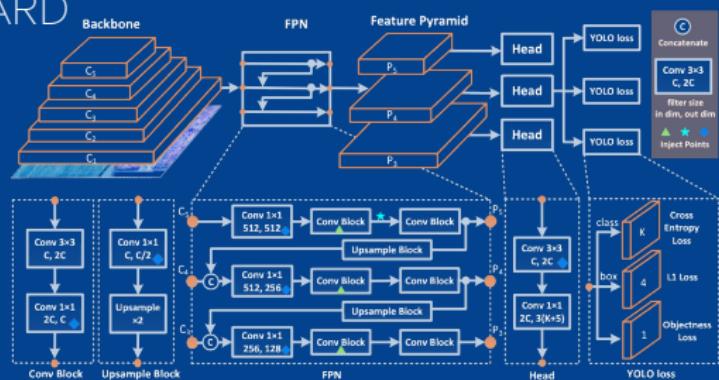




Lesson 10: Introduction to Advanced CNNs

CARSTEN EIE FRIGAARD

FALL 2023



Agenda

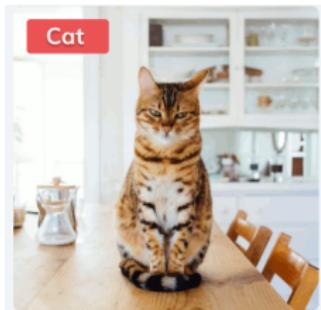
- ▶ Introduction to Advanced CNNs,
- ▶ In this-or-the-next lesson:
 - ▶ Real-world Examples,
 - ▶ Dimensionality Reduction,
 - ▶ Unsupervised Learning,
 - ▶ Reinforced Learning.

INTRO. TO ADVANCED CNNS

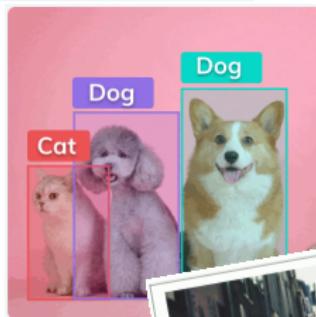


Computer Vision Tasks in ML

- ▶ classification,



- ▶ object detection,



- ▶ segmentation,
 - ▶ semantic segmentation
 - ▶ instance segmentation
 - ▶ panoptic segmentation



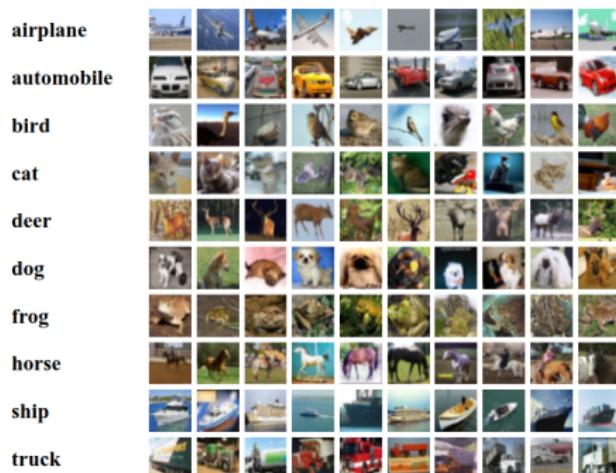
Figure 14-26. Semantic segmentation

Computer Vision Related Datasets

CICAR-10/100 (Canadian Institute For Advanced Research):

- ▶ "The CIFAR-10 and CIFAR-100 are labeled subsets of the 80 million tiny images dataset [...] The CIFAR-10 dataset consists of **60000 32x32 colour images in 10 classes**, with **6000 images per class**. There are **50000 training images** and **10000 test images**."

[<https://www.cs.toronto.edu/~kriz/cifar.html>]



Computer Vision Related Datasets

coco (Common Objects in Context):

- ▶ "COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features: Object segmentation, Recognition in context, [...] 330K images (>200K labeled), 1.5 million object instances, 80 object categories... [...]"

[<https://cocodataset.org/#home>]



Dataset examples

Computer Vision Related Datasets

ImageNet: ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

- "The ImageNet project is a large visual database designed for use in visual object recognition software research. More than 14 million images have been hand-annotated by the project to indicate what objects are pictured and in at least one million of the images, bounding boxes are also provided. ImageNet contains more than 20,000 categories [...]"

[<https://en.wikipedia.org/wiki/ImageNet>]



Data Augmentation

Data Augmentation

Data augmentation artificially increases the size of the training set by generating many realistic variants of each training instance. This reduces overfitting, making this a regularization technique. The generated instances should be as realistic as possible: ideally, given an image from the augmented training set, a human should not be able to tell whether it was augmented or not. Moreover, simply adding white noise will not help; the modifications should be learnable (white noise is not).

For example, you can slightly shift, rotate, and resize every picture in the training set by various amounts and add the resulting pictures to the training set (see Figure 14-12). This forces the model to be more tolerant to variations in the position, orientation, and size of the objects in the pictures. If you want the model to be more tolerant to different lighting conditions, you can similarly generate many images with various contrasts. In general, you can also flip the pictures horizontally (except for text, and other non-symmetrical objects). By combining these transformations you can greatly increase the size of your training set.

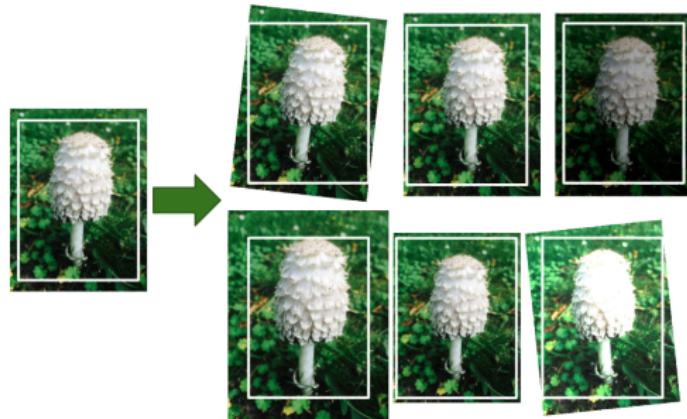


Figure 14-12. Generating new training instances from existing ones

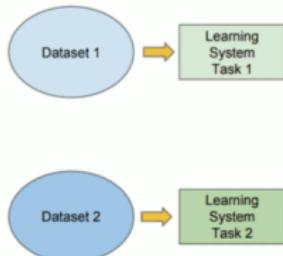
Pretrained Models and Transfer Learning

Traditional ML

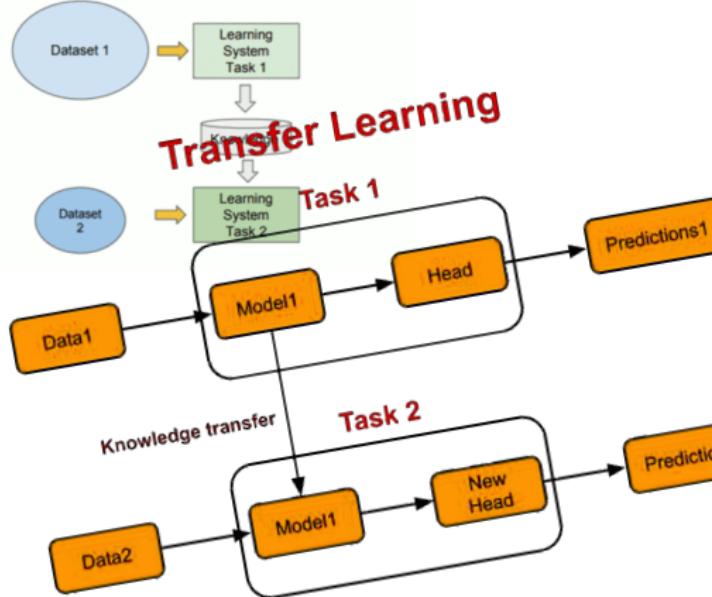
vs

Transfer Learning

- Isolated, Single task learning.
- Knowledge is not retained or accumulated. Learning is performed w/o consideration for knowledge learned from other tasks.



- Learning new tasks relies on previously learned tasks.
- Learning process can be faster, more accurate and/or need less training data.



[https://research.aimultiple.com/wp-content/uploads/2020/07/Transfer_Learning_Explanation-1920x943.png]
[https://www.topbots.com/wp-content/uploads/2019/12/cover_transfer_learning_1600px_web.jpg]

CNN Architectures

GoogLeNet, ResNet (Xception, SENet)

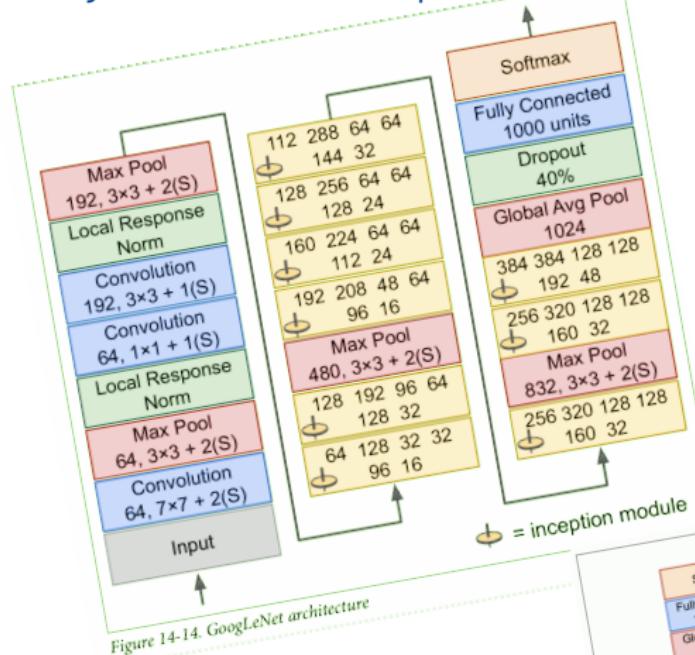
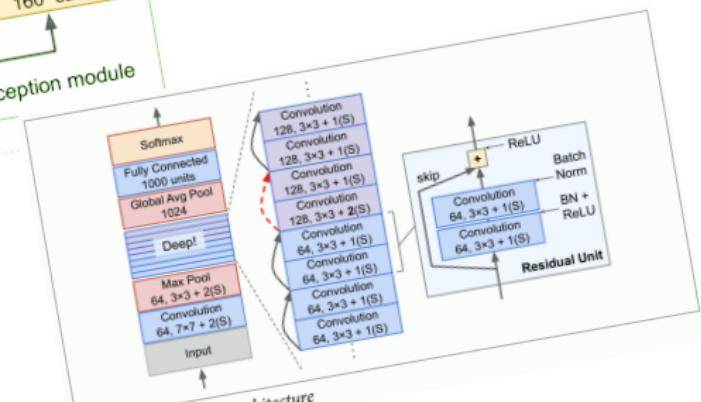


Figure 14-17. ResNet architecture



CNN Architectures

CNN Architectures:

- ▶ Fast R-CNN: .. (older)
- ▶ Faster/Mask R-CNN:
 - ▶ ResNet-50 (R50), ResNet-101 (R101)
 - ▶ ResNeXt-152 (X101)?

CNN 'Frameworks':

- ▶ Framework: Yolo (C, PyTorch)
 - Model Zoo:
 - ▶ Darknet53 nano, small, medium, large, xlarge
- ▶ Framework: Detectron2 (PyTorch)
 - Model Zoo:
 - ▶ Faster R-CNN: R50, R101, X101
 - ▶ RetinaNet: R50, R101 bases
 - ▶ RPN and Fast R-CNN: RPN R50-XX, Fast R-CNN R50-FPN..

CNN Architectures

YOLO version 3/4/5[..]/8

PyTorch

Get Started

Ecosystem

Mobile

Blog

Tutoria

YOLOv5

[View on Github](#)

[Open on Google Colab](#)



BEFORE YOU START

Start from a **Python>=3.8** environment with **PyTorch**:
<https://pytorch.org/get-started/locally/>. To install YOLOv5:

```
pip install -qr https://raw.githubusercontent.com/ultralytics/yolov5/master/requirements.txt
```

MODEL DESCRIPTION



Nano
YOLOv5n

4 MB_{FP16}
6.3 ms₁₀₀
28.4 mAP_{coco}

Small
YOLOv5s

14 MB_{FP16}
6.4 ms₁₀₀
37.2 mAP_{coco}

Medium
YOLOv5m

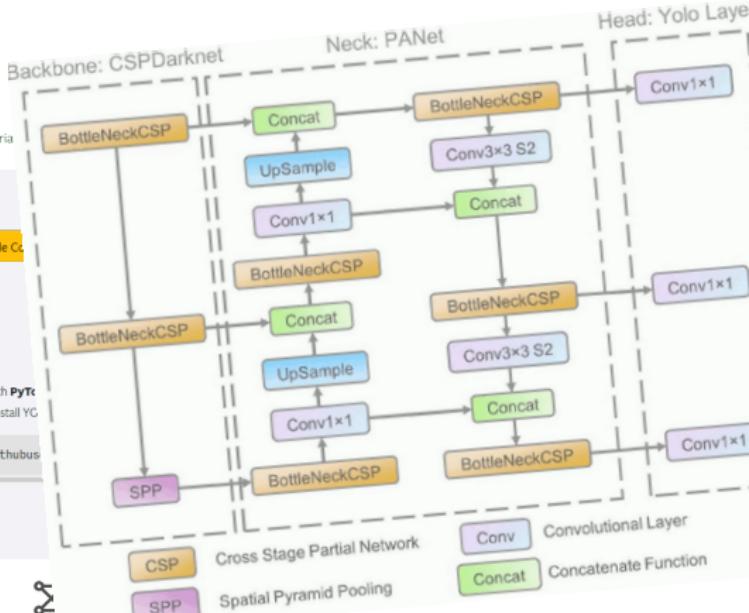
41 MB_{FP16}
8.2 ms₁₀₀
45.2 mAP_{coco}

Large
YOLOv5l

89 MB_{FP16}
10.1 ms₁₀₀
48.6 mAP_{coco}

XLarge
YOLOv5x

166 MB_{FP16}
12.1 ms₁₀₀
50.7 mAP_{coco}



- ▶ "All the YOLOv5 models are composed of the same 3 components: CSP-Darknet53 as a backbone, SPP and PANet in the model neck and the head used in YOLOv4 "

Pretrained Models for Transfer Learning

Detectron2 ImageNet Pretrained Zoo

ImageNet Pretrained Models

It's common to initialize from backbone models pre-trained on ImageNet classification tasks. The following backbone models are available:

- R-50.pkl: converted copy of MSRA's original ResNet-50 model.
- R-101.pkl: converted copy of MSRA's original ResNet-101 model.
- X-101-32x8d.pkl: ResNeXt-101-32x8d model trained with Caffe2 at FB.
- R-50.pkl (torchvision): converted copy of torchvision's ResNet-50 model.
More details can be found in the [conversion script](#).

Note that the above models have **different** format from those provided in Detectron: we do not fuse BatchNorm into an affine layer. Pretrained models in Detectron's format can still be used. For example:

- X-152-32x8d-IN5k.pkl: ResNeXt-152-32x8d model trained on ImageNet-5k with Caffe2 at FB (see ResNeXt paper for details on ImageNet-5k).
- R-50-GN.pkl: ResNet-50 with Group Normalization.
- R-101-GN.pkl: ResNet-101 with Group Normalization.

Performance Metric: IoU

Intersection over Union (IoU)

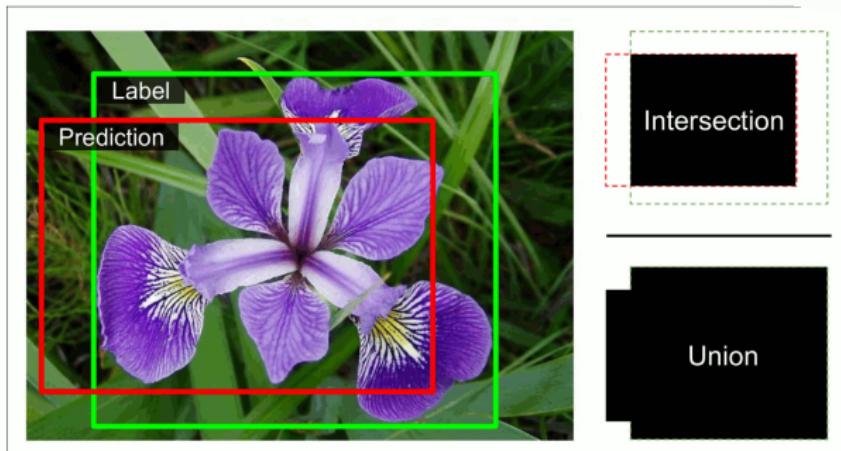
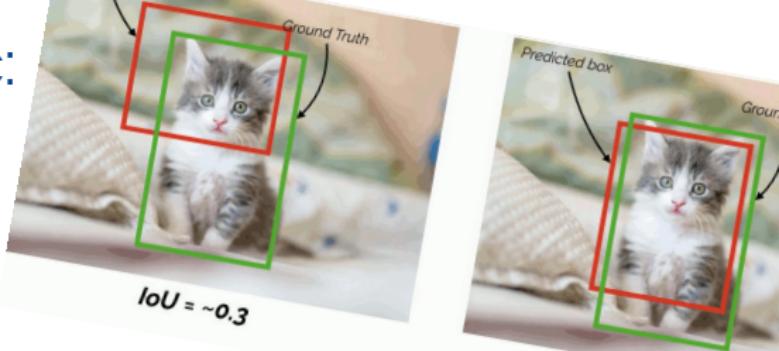


Figure 14-23. Intersection over Union (IoU) Metric for Bounding Boxes

Range: $0 \leq \text{IoU} \leq 1$ (best)

New hyperparameter:

threshold for 'correct' detection, say $\text{IoU} > 0.5$.

Performance Metric: mAP

Mean Average Precision (mAP)

Mean Average Precision (mAP)

A very common metric used in object detection tasks is the *mean Average Precision* (mAP). “Mean Average” sounds a bit redundant, doesn’t it? To understand this metric, let’s go back to two classification metrics we discussed in [Chapter 3](#): precision and recall. Remember the tradeoff: the higher the recall, the lower the precision. You can visualize this in a Precision/Recall curve (see [Figure 3-5](#)). To summarize this curve into a single number, we could compute its Area Under the Curve (AUC). But note that the Precision/Recall curve may contain a few sections where precision actually goes up when recall increases, especially at low recall values (you can see this at the top left of [Figure 3-5](#)). This is one of the motivations for the mAP metric.

Suppose the classifier has a 90% precision at 10% recall, but a 96% precision at 20% recall: there’s really no tradeoff here: it simply makes more sense to use the classifier at 20% recall rather than at 10% recall, as you will get both higher recall and higher precision. So instead of looking at the precision at 10% recall, we should really be looking at the *maximum* precision that the classifier can offer with *at least* 10% recall. It would be 96%, not 90%. So one way to get a fair idea of the model’s performance is to compute the maximum precision you can get with at least 0% recall, then 10% recall, 20%, and so on up to 100%, and then calculate the mean of these maximum precisions. This is called the *Average Precision* (AP) metric. Now when there are more than 2 classes, we can compute the AP for each class, and then compute the mean AP (mAP). That’s it!

However, in an object detection systems, there is an additional level of complexity: what if the system detected the correct class, but at the wrong location (i.e., the bounding box is completely off)? Surely we should not count this as a positive prediction. So one approach is to define an IOU threshold: for example, we may consider that a prediction is correct only if the IOU is greater than, say, 0.5, and the predicted class is correct. The corresponding mAP is generally noted mAP@0.5 (or mAP@50%, or sometimes just AP₅₀). In some competitions (such as the Pascal VOC challenge), this is what is done. In others (such as the COCO competition), the mAP is computed for different IOU thresholds (0.50, 0.55, 0.60, ..., 0.95), and the final metric is the mean of all these mAPs (noted AP@[.50:.95] or AP@[.50:0.05:.95]). Yes, that’s a mean mean average.

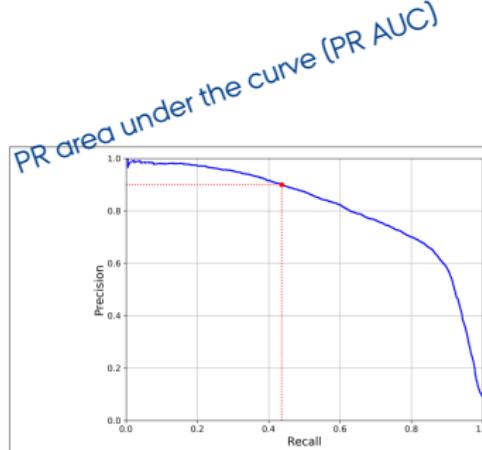


Figure 3-5. Precision versus recall

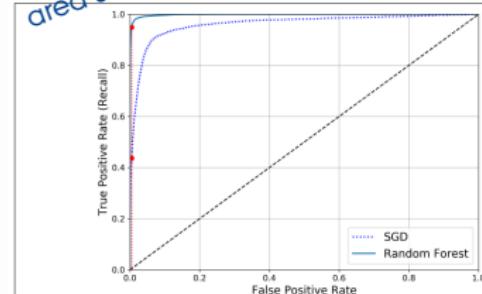


Figure 3-7. Comparing ROC curves

Performance Metric: mAP

Mean Average Precision (mAP)

Average Precision, AP summarizes a precision-recall curve

$$AP = \sum_n (R_n - R_{(n-1)})P_n$$

and threshold $n = 0.1, 0.2, \dots, 0.9$ for predicted scores

```
pred_scores = [0.7, 0.3, 0.5, 0.9, ..., 0.3]
```

mAP is the mean over all classes for AP.

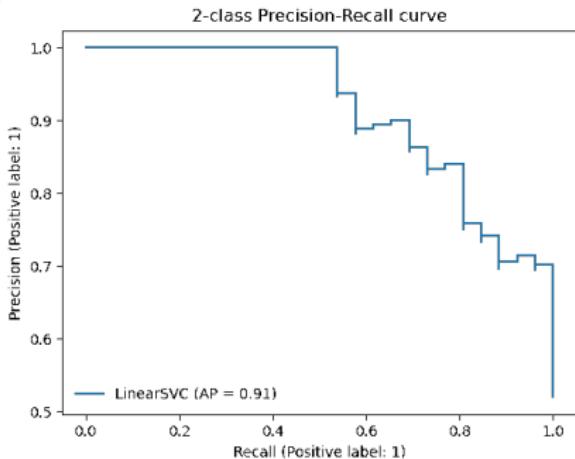
mAP@0.5 = mAP@50% = AP₅₀ =>

mAP at IoU > 0.5

and

mAP@[.50:.95] or mAP@[.50:0.05:.95]

is **mean mean average!**



Naïve Object Detection: Sliding Window

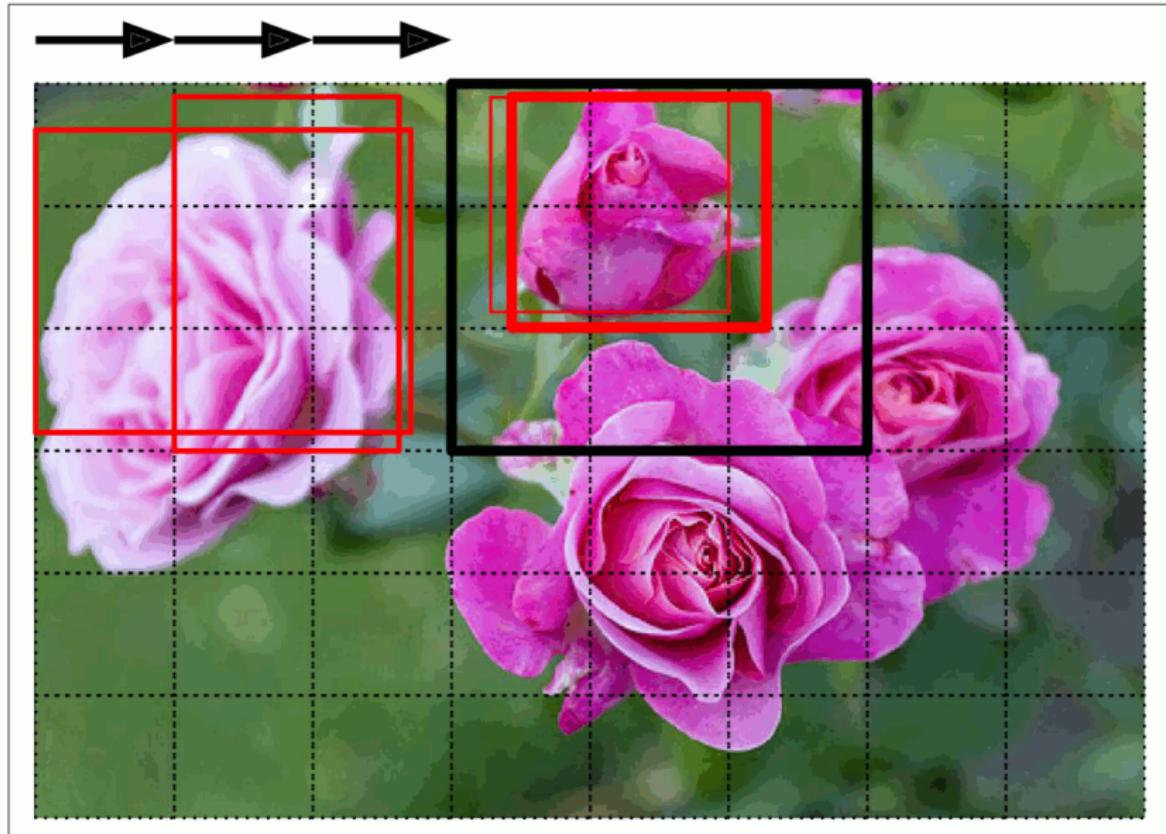


Figure 14-24. Detecting Multiple Objects by Sliding a CNN Across the Image

Fully Connected to Fully Convolutional Networks: YOLO

- It outputs two bounding boxes for each grid cell (instead of just one), which allows the model to handle cases where two objects are so close to each other that their bounding box centers lie within the same cell. Each bounding box also comes with its own objectness score.
- YOLO also outputs a class probability distribution for each grid cell, predicting 20 class probabilities per grid cell since YOLO was trained on the PASCAL VOC dataset, which contains 20 classes. This produces a coarse *class probability map*. Note that the model predicts one class probability distribution per grid cell, not per bounding box. However, it's possible to estimate class probabilities for each bounding box during postprocessing, by measuring how well each bounding box matches each class in the class probability map. For example, imagine a picture of a person standing in front of a car. There will be two bounding boxes: one large horizontal one for the car, and a smaller vertical one for the person. These bounding boxes may have their centers within the same grid cell. So how can we tell which class should be assigned to each bounding box? Well, the class probability map will contain a large region where the "car" class is dominant, and inside it there will be a smaller region where the "person" class is dominant. Hopefully, the car's bounding box will roughly match the "car" region, while the person's bounding box will roughly match the "person" region: this will allow the correct class to be assigned to each bounding box.

YOLO was originally developed using Darknet, an open source deep learning framework initially developed in C by Joseph Redmon, but it was soon ported to TensorFlow, Keras, PyTorch, and more. It was continuously improved over the years, with YOLOv2, YOLOv3, and YOLO9000 (again by Joseph Redmon et al.), YOLOv4 (by Alexey Bochkovskiy et al.), YOLOv5 (by Glenn Jocher), and PP-YOLO (by Xiang Long et al.).

Each version brought some impressive improvements in speed and accuracy, using a variety of techniques; for example, YOLOv3 boosted accuracy in part thanks to *anchor priors*, exploiting the fact that some bounding box shapes are more likely than others, depending on the class (e.g., people tend to have vertical bounding boxes while cars usually don't). They also increased the number of bounding boxes per grid cell, they trained on different datasets with many more classes (up to 9,000 classes organized in a hierarchy in the case of YOLO9000), they added skip connections to recover some of the spatial resolution that is lost in the CNN (we will discuss this shortly, when we look at semantic segmentation), and much more. There are many variants of these models too, such as YOLOv4-tiny, which is optimized to train on less powerful machines and which can run extremely fast (at over 1,000 frames per second!), but with a slightly lower *mean average precision* (mAP).

Mean Average Precision

A very common metric used in object detection tasks is the mean average precision. "Mean average" sounds a bit redundant, doesn't it? To understand this metric, let's go back to two classification metrics we discussed in Chapter 3: precision and recall. Remember the trade-off: the higher the recall, the lower the precision. You can visualize this in a precision/recall curve (see Figure 3-6). To summarize this curve into a single number, we could compute its area under the curve (AUC). But note that the precision/recall curve may contain a few sections where precision actually goes up when recall increases, especially at low recall values (you can see this at the top left of Figure 3-6). This is one of the motivations for the mAP metric.

Suppose the classifier has 90% precision at 10% recall, but 96% precision at 20% recall. There's really no trade-off here: it simply makes more sense to use the classifier at 20% recall rather than 10% recall, as you will get both higher recall and higher precision. So instead of looking at the precision *at* 10% recall, we should really be looking at the *maximum* precision that the classifier can offer *with at least 10% recall*. It would be 96%, not 90%. Therefore, one way to get a fair idea of the model's performance is to compute the maximum precision you can get with at least 0% recall, then 10% recall, 20%, and so on up to 100%, and then calculate the mean of these maximum precisions. This is called the *average precision* (AP) metric. Now when there are more than two classes, we can compute the AP for each class, and then compute the mean AP (mAP). That's it!

In an object detection system, there is an additional level of complexity: what if the system detected the correct class, but at the wrong location (i.e., the bounding box is completely off)? Surely we should not count this as a positive prediction. One approach is to define an IoU threshold: for example, we may consider that a prediction is correct only if the IoU is greater than, say, 0.5, and the predicted class is correct. The corresponding mAP is generally noted mAP@0.5 (or mAP@50%, or sometimes just AP₅₀). In some competitions (such as the PASCAL VOC challenge), this is what is done. In others (such as the COCO competition), the mAP is computed for different IoU thresholds (0.50, 0.55, 0.60, ..., 0.95), and the final metric is the mean of all these mAPs (noted mAP@[.50:.95] or mAP@[.50:0.05:.95]). Yes, that's a mean mean average.

YOLO, You Only Look Once

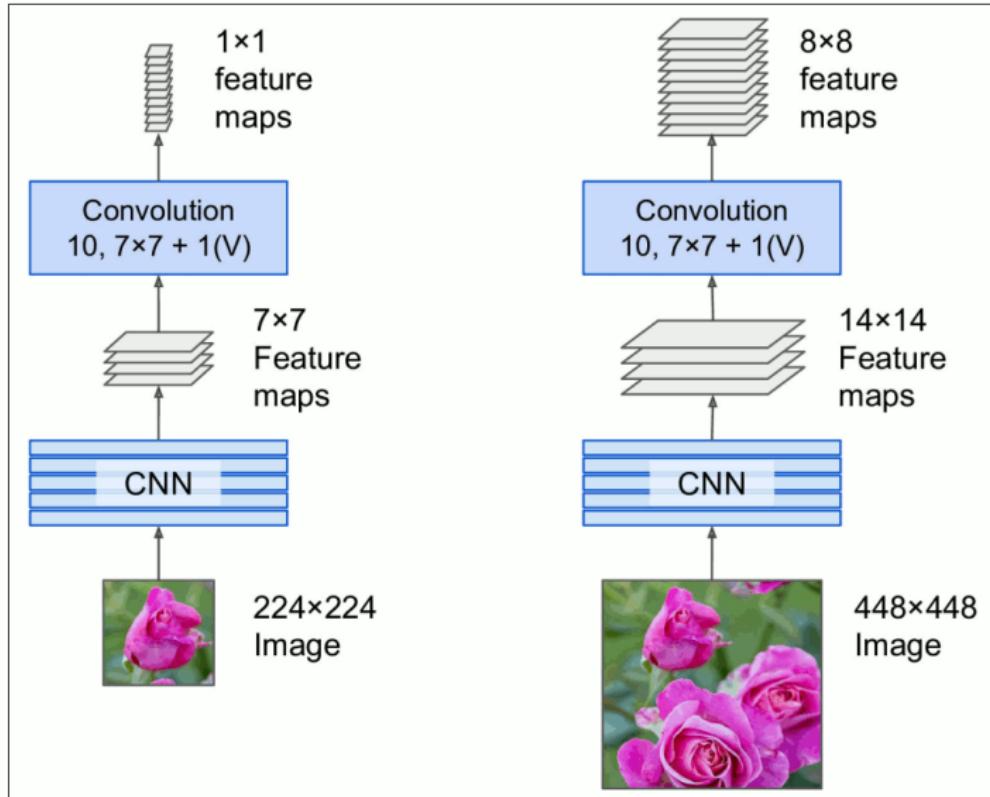
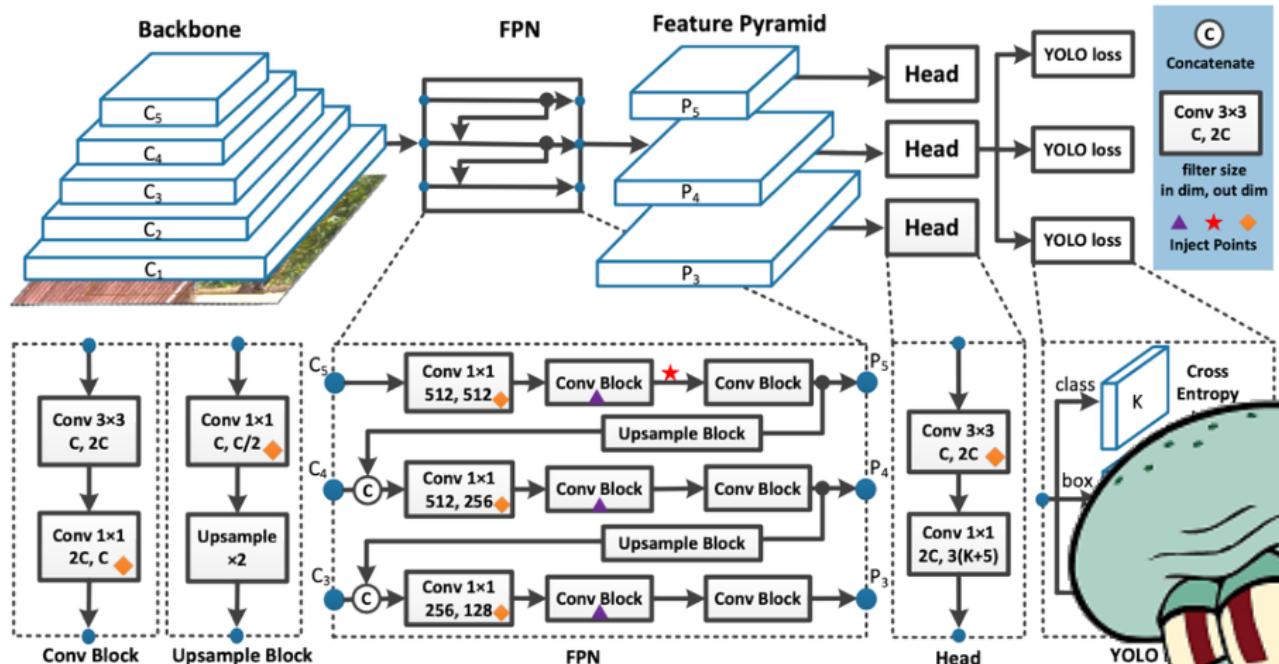


Figure 14-25. A Fully Convolutional Network Processing a Small Image (left) and a Large One (right)

YOLO, You Only Look Once



[<https://blog.roboflow.com/yolov7-breakdown/>]

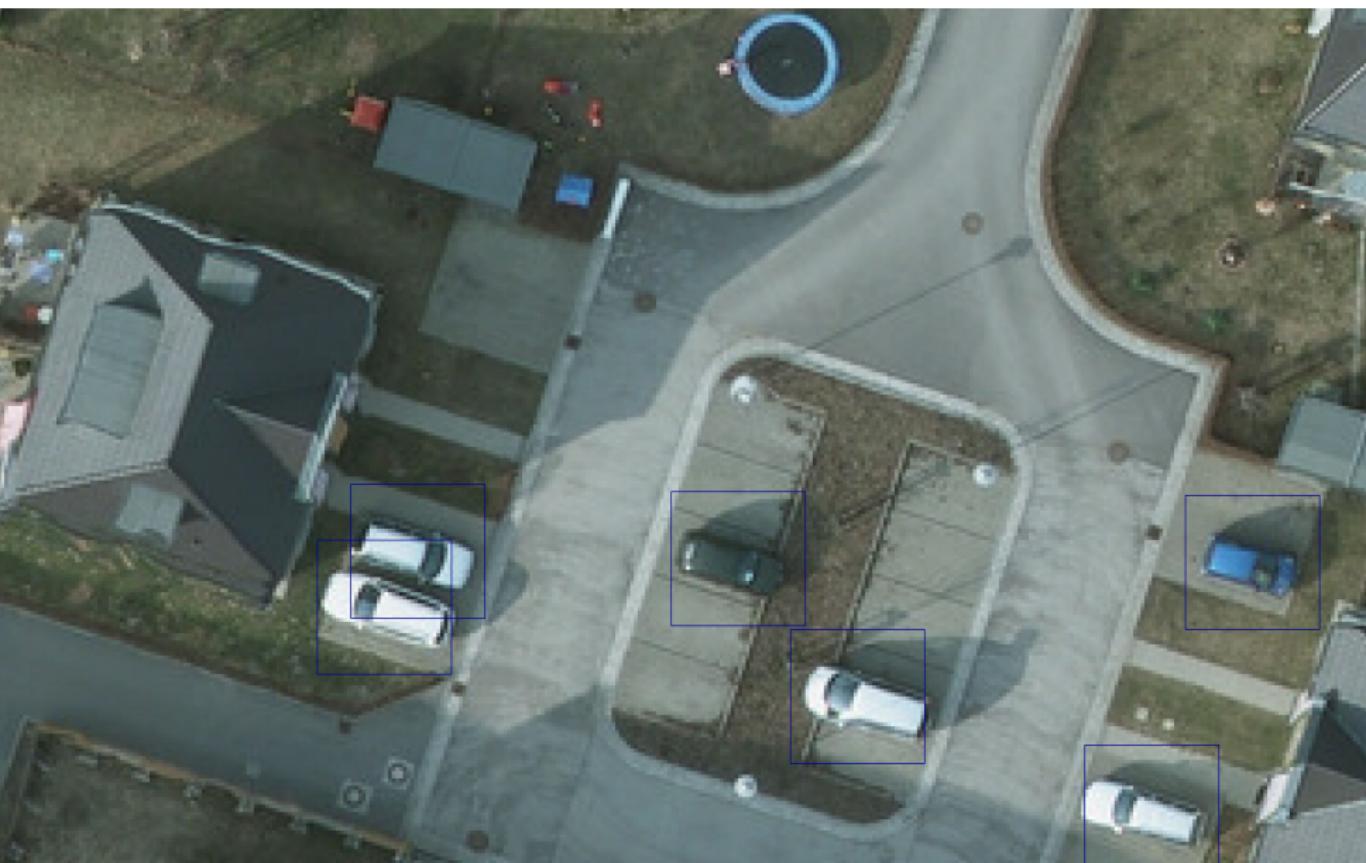
YOLO, You Only Look Once

My Postdam car recognition on 26640x26640 sized images



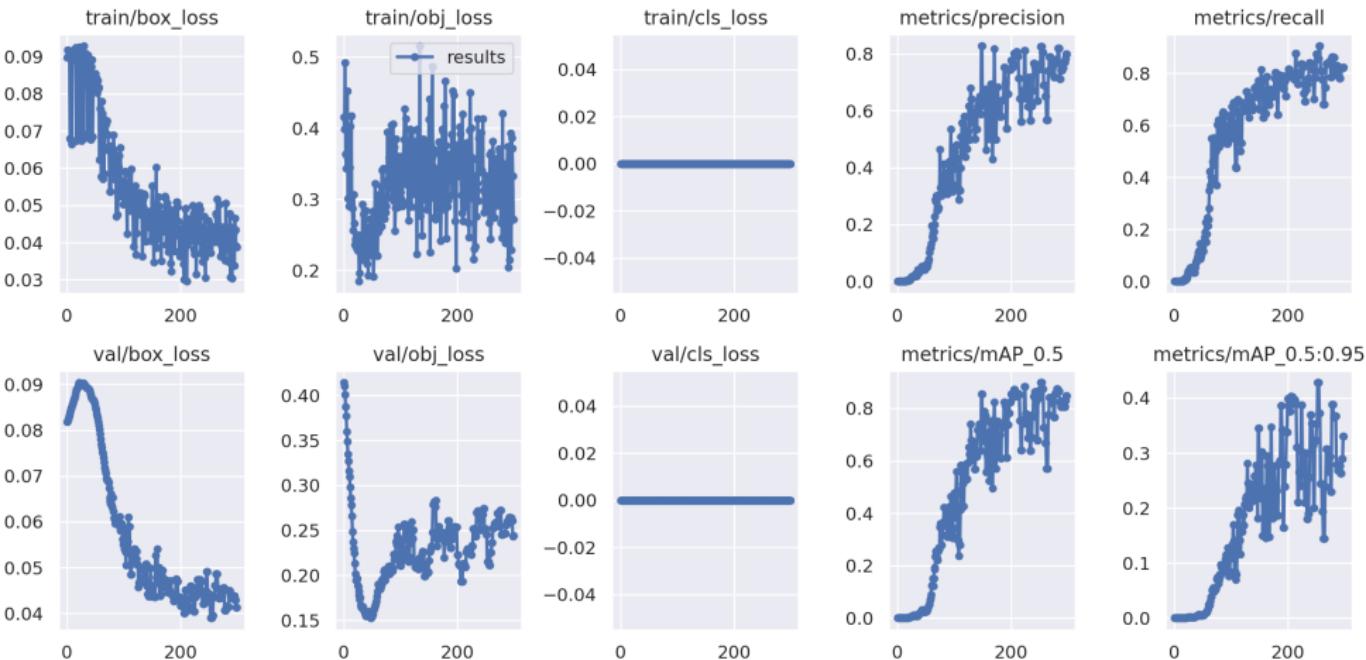
YOLO, You Only Look Once

My Postdam car recognition on 26640x26640 sized images



YOLO, You Only Look Once

Generated learning curves



YOLO, You Only Look Once

Generated learning curves

