# Detection and Segmentation
## CS60010: Deep Learning

Abir Das

IIT Kharagpur
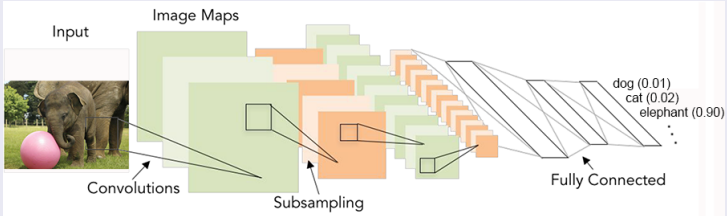
Feb 28, 2020

Introduction
oooooooooooo

Datasets
oooo

Localization
ooooooooooooooooooooooooo

## Agenda

To get introduced to two important tasks of computer vision - detection and segmentation along with deep neural network's application in these areas in recent years.

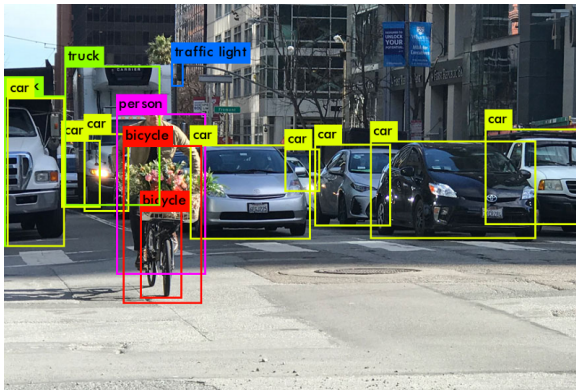# From Classification to Detection

## Classification



## Detection

Introduction

Datasets

Localization

○●○○○○○○○○○○

○○○○

○○○○○○○○○○○○○○○○○○○○○○○○

# Challenges of Object Detection

§ Simultaneous recognition and localization

§ Images may contain objects from more than one class and multiple instances of the same class
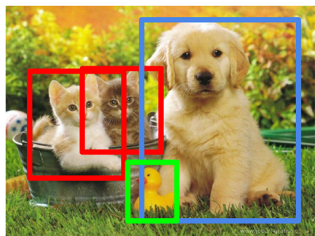
§ Evaluation

Introduction
0000000000000
Datasets
0000
Localization
000000000000000000000000

# Localization and Detection



**Classification**

**Classification + Localization**

**Object Detection**

CAT

CAT

CAT, DOG, DUCK

Single object

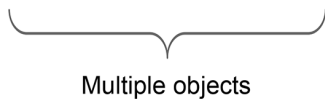Multiple objects

Introduction
○○○○●○○○○○○○

Datasets
○○○○

Localization
○○○○○○○○○○○○○○○○○○○○○○○○
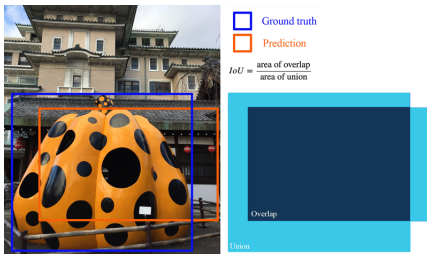
# Evaluation

§ At test time 3 things are predicted:- Bounding box coordinates, Bounding box class label, Confidence score

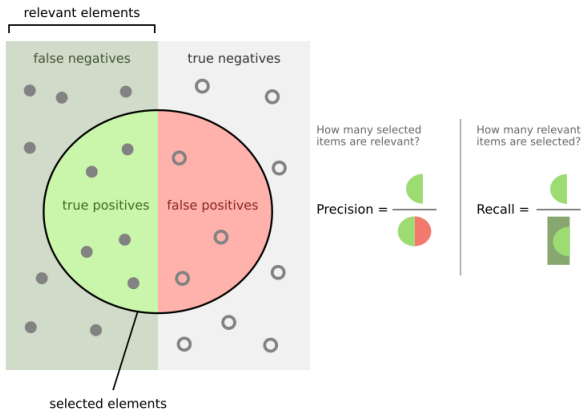§ Performance is measured in terms of IoU (Intersection over Union)



§ According to PASCAL criterion,

▶ a detection is correct if IoU $> 0.5$

▶ For multiple detections only one is considered **true positive**

by the (decreasing) confidence output. Multiple detections of the same object in an image were considered false detections e.g. 5 detections of a single object counted as 1 correct detection and 4 false detections—it was the responsibility of the participant's system to filter multiple detections from its output.

Image Source

Introduction
○○○○○●○○○○○○

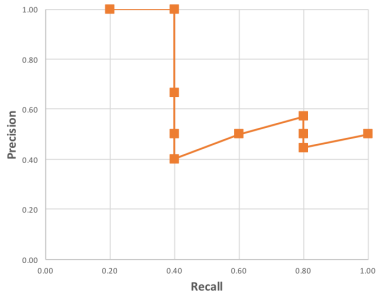Datasets
○○○○

Localization
○○○○○○○○○○○○○○○○○○○○○○○○

# Evaluation: Precision-Recall



§ precision $= \frac{tp}{tp+fp}$

§ recall $= \frac{tp}{tp+fn}$

Image Source

# Evaluation: Average Precision

Lets consider an image with 5 apples where our detector provides 10 detections.

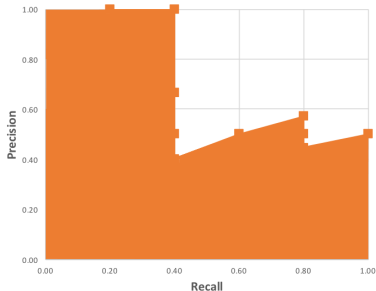| Rank | Correct | Precision | Recall |
|---|---|---|---|
| 1 | True Positive | 1.00 | 0.20 |
| 2 | True Positive | 1.00 | 0.40 |
| 3 | False Positive | 0.67 | 0.40 |
| 4 | False Positive | 0.50 | 0.40 |
| 5 | False Positive | 0.40 | 0.40 |
| 6 | True Positive | 0.50 | 0.60 |
| 7 | True Positive | 0.57 | 0.80 |
| 8 | False Positive | 0.50 | 0.80 |
| 9 | False Positive | 0.44 | 0.80 |
| 10 | True Positive | 0.50 | 1.00 |



Source: This medium post

# Evaluation: Average Precision

Area under curve is a measure of performance. This gives the average precision of the detector.

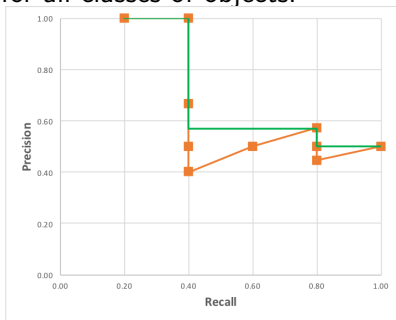| Rank | Correct | Precision | Recall |
|------|---------|-----------|--------|
| 1 | True Positive | 1.00 | 0.20 |
| 2 | True Positive | 1.00 | 0.40 |
| 3 | False Positive | 0.67 | 0.40 |
| 4 | False Positive | 0.50 | 0.40 |
| 5 | False Positive | 0.40 | 0.40 |
| 6 | True Positive | 0.50 | 0.60 |
| 7 | True Positive | 0.57 | 0.80 |
| 8 | False Positive | 0.50 | 0.80 |
| 9 | False Positive | 0.44 | 0.80 |
| 10 | True Positive | 0.50 | 1.00 |



Source: This medium post
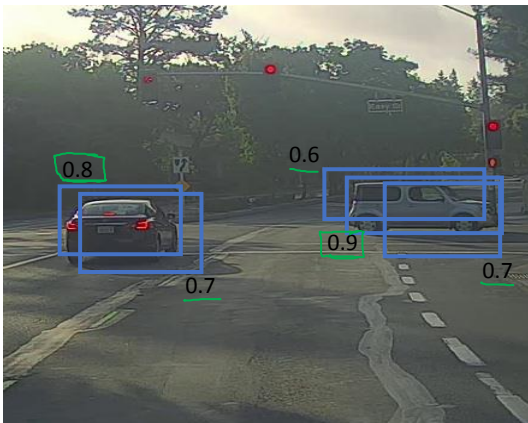
# Evaluation: mean Average Precision

A little more detail:

§ The curve is made smooth from the zigzag pattern by finding the highest precision value at or to the right side of the recall values.

§ Then the average is taken for 11 recall values (0, 0.1, 0.2, ... 1.0) - Average Precison (AP)

§ The mean average precision (mAP) is the mean of the average precisions (AP) for all classes of objects.



Source: This medium post

Introduction
○○○○○○○○○●○○

Datasets
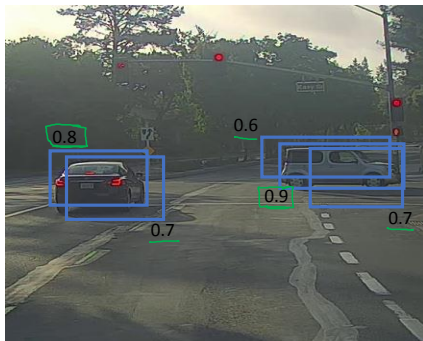○○○○

Localization
○○○○○○○○○○○○○○○○○○○○○○

## Non-max Suppression

What to do if there are multiple detections of the same object? Can you think its effect on precision-recall?
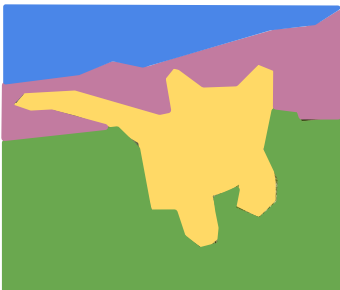


Source: deeplearning.ai

Introduction
⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙●⊙⊙

Datasets
⊙⊙⊙⊙

Localization
⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙⊙

# Non-max Suppression

§ Sort the predictions by the confidence scores

§ Starting with the top score prediction, ignore any other prediction of the same class and high overlap (*e.g.*, IoU > 0.5) with the top ranked prediction

§ Repeat the above step until all predictions are checked



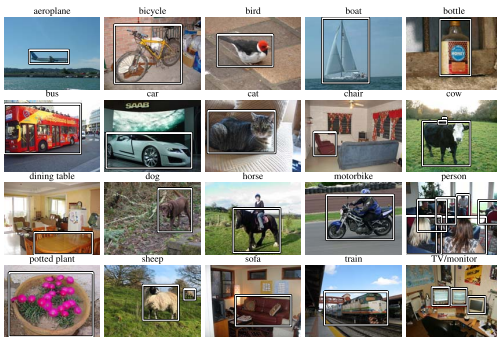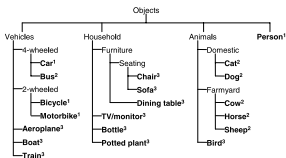Source: deeplearning.ai

## Segmentation

**Semantic Segmentation**

**Instance Segmentation**



**GRASS**, **CAT**, **TREE**, **SKY**



**DOG**, **DOG**, **CAT**

# PASCAL VOC



§ Dataset size (by 2012): 11.5K training/val images, 27K bounding boxes, 7K segmentations

Introduction
○○○○○○○○○○○

Datasets
○●○○

Localization
○○○○○○○○○○○○○○○○○○○○○○○○

## PASCAL VOC



Object detection renaissance (2013-present)

Source: ICCV '15, Fast R-CNN

Introduction
○○○○○○○○○○○○

Datasets
○○○●○

Localization
○○○○○○○○○○○○○○○○○○○○○○○

# COCO Dataset



## What is COCO?

COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features:

- ✔ **Object segmentation**
- ✔ **Recognition in context**
- ✔ **Superpixel stuff segmentation**
- ✔ **330K images (>200K labeled)**
- ✔ **1.5 million object instances**
- ✔ **80 object categories**
- ✔ **91 stuff categories**
- ✔ **5 captions per image**
- ✔ **250,000 people with keypoints**

http://cocodataset.org

Introduction
0000000000000

Datasets
0000

Localization
0000000000000000000000

# COCO Tasks



Image Classification

Semantic Segmentation

Object Detection

Instance Segmentation

Introduction
○○○○○○○○○○○

Datasets
○○○○

Localization
●○○○○○○○○○○○○○○○○○○○○○○○

## Classification + Localization

# Classification + Localization: Task

**Classification**: C classes
    **Input:** Image
    **Output:** Class label
    **Evaluation metric:** Accuracy

 ⟶ CAT

**Localization**:
    **Input:** Image
    **Output**: Box in the image (x, y, w, h)
    **Evaluation metric:** Intersection over Union

 ⟶ (x, y, w, h)

**Classification + Localization**: Do both

Source: cs231n course, Stanford University

Introduction
ooooooooooo

Datasets
oooo

Localization
oooooooooooooooooooooo

## Classification + Localization

# Idea #1: Localization as Regression

**Input**: image



Only one object,
simpler than detection

Neural Net $\longrightarrow$

**Output**:
Box coordinates
(4 numbers)

**Correct output**:
box coordinates
(4 numbers)

**Loss**:
L2 distance

Source: cs231n course, Stanford University

Introduction
oooooooooooo
Datasets
oooo
Localization
ooooooooooooooooooooooo

## Classification + Localization

# Simple Recipe for Classification + Localization

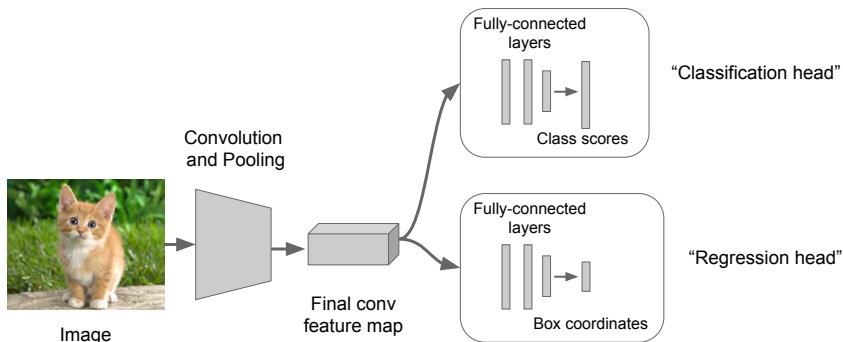**Step 1**: Train (or download) a classification model (AlexNet, VGG, GoogLeNet)



Convolution
and Pooling

Fully-connected
layers

Image

Final conv
feature map

Class scores

Softmax loss

## Classification + Localization

# Simple Recipe for Classification + Localization

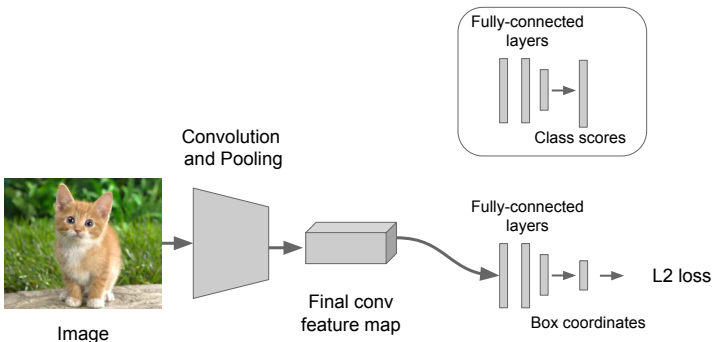**Step 2**: Attach new fully-connected "regression head" to the network



Source: cs231n course, Stanford University

## Classification + Localization

# Simple Recipe for Classification + Localization

**Step 3**: Train the regression head only with SGD and L2 loss
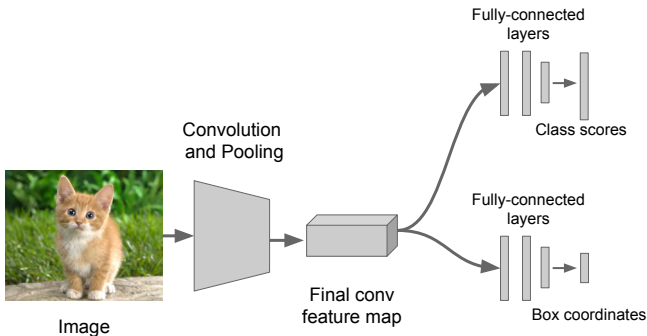


Source: cs231n course, Stanford University

## Classification + Localization

Simple Recipe for Classification + Localization
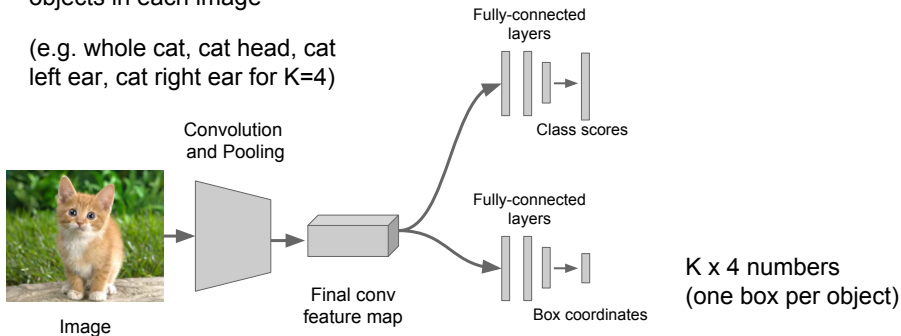
**Step 4**: At test time use both heads

Introduction
OOOOOOOOOOO
Datasets
OOOO
Localization
OOOOOOOOOOOOOOOOOOOOOOO

## Classification + Localization

# Aside: Localizing multiple objects

Want to localize **exactly** K objects in each image

(e.g. whole cat, cat head, cat left ear, cat right ear for K=4)



Image

Convolution and Pooling

Final conv feature map

Fully-connected layers

Class scores

Fully-connected layers

Box coordinates

K x 4 numbers
(one box per object)

Source: cs231n course, Stanford University

Introduction
○○○○○○○○○○○
Datasets
○○○○
Localization
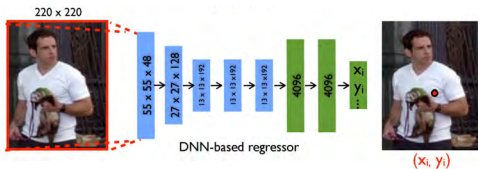○○○○○○○●○○○○○○○○○○○○○

## Classification + Localization

# Aside: Human Pose Estimation

Represent a person by K joints

Regress (x, y) for each joint from last fully-connected layer of AlexNet

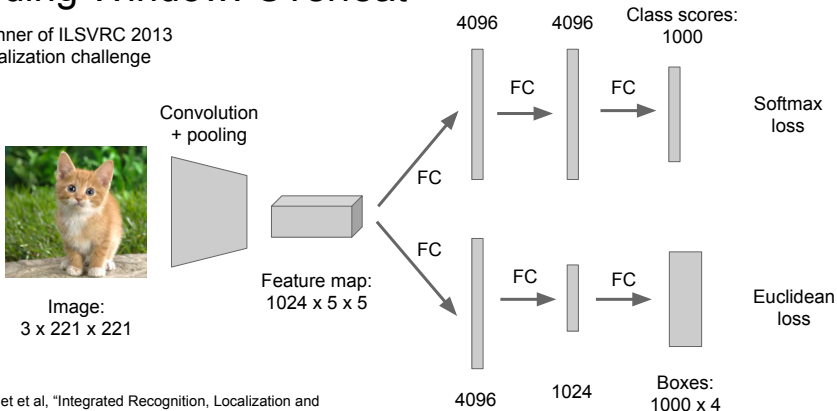(Details: Normalized coordinates, iterative refinement)



220 x 220

55 x 55 x 48 · 27 x 27 x 128 · 13 x 13 x 192 · 13 x 13 x 192 · 13 x 13 x 192 · 4096 · 4096 · $x_i$ $y_i$ :

DNN-based regressor

$(x_i, y_i)$

Toshev and Szegedy, "DeepPose: Human Pose Estimation via Deep Neural Networks", CVPR 2014

Source: cs231n course, Stanford University

## Classification + Localization

# Sliding Window: Overfeat

Winner of ILSVRC 2013
localization challenge



Image:
3 x 221 x 221

Convolution
+ pooling

Feature map:
1024 x 5 x 5

4096    4096    Class scores:
1000

FC      FC

Softmax
loss

FC

FC

4096    1024    Boxes:
1000 x 4

FC      FC

Euclidean
loss

Sermanet et al, "Integrated Recognition, Localization and
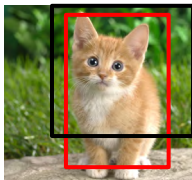Detection using Convolutional Networks", ICLR 2014
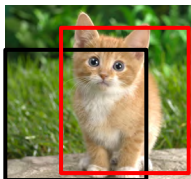
Source: cs231n course, Stanford University

## Classification + Localization

## Sliding Window: Overfeat



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

Source: cs231n course, Stanford University

## Classification + Localization

## Sliding Window: Overfeat



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

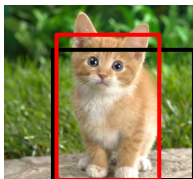| 0.5 | |
|---|---|
| | |

Classification scores:
P(cat)

Source: cs231n course, Stanford University

## Classification + Localization

# Sliding Window: Overfeat



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

| 0.5 | 0.75 |
|-----|------|
|     |      |

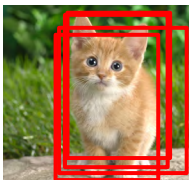Classification scores:
P(cat)

Source: cs231n course, Stanford University

## Classification + Localization

Sliding Window: Overfeat



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

| 0.5 | 0.75 |
|-----|------|
| 0.6 |      |

Classification scores:
P(cat)

Introduction
○○○○○○○○○○○

Datasets
○○○○

Localization
○○○○○○○○○○○○○○○○●○○○○○○○○

## Classification + Localization

### Sliding Window: Overfeat



Network input:
3 x 221 x 221



Larger image:
3 x 257 x 257

| 0.5 | 0.75 |
|-----|------|
| 0.6 | 0.8 |

Classification scores:
P(cat)

Source: cs231n course, Stanford University

Introduction
○○○○○○○○○○○○○

Datasets
○○○○

Localization
○○○○○○○○○○○○○○○●○○○○○○○

## Classification + Localization

# Sliding Window: Overfeat



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

| 0.5 | 0.75 |
|-----|------|
| 0.6 | 0.8 |

Classification scores:
P(cat)

## Classification + Localization

# Sliding Window: Overfeat

Greedily merge boxes and
scores (details in paper)



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

0.8

Classification score: P
(cat)

Source: cs231n course, Stanford University
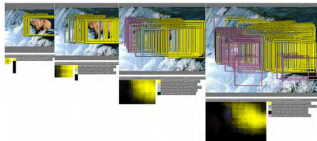
## Classification + Localization

# Sliding Window: Overfeat

In practice use many sliding window
locations and multiple scales
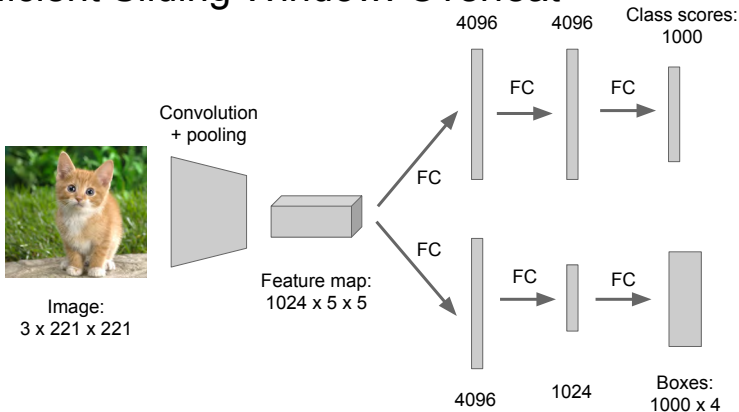
Window positions + score maps

Box regression outputs

Final Predictions

Source: cs231n course, Stanford University

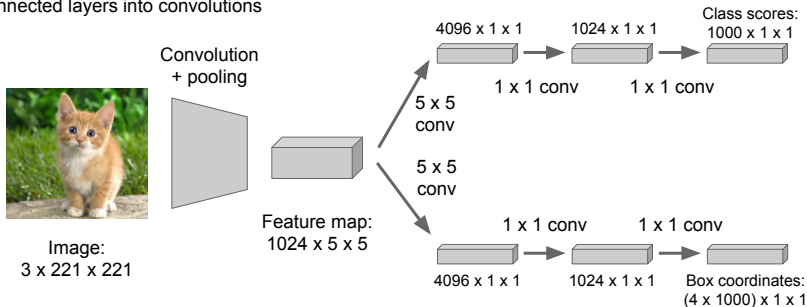## Classification + Localization

# Efficient Sliding Window: Overfeat

Introduction
○○○○○○○○○○○

Datasets
○○○○

Localization
○○○○○○○○○○○○○○○○○○○○●○○

## Classification + Localization

# Efficient Sliding Window: Overfeat

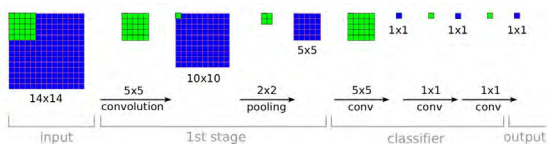Efficient sliding window by converting fully-connected layers into convolutions
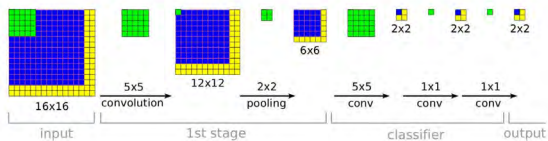


Source: cs231n course, Stanford University

## Classification + Localization

# Efficient Sliding Window: Overfeat

**Training time:** Small image, 1 x 1 classifier output



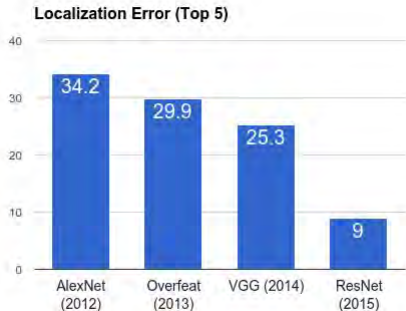**Test time:** Larger image, 2 x 2 classifier output, only extra compute at yellow regions



Sermanet et al, "Integrated Recognition, Localization and Detection using Convolutional Networks", ICLR 2014

Source: cs231n course, Stanford University

Introduction
OOOOOOOOOOO

Datasets
OOOO

Localization
OOOOOOOOOOOOOOOOOOOOOO●

## Classification + Localization

# ImageNet Classification + Localization



**AlexNet**: Localization method not published

**Overfeat**: Multiscale convolutional regression with box merging

**VGG**: Same as Overfeat, but fewer scales and locations; simpler method, gains all due to deeper features

**ResNet:** Different localization method (RPN) and much deeper features