
Seminar Computer Vision by Deep Learning

Classification with Localization and Single-stage object detectors

E. Geenjaar
4373480

P. de Rijk
4375017

S. Tak
4975154

Abstract

1 This chapter aims to establish an understanding of classification with localization,
2 and extend that understanding to single-stage detectors. The first part of this chapter
3 will cover the concept of classification as performed by contemporary computer
4 vision systems, and demonstrate how this can be extended to include localization.
5 Additionally, the concept of a sliding window detector will be elaborated on.
6 Building upon this knowledge, single-stage detectors are discussed with a focus on
7 YOLO. The basic principle of YOLO is to approach object detection as a single
8 regression problem, where image pixels are translated to bounding box coordinates
9 and class probabilities in one stage. During testing, only the best predictions are
10 kept using non-maximal suppression. The limitations of YOLO are discussed after
11 which solutions to these limitations are examined. The solutions include anchor
12 boxes and usage of multiple feature map scales. Finally, real-world applications
13 and controversies surrounding YOLO are covered.

14 1 Classification with Localization

15 We are already familiar with the concept of object detection, introduced in the previous chapter (5.1).
16 There are however two main tasks that precede object detection: **classification** and **classification**
17 **with localization**. The distinction between these tasks is illustrated in figure 1 and can be formulated
18 as follows: In a classification task the central question that needs answering is: *Is there a dog in the*
19 *image?* The classification with localization task extends this question to *Is there a dog and wherein*
20 *the image is it located?* Finally, with object detection, the question becomes: *Are there any objects*
21 *and wherein the image are they located?* This section will focus on answering the first two questions,
22 by first explaining how we can achieve a model that performs image classification, and subsequently
23 how this model can be extended to include localization. These concepts will form the basis for section
24 2, which discusses state-of the art single-stage object detection methods.

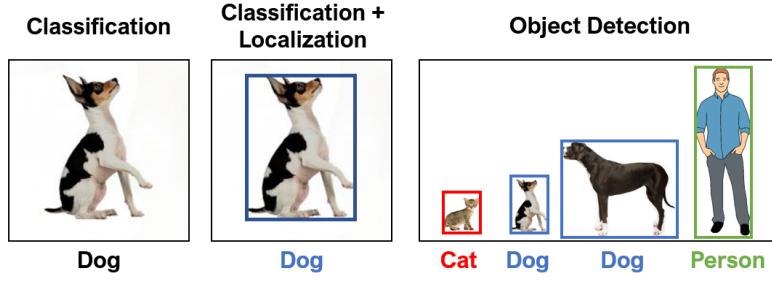


Figure 1: Overview comparison between Classification, Classification with localization and Object Detection

25 1.1 Classification

26 **Image classification** is one of the core tasks of computer vision. It forms the basis for other computer
27 vision tasks such as localization, detection, and segmentation [7]. Classification can be described as:

28 "The mapping from an input x to an output class in y , which consists of k categories, through a
29 function $y = f(x)$."

30 This definition is exemplified using figure 2: The input image x is represented as a $100 \times 100 \times 3$ array,
31 where the first two dimensions are the height and width, and the third dimension represents the three
32 color channels (red, green and blue) respectively. The array contains values for the pixel intensities
33 that range from 0 to 255. The output y is represented as a vector containing the probabilities for one
34 of $k = 4$ categories: $\{\text{dog}, \text{cat}, \text{horse}, \text{car}\}$. A good classifier f will be able to map the values in the
35 input image to the correct label for that image.

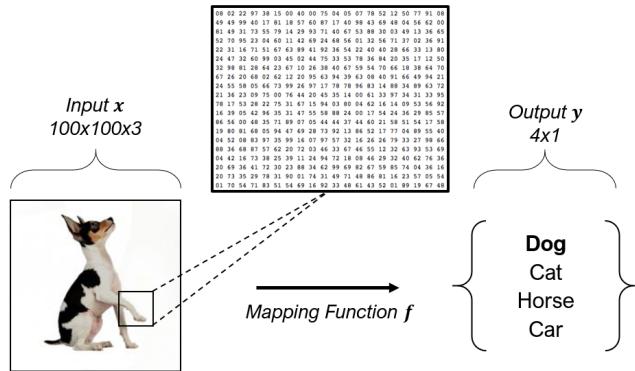


Figure 2: A good classifier can map an input image x in the form of an array to a correct label contained in the output vector y . Inspired by Karpathy [7]

36 Even though this classification task may seem trivial to a human, it can pose quite a lot of challenges
37 for a computer vision system. The algorithm only "sees" a large array of numbers, and it needs to
38 translate this array into the *semantic* idea of a, for example, a dog. Complications discussed in chapter
39 1.1, such as viewpoint-dependent object variability and the high in-class variability of object types
40 [3] make it hard for a classifier to *robustly* map each image of a certain object to its correct label:
41 There are many types of dogs; small dogs, large dogs, white dogs, brown dogs, etc. Dogs can also
42 be captured on an image from many viewpoints. All these variations make it hard for a classifier to
43 map the image array it receives as an input to the semantic label "dog". This problem of translating a
44 numerical array to the semantic classification of an object is also referred to as the **semantic gap** [7].

45 **1.1.1 Closing the Semantic Gap**

46 So how do we close the semantic gap? Traditionally, a dual-stage approach was used, where
47 handcrafted features were extracted from images using *feature descriptors* such as edge detectors and
48 corner detectors. These features were then used as the input to a trainable classifier [13], as shown in
49 figure 3. This approach suffered from the fact that the accuracy was highly dependent on the design
50 of the feature extraction stage, which in itself has shown to be a very challenging task [9].

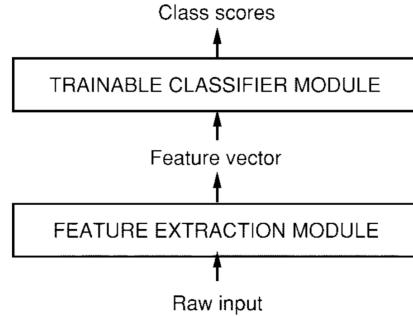


Figure 3: Traditional pattern recognition is performed with two modules: a fixed feature extractor and a trainable classifier. Taken from LeCun et al. [9]

51 The increase in computing power and the advent of larger datasets have enabled the adoption of
52 convolutional neural networks (ConvNets) in the field of computer vision. These ConvNets have
53 driven progress on classification tasks. The most significant advancement was achieved in the 2012
54 ImageNet Large Scale Visual Recognition Challenge [18]. A new ConvNet proposed by Krizhevsky
55 et al. [8] beat the competition by an error margin of over 10%, establishing the dominance of
56 ConvNets on classification tasks for the coming decade.

57 **1.1.2 Feature maps**

58 So how do ConvNets differ from the traditional methods which used manually crafted features to
59 classify an image? The main benefit arises from the fact that ConvNets can learn this feature map
60 by training on labelled data. **Feature maps** in ConvNets are 3-dimensional arrays, where the first
61 two dimensions are spatial dimensions and the third dimension contains the features for each spatial
62 location in the feature map. An example of a feature map is shown in figure 4. The network that
63 is used to produce the feature map is called a **backbone network**. The features in a feature map
64 are learned in such a way that they can be used to distinguish between the different class labels in a
65 classification task.

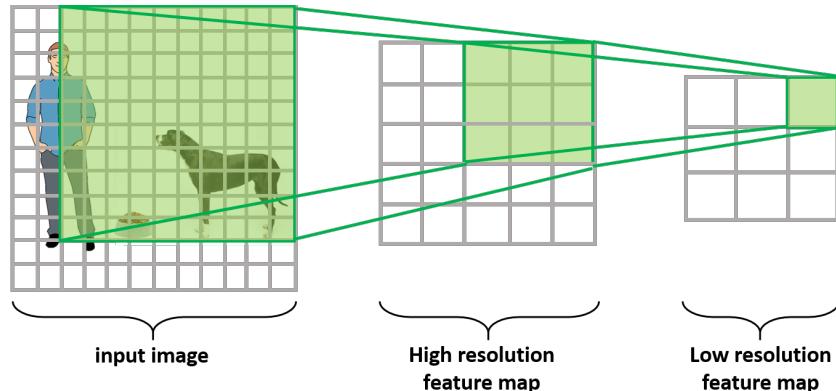


Figure 4: An example of a feature map, a simplified backbone network and the receptive field of the feature map

66 Throughout this chapter, we will refer to the spatial locations in a feature map as **pixels** in the feature
 67 map. A high spatial resolution of the feature map indicates a large number of pixels in the feature
 68 map and a low spatial resolution to a small number of pixels in the feature map. A high feature
 69 resolution refers to a large number of features (the third dimension) of a feature map, a small feature
 70 resolution refers to a small number of features in the feature map. The feature map has a lower spatial
 71 resolution than the input image and thus fewer pixels. These feature map pixels do however extract
 72 information from a very large part of the image. The part of the input image a pixel in the feature
 73 map can 'see' is referred to as the **receptive field** of that feature map pixel. The receptive field of a
 74 feature map pixel is often just shy of the whole image. Figure 4 shows an example of a feature map's
 75 receptive field. Two important factors that determine the receptive field of a feature map pixel are the
 76 *number of layers* in a backbone network and the *dilation rate* used in those layers. The backbone
 77 networks used for computer vision classification tasks are often ConvNets. An example of a feature
 78 map, a simplified backbone network and the receptive field of the feature map is shown in figure 4

79 To create classifications out of the feature map, a fully connected layer is used. A flattened version of
 80 the feature map is passed through a fully connected layer, which learns the mapping from the feature
 81 map to the class labels. The fully connected layer will push feature maps to contain information
 82 that is important for the classification of the image and these two parts can be trained end-to-end.
 83 **End-to-end** training is where all parameters of the network are trained in conjunction. The backbone
 84 network and the fully connected layer together are the classifier that can map an input x to an output
 85 class y . The number of output neurons of the fully connected layer is equal to the number of classes k
 86 we want to predict. Each neuron corresponds to a certain class label and will predict an unnormalized
 87 score for that class label. Figure 5 illustrates an example of the classification pipeline.

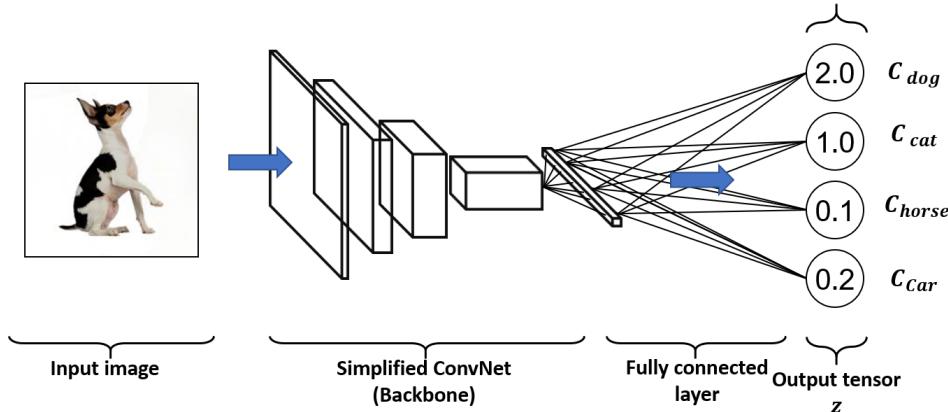


Figure 5: An example of the classification pipeline. The input image is fed through a backbone network containing several convolutional layers, which are followed by a flattening layer and a fully connected layer resulting in the class scores z

88 1.1.3 Measuring Classification Performance

89 The problem with the fully connected layer in our classifier is that it will provide predictions for each
 90 class that are not normalized and can range between $[-\infty, \infty]$. These unnormalized predictions are
 91 also referred to as **logit scores** and are the outputs of the fully connected layer. To make the outputs
 92 of the classifier more intuitive, we want to map them to a probability for each class. In order to do
 93 that we can use the **Softmax** function, which is defined as follows:

$$\text{Softmax}(\mathbf{z}) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} = \hat{\mathbf{y}} \quad (1)$$

94 Where \mathbf{z} is a vector of logit scores for each class, i is the index of the class the classifier is predicting.
 95 The sum in the denominator sums over each logit score to make sure that the probabilities for each
 96 class sum to 1. We refer to the output of the softmax function as $\hat{\mathbf{y}}$. Figure 6 demonstrates an example
 97 of a mapping between the logit scores and the probabilities that are assigned to each class.

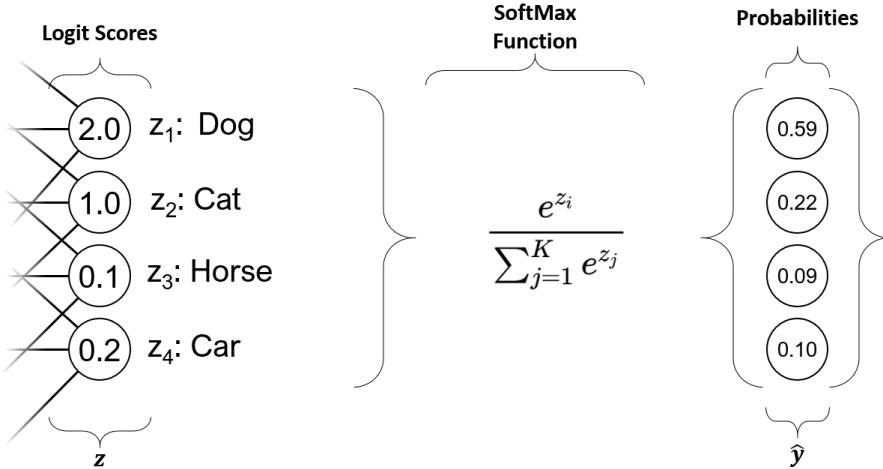


Figure 6: An example of how logit scores (left) are mapped to normalized probabilities (right) for each class using the softmax function (middle)

98 The loss function for classifiers can subsequently be defined by taking the difference between the
 99 ground truth class label y_i for an image and the discrete probability distribution that is given by the
 100 softmax. This is generally done using the cross entropy loss, which measures the difference between
 101 two probability distributions:

$$\mathcal{L}_{\text{cross-entropy}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_i y_i \log(\hat{y}_i) \quad (2)$$

102 Figure 7 illustrates the cross-entropy loss applied to our example network.

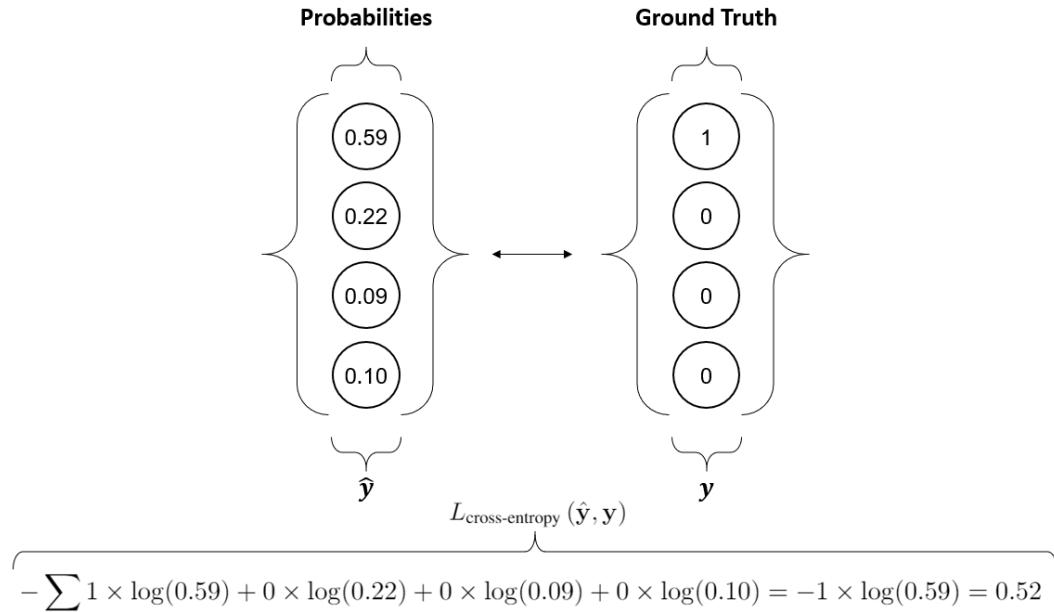


Figure 7: An example of how the cross-entropy loss is calculated on normalized probabilities for each class (left) and the ground truth class label (right)

103 The normalized probabilities that are predicted by the classifier are a discrete probability density
 104 function. The ground truth label is also a probability probability density function. The larger the
 105 difference between the predicted probability density function and the ground truth probability density
 106 function, the larger the gradients that flow back into the classifier.

107 **1.2 Extending to Localization**

108 In the previous subsection, we discussed image classification, where an image label is assigned to an
 109 input image. Assume we do not only want to classify an image as a "dog", but also want to know
 110 where it is located, we need to extend the classifier to be able to localize the dog. It turns out that the
 111 simplest way to localize an object in an image is by constructing a **bounding box** around the object.
 112 The bounding box is defined by 4 variables: b_x, b_y, b_w, b_h . These variables respectively represent the
 113 x and y position of the centre of the bounding box, the width, and the height. To make it easier to talk
 114 about the location of a bounding box in an image, we normalize the width and height of an image to
 115 $[0, 1]$. The left uppermost location of an image becomes $(0, 0)$ and the right bottom part becomes
 116 $(1, 1)$. Figure 8 illustrates an example of a bounding box.

117 Note that we are still only discussing classification and localization of *one* object in an image.

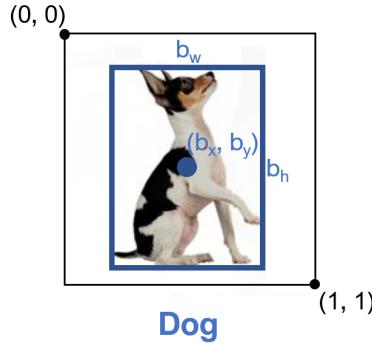


Figure 8: Bounding box around the object 'dog' including its centre and dimensions

118 So how do we add this additional localization task to our model? It turns out that a lot of the
 119 architecture used in the image classification model can be reused: Again, we feed the input image
 120 through our backbone network which in turn generates a feature map that summarizes the content
 121 of the image. We also use a fully-connected layer, but we add 5 output neurons that will learn to
 122 predict b_x, b_y, b_w, b_h and P_c . Where P_c is the probability that there is an object in the image. This
 123 probability is referred to as the **object score**. From now on we will refer to the classes that we can
 124 detect as C_i , where i indexes the k classes. Figure 9 illustrates what our example network would look
 125 like when it is extended to also handle localization.

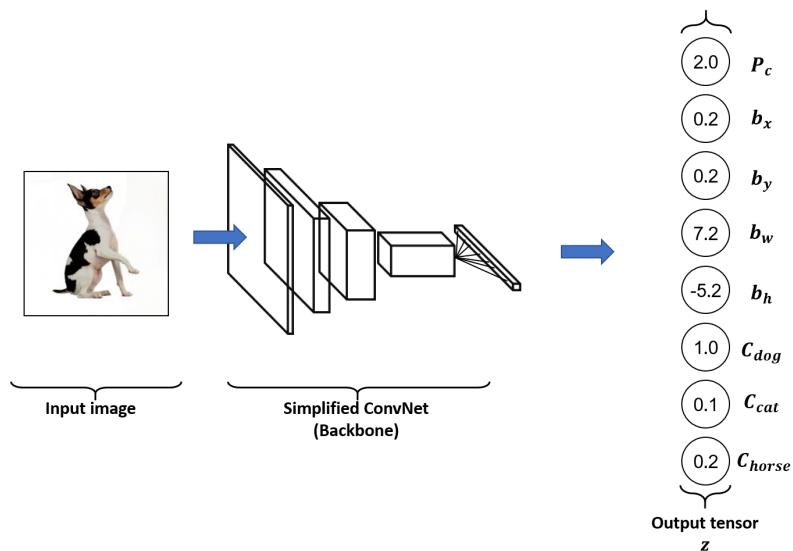


Figure 9: The classification with localization pipeline. The input image is fed into a backbone
 network and mapped onto an output tensor z by a fully-connected layer

126 **1.2.1 Measuring Localization Performance**

127 Classifiers are generally evaluated on their classification accuracy, which is the number of correctly
 128 classified samples divided by the total number of classified samples. This is not an adequate metric
 129 for classification and localization, because the metric does not measure how well the object detector
 130 localizes the object. It also does not take into account that there could be multiple objects in an image.
 131 **Intersection over Union** (IoU) $\in [0, 1]$ is a metric that can be used to determine the quality of the
 132 bounding box in localizing the object. The IoU is defined as follows:

$$\text{IoU} = \frac{\text{Pixels in ground truth bounding box} \cap \text{Pixels in predicted bounding box}}{\text{Pixels in ground truth bounding box} \cup \text{Pixels in predicted bounding box}} \quad (3)$$

133 Where \cap indicates the intersection. Hence, the numerator in equation 3 refers to the intersection
 134 between the pixels in the ground truth bounding box and the pixels in the predicted bounding box. \cup
 135 refers to the union, so the denominator in equation 3 refers to the union of the pixels in the ground
 136 truth bounding box and the pixels in the predicted bounding box. The formula is visualized in figure
 137 10

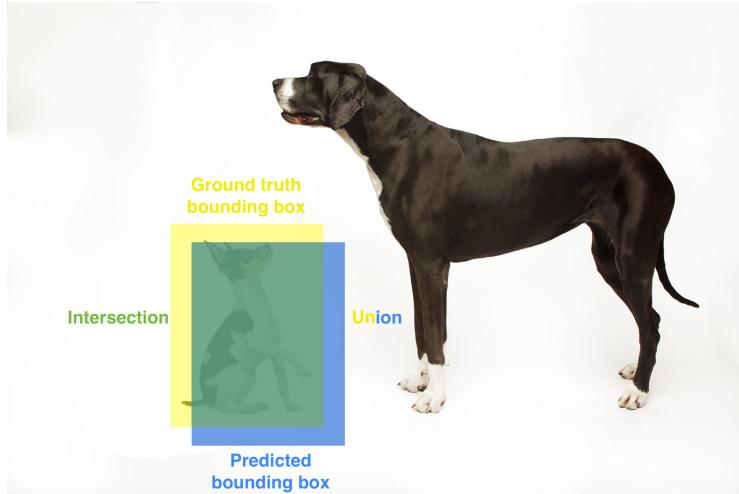


Figure 10: An example of the intersection over union metric. The ground truth bounding box is displayed in yellow, the predicted bounding box is shown in blue. The intersection between the two boxes is shown in green and the union is the area of the yellow, the blue and the green box

138 The IoU metric can then be used to determine whether a predicted bounding box is correct using a
 139 threshold $\tau \in [0, 1]$. A predicted bounding box is then considered correct if $\text{IoU} \geq \tau$. A common
 140 value for τ is 0.5, which means that at least 50% of the pixels of both the predicted bounding box
 141 and the ground truth bounding box have to overlap.

142 **1.2.2 Defining the Loss**

143 The localization of the bounding boxes can not be evaluated using a softmax activation on the
 144 fully-connected layers and a cross-entropy loss, as defined in section 1.1.3. We thus need to use
 145 an alternative loss function for the prediction of the bounding box. The fully-connected layer will
 146 provide logit scores for the centre, width, and height of the bounding box. These logit scores are
 147 however again unconstrained and not between 0 and 1. The values that represent the bounding box
 148 are therefore mapped between 0 and 1 using a sigmoid function to mitigate this. The sigmoid function
 149 is defined as follows.

$$S(z_i) = \frac{1}{1 + e^{-z_i}} \quad (4)$$

150 Where z is a vector with the logit scores for b_x, b_y, b_w and b_h . $S(z_i)$ is the value for those logit scores
 151 normalized to one, where i indexes each logit score. Figure 11 illustrates how a sigmoid layer is
 152 applied on our example network.

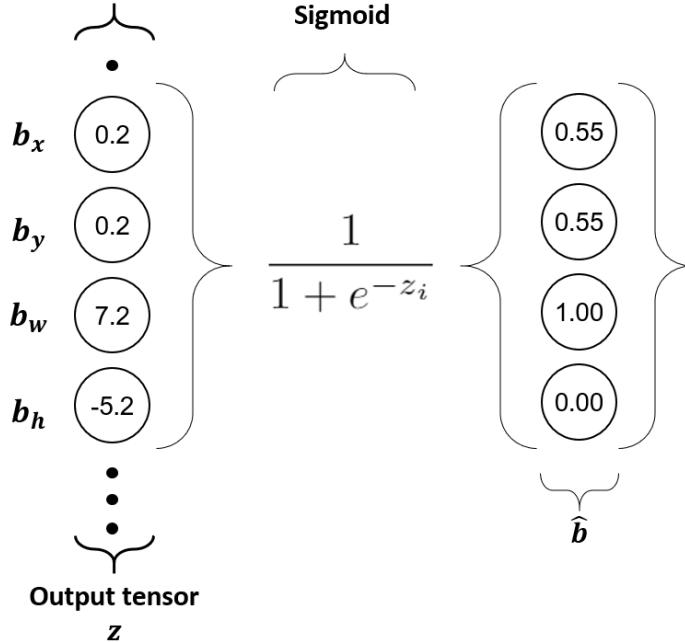


Figure 11: An example of how a sigmoid layer is used to transform b_x , b_y , b_w and b_h to remain within image bounds. Logit scores (left) are mapped to values 'squashed' between (0,1).

153 After applying the sigmoid function, we end up with the networks estimated values for $\hat{b}_x, \hat{b}_y, \hat{b}_w, \hat{b}_h$,
 154 which we will call the $\hat{\mathbf{b}}$ vector. Using $\hat{\mathbf{b}}$ and the ground truth vector \mathbf{b} we can now calculate the
 155 difference between those vectors. Some conventional functions used to accomplish this are the L1
 156 Loss, L2 loss, or the smooth L1 loss, defined in equations 5a, 5b and 5c respectively.

$$\mathcal{L}_{\text{L1}}(\hat{\mathbf{b}}, \mathbf{b}) = \sum_{i=1}^4 |b_i - \hat{b}_i| \quad (5a)$$

$$\mathcal{L}_{\text{L2}}(\hat{\mathbf{b}}, \mathbf{b}) = \sum_{i=1}^4 (b_i - \hat{b}_i)^2 \quad (5b)$$

$$\mathcal{L}_{1,\text{smooth}} = \begin{cases} |b_i - \hat{b}_i| & \text{if } |b_i - \hat{b}_i| > 1 \\ (b_i - \hat{b}_i)^2 & \text{if } |b_i - \hat{b}_i| \leq 1 \end{cases} \quad (5c)$$

159 A much more common measure for single-stage detectors however is to use the earlier explained IoU
 160 as a loss function. The IoU is scale invariant, unlike L1 and L2 loss. We can use network prediction
 161 $\hat{\mathbf{b}}$ and the ground truth bounding box \mathbf{b} to calculate the IoU, and consequently calculate the loss as
 162 follows:

$$\mathcal{L}_{\text{IoU}} = 1 - \text{IoU} \quad (6)$$

163 To calculate the total loss of the network, we calculate the cross-entropy loss for P_c and our class
 164 scores C_i and add them to the IoU loss as follows:

$$\mathcal{L}_{\text{Total}} = \mathcal{L}_{\text{IoU}} + \mathcal{L}_{\text{CE}, P_c} + \sum_{i=1}^k \mathcal{L}_{\text{CE}, C_i} \quad (7)$$

165 1.2.3 Overall performance

166 The performance can now be generalized to the localization of each object in image by introducing
 167 two more metrics, as follows:

$$\text{Precision} = \frac{\text{Number of correctly predicted bounding boxes}}{\text{Number of ground truth bounding boxes}} \quad (8a)$$

168

$$\text{Recall} = \frac{\text{Number of correctly predicted bounding boxes}}{\text{Number of predicted bounding boxes}} \quad (8b)$$

169 Where **precision** measures the fraction of correctly predicted bounding boxes relative to the number
 170 of ground truth boxes. **Recall** measures the fraction of correctly predicted bounding boxes relative to
 171 the total number of predicted bounding boxes. A smaller number of ground truth bounding boxes
 172 will increase the precision. A smaller number of predicted bounding boxes by the object detector will
 173 increase the recall. This creates a trade-off between the precision and the recall of an object detector.
 174 This is why these two metrics are taken together and plotted against each other, like in figure 12.

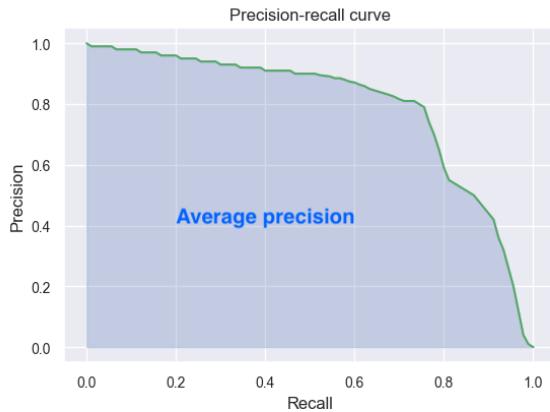


Figure 12: An example of a possible precision-recall curve. The precision is on the y-axis and the recall is on the x-axis. The precision-recall line is the green line in the figure and the transparent blue area under the green curve is the average precision

175 The area under the curve in figure 12 is called the **Average Precision** (AP) and is often used as a
 176 metric for object detectors. It takes the trade-off between the precision and recall into account. The
 177 AP only measures the performance of one class of the object detector. A version of the AP that is
 178 often used is the **mean Average Precision** (mAP), which is the average AP for all the classes the
 179 object detector can detect.

180 1.3 Sliding window detection

181 The original concept of classification and localization can only detect one object in an image, even if
 182 there are multiple objects present in that image. One way of trying to detect multiple objects in one
 183 image is to use a sliding window detector.

184 A sliding window detector uses the idea that a window with a certain shape and size can be slid over
 185 an image to detect objects in that image. Let's start simple, we have decided we want to detect the
 186 small dog in figure 13a. In this example, we again assume that we have a classifier $f(x)$ that can
 187 detect a dog in an image. We define a **sliding window** with a width w and a height h and place it
 188 in the top left corner of the image. The size of the window is generally *fixed*. The window 'slides'
 189 over the image to the right, with a step size dw until it hits the border on the right of the image. The
 190 window is then moved down with a certain step size dhs . The window is then slid horizontally again
 191 until it hits the border and this process is repeated until the window has been slid over the entire
 192 image. Each location of the window can be interpreted as a separate image that can be provided as
 193 input to the classifier. We get a predicted class and location for each window location and can thus
 194 classify and localize multiple objects in an image.

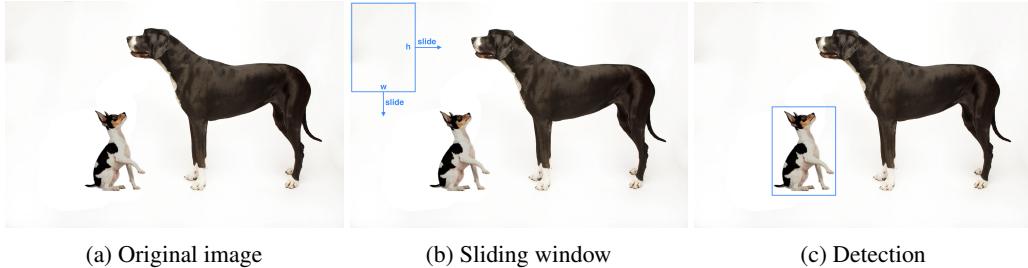


Figure 13: A sliding window detection example. (a) is the original image, (b) shows a sliding window and (c) shows a location of the sliding window that allows us to detect the smaller dog

195 There are a few disadvantages to sliding a fixed size window over an image, however. First, a small
 196 step size in the horizontal and vertical direction can lead to many different window locations. This
 197 is computationally expensive. Small step size can also lead to the classification and localization of
 198 one object across many different window locations. The step sizes dh and dw are therefore generally
 199 chosen to be the same as the width (w) and height (h) of the sliding window. This leads to a grid that
 200 is placed over the image, where each cell is a window, like in figure 14a. The size of the window in
 201 figure 14a makes it hard to detect the larger dog. If we increase the window size, however, like in
 202 figure 14b, we can detect the large dog, but it has now become impossible to detect more than 2 dogs
 203 because we only have 2 windows.

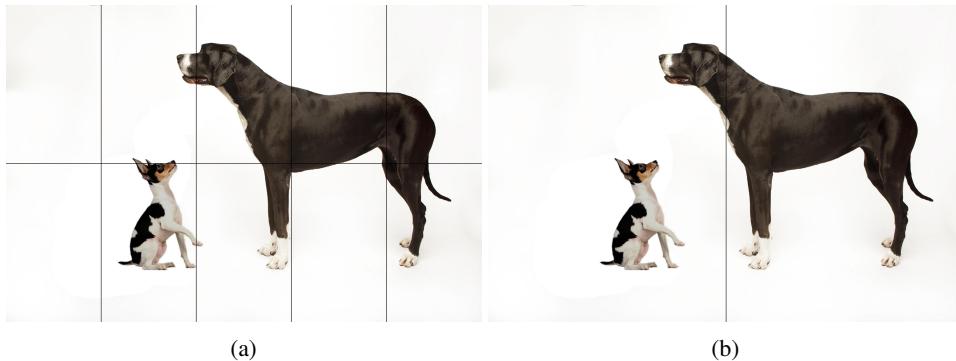


Figure 14: Windows as a grid, where the step size horizontally is equal to the width of the window. The vertical step size is equal to the height of the window. In image (a), a smaller window size is used, this may make it hard to detect the larger dog. In figure (b), a larger window size is used, but it is now impossible to detect more than 2 dogs in an image

204 More objects in the picture at different scales would mean that we need to place multiple window
 205 sizes over the image to detect each of the objects. We can keep the window size fixed and slide it
 206 over downsampled versions of the image to do that. This would mean, however, that we have to
 207 run the detection algorithm over multiple windows, over multiple scales of the image. There are
 208 some obvious downsides to this approach and the main disadvantages have been summarized in the
 209 following bullet points.

- 210 • Sliding a window over an image, even if multiple windows are slid over an image in parallel,
 211 is computationally inefficient;
- 212 • The size(s) of the windows relative to the image size and the number of different window
 213 sizes have to be determined using prior knowledge;
- 214 • Objects with different types of shapes can not be detected using the same window, for
 215 example a wide object or a tall object;
- 216 • A sliding window may not always line up well with a or multiple objects in an image, which
 217 may lead to that object not being detected.

218 **2 Single-Stage Object Detectors**

219 The previous section finished with a description of a sliding window detector, which is one of the
220 ways in which a classifier can detect multiple objects in an image. It turns out however that the
221 sliding window detector faces several challenges, amongst others that it is inefficient and might not be
222 feasible for all applications. One method to mitigate these challenges is to predict regions of interest
223 in the image and apply a classifier to this region. This is called a two-stage method. Stage one is
224 to predict the region of interest and stage two is to classify and localize objects within this region.
225 Two-stage methods are discussed in-depth in chapter 5.1.

226 Another method to mitigate these challenges is by making each feature map pixel responsible for the
227 classification and localization of a local region in the input image. The localization and classification
228 can be done using a single forward pass through the network, which is why it is called single-stage
229 object detection. Examples of single-stage detectors include YOLO, SSD [11], DSSD [5], RetinaNet
230 [10], M2Det [22], RefineDet [20]. Out of these detectors, YOLO is the most well known. We will
231 explain and explore this single-stage detector in the next section. A lot of the concepts that are true
232 for YOLO will generalize to the other single-stage detectors, but these other single-stage detectors are
233 outside of the scope of this chapter. We encourage you to take a look at the papers that are referenced.

234 **2.1 YOLO**

235 This subsection will elaborate on YOLO, or You Only Look Once, which was presented as a novel
236 approach to object detection by Redmon et al. [16]. Two official YOLO iterations have followed
237 since, namely YOLO9000 [14] and YOLOv3 [15]. These iterations were all created by the original
238 YOLO author, Joseph Redmon. YOLOv4 [1] is the most recent iteration and was not created by the
239 original author. Even though YOLOv4 offers a significant improvement over YOLOv3, the methods
240 that are used to improve upon YOLOv3 are outside the scope of this chapter. This is why we have
241 decided against discussing its improvements here.

242 YOLO is one of the fastest object detectors in terms of inference time. The inference time is the time
243 it takes to do one forward pass. It approaches object detection as a single regression problem, where
244 image pixels are translated to bounding box coordinates and class probabilities in one stage.

245 **2.1.1 Basic principle**

246 Single-stage detectors such as YOLO approach object detection based on feature maps. They map a
 247 2-dimensional grid over the input image corresponding to the spatial dimensions of the feature map.
 248 Take for example a 3-dimensional feature map with a width and height of 3 pixels, and a depth of
 249 1024 features. In this case, a 3 by 3 grid is mapped onto the input image. Each pixel in the feature
 250 map corresponds to a grid cell, which is a *rectangular area of pixels* in the input image. A visual
 251 understanding of this notion is given in figure 15. This mapping is direct, so the right upmost pixel in
 252 the feature map corresponds to the right upmost grid cell in the input image, etc.

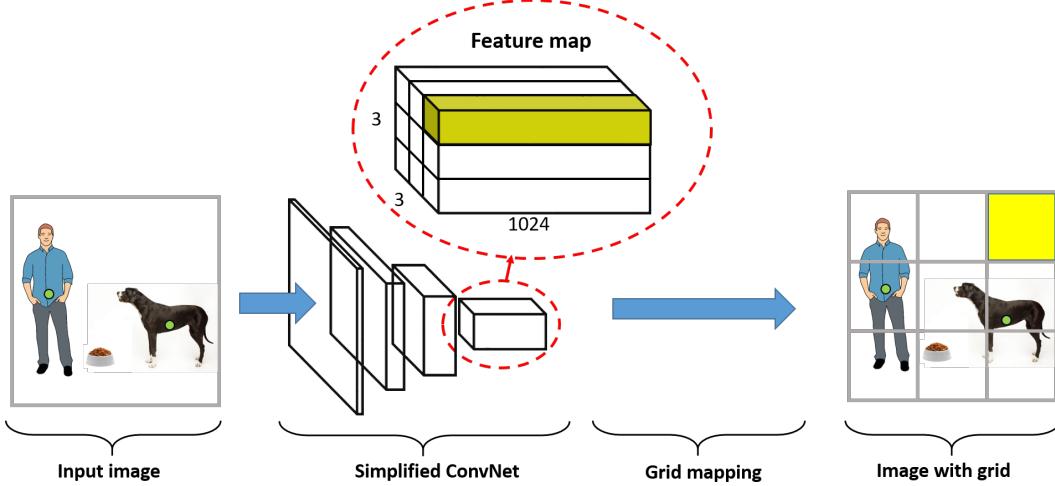


Figure 15: Input image is fed into a simplified ConvNet backbone which outputs a 3-dimensional feature map, which is $3 \times 3 \times 1024$ in this example. The spatial resolution of the feature map is mapped as a 3×3 grid onto the original input image. The yellow parts show the correspondence of one of the feature map pixels to its image grid cell

253 Each feature map pixel has a number of features attributed to it. These features are learnt from a
 254 rather large receptive field (section 1.1.2) of the input image. This receptive field exceeds the size of
 255 the grid cell in the input image that corresponds to this pixel. A feature map pixel is made *responsible*
 256 for the detection of objects in its corresponding grid cell. The features it uses can, however, contain
 257 information from a large part of the image to provide context. This is an important property because
 258 an object can span multiple grid cells. An example is shown in figure 16, where the human covers
 259 three grid cells.

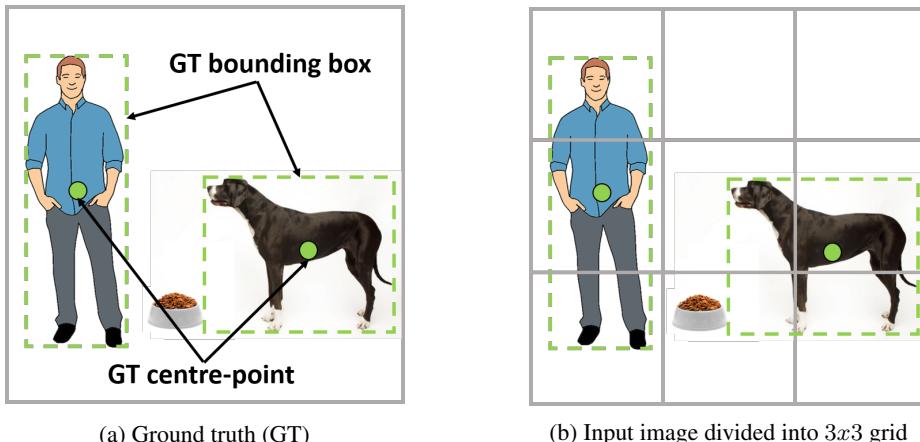
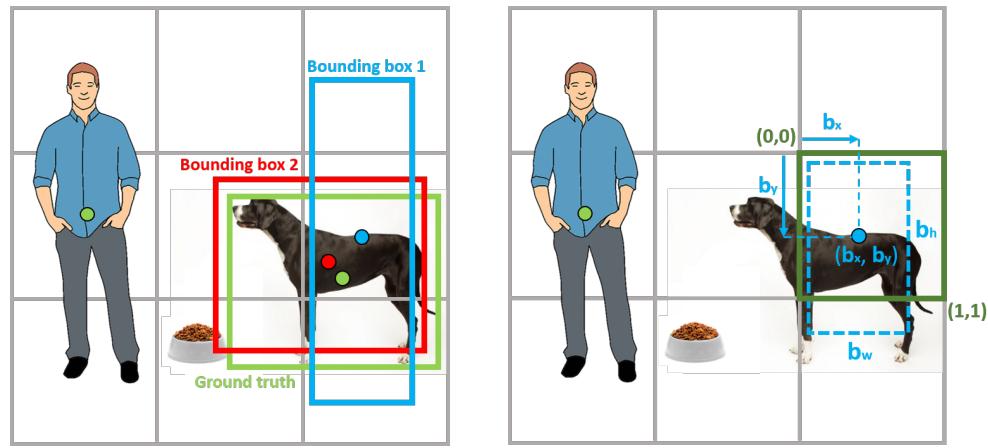


Figure 16: Input image with 2 identifiable objects including their ground truth (GT) is subdivided into an $S \times S$ grid. In this case the grid is 3×3

260 But how do we determine which feature map pixel is made responsible for an object that covers
 261 multiple grid cells? During training, YOLO deals with this by making the grid cell that contains the
 262 centre point of the object responsible for detecting that object. The grid cells on the input image
 263 are a direct mapping to and from the feature map pixels. The feature map pixel that corresponds
 264 to the grid cell with the centre of the object in it, will thus be responsible to detect that object. In
 265 figure 16b, the left-middle grid cell contains the centre of the human, the corresponding left-middle
 266 feature map pixel is thus made responsible for detecting this object. YOLO splits up the detection of
 267 objects into grid cells and assigns a responsibility to the feature map pixels to penalize the network
 268 when it detects a bounding box for an object in another cell. It also makes it easier for the network to
 269 converge on predictions because it splits the problem up into smaller parts. This allows the network
 270 to focus on object predictions for smaller sections of the image, which allows each part of the object
 271 detector (each feature map pixel) to focus on a small number of predictions. Figure 17a illustrates 2
 272 example bounding box predictions for one grid cell. The green bounding box represents the ground
 273 truth and the red and blue bounding boxes are predictions.



(a) ground truth bounding box and centre point (green) (b) Bounding box parameterization of centre point and predicted boxes and centre points (blue & red) dimensions

Figure 17: Bounding box predictions and parameterization

274 The final link in understanding the underlying principle of YOLO is the fully-connected layer that is
 275 used to map the information of the feature map to an output tensor. The **output tensor** is an extended
 276 version of the tensor introduced in section 1.2, and contains the class probabilities, object score and
 277 bounding box predictions for each grid cell in a vector format. Because the feature map is trained in
 278 conjunction with the fully-connected layer, each feature map pixel will learn which parts of the image
 279 are important as context for the cell it is responsible for. This allows YOLO to make predictions of
 280 multiple objects located around the entire image in one forward pass; hence the name You Only Look
 281 Once.

282 Summarizing the most important properties of YOLO:

- 283 1. The backbone network outputs a 3-dimensional feature map with a certain spatial and feature
 284 resolution, which captures important features in the input image;
- 285 2. A 2D grid is mapped onto the input image that is equivalent to the spatial resolution of the
 286 feature map;
- 287 3. During training, a feature map pixel that corresponds to a grid cell containing the ground
 288 truth centre coordinates of an object is made responsible for detecting that object;
- 289 4. Contextual information around a grid cell is also encoded in the corresponding feature map
 290 pixel due to the large receptive field the feature map has over the input image;
- 291 5. The feature map is trained in conjunction with the fully-connected layer and hence, each
 292 pixel in the feature map will learn which parts of the image are important as context for the
 293 cell it is responsible for.

294 **2.1.2 Output tensor**

295 For each feature map pixel, YOLO predicts an output label y . This label encodes information about
 296 the predicted bounding boxes and class corresponding to its grid cell. The length of the label mainly
 297 depends on the number of bounding boxes B predicted per pixel; each label has shape $B \cdot 5 + C$ and
 298 contains the following variables **per bounding box**; the object score P_c , the centre-point coordinates
 299 of the object, b_x and b_y , the width and height of the bounding box, b_w and b_h . Additionally, the
 300 output label contains a set of conditional class probabilities **C per grid cell**.

301 The meaning of these variables has been discussed previously in section 1.2. However, there are
 302 some important things to note; instead of the image spatial dimensions being normalized to [0, 1] as
 303 discussed in section 1.2.2, the spatial dimensions of the grid cells are all normalized to [0, 1]. This is
 304 illustrated in figure 17b. This normalization constrains the prediction of the centre of a bounding
 305 box to a certain grid cell. If we would use a sigmoid for the width and height, as discussed in section
 306 1.2.2, the width, and height of a bounding box would be constrained to the width and height of that
 307 grid cell as well (between 0 and 1). YOLO thus uses the exponential function $\exp z$ to map the logit
 308 scores for the width and height of a bounding box to [0, inf]. Where z corresponds to the logit scores
 309 for the width and height at the output of the fully-connected layer. This allows a bounding box to be
 310 larger than the grid cell that contains the centre of the object.

311 The final output tensor is a concatenation of all output labels and consequently has the shape
 312 $S \times S \times (B \cdot 5 + C)$, where $S \times S$ is the spatial resolution of the feature map, B the number of
 313 predicted bounding boxes per pixel and C the total number of classes. The notation of this output
 314 tensor differs between literature and is sometimes denoted as $S \times S \times B \times (5 + C)$. We will stick to
 315 the first notation. Figure 18 shows the output tensor for a YOLO algorithm that has a spatial feature
 316 map resolution of 3x3, predicts 2 bounding boxes per feature map pixel, and can predict a total of 3
 317 different classes.

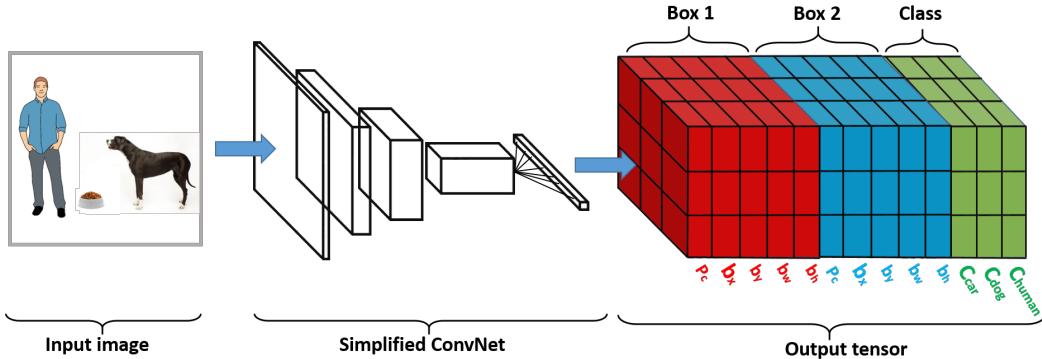


Figure 18: Simplified illustration of YOLO pipeline. Right: input image of a dog. Middle: simplified CNN. Left: YOLO output tensor of shape $S \times S \times (B \cdot 5 + C)$ corresponding to 3x3 gridcells S , 2 bounding boxes per cell B and 3 class labels C

318 **2.1.3 Non-Maximal Suppression**

319 Assume that we have an input image onto which a 3x3 grid is mapped. We will predict 2 bounding
 320 boxes per grid cell and thus per feature map pixel. This is illustrated in figure 19a, where the thicker
 321 bounding boxes reflect predictions with a higher confidence score. As you can see, multiple bounding
 322 boxes are trying to predict the same object. *During testing* (inference), we would like to end up
 323 with only one bounding box per object, which is illustrated in figure 19c. In order to get rid of the
 324 redundant bounding boxes we will use **Non-Maximal Suppression (NMS)**. The 3 steps for NMS
 325 are provided below.

- 326 1. Discard all bounding boxes with a confidence score below a certain threshold ϵ ; $P_c \leq \epsilon$;
- 327 2. Pick the bounding box with the largest confidence score and use that as predicted bounding
 box B_{pred} ;
- 328 3. Remove any remaining bounding box B_{remain} that has an IoU with the previously predicted
 bounding box (B_{pred}) that is above a certain threshold τ ; $IoU(B_{pred}, B_{remain}) \geq \tau$.

331 Step 2 and 3 are performed iteratively until there are no more remaining bounding boxes. The
 332 thresholds ϵ and τ can be set to any desired value.

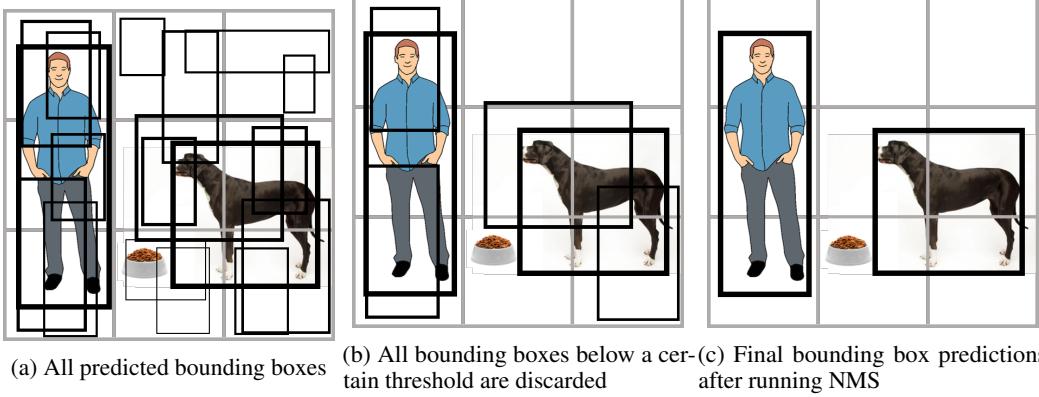


Figure 19: Removing redundant bounding box predictions using Non-Maximal Suppression

333 2.2 Limitations of YOLO

334 The ideas described in the previous subsections were used in YOLO, but there are a few downsides
 335 to this method. This motivated the next iterations YOLO9000 [14] and YOLOv3 [15]. The most
 336 important limitations of YOLO are the following.

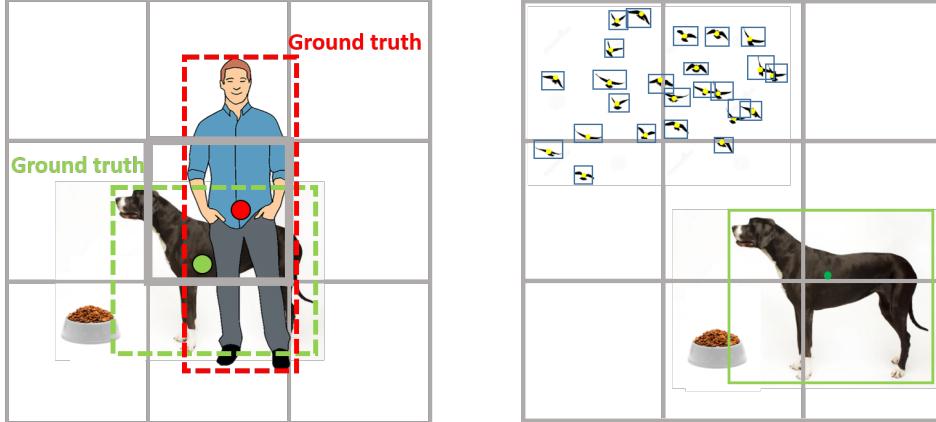
- 337 • **Each grid cell can only be assigned one class**, which is unrealistic because a grid cell can
 338 contain multiple different types of objects. The problem is visualized in figure 20a, where
 339 the centre of the human is in the same grid cell as the centre of the dog. The problem is
 340 caused by the output tensor, because the output tensor only encodes the classes on a grid
 341 cell level, instead of encoding the classes for each bounding box.;
- 342 • **Generalization to unusually shaped objects** is also a problem for the original YOLO
 343 network. The network has a hard time converging to bounding boxes with extreme aspect
 344 ratios, probably because the large deformation can take a long time to learn. Most bounding
 345 boxes will have a certain shape, so the deformations from the common bounding boxes to
 346 the ones needed for the outliers could be hard to learn;
- 347 • **Detection of small-scale objects** is a problem that is visualized in figure 20b. An important
 348 reason for this is that the feature map for YOLO is not fine enough (7×7). Therefore, each
 349 feature map pixel is responsible for a relatively large area in the input image, increasing
 350 difficulty in detecting small-scale objects.

351 2.3 YOLO9000 and YOLOv3

352 The limitations of the original YOLO have been tackled by more recent YOLO iterations. We will
 353 discuss two of the most prominent improvements introduced in YOLO9000 and YOLOv3, namely
 354 anchor boxes and different feature map scales.

355 2.3.1 Anchor boxes

356 In YOLO9000 [14] the concept of **anchor boxes**, or **priors**, was introduced. This improved per-
 357 formance with regards to generalization to unusually shaped objects and introduced multiple class
 358 predictions per grid cell. Anchor boxes are a set of predefined bounding boxes with a certain height
 359 and width. The shapes of the anchor boxes are defined to capture scale and aspect ratios of specific
 360 object classes within a dataset, e.g. a human is generally a tall shaped object and a car is generally a
 361 widely shaped object. We can incorporate this knowledge into YOLO by predefining anchor boxes
 362 corresponding to common shapes in the dataset. This makes it easier for the network to learn better
 363 predictions, as each anchor box specializes in detecting objects of a specific size and aspect ratio.
 364 Each anchor box corresponds to a certain class label, so the network can now learn multiple classes
 365 in a grid cell. Determining suitable shapes for the anchor boxes can be done by hand, but a more



(a) Centre-points of multiple objects of a different class share the same grid cell.
 (b) Many small objects, or birds, sharing the same grid cell

Figure 20: Removing redundant bounding box predictions using Non-Maximal Suppression

366 clever approach is preferred. YOLO9000 and YOLOv3 use k-means clustering on the ground truth
 367 bounding boxes in the training set to determine a good set of anchor boxes. The number of means
 368 that are optimal to use for YOLO9000 (5) is different from YOLOv3 (9). The clustering can be done
 369 before training the network.

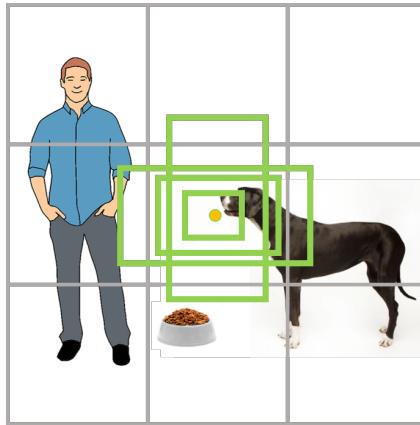


Figure 21: 4 predefined anchor boxes with different aspect ratios. Although not illustrated these anchor boxes are present in each of the grid cells

370 With anchor boxes, YOLO no longer predicts the centre coordinates of a bounding box or its height
 371 and width, but rather the offset from the predefined anchor box location and size to the predicted
 372 bounding box location and size. The anchor box is thus used as a prior for the shape of the object
 373 classes the object detector is trying to predict. Each anchor box can be assigned one object of any
 374 class. That means that each grid cell can detect as many different objects as there are anchor boxes.
 375 Essentially, the classification level has moved from the grid cell level to the bounding box level.

376 2.3.2 Different feature map scales

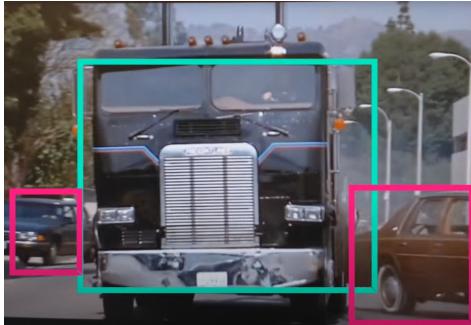
377 To overcome the problem of not being able to detect small-scale objects, YOLOv3 uses intermediate
 378 outputs of the backbone network as additional feature maps. Specifically, it uses feature maps
 379 produced at three locations in the network, which have a higher resolution. Because the intermediate
 380 feature maps are coming from lower layers, the receptive fields are smaller. To have meaningful
 381 intermediate feature maps, the number of layers in YOLOv3 is increased from 19 layers in YOLO9000
 382 to 53 layers in YOLOv3. The 3 feature maps with different sizes are each used to map a grid onto the
 383 input image. This results in 3 different grids; one very fine grid, a coarse grid, and an intermediately

384 coarse grid. The network uses 3 anchor boxes per grid cell to make predictions in each of the 3 grids.
 385 YOLO9000 does not use multiple feature maps and only uses 5 anchor boxes per grid cell. This
 386 results in an increase in inference time for YOLOv3, because of the number of anchor boxes (9 vs 5)
 387 it has to predict and the increase in network size (53 vs 19 layers).

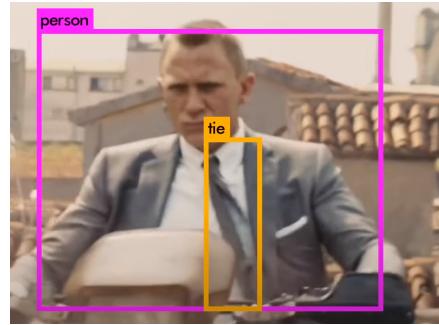
388 2.4 Real world usage

389 So what does YOLO look like when used on a real-world video? Along with the paper, the authors
 390 released a demonstrative showcase of YOLO on several videos and well-known movies. Additionally,
 391 we can use a very recent implementation of YOLOv4 on the Deepdrive PL Dataset by Karol Majek
 392 [12]. All videos are linked here on the PDF version: YOLO, YOLO9000, YOLOv3, YOLOv4.

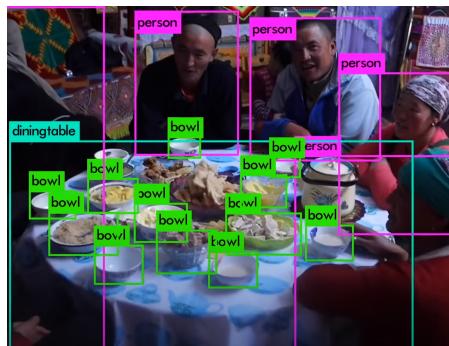
393 Figure 22 showcases selected frames of the YOLO, YOLO9000, YOLOv3 and YOLOv4 demon-
 394 stration videos. There is a clear increase in performance with each iteration producing more accurate
 detection and being able to capture more detail in the image.



(a) YOLO: The Algorithm is able to detect multiple separate vehicles



(b) YOLO9000: The algorithm is able to detect attributes that a person is wearing



(c) YOLOv3: the algorithm is able to detect several varying, complicated objects on top of each other



(d) YOLOv4: The algorithm is able to detect vehicles that are far away and even people present in the vehicle in front

Figure 22: Showcase of the current YOLO implementations

395

396 2.4.1 Current applications

397 The advantage of YOLO over two-stage systems is that it looks at the whole image at test time,
 398 ensuring that its predictions are informed by a global context. Furthermore, as it only uses just one
 399 network evaluation to make predictions it is much faster (1000x faster compared to RCNN and 100x
 400 Faster compared to Fast R-CNN [15]). It is therefore perfectly suited for applications that require
 401 speed such as autonomous cars [2]. The size of the network and the low inference time also makes
 402 it possible to translate a YOLO version to an embedded environment, like YOLO NANO [19] and
 403 YOLO-LITE [6].

404 **2.4.2 Controversy**

405 YOLO-based algorithms are known to be actively used in military applications, such as autonomous
406 **unmanned aerial vehicles** (UAV's) [17] and **surveillance tracking** by governments [21]. As a
407 result, it even caused the original creator of the algorithm, Joe Redmon, to stop working in the
408 computer vision field, because of his concerns regarding the usage of his work. As stated by him in a
409 Twitter post: "I stopped doing CV research because I saw the impact my work was having. I loved
410 the work but the military applications and privacy concerns eventually became impossible to ignore."

411 **Reference image**

412 The image that is used as a reference image throughout this chapter [4] was taken from Psychology
413 Today. The official source is SC Psychological Enterprises Ltd.

414 **References**

- 415 [1] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao. Yolov4: Optimal speed and accuracy of object
416 detection. *arXiv preprint arXiv:2004.10934*, 2020.
- 417 [2] J. Choi, D. Chun, H. Kim, and H.-J. Lee. Gaussian yolov3: An accurate and fast object detector
418 using localization uncertainty for autonomous driving. In *Proceedings of the IEEE International
419 Conference on Computer Vision*, pages 502–511, 2019.
- 420 [3] D. C. Ciresan, U. Meier, J. Masci, L. M. Gambardella, and J. Schmidhuber. Flexible, high perfor-
421 mance convolutional neural networks for image classification. In *Twenty-Second International
422 Joint Conference on Artificial Intelligence*, 2011.
- 423 [4] S. Coren. Behavior differences between smaller and larger dogs, Jul 2015.
424 URL [https://www.psychologytoday.com/us/blog/canine-corner/201507/
425 behavior-differences-between-smaller-and-larger-dogs](https://www.psychologytoday.com/us/blog/canine-corner/201507/behavior-differences-between-smaller-and-larger-dogs).
- 426 [5] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg. Dssd: Deconvolutional single shot
427 detector. *arXiv preprint arXiv:1701.06659*, 2017.
- 428 [6] R. Huang, J. Pedoeem, and C. Chen. Yolo-lite: a real-time object detection algorithm optimized
429 for non-gpu computers. In *2018 IEEE International Conference on Big Data (Big Data)*, pages
430 2503–2510. IEEE, 2018.
- 431 [7] A. Karpathy. Cs231n: Convolutional neural networks for visual recognition.
432 <https://cs231n.github.io/classification/>, 2016.
- 433 [8] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional
434 neural networks. In *Advances in neural information processing systems*, pages 1097–1105,
435 2012.
- 436 [9] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document
437 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 438 [10] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In
439 *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- 440 [11] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot
441 multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- 442 [12] K. Majek. Deep drive pl warsaw. <https://archive.org/details/0002201705192>, 2017.
- 443 [13] W. Rawat and Z. Wang. Deep convolutional neural networks for image classification: A
444 comprehensive review. *Neural computation*, 29(9):2352–2449, 2017.
- 445 [14] J. Redmon and A. Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE
446 conference on computer vision and pattern recognition*, pages 7263–7271, 2017.

- 447 [15] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *arXiv preprint*
448 *arXiv:1804.02767*, 2018.
- 449 [16] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time
450 object detection. In *Proceedings of the IEEE conference on computer vision and pattern*
451 *recognition*, pages 779–788, 2016.
- 452 [17] Á. Rocha and T. Guarda. *Developments and Advances in Defense and Security: Proceedings of*
453 *the Multidisciplinary International Conference of Research Applied to Defense and Security*
454 (*MICRADS 2018*), volume 94. Springer, 2018.
- 455 [18] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy,
456 A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition
457 Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi:
458 10.1007/s11263-015-0816-y.
- 459 [19] A. Wong, M. Famuori, M. J. Shafiee, F. Li, B. Chwyl, and J. Chung. Yolo nano: a highly
460 compact you only look once convolutional neural network for object detection. *arXiv preprint*
461 *arXiv:1910.01271*, 2019.
- 462 [20] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li. Single-shot refinement neural network for object
463 detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
464 pages 4203–4212, 2018.
- 465 [21] X. Zhang, X. Hao, S. Liu, J. Wang, J. Xu, and J. Hu. Multi-target tracking of surveillance
466 video with differential YOLO and DeepSort. In J.-N. Hwang and X. Jiang, editors, *Eleventh*
467 *International Conference on Digital Image Processing (ICDIP 2019)*, volume 11179, pages 701
468 – 710. International Society for Optics and Photonics, SPIE, 2019. doi: 10.1117/12.2540269.
469 URL <https://doi.org/10.1117/12.2540269>.
- 470 [22] Q. Zhao, T. Sheng, Y. Wang, Z. Tang, Y. Chen, L. Cai, and H. Ling. M2det: A single-shot object
471 detector based on multi-level feature pyramid network. In *Proceedings of the AAAI Conference*
472 *on Artificial Intelligence*, volume 33, pages 9259–9266, 2019.