# Introduction to
# Deep Convolutional Neural Networks

Ye Li (yli192@jhu.edu)
[1]Dept. of Electrical and Computer Engineering, Whiting School of Engineering
[2]Division of Medical Imaging Physics, Dept. of Radiology, School of Medicine
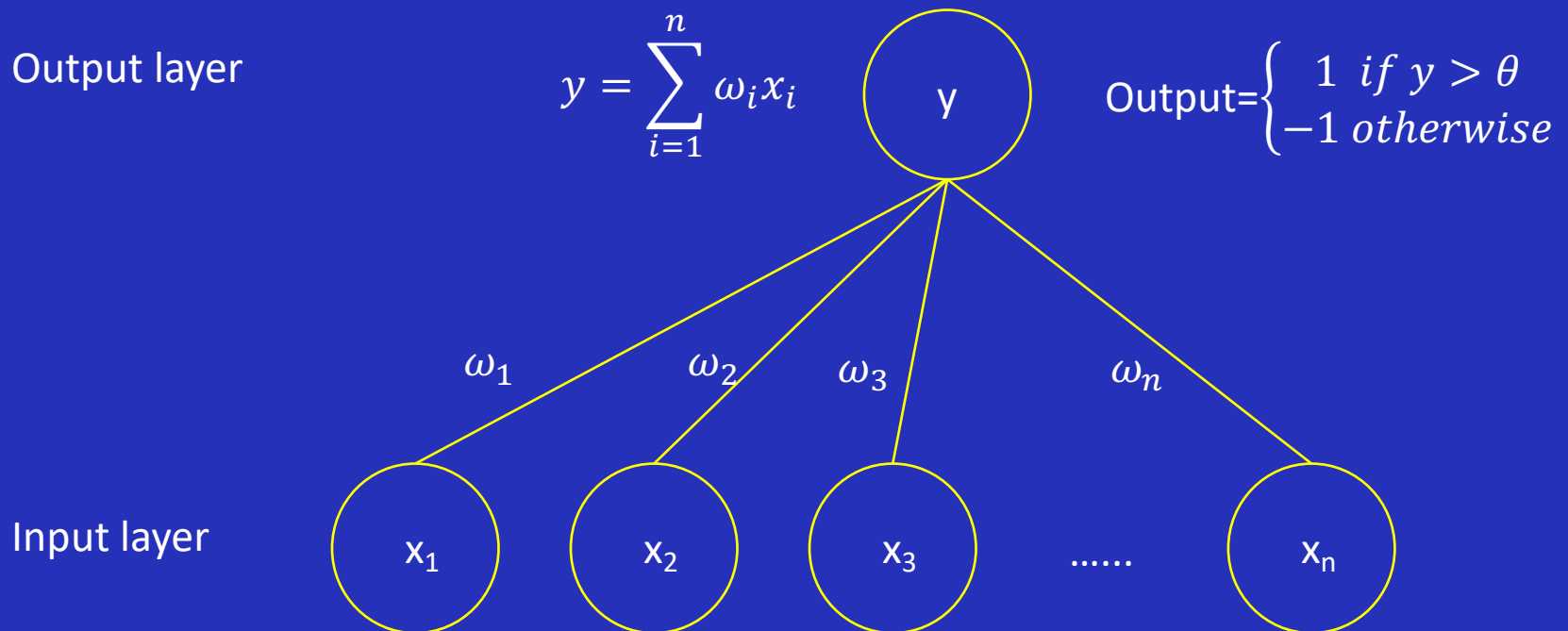Johns Hopkins University

# Context

- Introduction to Neural Networks
- Introduction to Deep Convolutional Neural Networks (DCNN)
- Deep Learning in Medical Image Segmentation
- DCNN Layer Functionality
- DCNN Architecture functionality
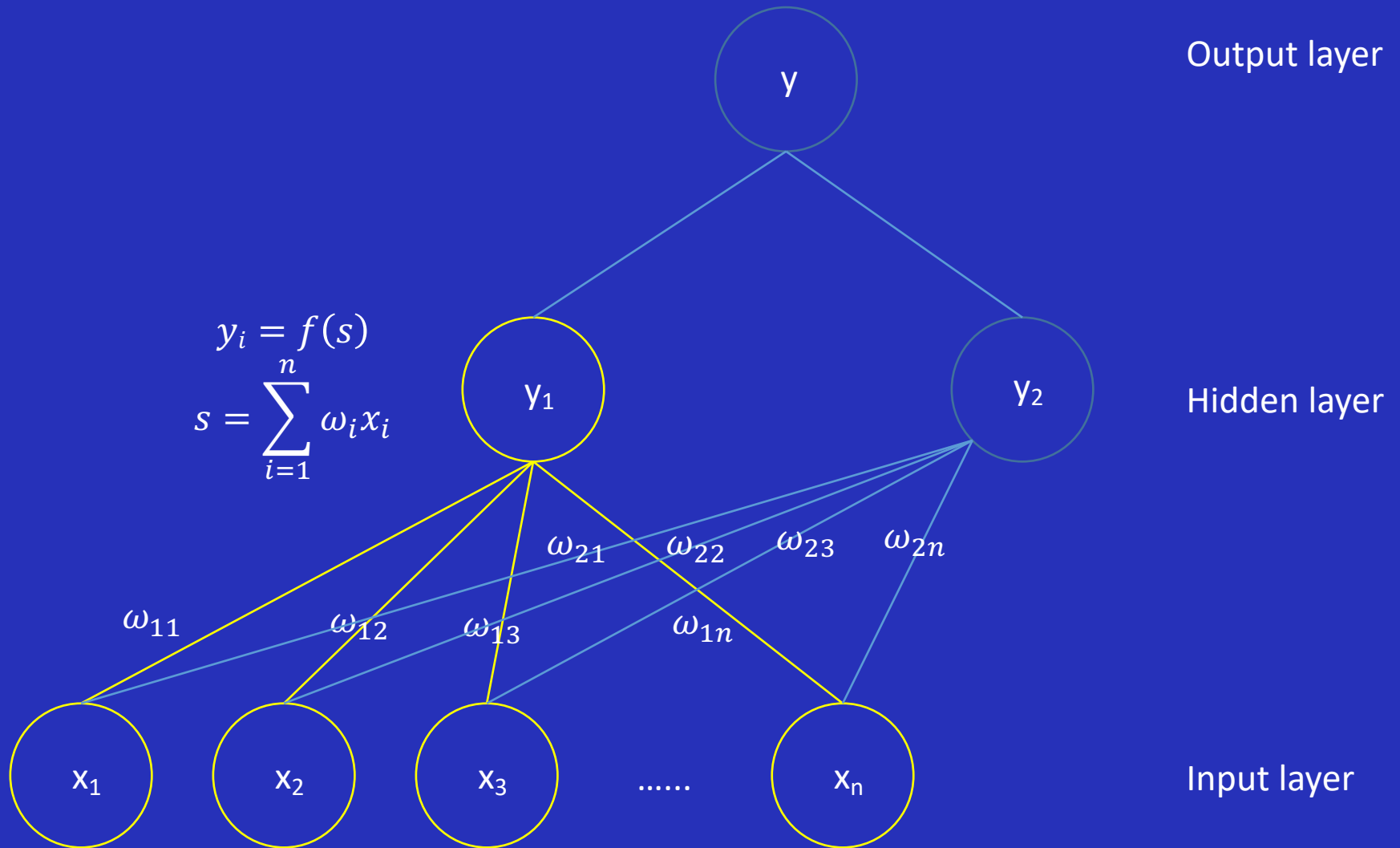
# Conventional Neural Networks

# Single-layer Perceptrons (SLP)

- Can classify linearly separable data into binary classes: -1 and 1.

- A feed-forward network based on a threshold transfer function

Output layer

$$y = \sum_{i=1}^{n} \omega_i x_i$$

$y$

$$\text{Output} = \begin{cases} 1 & \text{if } y > \theta \\ -1 & \text{otherwise} \end{cases}$$

$\omega_1$    $\omega_2$    $\omega_3$    $\omega_n$

Input layer    $x_1$    $x_2$    $x_3$    ……    $x_n$

# Multi-layer Perceptrons (MLP)

Output layer

$$y_i = f(s)$$
$$s = \sum_{i=1}^{n} \omega_i x_i$$

y

$y_1$ $y_2$

Hidden layer

$\omega_{21}$ $\omega_{22}$ $\omega_{23}$ $\omega_{2n}$

$\omega_{11}$ $\omega_{12}$ $\omega_{13}$ $\omega_{1n}$

$x_1$ $x_2$ $x_3$ ...... $x_n$

Input layer

# About MLP

- Differs from SLP by two things:
  - A soft thresholding function after each summation
  - Introduction of hidden layers

- Many levels can be specified to model non-linear relationship

- The number of hidden units is related to the capacity of the perceptron
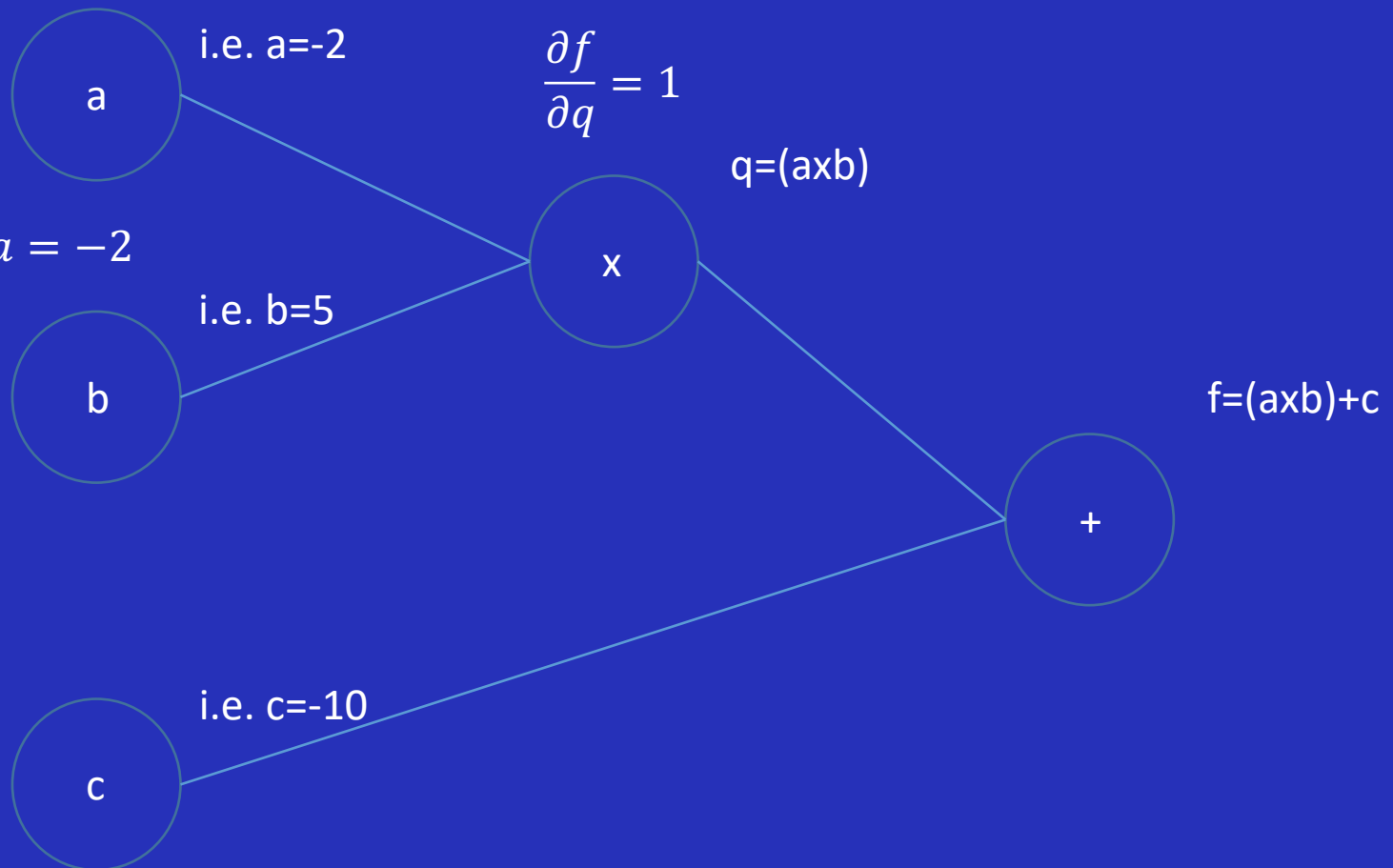
# Backpropagation

- To apply the chain rule many many times to calculate the gradient of a loss function with respect to all the inputs (weights, input data) in the network.
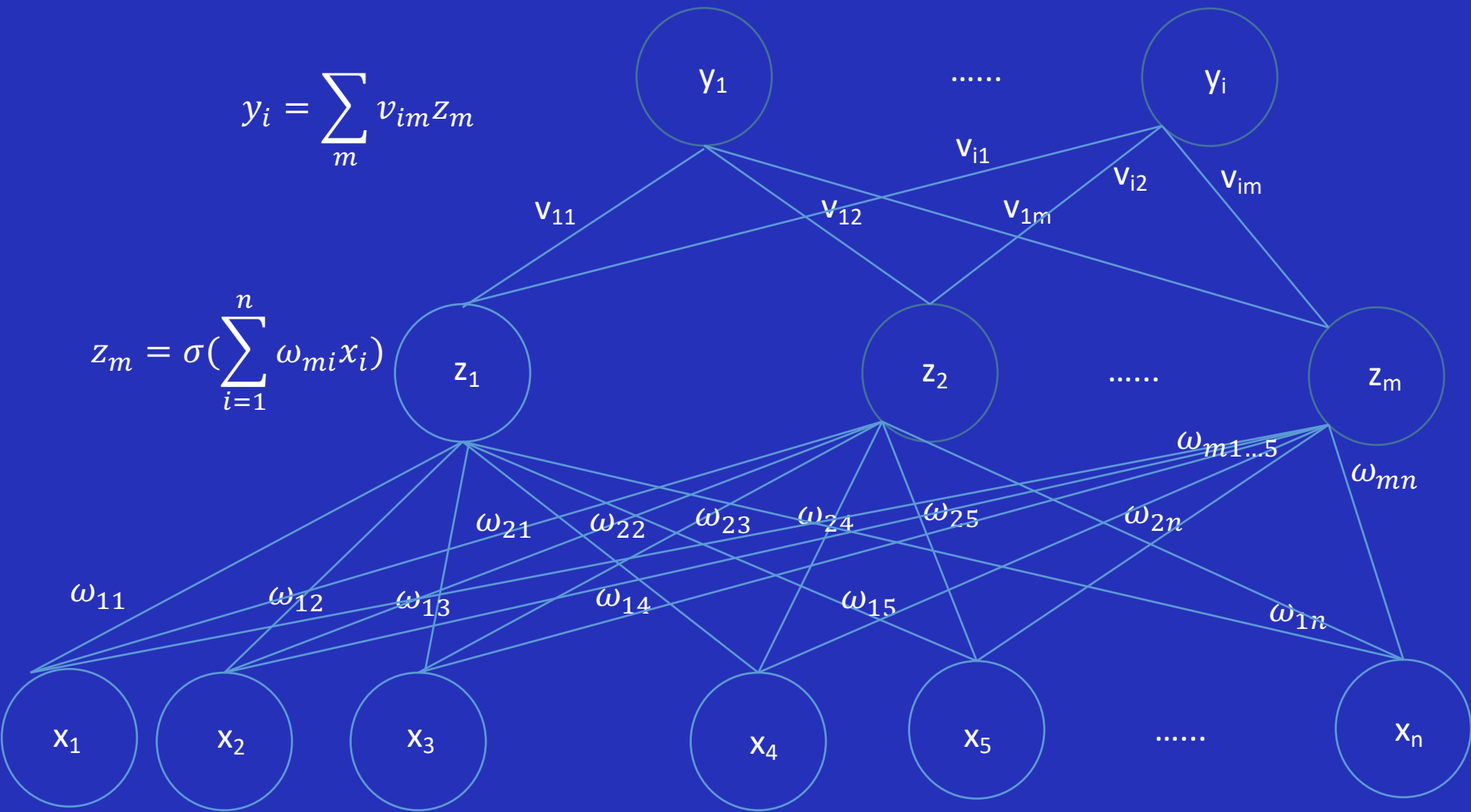
# Backpropagation (continued)

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial q}\frac{\partial q}{\partial a} = b = 5$$

i.e. a=-2

$$\frac{\partial f}{\partial q} = 1$$

q=(axb)

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial q}\frac{\partial q}{\partial b} = a = -2$$

i.e. b=5

f=(axb)+c

$$\frac{\partial f}{\partial c} = 1$$

i.e. c=-10

a

b

x

c

+

# Backpropagation for MLP



$$y_i = \sum_m v_{im} z_m$$

$$z_m = \sigma(\sum_{i=1}^{n} \omega_{mi} x_i)$$

# Backpropagation for MLP (cont'd)

- Loss function

$$E[\omega, \nu] = \sum_i \{y_i - \sum_m \nu_{im}\sigma(\sum_n \omega_{mn}x_n)\}^2$$

- Update terms are negative derivatives of the loss w.r.t the local parameters (weights)

$$\Delta\omega_{mn} = -\frac{\partial E}{\partial \omega_{mn}} \qquad \Delta\nu_{im} = -\frac{\partial E}{\partial \nu_{im}}$$
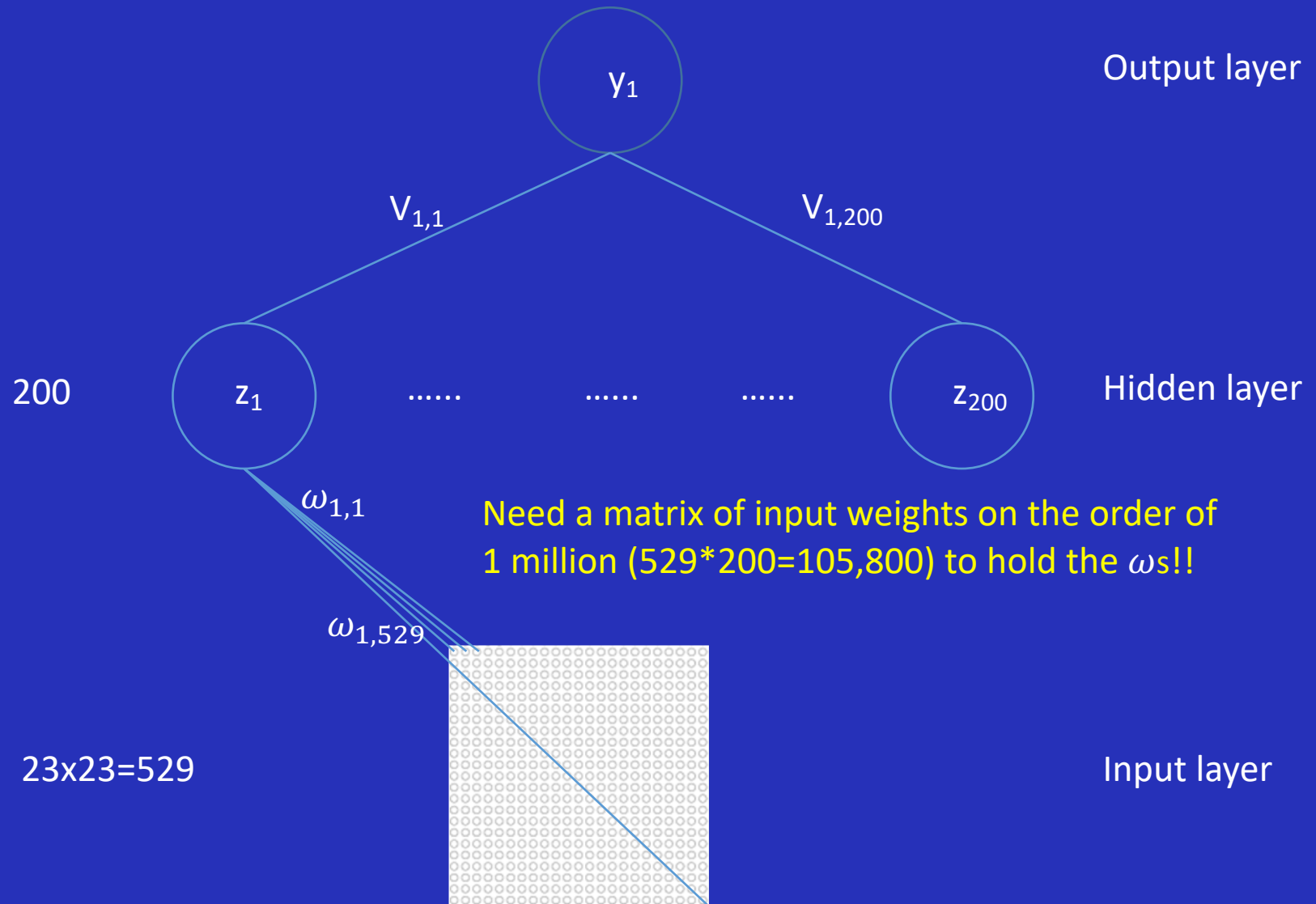
- By defining $z_m = \sigma(\sum_{i=1}^n \omega_{mi}x_i)$ and $E = \sum_{i=1}^n (y_i - \sum_m \nu_{im}z_m)^2$ … … …

$$\frac{\partial E}{\partial \omega_{mn}} = 2\sum_i (y_i - \sum_m(\nu_{im}z_m))\nu_{im}x_n\sigma(\sum_n \omega_{mn}x_n)\{1 - \sigma(\sum_n \omega_{mn}x_n)\}$$
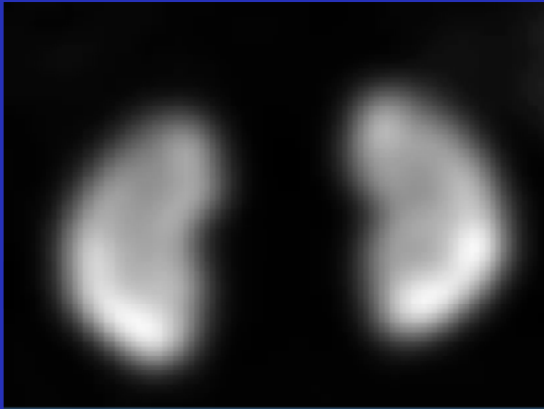
Detailed derivations are available at:
http://garyliye.com/Multilayer_perceptron_and_backpropagration.pdf

# What about a Real-world Image?

Output layer

$y_1$

$V_{1,1}$ $V_{1,200}$

200 $z_1$ …… …… …… $z_{200}$ Hidden layer

$\omega_{1,1}$

Need a matrix of input weights on the order of 1 million (529*200=105,800) to hold the $\omega$s!!

$\omega_{1,529}$

23x23=529 Input layer

# Spatial Structure



What we see



We computers see

Vectorizing an image completely ignores the complex 2D spatial structure of an image
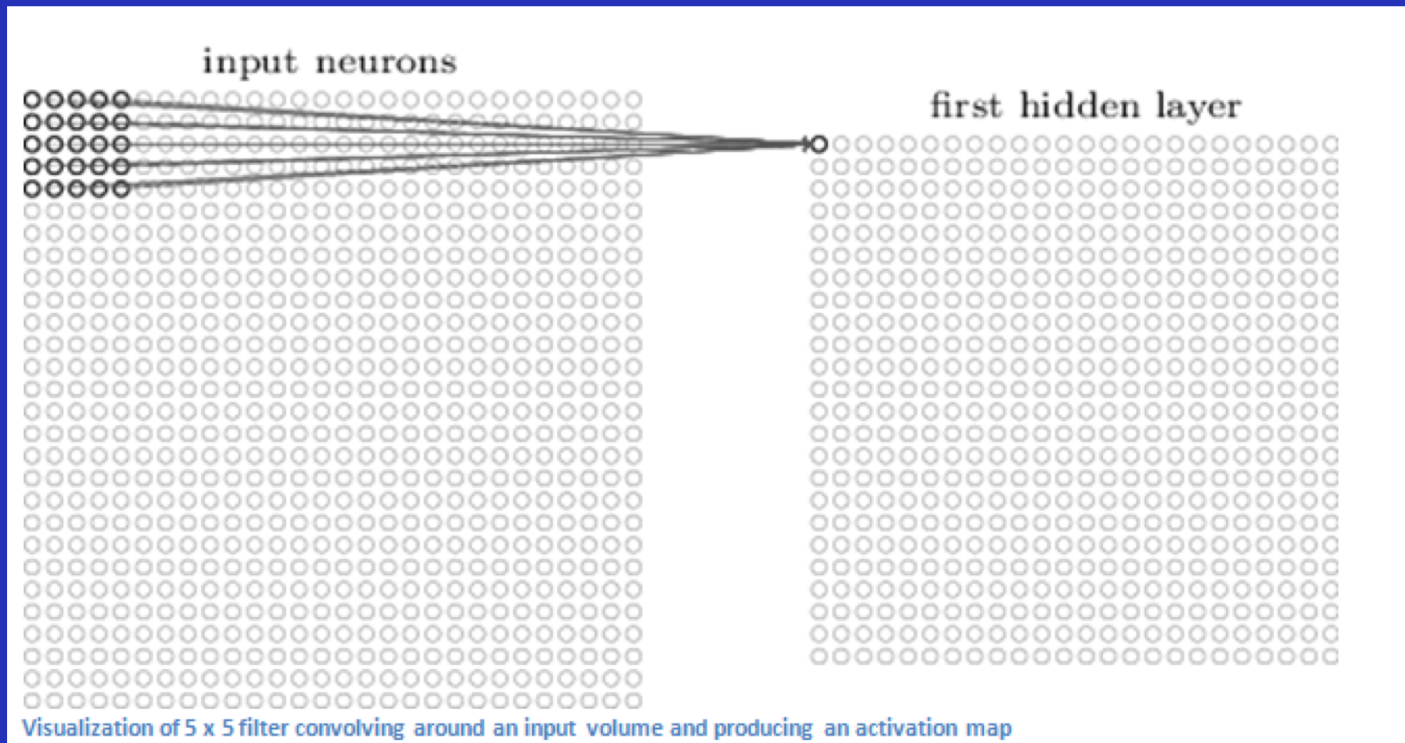
# Limitations of Conventioanl Neural Networks

- Impractical for real-world image classification

- Ignores 2D/3D spatial structure in image

- Solution to overcome both these disadvantages?

# One solution: Convolution

Use 2D convolution instead of matrix multiplications:
-Learning a set of convolutional filters (each of 5x5,say) is much more tractable than learning a large matrix (529x200)



input neurons

first hidden layer

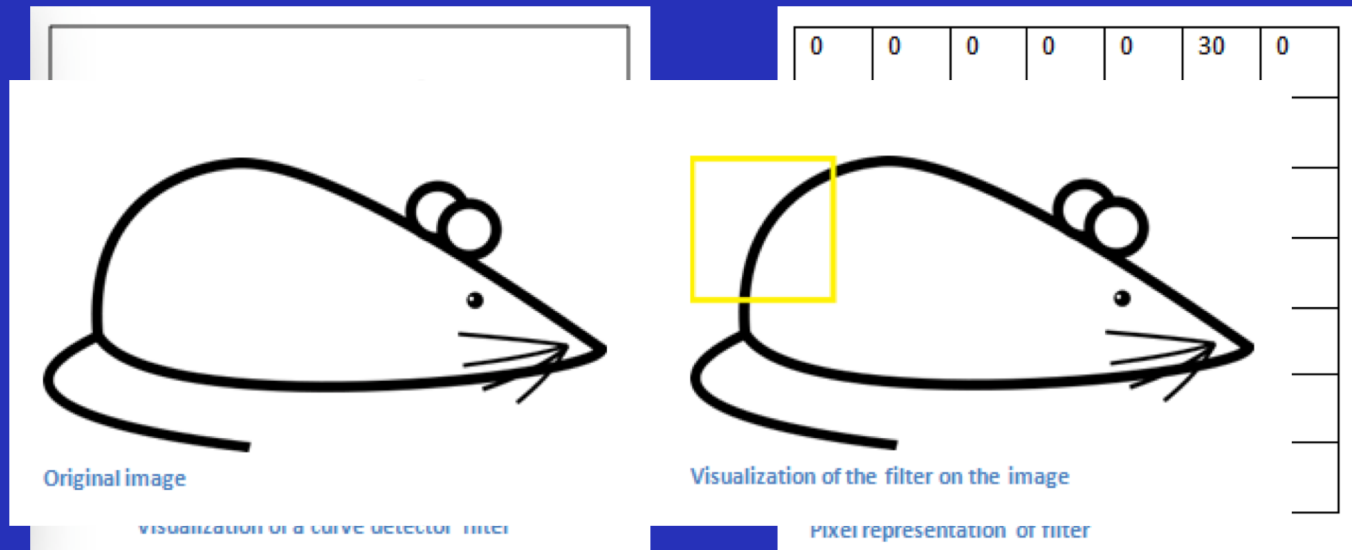Visualization of 5 x 5 filter convolving around an input volume and producing an activation map

# Convolutional Neural Networks

# Convolutional Neural Networks (CNN)

- CNN has proven very powerful
  -Retains structural or configural information in neighboring pixels or voxels in a (medical) image
  -Exploits extensive weight-sharing to reduce the degrees of freedom of models
  -Composed of convolution layers interspersed with pooling (sub-sampling) layers
  -Highly parallelizable
  -GPU implementations can accelerate 40 times or more
  -Trained using backpropagation algorithm and lots of labeled data
  -First uses in medical imaging in 1990's

- Improvement of artificial neural networks
  - More layers, higher levels of abstraction, improved predictions
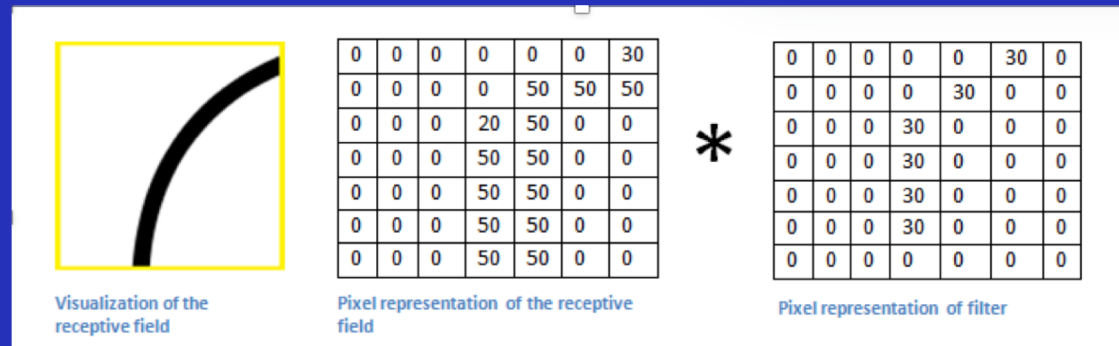
# High Level Perspective of Convolution

- Convolutional filters are essentially feature identifiers
- Features can be high-level (abstract) and low-level such as straight edges, simple colors, and curves.



Original image

Visualization of a curve detector filter

Visualization of the filter on the image

Pixel representation of filter

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |

# High Level Perspective of Convolution (cont'd)

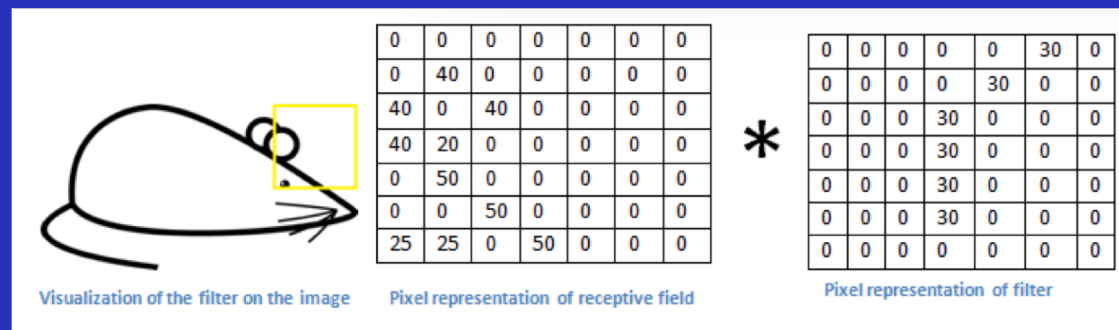The output of the filter has a high activation value. Or say, the neuron is fired/excited!

High activation



Visualization of the receptive field

Pixel representation of the receptive field

Pixel representation of filter

Multiplication and Summation = (50*30)+(50*30)+(50*30)+(20*30)+(50*30) = 6600 (A large number!)

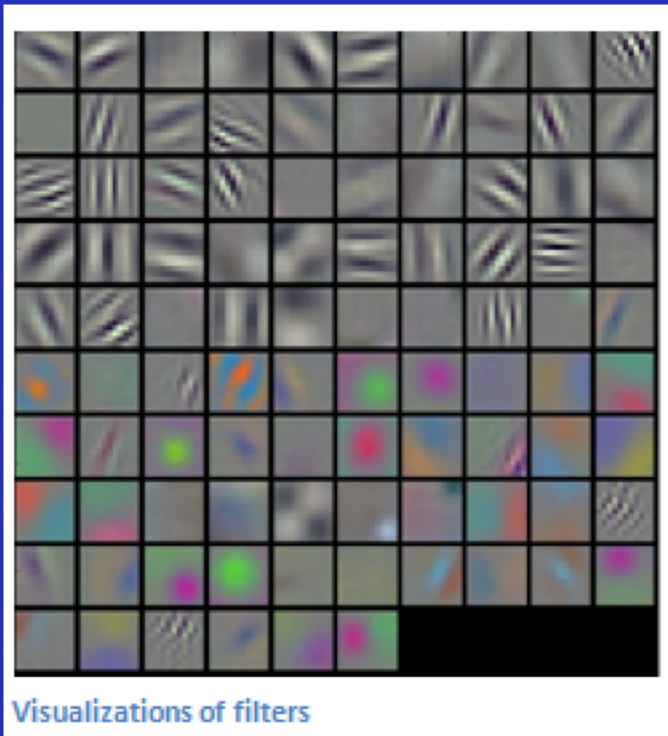The output of the filter has a low activation value.

No activation



Visualization of the filter on the image

Pixel representation of receptive field

Pixel representation of filter

Multiplication and Summation = 0

https://adeshpande3.github.io/A-Beginner's-Guide-To-Understanding-Convolutional-Neural-Networks/
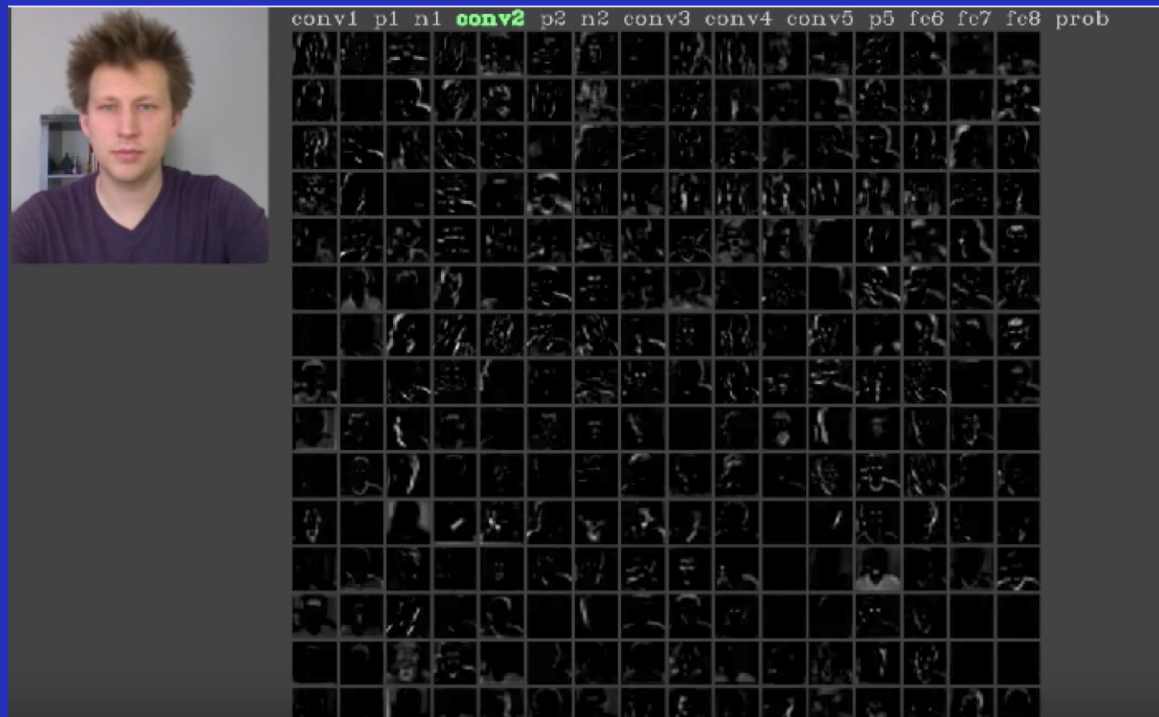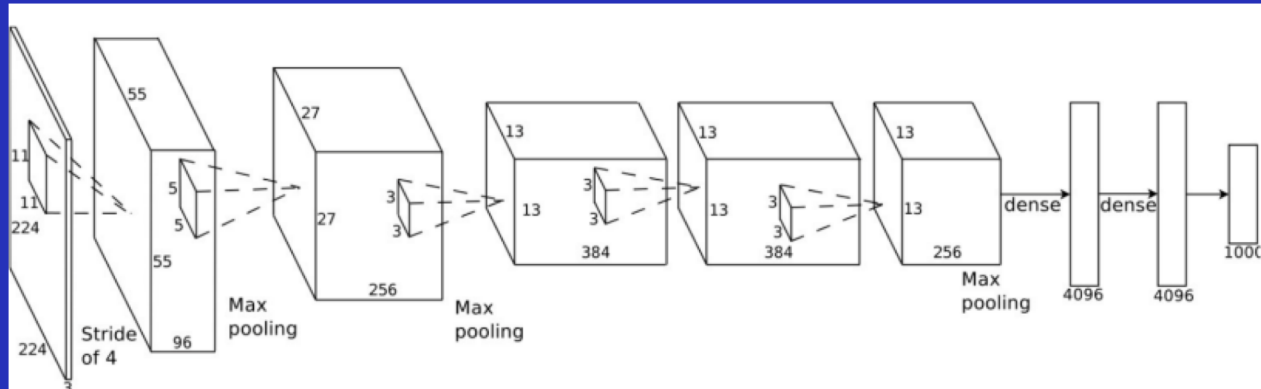
# 1st Conv Layer **Filters** Learned in AlexNet

Example filters learned by Krizhevsky et al. Each of the 96 filters shown below is of size 11x11x3.

**Each layer of the activation map(s) is basically describing the locations in the original image for where certain low level features appear.**
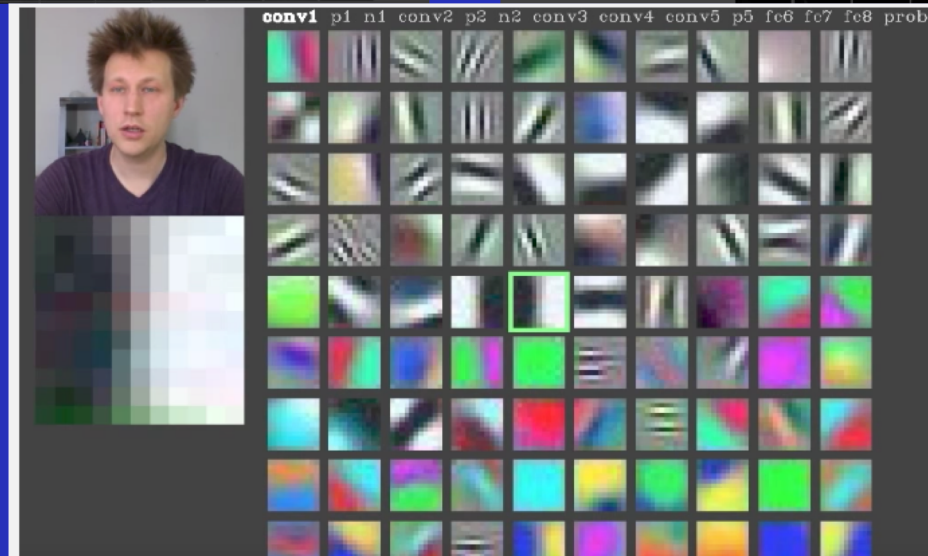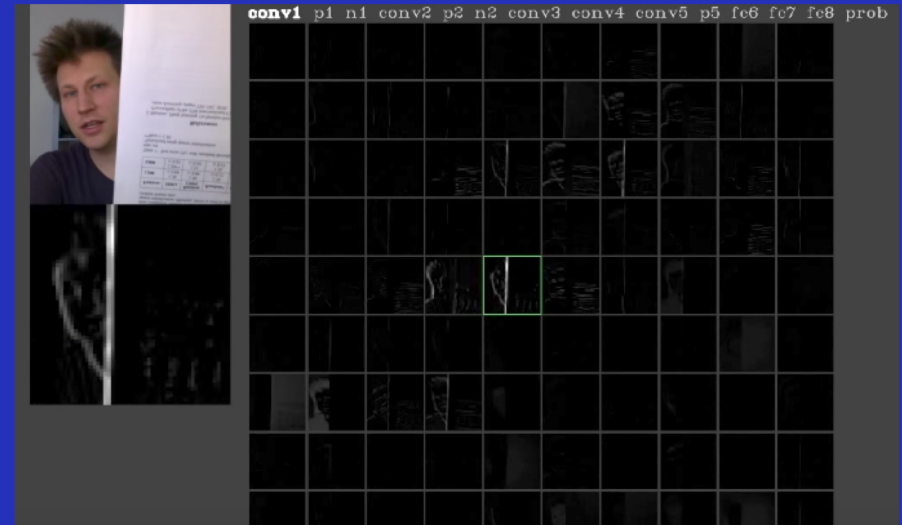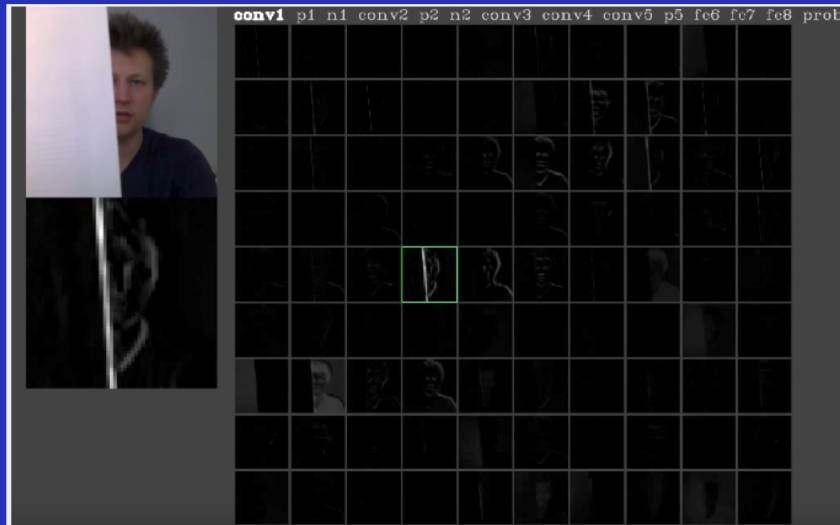


Visualizations of filters

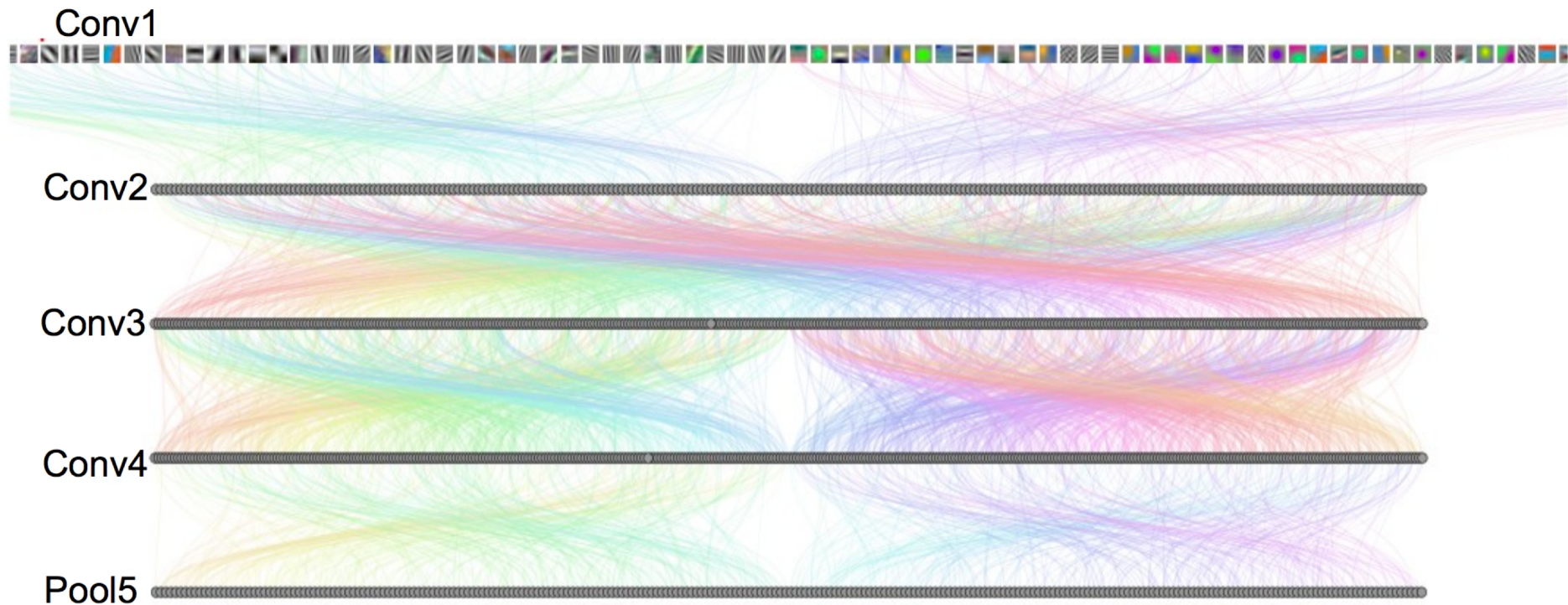# The 2$^{nd}$ Conv Layer Activation Map

# Filters and Activation Maps

# Connection Weights Between Convolutional Layers

- Let the learnable connection weights connecting feature map i at layer $l-1$ and the feature map j at the layer $l$ be $k_{ij}^l$. Specifically, the units of the convolutional layer $l$ compute their activations $A_j^l$ based only on a spatially contiguous subset of units in the feature maps $A_i^{l-1}$ of the preceding layer $l-1$ by convolving the kernels $k_{ij}^l$ as follows:

$$A_j^l = f(\sum_{i=1}^{M(l-1)} A_i^{l-1} * k_{ij}^l + b_j^l)$$

Say if there are 5 feature maps at layer $l-1$ and 4 feature maps at layer $l$, there would be 4 axax5(depth of the feature map at previous layer) connection weights

# How objects are represented in CNN?



Bolei Zhou, et al., arXiv.org
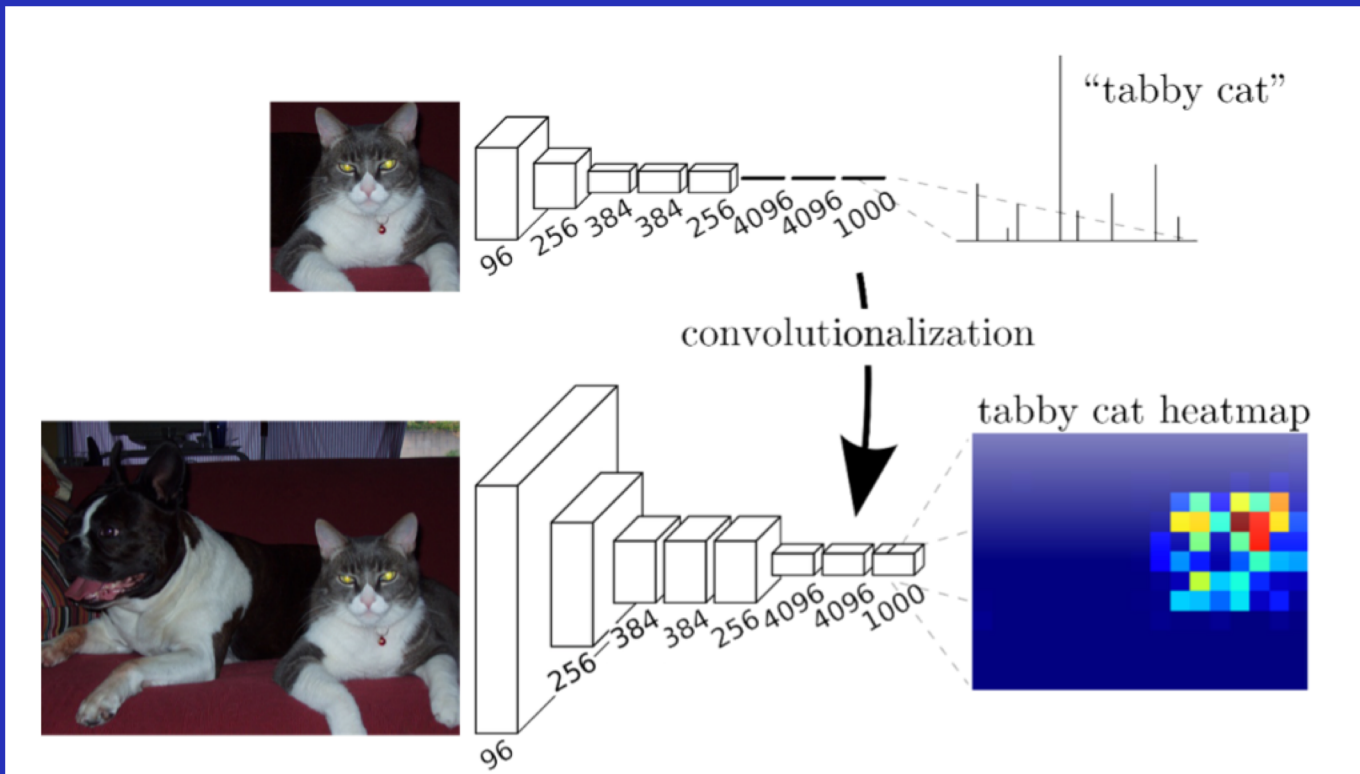
# Neuroscience connection

- Similar (convolution-like) computations within the human brain
- Primary visual cortex has simple and complex cells
- The simple cells responded primarily to oriented edges and gratings
- The complex cells were also sensitive to these edges and grating but exhibited spatial invariance

# Deep Learning in Medical Imaging

- Difficult to obtain large enough training data

- Some solutions to lack of "big data" in medical imaging

- What architecture to use?

H. Greenspan et al., IEEE TMI, May 2016

# Segmentation: pixel-wise classification

Transforming fully connected layers into convolution layers enables a classification net to output a spatial map.



tabby cat heatmap

Shelhamer, E. Long, J. Darrell, T. IEEE Trans Pattern Anal Mach Intell, 2016

# Network Depth and Receptive Field Size

- As you go deeper into the network, the filters begin to have a larger and larger receptive field, which means that they are able to consider information from a larger area of the original input volume (another way of putting it is that they are more responsive to a larger region of pixel space)

Layer functionality

# Convolutional layer

- Local connectivity
  -Because we use convolutional filter with size much smaller than the image it operates on. This contrasts with the global connectivity paradigm relevant to vectorized images

- Weight sharing
  -The same filter applied across the image

- Can be seen as a **local independent feature-detector;** To detect local features (local connectivity) at different position in the input feature maps
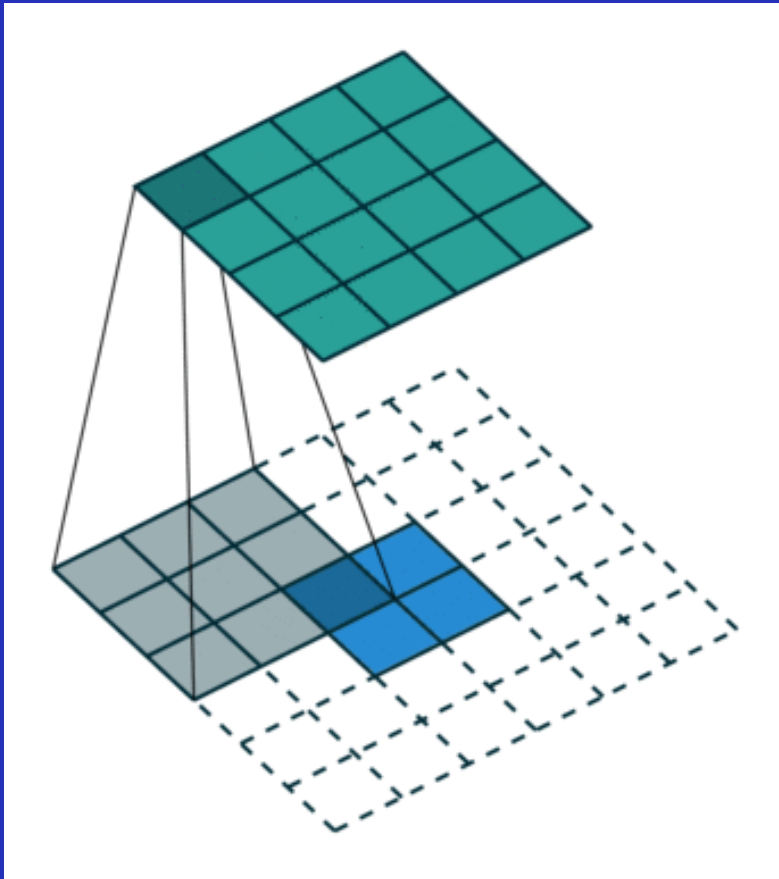
# Max-Pooling

- The Neocognitron model inspired the modeling of simple cells as convolutions.

- The complex cells can be modeled as a max-pooling operation, which can be thought as a max filter.

- Picks the highest activation in a local region, thus providing a small degree of spatial invariance, which is analogous to the operation of complex cells.
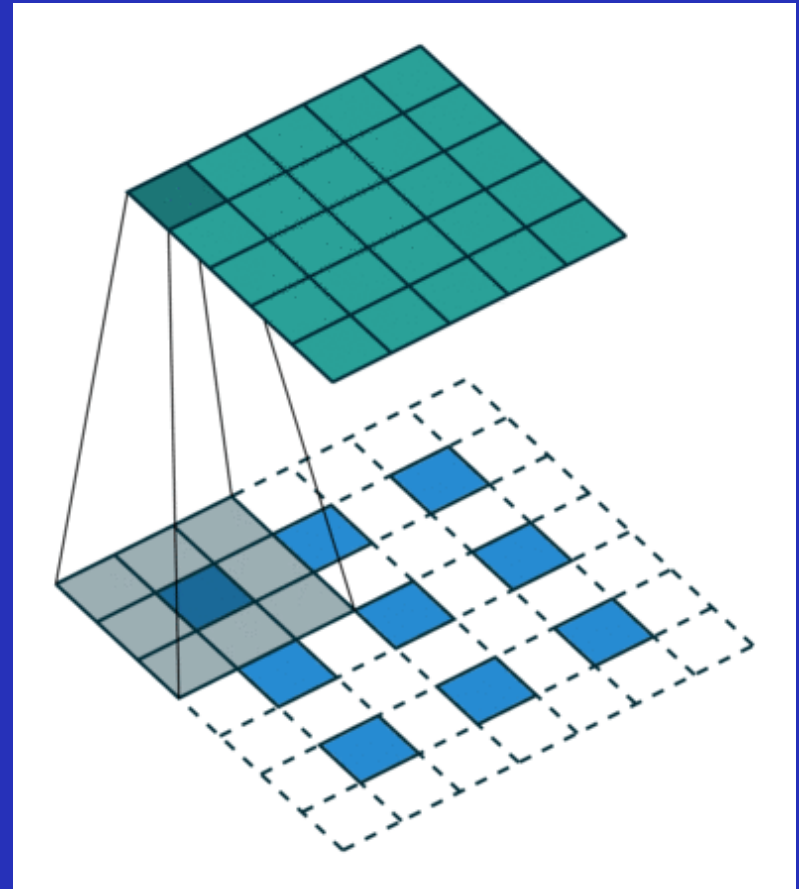
# Non-linearity layer

- Necessary because cascading linear (like convolution) systems is another linear system
- Non-linearity between layers ensure that the model is more expressive than a linear model
- In theory, no non-linearity has more expressive power than any other, as long as they are continuous, bounded, and monotonically increasing.
- Maas et al. introduced a new kind of nonlinearity, called the leaky-ReLU. ReLU(x)=max(0,x)+bmin(0,x)
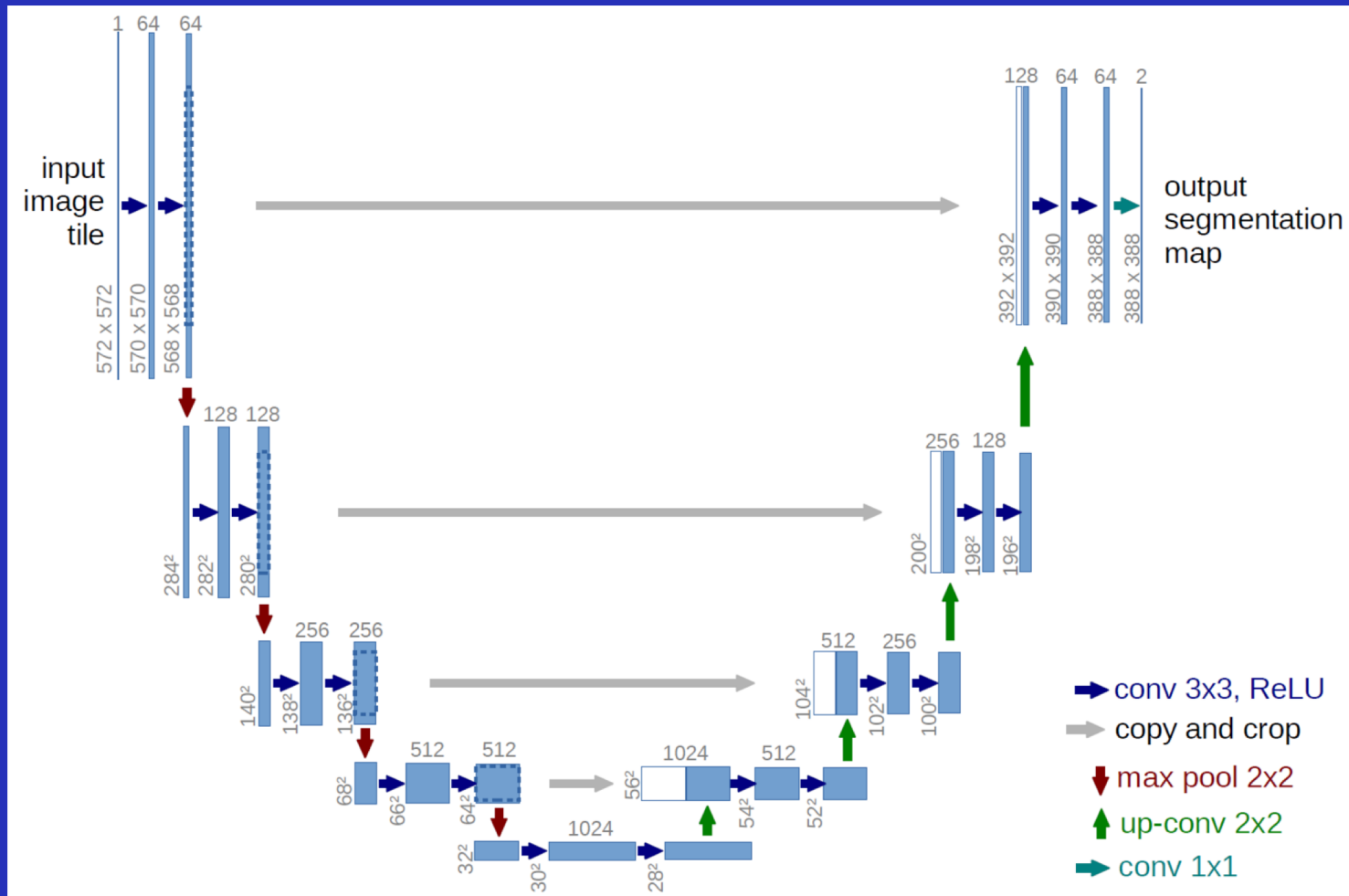
# Deconvolutional Layer
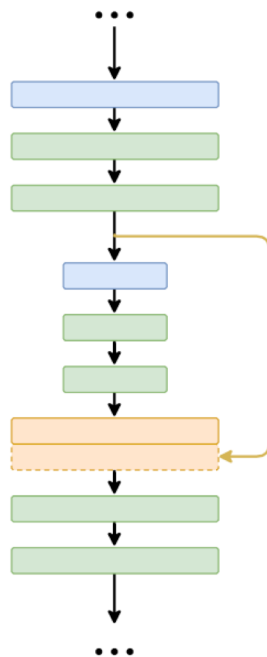


Without Padding

With Padding

# Architecture functionality (segmentation)

# Encoder-decoder architecture (U-net)

# Summation based skip architecture

U-net:
Copy
and past

FusionNet:
Copy and
add



Tran M.Q, et al., **arXiv.org**