

Understanding CNNs using visualisation and transformation analysis

Andrea Vedaldi

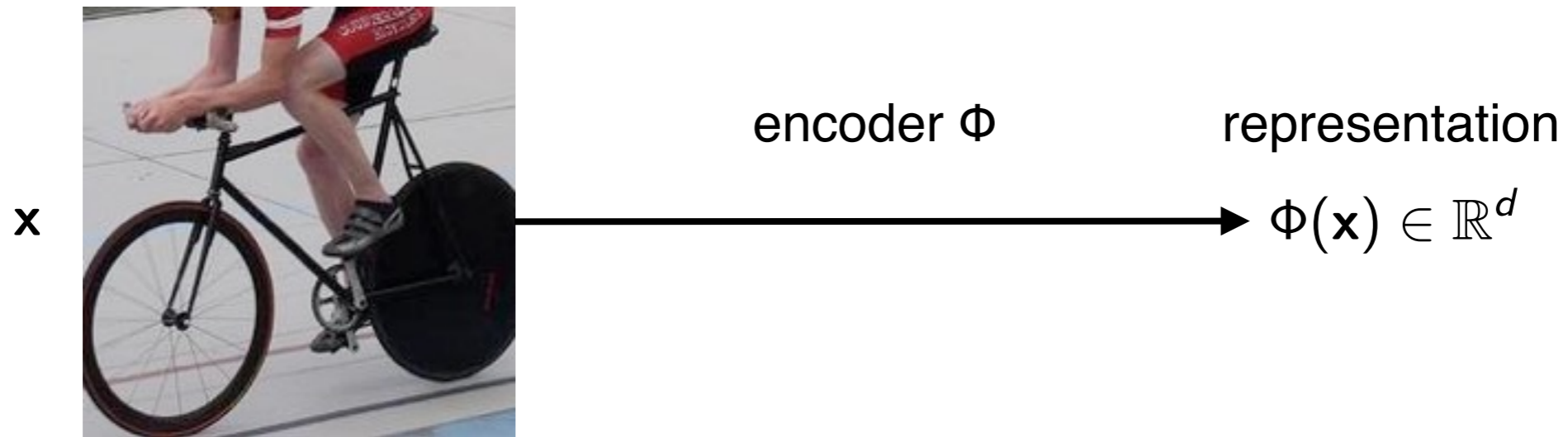
Medical Imaging Summer School

August 2016



UNIVERSITY OF
OXFORD

Image representations



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

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



Excellent **performance** in image understanding tasks

Millions of parameters learned from data

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

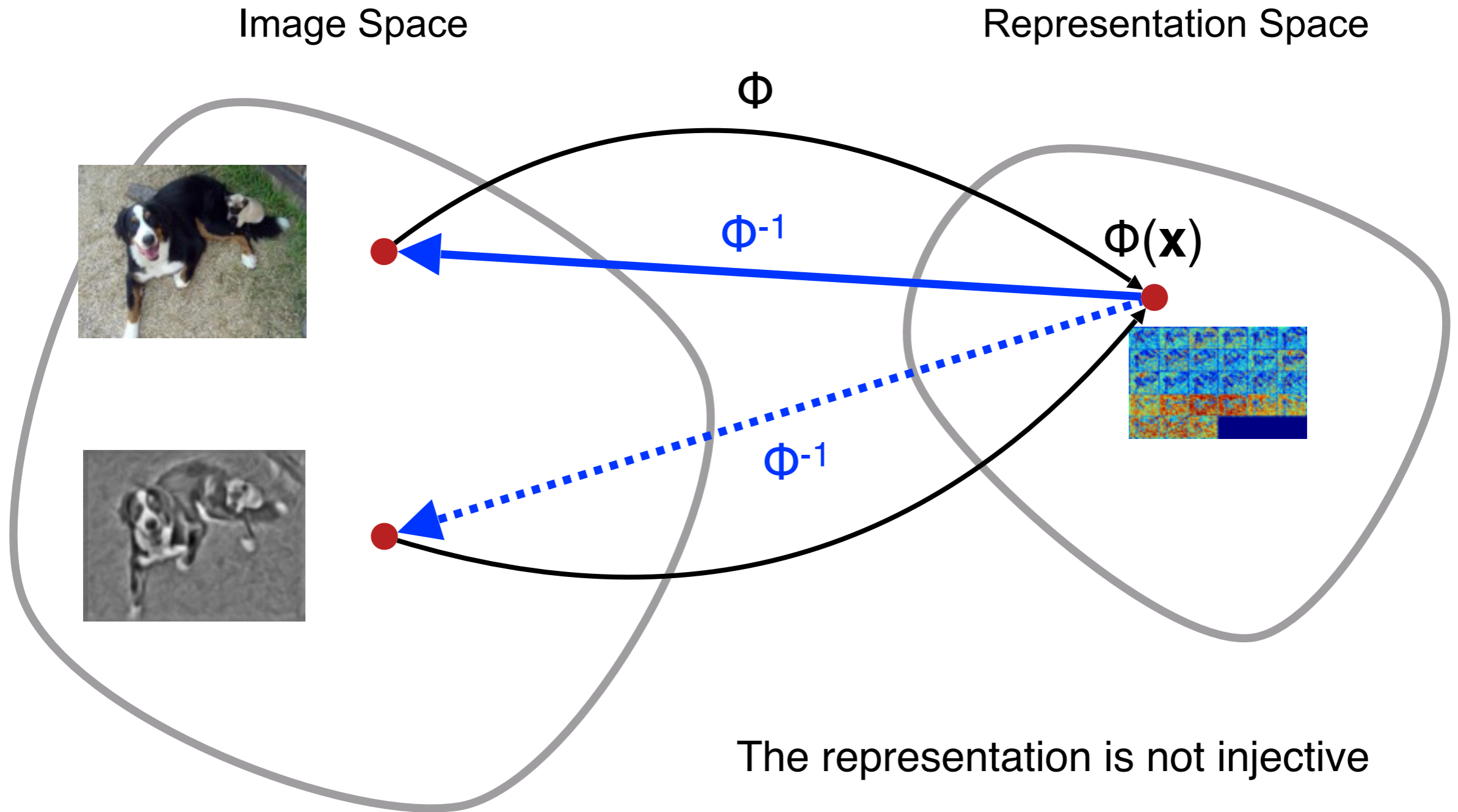
The “**meaning**” of the representation is **unclear**

Visualizing representations

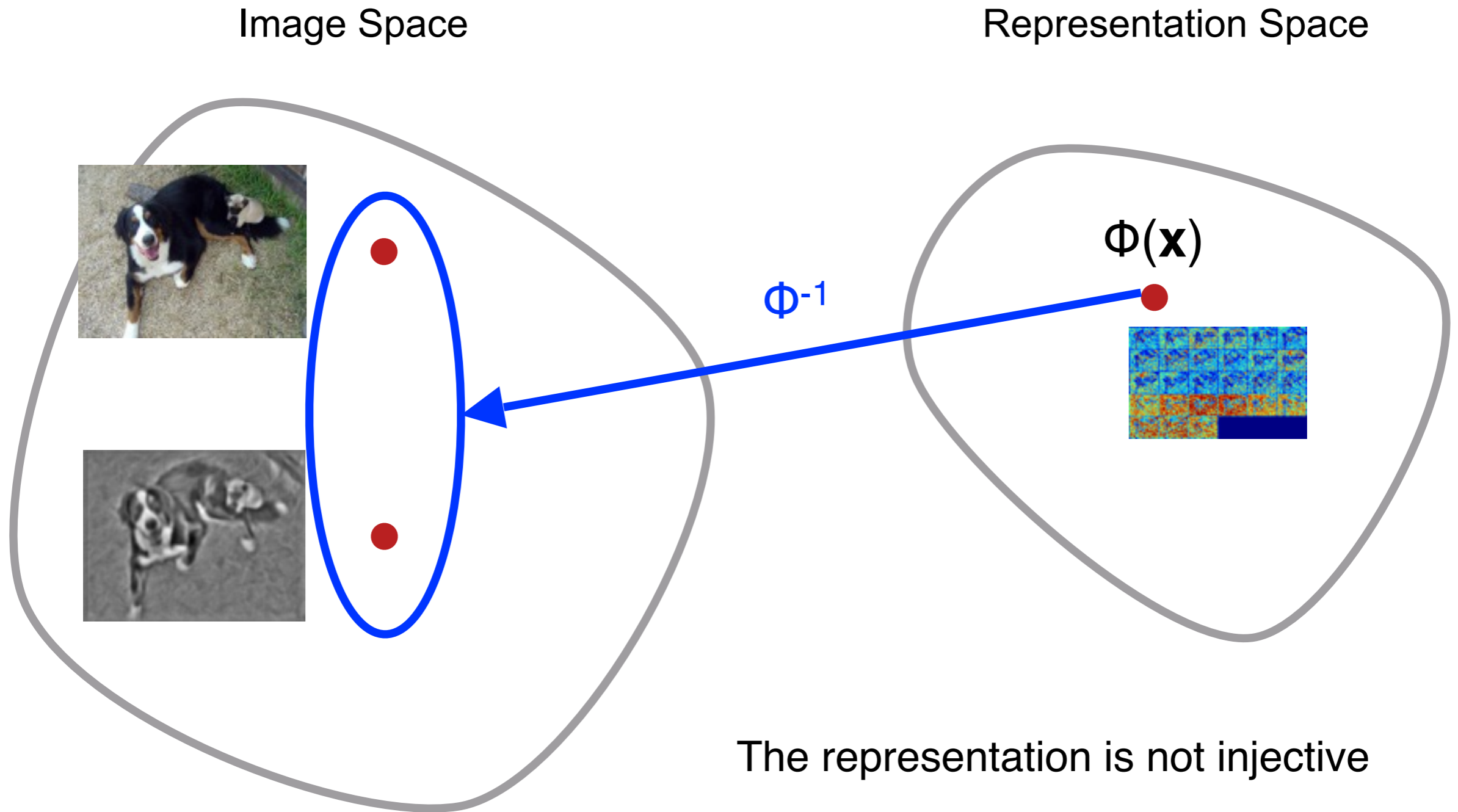
Backpropagation networks and “deconvolution”

Representations: equivalence & transformations

Visualization: Pre-Image



Visualization: Pre-Image



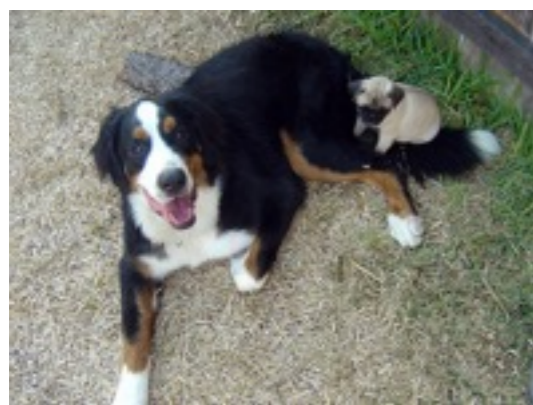
The representation is not injective

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

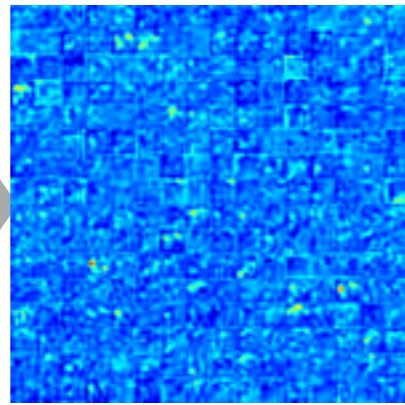
Finding a Pre-Image

A simple yet general and effective method

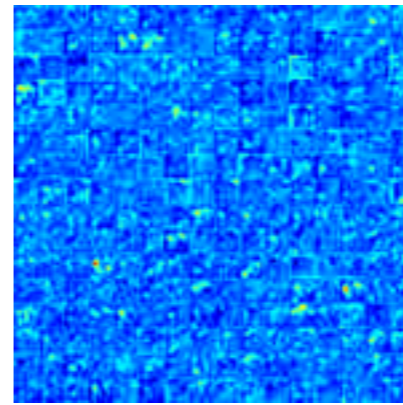
$$\min_{\mathbf{x}} \|\Phi(\mathbf{x}) - \Phi_0\|_2^2$$



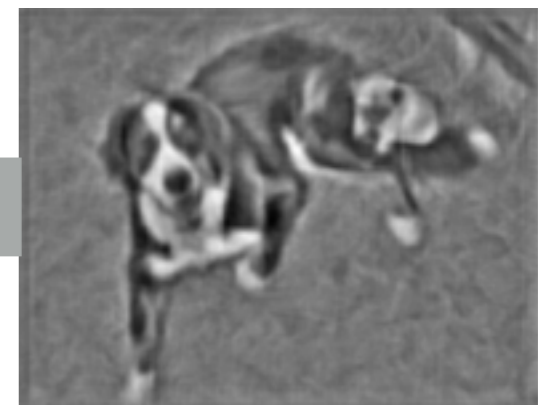
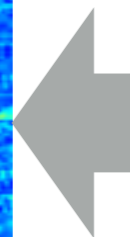
Image



Representation



Reconstruction



Pre-Image

Start from **random noise**

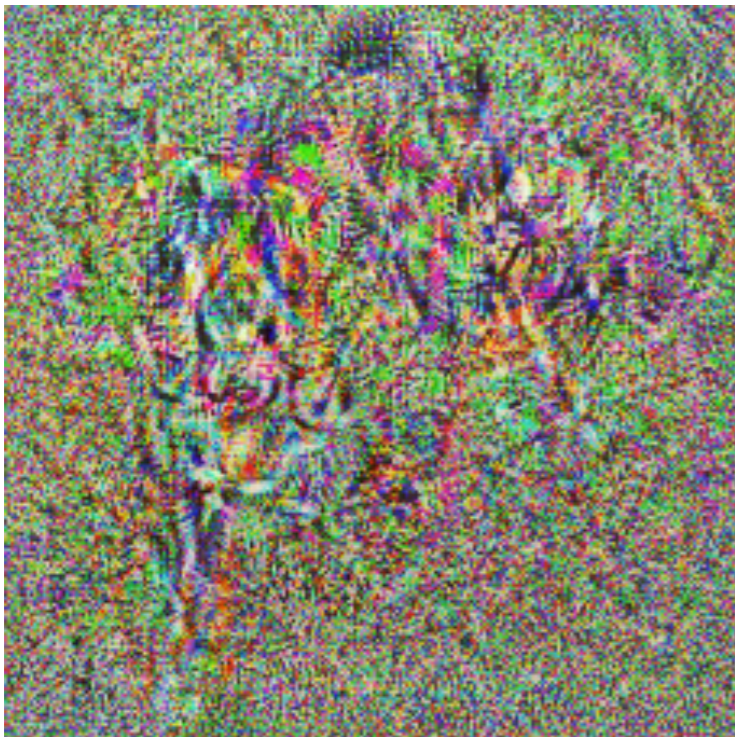
Optimize using stochastic **gradient descent**

Finding a Pre-Image

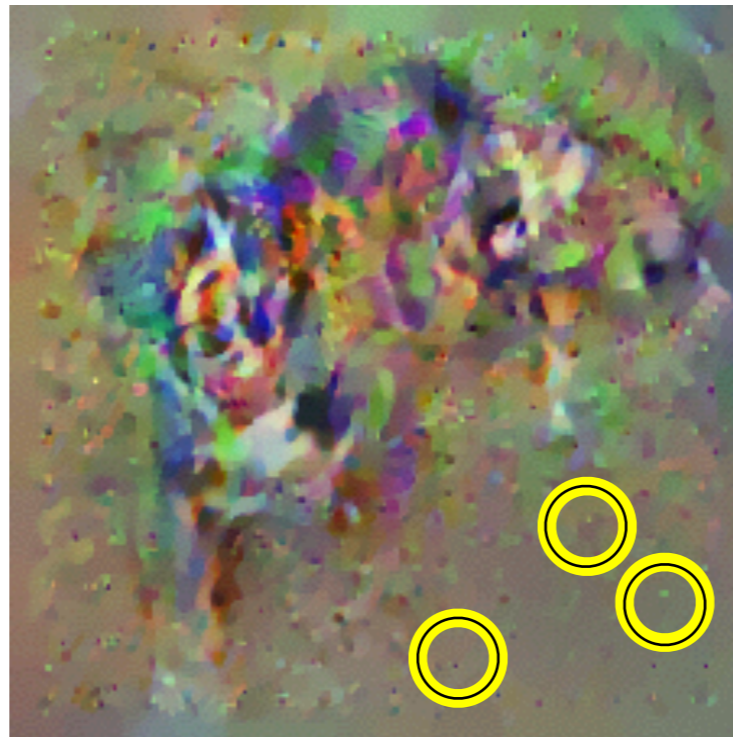
A simple yet general and effective method

$$\min_{\mathbf{x}} \|\Phi(\mathbf{x}) - \Phi_0\|_2^2$$

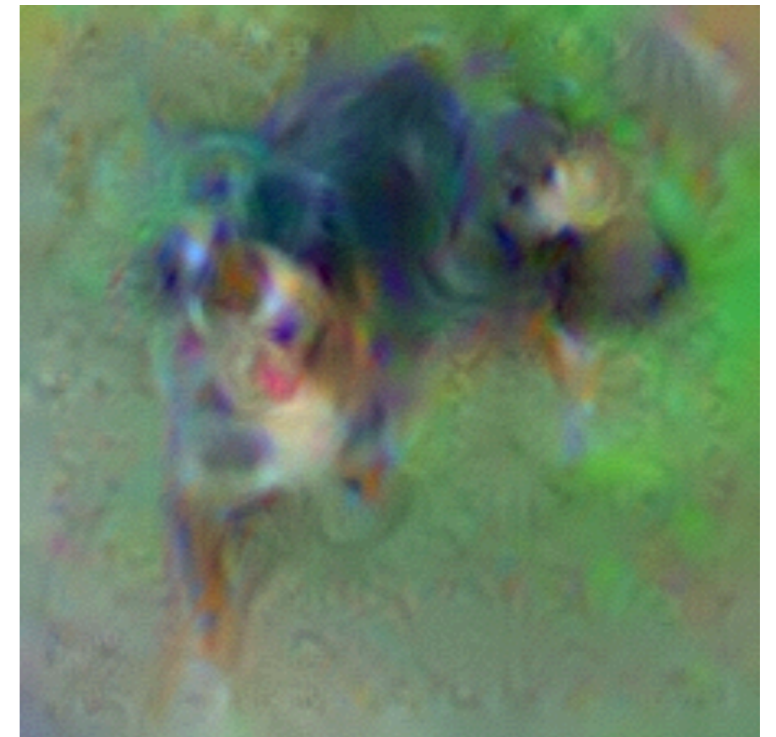
No prior



TV-norm $\beta = 1$



TV-norm $\beta = 2$



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

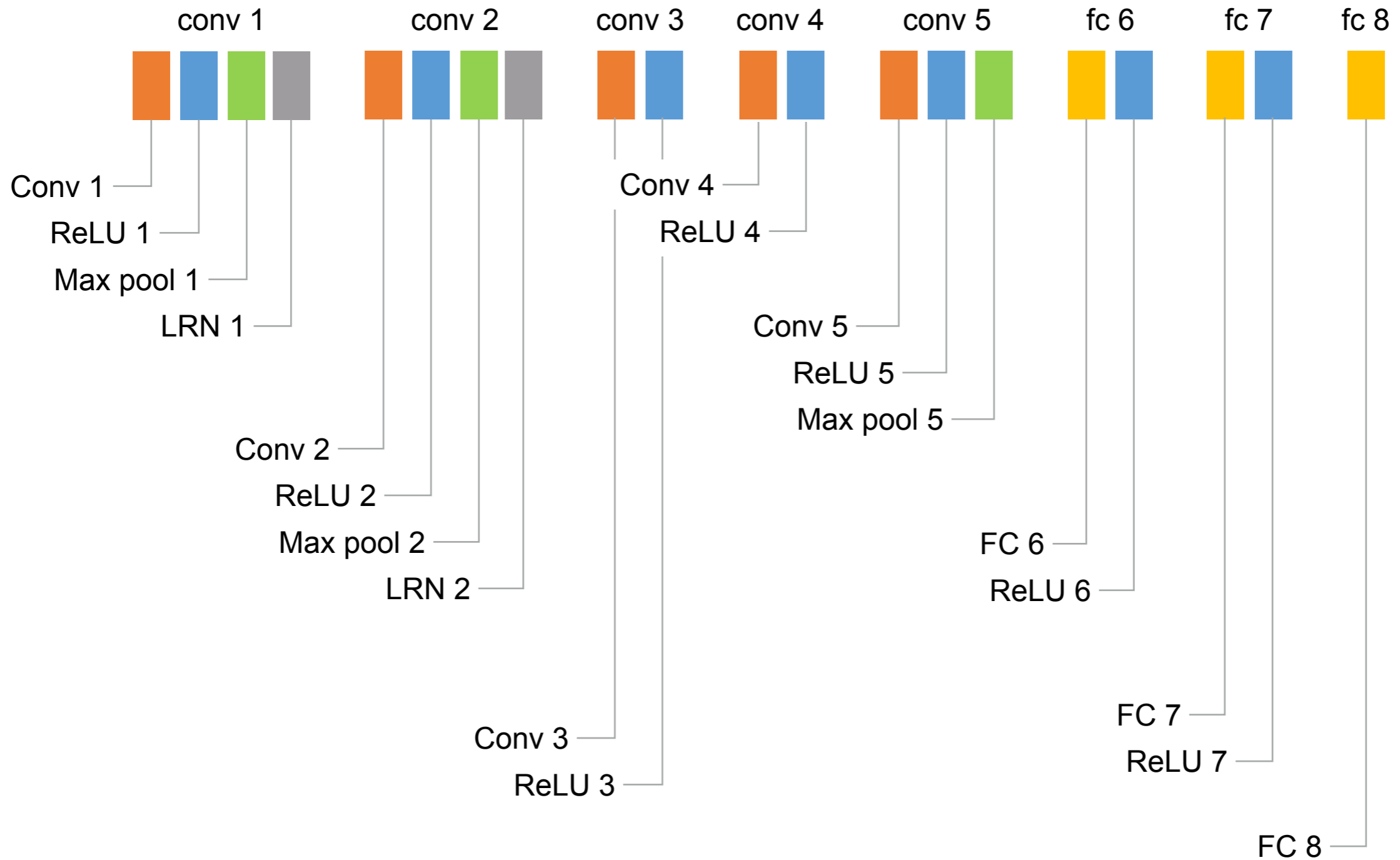
Google’s “inceptionism”

Mordvintsev et al. 2015

Style synthesis and transfer

Gatys et al. 2015

Inversion



AlexNet
[Krizhevsky et al. 2012]

Inversion



Original
Image



Inversion



Original Image



Inversion



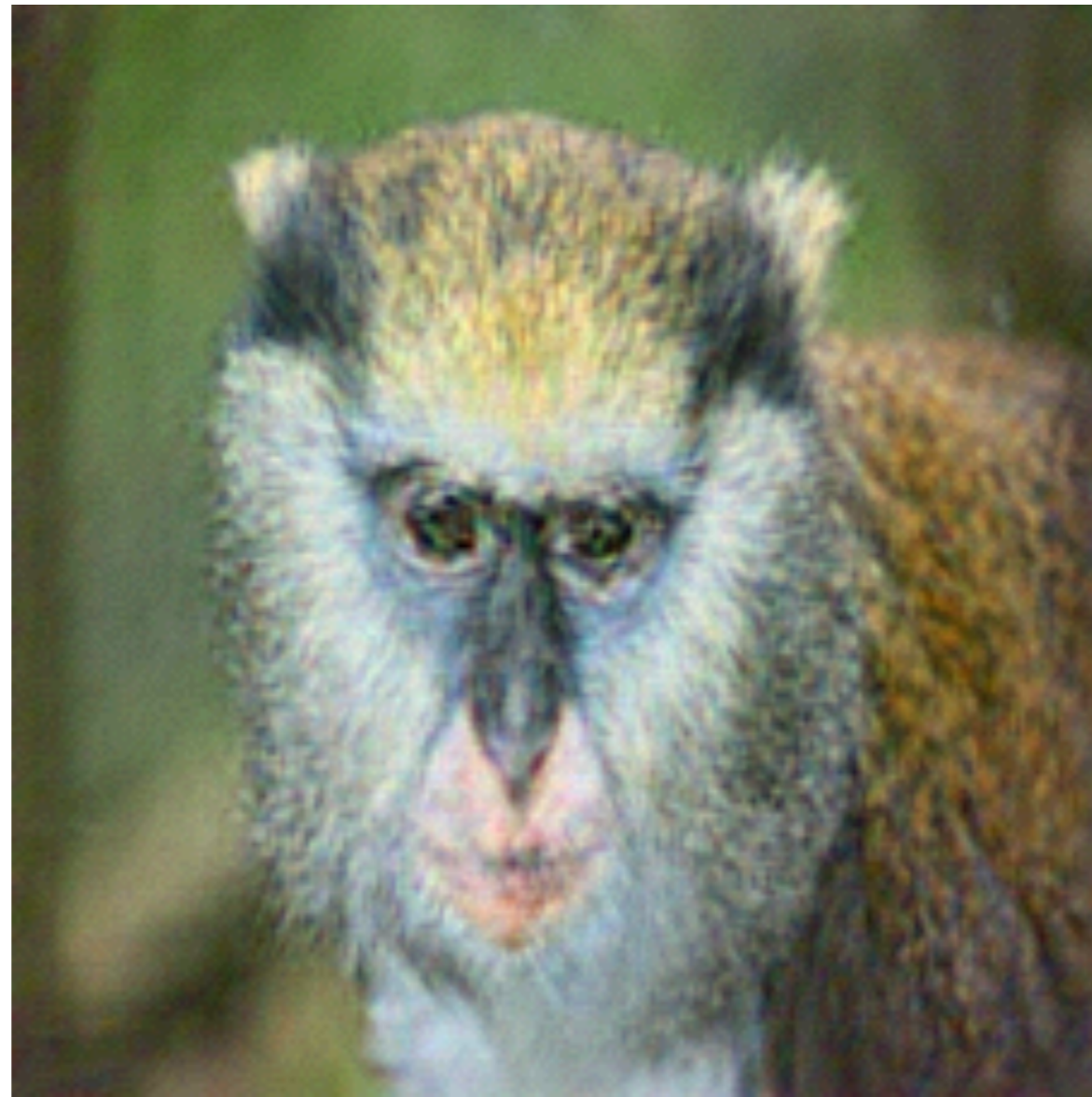
Original Image



Inversion



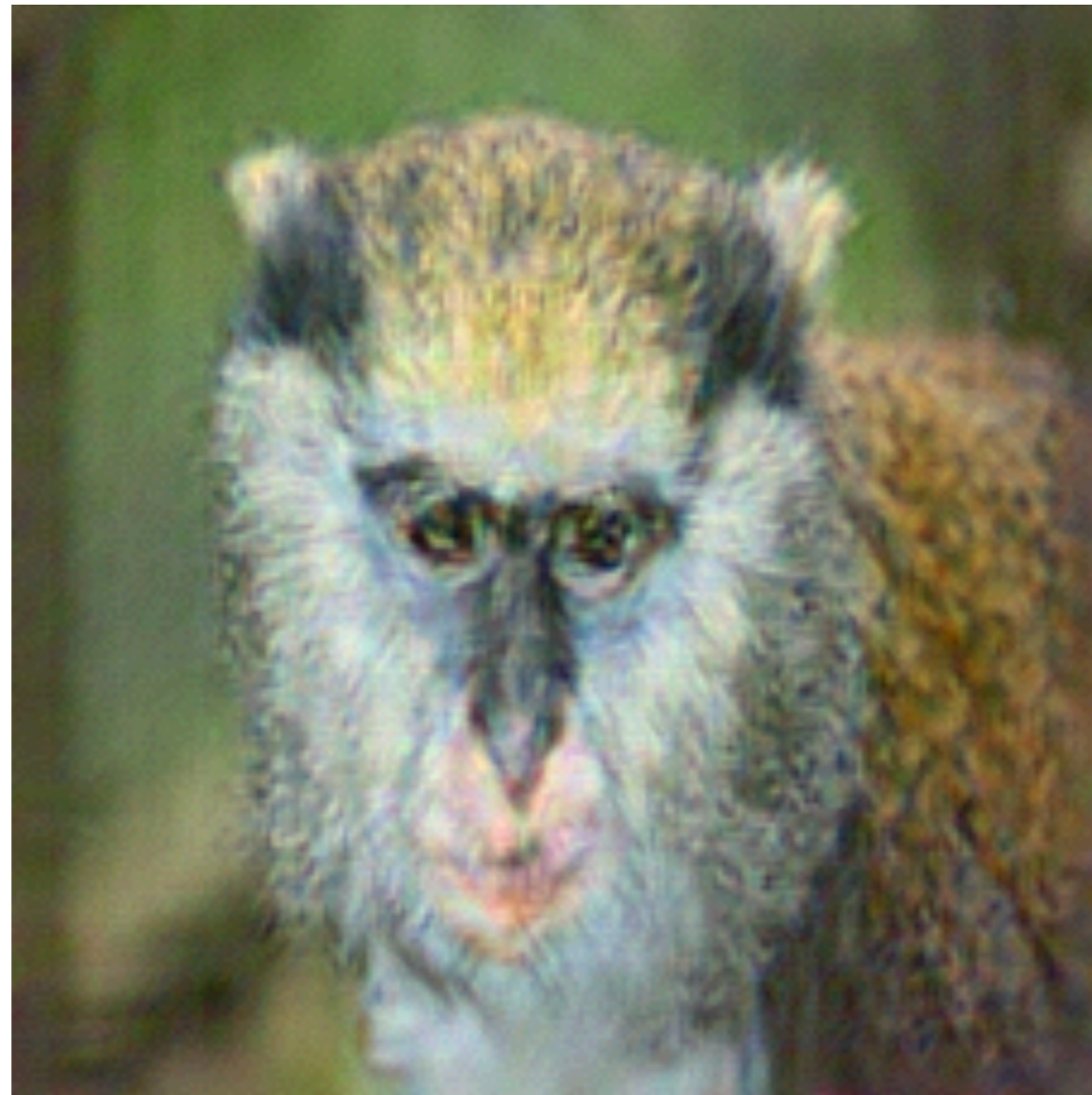
Original Image



Inversion



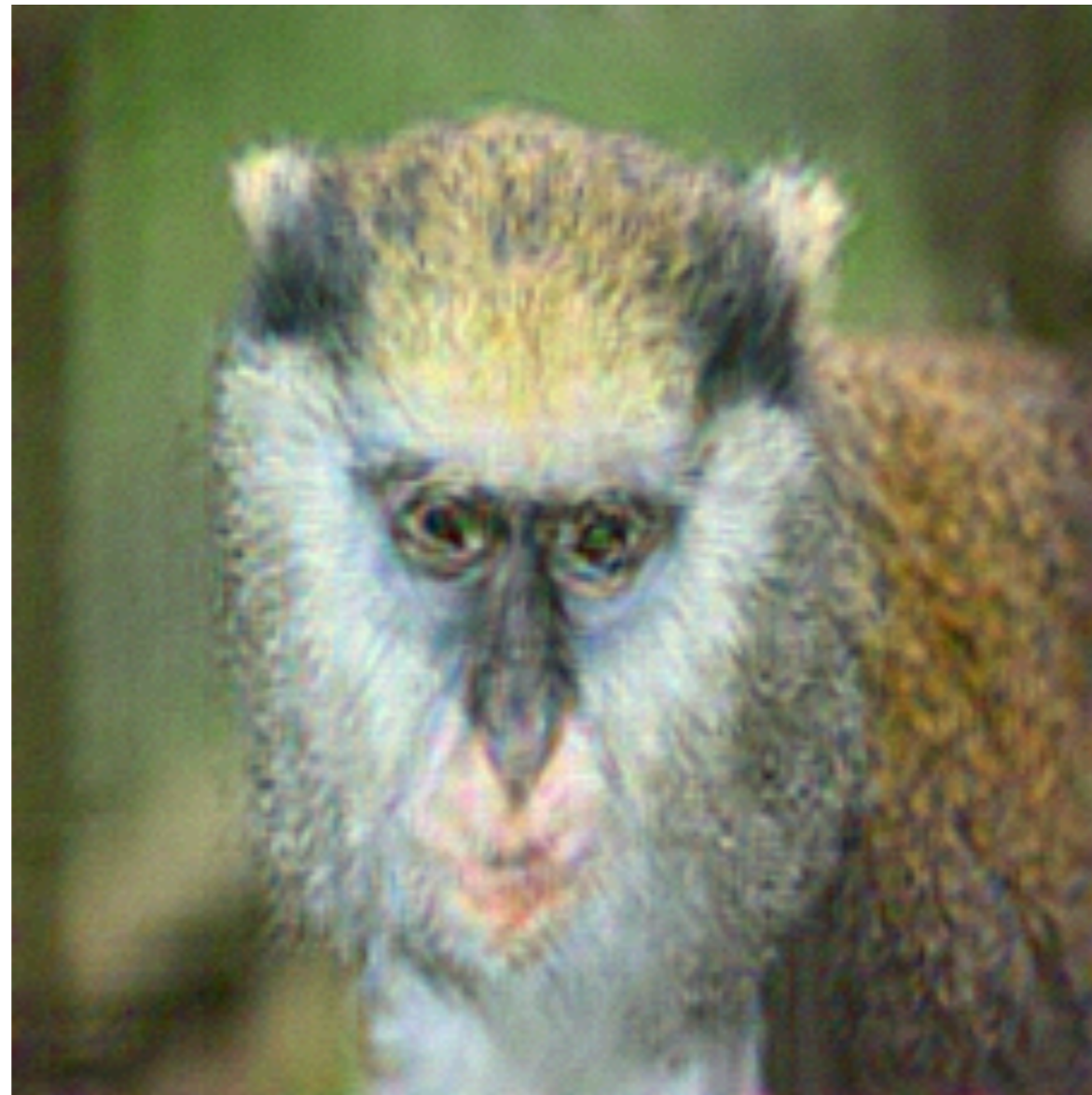
Original Image



Inversion



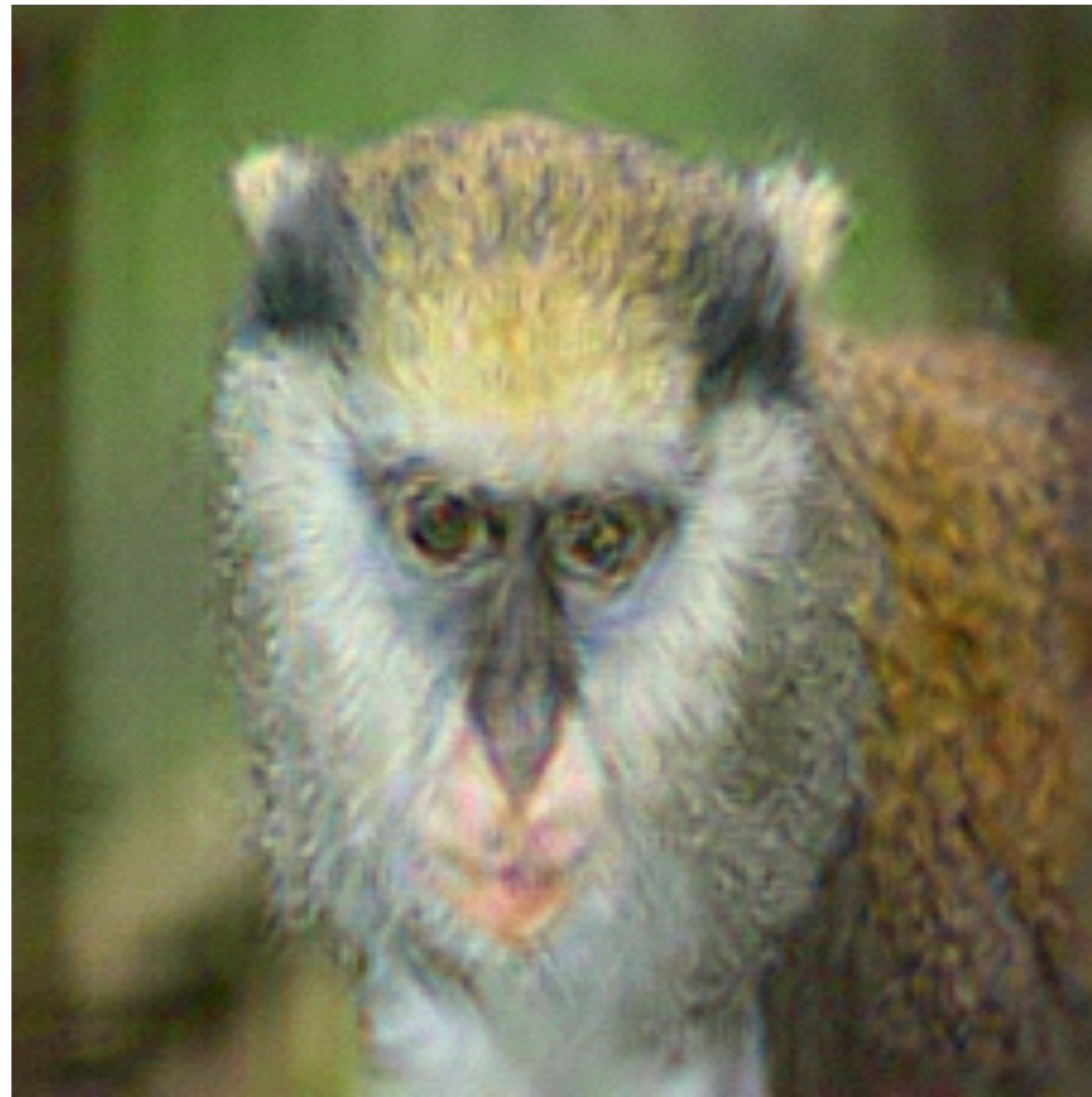
Original Image



Inversion



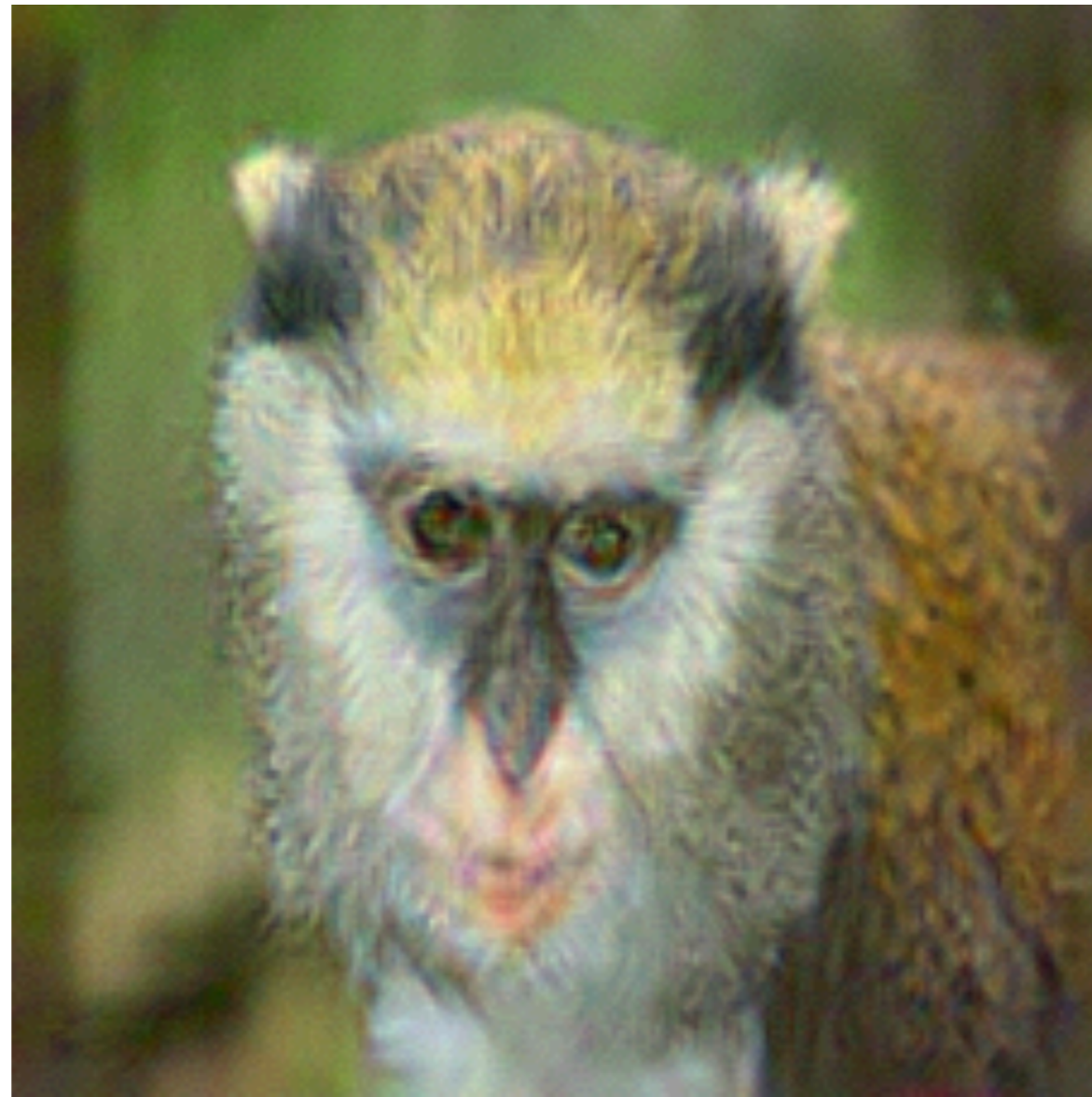
Original Image



Inversion



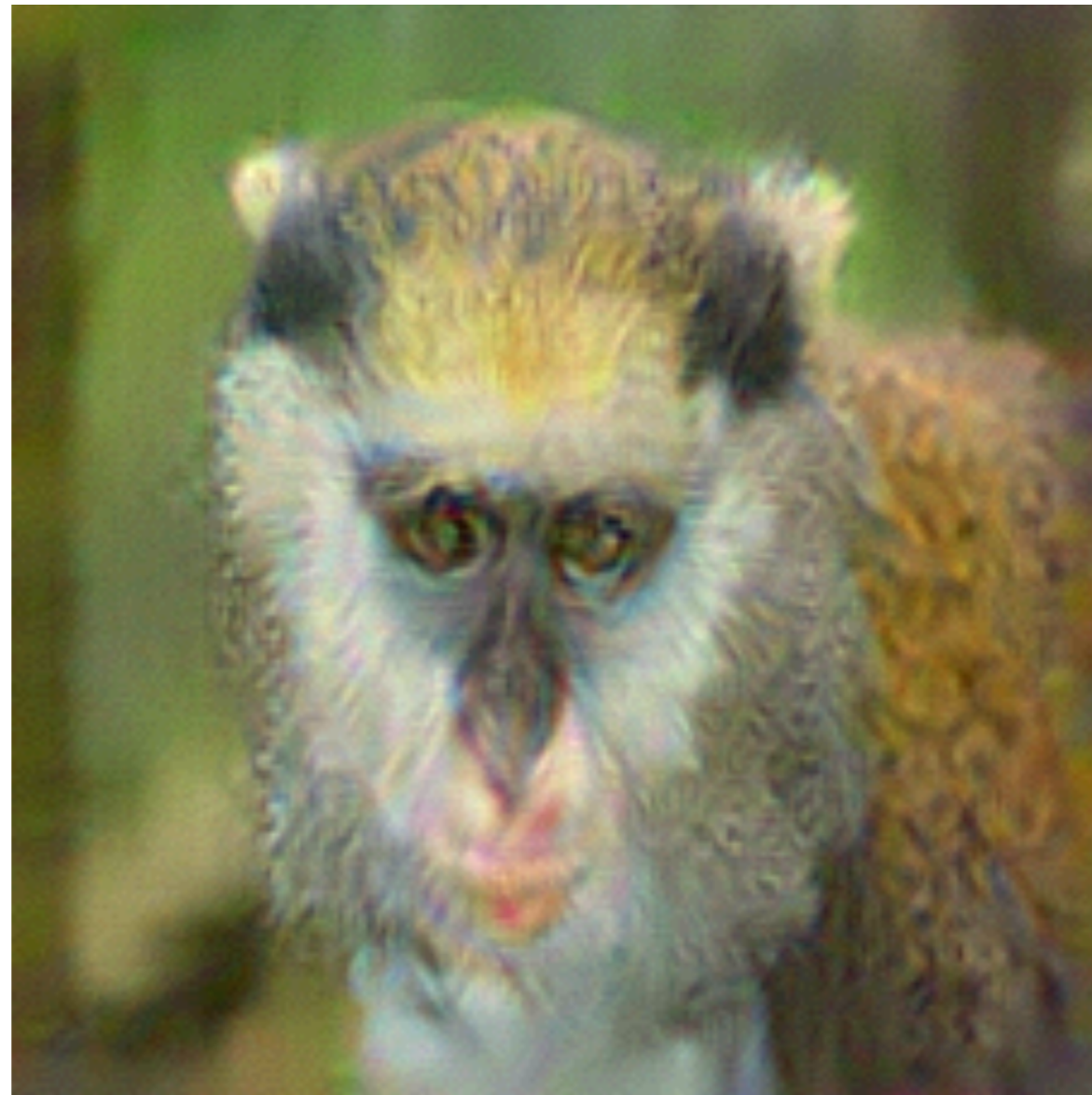
Original Image



Inversion



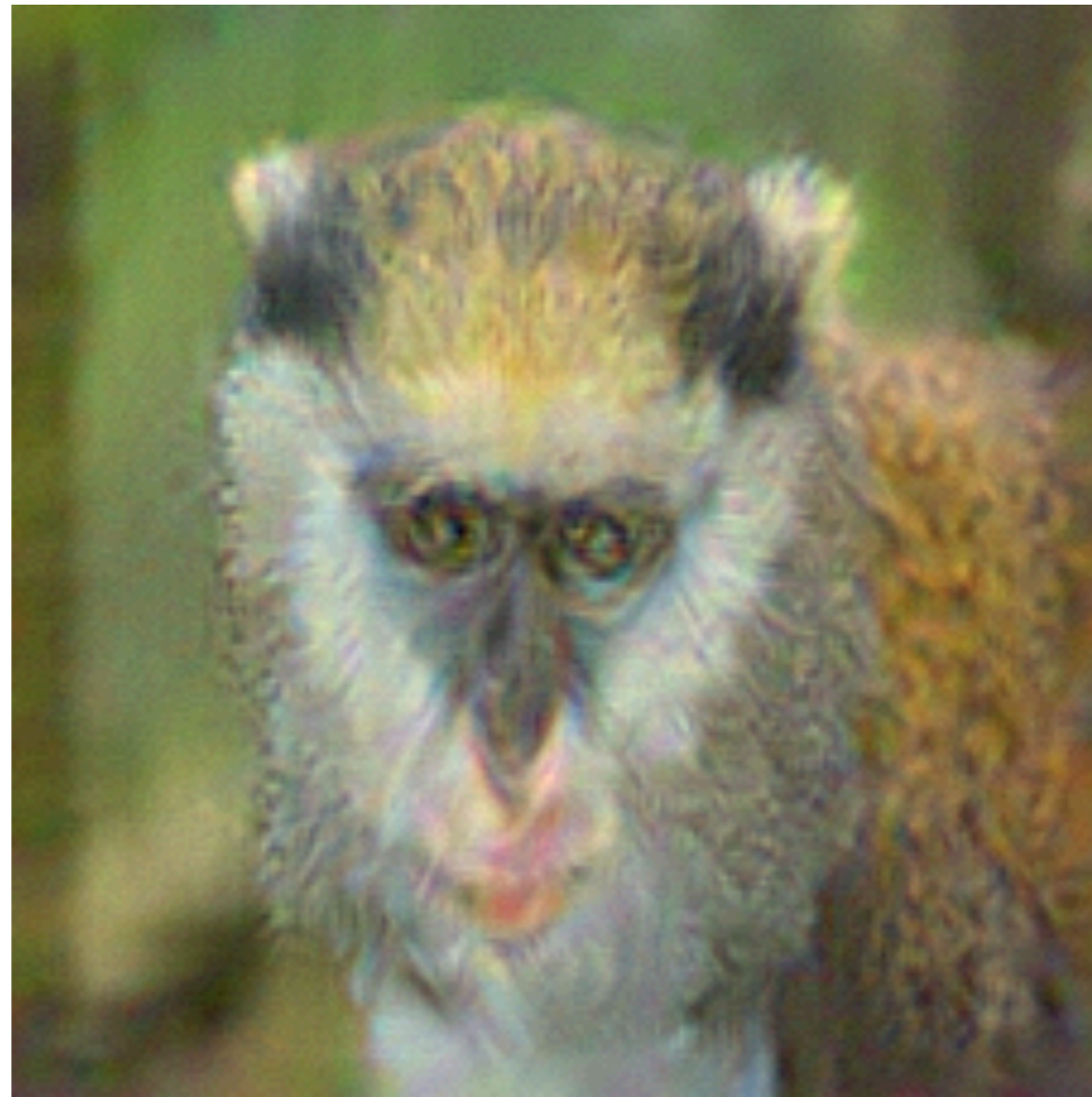
Original Image



Inversion



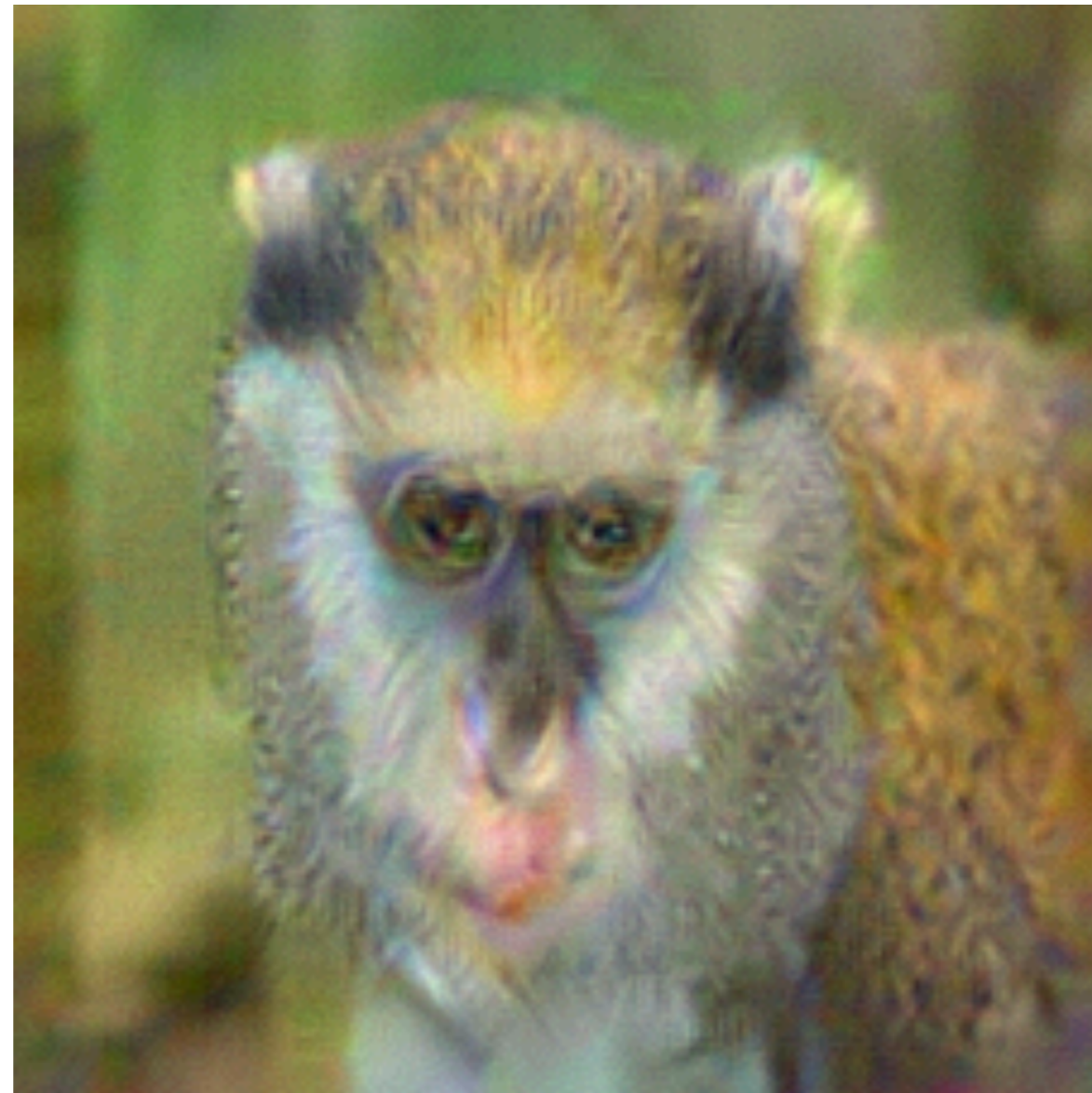
Original Image



Inversion



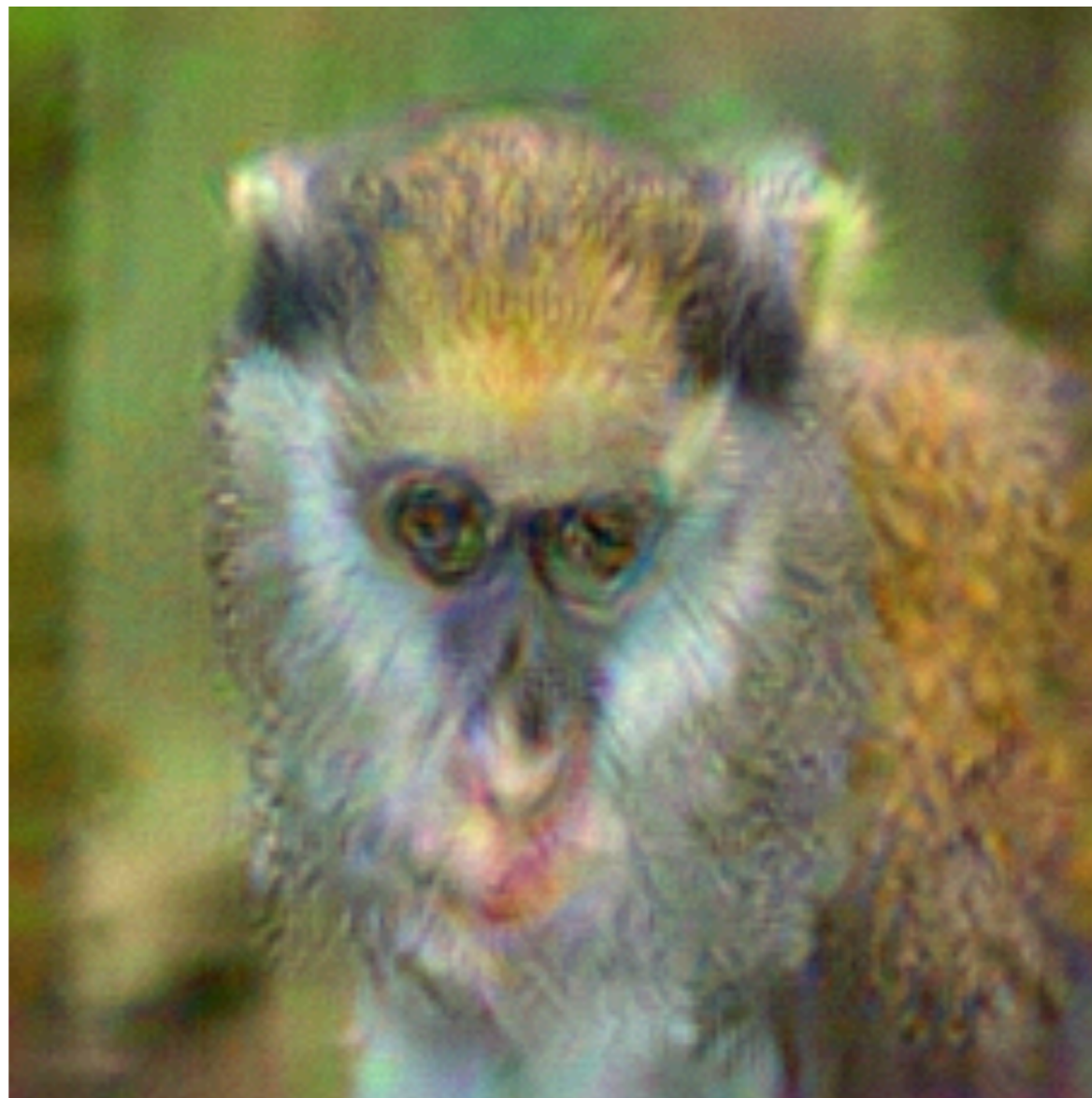
Original
Image



Inversion



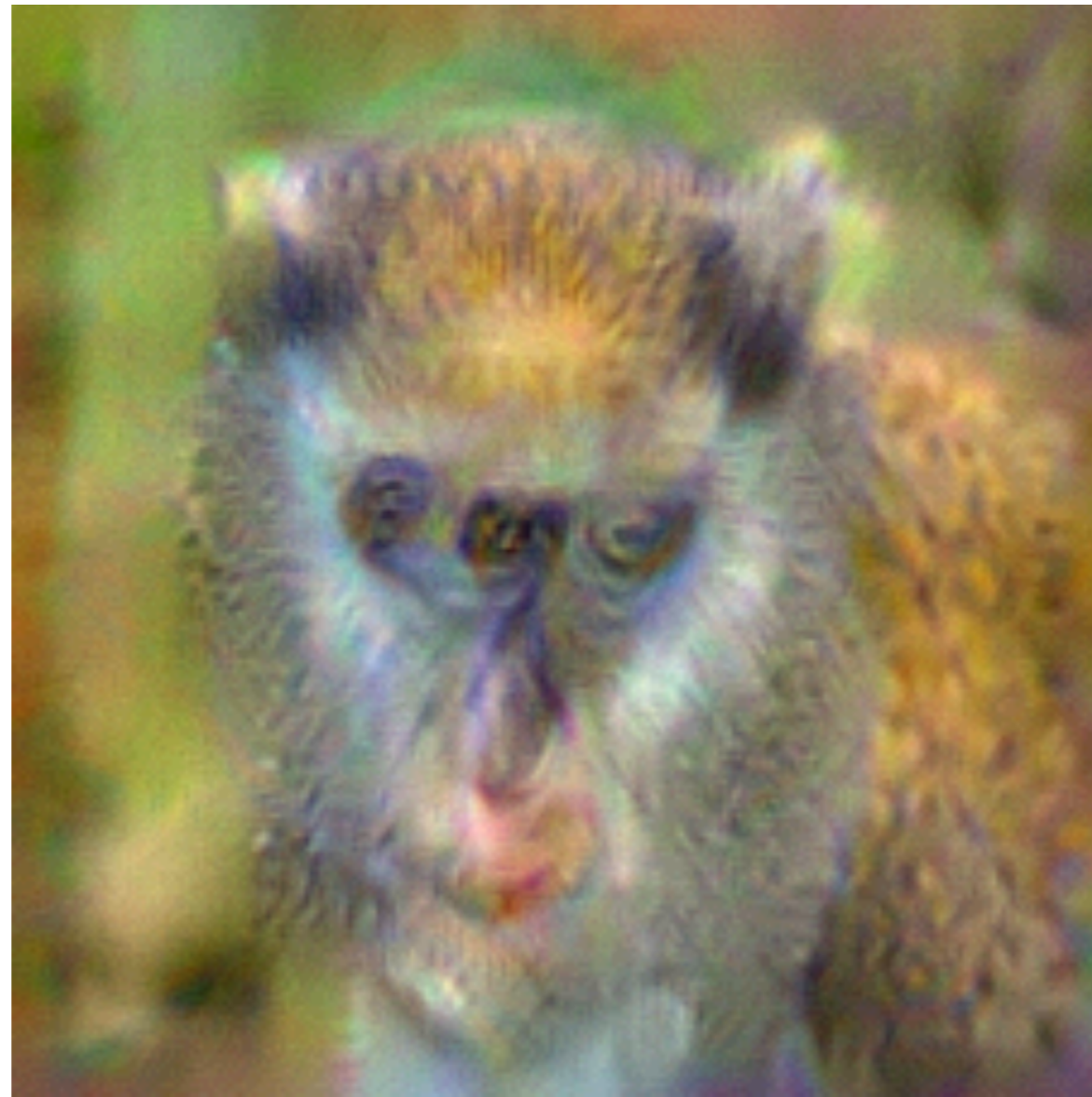
Original Image



Inversion



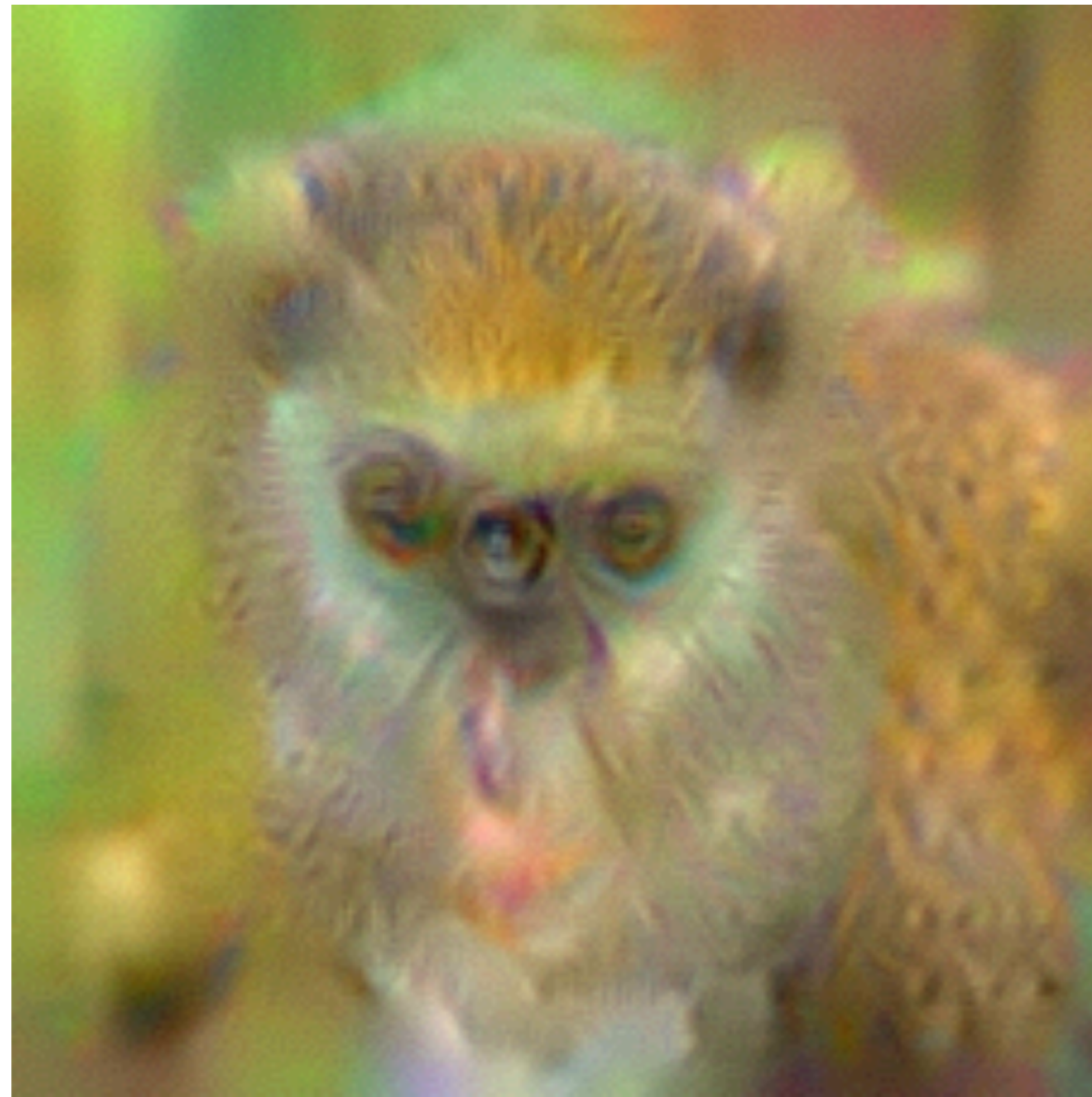
Original Image



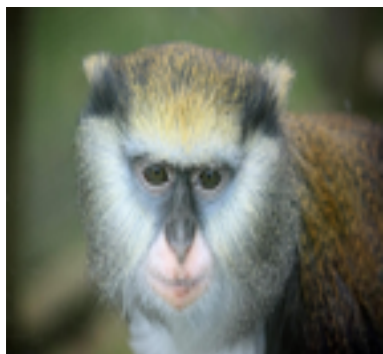
Inversion



Original Image



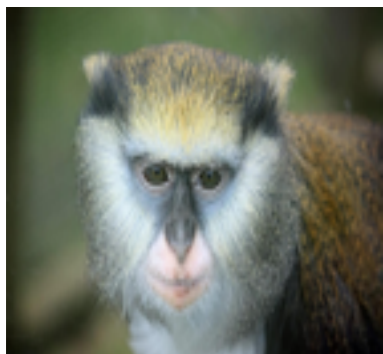
Inversion



Original
Image



Inversion



Original
Image



Inversion



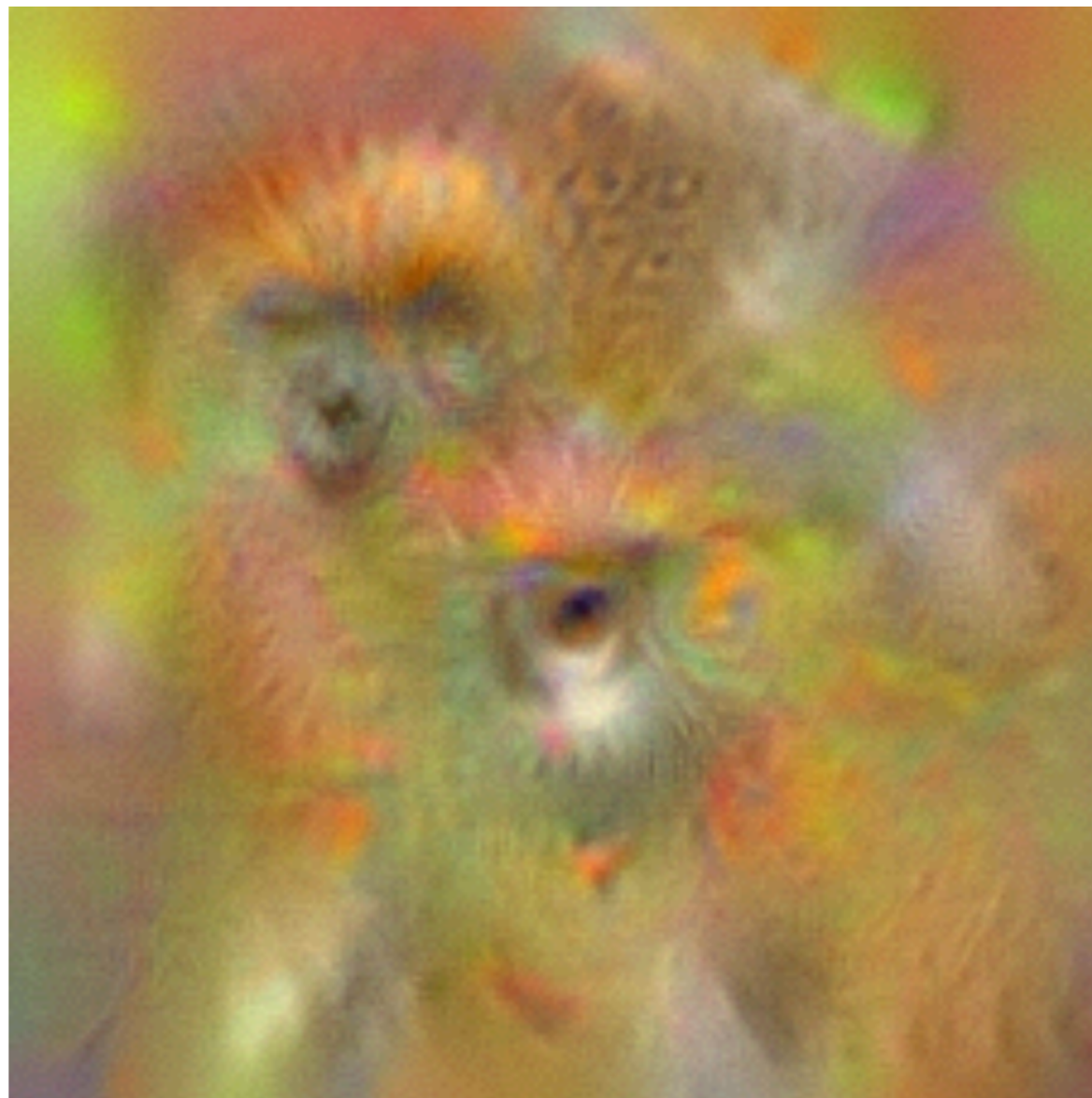
Original Image



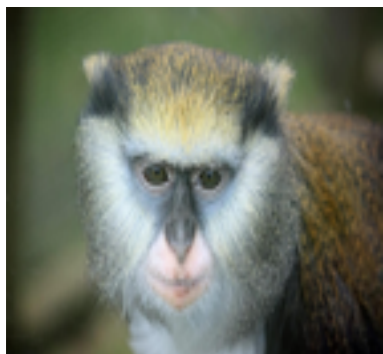
Inversion



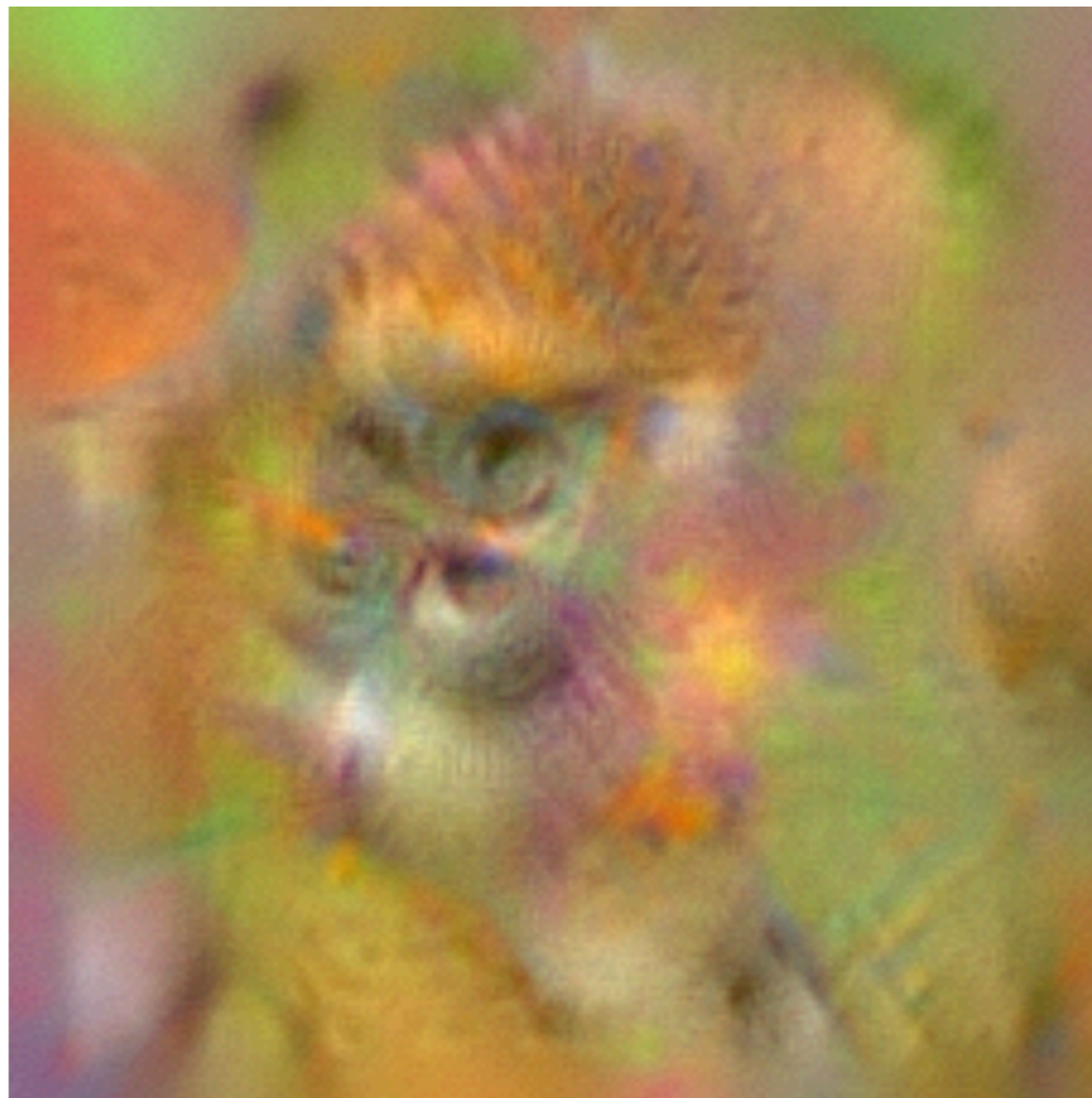
Original
Image



Inversion



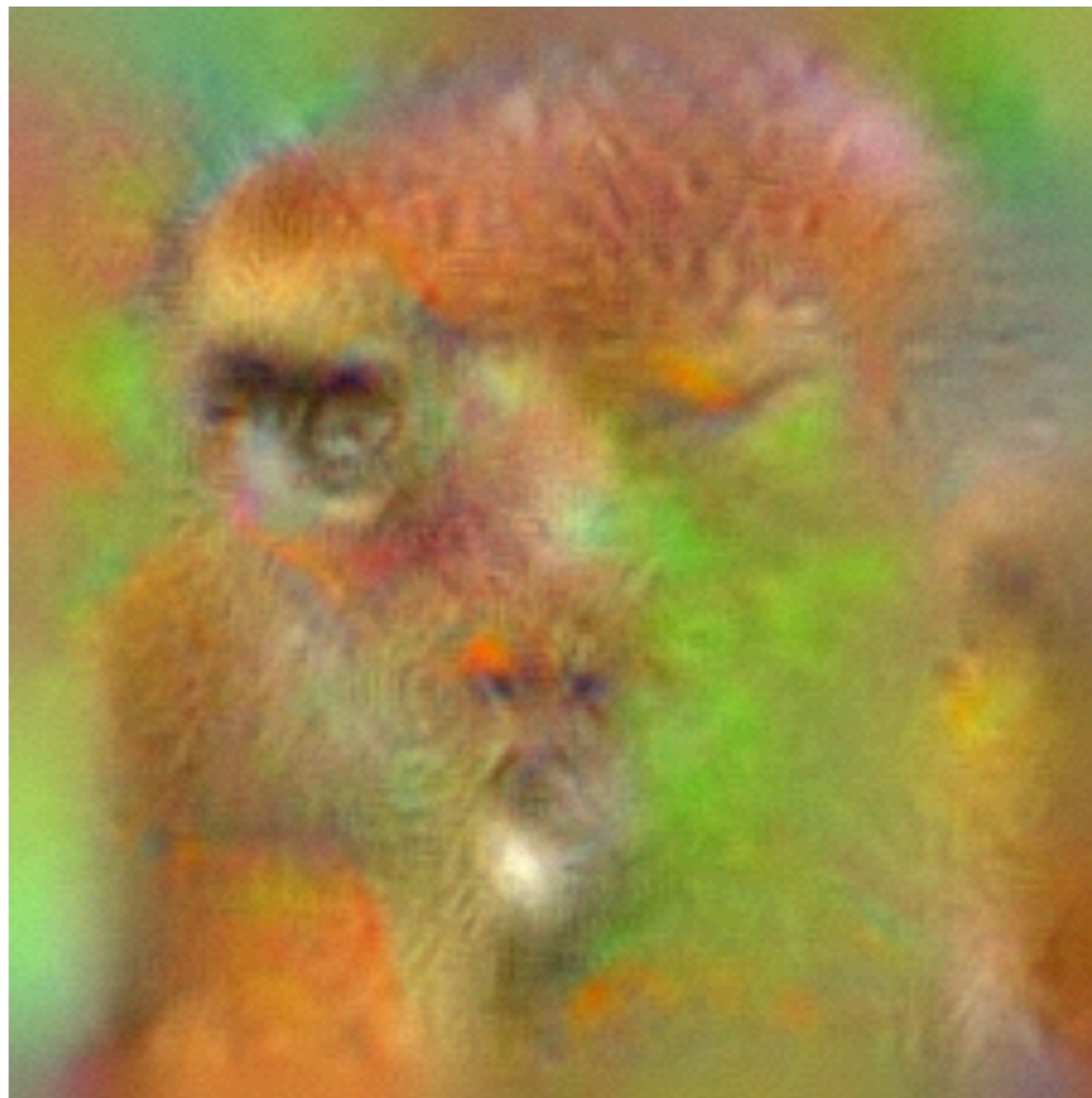
Original
Image



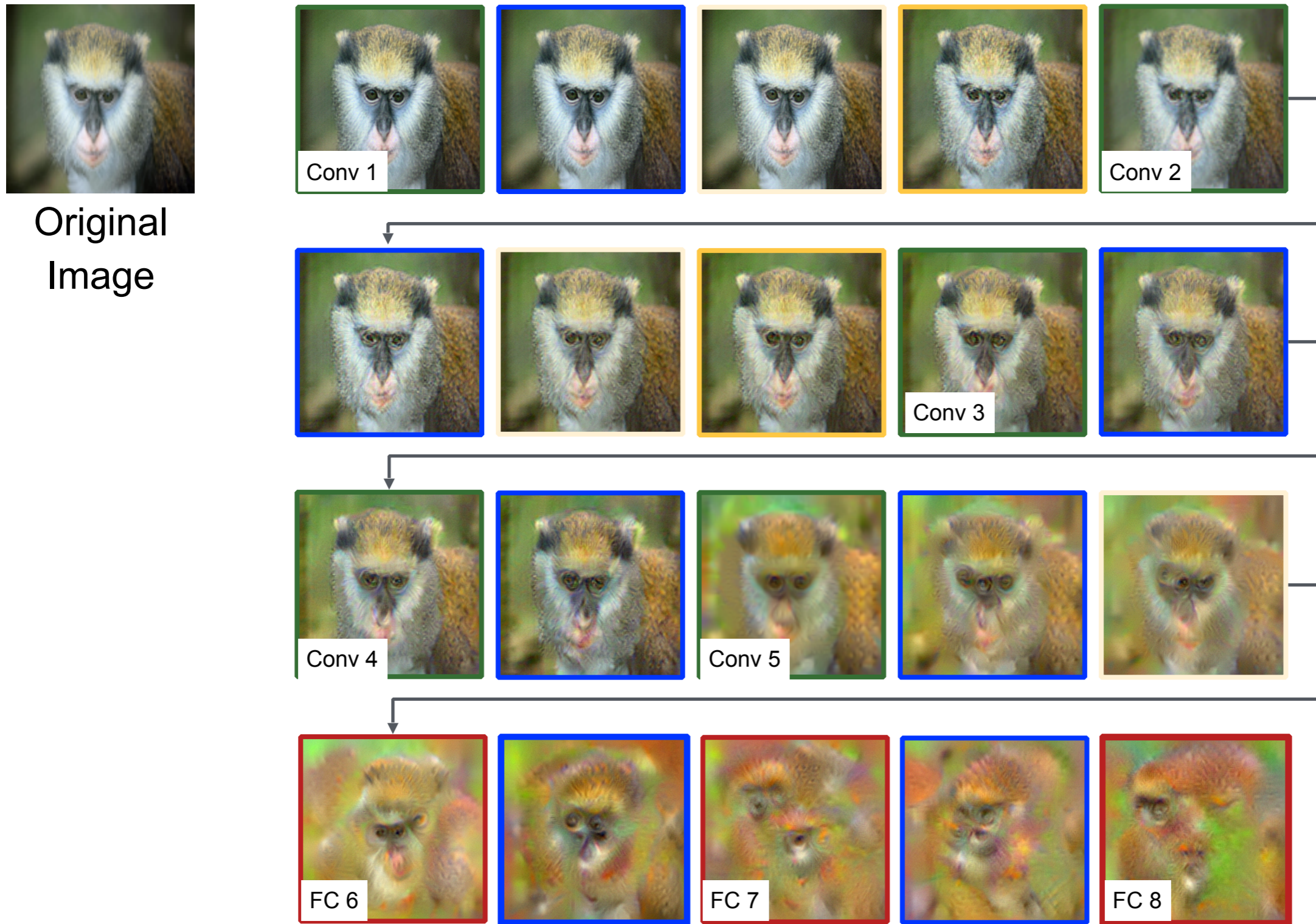
Inversion



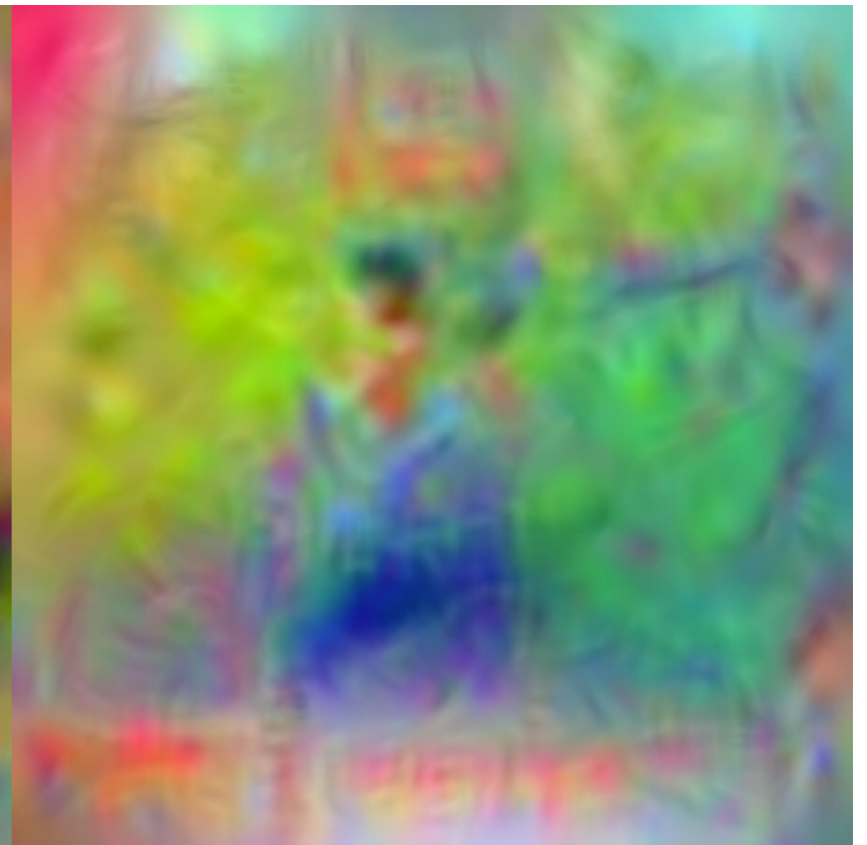
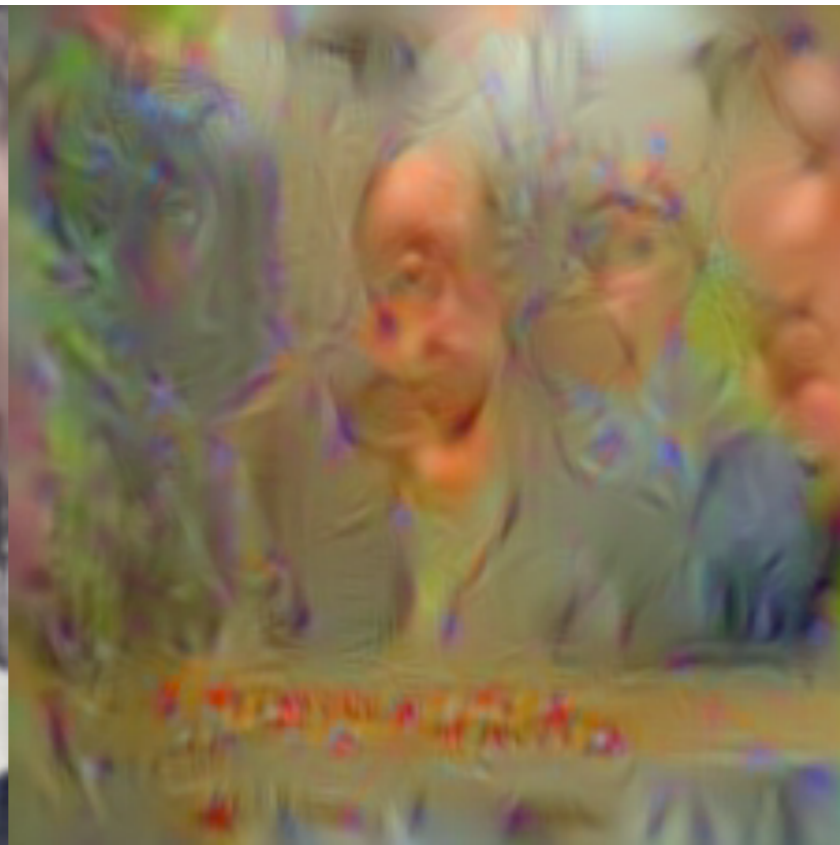
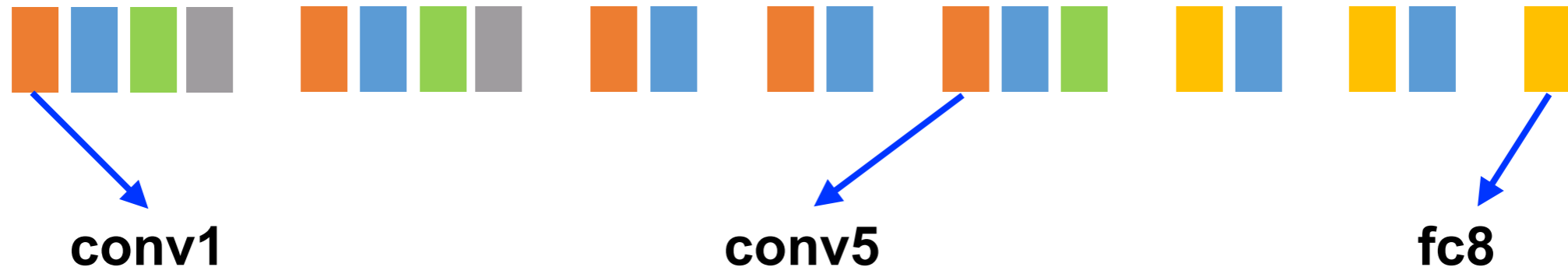
Original
Image



Inverting a Deep CNN



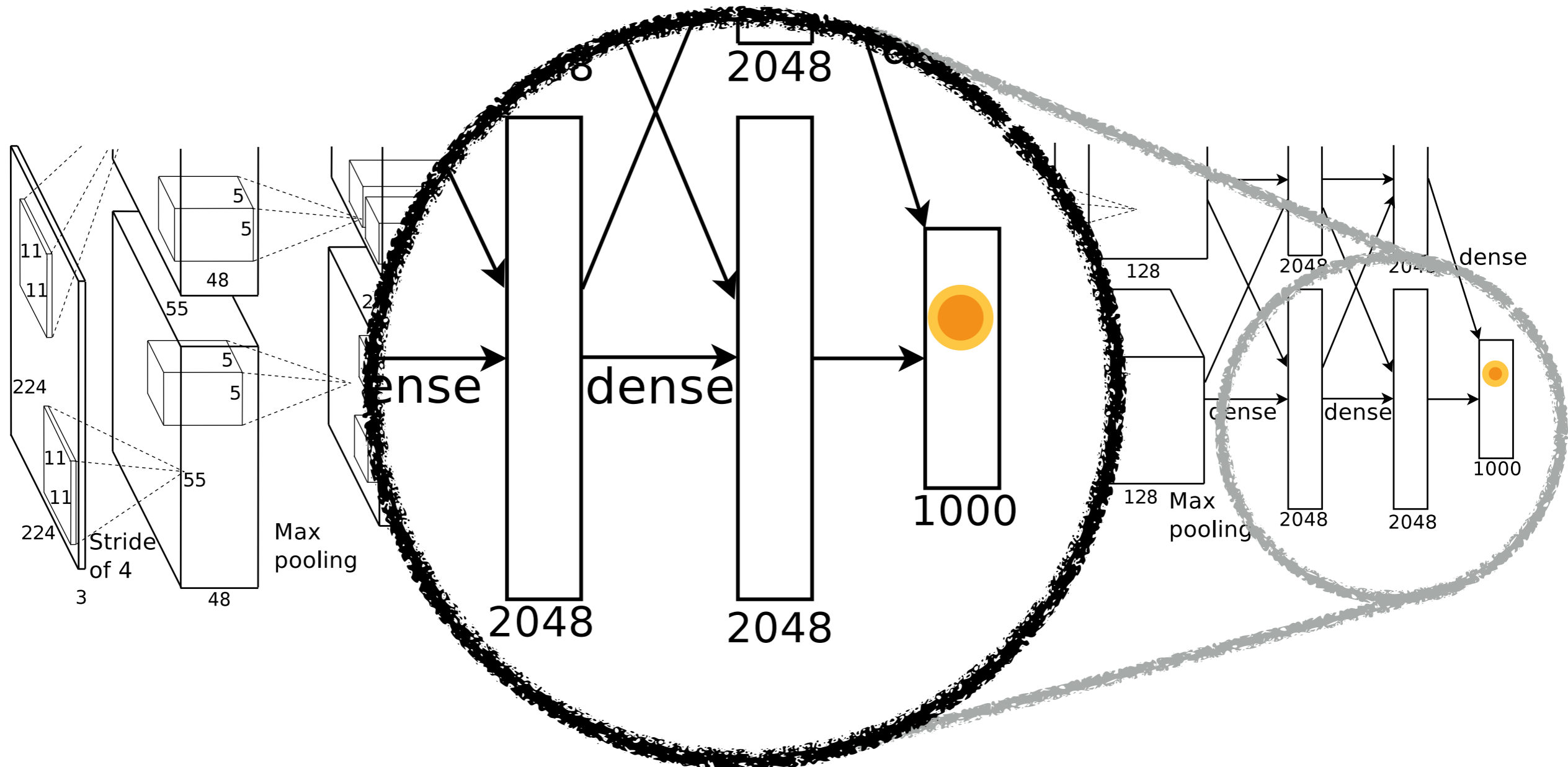
CNNs = visual codes?

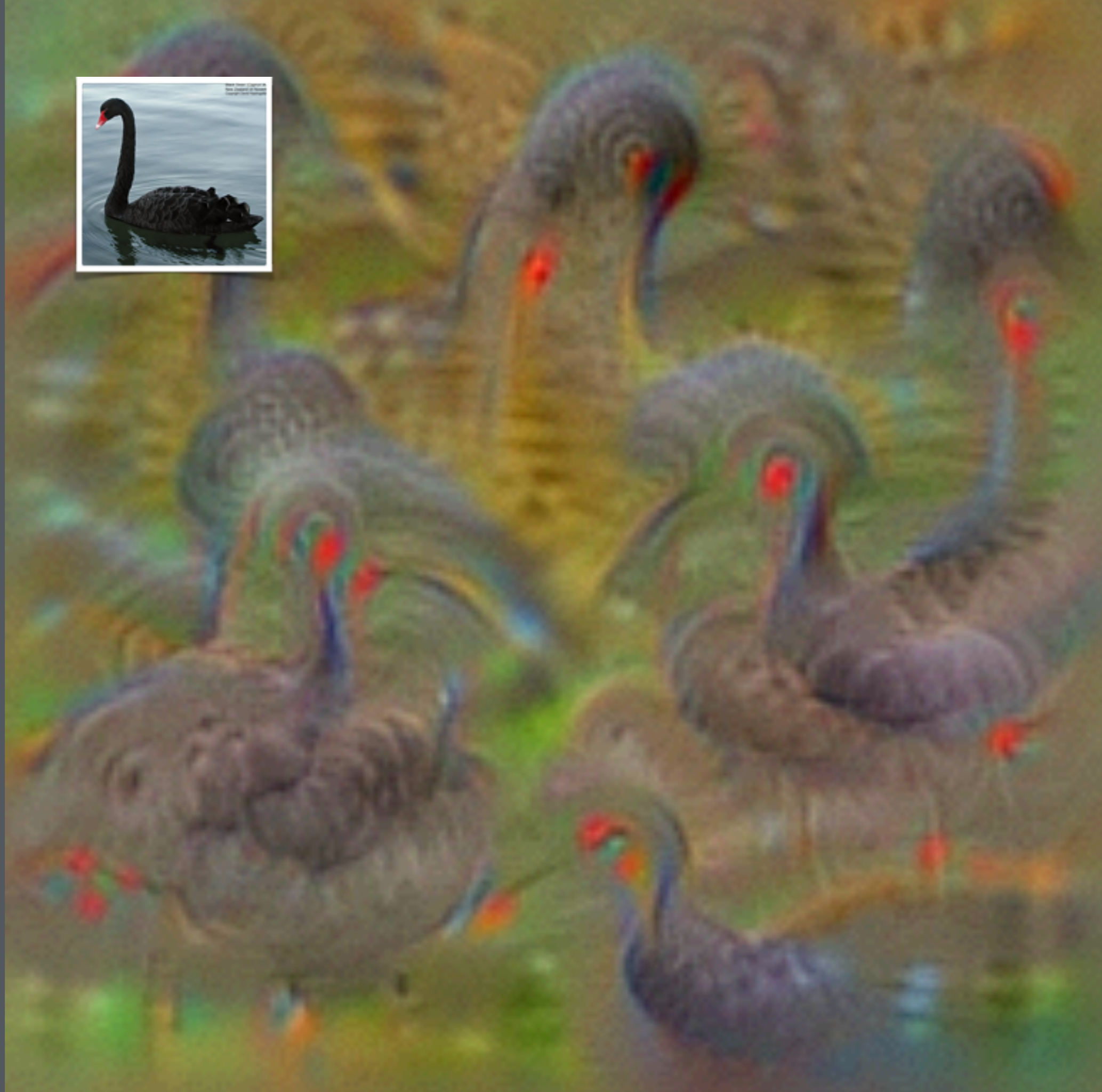
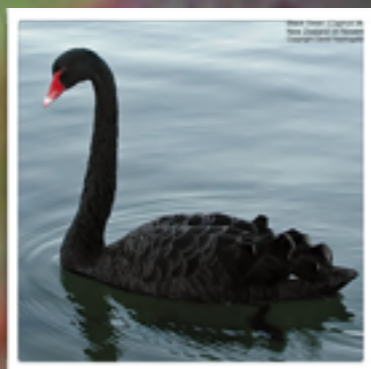


Activation Maximization

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

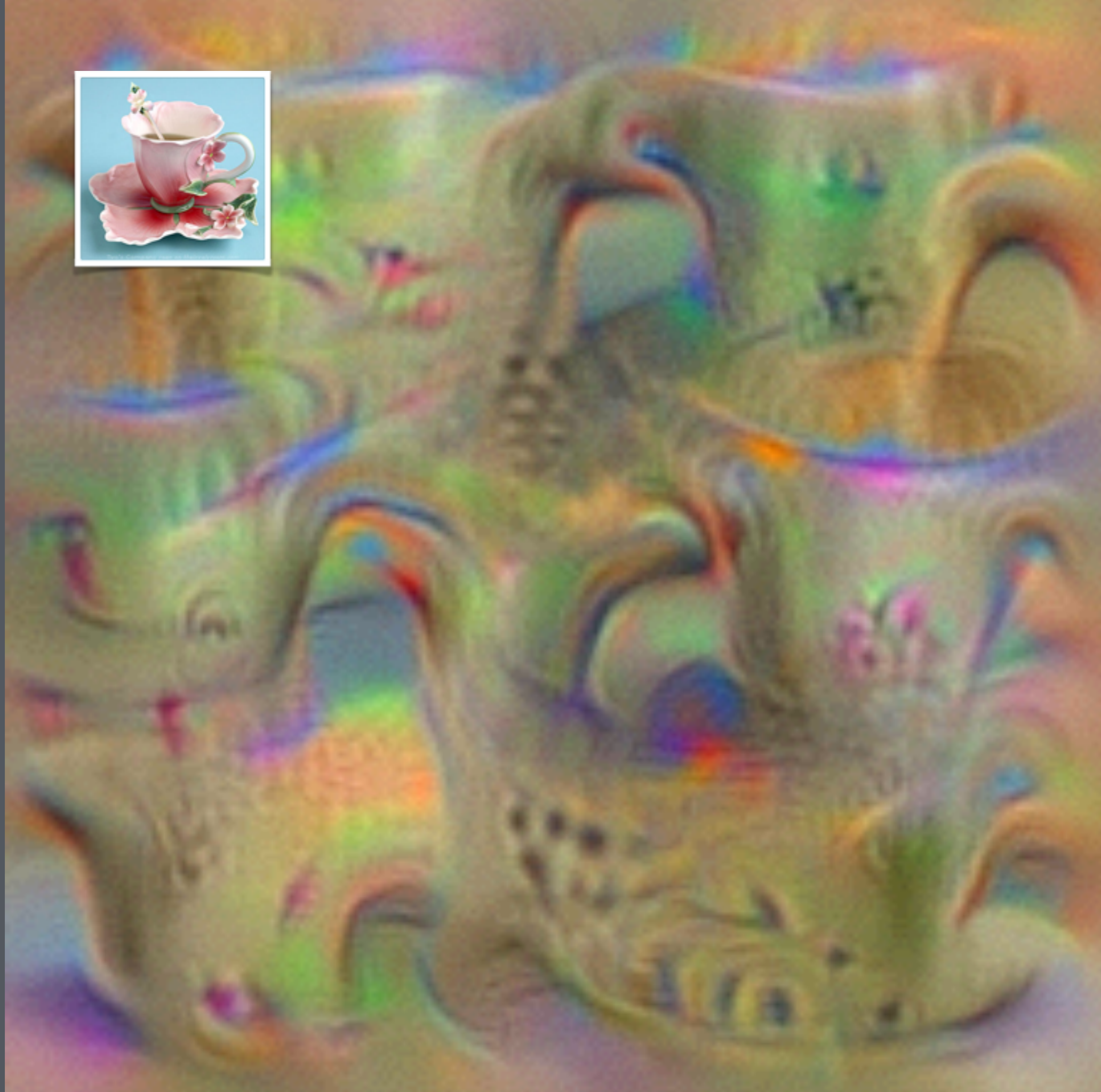
$$\min_{\mathbf{x}} -\langle \mathbf{e}_k, \Phi(\mathbf{x}) \rangle + R_{TV}(\mathbf{x}) + R_{\alpha}(\mathbf{x})$$

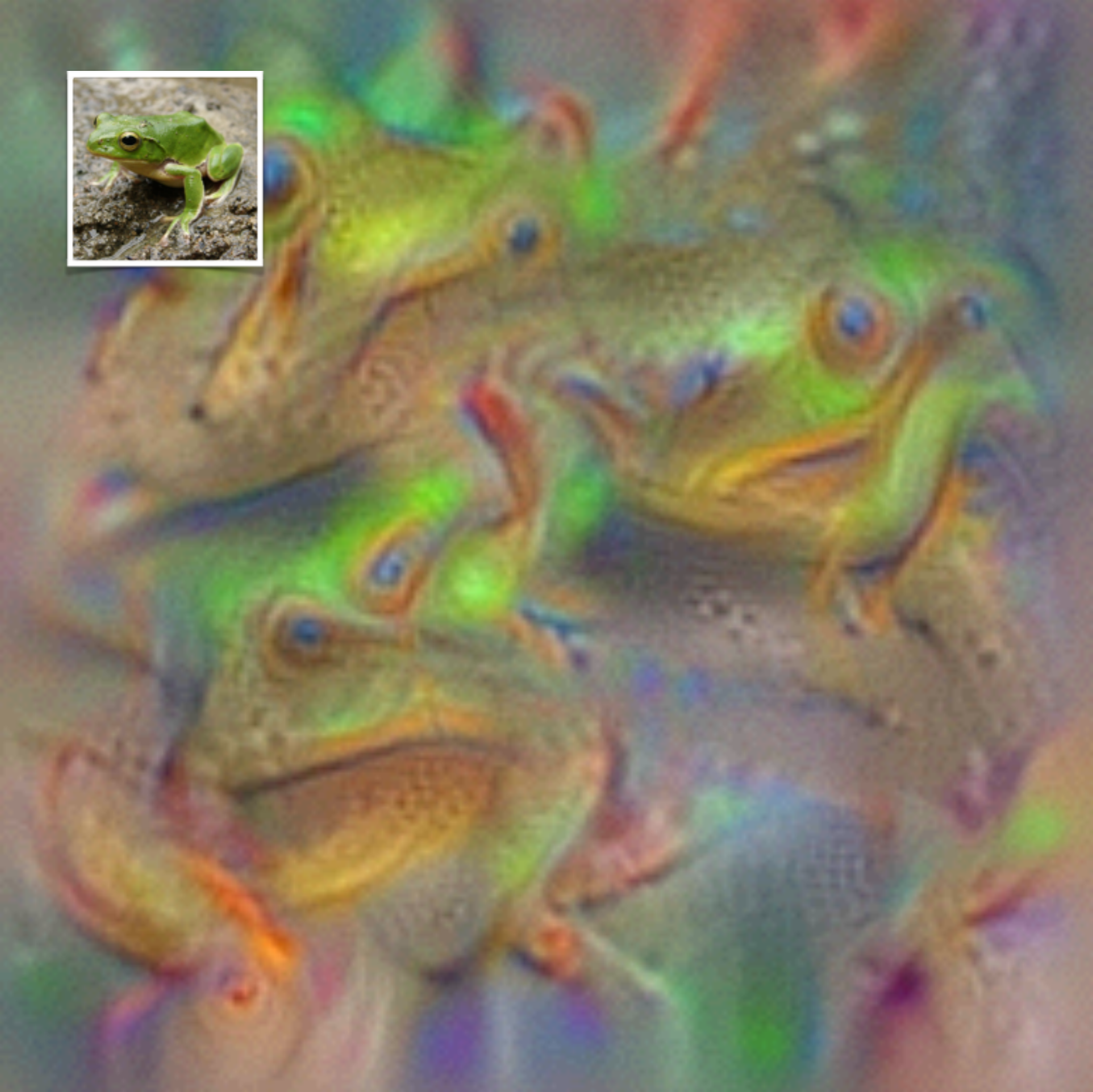


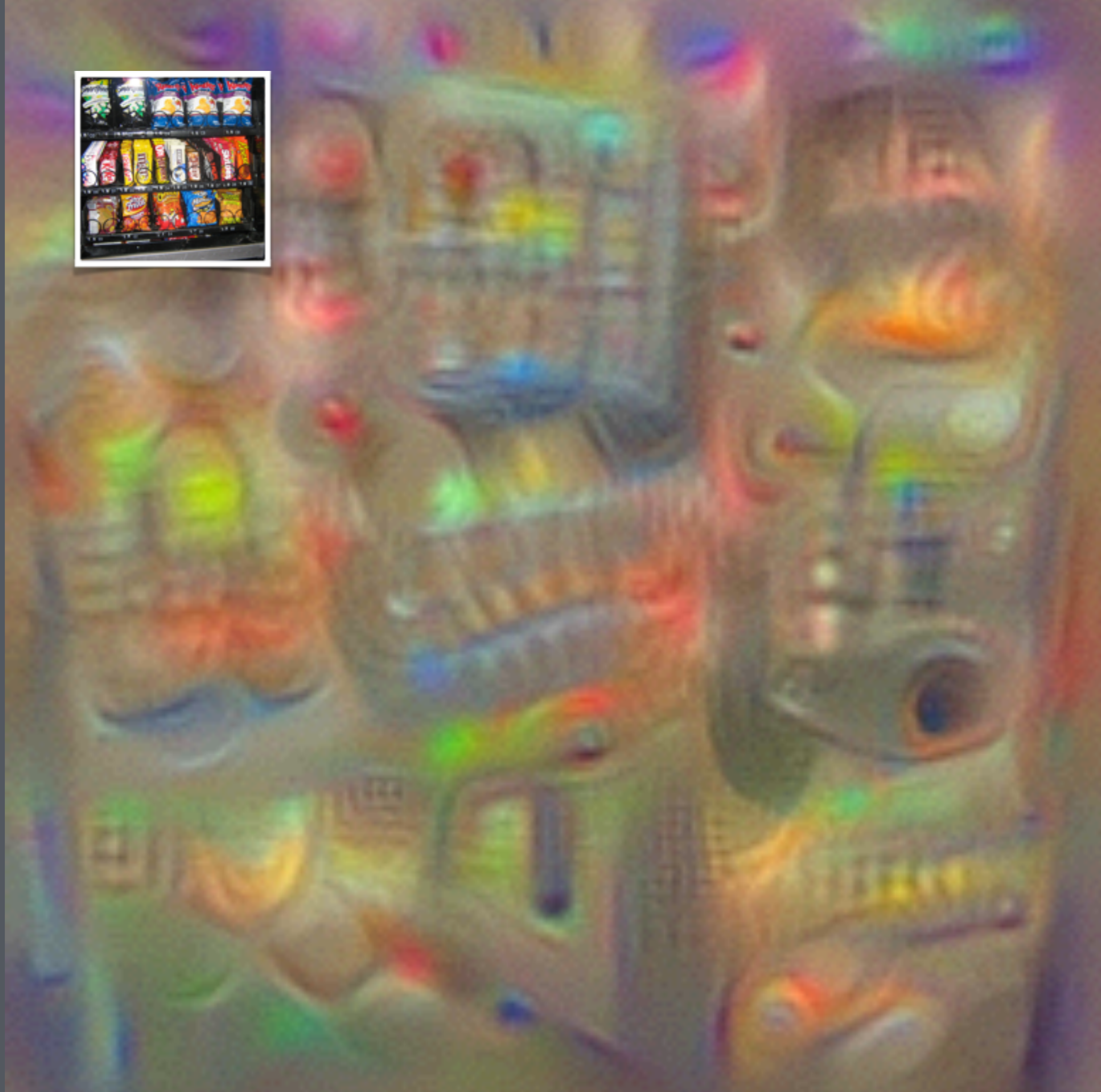






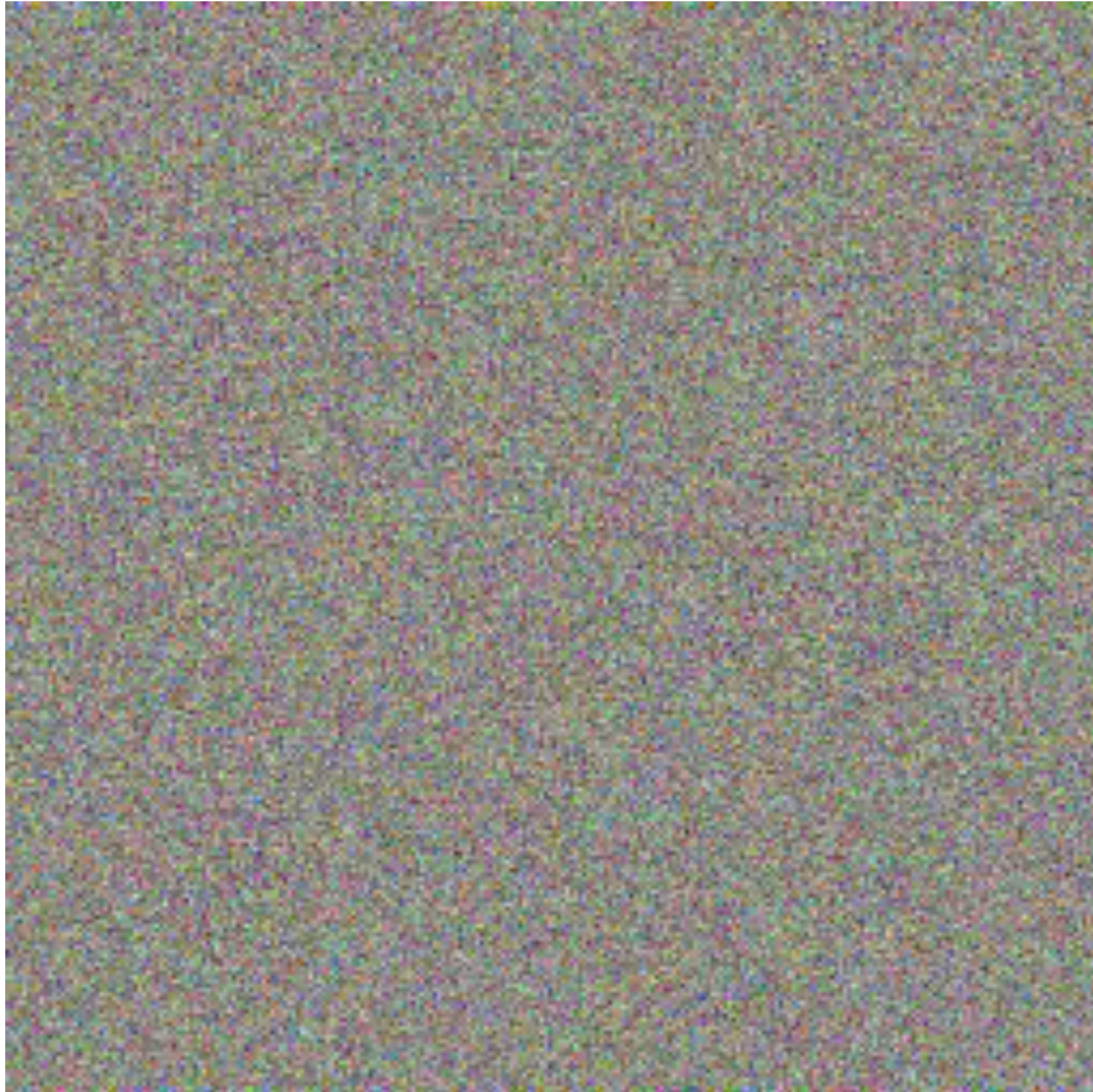


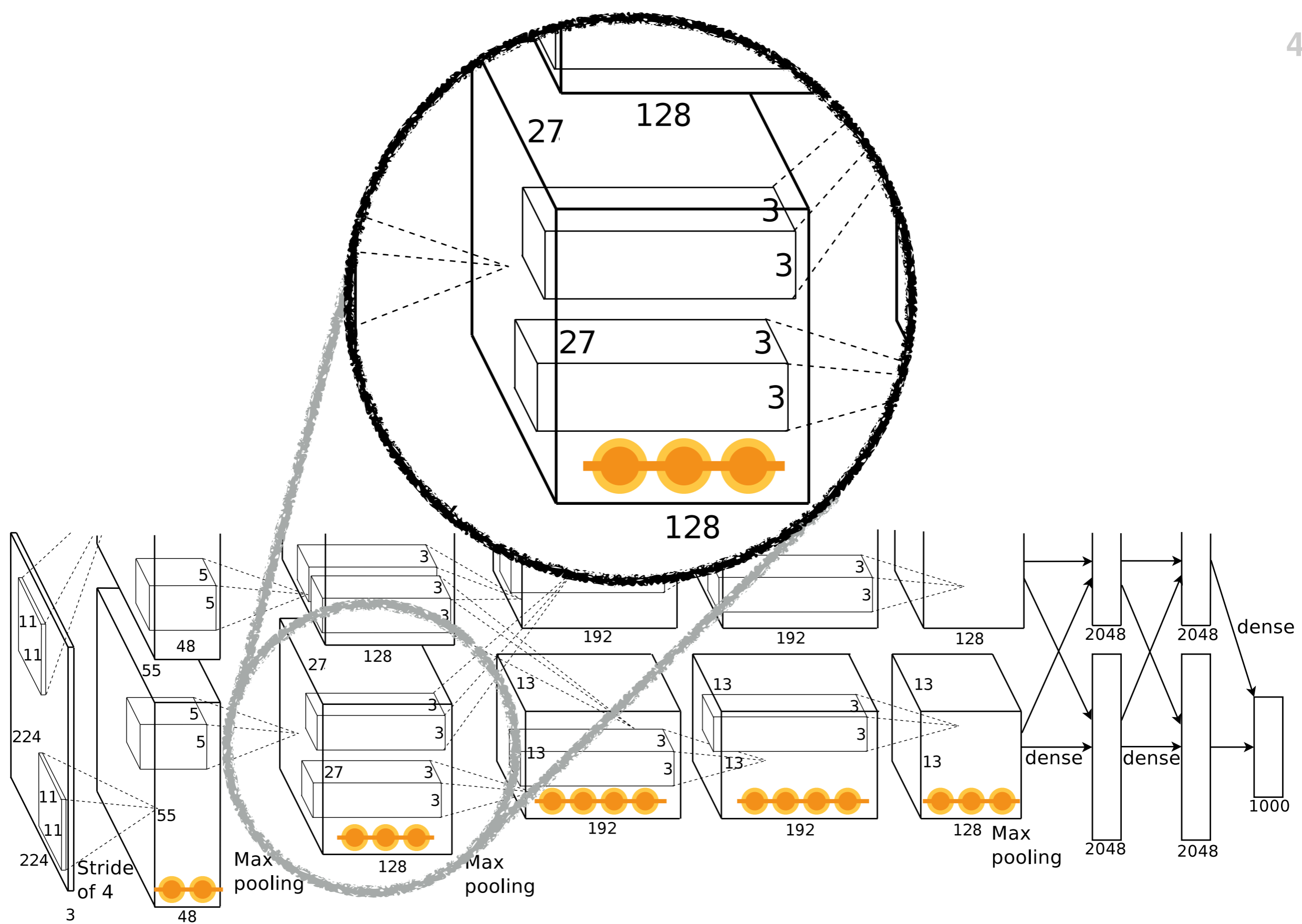




Remember: the starting point is white noise

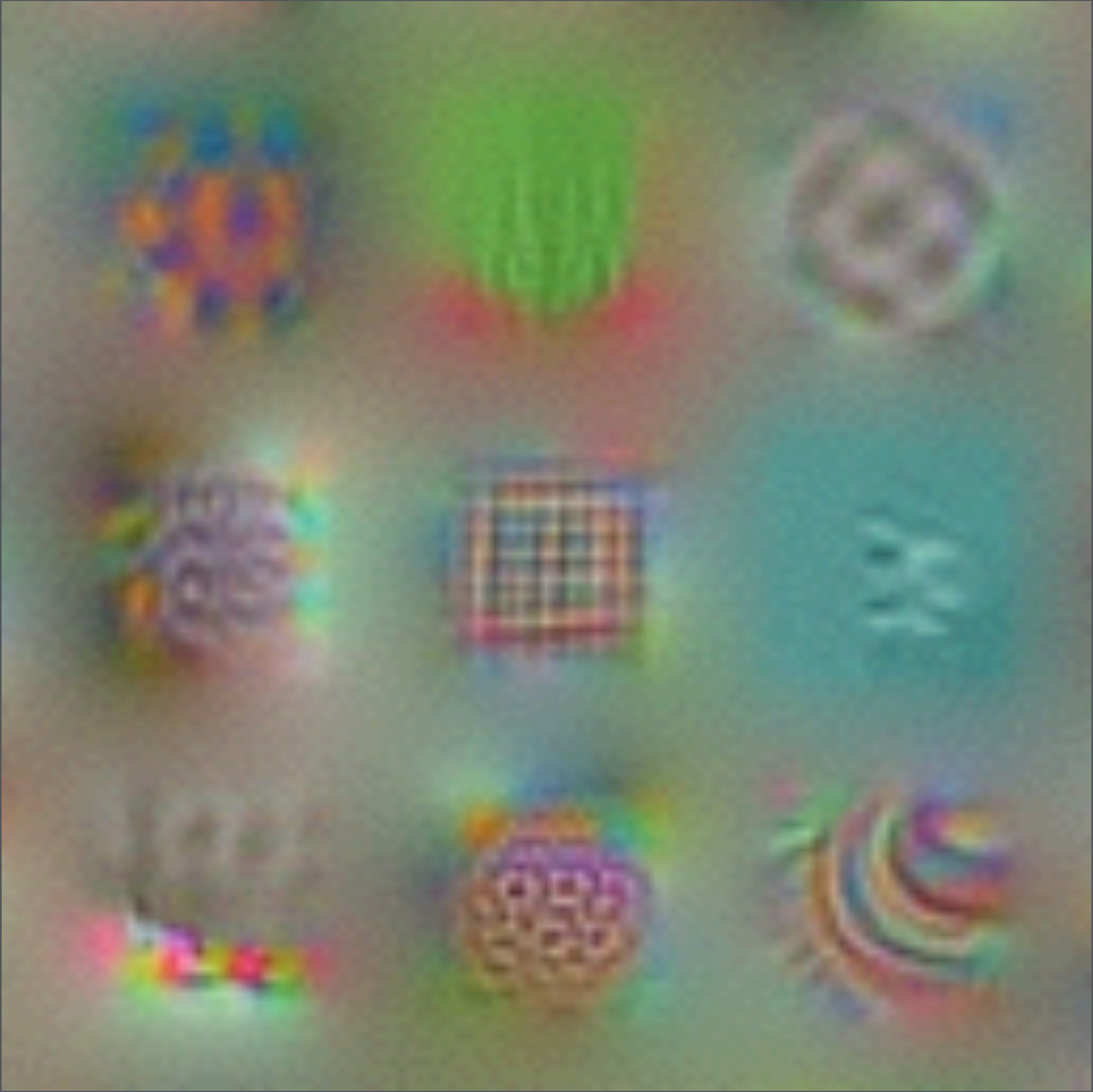
Not an image!







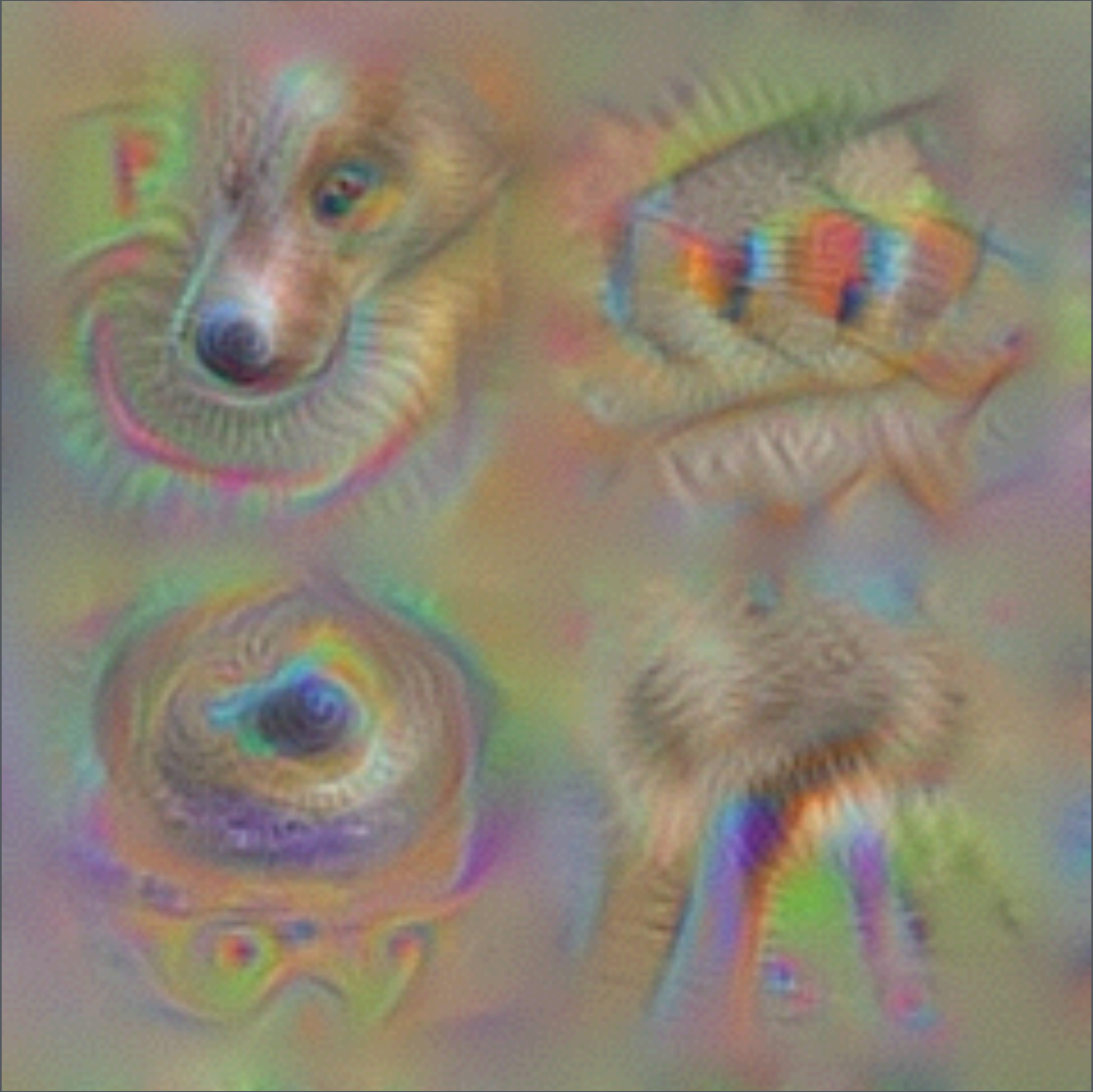
conv1



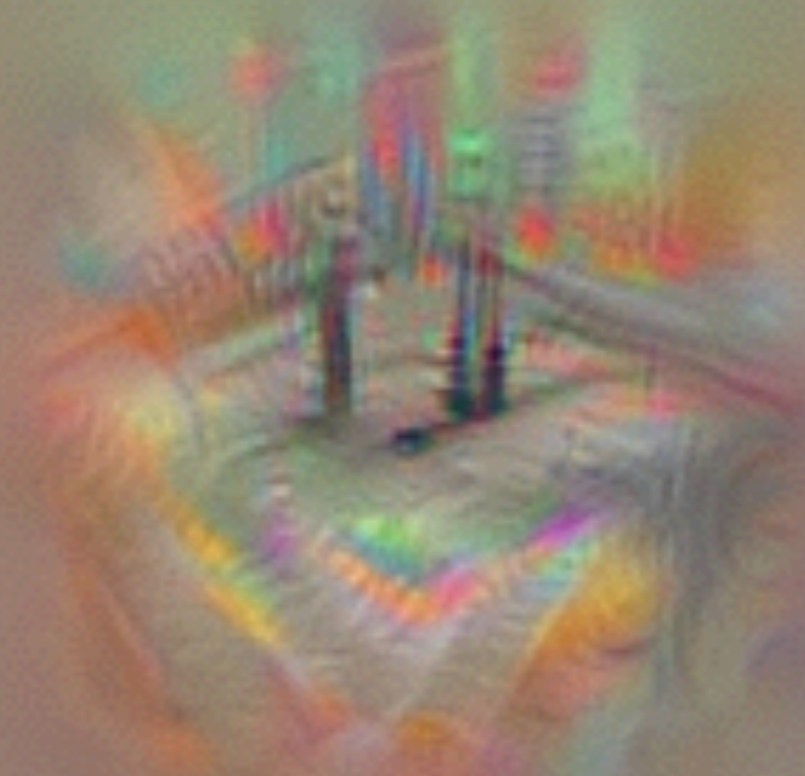
conv2



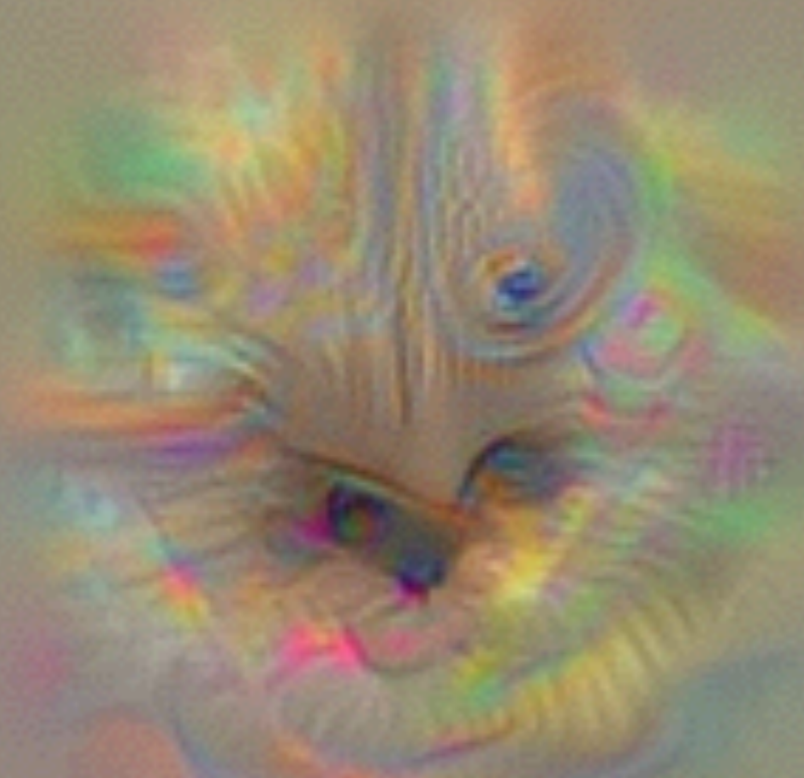
conv3



conv4



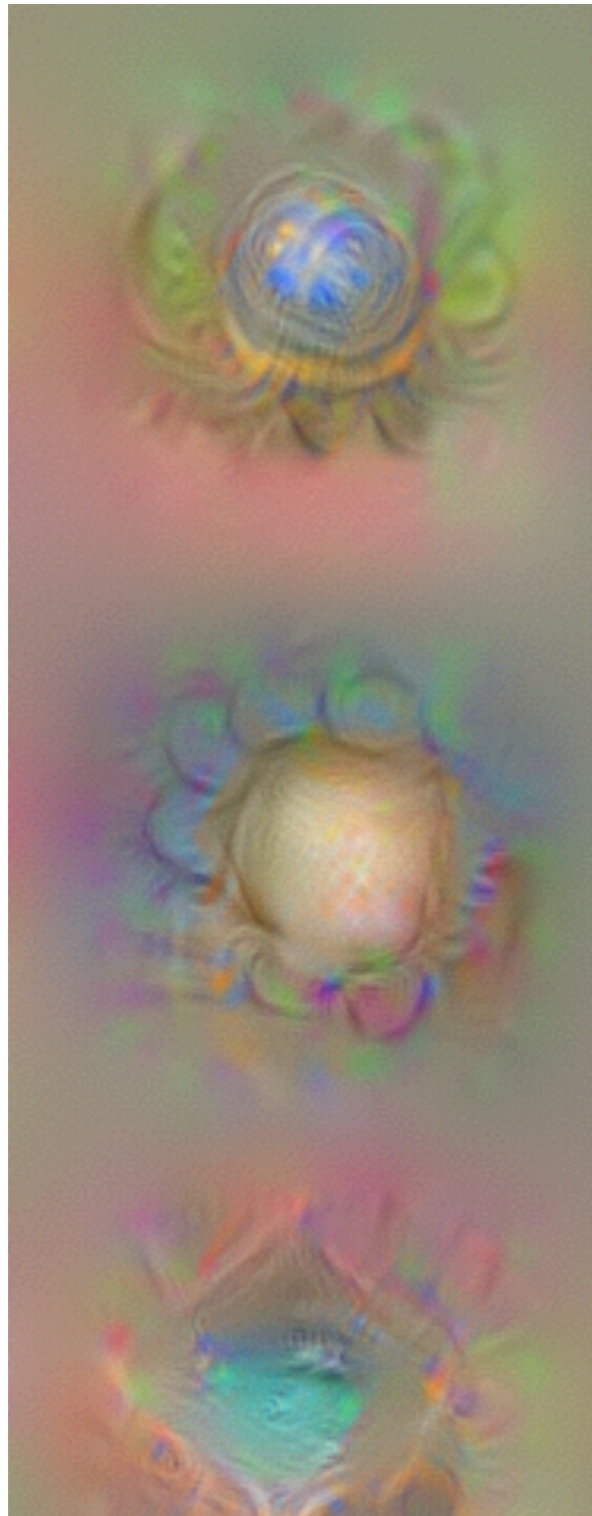
conv5



Network comparison

“conv5” features

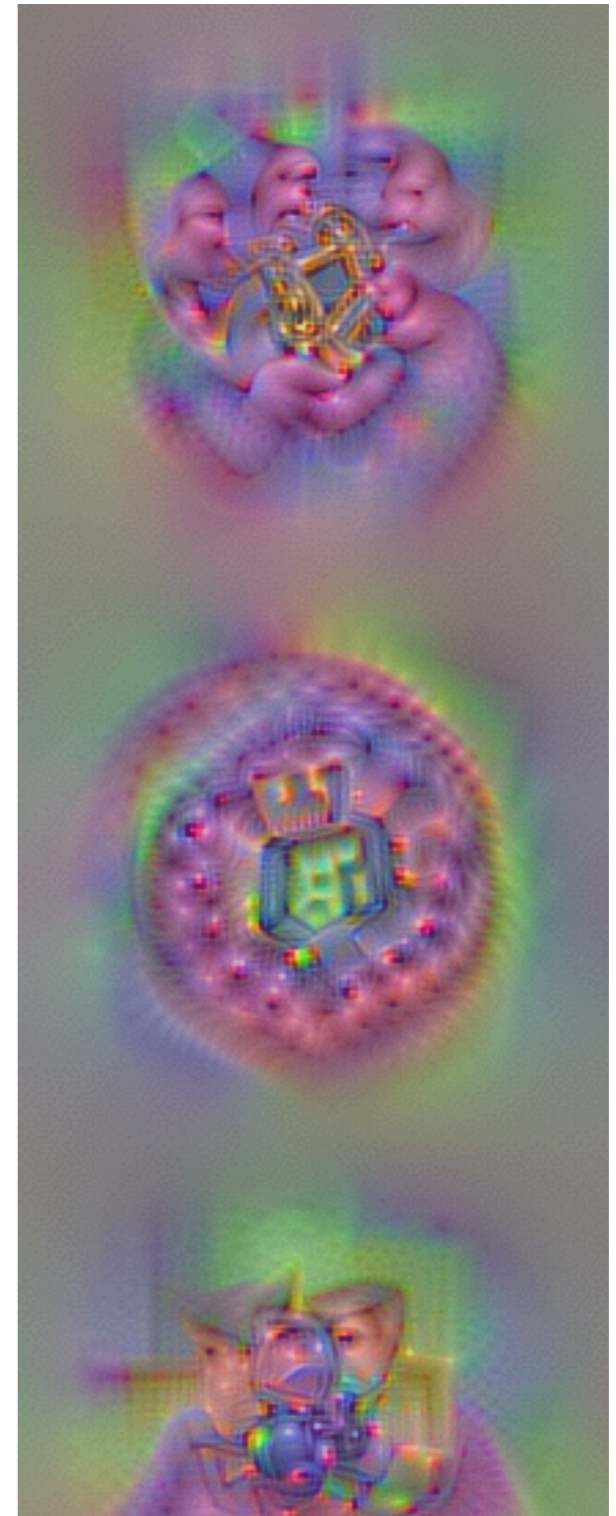
AlexNet



VGG-M



VGG-VD



[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)

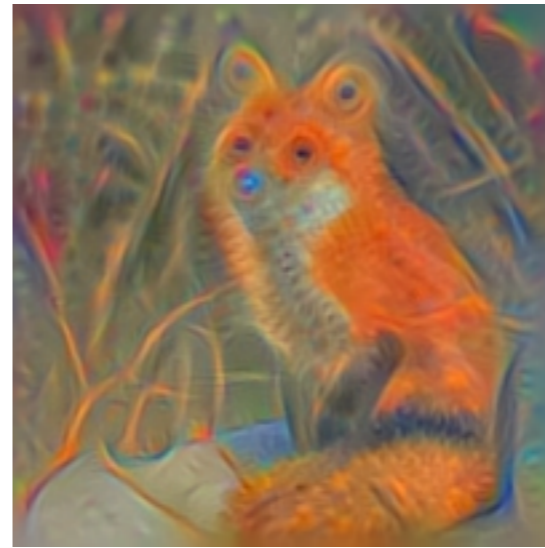
input



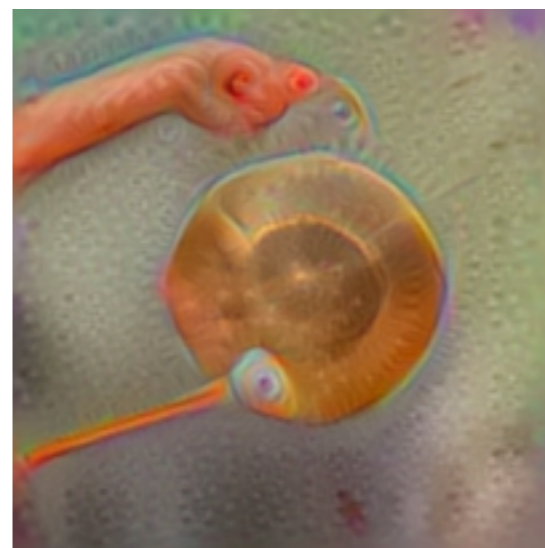
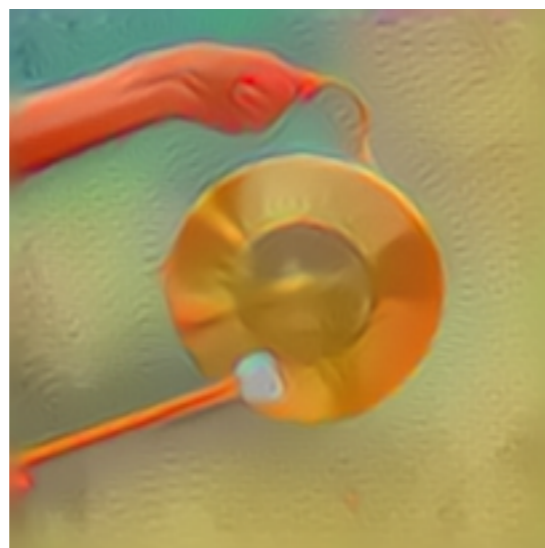
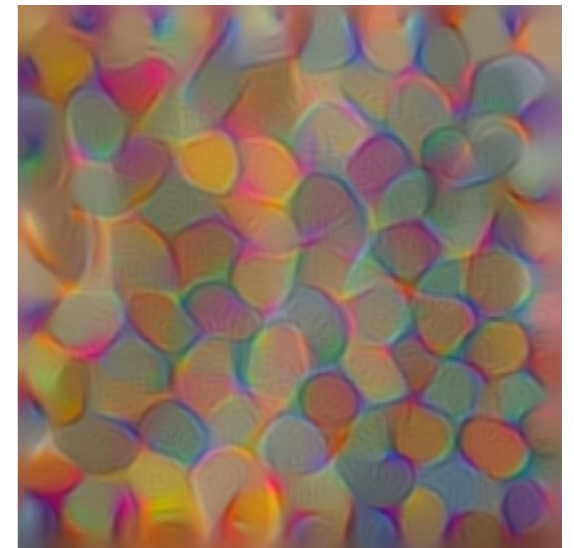
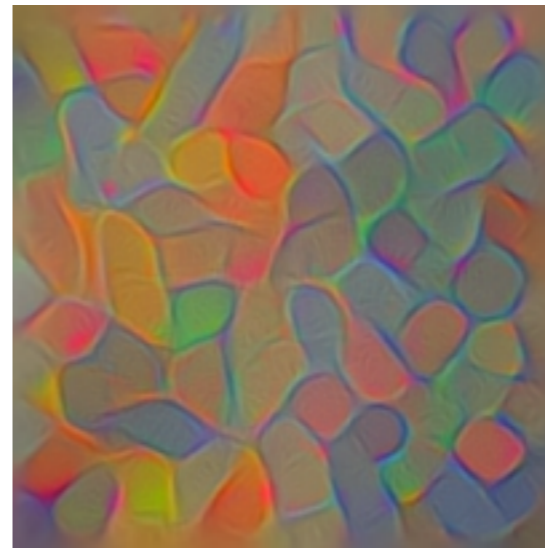
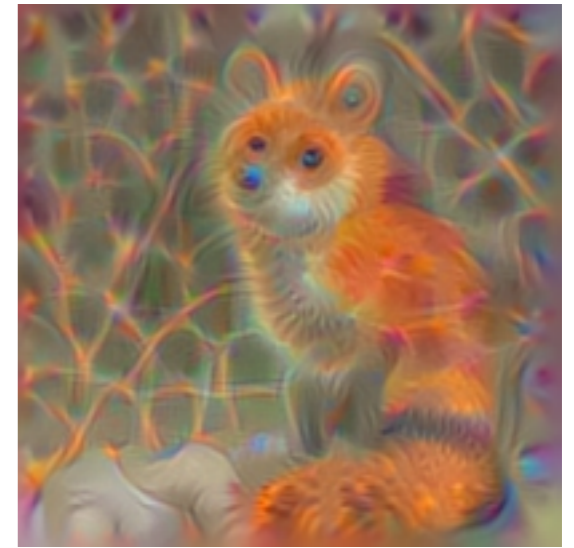
conv2



conv3



conv4



Caricaturization (VGG-M)

conv5



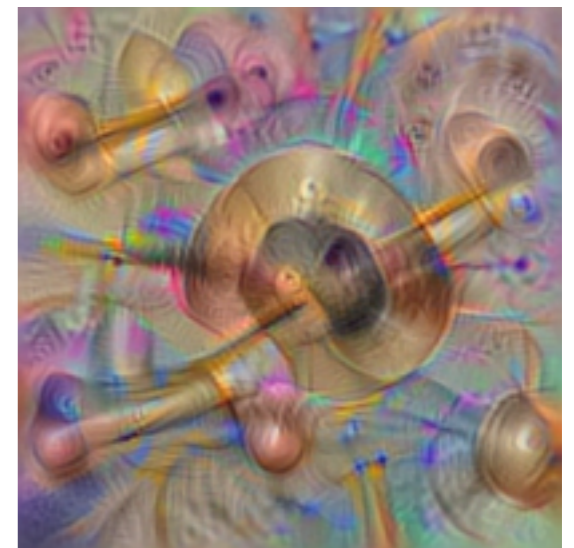
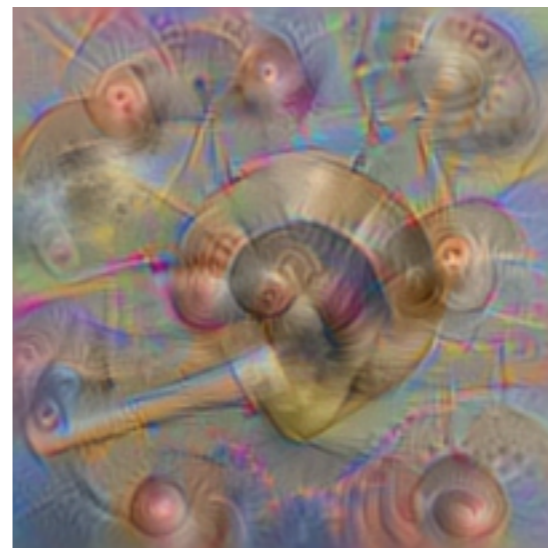
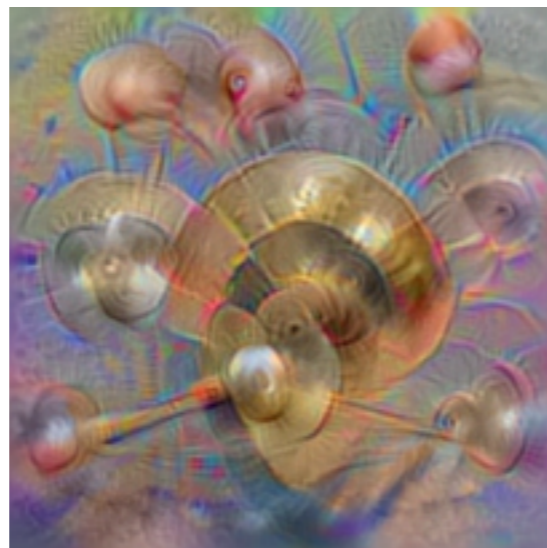
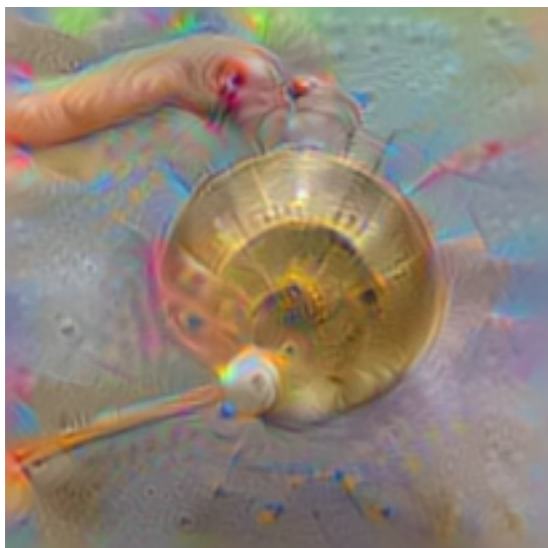
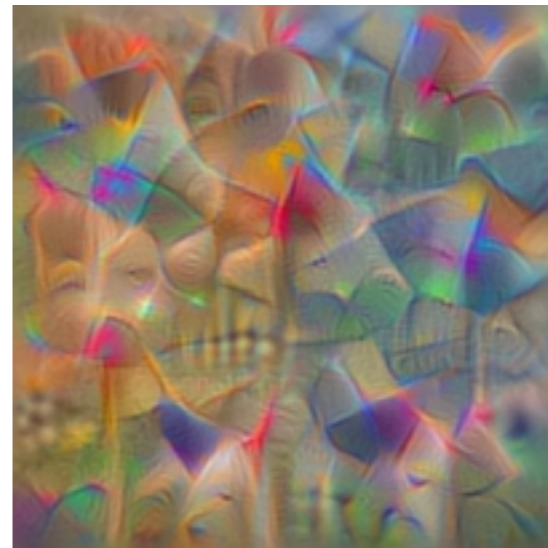
fc6



fc7



fc8





Surprisingly, the filters learned by discriminative neural networks capture well the “style” of an image.

This can be used to transfer the style of an image (e.g. a painting) to any other.

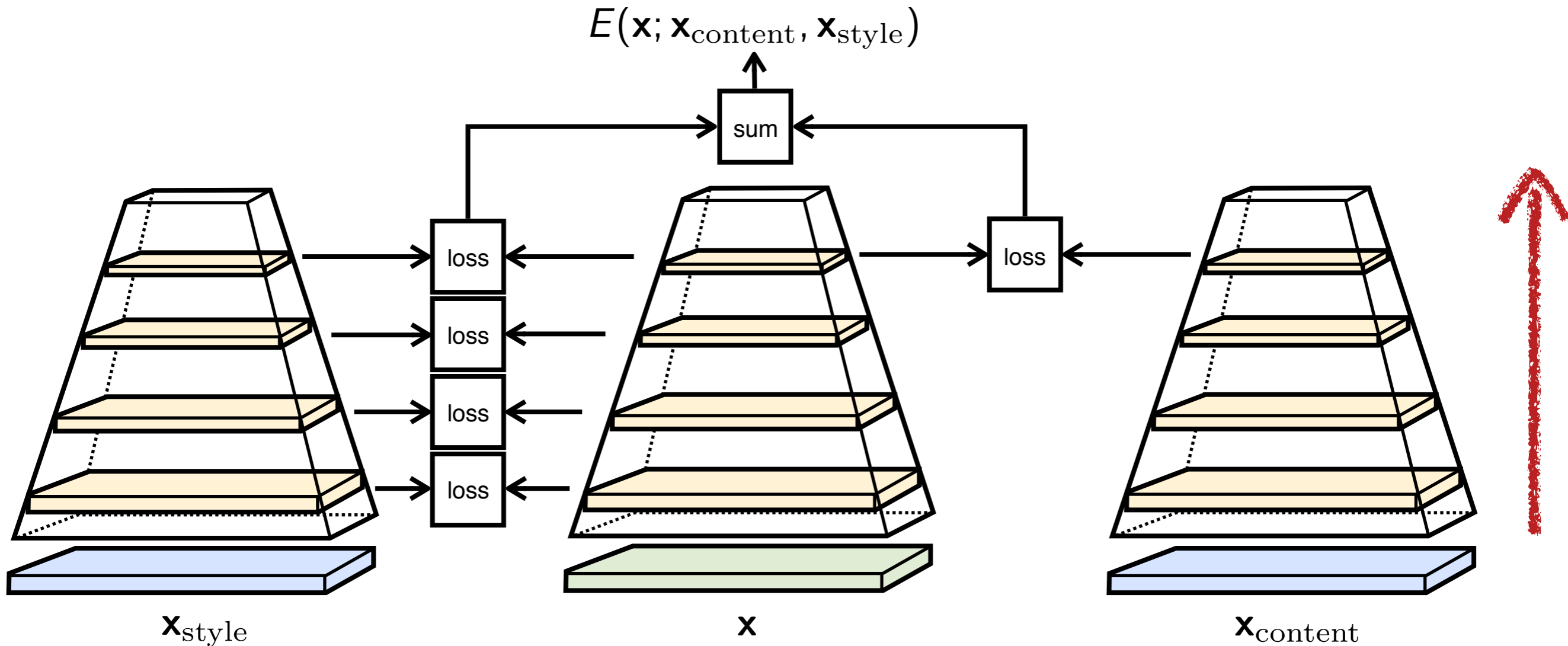
Optimisation based

L. A. Gatys, A. S. Ecker, and M. Bethge. Texture synthesis and the controlled generation of natural stimuli using convolutional neural networks. In Proc. NIPS, 2015.

Feed-forward neural network equivalents

D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. Proc. ICML, 2016.

J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In Proc. ECCV, 2016.



Moment matching

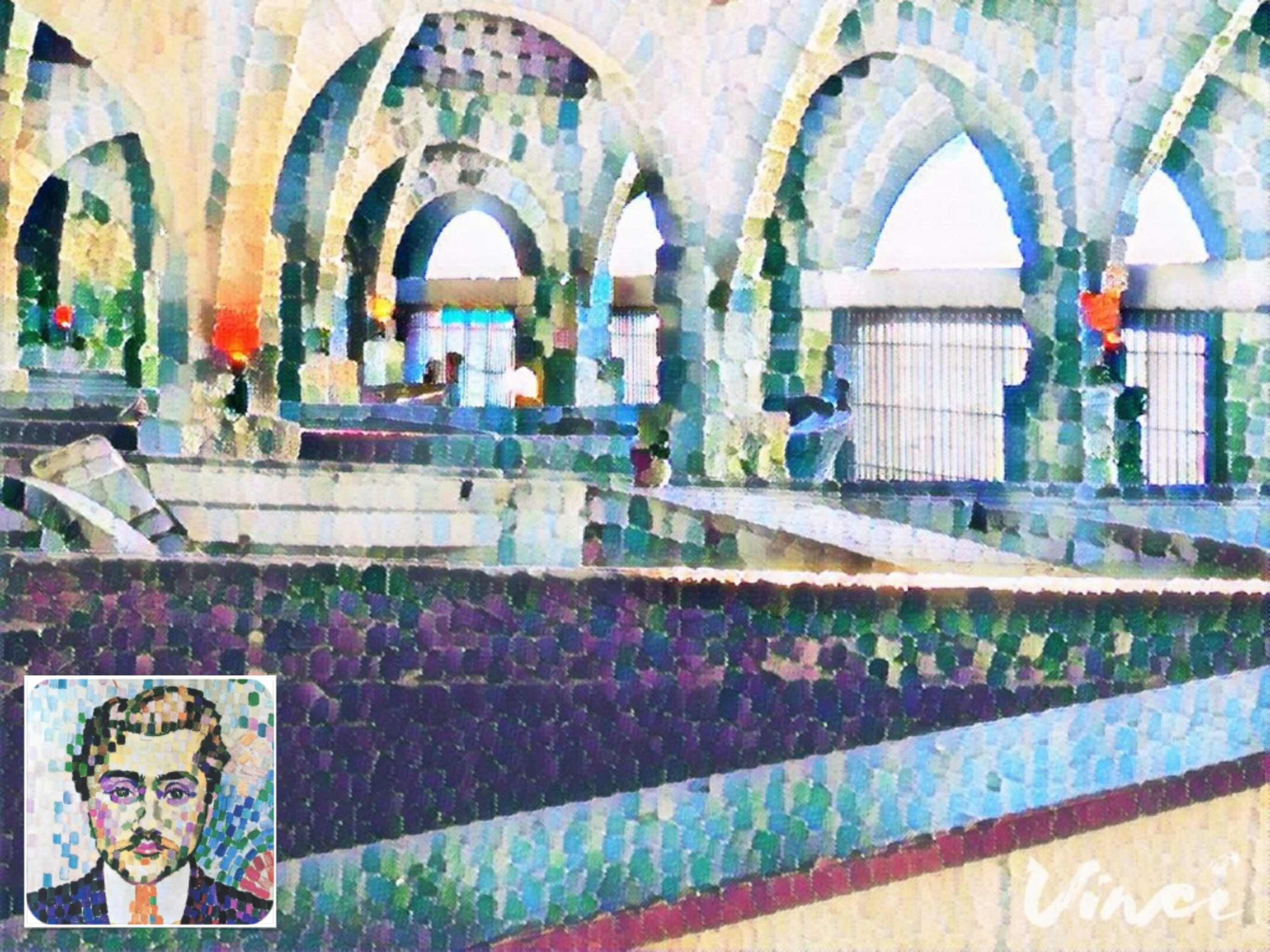
- ▶ Content statistics: same as inversion
- ▶ Style statistics: cross-channel correlations

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} E(\mathbf{x}; \mathbf{x}_{\text{content}}, \mathbf{x}_{\text{style}})$$





Vinci



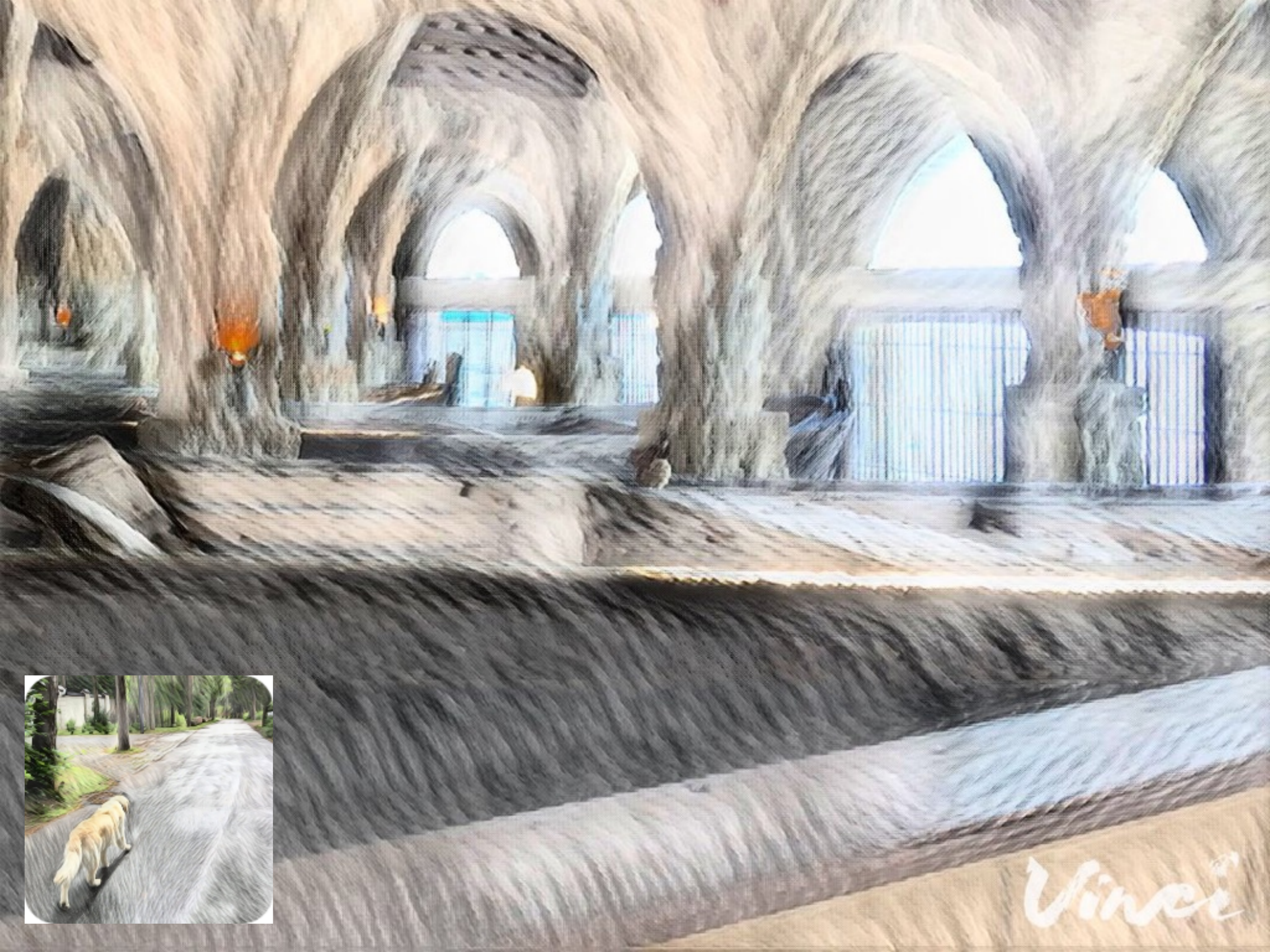
Vincci



Vincci



Vinci



Vinci



Vinci



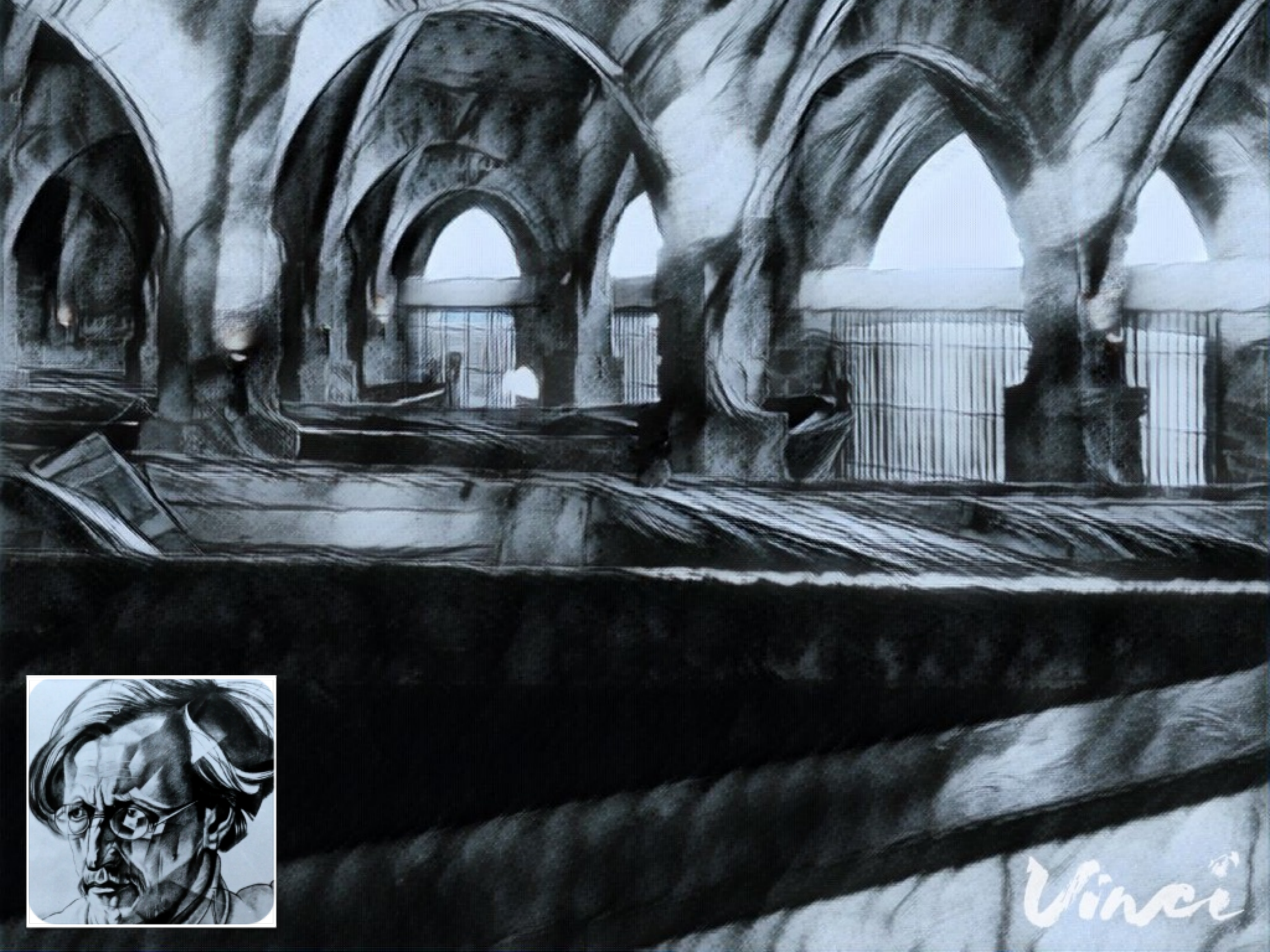
Vince



Vinci



Vinci



Vinci



Vinci





Vincci



Vincci

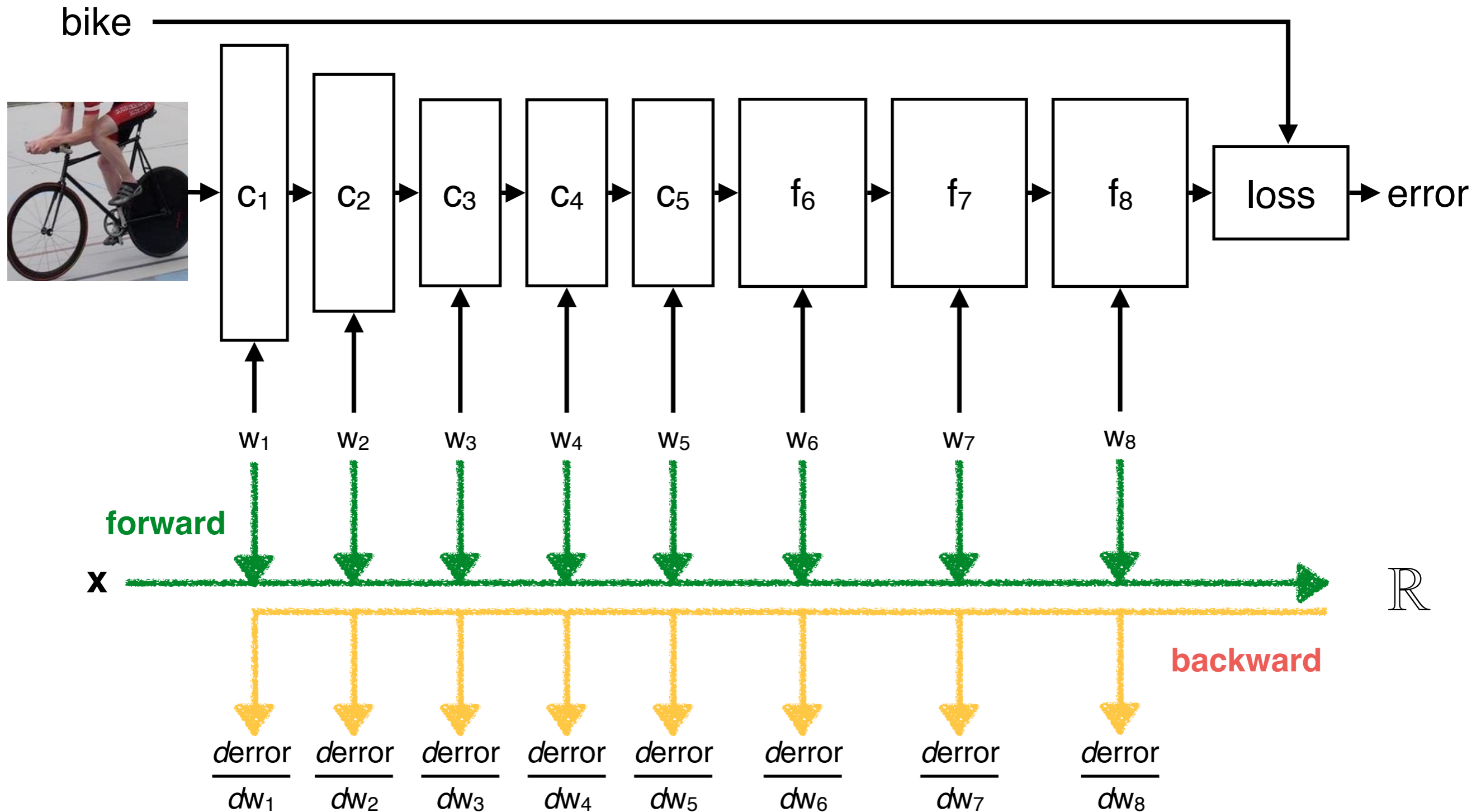
Visualizing representations

Backpropagation networks and “deconvolution”

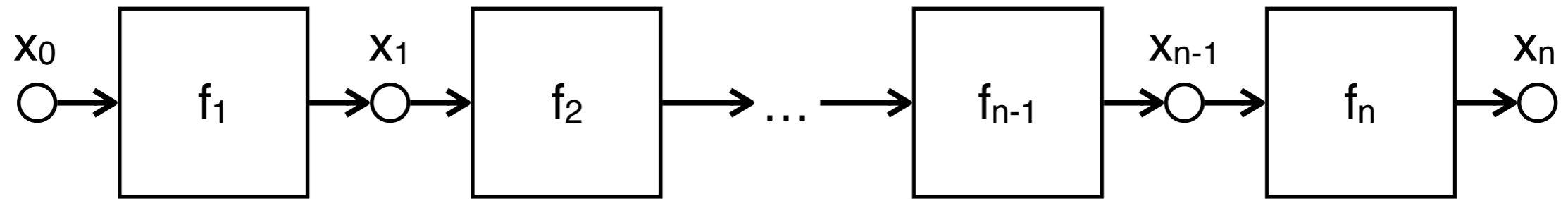
Representations: equivalence & transformations

Backpropagation

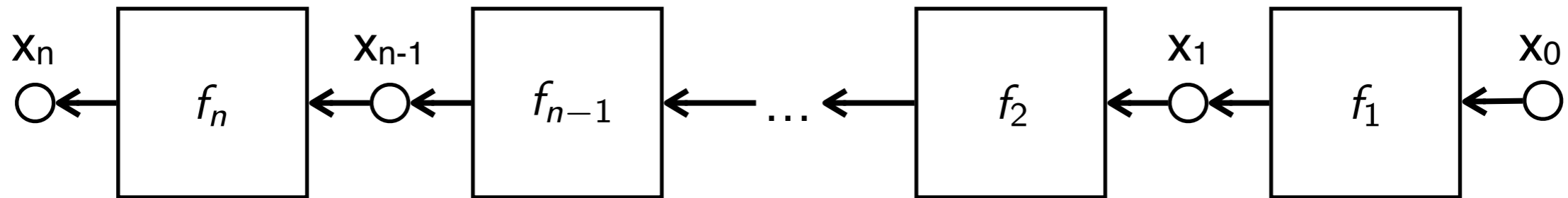
Compute derivatives using the chain rule



Chain rule: scalar version



Chain rule: scalar version



A composition of n functions

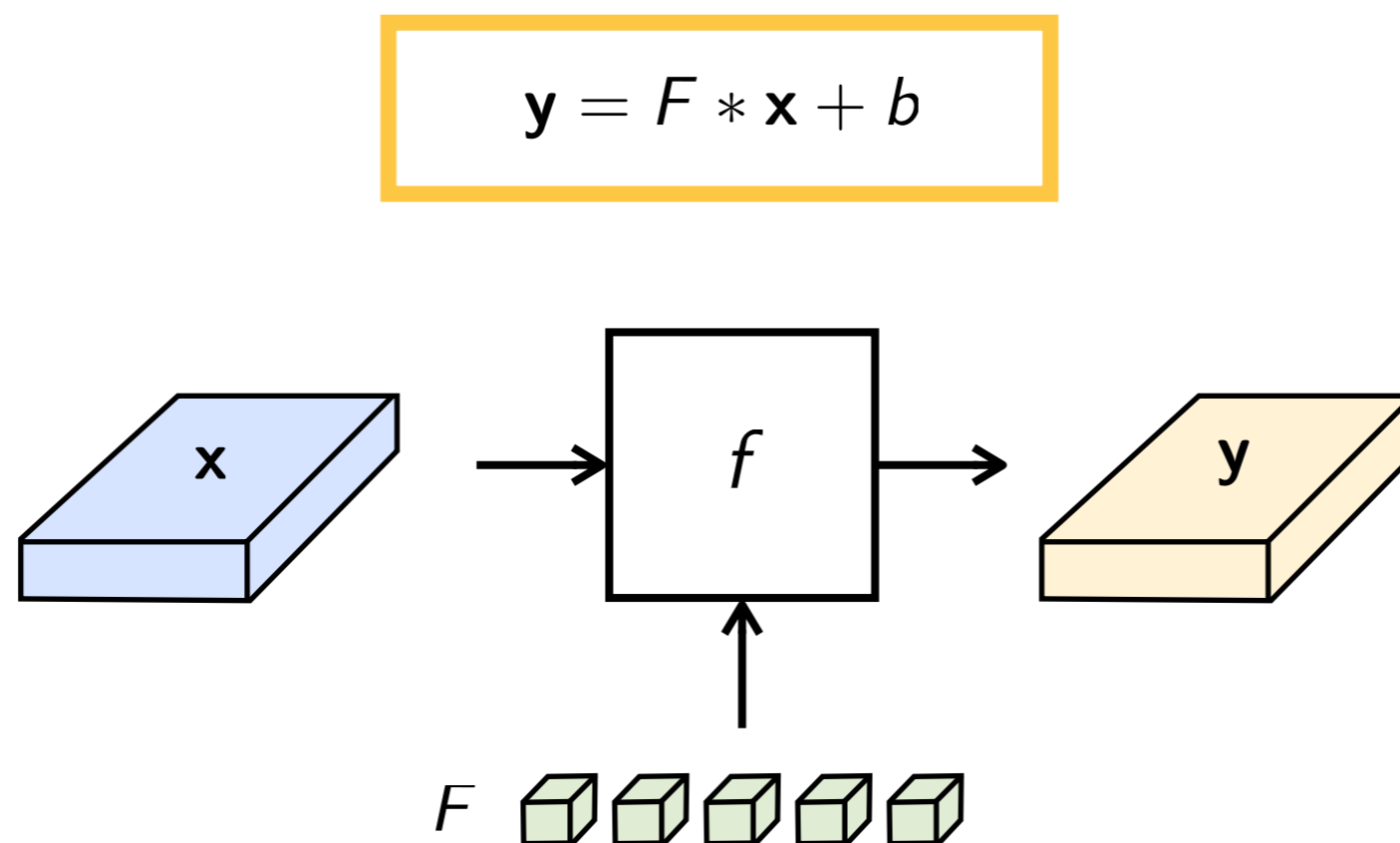
$$\mathbf{x}_n = (f_n \circ f_{n-1} \circ \dots \circ f_2 \circ f_1) (\mathbf{x}_0)$$
$$\frac{dx_n}{dx_0} = \frac{df_n}{dx_{n-1}} \times \frac{df_{n-1}}{dx_{n-2}} \times \dots \times \frac{df_2}{dx_1} \times \frac{df_1}{dx_0}$$

The diagram shows the relationship between the function composition and the chain rule. Dashed arrows point from each function in the composition to its corresponding derivative term in the chain rule equation.

Derivative ← chain rule

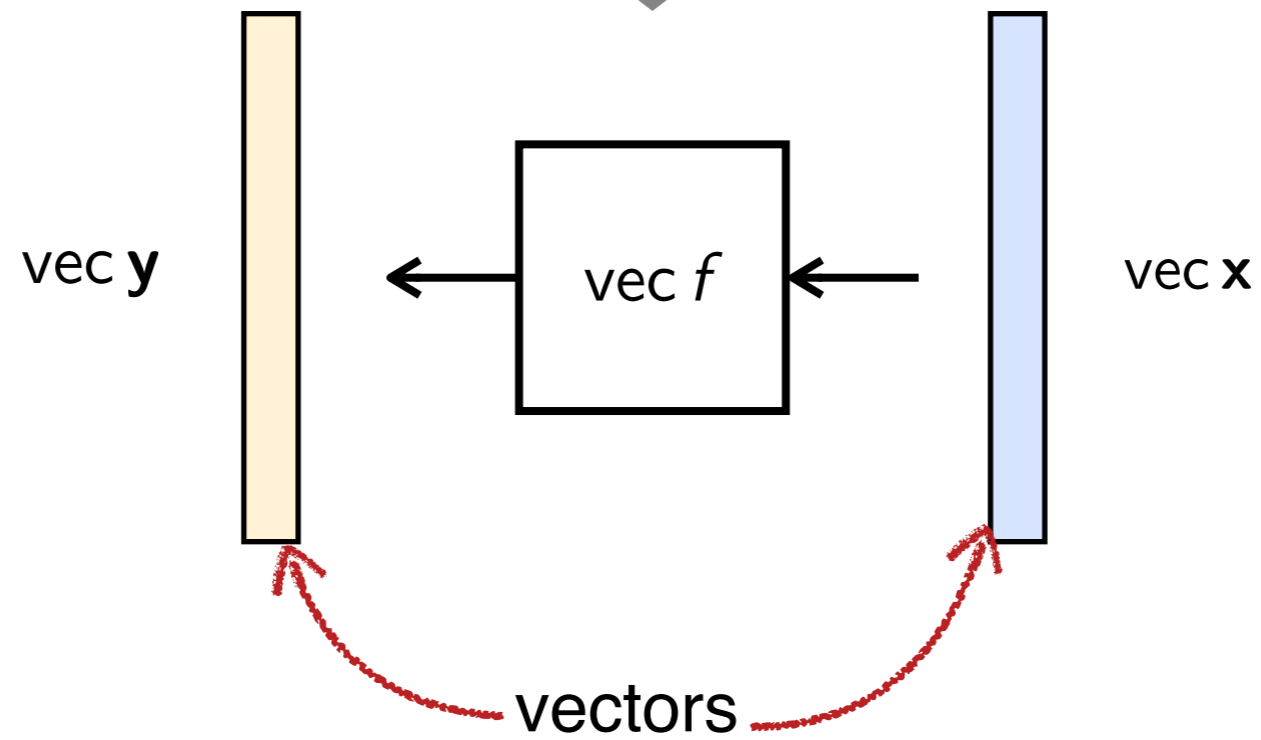
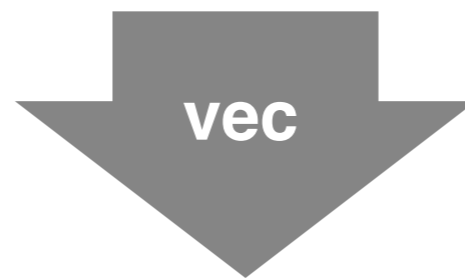
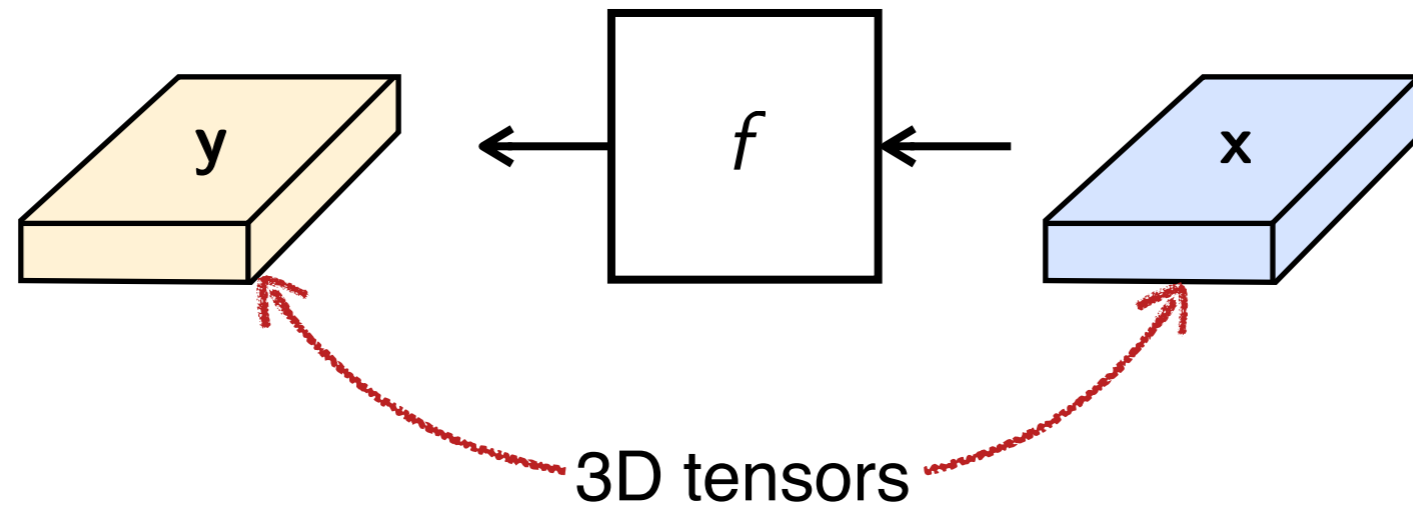
Tensor-valued functions

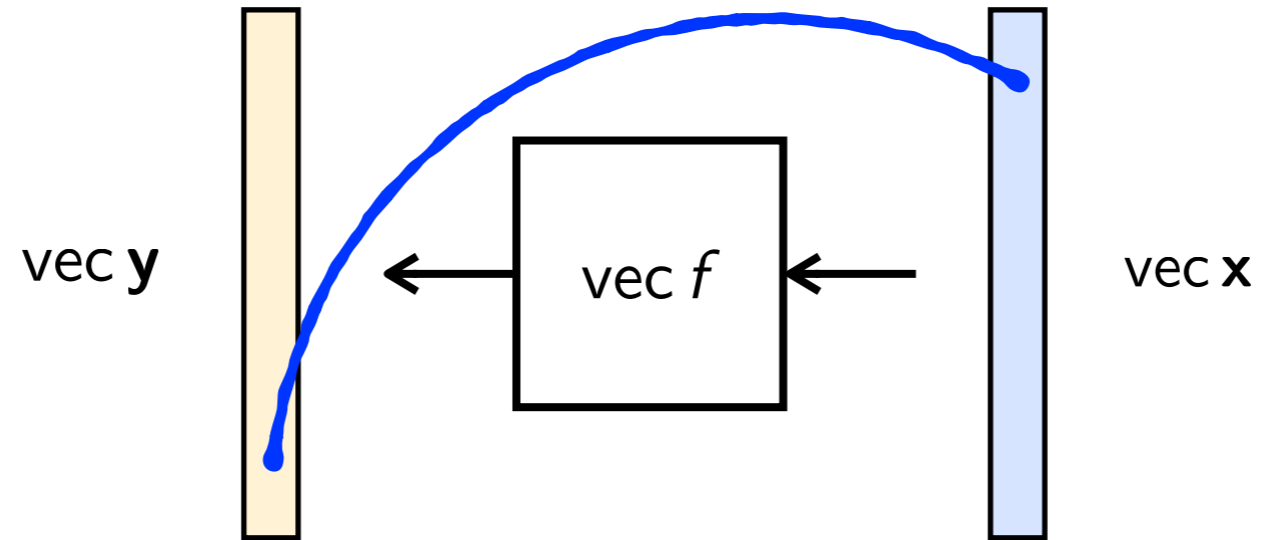
E.g. linear convolution = bank of 3D filters



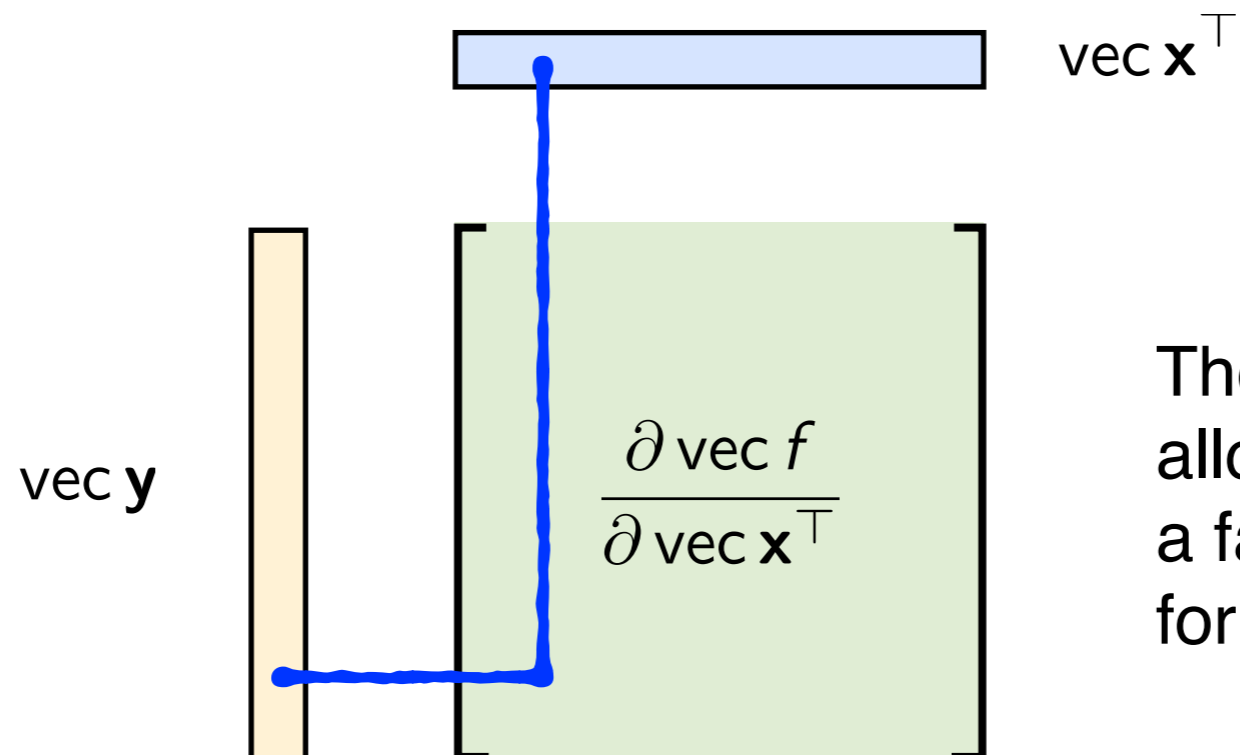
	height	width	channels	instances
input x	H	W	C	1 or N
filters F	H_f	W_f	C	K
output y	$H - H_f + 1$	$W - W_f + 1$	K	1 or N

Vector representation





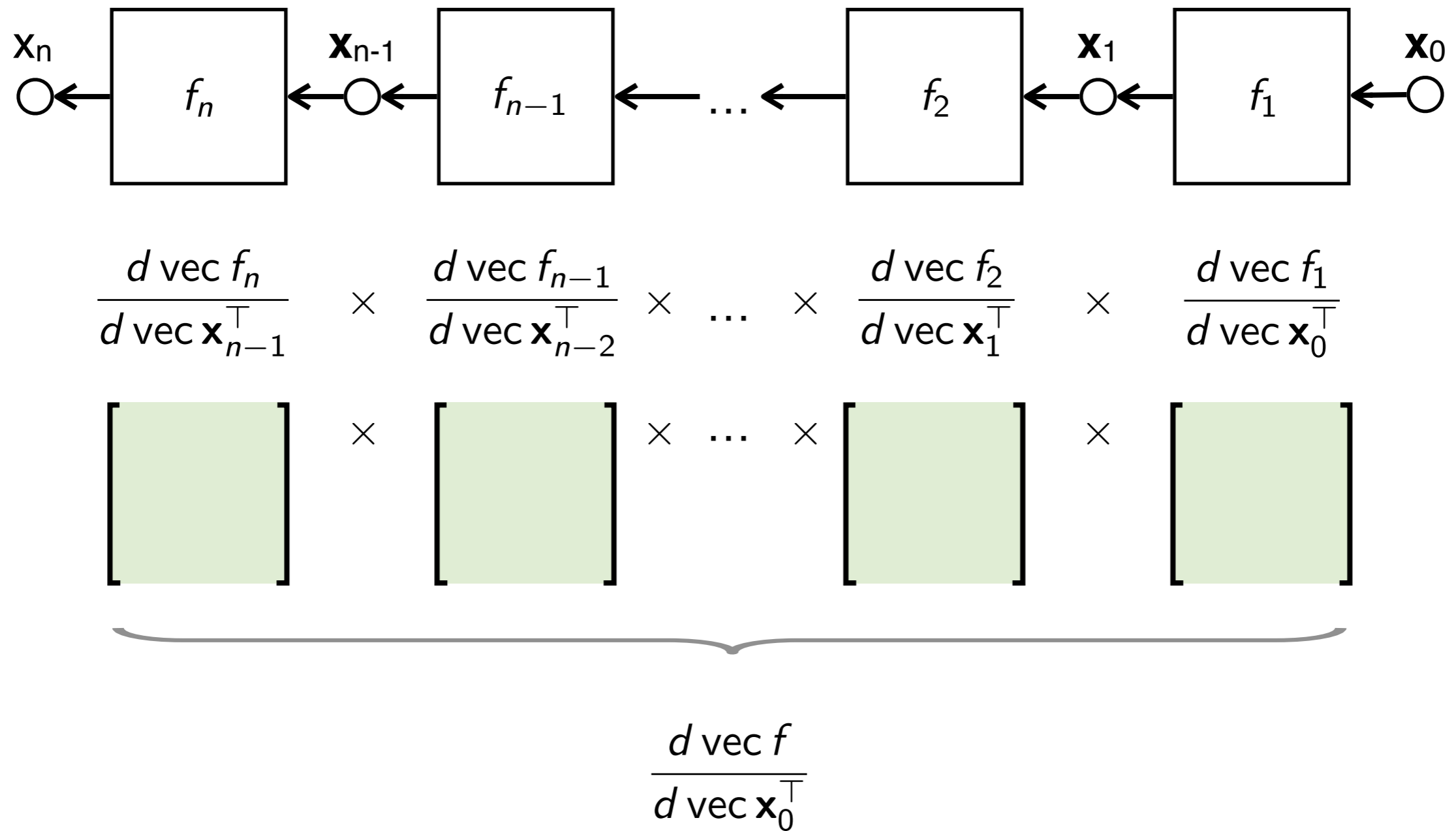
Derivative (Jacobian): every output element w.r.t. every input element!



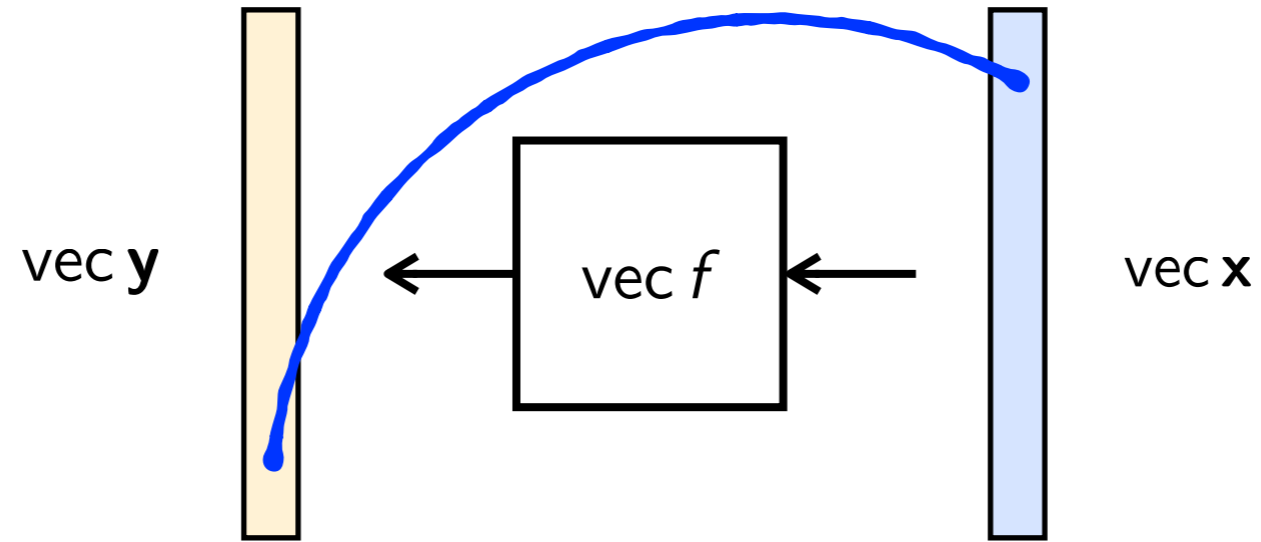
The **vec** operator allows us to use a familiar matrix notation for the derivatives

Chain rule: tensor version

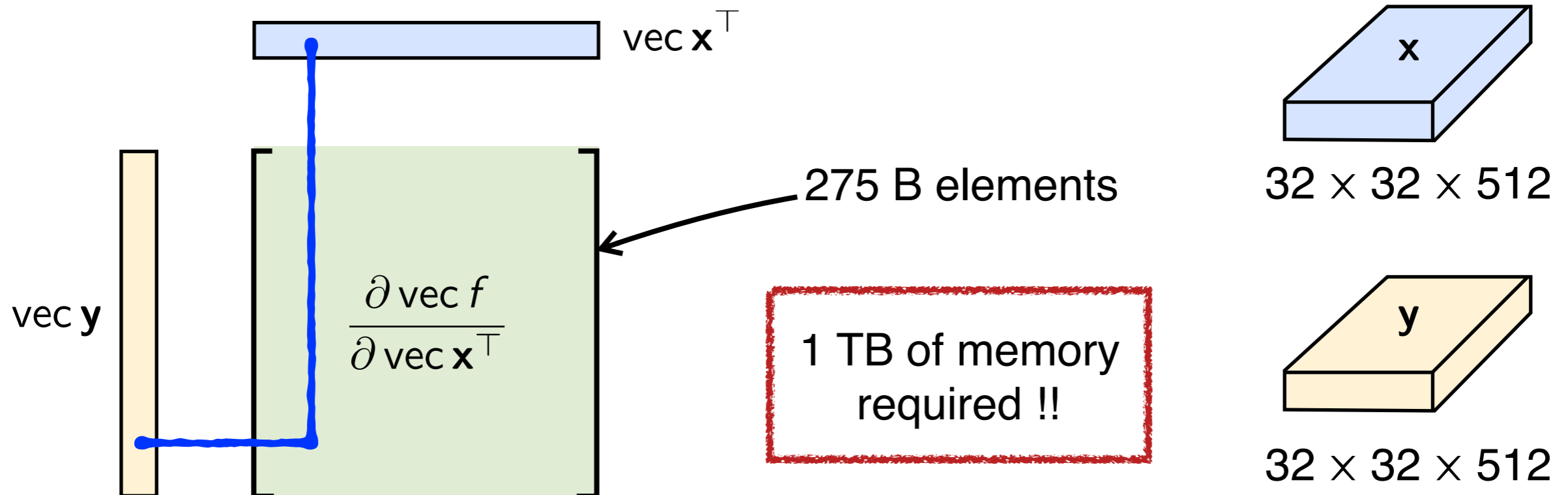
Using $\text{vec}()$ and matrix notation



The (unbearable) size of tensor derivatives



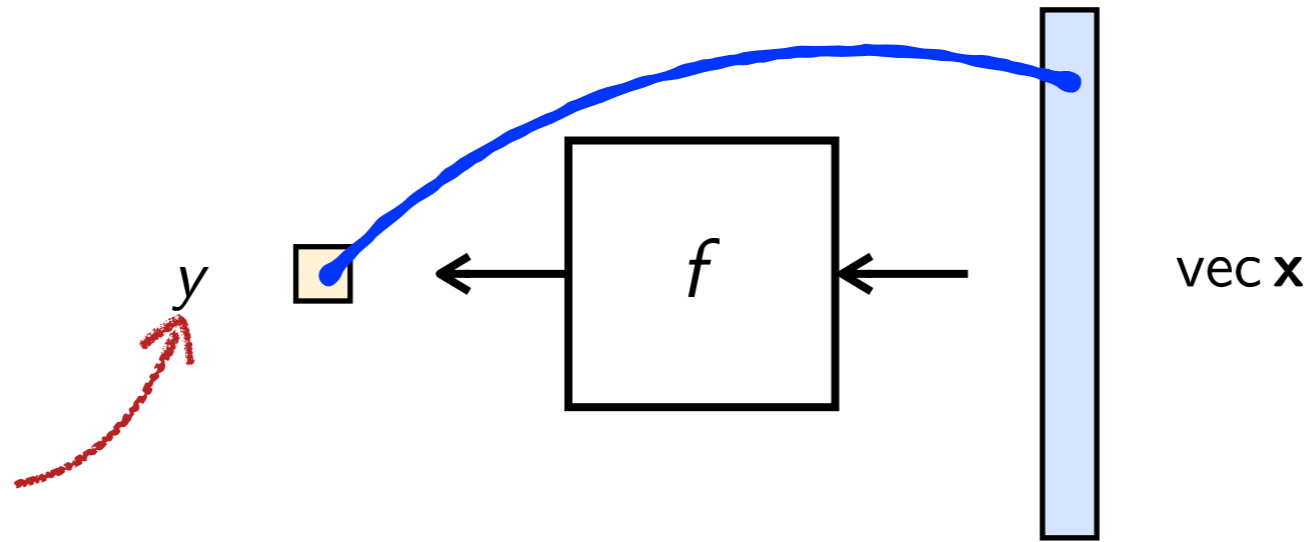
The size of these Jacobian matrices is **huge**. Example:



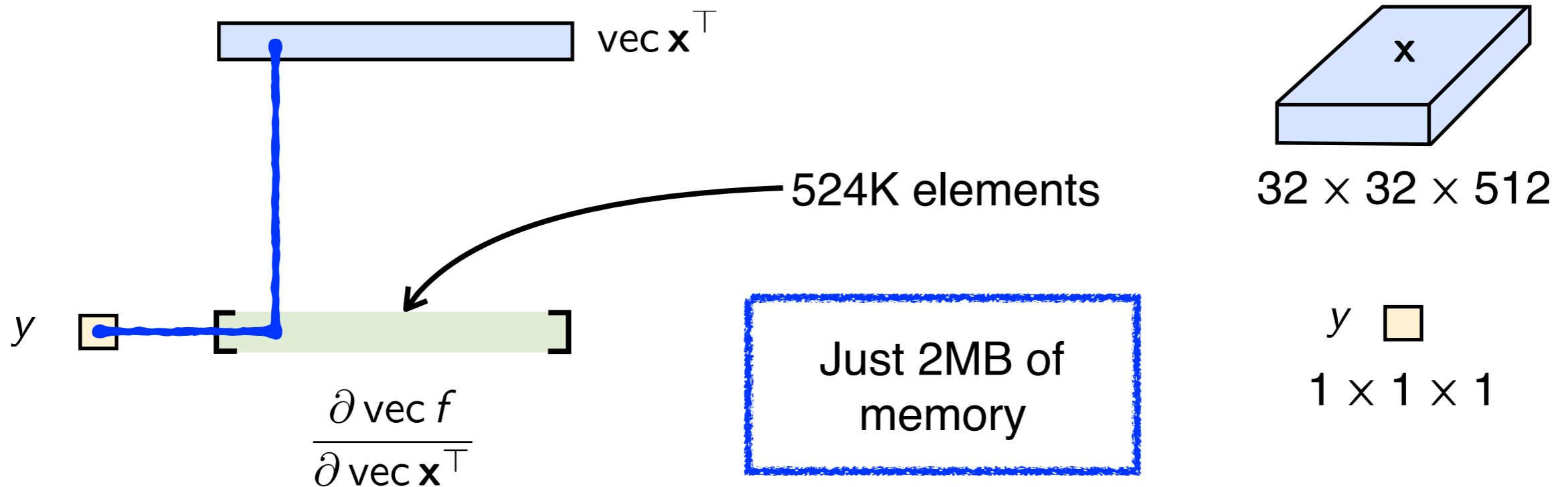
Unless the output is a scalar

Scalar

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

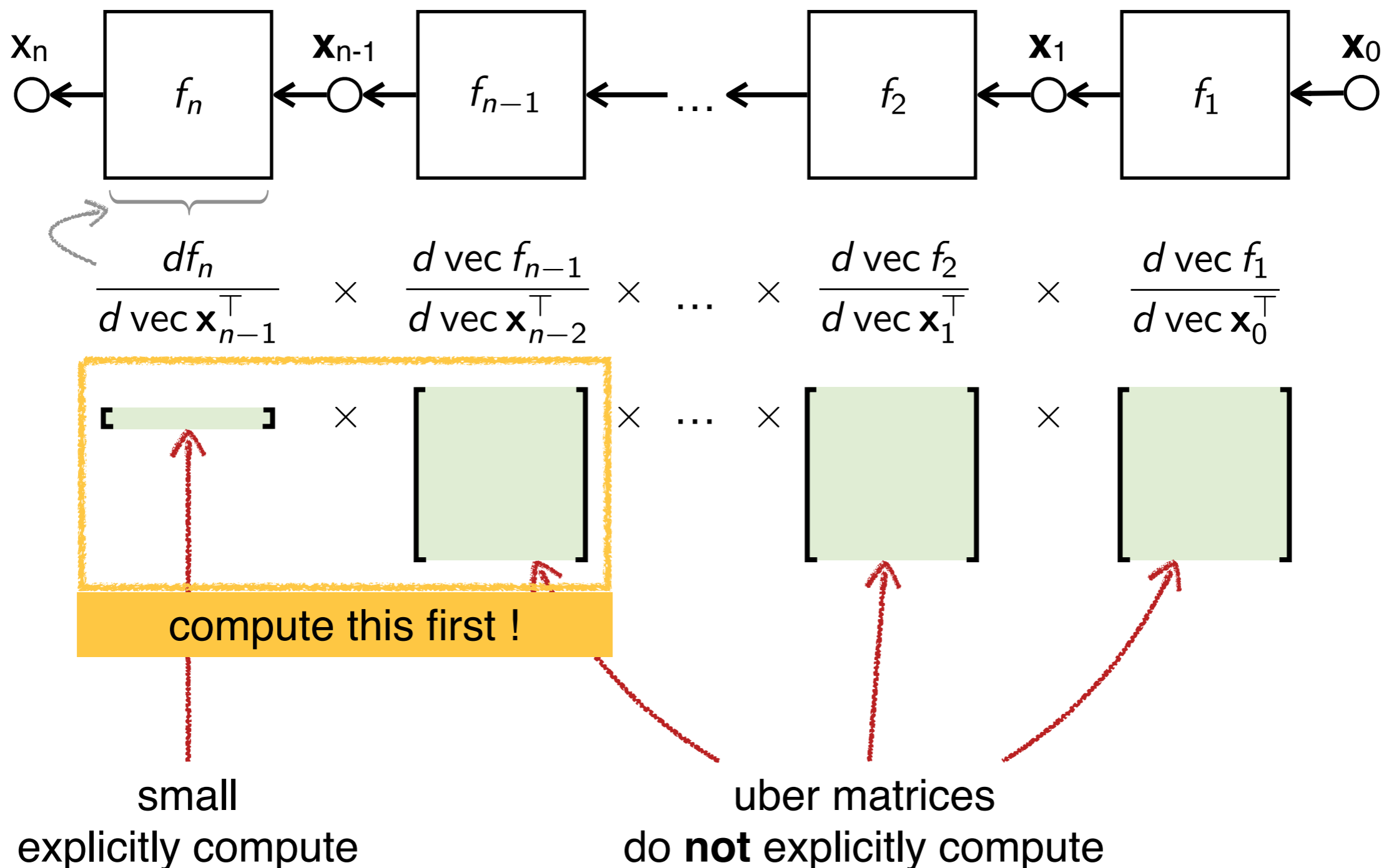


Now the Jacobian has the same size as \mathbf{x} . Example:



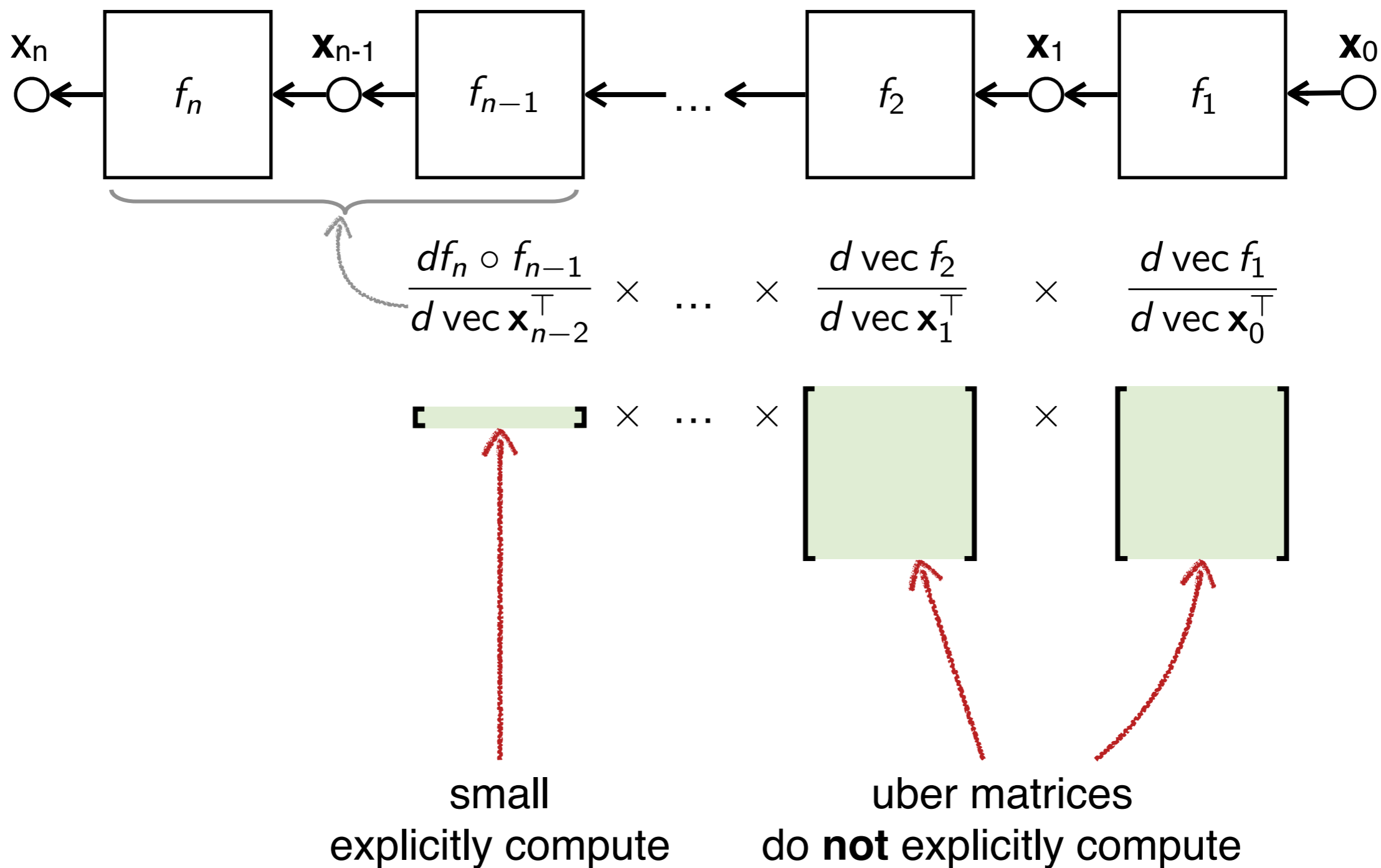
Backpropagation

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



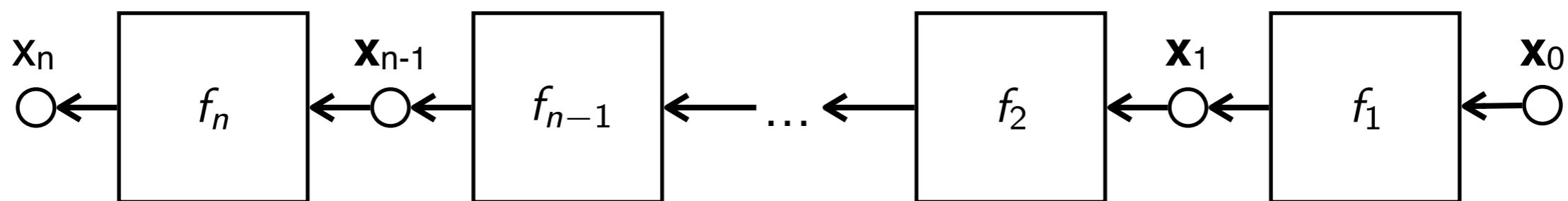
Backpropagation

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

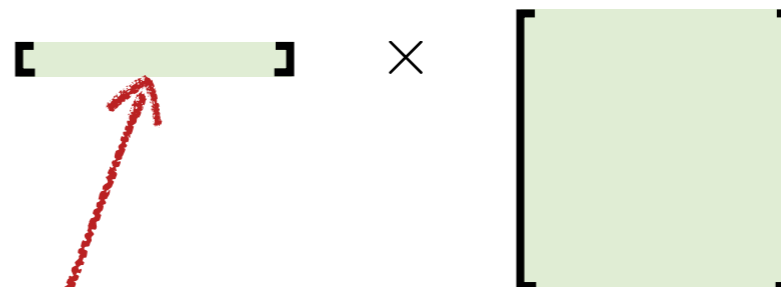


Backpropagation

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



$$\frac{df_n \circ \dots \circ f_2}{d \text{vec } \mathbf{x}_1^\top} \times \frac{d \text{vec } f_1}{d \text{vec } \mathbf{x}_0^\top}$$

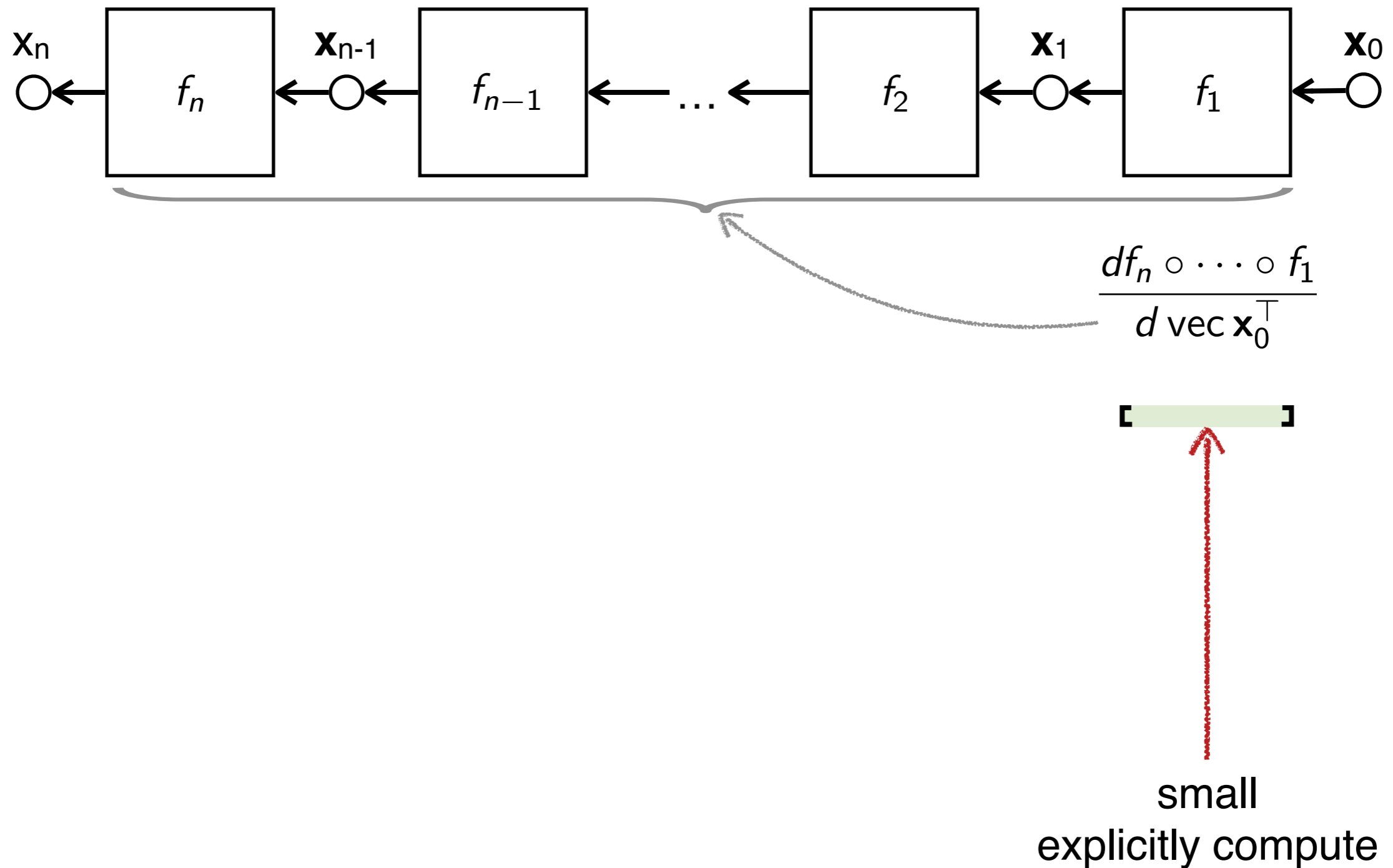


small
explicitly compute

uber matrix
do **not** explicitly compute

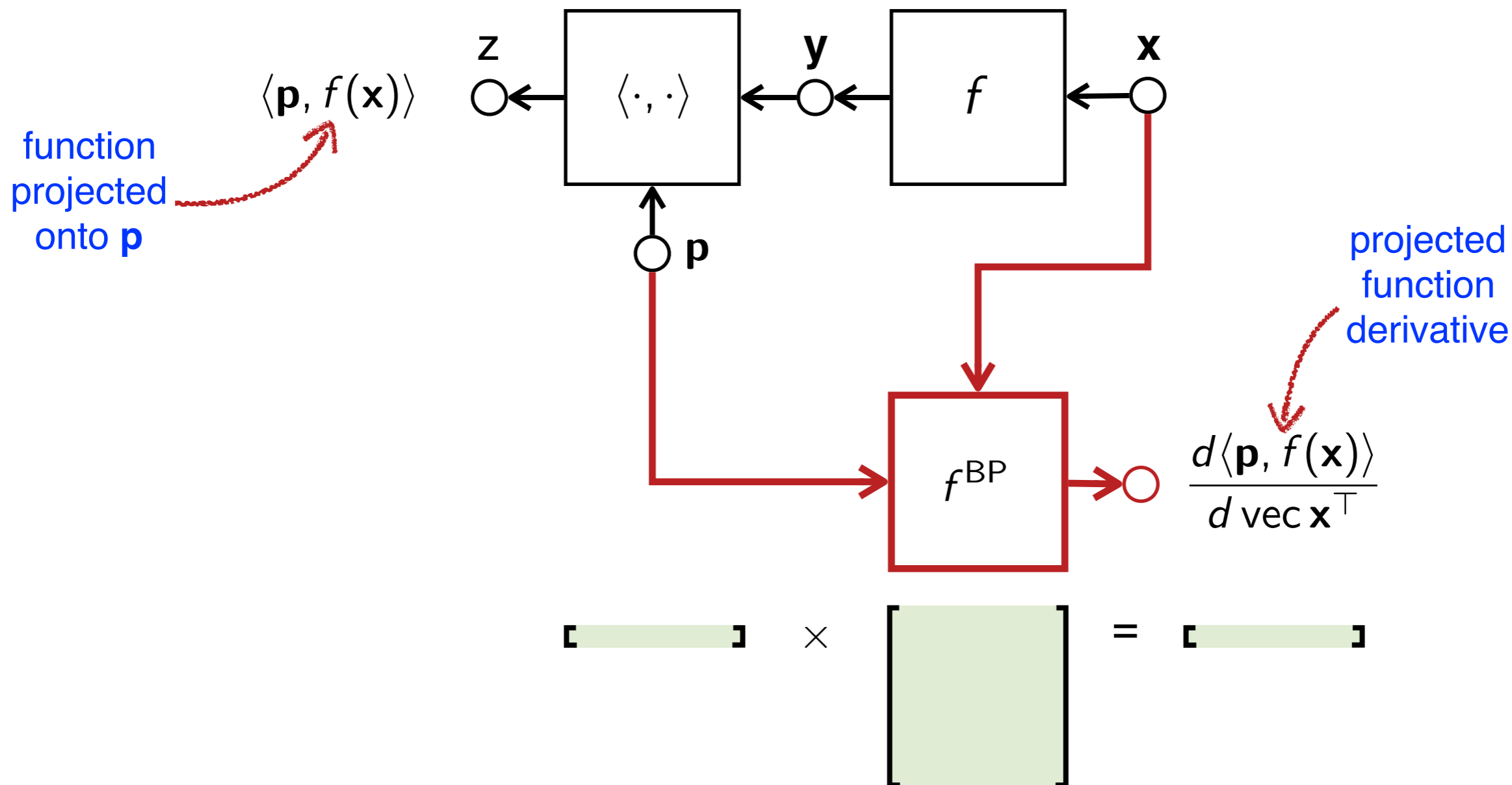
Backpropagation

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



Projected function derivative

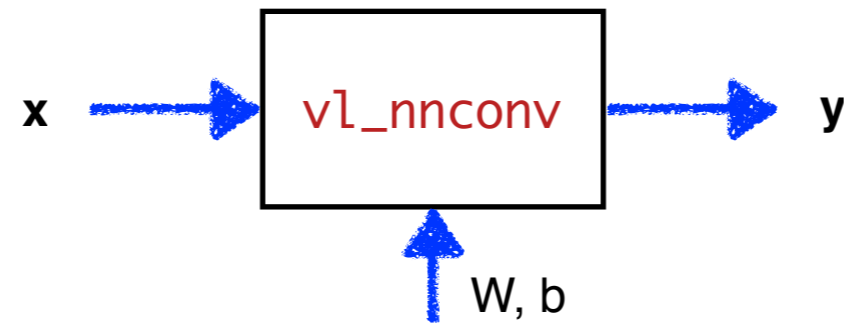
The “BP-reversed” layer



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

Anatomy of a building block

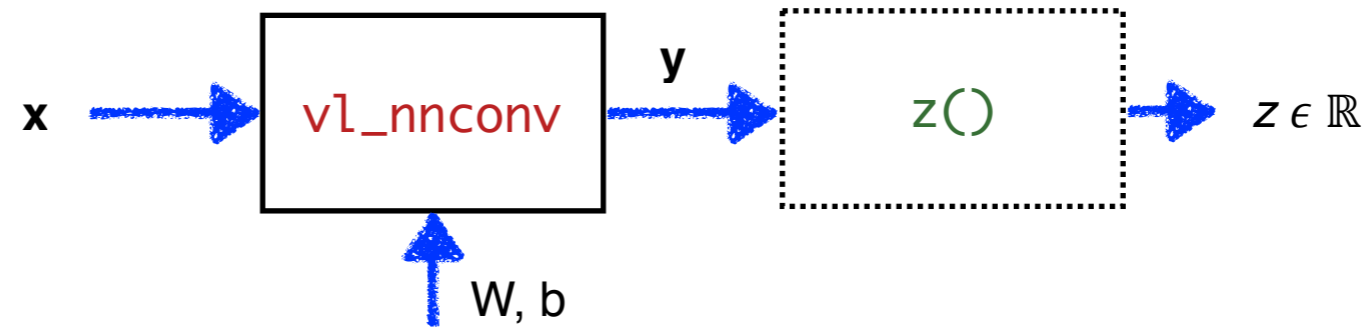
forward (eval)



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

Anatomy of a building block

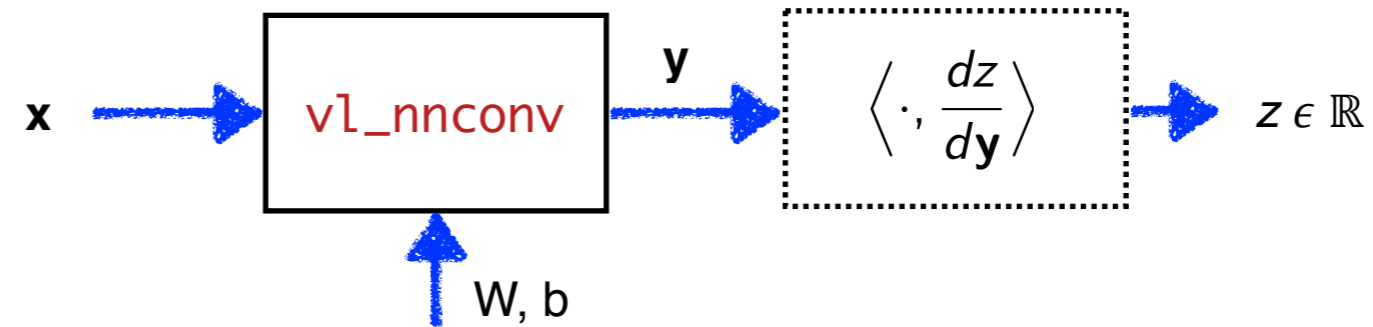
forward (eval)



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

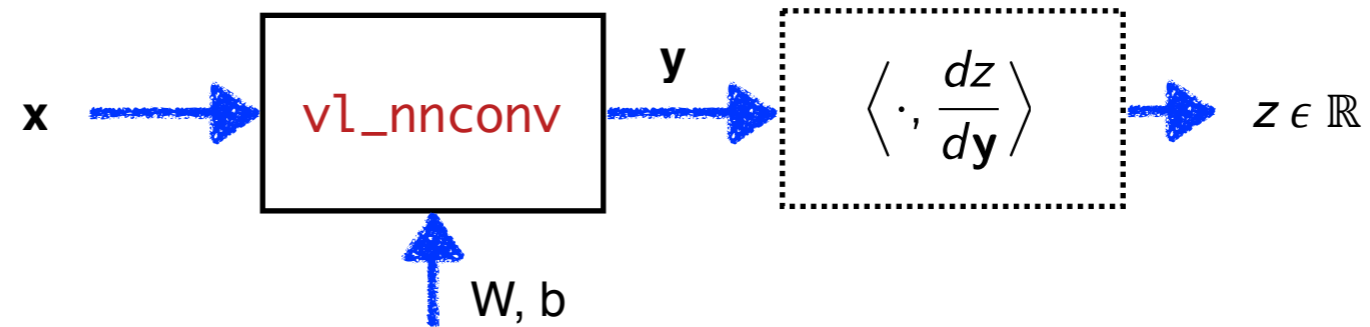
Anatomy of a building block

forward (eval)



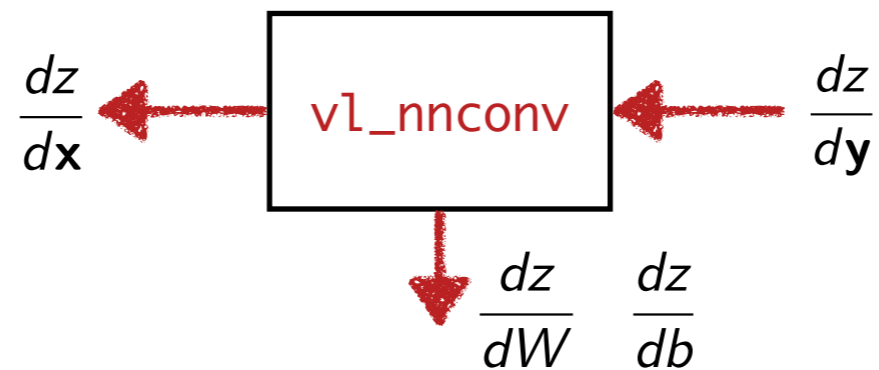
$$y = \text{vl_nnconv}(x, W, b)$$

forward (eval)



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

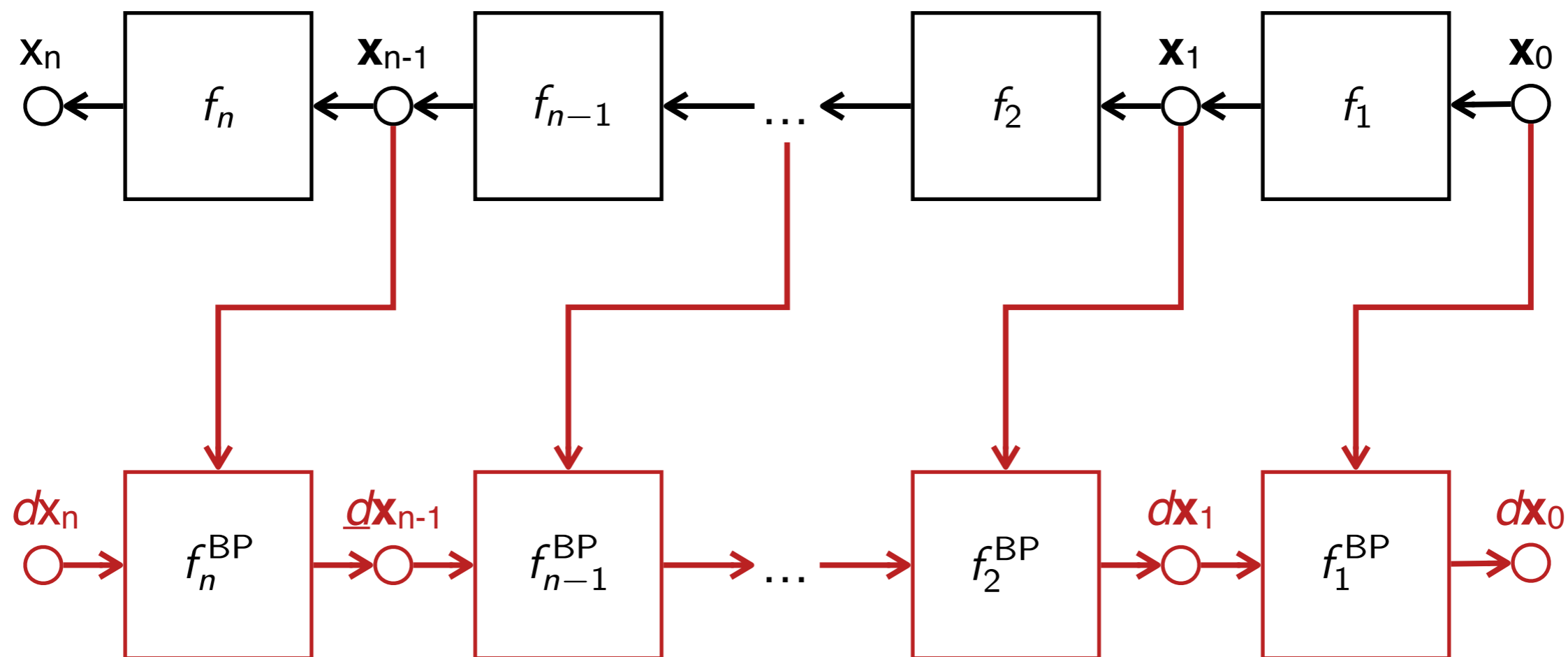
backward (backprop)



$$\text{dzdx} = \text{vL_nnconv}(x, W, b, \text{dzdy})$$

Backpropagation network

BP induces a “reversed” network

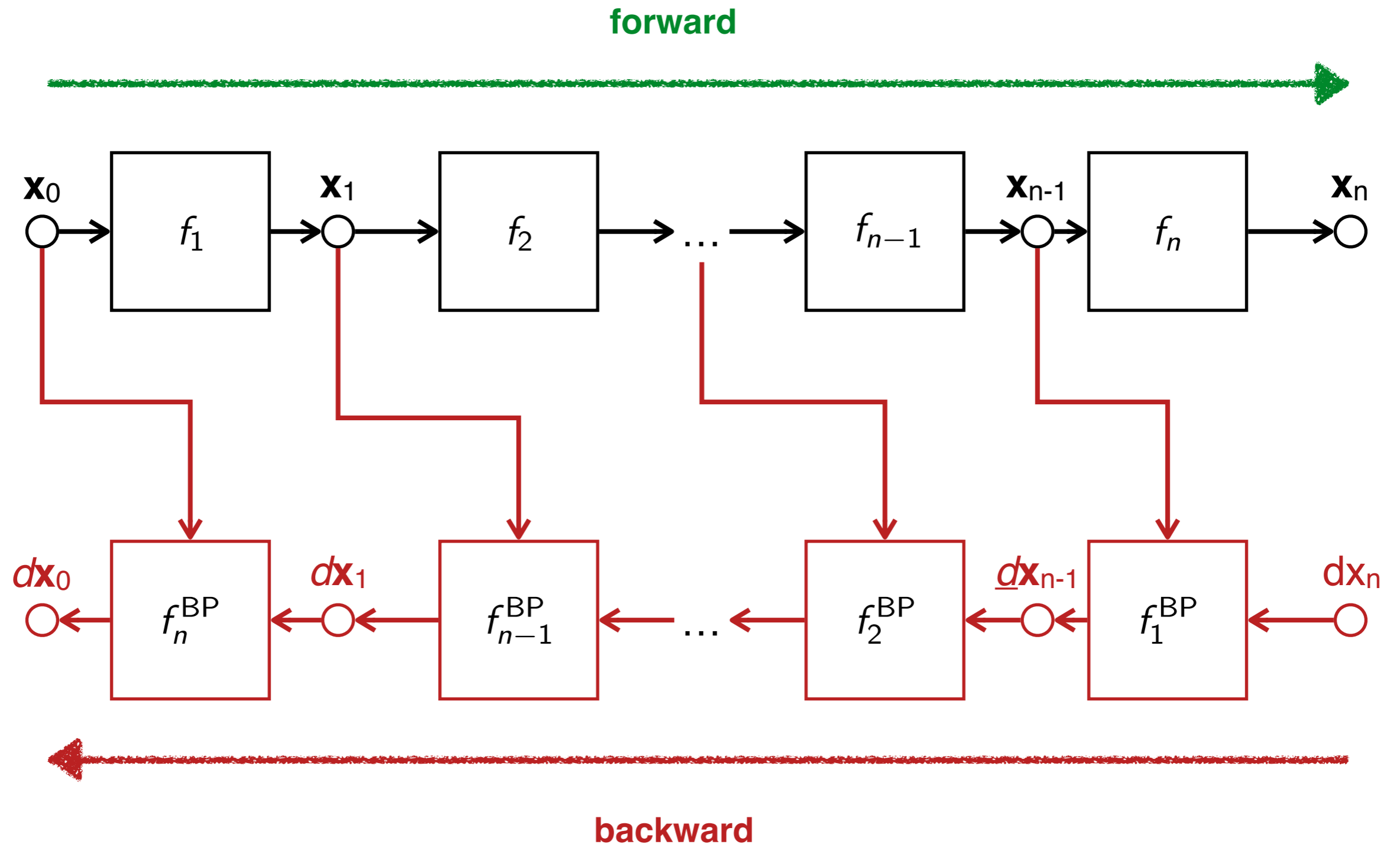


where $dx_i = \frac{df_n \circ \dots \circ f_{i+1}}{d \text{vec } x_i}$

Note: the BP network is linear in $dx_1, \dots, dx_{n-1}, dx_n$. Why?

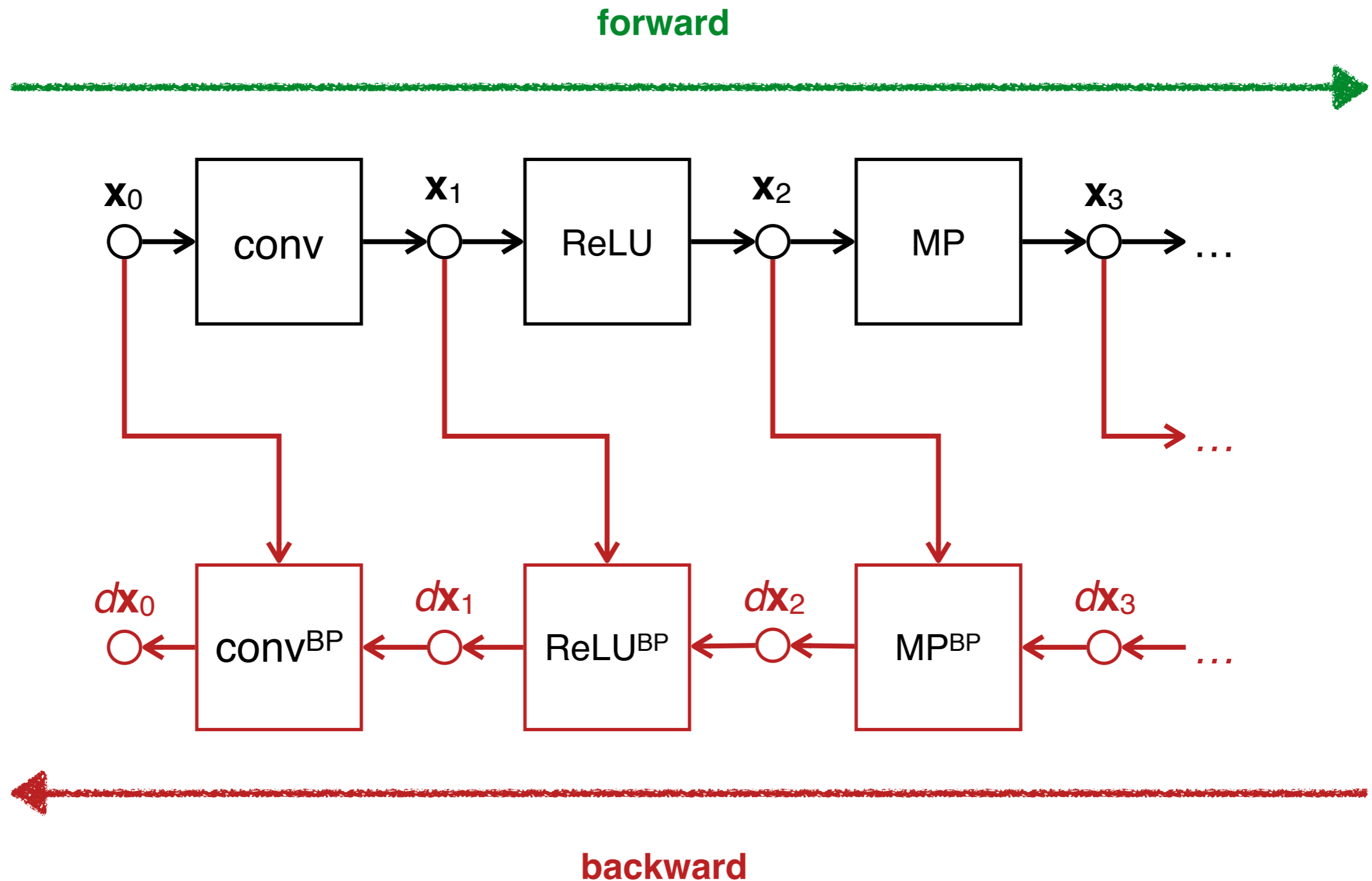
Backpropagation network

BP induces a “transposed” network



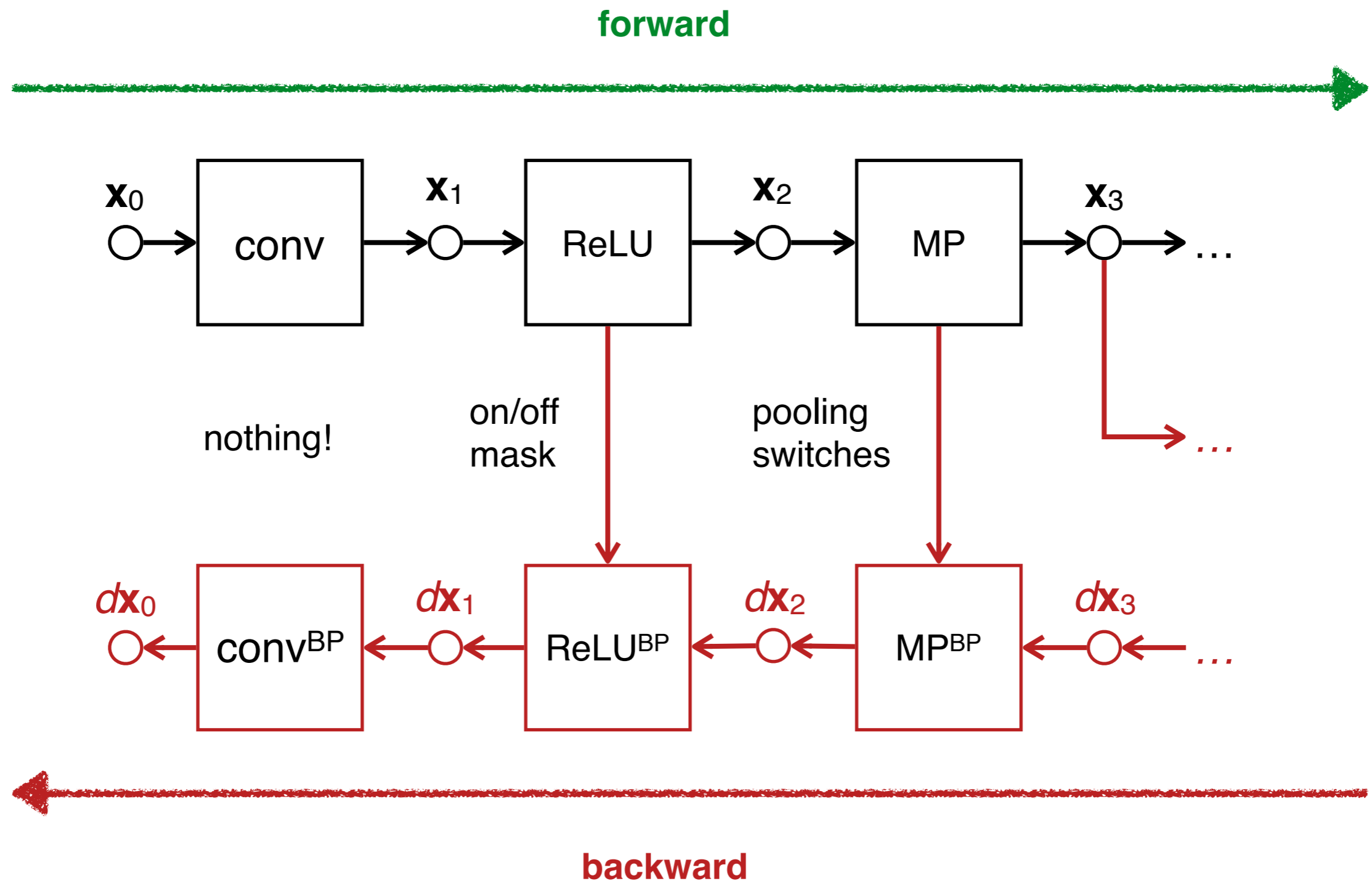
Backpropagation network

Conv, ReLU, MP and their transposed blocks



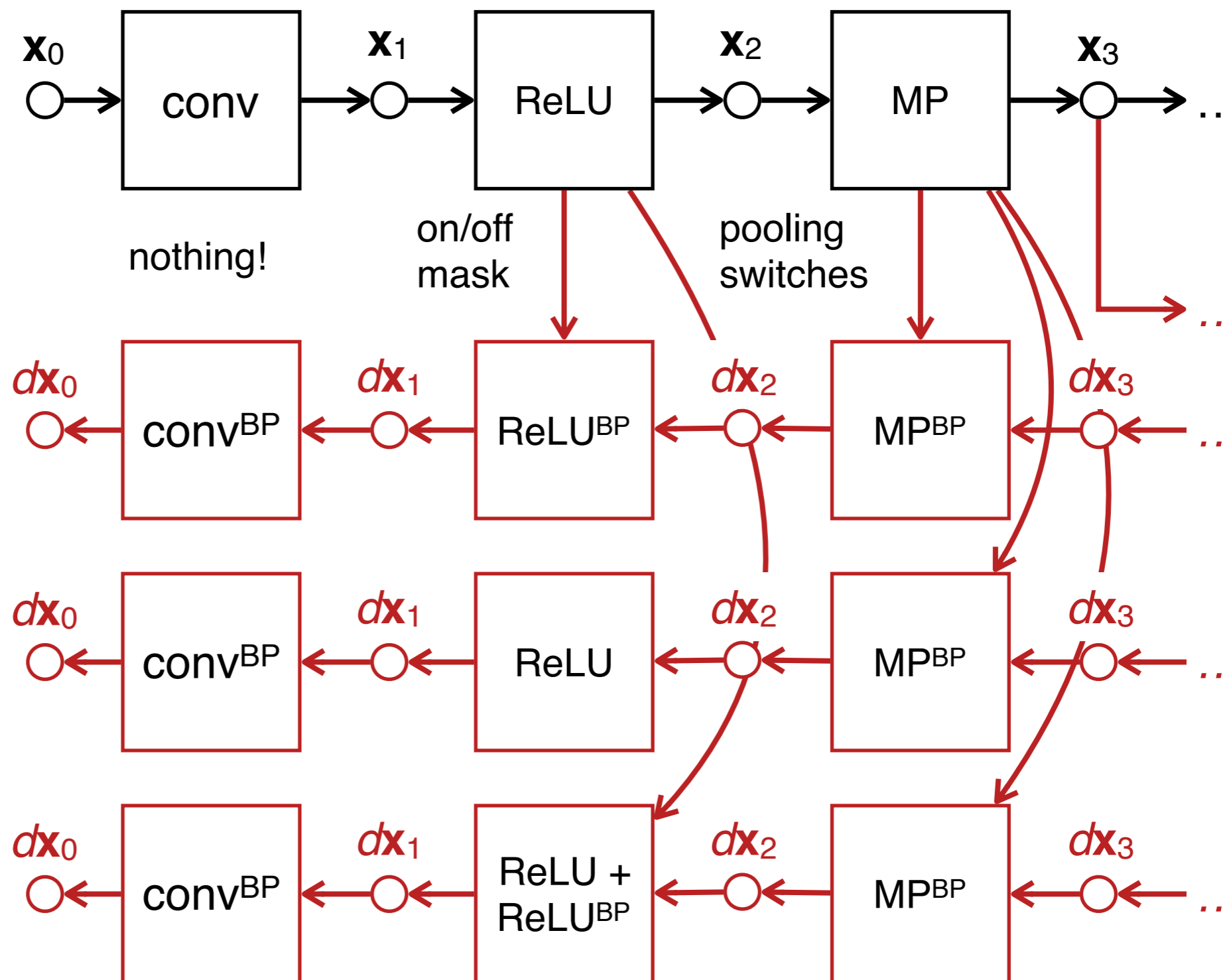
Sufficient statistics and bottlenecks

Usually much less information is needed



Three visualisation techniques

Modified backpropagation networks



SaliNet

Equiv. to
[Simonyan et al. 14]

DeConvNet

Same as
[Zeiler Fergus 14]

DeSaliNet

[Mahendran Vedaldi 16,
Springenberg et al. 15]

Results

———— VGG-VD Pool5_3 ————

———— VGG-VD FC8 ————

Image

DeConvNet

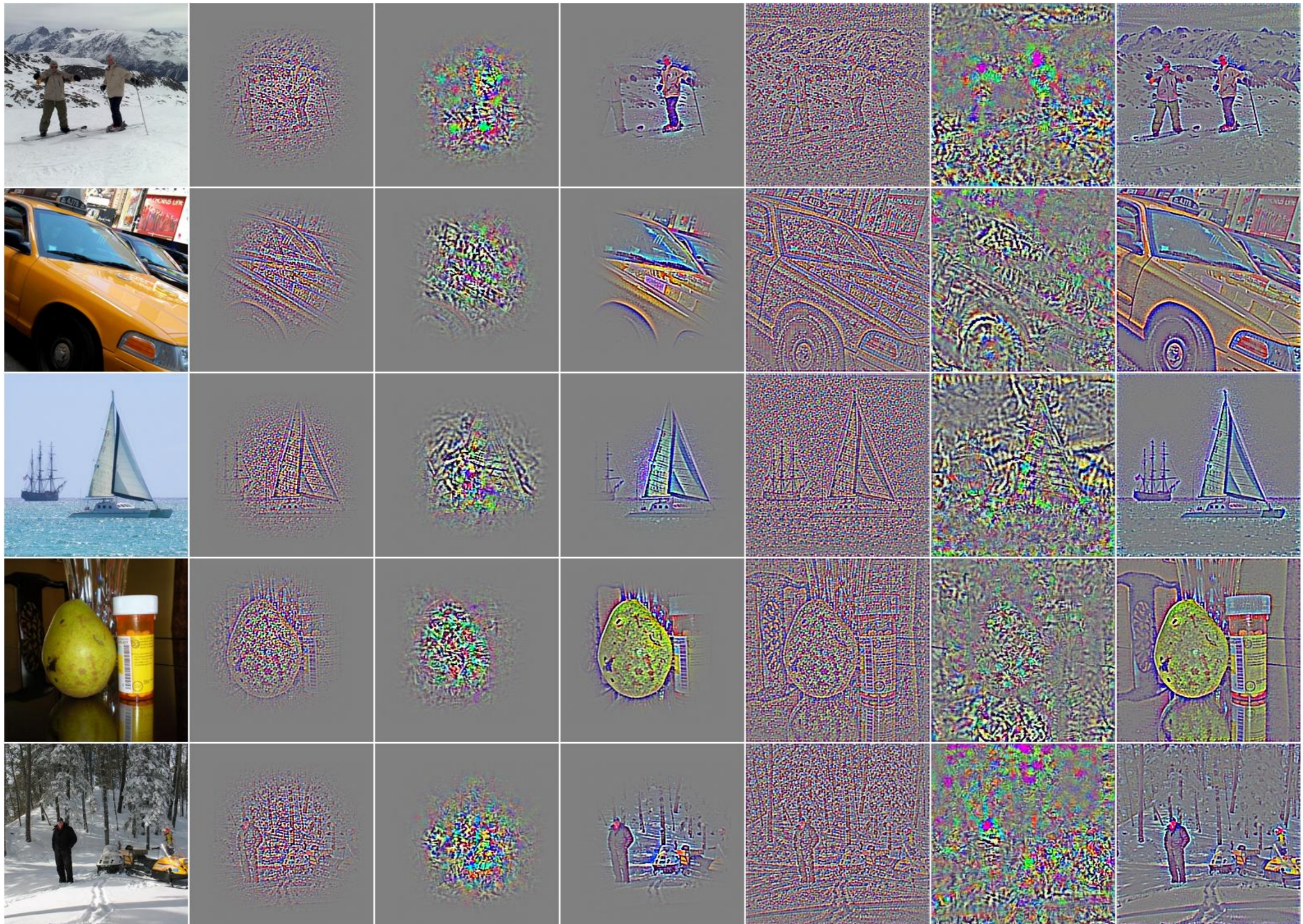
SaliNet

DeSaliNet

DeConvNet

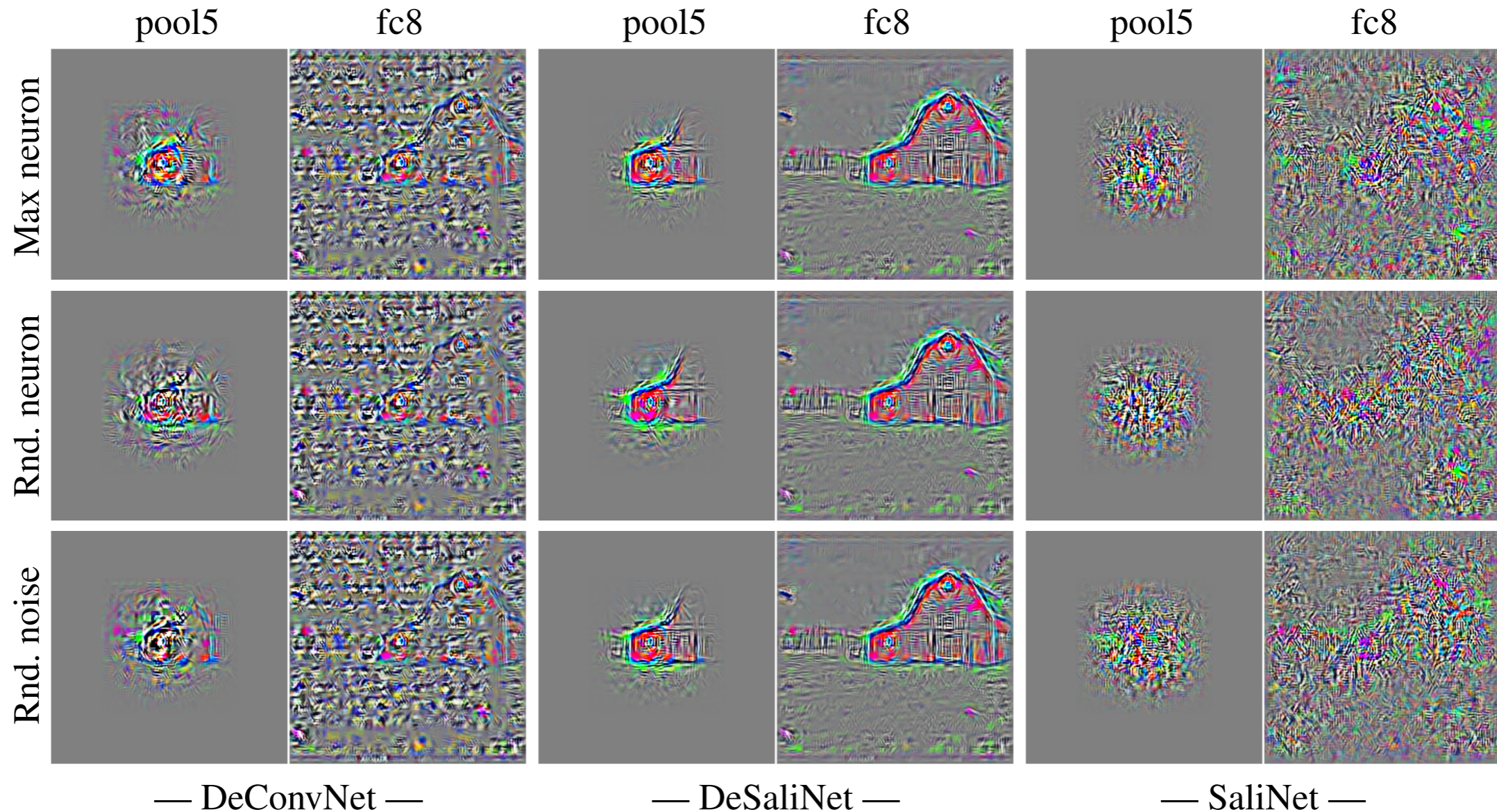
SaliNet

DeSaliNet



Key limitation of DeConvNets and similar

They are largely *not* neuron selective



Good for saliency

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

Visualizing representations

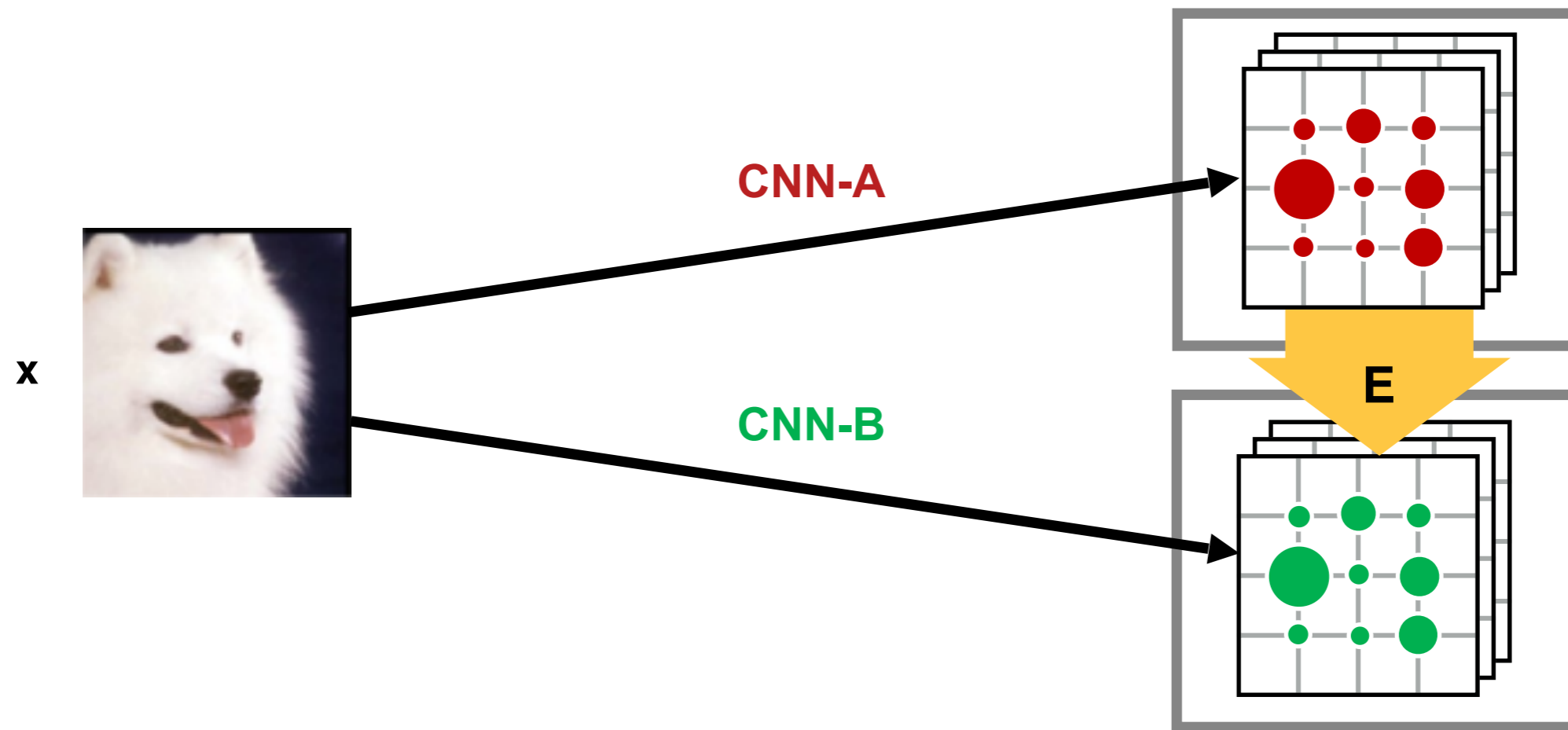
Backpropagation networks and “deconvolution”

Representations: equivalence & transformations

When are two representations the same?

Learning representations means that there is an endless number of them

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

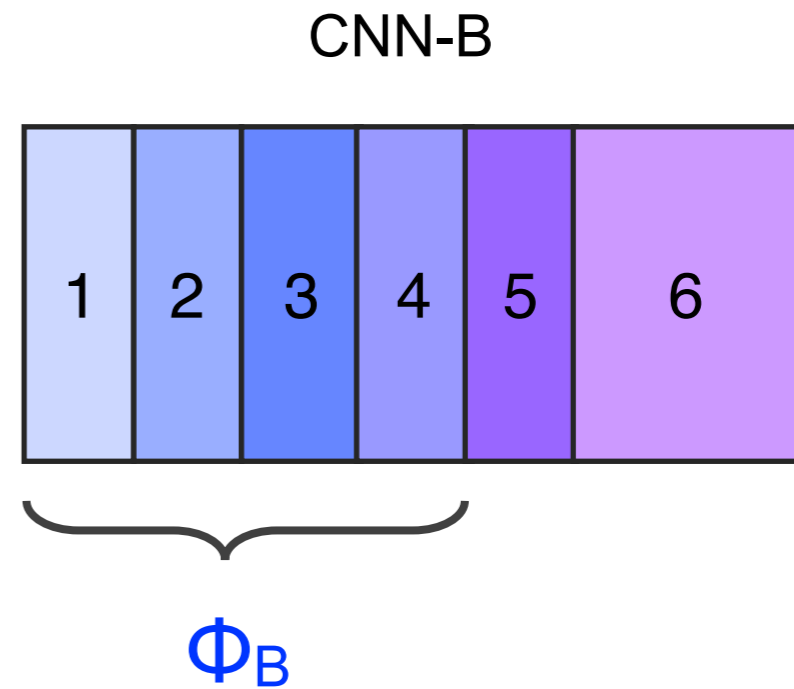
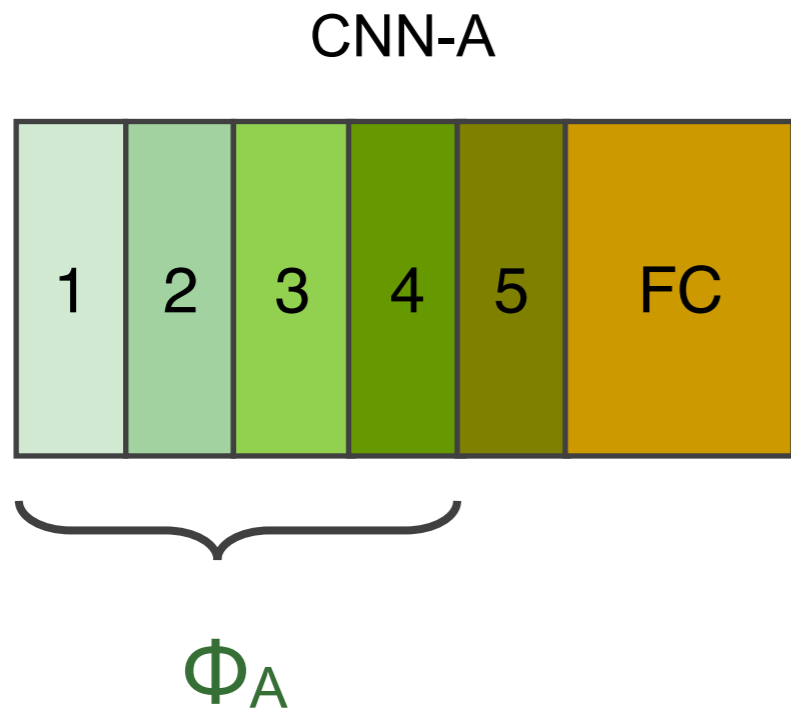


Equivalence

$$\Phi_B(\mathbf{x}) = E \Phi_A(\mathbf{x})$$

Equivalence

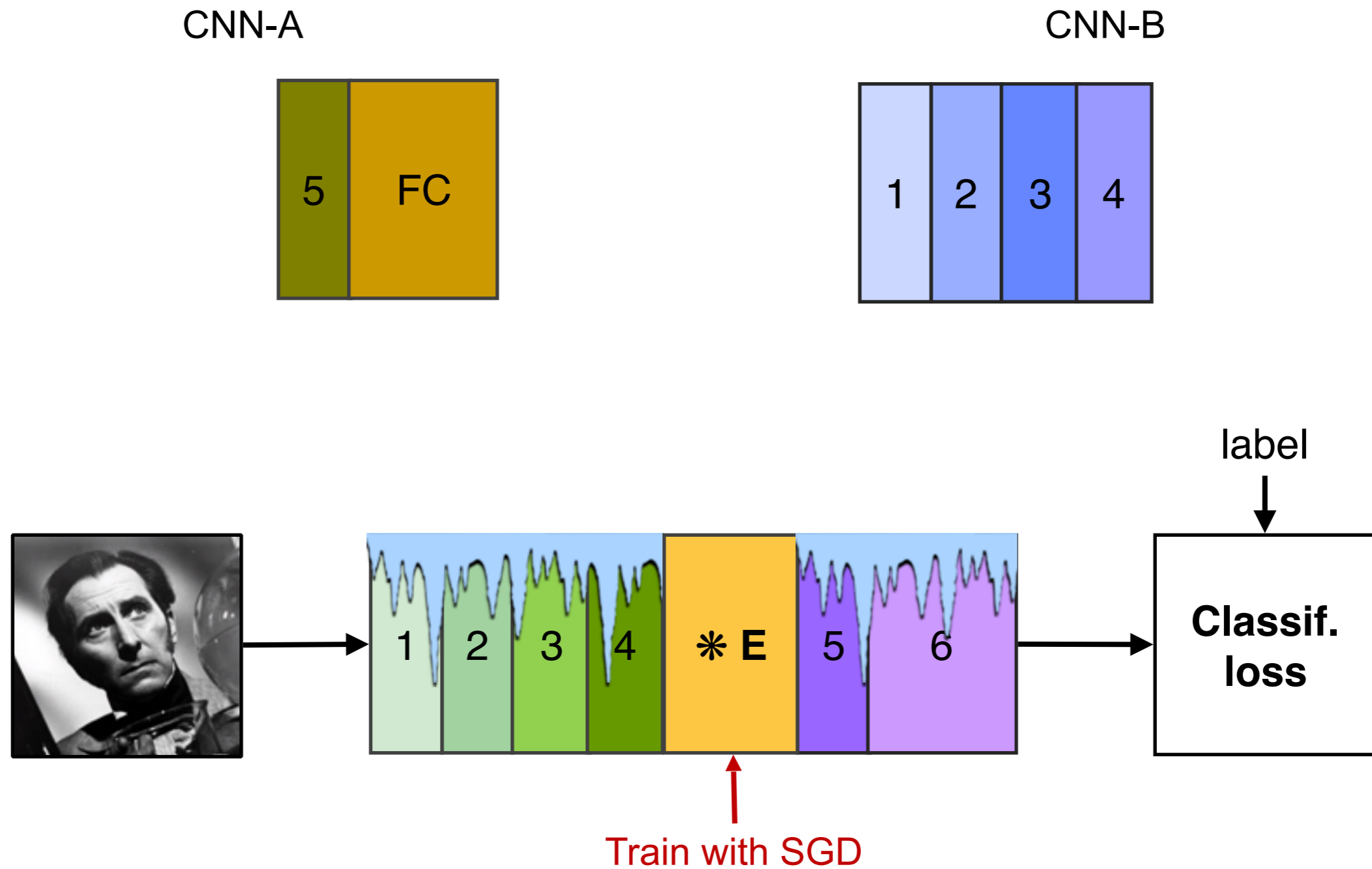
AlexNet, same training data, different parametrization:



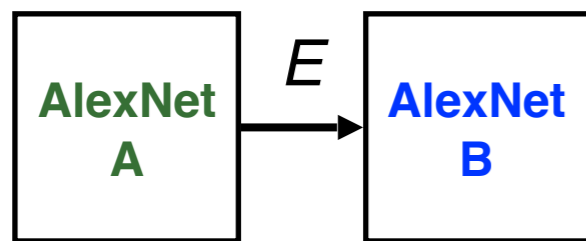
Are Φ_A and Φ_B equivalent ?

Equivalence

AlexNet, same training data, different parametrization:



Equivalence with different random seeds



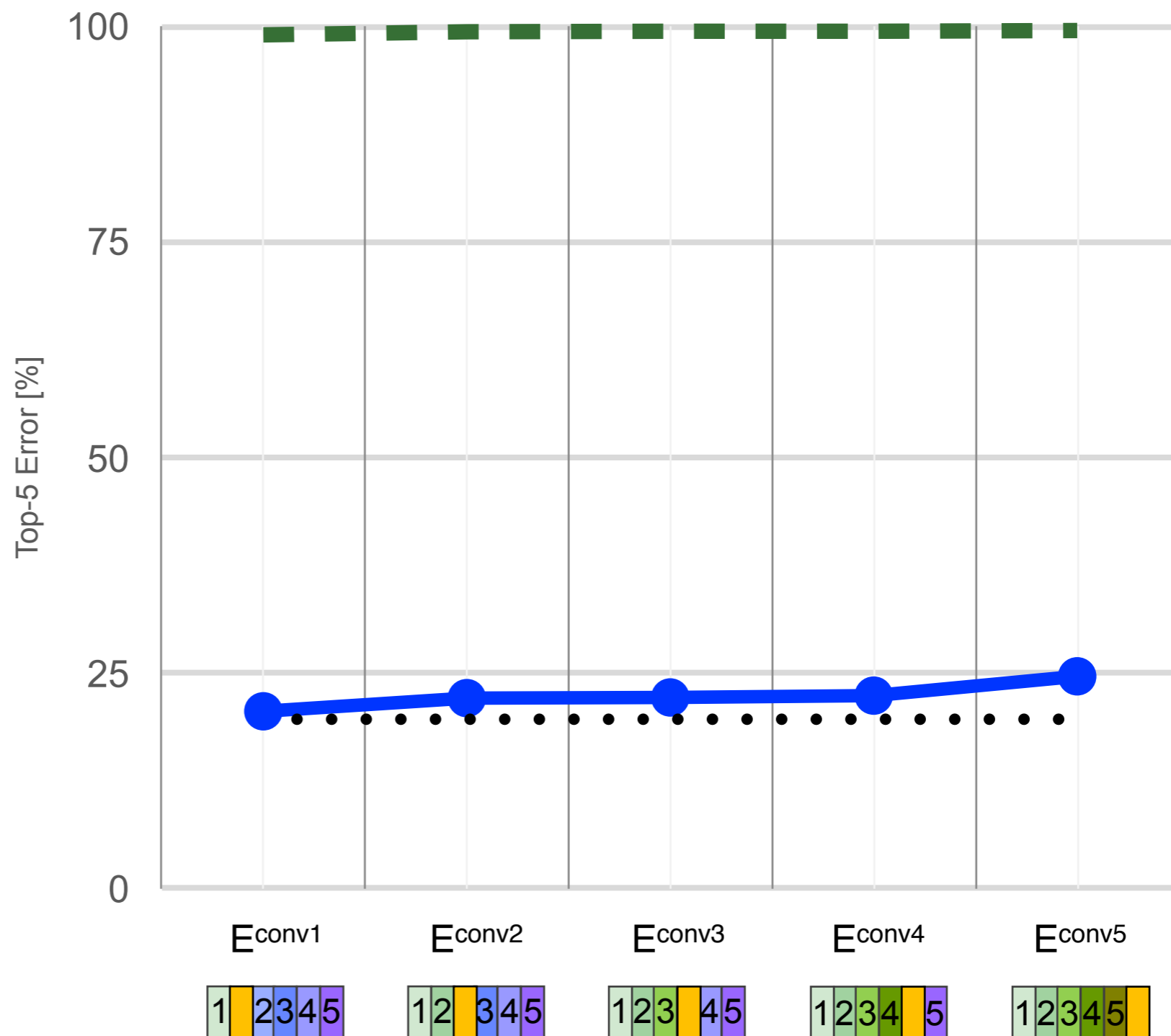
Baseline



Before training



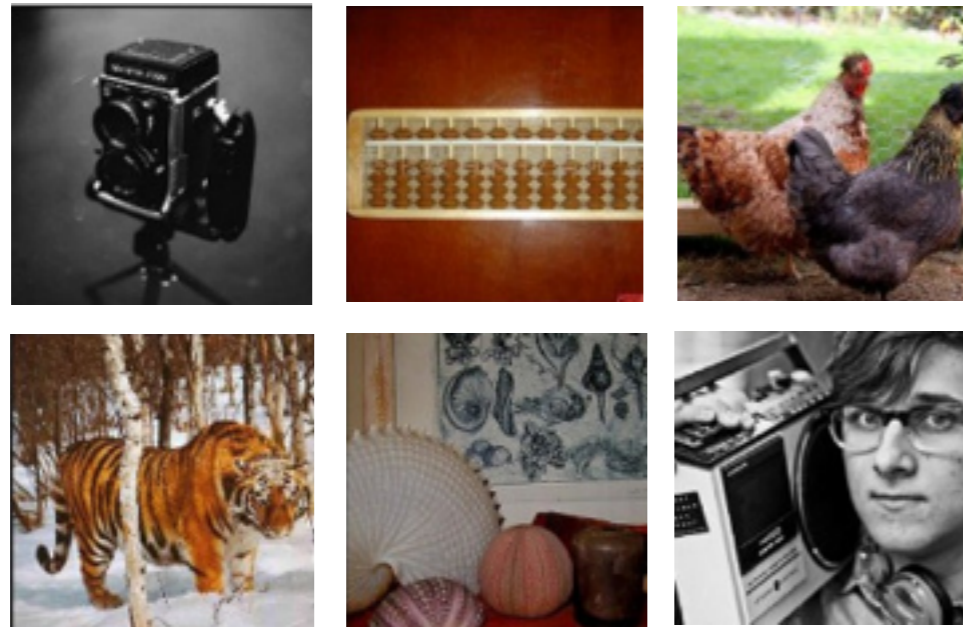
After training



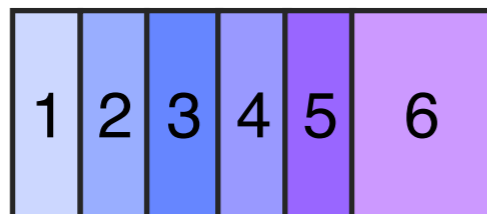
Equivalence of similar architecture

Train on two different datasets

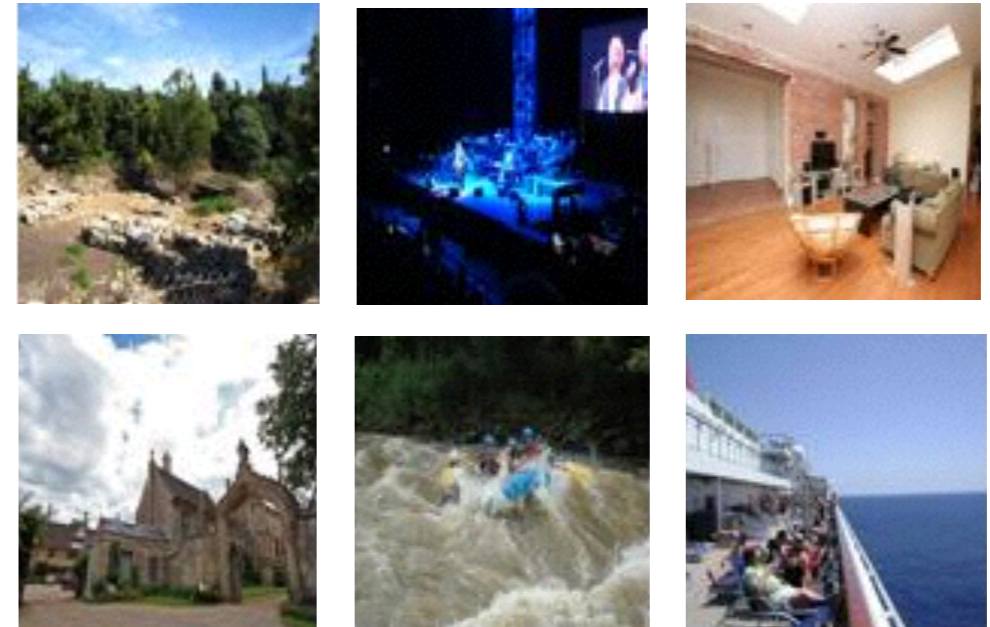
ILSVRC12 dataset



CNN-IMNET



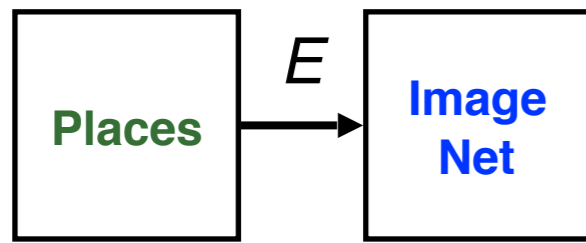
Places dataset



CNN-PLACES



Equivalence with different training data



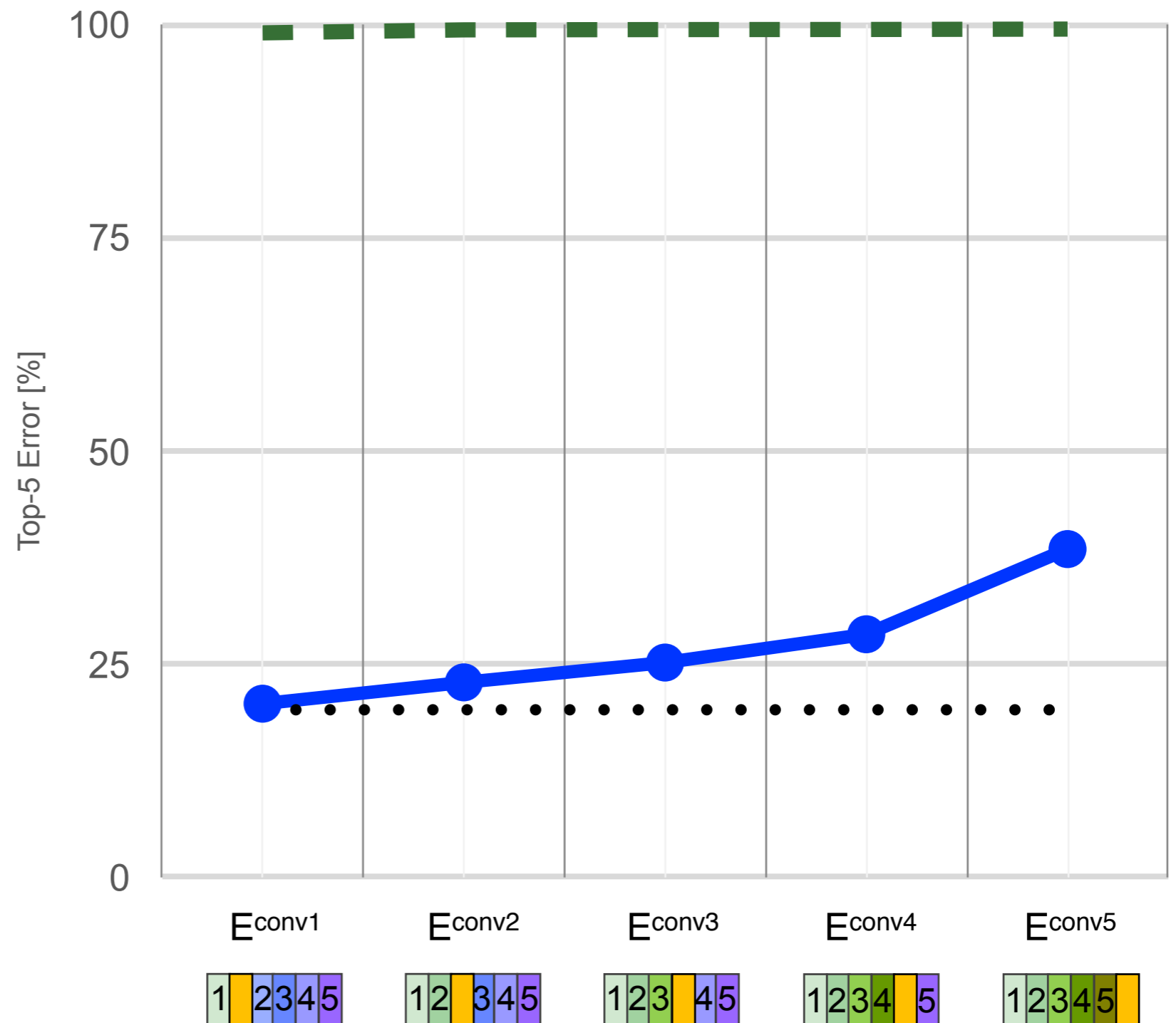
Baseline



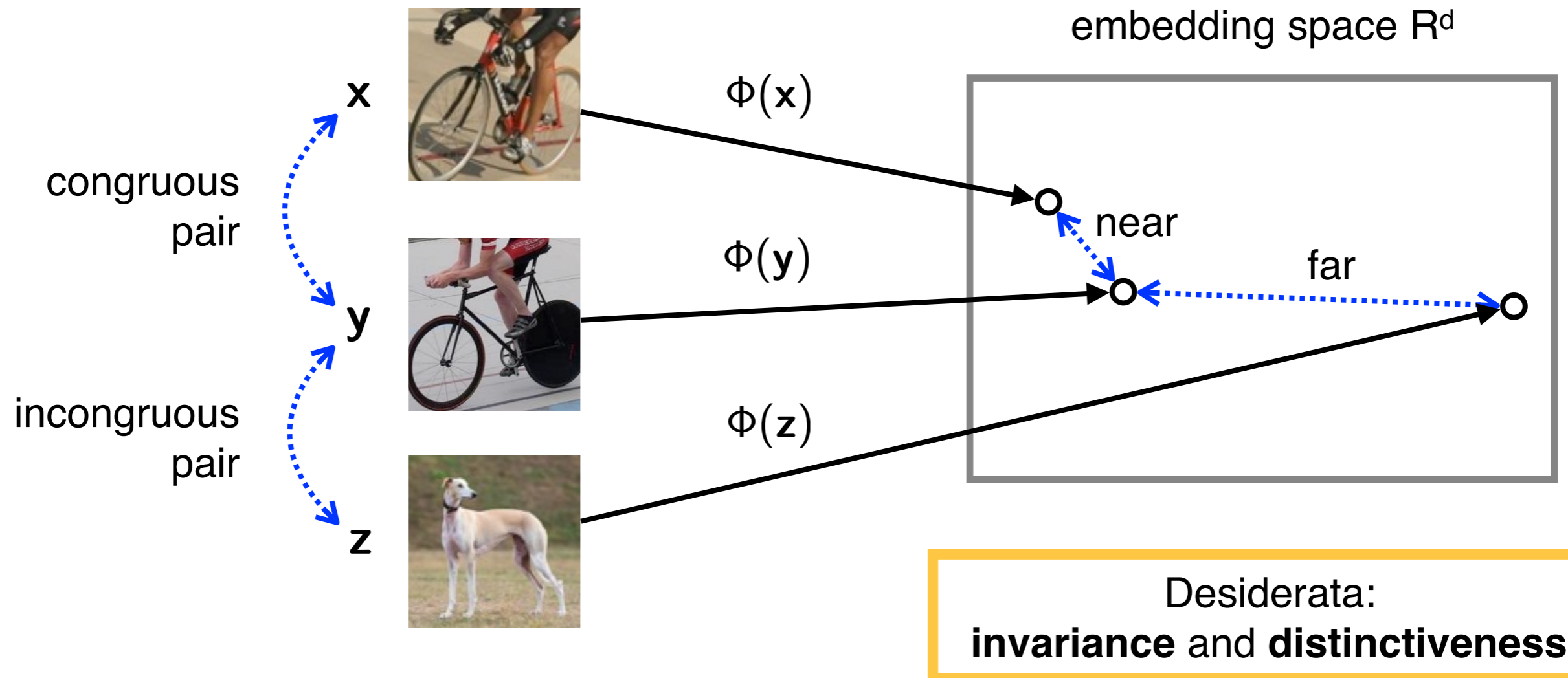
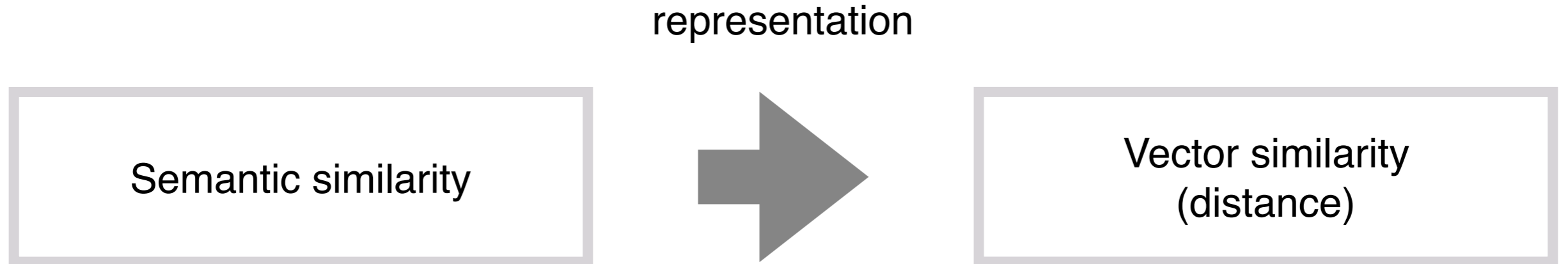
Naive stitching



Learned stitches

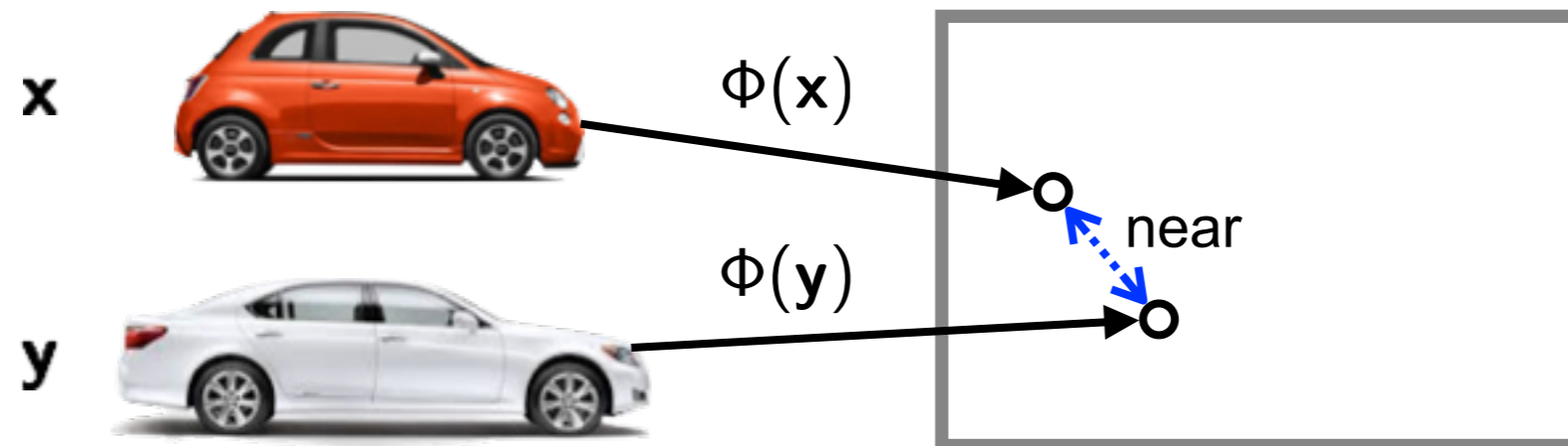


Meaningful representations



Invariance is task dependent

Are x and y the same category?



Are x and y the same colour?

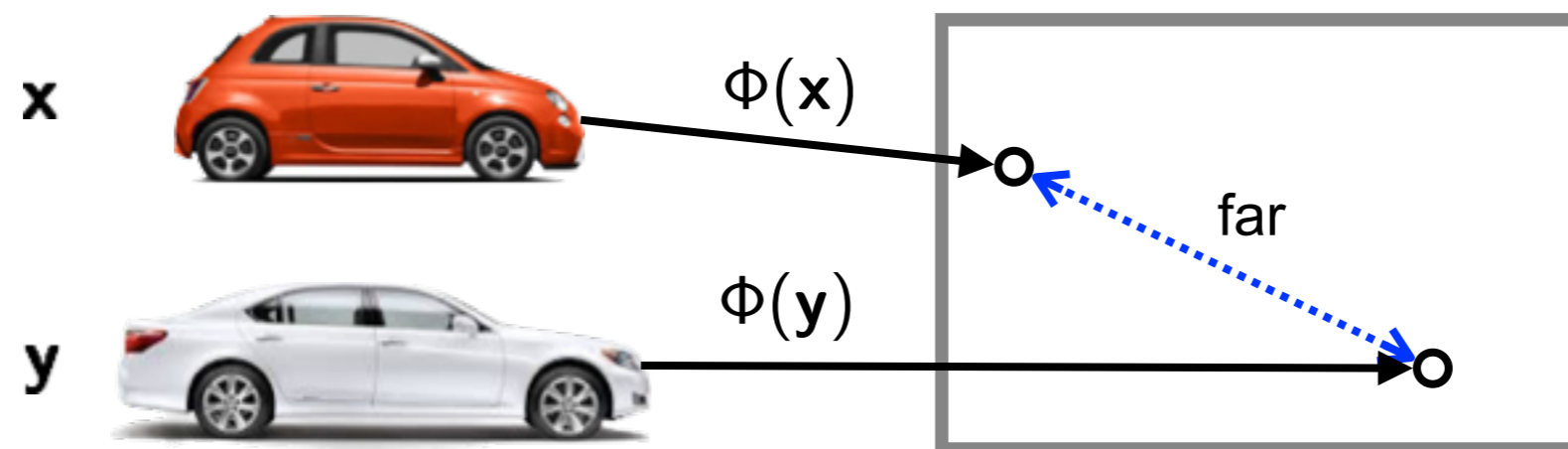
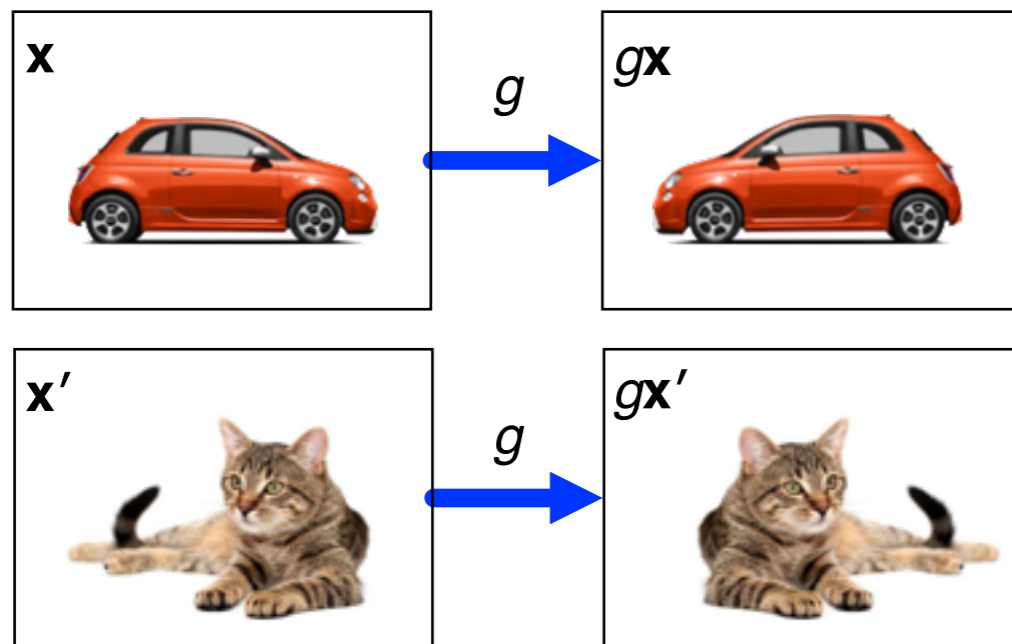
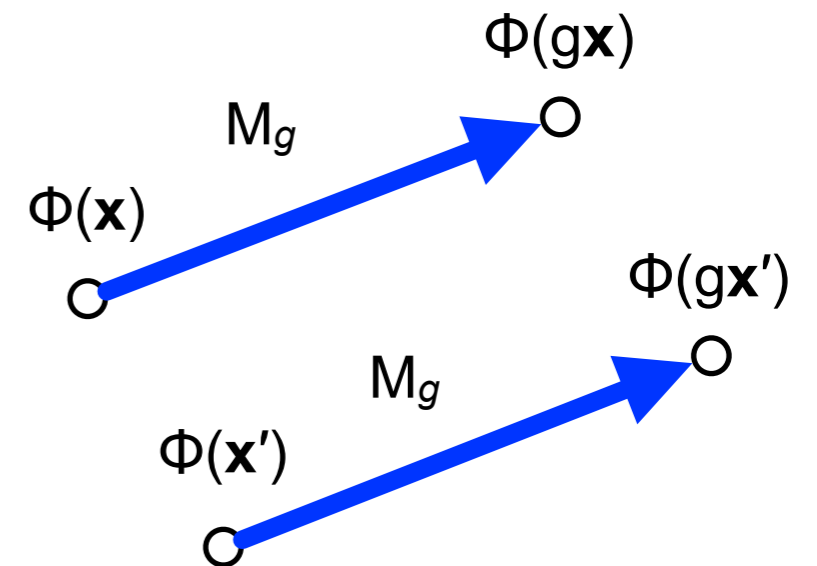


image space



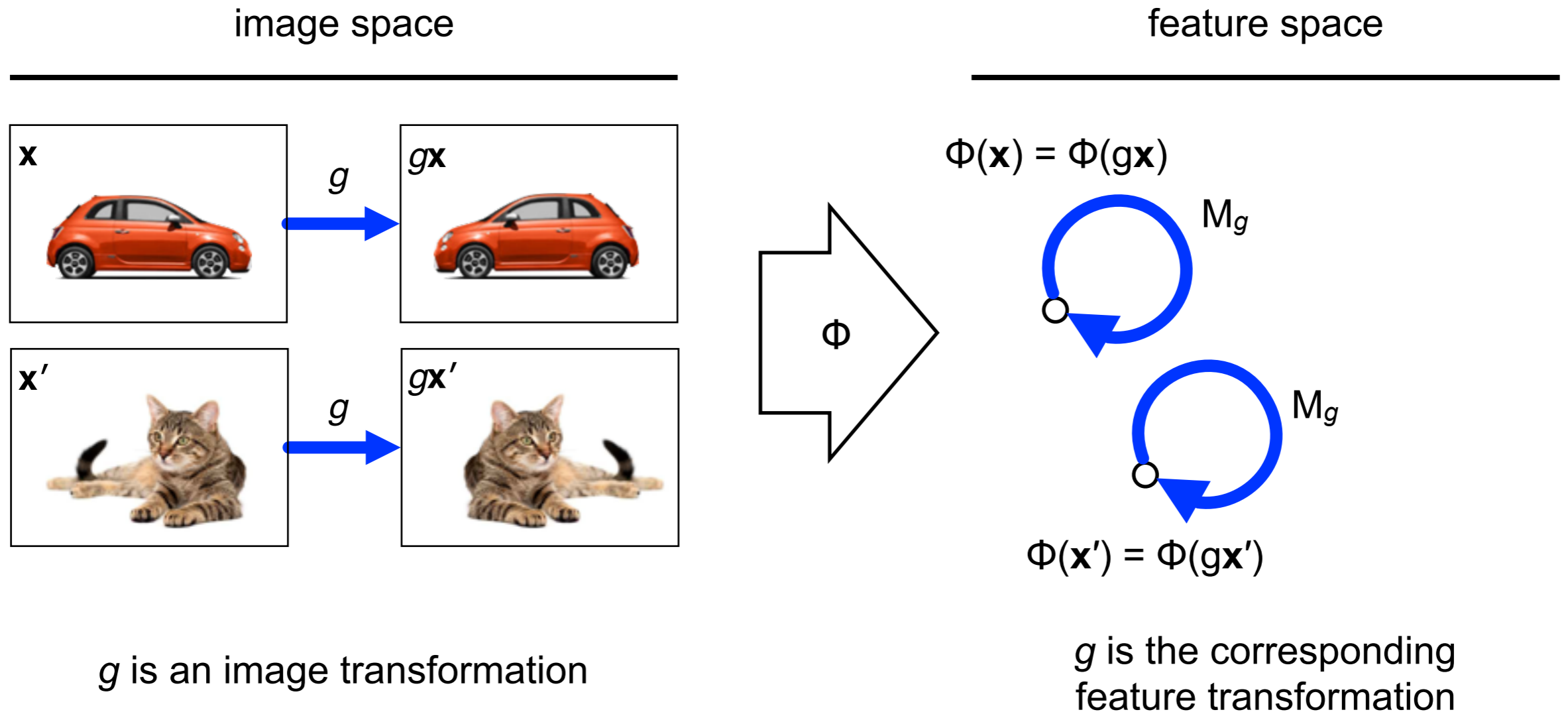
g is an image transformation

feature space



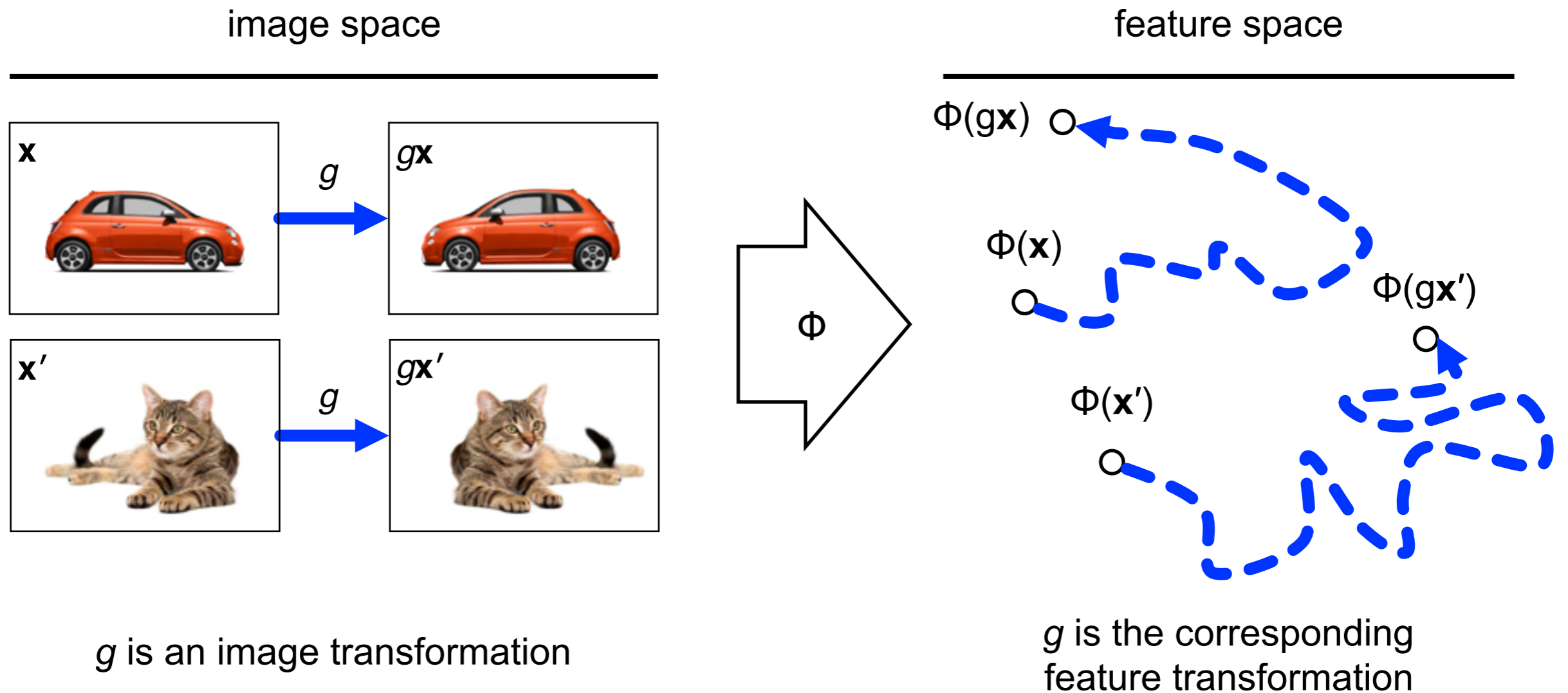
g is the corresponding feature transformation

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



the map M_g is the identity

No equivariance



the map M_g does not exist (or is intractable)

image space

feature space

g is an image transformation

Equivariance

There is a (simple) M_g

Invariance

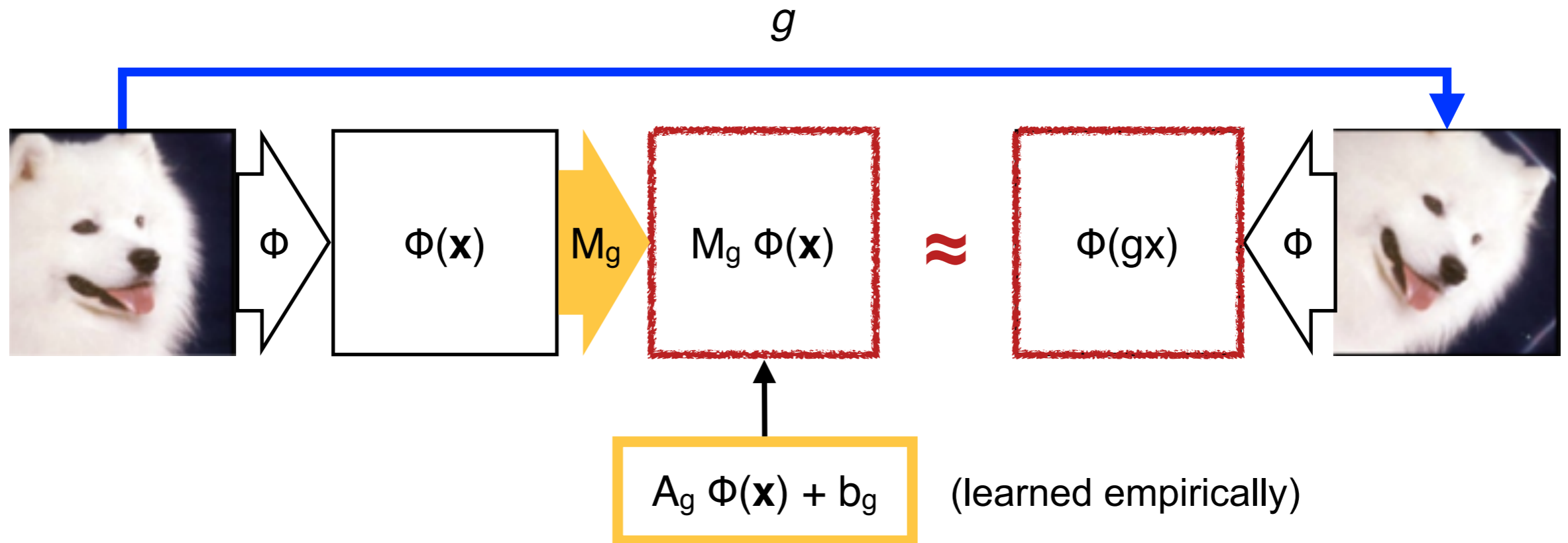
M_g is the identity

Neither

There is no (simple) M_g

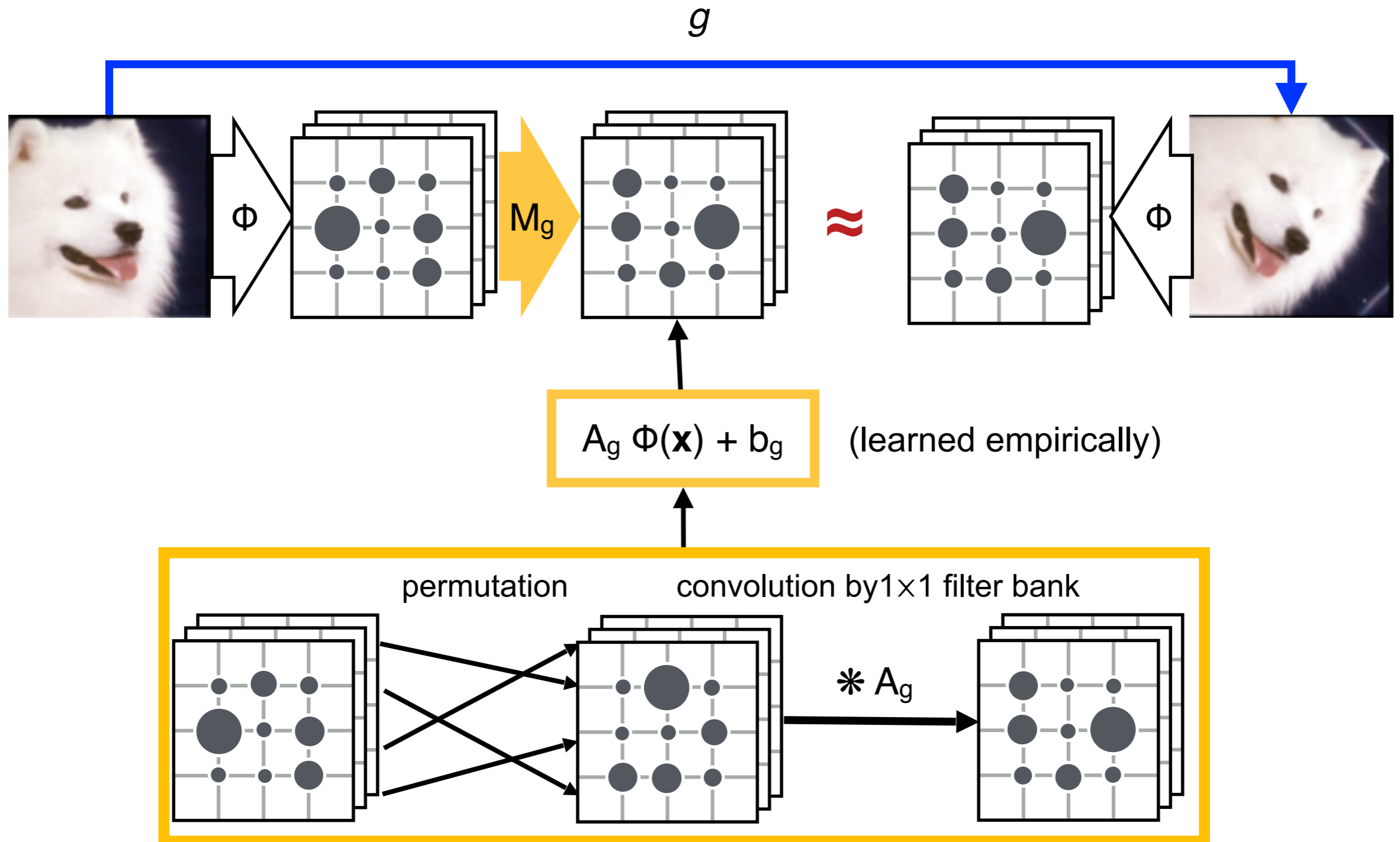
$$\Phi(g\mathbf{x}) = M_g \Phi(\mathbf{x})$$

Regularized linear regression



An empirical test of equivariance

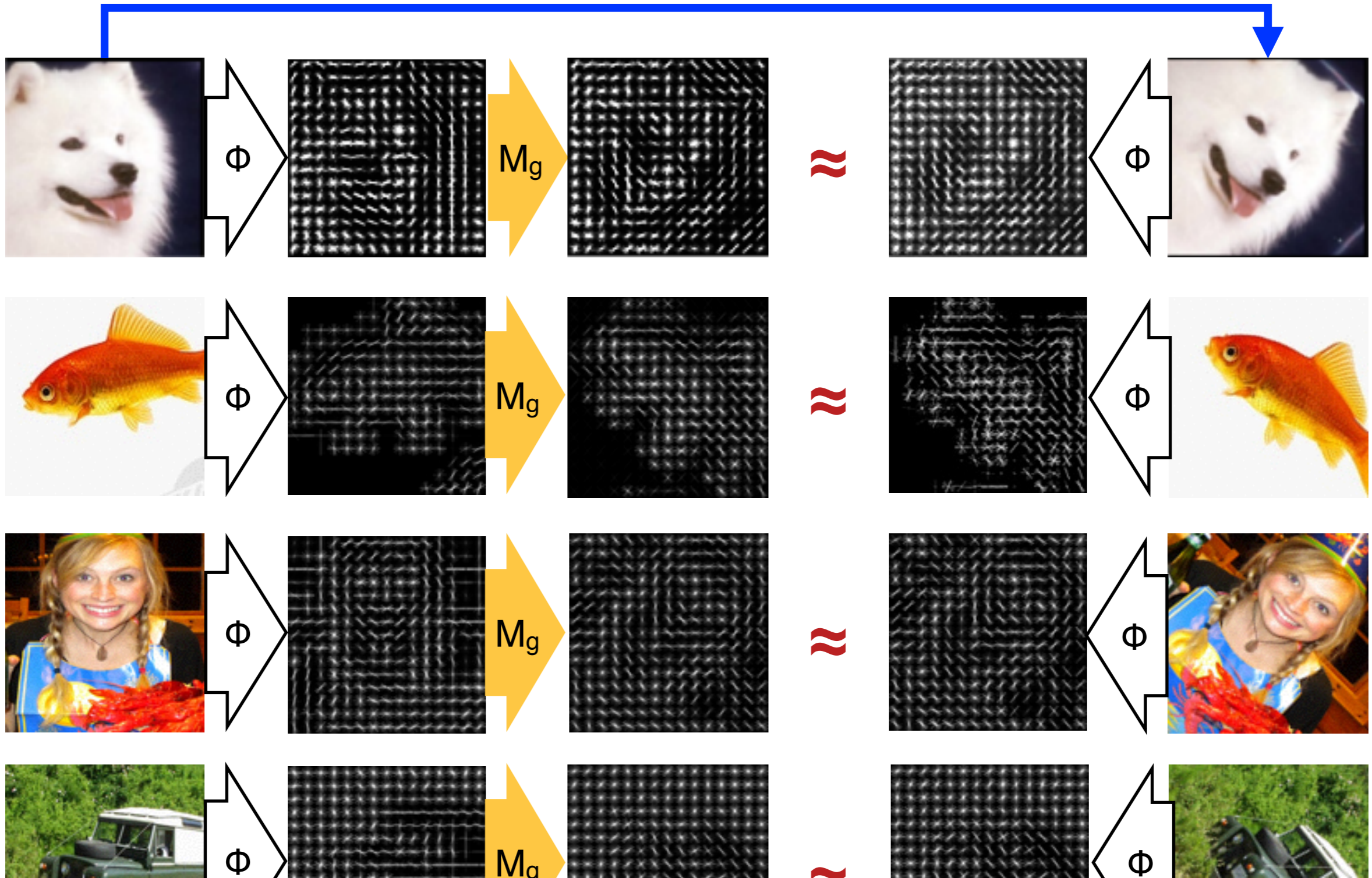
Regularized linear regression



An empirical test of equivariance

HOG features

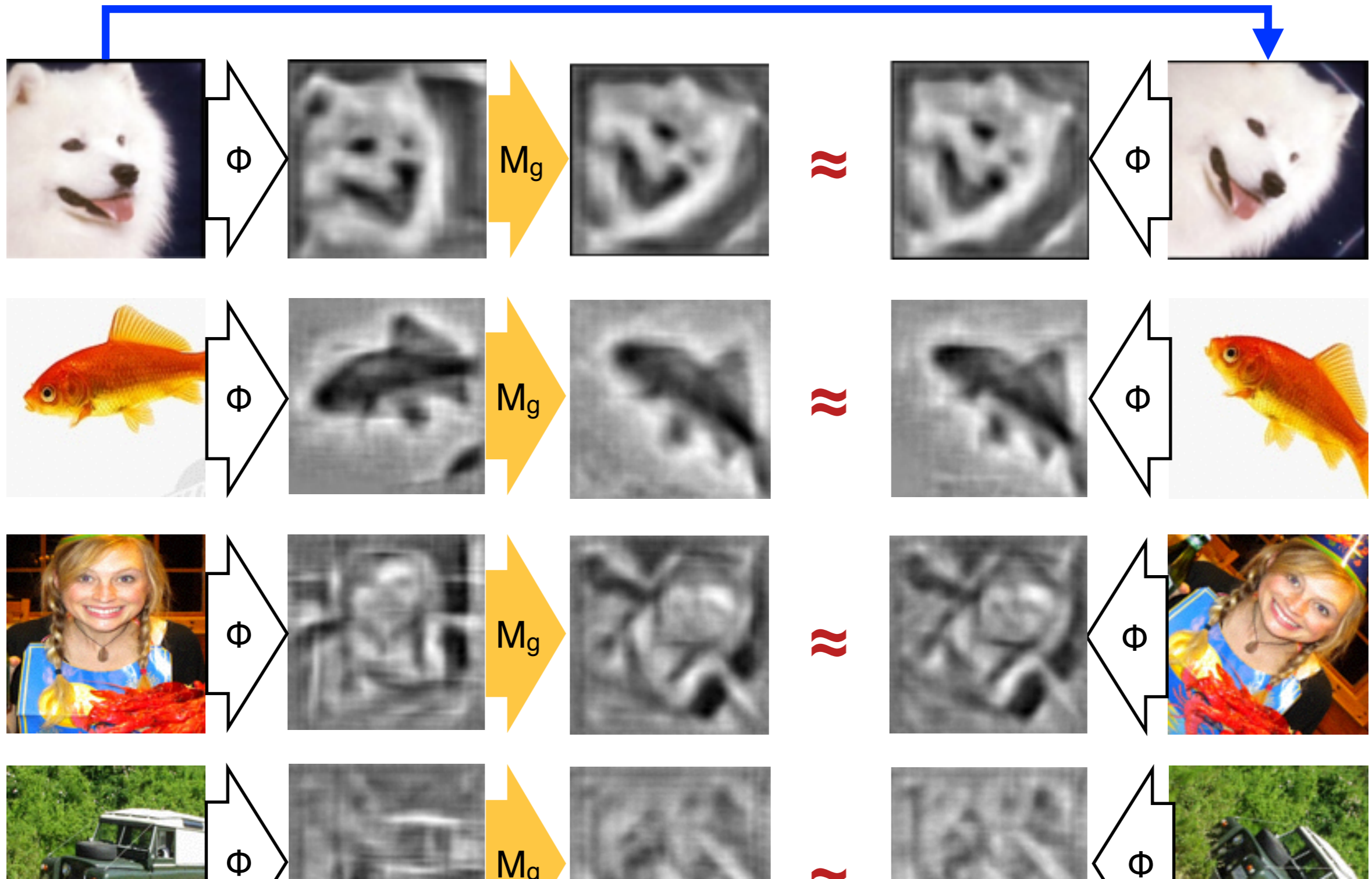
rotation 45 deg



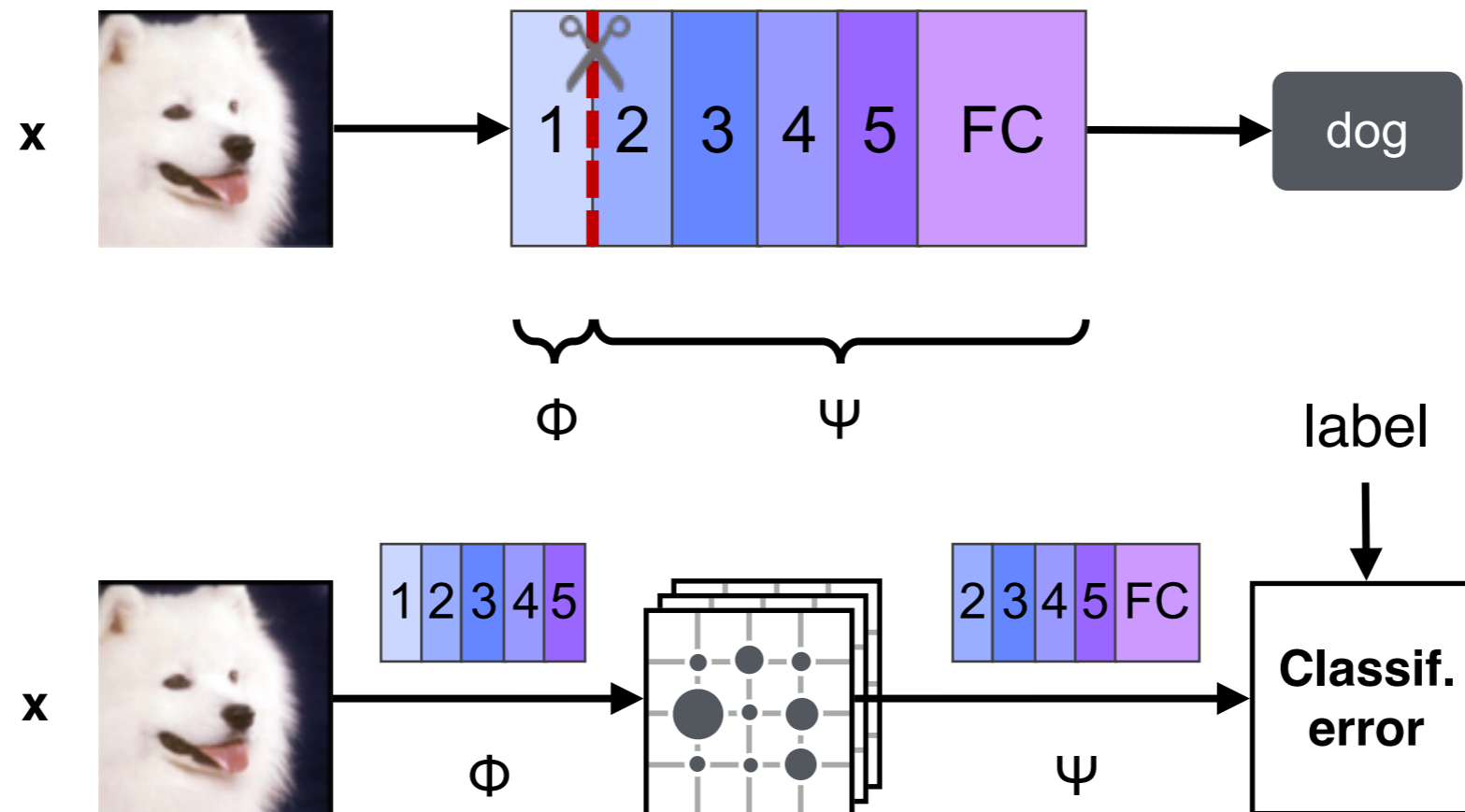
An empirical test of equivariance

HOG features

rotation 45 deg



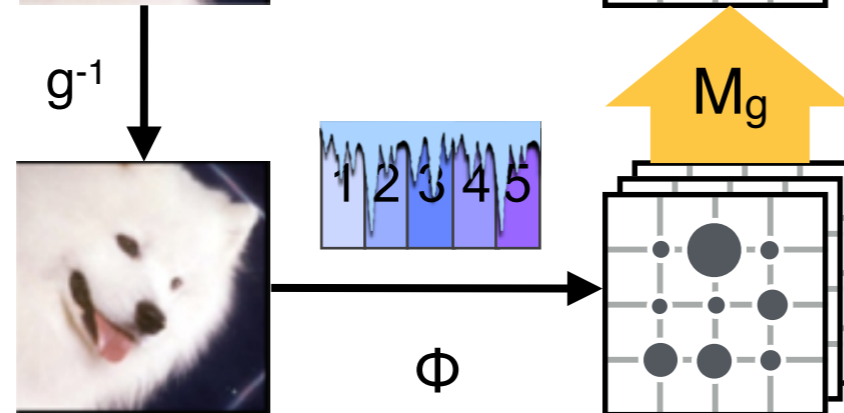
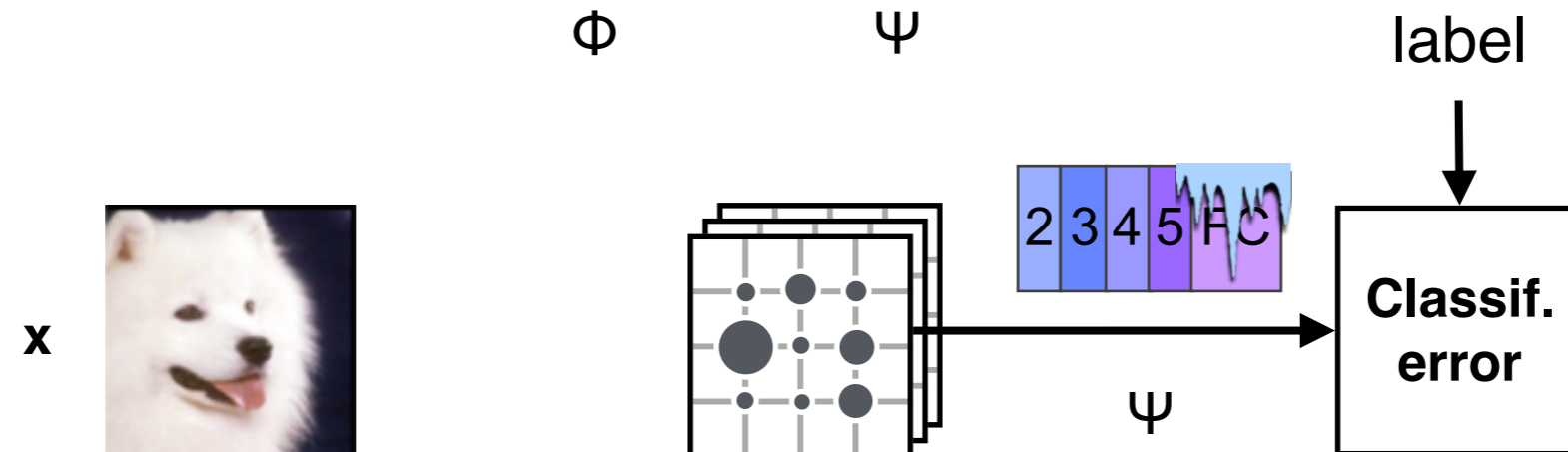
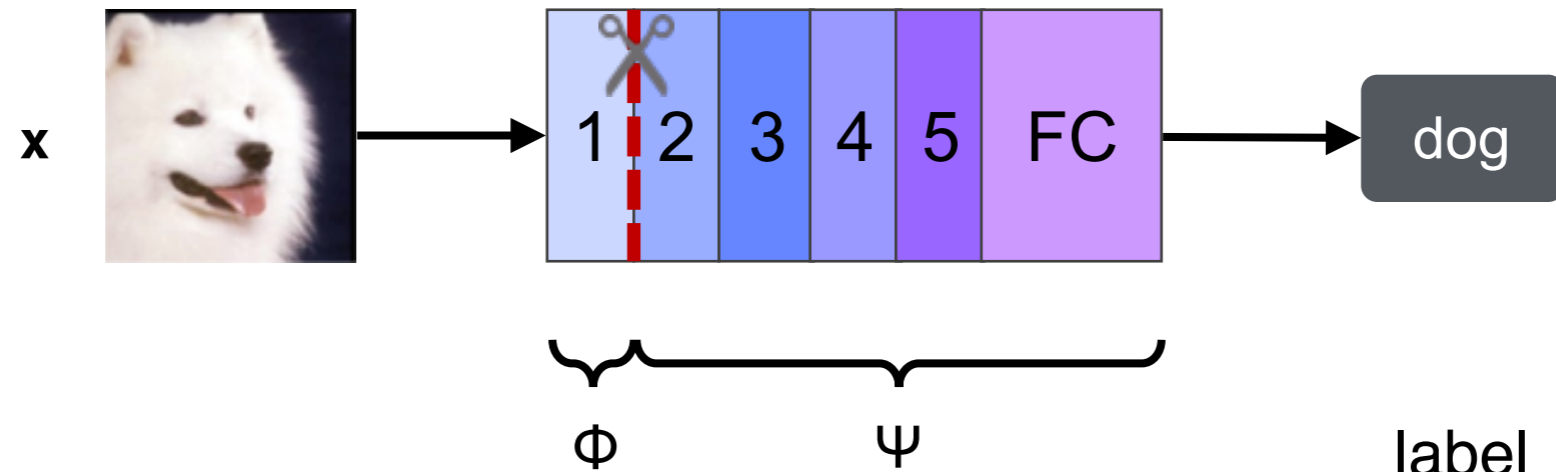
CNN: a sequence of representations



We run the same analysis on a typical CNN architecture

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

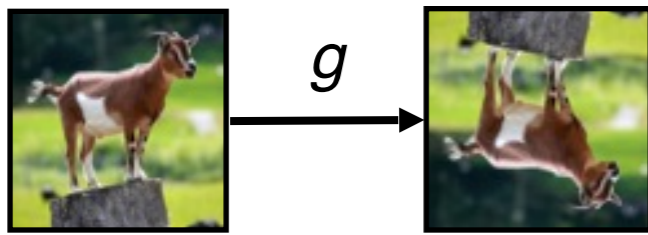
A discriminative goal to learn equivariance



M_g is learned empirically

All the other layers
(representation and classifier)
are frozen

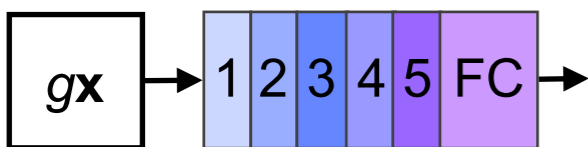
Vertical flips



Uncompensated, no TF



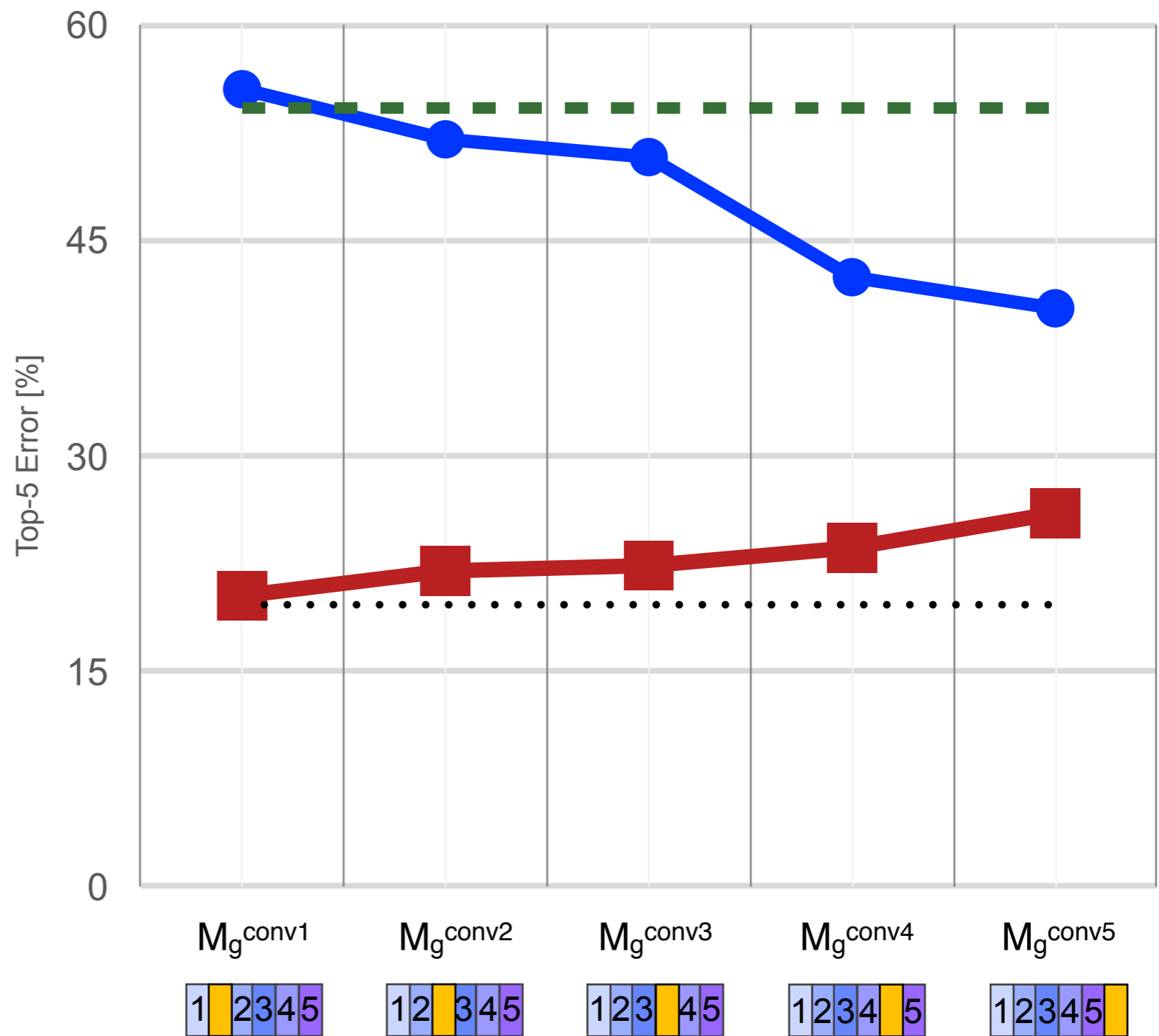
Uncompensated, TF



Compensated TF w perm.



Compensated TF, learned



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