

April 12, 2023

DSML: Computer Vision.

Class starts  
@ 9:05 pm.

## Siamese Networks: One-shot learning.



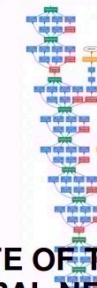
What normal people see  
when they walk on street



What Computer Vision  
folks see



**WHO WOULD WIN?**



STATE OF THE ART  
NEURAL NETWORK



ONE NOISY BOI

## Recap:

### \* Image segmentation

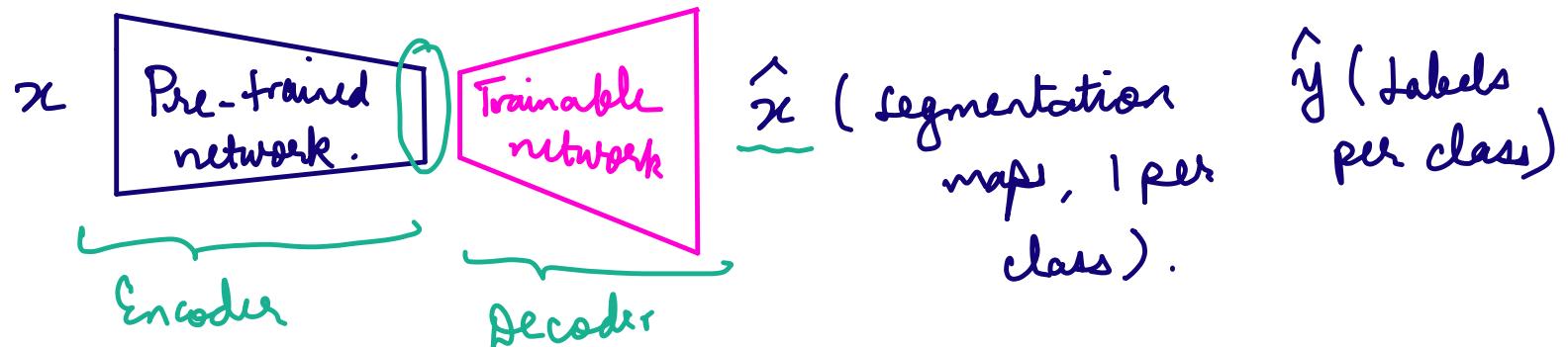
(a) Differences between segmentation & detection.

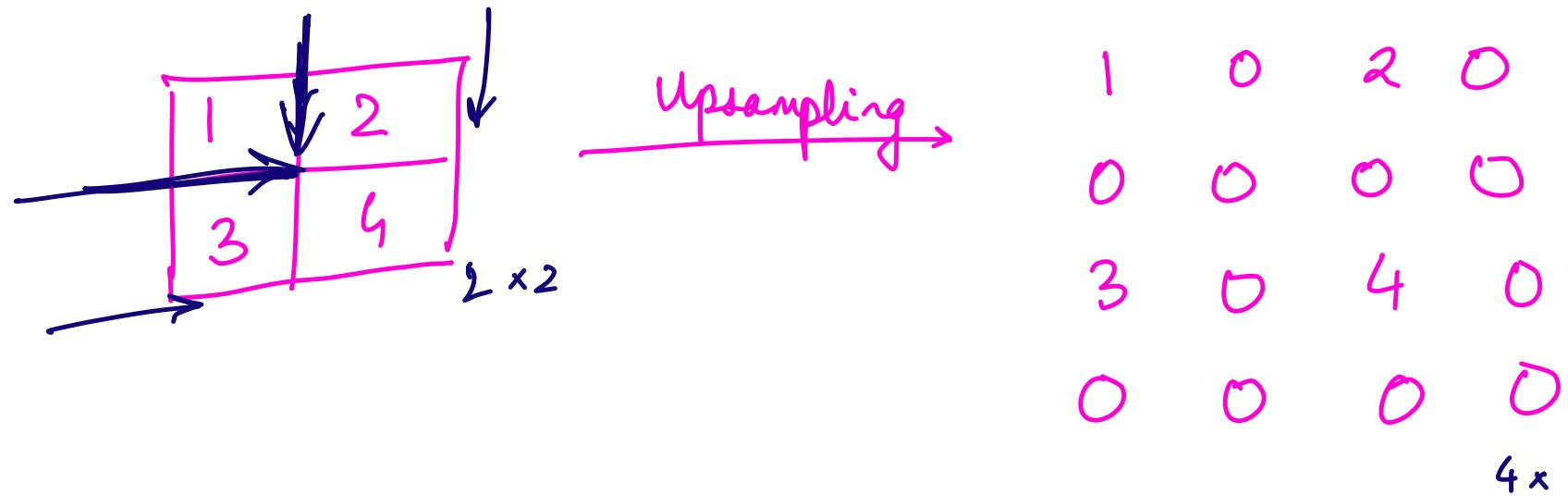
\* CNN architecture (Output layer): 2 heads  $\rightarrow$  classification  
regression  
 $\hookrightarrow$  same size masks as the input image.  $\approx 5$ .

\* Evaluation metric: IoU. DICE coefficient.

\* Loss: Combination of classification & regression.  
CCE loss

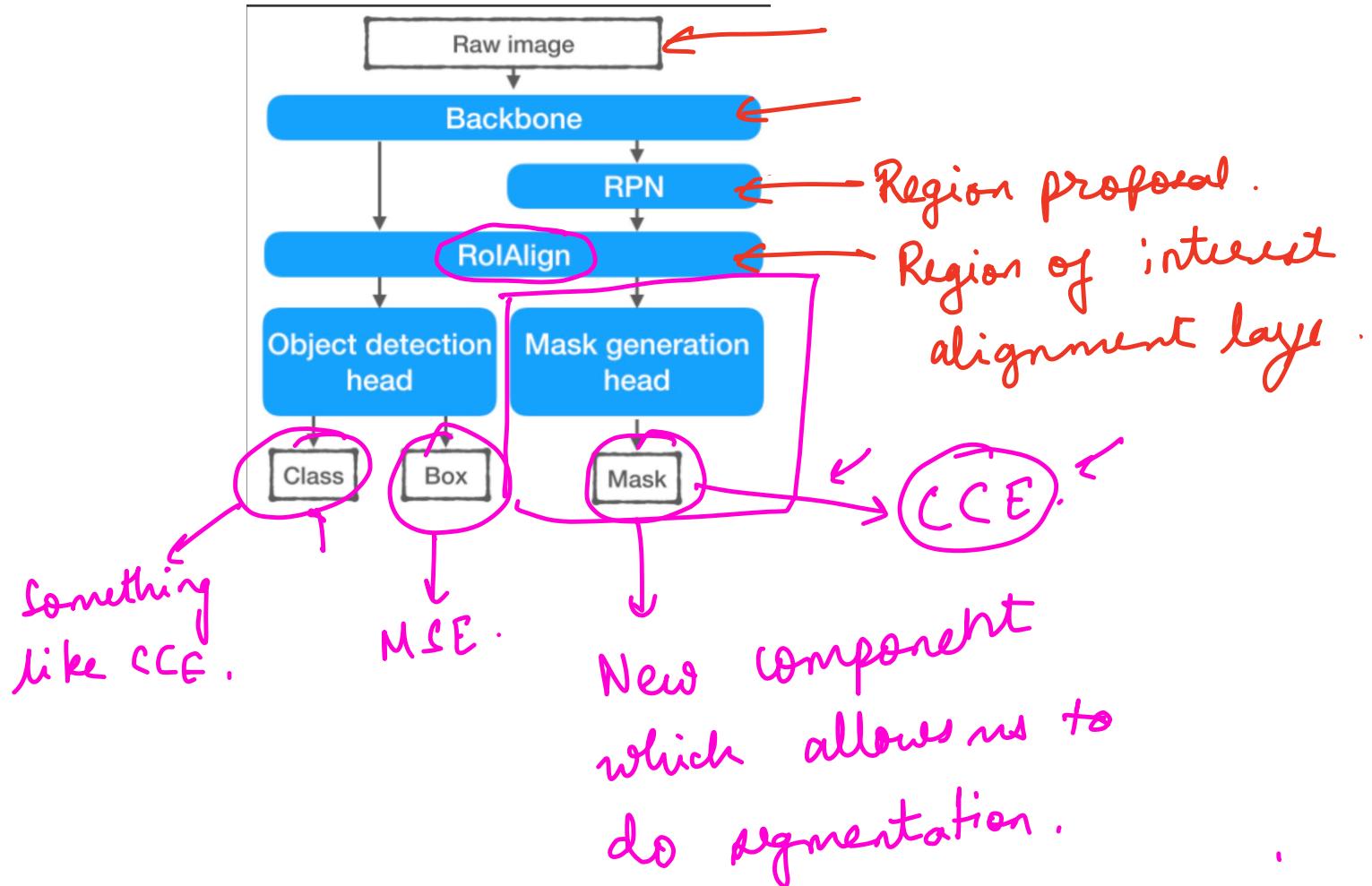
(b) FCN 8 architecture.

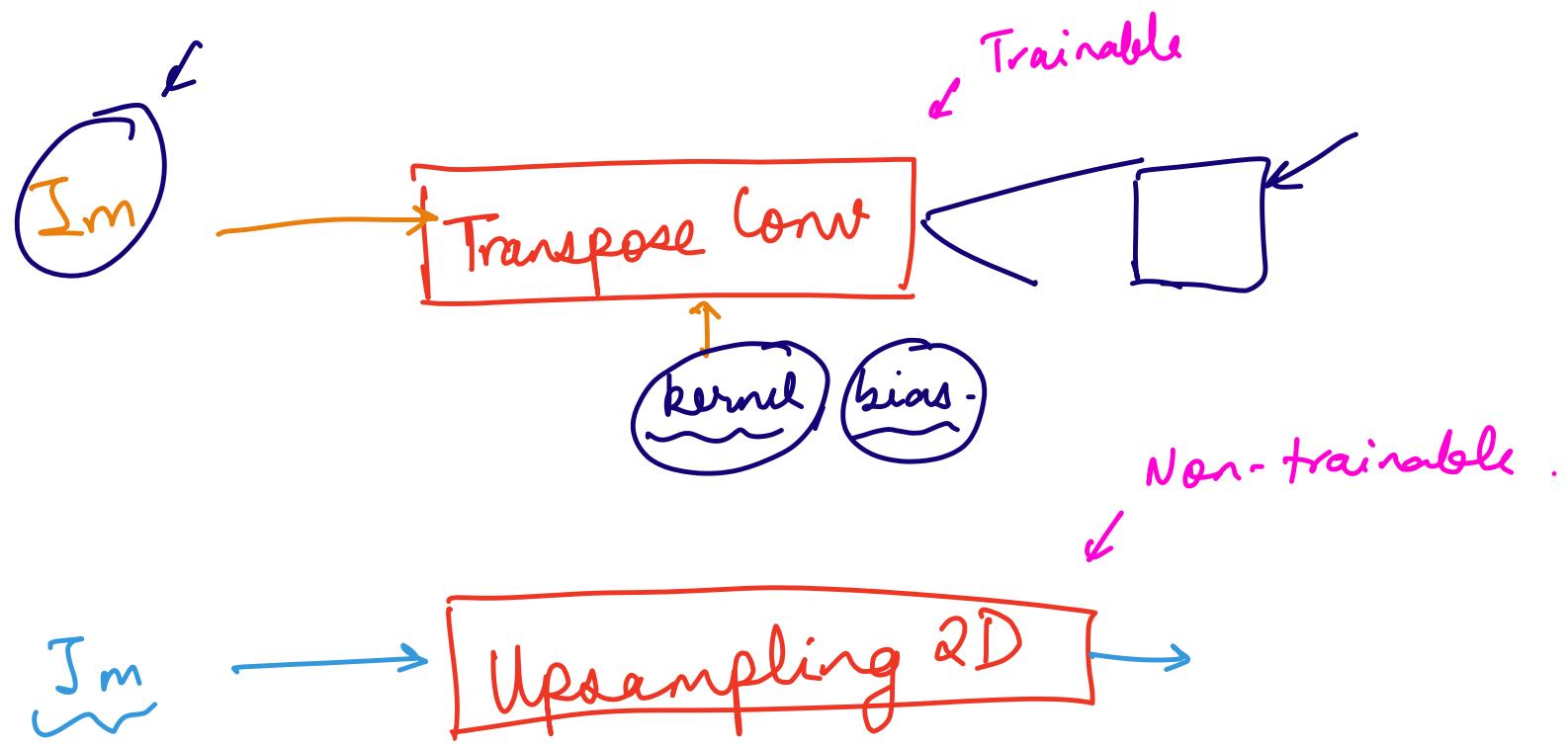




batch  $\times$   $\underbrace{h \times w}_{\text{no change}} \times \underline{\text{channels}}$

## Faster R-CNN.





## Object Classification



## One-shot learning

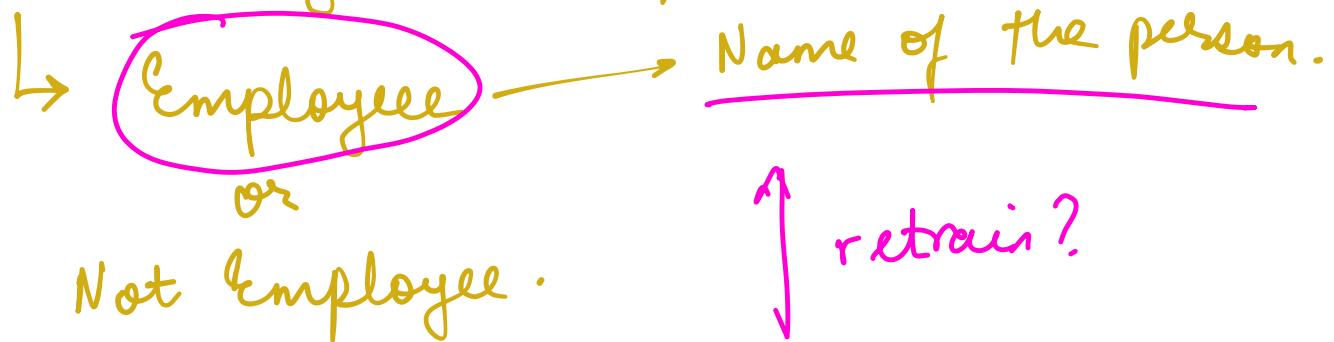


Are these the  
same objec<sup>t</sup>?

## Business Case:

→ Start a startup with 26 employees

→ Facial recognition software

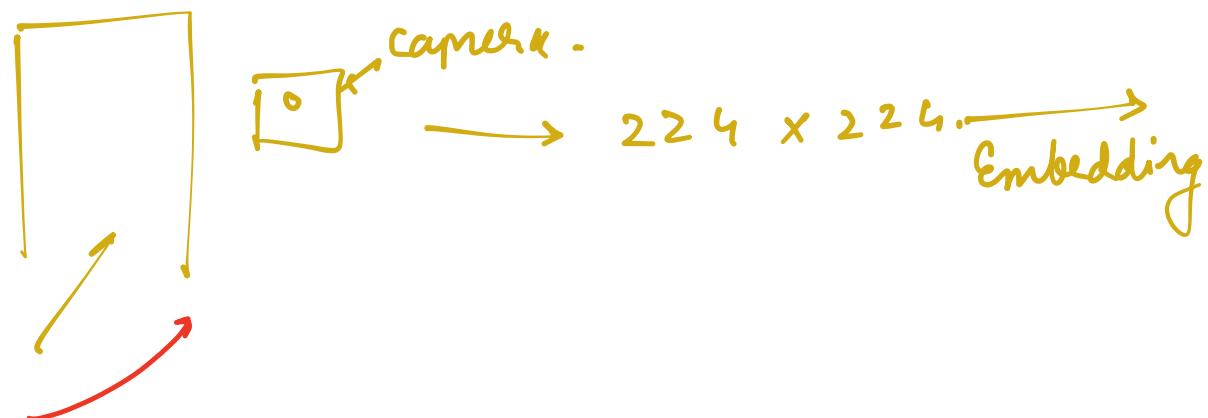
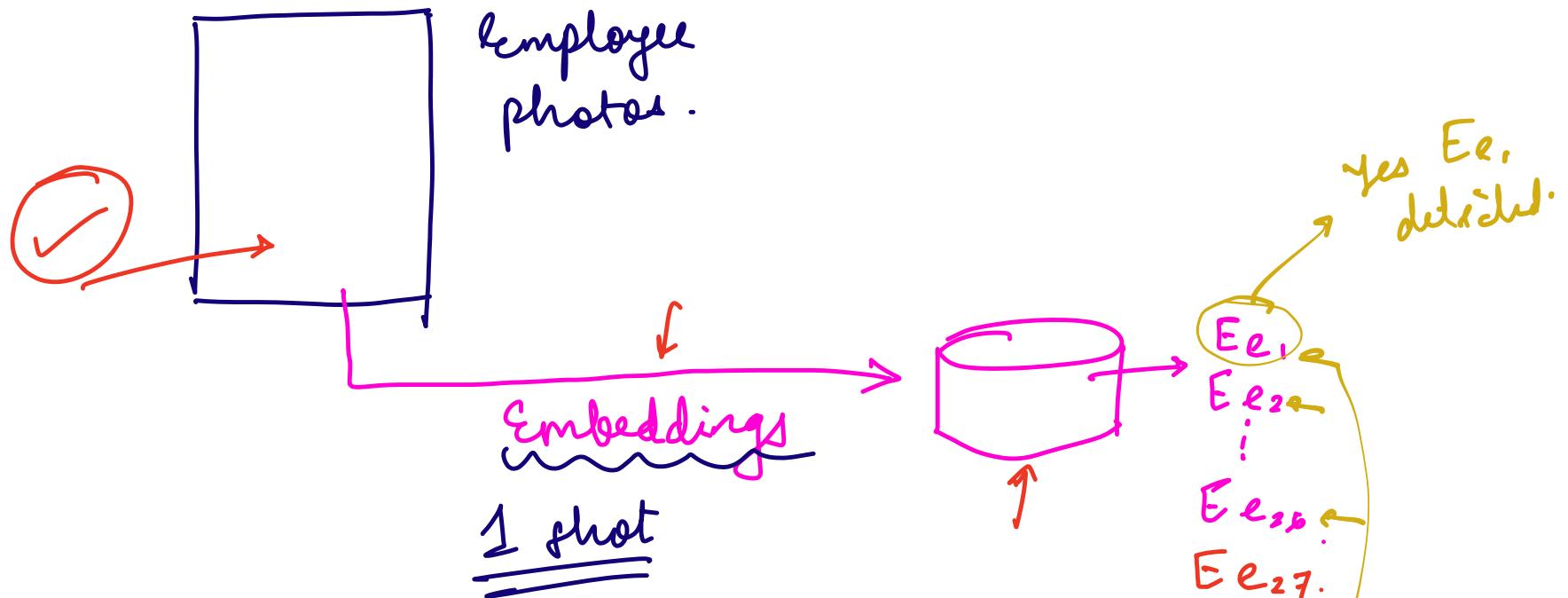


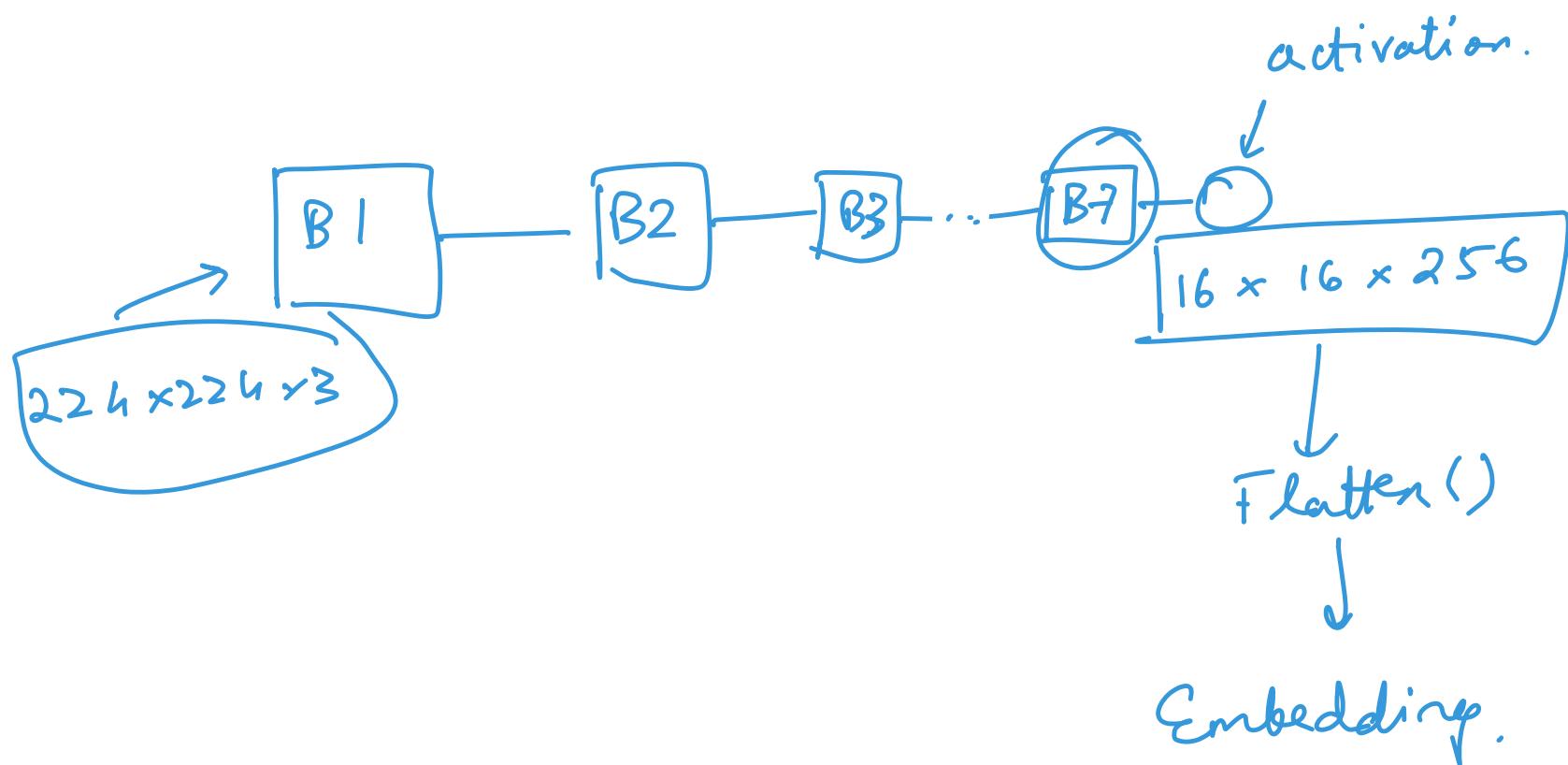
→ 1 new employee.

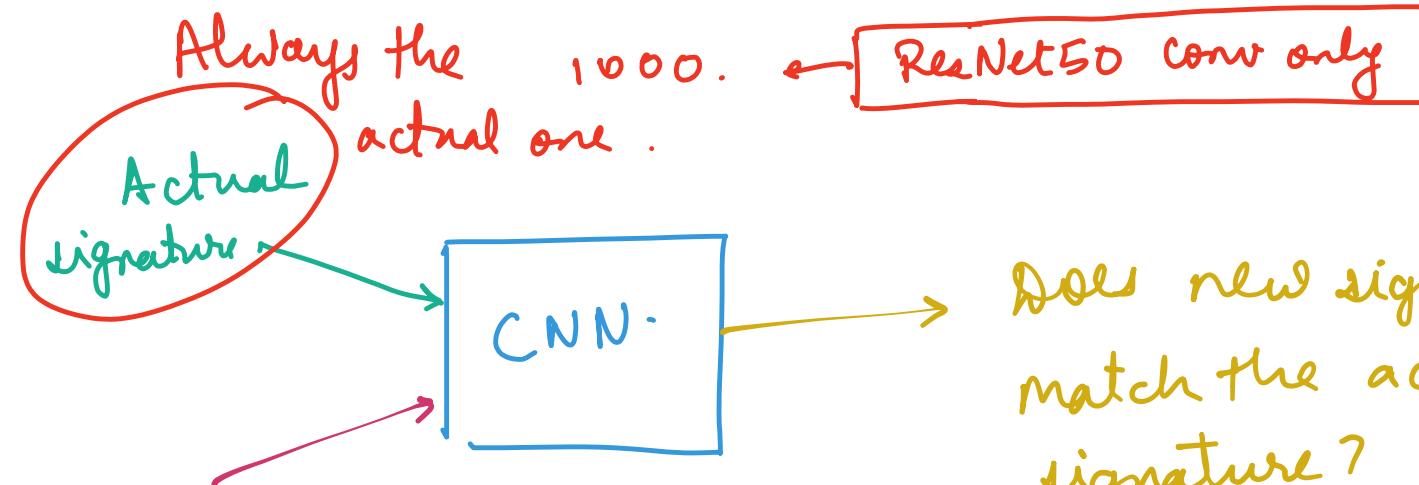
Problem: We will have to retrain the network  
every time we get new employees.

New system:

Database of all  
employee  
photos.



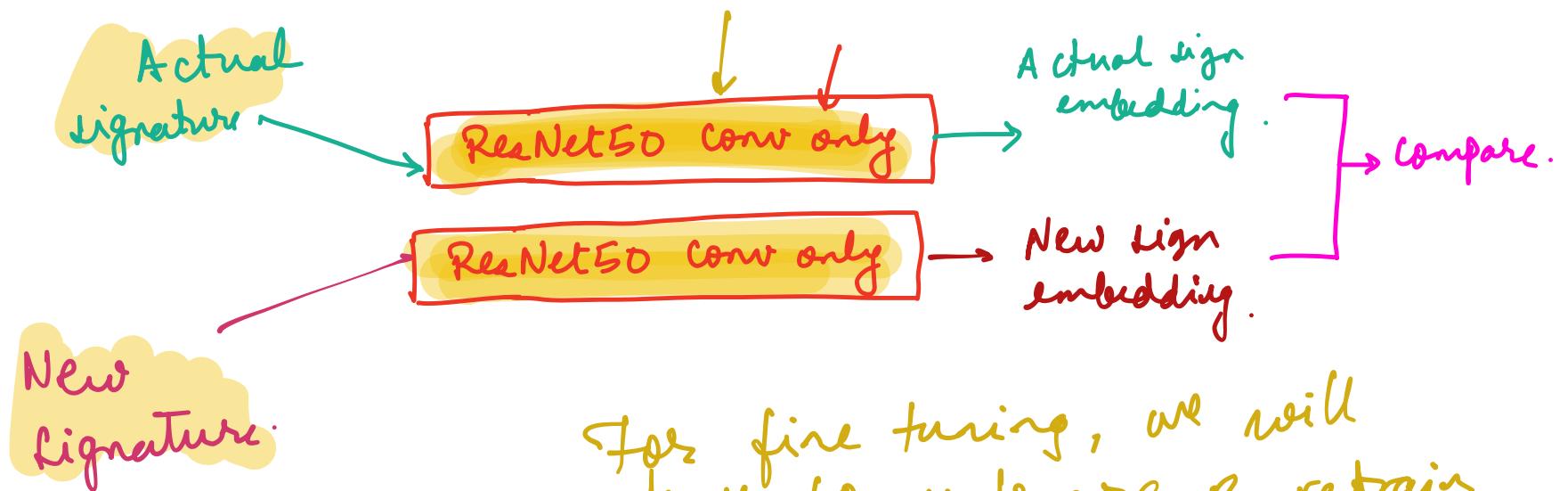




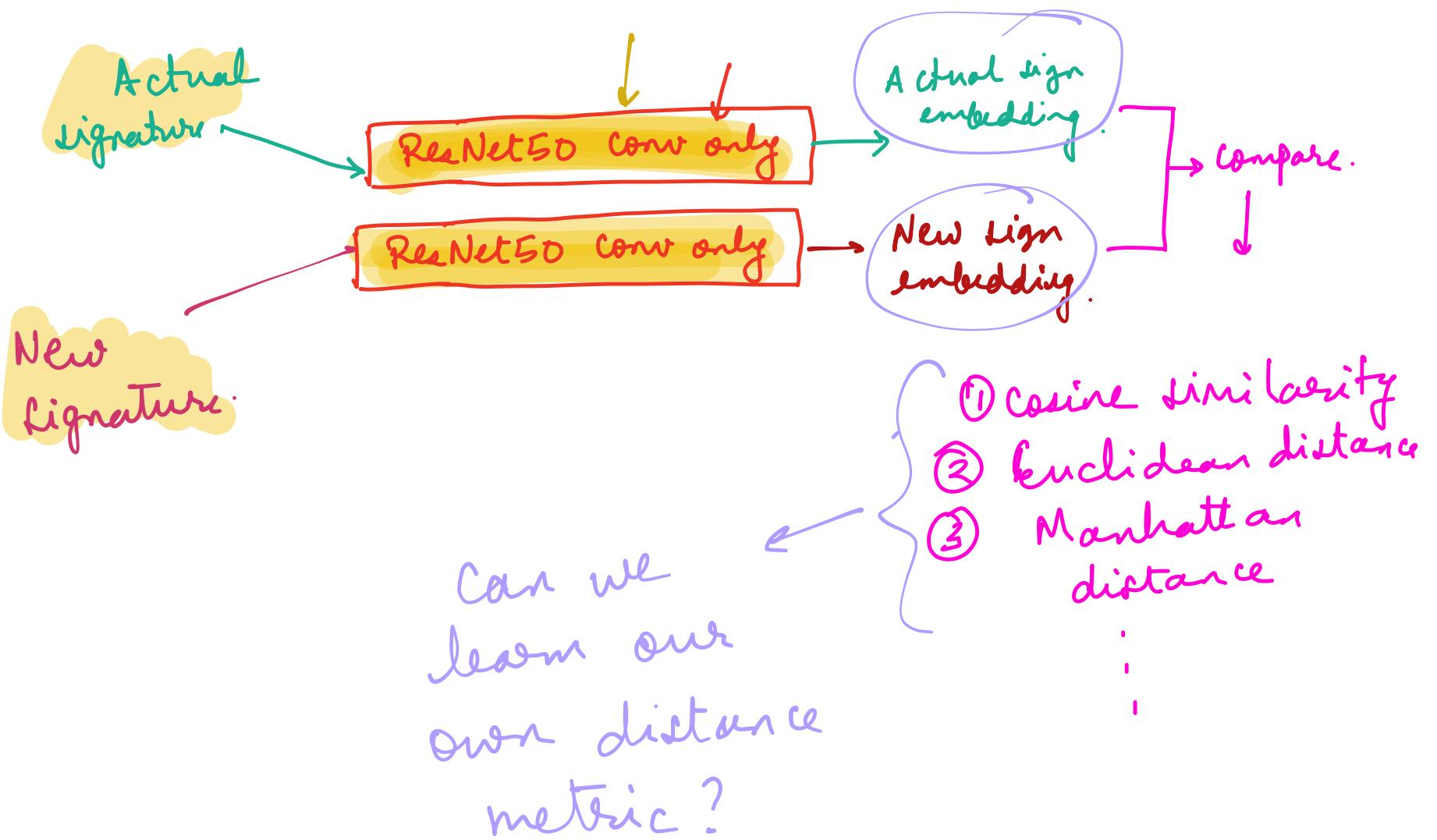
New  
signature:

"Siamese Network"

1 or 0.



For fine tuning, we will  
have to unfreeze & retrain  
the network!!



## Cases for One-shot learning:

Case 1: Genuine vs. Genuine. } We want to keep these close together

$$\frac{1 \cdot 7}{}, \frac{2 \cdot 2}{}, \frac{3 \cdot 8}{}, \frac{3 \cdot 6}{}, \dots$$

Case 2: Genuine vs. Fake. } Keep these far apart.

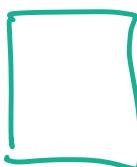
$$\dots, \frac{22 \cdot 7}{}, \frac{36 \cdot 5}{}, \frac{44 \cdot 8}{}, \frac{50 \cdot 3}{}, \dots$$

## Contrastive loss.



Input image  
(genuine)

$f v - a$



Query Image

$f v b$

label : ① -

Query Image  
is genuine.

0 - Query Image  
is fake -

Euclidean  
distance :

$D$ .

$$\rightarrow \text{Contrastive loss} = \gamma \cdot D^2 +$$

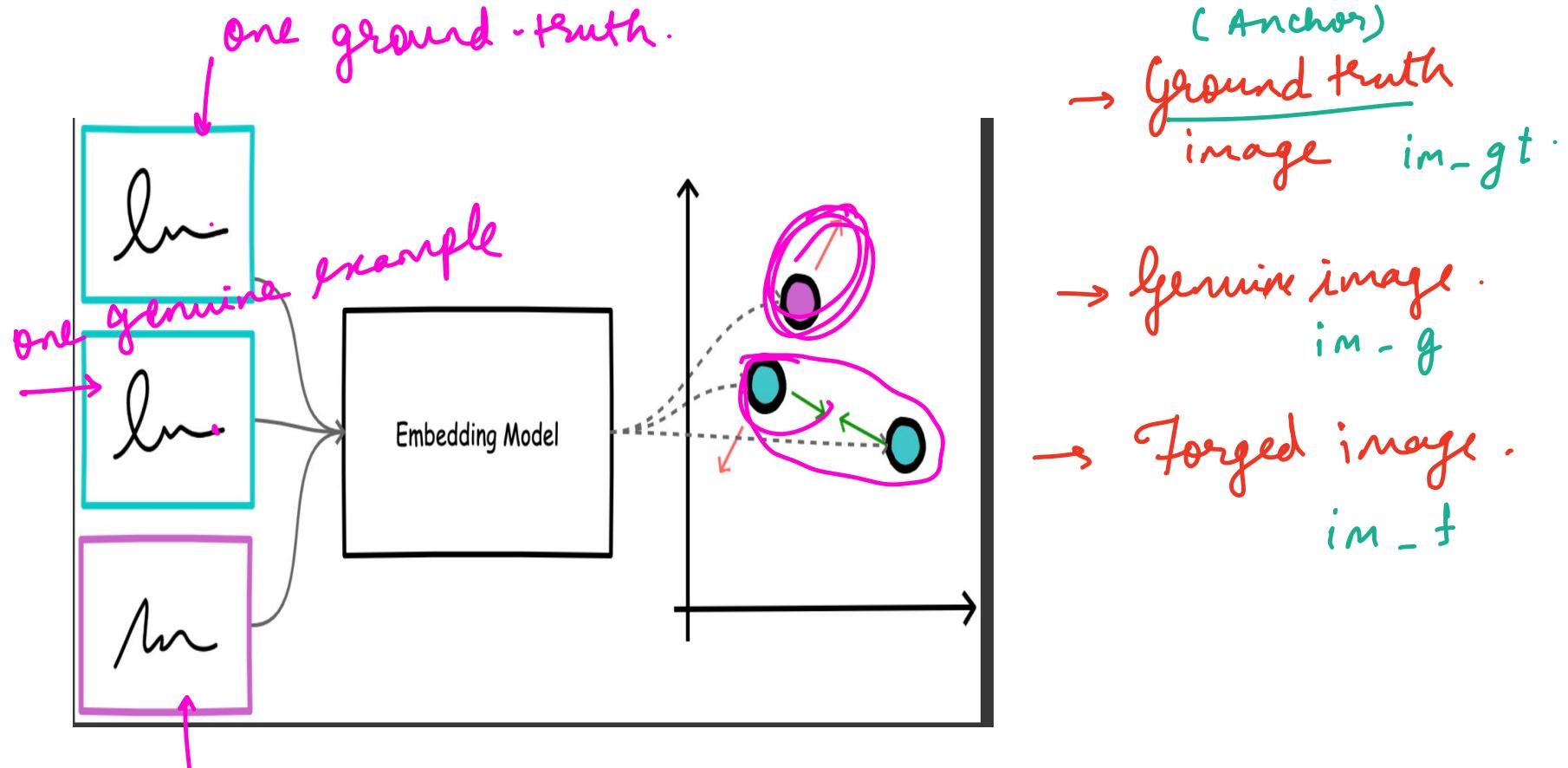
$$(1 - \gamma) \cdot \frac{\max((\text{margin} - D)^2, 0)}{ }$$

$$+ (1 - D)^2 \rightarrow \text{This is getting}$$

$$\min. D^2 \rightarrow (0.02)^2, (0.04)^2, (0.06)^2$$

↓

$$\min. (1-D)^2 \rightarrow (0.98)^2, (0.96)^2, (0.94)^2$$



$$\max \left( \text{Distance}(\text{im-gt}, \text{im-g}) - \text{Distance}(\text{im-gt}, \text{im-f}) \right) \geq 0.$$

$\leq$

$\text{Distance}(\text{im-gt}, \text{im-g}) \leq \text{Distance}(\text{im-gt}, \text{im-f})$

$\text{Distance}(\text{im-gt}, \text{im-g}) - \text{Distance}(\text{im-gt}, \text{im-f}) \geq 0$

## Triplet loss:

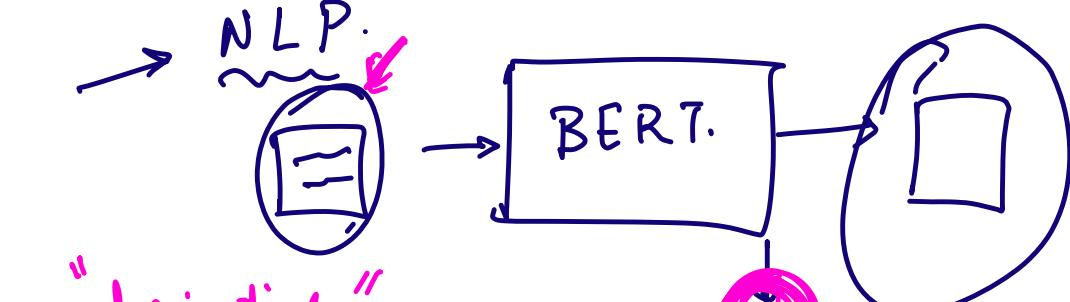
- A → Anchor embedding, ground truth embedding.
- P → Positive, genuine image.
- N → Negative, forged image.
- D(x, y) → Euclidean distance between  $x$  &  $y$ .

loss :  $\max \left( \underbrace{D(A, P) - D(A, N)}_{\text{negative}} , 0 \right)$

positive - Good, do nothing.

positive - Bad, do an update.

→ Zero-Shot learning →



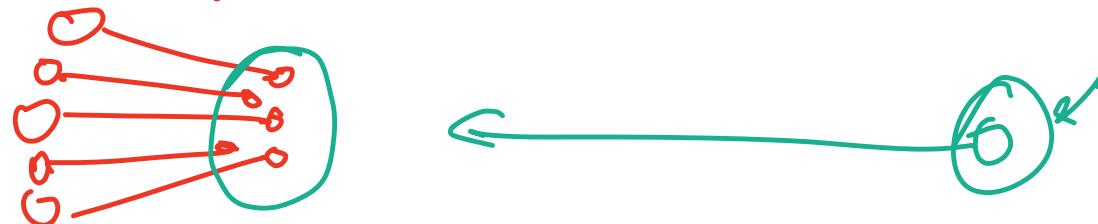
"derivative"

K shot learning.

No. of training examples of the  
positive class.

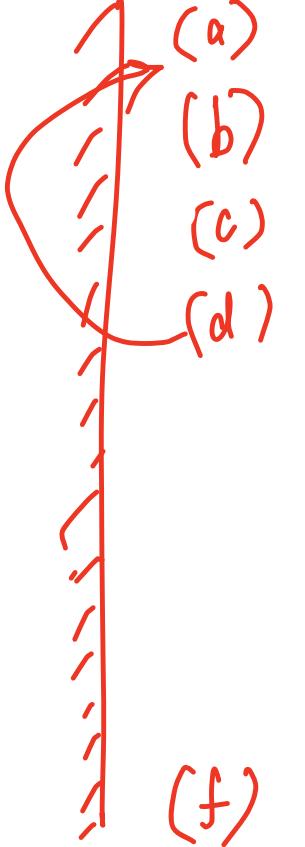
$K=5$

⑤ photos of  cat,  dog,  zebra.



I may see classes  
not encountered before,  
and must classify  
them anyway.

## Discussion Requests:

- 
- (a) Chat GPT.
  - (b) Image + Text Multimodal. (Diffusion, Dall-e)
  - (c) Lang Chain
  - (d) RL HF
  - ↓  
Re-inforcement learning through Human feedback.
  - (f) 3D related networks