

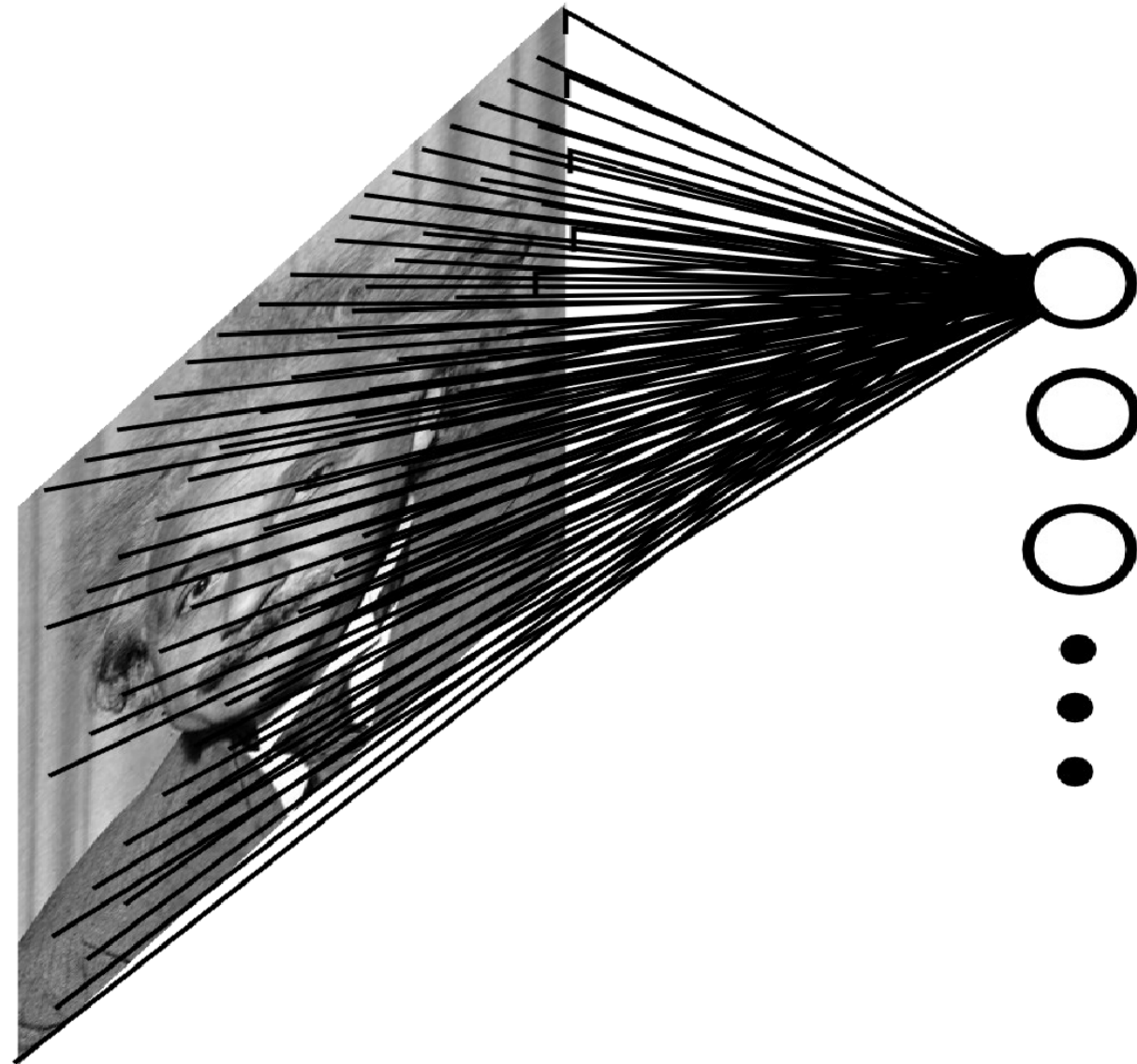
CSE 152: Computer Vision

Hao Su

Convolutional Neural Network



Images as input to neural networks

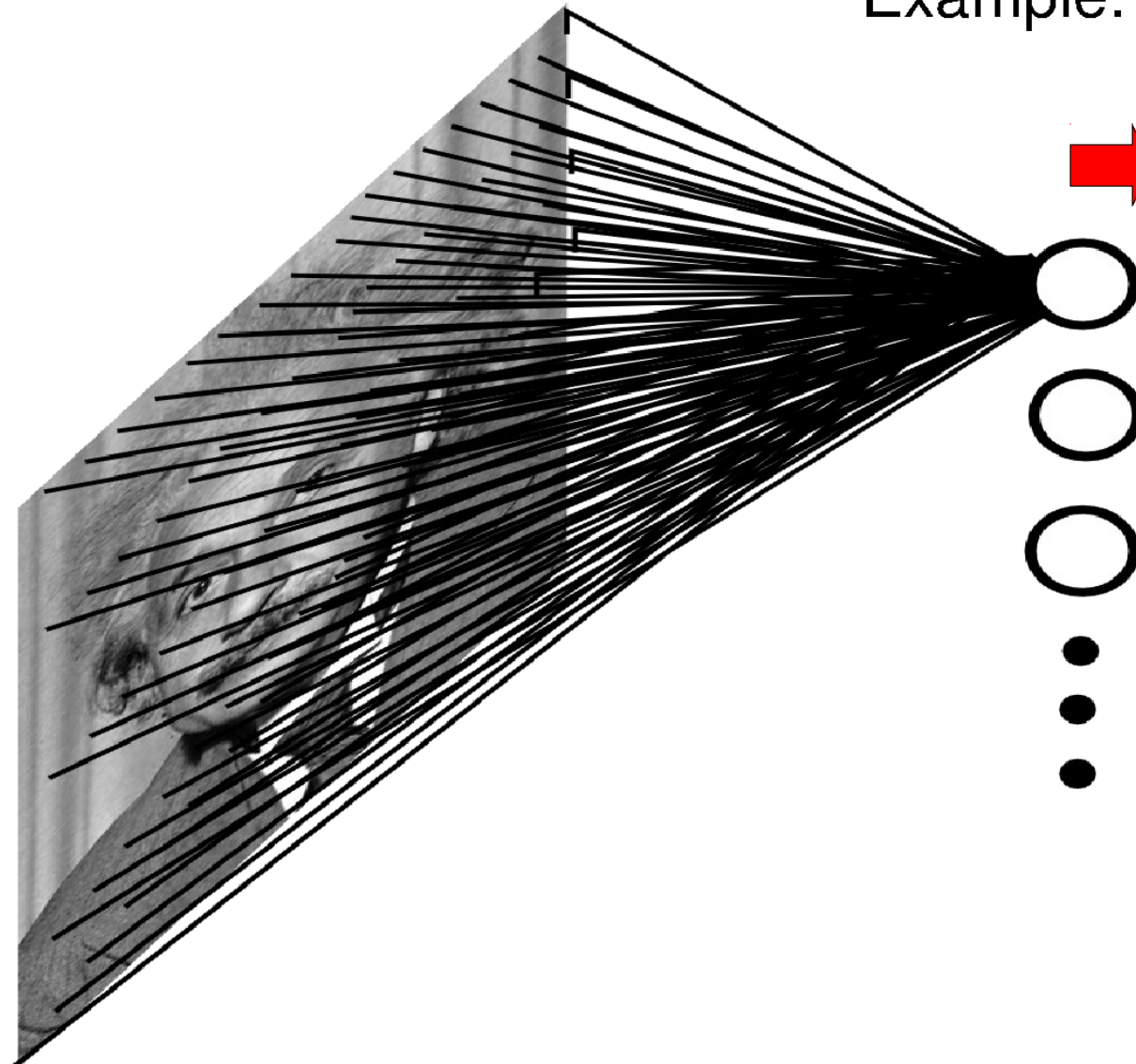


Images as input to neural networks

Example: 200x200 image

40K hidden units

→ **~2B parameters!!!**

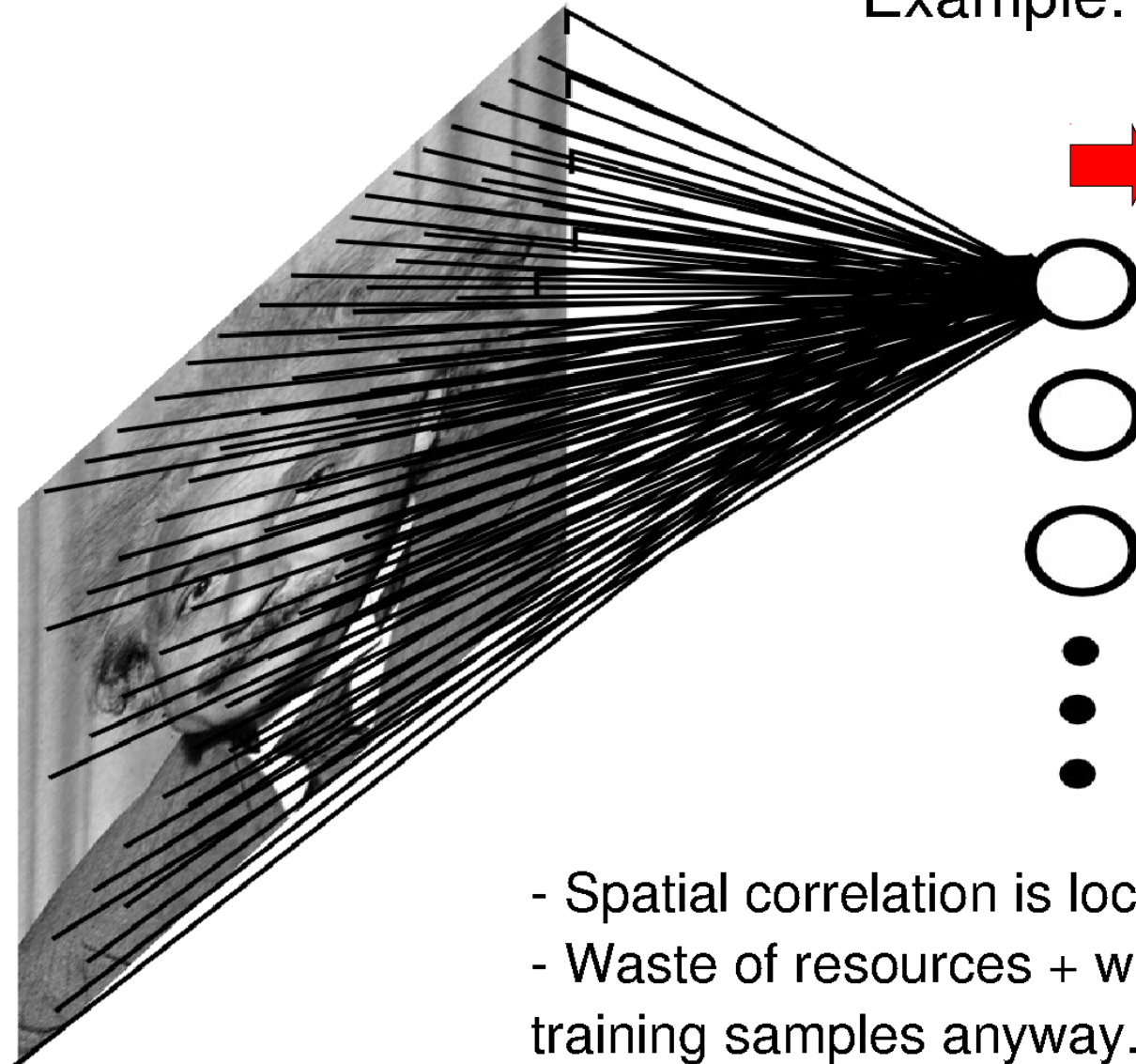


Images as input to neural networks

Example: 200x200 image

40K hidden units

➔ **~2B parameters!!!**



- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..

Convolutional Neural Networks

- CNN = a multi-layer neural network with
 - **Local** connectivity:
 - Neurons in a layer are only connected to a small region of the layer before it
 - **Share** weight parameters across spatial positions:
 - Learning shift-invariant filter kernels

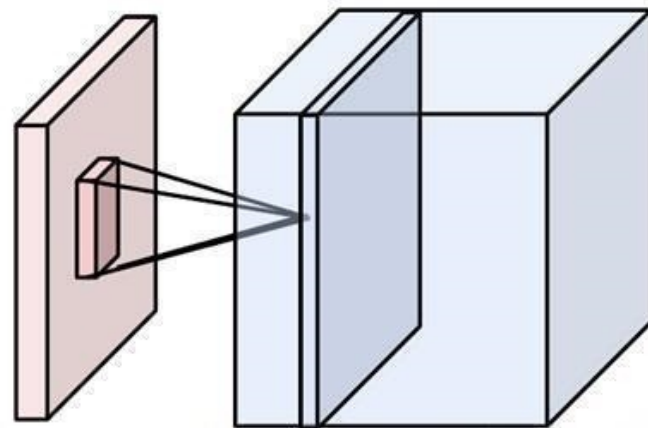
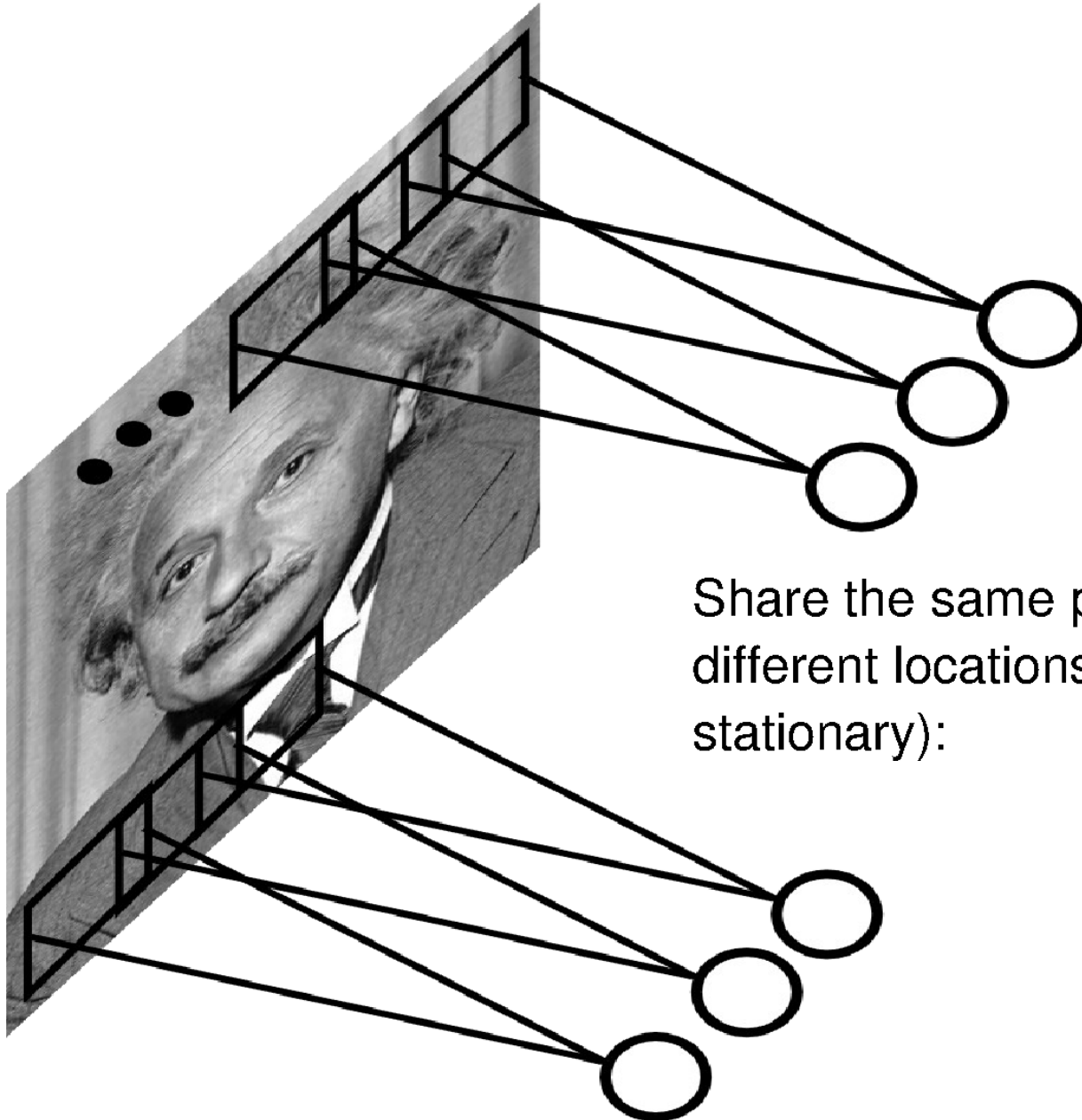
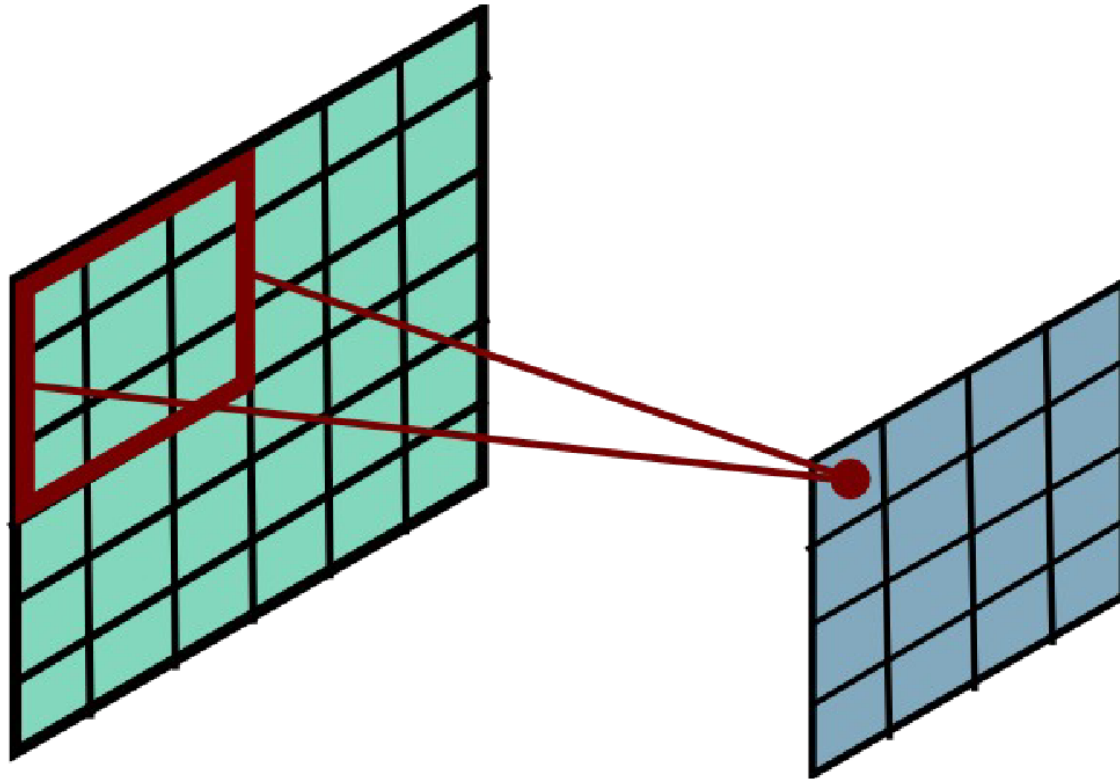


Image credit: A. Karpathy

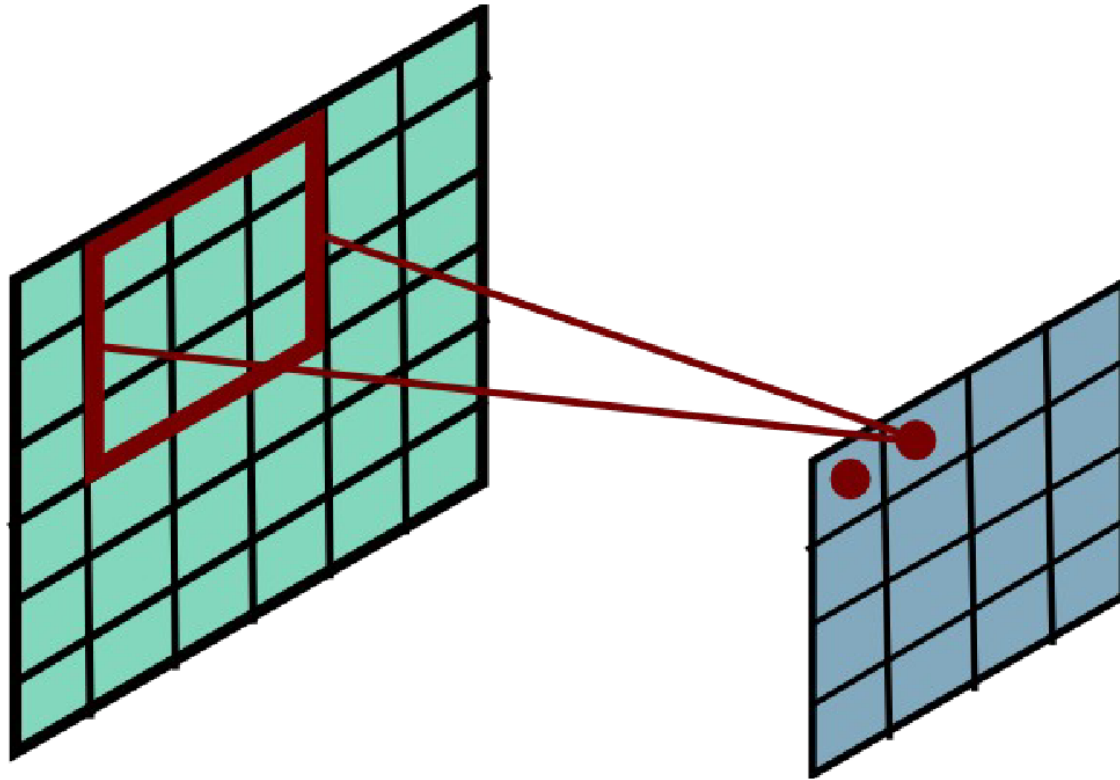


Share the same parameters across different locations (assuming input is stationary):

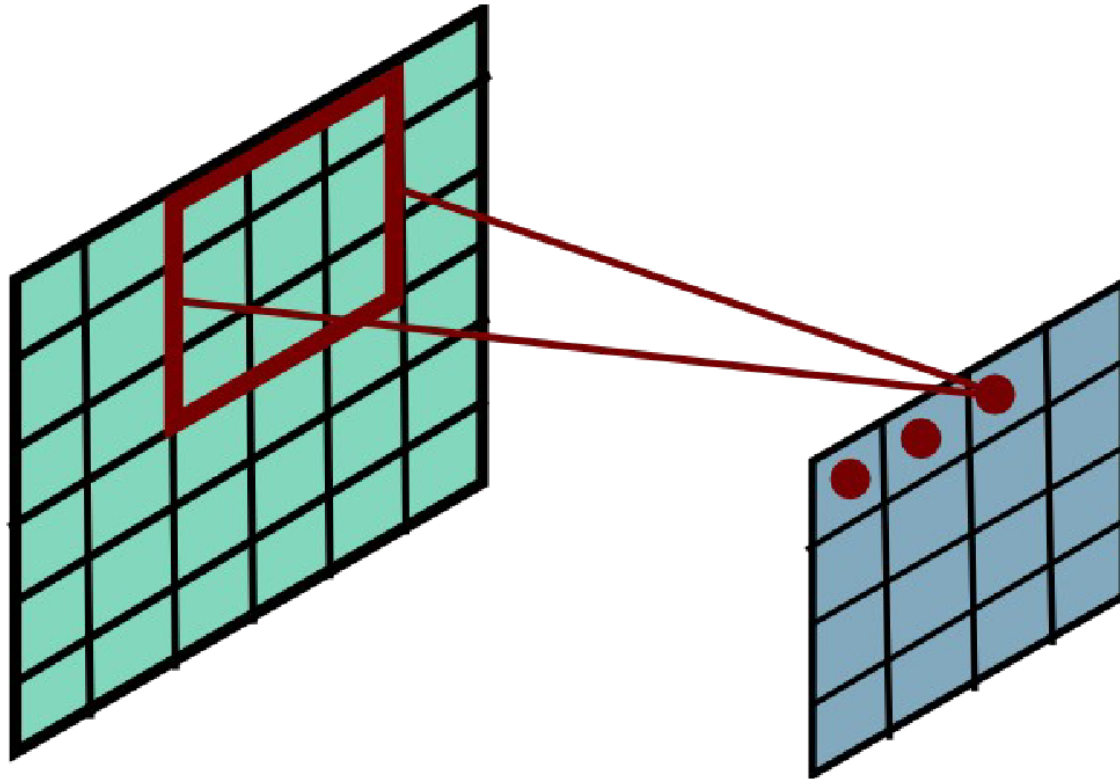
Convolutional Layer



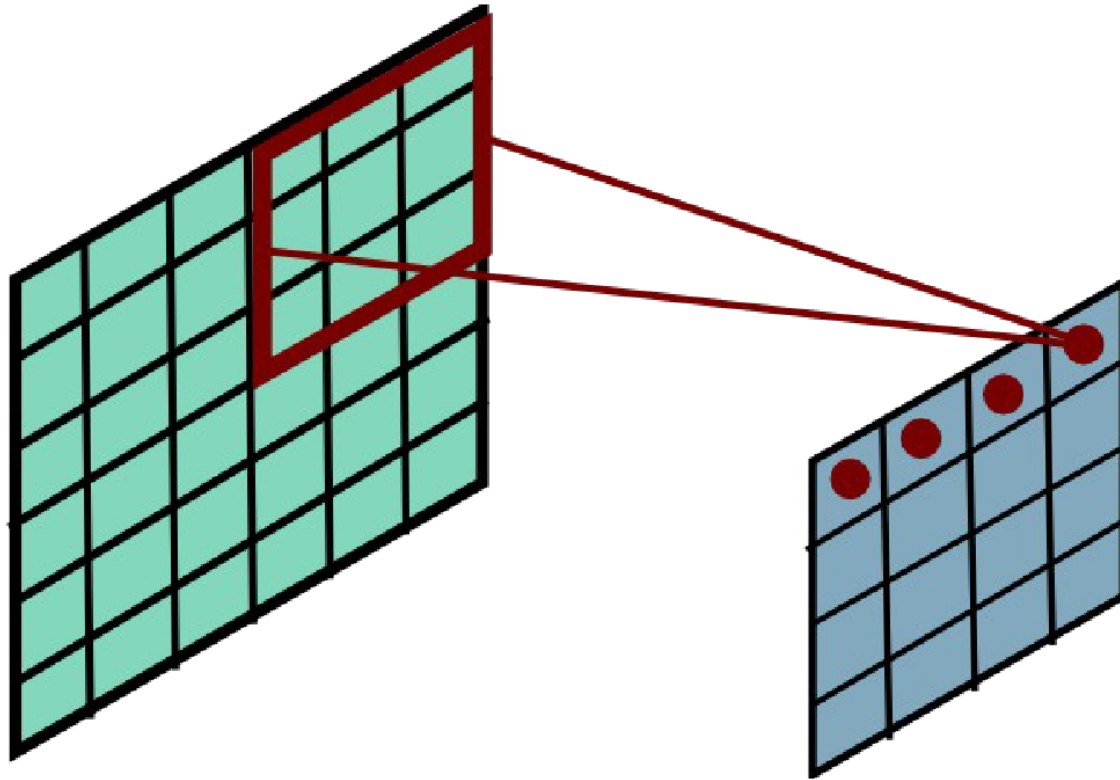
Convolutional Layer



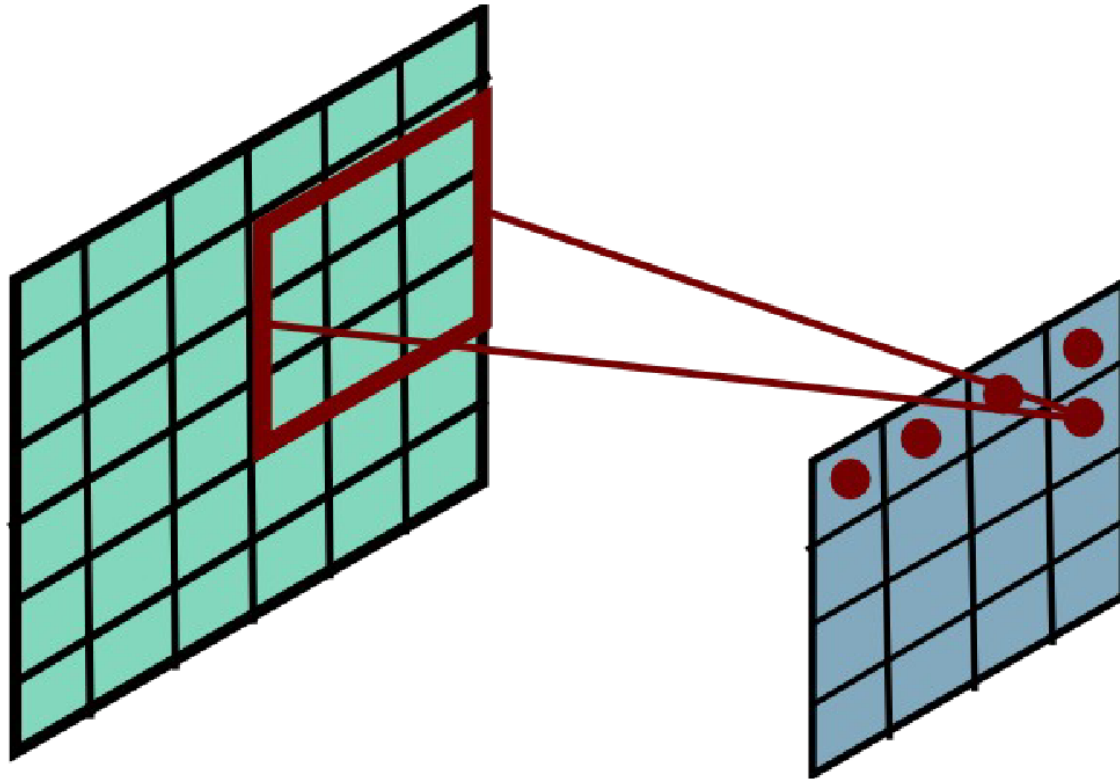
Convolutional Layer



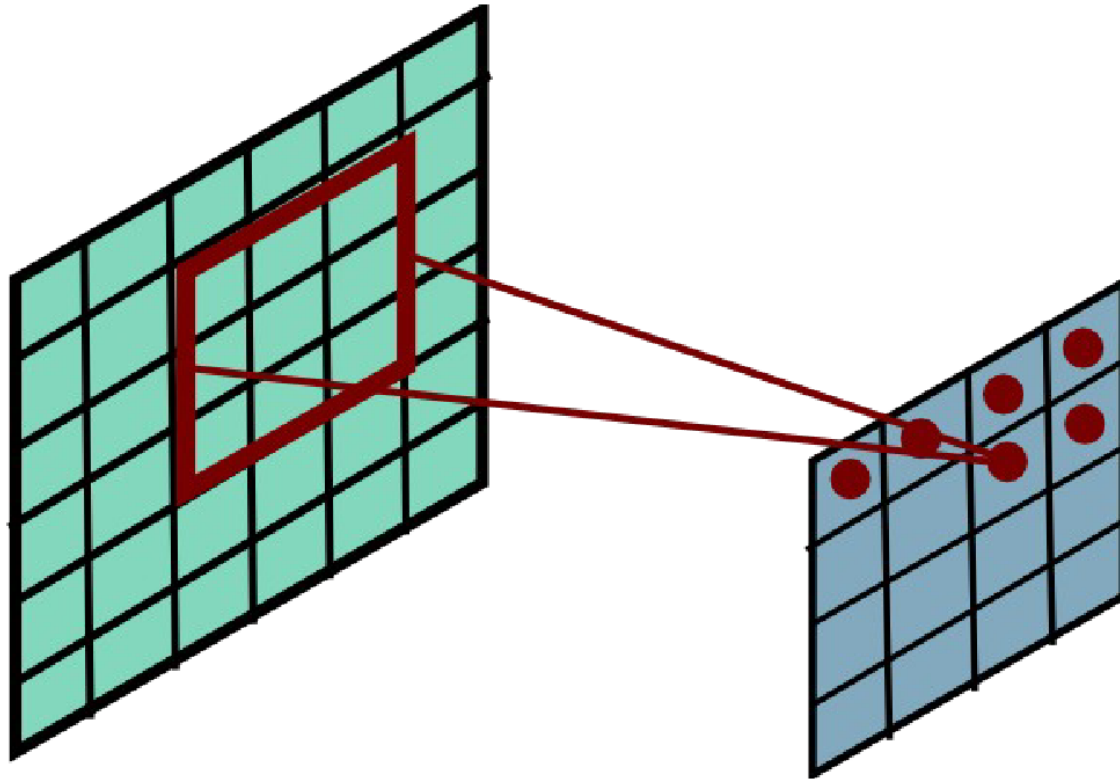
Convolutional Layer



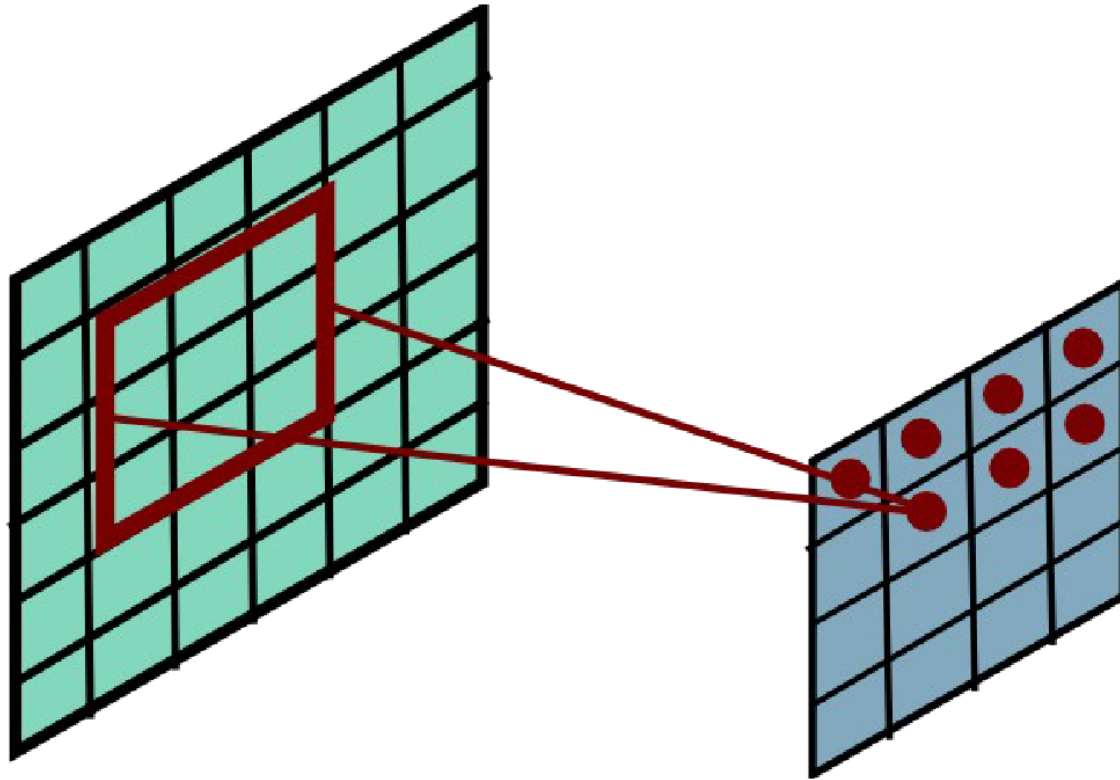
Convolutional Layer



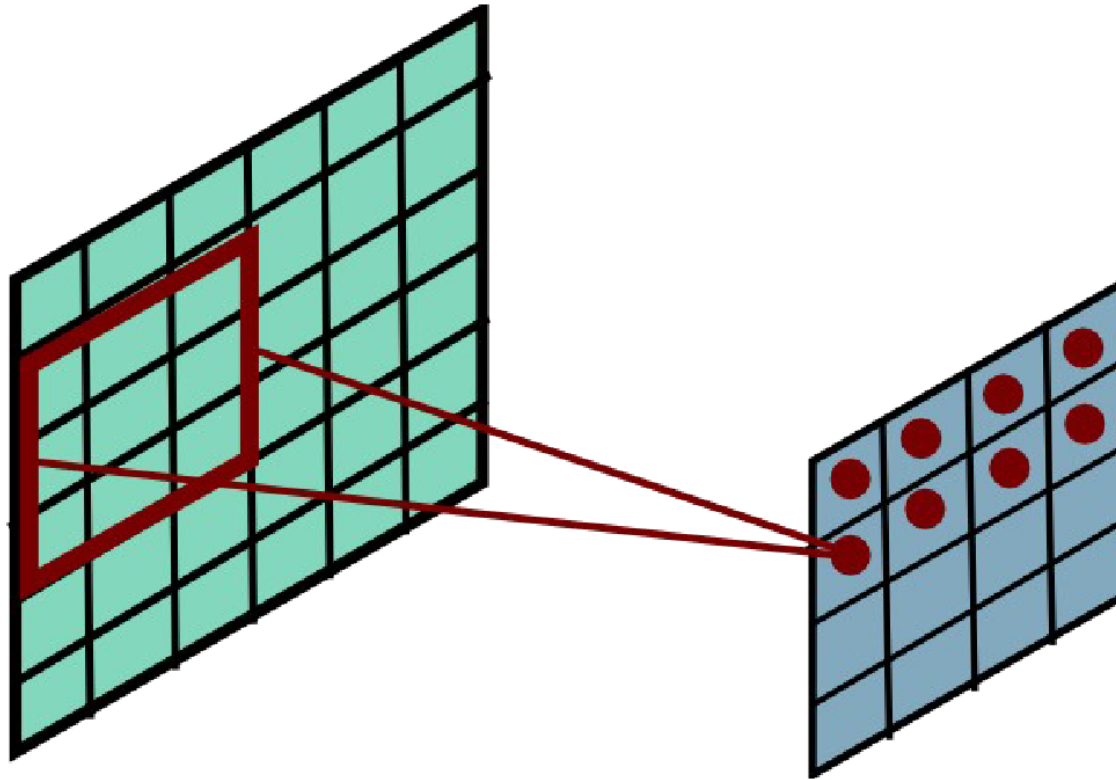
Convolutional Layer



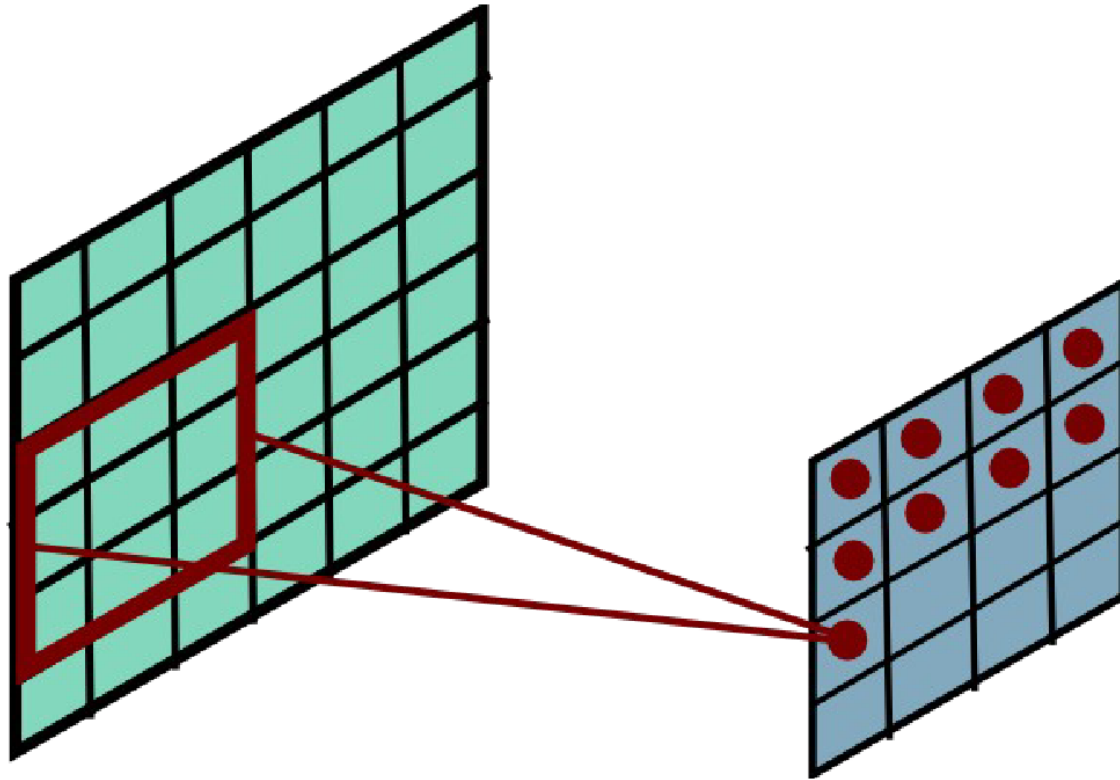
Convolutional Layer



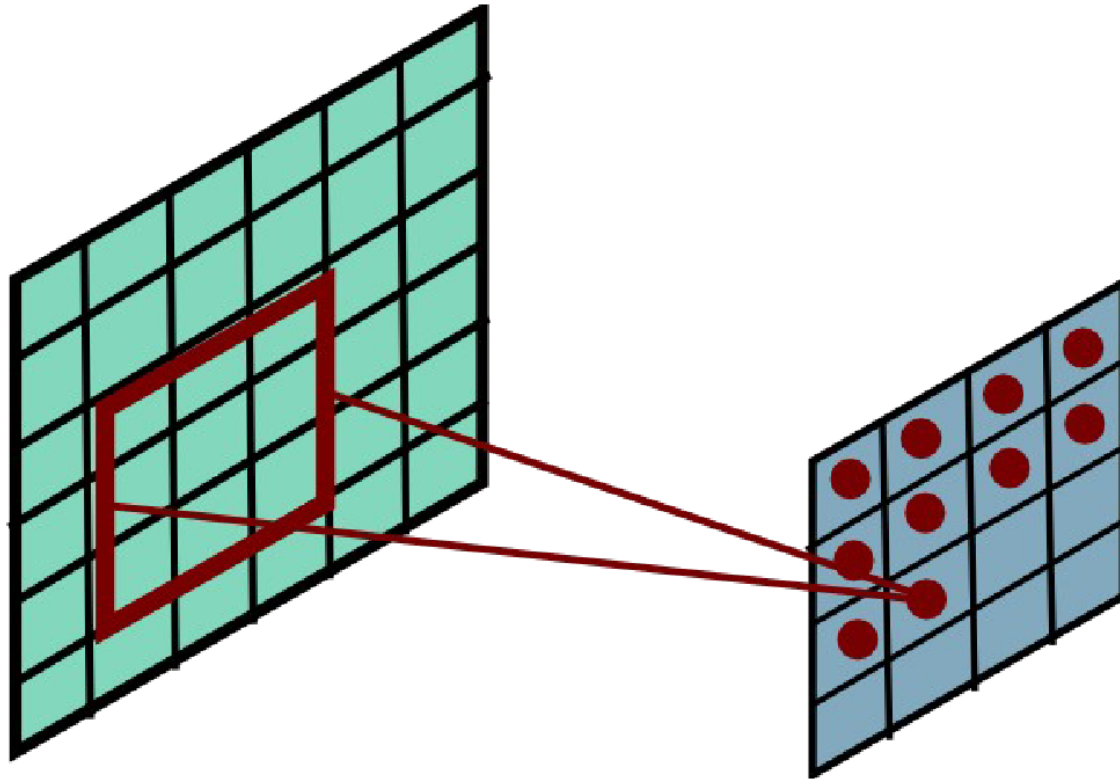
Convolutional Layer



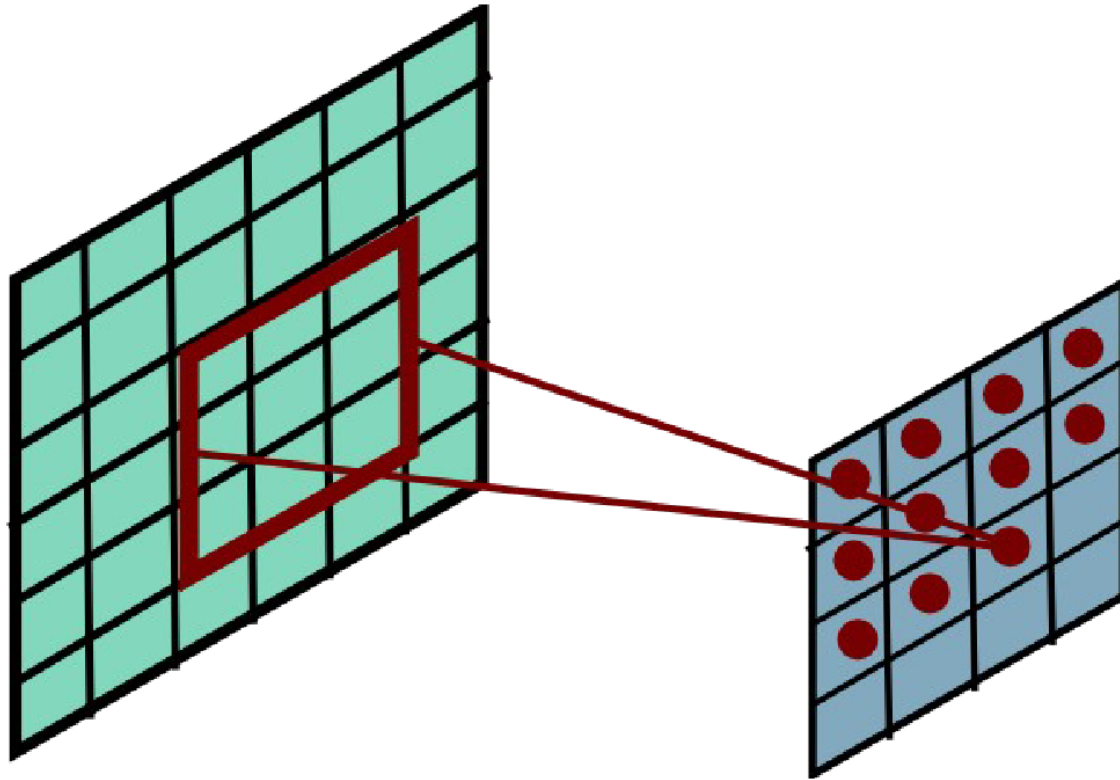
Convolutional Layer



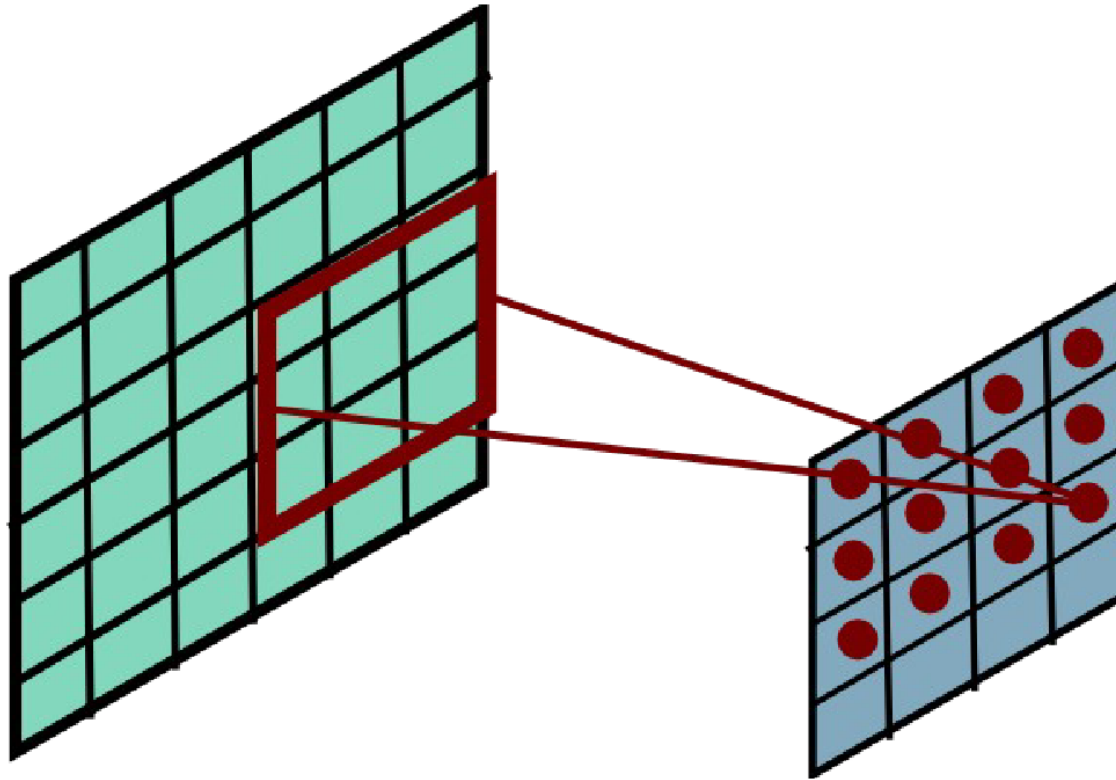
Convolutional Layer



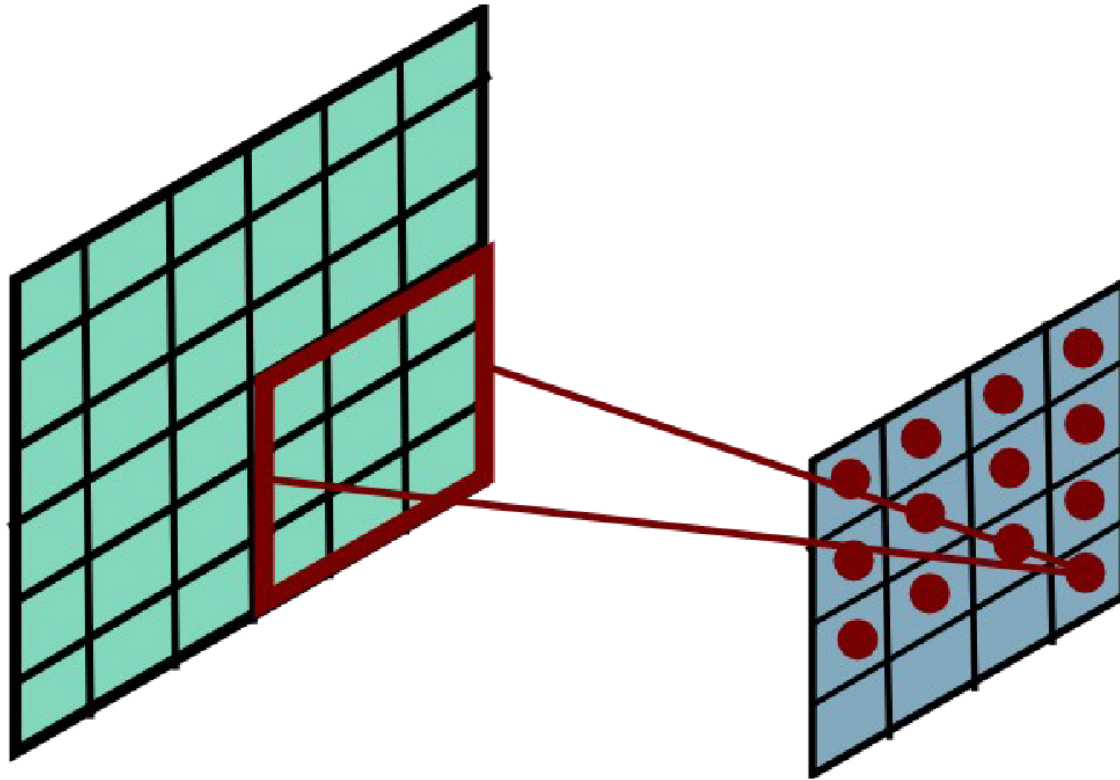
Convolutional Layer



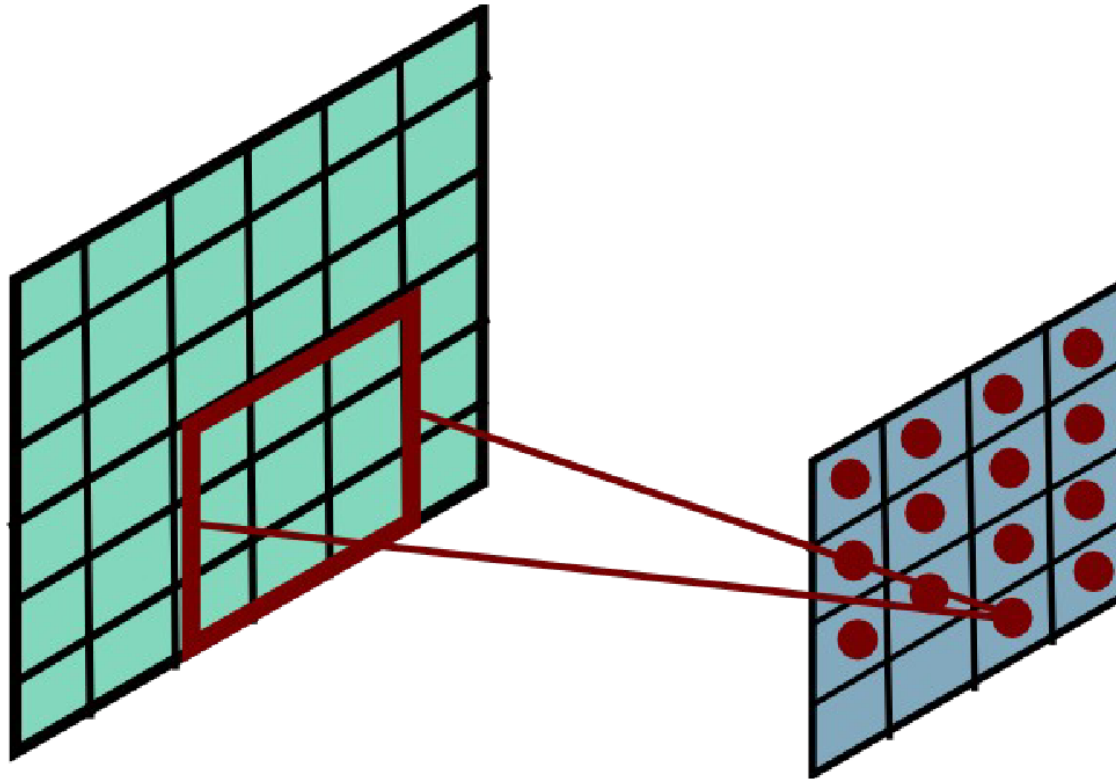
Convolutional Layer



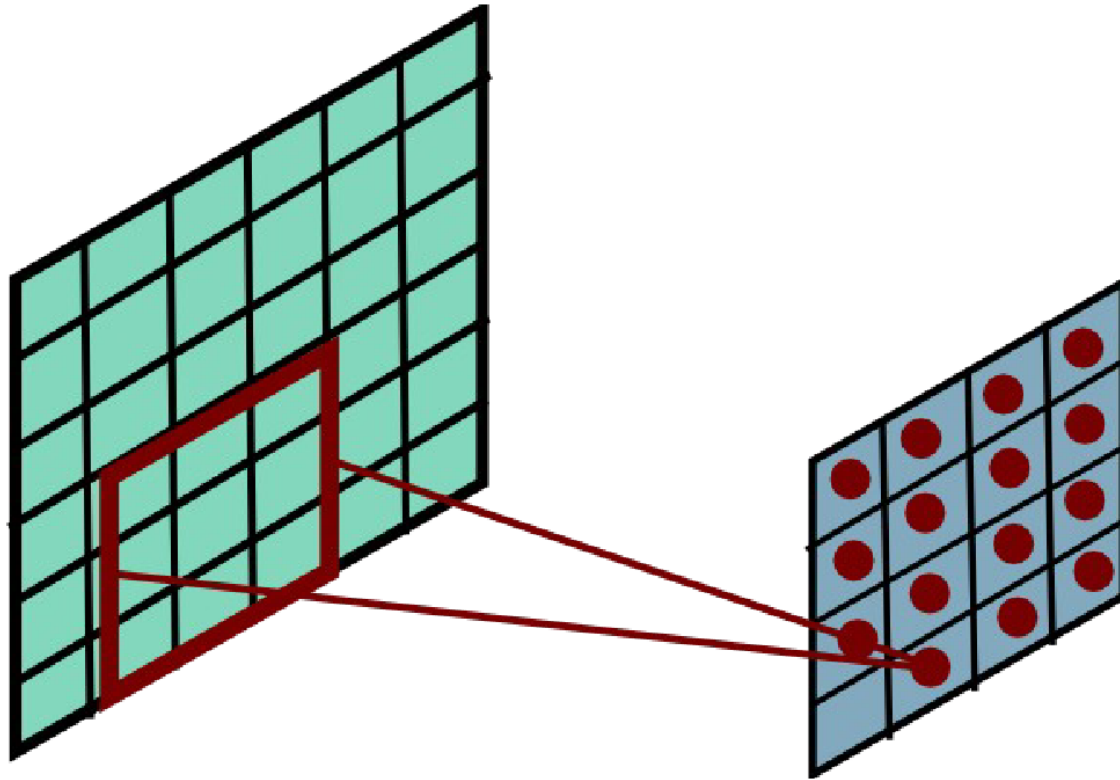
Convolutional Layer



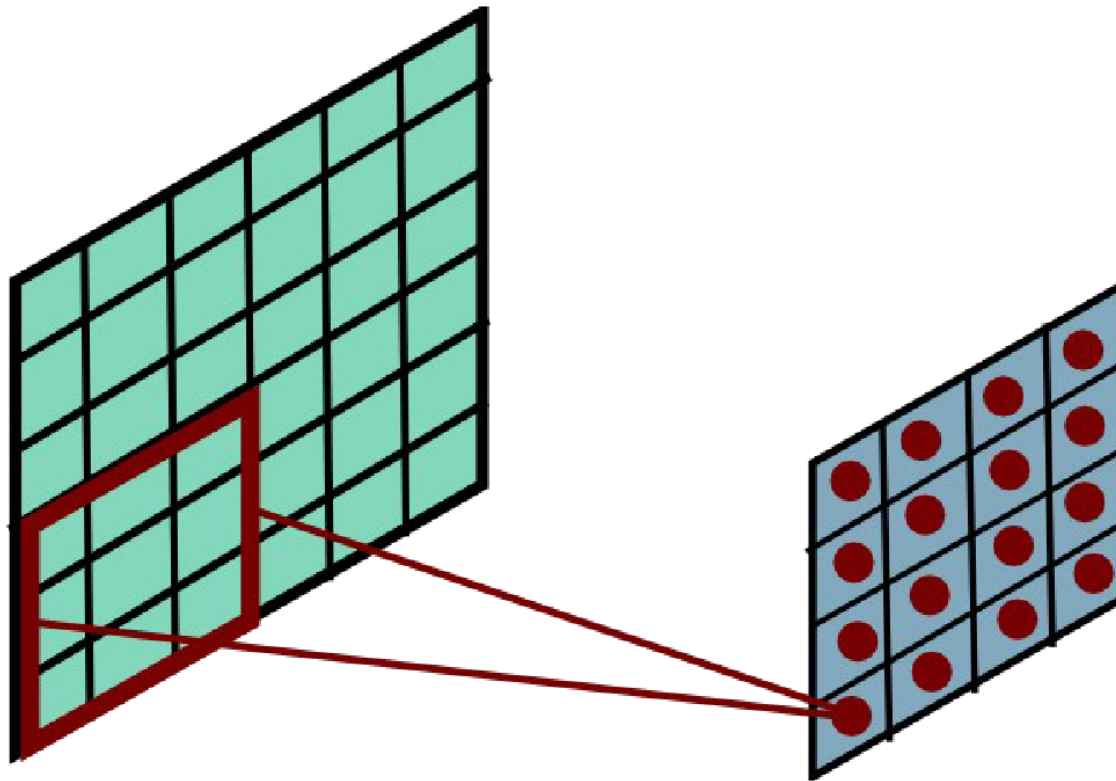
Convolutional Layer



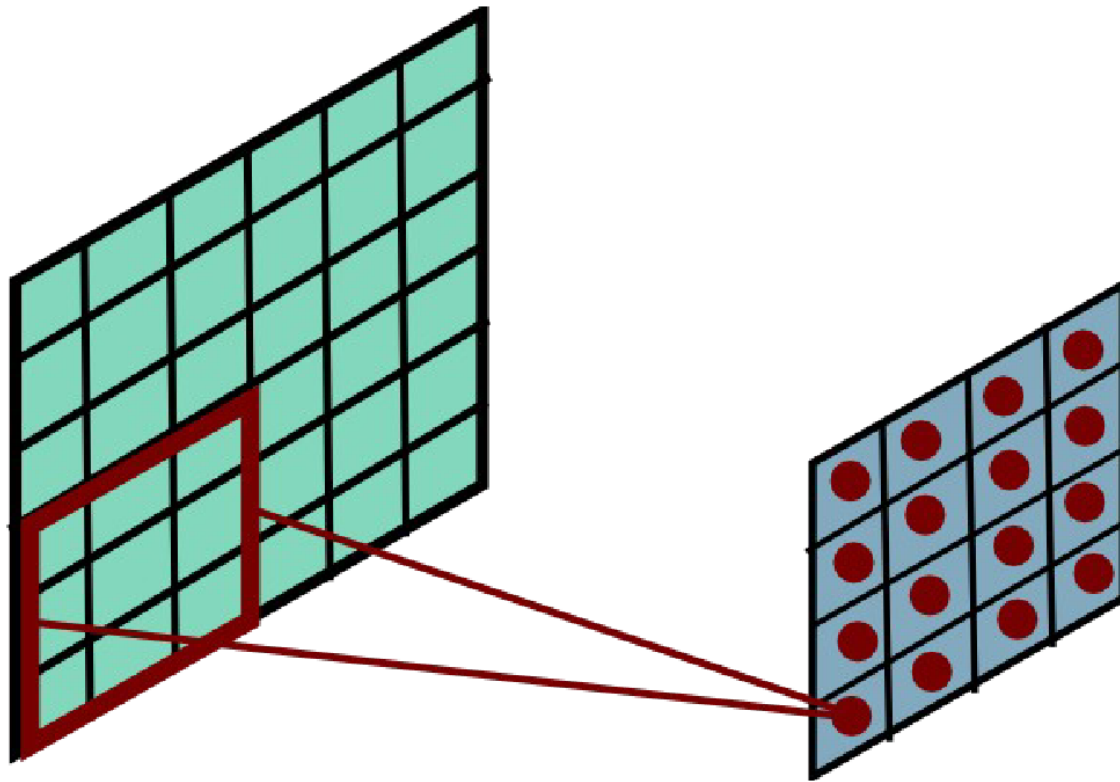
Convolutional Layer



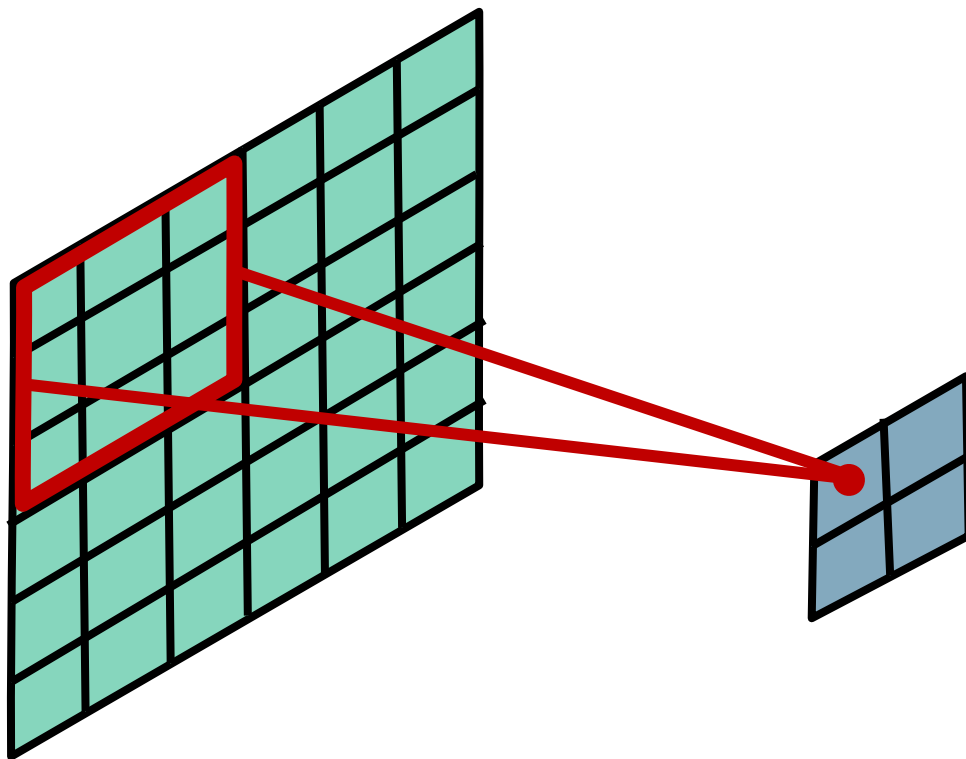
Convolutional Layer



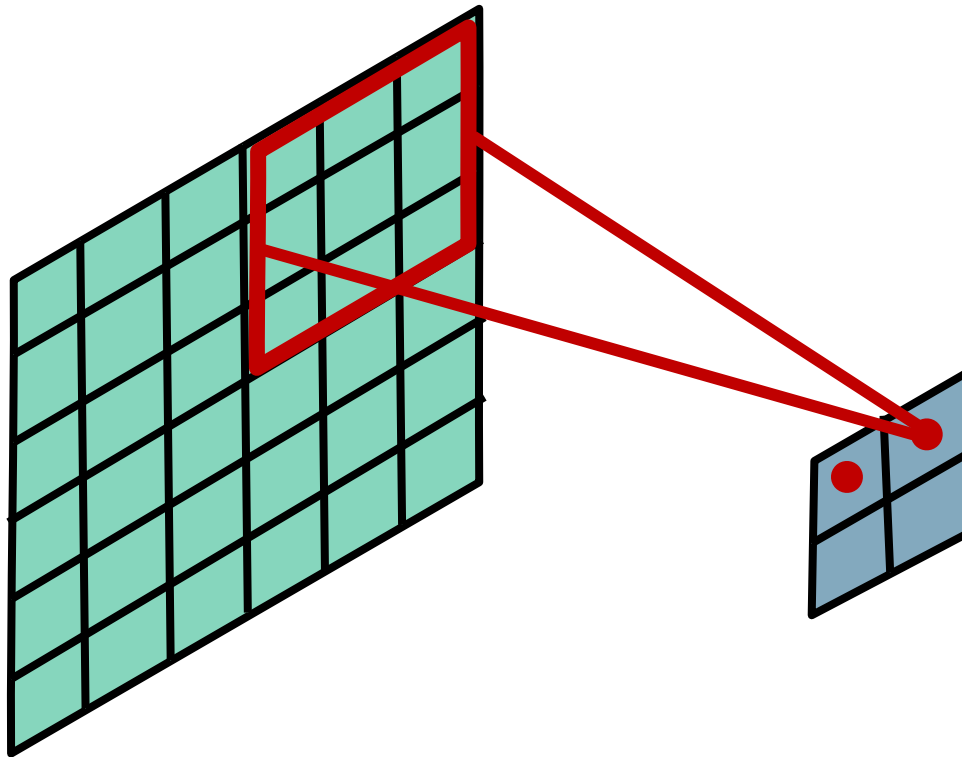
Stride = 3



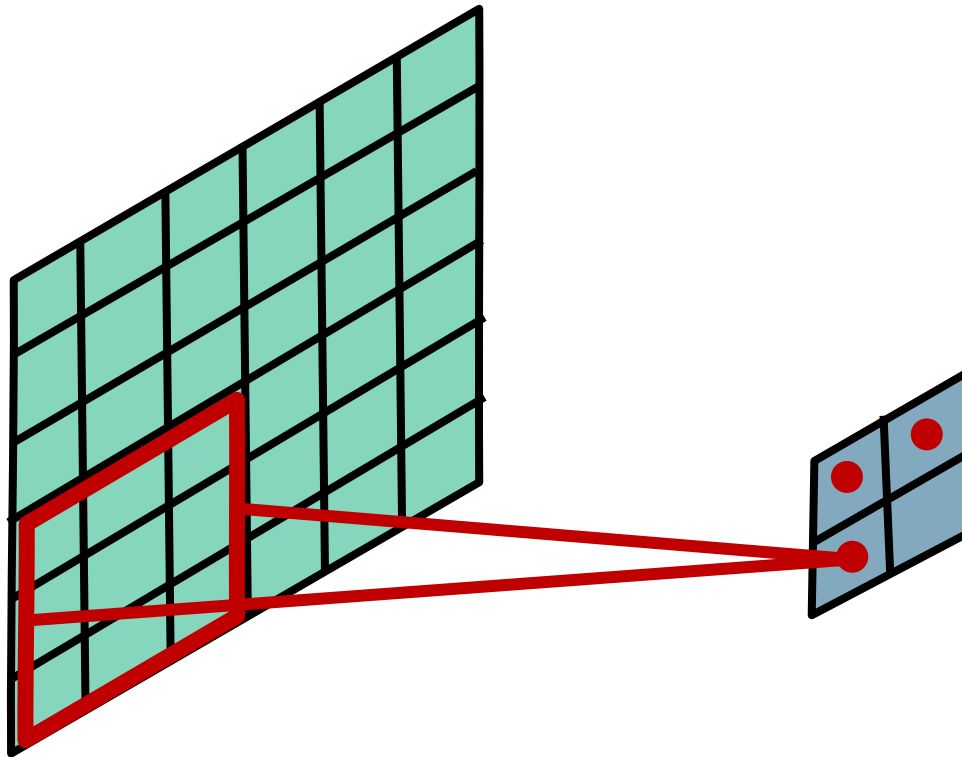
Stride = 3



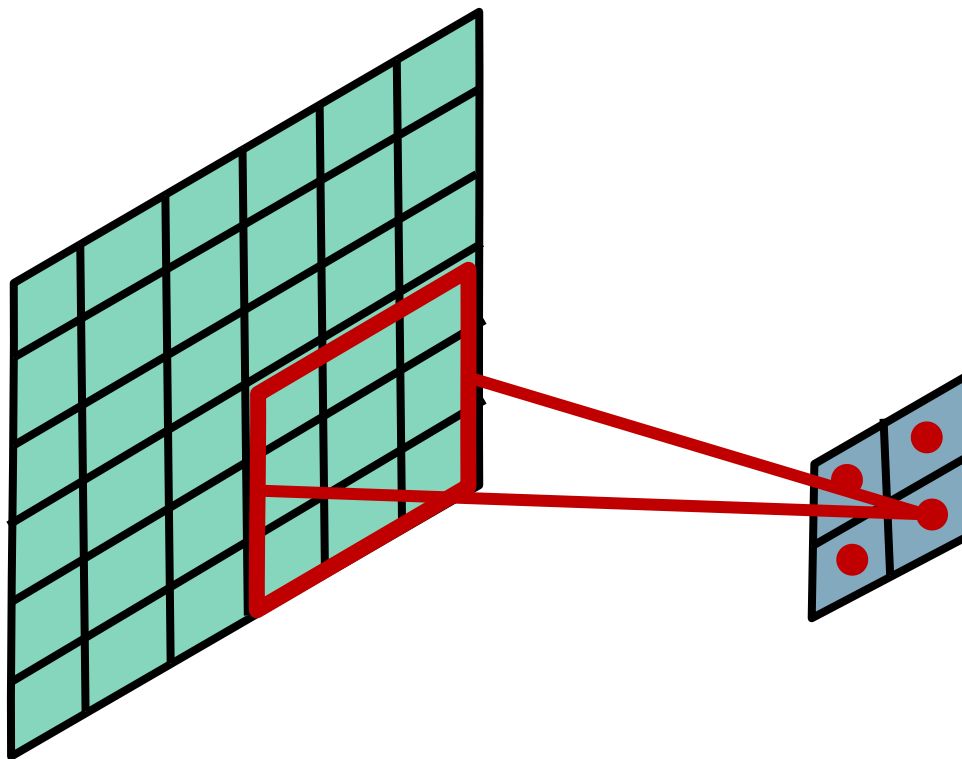
Stride = 3



Stride = 3

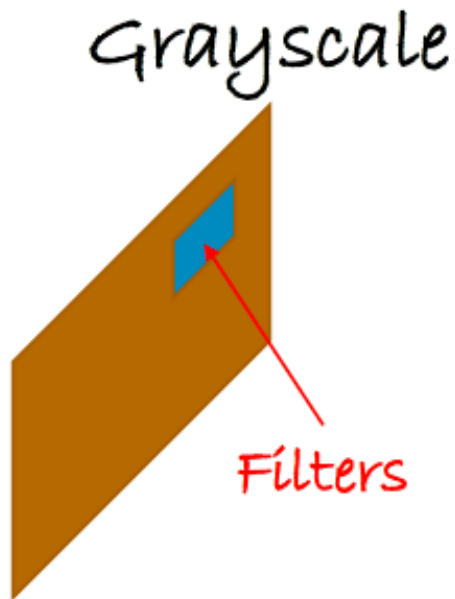


Stride = 3



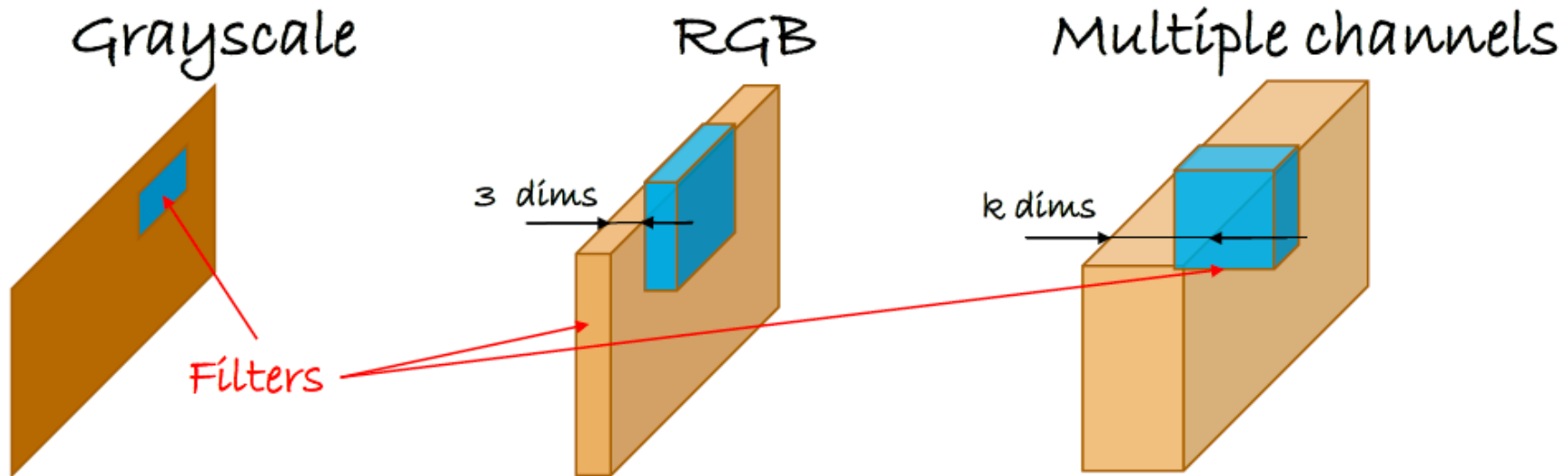
2D spatial filters

- If images are 2-D, parameters should also be organized in 2-D
 - That way they can learn the local correlations between input variables
 - That way they can “exploit” the spatial nature of images

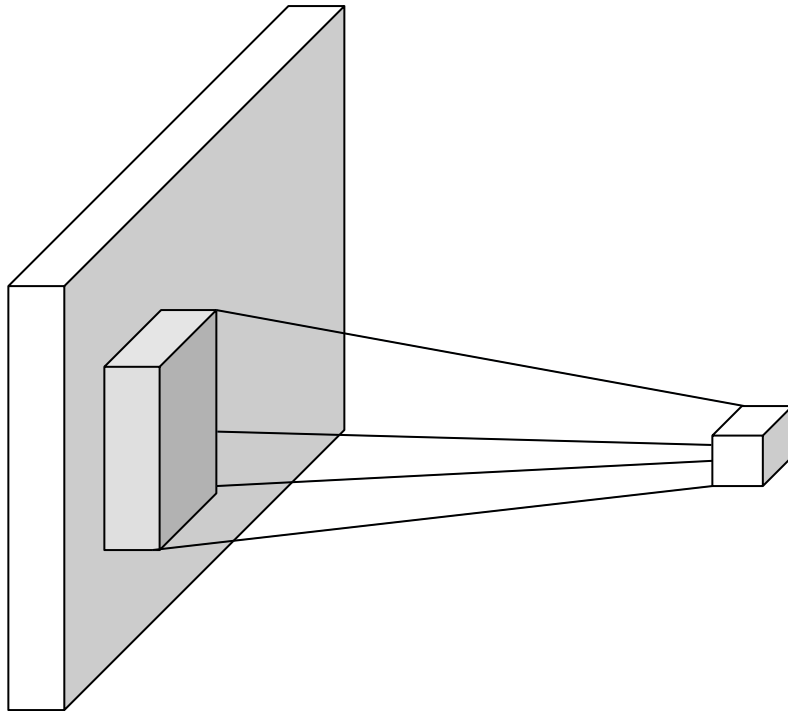


k-D spatial filters

- Similarly, if images are k-D, parameters should also be k-D



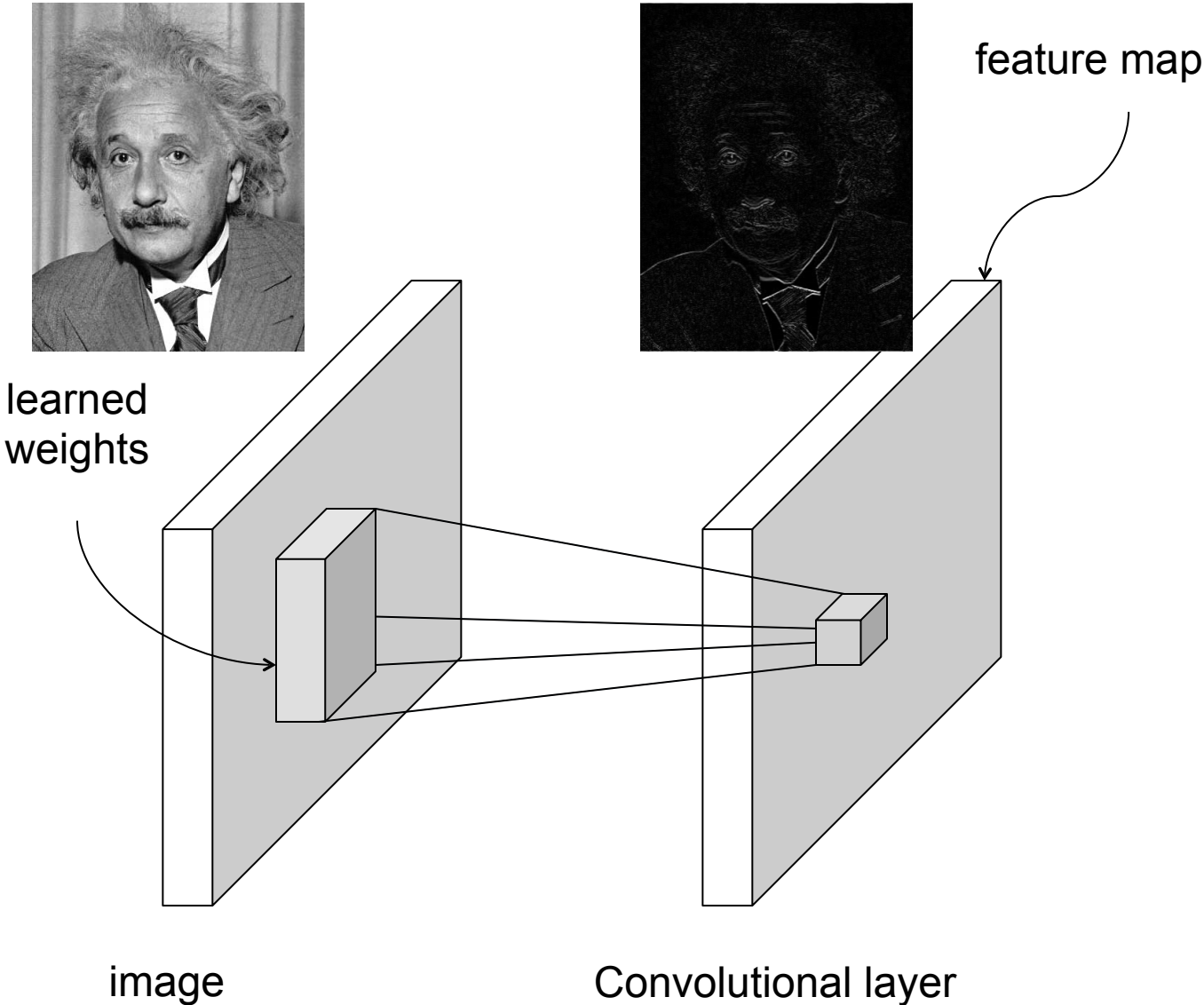
Dimensions of convolution



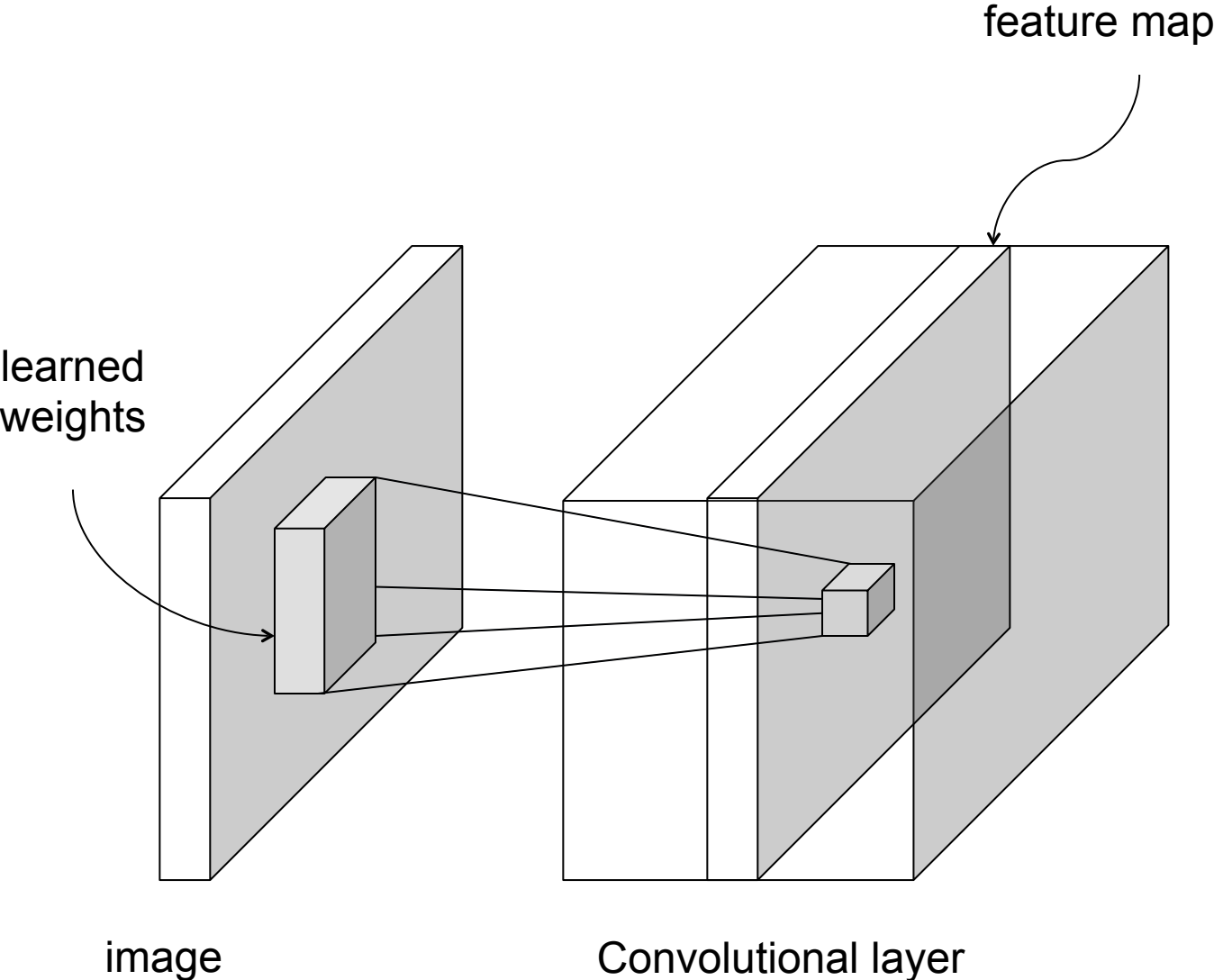
image

Convolutional layer

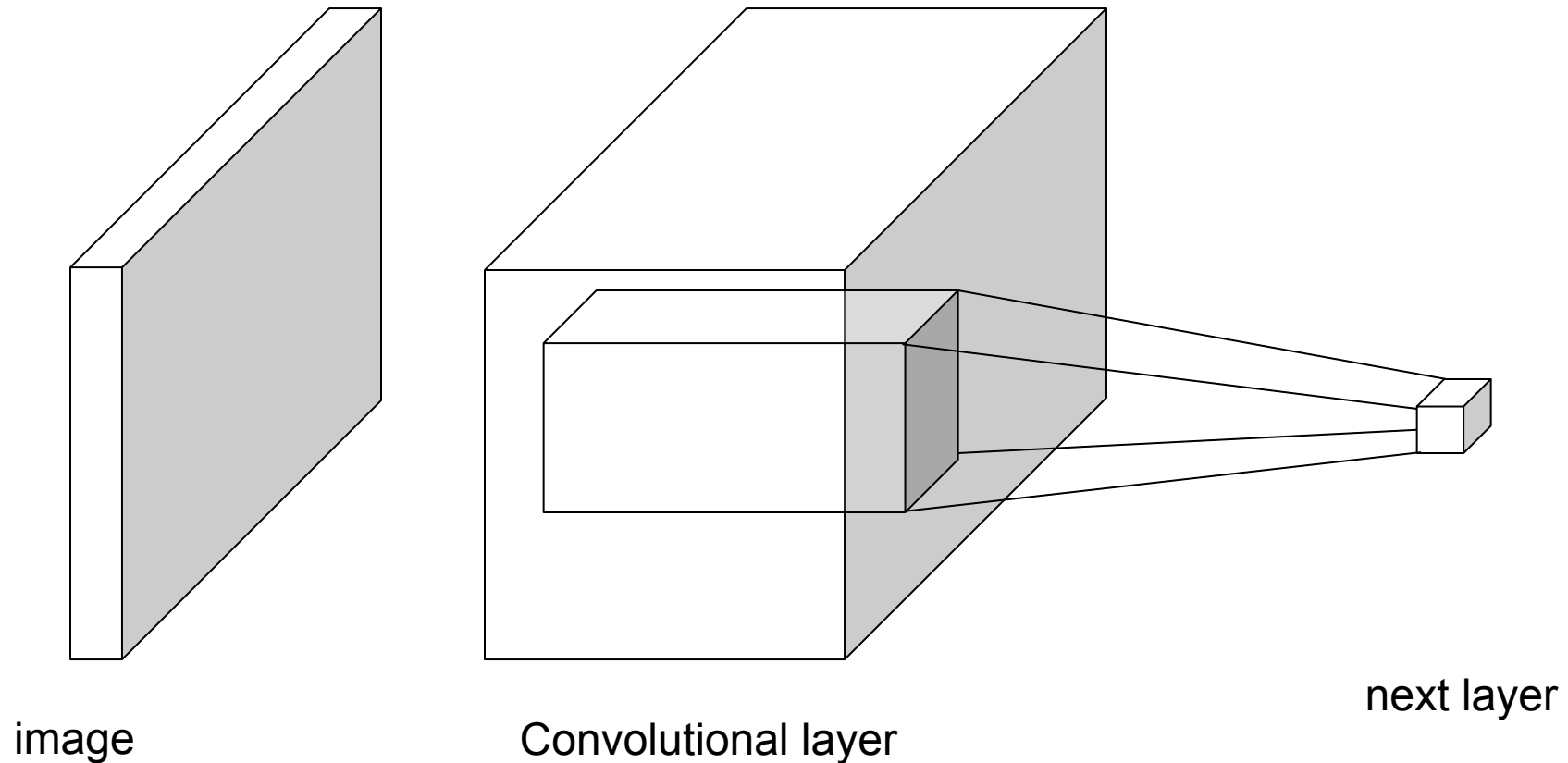
Dimensions of convolution



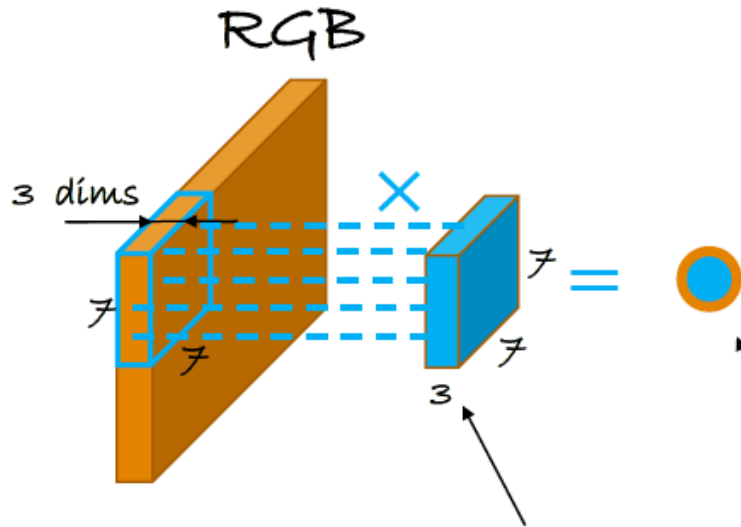
Dimensions of convolution



Dimensions of convolution



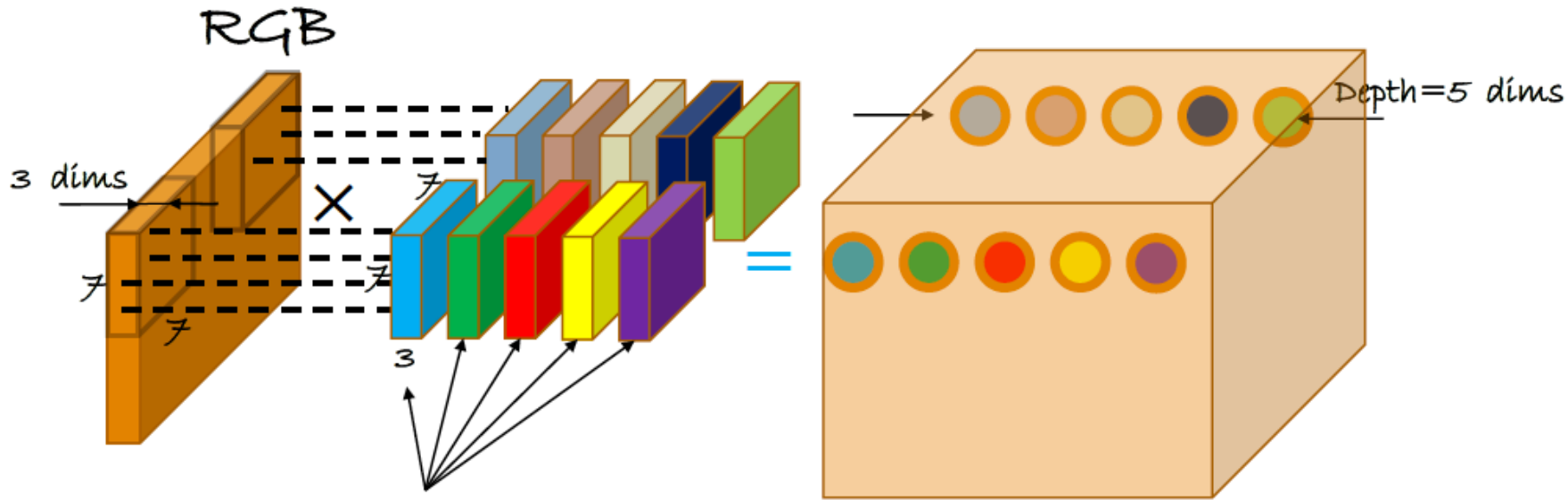
Number of weights



How many weights for this neuron?

$$7 \cdot 7 \cdot 3 = 147$$

Number of weights

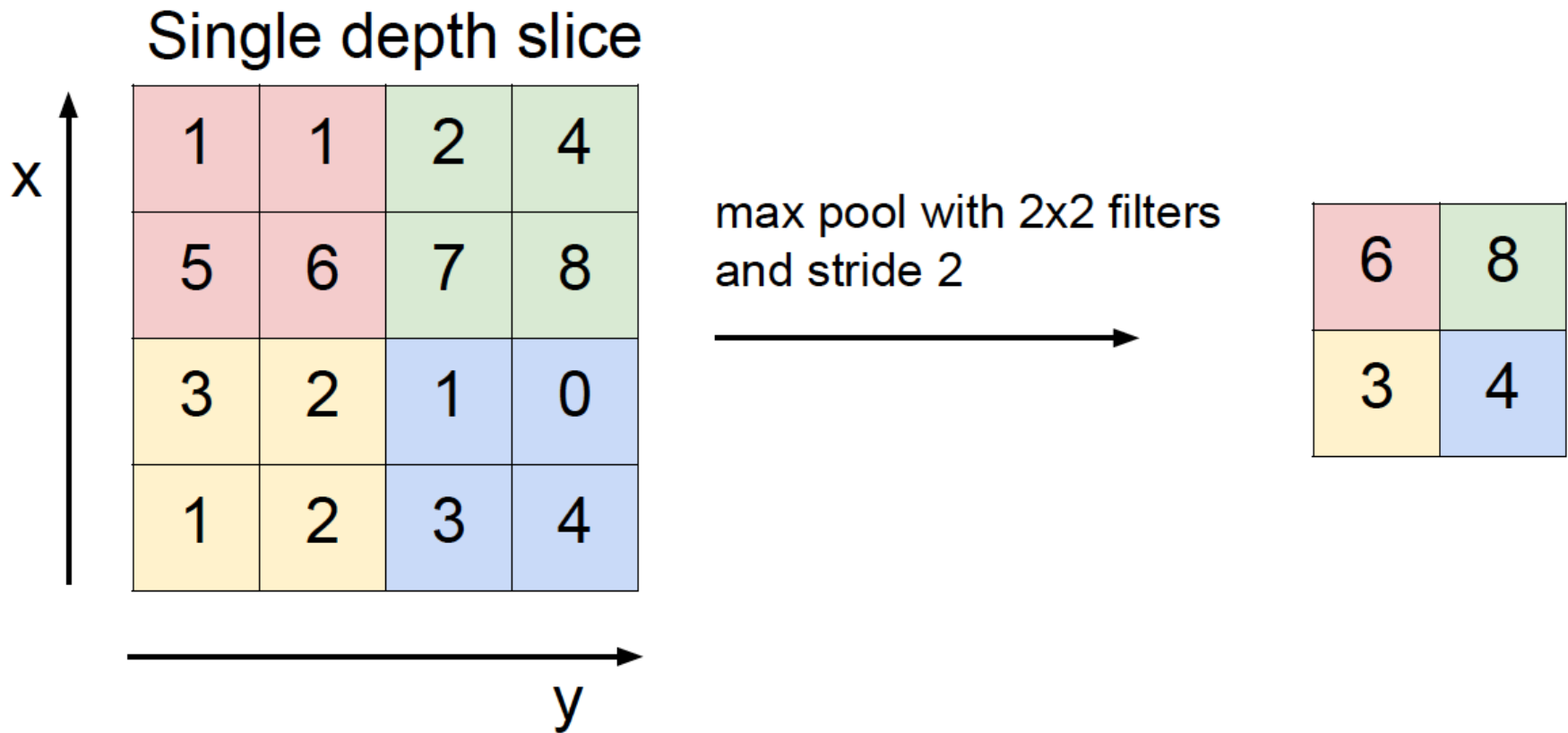


How many weights for these 5 neurons?

$$5 \cdot 7 \cdot 7 \cdot 3 = 735$$

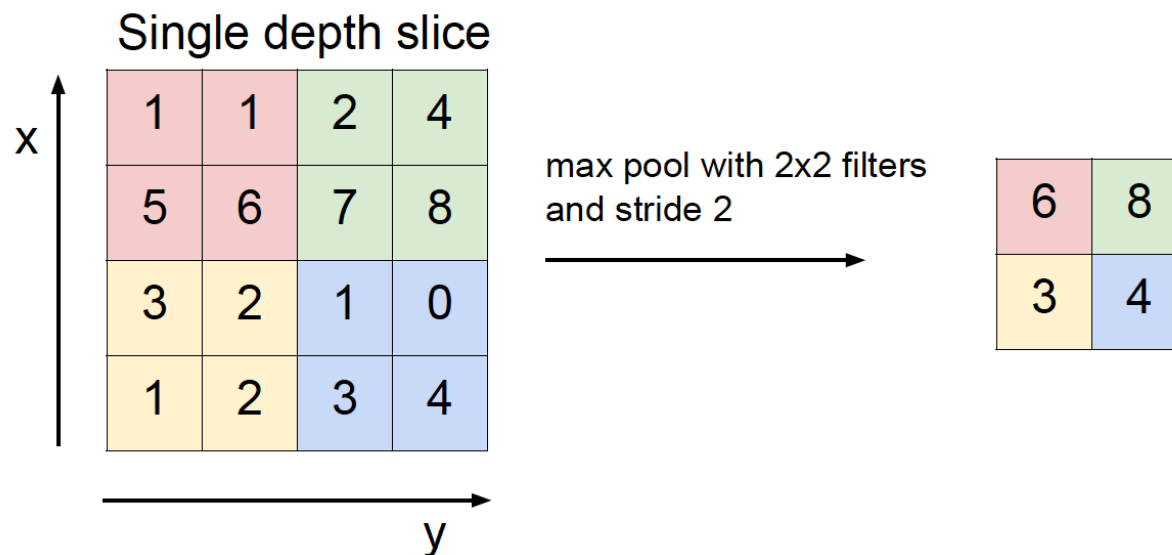
Pooling: Downsample feature maps

- Aggregate multiple values into a single value

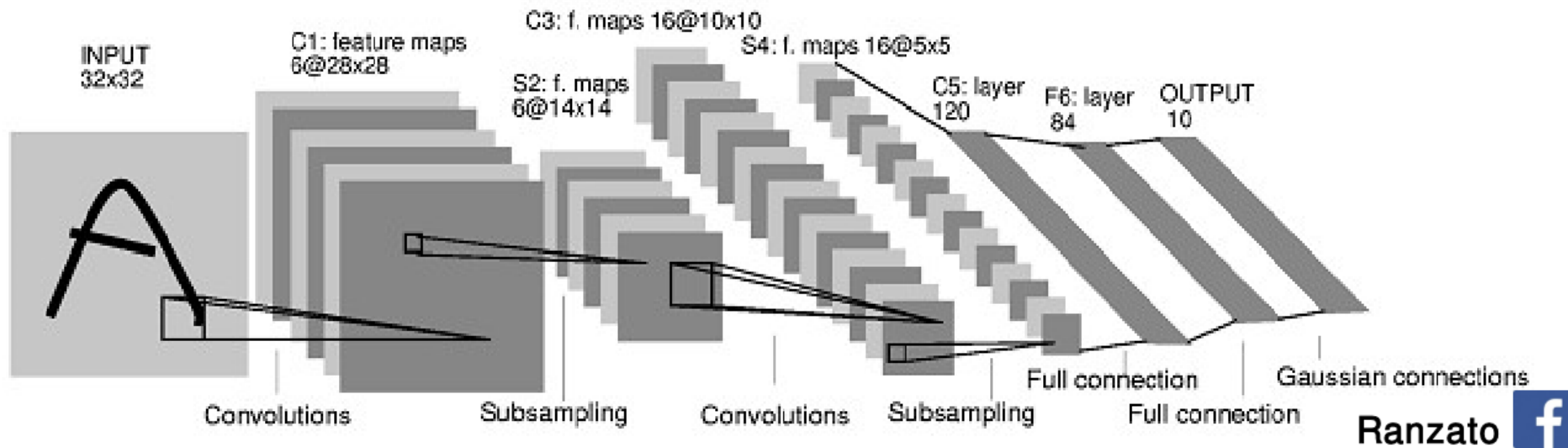


Pooling: Downsample feature maps

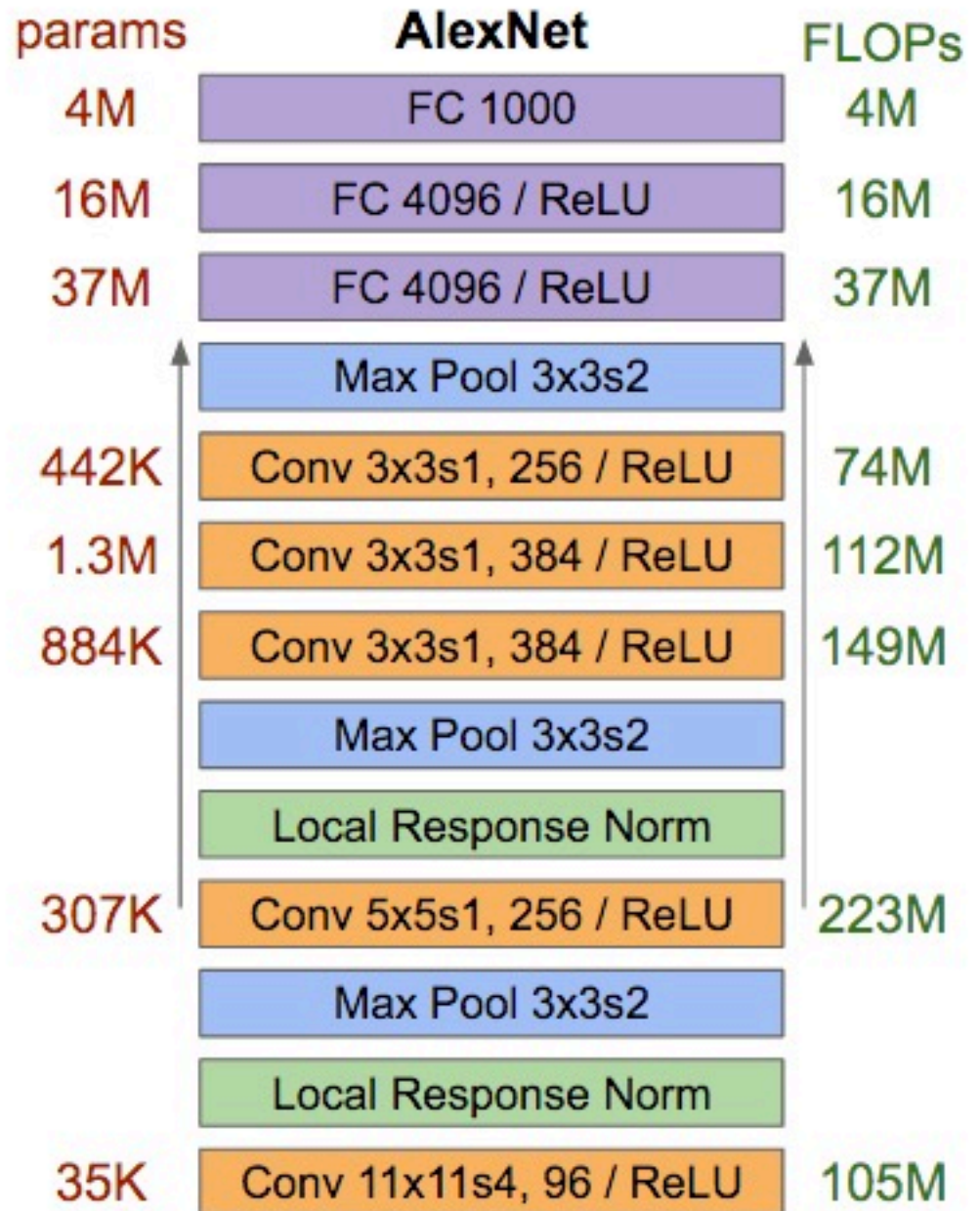
- Aggregate multiple values into a single value
- Invariance to small transformations
 - Keep only most important information for next layer
- Reduces the size of the next layer
 - Fewer parameters, faster computations
- Observe larger receptive field in next layer
 - Hierarchically extract more abstract features



Yann LeCun's MNIST CNN architecture



AlexNet for ImageNet

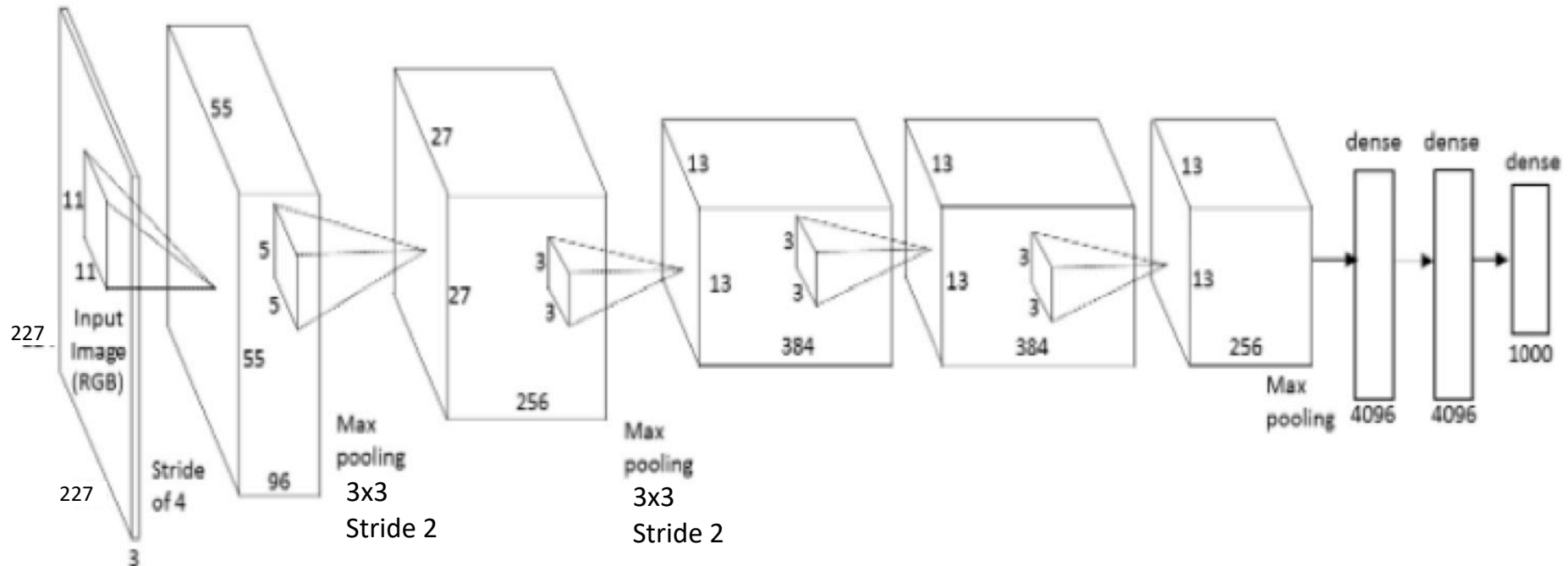


Layers

- Kernel sizes
- Strides
- # channels
- # kernels
- Max pooling

AlexNet diagram (simplified)

Input size
227 x 227 x 3



Conv 1

11 x 11 x 3

Stride 4

96 filters

Conv 2

5 x 5 x 48

Stride 1

256 filters

Conv 3

3 x 3 x 256

Stride 1

384 filters

Conv 4

3 x 3 x 192

Stride 1

384 filters

Conv 5

3 x 3 x 192

Stride 1

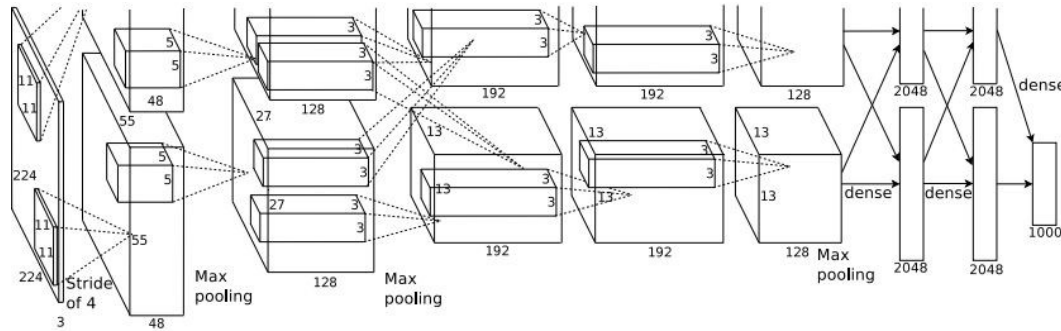
256 filters

Convolutional Neural Networks

- Question: Spatial structure?
 - Answer: Convolutional filters
- Question: Huge input dimensionalities?
 - Answer: Parameters are shared between filters
- Question: Local variances?
 - Answer: Pooling

What's going on inside ConvNets?

[This image is CC0 public domain](#)



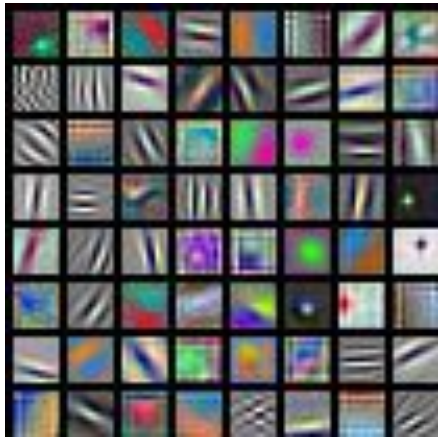
Class Scores:
1000 numbers

Input Image:
3 x 224 x 224



What are the intermediate features
looking for?

First Layer: Visualize Filters



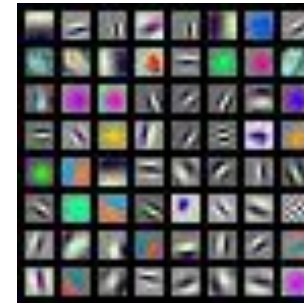
AlexNet:
64 x 3 x 11 x
11



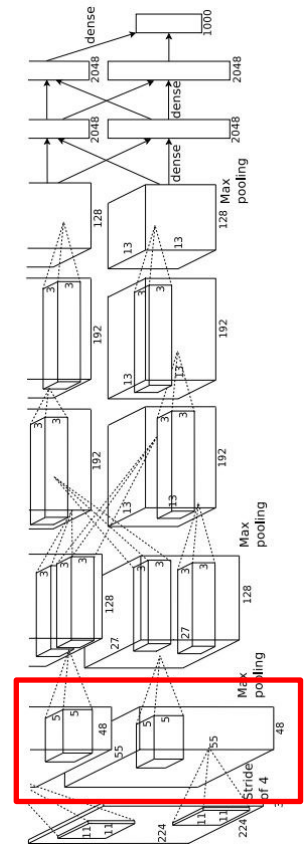
ResNet-
18: 64 x 3
x 7 x 7



ResNet-1
01: 64 x 3
x 7 x 7



DenseNet-
121: 64 x
3 x 7 x 7



Visualize the filters/kernels (raw weights)

We can visualize filters at higher layers, but not that interesting

(these are taken from ConvNetJS CIFAR-10 demo)

Weights:


layer 1 weights

16 x 3 x 7 x 7

Weights:


layer 2 weights

20 x 16 x 7 x 7

Weights:


layer 3 weights

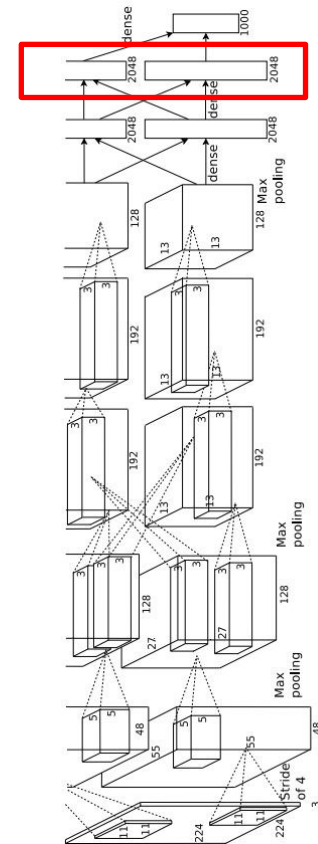
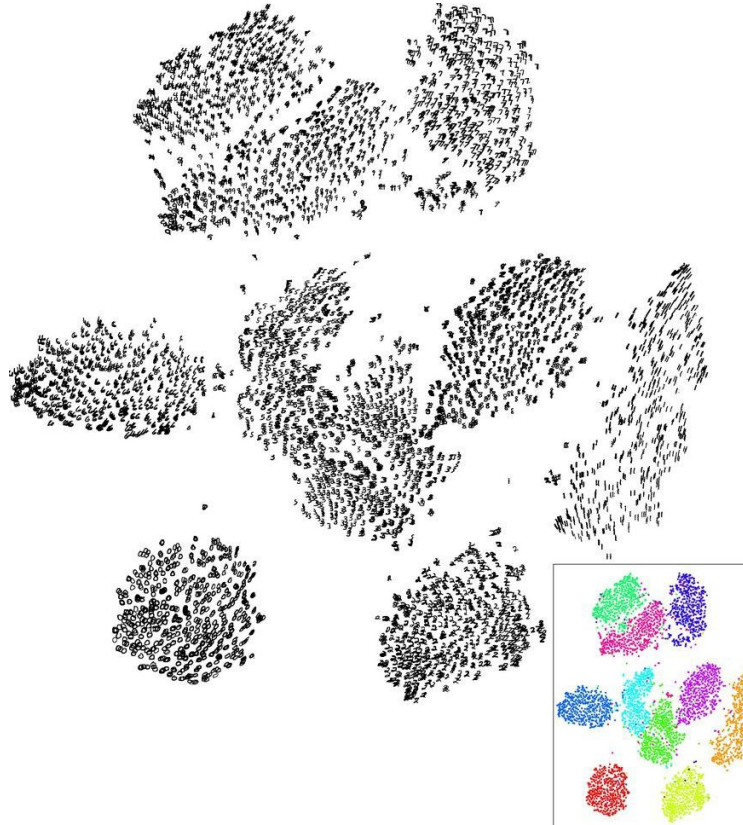
20 x 20 x 7 x 7

Last Layer: Dimensionality Reduction

Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

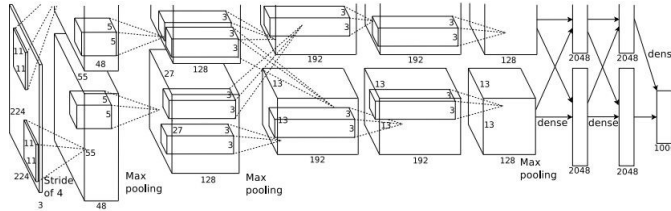
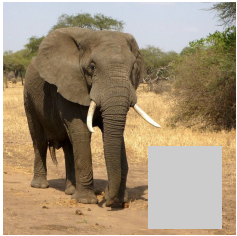
Simple algorithm: Principal Component Analysis (PCA)

More complex: **t-SNE**

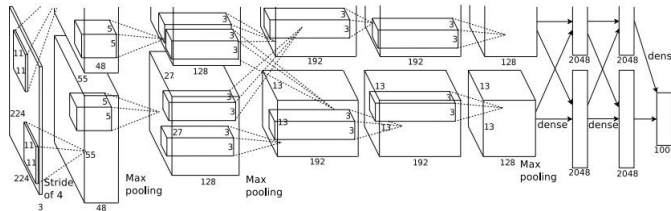
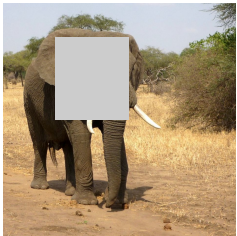


Which pixels matter: Saliency vs Occlusion

Mask part of the image before feeding to CNN, check how much predicted probabilities change



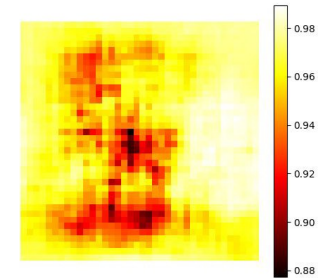
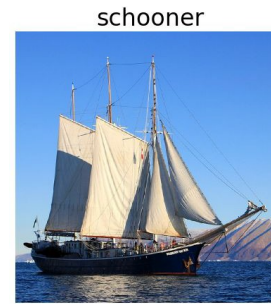
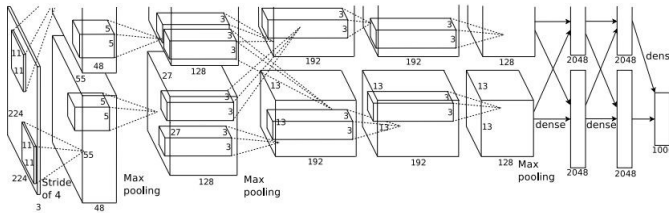
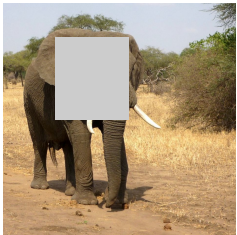
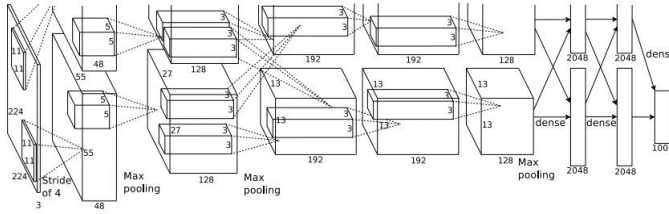
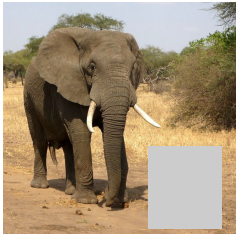
$P(\text{elephant}) = 0.95$



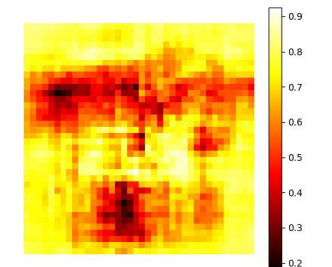
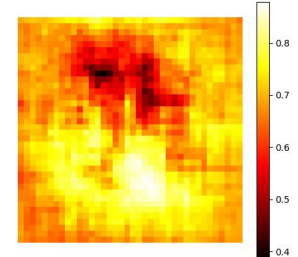
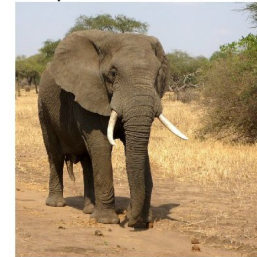
$P(\text{elephant}) = 0.75$

Which pixels matter: Saliency vs Occlusion

Mask part of the image before feeding to CNN, check how much predicted probabilities change

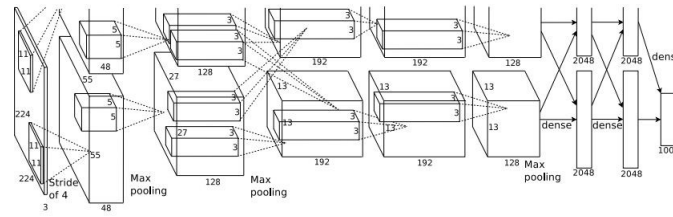


African elephant, *Loxodonta africana*



Which pixels matter: Saliency via Backprop

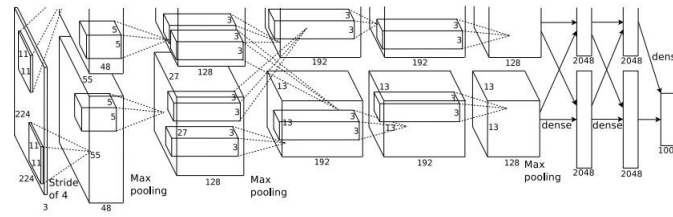
Forward pass: Compute probabilities



Dog

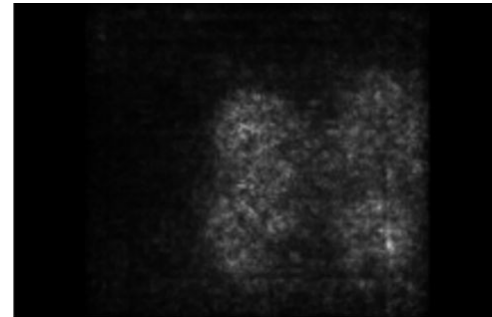
Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



Dog

Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels



Saliency Maps



Fooling Images / Adversarial Examples

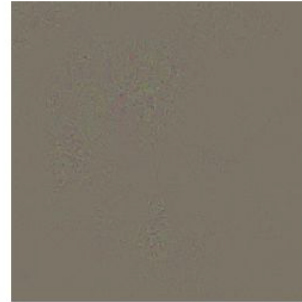
African elephant



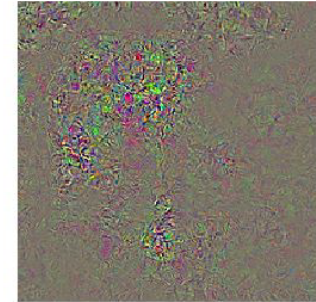
koala



Difference



10x Difference



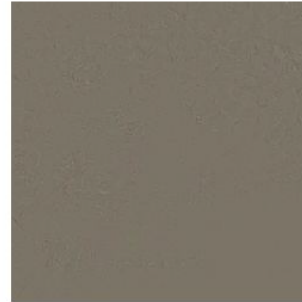
schooner



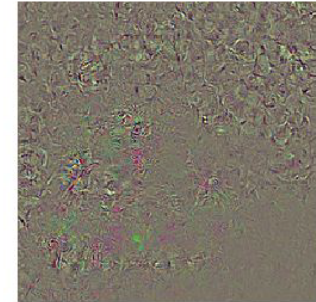
iPod



Difference



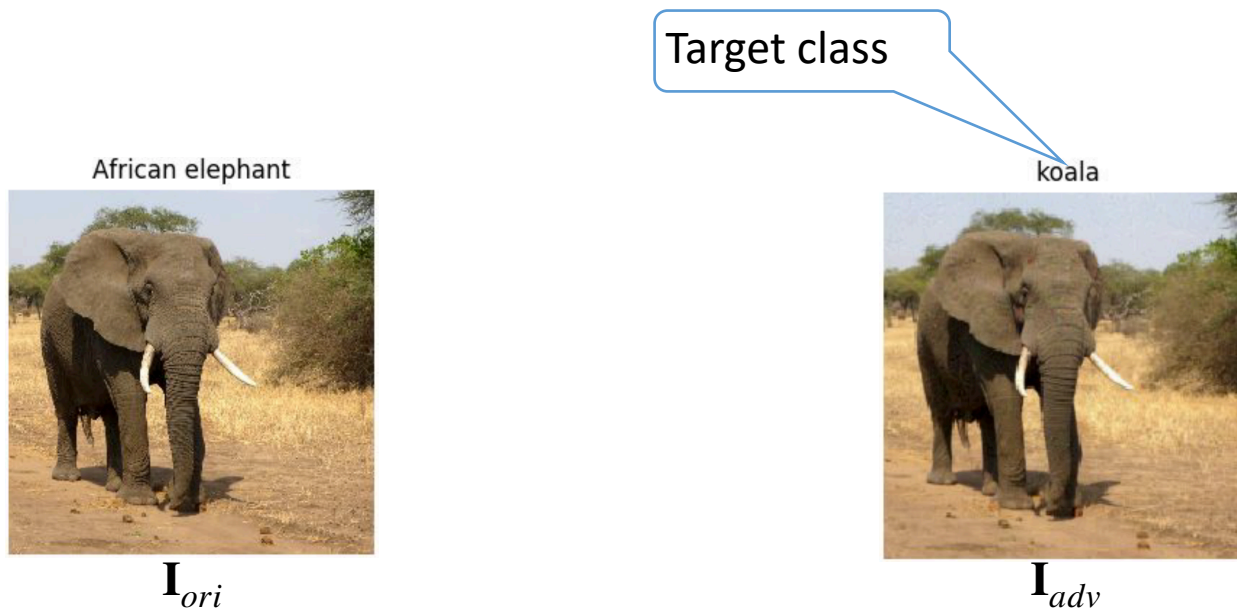
10x Difference



Fooling Images / Adversarial Examples

- (1) Start from an arbitrary image
- (2) Pick an arbitrary class
- (3) Modify the image to maximize the class
- (4) Repeat until network is fooled

Optimization Formulation



$Score_c(\mathbf{I}; \theta)$: the confidence score of an image belonging to class c , using a network of parameters θ

Attack: Modify the image \mathbf{I} to increase $Score_{target\ class}(\mathbf{I}; \theta)$

$$\begin{array}{ll} \underset{\mathbf{I}_{adv}}{\text{maximize}} & Score_{target\ class}(\mathbf{I}_{adv}) \\ \text{subject to} & \|\mathbf{I}_{adv} - \mathbf{I}_{ori}\| \leq \epsilon \end{array}$$

Gradient-based Attack

Fast Gradient Sign Method:

$$\mathbf{I}^{adv} = \mathbf{I} + \epsilon \text{sign}(\nabla_{\mathbf{I}} \text{Score}_{target\ class}(\mathbf{I}))$$



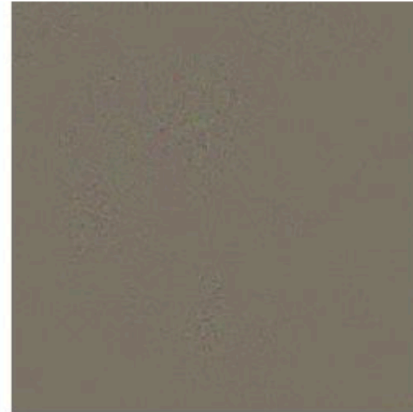
African elephant



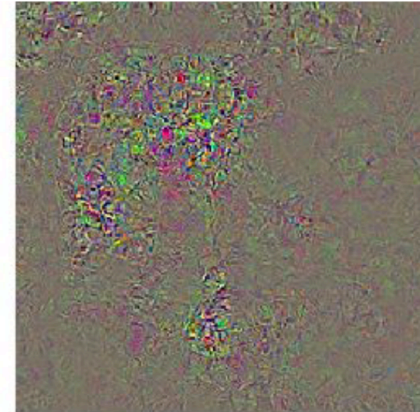
koala



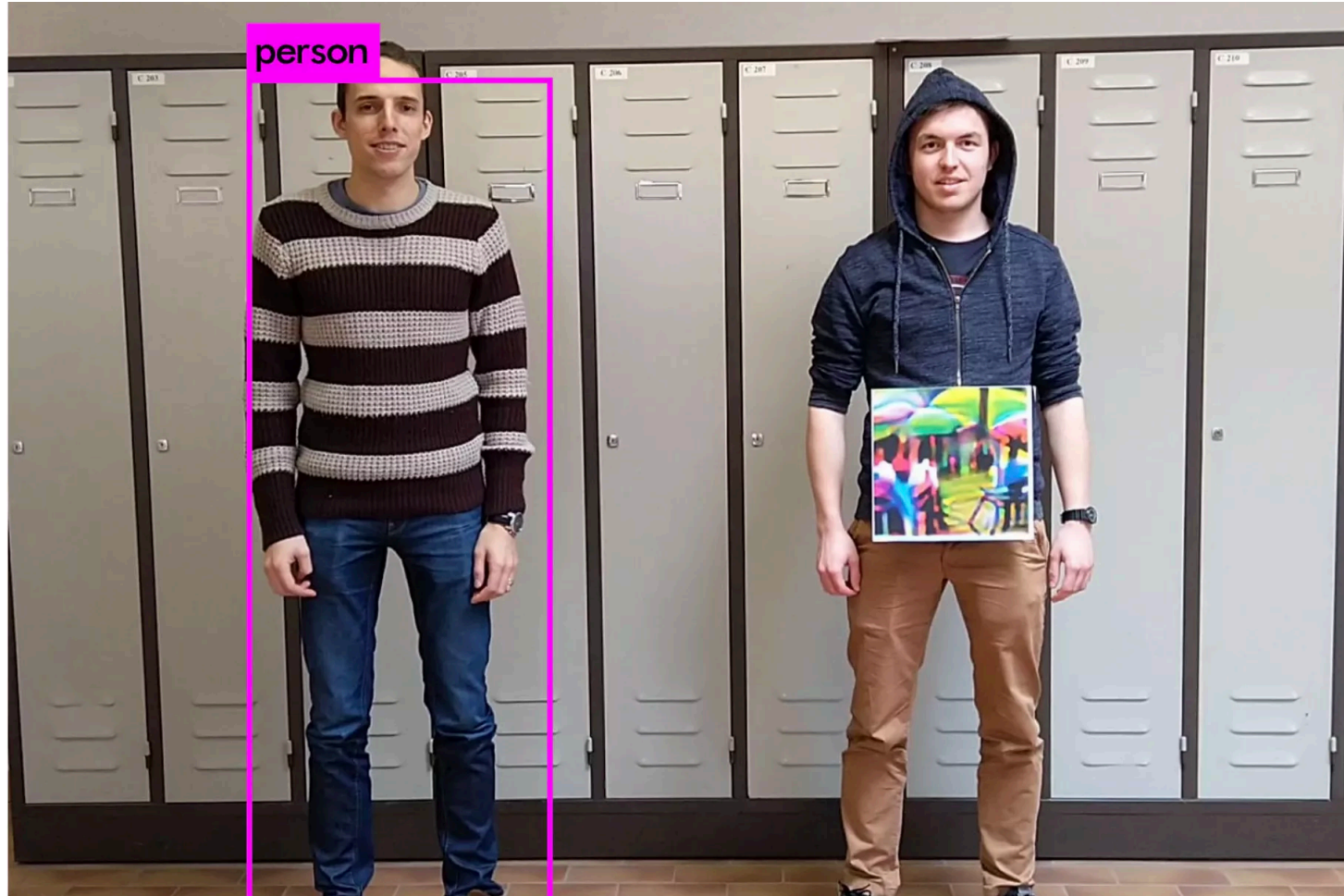
Difference



10x Difference



Patch-based Attack (Spatially Localized)



Dangerous!



(a)



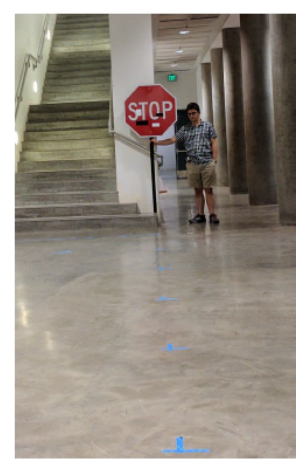
(b)



(c)



(d)



(e)

Neural Style Transfer

Content Image



This image is licensed under [CC-BY3.0](#)

+

Style Image



Starry Night by Van Gogh is in the public domain

Neural Style Transfer

Content Image



[This image](#) is licensed under [CC-BY3.0](#)

+

Style Image



[Starry Night](#) by Van Gogh is in the public domain

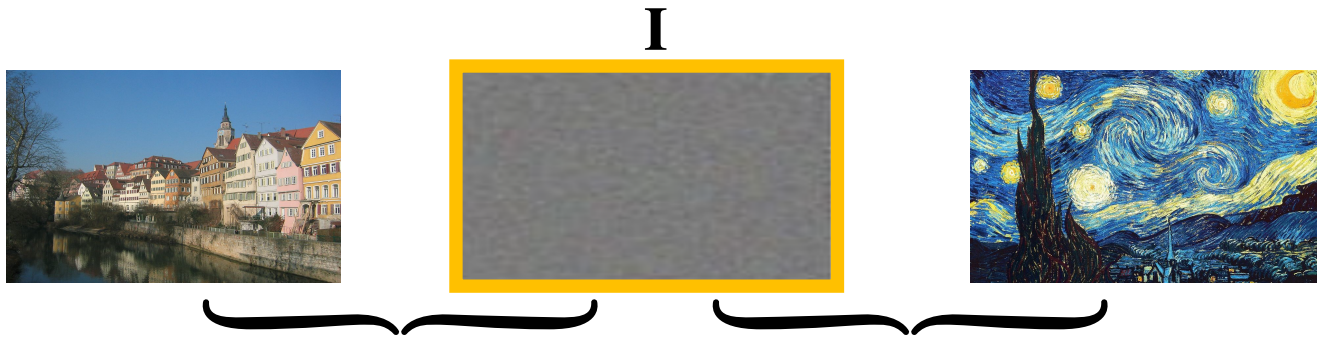
=

Style Transfer!



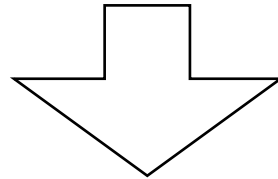
[This image](#) copyright Justin Johnson, 2015. Reproduced with permission.

Total Loss



$$\mathcal{L}_{total}(\mathbf{I}) = \alpha \mathcal{L}_{content}(\mathbf{I}) + \beta \mathcal{L}_{style}(\mathbf{I})$$

Minimize total loss



Neural Style Transfer



Example outputs

Neural Style Transfer



More weight to
content loss



More weight to
style loss