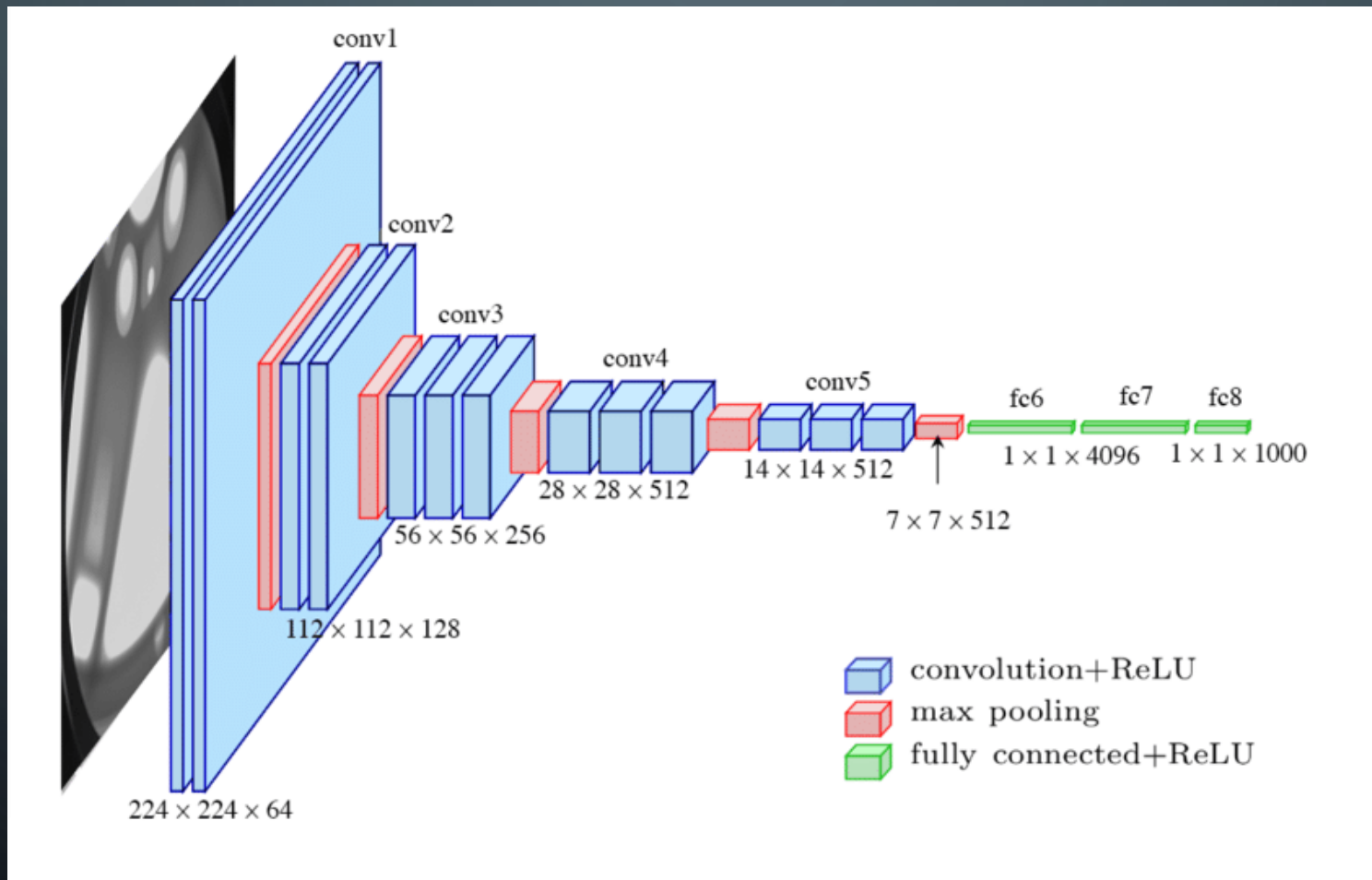




# POPULAR MODELS

MOHAMMAD GHODDOSI

# VGG16

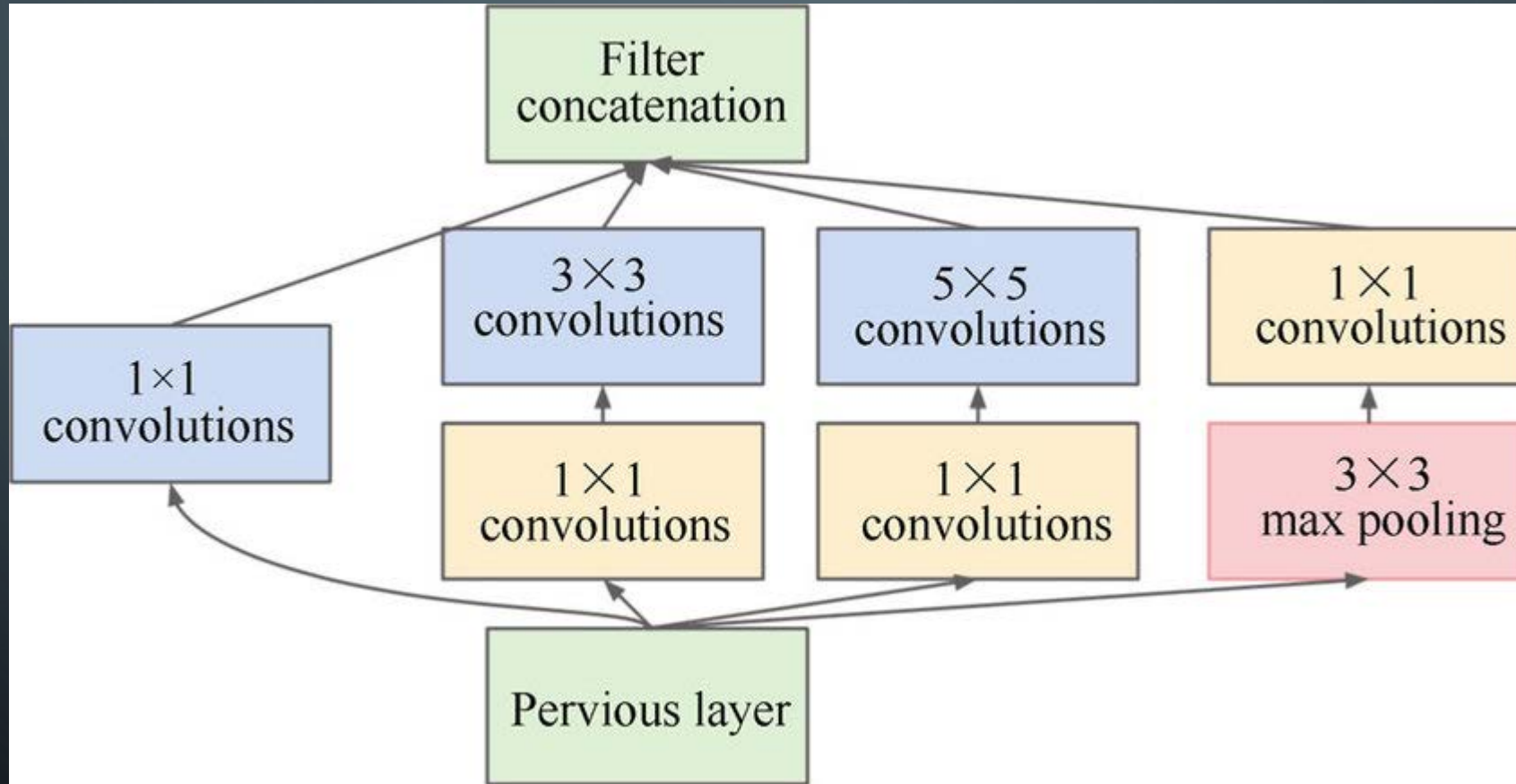


# GOOLENET (IDEAS)

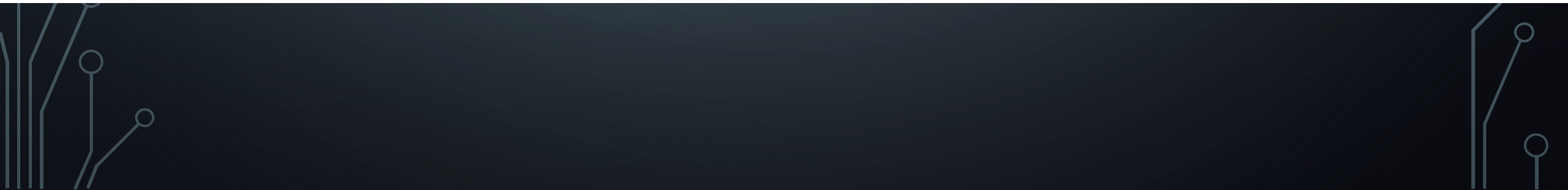


- We need to go deeper
  - Vanishing gradient
- 1x1 Convolution
  - Dimension reduction
- Inception module
- Global Average Pooling
- Auxiliary Classifiers

# GOOGLNET (INCEPTION MODULE)



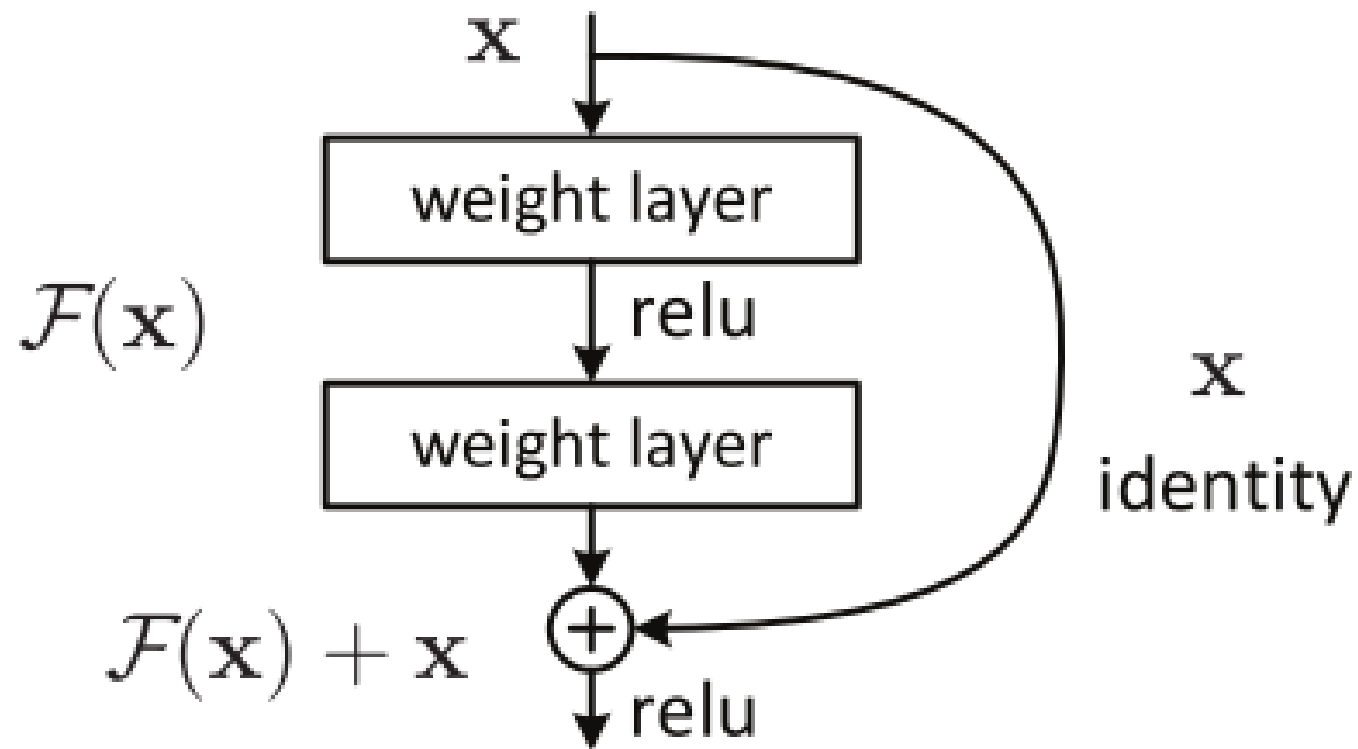






type	patch size/ stride	output size	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	params	ops
convolution	7×7/2	112×112×64	1							2.7K	34M
max pool	3×3/2	56×56×64	0								
convolution	3×3/1	56×56×192	2		64	192				112K	360M
max pool	3×3/2	28×28×192	0								
inception (3a)		28×28×256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28×28×480	2	128	128	192	32	96	64	380K	304M
max pool	3×3/2	14×14×480	0								
inception (4a)		14×14×512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14×14×512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14×14×512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14×14×528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14×14×832	2	256	160	320	32	128	128	840K	170M
max pool	3×3/2	7×7×832	0								
inception (5a)		7×7×832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7×7×1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7×7/1	1×1×1024	0								
dropout (40%)		1×1×1024	0								
linear		1×1×1000	1							1000K	1M
softmax		1×1×1000	0								

# RESNET (RESIDUAL BLOCK)



# RESNET



layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112	7×7, 64, stride 2				
conv2_x	56×56	3×3 max pool, stride 2				
		$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				



# MOBILENET (IDEAS)



- Significant reduction in number of parameters. (same depth)
- Lightweight deep neural network
- Depthwise separable convolution
  - Depthwise convolution
  - Pointwise convolution



# MOBILENET (CONVOLUTION)

- We can separate dimensions:

$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \cdot [-1 \quad 0 \quad 1]$$

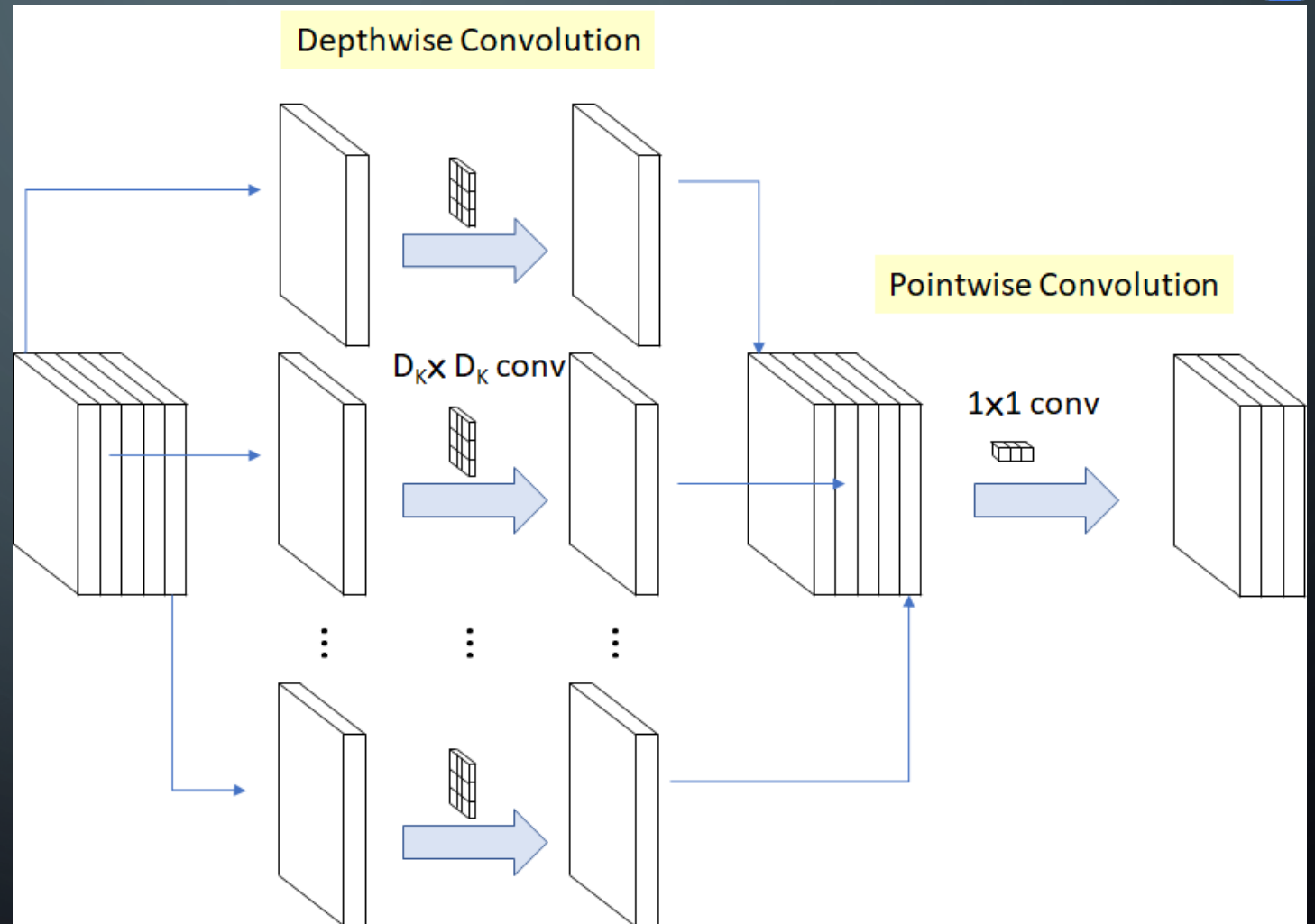
- From 9 parameters to 6 parameters
- Lower operations
- We can apply this for depth to

# MOBILENET (CONVOLUTION)

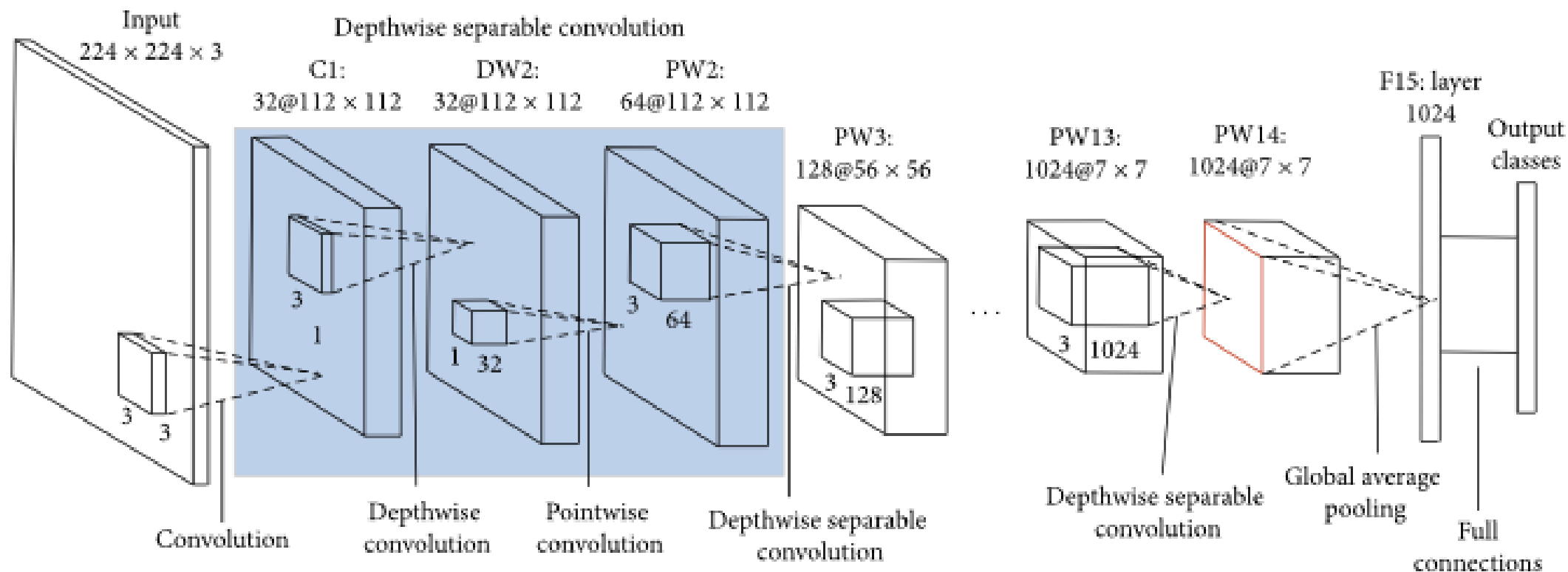


Example:

- 512\*512 input image
- 3\*3 filters
- 64 input depth
- 128 output depth
- Params:
  - Normal => 73,728
  - MobileNet => 704
- Operations:
  - Normal => 38.65e9
  - MobileNet => 4.59e9

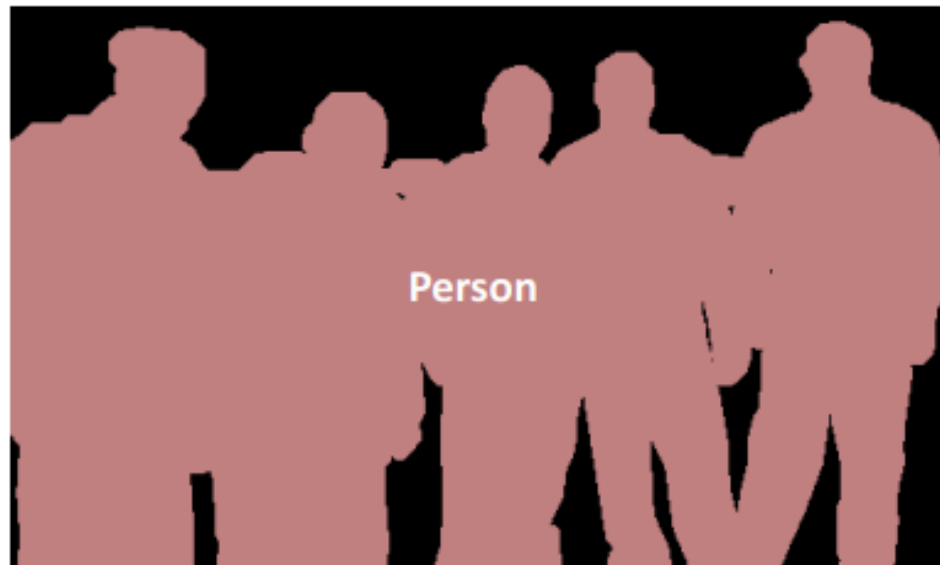


# MOBILENET





# SEGMENTATION



Semantic Segmentation



# U-NET OPERATIONS

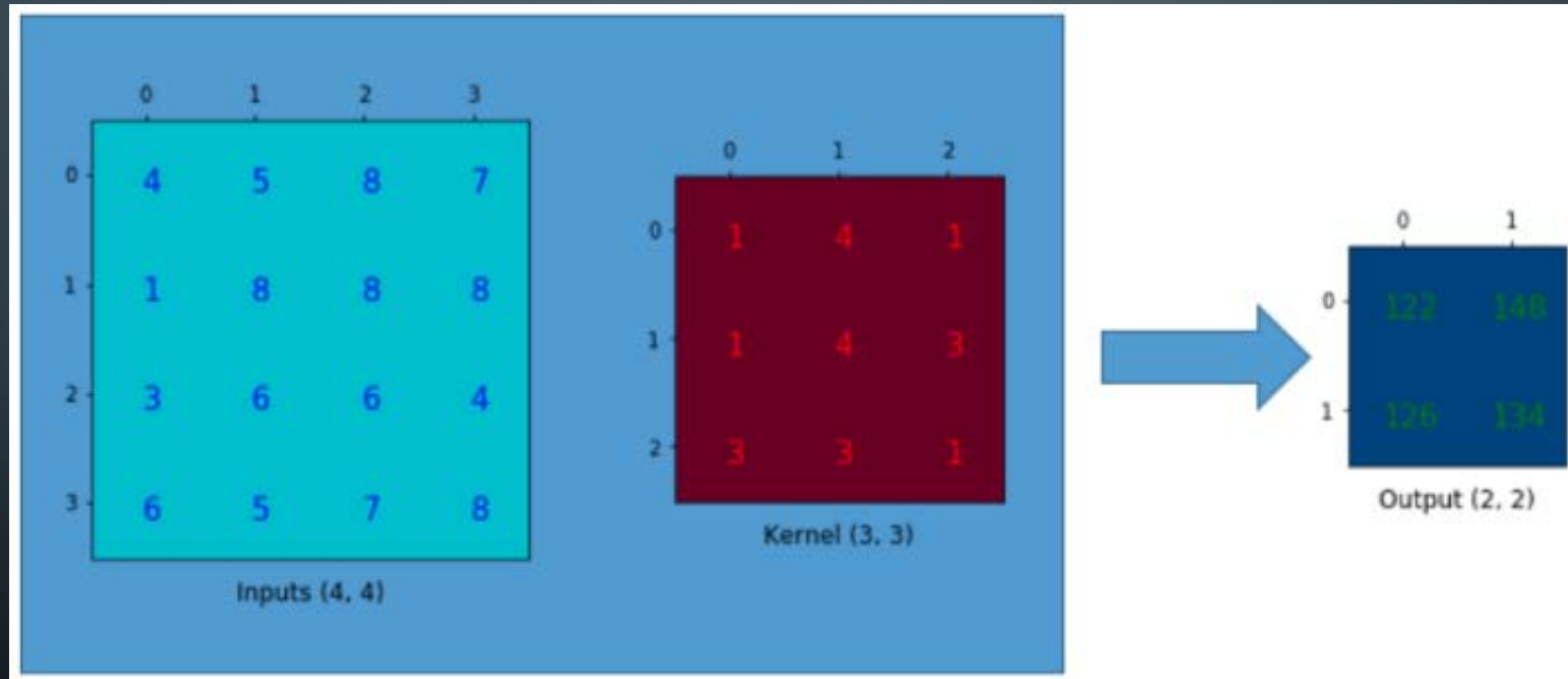
- Convolution
- Pooling (down sampling)
  - Lose “WHERE” information
  - We need “WHERE” and “WHAT” information both
- Up sampling
  - Transpose Convolution (deconvolution)
    - Convolution with fraction stride
- Skip connection

# TRANSPOSE CONVOLUTION



- For matrices we know:
  - $A_{1*100} \cdot B_{100*3} = C_{1*3}$
  - $C_{1*3} \cdot B_{3*100}^T = D_{1*100}$
- We write convolution as dot product and
- Then write it the other way around

# CONVOLUTION





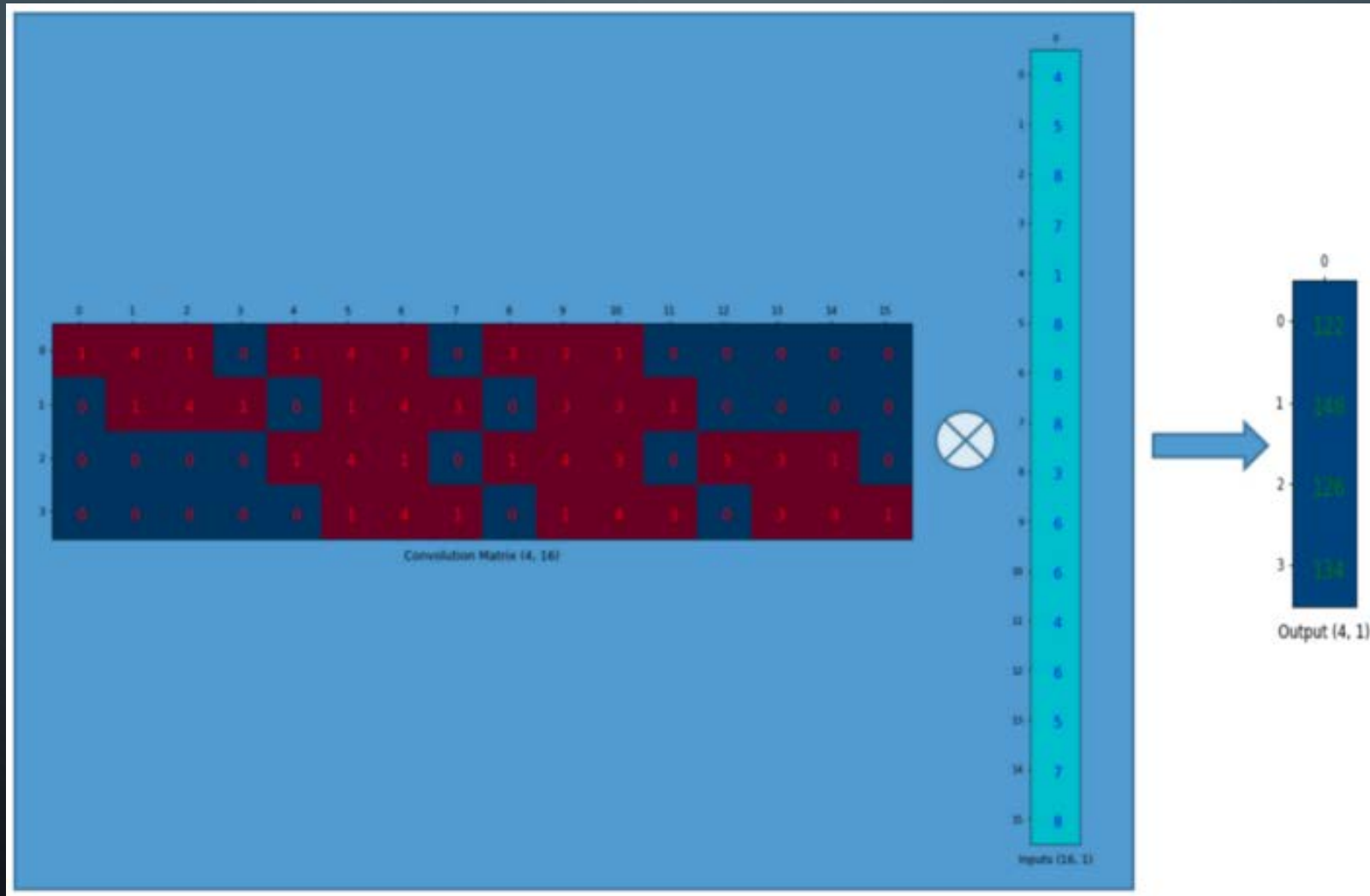
# CONVOLUTION AS DOT PRODUCT



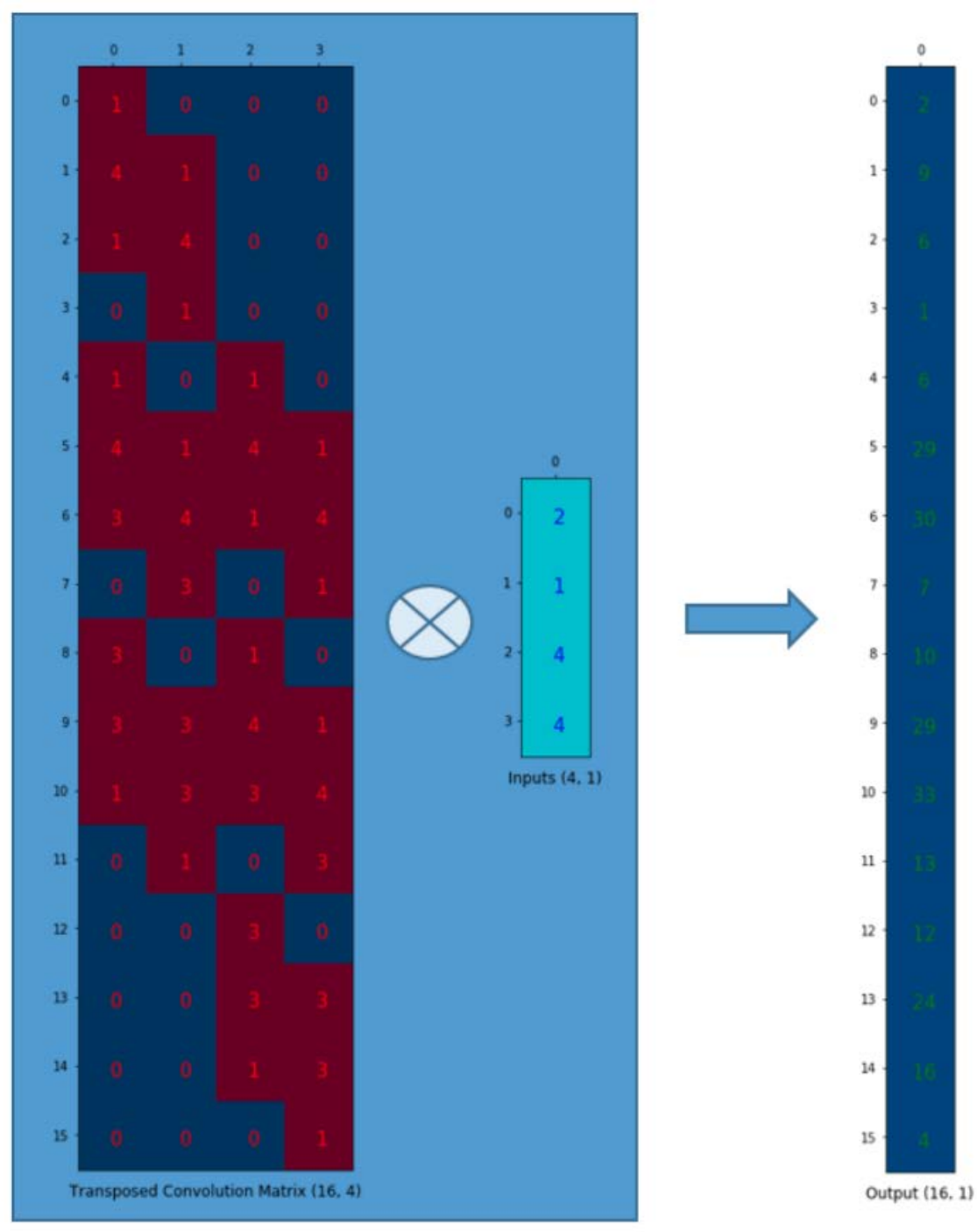
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	1	4	1	0	1	4	3	0	3	3	1	0	0	0	0	0
1	0	1	4	1	0	1	4	3	0	3	3	1	0	0	0	0
2	0	0	0	0	1	4	1	0	1	4	3	0	3	3	1	0
3	0	0	0	0	0	1	4	1	0	1	4	3	0	3	3	1

Convolution Matrix (4, 16)

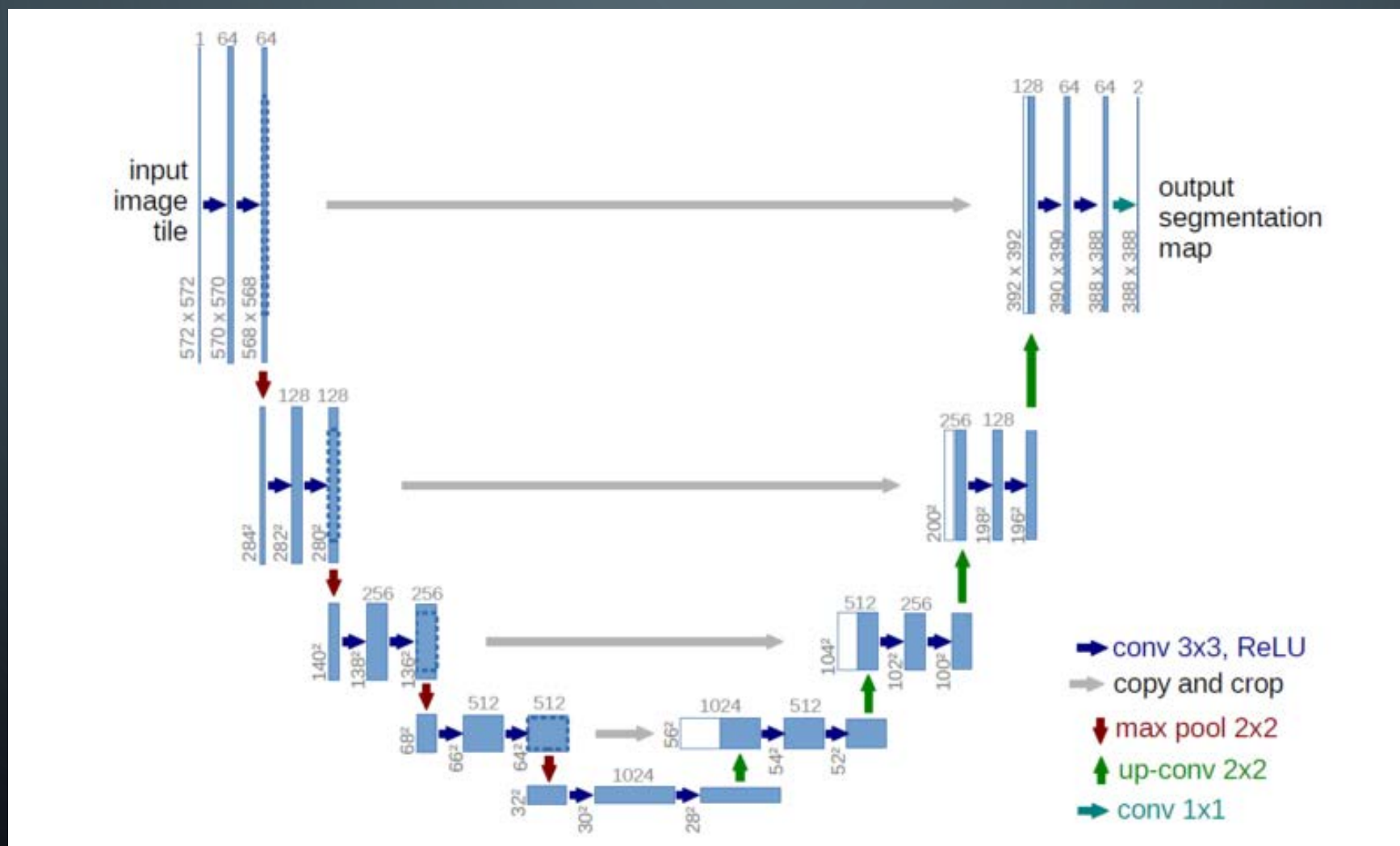
# CONVOLUTION AS DOT PRODUCT



# TRANSPOSE CONVOLUTION



# U-NET ARCHITECTURE





# OBJECT DETECTION



Object Detection

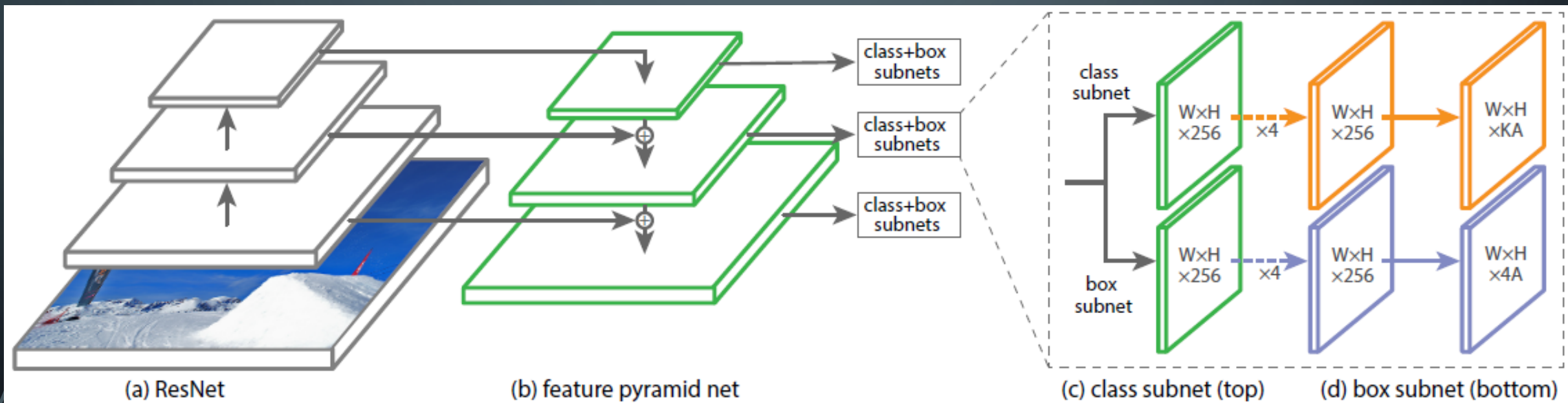
# RETINANET



- Focal loss
  - $FL(p_t) = -(1 - p_t)^\gamma \log(p_t)$
  - Lower weight to easy samples
  - Focus training on hard negatives
- Non maximum suppression

$$p_t = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases}$$

# RETINANET



# FURTHER READING

- LeNet
- AlexNet
- ZFNet
- Mask R-CNN
- YOLO (You Only Look Once)
- SSD (Single Shot Detector)





# IDENTIFICATION

- Verification
- Identification

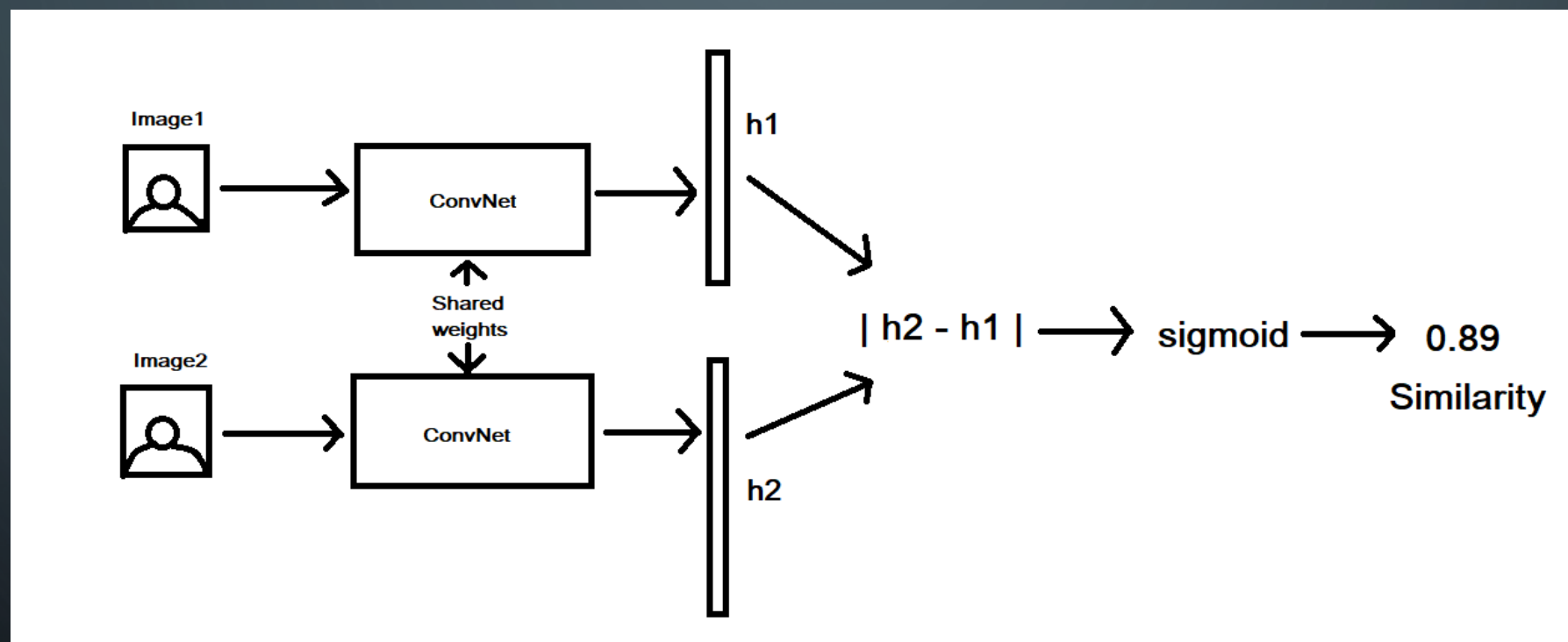


# IDENTIFICATION AS CLASSIFICATION

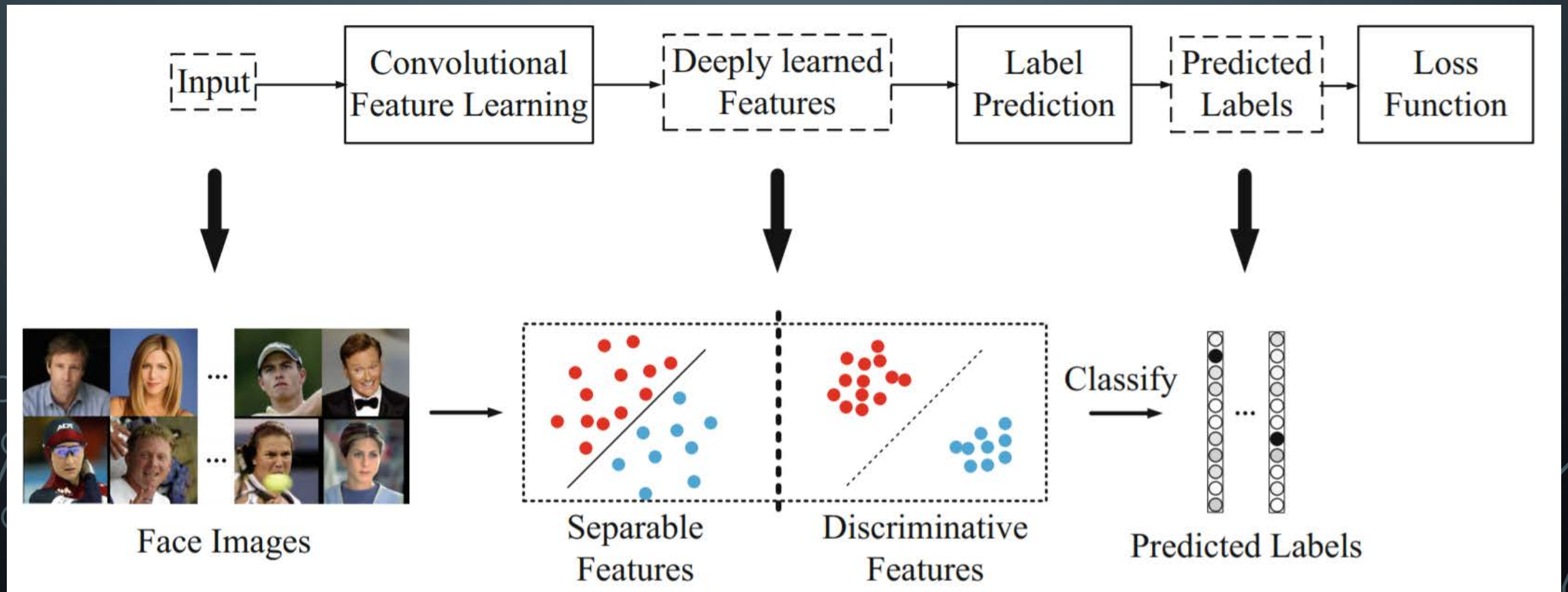


- Each identity is a class
- We need lots of images for each person
- Its not easy to add one person
- Solution -> one shot learning

# SIAMESE NETWORK

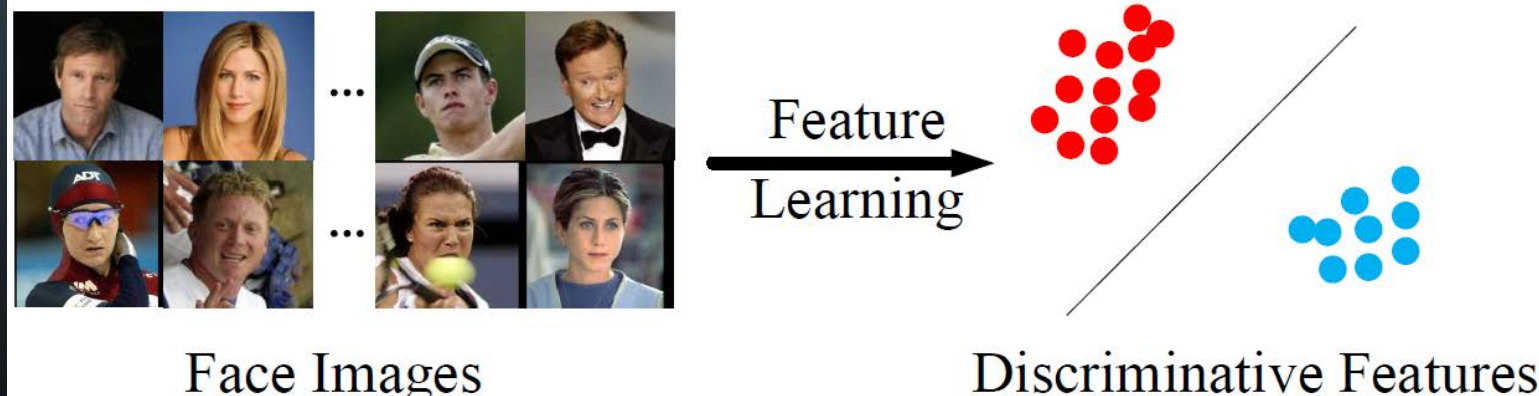
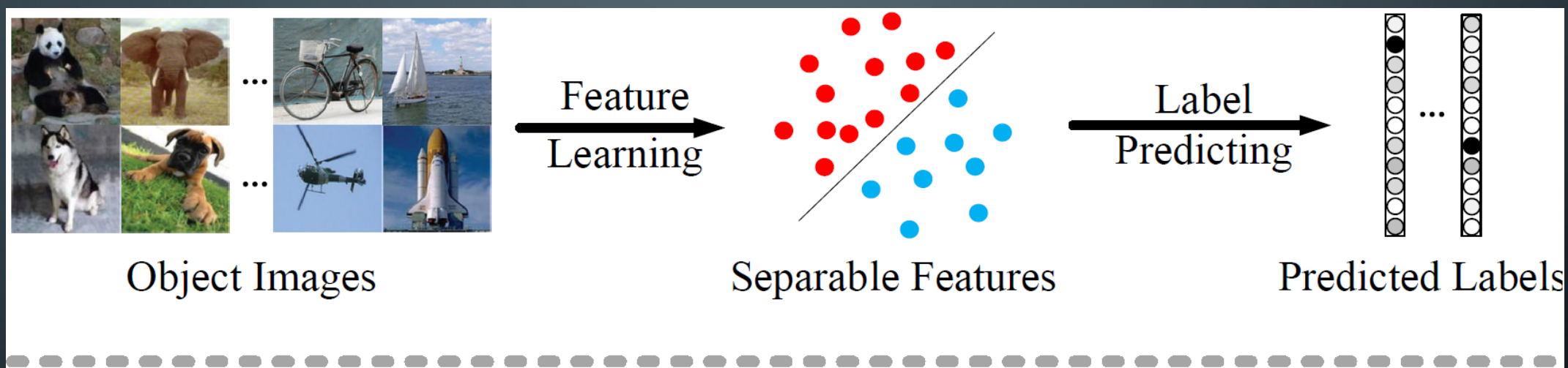


# DISCRIMINATIVE VS SEPARABLE

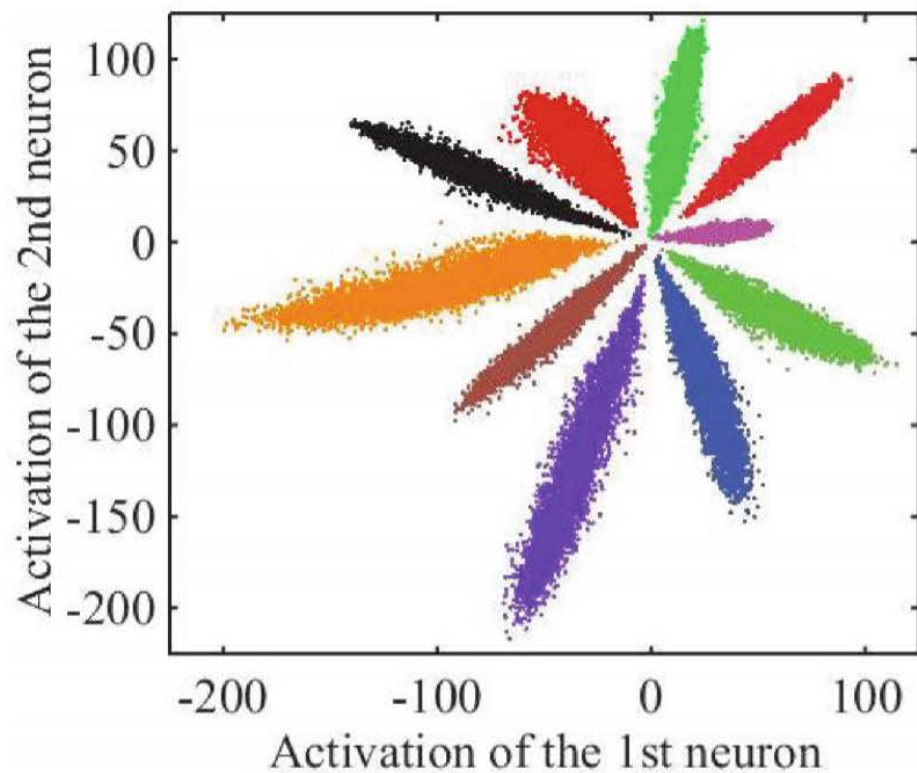




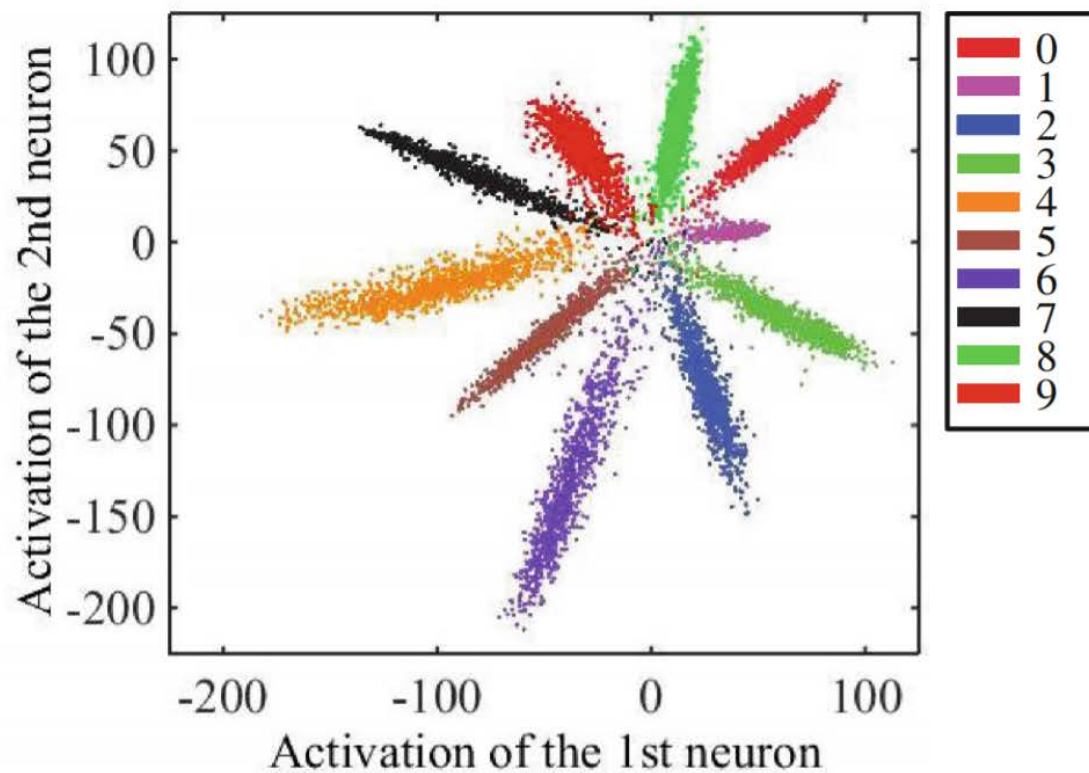
# DISCRIMINATIVE VS SEPARABLE



# CROSS ENTROPY 2D FEATURE SPACE

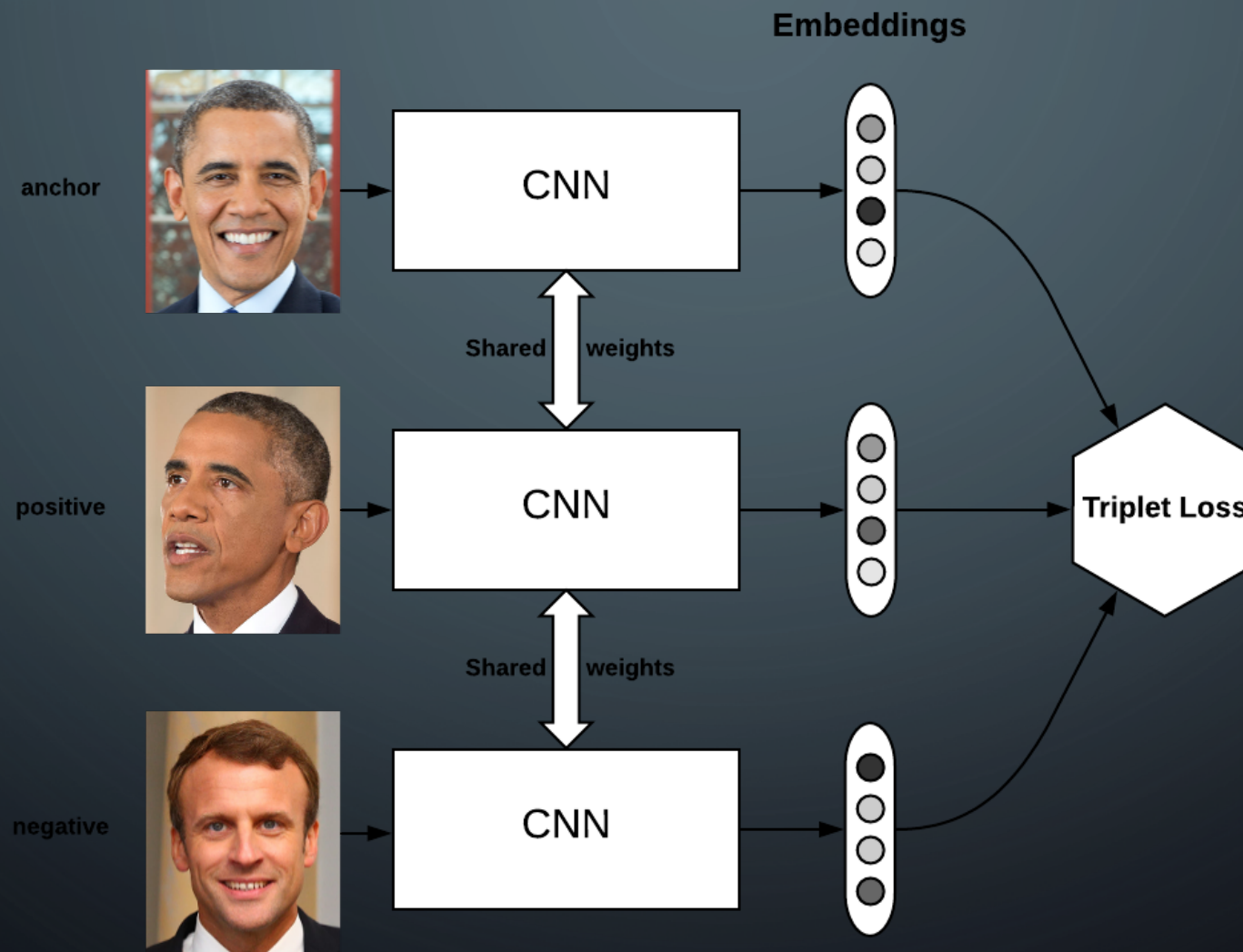


(a)



(b)

# TRIPLET LOSS





# TRIPLET LOSS



- $L(a, p, n) = \max(0, D(a, p) - D(a, n) + \textit{margin})$
- D is distance
- D could be L2 distance
- D could be (1 – cosine similarity)





# PROBLEMS WITH TRIPLET LOSS

- Most triplet has no information
- We can choose hard negative
- Noisy data and mislabeled data
- Computation is time consuming
- We can choose farthest positive and nearest negative from batch
- We need large batch-size (for example 1000)
- Memory bottleneck

# CENTER LOSS



- We can use classification to train model
- Try to hold feature vector of each class near each other
- Force model to shape feature vectors as sphere
- This feature space will be discriminative

# CENTER LOSS



$$\mathcal{L}_S = - \sum_{i=1}^m \log \left( \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}} \right)$$

$$\mathcal{L}_C = \frac{1}{2} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{c}_{y_i}\|_2^2$$



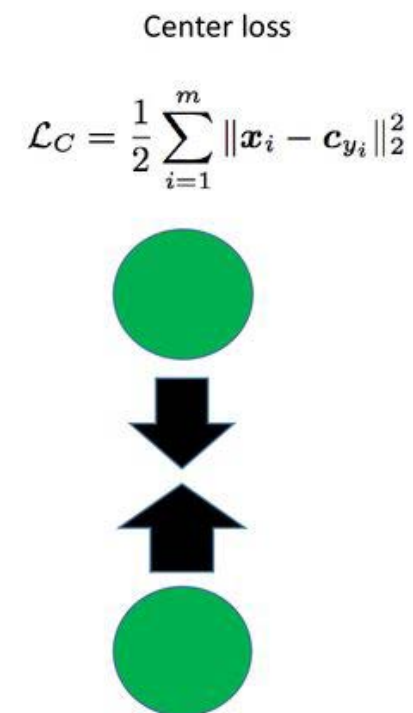
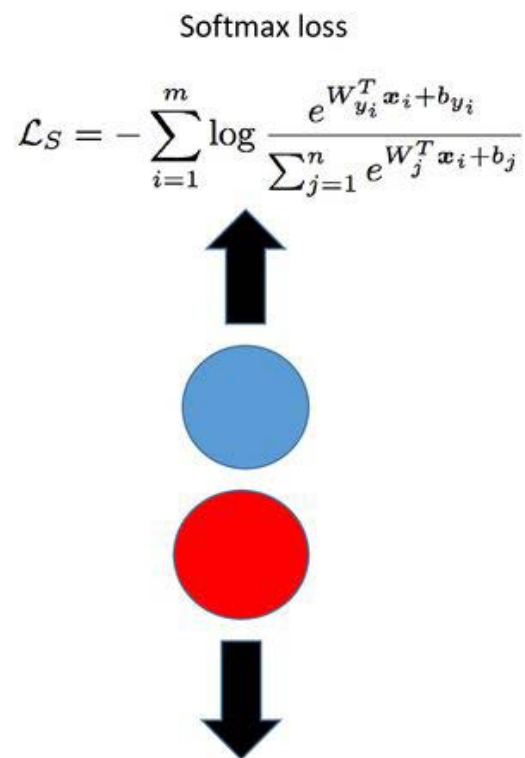
$$\mathcal{L} = \mathcal{L}_S + \lambda \mathcal{L}_C$$

$$= - \sum_{i=1}^m \log \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}} + \frac{\lambda}{2} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{c}_{y_i}\|_2^2$$

# SOFTMAX LOSS VS CENTER LOSS

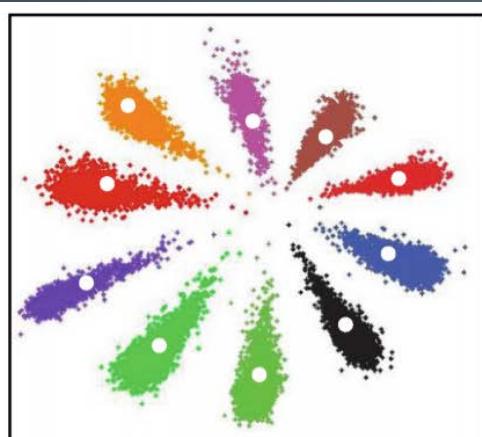


- Just Softmax loss
  - Separable
  - not discriminative
- Just Center loss
  - all feature vectors will be same
- Softmax loss => between-class variation
- Center loss => inter-class similarity

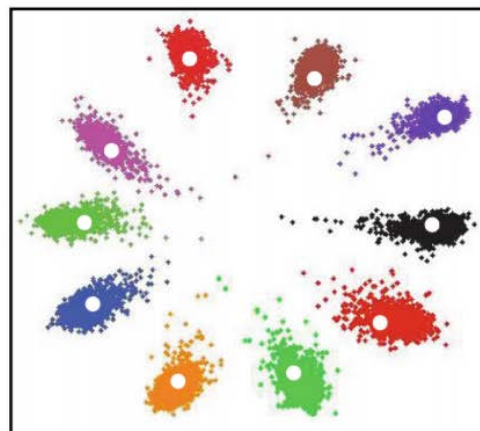




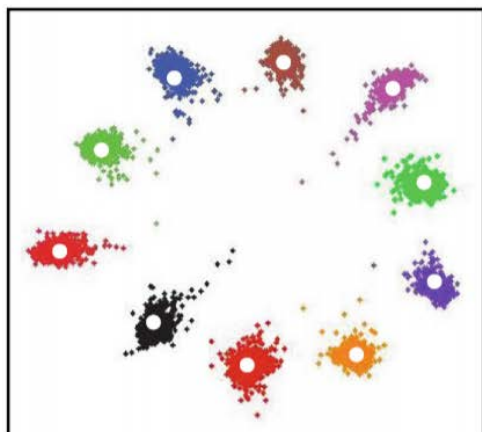
# CENTER LOSS + CROSS ENTROPY FEATURE SPACE



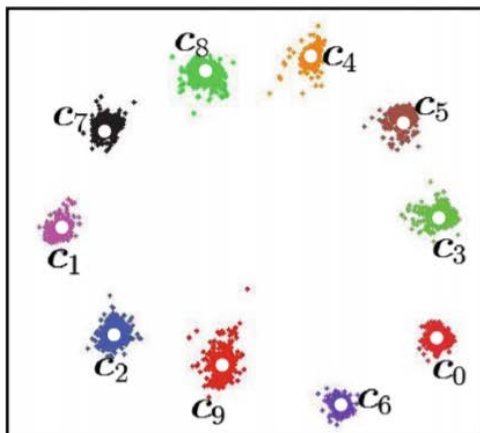
(a)  $\lambda = 0.001$



(b)  $\lambda = 0.01$



(c)  $\lambda = 0.1$



(d)  $\lambda = 1$

