## Two shot detectors

* use two passes on the data:
  - detect candidates (region proposal) (region proposals)
  - classify and regress objects in proposal (including a no-object class)

* Methods:

More complex →
  - Fast-RcNN - cNN classifier
  - R-FcN - Fully convolutional RcNN
  - Faster RcNN - use region proposal network
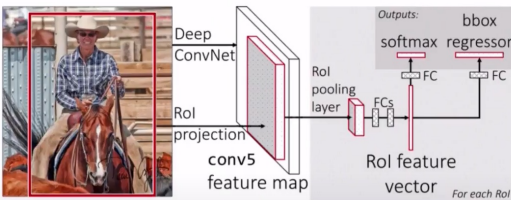  - Mask RcNN - add object segmentation (instance...)

## Two shot detectors

* RcNN:
- Instead of processing candidates on a grid use Super Pixels (Selective Search alg.) to find candidates
- Extract 2000 proposals

- Advantage:
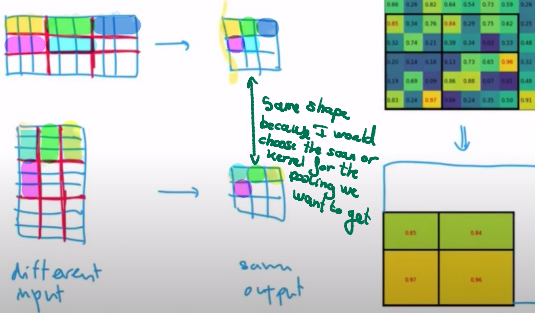  • candidates have different aspect ratios and scales
  • Fewer candidates

* Fast-RcNN architecture:



* RoI-pooling - use max pooling to convert the region of interest to a fixed size (e.g. 7x7) to the FC classifier

This page does not seem important

* ROI pooling:



Save shape
because I would
choose the same or
kernel for the
pooling we
want to get

different
input

same
output

---

* Region proposal Network (RPN)
- Part of faster - RCNN
- train CNN to produce region proposals
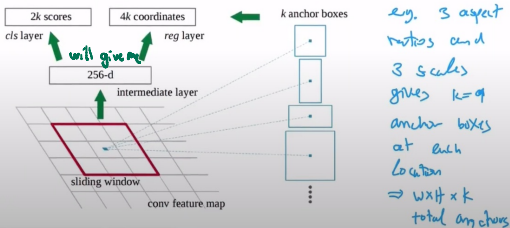- Light weight CNN (relatively simpler task)

* Algorithm:

main idea

- slide window and classify (object / no-object)
  use fast convolution implementation (sliding window)

- use k anchor boxes at each location
  (different size and aspect ratio)
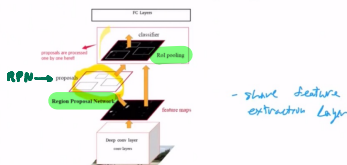  ⟹ around 200k boxes

---

* RPN algorithm contd.
- classify / regress k - boxes
  cls - output 2k classification scores (obj / no-obj)
  reg - output 4k regression boxes



| 2k scores | 4k coordinates |
|---|---|
| cls layer | reg layer |

will give me

256-d
intermediate layer

sliding window

conv feature map

k anchor boxes

e.g. 3 aspect
ratios and
3 scales
gives k = 9
anchor boxes
at each
location
⟹ W × H × k
total anchors

---

* Faster RCNN (fast RCNN + RPN):
- use RPN for region proposals
- use ROI - pooling then classify using dense layers



FC Layer

classifier

proposals are processed
one by one here!!

RoI pooling

RPN → proposals

Region Proposal Network

feature maps

Deep conv layer
conv layers

- share feature
  extraction layer

# Two shot detectors

* **RFCN** (Region based fully convolutional Network)

 - use Convolutions instead of costly dense layers
 that are applied for each region)

 - create position sensitive score maps
 where each score map is sensitive to
 another region.

 k×k grid ⟹ $k^2$ score maps

 each score map has C+1 channels
 (c classes + one no-obj class)

# Two shot detectors

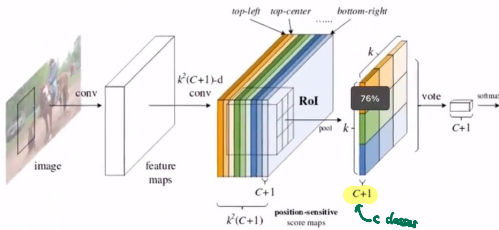* Position sensitive RoI Pooling:



Figure 1: Key idea of **R-FCN** for object detection. In this illustration, there are $k \times k = 3 \times 3$ position-sensitive score maps generated by a fully convolutional network. For each of the $k \times k$ bins in an RoI, pooling is only performed on one of the $k^2$ maps (marked by different colors).

* RFCN architecture:
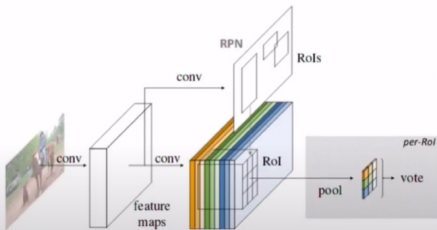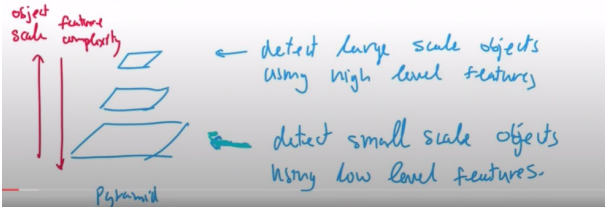
 - use RPN for proposals



Figure 2: Overall architecture of R-FCN. A Region Proposal Network (RPN) [18] proposes candidate RoIs, which are then applied on the score maps. All learnable weight layers are convolutional and are computed on the entire image; the per-RoI computational cost is negligible.

* Feature pyramid network (FPN)
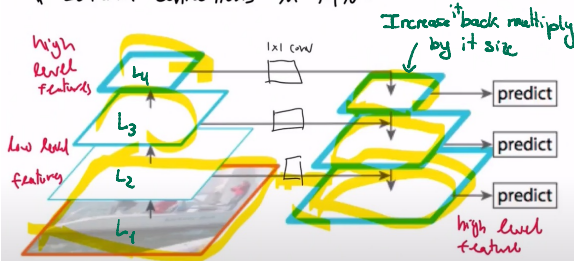
  - used by faster-RCNN

  - Goal: detect objects at different scales using higher level features

object scale | feature complexity

← detect large scale objects using high level features

← detect small scale objects using low level features.

Pyramid

* Lateral connections in FPN:

high level features

$L_4$

$L_3$

low level features

$L_2$

$L_1$

1x1 conv

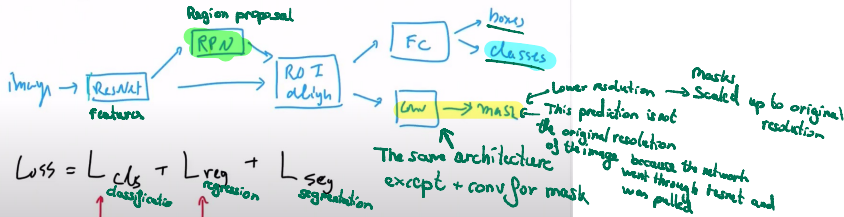"Increase" back multiply by it size

predict

predict

predict

high level feature

- maintain strong semantic features at each level
- lateral connections combine features at

* Mask RCNN: → Try to solve the problem of instant segmentation

- Add segmentation to faster RCNN → instance Segmentation

Region proposal
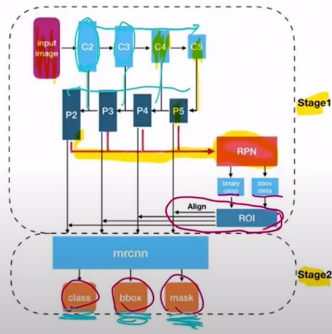RPN

FC → boxes, classes

image → ResNet features → ROI align → conv → MASK

The same architecture except + conv for mask

Lower resolution → This prediction is not the original resolution of the image because the network went through & reset and was pulled

Masks → Scaled up to original resolution

$$Loss = L_{cls} + L_{reg} + L_{seg}$$

classification    regression    segmentation

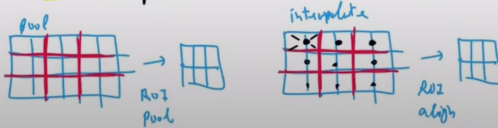$L_{cls}^{RPN} + L_{cls}^{obj}$    $L_{reg}^{RPN} + N_{reg}^{obj}$

* Mask - RcNN uses **FPN** as backbone:



## Two shot detectors

* ROI align: ← Try to solve the problem of **max** pulling → You're sensitive to alignment

- RoI - pooling extracts small feature maps of fixed size.

- RoI - pooling is sensitive to alignment ← Because regions are fixed. shifting the grid will change results.

- RoI - align solves this by sampling with interpolation



```python
import tensorflow as tf
import tensorflow_datasets as tfds
from tensorflow.keras.applications import MobileNetV2
from tensorflow.keras.layers import Dense, Flatten, Conv2D, Reshape, Input
from tensorflow.keras.models import Model
import matplotlib.pyplot as plt

# Load the Oxford-IIIT Pet dataset        breeds of dogs
dataset, info = tfds.load('oxford_iiit_pet', split='train', with_info=True,
as_supervised=False)

# Preprocessing function
def preprocess_data(sample):
    image = sample['image']
    bbox = sample['objects']['bbox']           Three things that I need
    label = sample['objects']['label']

    # Resize image to a fixed shape (224x224)
    image = tf.image.resize(image, (224, 224))
    image = image / 255.0  # Normalize pixel values

    # Convert bbox to [xmin, ymin, xmax, ymax] format and normalize coordinates
    bbox = tf.stack([bbox[:, 1], bbox[:, 0], bbox[:, 3], bbox[:, 2]], axis=1)  #
Convert to [xmin, ymin, xmax, ymax]

    return image, {'bbox': bbox, 'class': tf.one_hot(label, depth=37)}  # 37 classes
```