

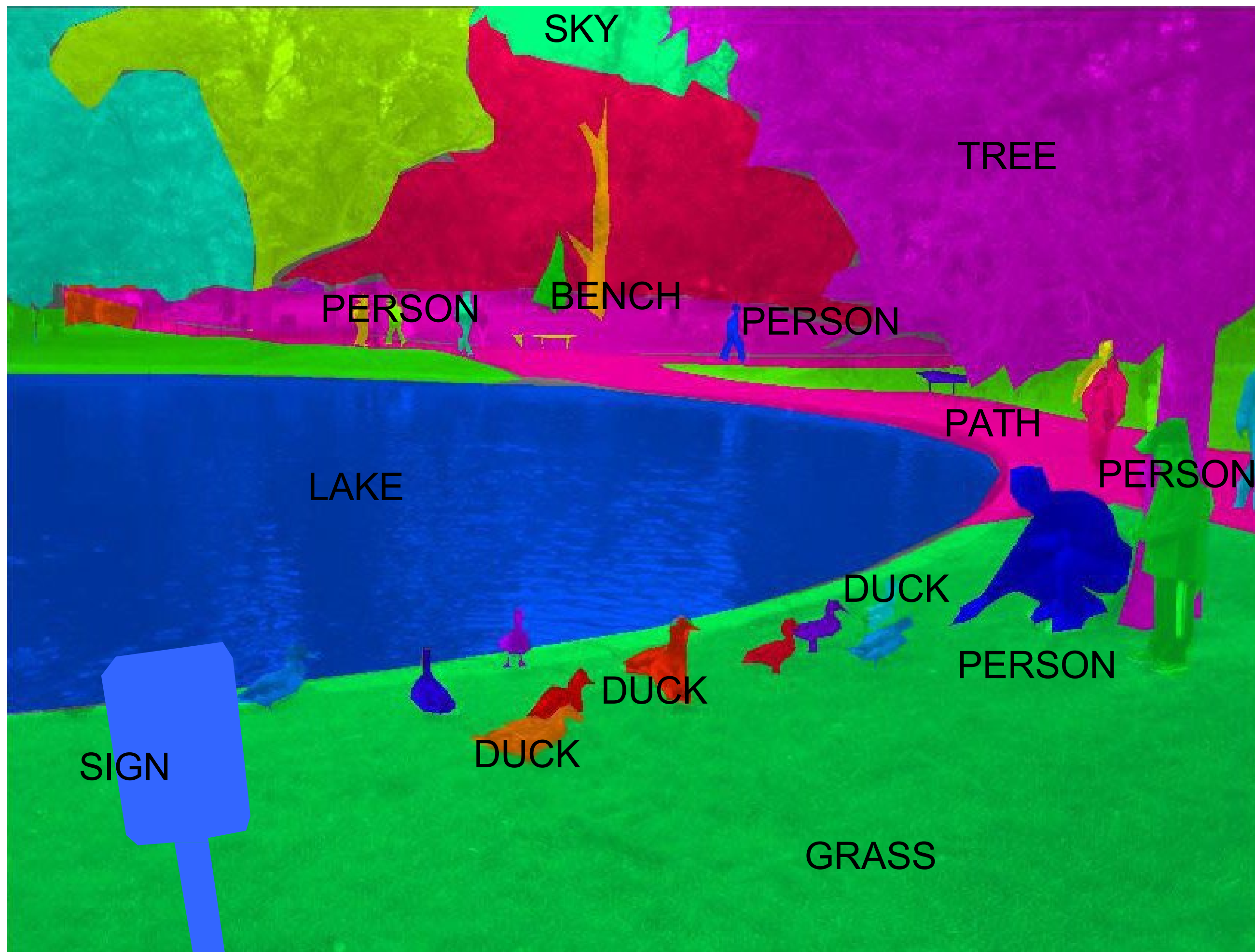
Visual Recognition with Deep Neural Networks

CMSC 491/691 Robust Machine Learning



Announcements

- **Team formation is due by next Monday**
 - Announcement on Blackboard has sign-up link (same link as Presentation/Survey)
- **Project Proposals will be due around Sept 23 (tentative)**
 - Announcement coming soon. 2-page proposal
 - Rough outline:
 - *Title of the project*
 - *Names and emails of group members*
 - *Problem you wish to tackle (and why)*
 - *Proposed approach and methods*
 - *Novelty (how is your approach better or different than previous work)*
 - *Timeline*
 - *Plan for Individual Contributions (What each student in the group will do)*
 - *Expected Outcome and Worst-Case Outcome: tell us the what you ideally want to achieve through this project*



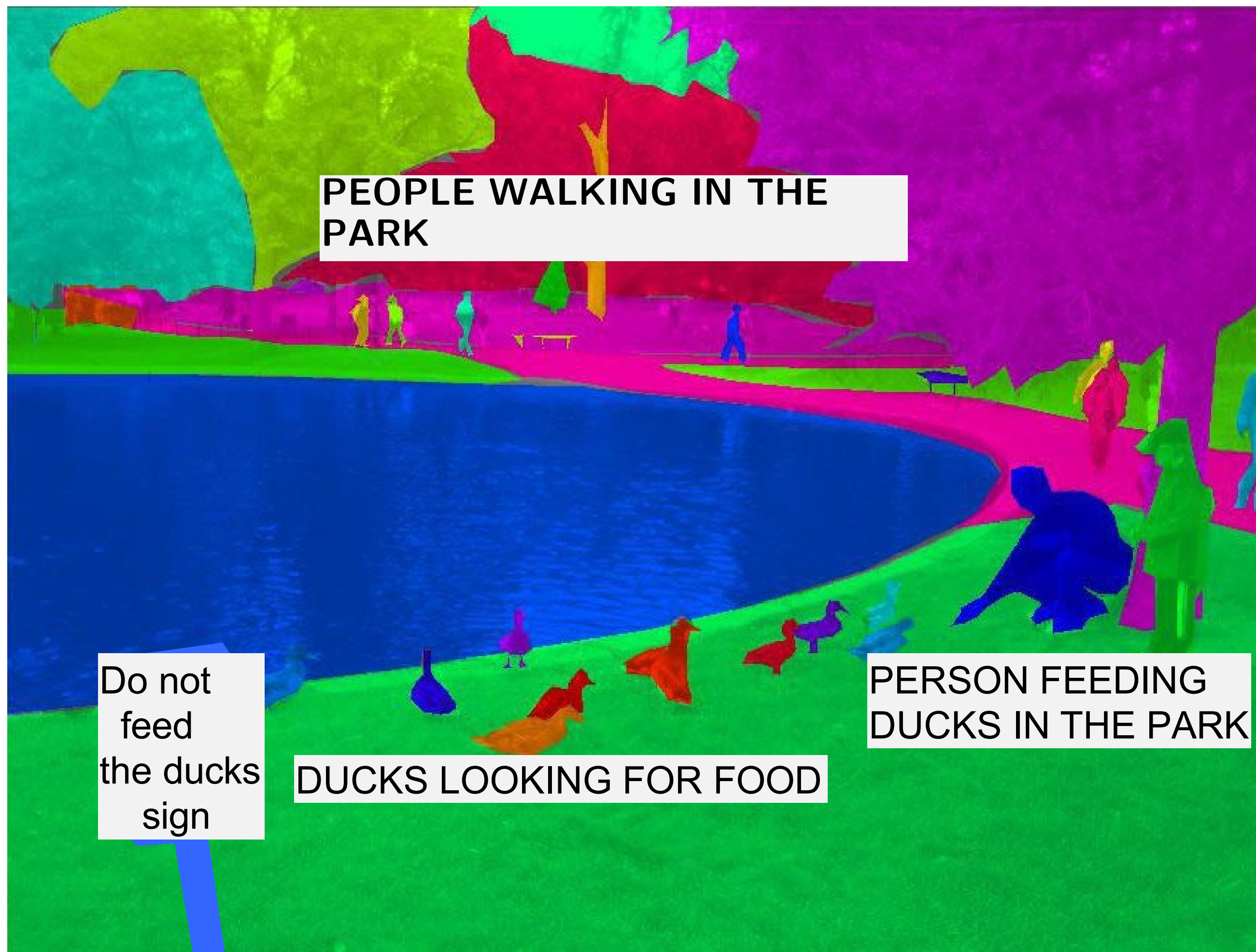
Label each pixel as a category. Each category has a unique color.



Scene-Level Classification: This is a "PARK"



Image Captioning: Describe the image in human language (e.g. English)



Dense Image Captioning: Describe several parts of the image

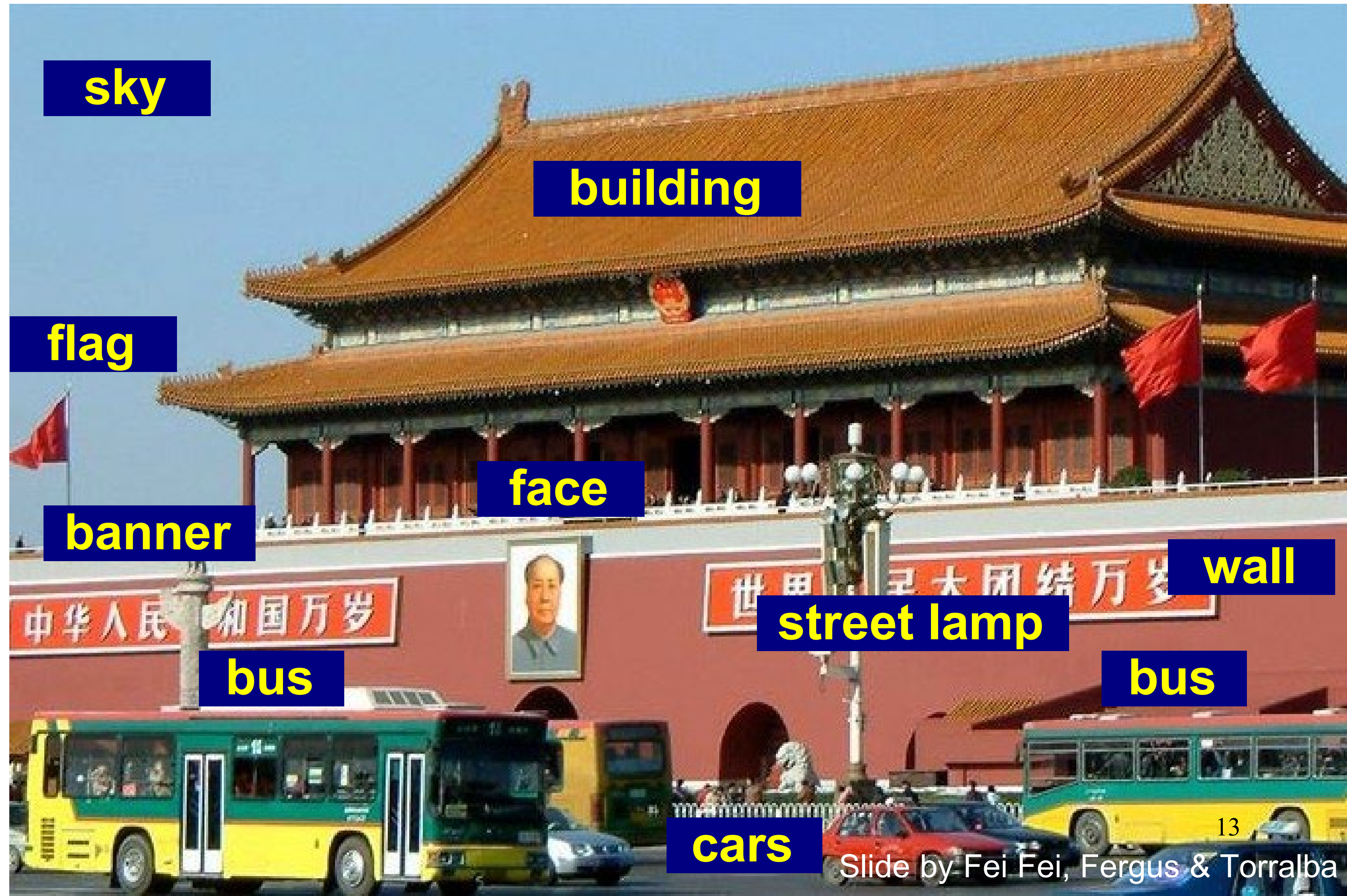
What makes this challenging?

Why do we care about recognition?



- The concept of “**categories**” encapsulates semantic information that humans use when communicating with each other.
- Categories are also linked with what can we do with those objects.

Object categories aren't everything



Object categories aren't everything

sky

A picture is worth a 1000 words...

Or just 10?

building

flag

face

banner

wall

street lamp

bus

bus

cars

How finegrained should categories be ?



A Beijing City Transit Bus #17, serial number 43253?

Need more general (useful) information



What can we say the very first time we see this thing?

Functional:

- A large vehicle that may be moving fast, probably to the right, and will hurt you if you stand in its way.
- However, at specified places, it will allow you to enter it and transport you quickly over large distances.

Communicational:

- bus, autobus, λεωφορείο, ônibus, автобус, 公共汽车, etc.

Visual challenges with categories

- A lot of categories are functional
- Categories are 3D, but images are 2D
- World is highly varied

Chair



car



train



Limits to direct perception



Importance of Context



We might think seeing is believing ...



Video by Antonio Torralba (starring Rob Fergus)

But is it ?



Video by Antonio Torralba (starring Rob Fergus)

Image classification



image x



Classifier



"Fish"

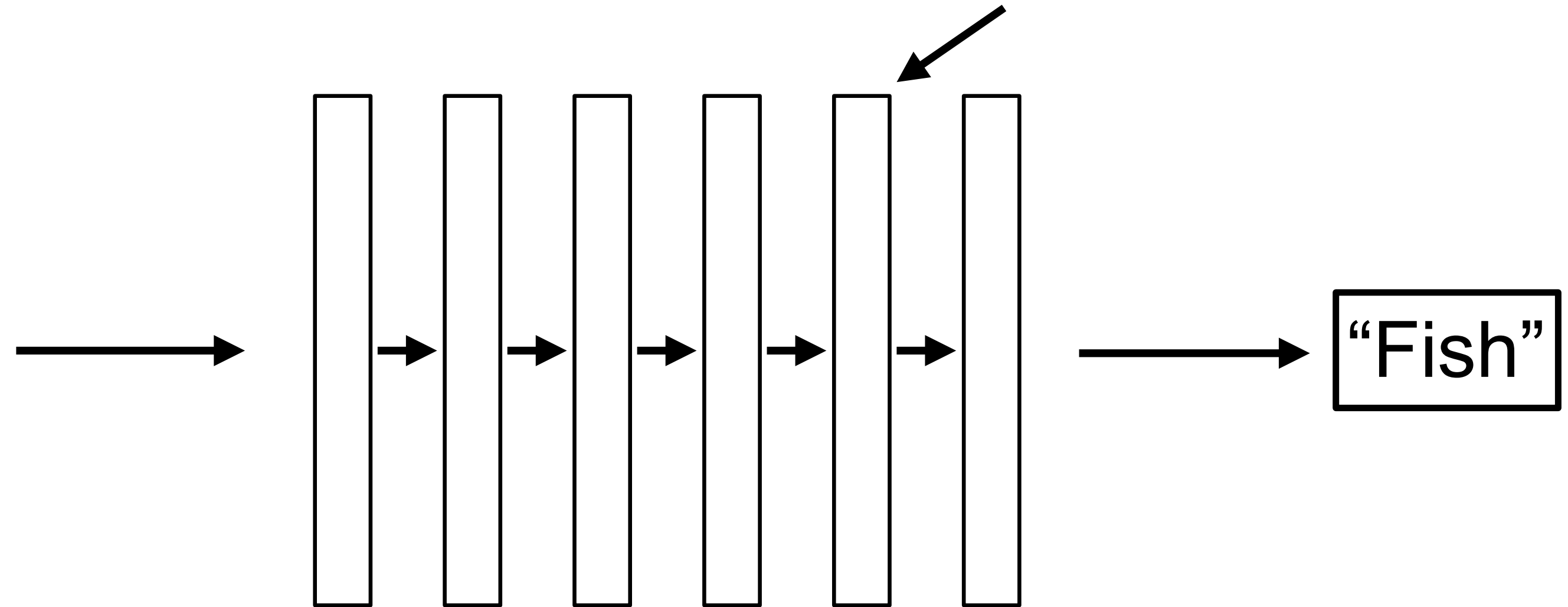
label y

Image classification

What should these be?



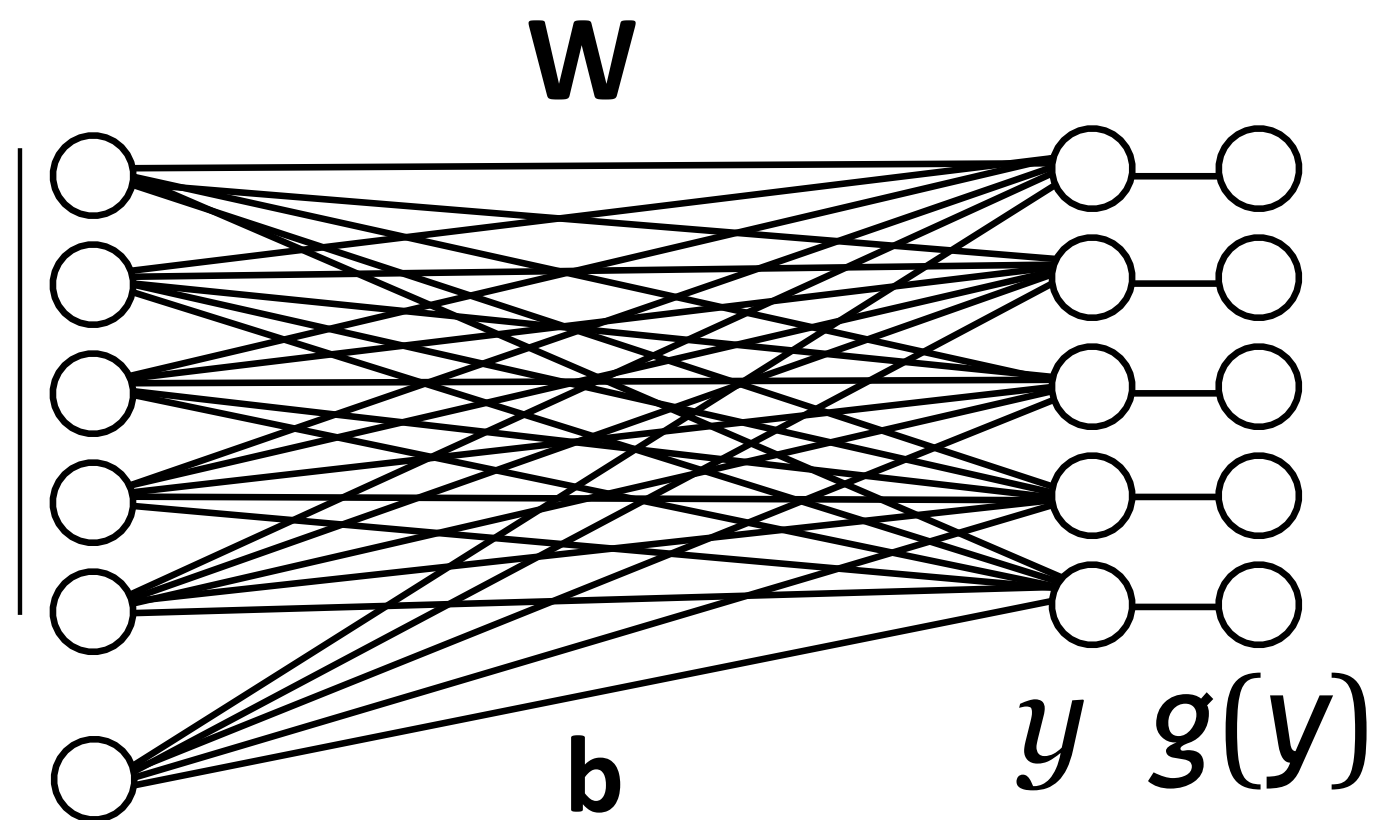
image x



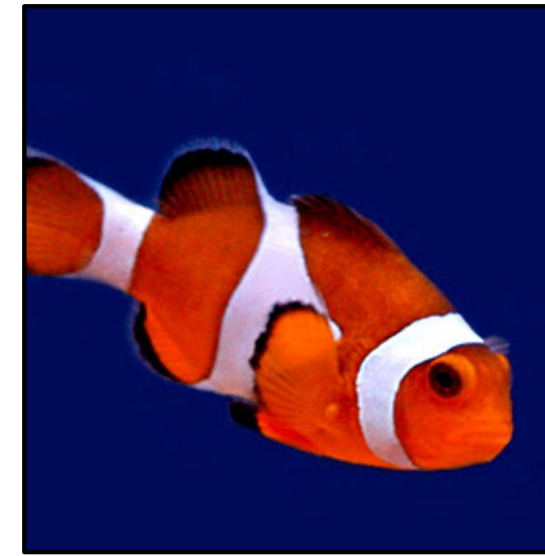
label y

Idea #1: Fully-connected network

Fully-connected (linear) layer

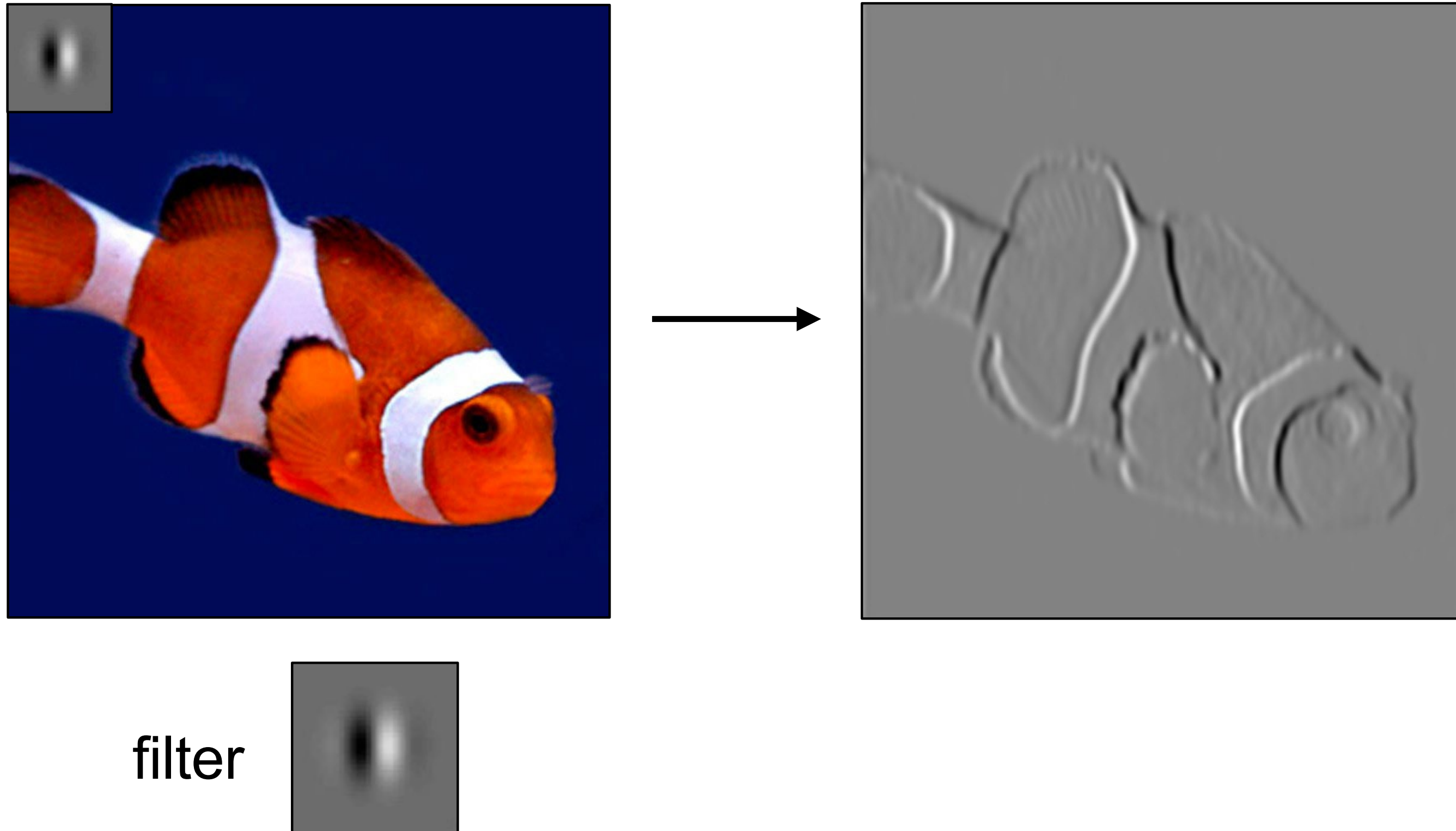


=

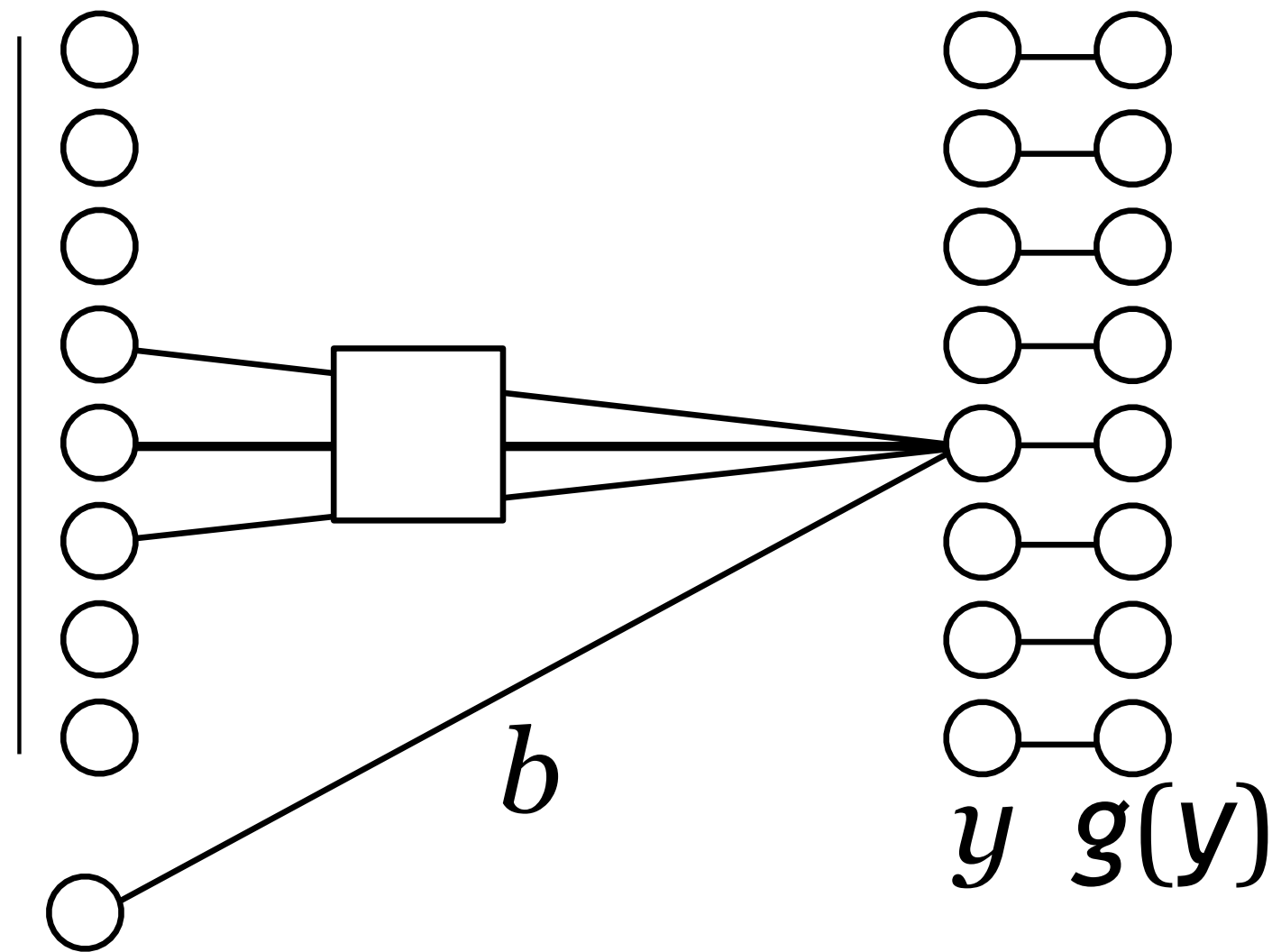
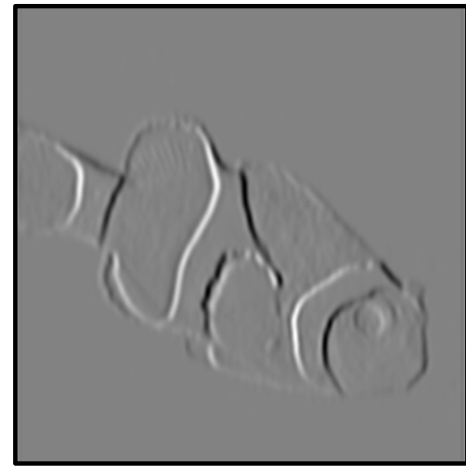


- is really big! E.g. 256×256

Can we use convolution in a neural network?



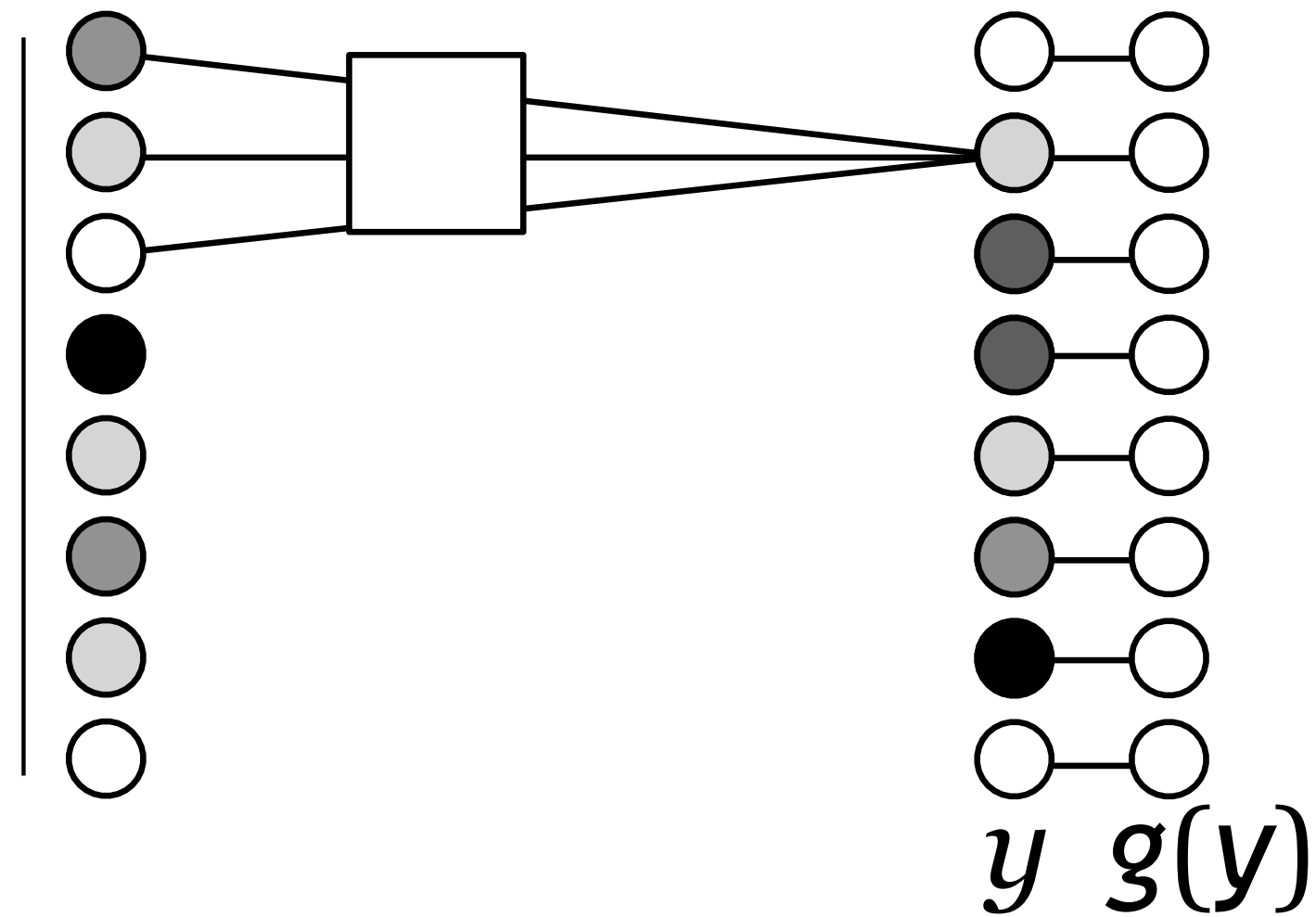
Recall: Sparsely connected network



Each unit is connected to a subset of the units in the previous layer.

Convolutional neural network

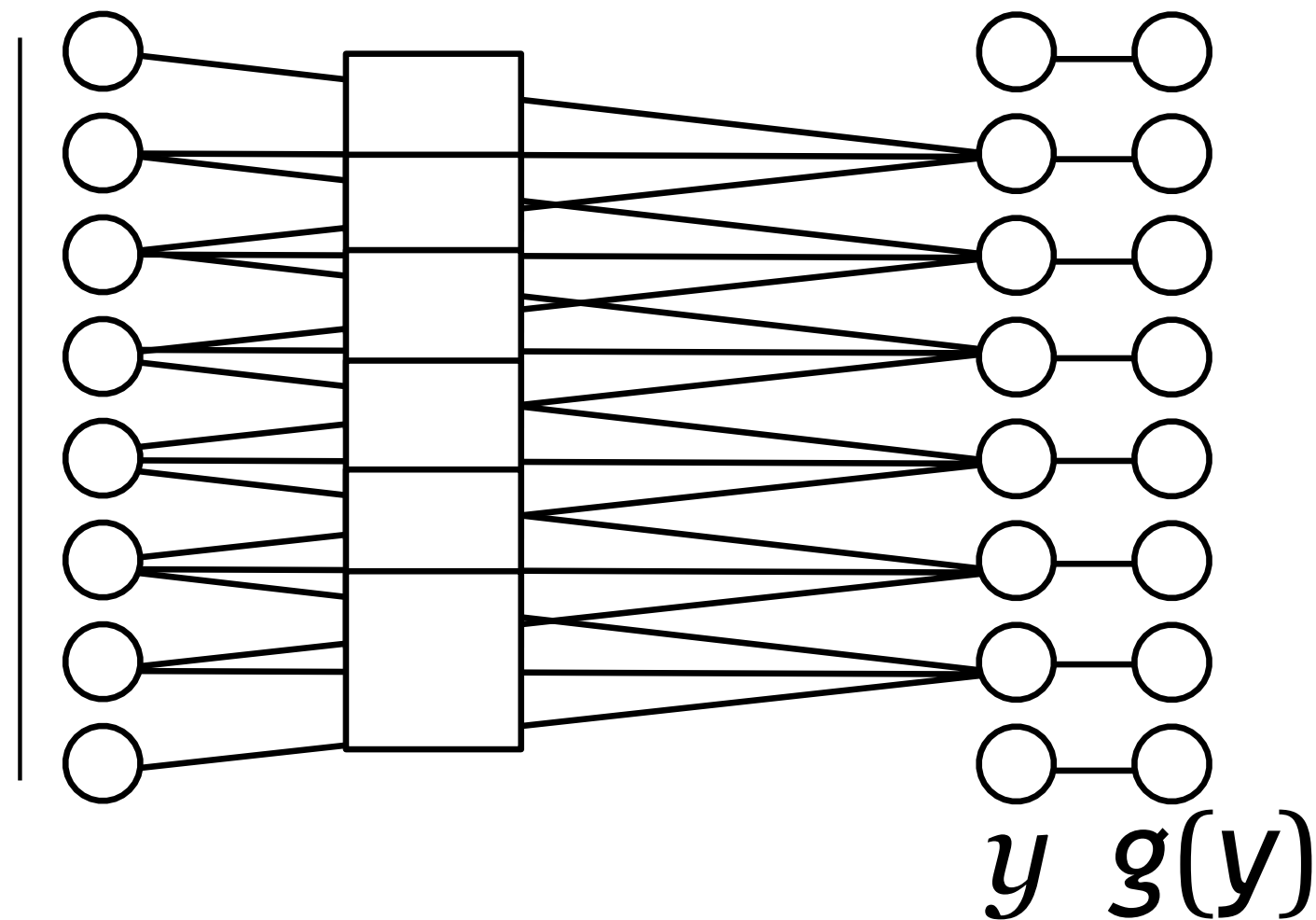
Conv layer



Each output unit is computed from an image patch.

Weight sharing

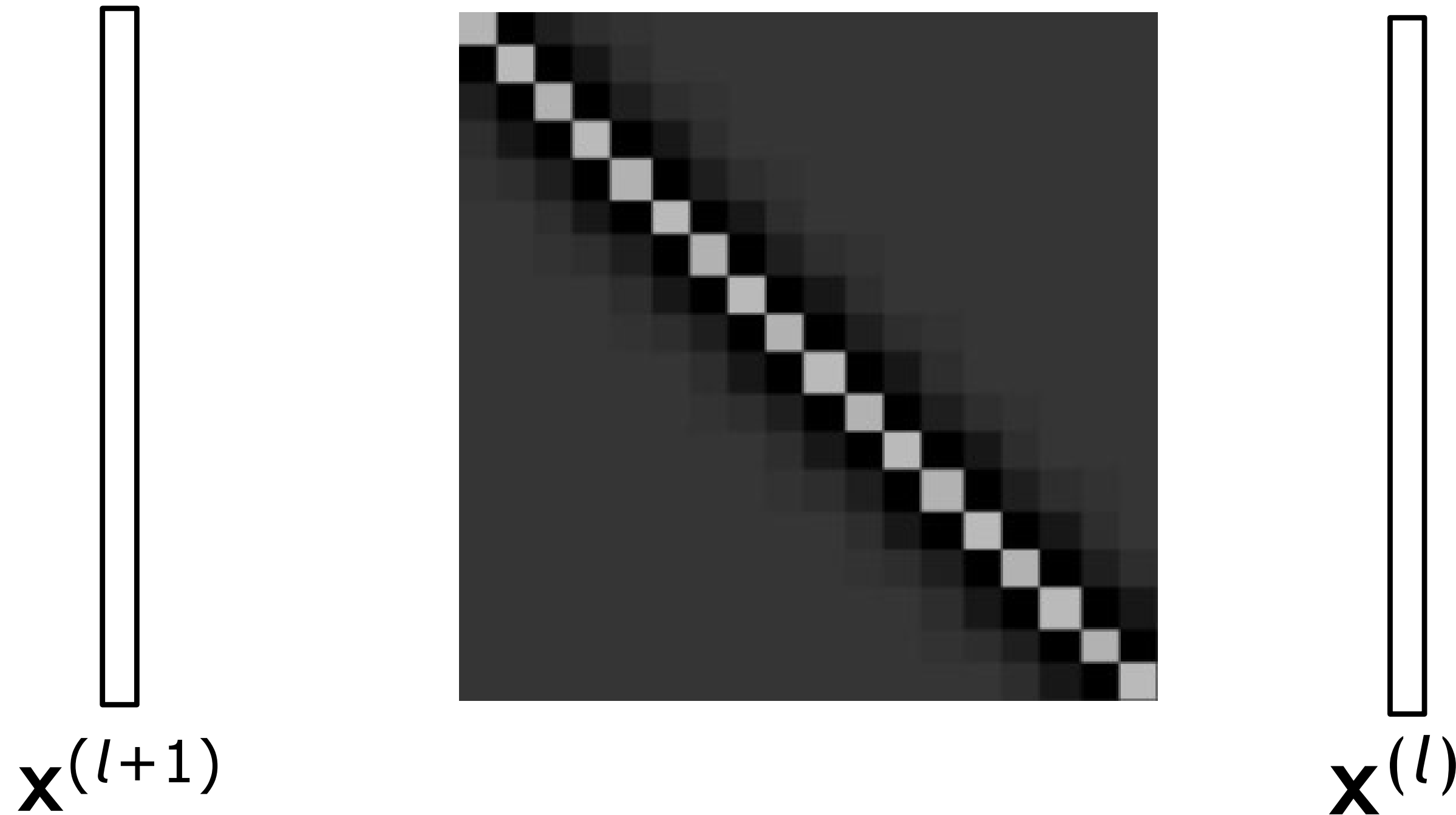
Conv layer



We “share” weights for each patch.

If a feature is useful in one position,
it should be useful in others, too.

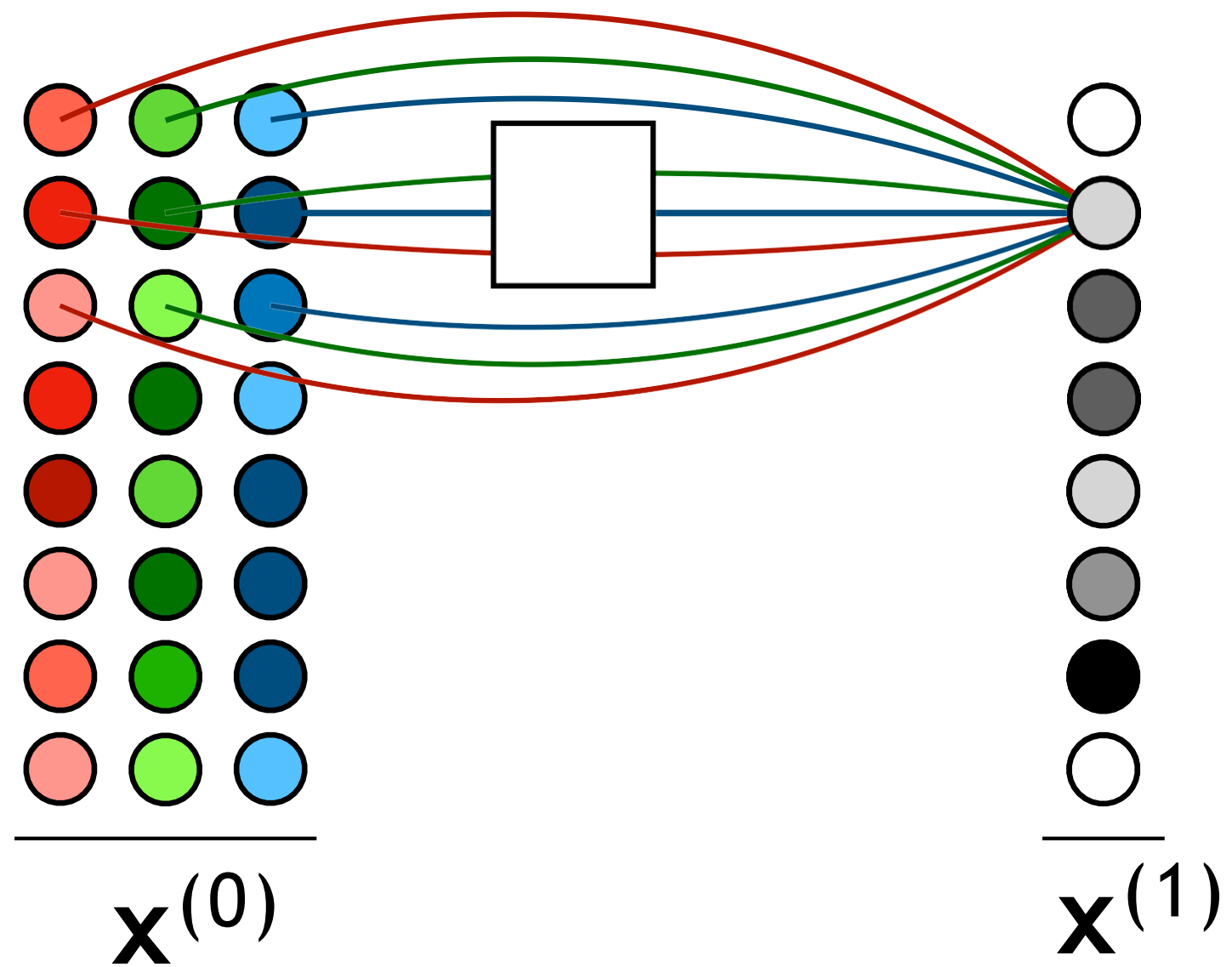
Convolution is a linear function



- Constrained linear layer
 - Fewer parameters: easier to learn, less overfitting
 - Usually use zero padding
- e.g., image

Multiple channels

Conv layer

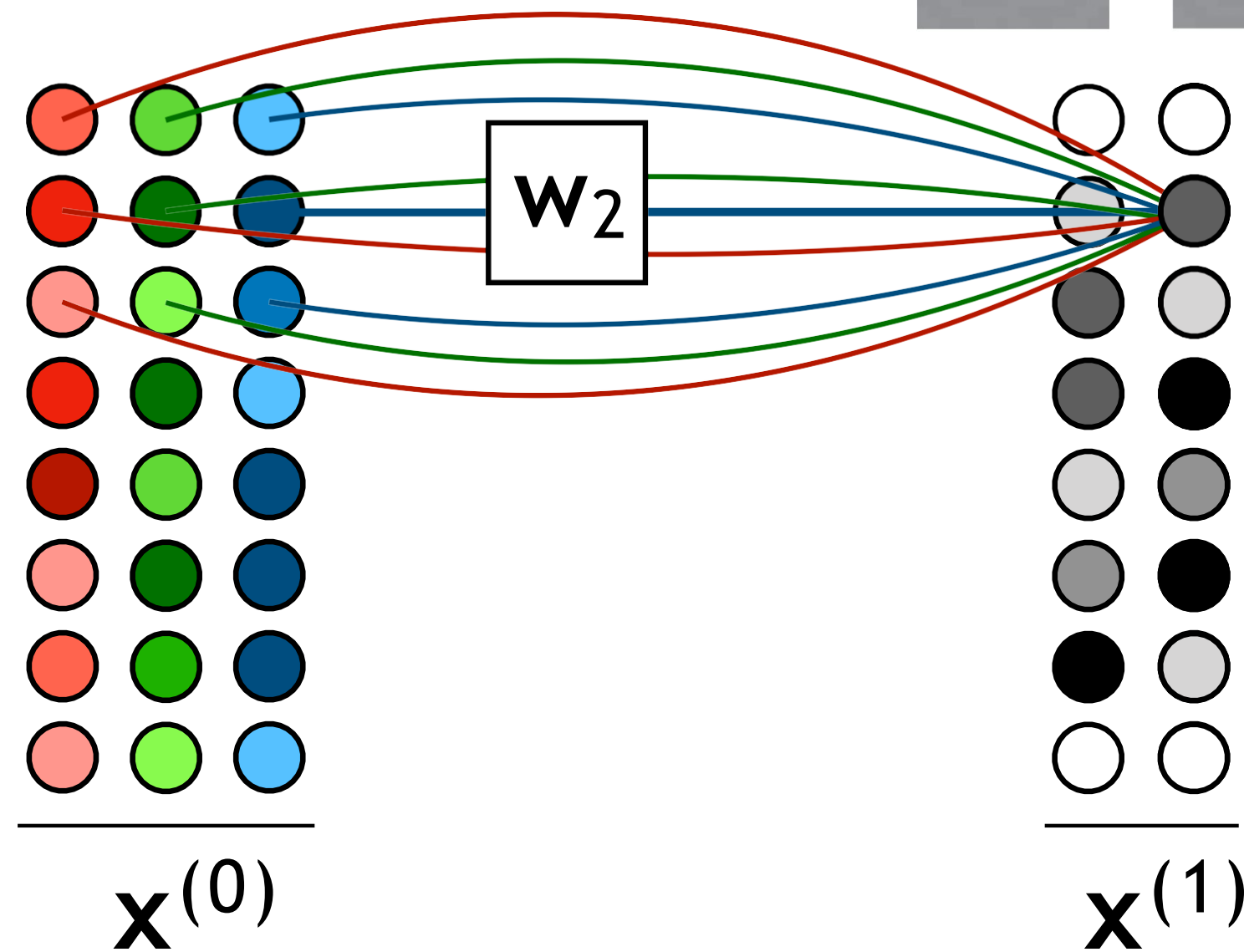
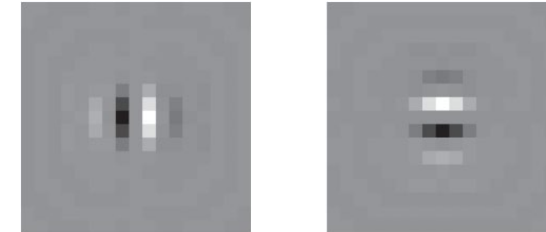


$$\mathbf{R}^{N \times C} \neq \mathbf{R}^{N \times 1}$$

Multiple channels

E.g.:

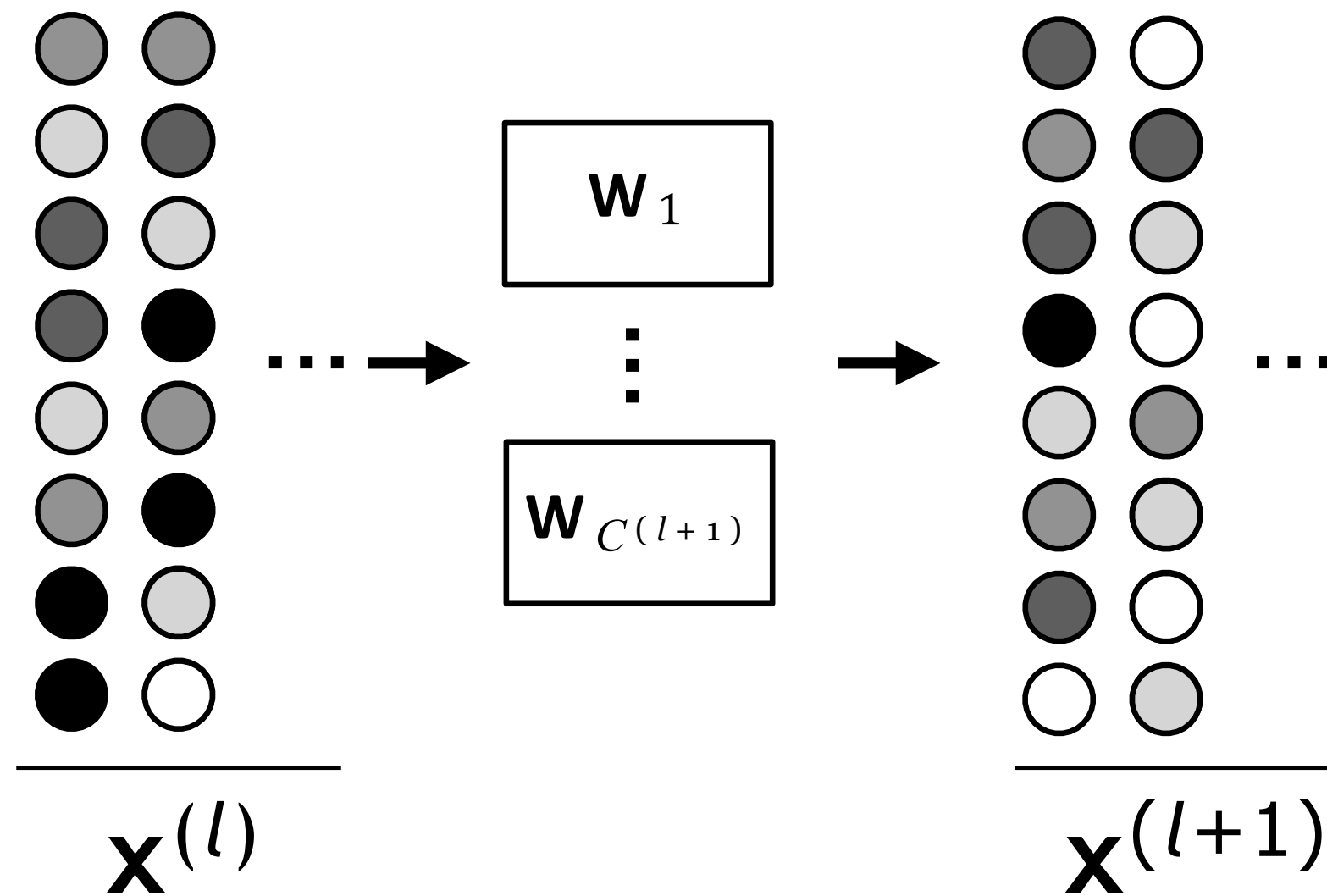
Conv layer



$$\mathbf{R}^{N \times C^{(0)}} \quad ! \quad \mathbf{R}^{N \times C^{(1)}}$$

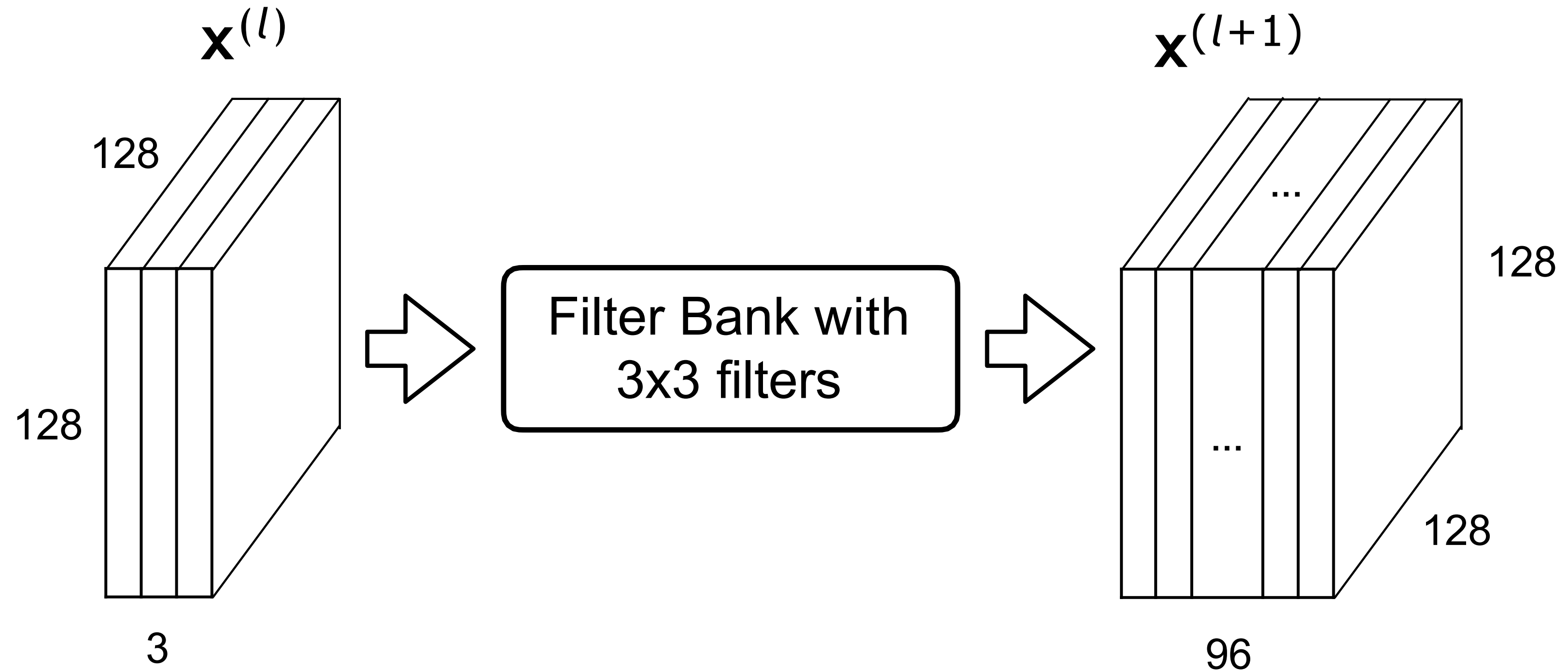
Multiple channels

Conv layer



$$\mathbf{R}^{N \times C^{(l)}} \quad ! \quad \mathbf{R}^{N \times C^{(l+1)}}$$

Multiple channels: Example



How many parameters does *each filter* have?

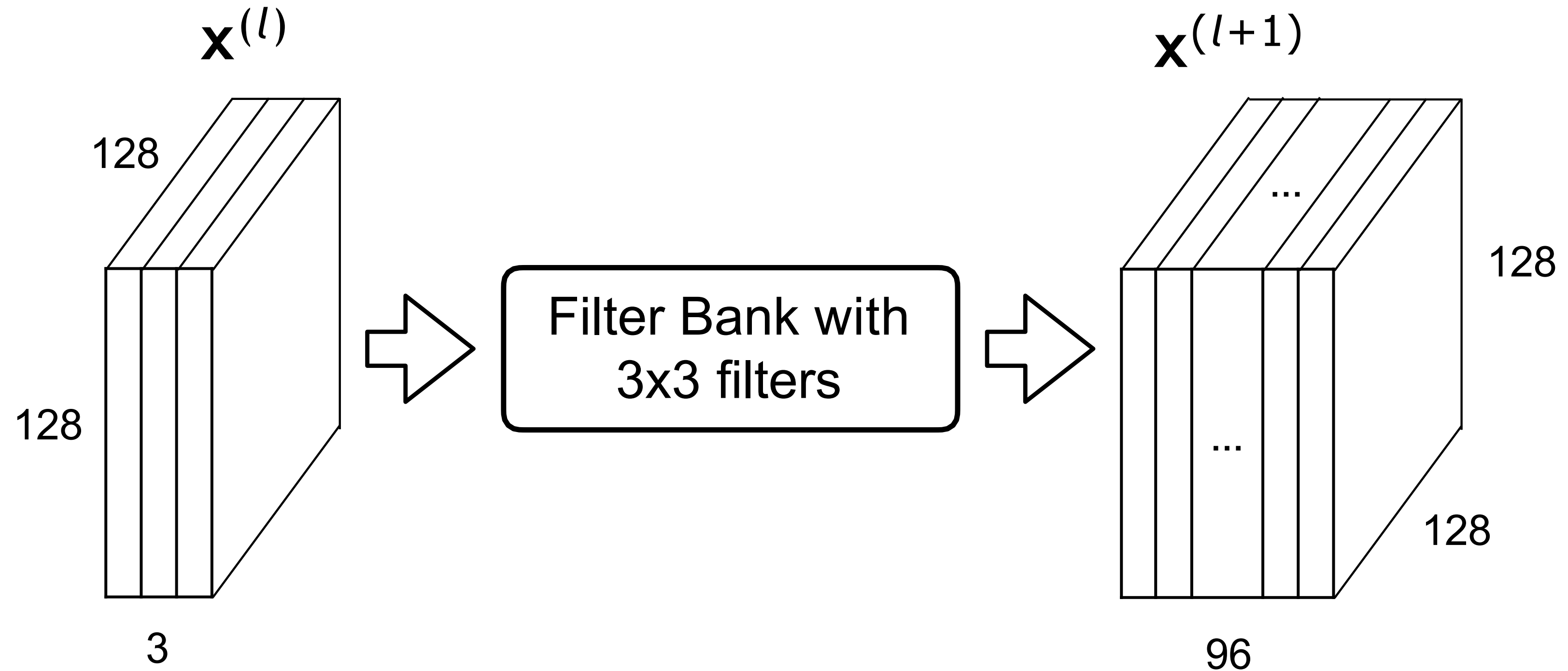
(a) 9

(b) 27

(c) 96

(d) 2592

Multiple channels: Example



How many parameters *total* does this layer have?

(a) 9

(b) 27

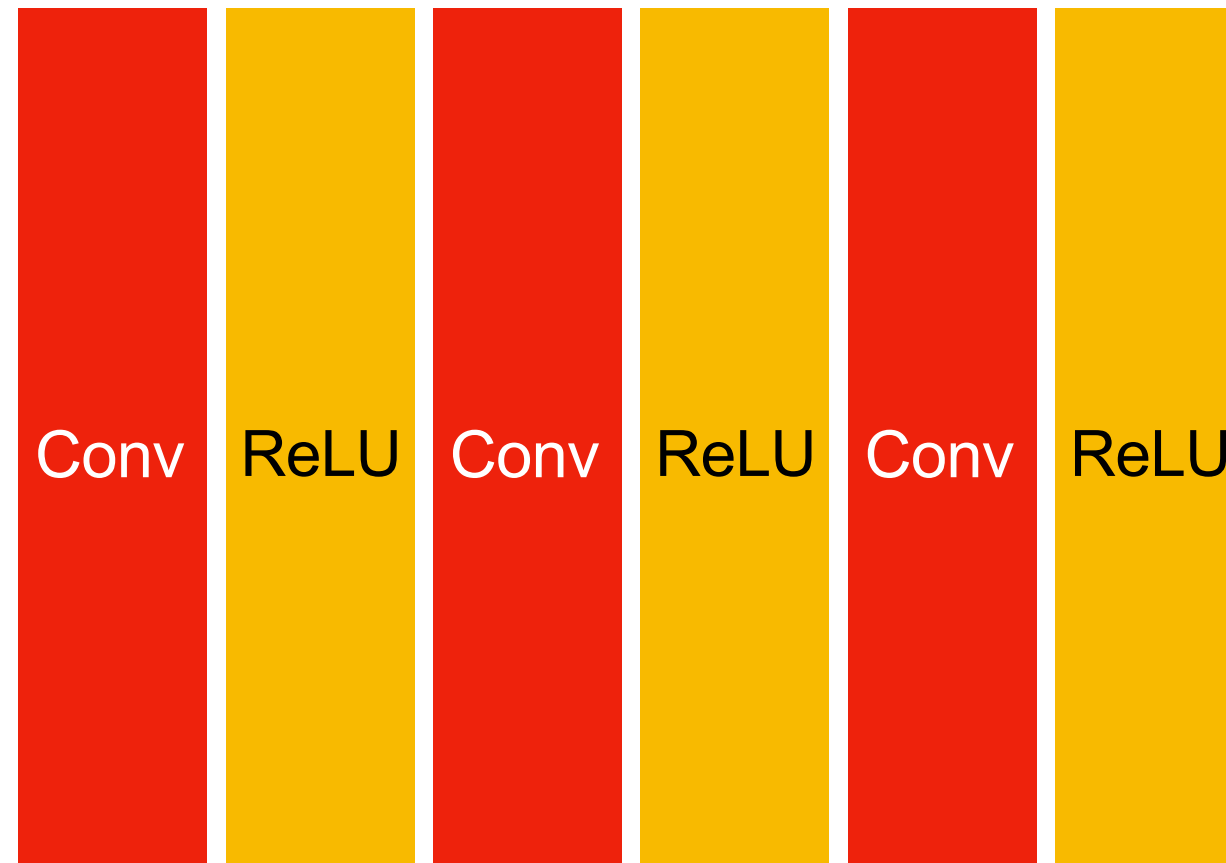
(c) 96

(d) 2592

Image classification



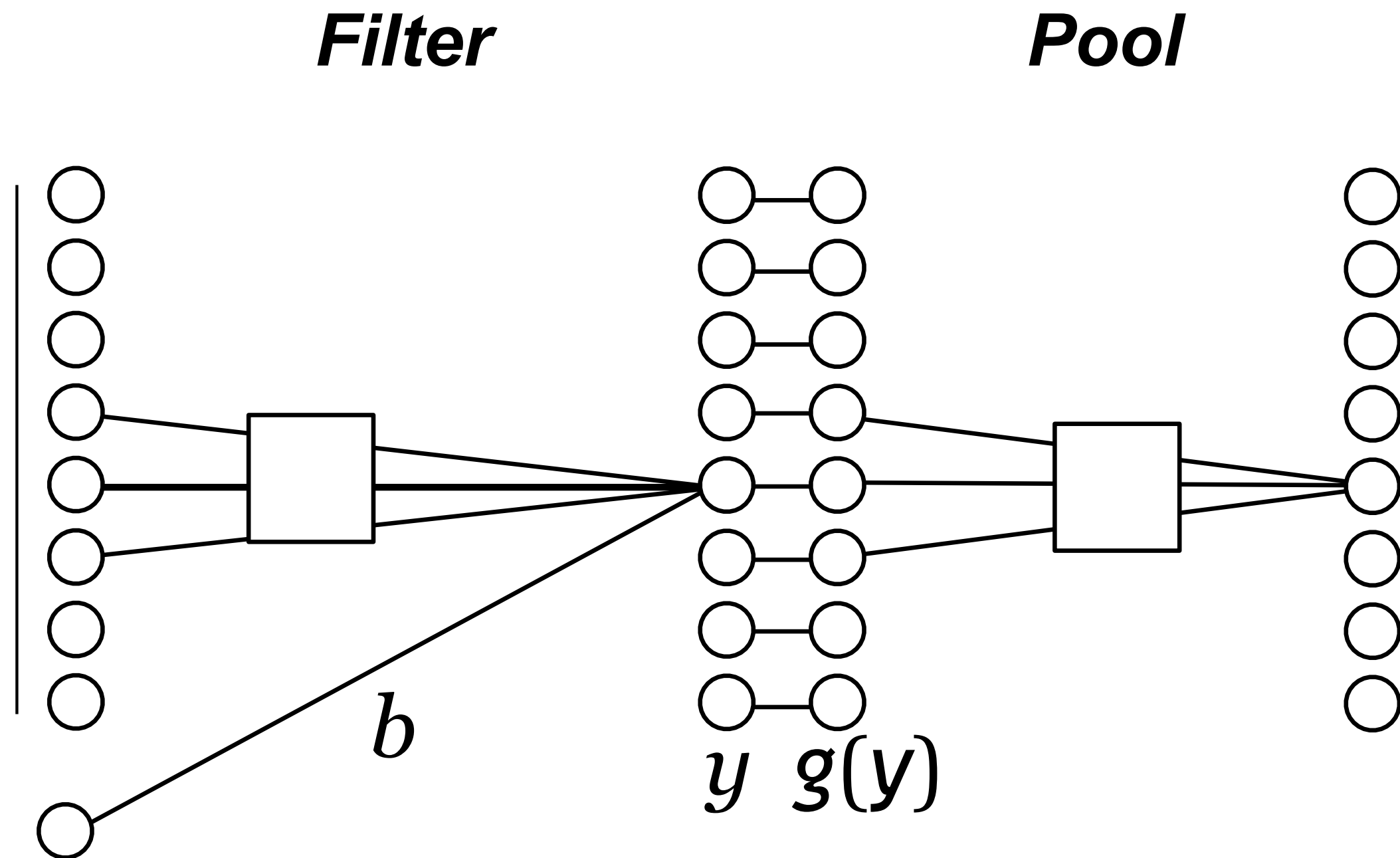
image x



“Fish”

label y

Pooling

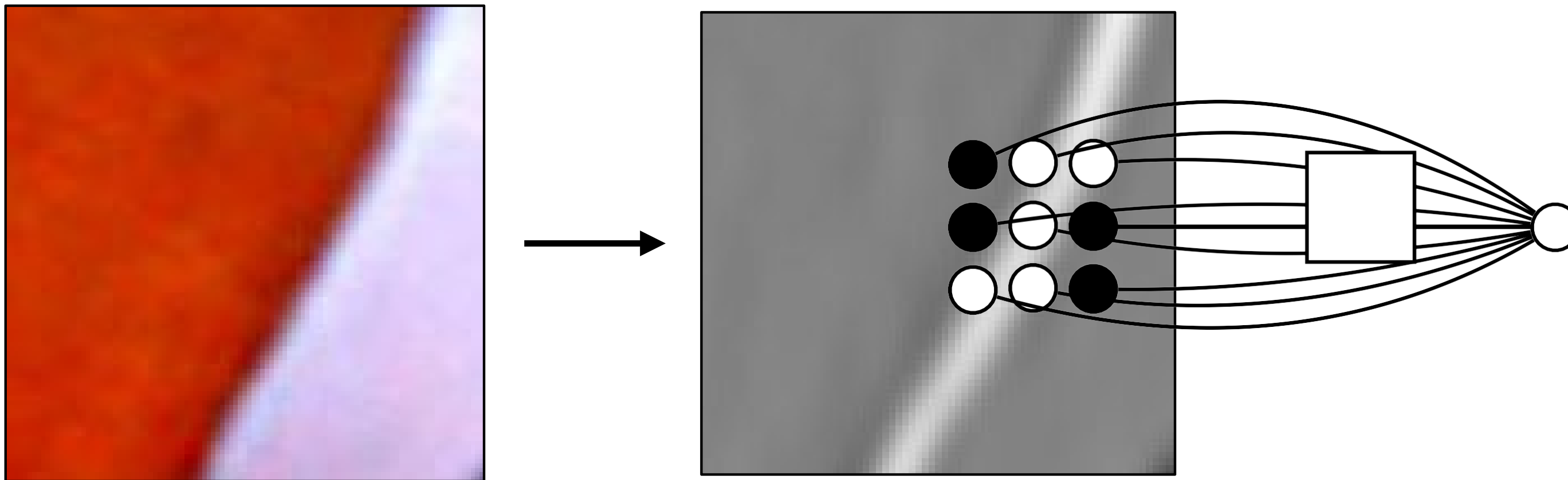


Max pooling

$$z_k = \max_{j \in N(k)} g(y_j)$$

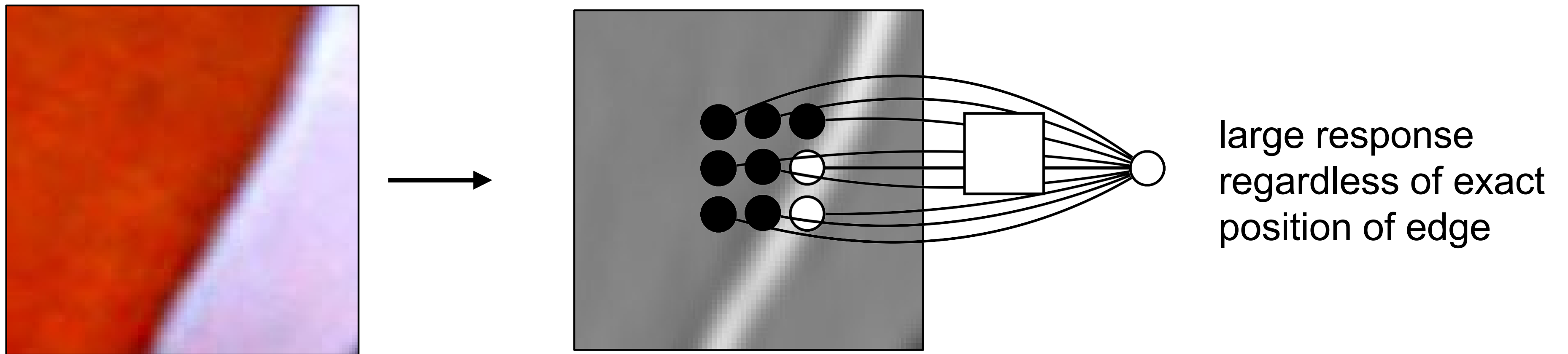
Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:



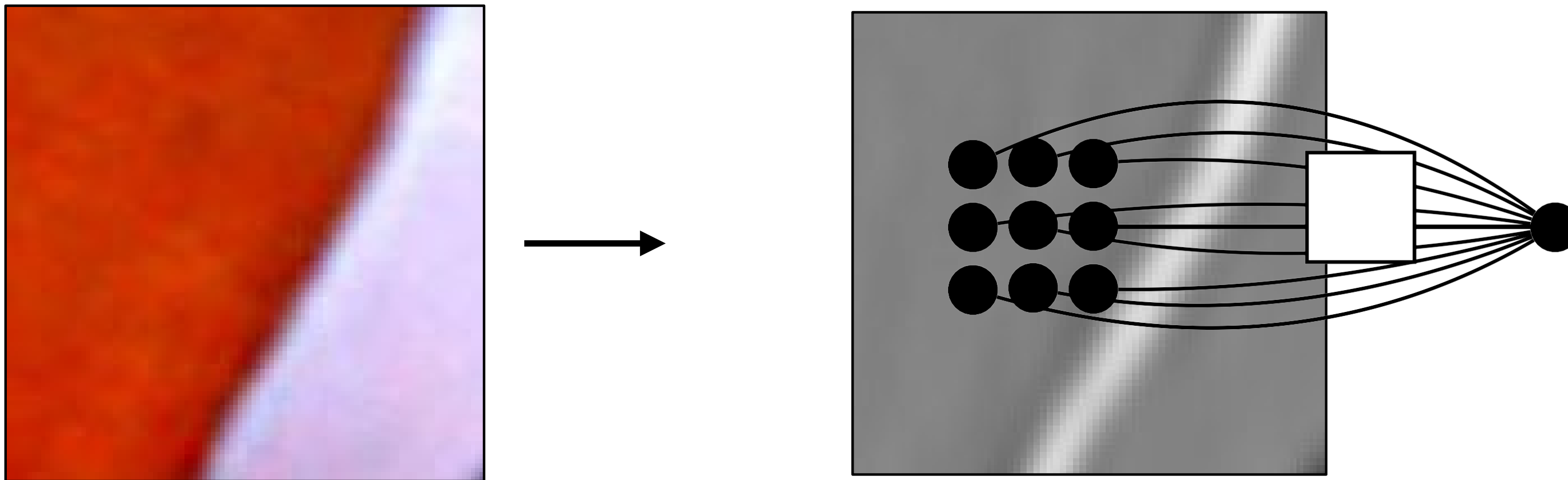
Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:

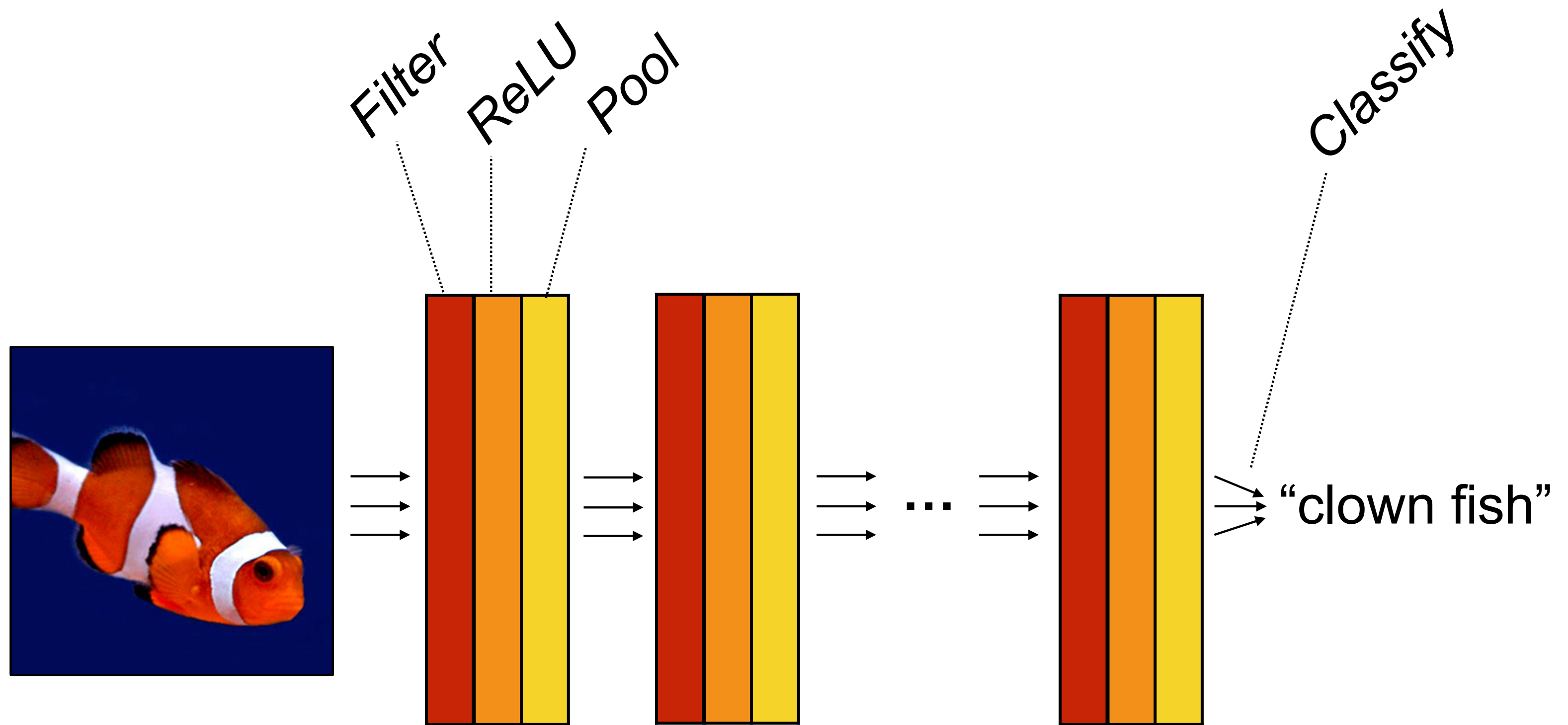


Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:

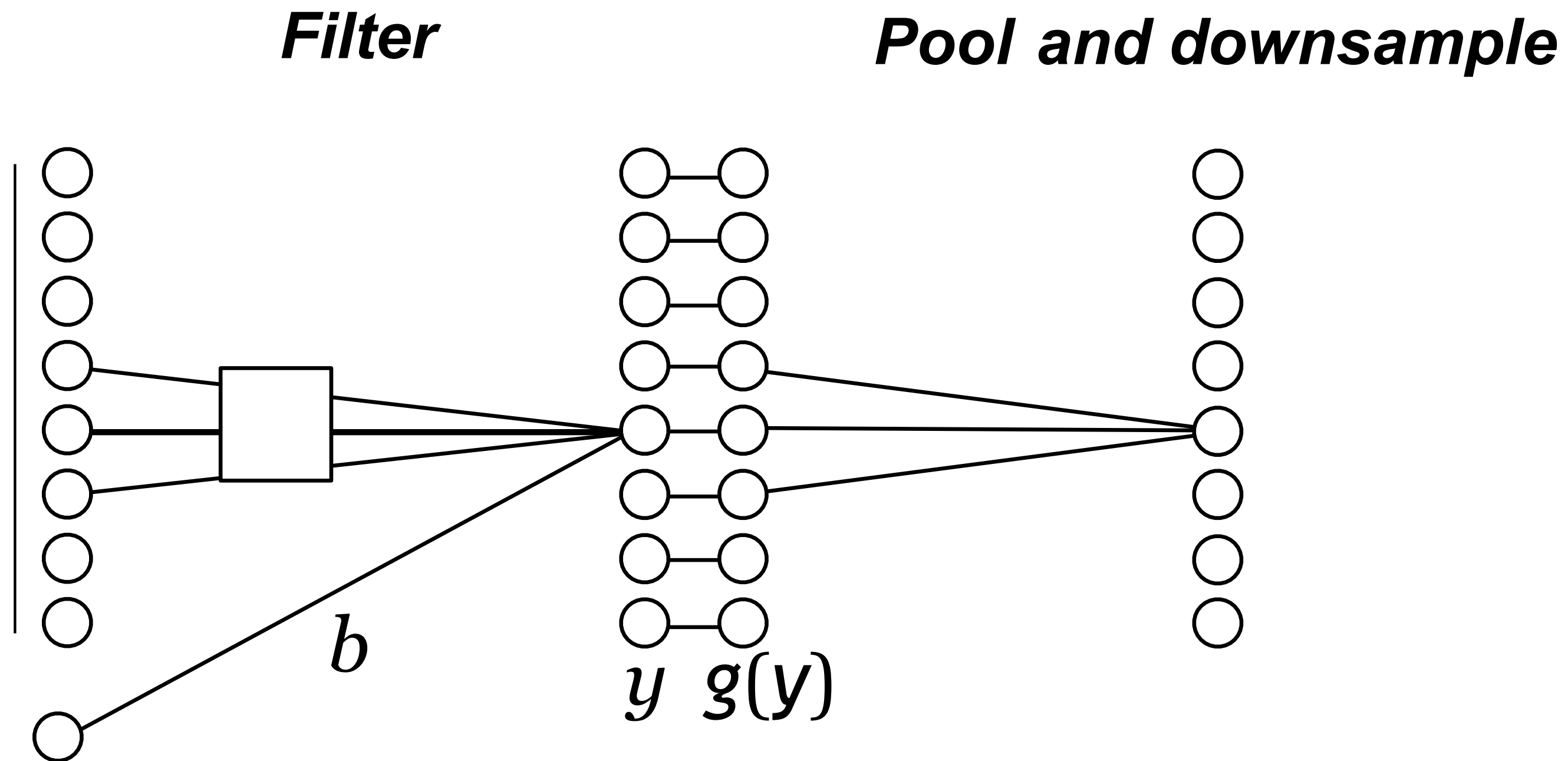


Computation in a neural net



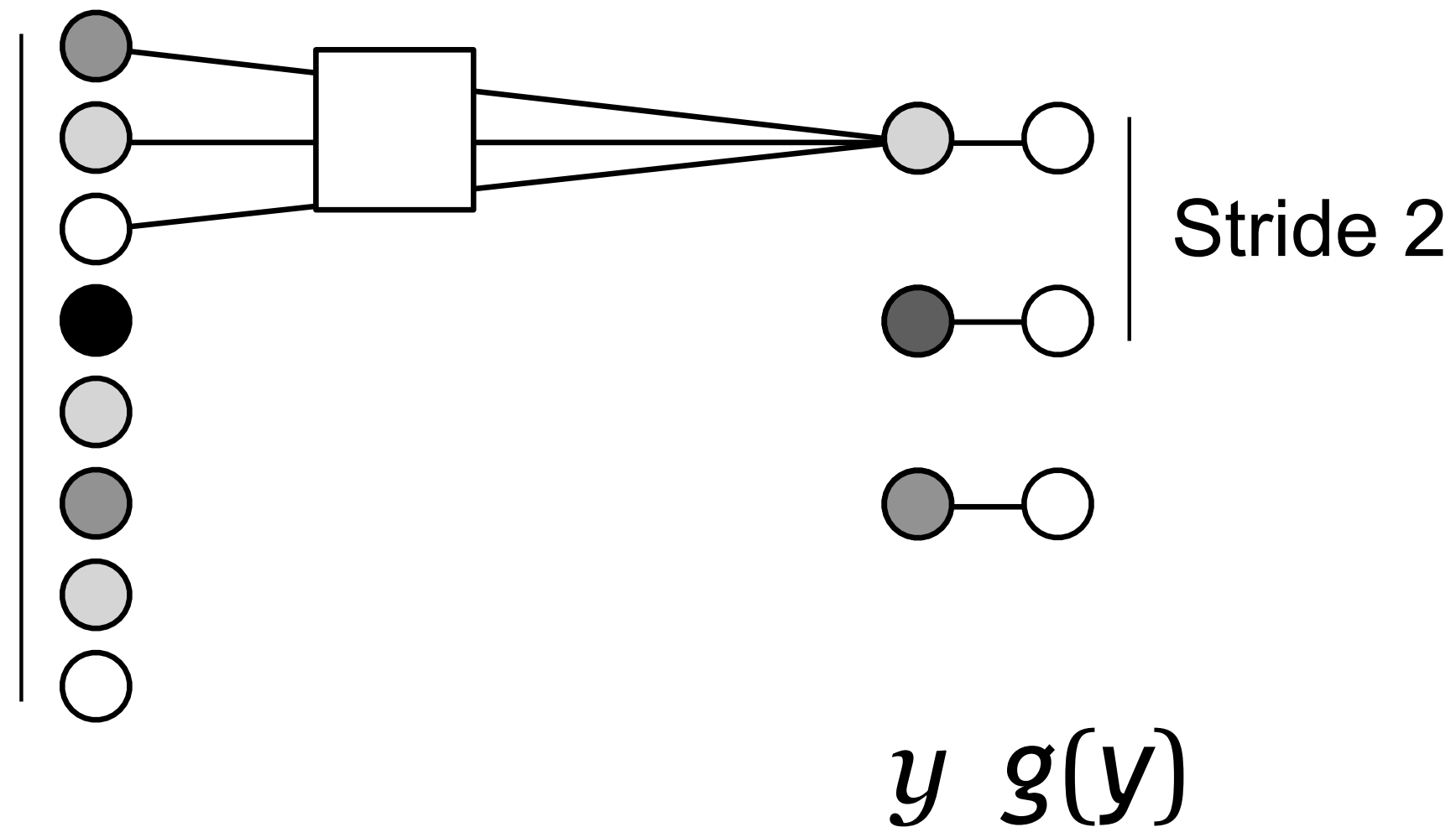
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Downsampling



Strided operations

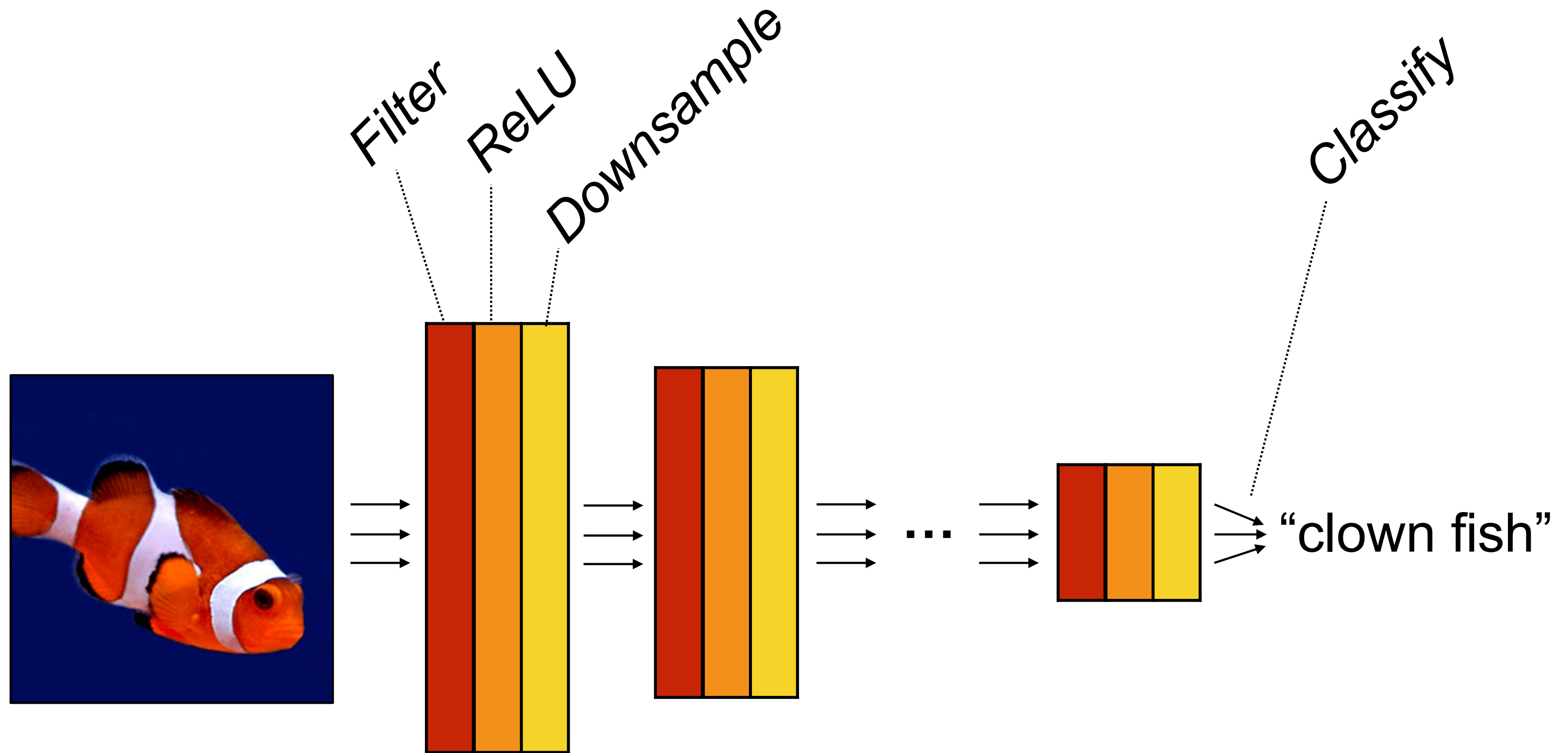
Conv layer



Strided operations combine a given operation (convolution or pooling) and downsampling into a single operation.

Strided convolution is an alternative to pooling layers:
just do a strided convolution!

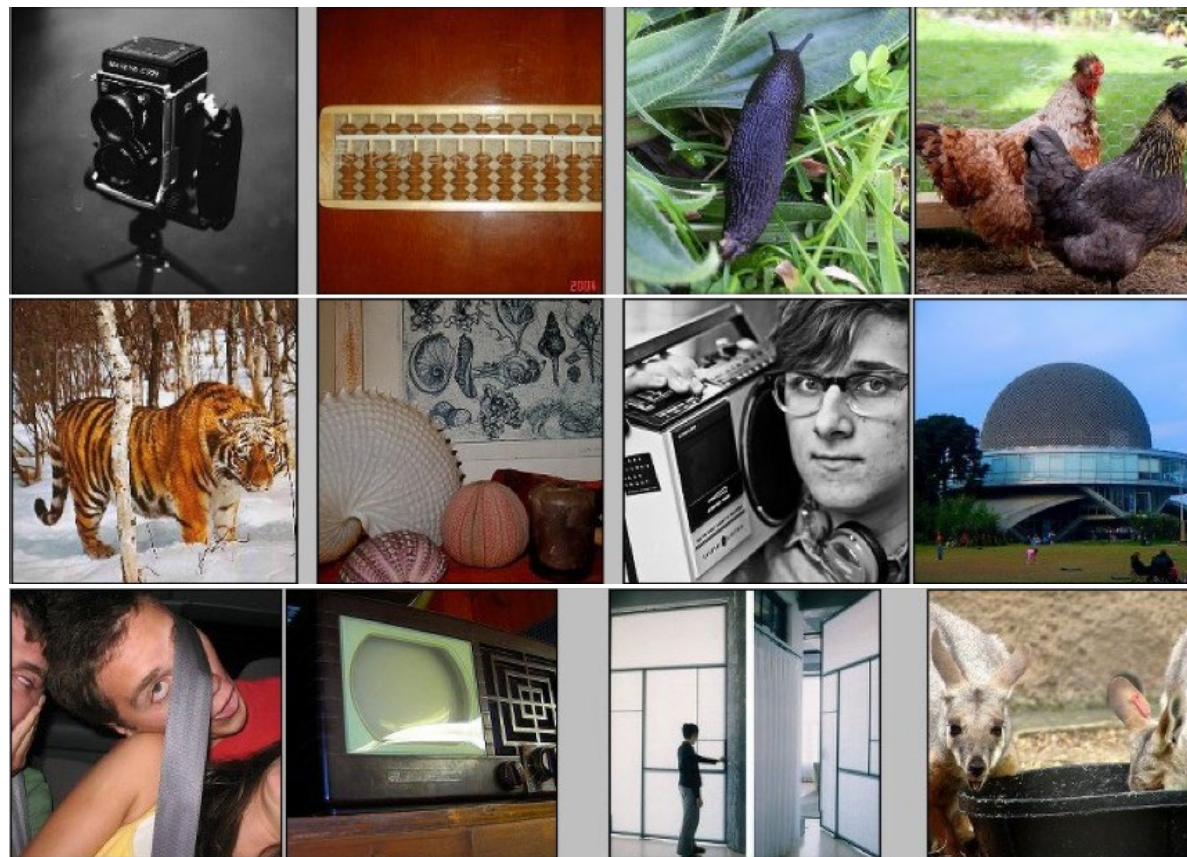
Computation in a neural net



$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

ImageNet Challenge

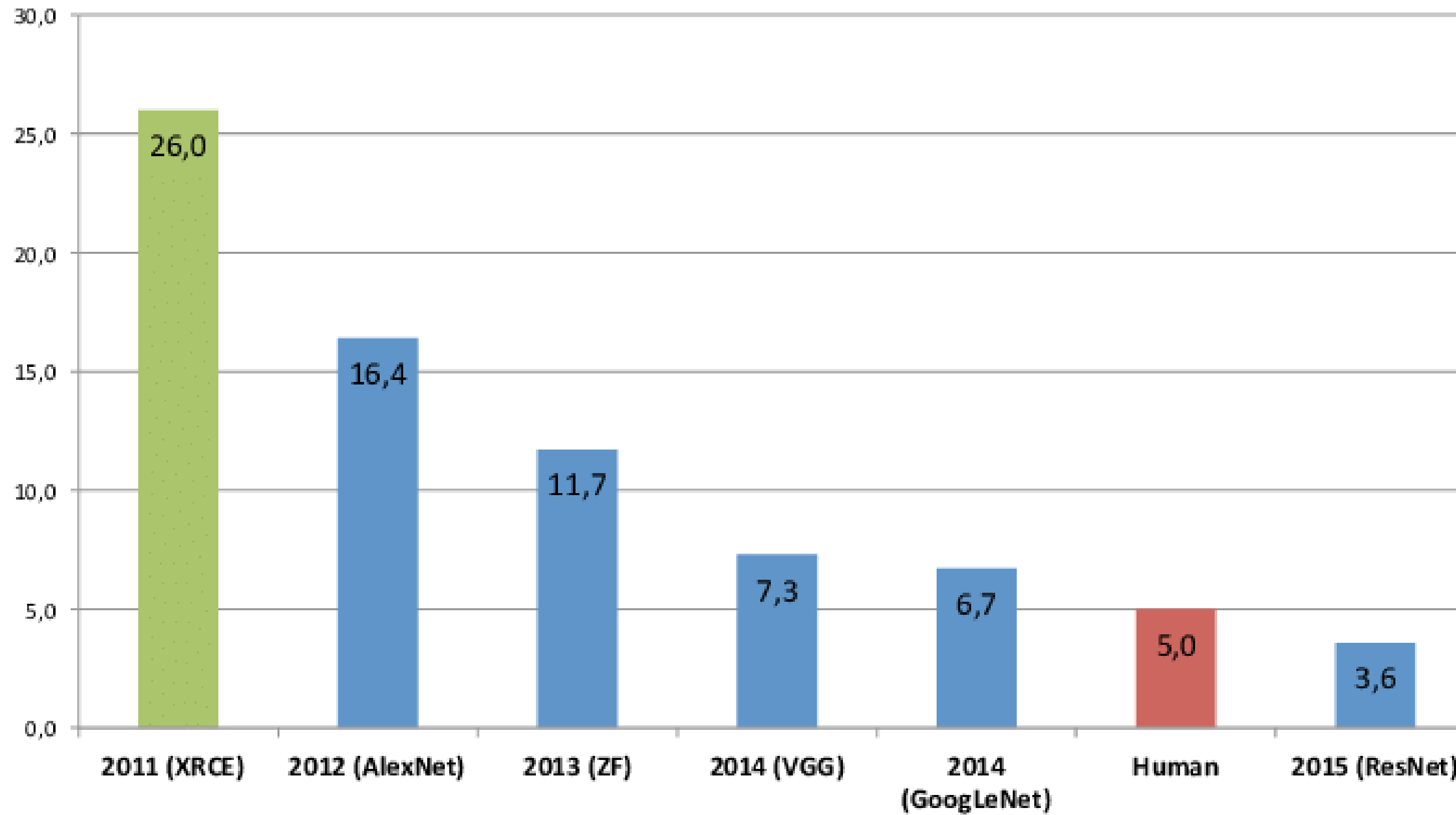
IMAGENET



- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon MTurk
- ImageNet Large-Scale Visual Recognition Challenge (ILSVRC):
1.2 million training images, 1000 classes

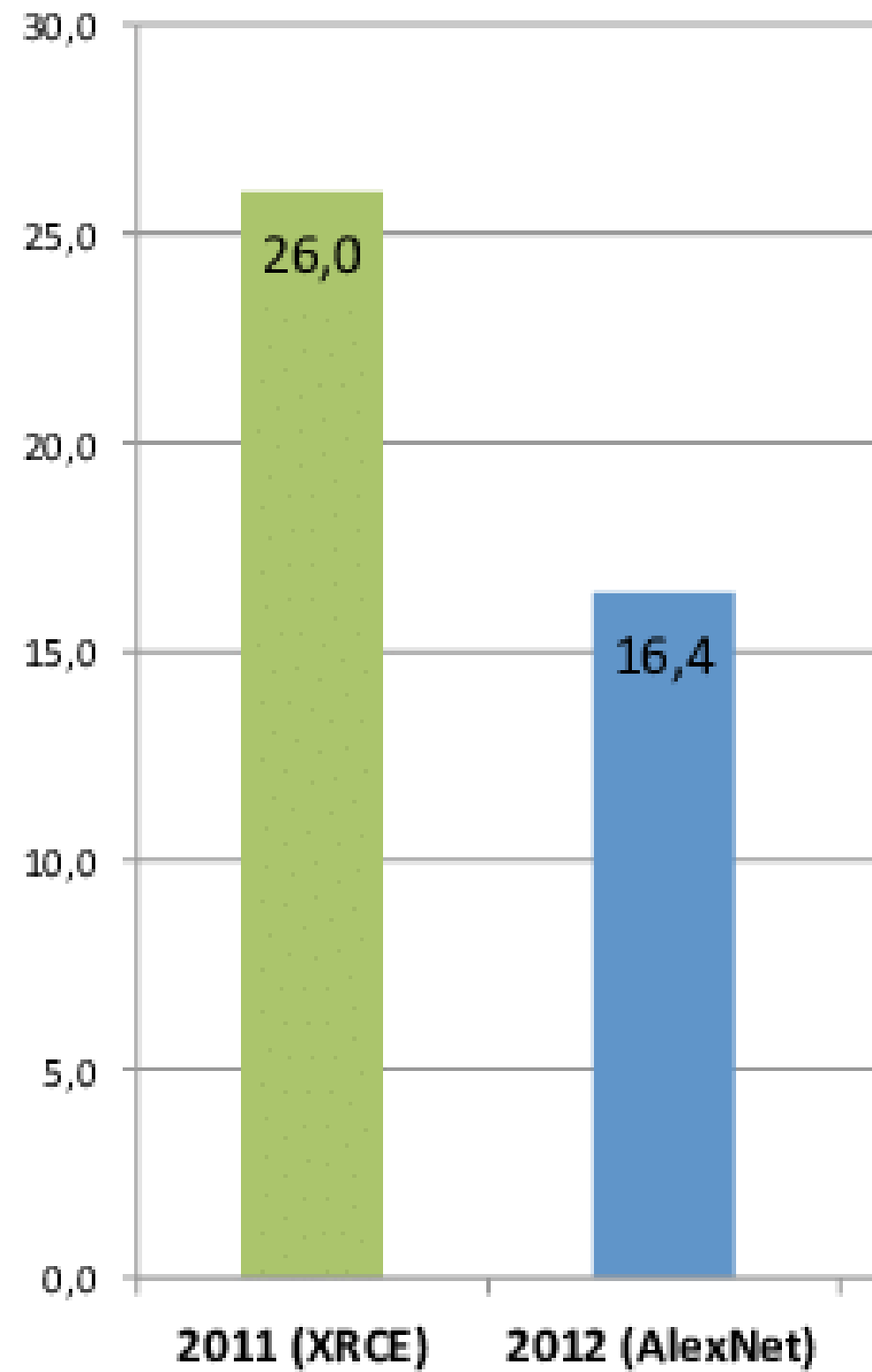
[Russakovsky, Deng, Su, Krause, Satheesh, Ma, Huang, Karpathy, Khosla, Bernstein, Berge, Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge", 2015]

ImageNet Classification Error (Top 5)

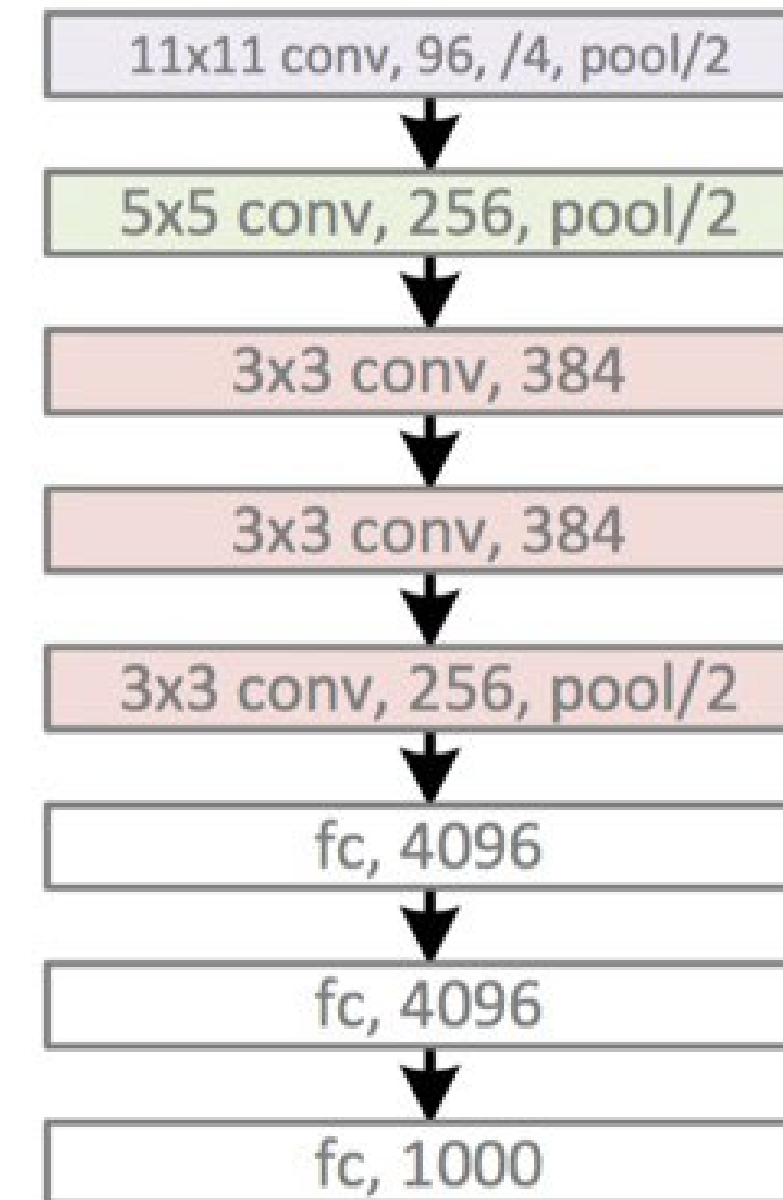


Network designs

ImageNet classification (top 5)



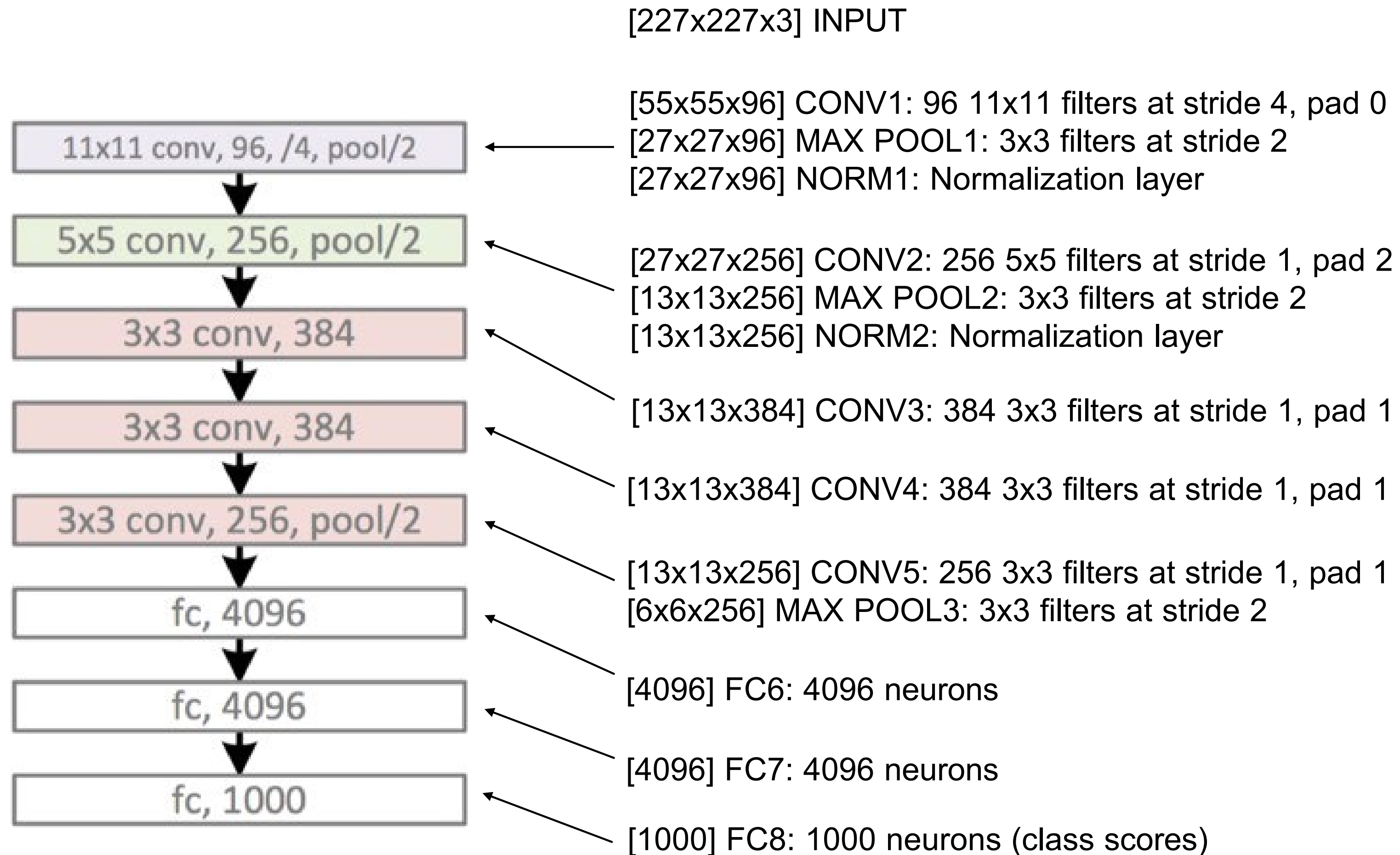
2012: AlexNet
5 conv. layers

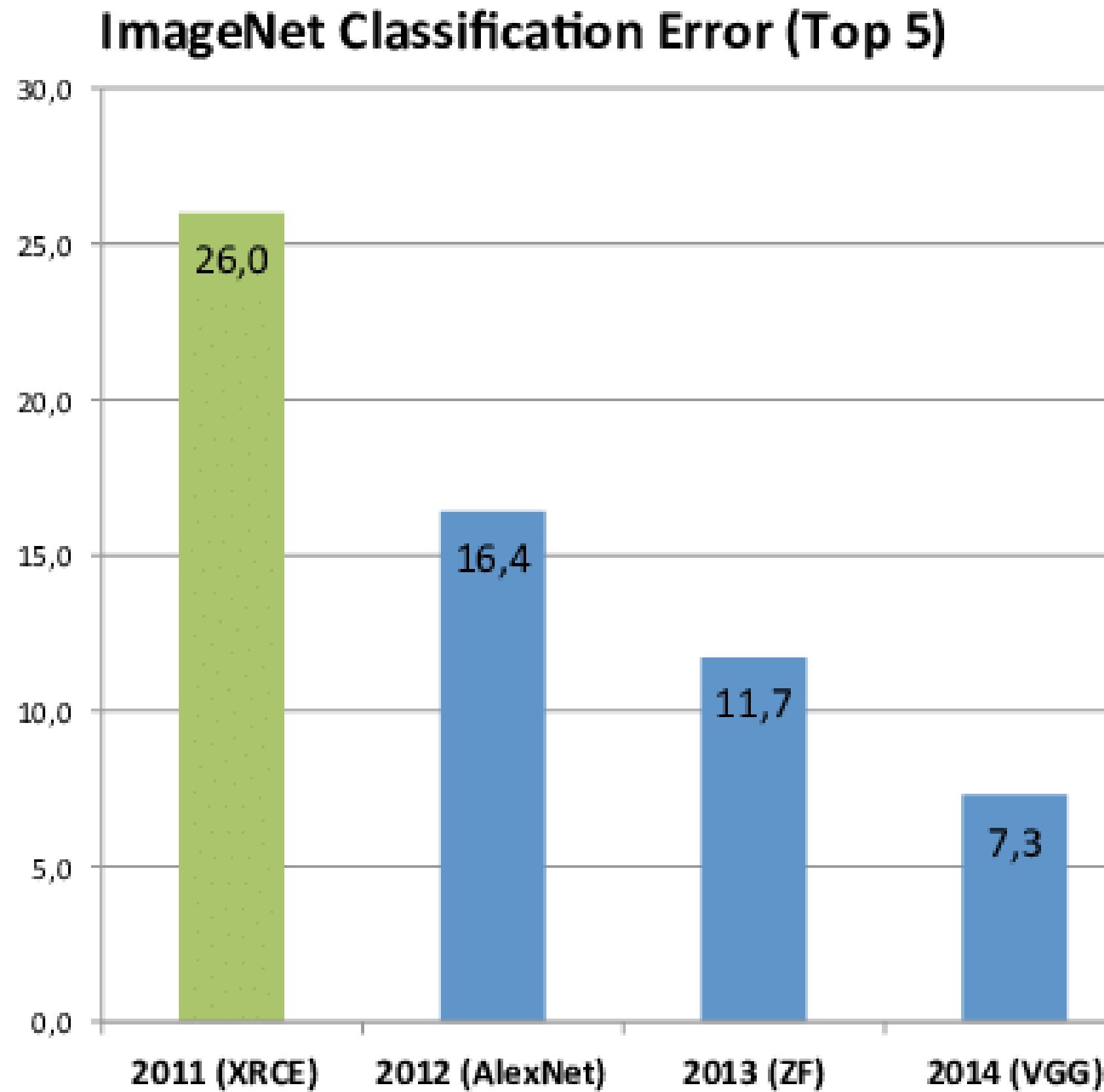


Error: 16.4%

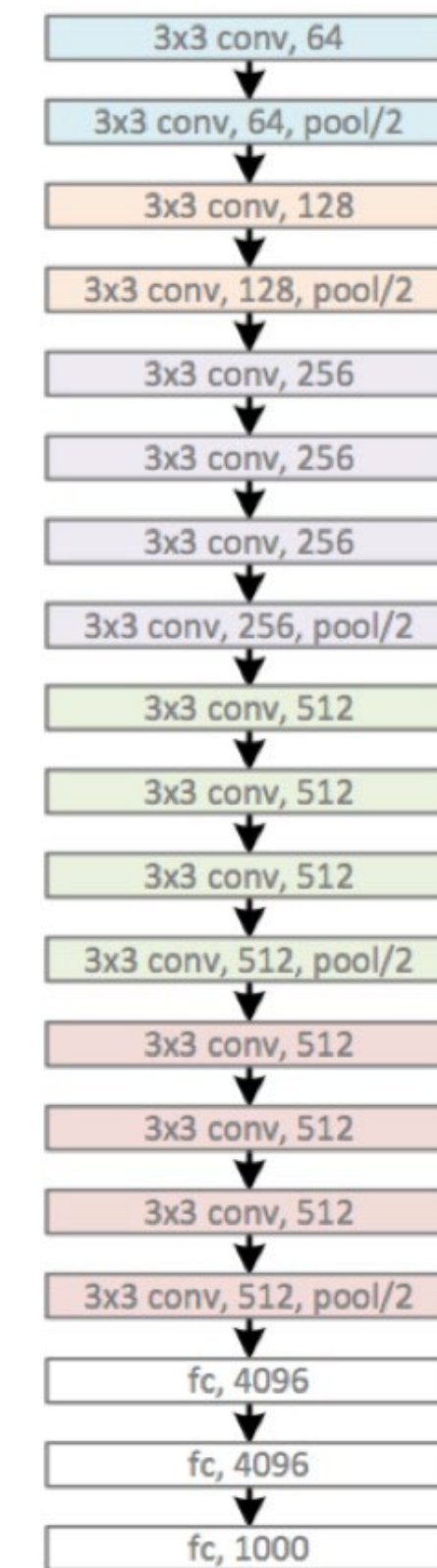
[Krizhevsky et al: ImageNet Classification with Deep Convolutional Neural Networks, NeurIPS 2012]

Alexnet — [Krizhevsky et al. NIPS 2012]





2014: VGG
16 conv. layers

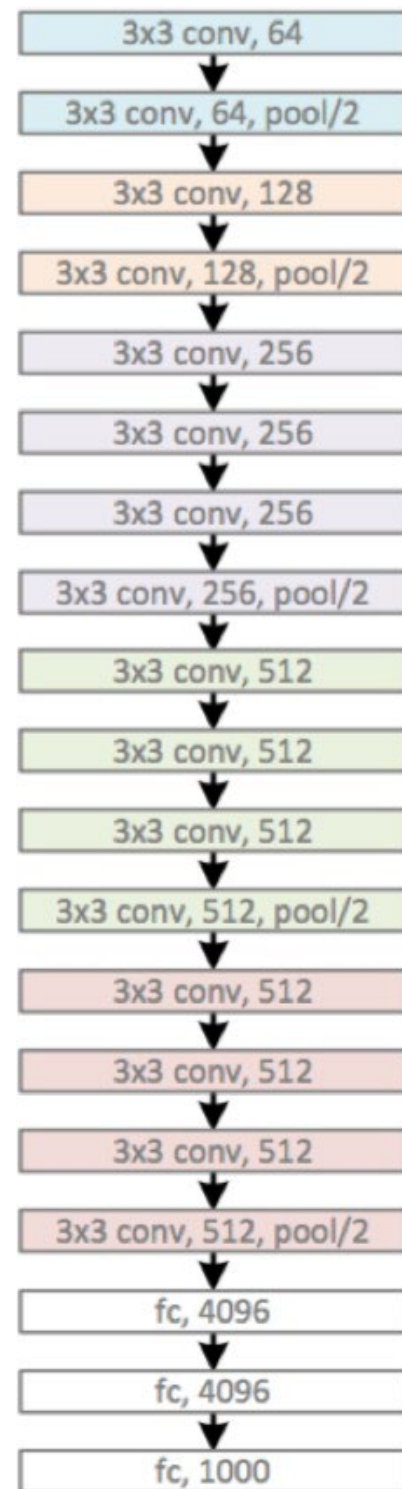


Error: 7.3%

[Simonyan & Zisserman: Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR 2015]

VGG-Net [Simonyan & Zisserman, 2015]

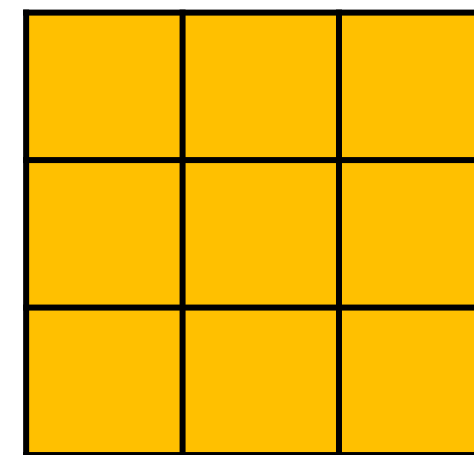
2014: VGG
16 conv. layers



Error: 7.3%

Main developments

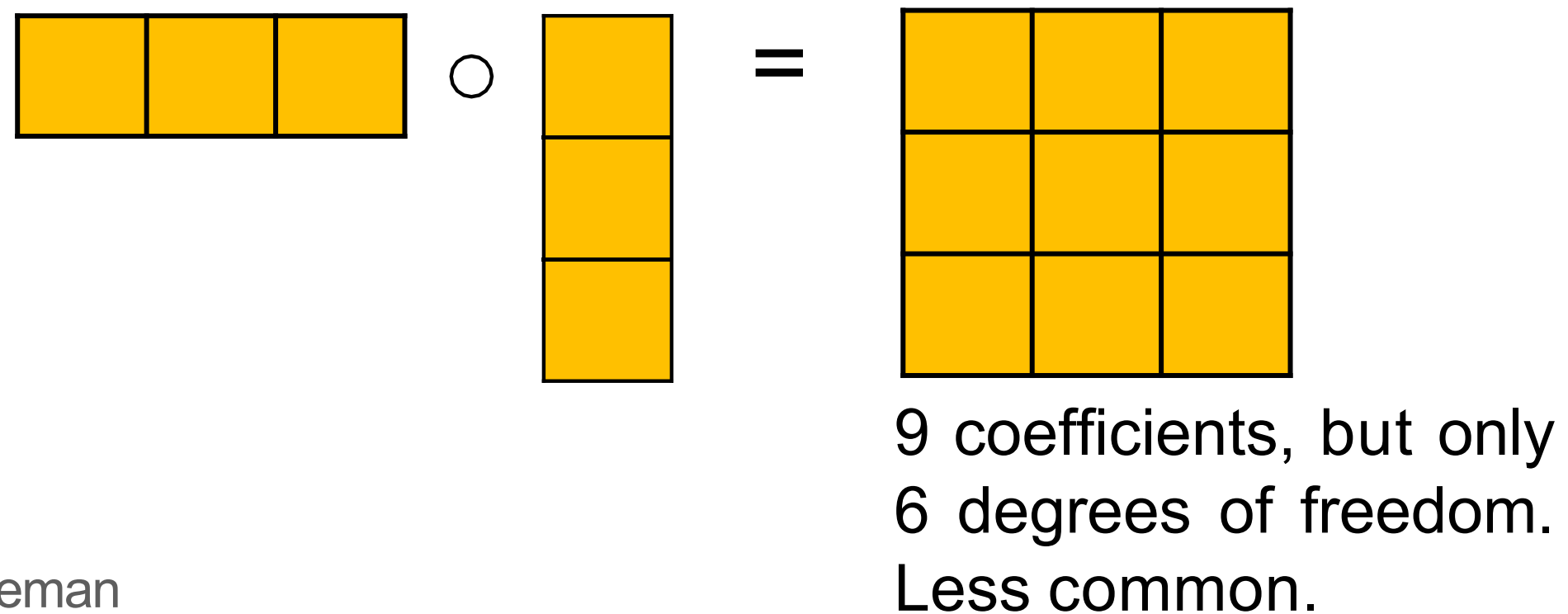
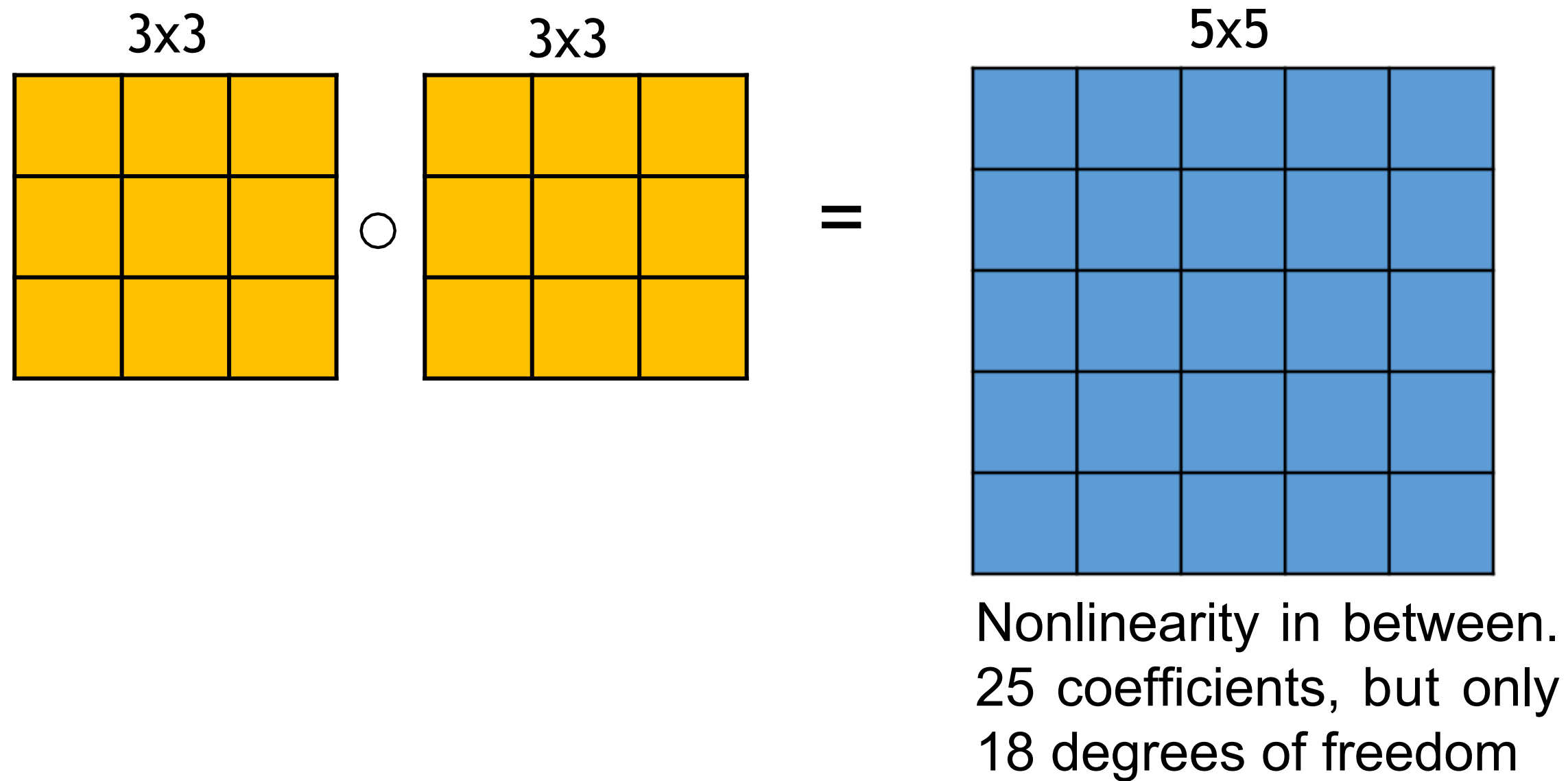
- Small convolutional kernels: only 3x3



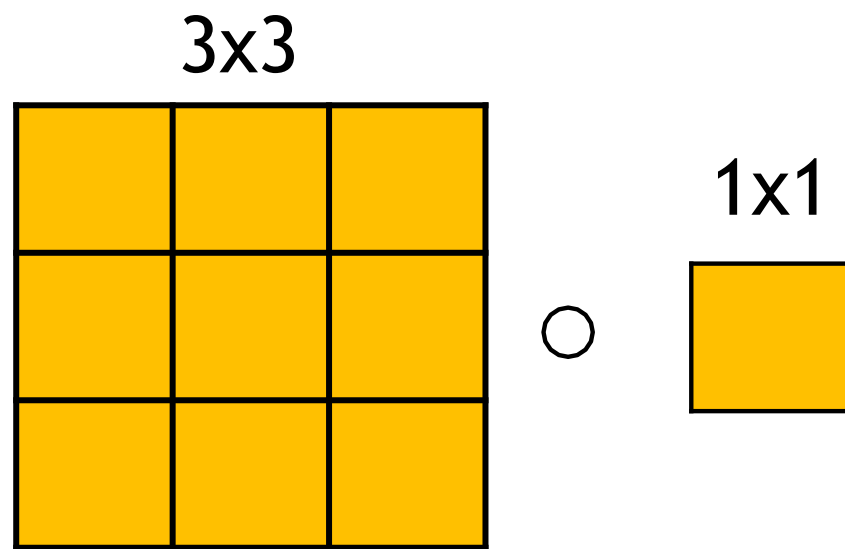
- Increased depth (5 -> 16/19 layers)

Other tricks for designing convolutional nets

Chaining convolutions



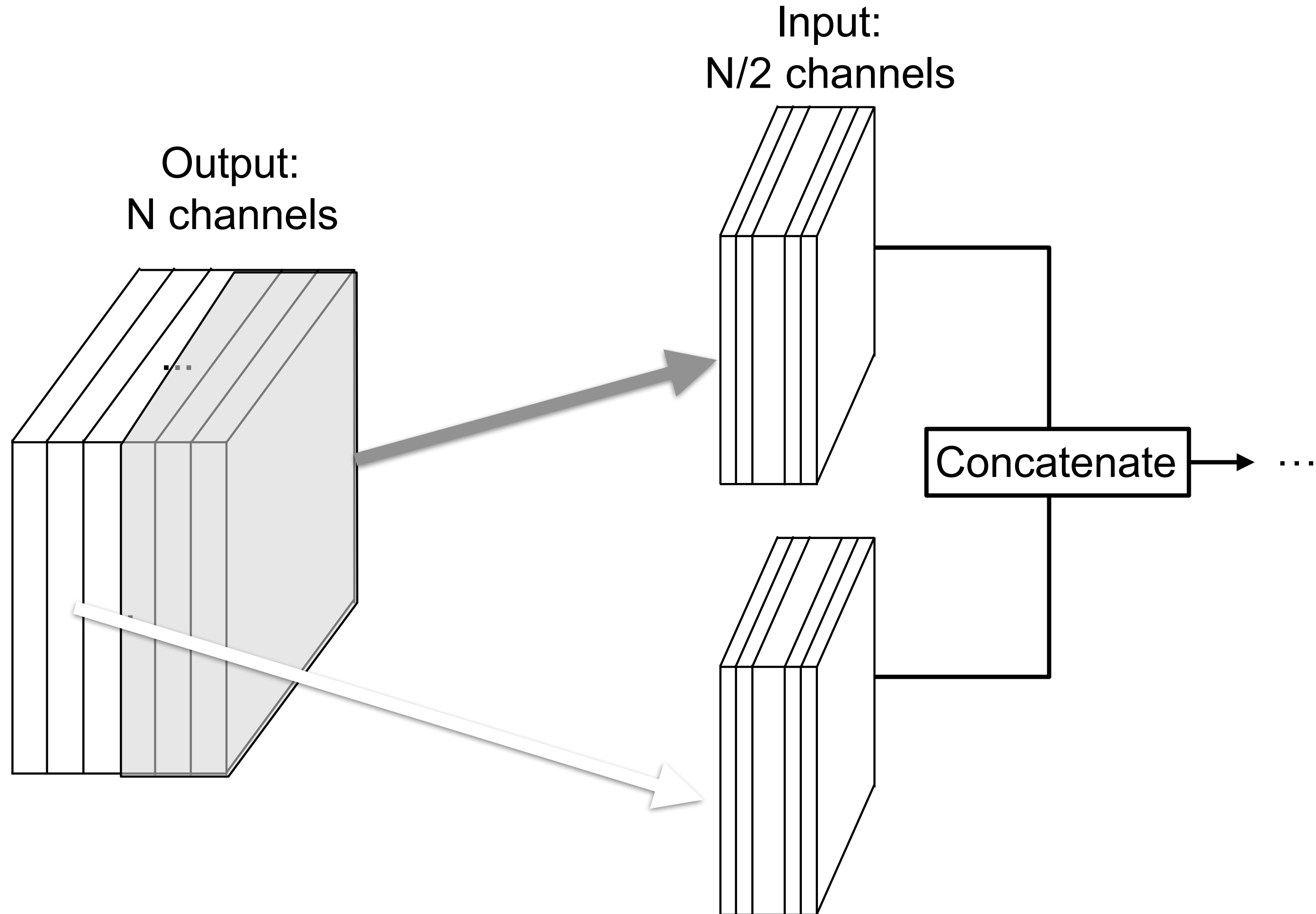
1x1 convolutions



Why do this?

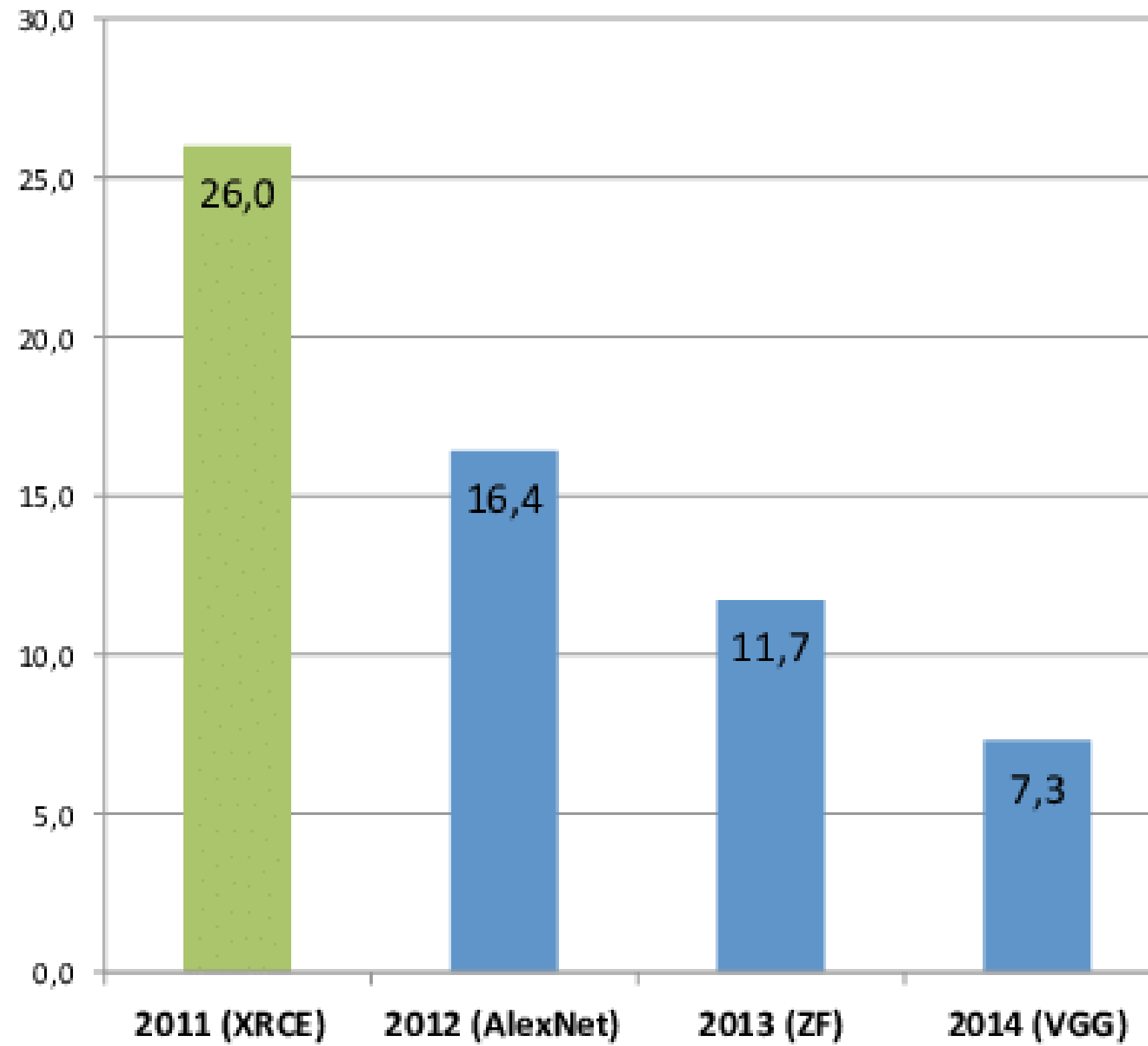
(nonlinearity in between)

Grouped Convolutions

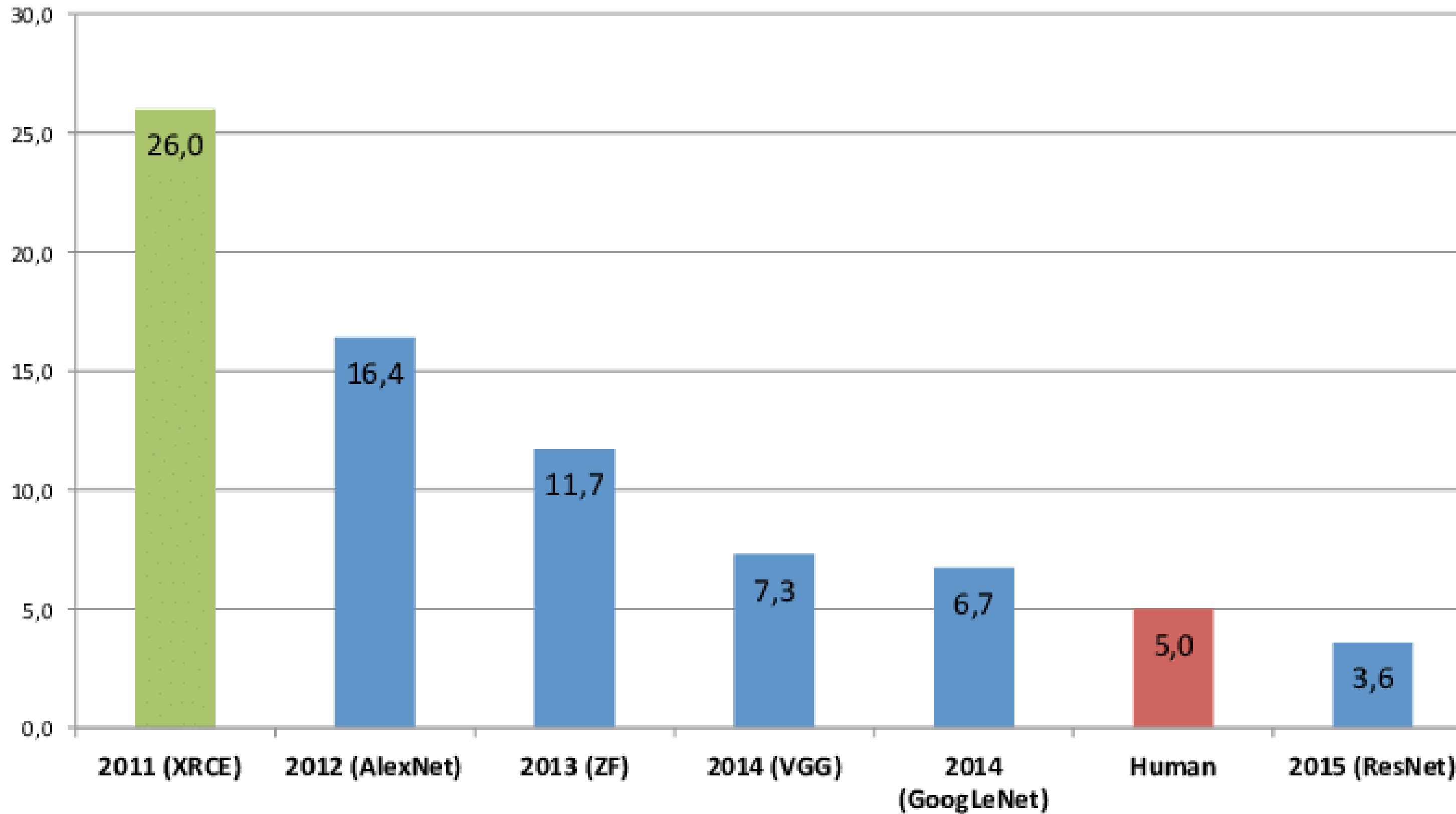


Split channels into N groups, and process separately with N convolution⁵⁸ layers.

ImageNet Classification Error (Top 5)



ImageNet Classification Error (Top 5)



2016: ResNet
>100 conv. layers

Error: 3.6%

[He et al: Deep Residual Learning for Image Recognition, CVPR 2016]

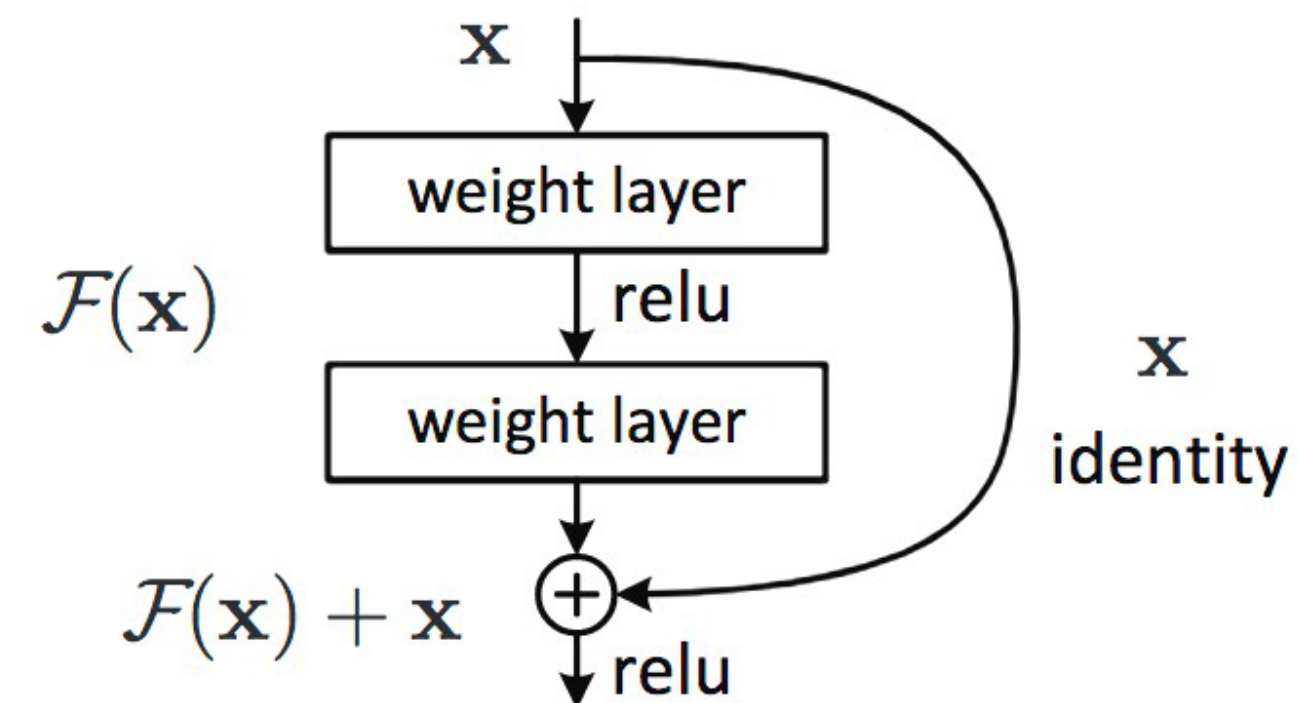
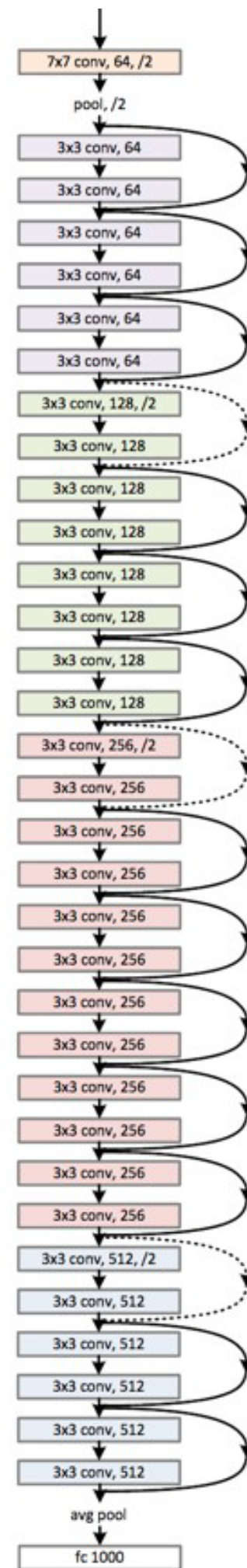
2016: ResNet
>100 conv. layers

ResNet [He et al, 2016]

Main developments

- Increased depth possible through residual blocks

Error: 3.6%



Residual Blocks

Problem: Hard to train very deep nets (50+ layers). This is an optimization issue, not overfitting: shallow models often get higher *training* accuracy than deep ones!

Idea: Make it easy to represent for the network to implement the identity.

Normal convolution + relu:

$$x_{i+1} = \text{relu}(x_i \text{ of})$$

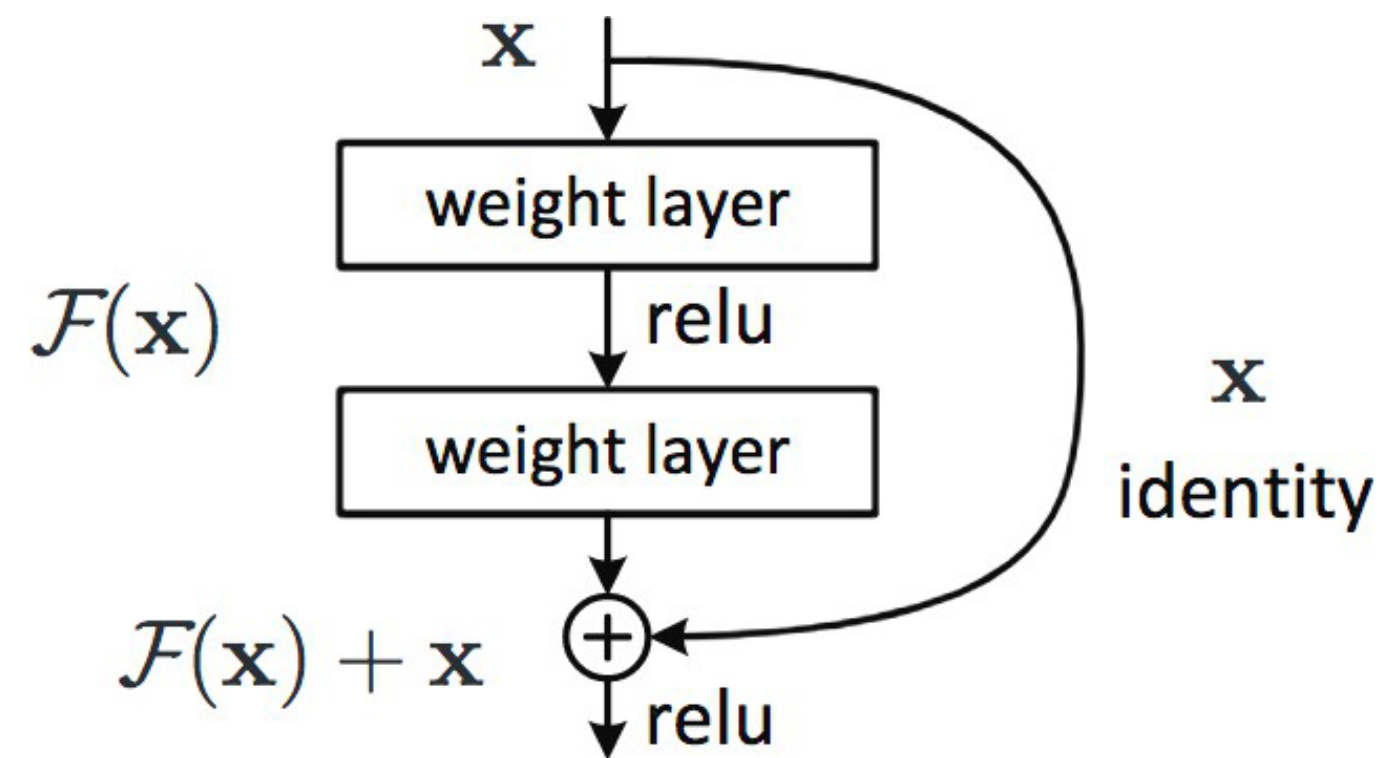
Residual connection

$$x_{i+1} = \text{relu}((x_i \text{ of}) + x_i)$$

In general, do multiple convolutions (with nonlinearities) before summing:

$$x_{i+1} = \text{relu}(F(x_i) + x_i)$$

Residual Blocks



Why do they work?

- Gradients can propagate faster (via the identity mapping)
- Within each block, only small residuals have to be learned

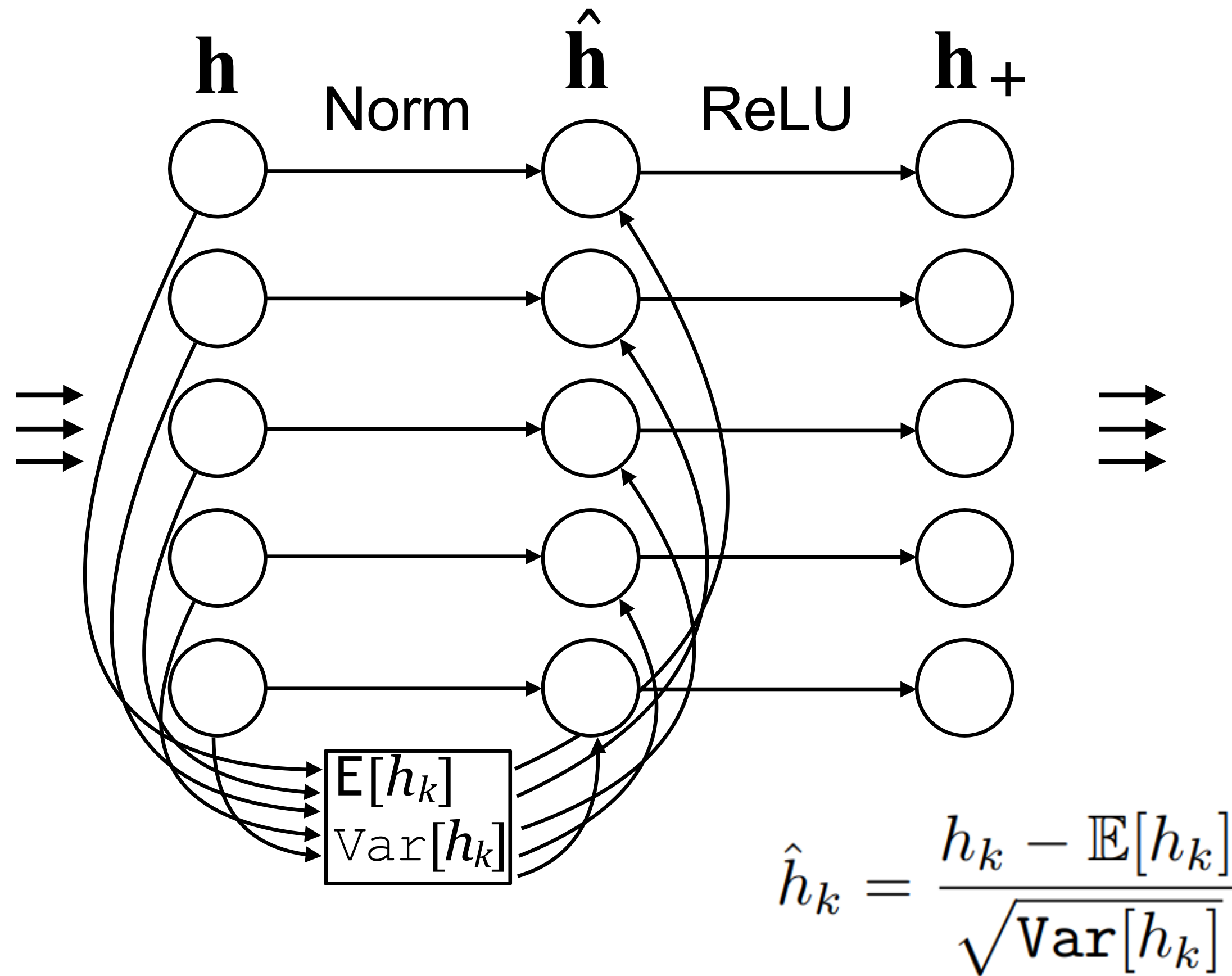
Normalization layers

Standardize activations by subtracting mean and dividing by standard deviation (averaged over all spatial locations).

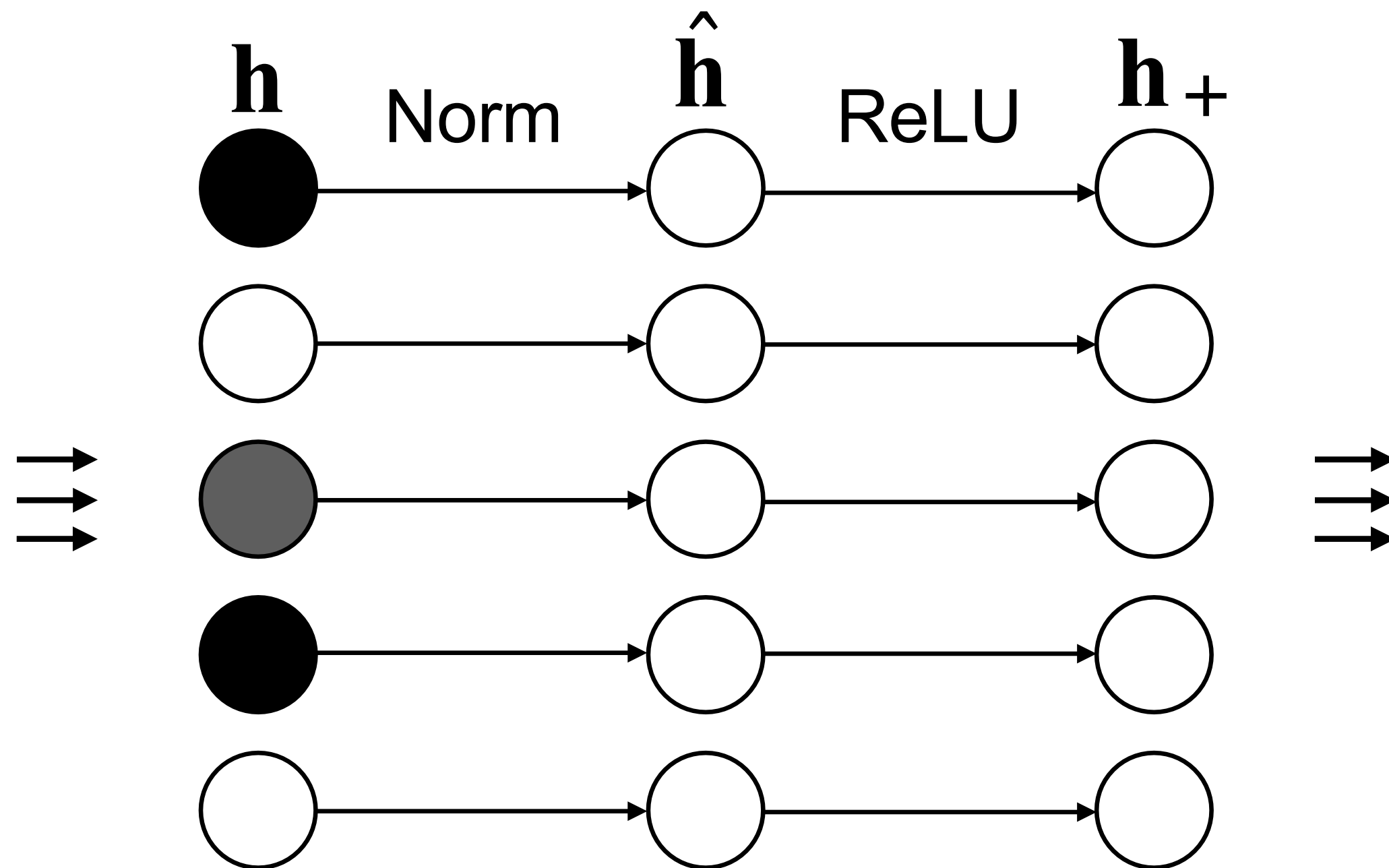
This provides a constant “interface” for later layers of the networks.
Ensures that the previous layer will have zero mean,

Obtains invariance to mean and variance.

Normalization layers

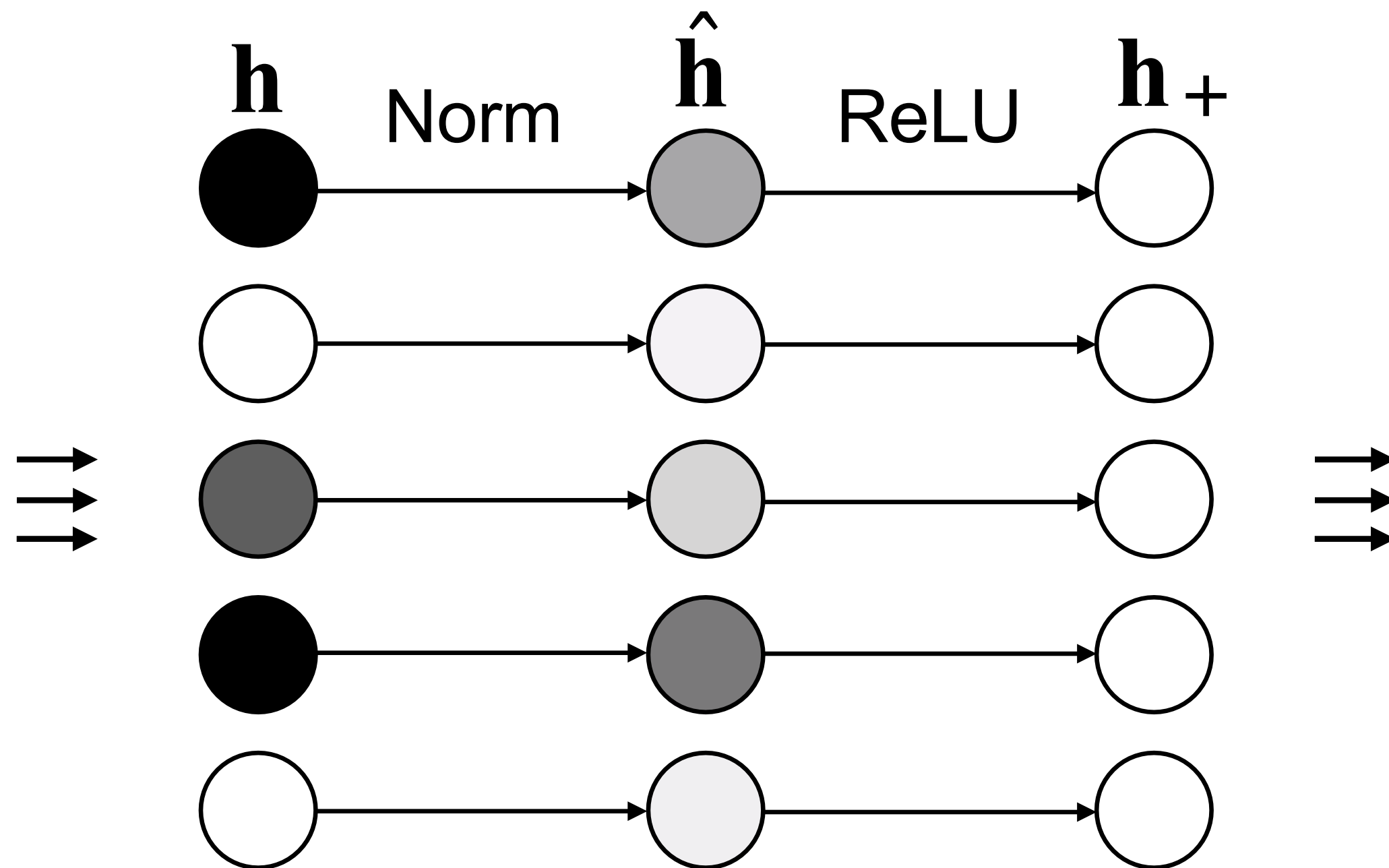


Normalization layers



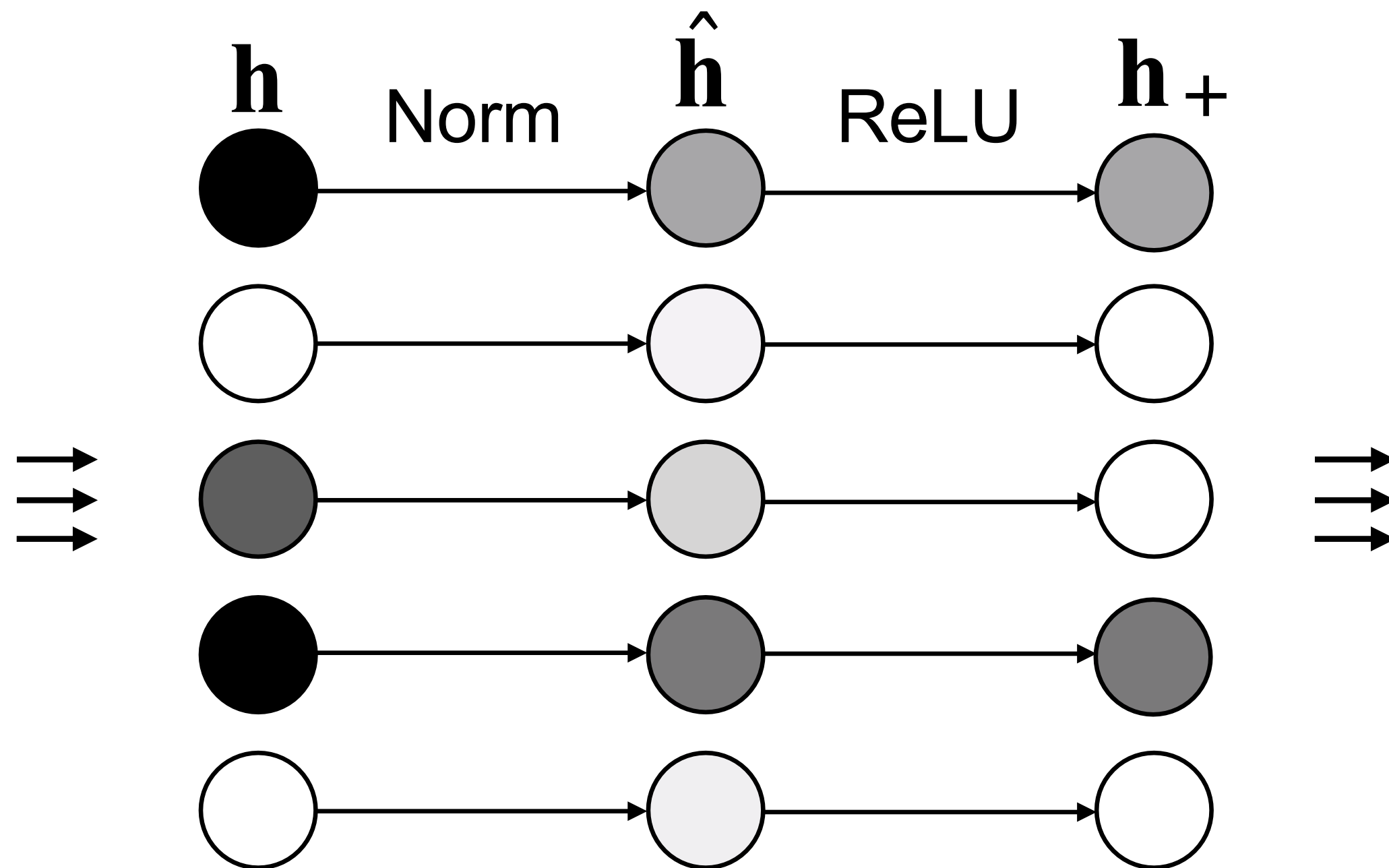
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers



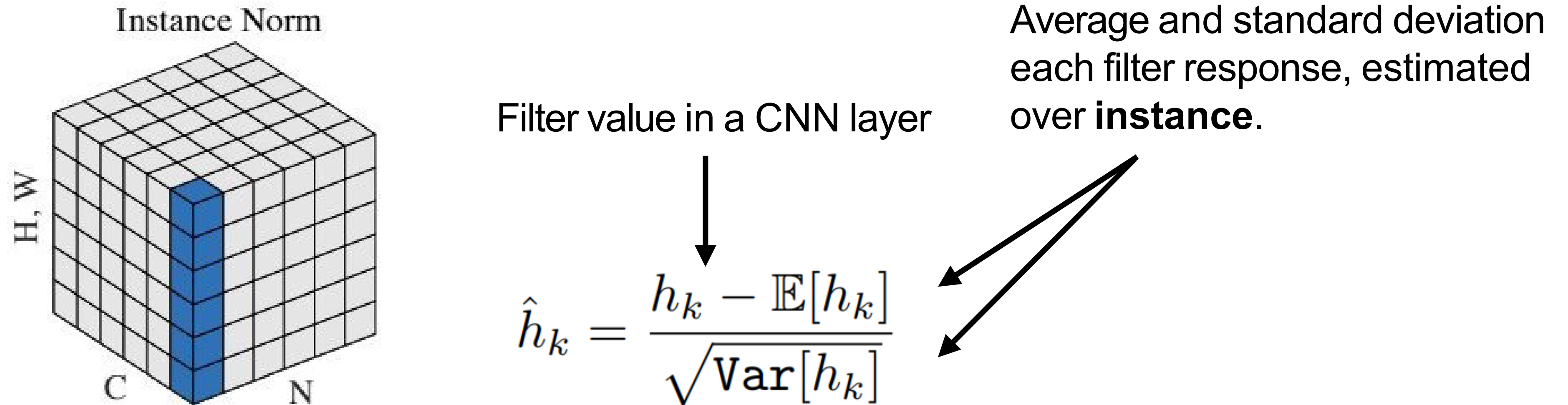
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers



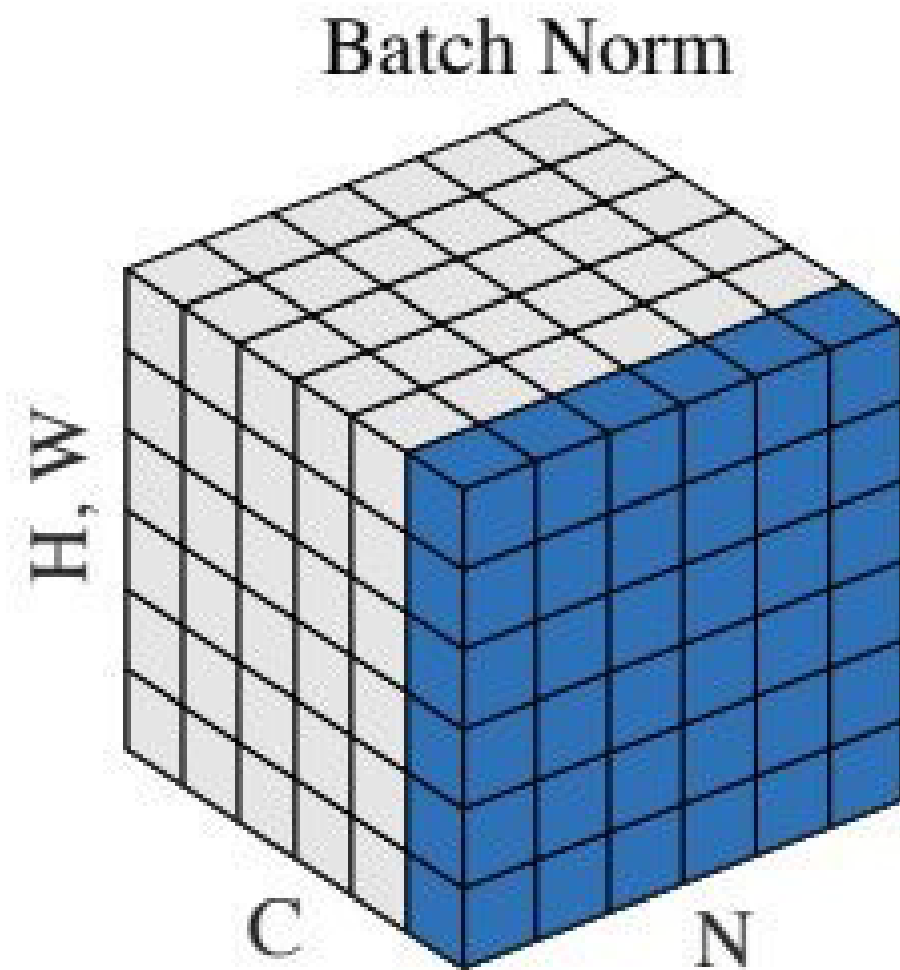
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Instance normalization



Normalize a single hidden unit's activations to be mean 0, standard deviation 1.

Batch normalization



Filter value in a CNN layer

$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Average and standard deviation
each filter response, estimated
over **whole batch**.

Normalize a single hidden unit's activations to be mean 0, standard deviation 1.

At test time, remember the mean and standard deviation seen during training.

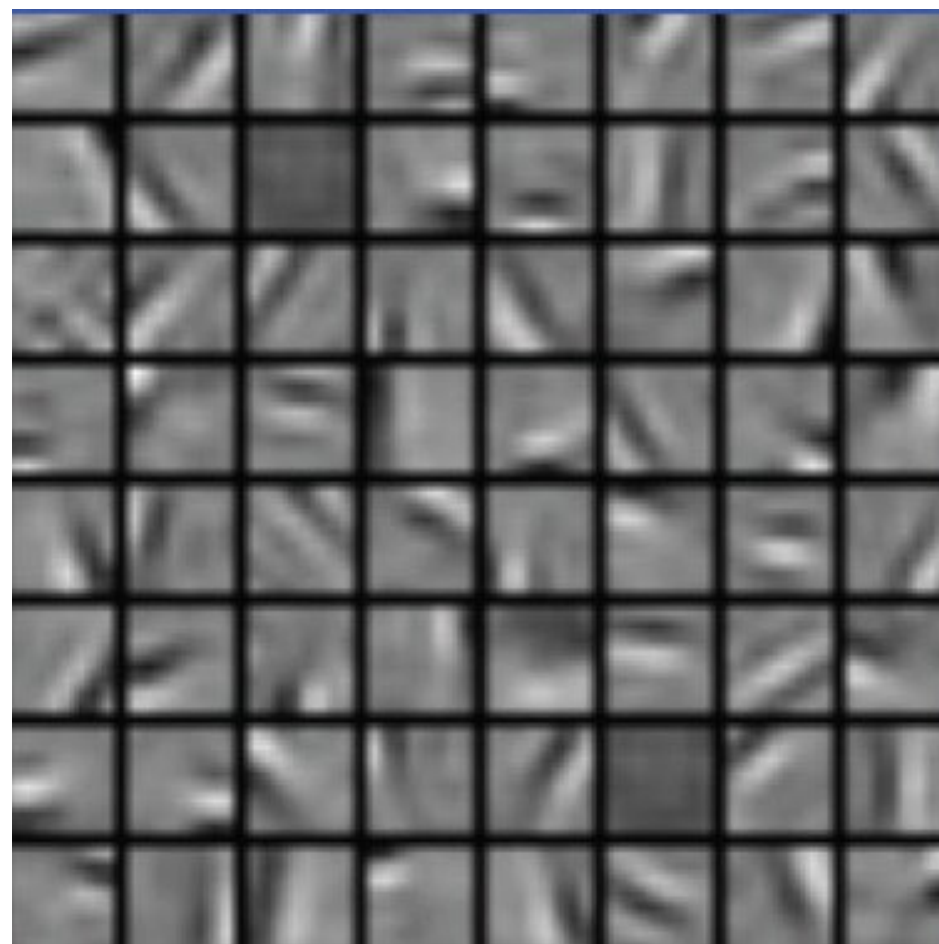
Can allow you to train with larger learning rate and significantly speed up training!

[Figure from Wu & He, arXiv 2018] [Ioffe & Szegedy, 2015]

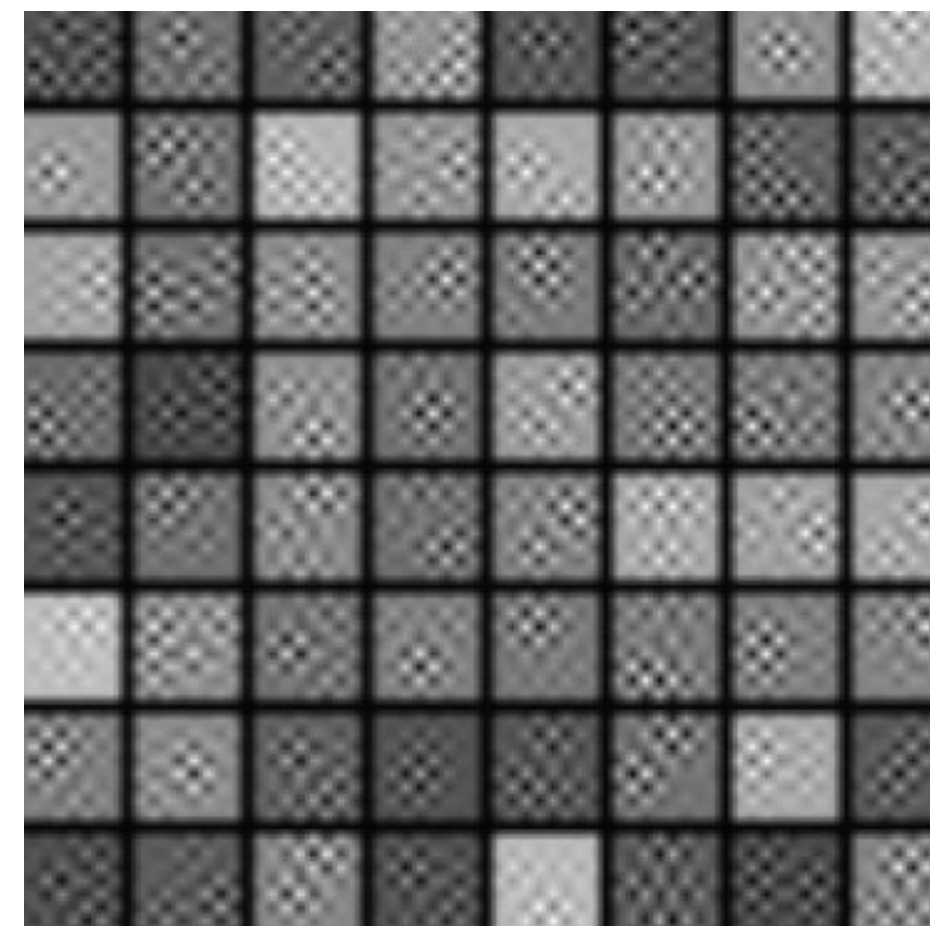
What filters are learned?

What filters are learned?

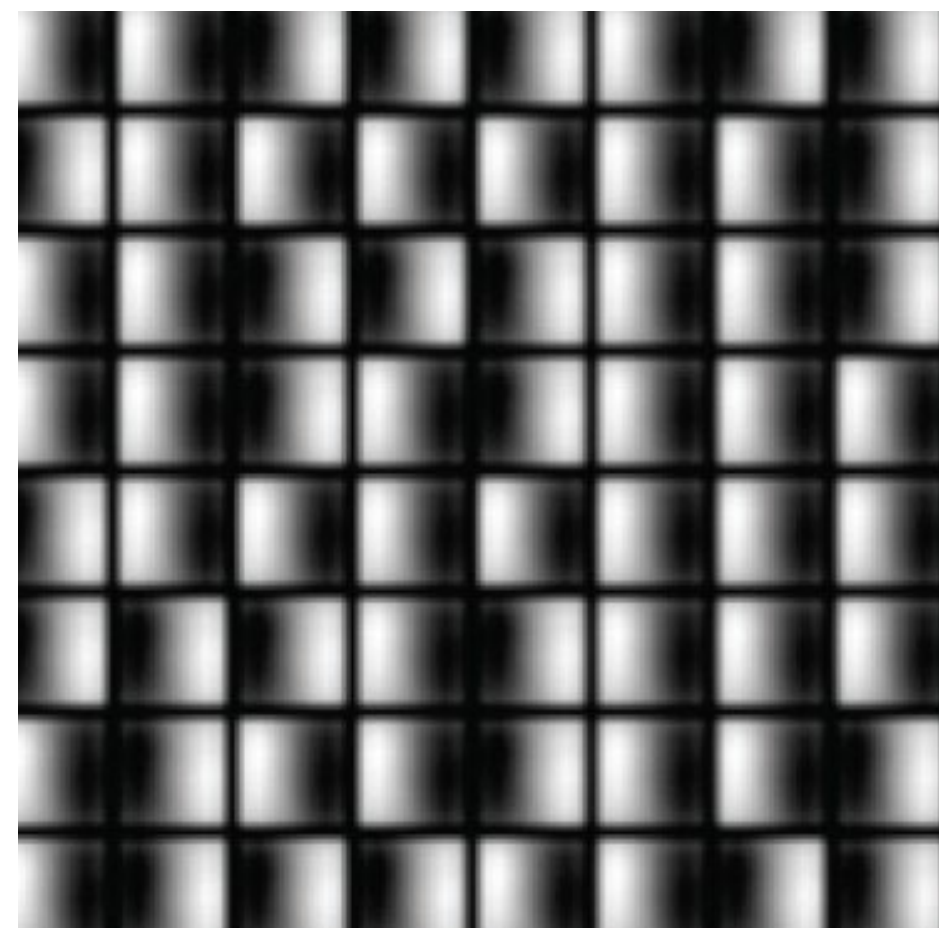
A



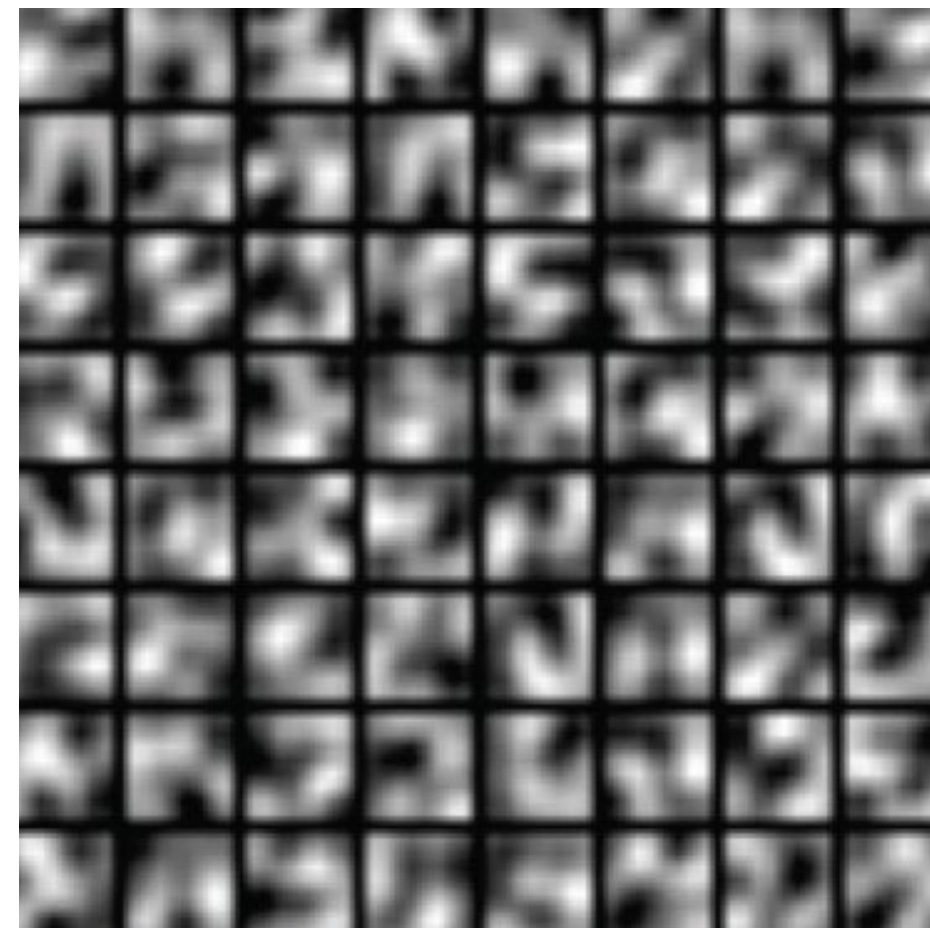
B



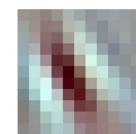
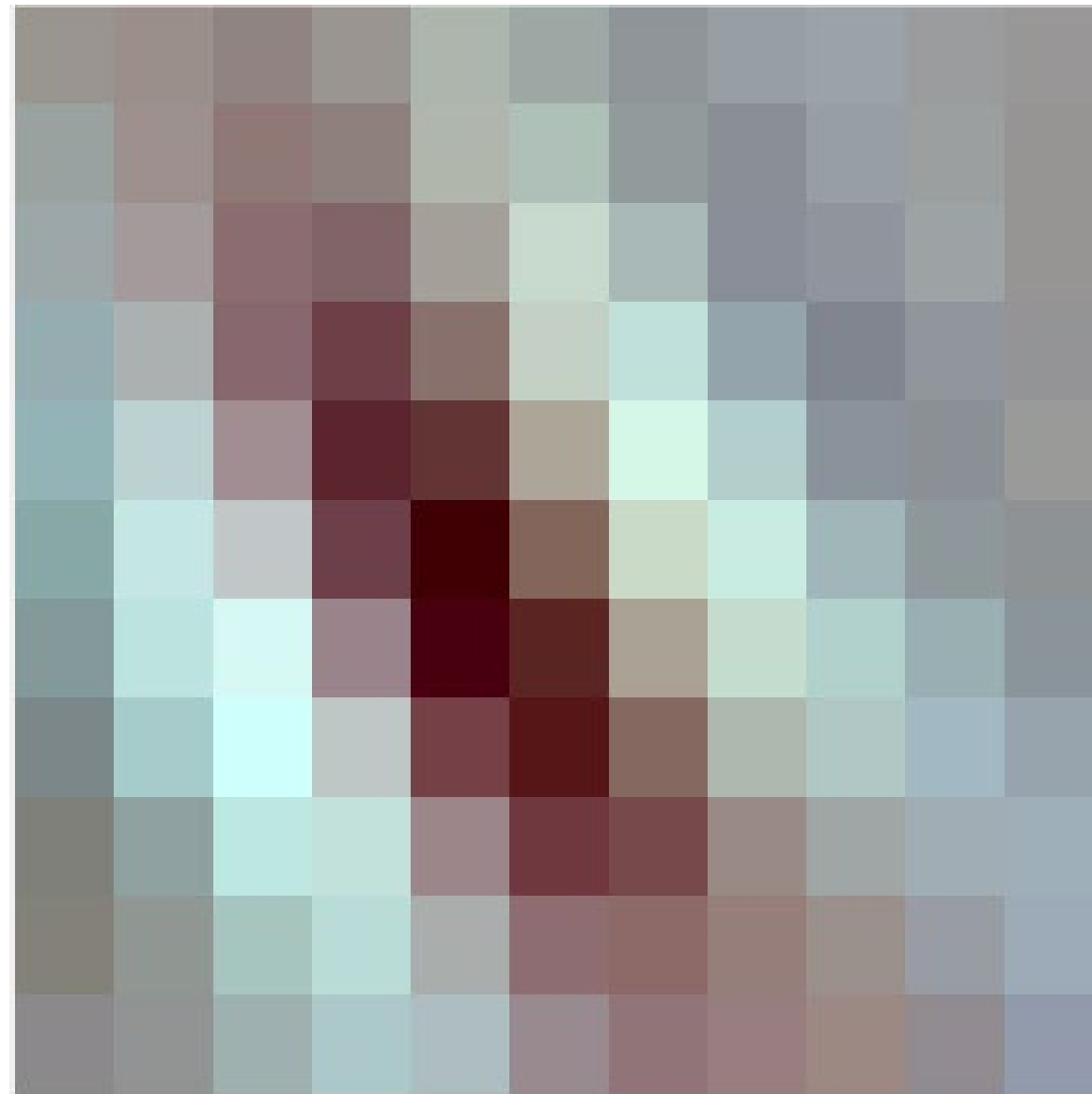
C



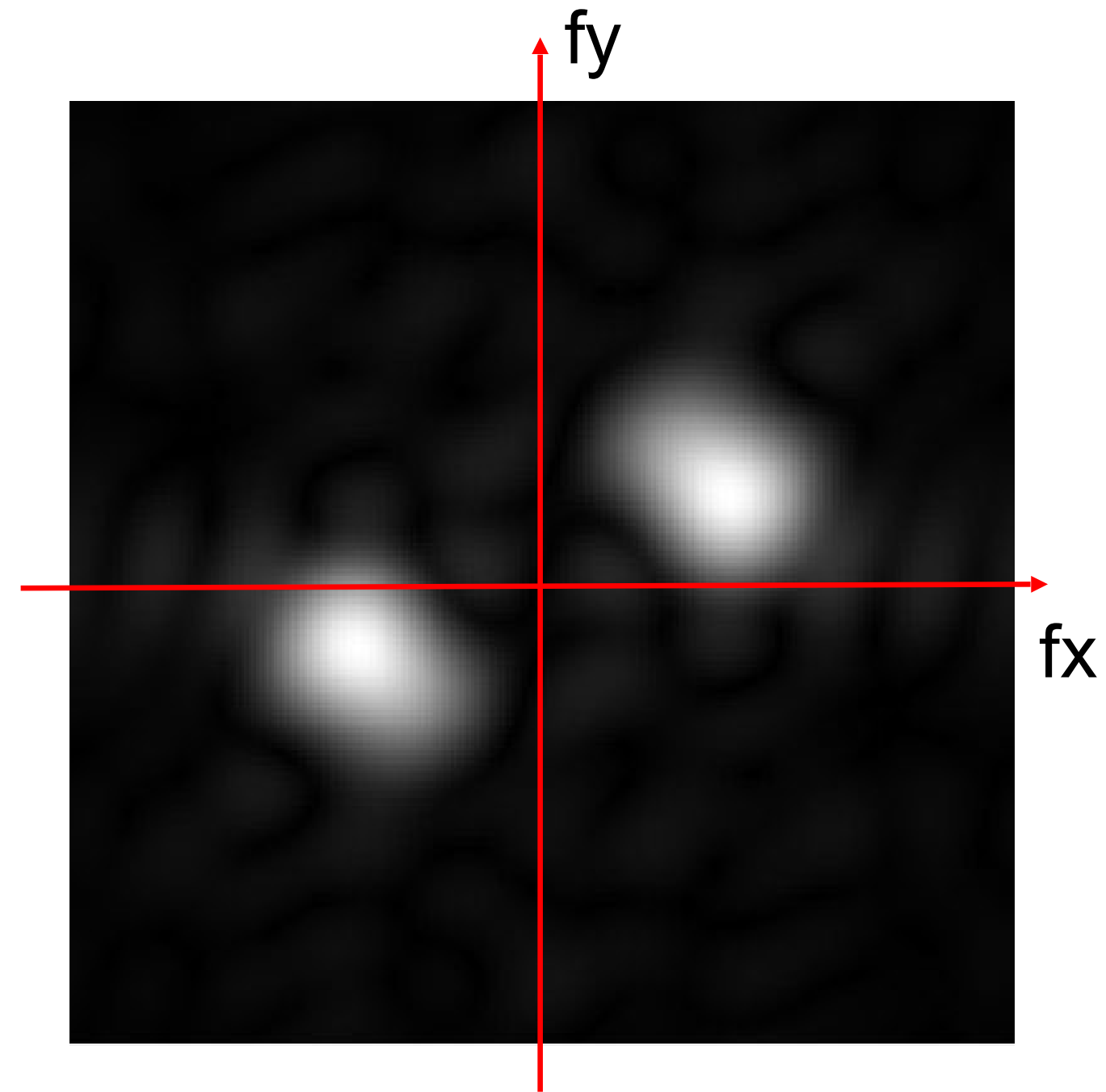
D



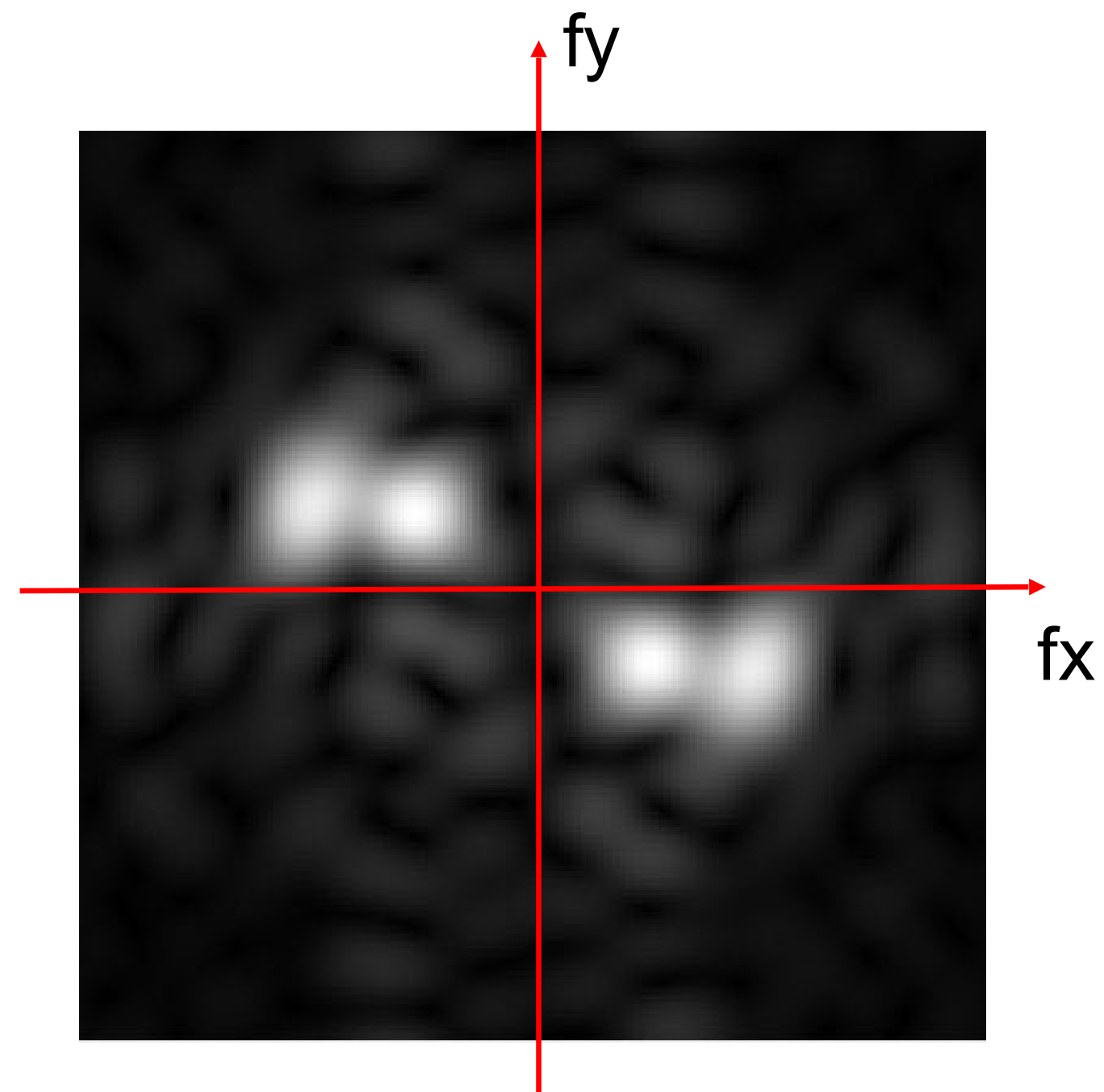
Filters in first layer



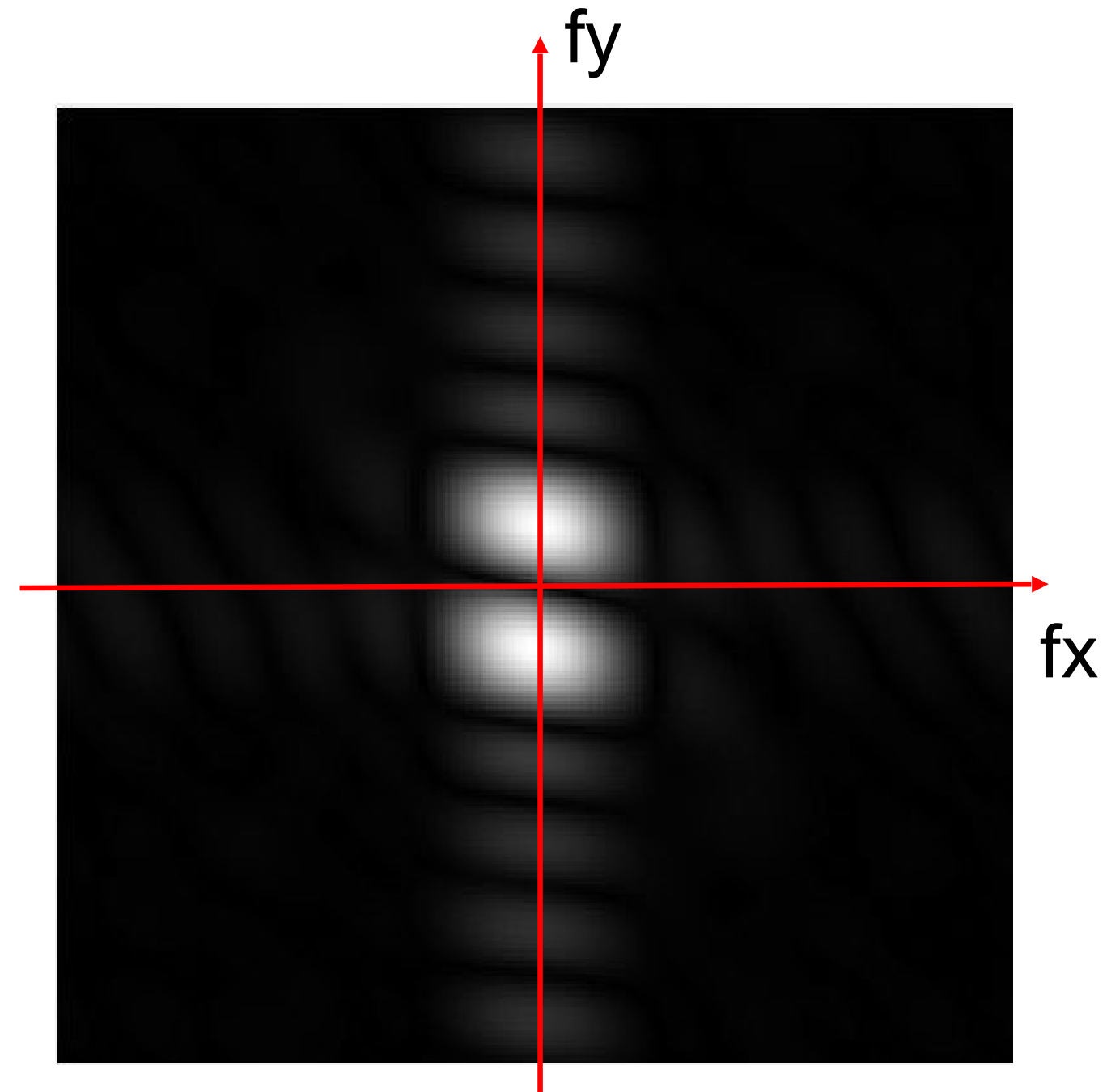
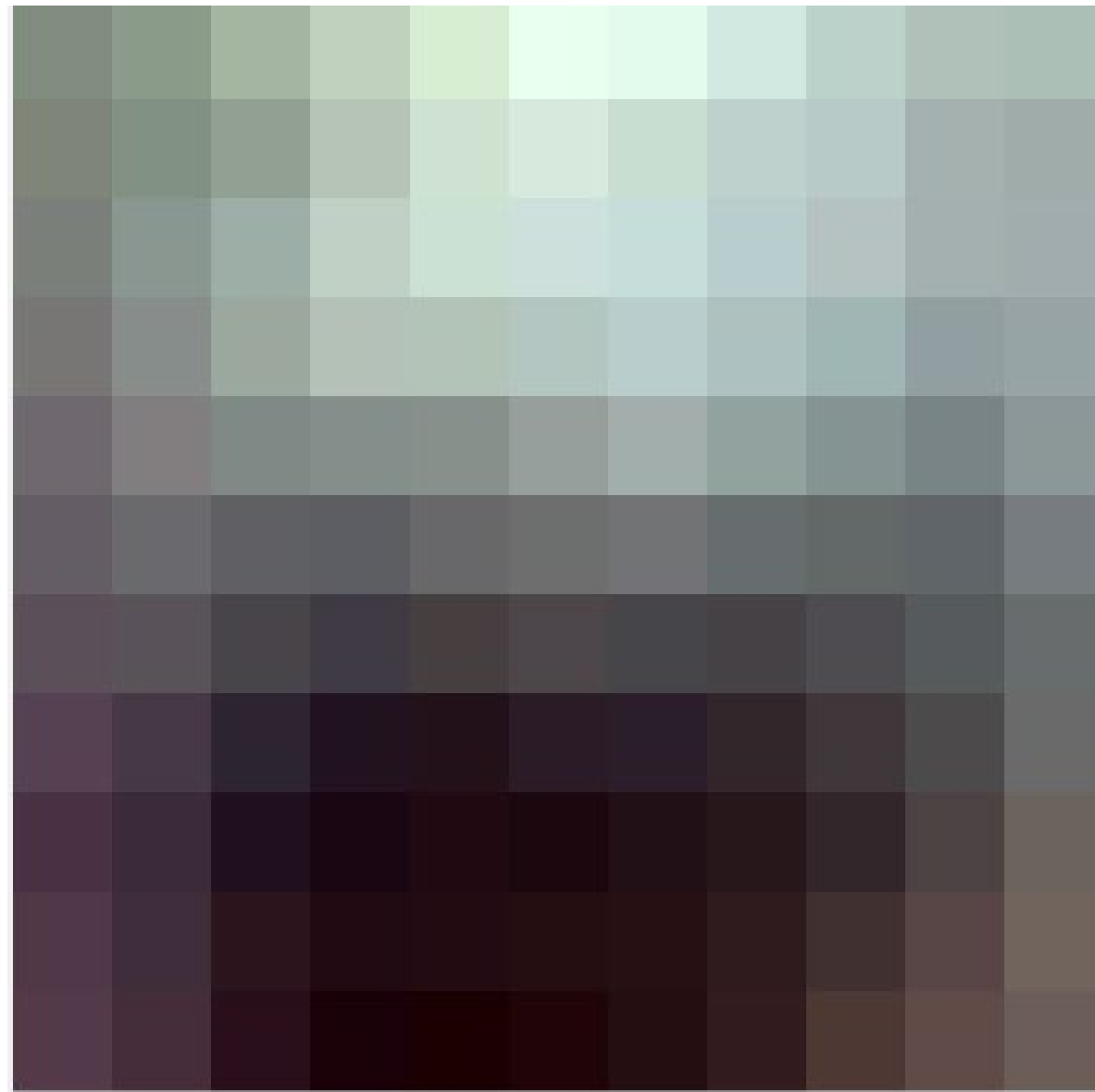
11x11 convolution kernel
(3 color channels)



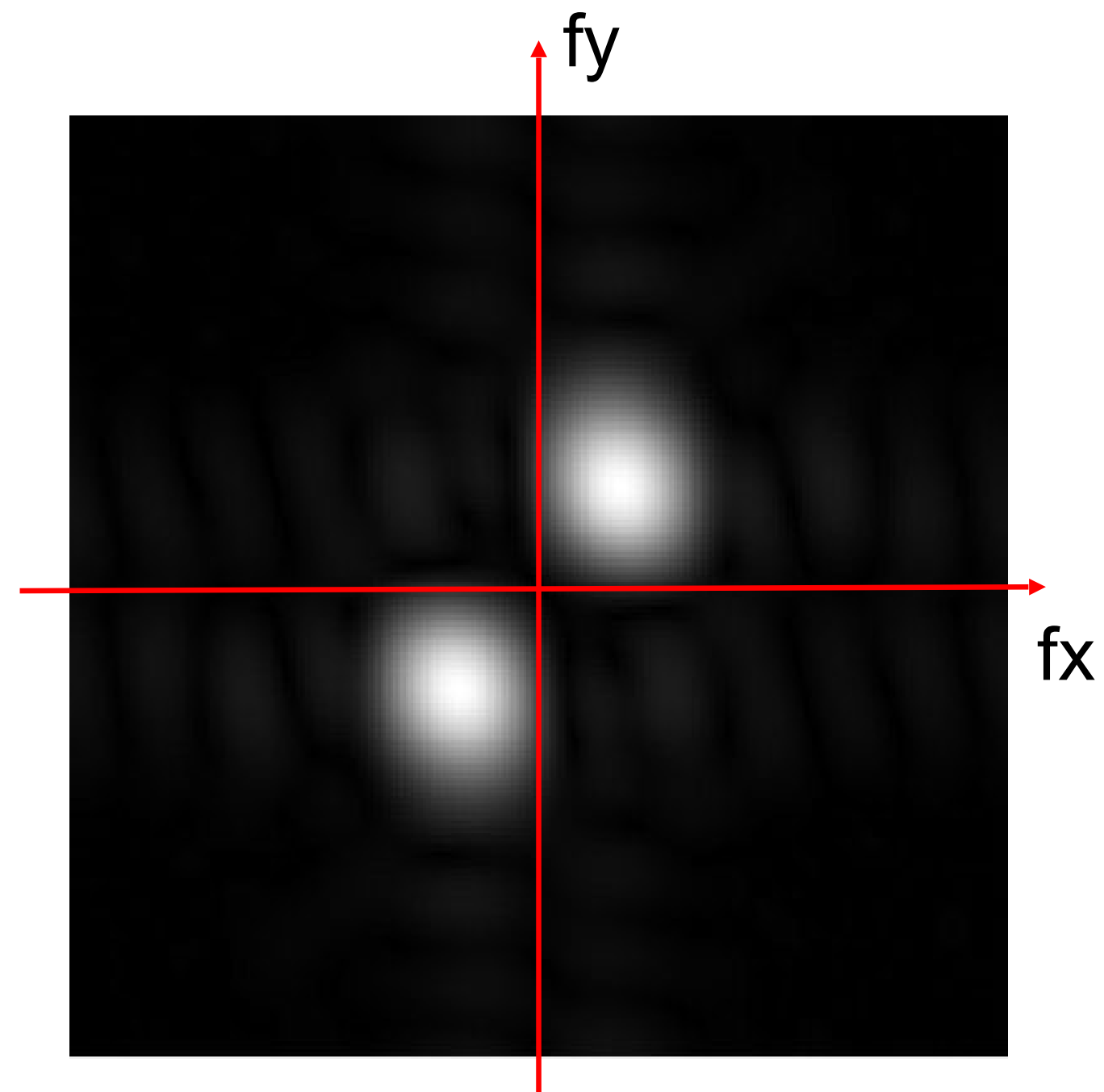
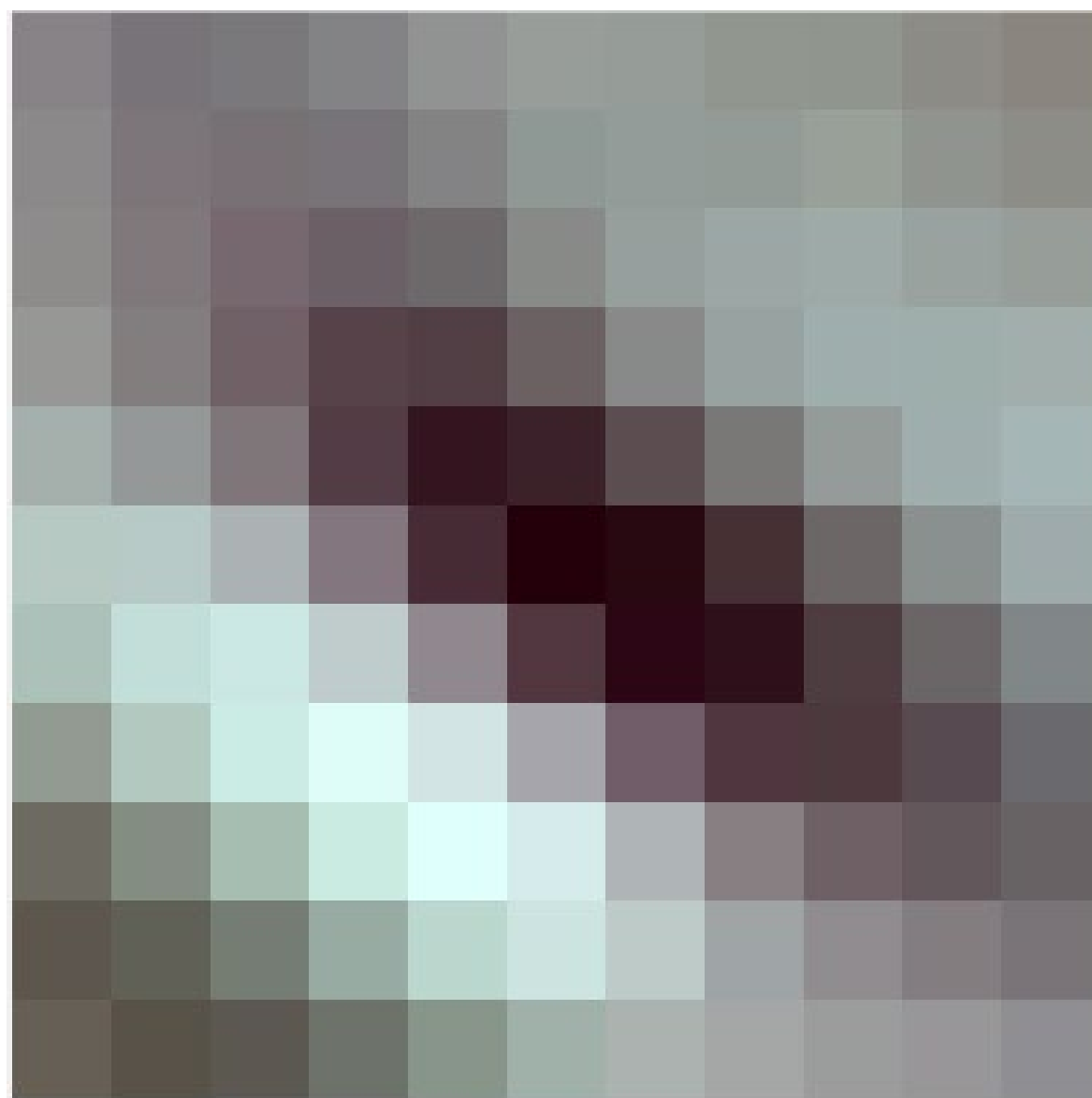
Filters in first layer



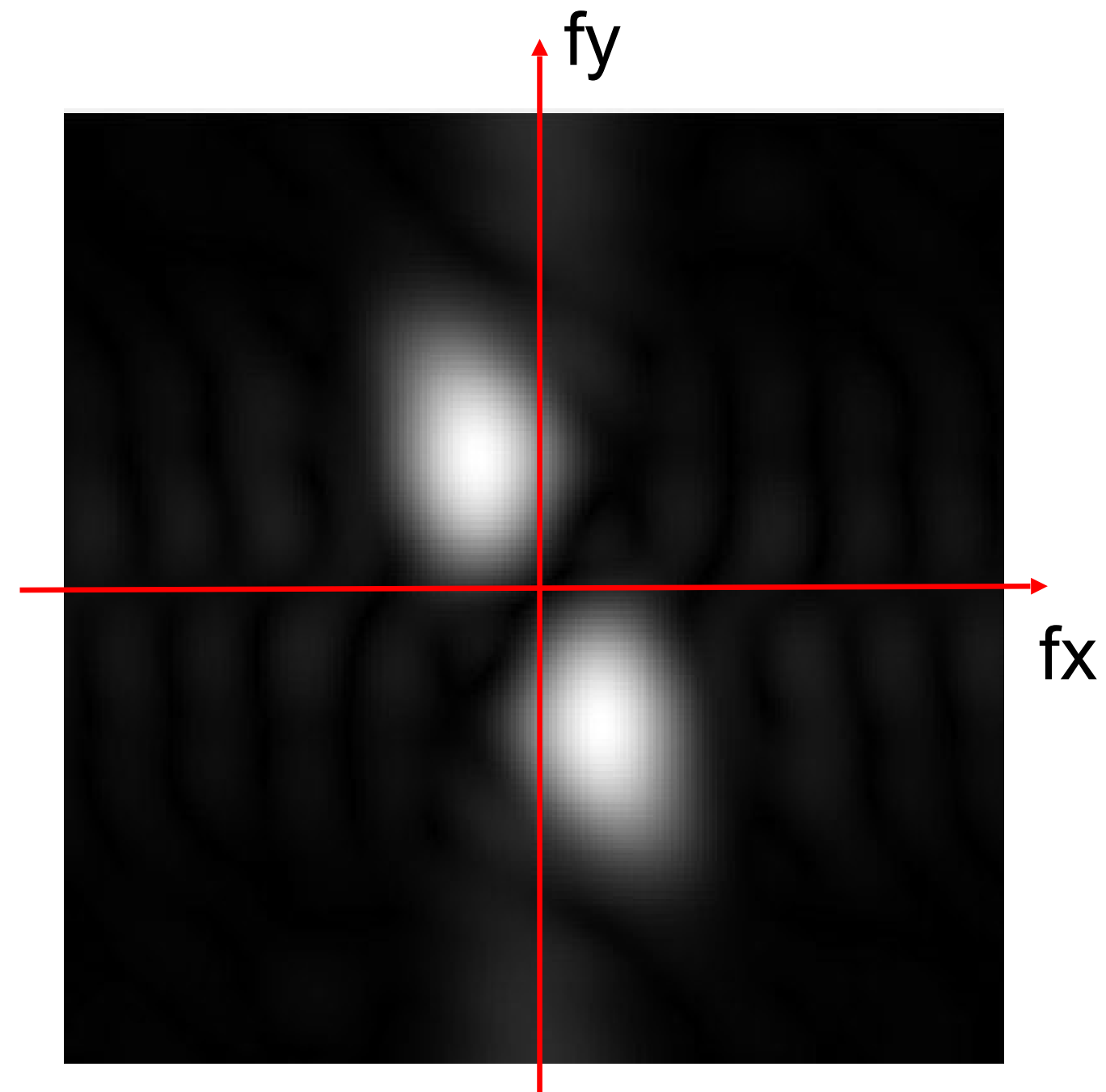
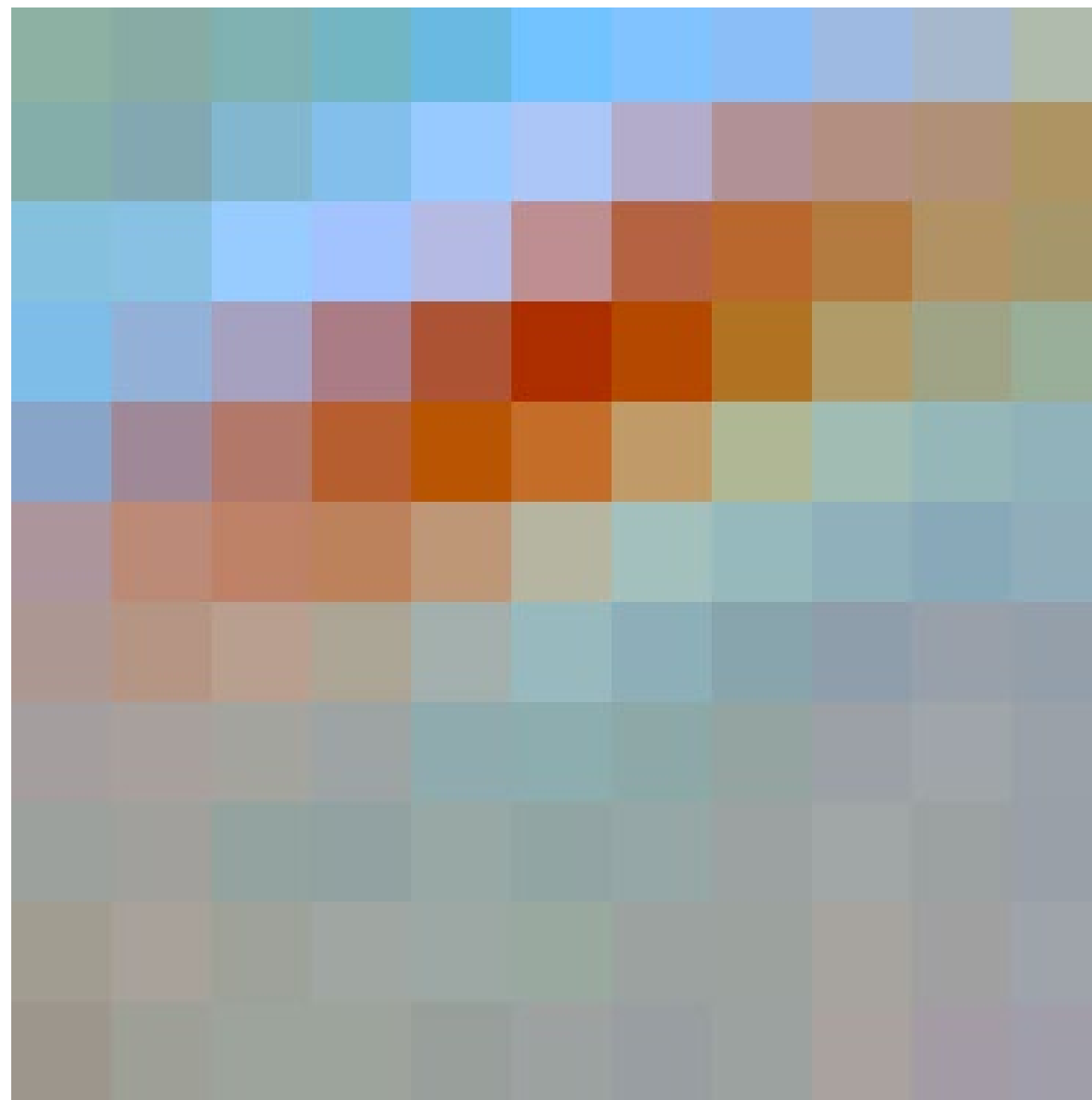
Filters in first layer



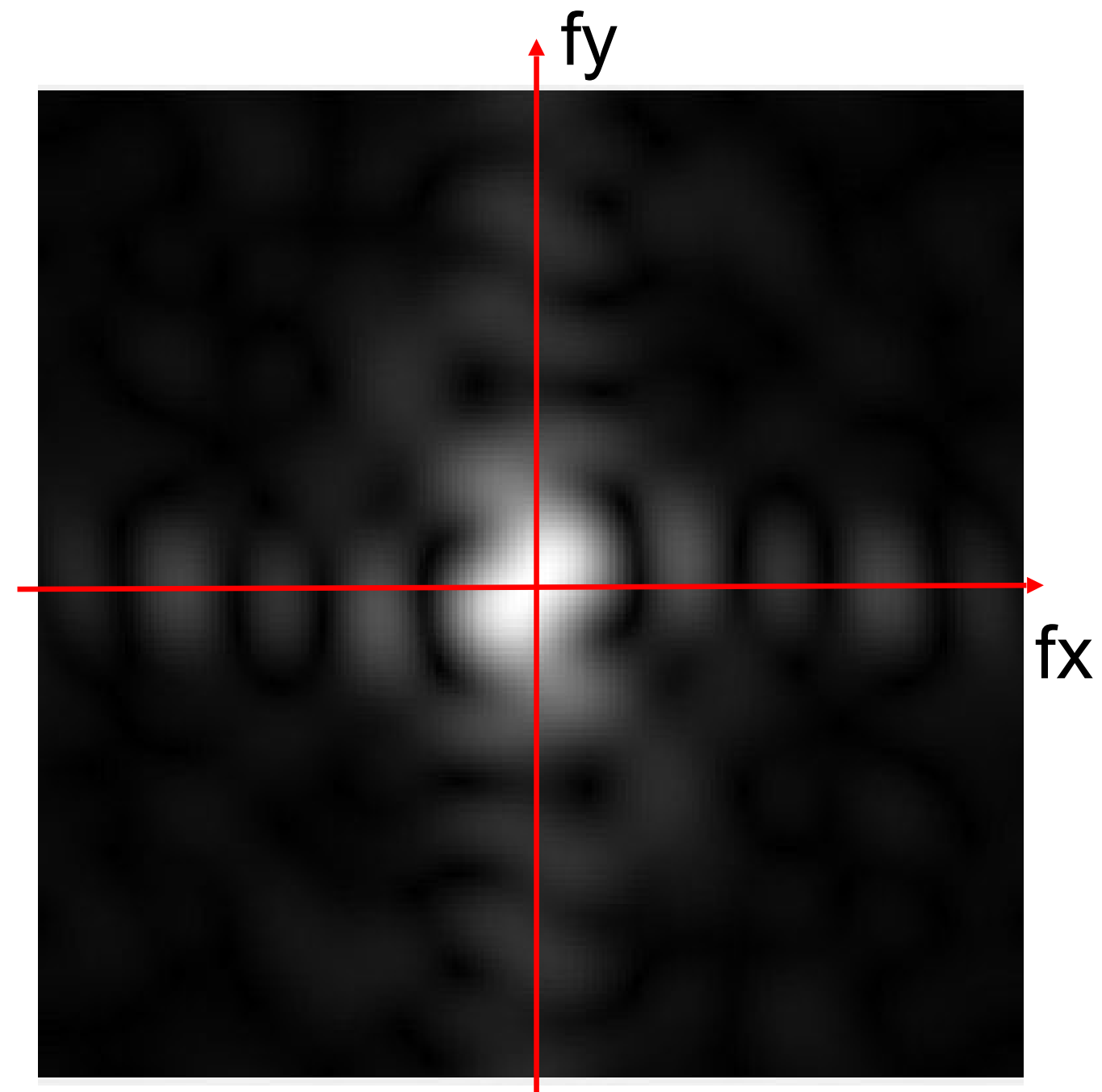
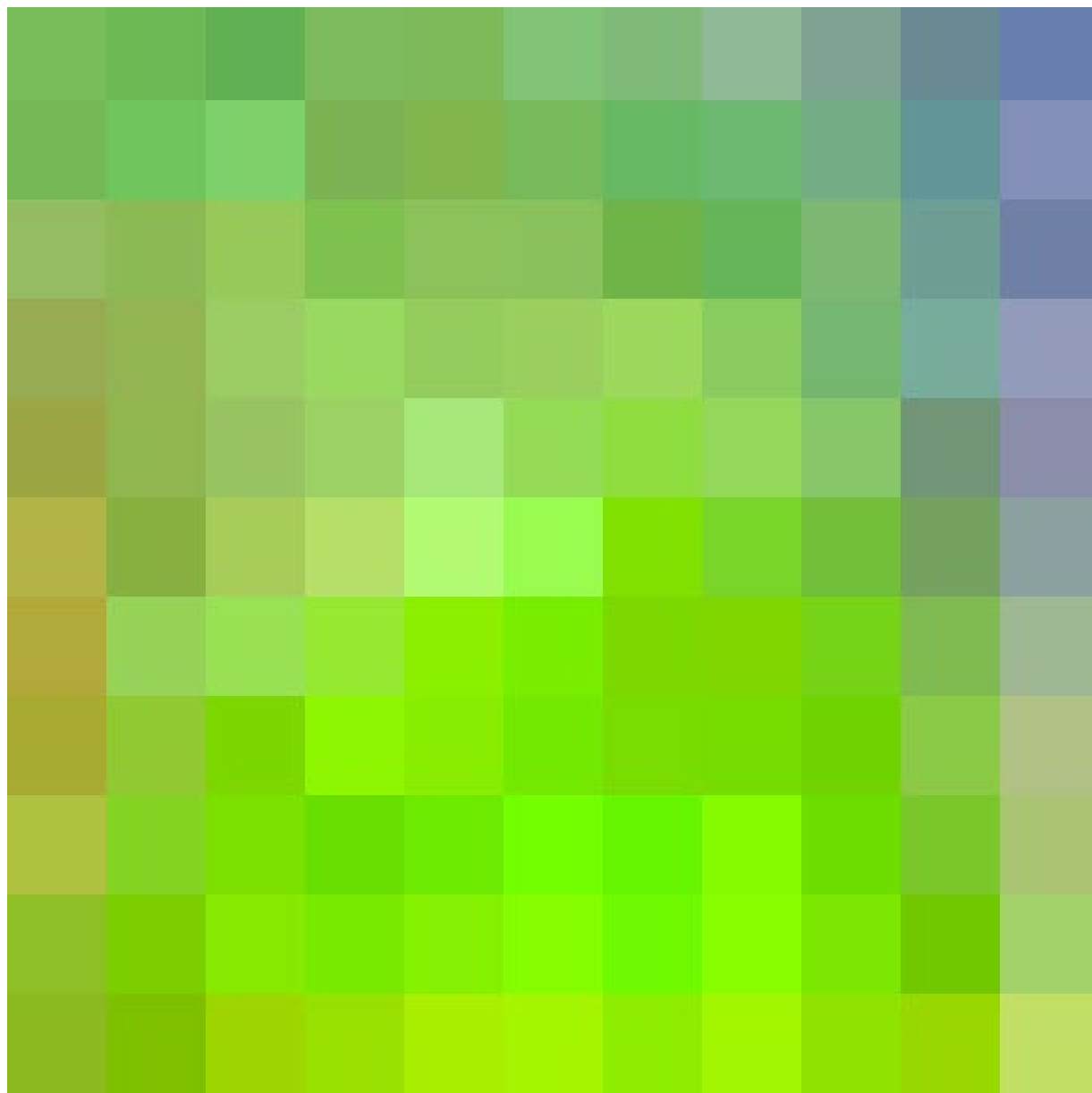
Filters in first layer



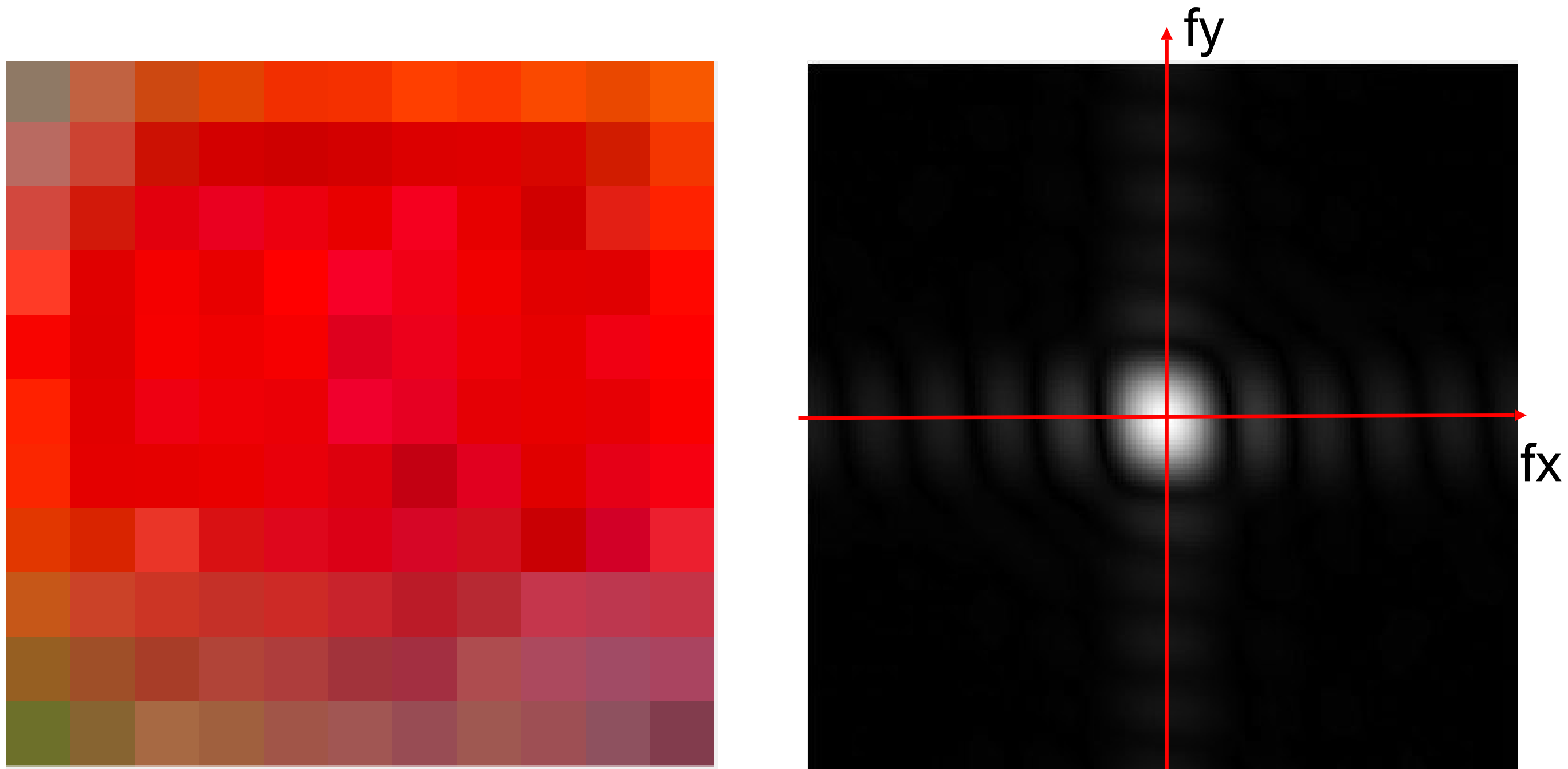
Filters in first layer



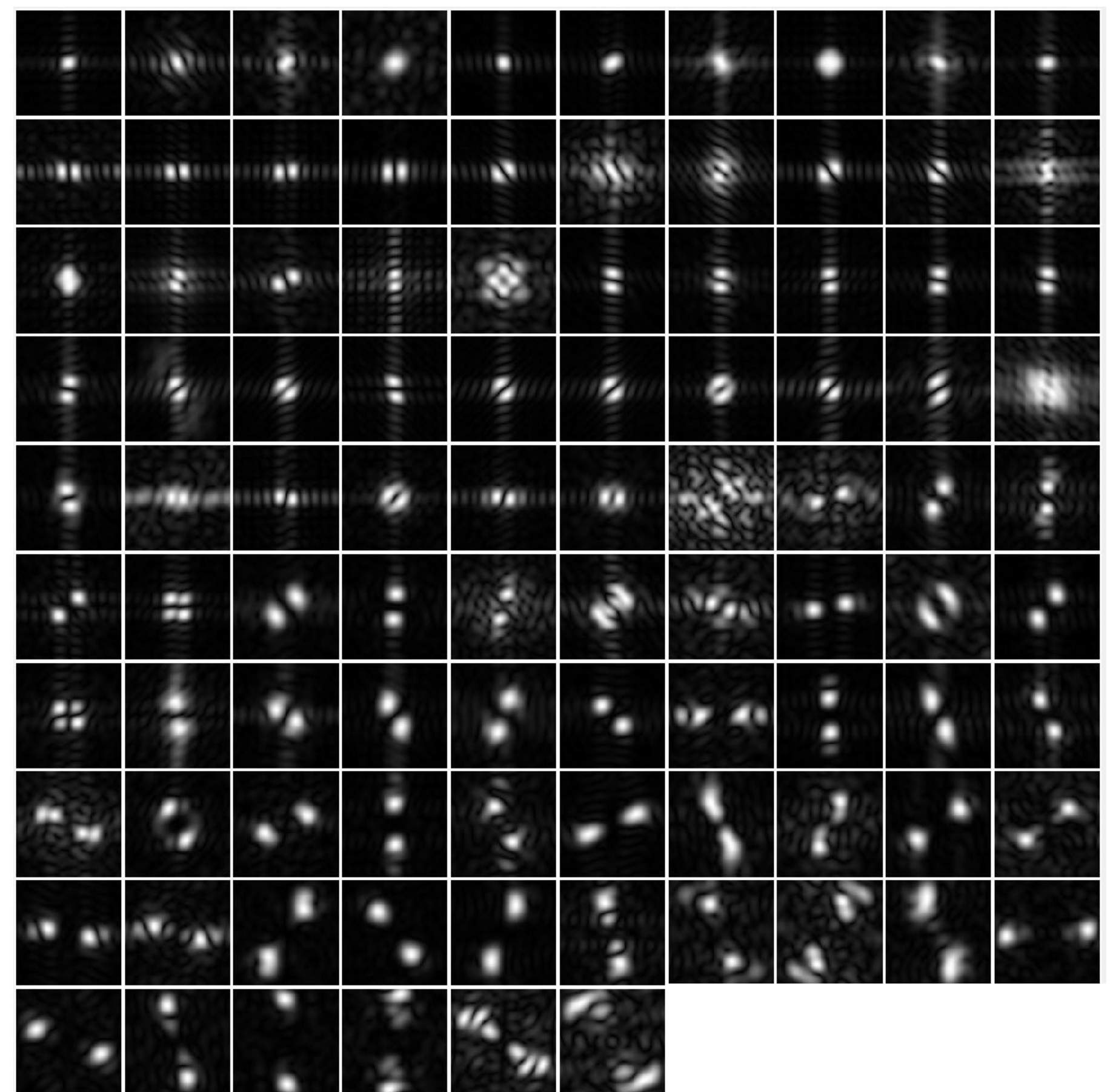
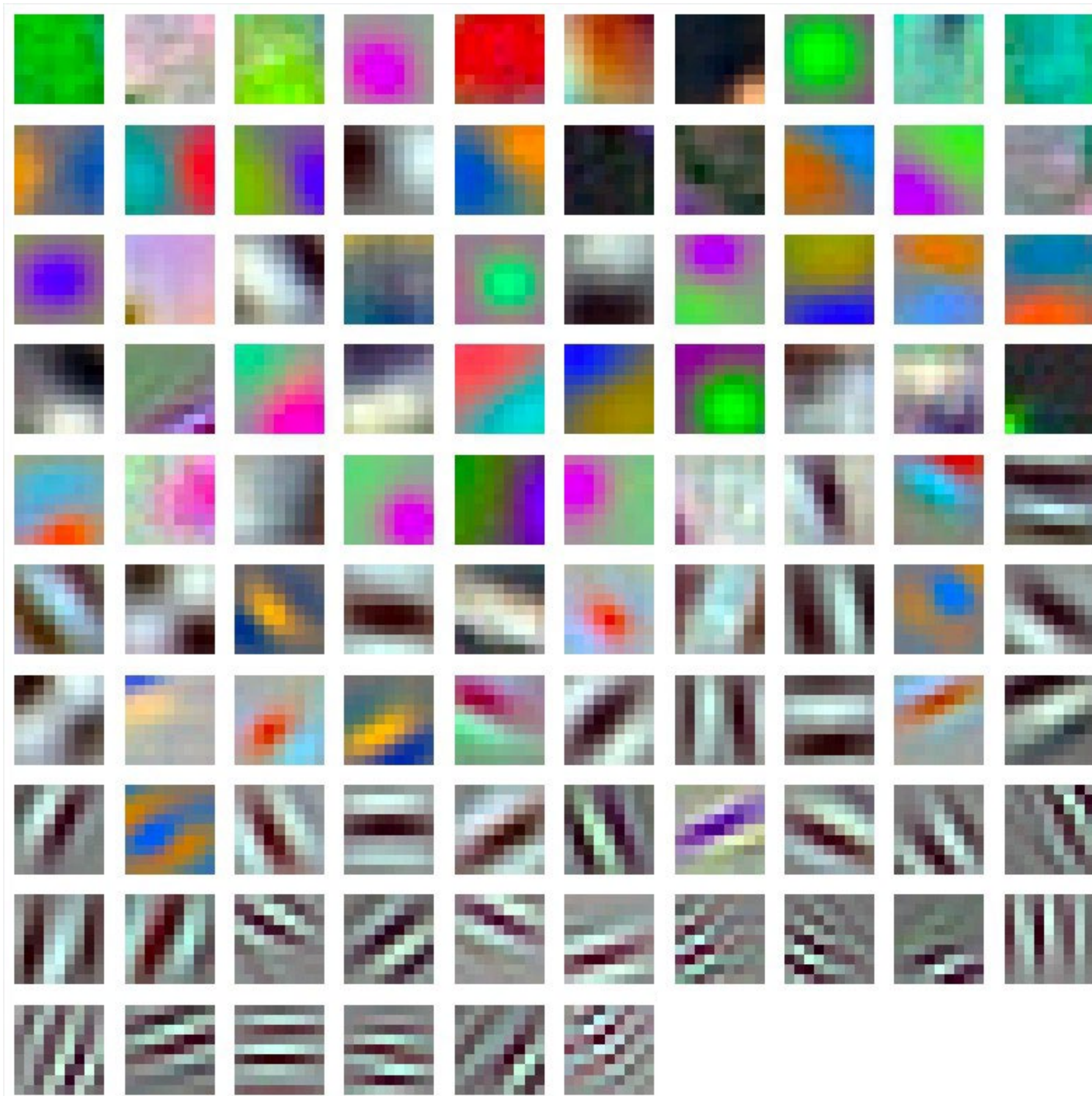
Filters in first layer



Filters in first layer



Filters in first layer



96 Units in conv1

A few practical issues

Dealing with rectangular images



Rectangular
images



Resize, then take square
crop from center

Training with data augmentation



Original image



Flipping

- Less susceptible to overfitting.
- Improves performance by simulating examples.

Training with data augmentation



Original image



Flipping



Cropping



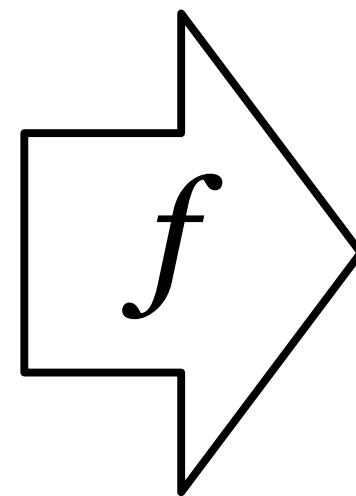
Scaling



Color jittering

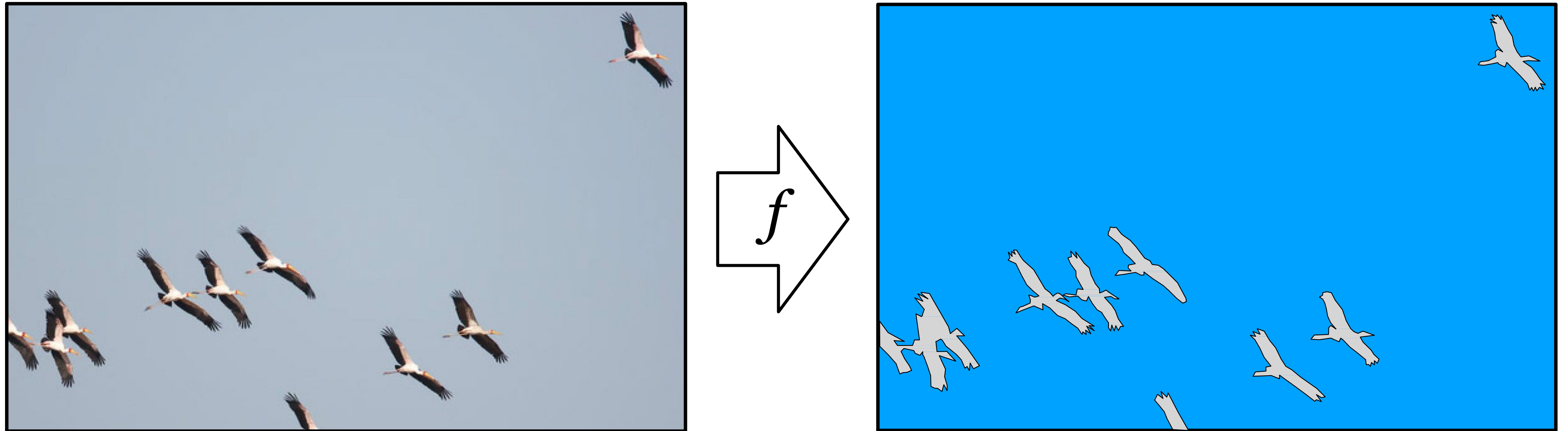
Beyond image labeling

Object recognition: what objects are in the image?



“Birds”

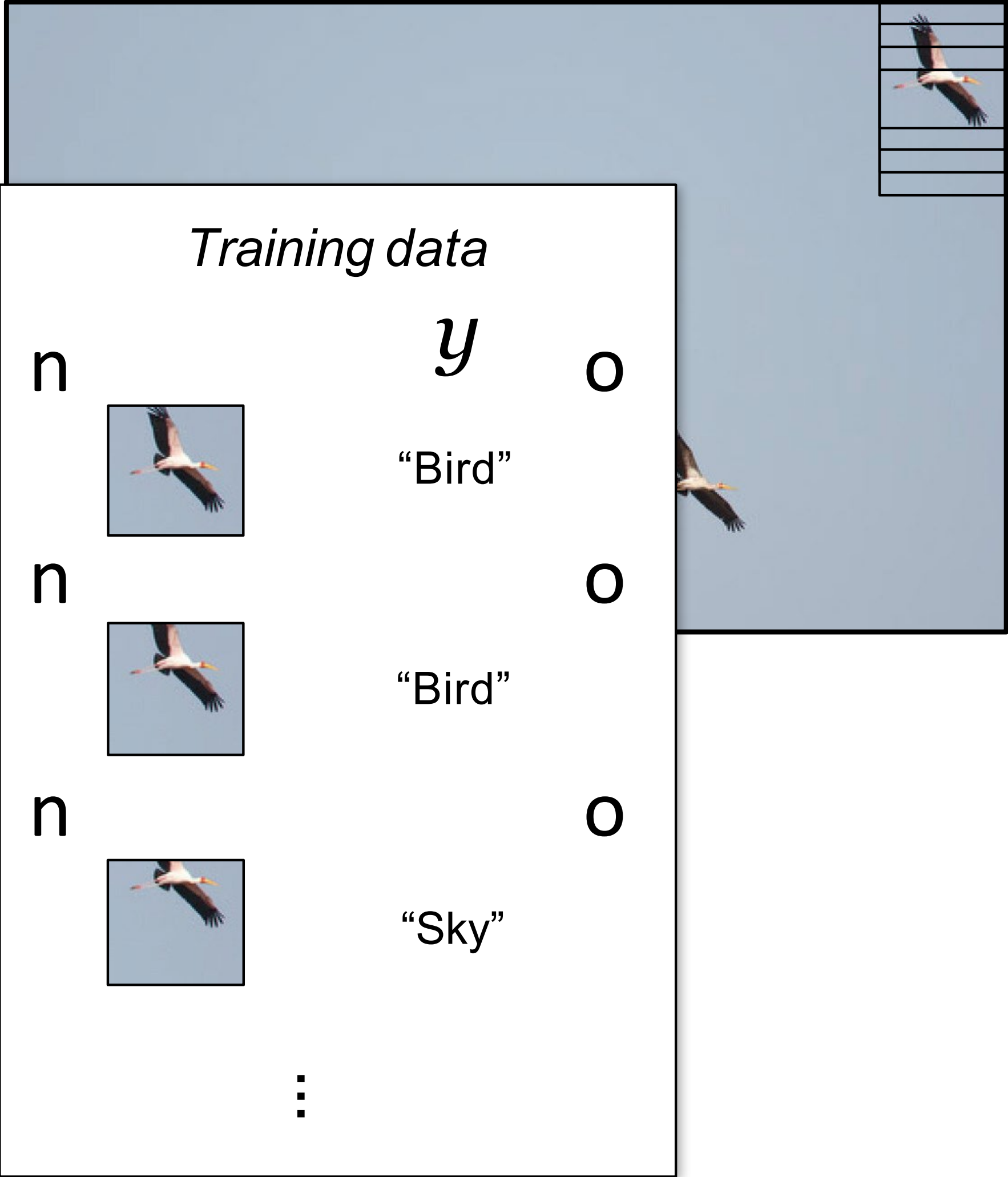
Semantic segmentation



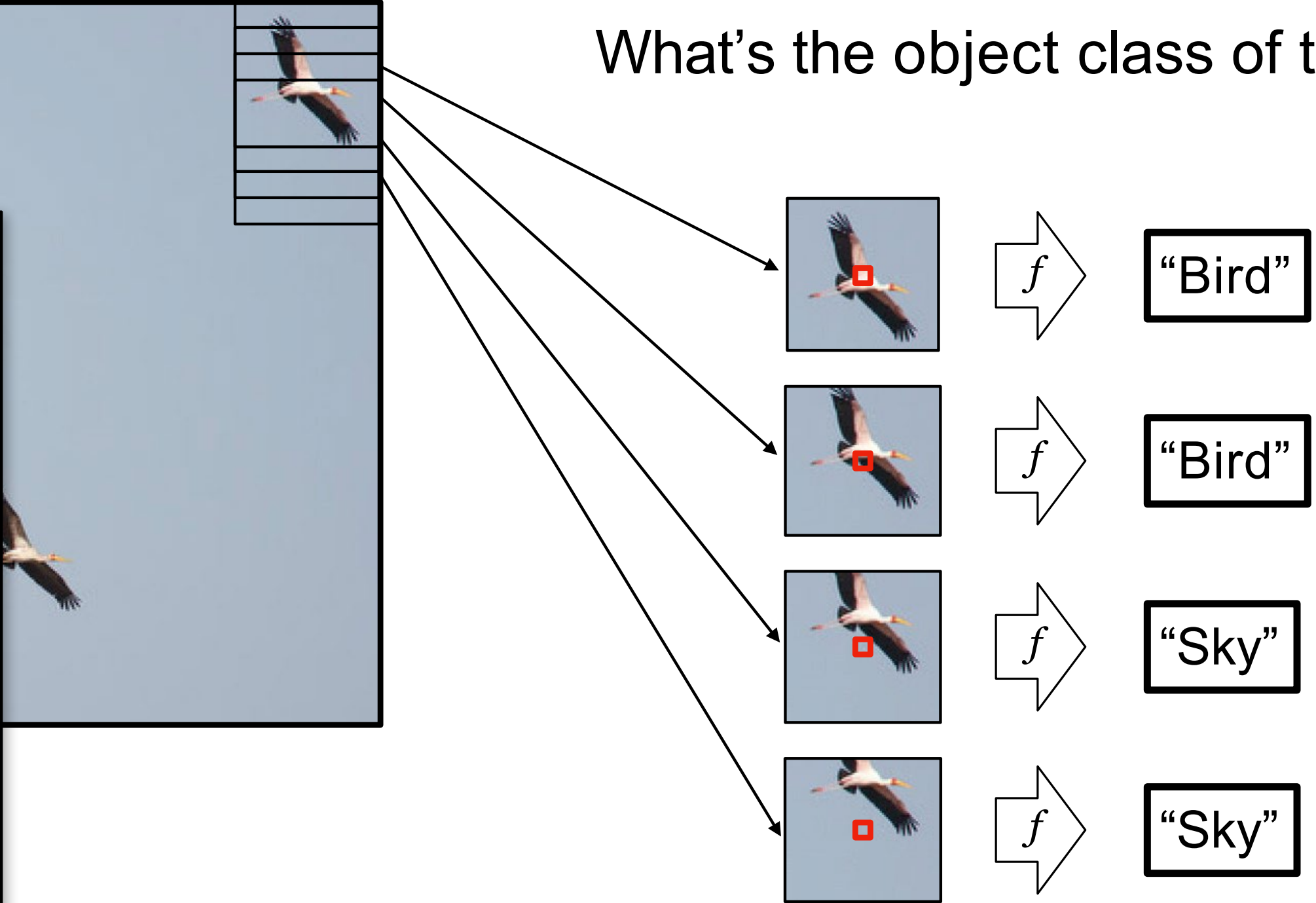
(Colors represent categories)

General technique: predict something at every pixel!

Idea #1: Independently classify windows

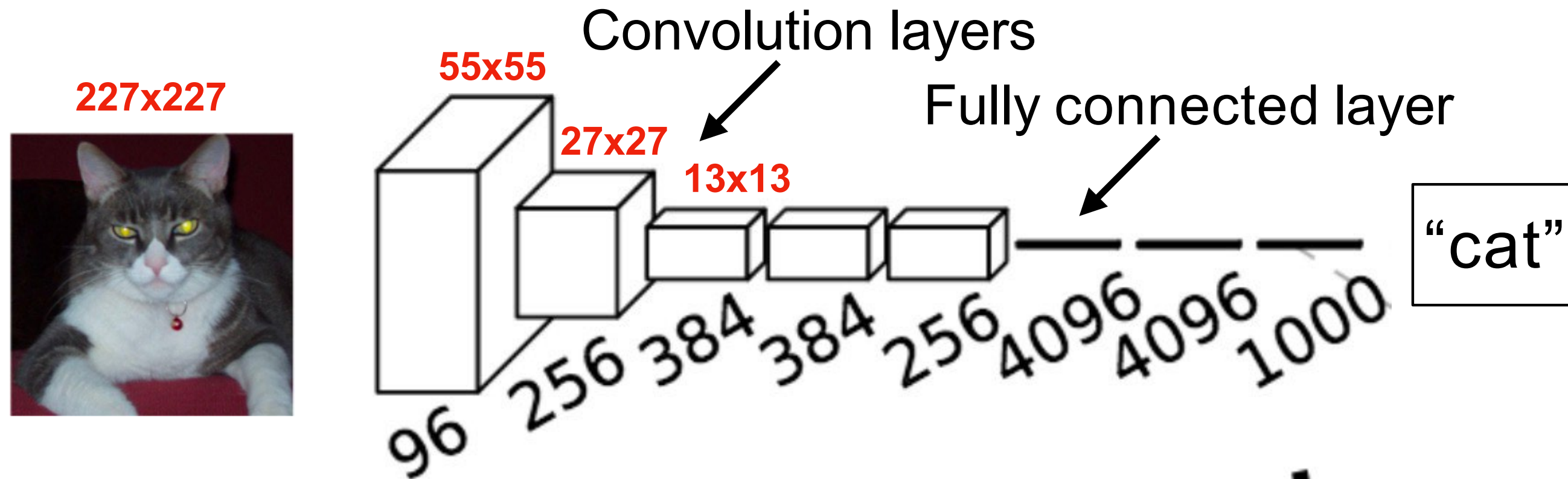


What's the object class of the center pixel?

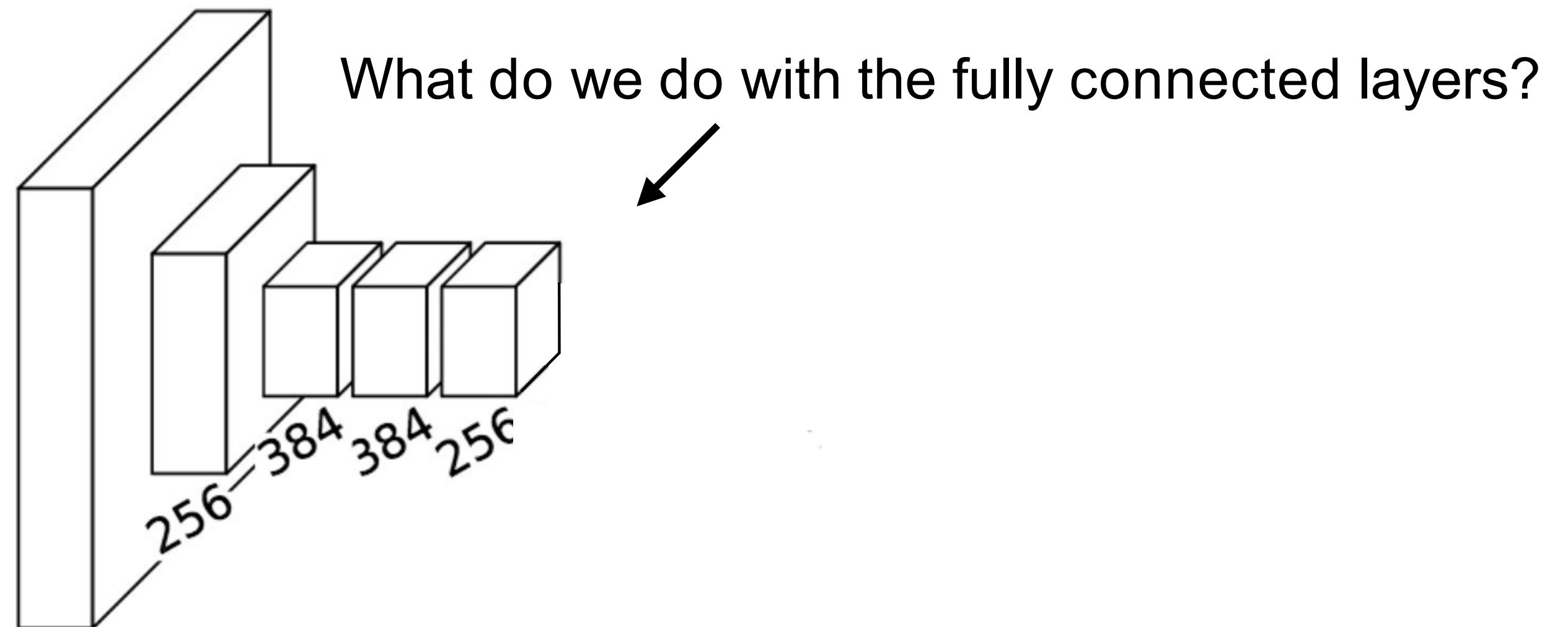


Idea #2: Fully convolutional networks

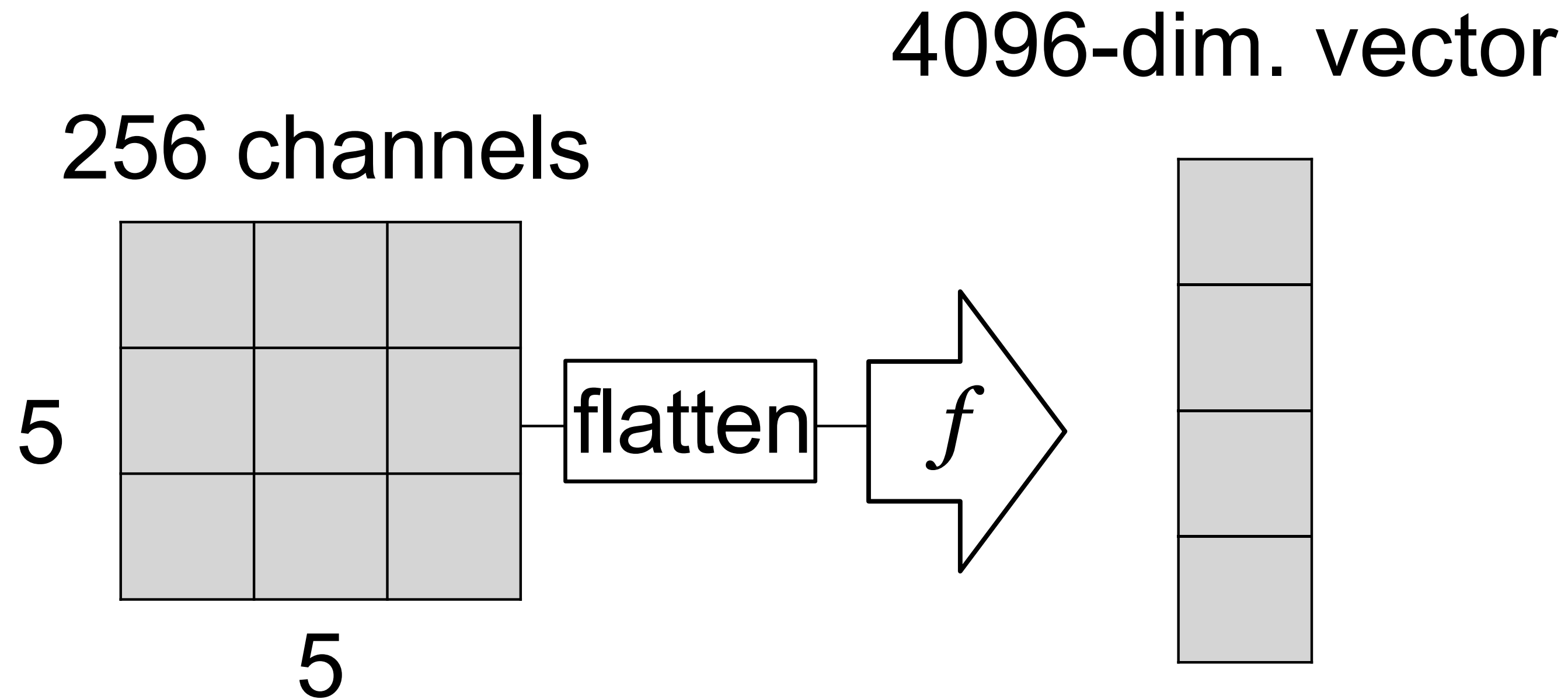
Reuse features across windows



HxW

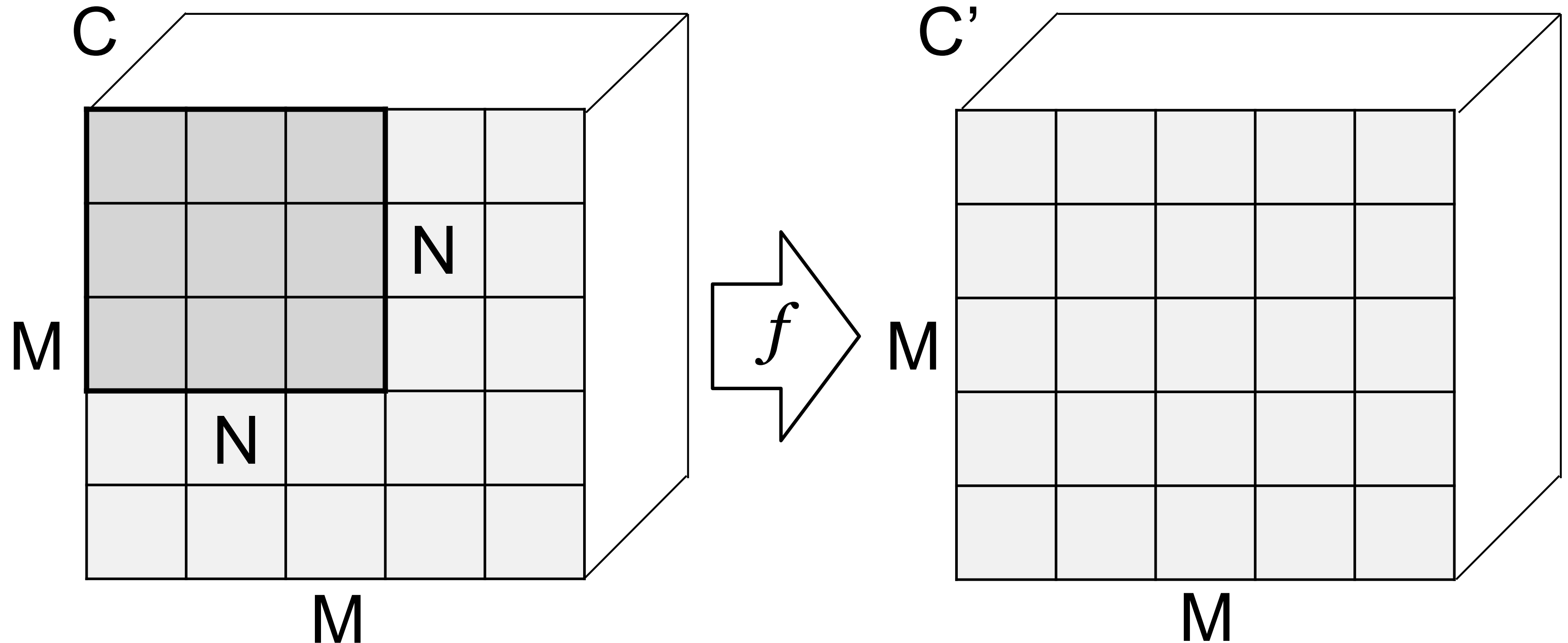


Converting fully connected layer to convolution



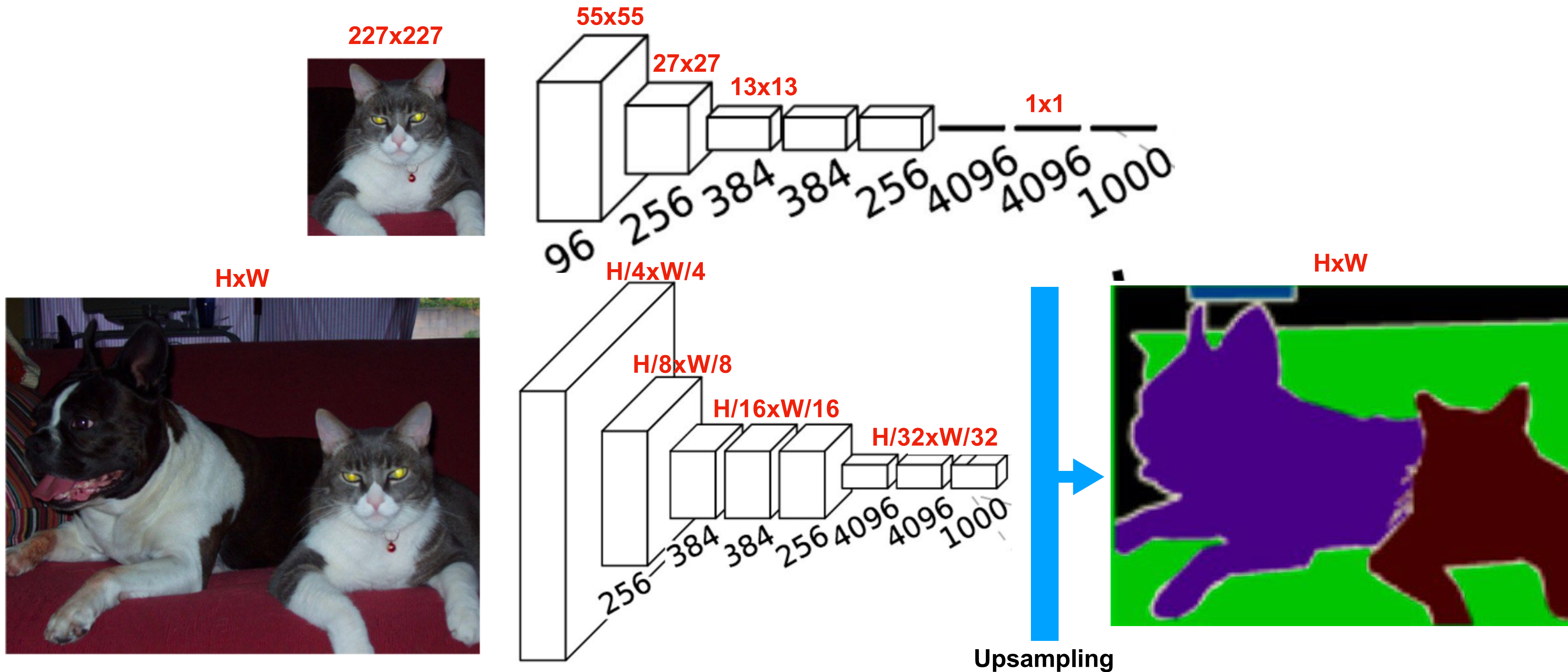
Fully connected layer: $\mathbb{R}^{N \times N \times C} \rightarrow C'$

Converting fully connected layer to convolution

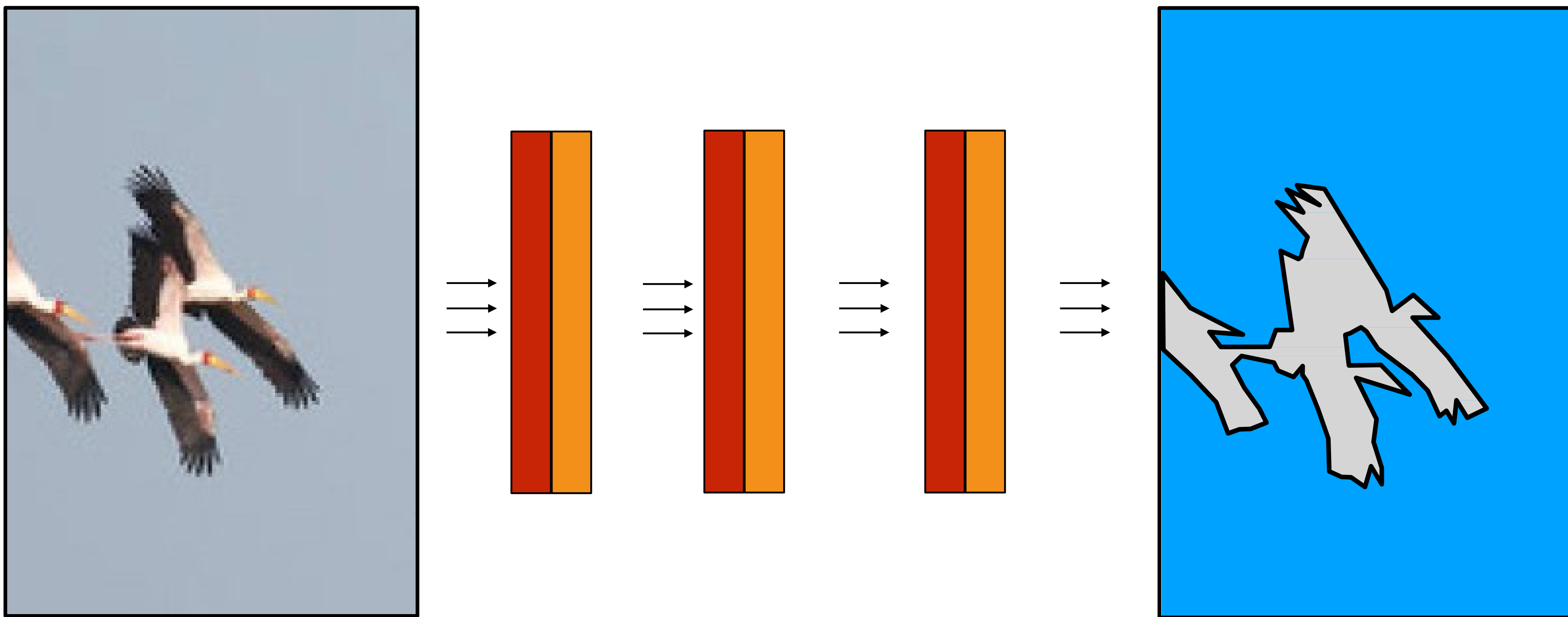


Convert into an $N \times N$ convolution!

Fully Convolutional Networks



What if we remove subsampling layers?



Problems: small receptive fields (and expensive)

Idea #3: Dilated convolutions

Dilated convolutions

3x3

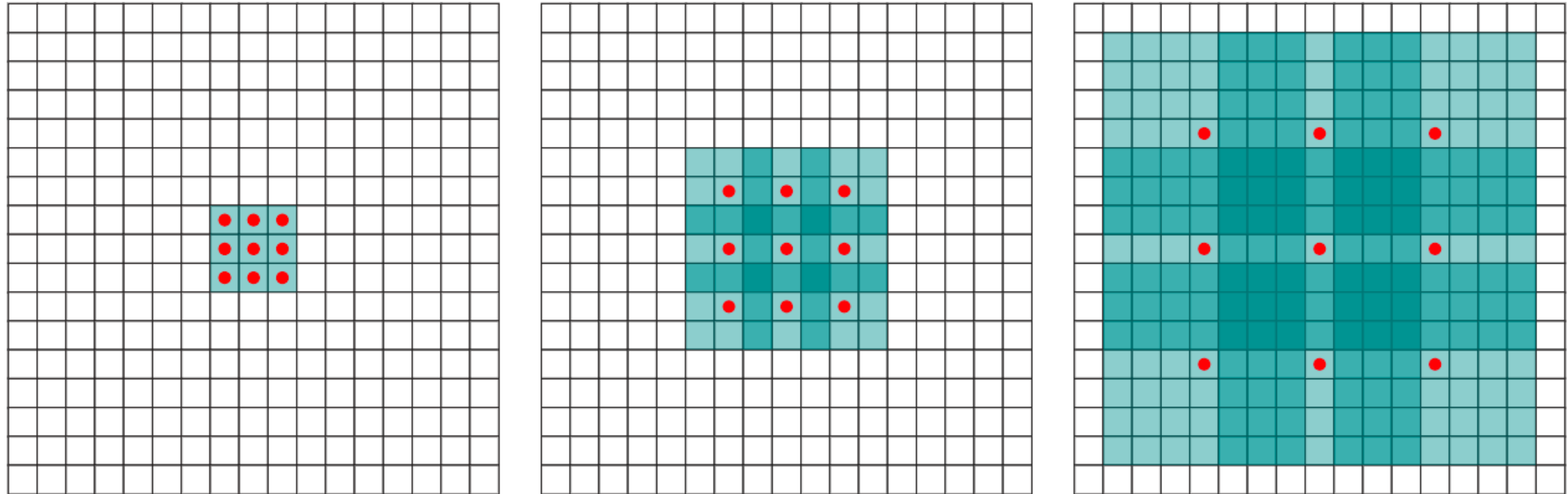
| | | |
|---|---|---|
| a | b | c |
| d | e | f |
| g | h | i |

5x5

| | | | | |
|---|---|---|---|---|
| a | 0 | b | 0 | c |
| 0 | 0 | 0 | 0 | 0 |
| d | 0 | e | 0 | f |
| 0 | 0 | 0 | 0 | 0 |
| g | 0 | h | 0 | i |

- Alternative to pooling that preserves input size
- 9 degrees of freedom
- 5x5 receptive field

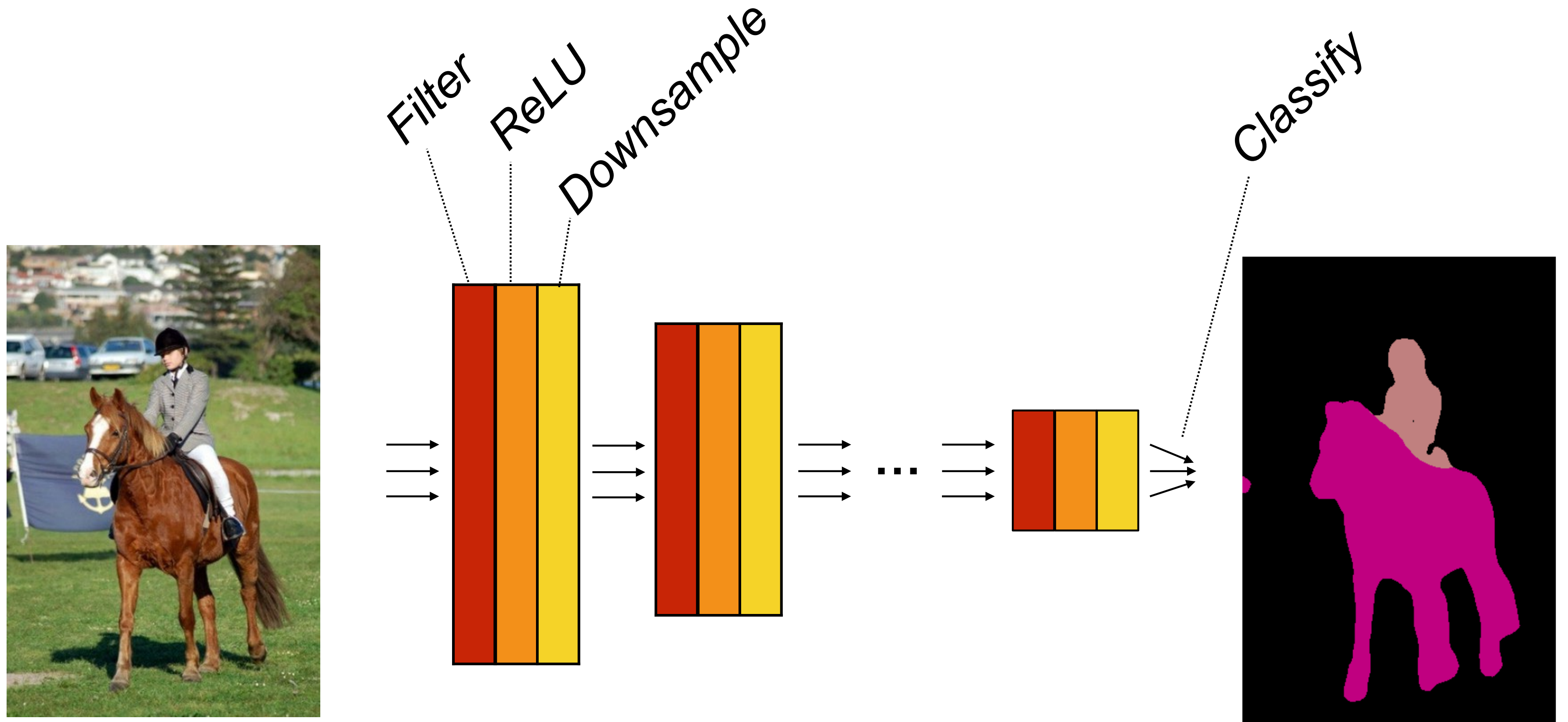
Architectures with dilated convolutions



- Architecture design: dilation by 2^L instead of striding L times
- Obtains comparable receptive field to CNN with strides.

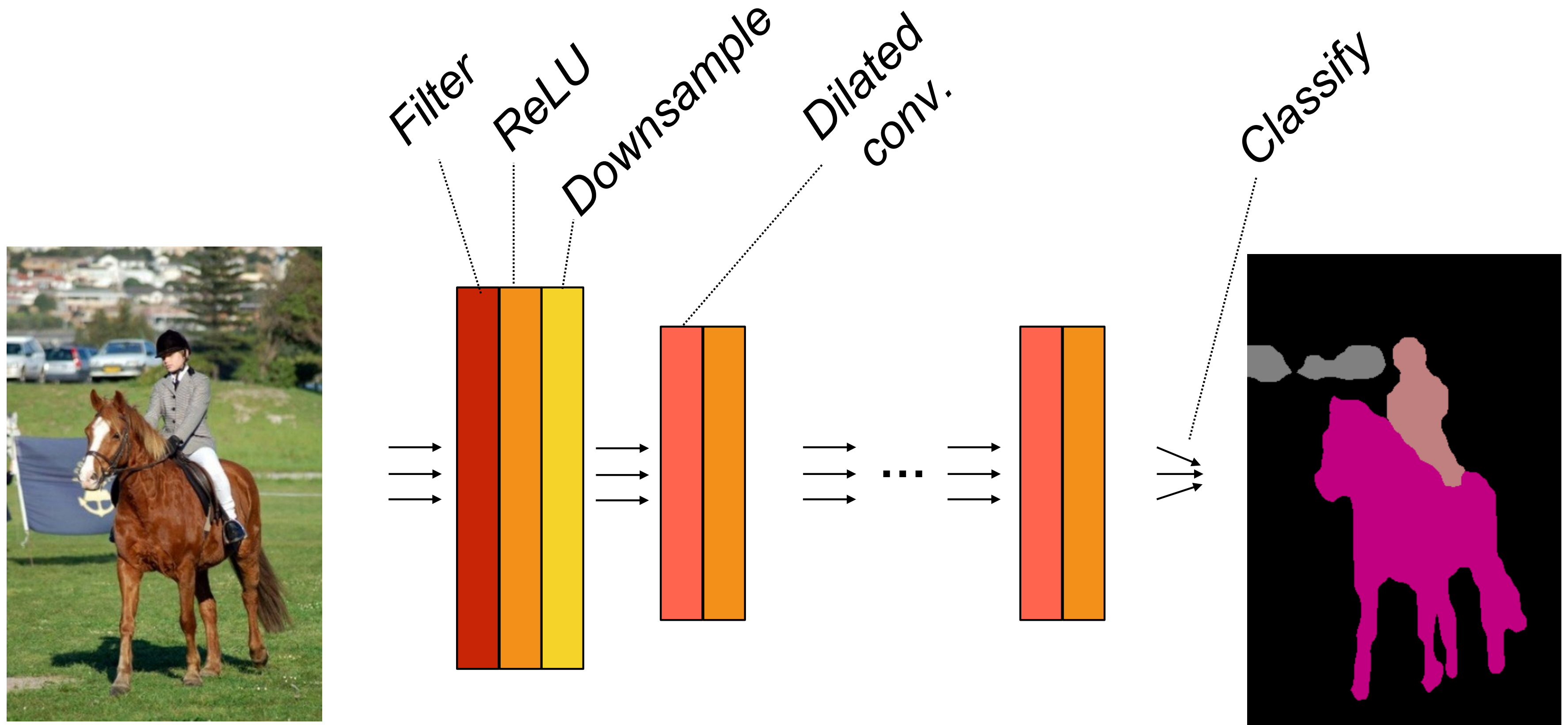
[Yu and Koltun 2016, <https://arxiv.org/pdf/1511.07122.pdf>]

CNN without dilated convolutions



Apply CNN convolutionally

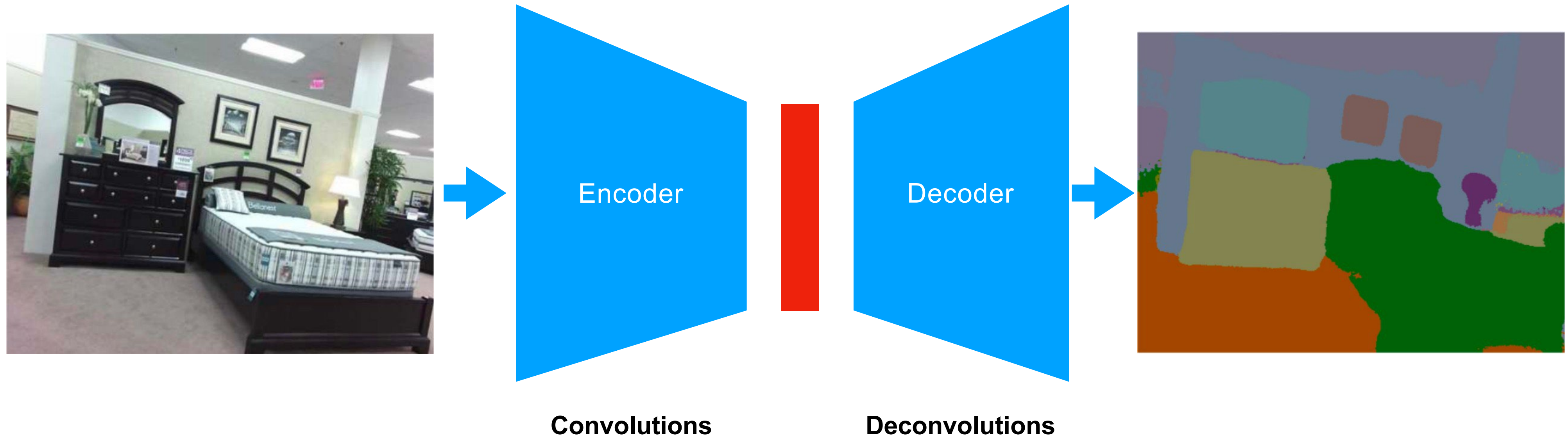
CNN with dilated convolutions



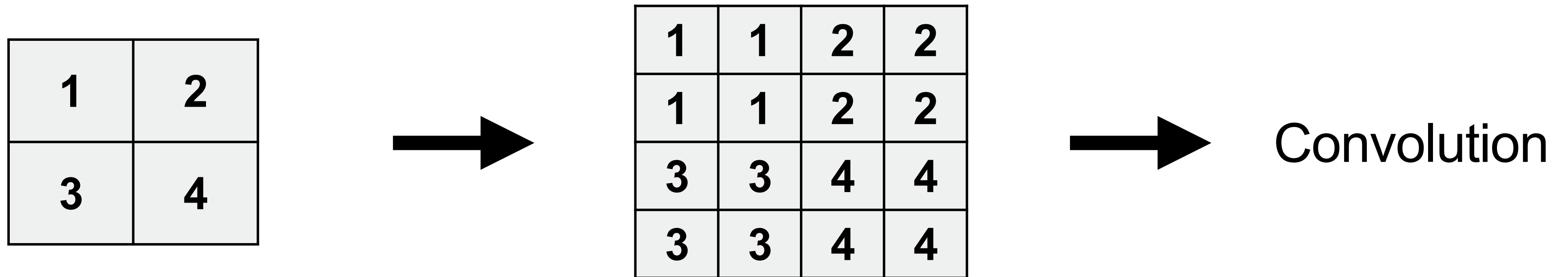
Even with dilated convolutions, still not full resolution

Idea #4: Encoder-decoder models

Encoder-decoder architectures

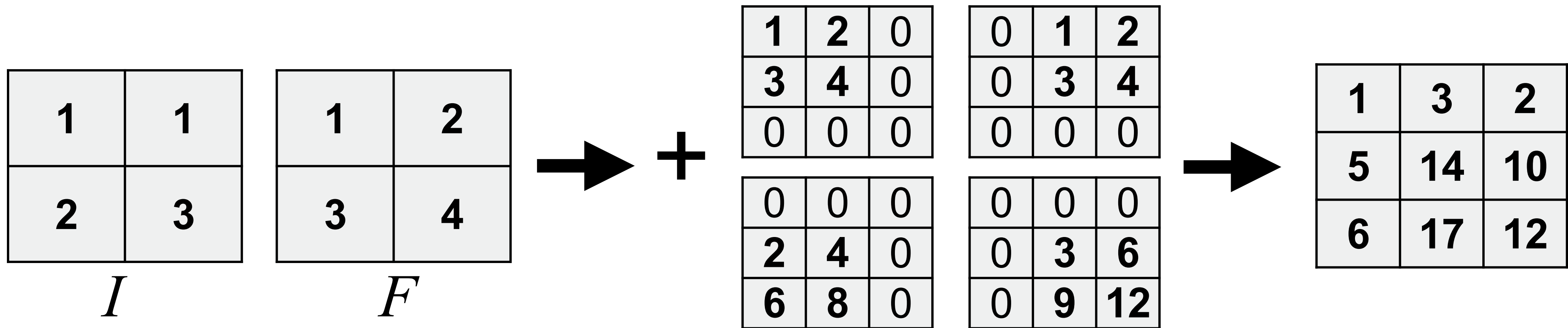


Upsampling



- Often using nearest-neighbor upsampling
- Can also use interpolation.
- Produces fewer “checkerboard” artifacts

Transposed convolution



- Weight the filter by the image coefficient and sum.
- Also sometimes called “upconvolution” or “deconvolution”.

Transposed convolution

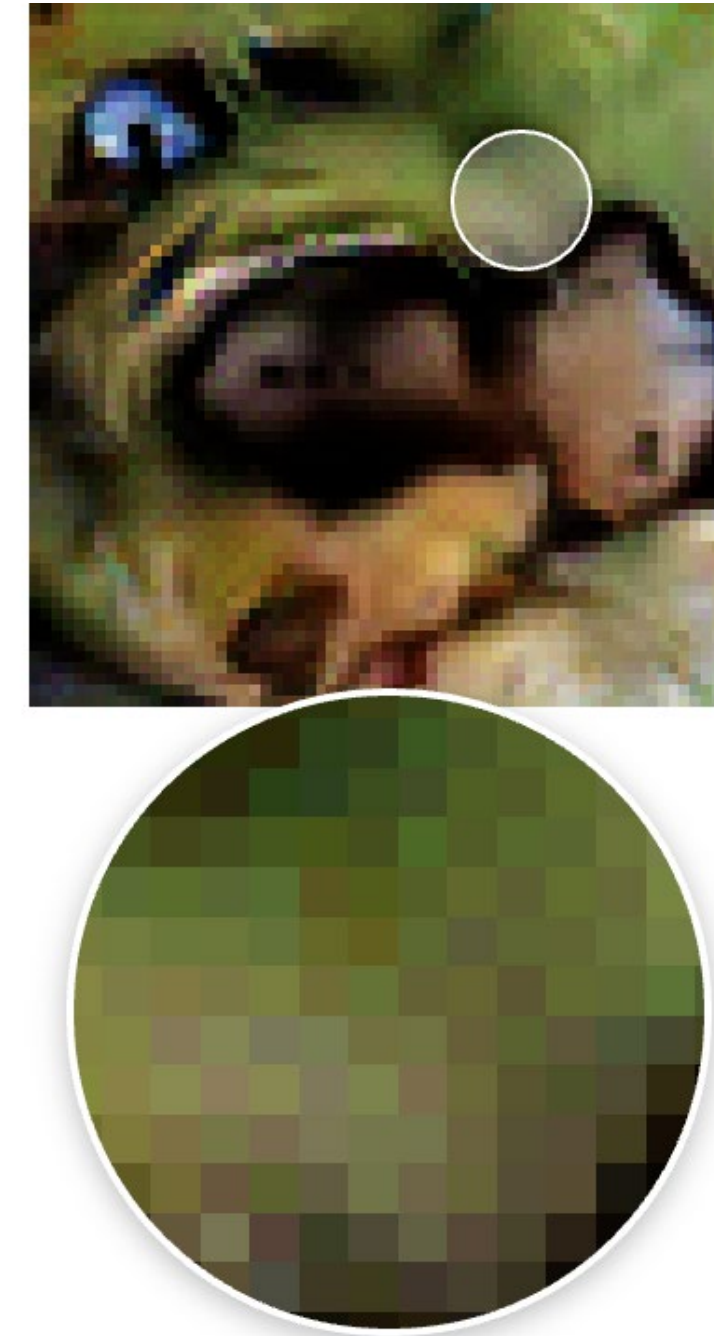
| | |
|---|---|
| 1 | 1 |
| 2 | 3 |

I

| | |
|---|---|
| 1 | 2 |
| 3 | 4 |

F

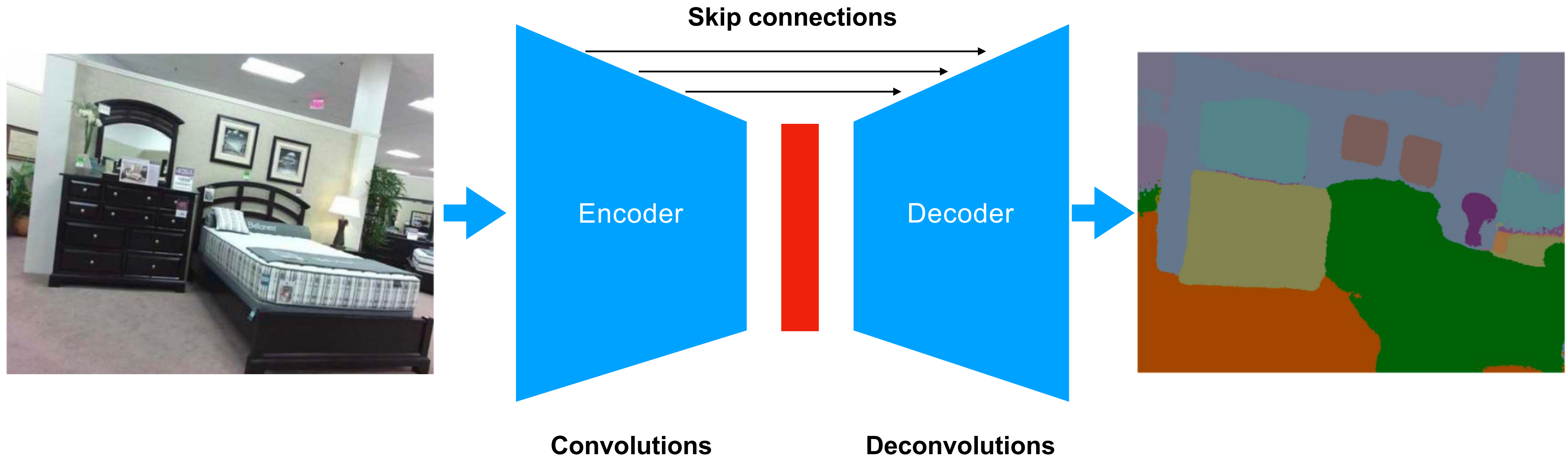
Can lead to “checkerboard” artifacts.



Donahue, et al., 2016 [3]

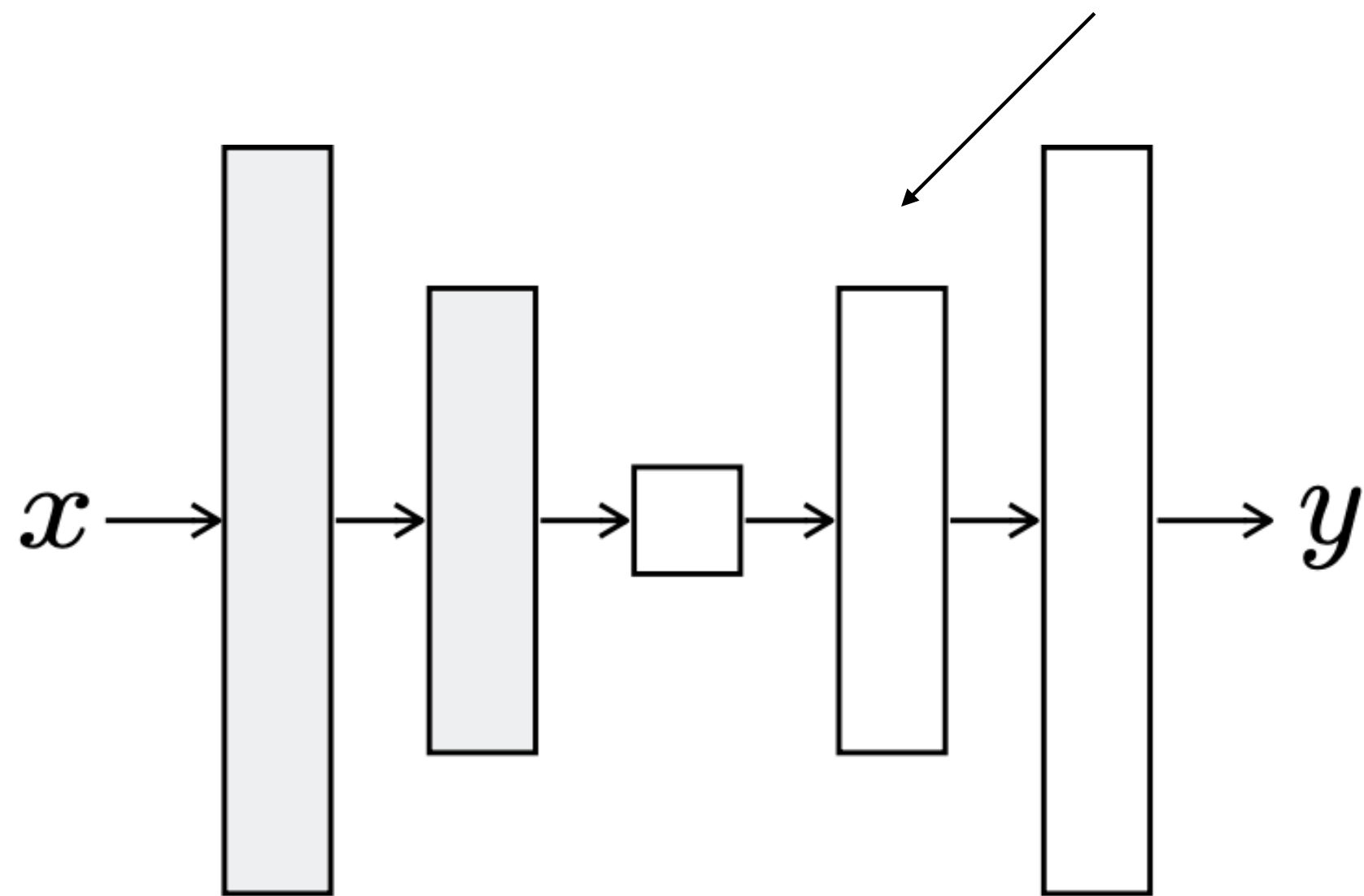
[Odena et al. Distill article]

Encoder-decoder architectures



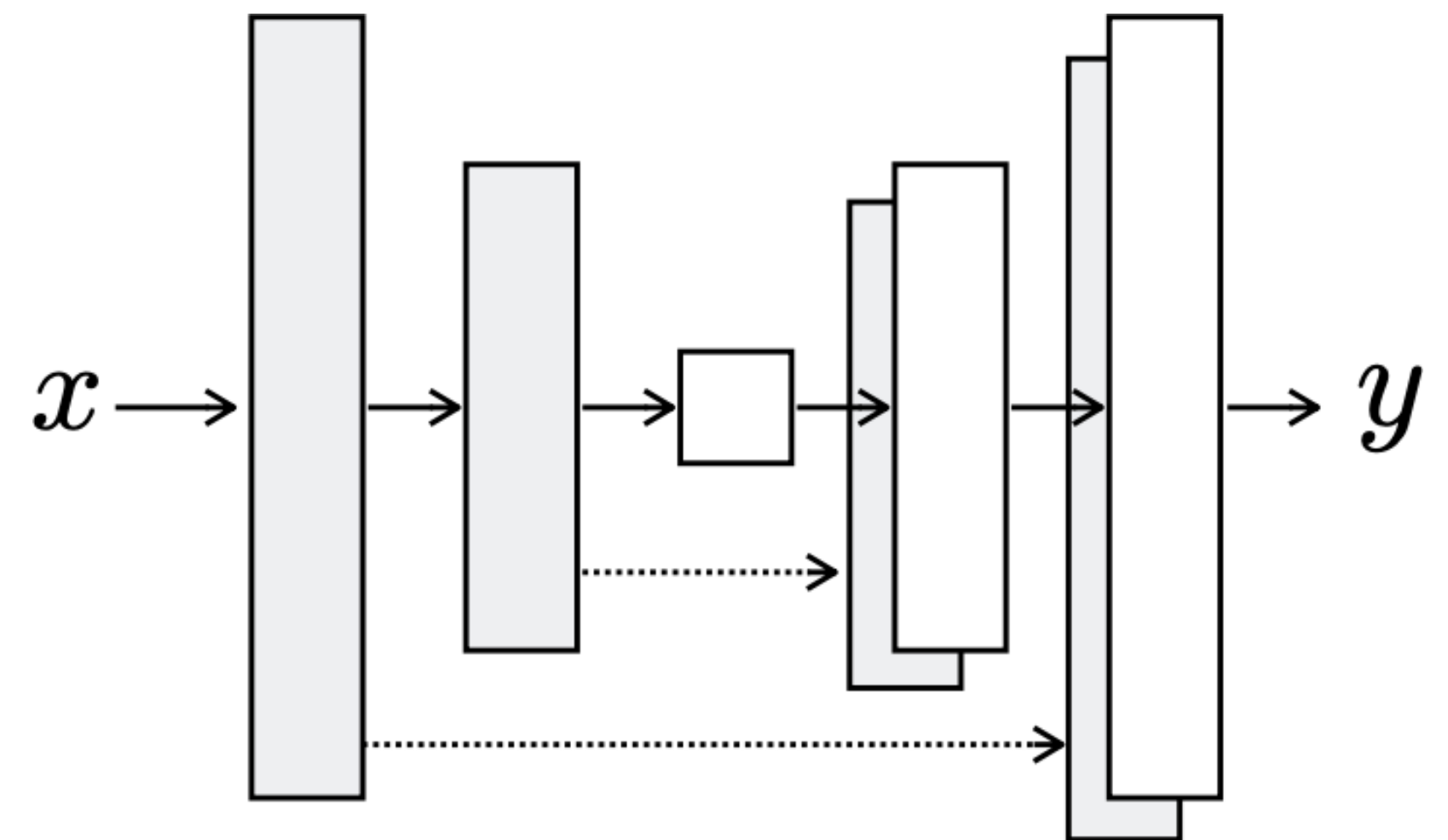
Encoder-decoder architectures

Transposed convolution



“Vanilla” encoder-decoder architecture

Early layers and late layers have same shape. Concatenate channel-wise!



U-Net

Encoder-decoder architectures

SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation

V. Badrinarayanan, A. Kendall, R. Cipolla, *Senior Member, IEEE*,

This paper presents a deep convolutional encoder-decoder architecture for semantic pixel-wise segmentation. The encoder network, a corresponding decoder network followed by a skip connection, maps the input image to a low resolution feature map. The decoder network upsamples its lower resolution input feature map to the original resolution of the input image. The upsampling is performed by a max-pooling step of the corresponding encoder to produce sparse maps. The upsampling maps are sparse and are then combined with the proposed architecture with the widely adopted FCN [2]. This comparison reveals the memory versus

performance of SegNet. SegNet is designed to be efficient both in terms of memory and number of trainable parameters than other competing architectures. We also performed a controlled benchmark of SegNet on segmentation tasks. These quantitative assessments and most efficient inference memory-wise as compared to other architectures are available at <http://mi.eng.cam.ac.uk/projects/segnet/>.

Keywords: Segmentation, Indoor Scenes, Road Scenes, Encoder,

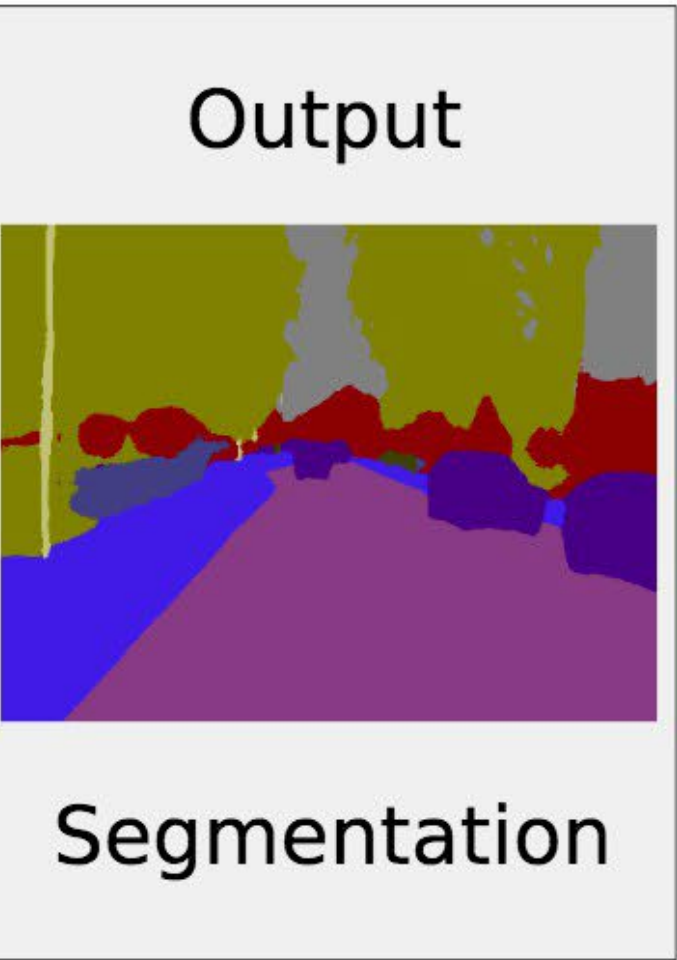
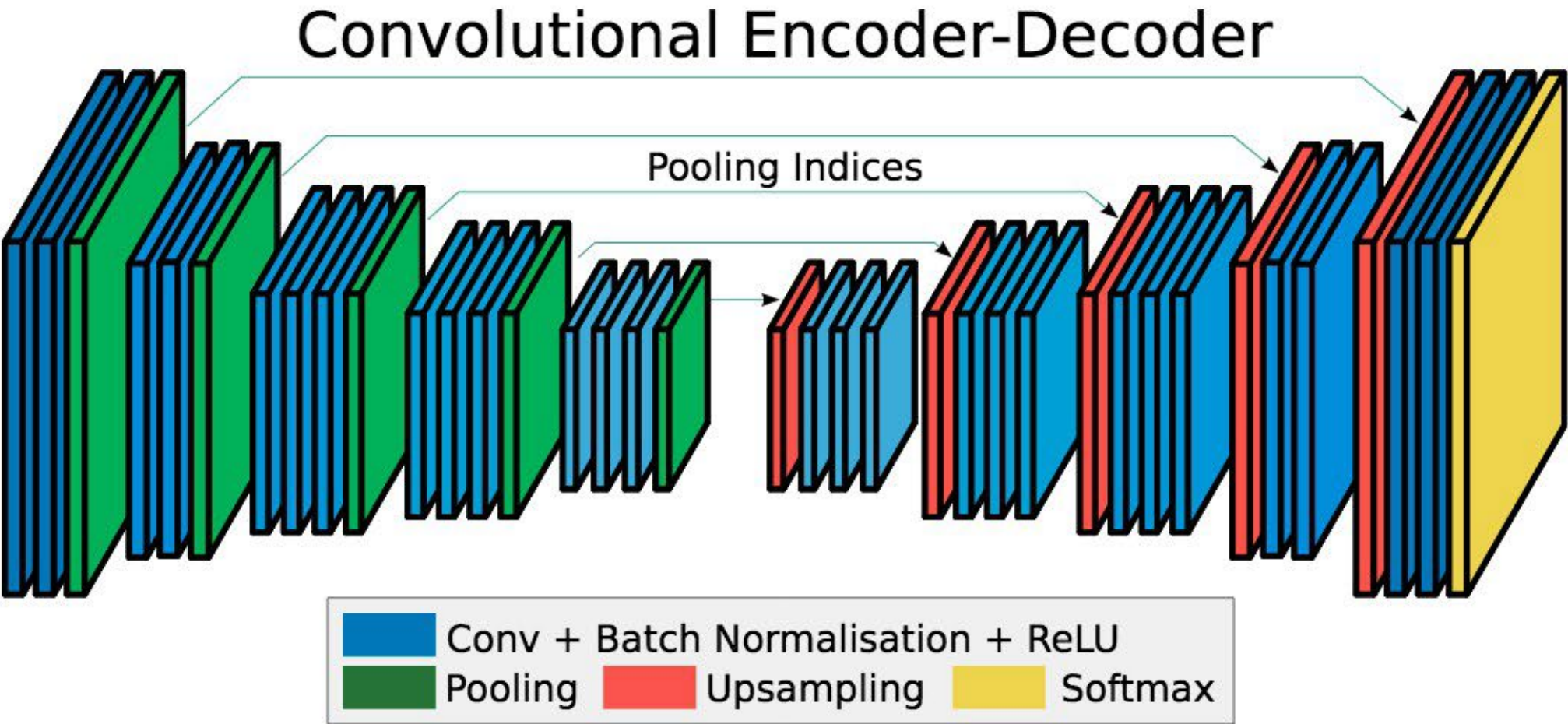
and understand the spatial-relationship (context) between different classes such as road and side-walk. In typical road scenes, the majority of the pixels belong to large classes such as building and hence the network must produce smooth boundaries. The engine must also have the ability to delineate boundaries based on their shape despite their small size. Hence it is necessary to retain boundary information in the extracted image features. From a computational perspective, it is necessary for the network to be efficient in terms of both memory and inference time during inference. The ability to train end-to-end to jointly optimise all the weights in the network using a single weight update technique such as stochastic gradient descent (SGD) [17] is an additional benefit since it is more easily repeatable. The design of SegNet arose from a need to match these criteria.

The encoder network in SegNet is topologically identical to the convolutional layers in VGG16 [11]. We remove the fully connected layers of VGG16 which makes the SegNet encoder network significantly smaller and easier to train than many other recent architectures [2], [4], [11], [18]. The key component of SegNet is the decoder network which consists of a hierarchy of decoders one corresponding to each encoder. Of these, the appropriate decoders use the max-pooling indices received from the corresponding encoder to perform non-linear upsampling of their input feature maps. This idea was inspired from an architecture designed for unsupervised feature learning [19]. Reusing max-pooling indices in the decoding process has several practical

This is primarily because max pooling and sub-sampling reduce feature map resolution. Our motivation to design SegNet arises from this need to map low resolution features to input resolution for pixel-wise classification. This mapping must produce features which are useful for accurate boundary localization.

Our architecture, SegNet, is designed to be an efficient architecture for pixel-wise semantic segmentation. It is primarily motivated by road scene understanding applications which require the ability to model appearance (road, building), shape (cars,

• V. Badrinarayanan, A. Kendall, R. Cipolla are with the Machine Intelligence Lab, Department of Engineering, University of Cambridge, UK. E-mail: vb292,agk34,cipolla@eng.cam.ac.uk



Input



Segnet



FCN

