

# Convolutional Neural Networks

\* Edge detection →

$$\begin{array}{|c c c|} \hline 3 & 0 & 1 \\ \hline 1 & 5 & 8 \\ \hline 2 & 7 & 2 \\ \hline \end{array}
 \begin{array}{|c c c|} \hline 2 & 7 & 4 \\ \hline 9 & 3 & 1 \\ \hline 5 & 1 & 3 \\ \hline \end{array}
 \begin{array}{|c c c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}
 = \begin{array}{|c c c c|} \hline -5 & -4 & 0 & 8 \\ \hline -10 & -2 & 2 & 3 \\ \hline 0 & -2 & -4 & -7 \\ \hline -3 & -2 & -3 & -16 \\ \hline \end{array}$$

Image  $\quad * \quad$  Filter  $\quad$  Vertical

$a_{11} = 3 \cdot 1 + 1 \cdot 1 + 2 \cdot 1$   
 $+ 0 \cdot 0 + 5 \cdot 0 + 7 \cdot 0$   
 $+ 1 \cdot (-1) + 8 \cdot (-1) + 2 \cdot (-1)$

Horizontal  
Filter  $3 \times 3$

$$\begin{matrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{matrix}$$

Sheard Vert  
Filter  $3 \times 3$

$$\begin{matrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{matrix}$$

Scharr Vert  
Filter  $3 \times 3$

$$\begin{matrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{matrix}$$

Learned Plane  
NN Edge detect.

$$w_1, w_2, w_3$$

$$w_4, w_5, w_6$$

$$w_7, w_8, w_9$$

Not necessarily  $0^\circ, 90^\circ$

but also  $45^\circ, 70^\circ, 85^\circ$

\* Padding →

sym vert vs horiz.

- If padding is not used, pixel at borders are not used as much.

- Prevents image shrinking, also keeps information of borders.

- Valid convolution: "No padding"  $\rightarrow n \times n * f \times f \rightarrow (n-f+1) \times (n-f+1)$

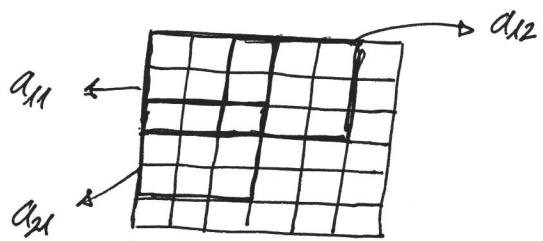
- Same resolution  $\rightarrow$  Out size = In size.

$$(n+2p-f+1) \times (n+2p-f+1) \rightarrow n+2p-f+1 = n \rightarrow p = \frac{f-1}{2}$$

$f$  = usually odd number.

## \* Stride $\rightarrow$

Stride = 2



Receptive dimension  $\rightarrow$

$$\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

$$\lfloor z \rfloor = \text{floor}(z)$$

$\hookrightarrow$  If the convolution doesn't fit the matrix, convolution is not to do.

## \* Mathematical convolution $\rightarrow$

A. The filter is flipped horizontally.

B. In here we skip this step, which technically is cross-correlation.

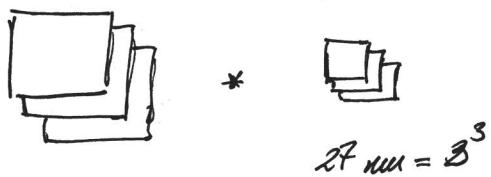
## \* Convolution over volume $\rightarrow$

RGB images

$6 \times 6 \times 3$

$3 \times 3 \times 3$

$4 \times 4$



$$a_{ii} = \text{Red conv} + \text{Green conv} + \text{Blue conv}$$

$$27 \text{ num} = 3^3$$

## - Multiple filters $\rightarrow$

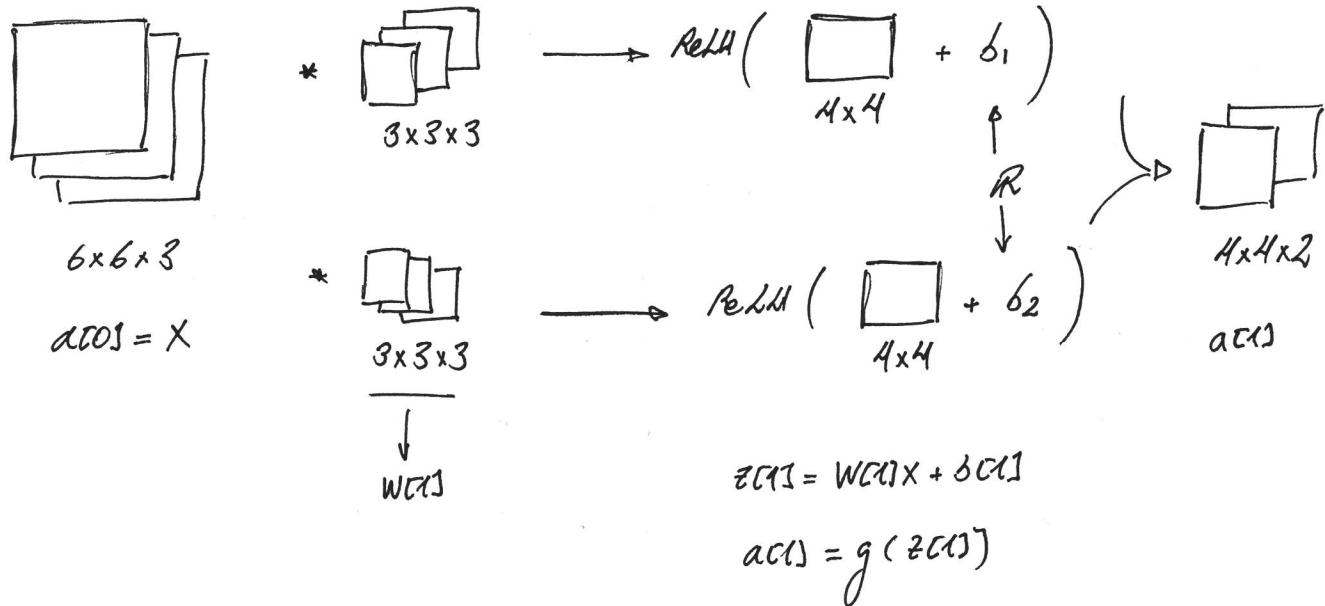
$$6 \times 6 \times 3 * 3 \times 3 \times 3 = 4 \times 4$$

$$* 3 \times 3 \times 3 = 4 \times 4$$

Number of filters.

$4 \times 4 \times 2$

- One layer CNN →



2 filters ; each  $3 \times 3 \times 3 \rightarrow$  Parameters?

$$3 \times 3 \times 3 + \text{bias/filter} = 28 \text{ param/filter} \Rightarrow 28 \cdot 2 = 56 \text{ parameters.}$$

\* Notation →

- layer  $l$  = convolutional layer.
- $f_{CL}$  = filter size
- $p_{CL}$  = padding
- $s_{CL}$  = stride
- $n_{CL}$  = # filters.

Filter →

$$f_{CL} \times f_{CL} \times n_{CL-1}$$

Input →

$$n_{HCL-1} \times n_{WCL-1} \times n_{CCL-1}$$

Output →

$$n_{HCL} \times n_{WCL} \times n_{CCL}$$

$$n_{HCL} = n_{WCL} = \left\lfloor \frac{n_{CL-1} + p_{CL} - f_{CL}}{s_{CL}} + 1 \right\rfloor$$

Weights →

$$f_{CL} \times f_{CL} \times n_{CL-1} \times n_{CCL}$$

Activations →

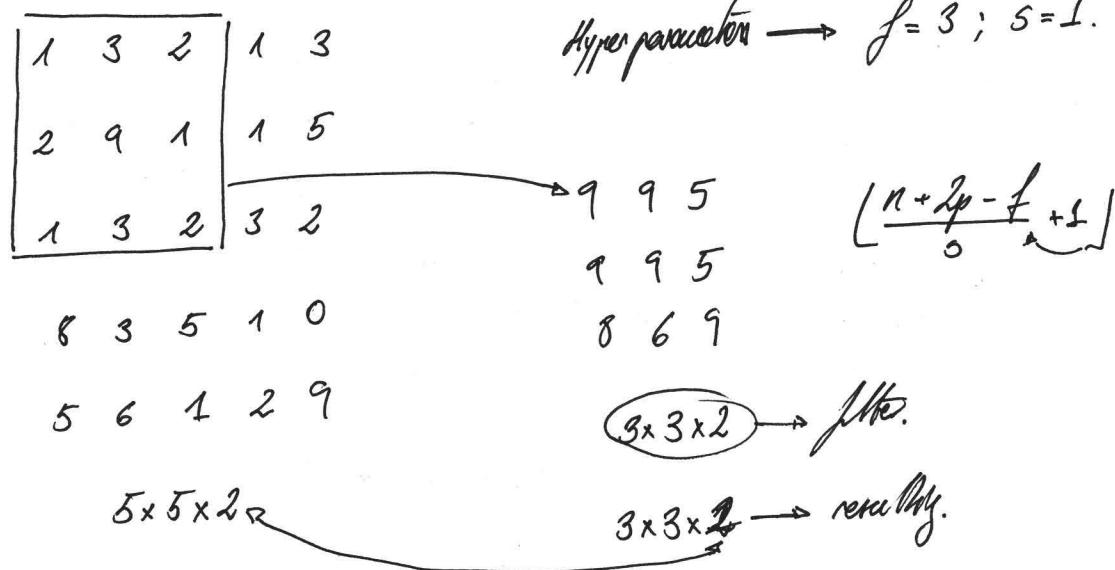
$$act \rightarrow n_{HCL} \times n_{WCL} \times n_{CCL}$$

$$ACCL \rightarrow n_{HCL} \times n_{WCL} \times n_{CCCL} \times M$$

Bias →

$$n_{CCCL}$$

- Pooling layer →

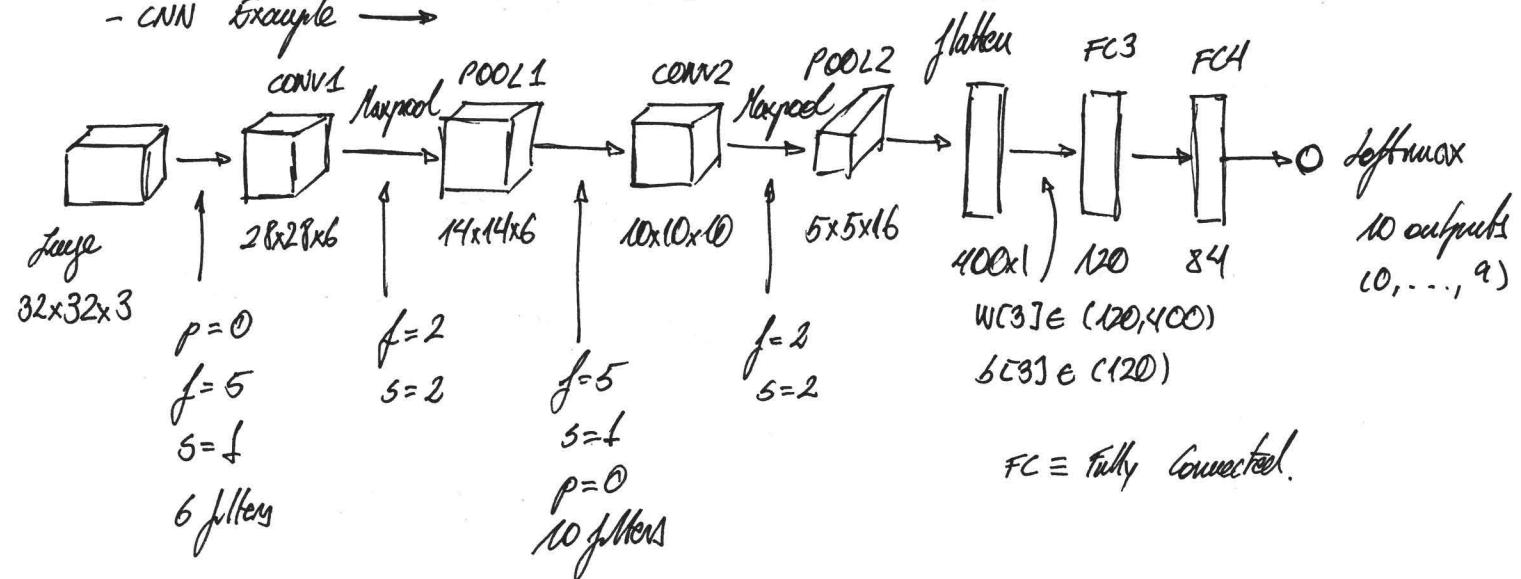


- \* Good for detecting distinctive features.
- \* No parameters to learn. → But it affects backpropagation calculation.
- \* Average pooling → Take average instead of max.

$$n_H \times n_W \times n_C \rightarrow \left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \times n_C$$

$\left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \quad \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor$

- CNN Example →



\*  $n_H, n_W, n_C$

\*  $n_C \uparrow$

	Activation shape	Activation size	# parameters.
Input	(32, 32, 3)	3,072	0
conv1	(28, 28, 8)	6,272	208
Pool 1	(14, 14, 8)	1,561	0
Conv 2	(10, 10, 16)	1600	416
Pool 2	(5, 5, 16)	400	0
FC 3	(120, 1)	120	48,001
FC 4	(84, 1)	84	10,081
Softmax	(10, 1)	10	841

↓ Decreases.

→ Most of the parameters.

- Parameter sharing →

- \* Convolutional NN farms reduce the total number of parameters, reducing overfitting.

- \* It also allows a feature detector to be used in multiple locations throughout the whole input image/Volume.

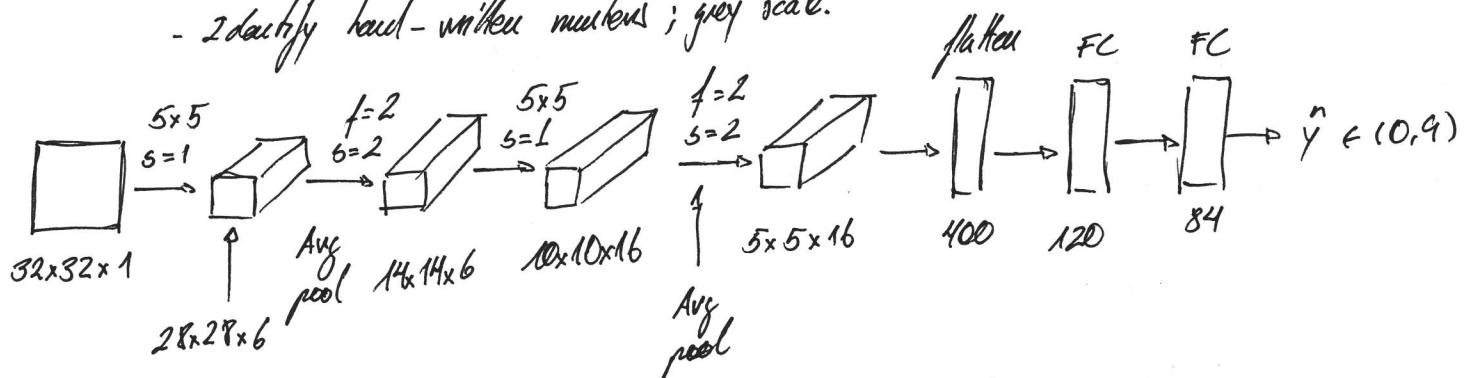
- Sparsity of connections →

- \* Each activation in the next layer depends only a small number of activations from the prev. layer.

- Case Studies →

+ de-Net 5 →

- 2 classify hand-written numbers; grey scale.



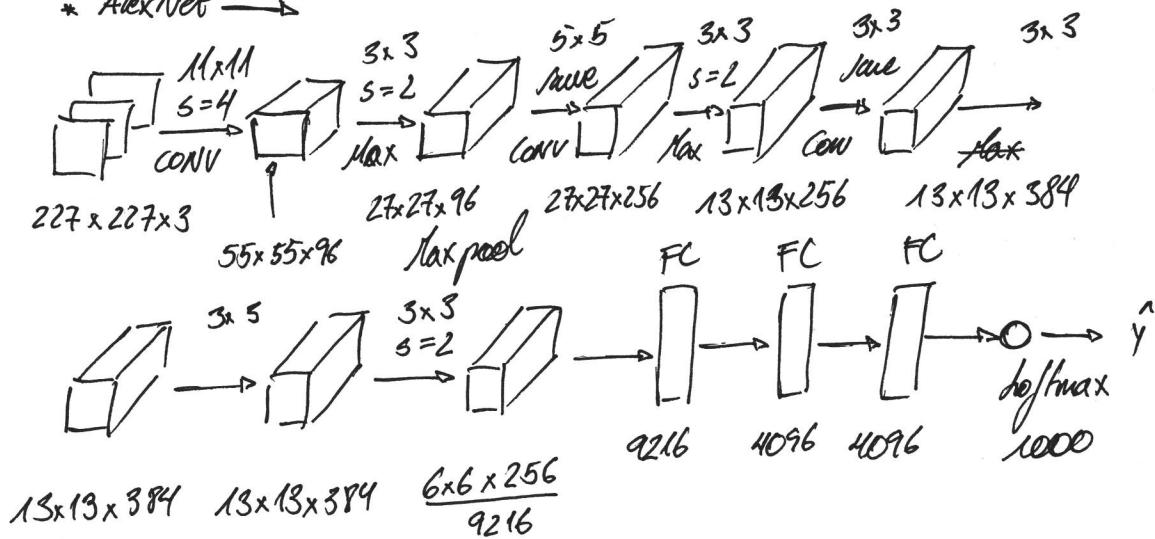
- 60K parameters. → much compare to now (10M).

-  $n_H, n_W \downarrow$ ;  $K_c \uparrow$

- One or more → conv + pool

- Activation function → sigmoid / tanh  $\xrightarrow{\text{ReLU (now)}}$

\* AlexNet →



-  $16 \cdot 10^6$  parameters.

- Similar last layer to Net.

- Activation  $\rightarrow$  ReLU.

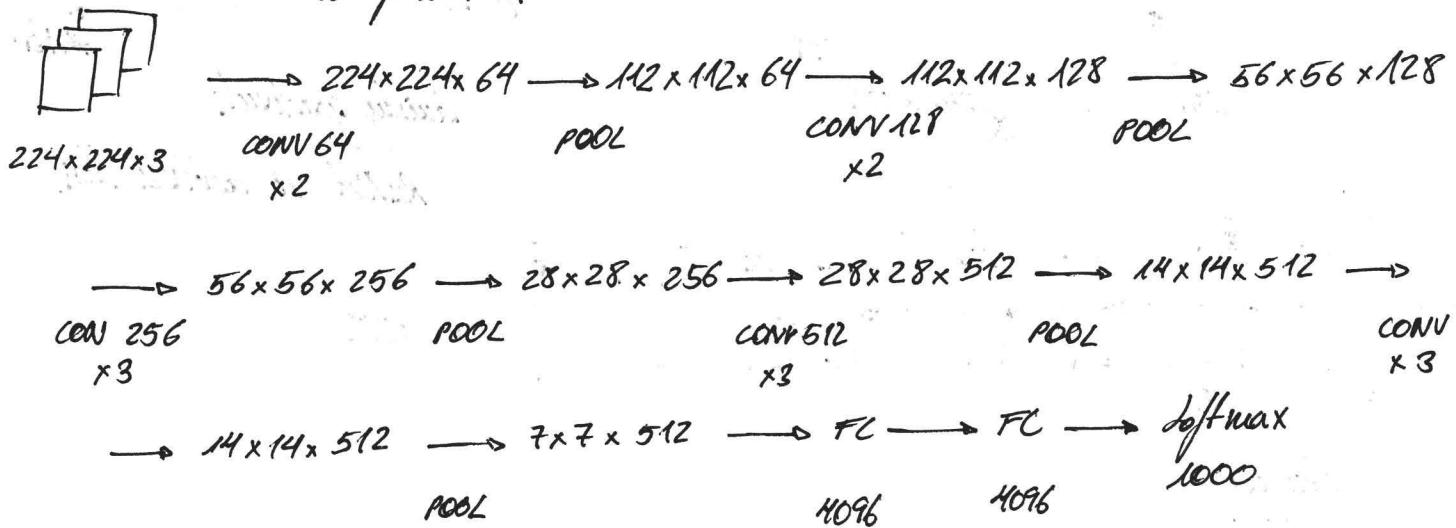
\* VGG - 16 →

- 16 layers with weights.

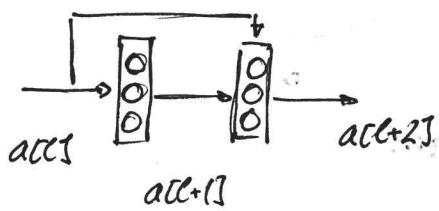
- Conv →  $3 \times 3$  filter,  $s=1$ , same

- Max pool →  $2 \times 2$ ,  $s=2$

-  $138 \cdot 10^6$  parameters.



- ResNet = Residual Network →



$$accl+2 = g(accl+1) + accl$$

$$accl+2 = W[accl+2] accl+1 + b[accl+2] =$$

$$= W[accl+2] g(W[accl+1] accl + b[accl+1]) + b[accl+2]$$

(Other layer skips 3)

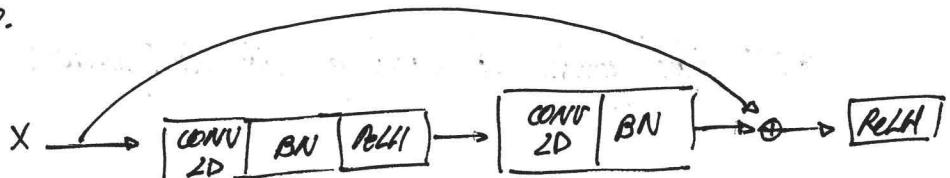
\* if we choose  $W[accl+2] = b[accl+2] = 0 \rightarrow accl+2 = g(accl)$

- easy for the ResNet to learn the identity function.

- it doesn't impact much the performance.

- it is not more difficult to learn the identity function

as deep as you go.

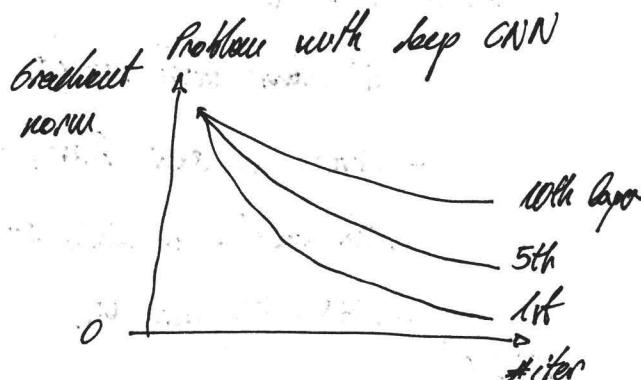


- Assumption  $z_{[l+2]} \& a_{[l]}$  have same dimension  $\rightarrow$  same convolution.

- In any other case  $\rightarrow$

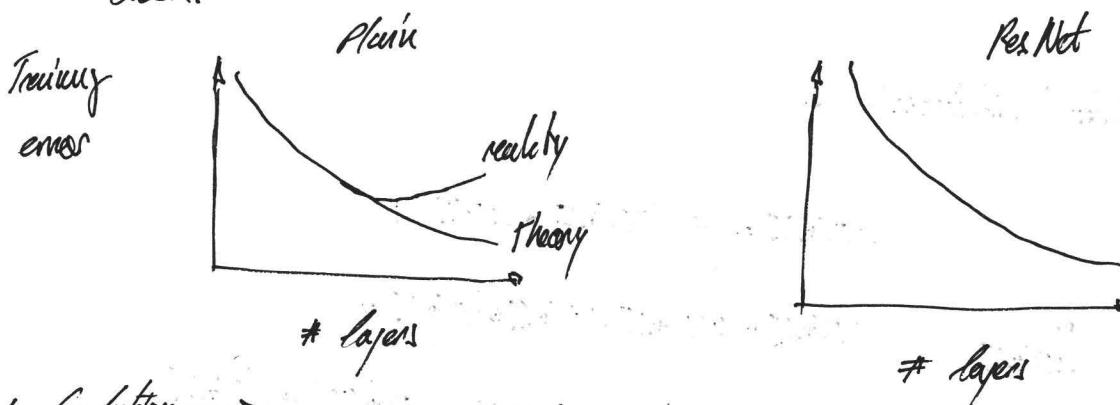
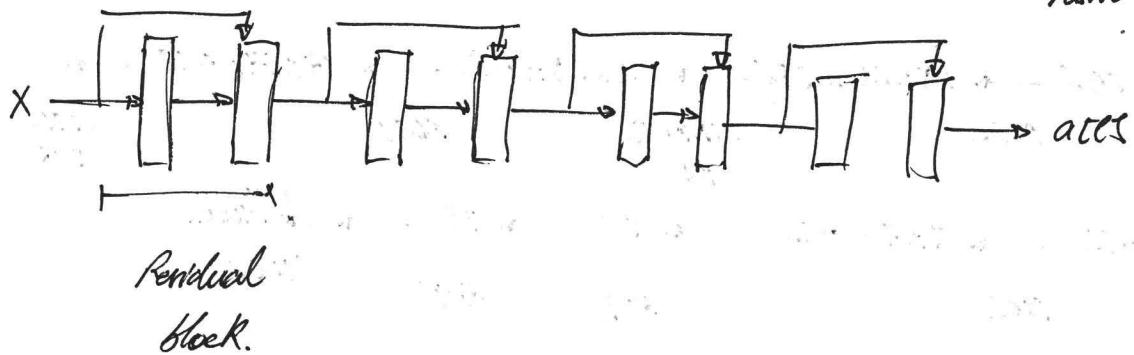
$$g(z_{[l+2]} + w_s a_{[l]})$$

$256$                            $R^{256 \times 128}$                            $128$

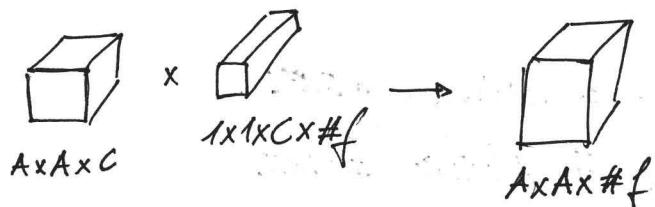


Vanishing gradient.

ResNet to address this.



- 1x1 convolution  $\rightarrow$



ReLU (but 1x1 per each of the channels).

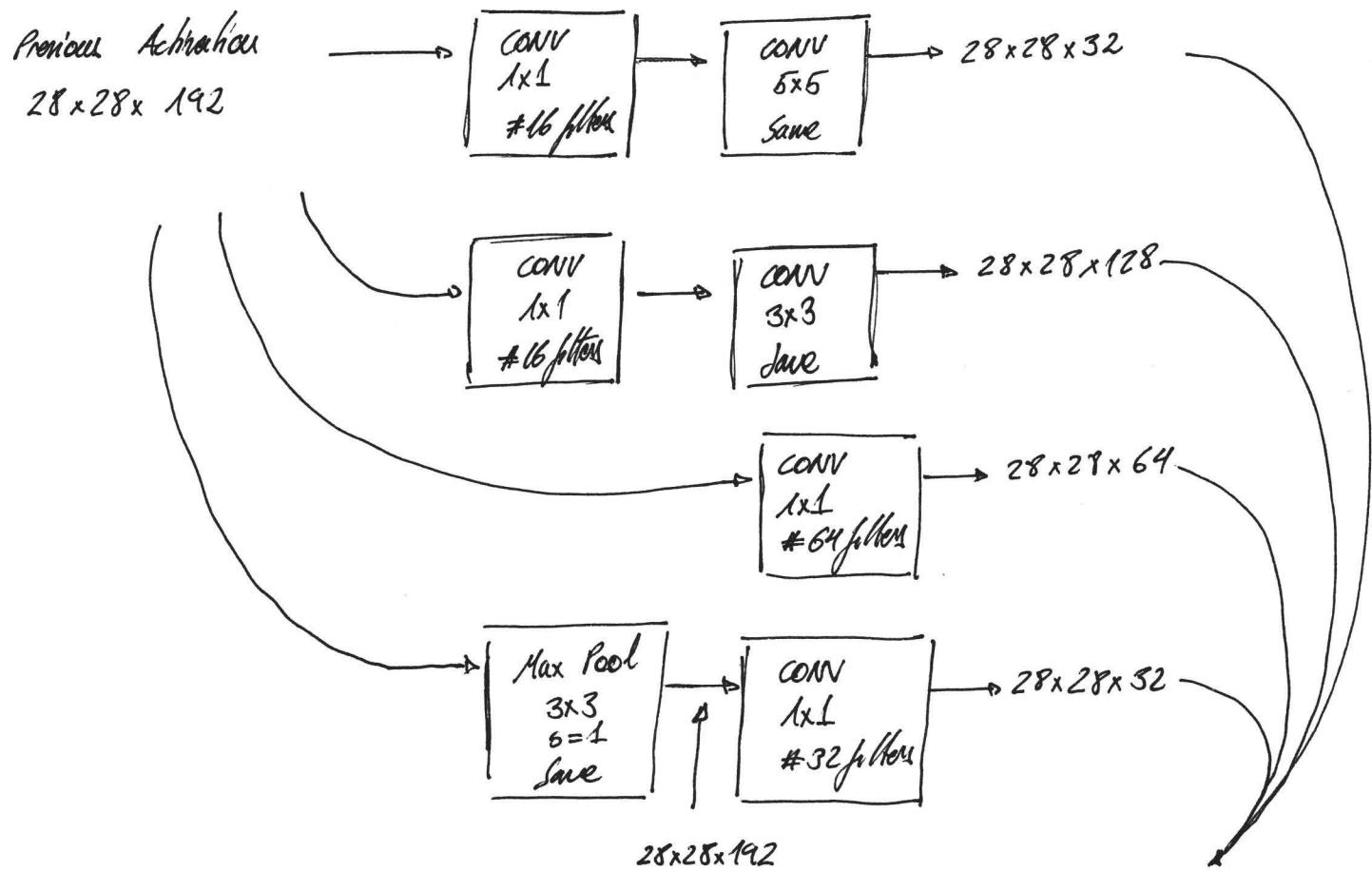
\* Keeps  $N_H$ ,  $N_W$  and can shrink # channels.

\* If you keep # channels, it allows you to add non-linearity to the model.

\* FC that applies to each of the  $A \times A$  patches, by the # filters.

- Inception Network →

\* Inception Block →



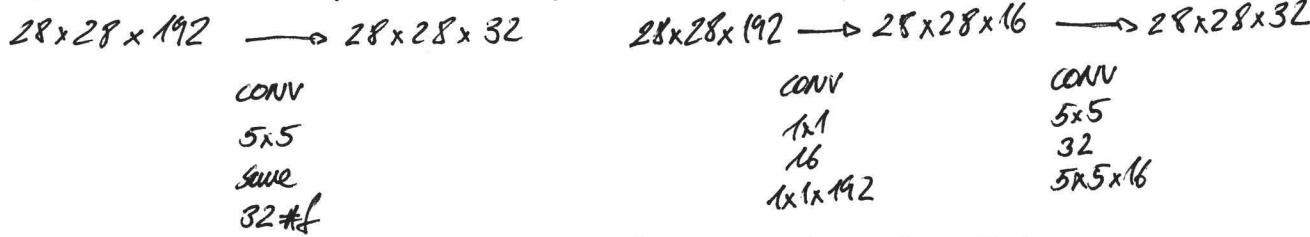
\* Inception network as a NN based on these.

$$28 \times 28 \times (128 + 64 + 32 + 32) \times 256$$

\* This network may have different branches with softmax →  
helps predict an output cost and prevent overfitting.

\* GoogLeNet

\* Computation cost very and not very  $1 \times 1 \rightarrow$



$$(28 \times 28 \times 16) \times (192) = 2.4M$$

$$(28 \times 28 \times 192) \times (5 \times 5 \times 192) = 120M$$

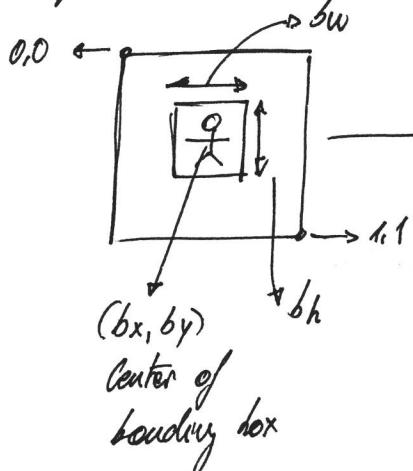
$$12.4M$$

$$(28 \times 28 \times 32) \times (5 \times 5 \times 16) = 104M$$



## Detection Algorithms

### \* Object detection →



- If we define the loss function as squared error →

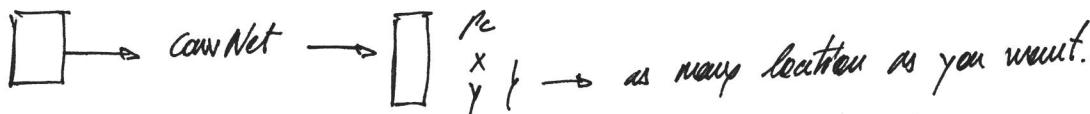
$$\text{If } y_i = 1 \rightarrow \sum_i^8 (y_i - \hat{y}_i)^2$$

$$\text{if } y_i = 0 \rightarrow (y_i - \hat{y}_i)^2$$

### \* Landmark detection →

- Use a NN to extract  $x, y$  locations in an image.

- Basic node for face recognition → Need to label training set by hand.



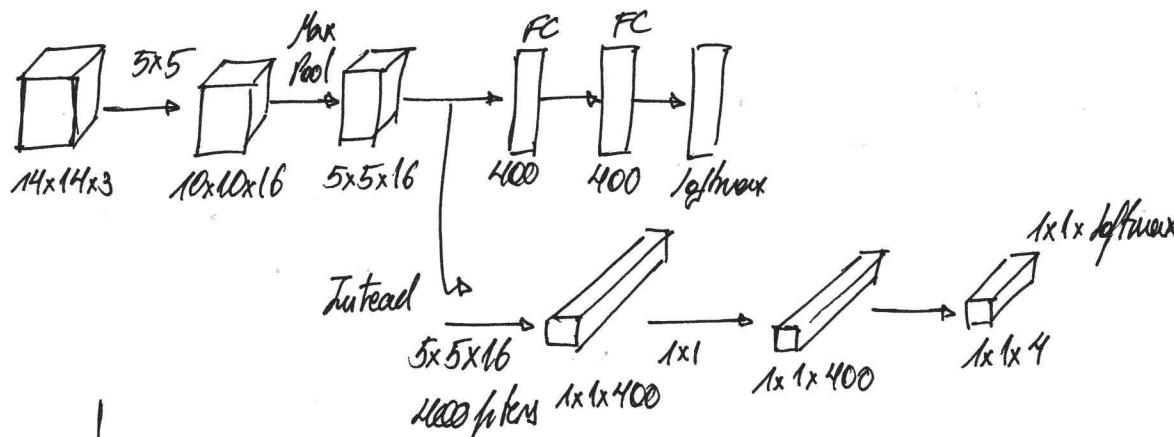
### \* Object detection →

- On training set feed done cropped images, mostly large of object.

- Then we do full image detection → This has a huge computational cost.

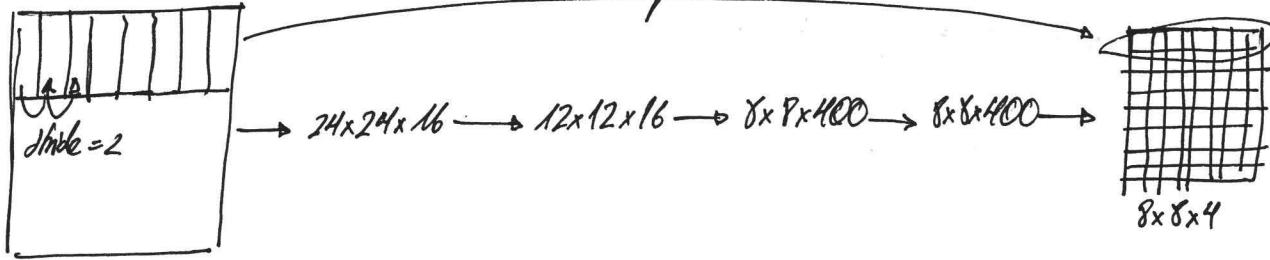
\* Conventional implementation of visibility windows →

- 1st idea → Turn FC into a conv layer.



→ Memory to a big nudge

save correspondence.



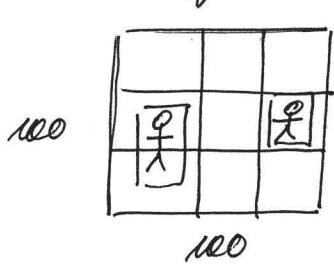
28x28x3

\* Most of the calculation has already been done, highly repetitive.

\* Instead process the whole image.

\* Doing box prediction →

- YOLO Algorithm →



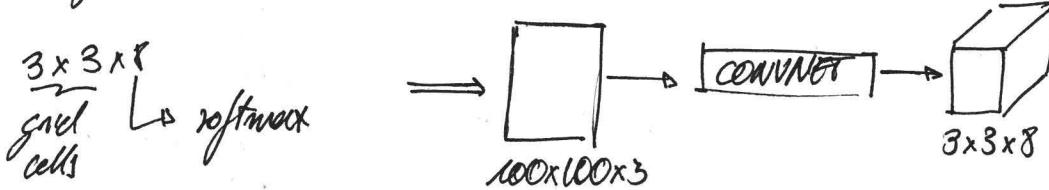
\* Divide image in grid of cells and object localisation for each cell

\* if anyone objects to a cell

\* for each grid cell label  $y$  →

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_H \\ b_W \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

- \* do the resulting boxes should be →



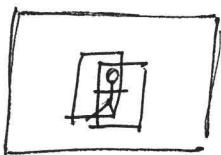
- \* Encoding the bounding boxes →

$$y = \begin{bmatrix} 1 \\ bx \\ by \\ bw \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$b_w = 0.5, b_H = 1.1$  can be greater like in this case.  
 $\left. \begin{array}{l} b_x \approx 0.5 \\ b_y \approx 0.1 \end{array} \right\} \rightarrow \text{always between } 0.1$

- Intersection over Union →

- \* used in YOLO as a way to evaluate object detection.



\* For each of the anchor boxes, it will be used with the one of highest prob → Recalling all boxes that have a large overlap.



Union      Intersection  $\rightarrow I_{OU} = \frac{\text{Area Intersection}}{\text{Area Union}}$

If  $I_{OU} > \text{iou\_threshold} \rightarrow \text{Recall box}$

- Non-max Suppression →

- \* scenario → one object with several bounding boxes

- \* Algorithm →

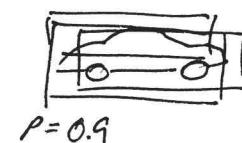
1. Takes the bbox with highest prob.

2. For each remaining →

- calculates  $I_{OU}$ , discard if  $>$  threshold

- \* each cell of the grid → ~~(one object)~~ output  $(p_c, b_x, b_y, b_H, b_w)$

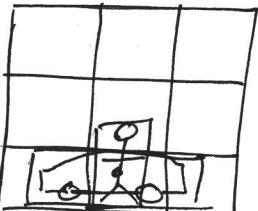
$P = 0.6$



$P = 0.9$

- Anchor boxes →

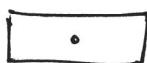
\* Overlapping objects →



Anchor box 1



Anchor box 2



center both objects  
in same grid cell

\* Now each grid can have more than  
one output associated.

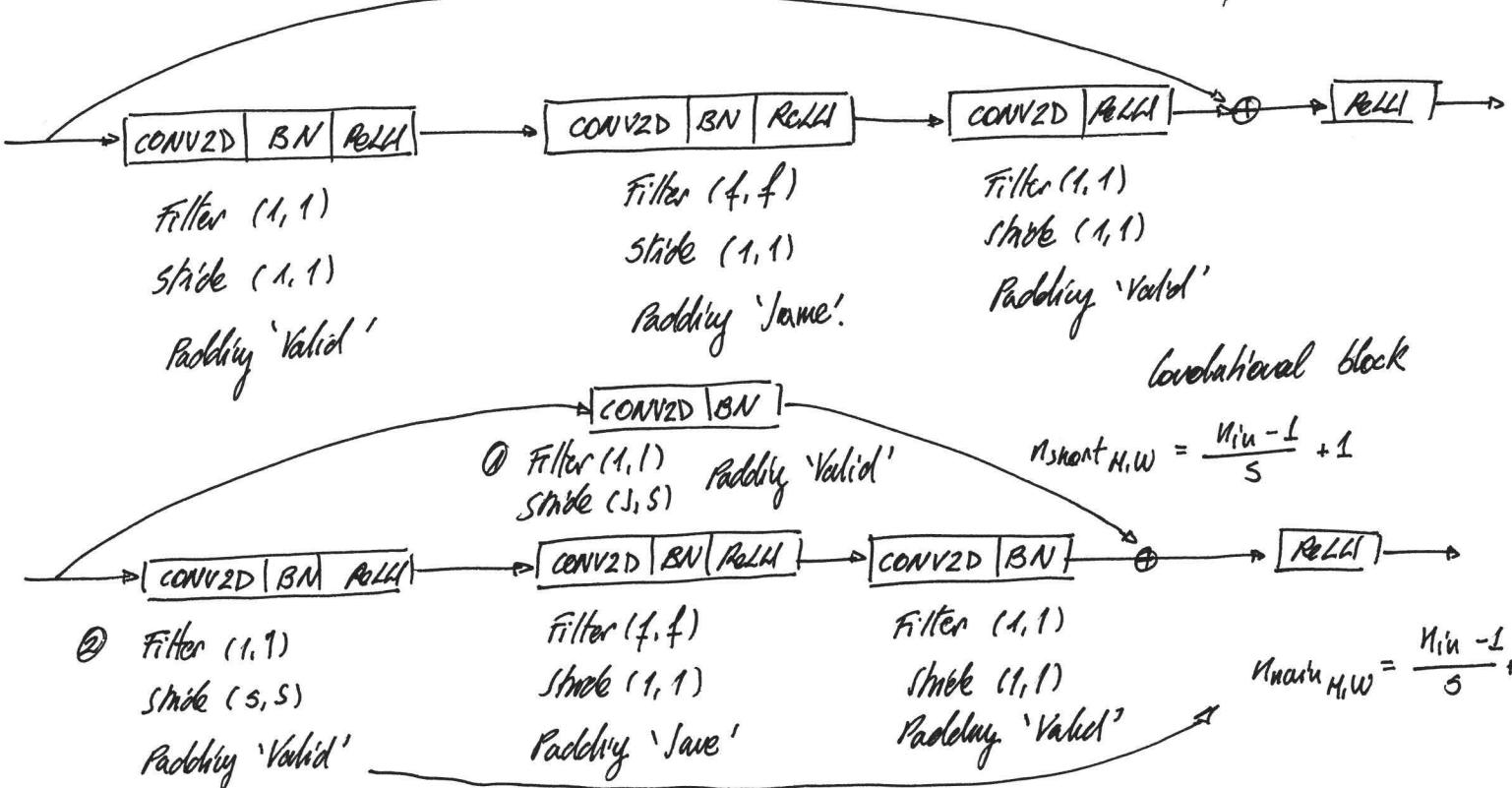
\* Only works if 2 different object, each associated to different anchor boxes.

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_H \\ b_W \\ C_1 \\ C_2 \\ C_3 \\ p_{c'} \\ b'_{x} \\ b'_{y} \\ b'_{H} \\ b'_{W} \\ C'_1 \\ C'_2 \\ C'_3 \end{bmatrix}$$

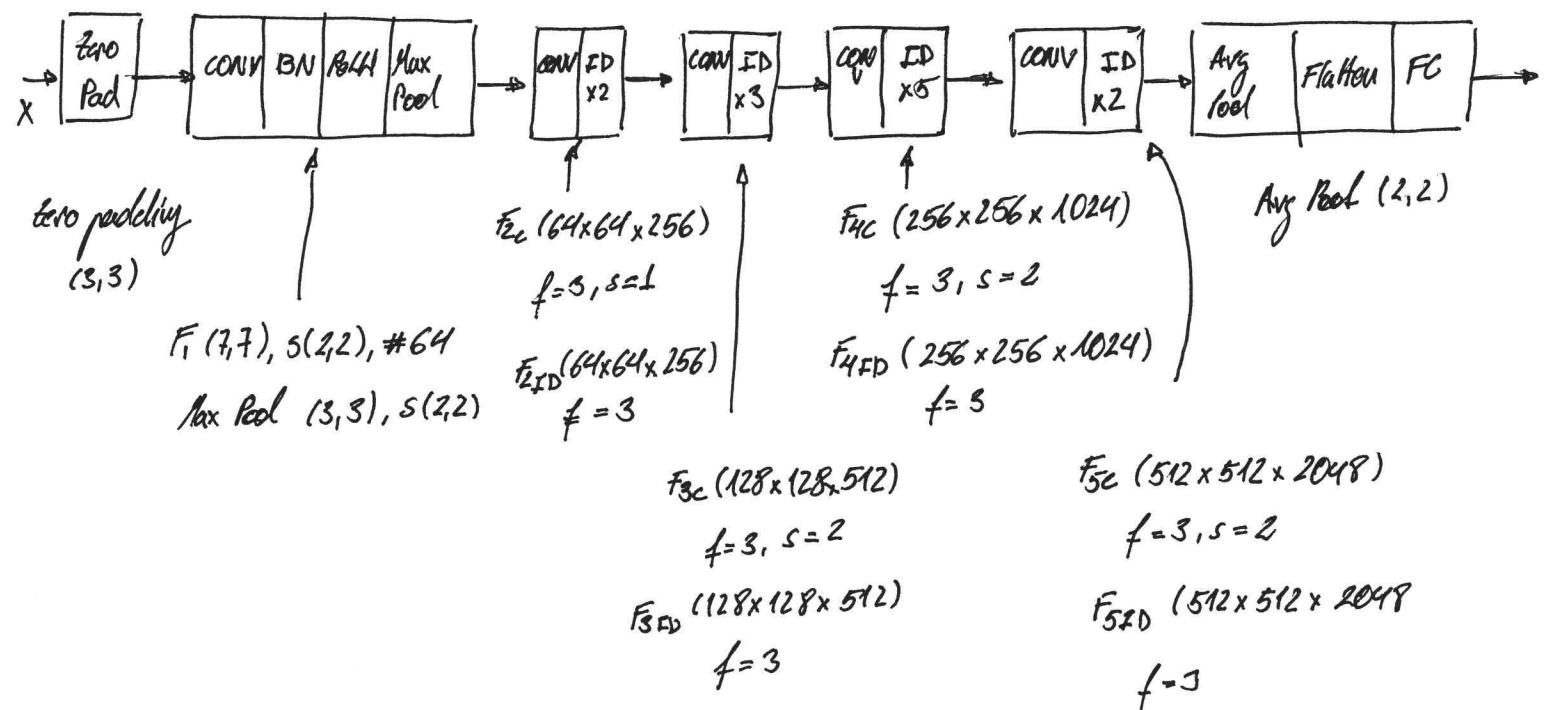
→ Can provide  
information  
of 2  
pedestrian & Car



## Identity block



② Used to match up the addition, since filter path ② changes shape.



ResNet 50 layers.



## Face Recognition

\* Verification vs recognition →

### Verification

Input image vs name,  $ID \Rightarrow 1:1$

Output whether the image is of that person

### Recognition

DB with  $K$  persons

Get input image  $\xrightarrow{1:K}$

More prob of making mistake.

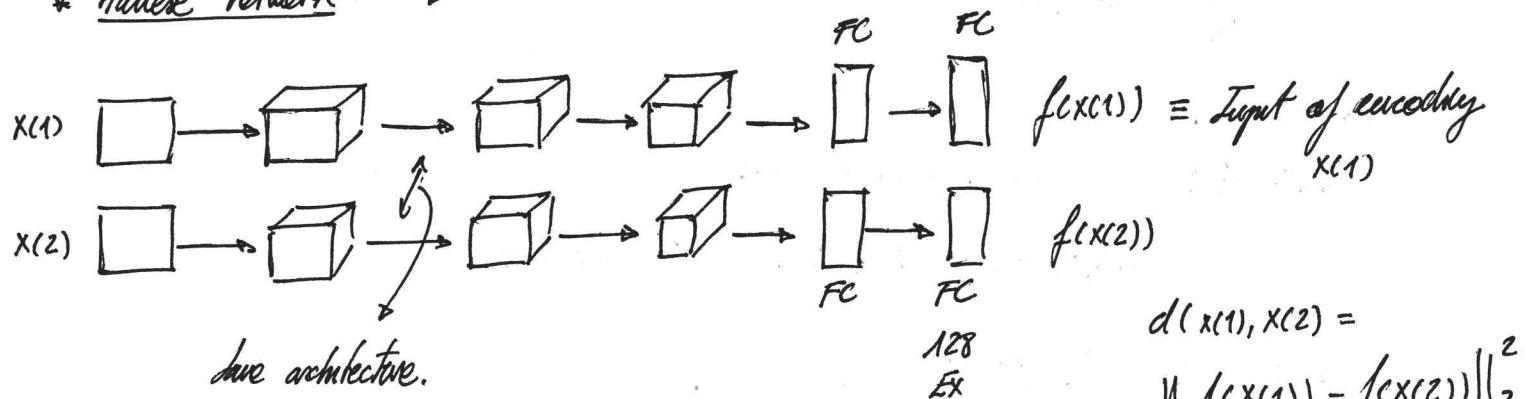
\* One shot learning →

- Objective → Learn from one instance to recognise the person again.
- Use a similarity function →

$d(\text{img1}, \text{img2}) \equiv$  degree of difference between images

$$d(\text{img1}, \text{img2}) \leq \epsilon$$

\* Feature network →



$$d(x(1), x(2)) = \|f(x(1)) - f(x(2))\|_2^2$$

\* Train network so it encodes the similarity function.

\* Parameters of NN define the encoding. →

If  $x(i), x(j)$  are same person →  $\|f(x(i)) - f(x(j))\|_2^2$  is small.

\* Paper → DeepFace, closing the gap to human level performance.

\* Triplet loss function →

- Uses Anchor, Positive, Negative →

$\| \cdot \|_2^2 \rightarrow$  square Euclidean distance

$$\frac{\|f(A) - f(P)\|_2^2}{d(A, P)} \leq \frac{\|f(A) - f(N)\|_2^2}{d(A, N)} \quad L = (f(A_1) - f(P_1))^2 + (f(A_2) - f(P_2))^2 + \dots$$

- To prevent the network to just  $f(\text{image}) = \vec{0} \rightarrow$

$\alpha = \text{margin}$ ; forces the network to learn to fix a gap between P, N.

$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

The NN shouldn't care how much negative it is for no match cases, force 0.

- Training set → 10K pictures of 1K persons.

$$J = \sum_{i=1}^M L(A^{(i)}, P^{(i)}, N^{(i)})$$

- Triplets →

\* If  $A, P, N$  are chosen randomly →

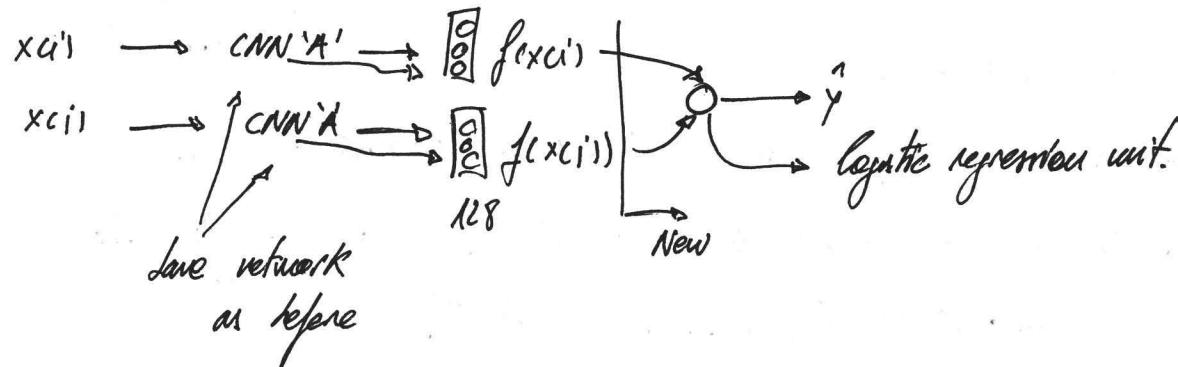
$d(A, P) + \alpha \leq d(A, N)$  is easily satisfied.

+ There triplets 'hard' to train on →

$d(A, P) \geq d(A, N) \rightarrow$  Network forced to learn this case.

## Face verification & Binary classification →

- Another way of learning the similarity function.



$$\hat{y} = \tau \left( \underbrace{\sum_{k=1}^{128} w_i |f(x_{ci})_k - f(x_{cj})_k| + b}_{\frac{(f(x_{ci}))_k - (f(x_{cj}))_k)^2}{f(x_{ci})_k + f(x_{cj})_k}} \right)$$

- If the 2nd stage belongs to a db, you can have those already computed.  $f(x_{cj})$

## Neural style Transfer →

### \* Deep ConvNet learning →

Image →  $110 \times 110 \times 96 \rightarrow 55 \times 55 \times 96 \rightarrow 26 \times 26 \times 256 \rightarrow 13 \times 13 \times 384 \rightarrow 224 \times 224 \times 3$

→  $13 \times 13 \times 384 \rightarrow 6 \times 6 \times 256 \rightarrow 4096 \xrightarrow{\text{FC}} 4096 \xrightarrow{\text{FC}} \hat{y}$  AlexNet like network.

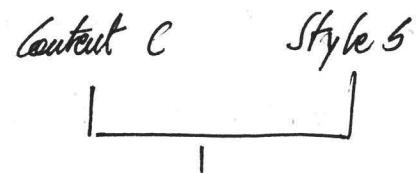
- Paper → Visualizing and understanding convolutional networks. 2013, Zeiler and Fergus.

- When you visualize activation function outputs →

- \* On the first layer, it look for edges, small features.
- \* As you go deeper it identifies larger portions.

\* loss function  $\rightarrow$

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$



- Find generated  $G \rightarrow$

1. Initialize randomly  $\rightarrow G: 100 \times 100 \times 3$  (Or whatever shape).

2. Use gradient descent to minimize  $J(G) \rightarrow G = G - \frac{\lambda}{2G} J(G)$

Update on  $G$  values (Image output).

- Content cost function  $\rightarrow$

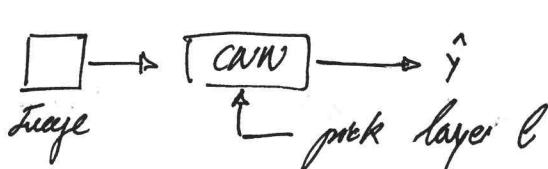
\* Use hidden layer  $l$  to compute cost  $\rightarrow$  neither shaller, too deep.

\* Pretrained CNN  $\rightarrow$  E.g. VGG

\* Let  $\text{act}(c)$  and  $\text{act}(G)$  be the activation function of layer  $l$  activations

$$J_{\text{content}}(C, G) = \frac{1}{2} \| \underbrace{\text{act}(c)}_{\text{activations}} - \text{act}(G) \|_2^2 \rightarrow \text{Measures how similar the a few element-wise, unroll them into vectors.}$$

- Style cost function  $\rightarrow$



layer  $l$   
 $n_H$   $n_W$   $n_C$   $\rightarrow$  more likely to find correlated areas  
in later layers. (Easier).

\* Define the correlation between activation from different channels.

\* Style matrix  $\rightarrow a_{ijk, lcl} = \text{Activation at } (i, j, k). G_{ll} \text{ is } n_{CCS} \times n_{CCS}$

$$G_{KK, CCS} = \sum_i^{n_{CCS}} \sum_j^{n_{CCS}} a_{ijk, lcl} a_{ijk, lcl}$$

$N \times 1 \#C$

$G_{KK, CCS} = \text{Measures how correlated are activations in } K \text{ with } K!$

Point of this  $\rightarrow$  see how correlated are the style in the generated map.

\* Construct this matrix for style and generated  $\rightarrow$   
 $G_{KK^T} CCS(S)$ ;  $G_{KK^T} CCS(G)$ ;  $G \equiv$  Gram Matrix.

$$J_{Style}(CCS(S, G)) = \|G CCS(S) - G CCS(G)\|_F^2 = \frac{1}{(2n_W n_C n_K n_{CCS})^2} (G_{KK^T} CCS(S) - G_{KK^T} CCS(G))^2$$

\* For multiple layers  $\rightarrow$

$$J_{Style}(J, G) = \sum_e d_{CCS} J_{Style}(S, G)$$

## Neural style transfer example

- A. Neural style starts from a pretrained network and builds on top  
 B. Paper → VGG-19, VGG with 19 layers.

\* Reuse from model → dictionary with model in python  
 'layer name' = object.

\* feed image → model ['input']. assign (image); image 600x800

\* access activations of particular layer →

res. var (model ['conv2\_1']) → same convention as in prev array.

- c. Our algorithm on this → hyperparameters that control relative weight of content and style ↴

$$J_{content}(C, G) ; J_{style}(S, G) \rightarrow J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$

\* Ensure generated image G match content image C (feeded image).

- choose some layer in between, not too shallow or deep.

- Layer C → VGG forward.

activation function of layer a<sup>(C)</sup>CS

- Layer G → VGG forward.

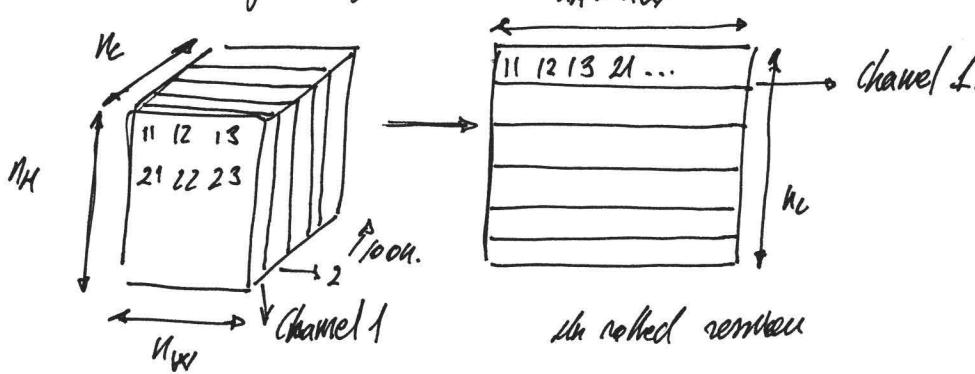
a<sup>(G)</sup>CS

$$J_{content}(C, G) = \frac{1}{N \times H_w \times H_h \times N_C} \sum_{\text{all entries.}} (a^{(C)}_{CS} - a^{(G)}_{CS})^2$$

1. measures how

- Unrolling → good practice.

different aCS(C) &  
aCS(G) are.



un rolled volume

2. Minimizing cost →  
 makes G to have  
 similar content as C.

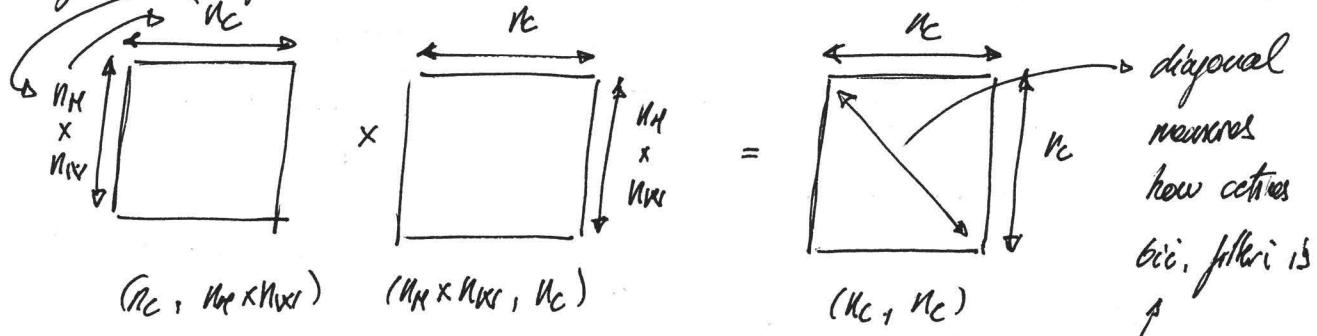
### \* Style Matrix.

- Gram matrix  $\rightarrow$  6 of vectors  $(v_1, \dots, v_n)$  is the matrix of the dot products.

$$G_{ij} = v_i^T v_j = \text{np.dot}(v_i, v_j) \rightarrow G \text{ from Gram, not to take for Generated } G.$$

$\rightarrow$  leaves how similar are the vectors. As similar, as large.

- we generate it from the another matrices,  $G$  and  $G^T$  (transposed of  $G$ ).



- leaves how similar the activation functions are.

### \* Style cost $\rightarrow$

- goal is to minimize the distance between Gram matrix of the style matrix and the generated image.

Gram  
Style  
 $\downarrow$   
+  
Gram  
Generated.  
 $\downarrow$

$$J_{\text{Style TCS}}(S, G) = \frac{1}{4 \cdot n_c^2 \cdot (n_H \cdot n_W)^2} \sum_{i=1}^{n_c} \sum_{j=1}^{n_c} (G_{ij}(S) - G_{ij}(G))^2$$

layer:

### D. Style Weight $\rightarrow$

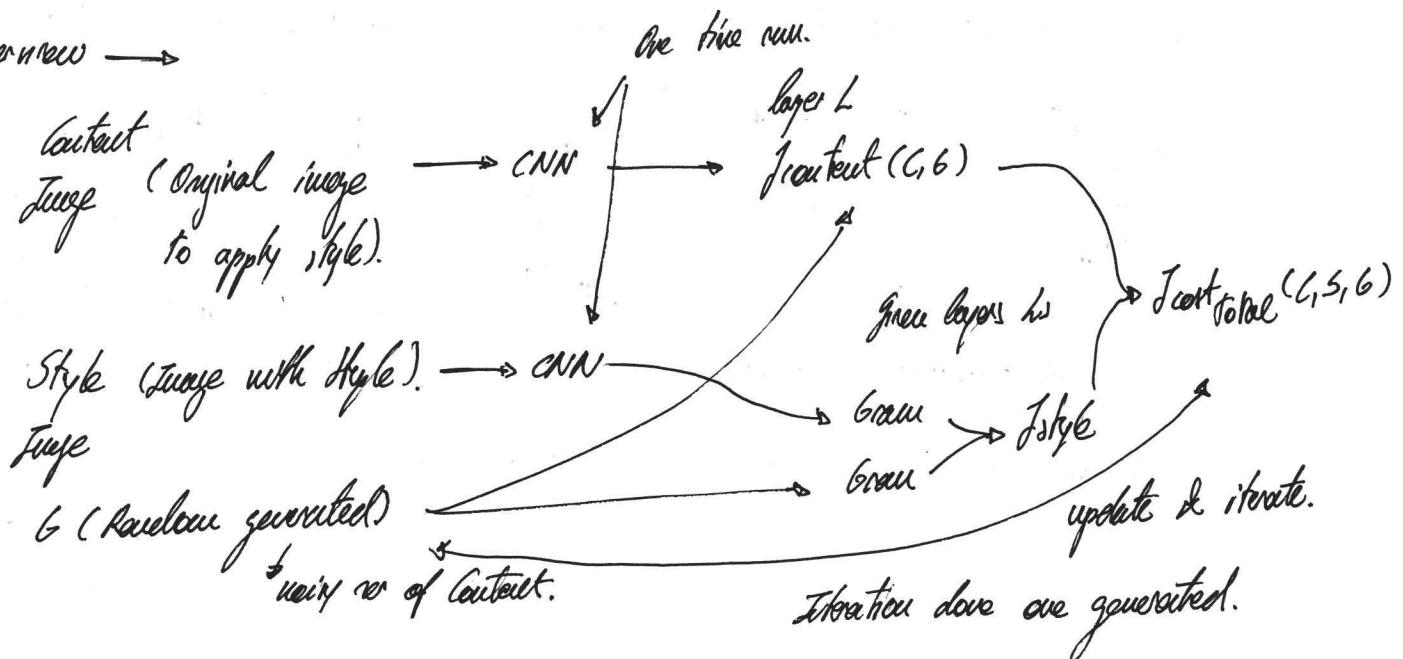
- $\rightarrow$  so far, only for one layer. Better results for several different layers.

$$J_{\text{Style}}(S, G) = \sum_{l} \alpha_l J_{\text{Style TCS}}(S, G_l)$$

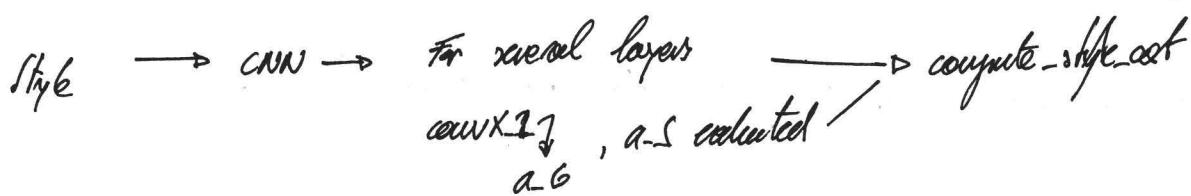
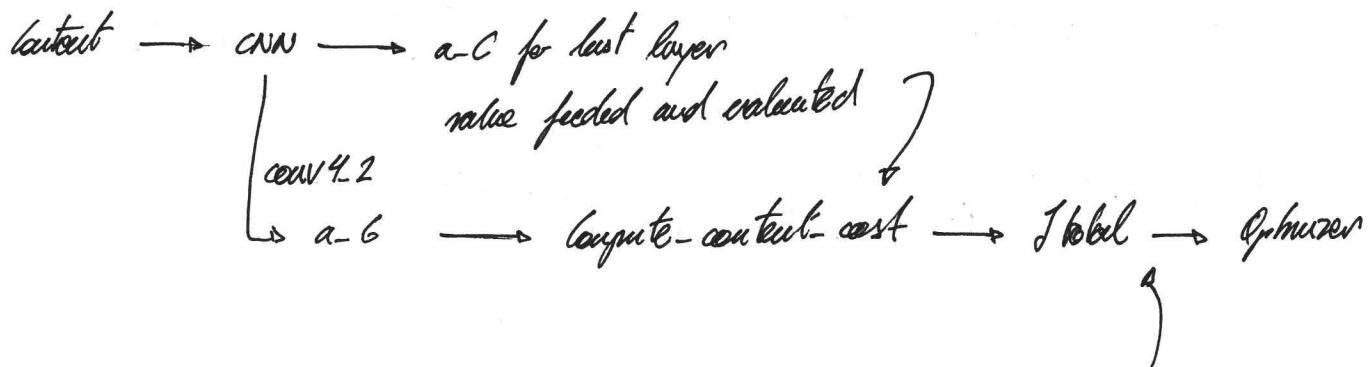
- $\rightarrow$  Any weights per layer

- \* Style of an image can be represented by the gram of a hidden layer function. (using a few gives a better result.)
- \* For the content generation, just one layer is sufficient.
- \* Minimizing the style function  $\rightarrow$  makes it to follow the style of S.

D. Overview  $\rightarrow$



E. Tensorflow graph  $\rightarrow$



Feed converted image  $\rightarrow$  model  $\rightarrow$  Gets values at the layers  $\rightarrow$  **l.n.e.**

for  $X$  iterations  $\rightarrow$  **run(Optimizer)**; runs gradient descent and updates generated.

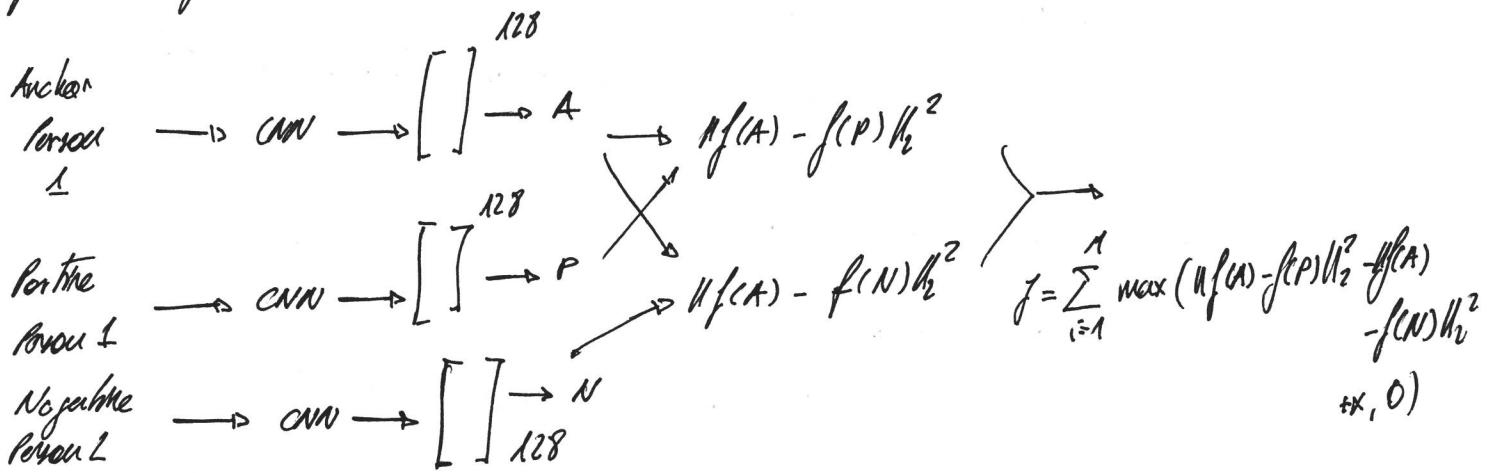
## Face Recognition

\* Using FaceNet model → 128 Bit

$$\text{Face} \rightarrow \begin{bmatrix} \text{CNN} \\ \text{inception model} \end{bmatrix} + \begin{bmatrix} \cdot \end{bmatrix}$$

- This model recognizes features from faces

\* Triplet loss function →



$$\|f(A) - f(P)\|_2^2 = \sum_{i=1}^m (f(A(i)) - f(P(i)))^2 = \text{some euclidean distance.}$$

\* Once define the loss function → Trivially.

- Objective is to pick images that are similar (P, N) so the network is trained on the hard cases.

Notes →

A. Tensorflow conv2d →

padding = 'same' →

$$\text{new\_height} = \text{ceil}(\text{float(in-height)} / \text{float(strides[1])})$$

$$\text{new\_width} = \text{ceil}(\text{float(in-width)} / \text{float(strides[2])})$$

padding = 'valid' →

$$\text{new\_height} = (\text{input-height} - \text{filter-height} + 2 \cdot P) / S + 1$$

$$\text{new\_width} = (\text{input-width} - \text{filter-width} + 2 \cdot P) / S + 1$$

B. Tensorflow conv2d\_maxpool →

padding = 'same' →

$$\text{out\_height} = \text{ceil}(\text{float(in-height)} * \text{float(strides[1])})$$

$$\text{out\_width} = \text{ceil}(\text{float(in-width)} * \text{float(strides[2])})$$

padding = 'valid' →

$$\text{out\_height} = \text{in\_height} * \text{stride} + \text{filter\_height} - 1$$

$$\text{out\_width} = \text{in\_width} * \text{stride} + \text{filter\_width} - 1.$$