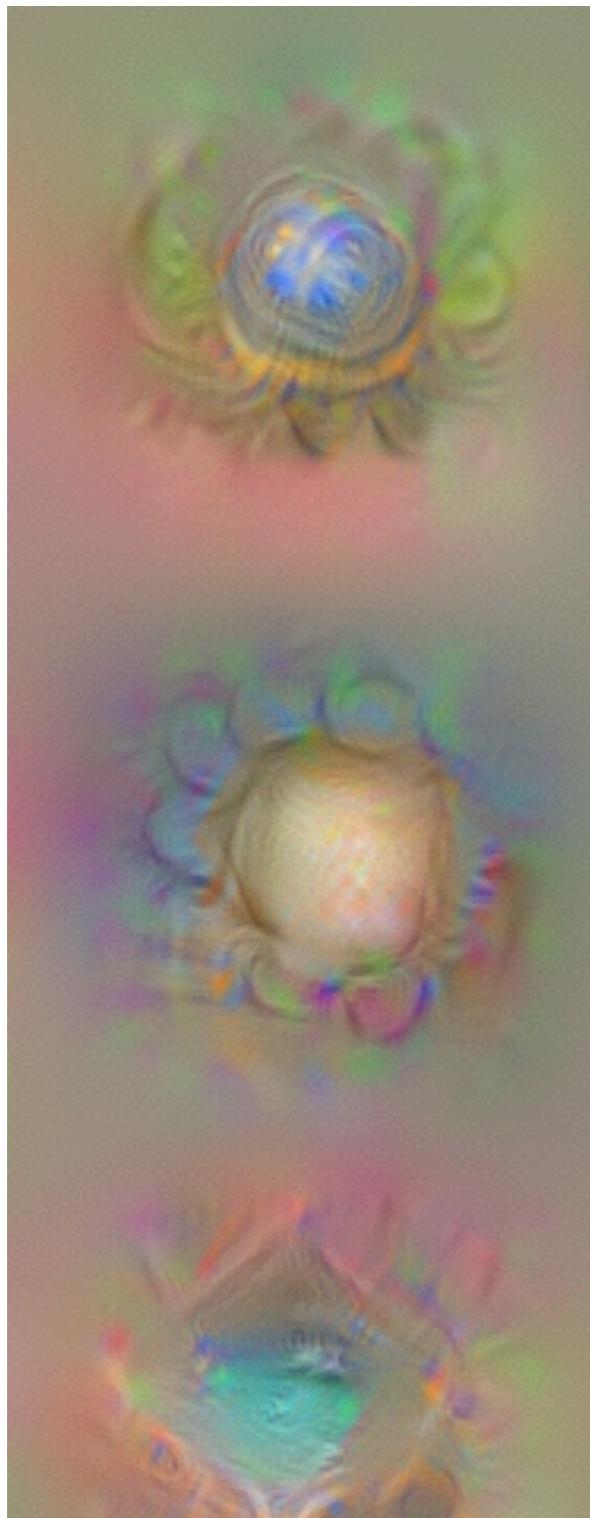
conv5



Network comparison

“conv5” features

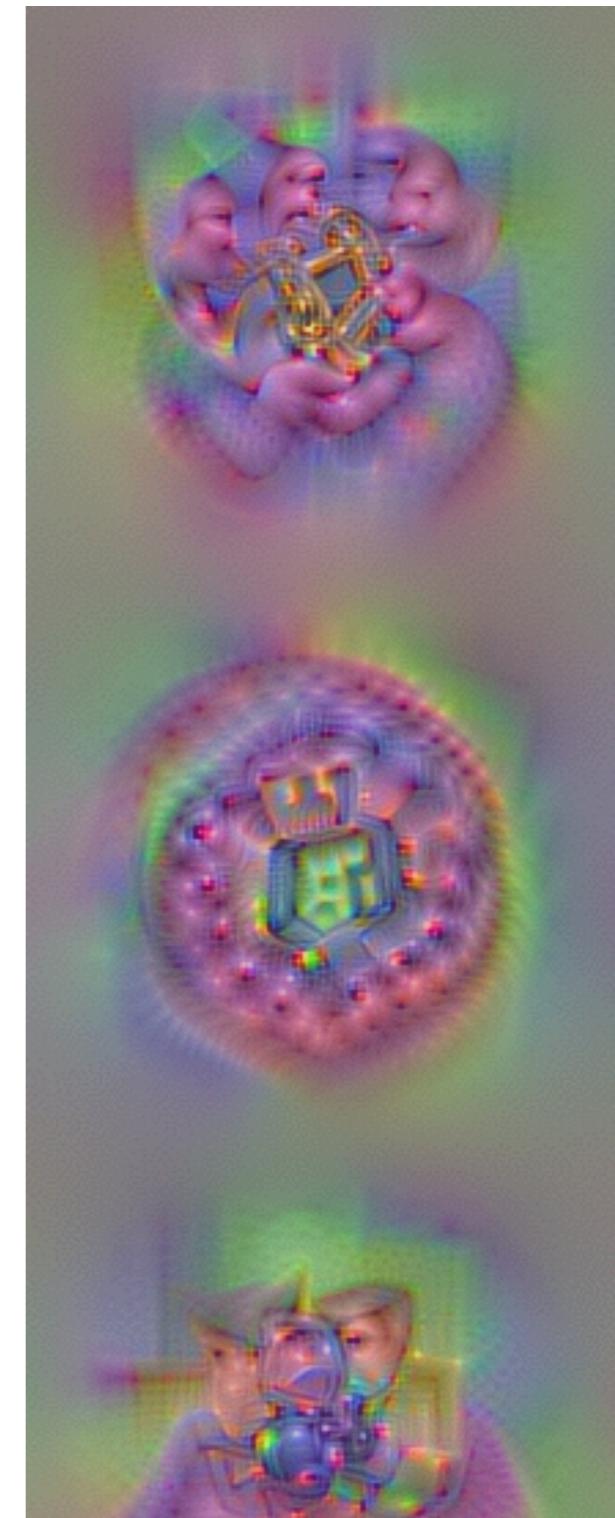
AlexNet



VGG-M



VGG-VD



Caricaturization

[Google Inceptionism 2015, Mahendran et al. 2016]

Emphasise patterns that are detected by a certain representation

$$\min_{\mathbf{x}} -\langle \Phi(\mathbf{x}_0), \Phi(\mathbf{x}) \rangle + R_{TV}(\mathbf{x}) + R_\alpha(\mathbf{x})$$

Key differences:

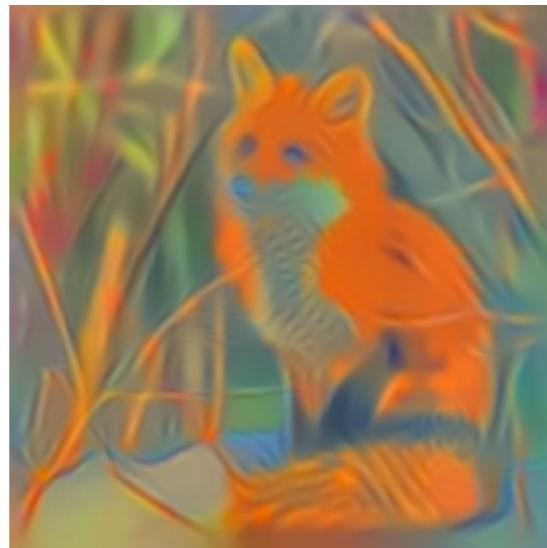
- ▶ the starting point **is** the image \mathbf{x}_0
- ▶ particular configurations of features are emphasized, not individual features

Caricaturization (VGG-M)

input



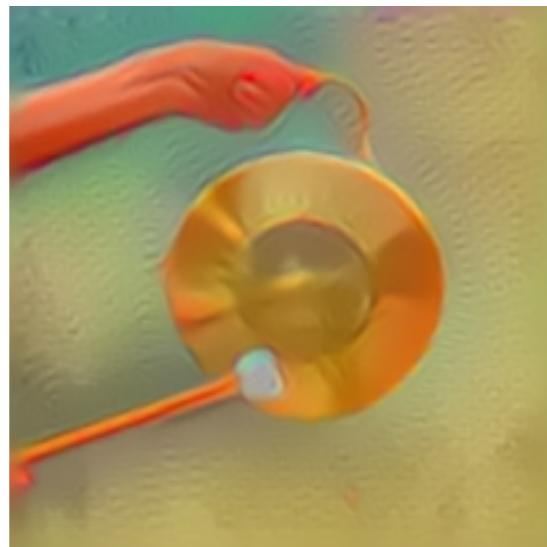
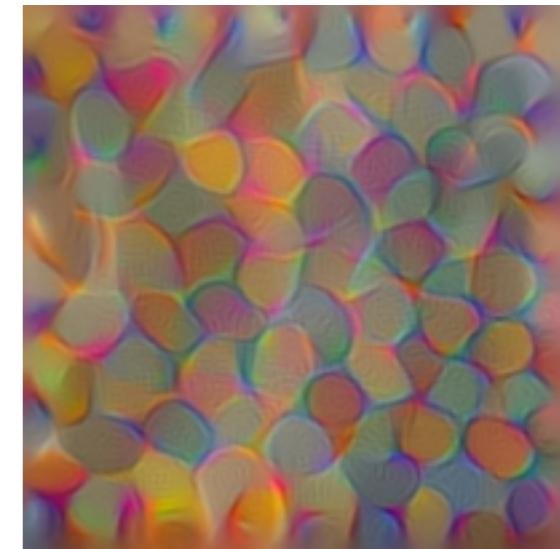
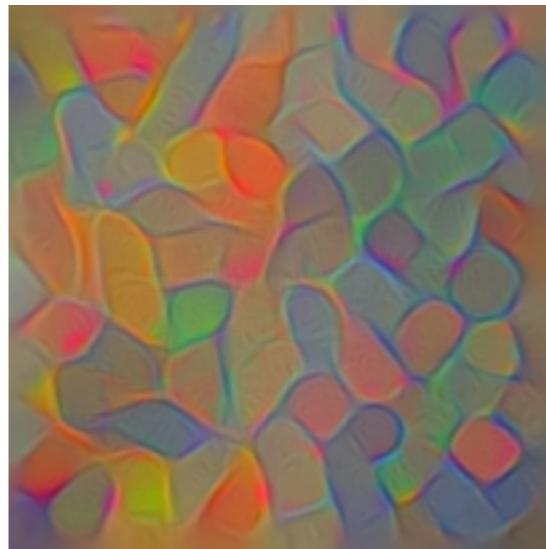
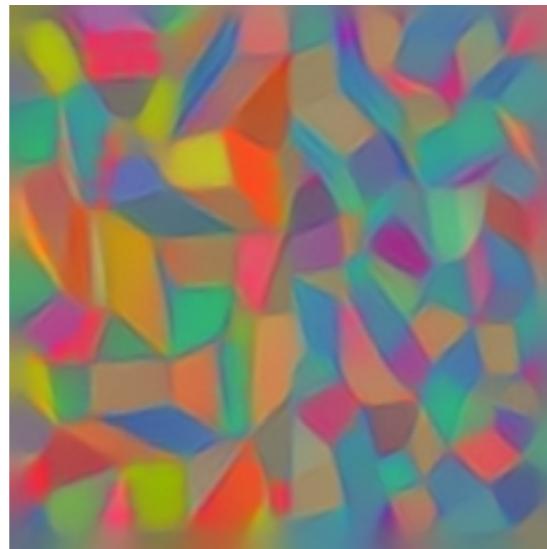
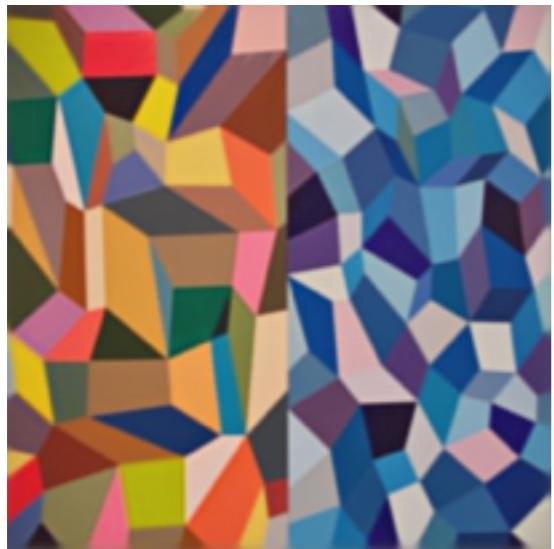
conv2



conv3



conv4



Caricaturization (VGG-M)

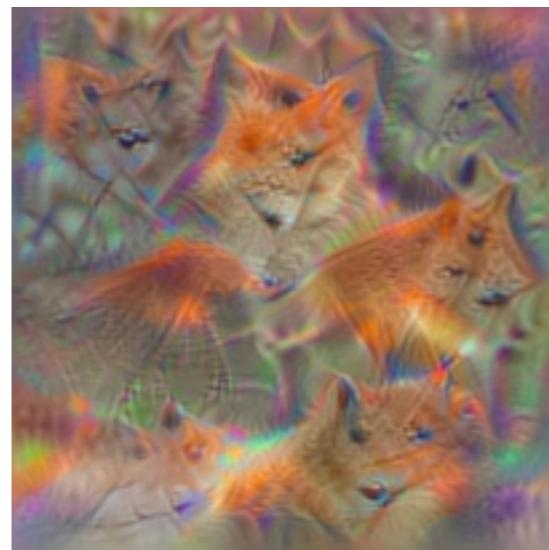
conv5



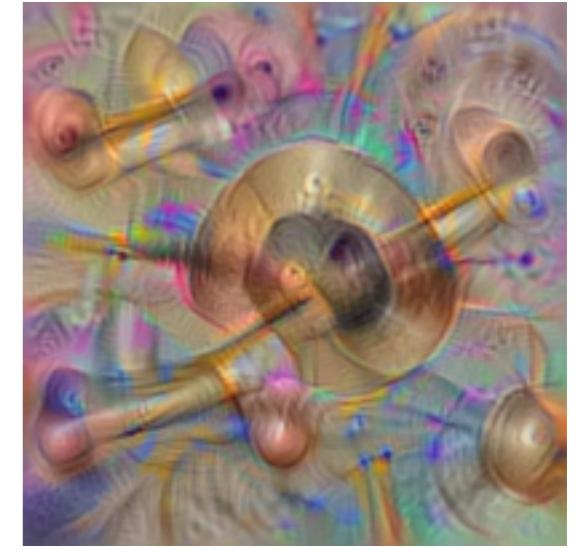
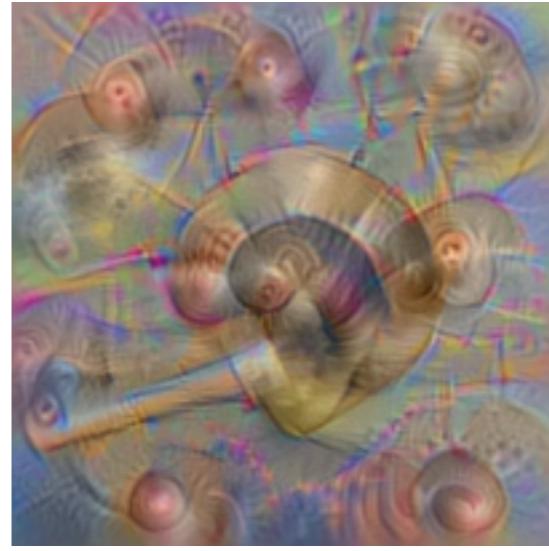
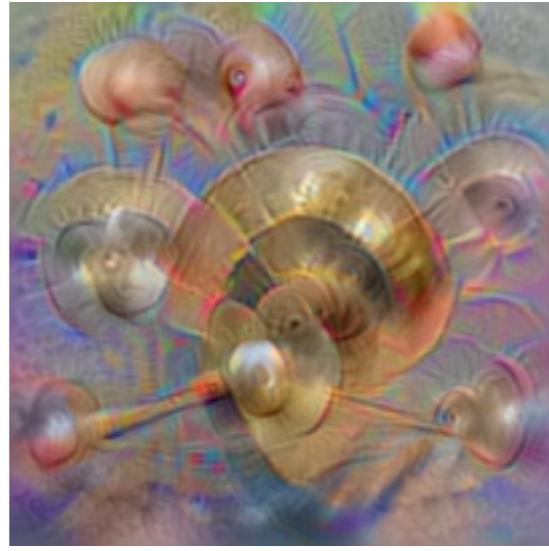
fc6



fc7



fc8





Intro

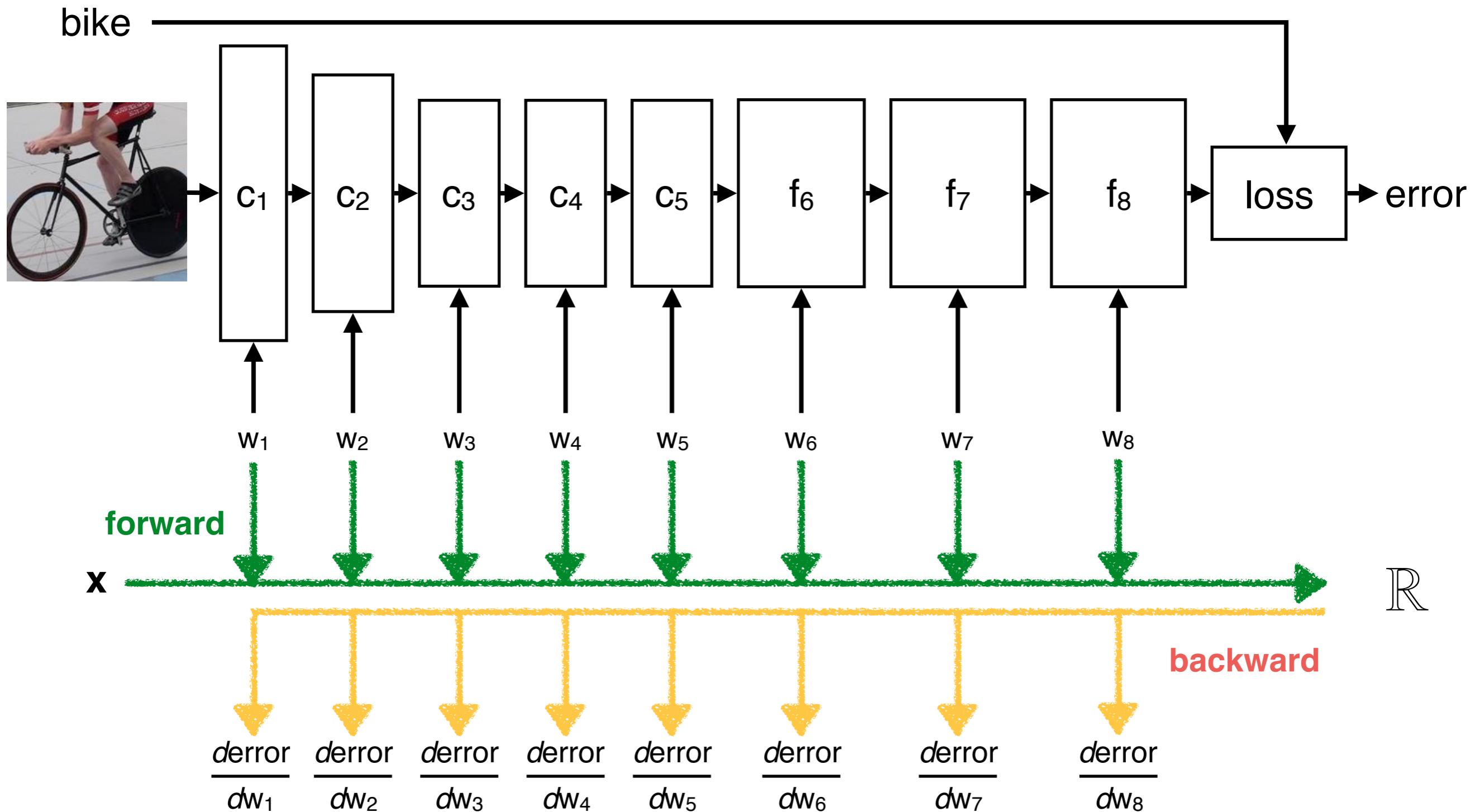
Visualizing representations

Backpropagation networks and “deconvolution”

Representations: equivalence & transformations

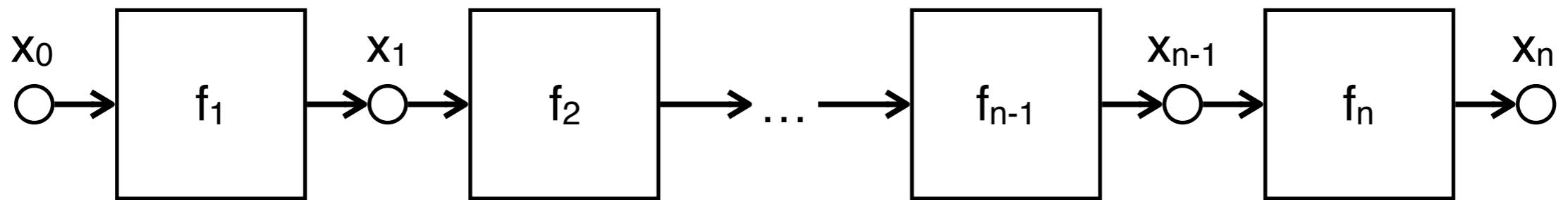
Backpropagation

Compute derivatives using the chain rule

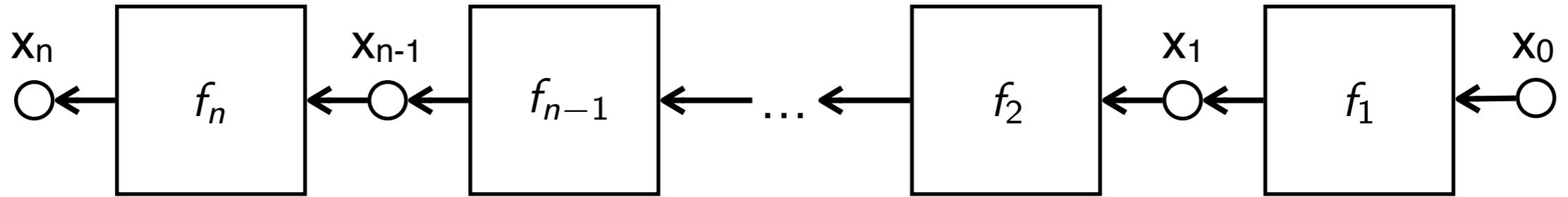


Chain rule: scalar version

99



Chain rule: scalar version



A composition of n functions

$$x_n = (f_n \circ \dots \circ f_2 \circ f_1)(x_0)$$

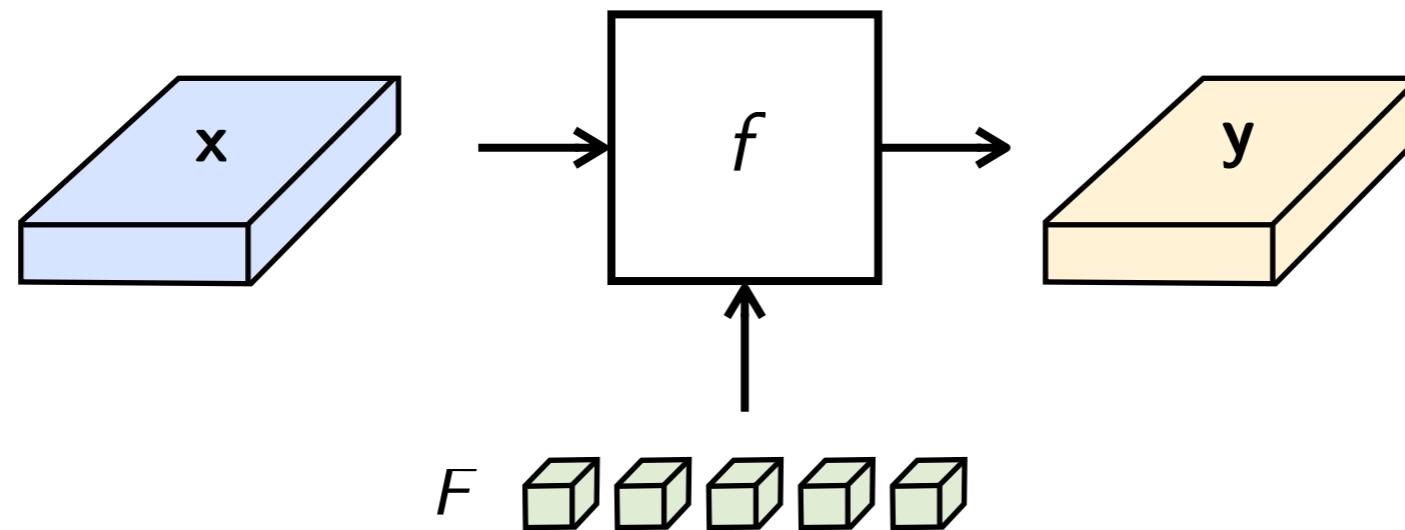
$$\frac{dx_n}{dx_0} = \frac{df_n}{dx_{n-1}} \times \frac{df_{n-1}}{dx_{n-2}} \times \dots \times \frac{df_2}{dx_1} \times \frac{df_1}{dx_0}$$

Derivative \leftarrow chain rule

Tensor-valued functions

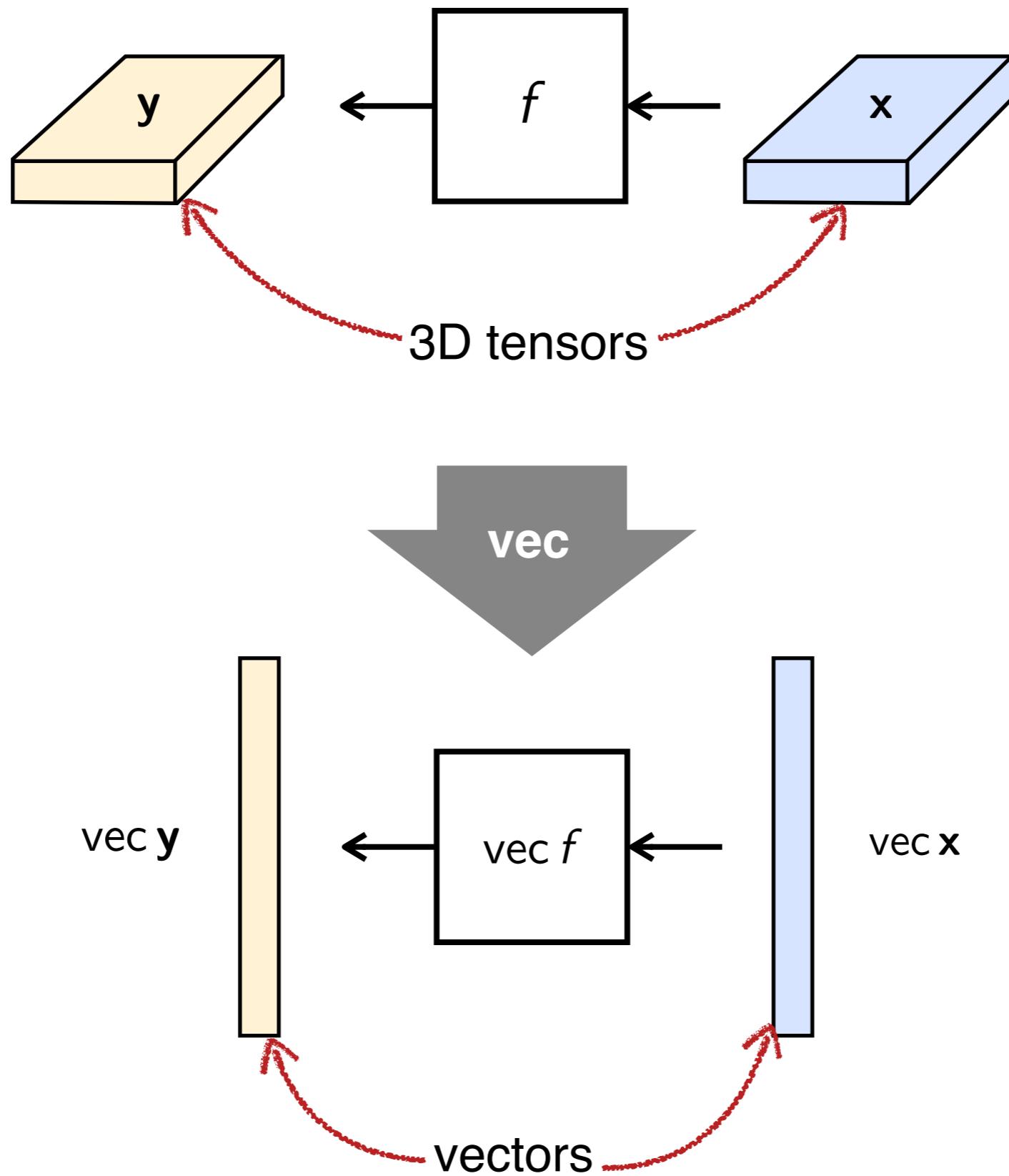
E.g. linear convolution = bank of 3D filters

$$\mathbf{y} = \mathcal{F} * \mathbf{x} + b$$



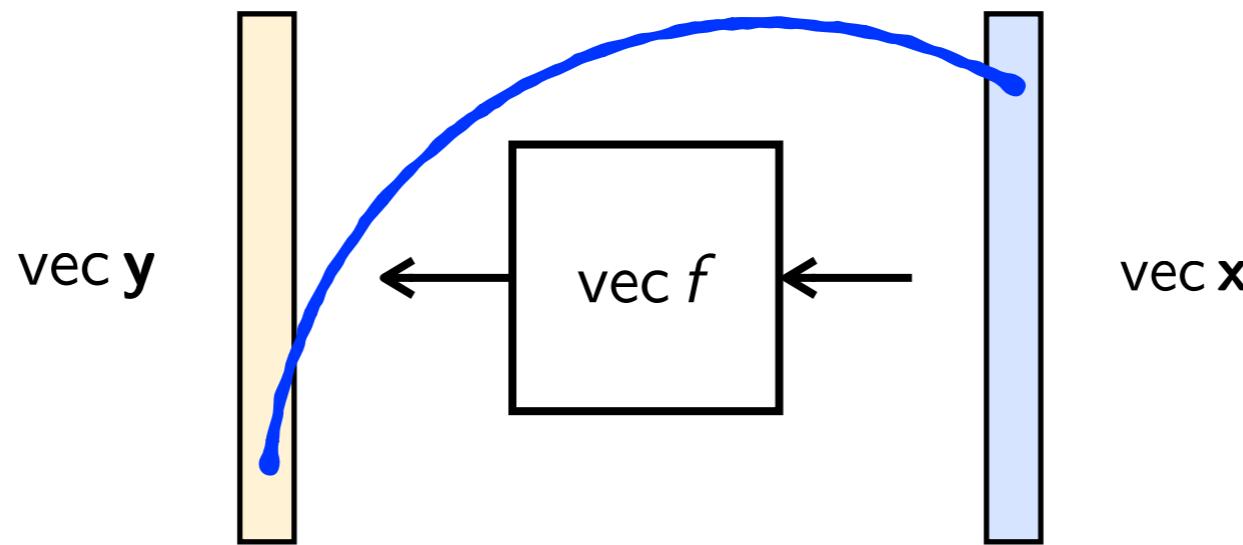
| | height | width | channels | instances |
|-----------------------|---------------|---------------|----------|-----------|
| input \mathbf{x} | H | W | C | 1 or N |
| filters \mathcal{F} | H_f | W_f | C | K |
| output \mathbf{y} | $H - H_f + 1$ | $W - W_f + 1$ | K | 1 or N |

Vector representation

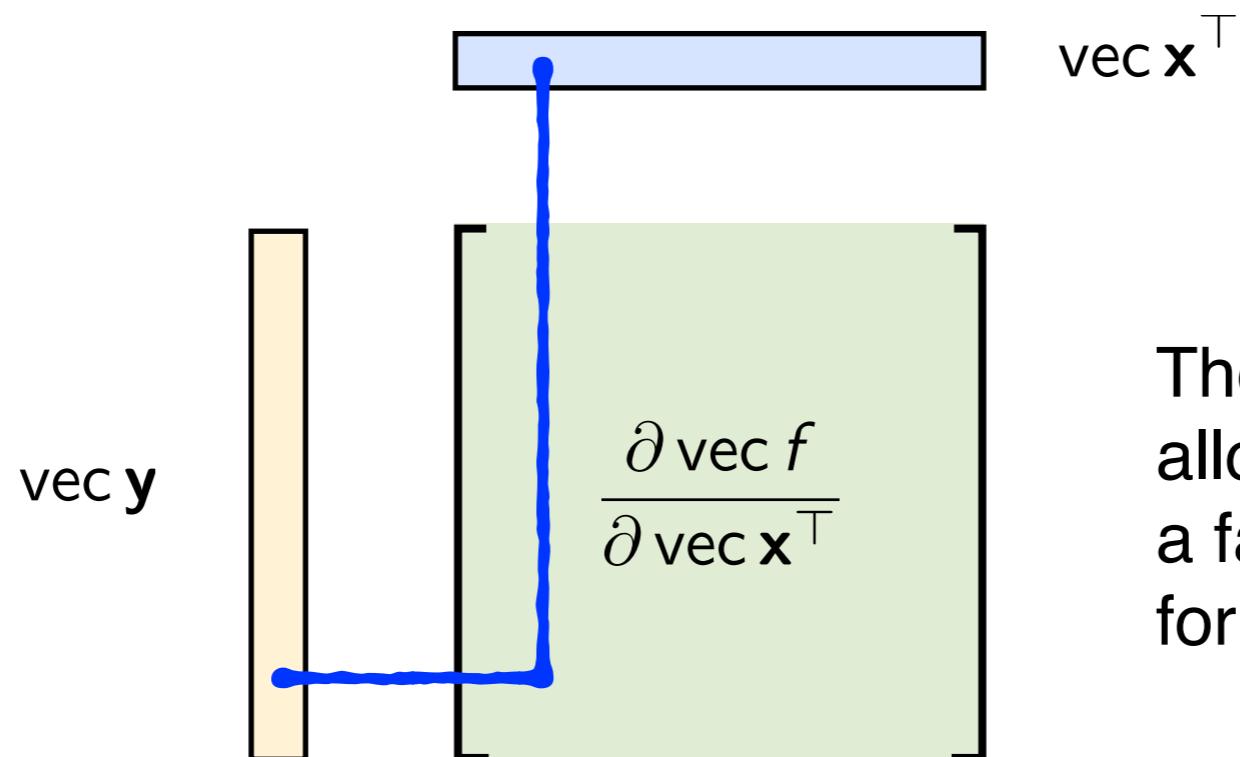


Derivative of tensor-valued functions

103



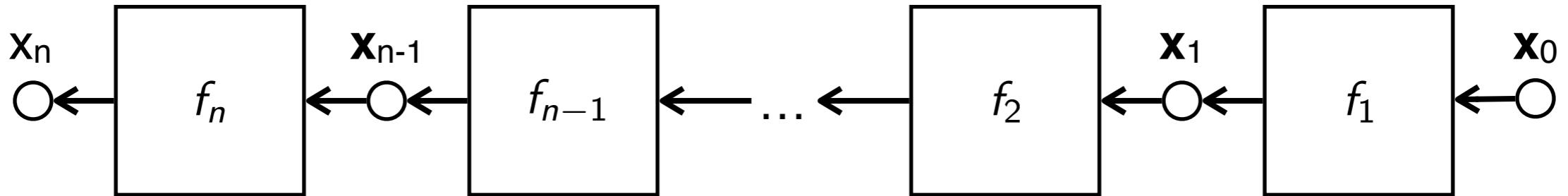
Derivative (Jacobian): every output element w.r.t. every input element!



The **vec** operator allows us to use a familiar matrix notation for the derivatives

Chain rule: tensor version

Using $\text{vec}()$ and matrix notation



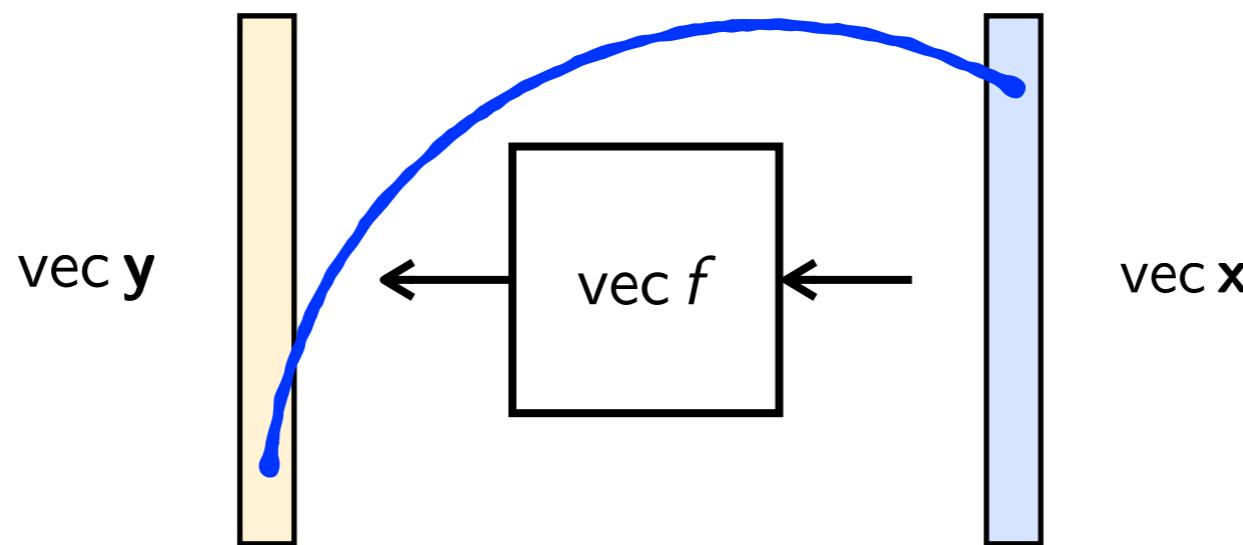
$$\frac{d \text{vec } f_n}{d \text{vec } \mathbf{x}_{n-1}^\top} \times \frac{d \text{vec } f_{n-1}}{d \text{vec } \mathbf{x}_{n-2}^\top} \times \dots \times \frac{d \text{vec } f_2}{d \text{vec } \mathbf{x}_1^\top} \times \frac{d \text{vec } f_1}{d \text{vec } \mathbf{x}_0^\top}$$

$$\left[\begin{array}{c} \text{green box} \end{array} \right] \times \left[\begin{array}{c} \text{green box} \end{array} \right] \times \dots \times \left[\begin{array}{c} \text{green box} \end{array} \right] \times \left[\begin{array}{c} \text{green box} \end{array} \right]$$

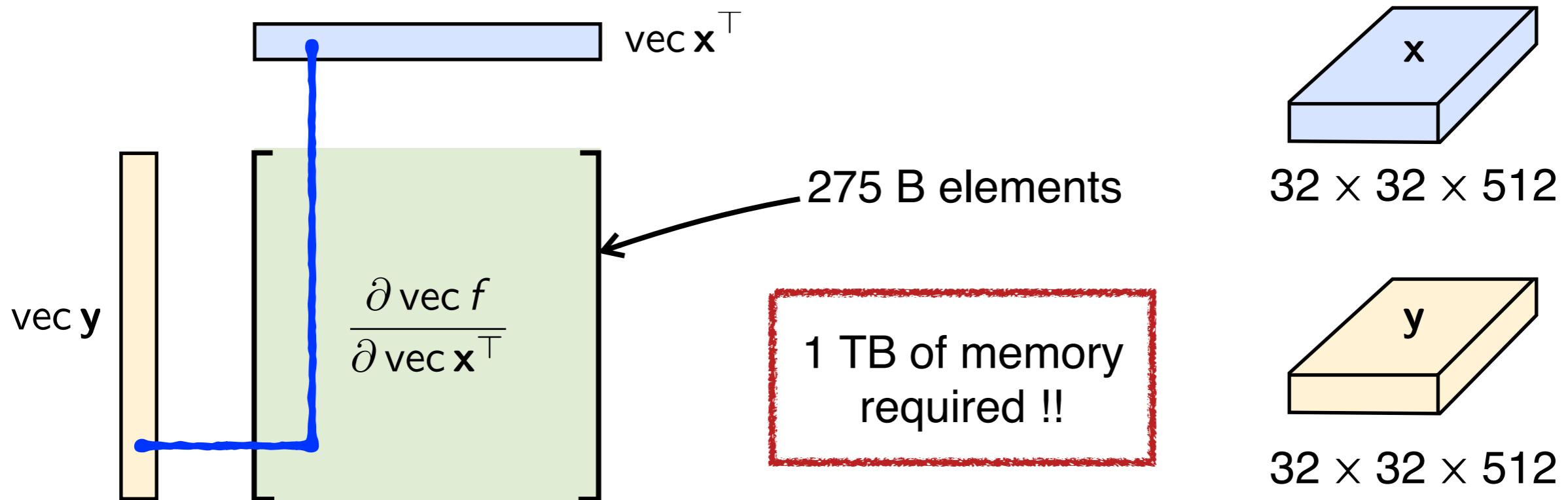
$$\underbrace{\quad\quad\quad}_{d \text{vec } f} \frac{d \text{vec } f}{d \text{vec } \mathbf{x}_0^\top}$$

The (unbearable) size of tensor derivatives

105



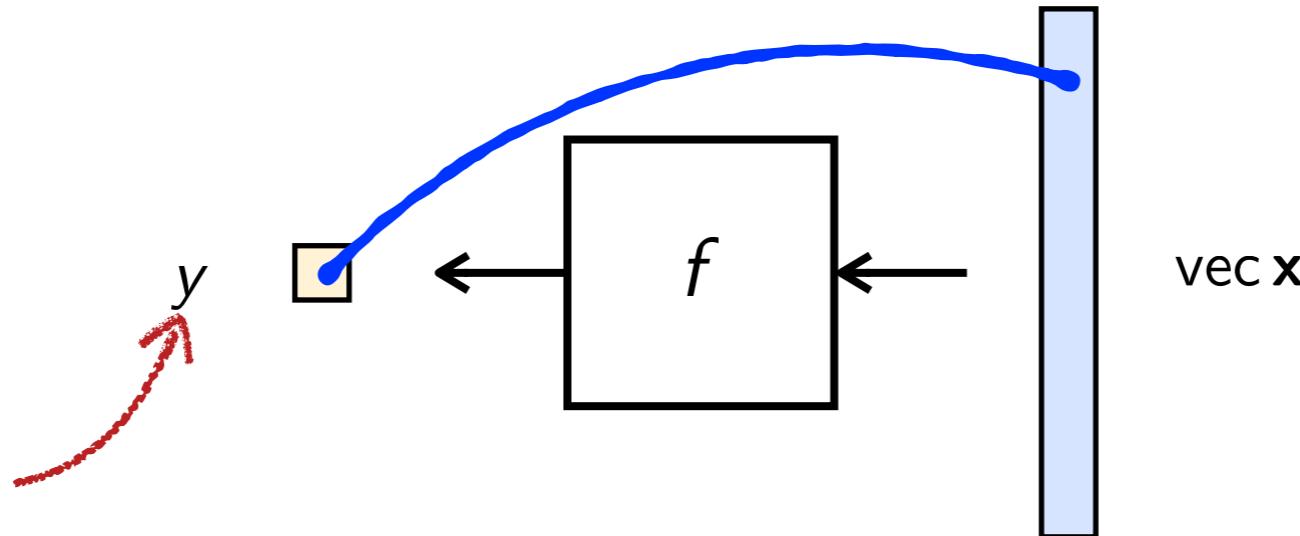
The size of these Jacobian matrices is **huge**. Example:



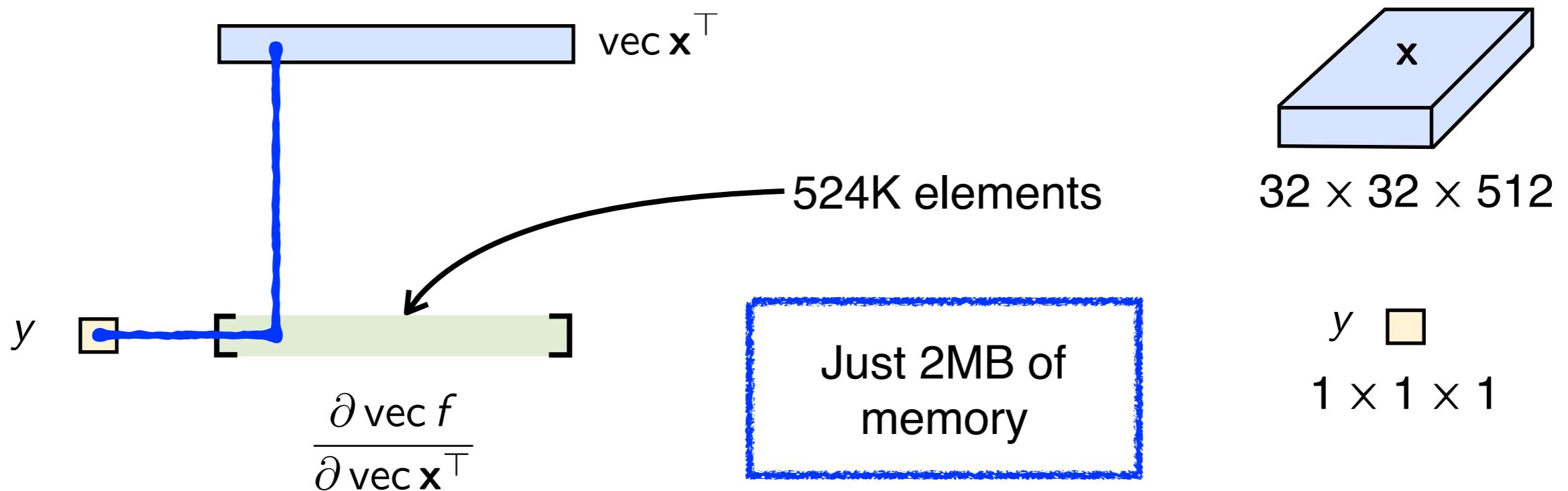
Unless the output is a scalar

Scalar

This is always the case
if the last layer
is the **loss function**

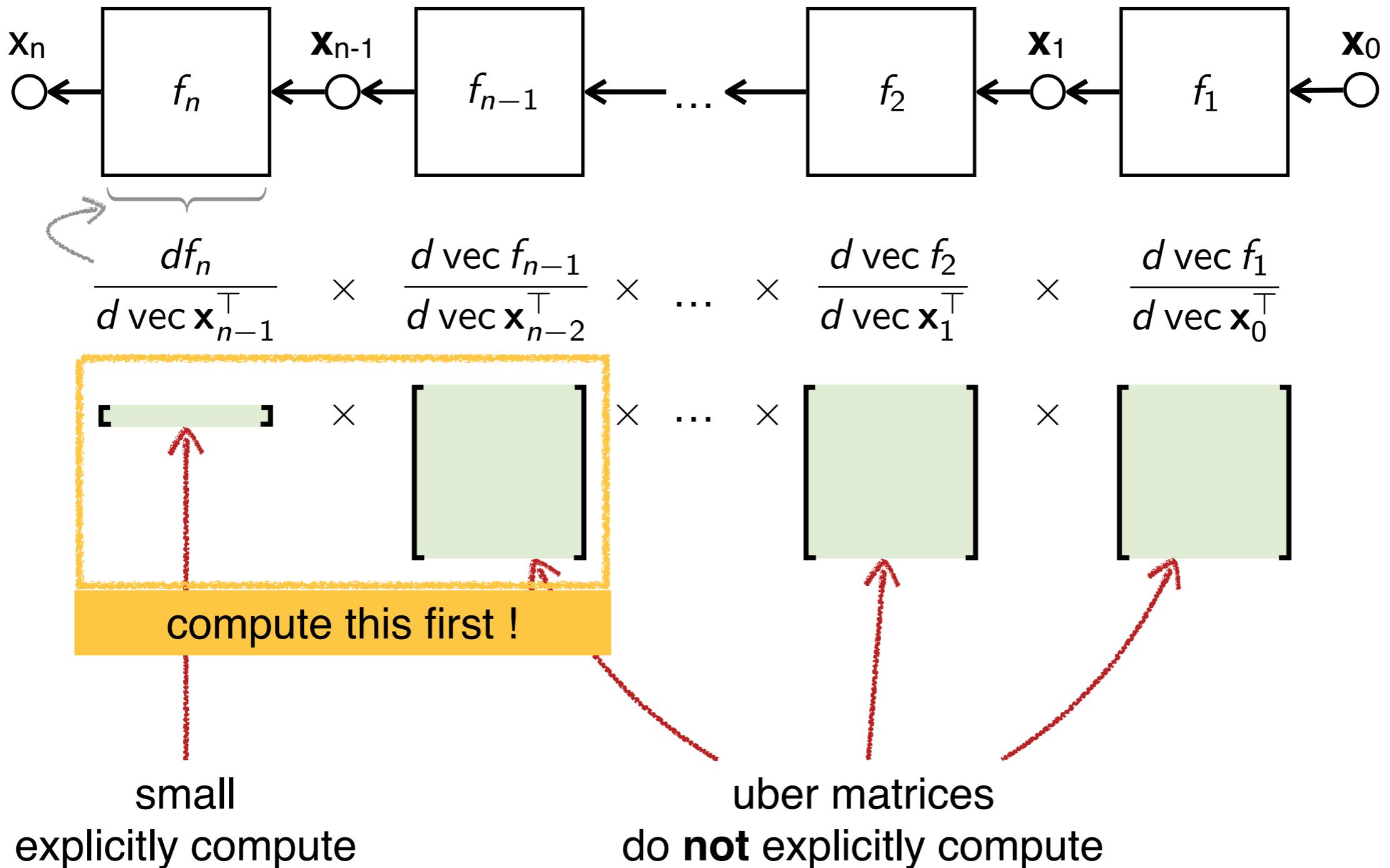


Now the Jacobian has the same size as x . Example:



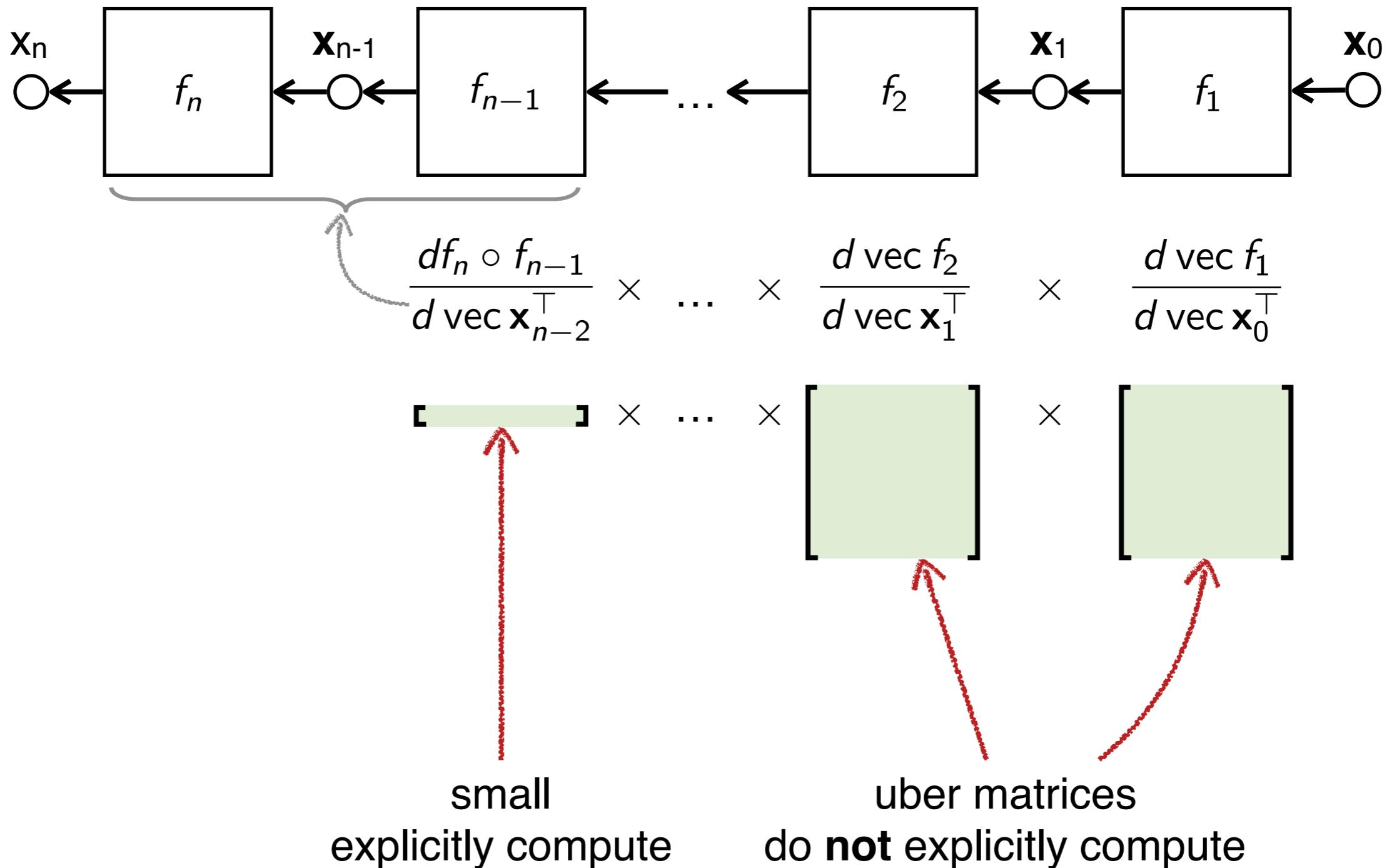
Backpropagation

Assume that x_n is a scalar (e.g. loss)



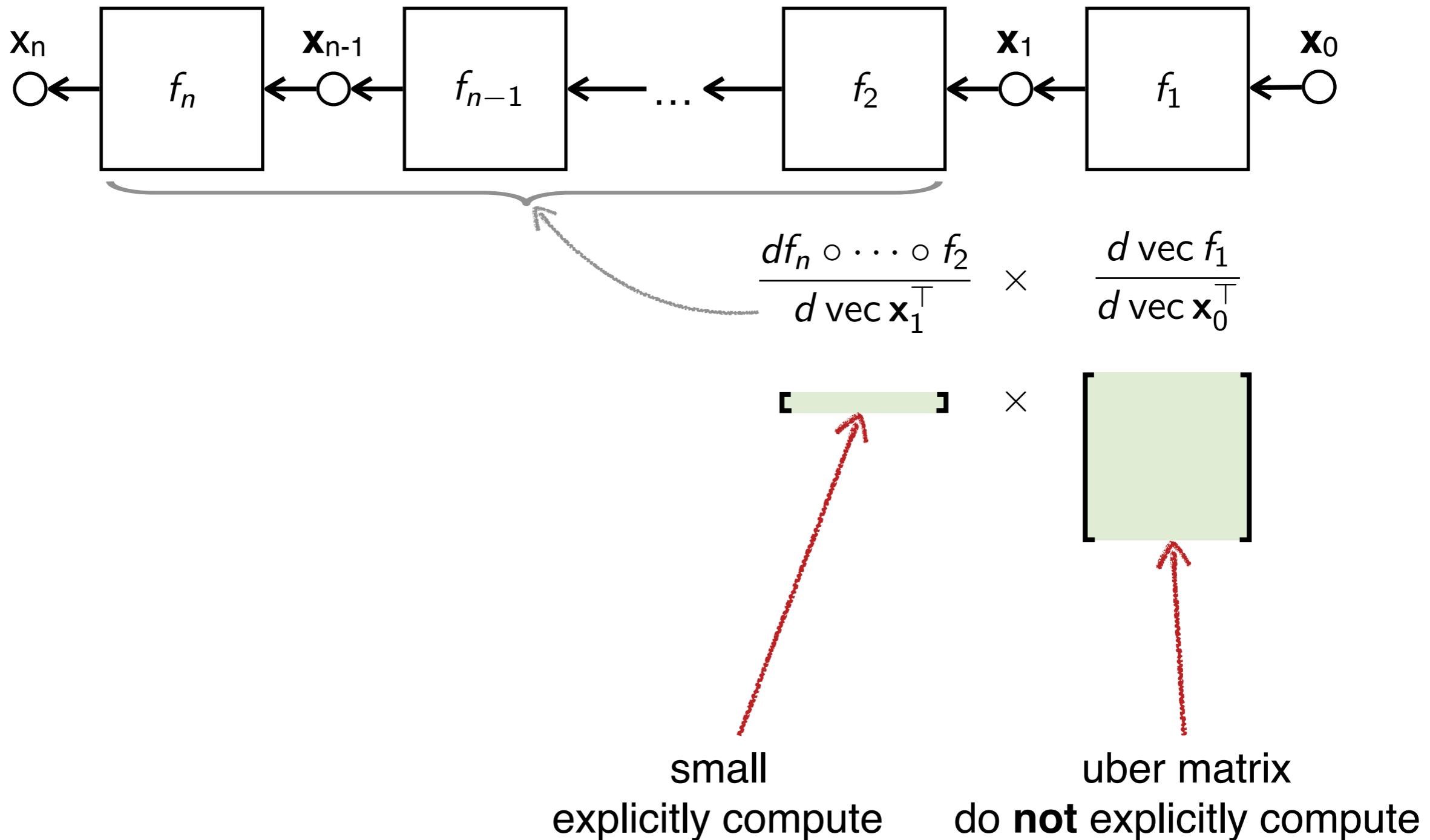
Backpropagation

Assume that x_n is a scalar (e.g. loss)



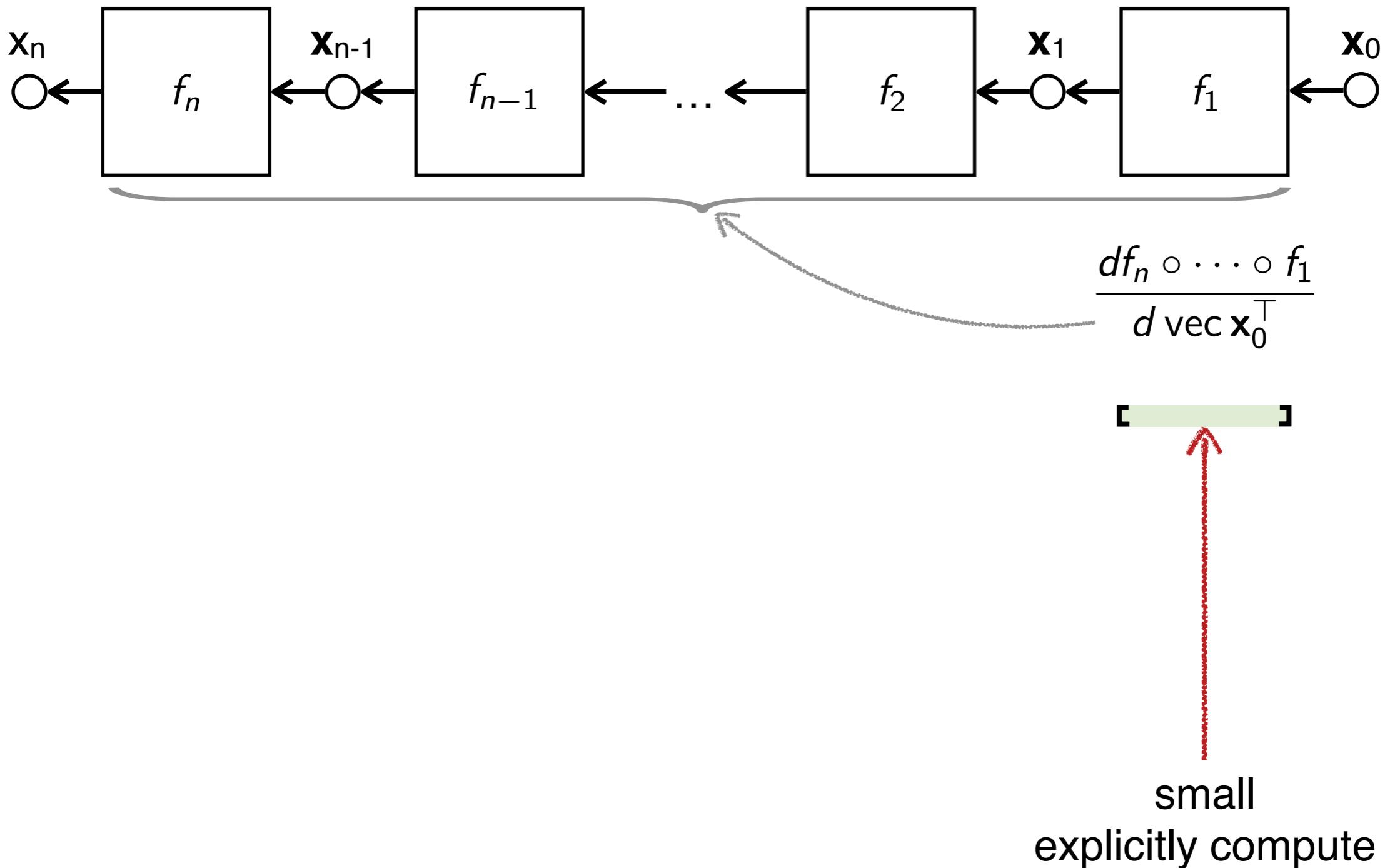
Backpropagation

Assume that x_n is a scalar (e.g. loss)



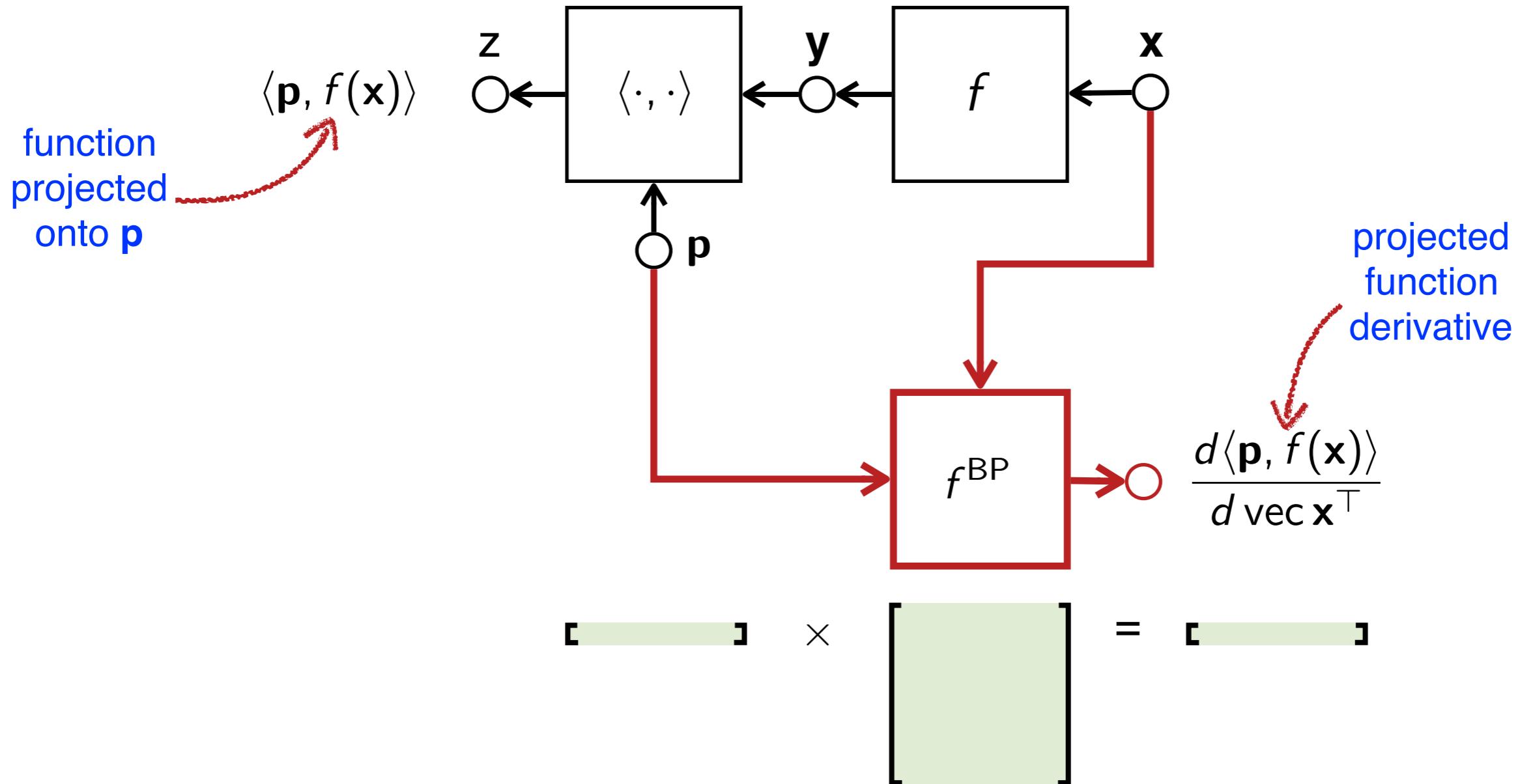
Backpropagation

Assume that x_n is a scalar (e.g. loss)



Projected function derivative

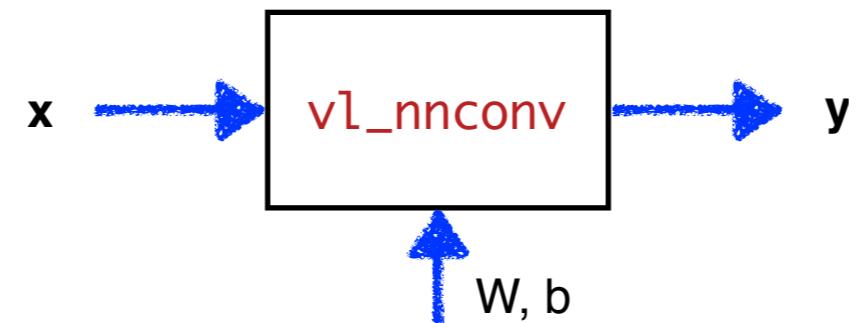
The “BP-reversed” layer



An “equivalent circuit” is obtained by introducing a transposed function f^\top

Anatomy of a building block

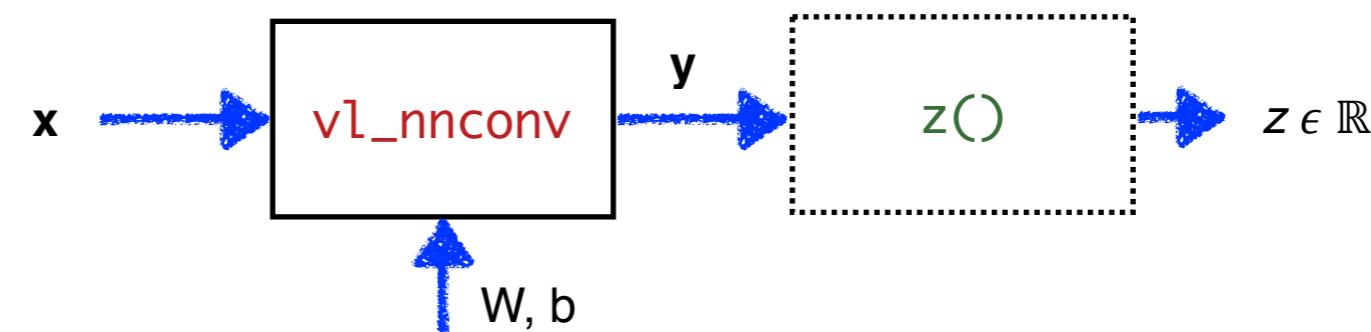
forward (eval)



$$y = \text{vl_nnconv}(x, w, b)$$

Anatomy of a building block

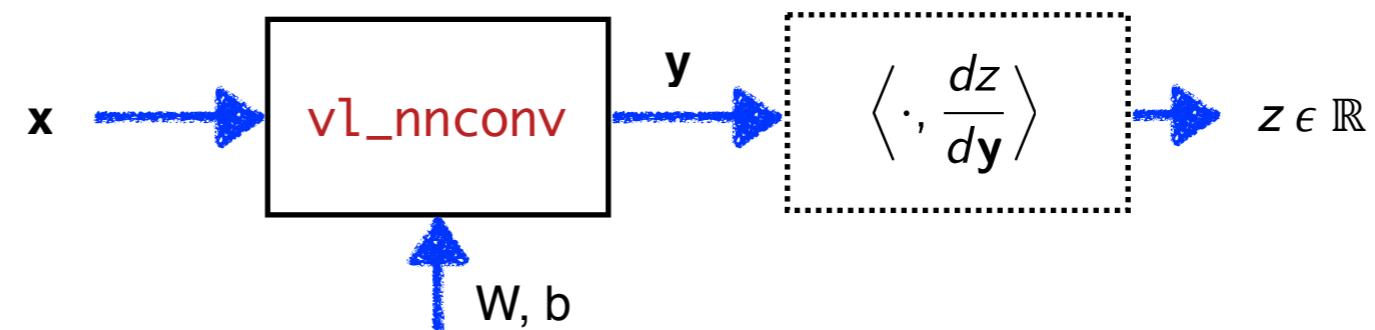
forward (eval)



$$y = \text{vl_nnconv}(x, W, b)$$

Anatomy of a building block

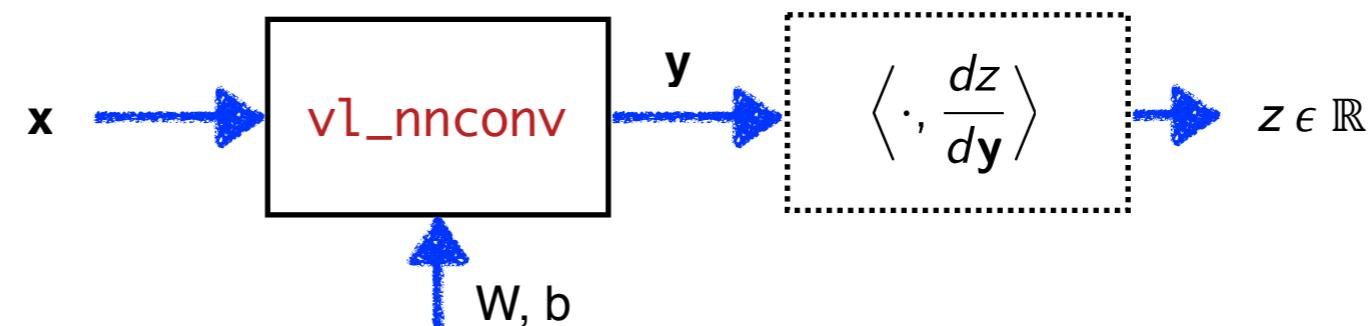
forward (eval)



$$y = \text{vl_nnconv}(x, W, b)$$

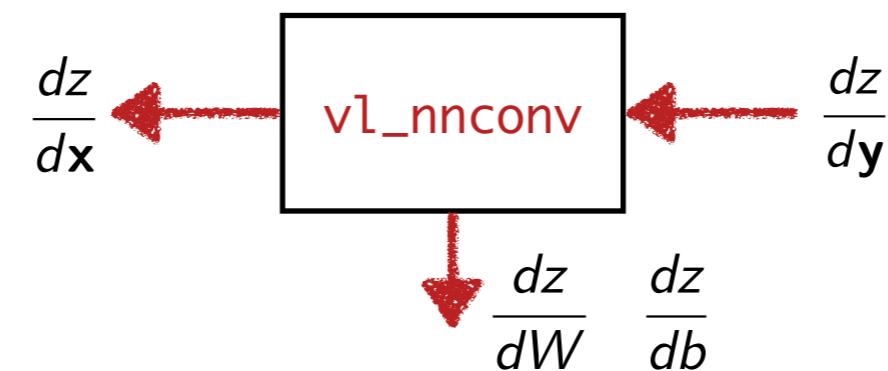
Anatomy of a building block

forward (eval)



$$y = \text{vl_nnconv}(x, W, b)$$

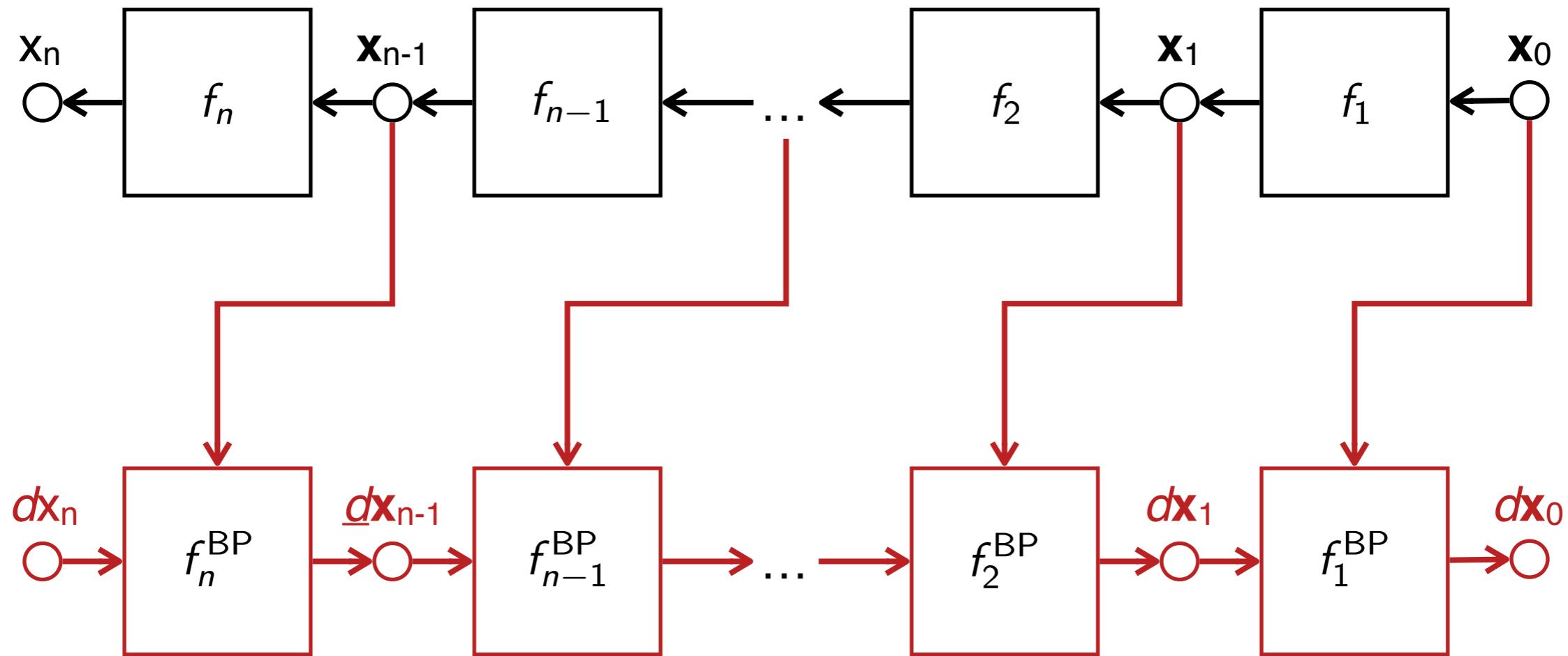
backward (backprop)



$$\frac{dz}{dx} = \text{vl_nnconv}(x, W, b, \frac{dz}{dy})$$

Backpropagation network

BP induces a “reversed” network

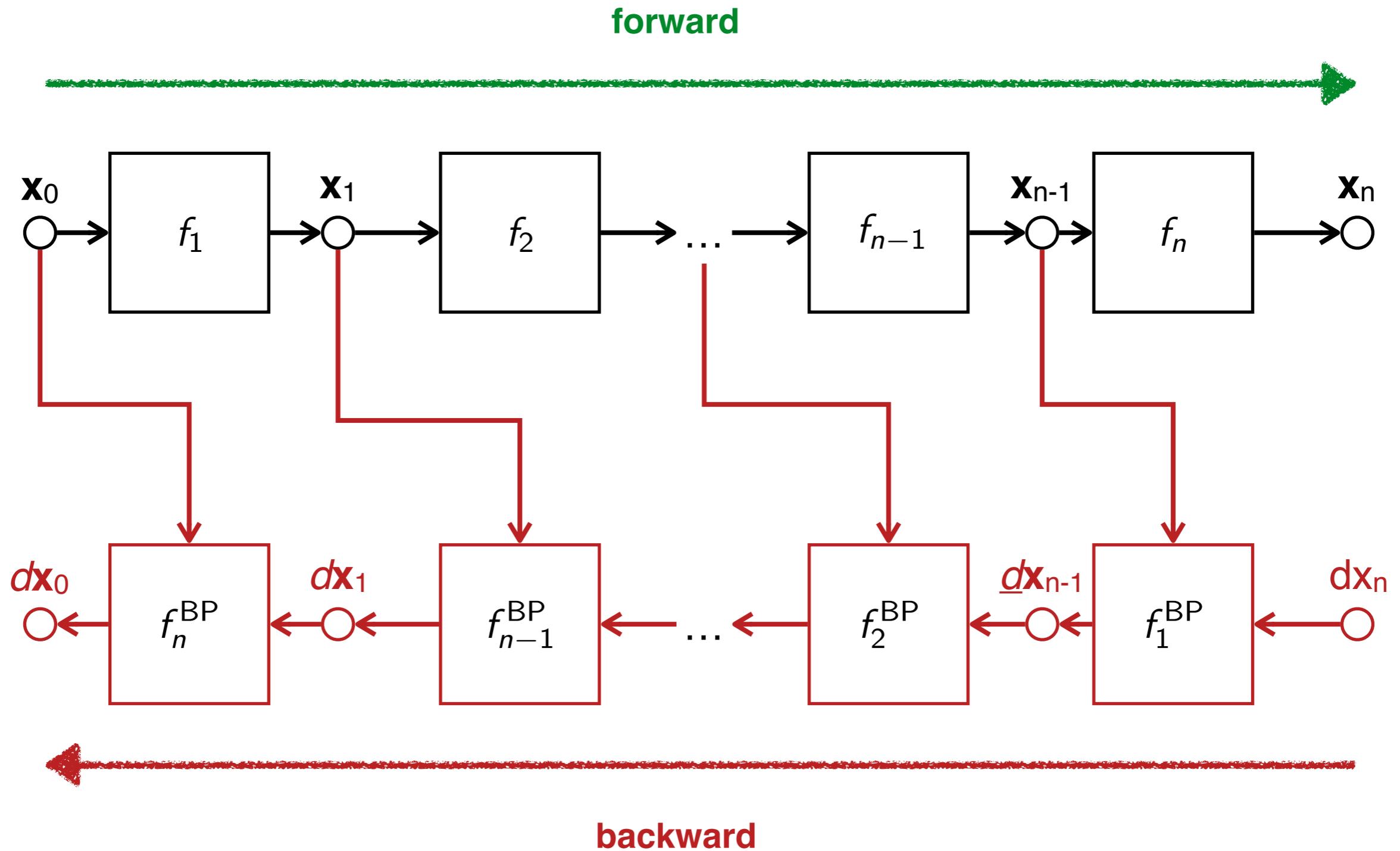


where $d\mathbf{x}_i = \frac{df_n \circ \cdots \circ f_{i+1}}{d \text{vec } \mathbf{x}_i}$

Note: the BP network is linear
in $d\mathbf{x}_1, \dots, d\mathbf{x}_{n-1}, d\mathbf{x}_n$. Why?

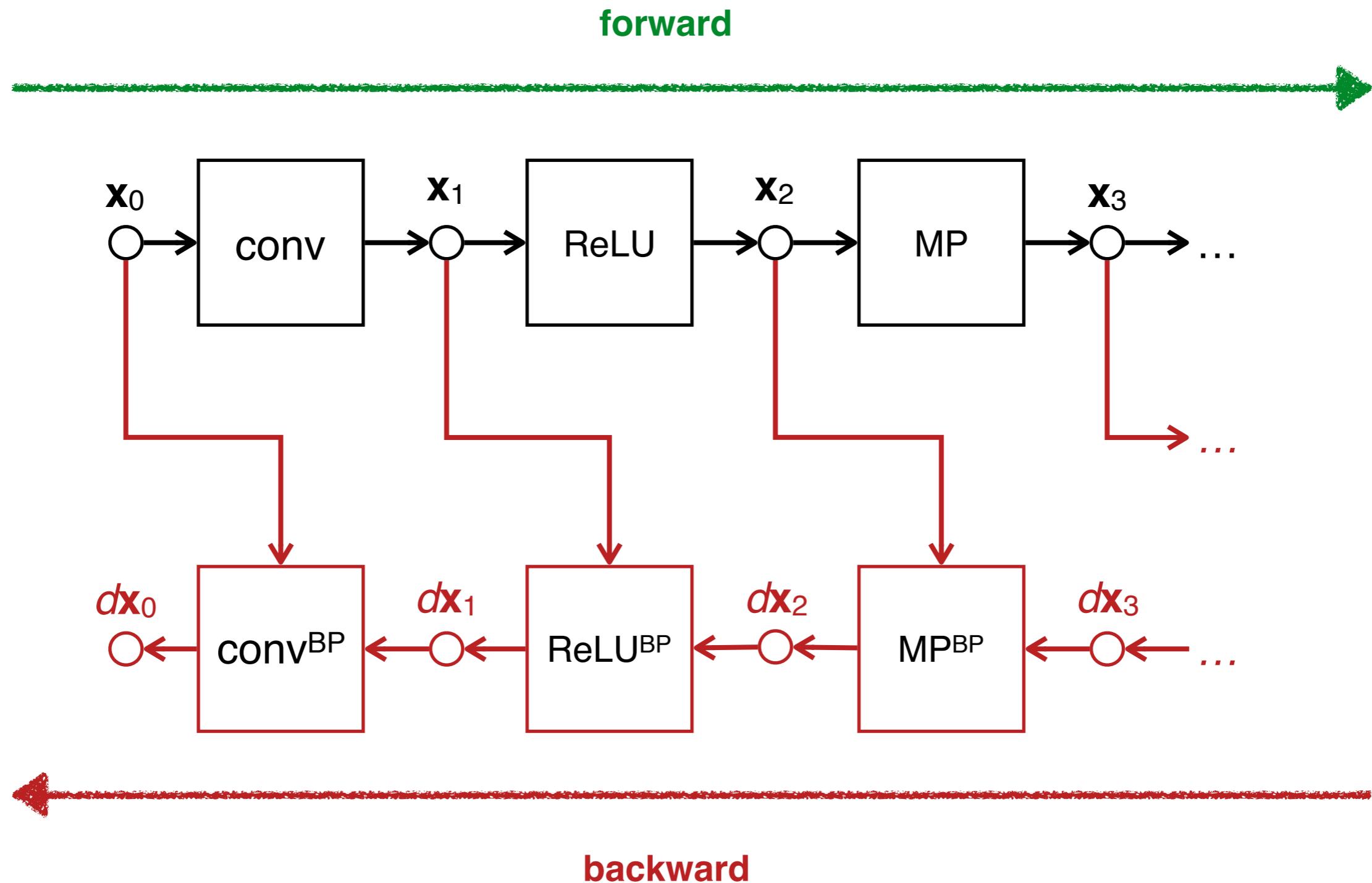
Backpropagation network

BP induces a “transposed” network



Backpropagation network

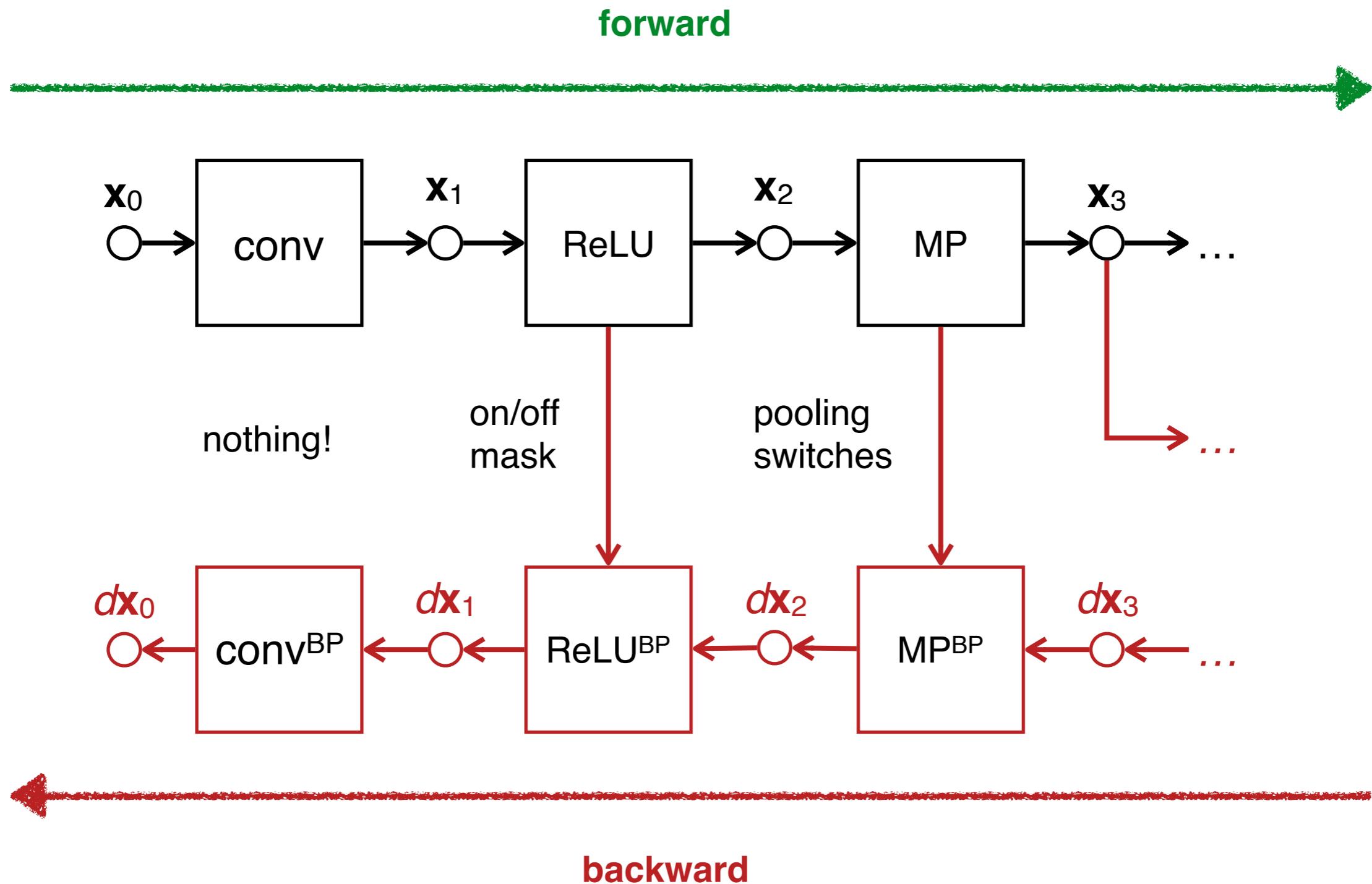
Conv, ReLU, MP and their transposed blocks



Sufficient statistics and bottlenecks

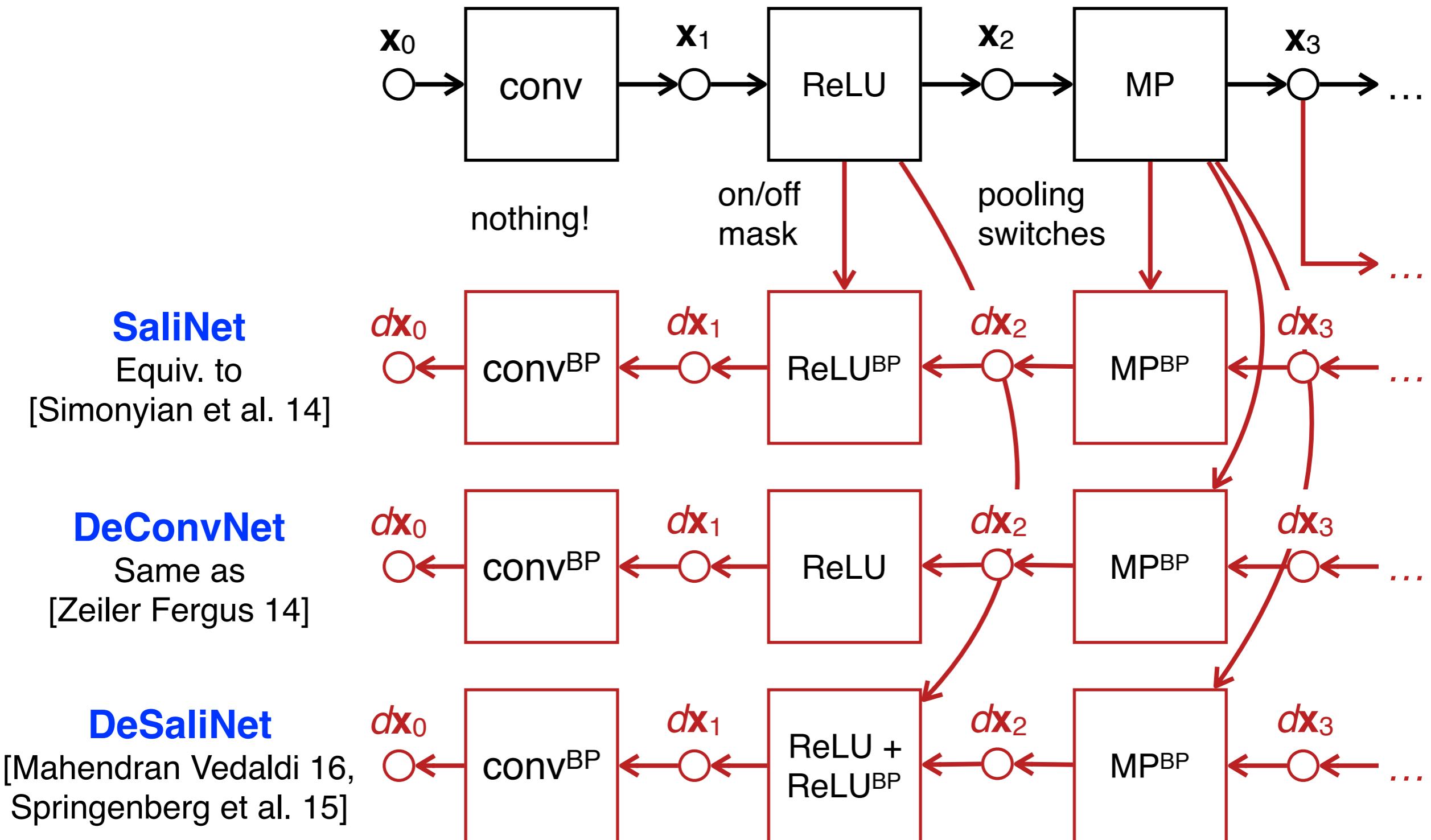
119

Usually much less information is needed

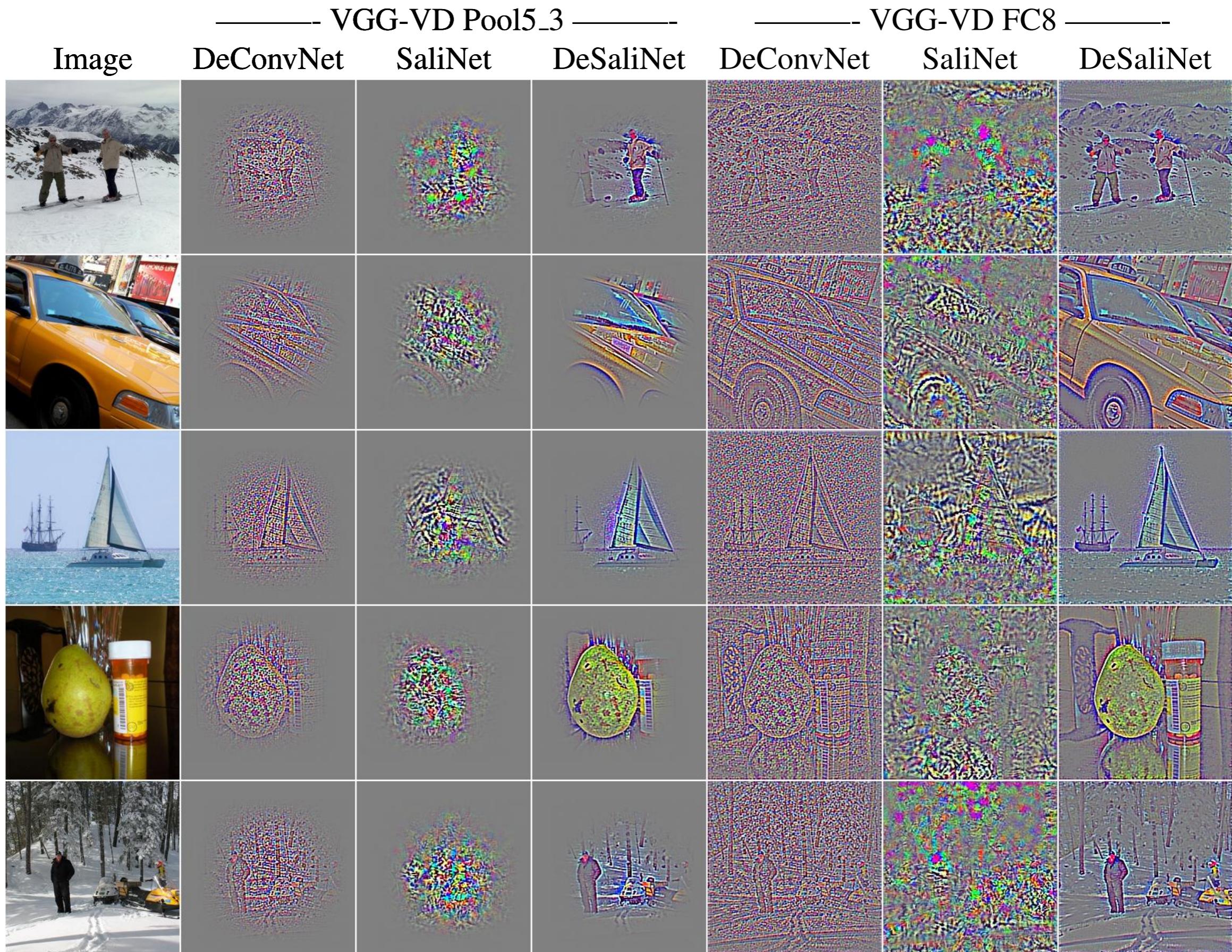


Three visualisation techniques

Modified backpropagation networks



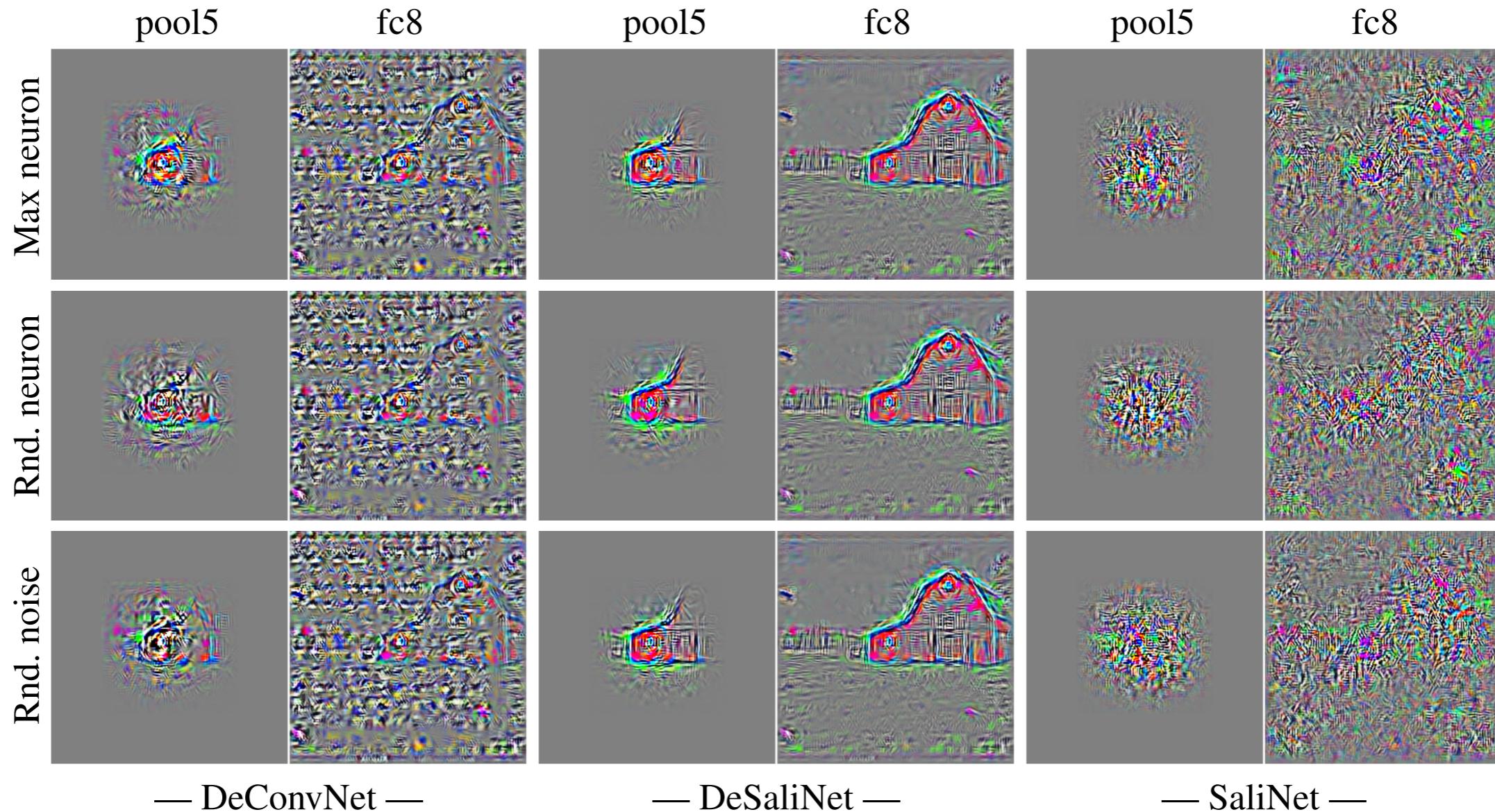
Results



Key limitation of DeConvNets and similar

122

They are largely *not* neuron selective



Good for saliency

Bad for studying individual neurons
(use inversion, act. max., etc. instead)

Intro

Visualizing representations

Backpropagation networks and “deconvolution”

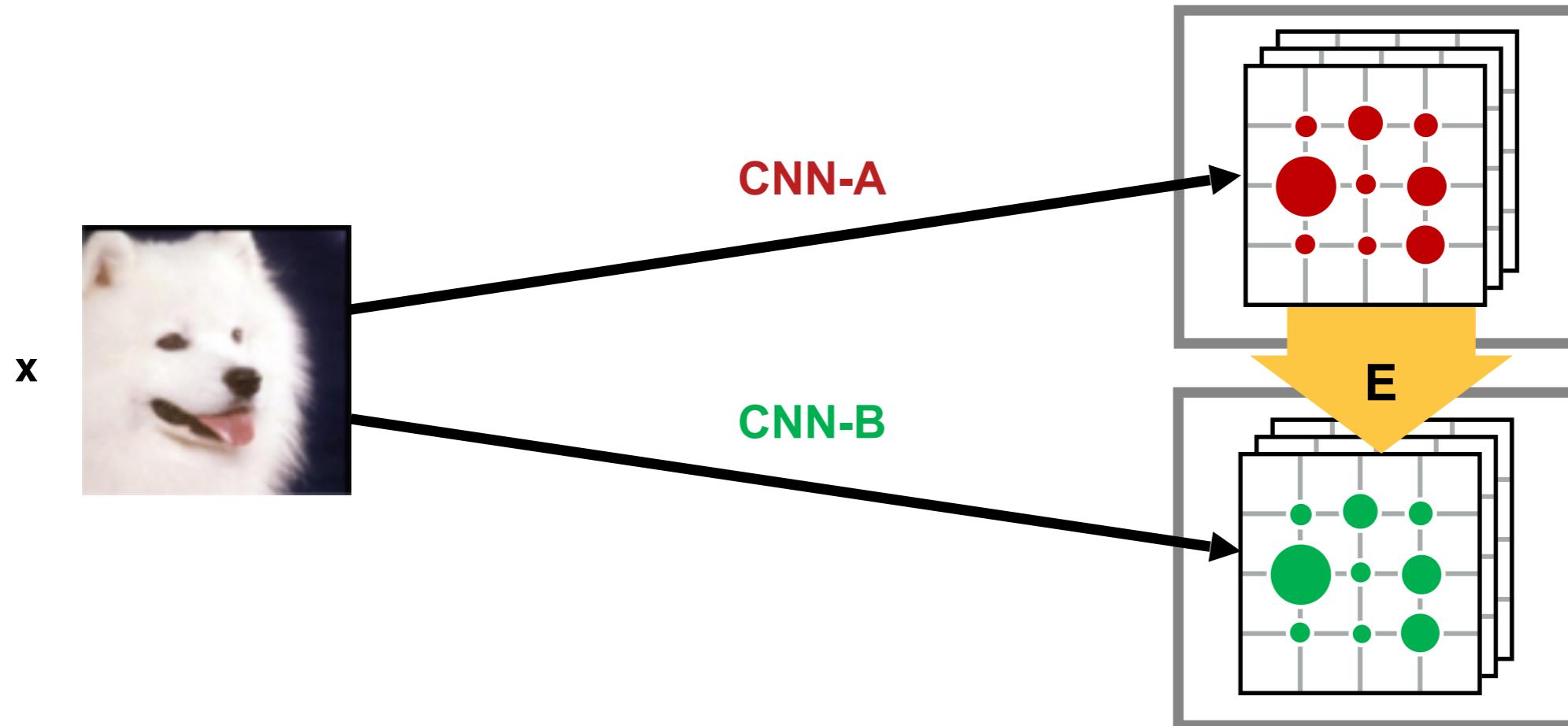
Representations: equivalence & transformations

When are two representations the same?

124

Learning representations means that there is an endless number of them

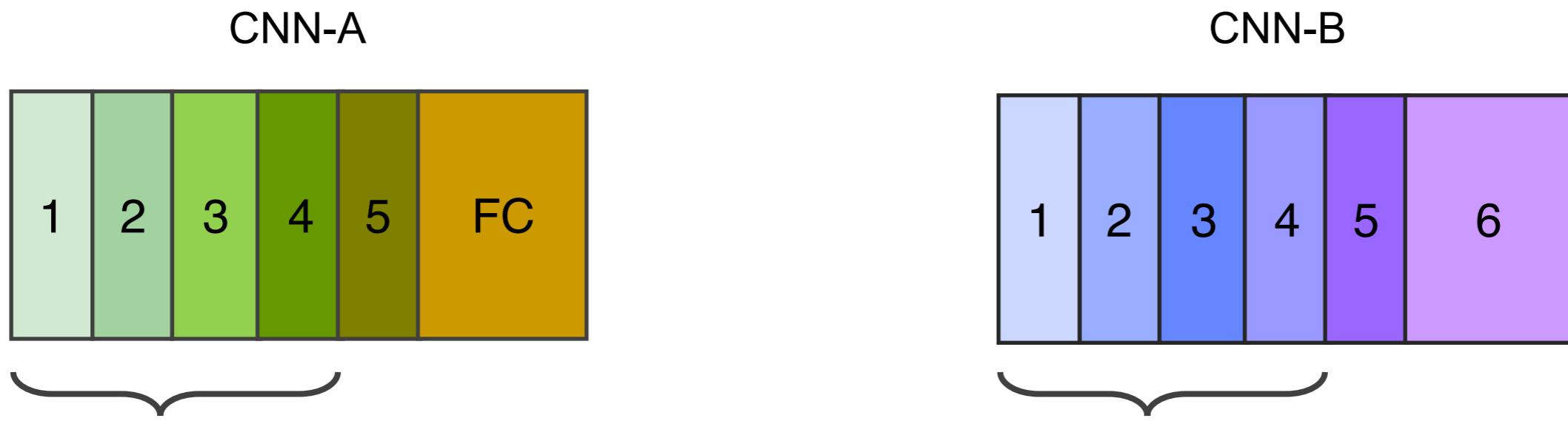
Variants obtained by learning on different datasets, or different local optima



Equivalence
 $\Phi_B(x) = E \Phi_A(x)$

Equivalence

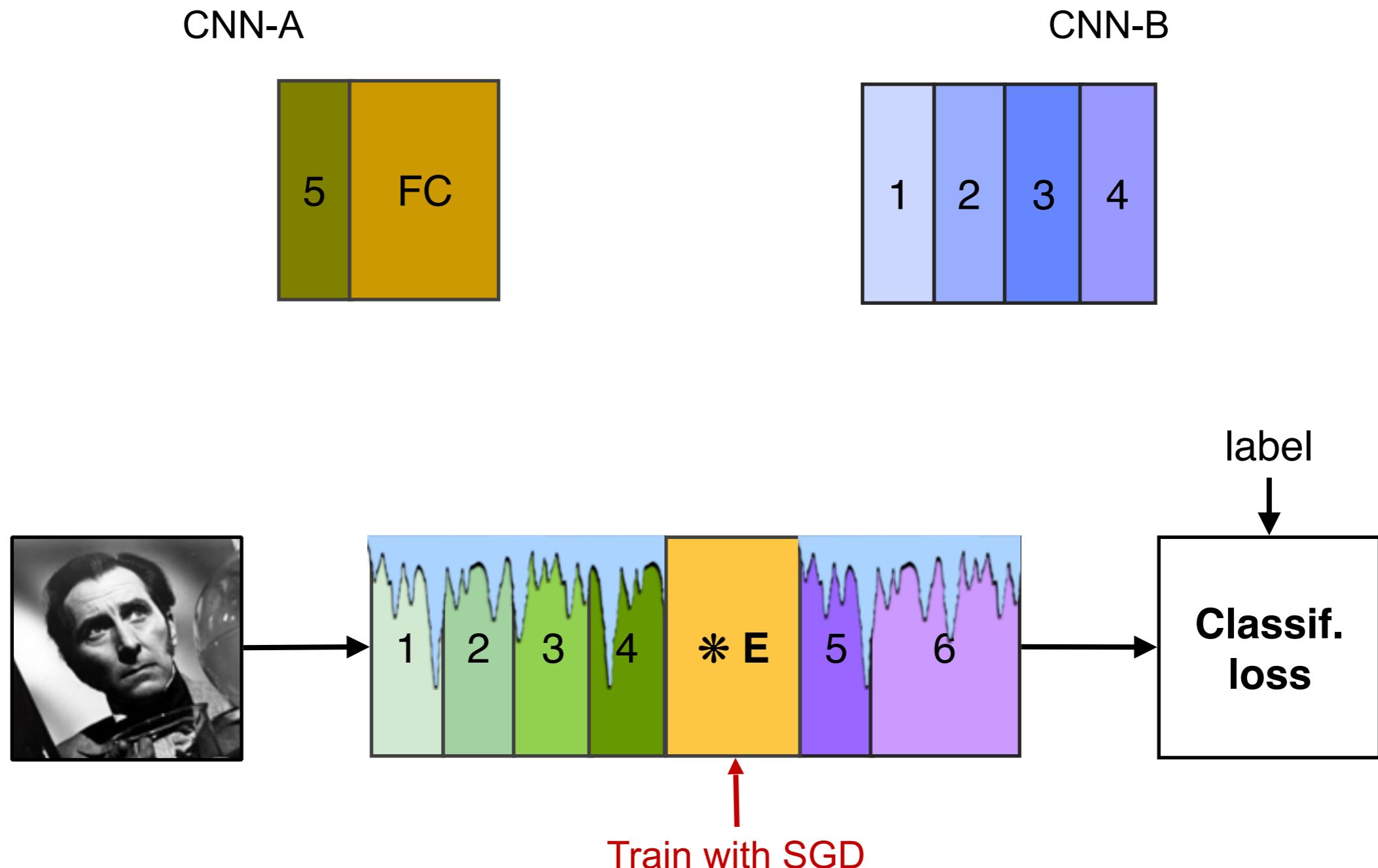
AlexNet, same training data, different parametrization:



Are Φ_A and Φ_B equivalent ?

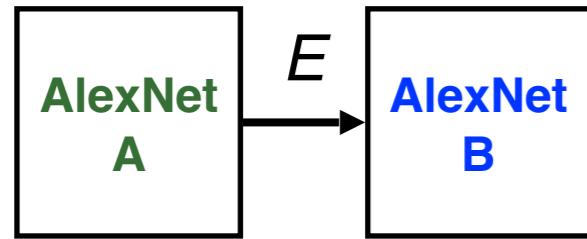
Equivalence

AlexNet, same training data, different parametrization:



Equivalence with different random seeds

127



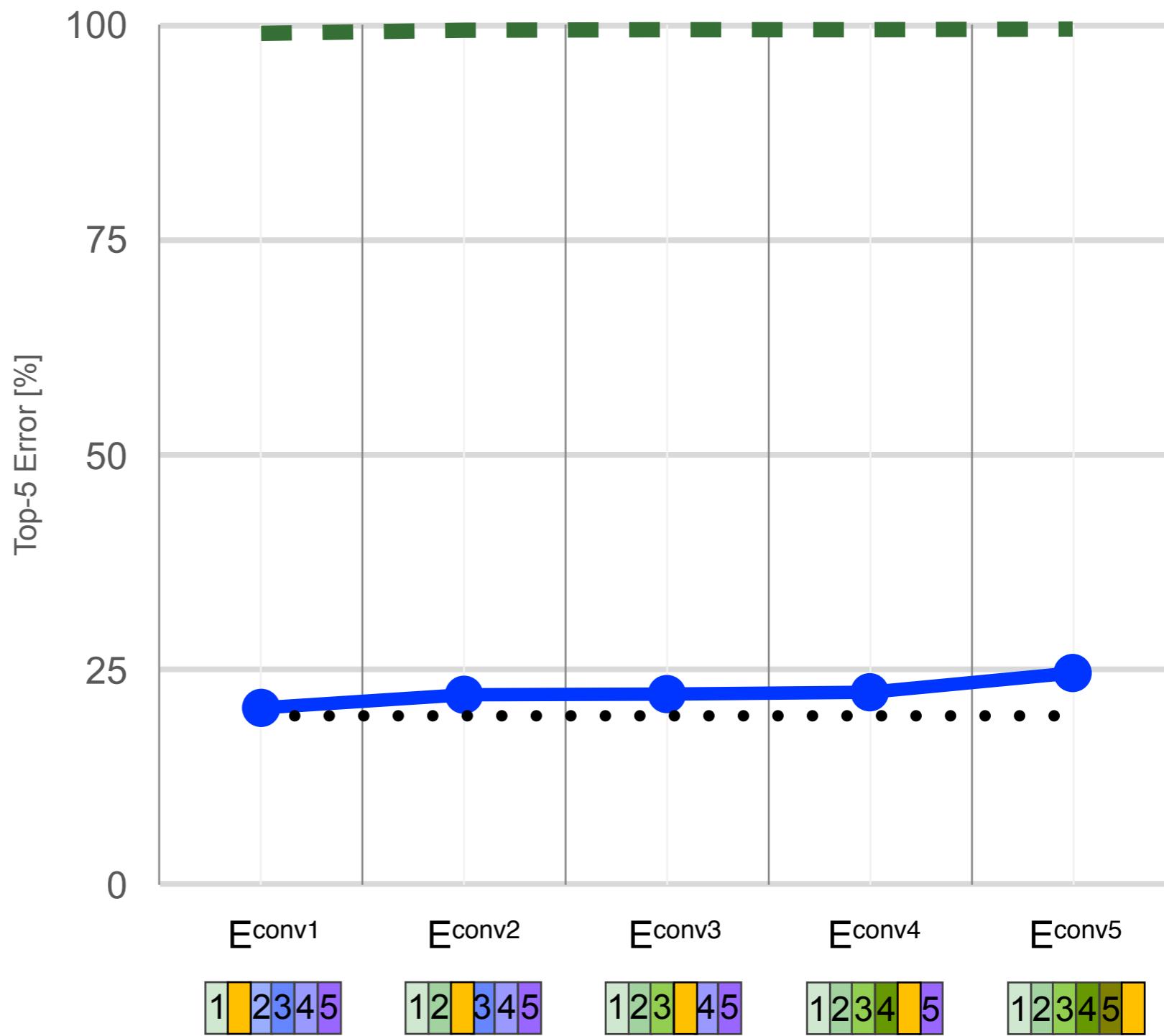
Baseline



Before training



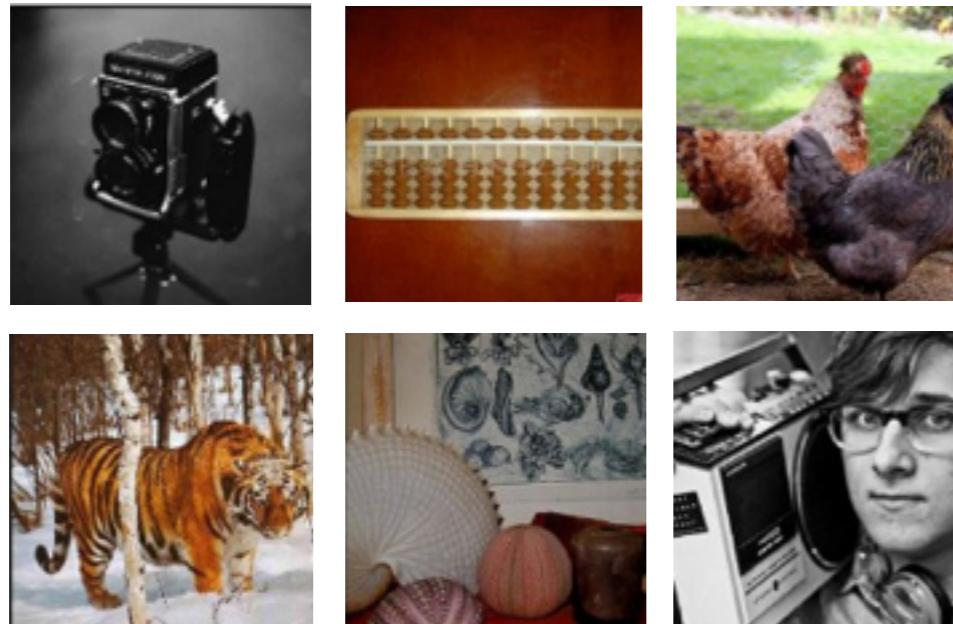
After training



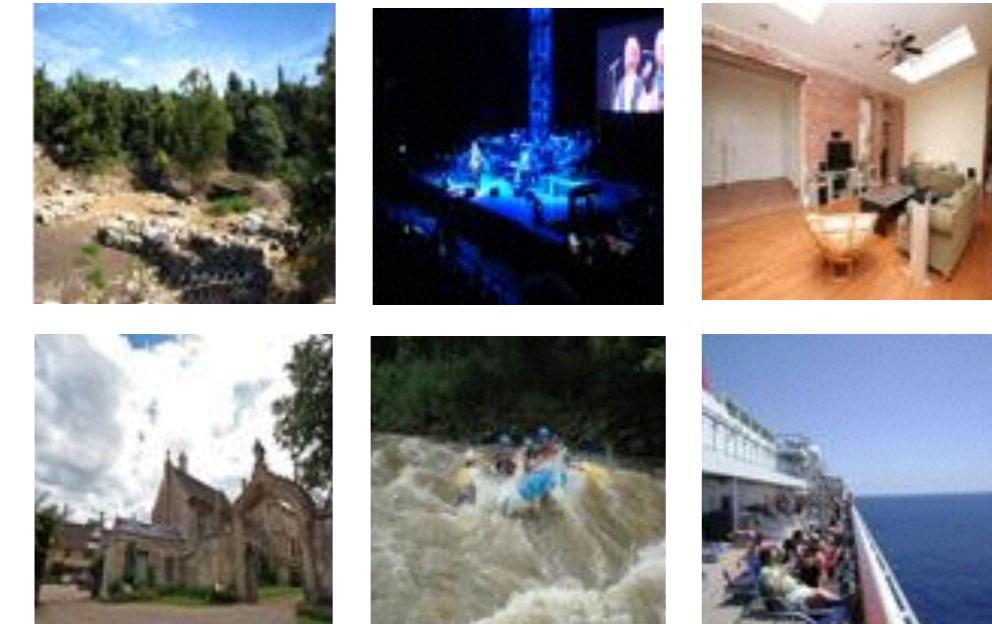
Equivalence of similar architecture

Train on two different datasets

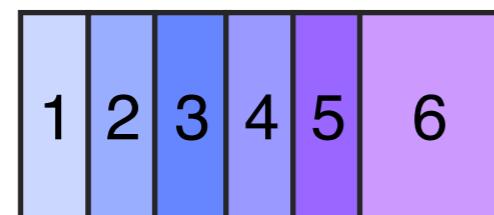
ILSVRC12 dataset



Places dataset



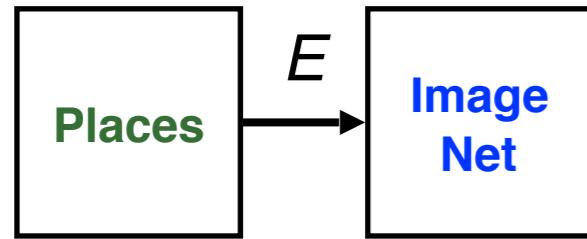
CNN-IMNET



CNN-PLACES



Equivalence with different training data



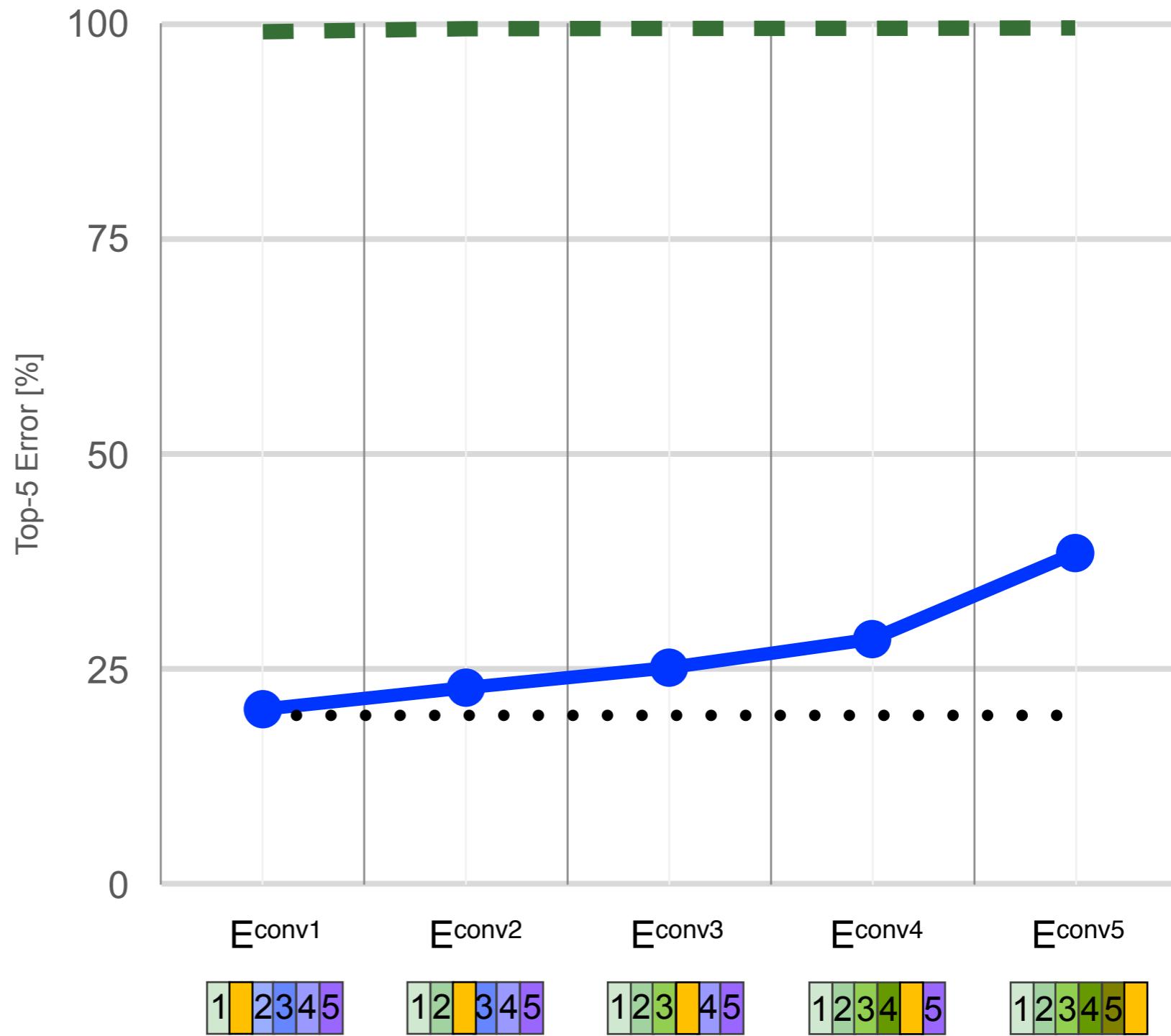
Baseline



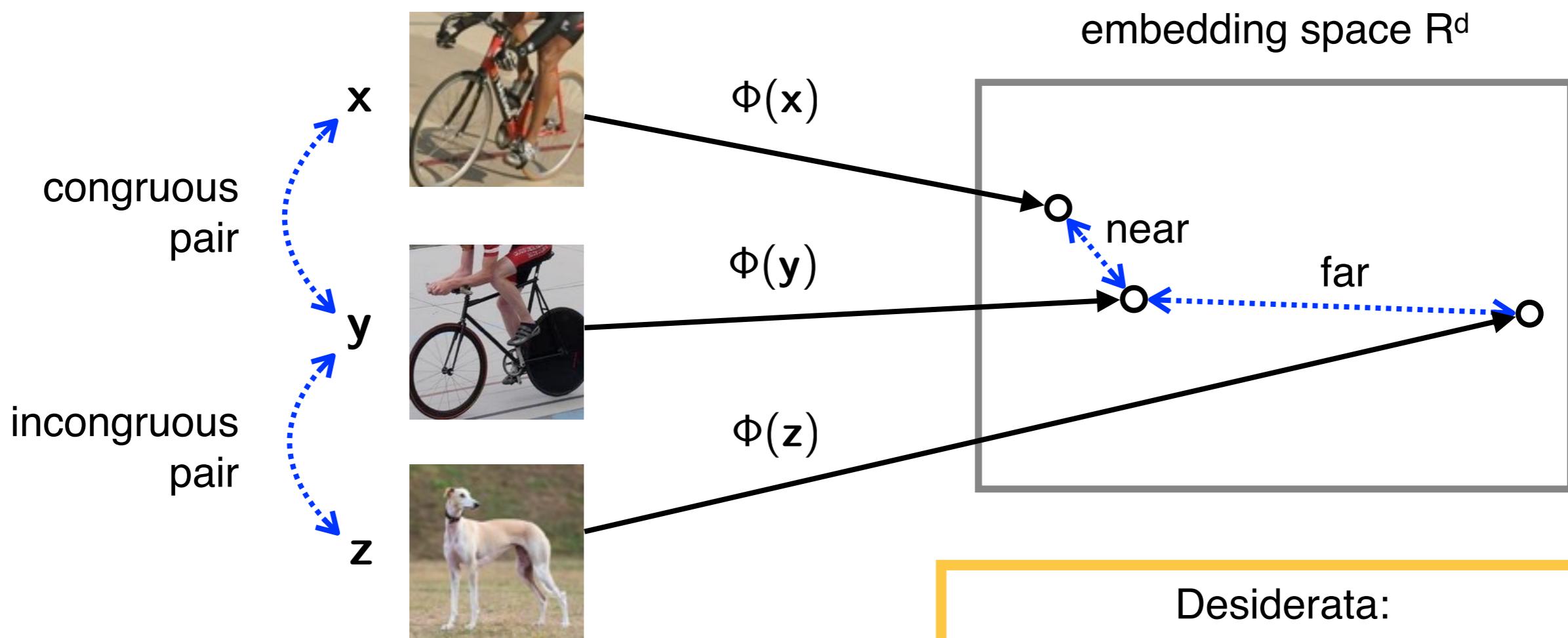
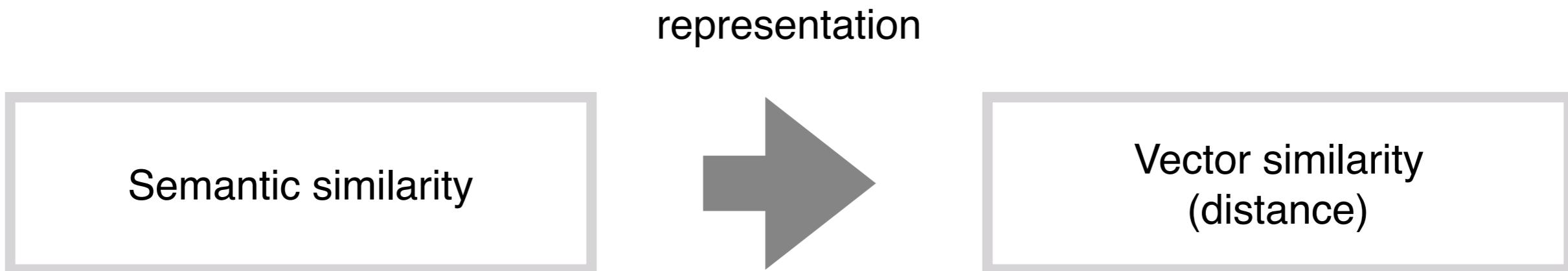
Naive stitching



Learned stitches

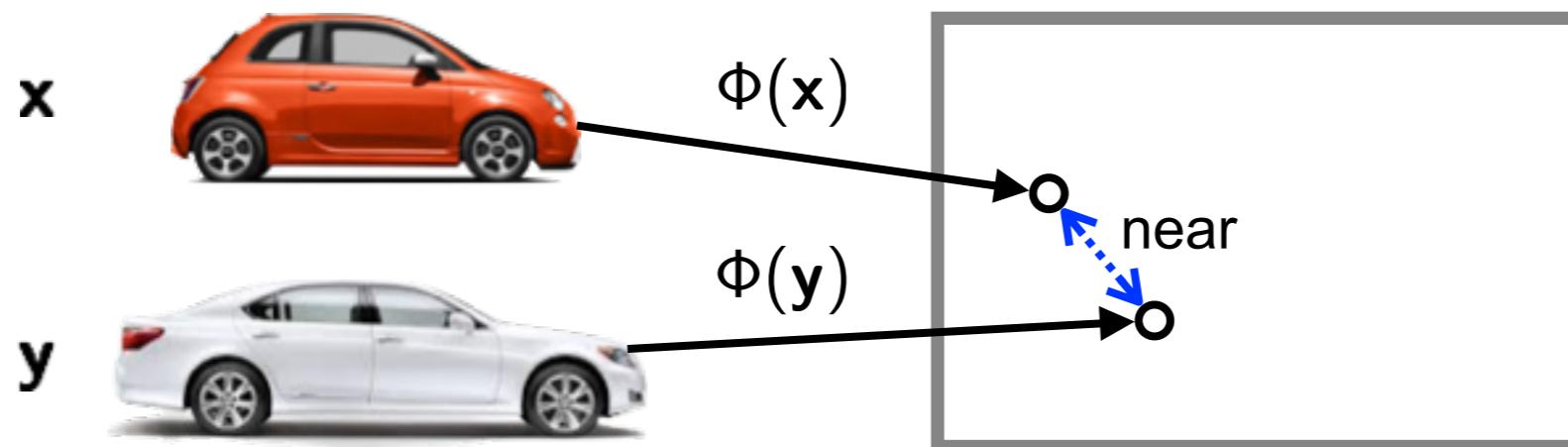


Meaningful representations

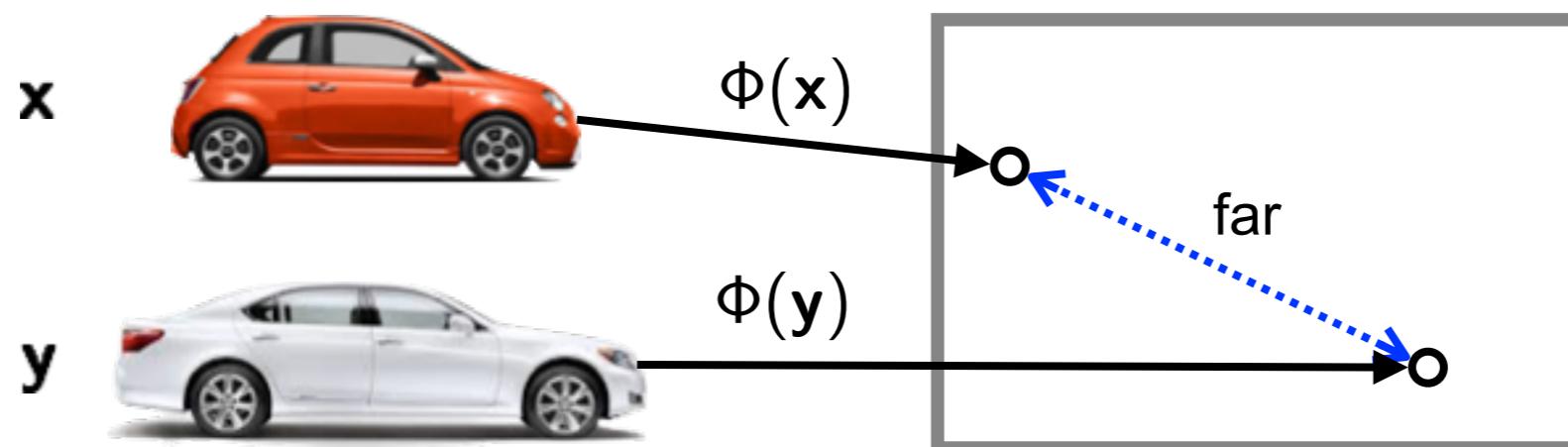


Invariance is task dependent

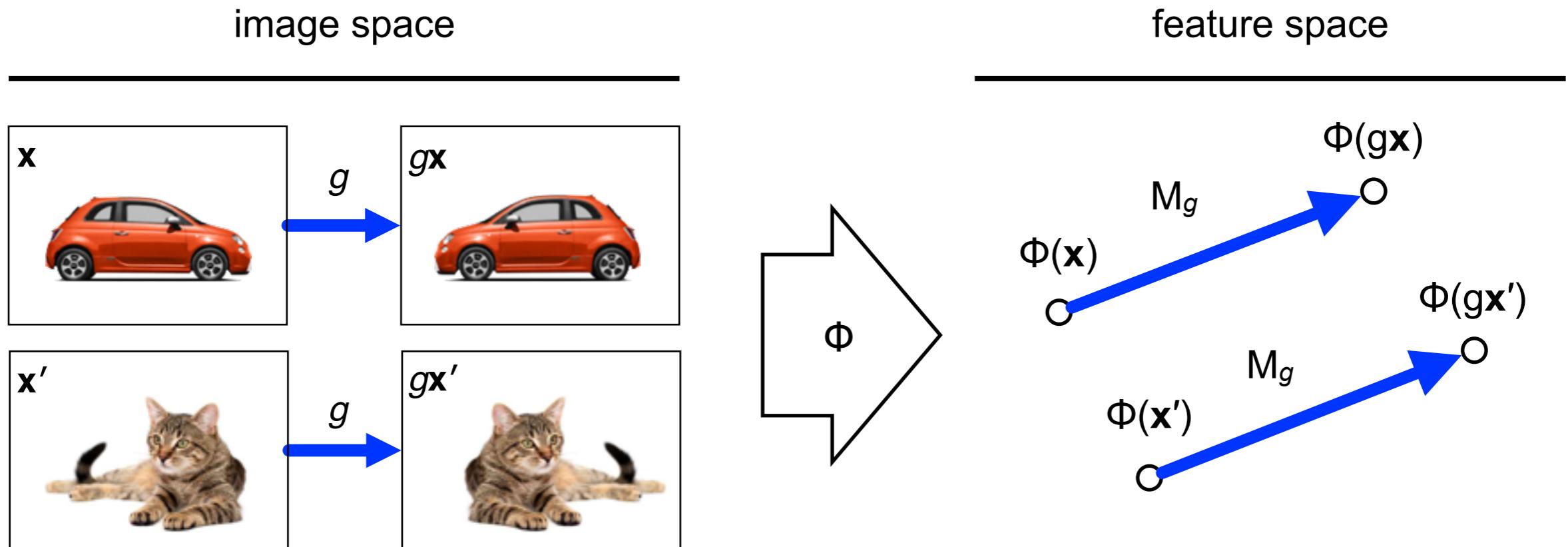
Are x and y the same category?



Are x and y the same colour?



Equivariance

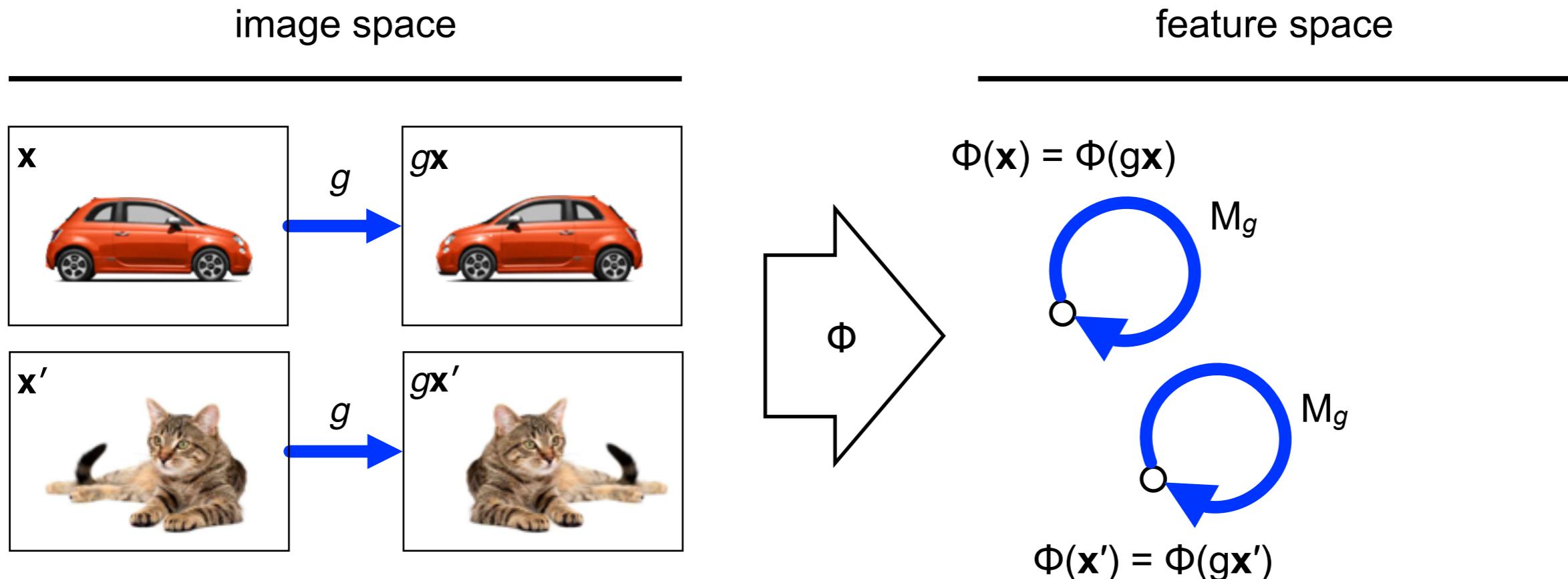


g is an image transformation

g is the corresponding
feature transformation

$$\Phi(g\mathbf{x}) = M_g \Phi(\mathbf{x})$$

Invariance

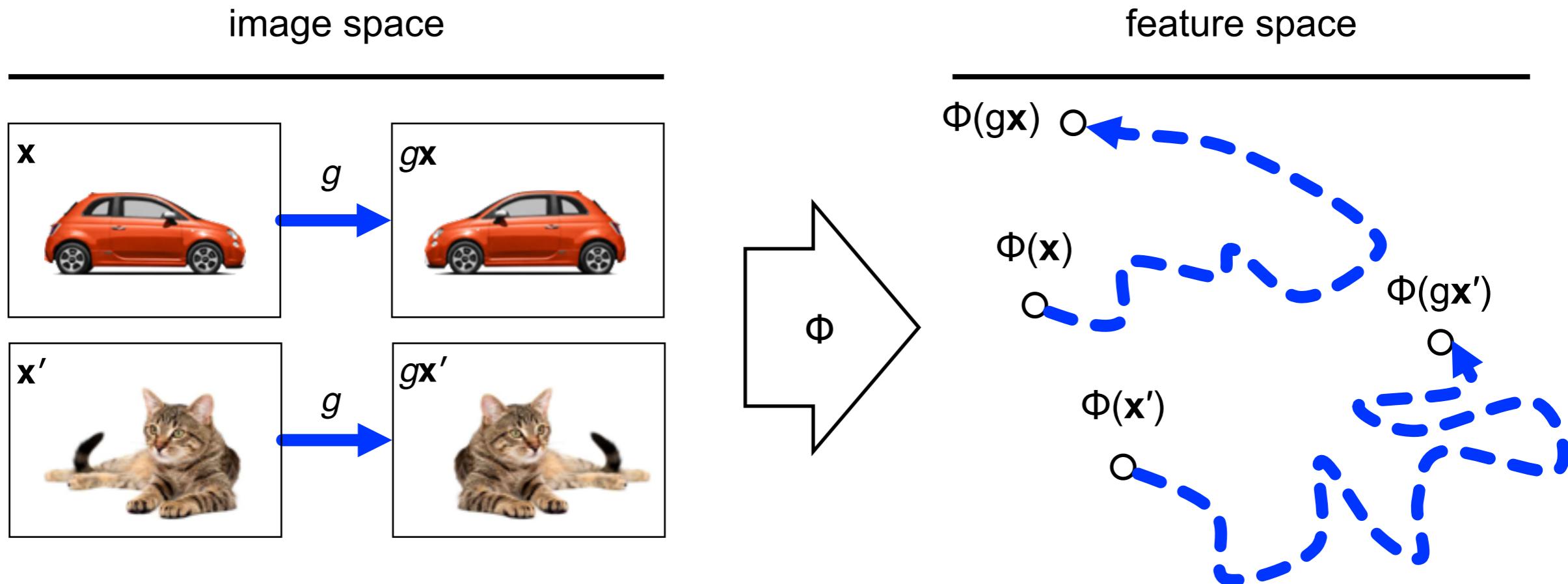


g is an image transformation

g is the corresponding
feature transformation

the map M_g is the identity

No equivariance



g is an image transformation

g is the corresponding
feature transformation

the map M_g does not
exist (or is intractable)

Representations and transformations

135

image space

g is an image transformation

feature space

Equivariance
There is a (simple) M_g

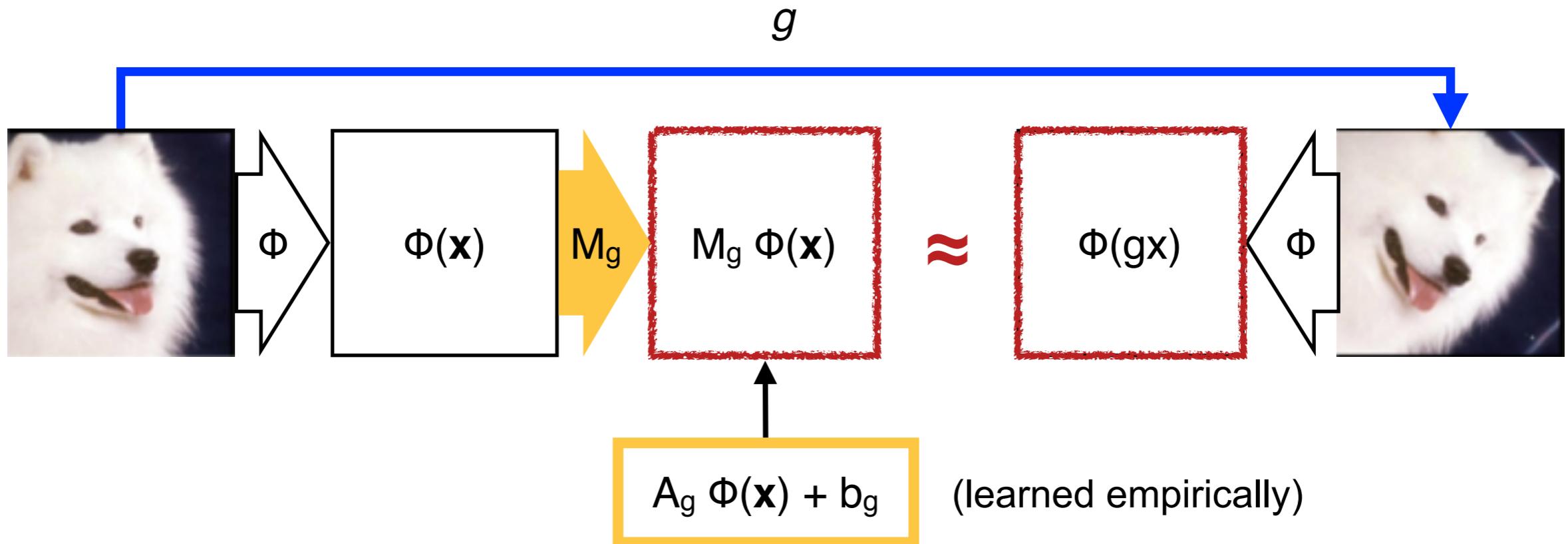
Invariance
 M_g is the identity

Neither
There is no (simple) M_g

$$\Phi(g\mathbf{x}) = M_g \Phi(\mathbf{x})$$

An empirical test of equivariance

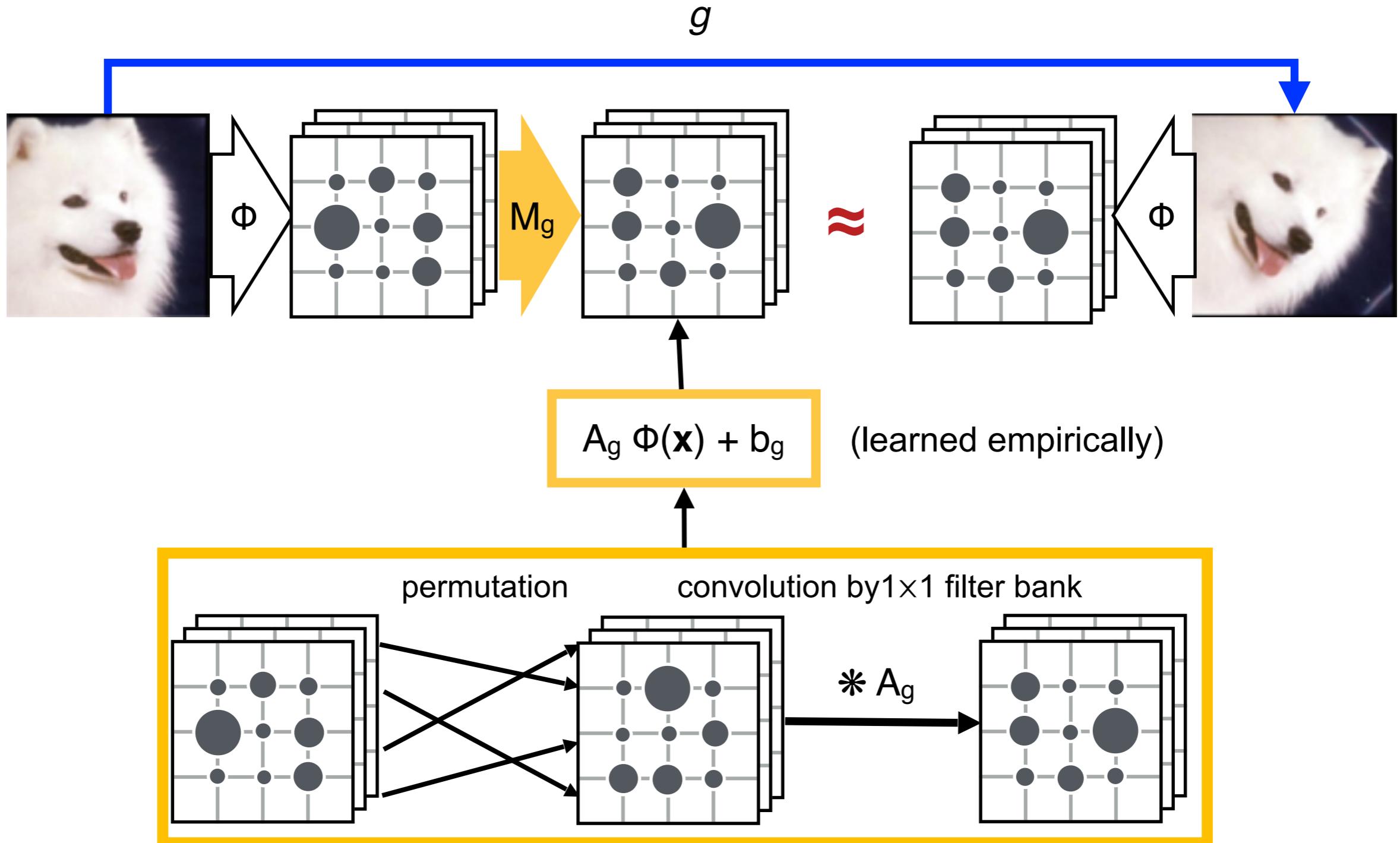
Regularized linear regression



An empirical test of equivariance

137

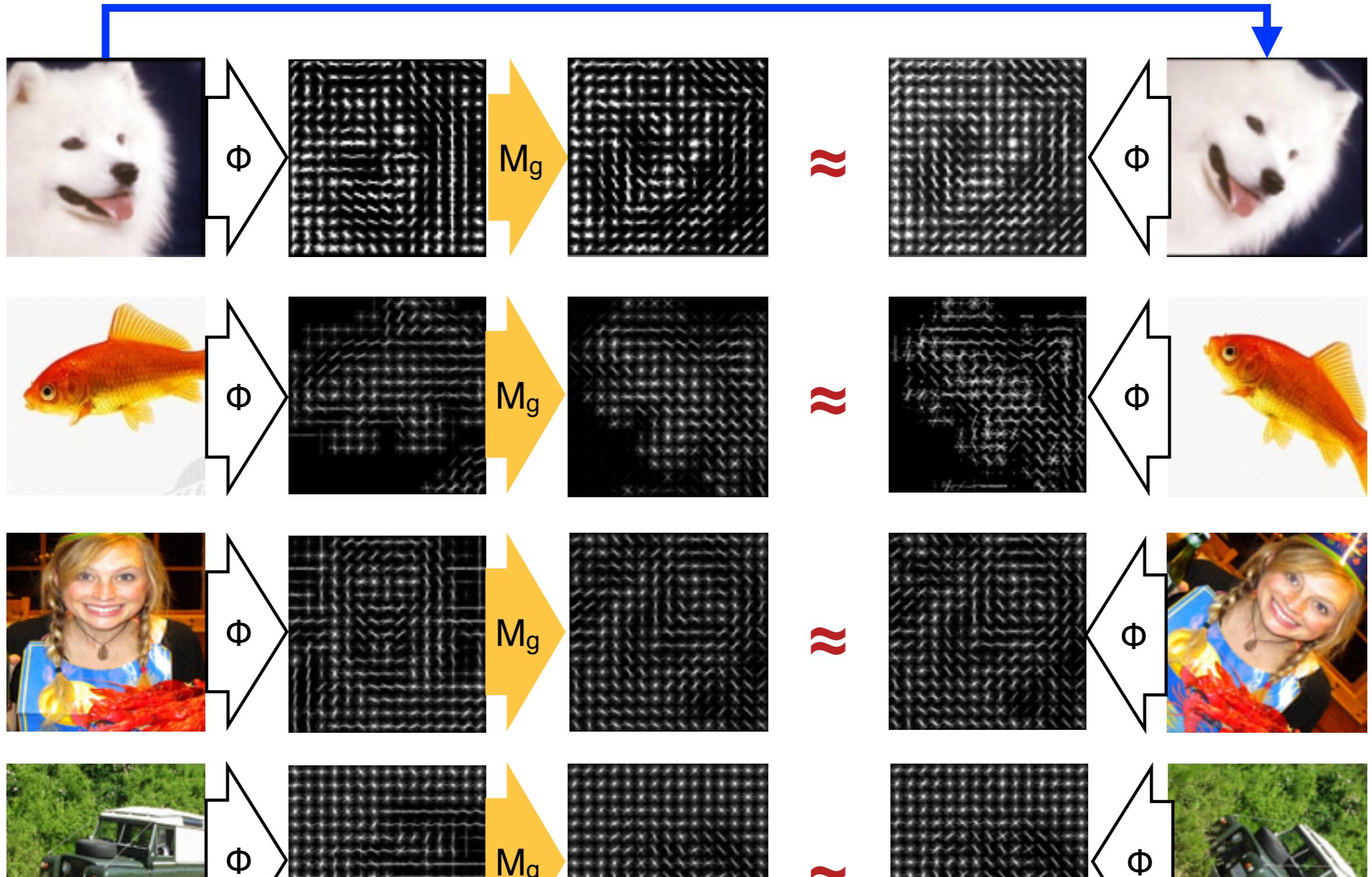
Regularized linear regression



An empirical test of equivariance

HOG features

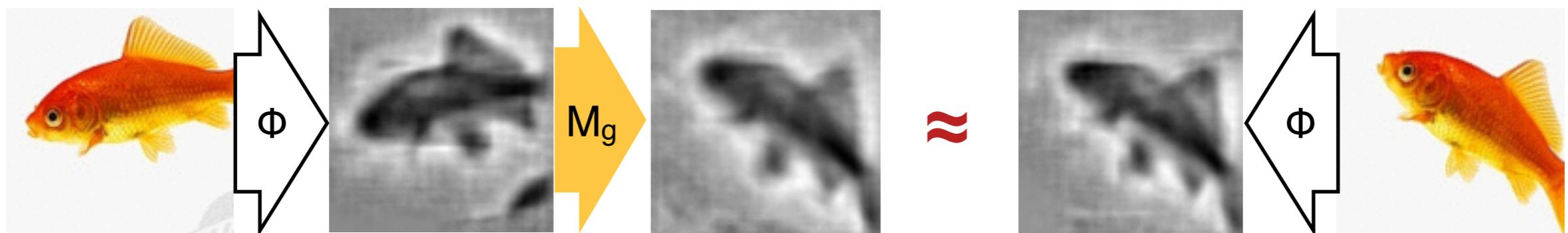
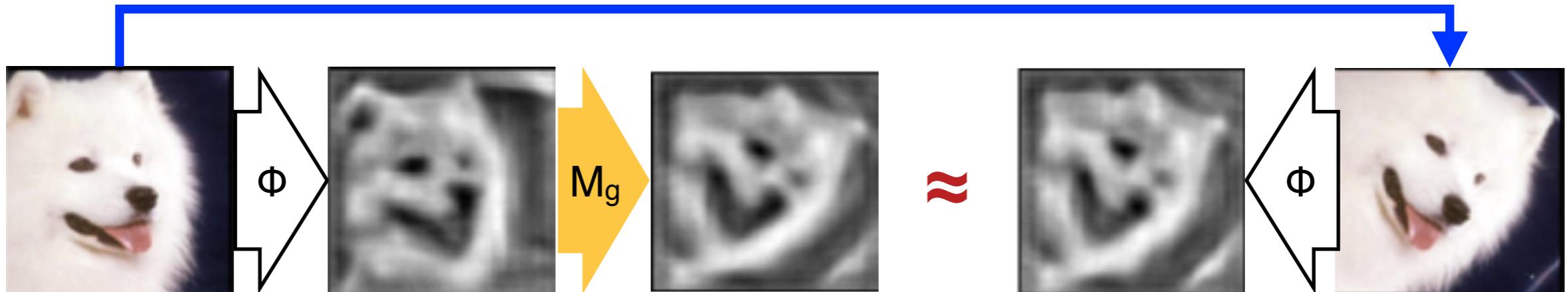
rotation 45 deg



An empirical test of equivariance

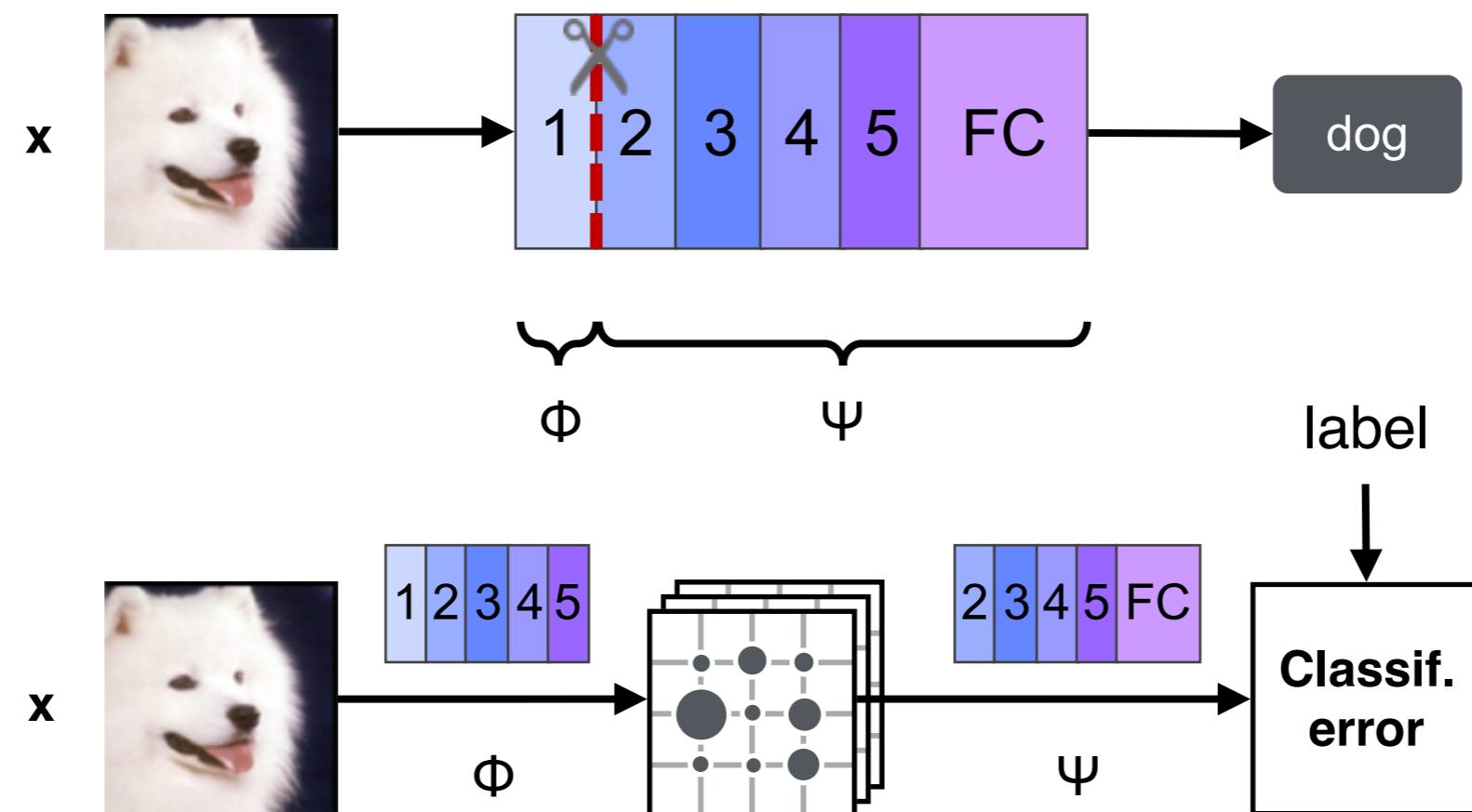
HOG features

rotation 45 deg



CNN: a sequence of representations

140

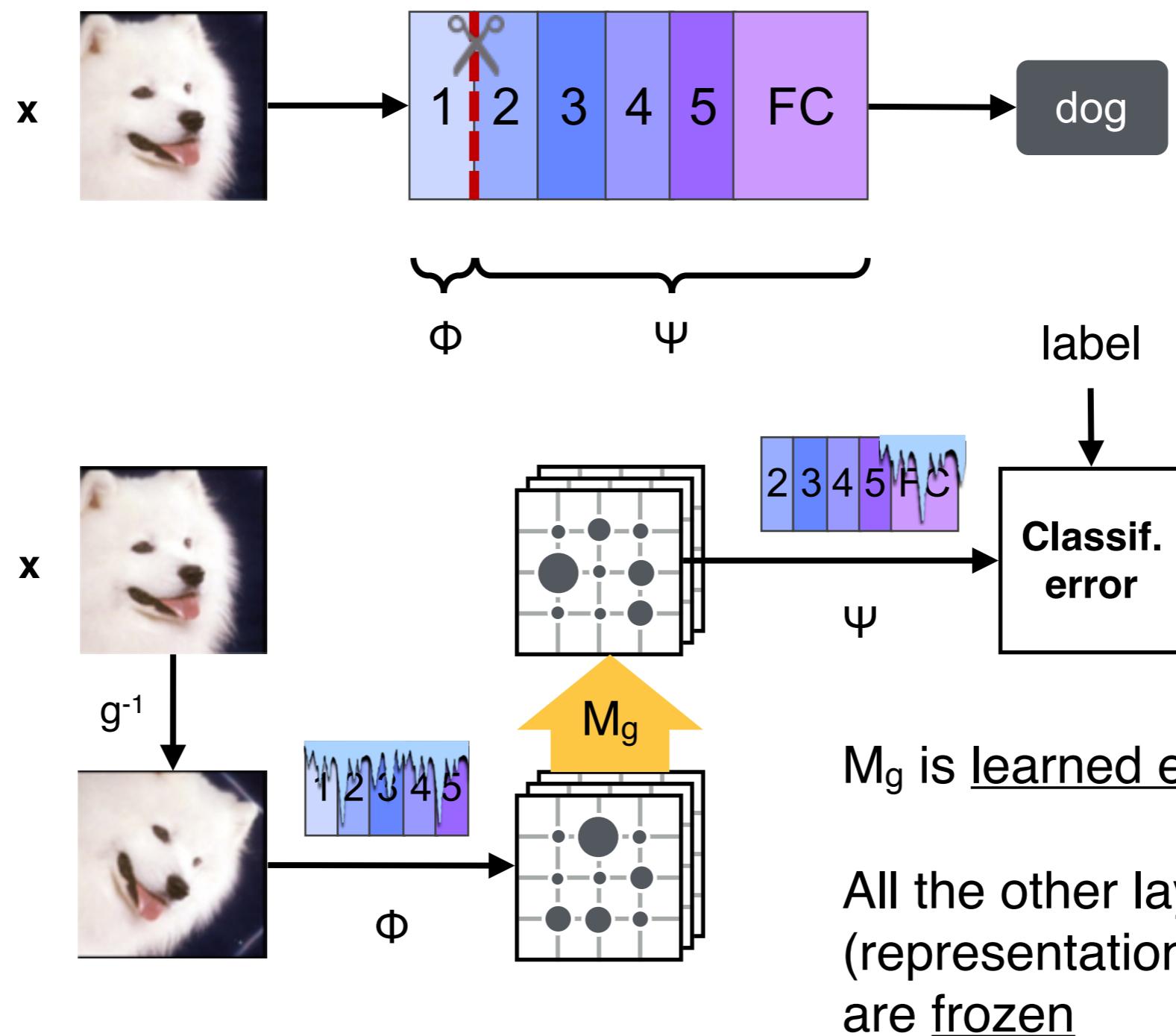


We run the same analysis on a typical CNN architecture

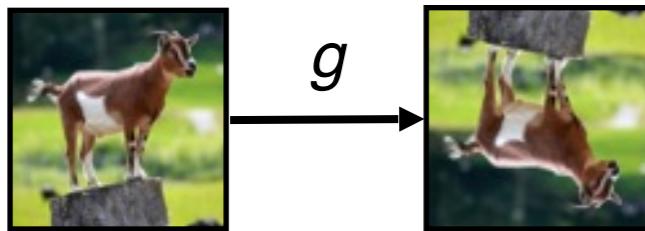
- ▶ AlexNet [Krizhevsky et al. 12]
- ▶ 5 convolutional layers + fully-connected layers
- ▶ Trained on ImageNet ILSVRC

A discriminative goal to learn equivariance

141



Vertical flips



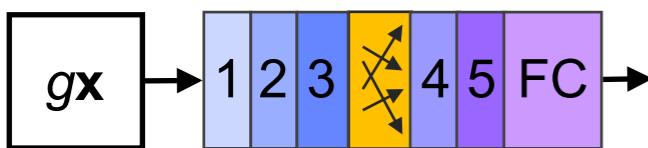
Uncompensated, no TF



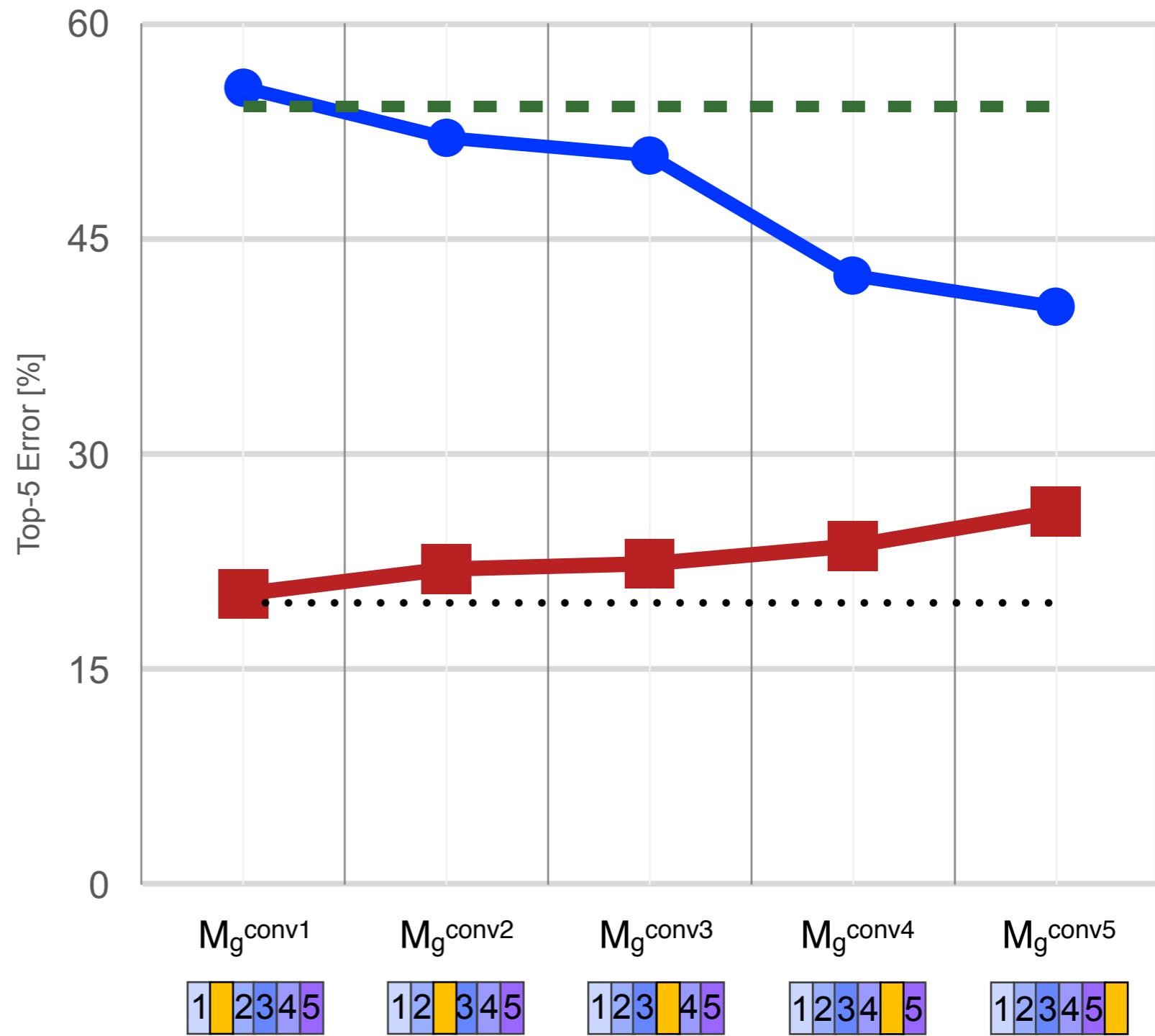
Uncompensated, TF



Compensated TF w perm.



Compensated TF, learned



Summary

Modern CNNs are still “new” technology

- ▶ Expect bigger and meaner CNNs to further improve performance
- ▶ CNNs are more than big balls of parameters
- ▶ We do not really understand what they do

Understanding deep nets

- ▶ What do deep networks do?
- ▶ Are there computational / statistical principles to be learned?
- ▶ How much do deep nets learn about the visual world?

Possible angles

- ▶ Visualisations
- ▶ Probing statistical properties

Are we there yet?

No!

Integrated computer vision

- ▶ One network to rule them all

Cognition

- ▶ Integrate perception and “the rest”
- ▶ Planning, memory, background knowledge, ...

No labels required

- ▶ Unsupervised
- ▶ Alternatively supervised (reinforcement, pseudo-tasks)

Spatial reasoning

- ▶ From image-centric to object-centric and scene-centric understanding
- ▶ Representing deformable 3D shape
- ▶ Dynamic environments, physics

Understanding visual representations

Learnable representations, and deep convolutional neural networks (CNNs) in particular, have become the preferred way of extracting visual features for image understanding tasks, from object recognition to semantic segmentation. Part of the reason for the power of these models is that they contain millions of parameters learned from millions of images. However, this also means that deep networks are black boxes of difficult interpretation.

In this talk, I will focus on the problem of understanding deep networks. First, I will discuss a number of techniques to visualise the feature learned in a CNN: inversion, exploring which visual information is retained in a deep representation, activation maximisation, which finds patterns that maximally excite representation components, and caricaturisation, which emphasise which visual properties are captured by different representation layers. I will also briefly discuss methods to study how image transformations are encoded by representations, and whether different representations are equivalent.

Bio

Andrea Vedaldi is Assistant Professor of Engineering Science at the University of Oxford. His main research interests include large scale machine learning and optimization applied to semantic image understanding. He is author of more than 50 peer-reviewed publications in the top machine vision and artificial intelligence conferences and journals (CVPR, ECCV, ICCV, NIPS; h-index 34). He is project leader of VLFeat (<http://www.vlfeat.org>), a popular computer vision software toolkit (ACM Open Source Award 2010), and of MatConvNet (<http://www.vlfeat.org/matconvnet>), a MATLAB-based toolbox for convolutional neural networks. He is a recipient of the ERC Starting Grant (similar to the NSF CAREER award), of the Mark Everingham Prize (ICCV 2015) for selfless contributions to the computer vision community, a Glasstone Research Fellow, a J. J. Spooner Fellow, a winner of the PASCAL Visual Object Classification competition, and a recipient of the Best Paper Award at the British Machine Vision Conference (BMVC) and of the Open Source Software Award by the ACM. He is supported by the ERC, EPSRC, BP, and XRCE.