

# The Integration of Thermal and Visual Images for Human Detection Using Deep Learning

GFM M.Sc. Mid-term Defence

*Qiao REN*

First Supervisor: Dr. Siavash Hosseinyalamdary

Second Supervisor: Dr. Francesco Nex

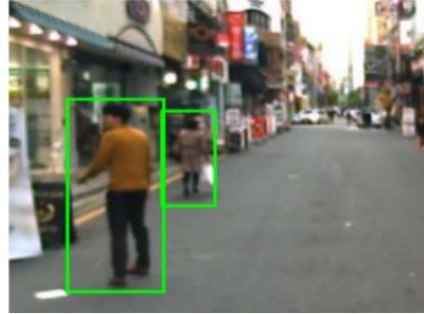
Nov 21, 2018

# Content

- 
1. Problem Statement
  2. Research Identification
  3. Methodology
  4. My Progress
  5. Intermediate Result
  6. What to do next

# 1. Problem Statement

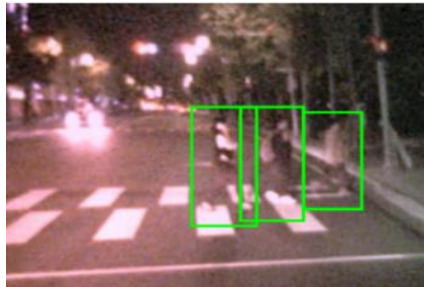
- ▶ Main information sources: Visual Images and Thermal Images
- ▶ Drawbacks of visual and thermal images
- ▶ Sensor fusion: RGBD (RGBT)



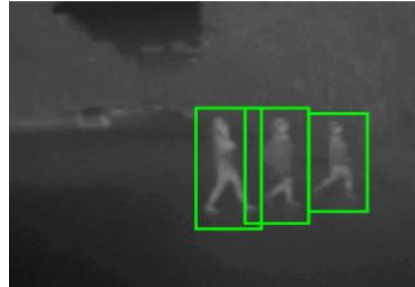
(a)



(b)



(c)



(d)

# Content

- 
- 1. Problem Statement
  - 2. Research Identification
  - 3. Methodology
  - 4. My Progress
  - 5. Intermediate Result
  - 6. What to do next

## 2. Research Identification

### - Research Objective

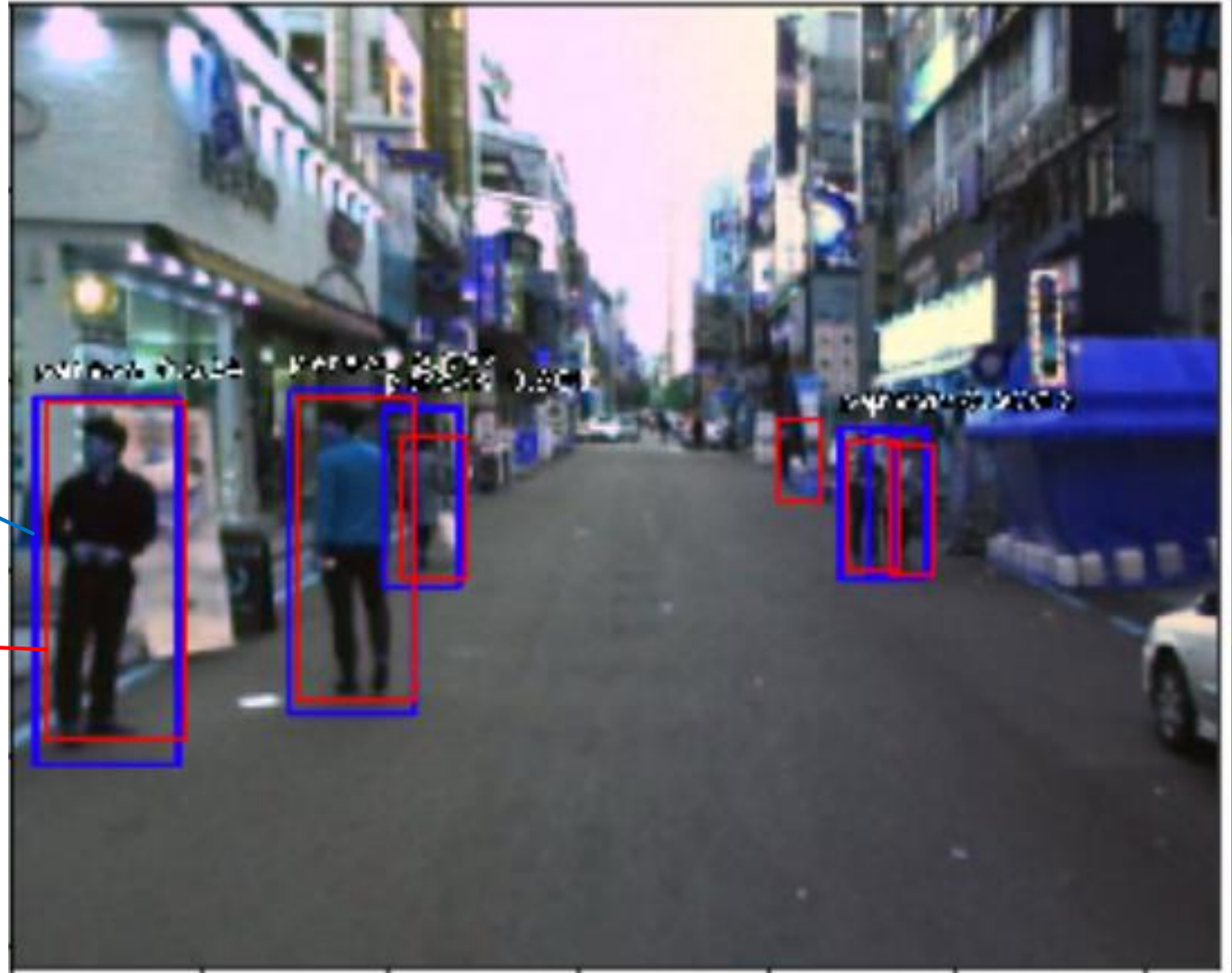
- ▶ Test: different sensor fusion architectures
- ▶ Approach: convolutional neural network (CNN)
- ▶ Applied in: human detection
- ▶ Input: thermal and visual images
- ▶ Provide: the most accurate architecture

## 2. Research Identification

### ► Output: bounding box

Blue bounding box:  
prediction

Red bounding box:  
annotation (ground  
truth)

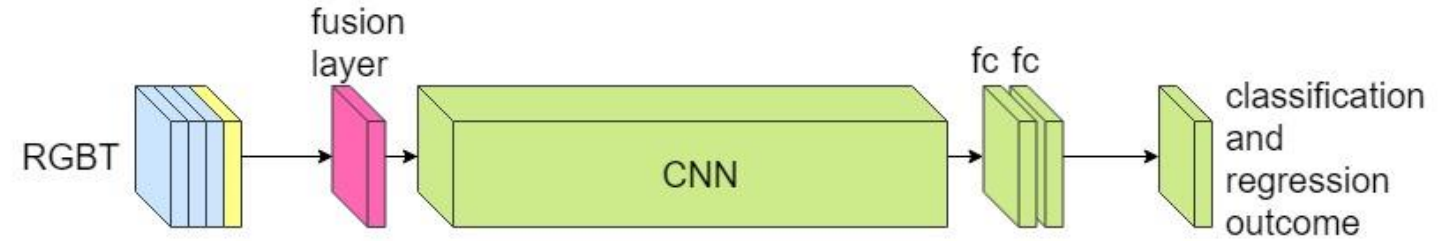


# Content

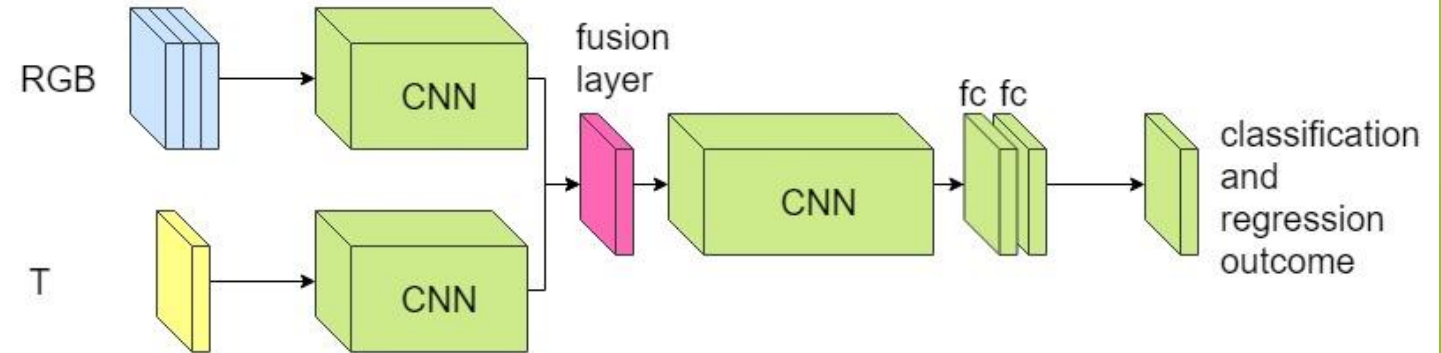
- 
- 1. Problem Statement
  - 2. Research Identification
  - 3. Methodology
  - 4. My Progress
  - 5. Intermediate Result
  - 6. What to do next

# 3. Methodology

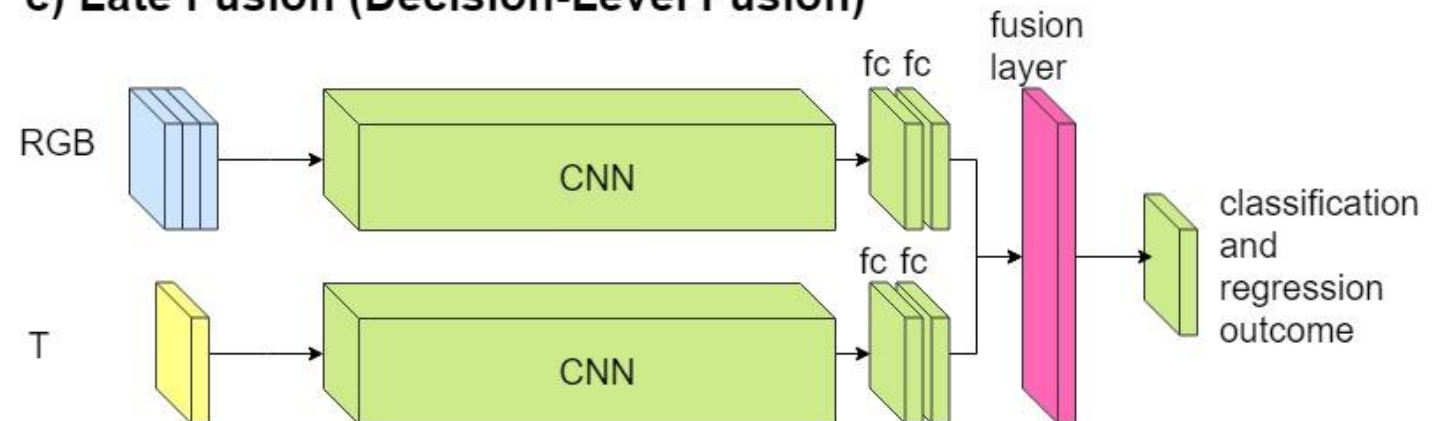
## a) Early Fusion (Pixel-Level Fusion)



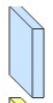



## b) Half-way Fusion (Feature-Level Fusion)



## c) Late Fusion (Decision-Level Fusion)

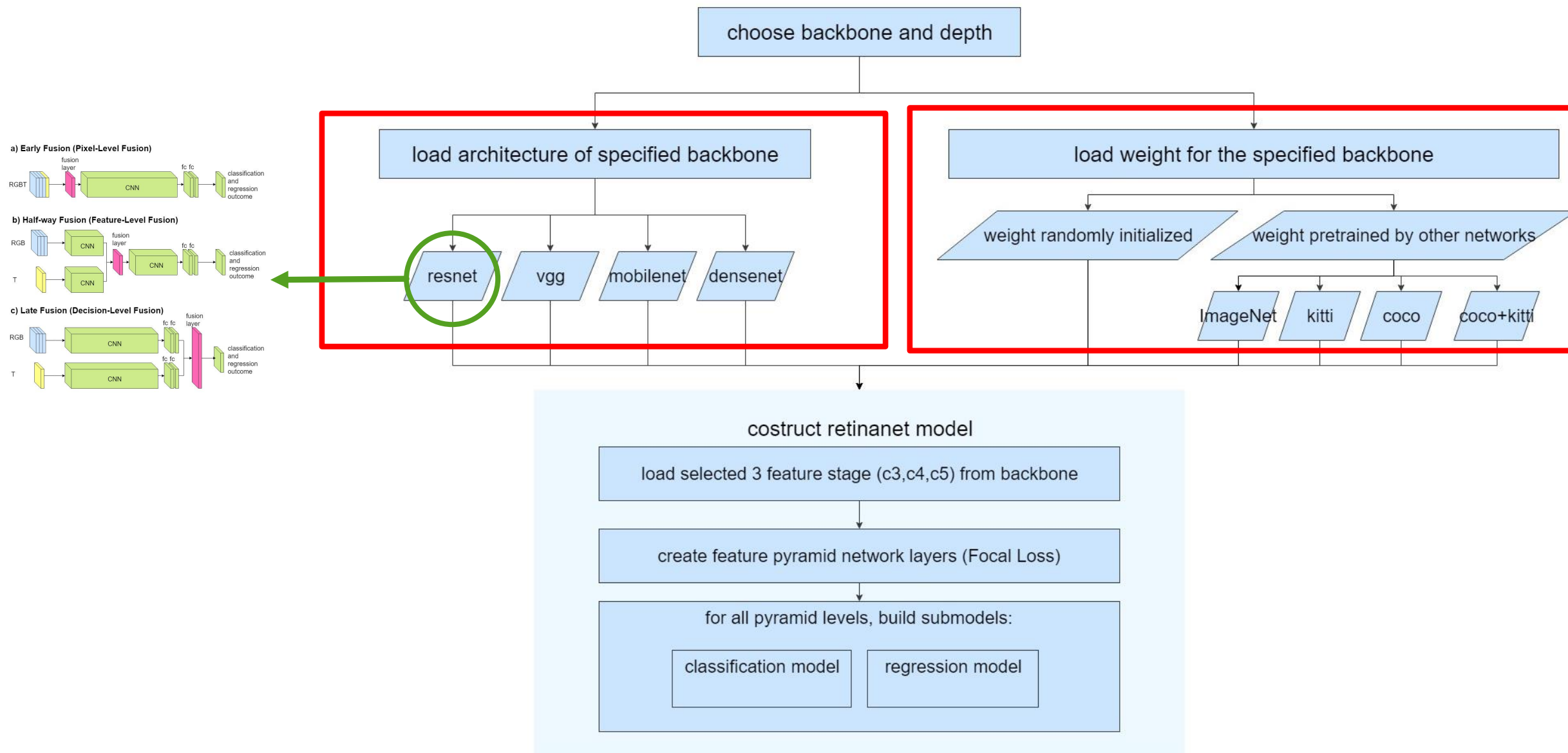


### legend

-  a channel in visual images
-  a channel in thermal images
-  a fusion layer (concatenate layer)
-  a convolutional layer, a fully connected layer (fc), a classification output layer or a regression output layer

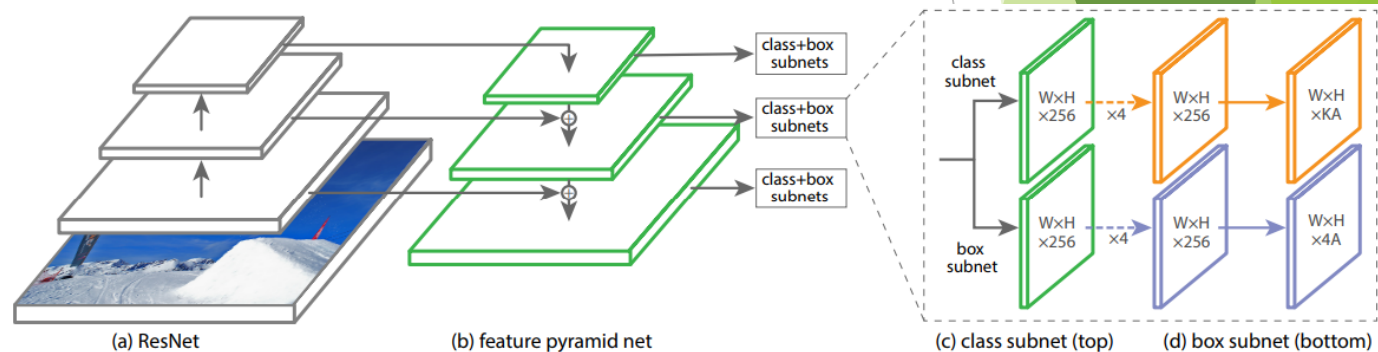


# 3. Methodology-Retinanet



# Why choose retinanet?

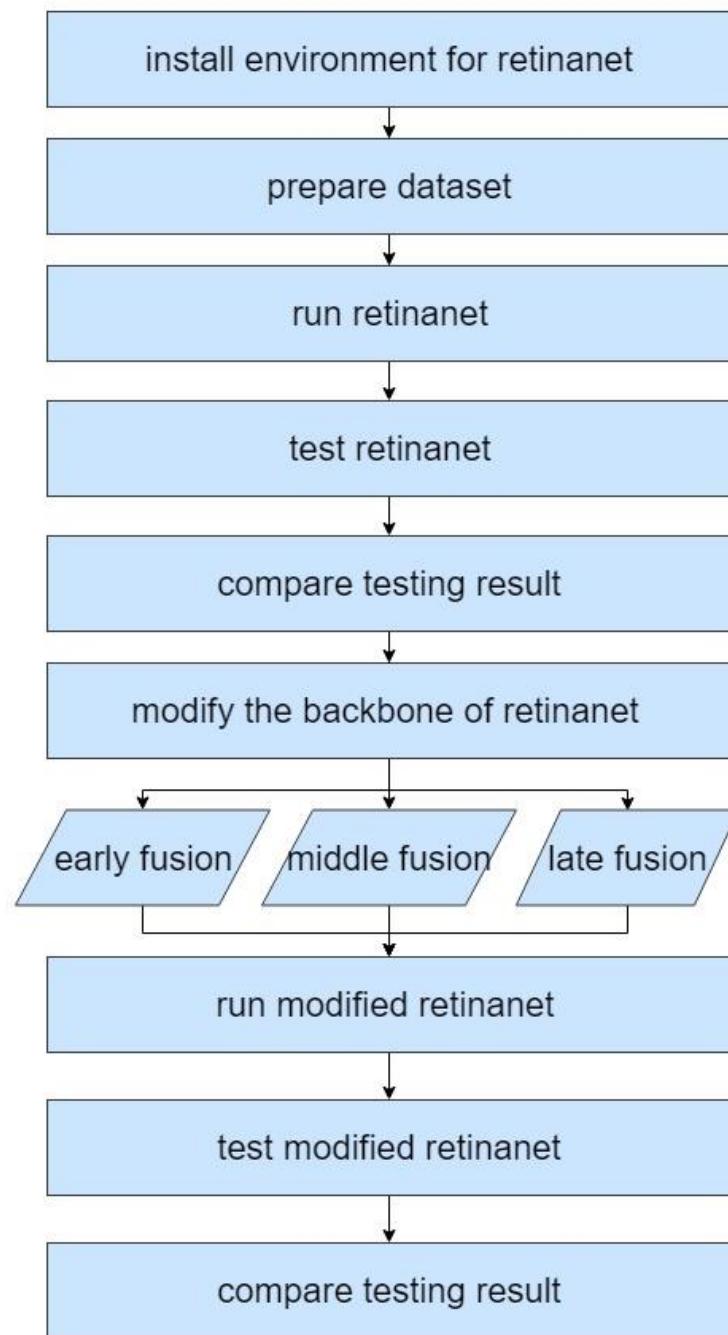
- ▶ solve two problems:
    - ▶ Problem 1: class imbalance
    - ▶ Solution: Focal Loss
  - ▶ Problem 2: high resolution maps (earlier layers) have low-level features
  - ▶ Solution: feature pyramid network
- 
- ▶ Speed vs Accuracy on Coco
    - ▶ surpass yolo, SSD, DSSD, R-FCN (Lin,T., Goyal P., Girshick,R, 2018)



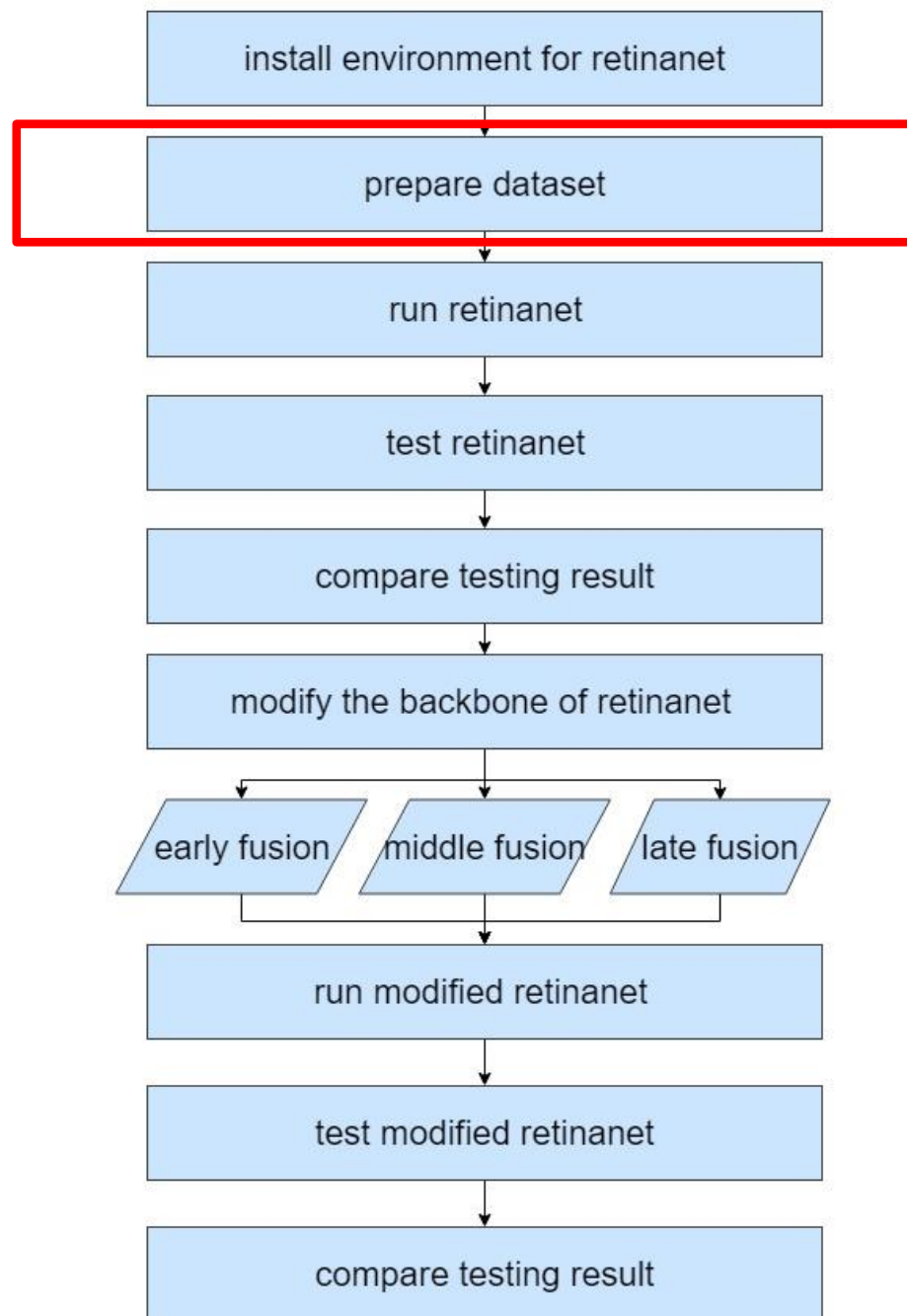
# Content

- 
- 1. Problem Statement
  - 2. Research Identification
  - 3. Methodology
  - 4. My Progress
  - 5. Intermediate Result
  - 6. What to do next

# Workflow



# Workflow



# Data Description

- ▶ KAIST (Korea Advanced Institute of Science and Technology)
- ▶ Image size: 512 pixels \* 640 pixels
- ▶ Total: 95k colour-thermal pairs
- ▶ Total: 1182 annotated pedestrians

## Train 50k

- Set 00** / Day / Campus / 17,498 images
- Set 01** / Day / Road / 8,035 images
- Set 02** / Day / Downtown / 7,866 images
- Set 03** / Night / Campus / 6,668 images
- Set 04** / Night / Road / 7,200 images
- Set 05** / Night / Downtown / 2,920 images

Validation: will be extracted

## Test 45k

- Set 06** / Day / Campus / 12,988 images
- Set 07** / Day / Road / 8,141 images
- Set 08** / Day / Downtown / 8,050 images
- Set 09** / Night / Campus / 3,500 images
- Set 10** / Night / Road / 8,902 images
- Set 11** / Night / Downtown / 3,560 images

Improved annotation for test: 2250 images

set06-11 1/20 in each set

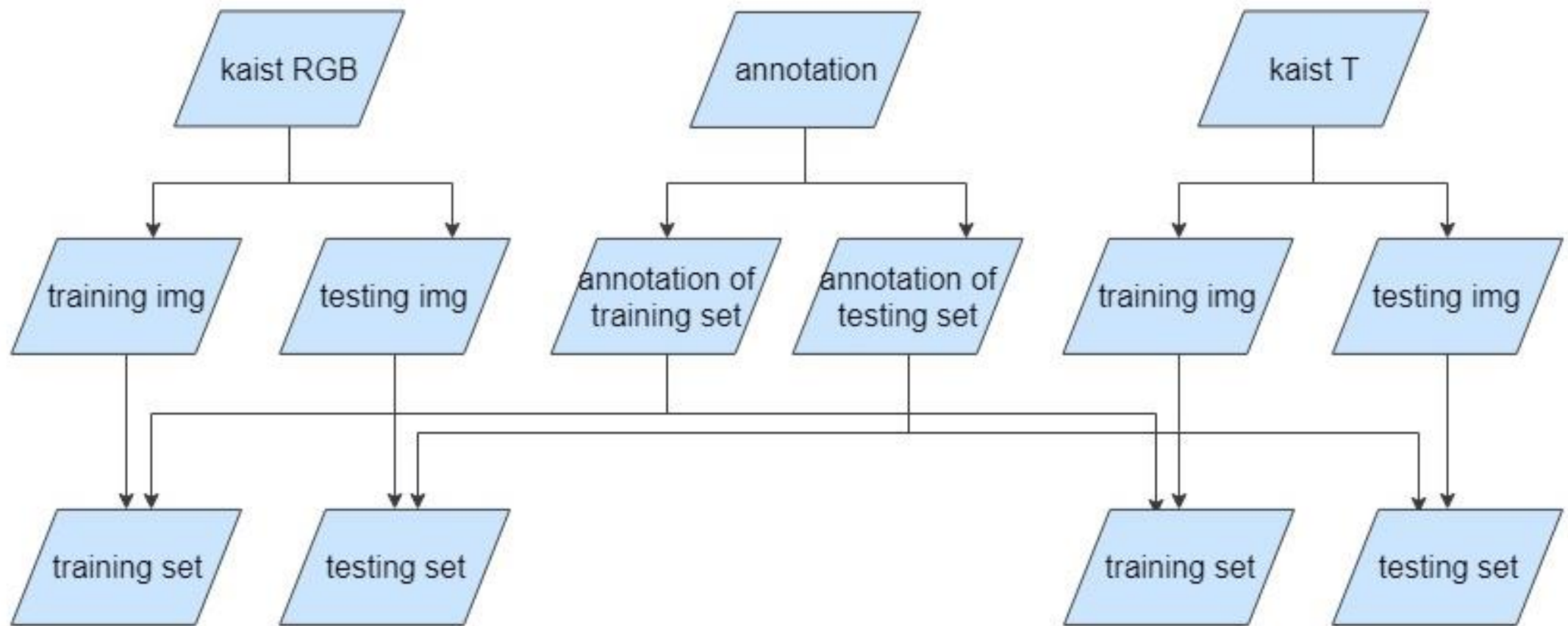


(a)



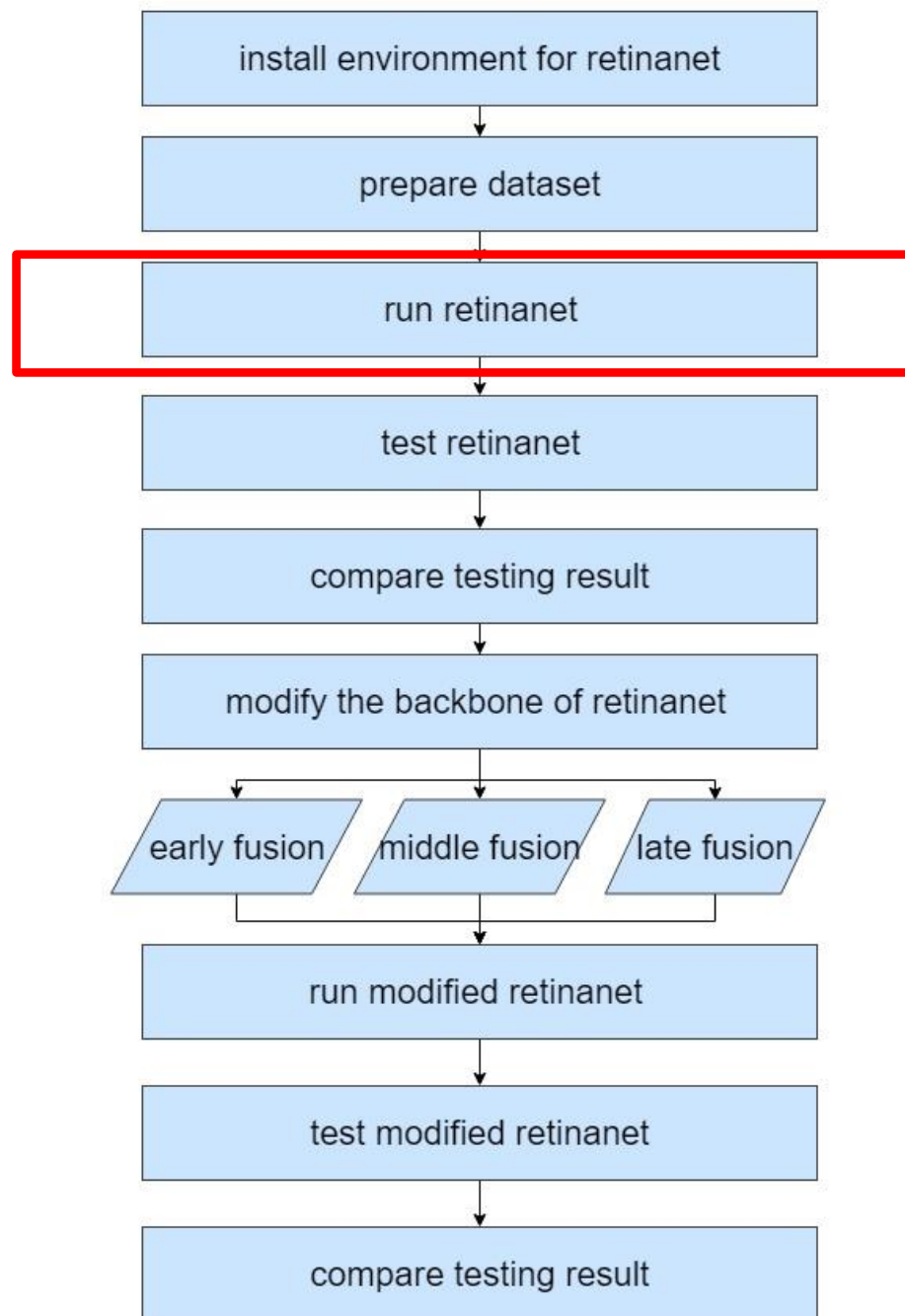
(b)

# Prepare dataset



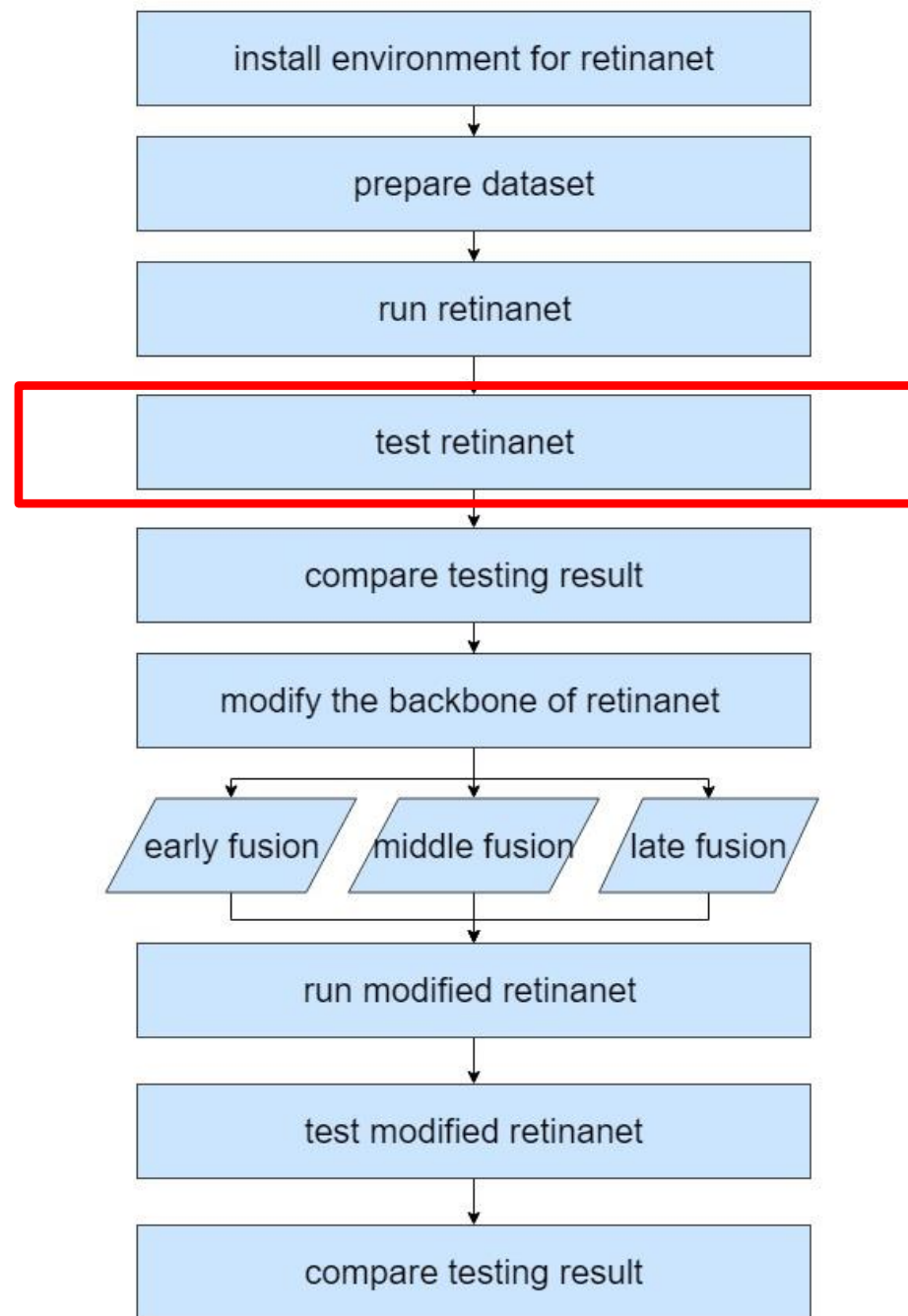


# Workflow



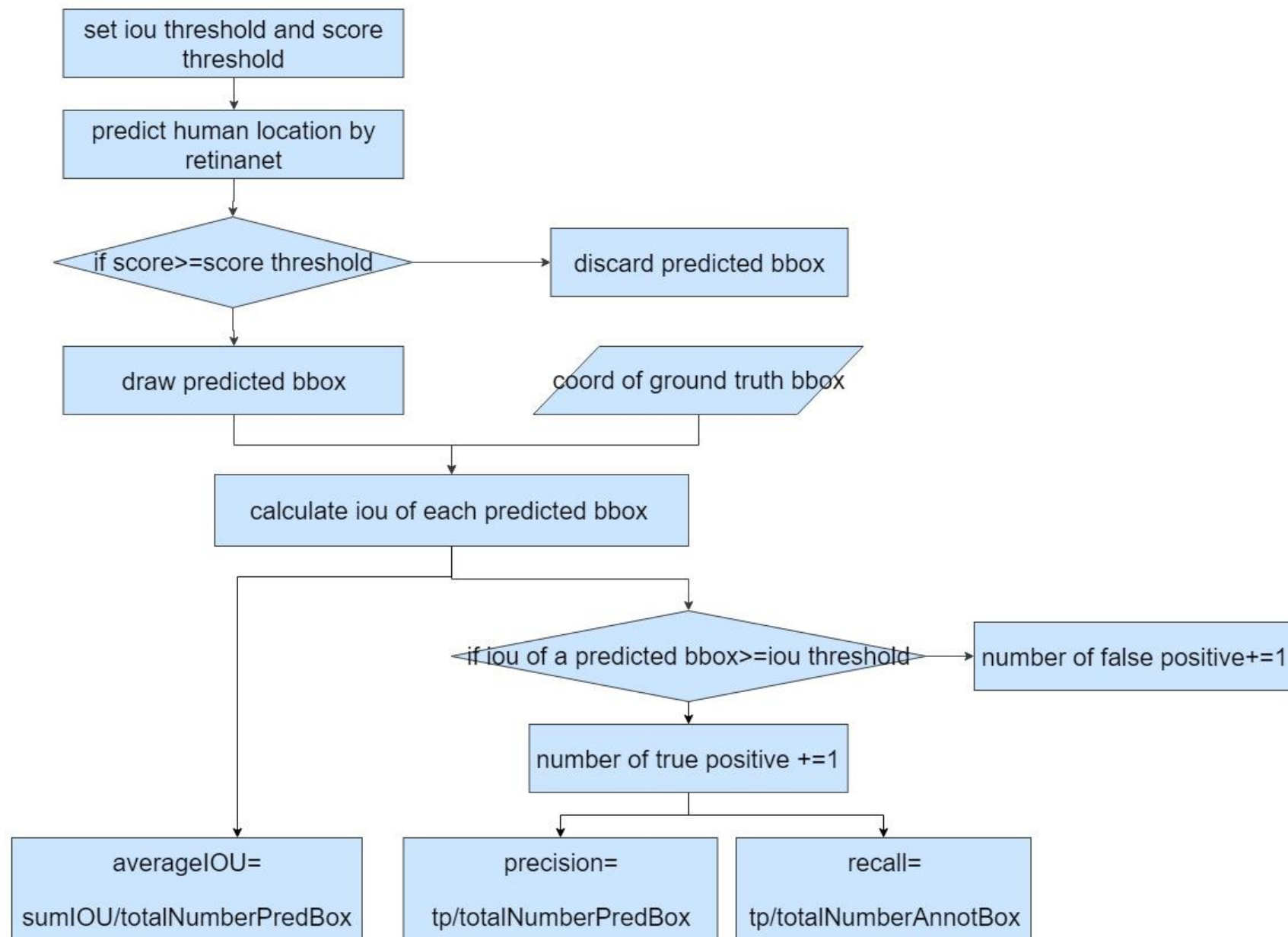


# Workflow

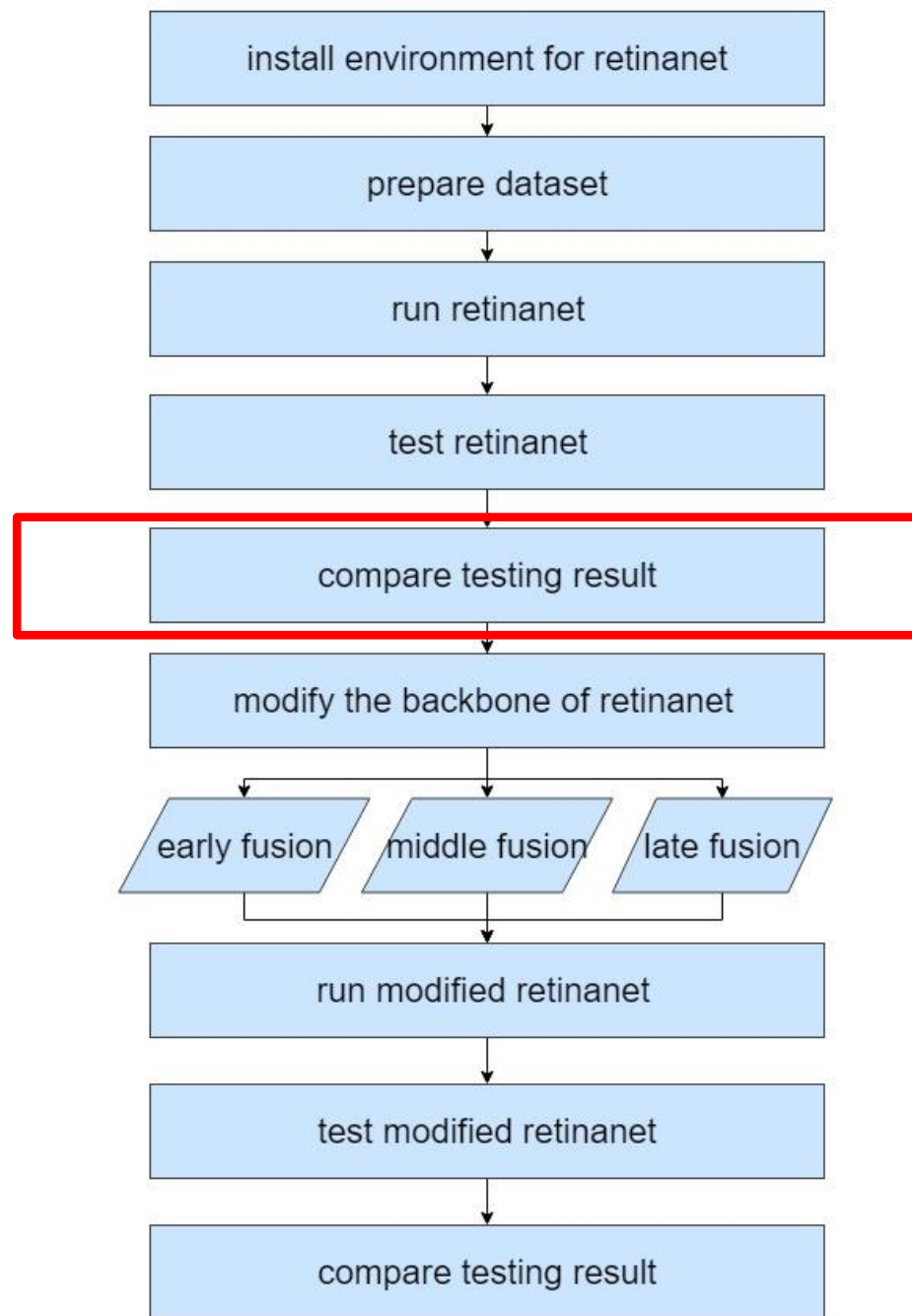


# Test Retinanet


- Manual test
  - Average IOU
  - Precision
  - recall
- Evaluate.py



# Workflow

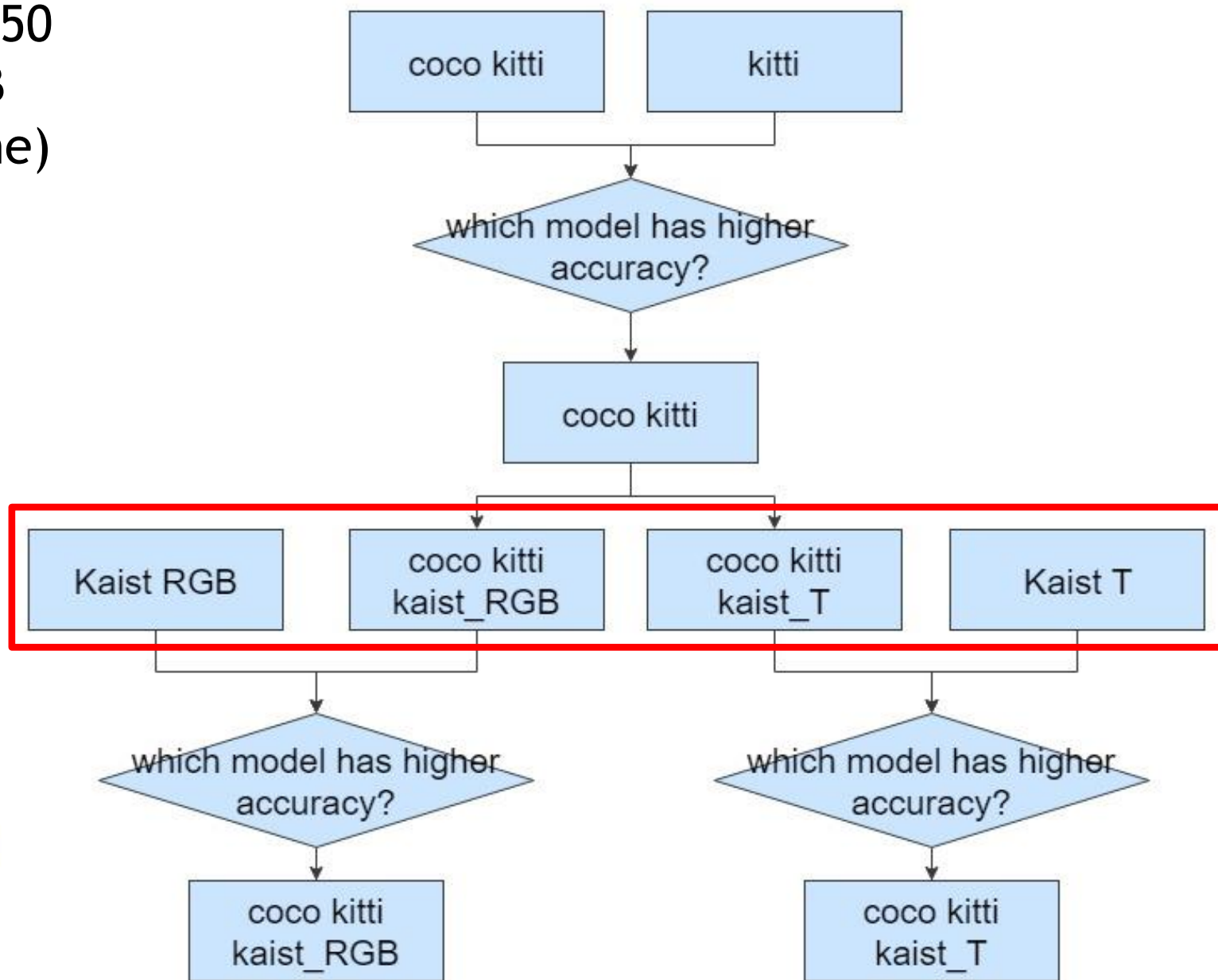
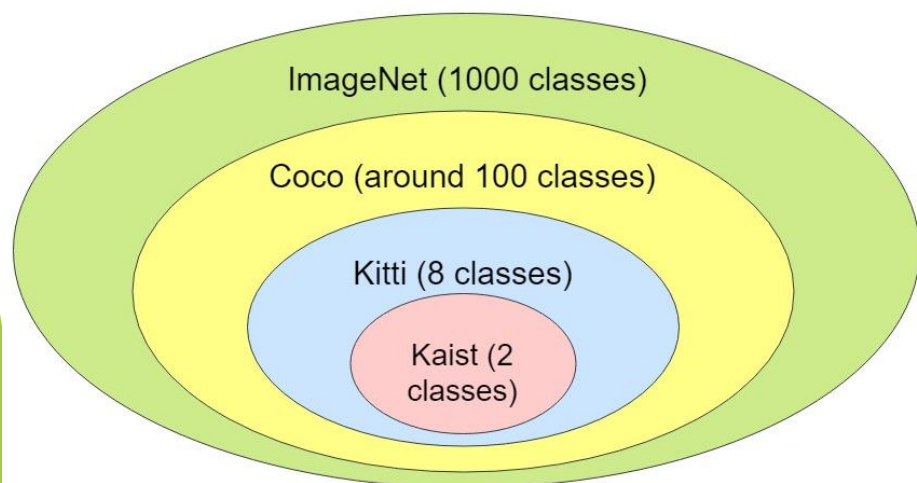


# Content

- 
1. Problem Statement
  2. Research Identification
  3. Methodology
  4. My Progress
  5. Intermediate Result
  6. What to do next

Backbone of all models: Resnet 50  
ImageNet (default)+KaistT/RGB  
Coco+Kitti+KaistT/RGB (finetune)

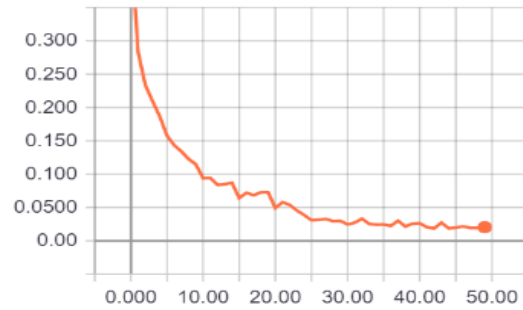
	cocokitti	kitti
<b>Pedestrian with average precision</b>	0.44	0.34
<b>Cyclist with average precision</b>	0.46	0.39



# Classification Loss and Regression Loss

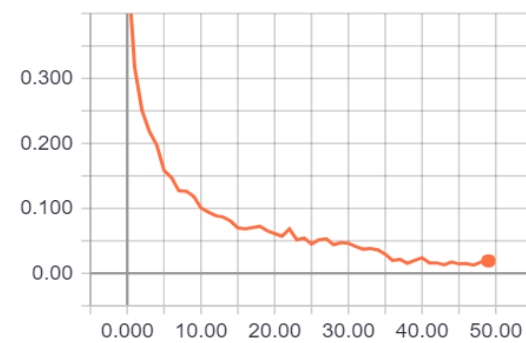
KaistT

classification\_loss



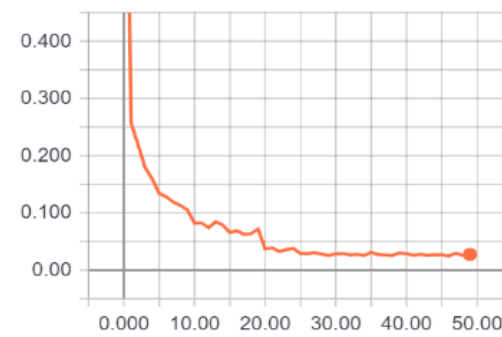
KaistRGB

classification\_loss



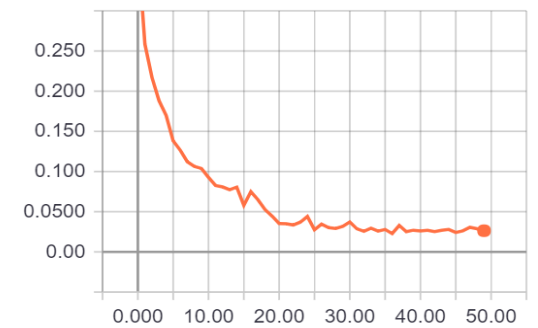
CocoKittiKaistT

classification\_loss

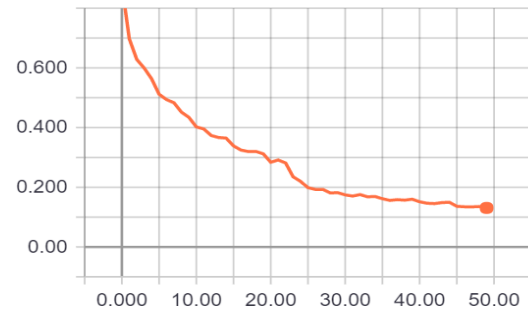


CocoKittiKaistRGB

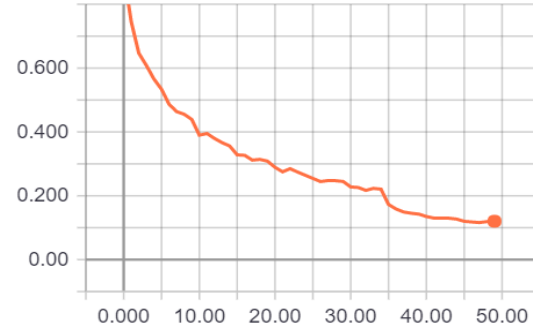
classification\_loss



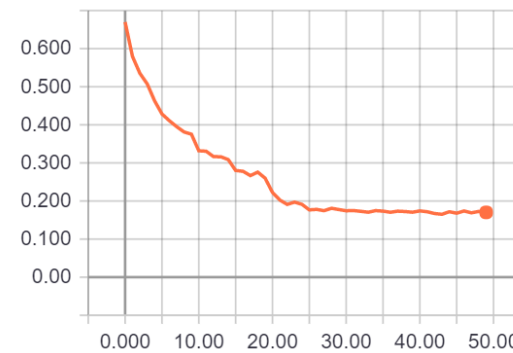
regression\_loss



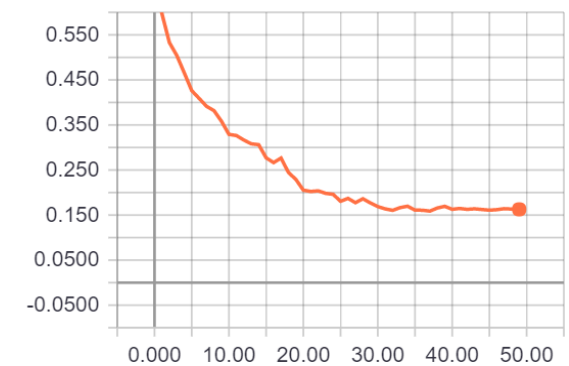
regression\_loss



regression\_loss



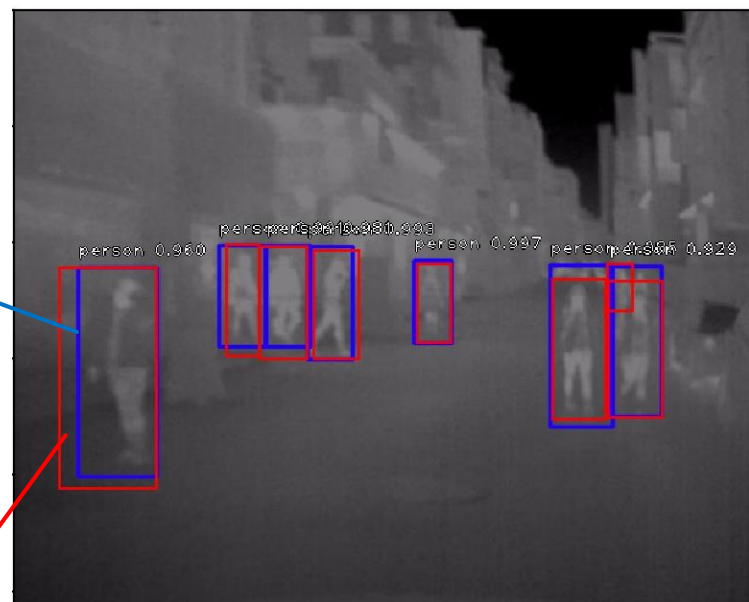
regression\_loss



Will include validation loss



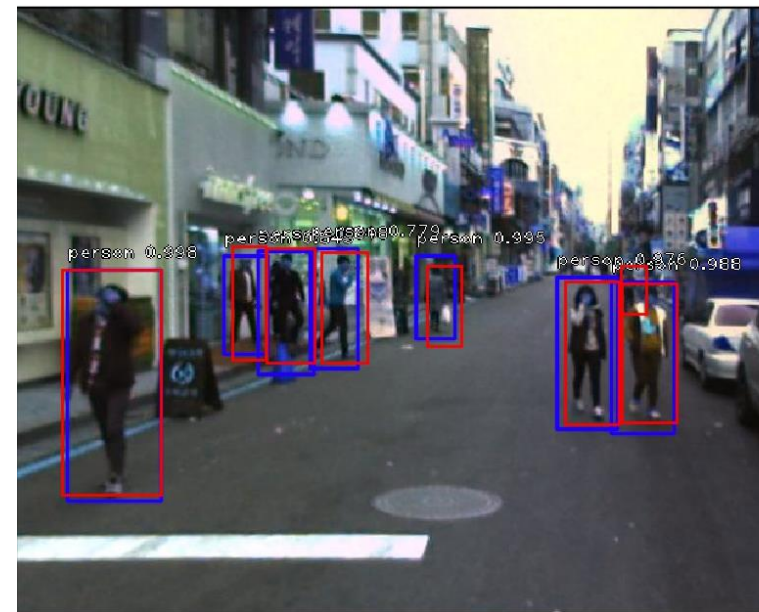
KaistT



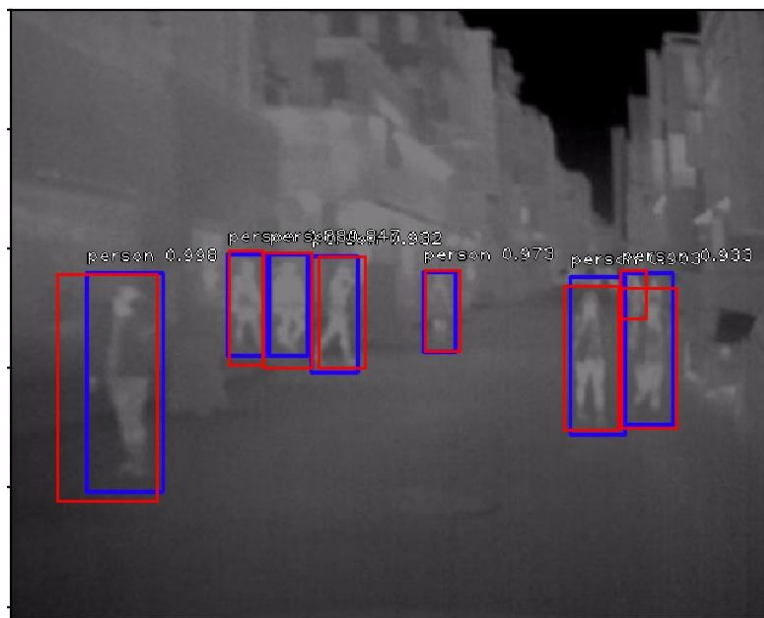
Blue bounding box:  
prediction

Red bounding box:  
annotation  
(ground truth)

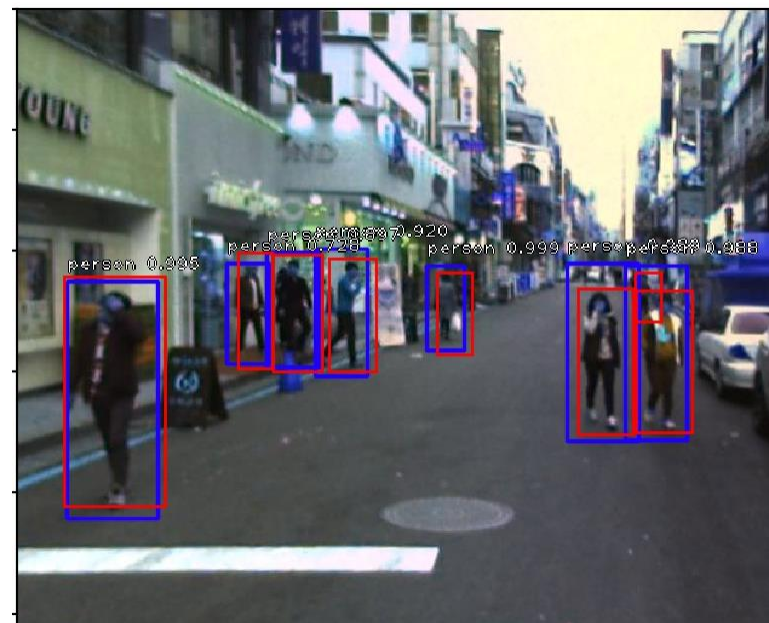
KaistRGB



# CocoKittiKaistT

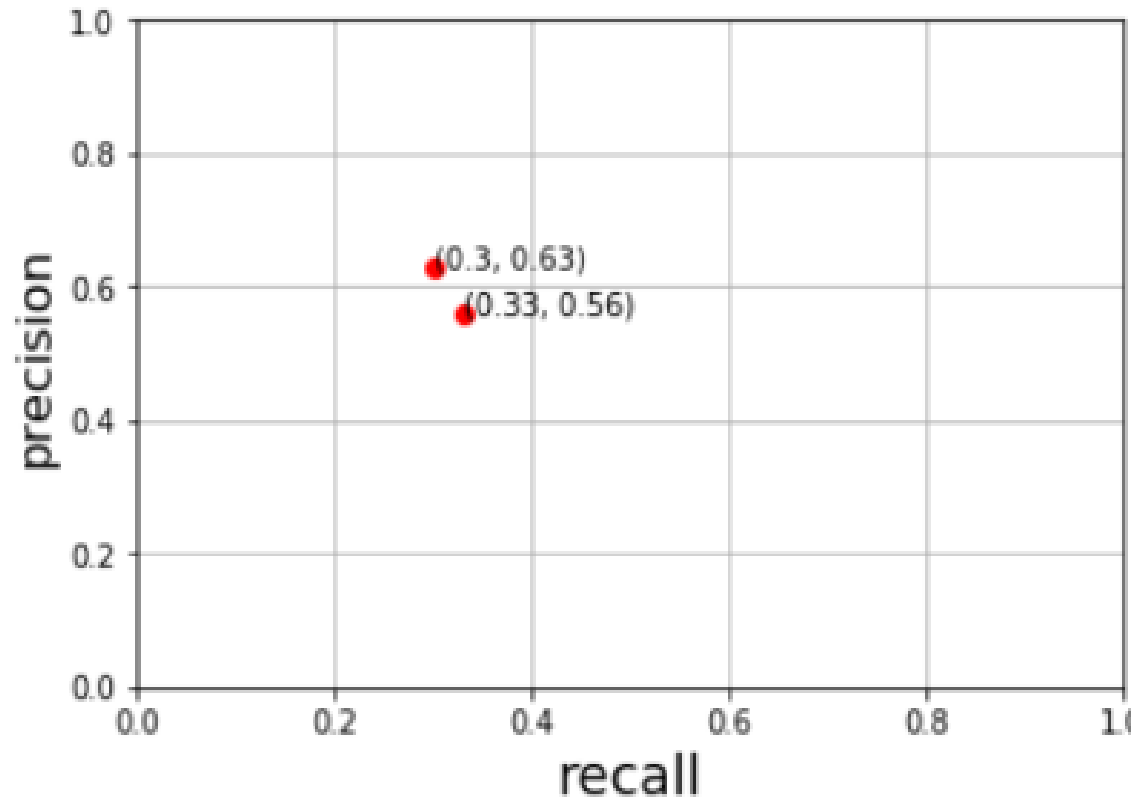


## CocoKittiKaistRGB



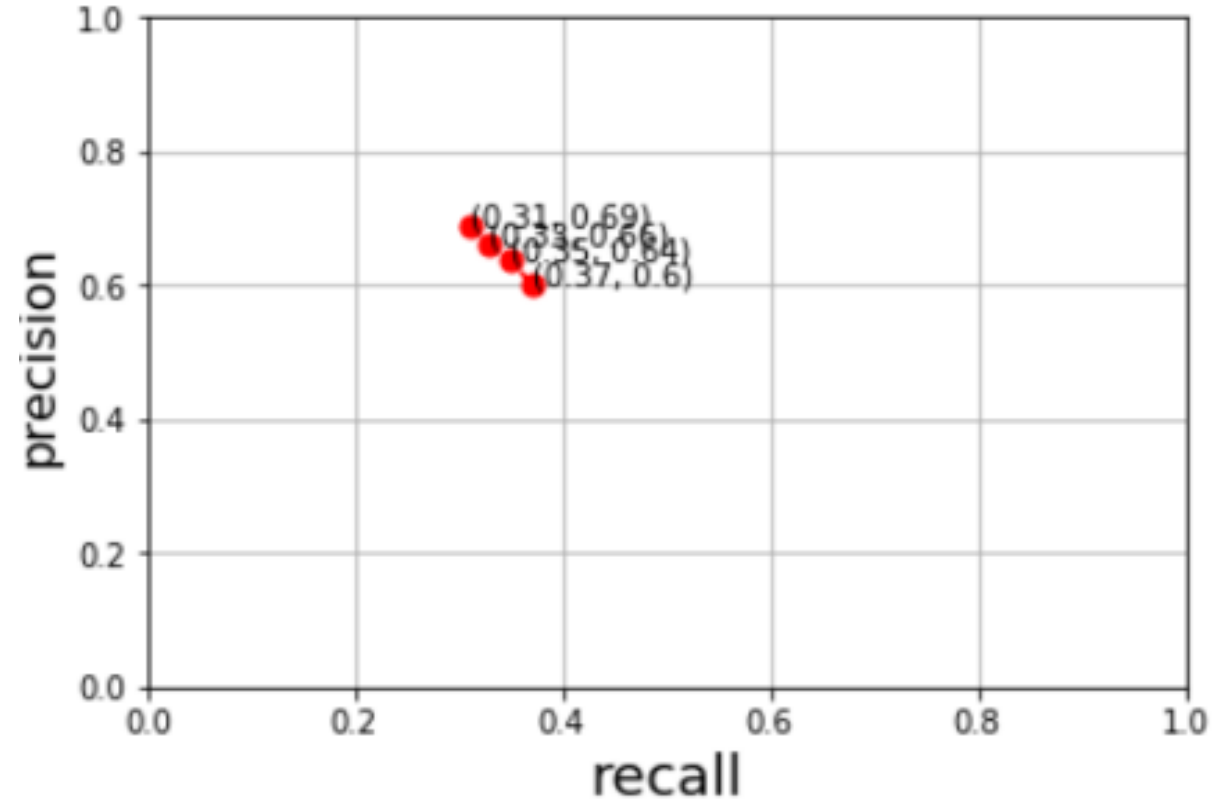
# Precision-Recall Graph when iou\_threshold=0.7

KaistT



Score>=0.5 recall:0.33 precision: 0.56  
Score>=0.6 recall: precision:  
Score>=0.7 recall: precision:  
Score>=0.8 recall:0.30 precision: 0.63

CocoKittiKaistT

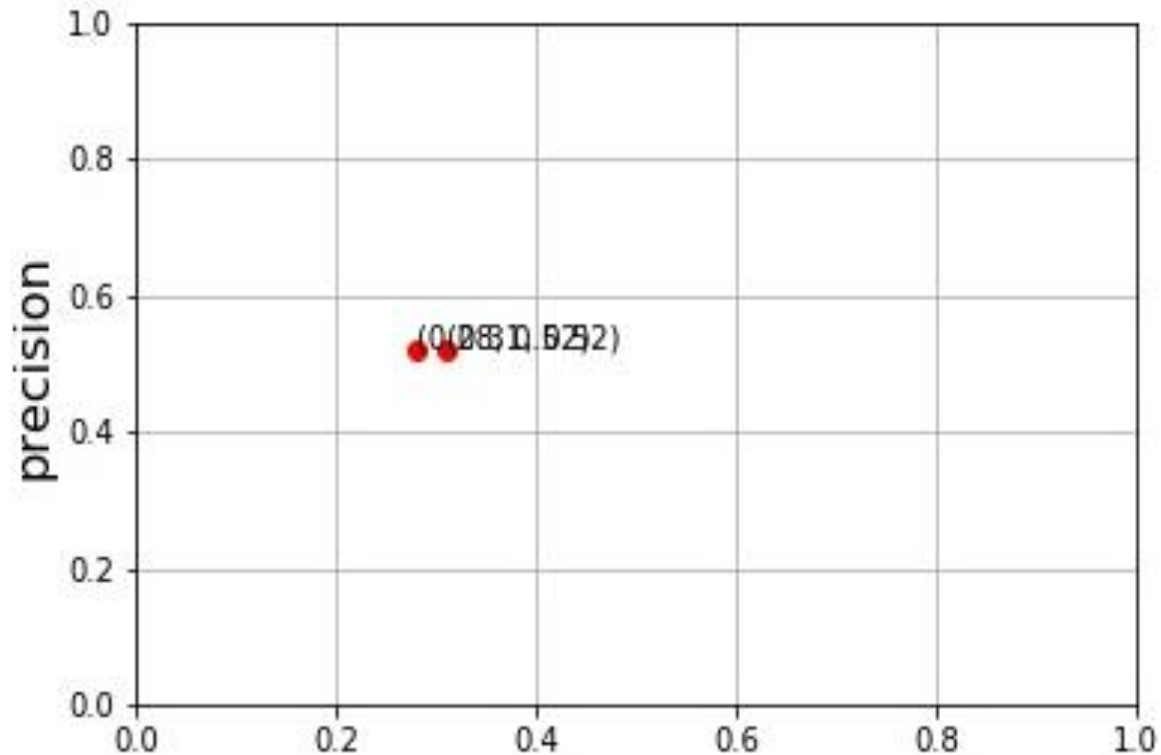


Score>=0.5 recall:0.37 precision: 0.60  
Score>=0.6 recall:0.35 precision: 0.64  
Score>=0.7 recall:0.33 precision: 0.66  
Score>=0.8 recall:0.31 precision: 0.69



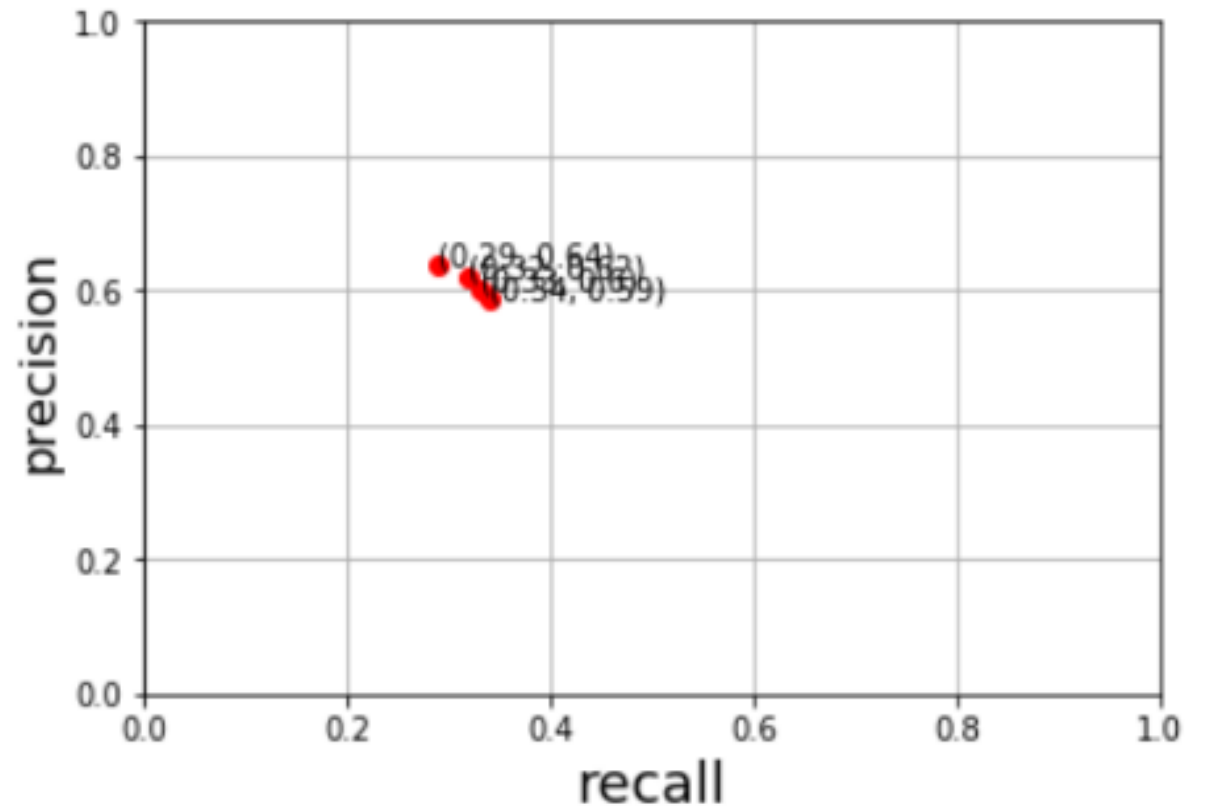
# Precision-Recall Graph when iou\_threshold=0.7

KaistRGB



Score $\geq$ 0.5	recall:0.31	precision: 0.52
Score $\geq$ 0.6	recall:	precision:
Score $\geq$ 0.7	recall:	precision:
Score $\geq$ 0.8	recall:0.28	precision: 0.52

CocoKittiKaistRGB



Score $\geq$ 0.5	recall:0.34	precision: 0.59
Score $\geq$ 0.6	recall:0.33	precision: 0.60
Score $\geq$ 0.7	recall:0.32	precision: 0.62
Score $\geq$ 0.8	recall:0.29	precision: 0.64

## Average IOU when IOU threshold=0.7

	kaistT	kaistRGB	cocokittikaistT	cocokittikaistRGB
score threshold=0.5	0.693	0.662	0.704	0.692
score threshold=0.8	0.724	0.690	0.739	0.715

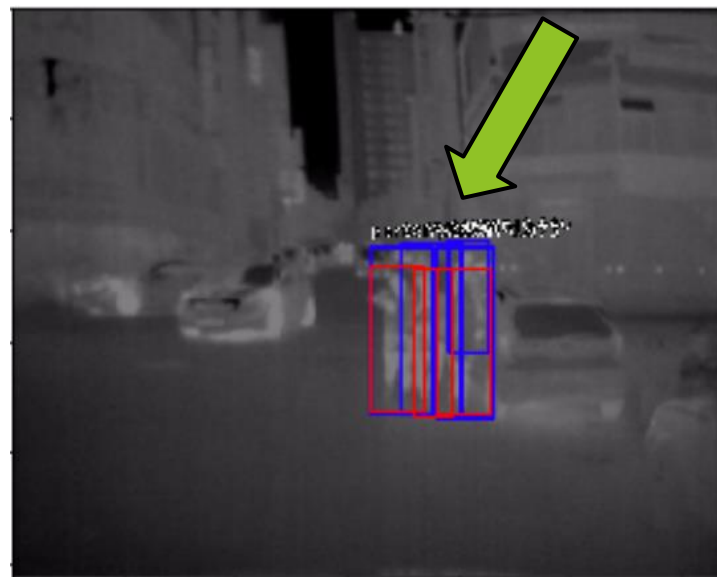
Train: set 0-5

Test: set 08/v000 with improved annotation

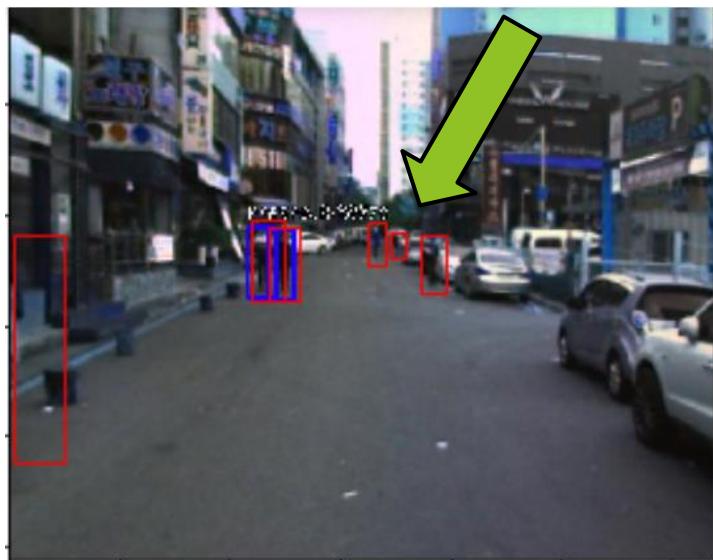
Half person



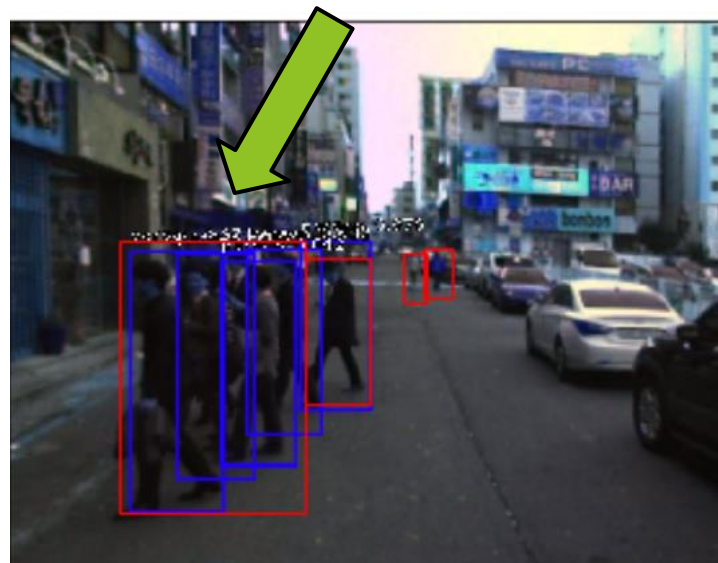
occluded person  
(able to detect but  $iou < \text{threshold}$ )



Persons are too small



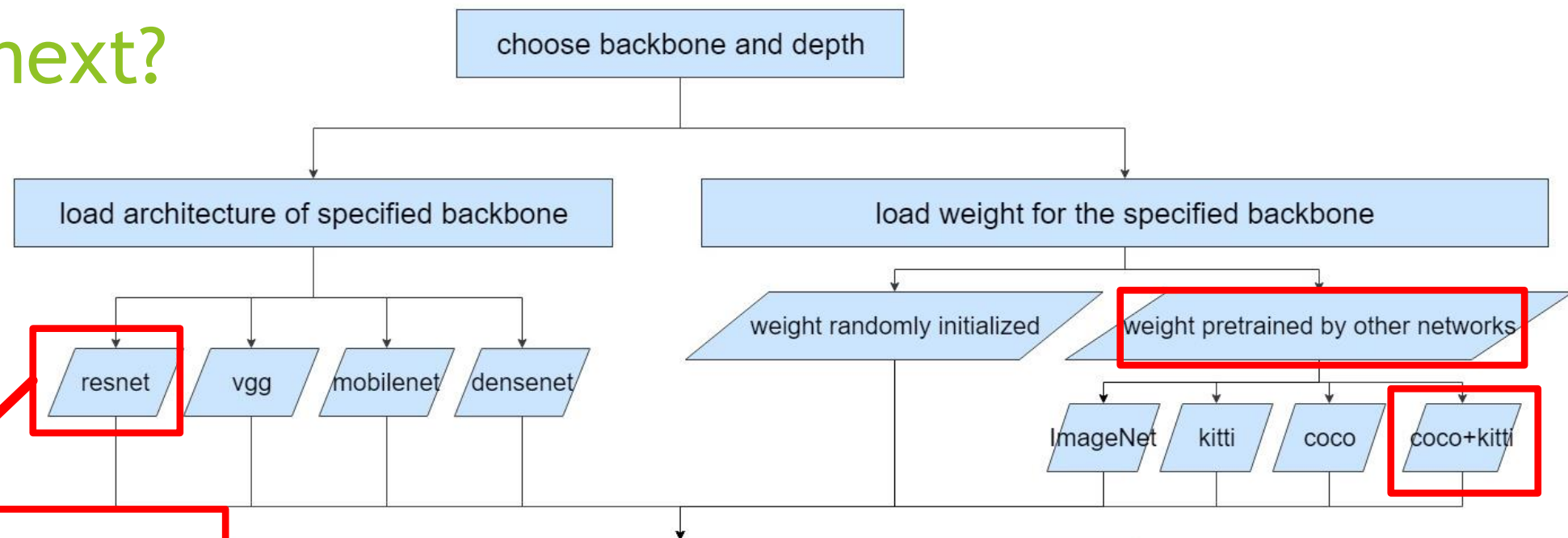
Annotation is imperfect



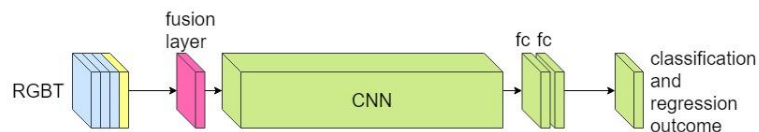
# Content

- 
1. Problem Statement
  2. Research Identification
  3. Methodology
  4. My Progress
  5. Intermediate Result
  6. What to do next

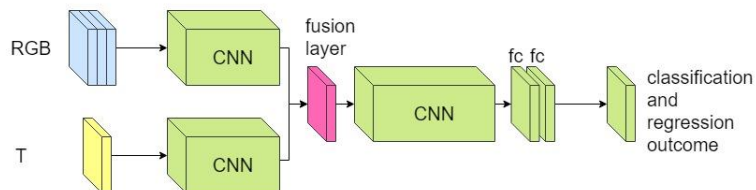
# What to do next?



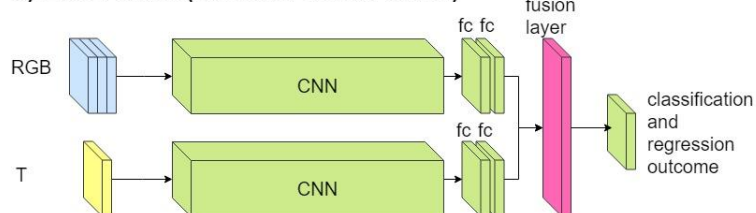
## a) Early Fusion (Pixel-Level Fusion)



## b) Half-way Fusion (Feature-Level Fusion)



## c) Late Fusion (Decision-Level Fusion)



## construct retinanet model

load selected 3 feature stage (c3,c4,c5) from backbone

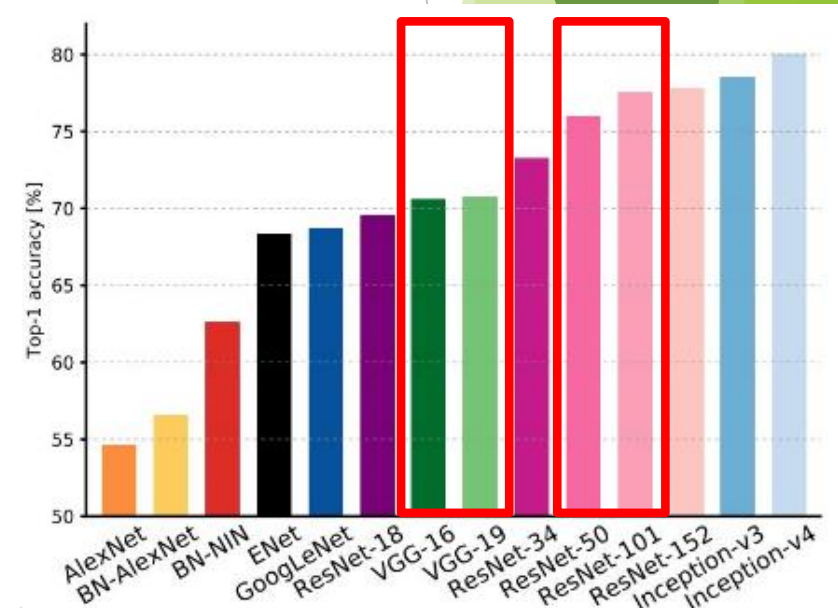
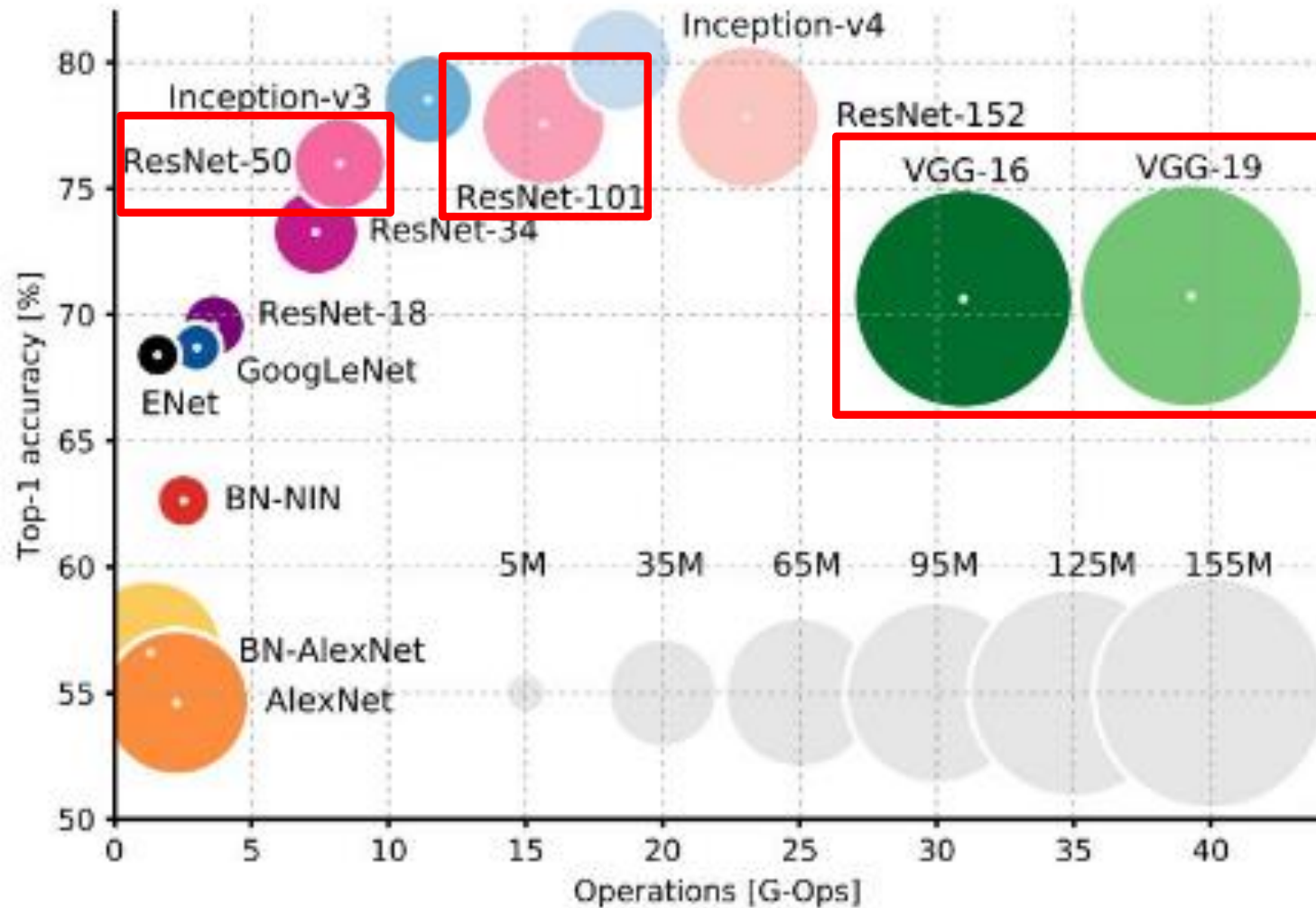
create feature pyramid network layers (Focal Loss)

for all pyramid levels, build submodels:

classification model

regression model

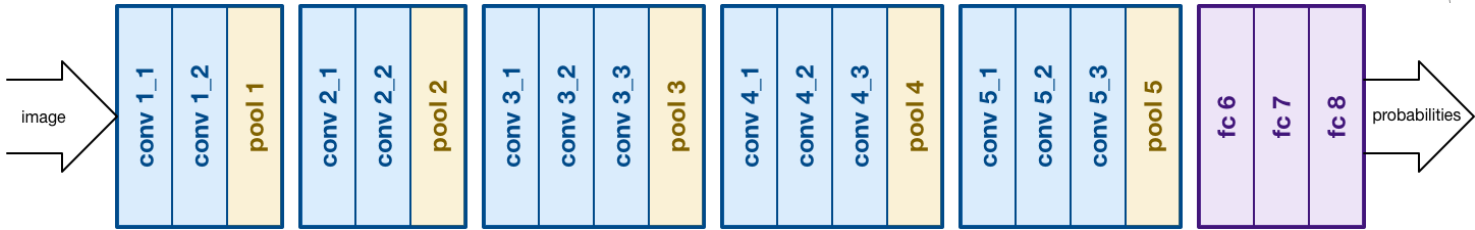
# Which backbone to choose? Resnet, VGG



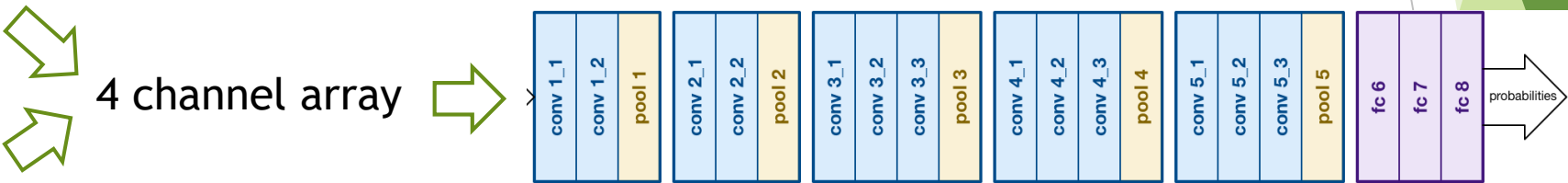


# I am working on:

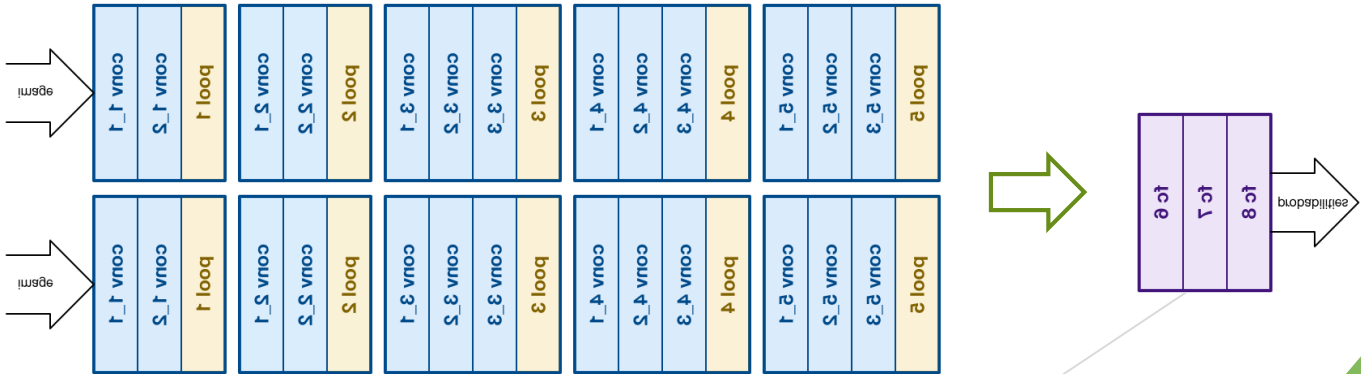
VGG original



VGG early



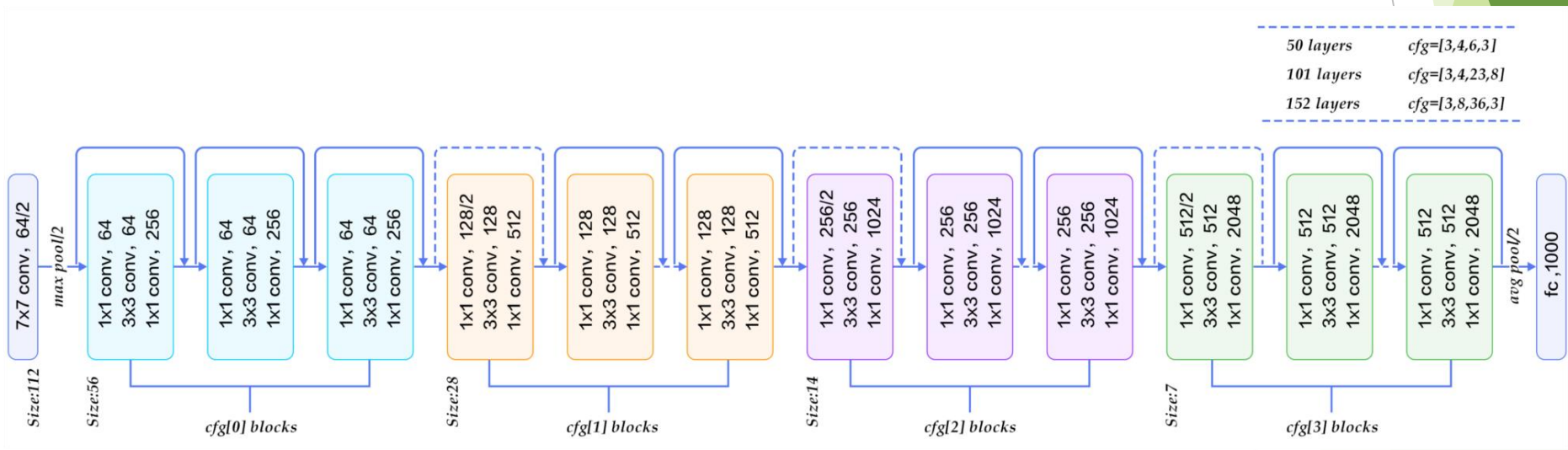
VGG late



VGG



resnet





# What to do next?

- ▶ Training:
  - ▶ Modify backbone architecture:
    - ▶ Firstly vgg, then resnet
    - ▶ Firstly early, then late, finally halfway
  - ▶ Find out optimal value of hyperparameters
- ▶ Validation: take 20% data from set 0-5
- ▶ Testing:
  - ▶ Whole testing dataset
  - ▶ Understand why Evaluate.py gives very low mAP
  - ▶ Evaluate computational complexity
- ▶ report



# hyperparameters

- ▶ hyperparameters for building up network:
  - ▶ 1) Number of hidden layers and units
  - ▶ 2) Weight initialization
- ▶ hyperparameters to train network:
  - ▶ 1) Learning rate: speed vs convergence
  - ▶ 2) Number of epochs: observe validation accuracy
  - ▶ 3) Batch size: speed vs accuracy