

# USING MOBILE NEURAL NETWORK FOR PETS CLASSIFICATION

Jie Wen 1155092211

*Supervised by Prof. Michael R. LYU*

## MAIN TOPIC IN THE 2ND SEMESTER

1. Similar Species Classification
2. Object Detection
3. Simple UI Improve & Explore

# MOTIVATION

## Why Mobile

General public

Number of User

Cost of time/money

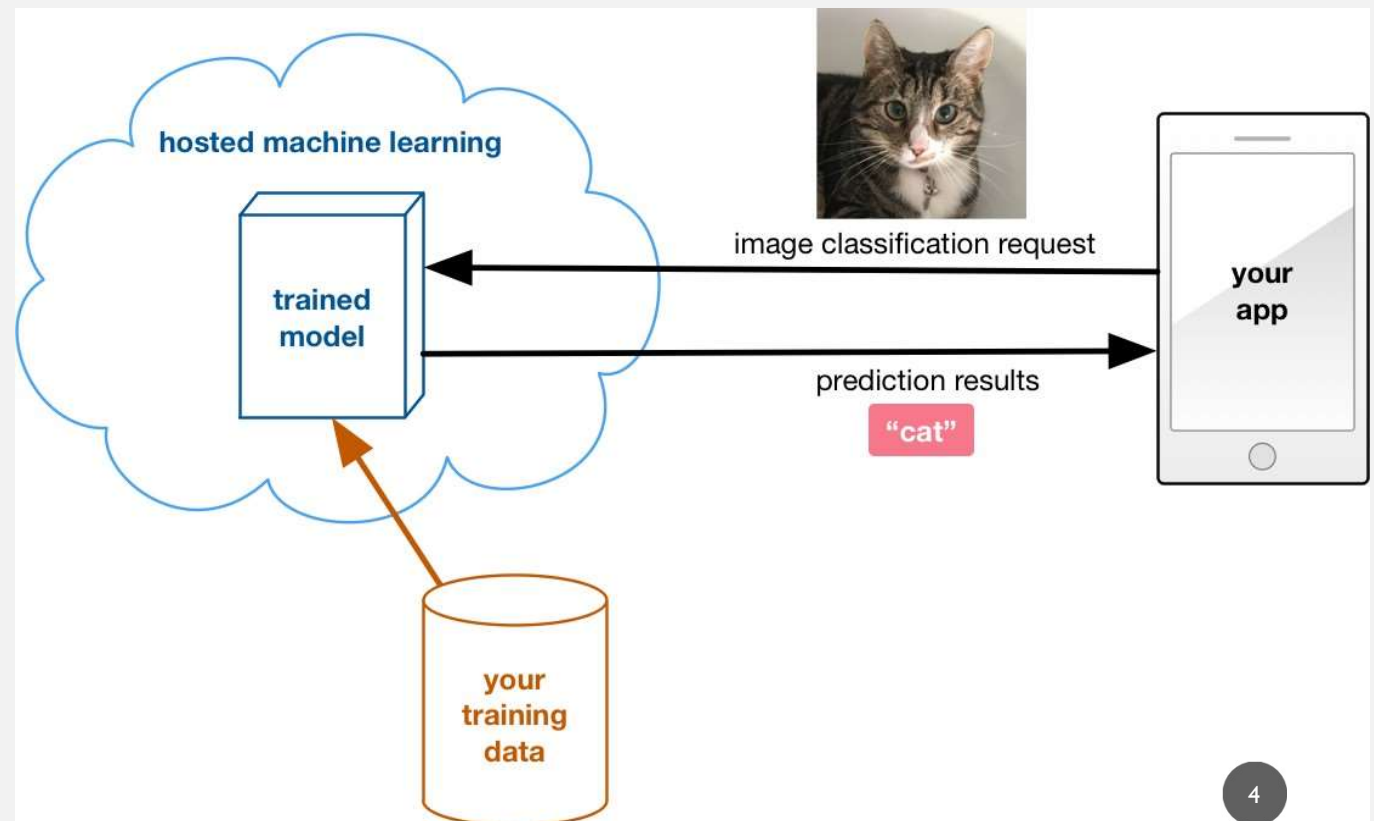
Wide usage scenario

# MOTIVATION

Mobile implementation

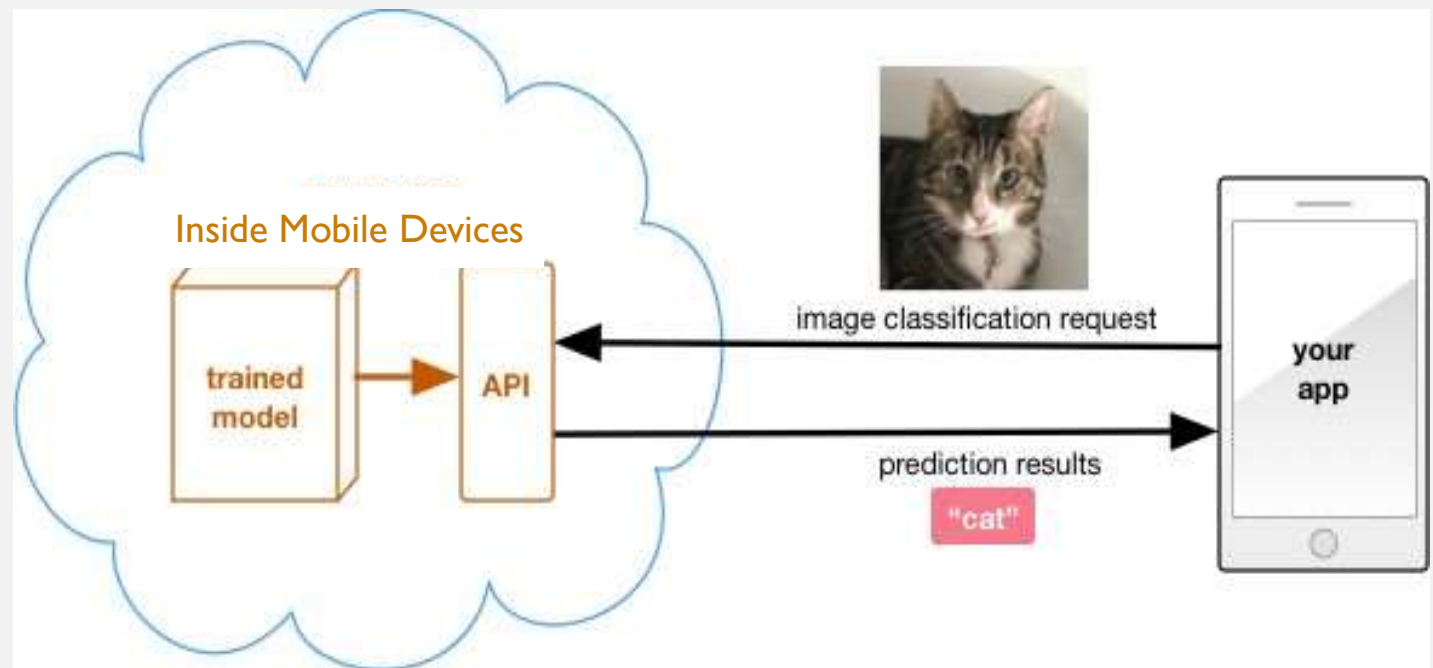
Need cloud server

Need Network



# MOTIVATION

## Mobile implementation



# DATASET

Last term:

BotanWiki -> AnimalWiki

(Like ImageNet: dog, cat, cow, bird...)

AnimalWiki -> PetWiki

(Shiba, Husky, Scottish fold...)

# DATASET

This term:

Crawled Dataset from Internet (22 species)

+ Stanford Dogs Dataset (120 species)

= Dataset 133 (133 species, for similar species classification)

Oxford-IIIT Pets Dataset (37 species, for object detection)

# DATASET 133

## About 23,500 images for 133 species (dogs)

129 dogs and 4 cats

Large dog (31): German Shepherd, Greyhound, Saint Bernard,  
Tibetan Mastiff, Samoyed, Scotch Collie, Husky...

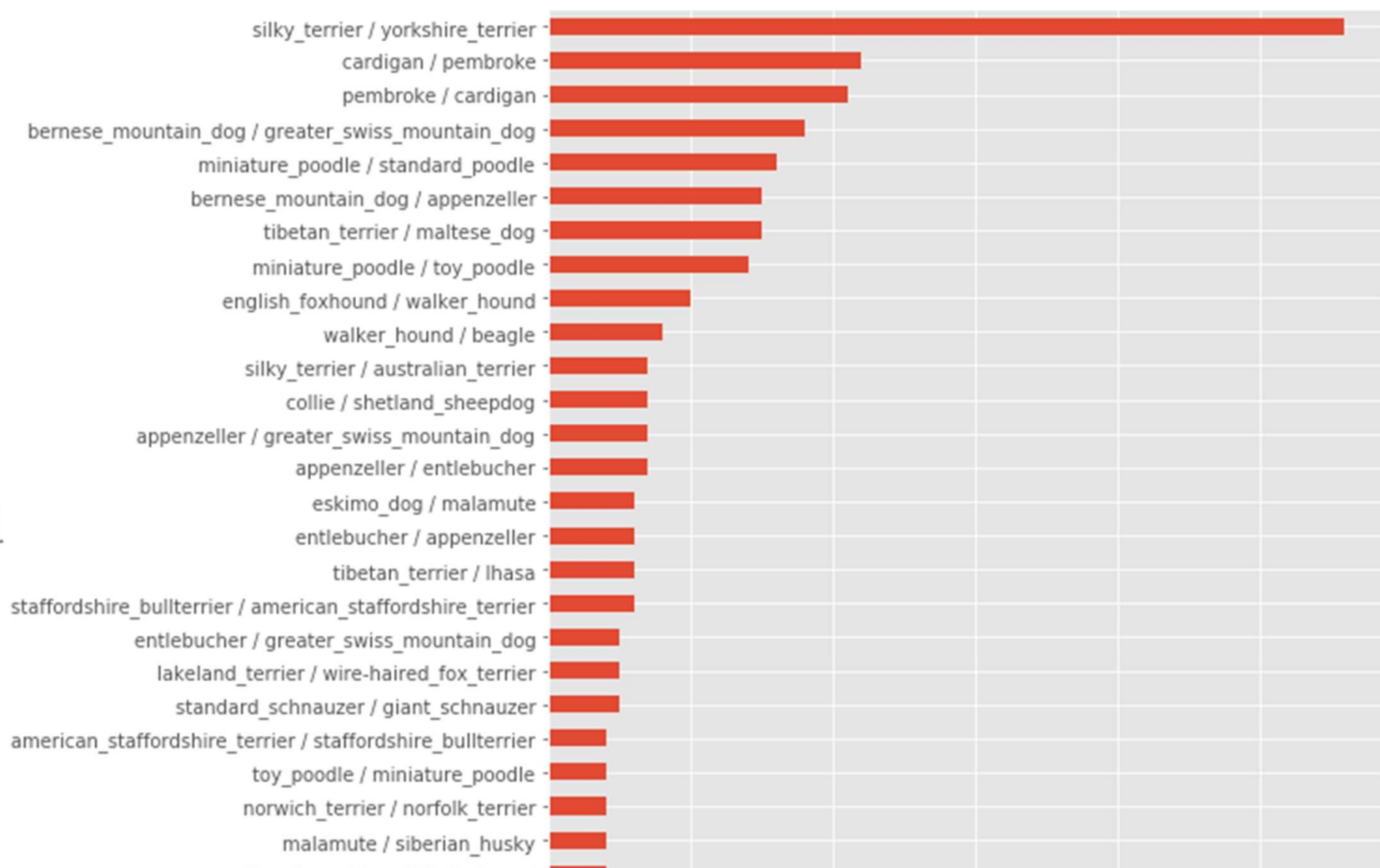
Mid-size dog (53): Shiba, Black Shiba, Border Collie, Dalmatian, Shar Pei, Pug. ..

Small dog (45): Bichon frise, Chihuahua, Corgi, Poodle, Schnauzer...

Cat (4): Bobcat, Persian Cat, Scottish Fold, Siamese Cat.

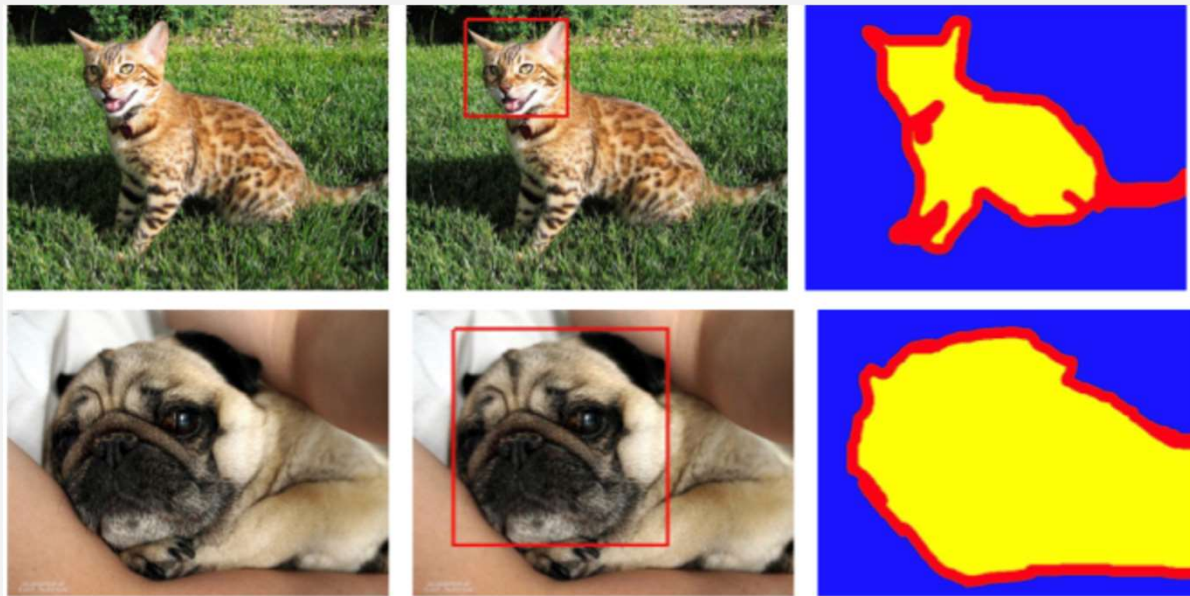






# OXFORD-IIIT PETS DATASET

37 species with roughly 200 images for each class



# OXFORD-IIIT PETS DATASET

Oxford VS Stanford

Why use Oxford Dataset in object detection:

I. Complexity of using TensorFlow-gpu in CSE server without root/sudo access:  
TensorFlow/CuDNN outdated -> Anaconda exceed disk quota -> Miniconda can't find \$PATH. Without multiple GPU, Stanford dataset would be painful.

Basically there'll be a lot of problems once TensorFlow/cuda/cudnn is outdated. Finally I moved all my data to Google Cloud.

# OXFORD-IIIT PETS DATASET

Oxford VS Stanford

Why use Oxford Dataset:

2. Limitation of computation power.

37 species with a simple MobileNet-v1 structure took me a whole day to get the result on Google Cloud. It would be too time-consuming to use Stanford dataset.

# MODEL

Which model to choose?

Inception (V3)

MobileNet

Faster-RCNN

Mask-RCNN

# MODEL

R-CNN(CVPR 2014) -> Fast R-CNN(2015)

-> Faster R-CNN(2016) -> Mask-RCNN(ICCV 2017)

# MODEL

## R-CNN

### R-CNN: *Regions with CNN features*

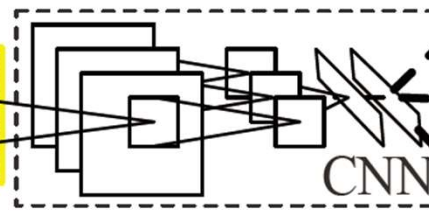


1. Input image



2. Extract region proposals (~2k)

warped region



3. Compute CNN features

aeroplane? no.  
:  
person? yes.  
:  
tvmonitor? no.

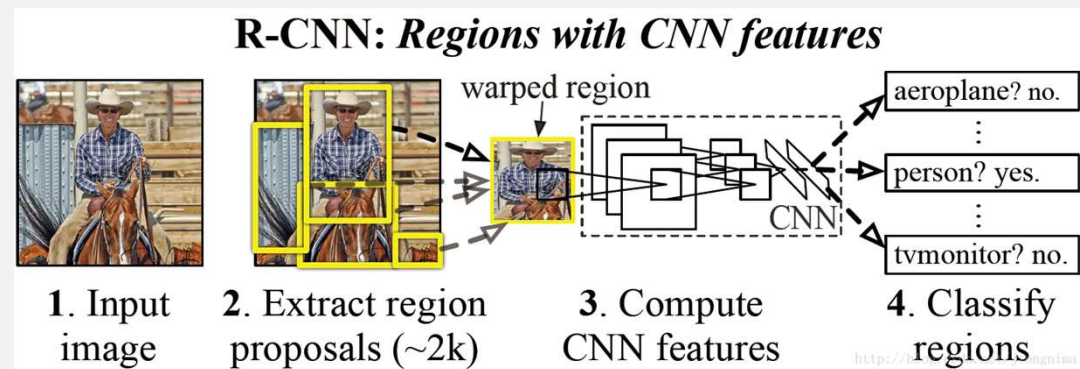
4. Classify regions

<http://blog.robots.ox.ac.uk/pongnima>

# MODEL

## R-CNN

1. Extract 2k regions
2. Each of them go through CNN one by one to extract features.
3. Use SVM to classify regions.
4. Adjust the region through bounding box regression



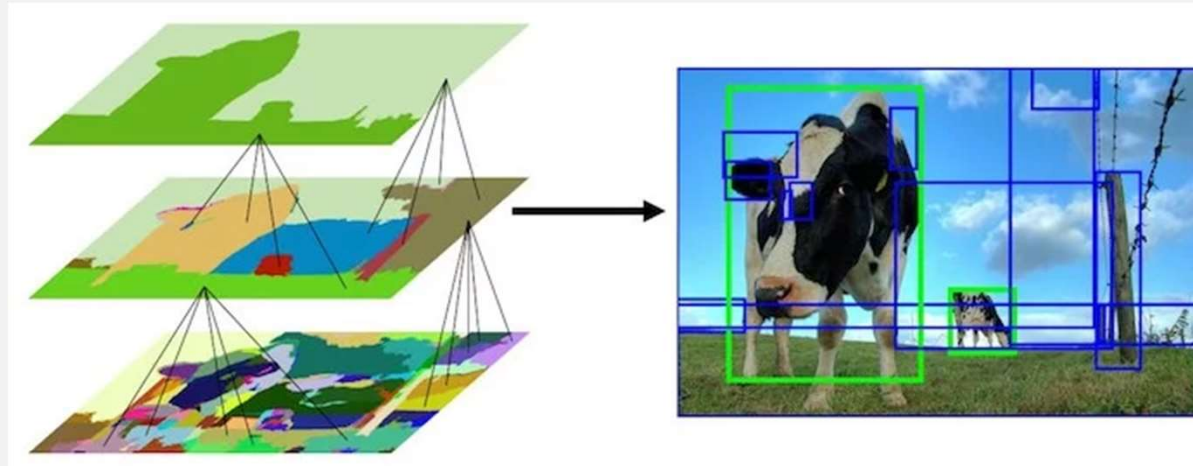


# MODEL

## R-CNN

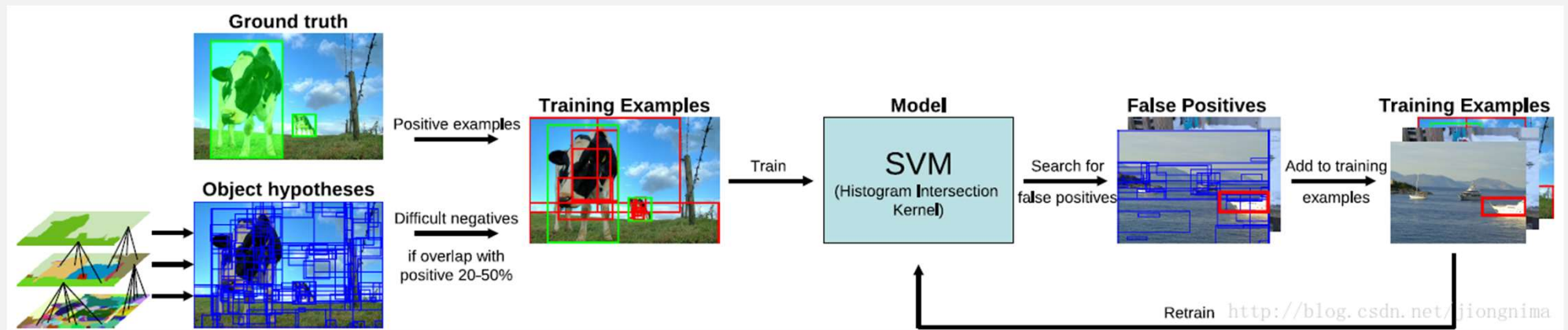
How to extract regions:  
**Selective Search**

- a) Based on traditional methods to segment images
- b) Combine segments based on similarity and then go back to a).
- c) Keep doing this and we will have the result at the right.



# MODEL

## R-CNN



# MODEL

## R-CNN

Pros: Use CNN to extract features.

Use bounding box regression to adjust final result.

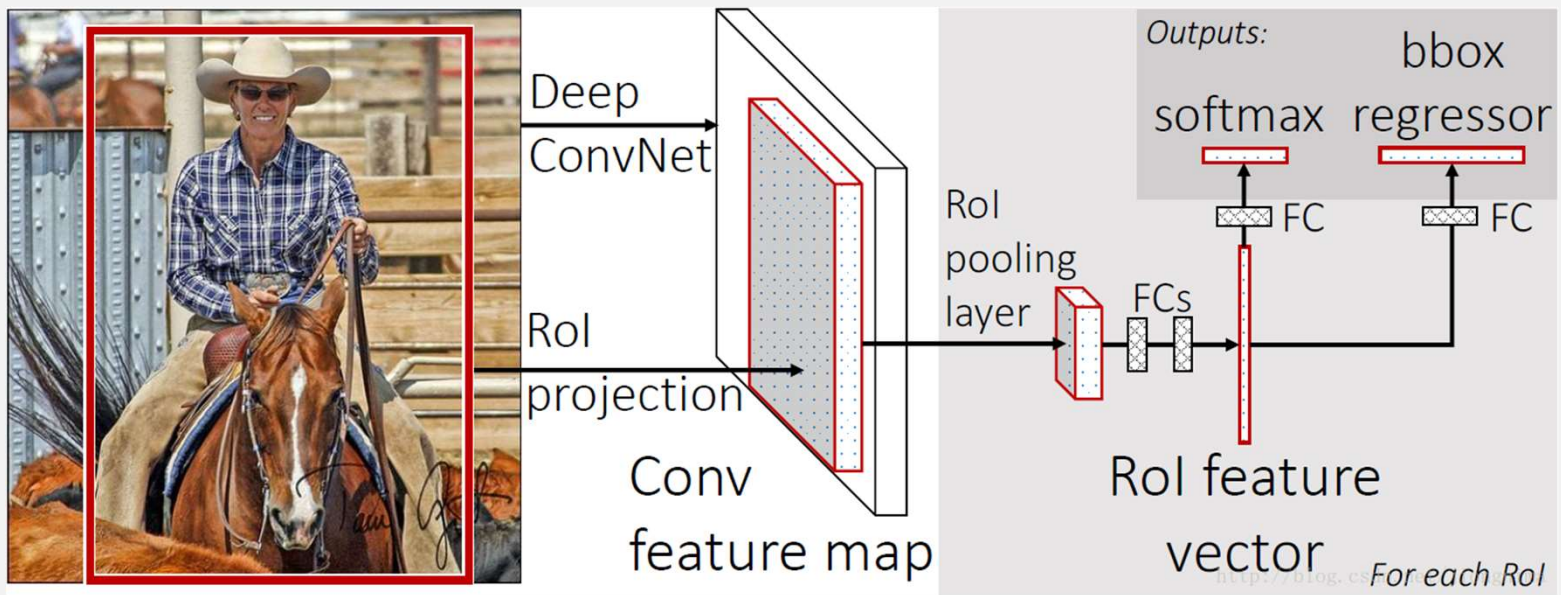
Cons: Selective Search is time-consuming.

(Series) CNN forward propagation is time-consuming.

Each parts trained separately, waste of time & space.

# MODEL

## Fast R-CNN

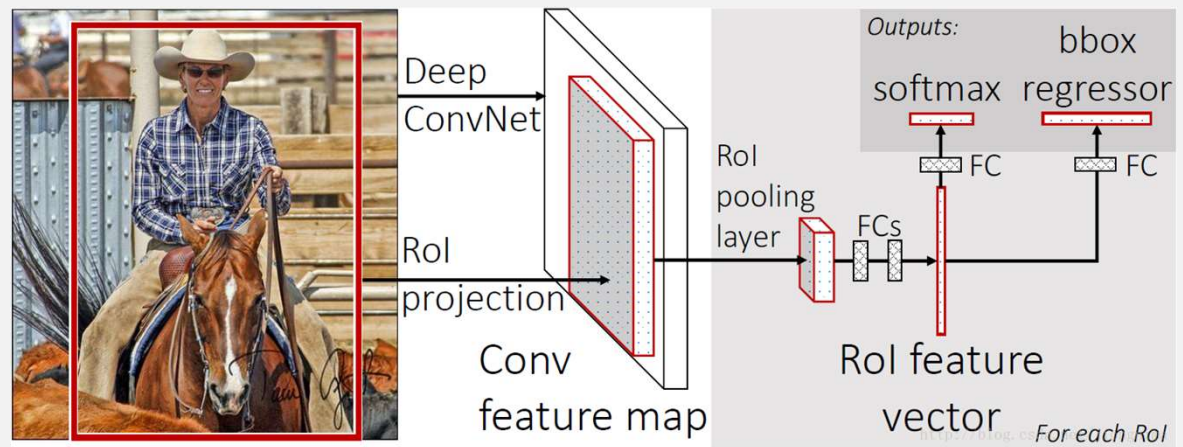


# MODEL

## Fast R-CNN

Still use selective search, and a neural network to extract features on the whole graph.

After that, an RoI Pooling Layer will be used to extract features from feature map and pass to FC Layer for correction.



# MODEL

## Fast R-CNN

Pros: Use a NN to extract features based on the whole image instead of doing it one by one.

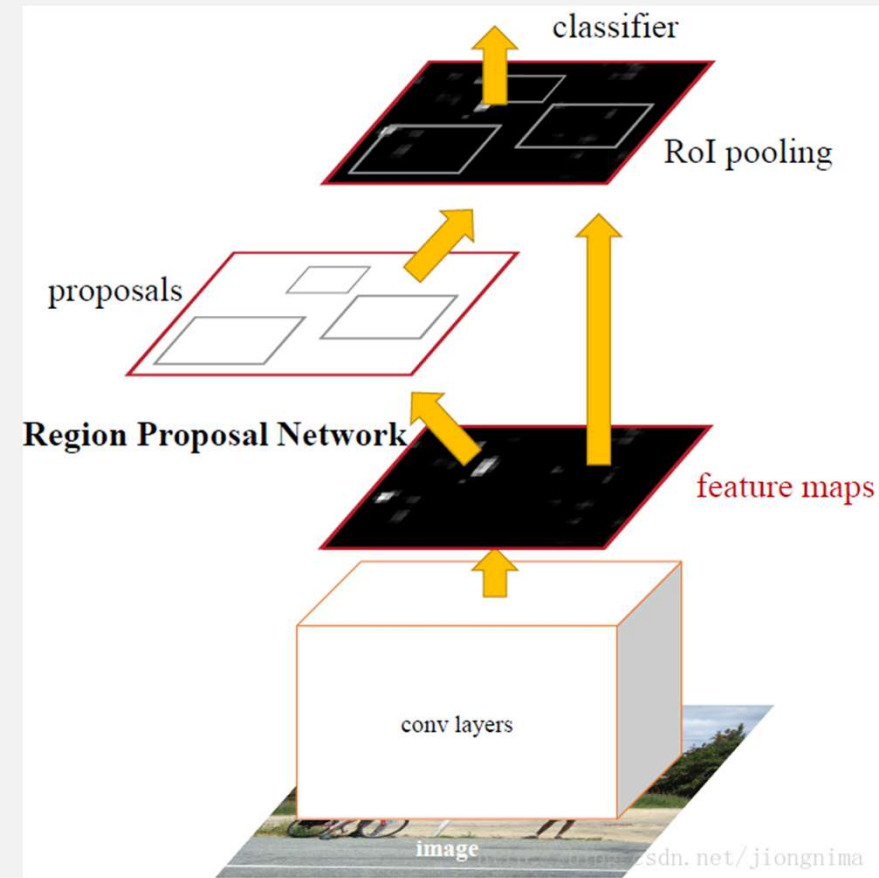
The other parts could be combined during training except for selective search.

Cons: Selective Search is still there.

# MODEL

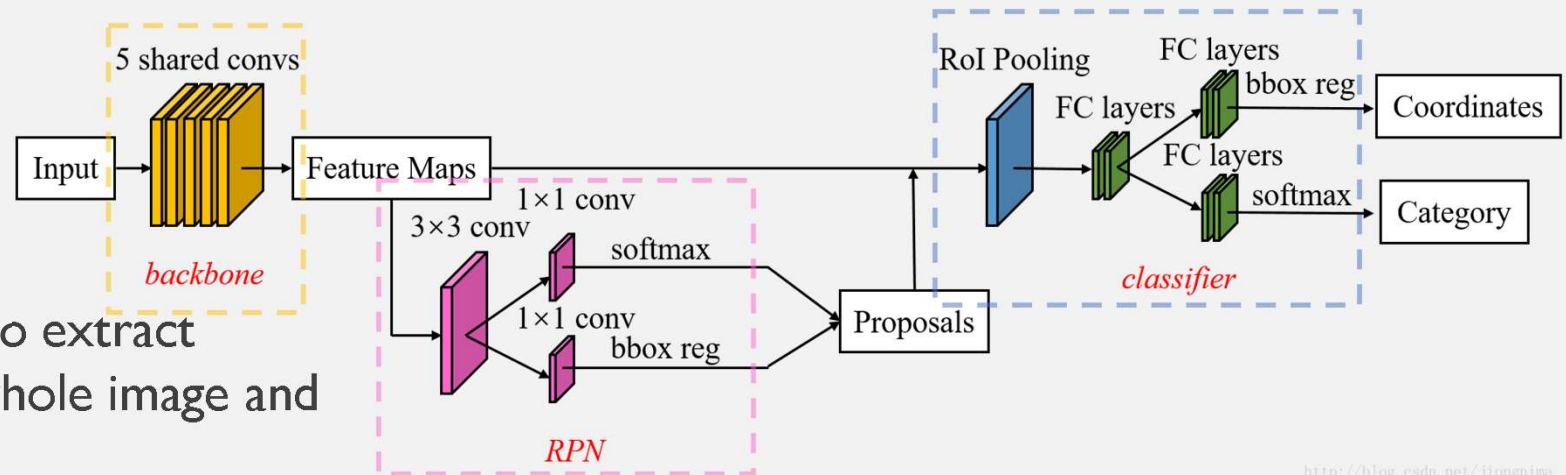
## Faster R-CNN

Use Regional Proposal Network(RPN) to replace selective search.



# MODEL

## Faster R-CNN



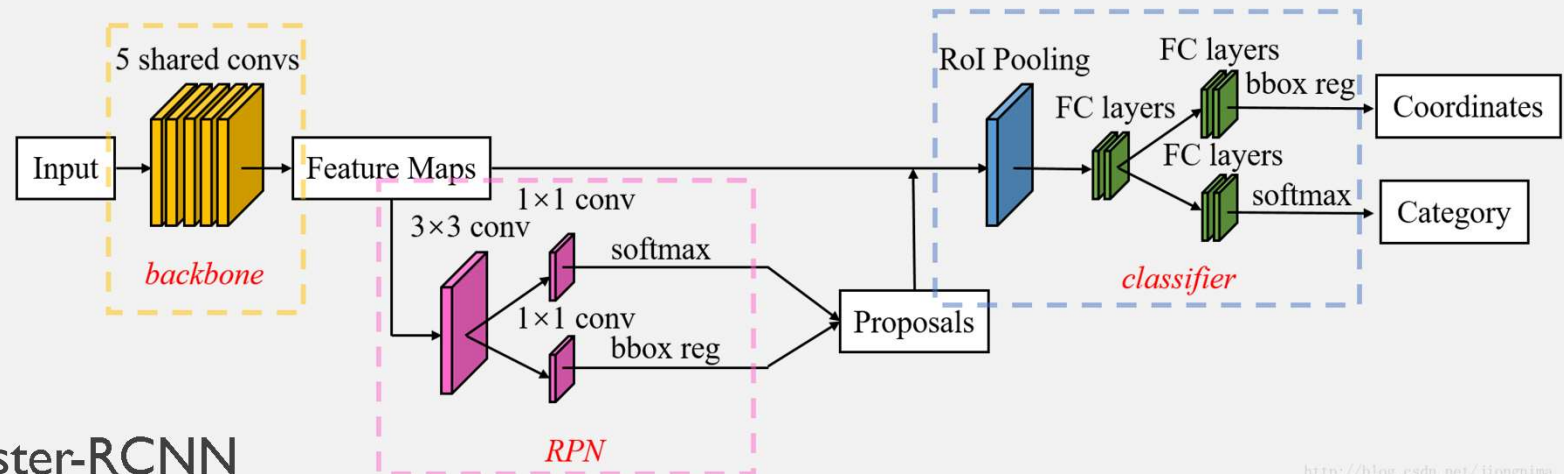
<http://blog.csdn.net/jiongningma>

1. Use backbone to extract features for the whole image and pass to RPN.
2. RPN will generate bounding box of ROI and slightly fix it.
3. ROI pooling layer will select feature for each ROI and FC layer will do classification.



# MODEL

## Faster R-CNN



As we can see, Faster-RCNN mainly consists of 3 parts: backbone, RPN and classifier. Let's look at RPN first.

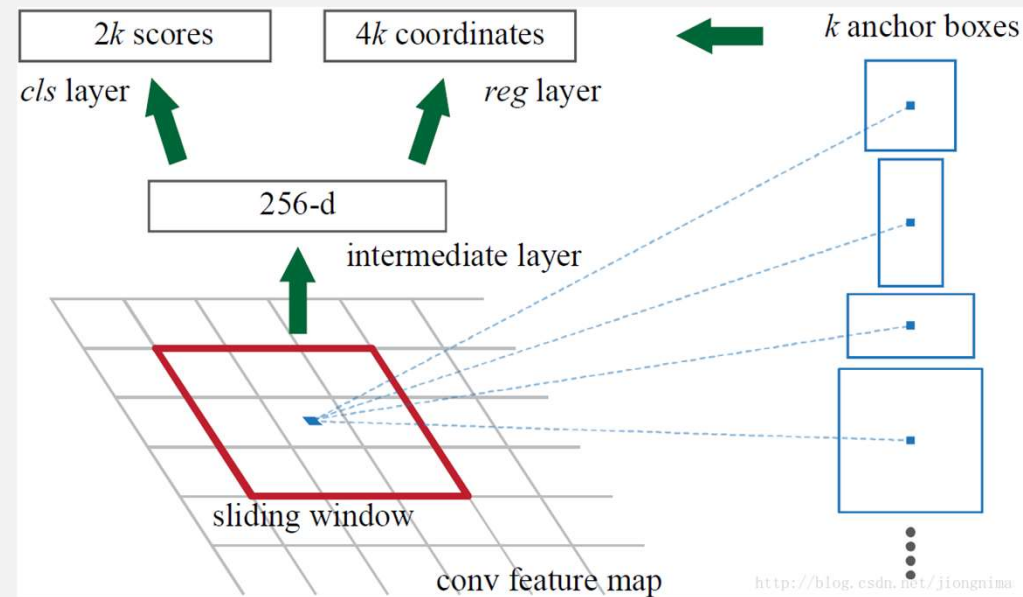
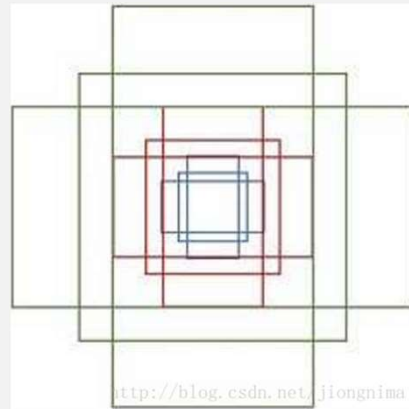
<http://blog.csdn.net/jiongningma>

# MODEL

## Faster R-CNN

RPN relies on the sliding window to generate 9 pre-defined anchors for each position.

The 9 pre-defined anchors can be shown as:  
(3 sizes, 128\*128, 256\*256, 512\*512.  
3 width-height ratios:  
1:1, 1:2, 2:1)



# MODEL

Faster R-CNN

RPN

There will be  $40 \times 60$  shared feature map, and thus  $40 \times 60 \times 9$  (~20000) anchors.

For each anchor, RPN needs to decide whether it is front or back(cover the object?) and adjust it if it's front.

For the 1<sup>st</sup> question, RPN use SoftmaxLoss to train.

For the 2<sup>nd</sup> question, RPN use SmoothL1 Loss to train.

# MODEL

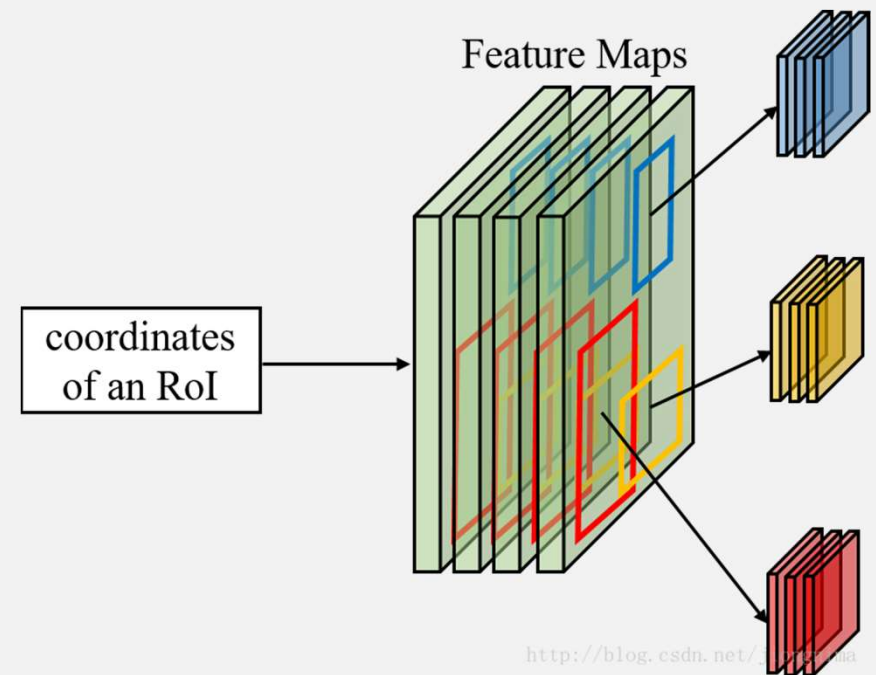
## Faster R-CNN

### ROI Polling

For each ROI, we need to get it from the combined feature map and send it to classifier.

ROI Polling will select feature for each ROI and convert the dimension to meet the FC layer requirement.

E.g.: ROI polling in the right will pick out the feature for each ROI and convert it to  $6 \times 6$ .



# MODEL

Faster R-CNN

Classifier and boulder adjust

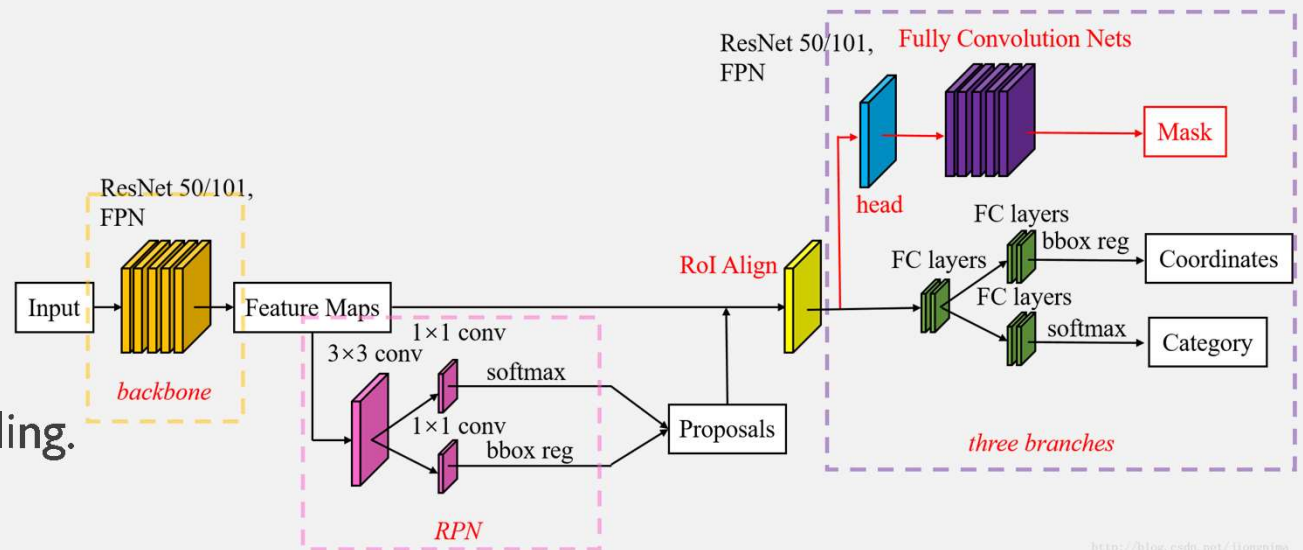
Classifier: see what exactly is this ROI (Human, car, flower)

boulder adjust: also use SmoothL1 Loss to adjust non-background ROI boundaries.

# MODEL

## Mask R-CNN

1. ROI polling -> ROI Align:  
Better align with the original  
ROI region compared to polling.



2. Mask branch

FCN SoftmaxLoss -> average binary cross-entropy loss of K Mask predicts

<http://blog.csdn.net/jiangniwo>

# MODEL

## Mask R-CNN

Great result with great cost.

	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
MNC [7]	ResNet-101-C4	24.6	44.3	24.8	4.7	25.9	43.6
FCIS [20] +OHEM	ResNet-101-C5-dilated	29.2	49.5	-	7.1	31.3	50.0
FCIS+++ [20] +OHEM	ResNet-101-C5-dilated	33.6	54.5	-	-	-	-
<b>Mask R-CNN</b>	ResNet-101-C4	33.1	54.9	34.8	12.1	35.6	51.1
<b>Mask R-CNN</b>	ResNet-101-FPN	35.7	58.0	37.8	15.5	38.1	52.4
<b>Mask R-CNN</b>	ResNeXt-101-FPN	<b>37.1</b>	<b>60.0</b>	<b>39.4</b>	<b>16.9</b>	<b>39.9</b>	<b>53.5</b>

# RESULT

Training Accuracy &  
Model Size

80% Training  
10% Validation  
10% Testing

<b>Inception v3</b>	<b>Accuracy: 0.98 Size: 88M</b>
<b>MobileNet 100, 224</b>	<b>Accuracy: 0.94 Size: 10M</b>
<b>MobileNet 050, 224</b>	<b>Accuracy: 0.92 Size: 3M</b>
<b>MobileNet 050, 128</b>	<b>Accuracy: 0.91 Size: 3M</b>
<b>MobileNet 035, 224</b>	<b>Accuracy: 0.94 Size: 2M</b>



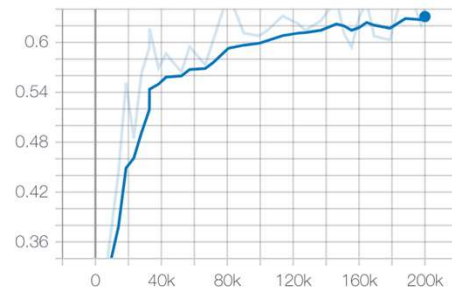
# RESULT

DetectionBoxes\_Precision

6

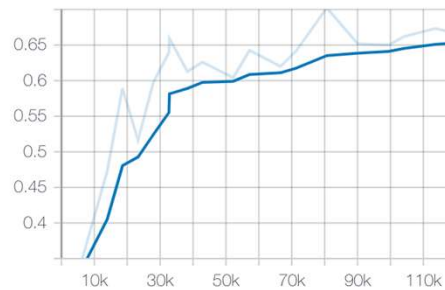
mAP

tag: DetectionBoxes\_Precision/mAP



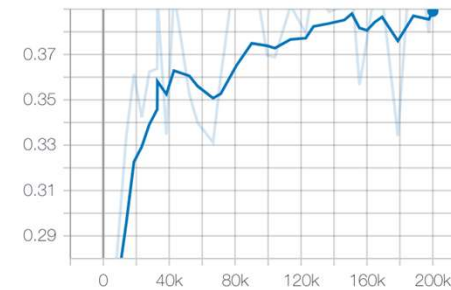
mAP (large)

tag: DetectionBoxes\_Precision/mAP (large)



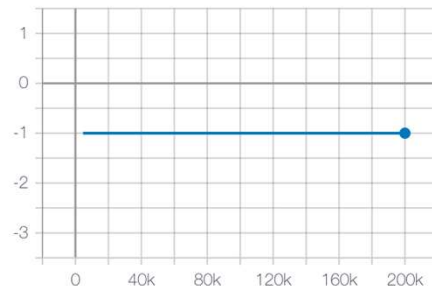
mAP (medium)

tag: DetectionBoxes\_Precision/mAP (medium)



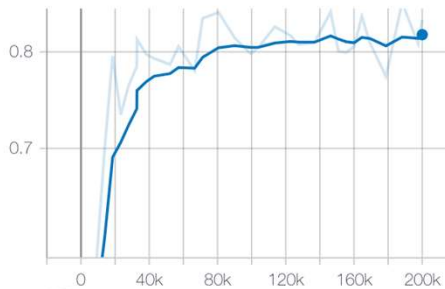
mAP (small)

tag: DetectionBoxes\_Precision/mAP (small)



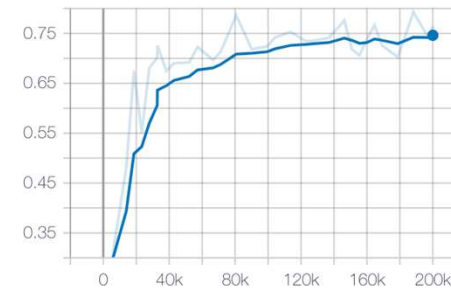
mAP@.50IOU

tag: DetectionBoxes\_Precision/mAP@.50IOU



mAP@.75IOU

tag: DetectionBoxes\_Precision/mAP@.75IOU



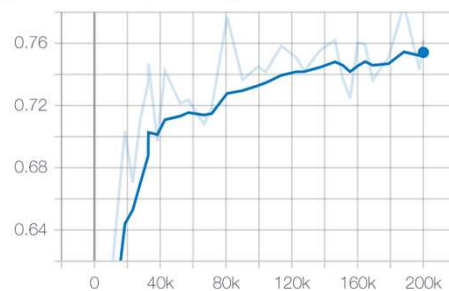
# RESULT

DetectionBoxes\_Recall

6

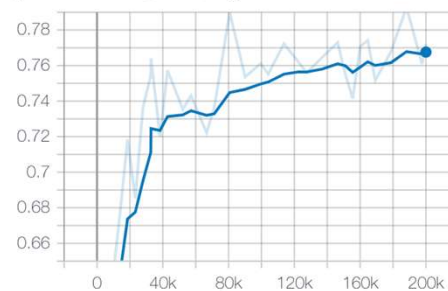
AR@1

tag: DetectionBoxes\_Recall/AR@1



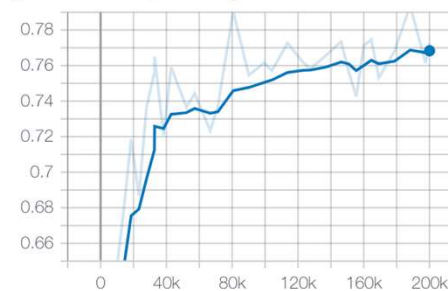
AR@10

tag: DetectionBoxes\_Recall/AR@10



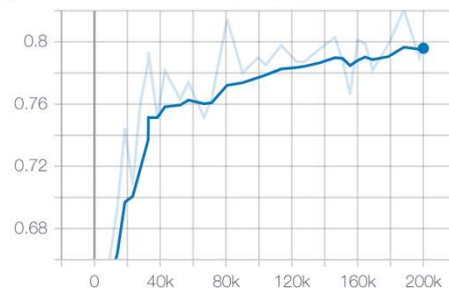
AR@100

tag: DetectionBoxes\_Recall/AR@100



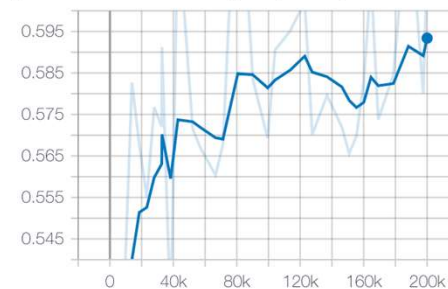
AR@100 (large)

tag: DetectionBoxes\_Recall/AR@100 (large)



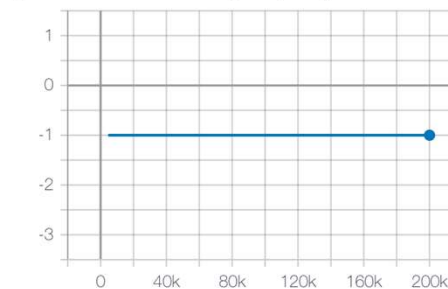
AR@100 (medium)

tag: DetectionBoxes\_Recall/AR@100 (medium)



AR@100 (small)

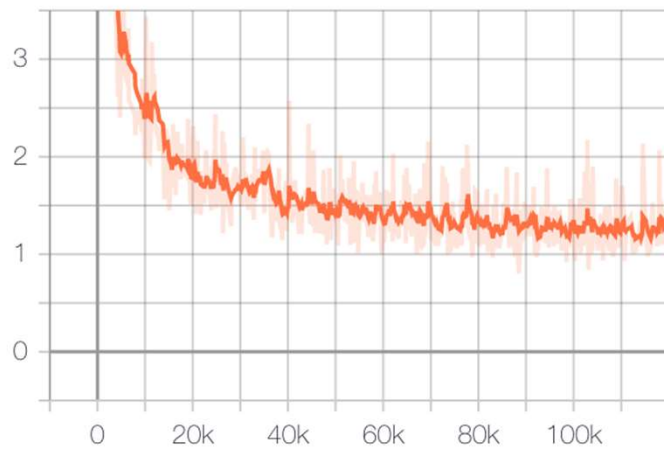
tag: DetectionBoxes\_Recall/AR@100 (small)



# RESULT

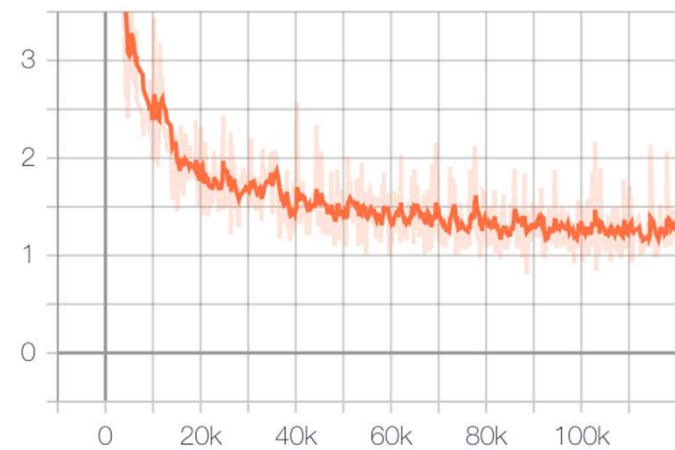
loss\_1

loss\_1



loss\_2

loss\_2



## DEMO - CLASSIFICATION

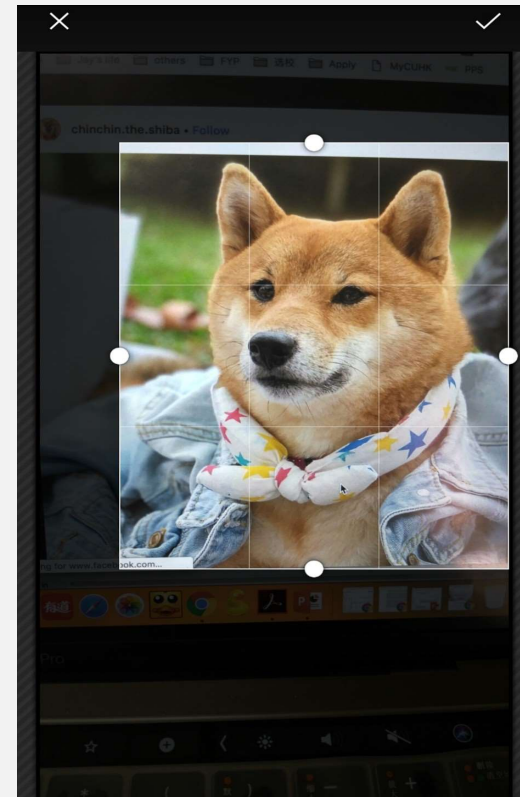
### Application

1. Take/Choose photo with Inception model
2. Real-time Classify with Inception model

## DEMO - CLASSIFICATION

Take/Choose photo  
with Inception model

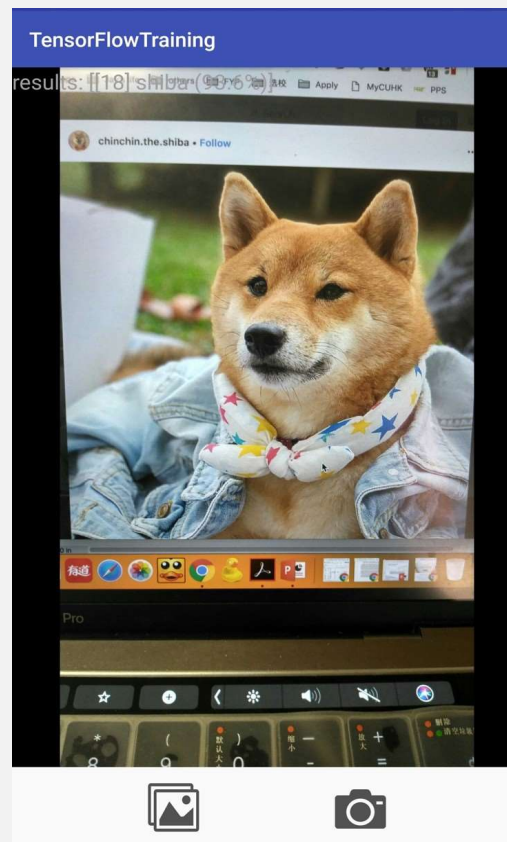
I. Take & Crop  
photo in APP





# DEMO - CLASSIFICATION

Take/Choose photo  
with Inception model

1. Take a photo in  
APP
2. Choose a photo  
from album:  
Crop & Not Crop



	SIBERIAN HUSKY	ALASKAN MALAMUTE
Picture		
Origin	Siberia	Alaska
Size	51 - 60 cm	58 - 71 cm
Weight	16 - 34 kg	39 - 57 kg
Function	To carry a light load at moderate speed over great distances	To carry a heavy load
Eyes	Blue or Brown	Only Brown
Ears	Set high on the head	Set wide apart on the head
Tail	Fox brush carried in a sickle	A waving plume
Personality Traits	Highly Active & Vocal	Laid Back
	Friendly towards other dogs	Gender aggressive towards dogs of the same sex
	No loyalty to one person - they love everyone & everything	Family orientated - Babysat the Mahlemut children in the tribe



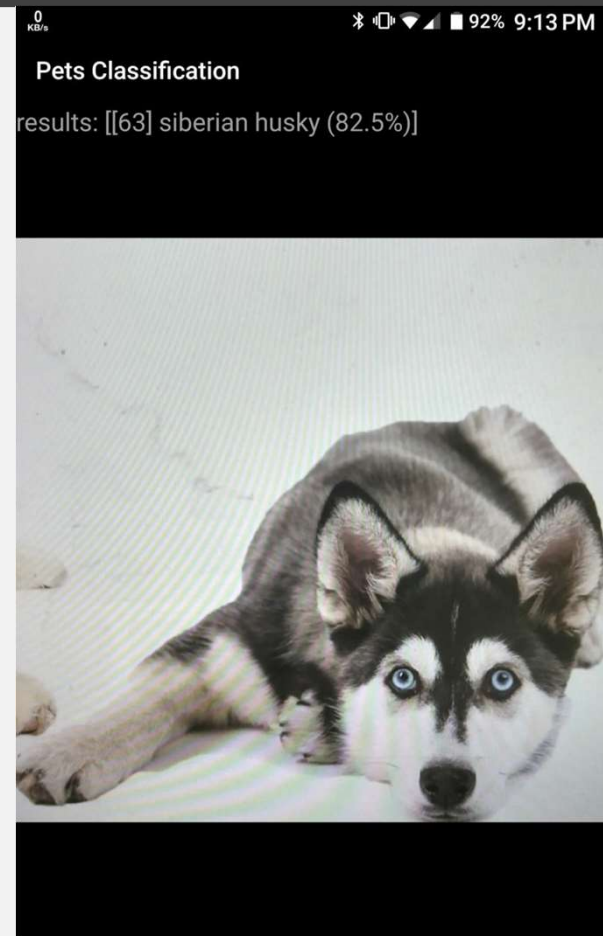
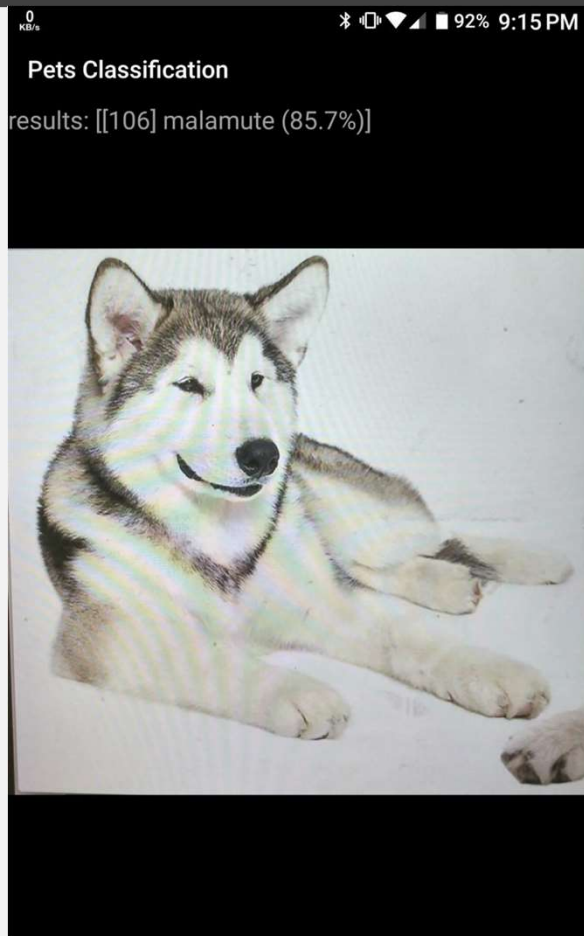
## DEMO - CLASSIFICATION

### 3. Similar Species (Husky vs Malamute, High Quality)





# DEMO - CLASSIFICATION

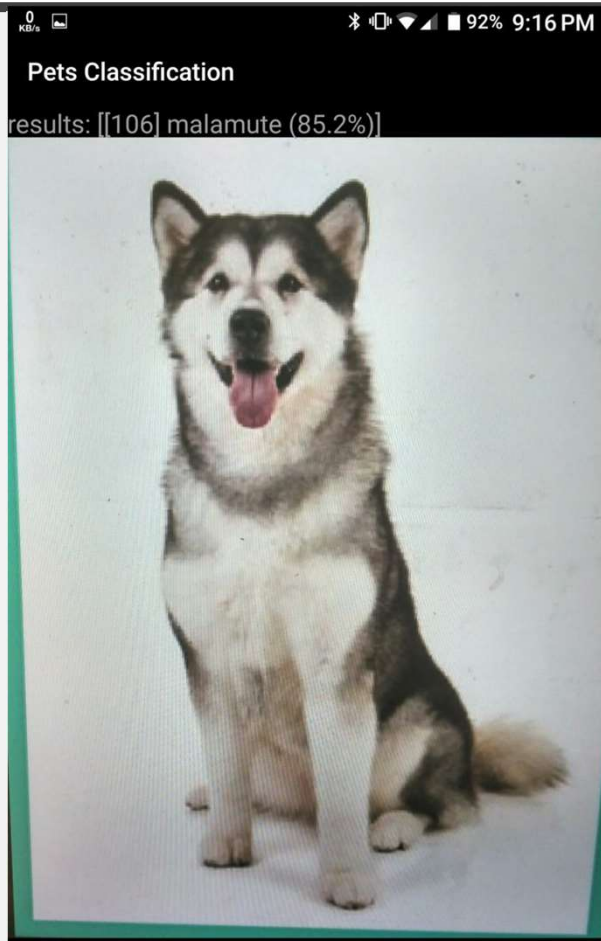


## DEMO - CLASSIFICATION

### 3. Similar Species (Husky vs Malamute, Normal Quality)



## DEMO - CLASSIFICATION

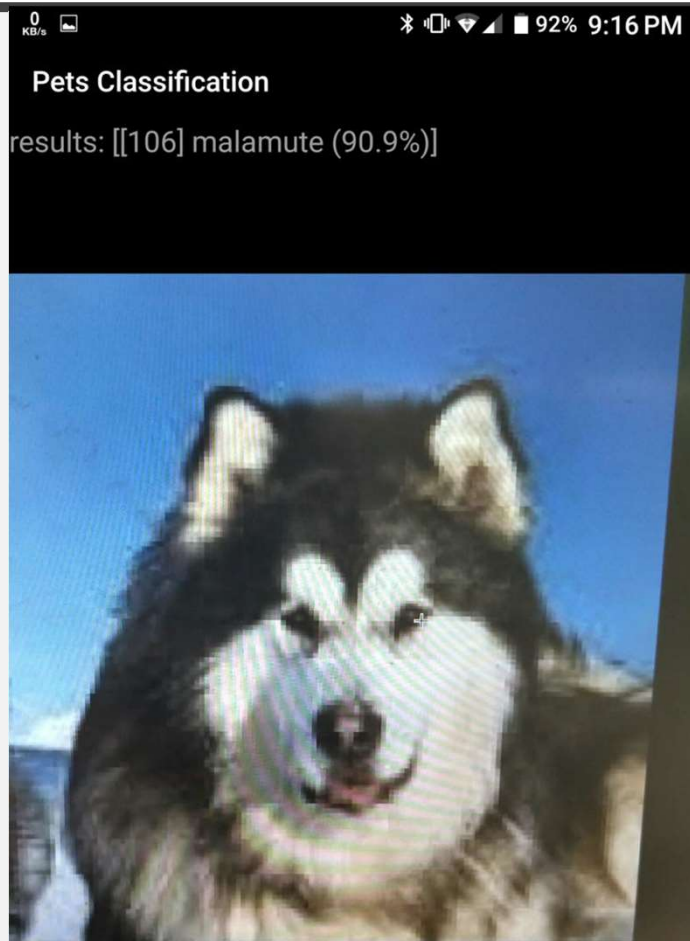


## DEMO - CLASSIFICATION

### 3. Similar Species (Husky vs Malamute, Low Quality)



## DEMO - CLASSIFICATION

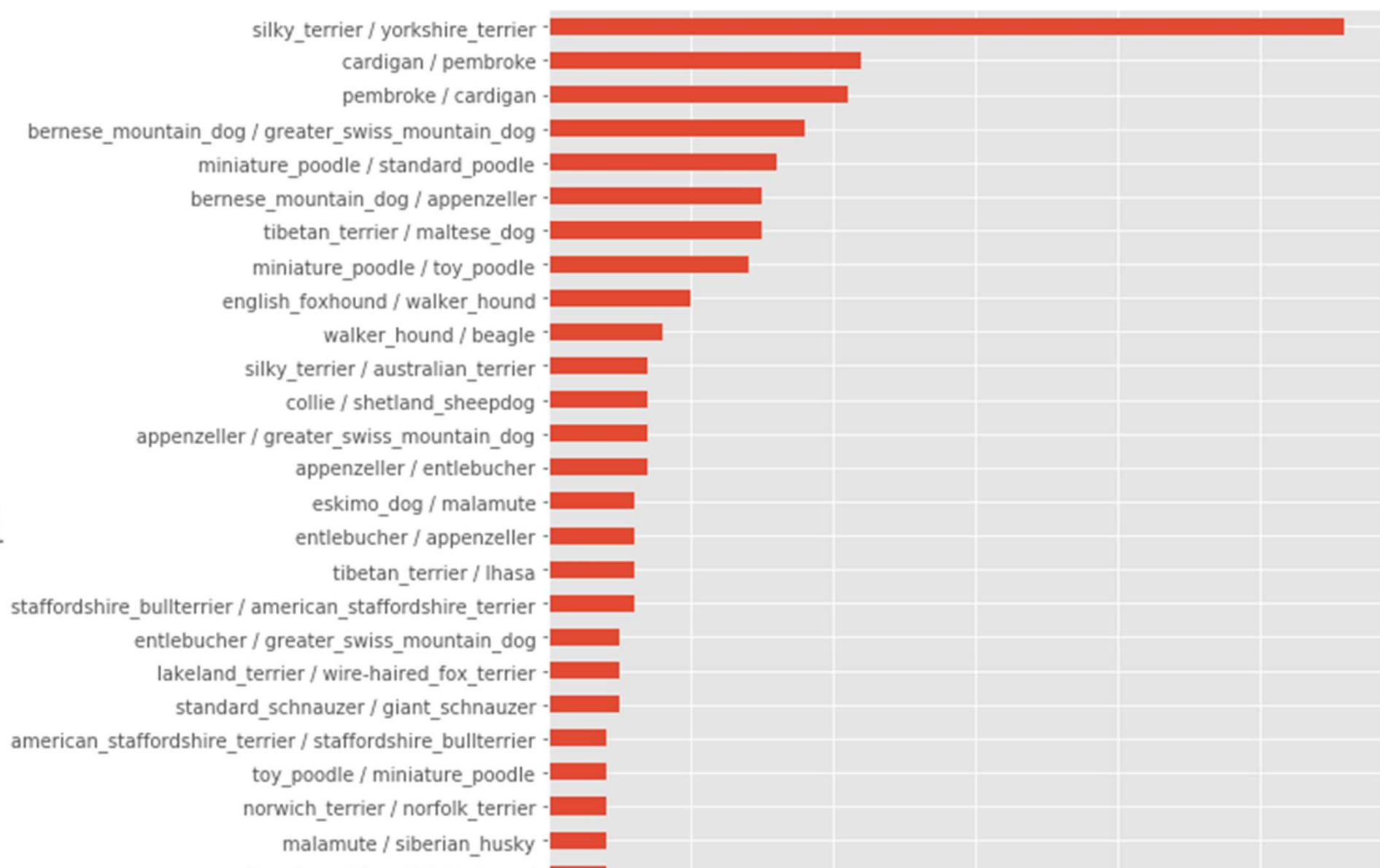




## DEMO - CLASSIFICATION

### 4. Other Similar Species (Cairn vs Affenpinscher vs Schnauzer)





## DEMO – DOG FACE DETECTION

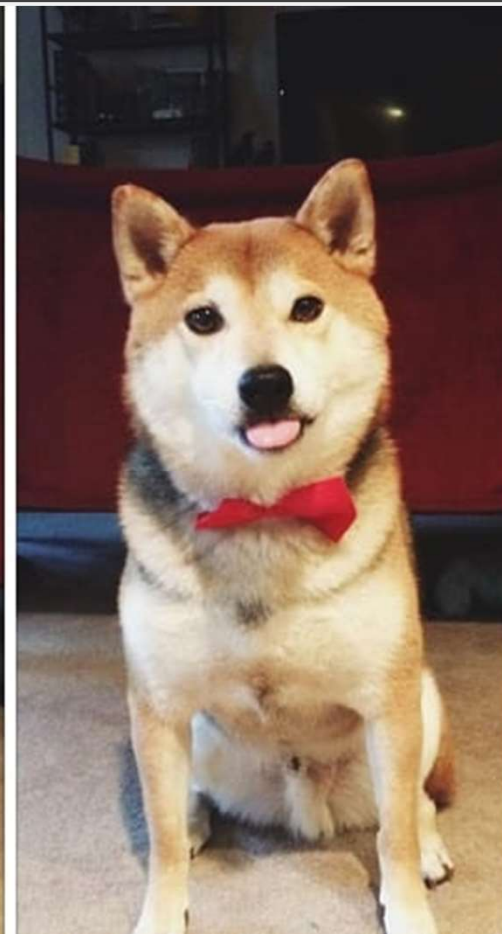
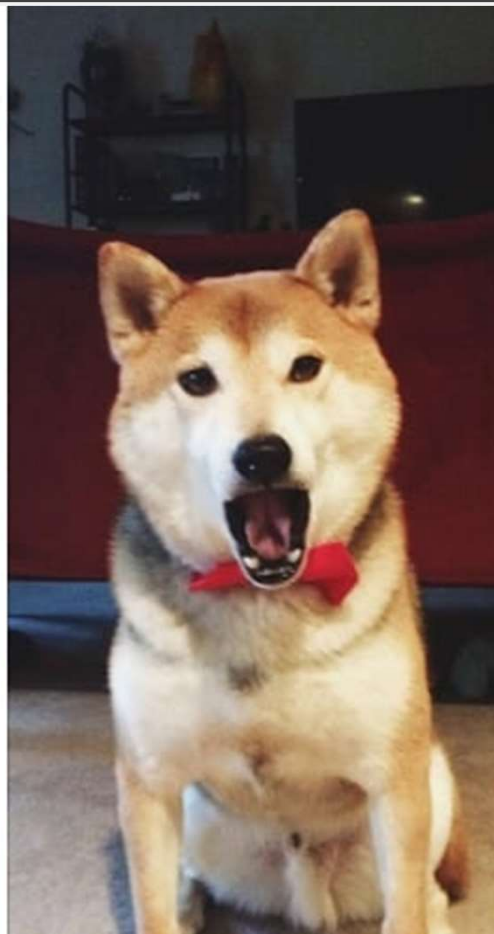




## DEMO – DOG FACE DETECTION



## DEMO – DOG FACE DETECTION





## DEMO – DOG FACE DETECTION

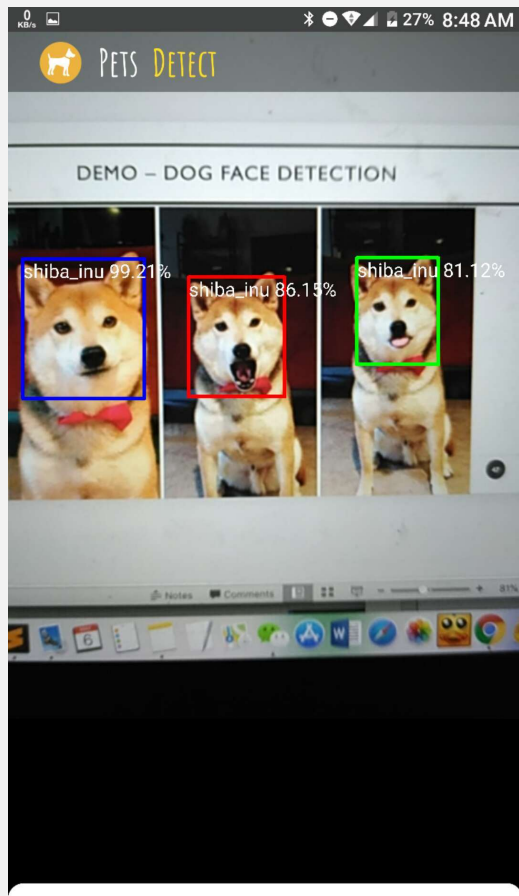


## SMALL UI IMPROVE

german shepherd: 0.88261706



## SMALL UI IMPROVE

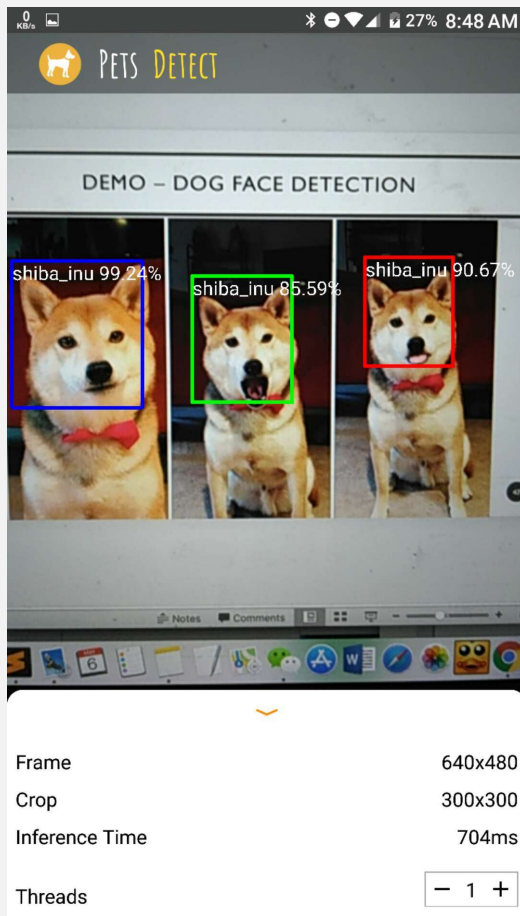


Top bar – transparent

Change icon, logo

Bottom – hidden menu

## SMALL UI IMPROVE



Swipe up -> hidden menu will show

Current frame size

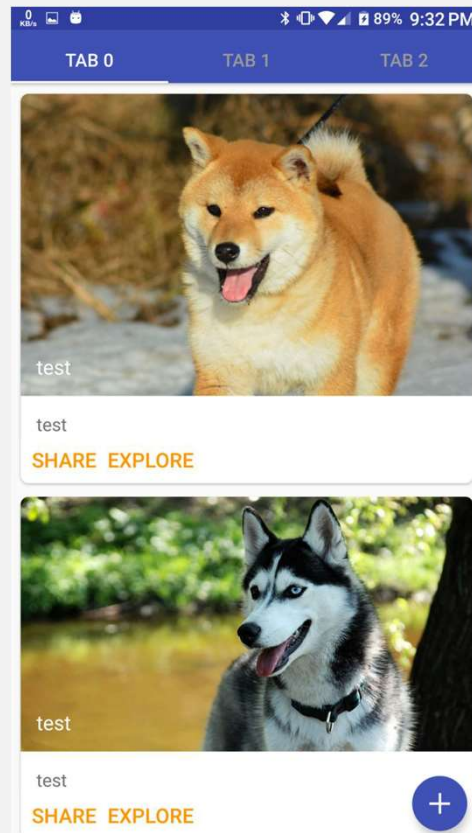
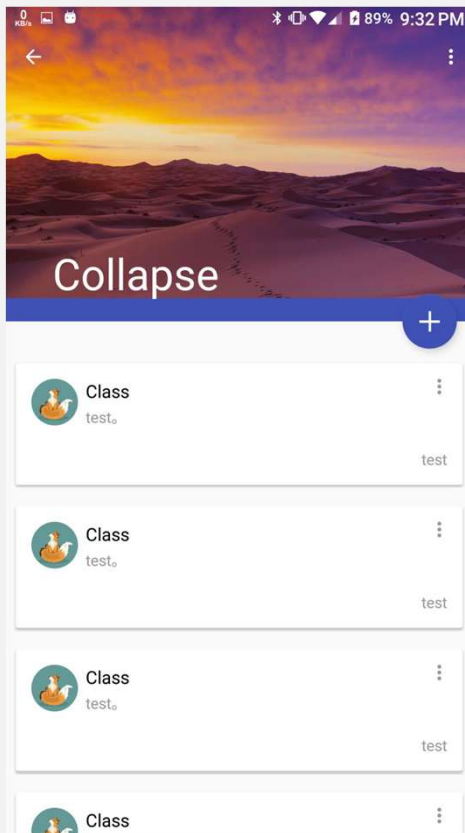
Crop size

Inference Time

Thread Num



# UI EXPLORE



The names & numbers are somehow packed and protected in this version.

Better not to change the library function or nested structure.

# CONCLUSION

## Improvements:

For classification:

More images for each class probably would be better.

UI could be improved.

For detection:

More models should be trained with gpu, especially faster one.

Stanford dataset would be better.



## CONCLUSION

Thank you!