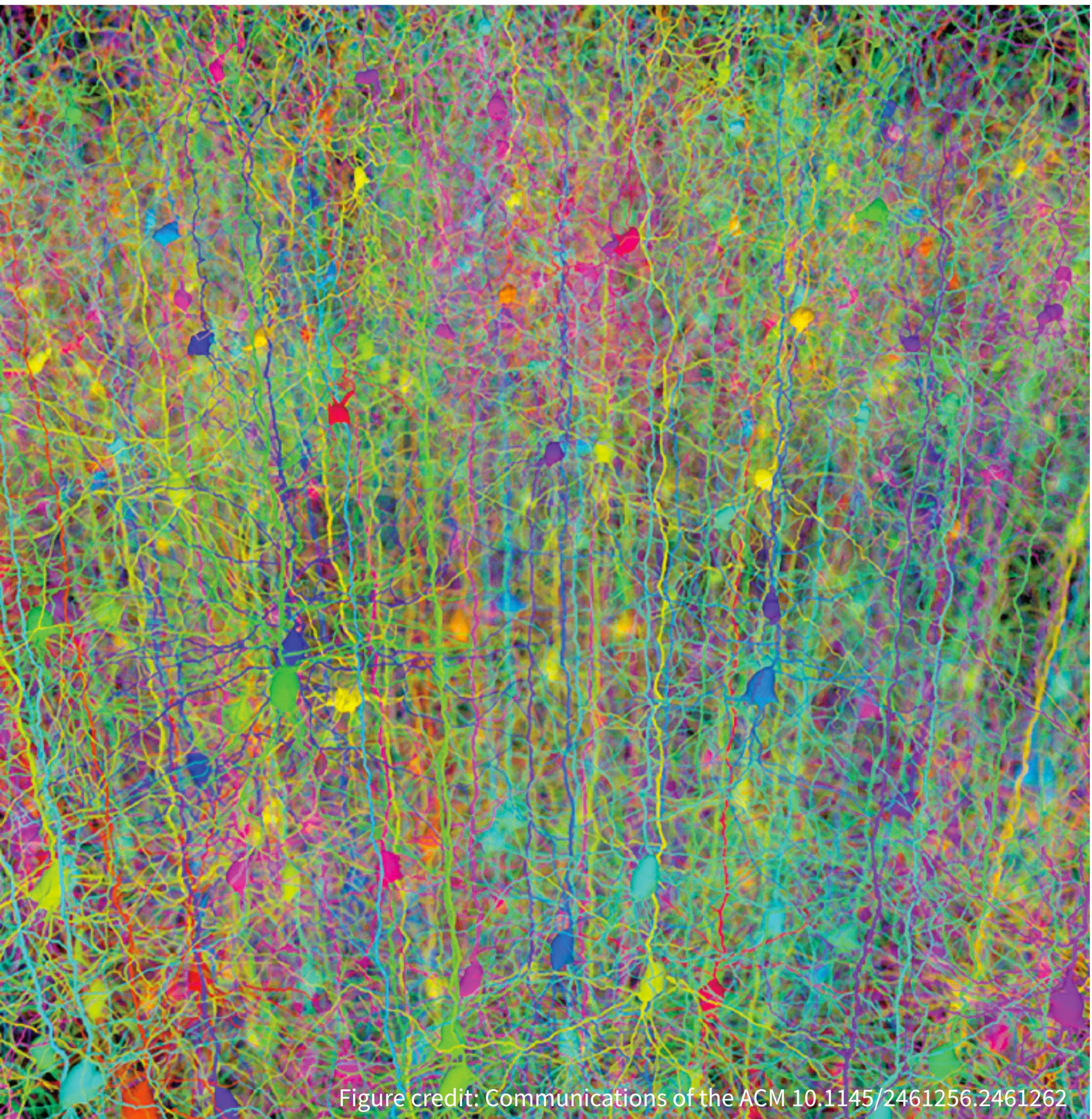


# Convolutional Neural Networks

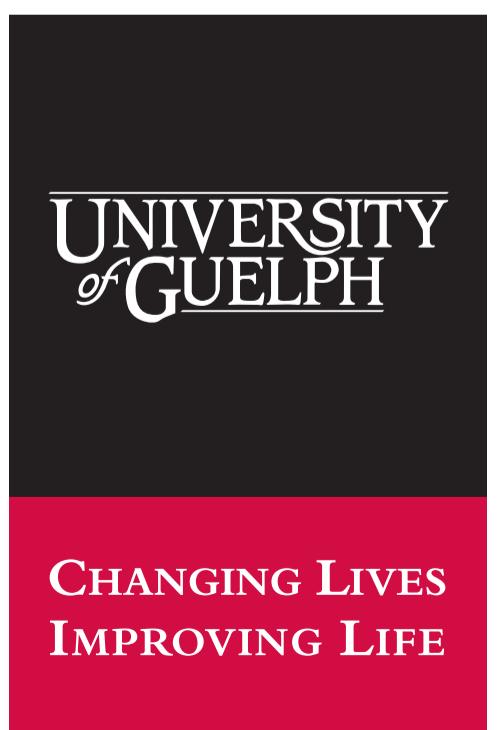


GRAHAM TAYLOR

VECTOR INSTITUTE

SCHOOL OF ENGINEERING  
UNIVERSITY OF GUELPH

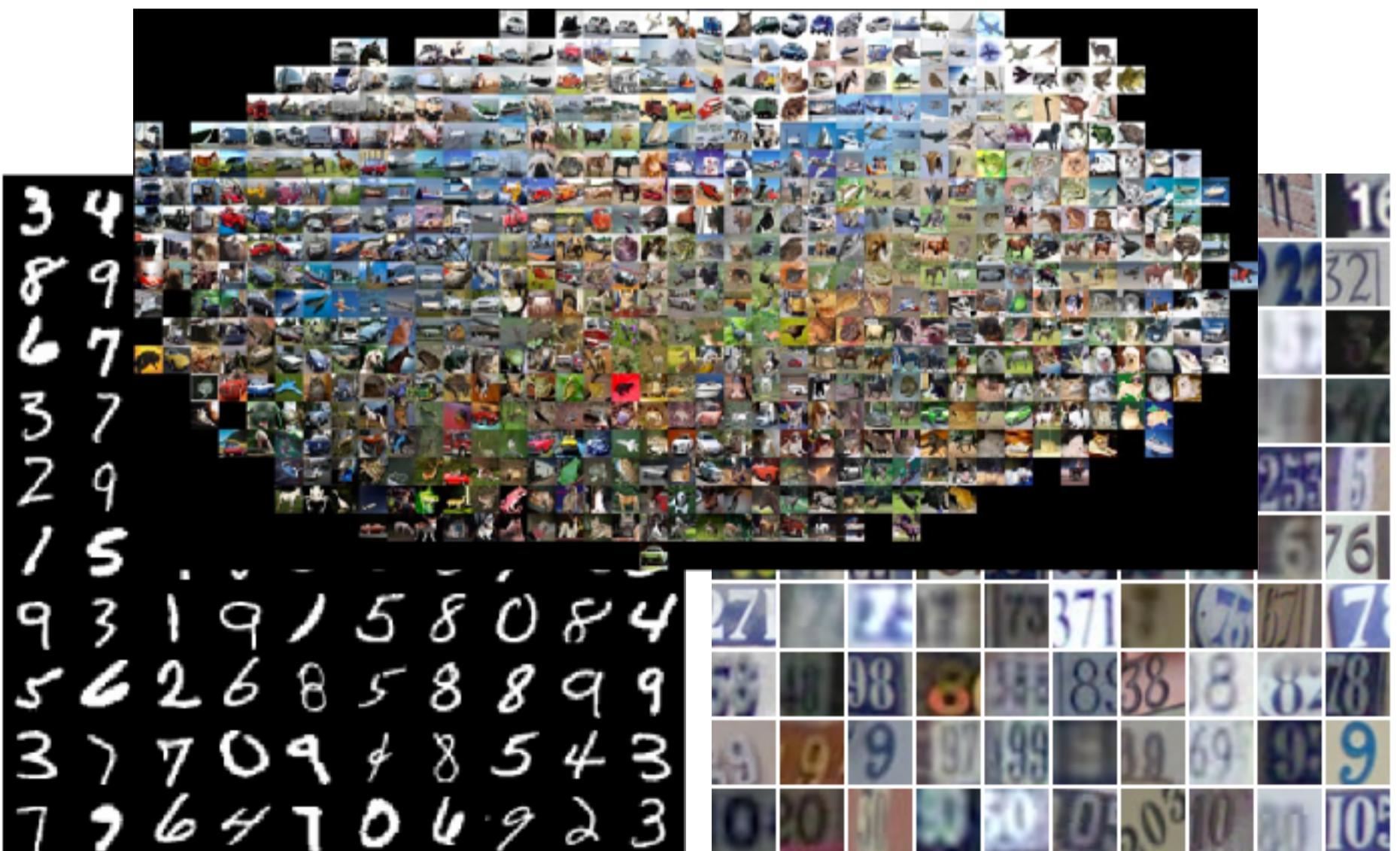
CANADIAN INSTITUTE  
FOR ADVANCED RESEARCH



**CIFAR**  
CANADIAN  
INSTITUTE  
FOR  
ADVANCED  
RESEARCH

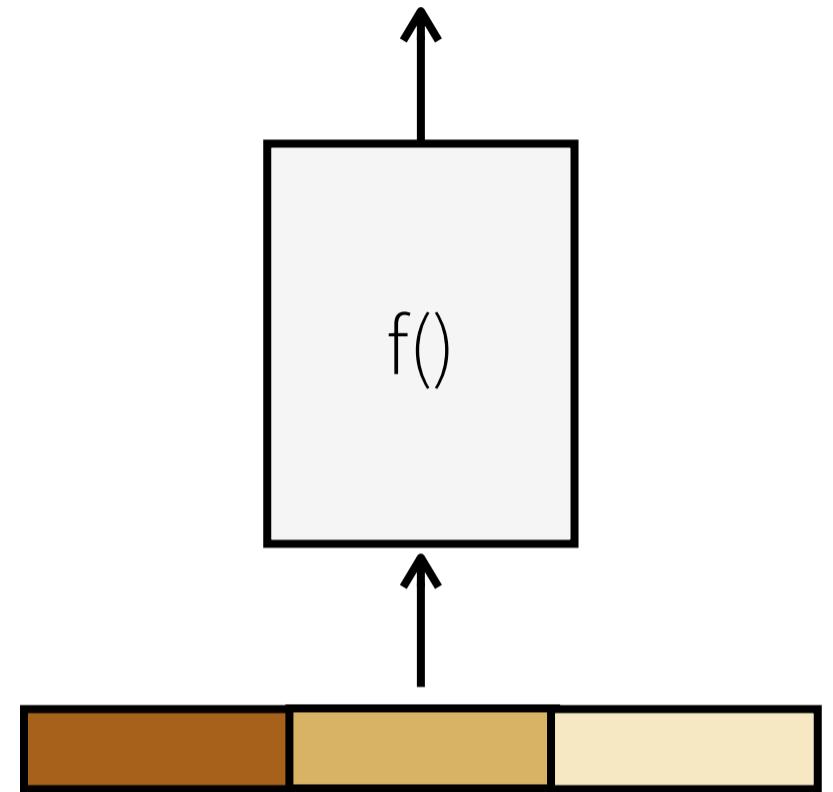
# DNNs for Visual Recognition

- For a long time, **images** have been a common “testbed” for deep learning
- Computer vision and machine learning have been intimately linked, moreso in recent years
- The application of feedforward neural networks to image data has typically been **patch-based**, or on very tiny images, e.g.
  - MNIST ( $28 \times 28$ )
  - CIFAR-10 ( $32 \times 32$ )
  - SVHN ( $32 \times 32$ )
  - etc...



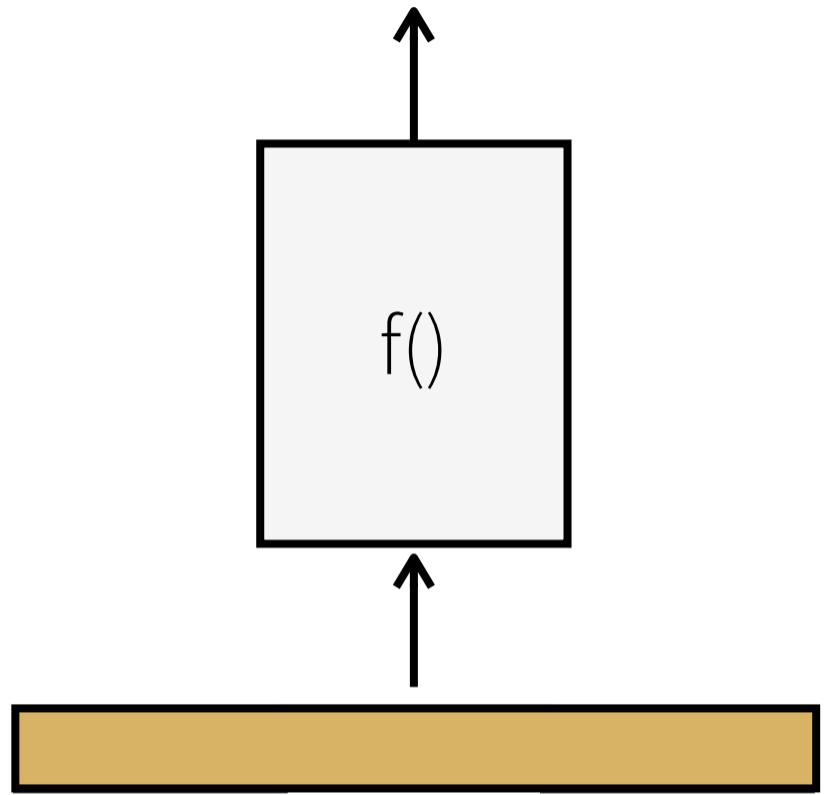
# Encoding Domain Knowledge

- We've so-far operated on vectors
  - Paying no heed to structured relationships among elements
- Images can also be processed as vectors
  - Ignoring that we are dealing with pixels
- But this ignores readily available domain knowledge about the structure of the problem



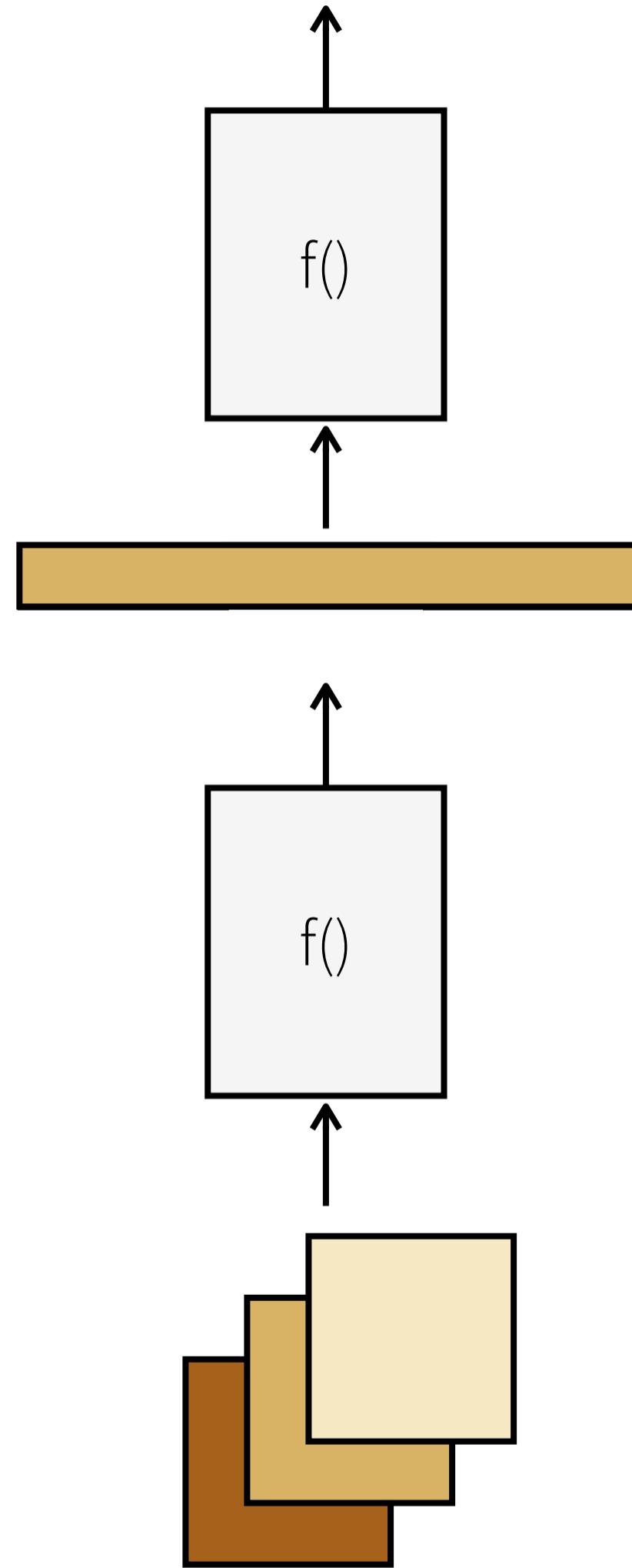
# Encoding Domain Knowledge

- We've so-far operated on vectors
  - Paying no heed to structured relationships among elements
- Images can also be processed as vectors
  - Ignoring that we are dealing with pixels
- But this ignores readily available domain knowledge about the structure of the problem

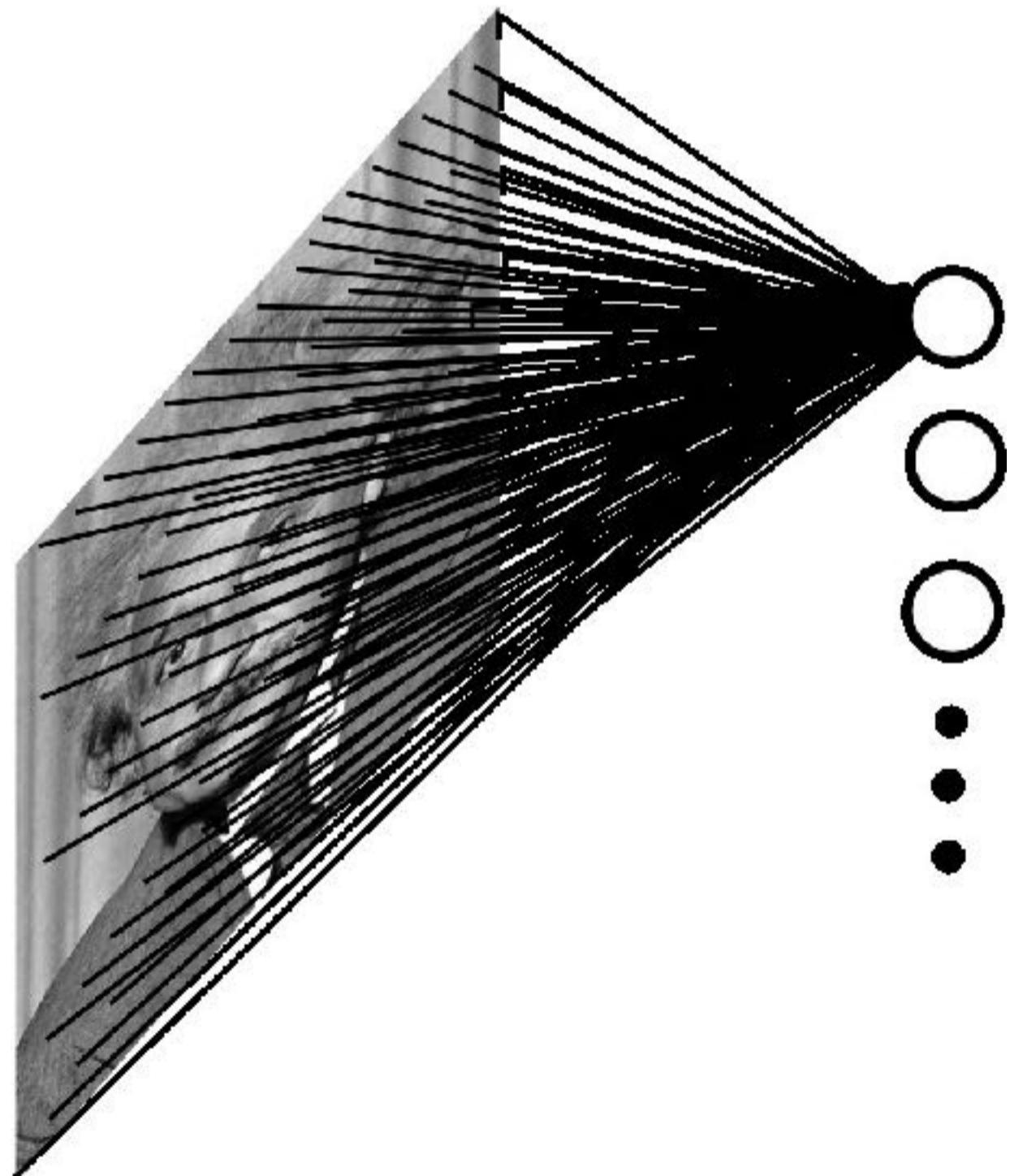


# Encoding Domain Knowledge

- We've so-far operated on vectors
  - Paying no heed to structured relationships among elements
- Images can also be processed as vectors
  - Ignoring that we are dealing with pixels
- But this ignores readily available domain knowledge about the structure of the problem



# Fully-Connected Layer

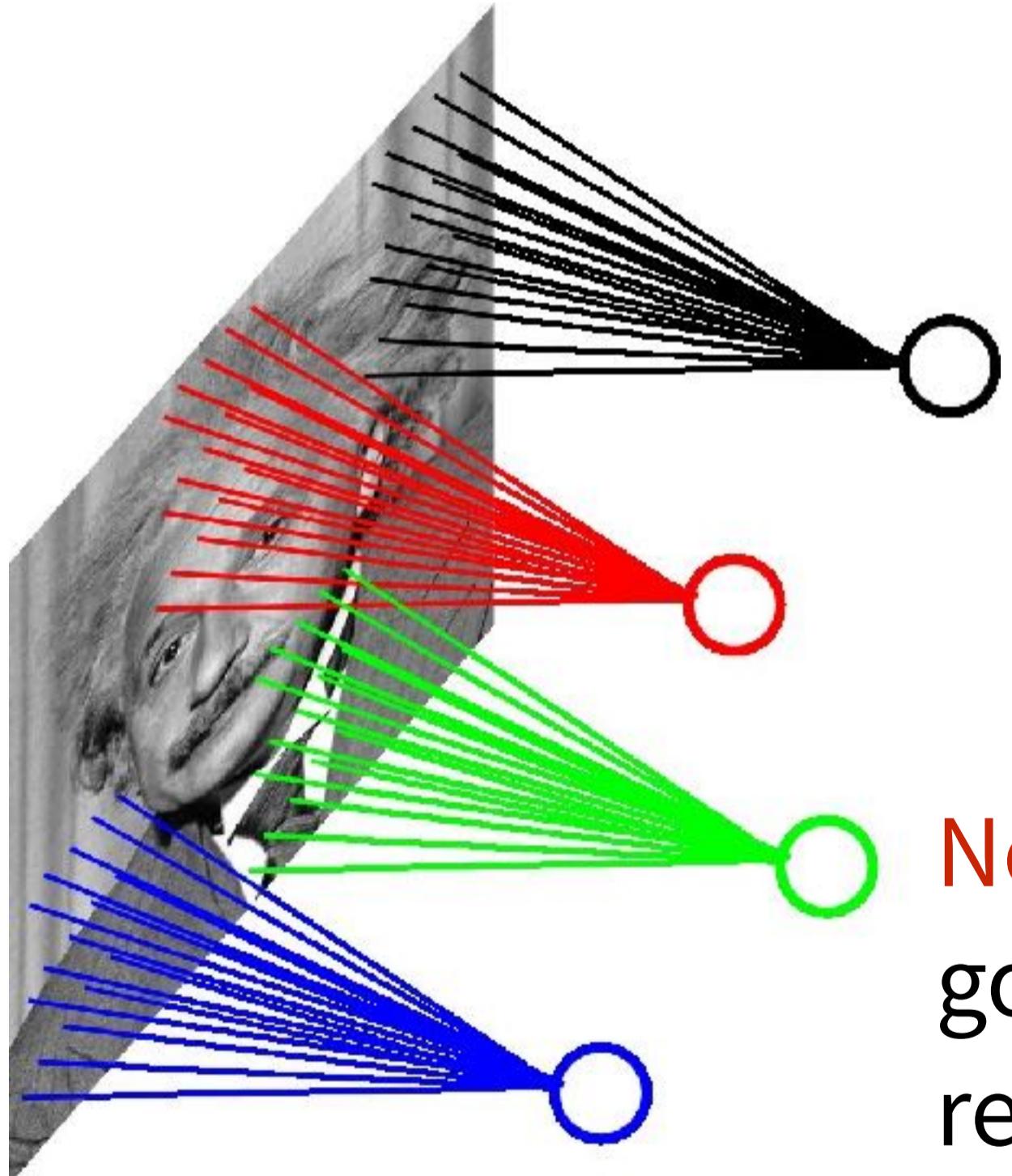


Example:  $200 \times 200$  image  
40k hidden units

→ ~2B parameters!

- Spatial correlation is local
- Waste of resources
- We don't have enough training examples to fit!

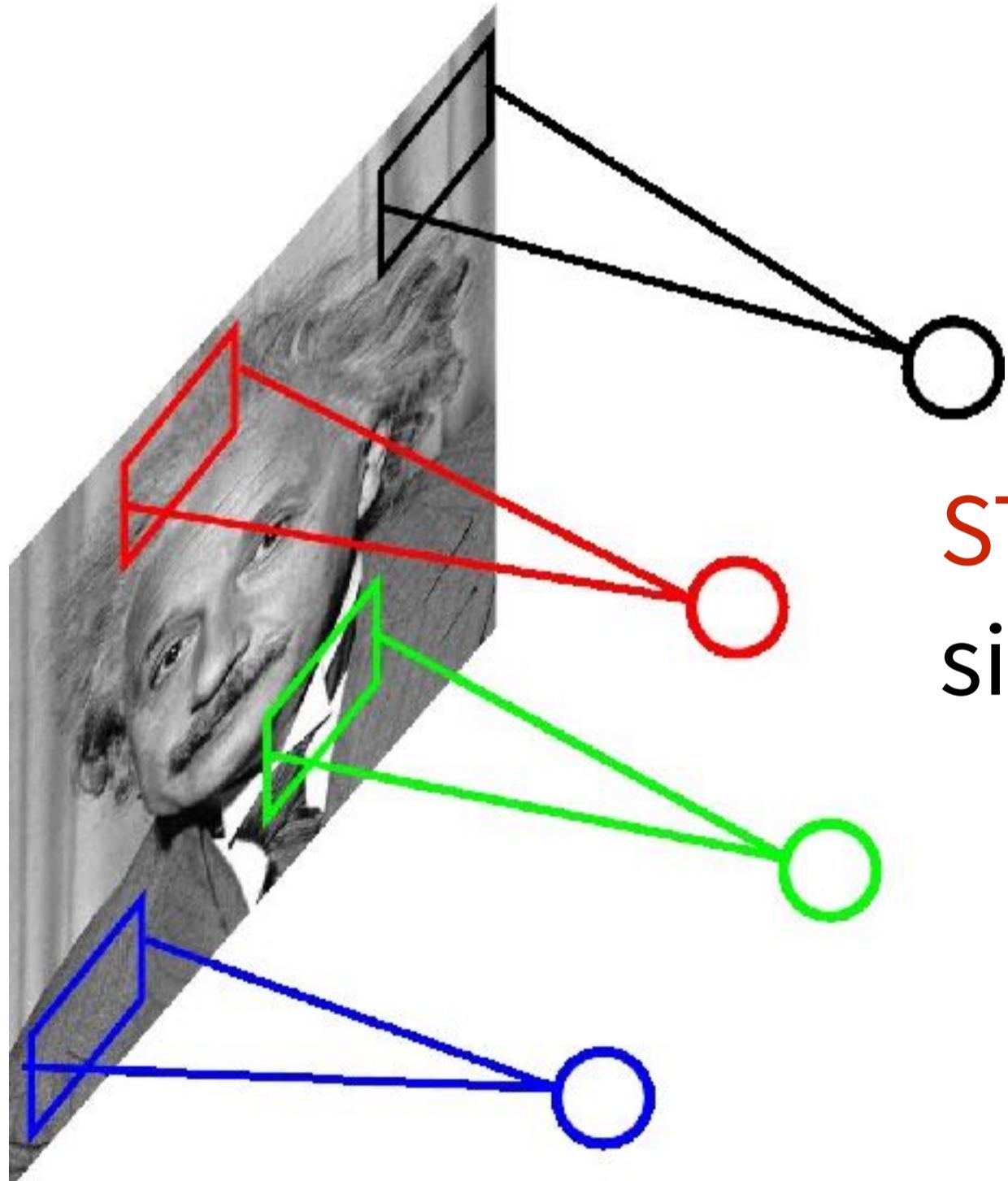
# Locally-Connected Layer



Example:  $200 \times 200$  image  
40k hidden units  
Filter size:  $10 \times 10$   
4M parameters

**Note:** this parameterization is good when the input images are registered (e.g. face recognition)

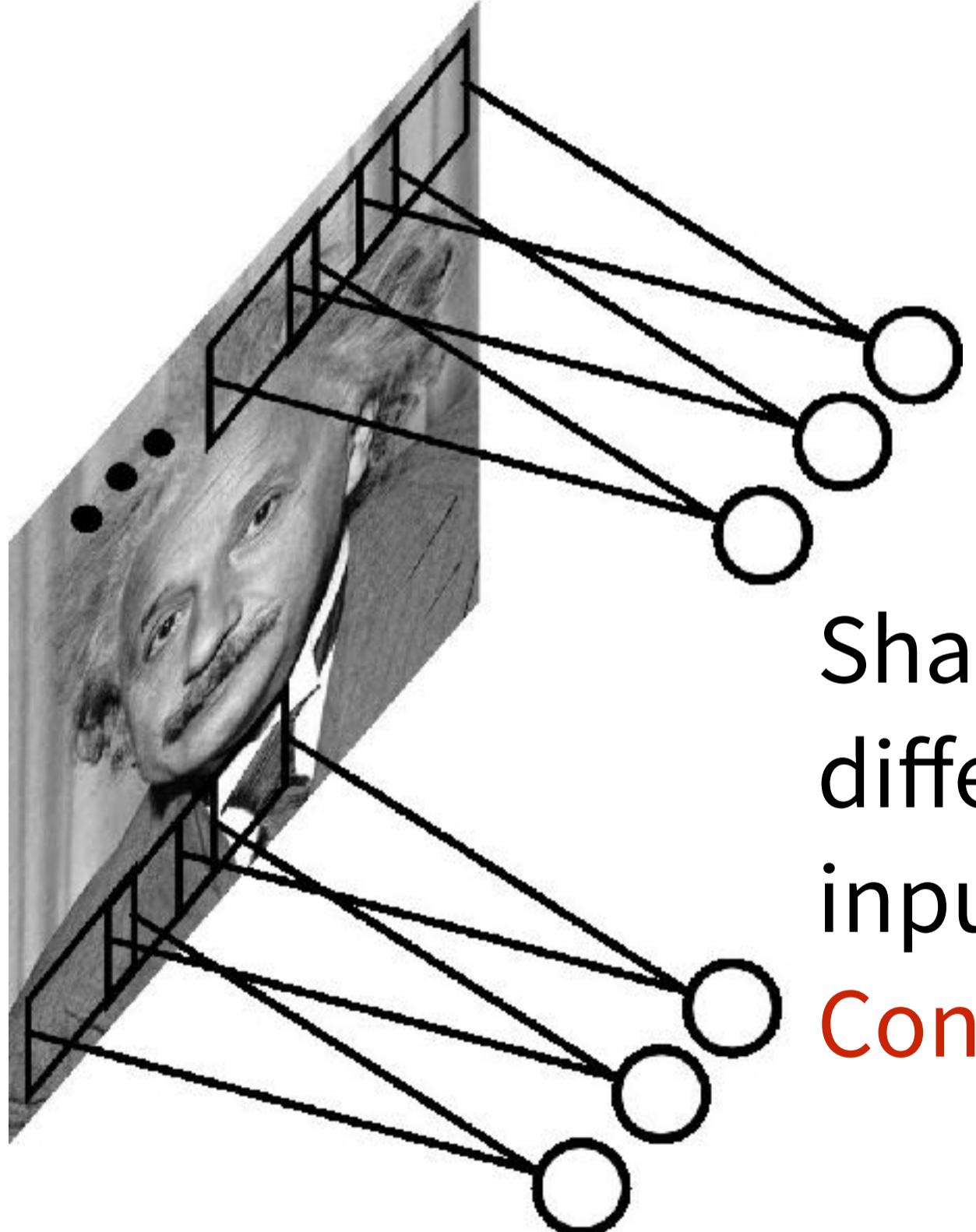
# Locally-Connected Layer



Example:  $200 \times 200$  image  
40k hidden units  
Filter size:  $10 \times 10$   
4M parameters

**STATIONARITY?** Statistics are similar at different locations.

# Weight Sharing

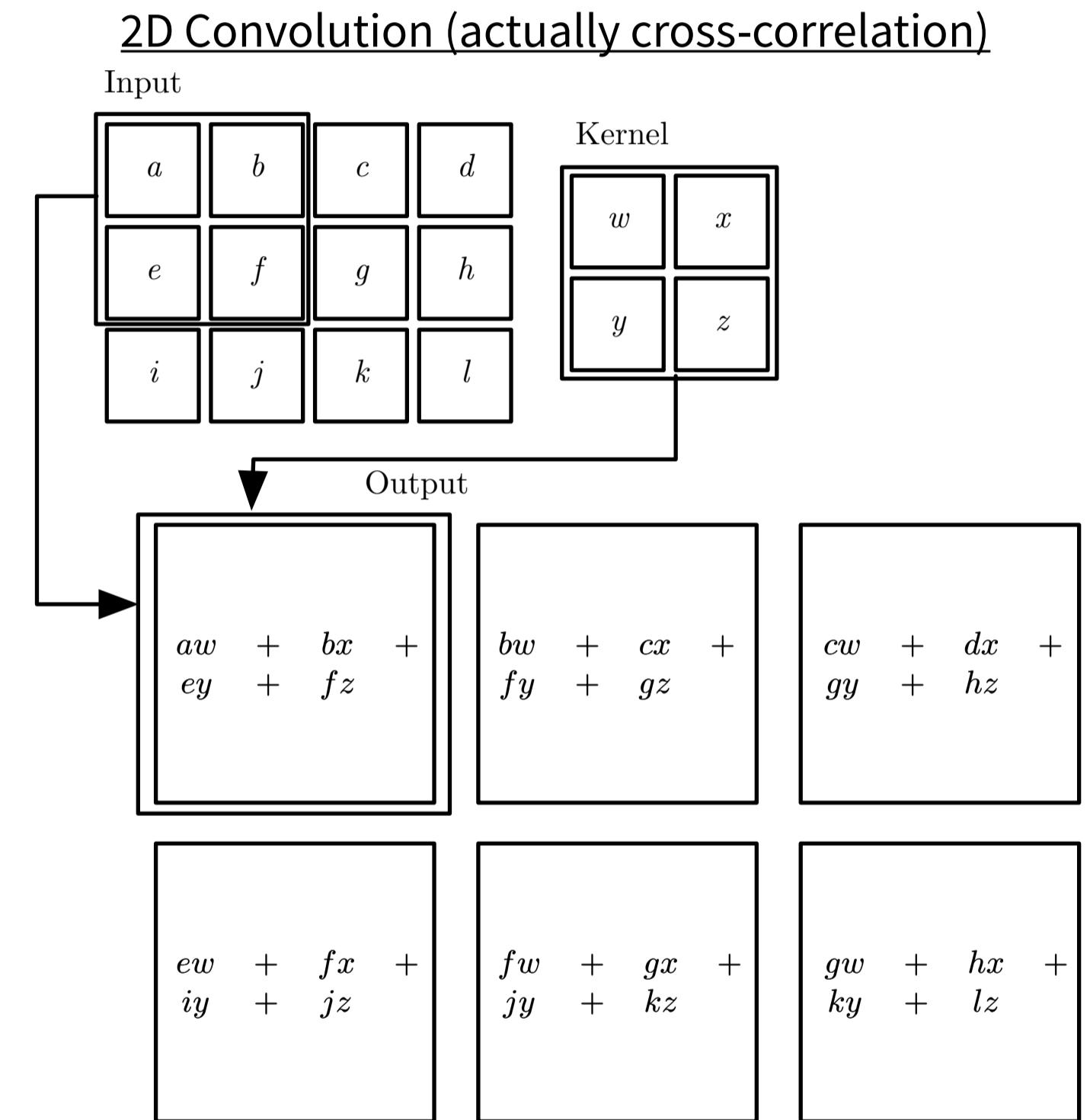


Share the same parameters across  
different locations (assuming  
input is stationary)  
**Convolutions with learned kernels**

# What is Convolution?

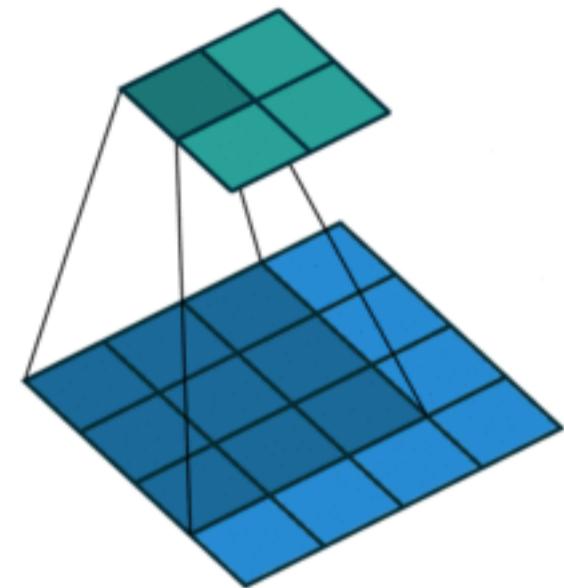
- Convolution is a specialized kind of **linear operation**
- Can think of it as a **weighted averaging** operation in time or space
- For vision, we “convolve” two-dimensional convolution kernels with two-dimensional images:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n)$$



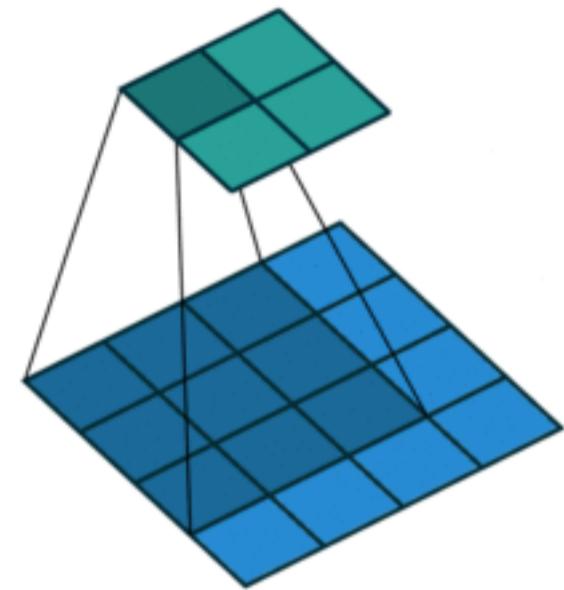
# Convolution (Animation)

- Note: hidden units are spatially organized into a “feature map”
- Each hidden unit only receives input from a local region of the input, its “receptive field”

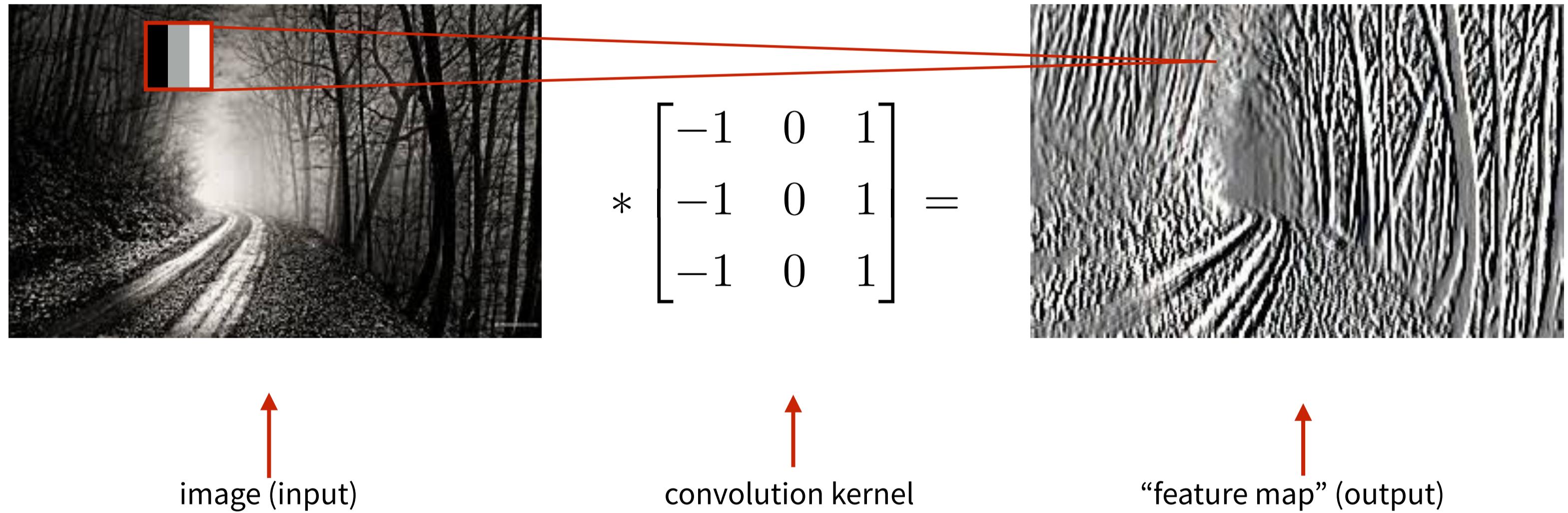


# Convolution (Animation)

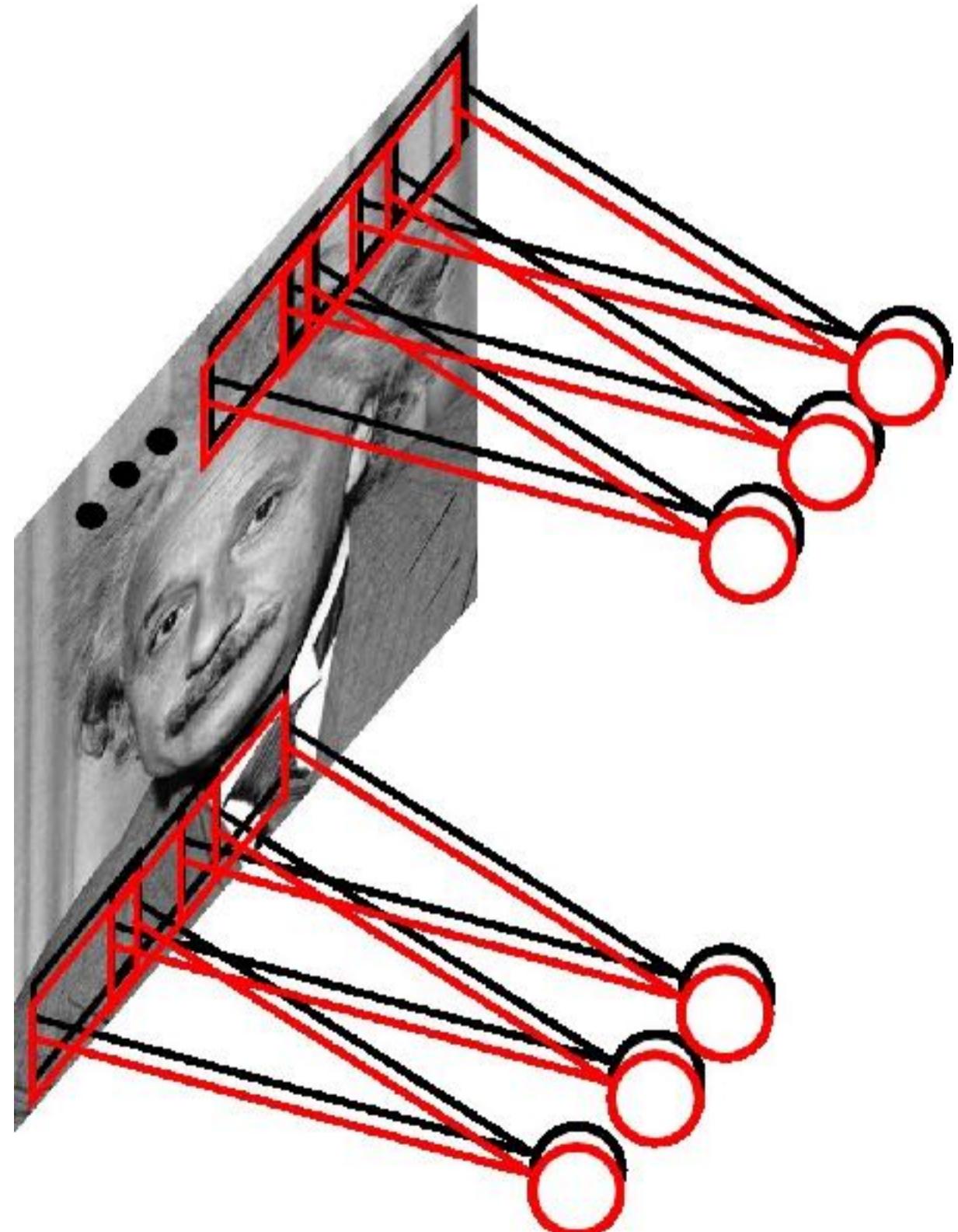
- Note: hidden units are spatially organized into a “feature map”
- Each hidden unit only receives input from a local region of the input, its “receptive field”



# Convolutional Layer



# Convolutional Layer



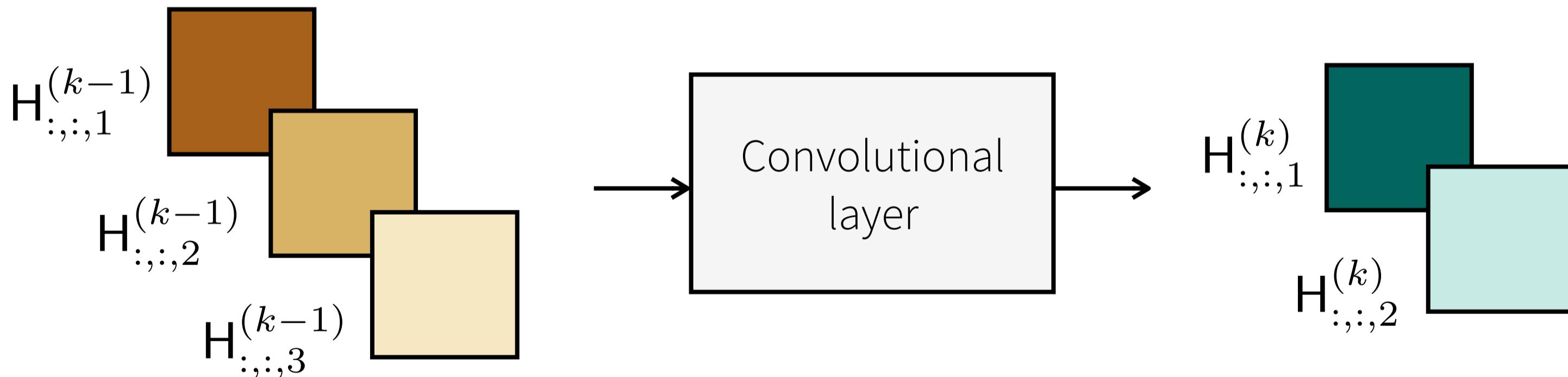
Learn multiple filters.

Example:  $200 \times 200$  image  
100 filters  
Filter size:  $10 \times 10$   
10k parameters

# Convolutional Layer

$$H_{:,:,i}^{(k)} = g \left( \sum_j H_{:,:,j}^{(k-1)} * W_{:,:,i,j}^{(k)} \right)$$

↑  
output  
feature map      ↑  
input  
feature map      ↑  
kernel

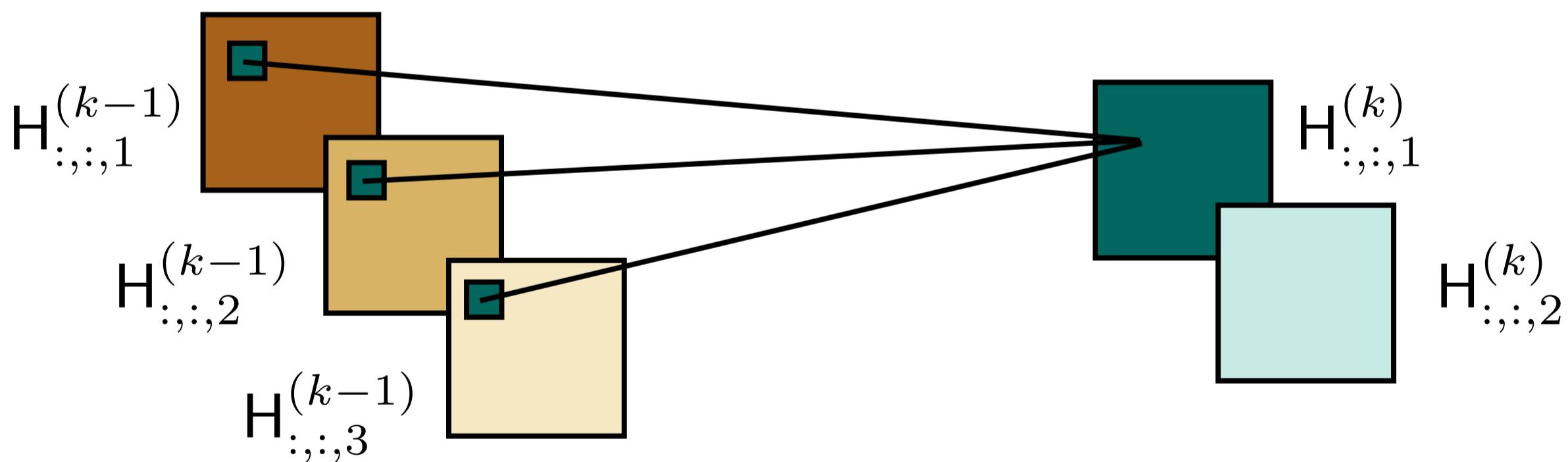


$k$  indexes layers  
 $i, j$  index maps (features)

# Convolutional Layer

$$H_{:,:,i}^{(k)} = g \left( \sum_j H_{:,:,j}^{(k-1)} * W_{:,:,i,j}^{(k)} \right)$$

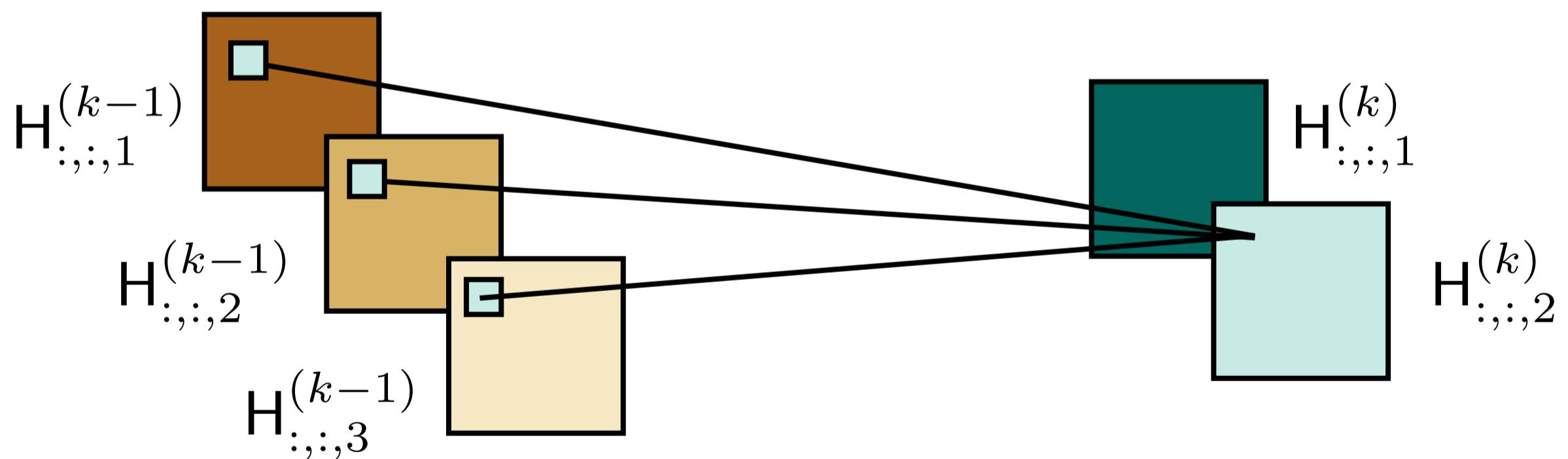
output feature map      input feature map      kernel



# Convolutional Layer

$$H_{:,:,i}^{(k)} = g \left( \sum_j H_{:,:,j}^{(k-1)} * W_{:,:,i,j}^{(k)} \right)$$

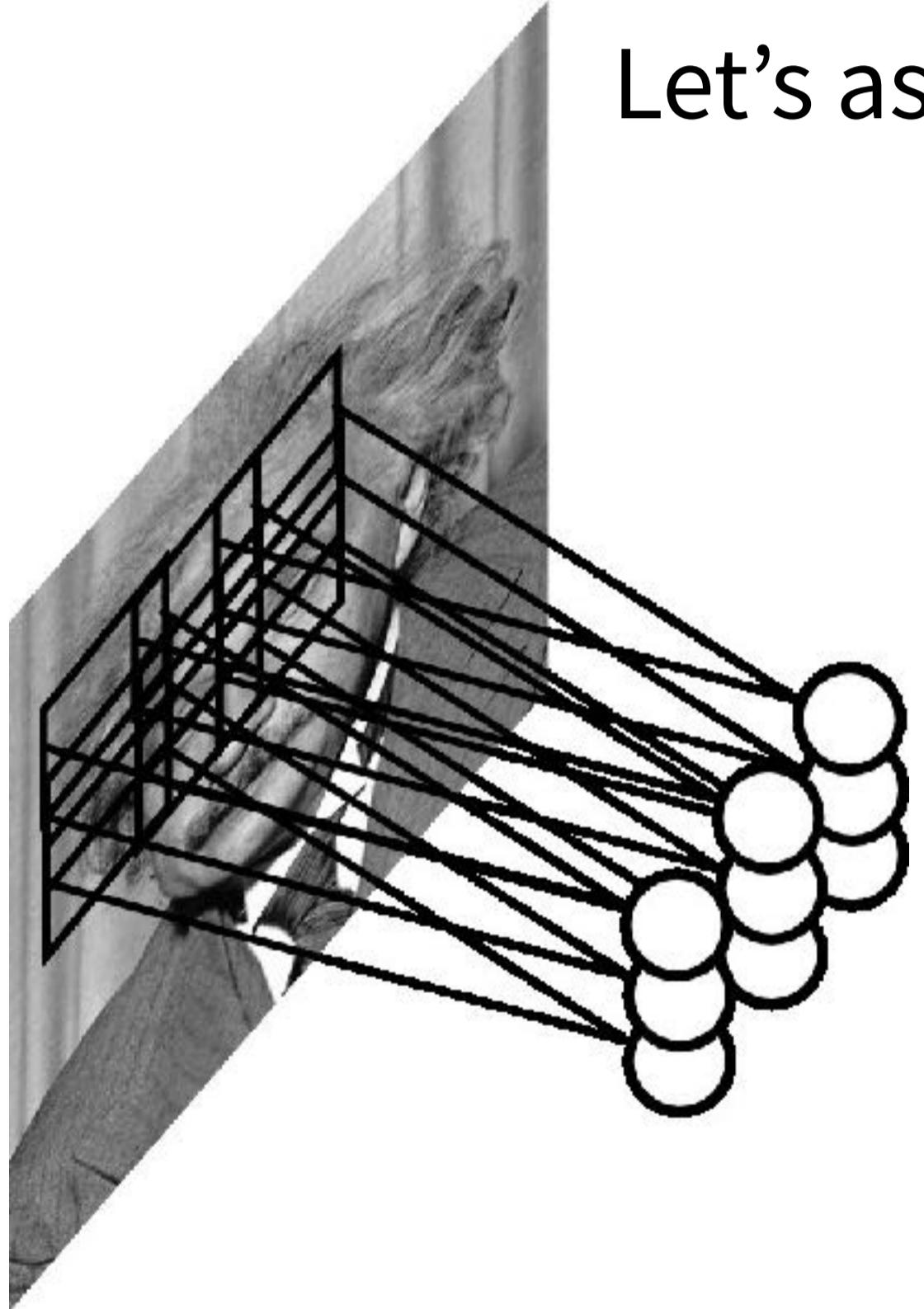
↑  
output  
feature map      ↑  
input  
feature map      ↑  
kernel



# Convolutional Net - Recap

- A **fully-connected** neural network applied to images
  - scales quadratically with the size of the input
  - does not leverage stationarity
- Solution
  - connect each hidden to a small patch of input
  - share the weights across space
- This is called a **convolutional layer**
- A network with convolutional layers is called a **convolutional net**

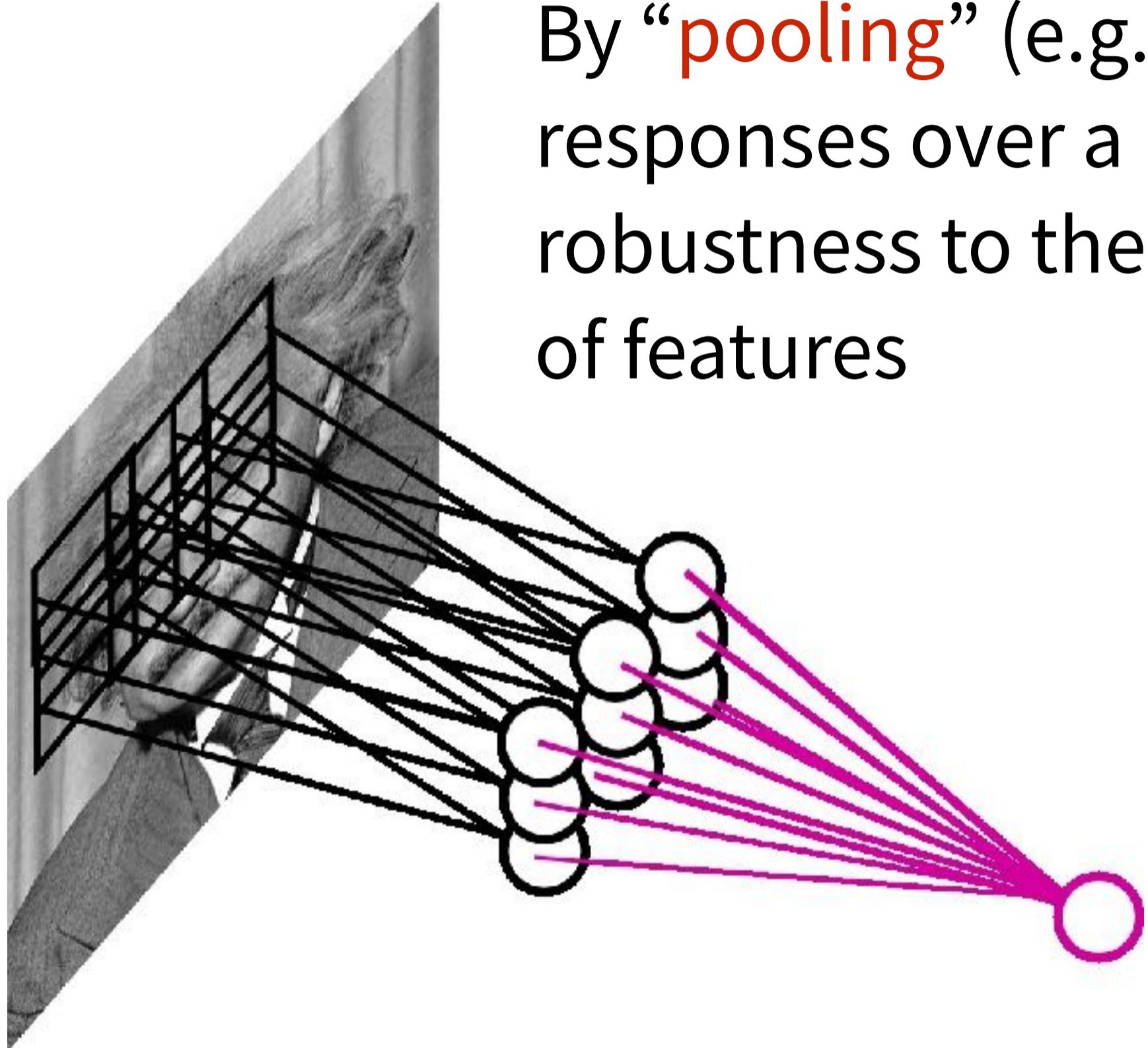
# Pooling Layer



Let's assume the filter is an “eye” detector

How can we make the detection  
robust to the exact location  
of the eye?

# Pooling Layer



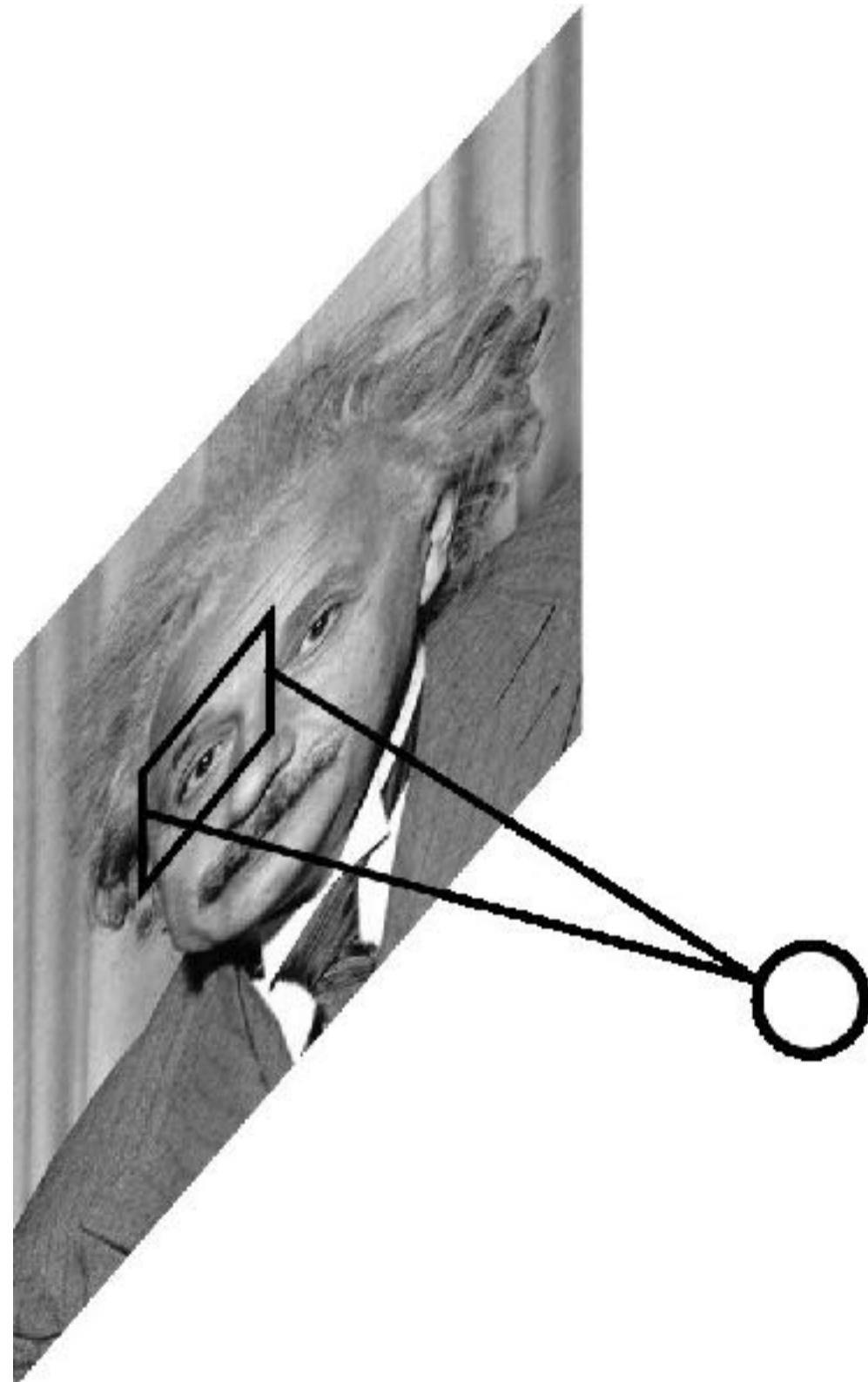
By “**pooling**” (e.g. taking max) the filter responses over a local region we gain robustness to the exact spatial location of features

“Shift invariance”

# Types of Pooling

Max-Pooling	$H_{m,n,i}^{(k)} = \max_{m' \in \mathbb{N}(m), n' \in \mathbb{N}(n)} H_{m',n',i}^{(k-1)}$
Average-Pooling	$H_{m,n,i}^{(k)} = \frac{1}{ \mathbb{N}(m)  \mathbb{N}(n) } \sum_{m' \in \mathbb{N}(m), n' \in \mathbb{N}(n)} H_{m',n',i}^{(k-1)}$
L2-Pooling	$H_{m,n,i}^{(k)} = \sqrt{\sum_{m' \in \mathbb{N}(m), n' \in \mathbb{N}(n)} (H_{m',n',i}^{(k-1)})^2}$
L2-Pooling Over Features	$H_{m,n,i}^{(k)} = \sqrt{\sum_{j \in \mathbb{N}(i)} (H_{m,n,j}^{(k-1)})^2}$

# Local Contrast Normalization



$$H_{m,n,i}^{(k)} = \frac{H_{m,n,i}^{(k-1)} - M_{m,n,i}^{(k-1)}}{S_{m,n,i}^{(k-1)}}$$

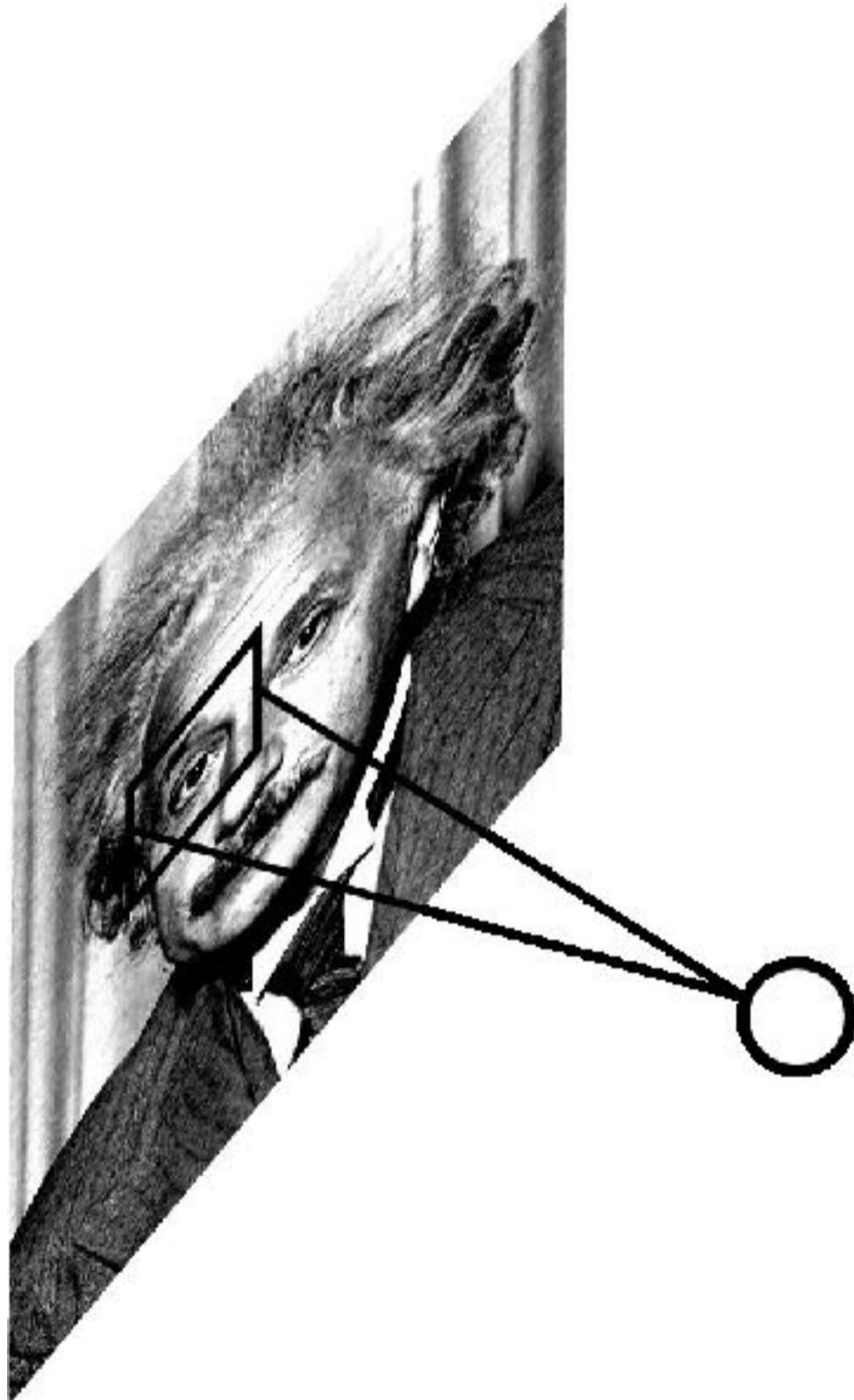
$M_{m,n,i}^{(k)}$

mean of a local neighbourhood  
around pixels  $m, n$

$S_{m,n,i}^{(k)}$

standard deviation of a local  
neighbourhood around pixels  $m, n$

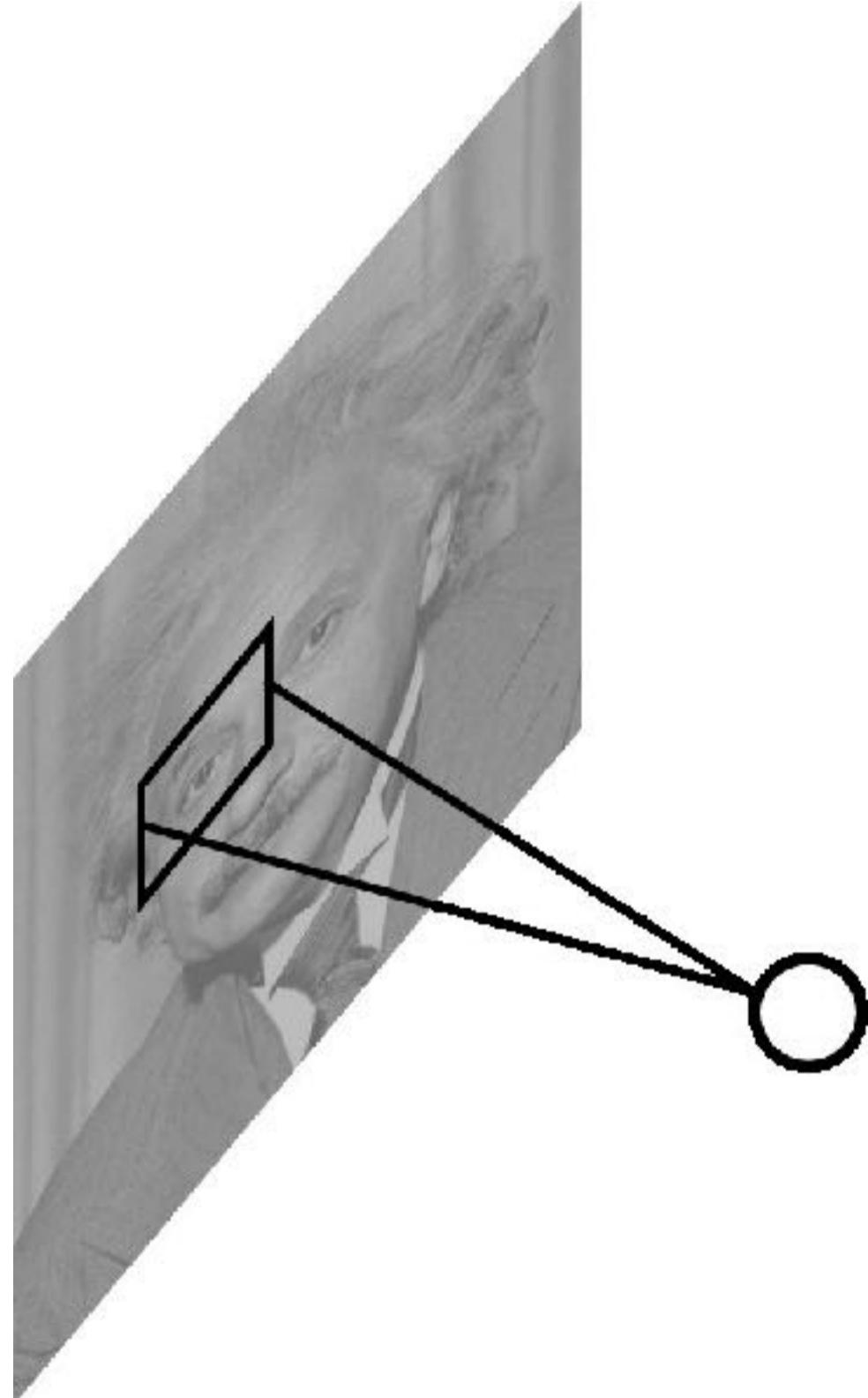
# Local Contrast Normalization



$$H_{m,n,i}^{(k)} = \frac{H_{m,n,i}^{(k-1)} - M_{m,n,i}^{(k-1)}}{S_{m,n,i}^{(k-1)}}$$

Controls the dynamic range  
of each feature map.

# Local Contrast Normalization



$$H_{m,n,i}^{(k)} = \frac{H_{m,n,i}^{(k-1)} - M_{m,n,i}^{(k-1)}}{S_{m,n,i}^{(k-1)}}$$

Performed also across features  
and in the higher layers...

Effects:

- improves invariance
- improves optimization
- increases sparsity

**Note:** computational cost is negligible  
compared to convolutional layer.

# Convnets: Single Stage

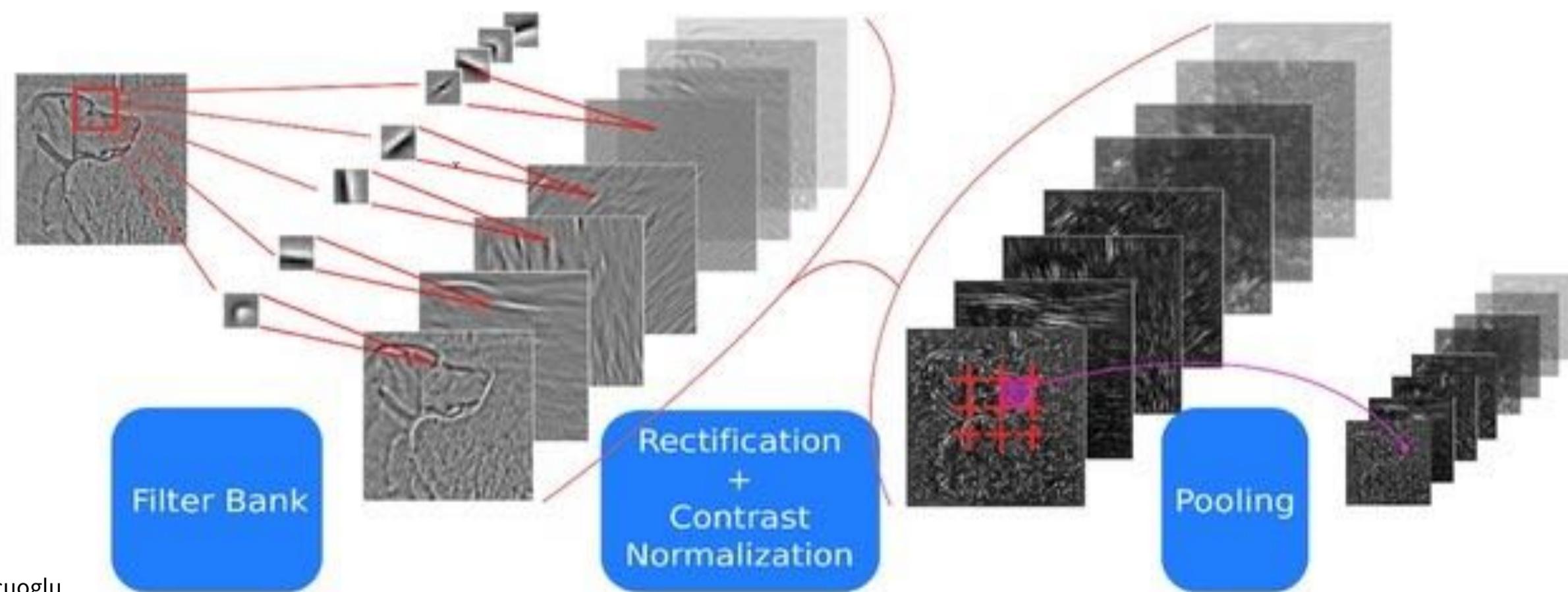


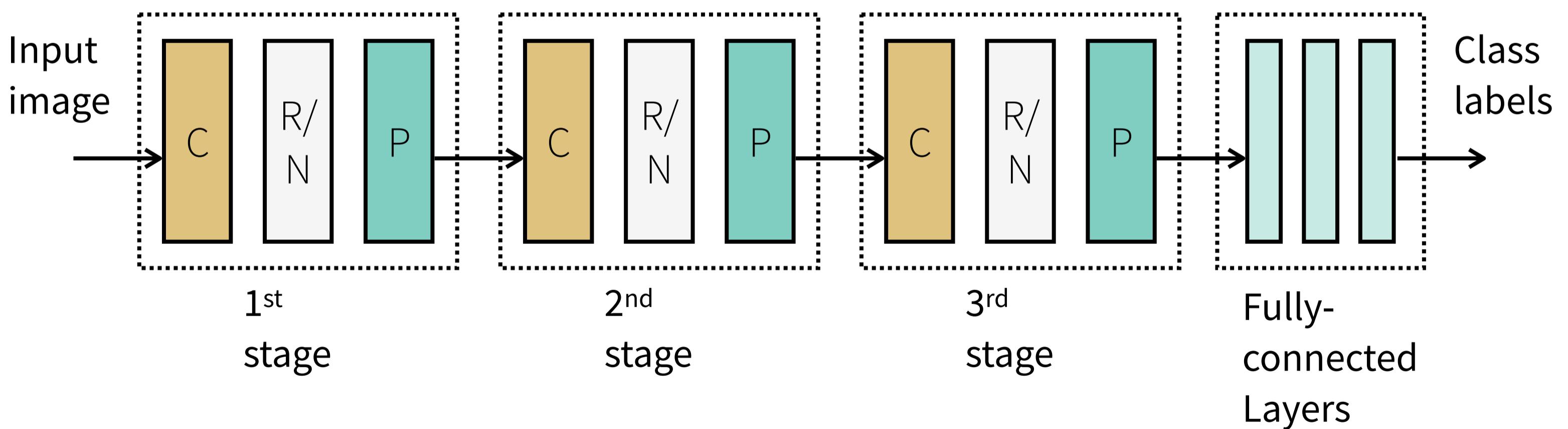
Image credit: Koray Kavukcuoglu

# Convnets: Typical Architecture

Single stage



Whole system

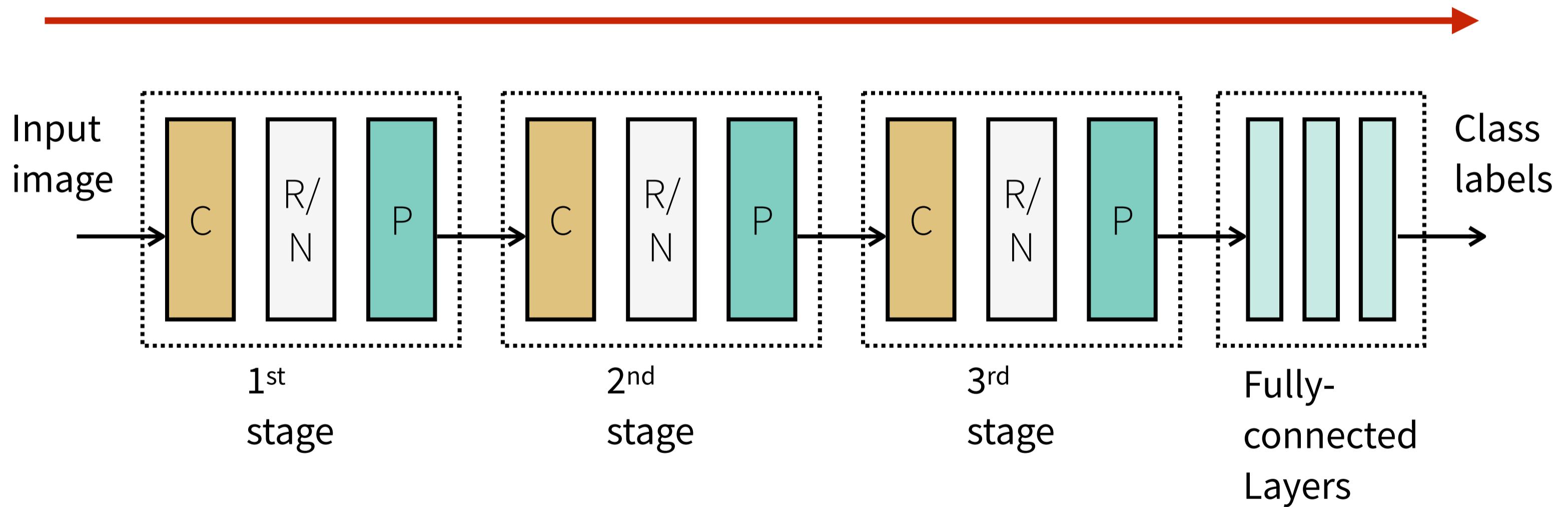


# Convnets: Training

- All layers are differentiable
- We can use standard back-propagation
- Algorithm:
  - Given a small mini-batch:
    - F-Prop
    - B-Prop
    - Parameter updates

# Convnets: Testing

At test time, only run forward propagation



Convnets can naturally process larger images at little cost.  
Traditional methods use inefficient sliding windows.

# Convnets: today

- Until 2012, the common wisdom was that training didn't work because we would “**get stuck in local minima**”
- Local minima are all similar, there are long plateaus, and it can take a long time to break symmetries
- Optimization is not the real problem, when:
  - the dataset is large
  - units do not saturate much
  - we use normalization layers

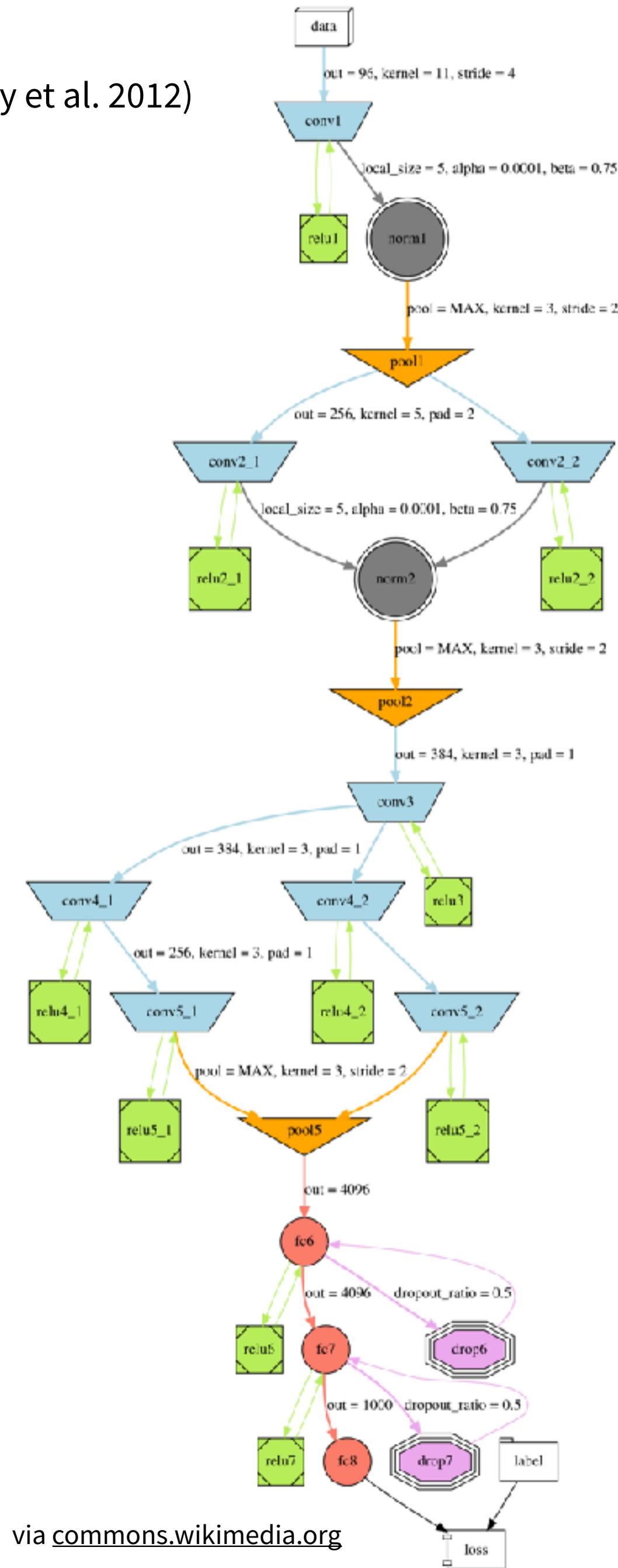
# Scalability

- The real challenges are:
  - Generalization
    - How many training examples to fit 1B params?
    - How many parameters/samples to model spaces with 1M dimensions?
  - Scalability

# Profile: AlexNet

(Krizhevsky et al. 2012)

- First work to popularize DNNs in vision
- Astonishing result on ILSVRC 2012 (16% top-5 error compared to 26% for runner-up)
- Similar to early LeNet architectures but deeper, bigger, and stacked convolutional layers (without pooling)
- Uses dropout regularization
- Optimized for 2 GPUs

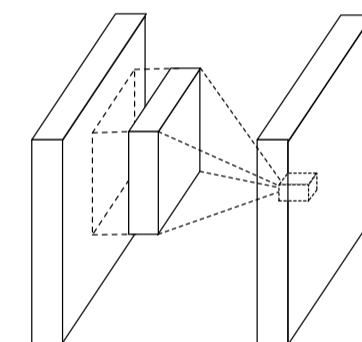


# Profile: Network-in-Network

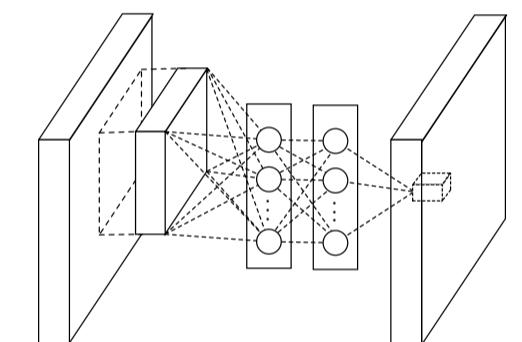
Not an ILSVRC competitor, this paper proposed a number of novel ideas which were **utilized by later architectures**:

- Replaces linear convolution layer with a “micro network” structure which is a general nonlinear operation
- Noting that fully connected layers are prone to overfitting, proposed “Global Average Pooling”

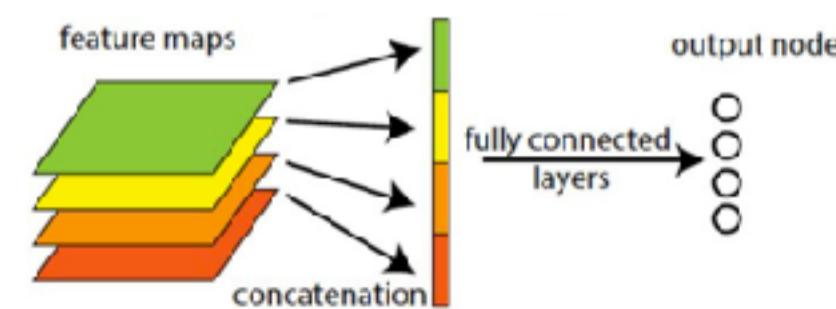
Linear convolution layer



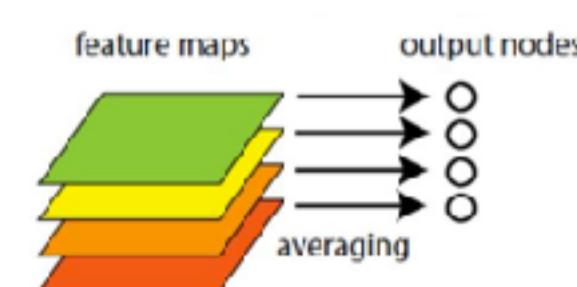
Mlpconv layer



CNN

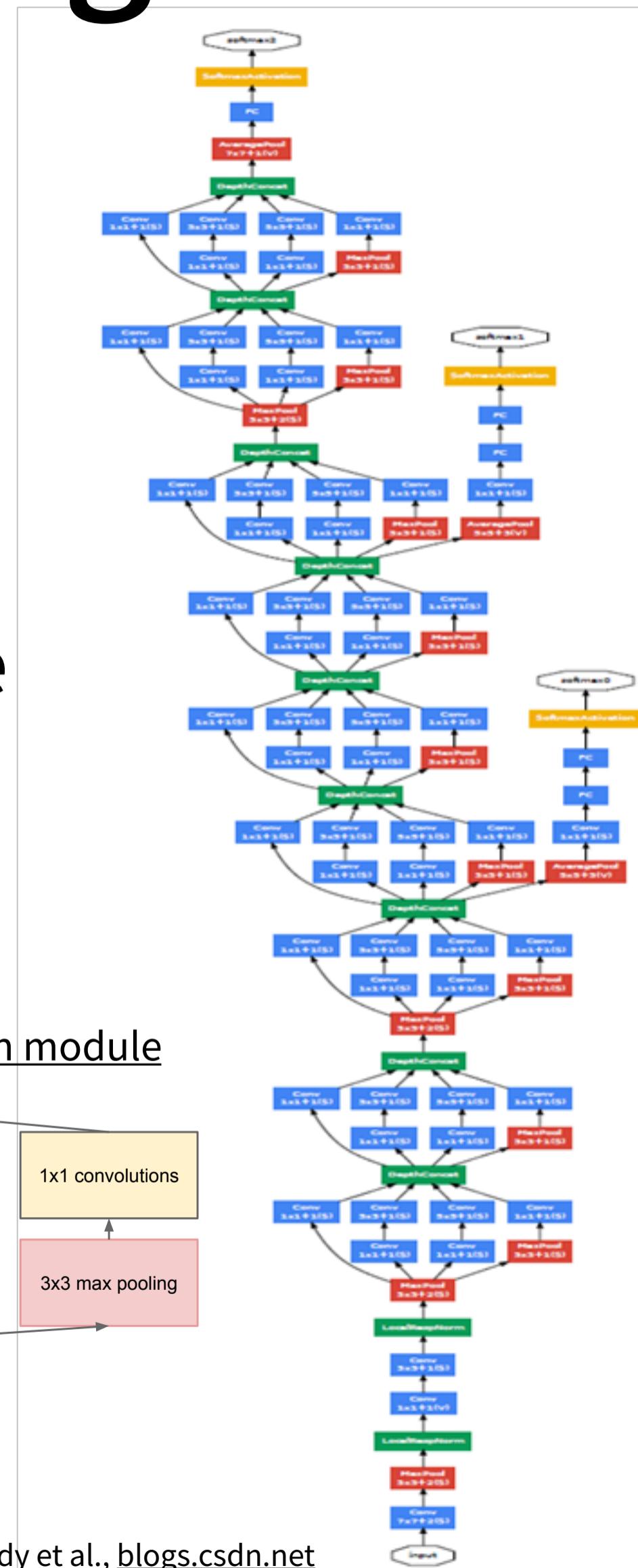
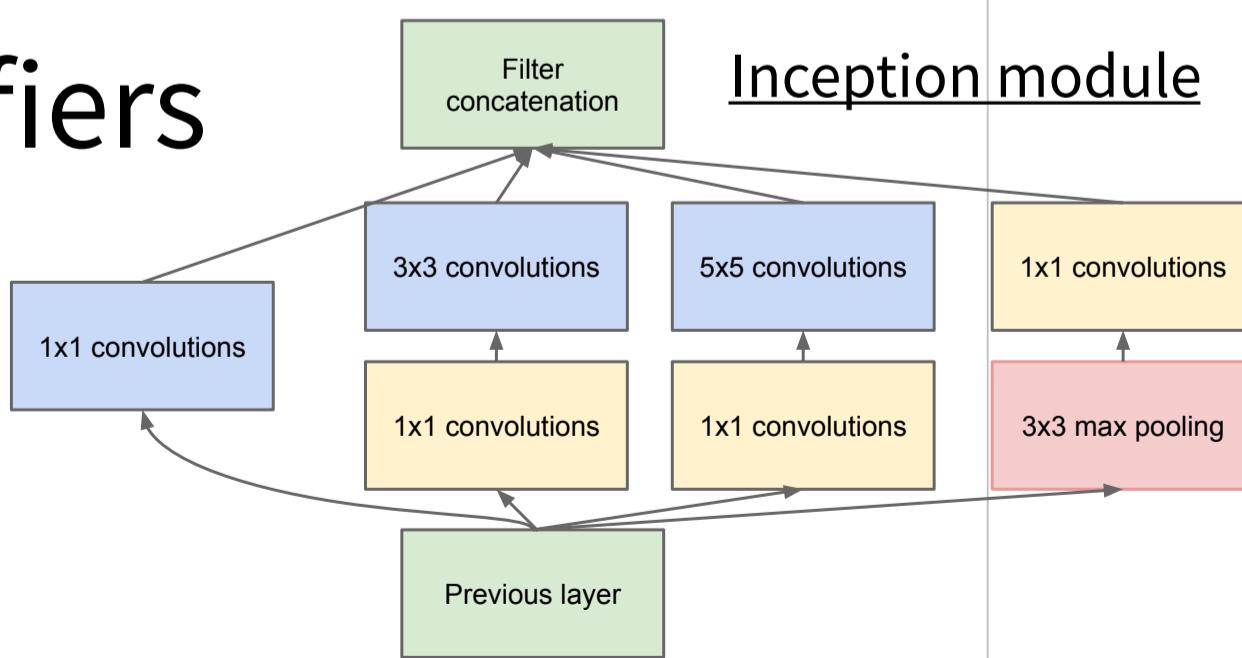


GAP



# Profile: Inception/GoogLeNet

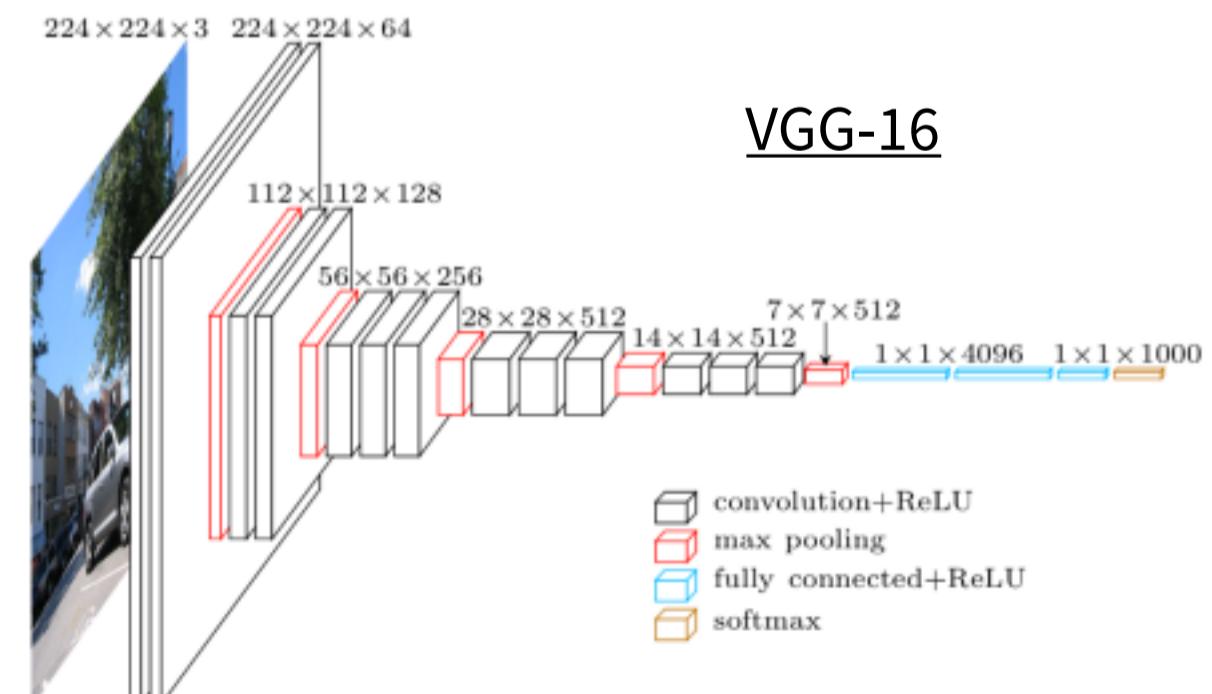
- Winner of ILSVRC 2014
- Main contribution is the “Inception Module” that dramatically reduced the number of parameters in the network (4M compared to 60M for AlexNet)
- Uses auxillary classifiers



# Profile: VGGNet

(Simonyan and Zisserman 2014)

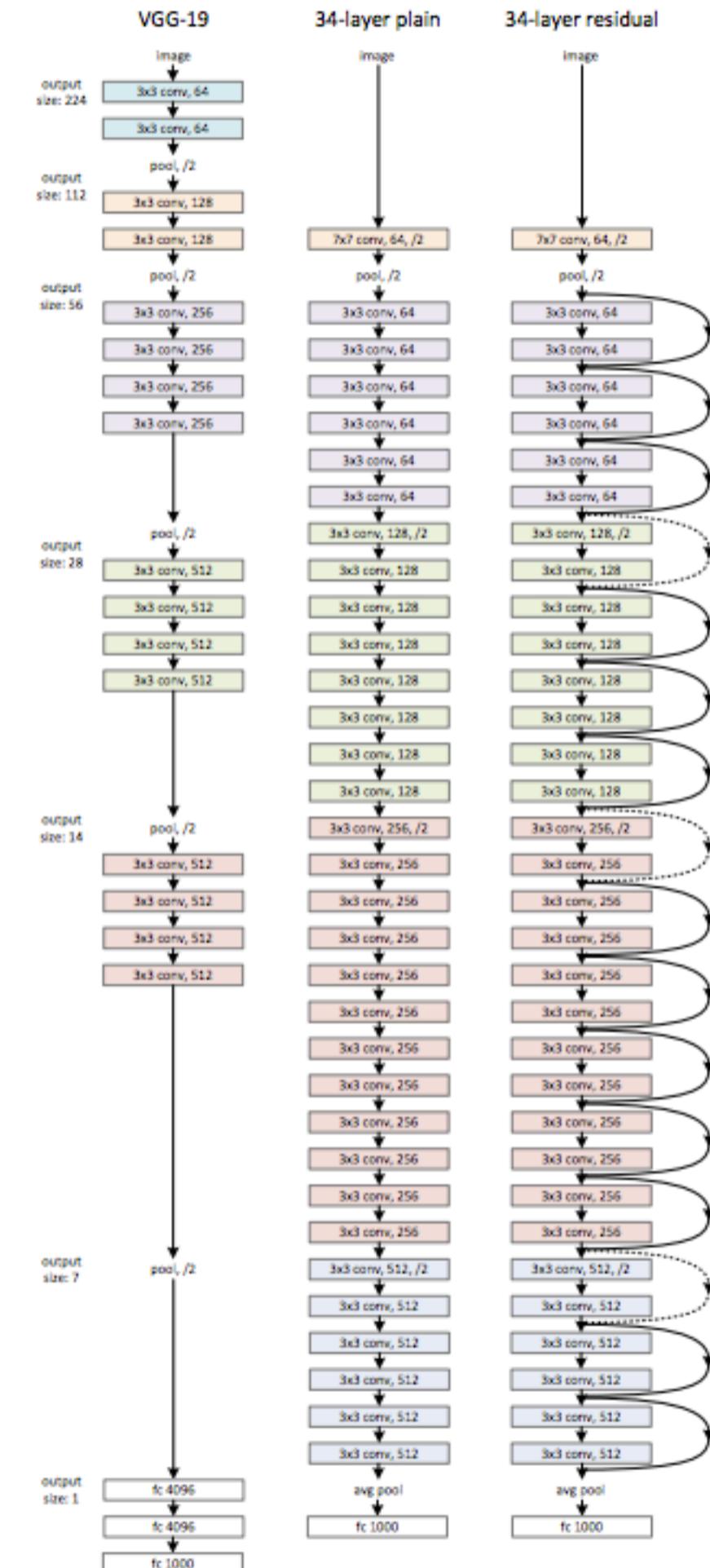
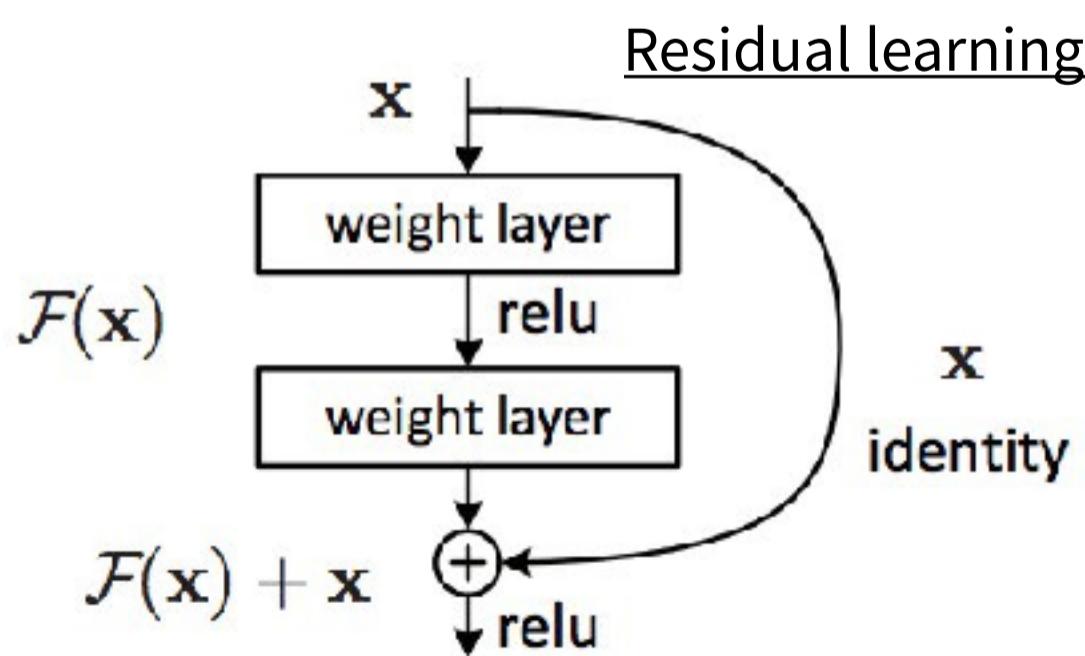
- Runner-up to ILSVRC 2014
- Main contribution: showing that depth of the network is a critical component for good performance (16 CONV/FC layers)
- Extremely homogeneous architecture that only performs 3x3 convolutions and 2x2 pooling from beginning to end
- Expensive to evaluate, uses a lot of memory and has 140M parameters



# Profile: ResNet

(He et al. 2015)

- Winner of ILSVRC 2015
  - It features special **skip connections** and heavy use of **batch normalization**
  - The architecture is also missing fully connected layers at the end of the network (note use of GAP)
  - Winning network had 152 layers; 1000-layer models have been trained



# Modern Architectures

