

## Face Recognition

Input Image size = 224

face\_detection = RetinaNet() (Resnet18 + FPN)

face\_recognition = SNFaceNetModel() (SqueezeNet1.1)

```
input_image = Image.open(opt.input_file_name)
retinanet = model.resnet18(num_classes=1)
model_checkpoint = torch.load(opt.face_detect_model,
                               map_location=cpu)
```

# Remove keys related to DataParallel

```
new_state_dict = OrderedDict()
```

```
for K, V in model_checkpoint['state_dict'].items():
```

```
    name = K[7:]
```

```
    new_state_dict[name] = V
```

```
retinanet.load_state_dict(new_state_dict)
```

If use-gpu:

```
retinanet = torch.nn.DataParallel(retinanet)
```

```
facenet = model.SNFaceNetModel()
```

## OpenVINO:

```
classification, regression = iter(net.output) # face-detection outputs
```

```
output_blob_recog, output_blob_recog_gunder = iter(net_recog.output) # face-recognition output
```

```
exec_net = plugin.load(network=net) # load the face detection network
```

```
exec_net_recog = plugin.load(network=net_recog) # load face recognition detection network
```



$$x = (512, 64, 256, 128)$$

$$x = 512, \dots$$

$$\text{genout} = (512, 64) \rightarrow \text{Conv} \\ + \text{BN}$$

$$+ \text{ReLU}$$

$$+ \text{genfca} (64, 2) \rightarrow$$

$$s^0, \text{genout} = [\text{male}, \text{female}]$$

$$x = (512, 128)$$

$$\text{return } x, \text{genout} \\ \Downarrow \quad \Downarrow \\ [128] \quad [2]$$



if FACE\_ALIGN:

if FACE\_ALIGN\_POINTS=68:

sp = dlib.shape\_predictor('shape\_predictor\_68\_face\_landmarks.dat')

else  
sp = \_\_\_\_\_ (5-face-\_\_\_\_\_)

# Finding embeddings of reference images:

folder\_list = sorted([f.name for f in os.listdir(opt.ref\_image\_dir) if f.name != '\_\_'])

ref\_img\_list = [[] for \_ in range(len(folder\_list))]

ref\_img\_encoded\_list = [[] for \_ in range(len(folder\_list))]

for i in range(len(folder\_list)):

for name in files:

ref\_img\_list[i].append(name)

# i are all folders which contains

images (2/3) for that person

ref\_img\_list = [[random, random], [pure], [suit, suits], ...]

# running the face recognition infer request for reference images

ref\_image\_encoded = exec\_net\_recog.infer(input={input\_image:  
batch})

ref\_image\_encoded\_list[i].append(ref\_image\_encoded)

# running infer request for face recognition model:

face\_image\_encoded = exec\_net\_recog.infer(input={input\_image:  
batch})



1 inference image  $\rightarrow$  detection + encoding  
[encodings]

reference images  $\rightarrow$  resize-center-crop

(  
only one time  
when we start  
the script

encoding (only encode  
model)  
not detect

[encodings]  $\forall$  ref images

$\rightarrow$  Calculate the dis b/w face-encodings

& reference images encoding.

$\rightarrow$  predict the class (using min. dist).

---

2 train  $\rightarrow$  RetinaNet  $\rightarrow$  (classification loss)  
 $\rightarrow$  Sigmoid Loss

$\rightarrow$  triplet loss.

(4/5 S.p.u for each class)



face\_image\_encoded = ~~face\_image\_encoded~~  
l2\_norm(face\_image\_encoded [output  
blob\_recog]).

face\_image\_encoded \* 10

face\_img\_encoded\_list.append(face\_image\_encoded)

for id in range(len(folder\_list)):

distances = []

for j in range(len(ref\_images\_list[id])):

dist = l2\_dist(face\_image\_encoded,  
face\_img\_encoded\_list[id][j]).

distances.append(min\_dist.detach().numpy().item())

top\_5\_idx = np.argsort(distances)[:5]

face net with gender classification - train

~~face~~  
model = FaceNetModel( ) → resNet-34  
gen\_model = SNFaceNetModel( ) → SqueezeNet 1.1

---

Face Recog



# • Face-Recognition- summary

① face detection

⇓  
RetinaNet

⇓  
(Base model is ResNet-18)

+  
(FPN)

+ focal loss.

+ face Recognition

⇓  
SiFaceNet

⇓  
(Base model is SqueezeNet)

⇓  
Generate 128 length embeddings

+ triplet loss for is used

② use OpenVINO (for inference).

③  $img \rightarrow$  inference image

we kept a set of reference images

• inference image ( $img$ )  $\rightarrow$  detection (RetinaNet) + encodings (SiFaceNet)

• reference images  $\rightarrow$  (resize  $\rightarrow$  center crop)   
 preprocessing

⇓  
encoding (only encoding, no detection)

⇓  
[embeddings] of reference images.

↳ then calculate the ~~distance~~ (L2 Norm) distance  
(L2 Norm) b/w face encoding ( $img$ ) &  
reference images encodings.



and, compute the dist from all images in reference - img list.

↳ take the min distance.

↳ min dist should be less than a threshold.

↳ + used face landmarks using dlib to calculate the distance b/w eyes (or b/w ear & reference images) for fine tuning the model.

➡ Based on the previous threshold & this threshold final class name was predicted.

↳ Name was taken from ref. images folder name.

• My Role:

- ① RetinaNet - training - setup. (& test setup) & integrated it into pipeline.
- ② Understand the SFaceNet model & integrated the embeddings to pipeline.



## Face Recognition

- ① ODI (RetinaNet)
- ② embeddings (SiameseNet)

## License Plate

- ① ODI (RetinaNet)

Wpoc. network

- ② Character Segmentation (RetinaNet)  
+ character classification (CNN)

## Vehicle

- ① ODI → YOLO

(CN, UN no., hazard symbols, seals etc)  
(Retrained on new dataset).

- ② UN no. ~~segment~~ character Segmentation  
(w/ RetinaNet).

+ UN no. character classification.

## ↳ post-processing

- ① K-means Clustering.

- ② take the max. occurrence of a character in Real time. (inference).



$$\|d(x_1, x_2)\|_2 = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$

$$x_1 = [x_{11}, x_{12}, x_{13}, \dots, x_{1n}]$$

$$x_2 = [x_{21}, x_{22}, x_{23}, \dots, x_{2n}]$$

Triplet loss:

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha$$

$$\forall (f(x_i^a), f(x_i^p), f(x_i^n)) \in T$$

$$\text{Loss} = \max\left(0, \sum_i \left[\|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha\right]\right)$$

ie,

$$\text{Loss}_{\text{triplet}} = \max(0, [\text{pos dist} - \text{neg dist} + \alpha])$$

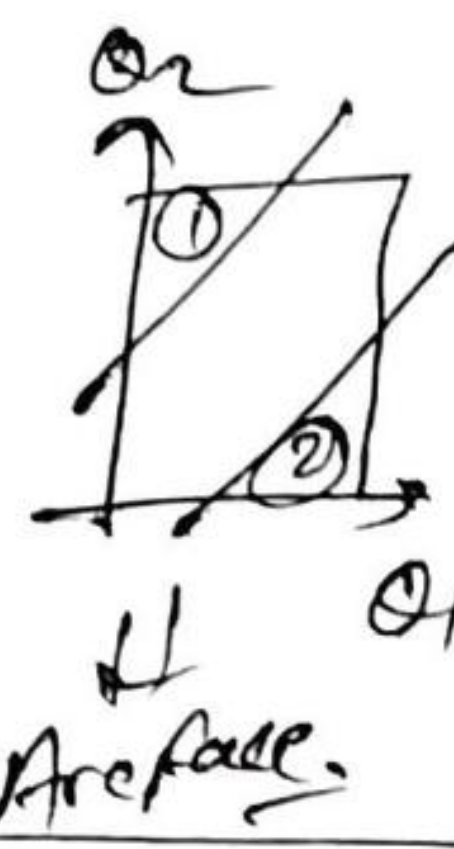
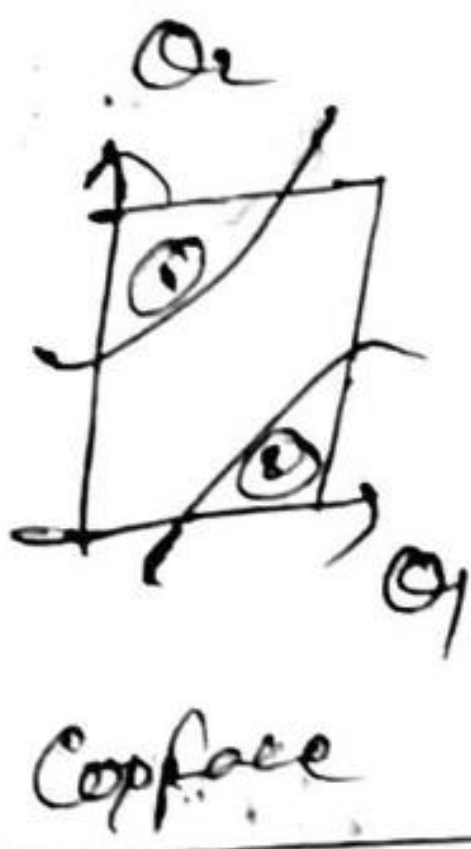
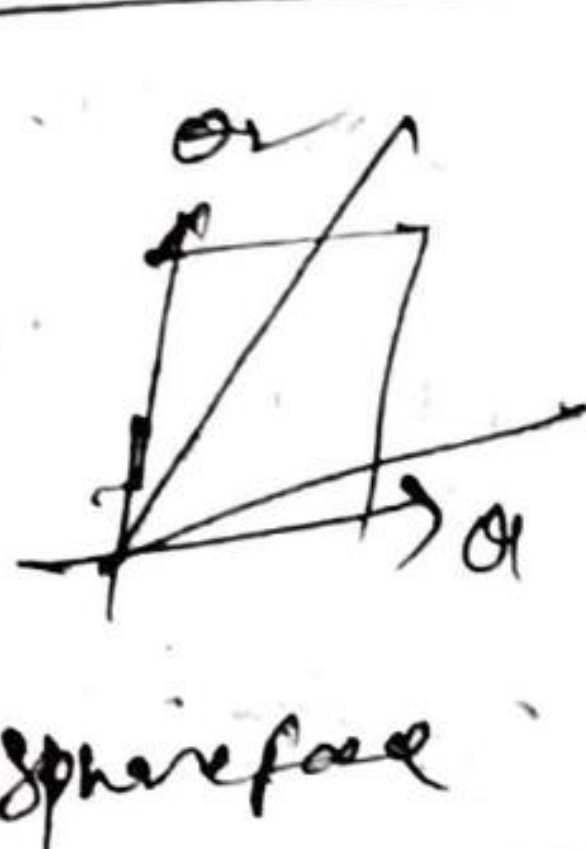
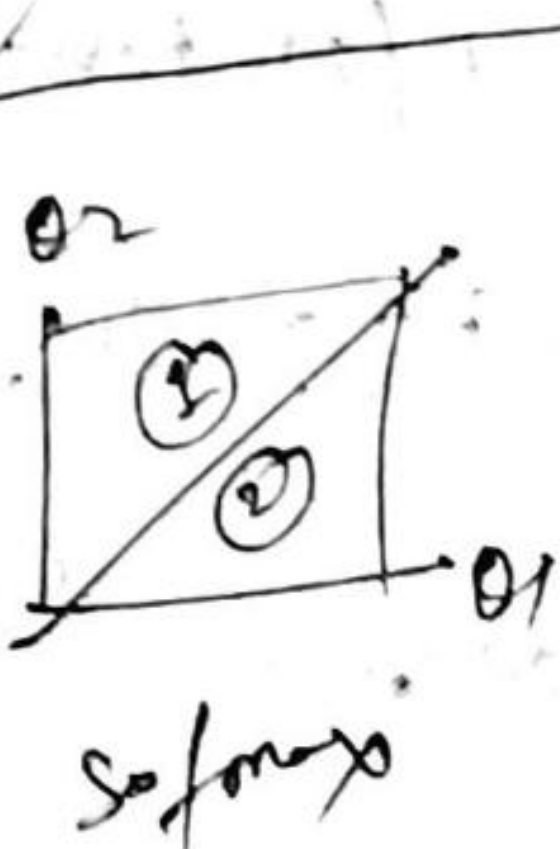
Proface  $\rightarrow$  Modified softmax loss

$$L_{\text{softmax}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{w_i^T x_i + b_i}}{\sum_{j=1}^N e^{w_j^T x_i + b_j}}$$



$$L_{\text{argane}} = \frac{1}{N} \sum \log e^{s(\cos(\theta y_i + m) + m)}$$

$$L_{\text{argane}} = \frac{1}{N} \sum \log \frac{e^{s(\cos(\theta y_i + m))}}{e^{s(\cos(\theta y_i + m))} + \sum_{\substack{j=1 \\ j \neq y_i}}^n e^{-s \cos \theta_j}}$$



# Cargo: Arcface

train.py \*optimizer = Adam(list(model.parameters()) + list(detach\_model.parameters()),  
loss = lr, amsgrad = True,  
weight\_decay = 1e-4)

class Arcface(nn.Module):  
 return s as output

model = S1NfaceNetModel()

detach\_model = Arcface()

~~crit~~ criterion = nn.CrossEntropyLoss()  $\Rightarrow$  training loss

def train(model, -- detach\_model):

for epoch ---:

anc\_emb = model(anc\_img)

anc\_cls = detach\_model(anc\_emb, anc\_cls)

loss = criterion(anc\_cls, anc\_label)



```
optimizer.zero_grad; loss.backward(); optimizer.step()
loss_sum += loss.item()
avg_loss = loss_sum / len(dataloaders['train']); print('CE-loss = avg_loss')
```

```
valid(model, dataloader, detach_model):
```

```
for batch_idx, batch in enumerate(dataloader[valid]):
```

```
anc_embed = model(anc_img)
```

```
pos_embed = model(pos_img)
```

```
neg_embed = model(neg_img)
```

```
pos_dist = l2_dist.forward(anc_embed, pos_embed)
```

```
neg_dist = l2_dist.forward(anc_embed, neg_embed)
```

```
all = (neg_dist - pos_dist < margin)
```

```
triplet_loss = TripletLoss(margin).forward(anc_hard_emb,
pos_hard_emb, neg_hard_emb).
```

```
dist = l2_dist.forward(anc_embed, pos_embed)
```

```
distances.append(dist.data.cpu().numpy())
```

```
label.append(np.ones(dist.size()))
```

```
dist = l2_dist.forward(anc_embed, neg_embed)
```

```
distances.append(dist.data.cpu().numpy())
```

```
labels.append(np.zeros(dist.size()))
```

```
triplet_loss_sum += triplet_loss.item()
```

```
avg_triplet_loss = triplet_loss_sum / len(dataloaders[valid])
```

```
labels = np.array([sublabel for label in labels for
sublabel in label])
```

```
distances = np.array([subdist for dist in distances
for subdist in dist])
```



$tp_r, fp_r, accuracy = evaluate(distances, labels)$

if phase == 'valid':

torch.save(model.state\_dict())

torch.save(detach\_model.state\_dict())

plot\_roc(fp\_r, tp\_r, figure\_name='acc.png')

val\_loss = avg\_val\_loss

return val\_loss

for  
dist.  
pos  $\rightarrow 1$   
Label  $\rightarrow$  and  
Neg  $\rightarrow 0$   
for dist.



# ALPR

(WPOD  $\rightarrow$  Wrapped planar Object detection)

ResNet o/p  $\rightarrow$  (128, )

self.clf  $\rightarrow$  (128, 2)

self.regression  $\rightarrow$  (128, 6)

code: 2

$x = \text{self.layer4}(x)$

# Clf  
 $c = \text{self.classification}(x)$

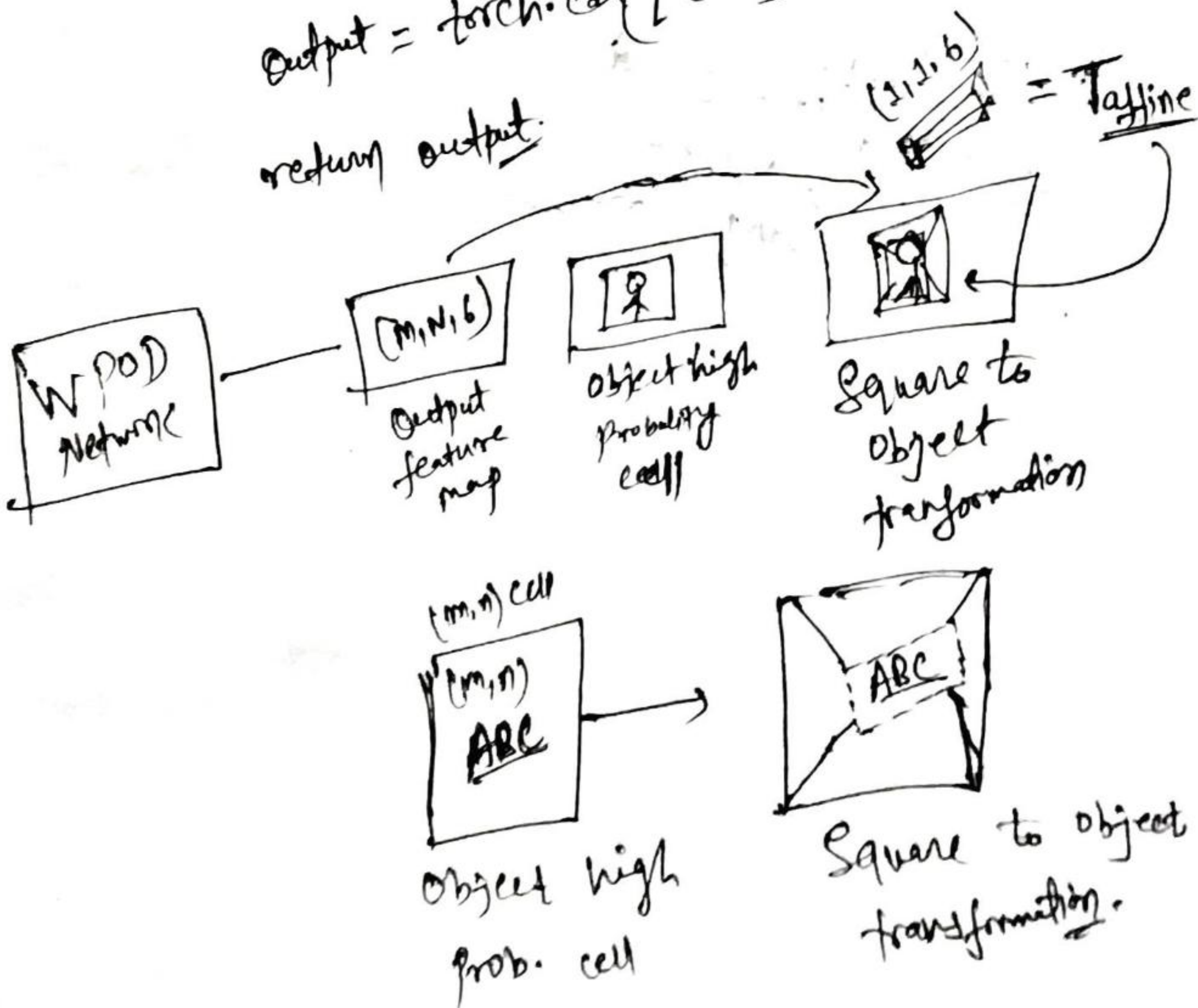
$c = \text{self.softmax}(c)$

# reg.

$r = \text{self.regression}(x)$

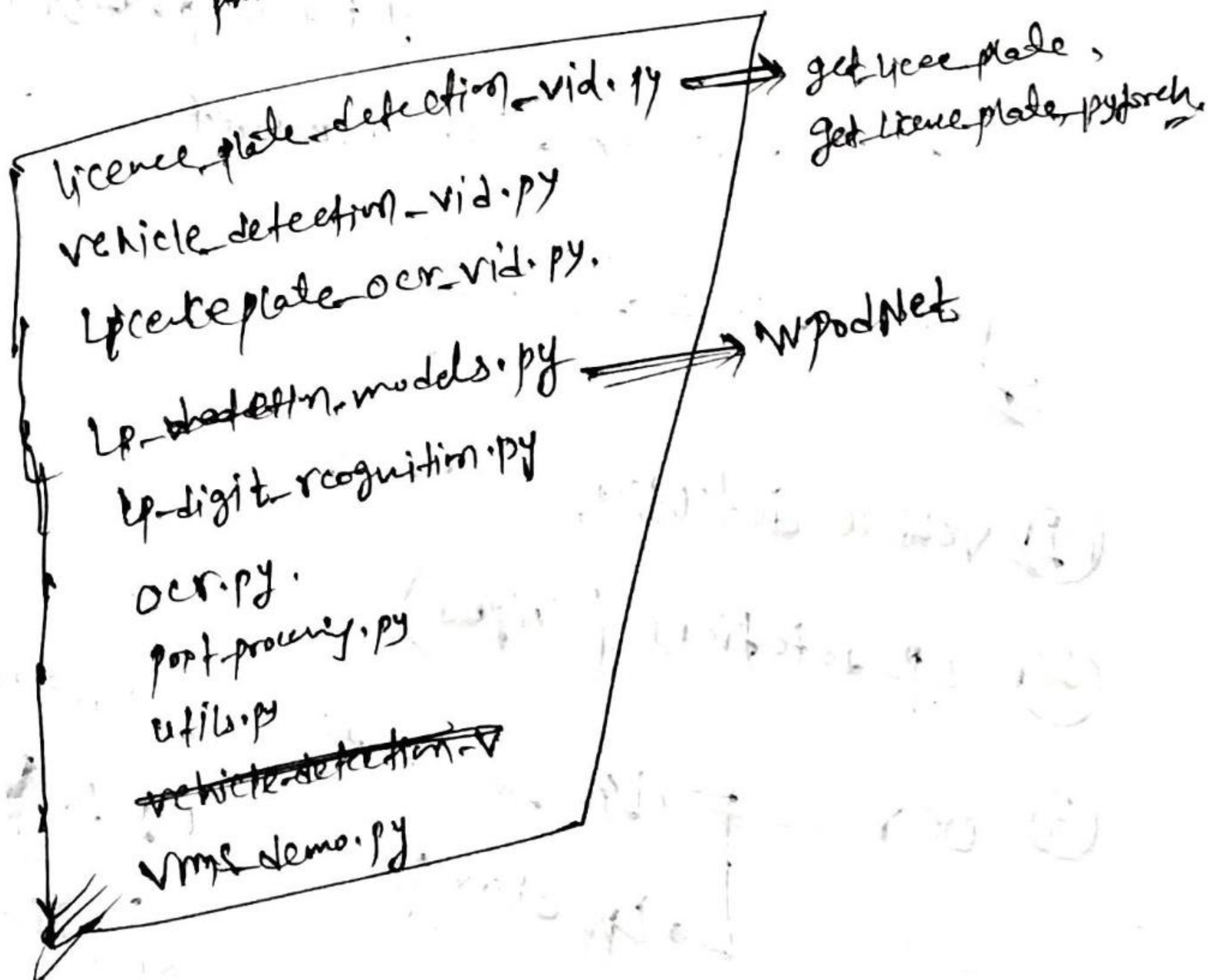
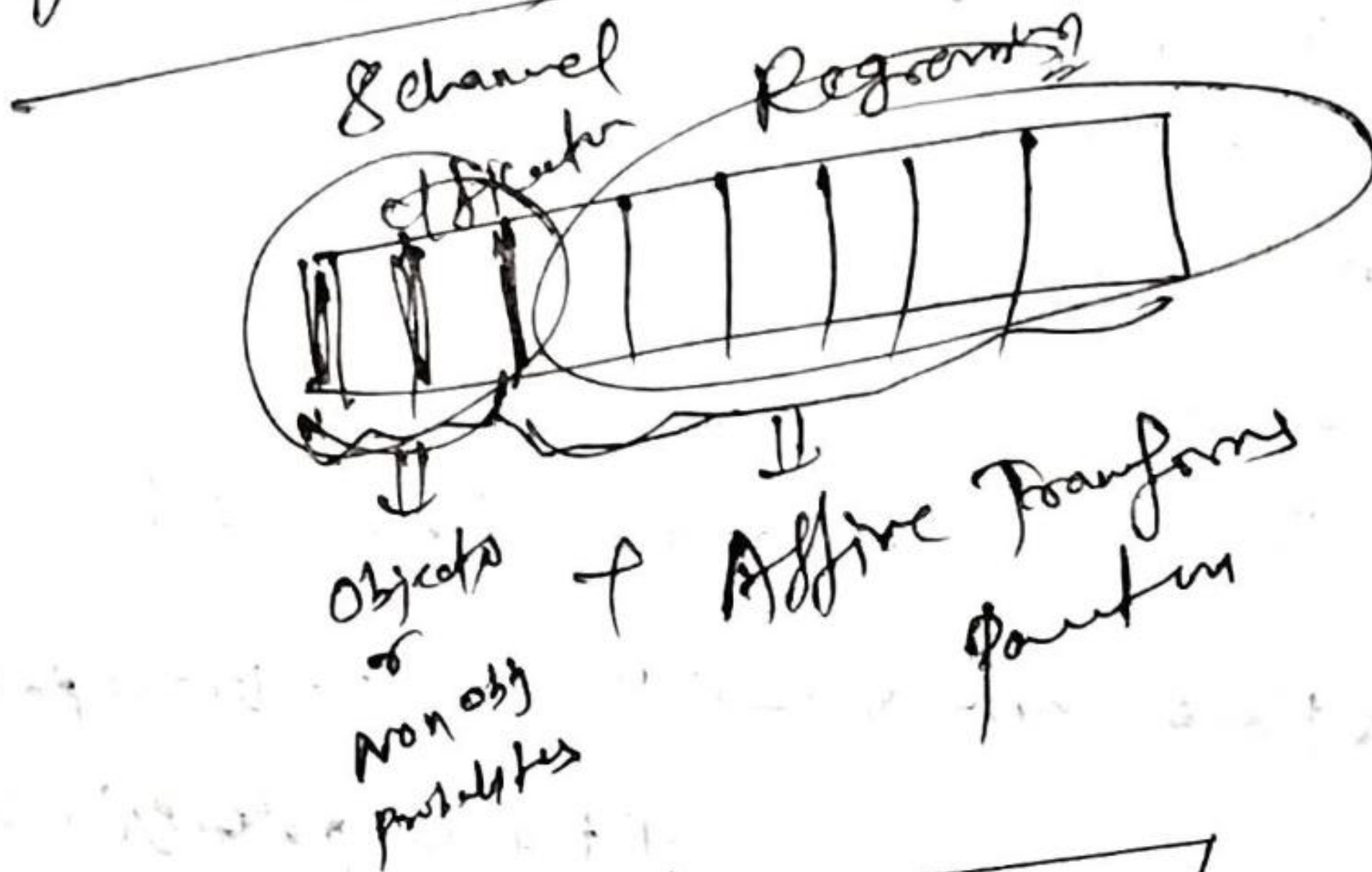
$\text{output} = \text{torch.cat}([c, r], \text{dim}=1)$

return output





feature map:





## Vehicle-detection.py

sys.path.append('..yolo')  
from yolo predict

import darknet as dn

postprocessing → lp-sort

ocr → character-recogtn

cfg ~~input~~, wpo2 net path, lp-detect-model

license plate detection vid ~~input~~ → get license plate  
get license plate torch

lp-detection-model ~~input~~, wpo2 net

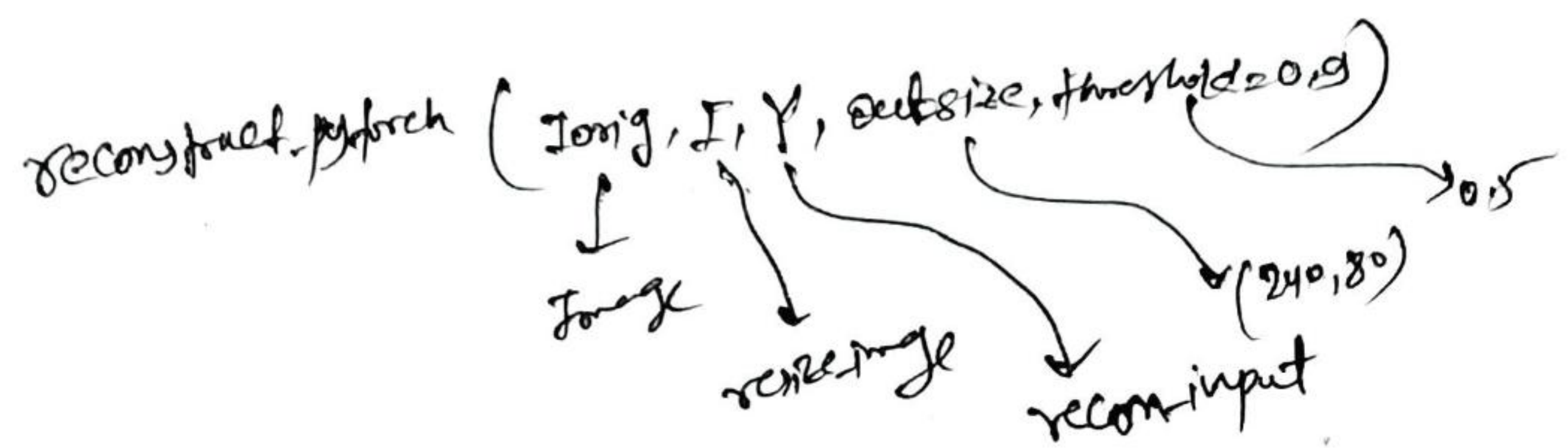
Yolo v3, lp-detection, ocr

① vehicle detection (Yolo v3)

② lp-detection (wpo2)

③ ocr → lp-segmentation (RefineNet)  
→ lp-classification (CNN)  
Resnet





$x = (2, 3, 4)$   
 $x[:, :, 0]$



## #Alpr summary

- ① vehicle-detection  $\rightarrow$  YOLOv3
- ② LP-detection  $\rightarrow$  WPOD Net + Prospective Transform.
- ③ LP-segmentation  $\rightarrow$  RetinaNet
- ④ LP-classification  $\rightarrow$  CNN (ResNet),  
based CIF model.

### • Post processing:

- ① d2c map
- ② c2d map



## # Note ODI & UN Number Summary

### ① ODI (YOLOV3)

- CN
- UN
- LP
- Seals
- LP sticker
- Hazard category
- hole
- dent

↳ Yolo detects each of the above classes  
↳ saved cropped image for each classes & stored it in respective input folder (in ~~np~~ .npy format) for further processing.

for some module  
& for some module  
it was saved  
in cv2 image format.

### ② UN seg. & clf.

- ↳ UN No. Seg : → RetinaNet
- ↳ UN No. Clf → CNN (ResNet) Based Model.
- ↳ UN No post processing.
  - ↳ ~~2d~~ 2d crop
  - ↳ c2d crop
  - ↳ K-Means clustering Algorithm.



## # KYCL summary

① Twist-lock Rot detection.  
↳ RetinaNet

② Twist-lock classification  
↳ CNN (ResNet) Based classification Model.

③ OpenVino conversion of PyTorch Model.



## # Training set up

### # Face Recognition

① face detection → RetinaNet → our exp. data

② FaceNet → (face recog.) → LFW dataset (14~16K)

### # VMS

→ vehicle detection → YOLO

→ LP-detection → WPOD

→ LP-segmentation → RetinaNet

→ LP-classification → CNN (ResNet) based Model.

our Car dataset

### # Wate:

→ ODI → YOLO (60k images)

→ UN No. → original, UN No. collected from internet

→ seg-RetinaNet  
→ clf-CNN (ResNet)  
→ Type-1 UN No.  
→ Type-2 UN No.

### # KYCL

→ Twistlock ROI Detection → RetinaNet

→ Twistlock light clf → CNN (ResNet)

### # Cargo

→ cargo-classification (Reg-Misc) → CNN (ResNet)

→ ~~ROI~~ (cargo box detection) → YOLO V3

→ ~~Tracking~~ (DeepSort, Kalman filter)

→ Feature extraction (ArcFace Model - ~~arcface~~)

→ Cargo dataset

→ Cross Entropy loss (train)  
→ Triplet loss (val).