

EE 628
Deep Learning
Spring 2020

Lecture 12
04/16/2020

Sergul Aydore

*Applied Scientist
Amazon Web Services*

Announcements

- Deadline for Projects **04/27/2020 Monday at 5pm ET**
 - project presentation on **04/30/2020** (40 %)
 - Details on grading: <https://github.com/sergulaydore/EE-628-Spring-2020>

Overview

- Last lecture we covered
 - Natural Language Processing: Pretraining & Applications
- Today, we will cover
 - Advanced Topics in Computer Vision
- Source material:
 - Dive into Deep Learning (<https://d2l.ai/>)
 - <https://github.com/sergulaydore/EE-628-Spring-2020>

Computer vision

- Here, we will introduce **image augmentation** and **fine tuning** methods
- Then, we will explore **various methods of object detection**.
- After that, we will learn how to use fully convolutional networks to perform **semantic segmentation** on images.
- Finally, we will introduce **Generative Adversarial Networks (GANs)**

Image Augmentation

- We mentioned that large-scale data sets are prerequisites for the successful application of deep neural networks
- Image augmentation technology expands the scale of training data sets by making a series of random changes to the training images to produce similar, but different, training examples.
- Another way to explain image augmentation is that randomly changing training examples can reduce a model's dependence on certain properties, thereby improving its capability for generalization.
- It can be said that image augmentation technology contributed greatly to the success of AlexNet.

Common Image Augmentation Methods

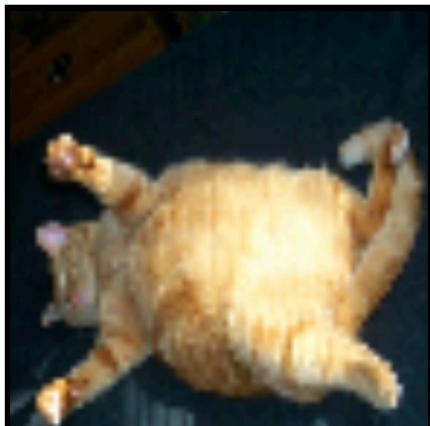
Flip



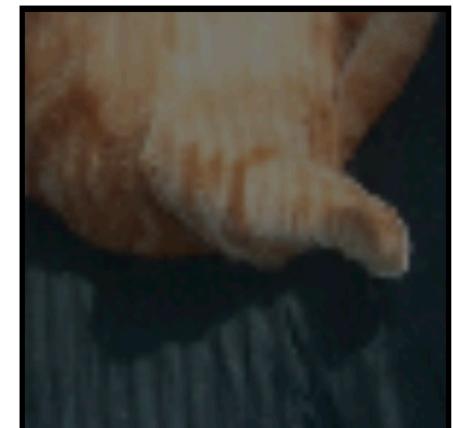
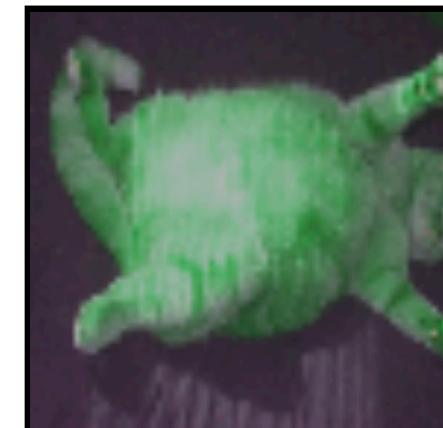
Change Color



Crop



Overlaying multiple methods



Fine Tuning

- Assume that we want to identify different kinds of chairs in images.
- One possible method is to first find a hundred common chairs, take one thousand different images with different angles for each chair, and then train a classification model on the collected image data set.
- Although this data set may be larger than Fashion-MNIST, the number of examples is still less than ImageNet.
- This may result in the overfitting of the complicated model applicable to ImageNet on this data set.
- In order to deal with the above problems, an obvious solution is to collect more data.
- However, collecting and labeling data can consume a lot of time and money.

Fine Tuning

- Another solution is to apply transfer learning to migrate the knowledge learned from the source data set to the target data set.
- For example, although the images in ImageNet are mostly unrelated to chairs, models trained on this data set can extract more general image features that can help identify edges, textures, shapes, and object composition.
- These similar features may be equally effective for recognizing a chair.
- Now, we introduce a common technique in **transfer learning**: **fine tuning**.

Fine Tuning

- Fine tuning consists of the following four steps:
 1. Pre-train a neural network model, i.e., the source model, on a source data set (e.g., the ImageNet dataset).
 2. Create a new neural network model, i.e., the target model. This replicates all model designs and their parameters on the source model, except the output layer.
 3. Add an output layer whose output size is the number of target data set categories to the target model, and randomly initialize the model parameters of this layer.
 4. Train the target model on a target data set, such as a chair data set.

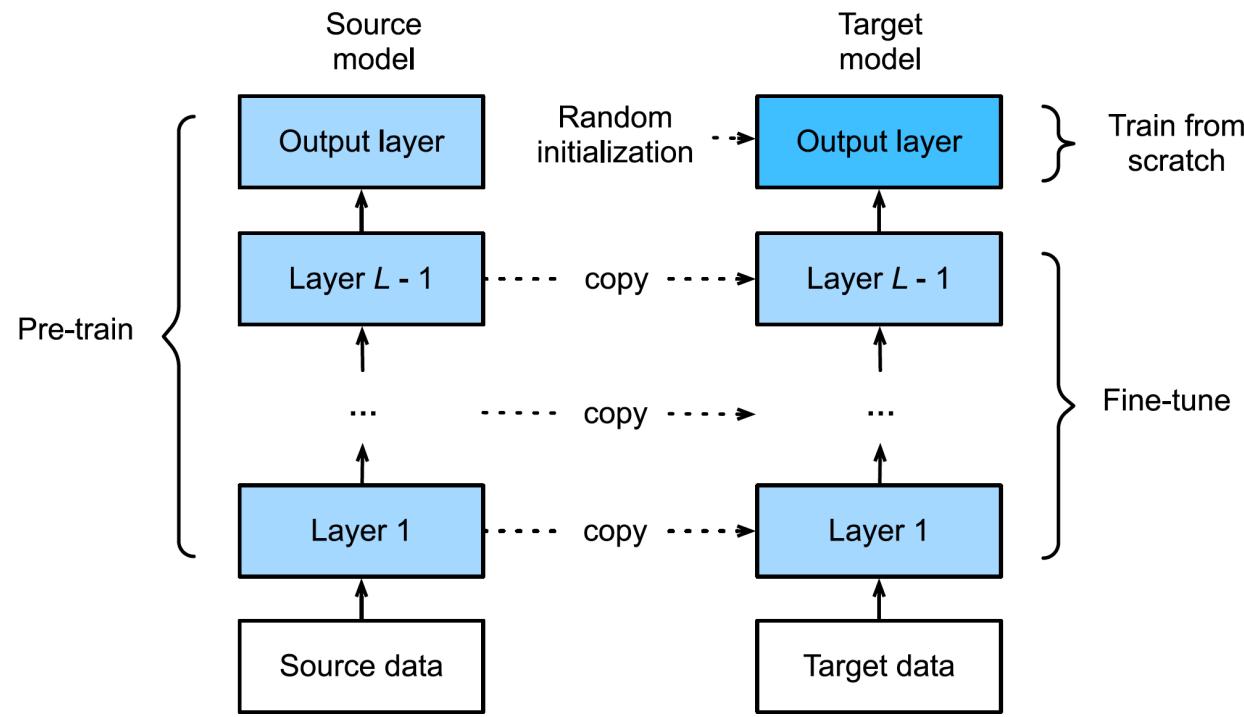


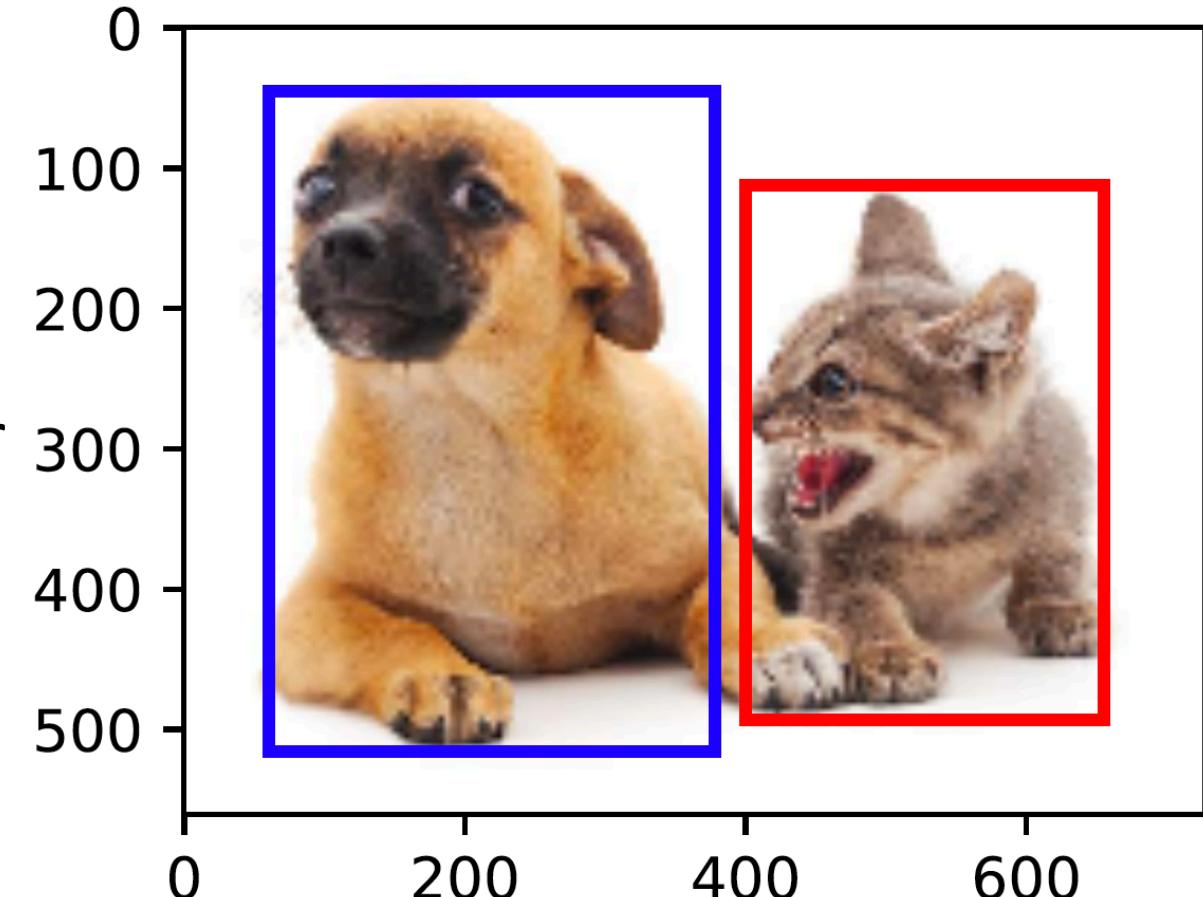
Fig. 14.2.1: Fine tuning.

Object Detection and Bounding Boxes

- In image classification tasks, we assume that there is only one main target in the image and we only focus on how to identify the target category.
- However, in many situations, there are multiple targets in the image that we are interested in.
- We not only want to classify them, but also want to obtain their specific positions in the image.
- In computer vision, we refer to such tasks as object detection (or object recognition).
- Object detection is widely used in many fields
 - in self-driving technology, we need to plan routes by identifying the locations of vehicles, pedestrians, roads, and obstacles in the captured video image.
 - Systems in the security field need to detect abnormal targets, such as intruders or bombs.

Bounding Boxes

- In object detection, we usually use a bounding box to describe the target location.
- The bounding box is a rectangular box that can be determined by the x and y axis coordinates in the upper-left corner and the x and y axis coordinates in the lower-right corner of the rectangle.
- The origin of the coordinates in the above image is the upper left corner of the image, and to the right and down are the positive directions of the x axis and the y axis, respectively.



bbox is the abbreviation for bounding box

```
dog_bbox, cat_bbox = [60, 45, 378, 516], [400, 112, 655, 493]
```

Anchor Boxes

- Object detection algorithms usually sample a large number of regions in the input image, determine whether these regions contain objects of interest, and adjust the edges of the regions so as to predict the ground-truth bounding box of the target more accurately.
- Different models may use different region sampling methods.
- Here, we introduce one such method: it generates multiple bounding boxes with different sizes and aspect ratios while centering on each pixel.
- These bounding boxes are called **anchor boxes**.

Generate Multiple Anchor Boxes

- Assume the input image has a height of h and width of w .
- We generate anchor boxes with different shapes centered on each pixel of the image.
- Assume the size is $s \in (0, 1]$, the aspect ratio is $r > 0$, and the width and height of the anchor box are ws/\sqrt{r} and hs/\sqrt{r} , respectively
- When the center position is given, an anchor box with known width and height is determined.
- Let's set a set of sizes s_1, \dots, s_n and a set of aspect ratios r_1, \dots, r_m . If we use a combination of all sizes and aspect ratios with each pixel as the center, the input image will have a total of $whmn$ anchor boxes.
- Although these anchor boxes may cover all ground-truth bounding boxes, the computational complexity is often excessive.

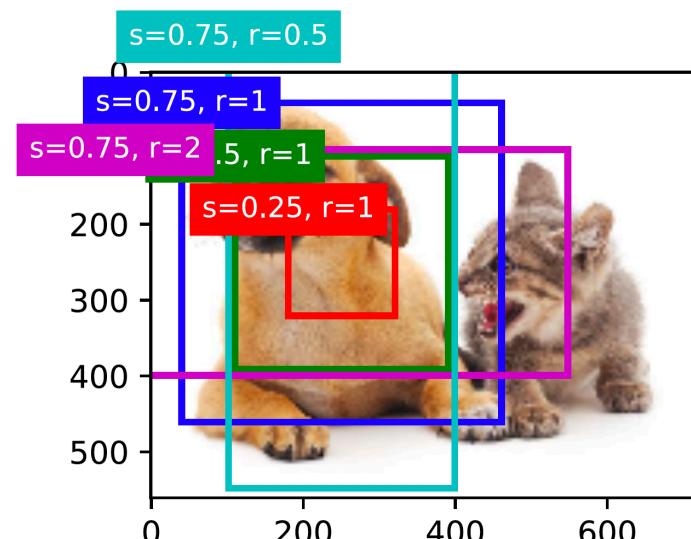
Generate Multiple Anchor Boxes

- Therefore, we are usually only interested in a combination containing s_1 or r_1 sizes and aspect ratios, that is:

$$(s_1, r_1), (s_1, r_2), \dots, (s_1, r_m), (s_2, r_1), (s_3, r_1), \dots, (s_n, r_1).$$

- That is, the number of anchor boxes centered on the same pixel is $n + m - 1$. For the entire input image, we will generate a total of $wh(n + m - 1)$ anchor boxes.

```
Y = contrib.nd.MultiBoxPrior(X, sizes=[0.75, 0.5, 0.25], ratios=[1, 2, 0.5])
```

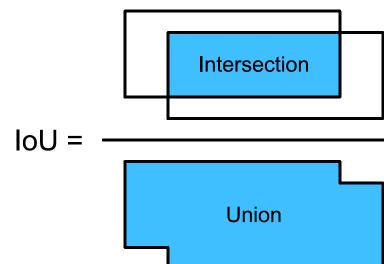


Intersection over Union

- How to measure how well the anchor box covers the target?
- We measure the similarity between anchor boxes and the ground-truth bounding box.
- Given set A and B , their Jaccard index is the size of their intersection divided by the size of their union

$$J(\mathcal{A}, \mathcal{B}) = \frac{|\mathcal{A} \cap \mathcal{B}|}{|\mathcal{A} \cup \mathcal{B}|}.$$

- In fact, we can consider the pixel area of a bounding box as a collection of pixels.
- When we measure the similarity of two bounding boxes, we usually refer the Jaccard index as intersection over union (IoU)



Labeling Training Set Anchor Boxes

- In the training set, we consider each anchor box as a training example.
- In order to train the object detection model, we need to mark two types of labels for each anchor box:
 - first, the category of the target contained in the anchor box (category) and,
 - second, the offset of the ground-truth bounding box relative to the anchor box (offset).
- In the training set, each image is labelled with the location of the ground truth bounding box and the category of the target contained.
- After the anchor boxes are generated, we primarily label anchor boxes based on the location and category information of the ground-truth bounding boxes similar to the anchor boxes.
- So how do we assign ground-truth bounding boxes to anchor boxes similar to them?

Labeling Training Set Anchor Boxes

- Assume the anchor boxes in the image are A_1, \dots, A_{n_a}
- And the ground truth bounding boxes are B_1, \dots, B_{n_b}
- And $n_a \geq n_b$
- Define a matrix $\mathbf{X} \in R^{n_a \times n_b}$, where element x_{ij} is the IoU of the anchor box A_i to the ground-truth bounding box B_j

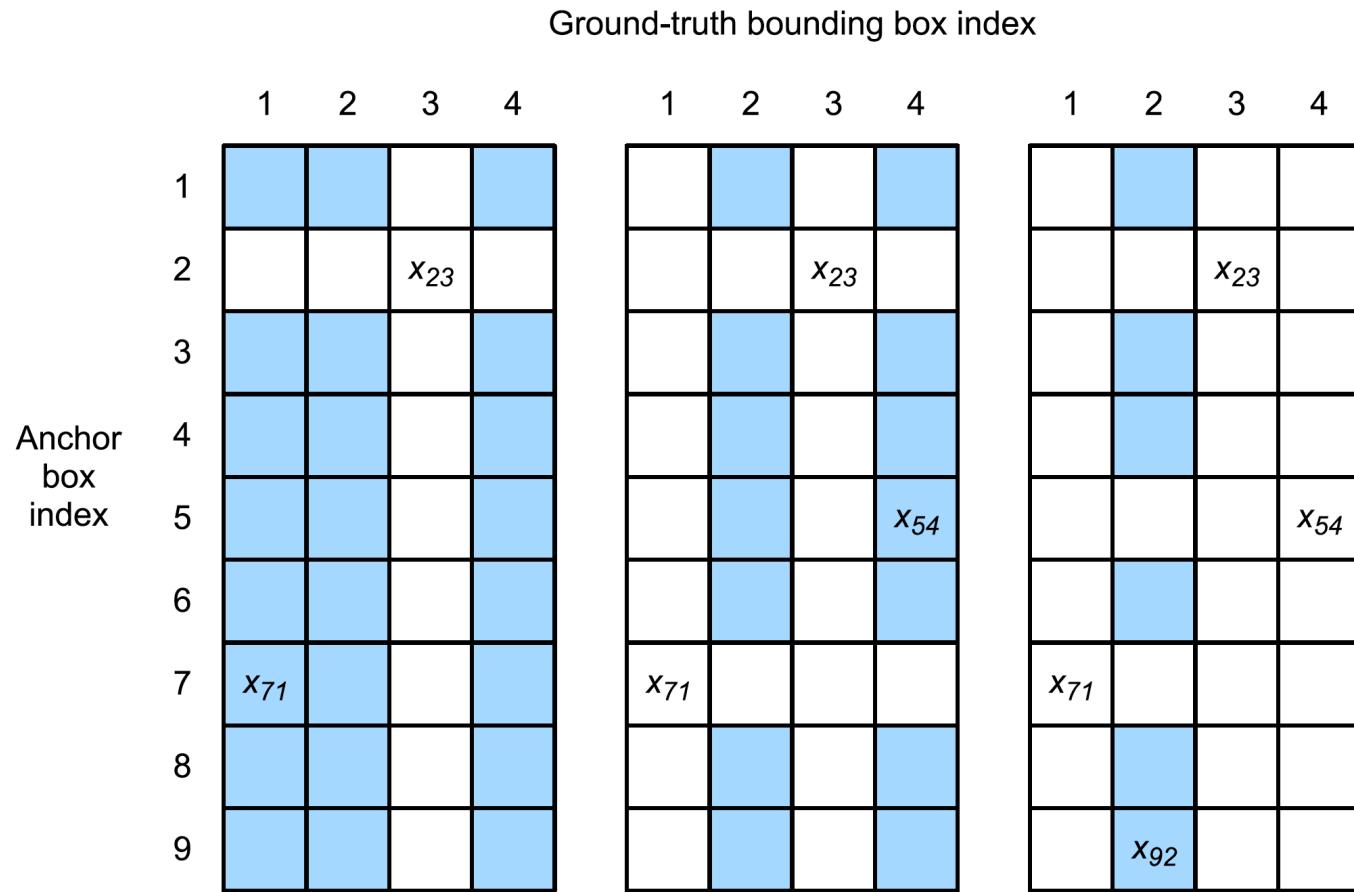
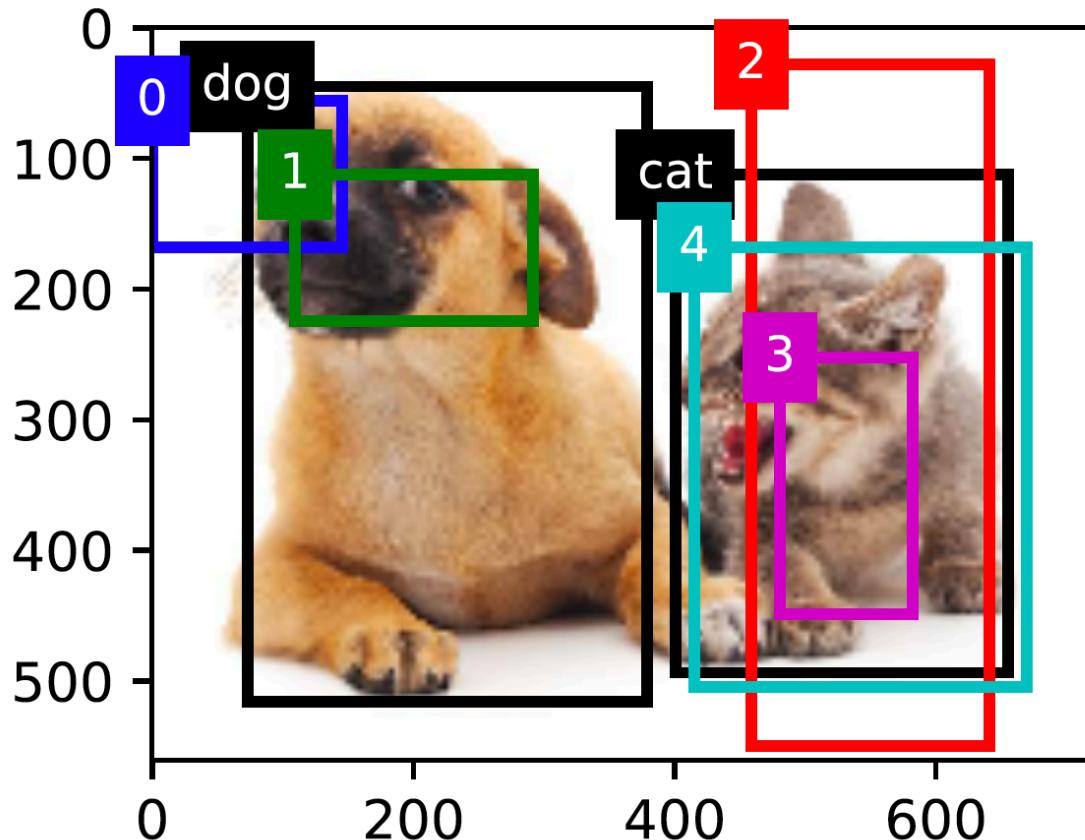


Fig. 14.4.2: Assign ground-truth bounding boxes to anchor boxes.

Labeling Training Set Anchor Boxes

1: dog
2: cat



```
[[0. 1. 2. 0. 2.]]
```

Labeling Training Set Anchor Boxes

- The offset of the anchor box A is set according to the relative position of the central coordinates of B and A and the relative sizes of the two boxes.
- Because the positions and sizes of various boxes in the data set may vary, these relative positions and relative sizes usually require some special transformations to make the offset distribution more uniform and easier to fit.
- Assume the center coordinates of anchor box A and its assigned ground-truth bounding box B are $(x_a, y_a), (x_b, y_b)$ the widths of A and B are w_a, w_b , and their heights are h_a, h_b , respectively.
- A common technique is to label the offset of A as

$$\left(\frac{\frac{x_b - x_a}{w_a} - \mu_x}{\sigma_x}, \frac{\frac{y_b - y_a}{h_a} - \mu_y}{\sigma_y}, \frac{\log \frac{w_b}{w_a} - \mu_w}{\sigma_w}, \frac{\log \frac{h_b}{h_a} - \mu_h}{\sigma_h} \right)$$

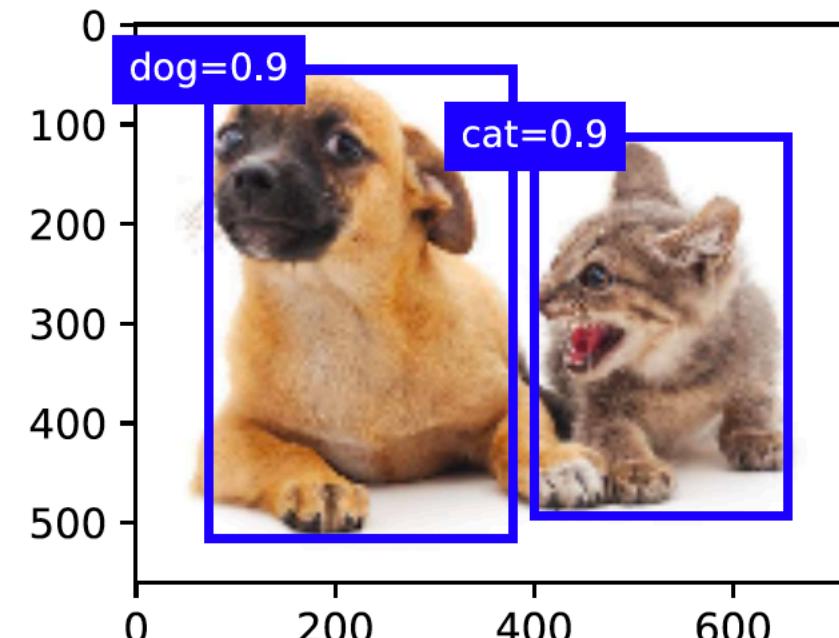
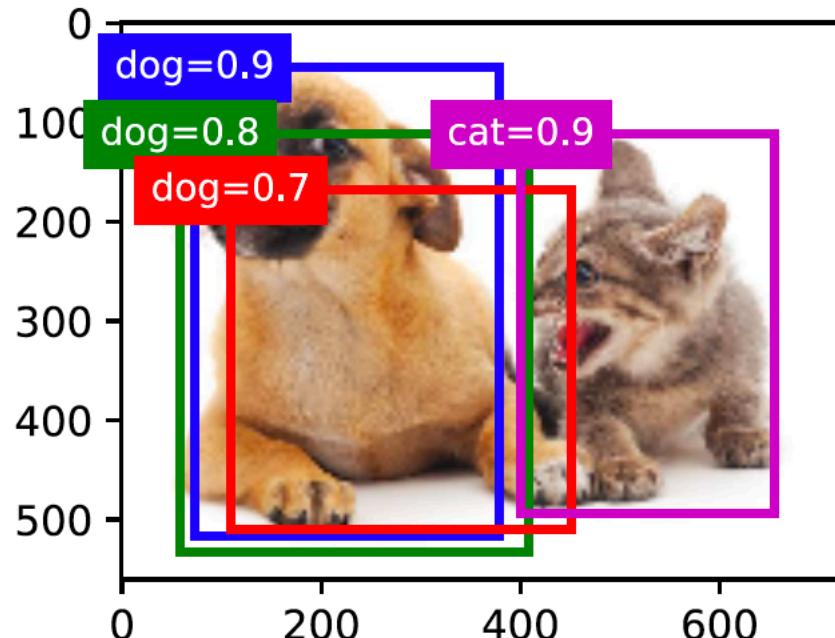
Output Bounding Boxes for Prediction

- During model prediction, we first generate multiple anchor boxes for the image and then predict categories and offsets for these anchor boxes one by one.
- Then, we obtain prediction bounding boxes based on anchor boxes and their predicted offsets.
- When there are many anchor boxes, many similar prediction bounding boxes may be output for the same target.
- To simplify the results, we can remove similar prediction bounding boxes.
- A commonly used method is called non-maximum suppression (NMS).

Output Bounding Boxes for Prediction

- Let us take a look at how NMS works.
- For a prediction bounding box B , the model calculates the predicted probability for each category.
- Assume the largest predicted category is p , the category corresponding to this probability is the predicted category of B .
- We also refer to p as the confidence level of prediction bounding box B .
- On the same image, we sort the prediction bounding boxes with predicted categories other than background by confidence level from high to low, and obtain the list L .
- Select the prediction bounding box B_1 with highest confidence level from L as a baseline and remove all non-benchmark prediction bounding boxes with an IoU with B_1 greater than a certain threshold from L .
- At this point, L retains the prediction bounding box with the highest confidence level and removes other prediction boxes similar to it.
- Next, we select the prediction bounding box B_2 with the second highest confidence level from L as a baseline, and remove all non-benchmark prediction bounding boxes with an IoU with B_2 greater than a certain threshold from L .
- Repeat this process until all prediction bounding boxes in L have been used as a baseline,
- At this time, the IoU of any pair of prediction bounding boxes in L is less than the threshold.
- Finally, output all prediction bounding boxes in the list L

Output Bounding Boxes for Prediction



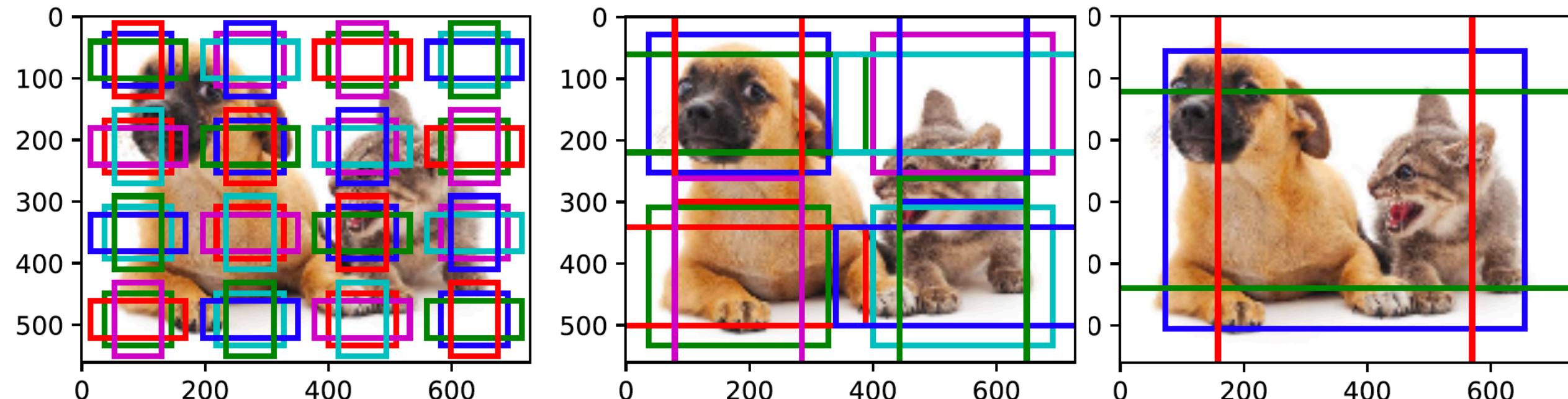
Multiscale Object Detection

- So far, we generated multiple anchor boxes centered on each pixel of the input image.
- These anchor boxes are used to sample different regions of the input image.
- However, if anchor boxes are generated centered on each pixel of the image, soon there will be too many anchor boxes to compute.
- It is not difficult to reduce the number of anchor boxes.
- We can apply uniform sampling on a small portion of pixels from the input image and generate anchor boxes centered on the samples pixels.

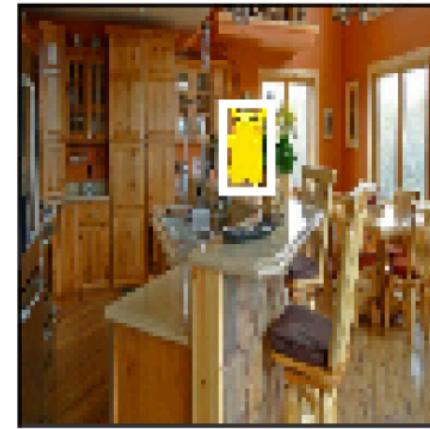
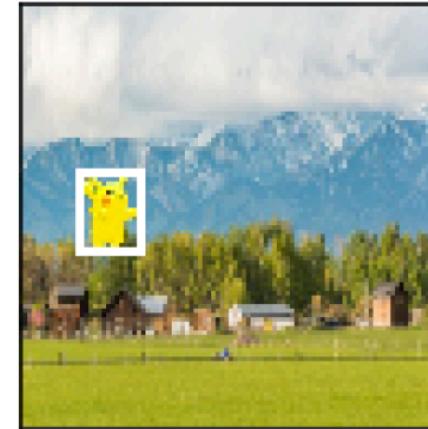
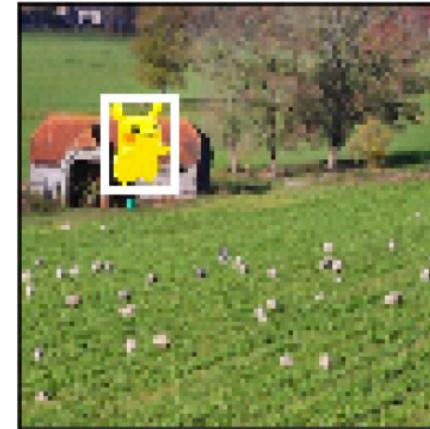
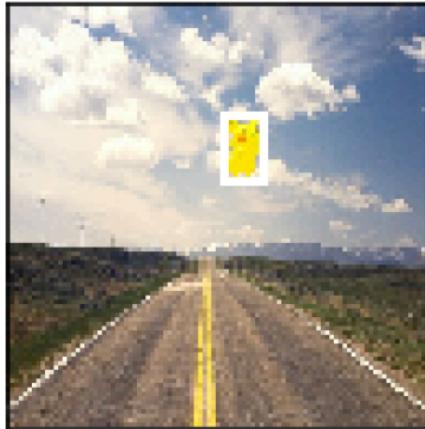
Multiscale Object Detection

- In addition, we can generate anchor boxes of varied numbers and sizes on multiple scales.
- Notice that smaller objects are more likely to be positioned on the image rather than larger ones.
- Consider an image of size 2x2:
 - Say, there are objects with shapes of 1x1, 1x2, and 2x2
 - Each of these objects 4, 2, and 1 possible positions on an image
- Therefore, when using smaller anchor boxes to detect smaller objects, we can sample more regions; when using larger anchor boxes to detect larger objects, we can sample fewer regions.

Multiscale object detection



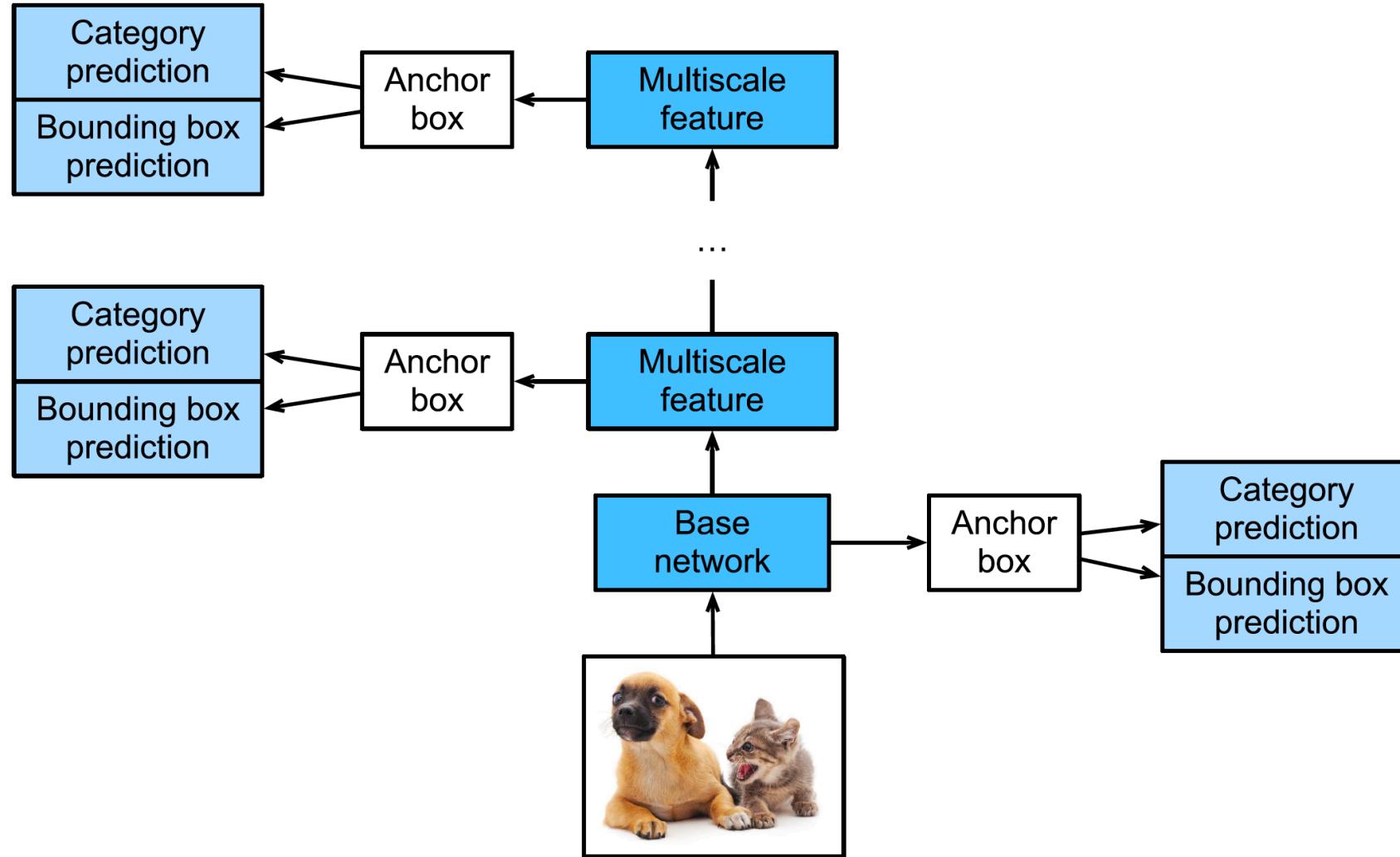
Object Detection Data Set (Pikachu)



Single Shot Multibox Detection (SSD)

- Now, we will use our knowledge on bounding boxes, anchor boxes and multiscale object detection to construct an object detection model: single shot multibox detection (SSD).
- This quick and easy model is already widely used.
- Some of the design concepts and implementation details of this model are also applicable to other object detection models.

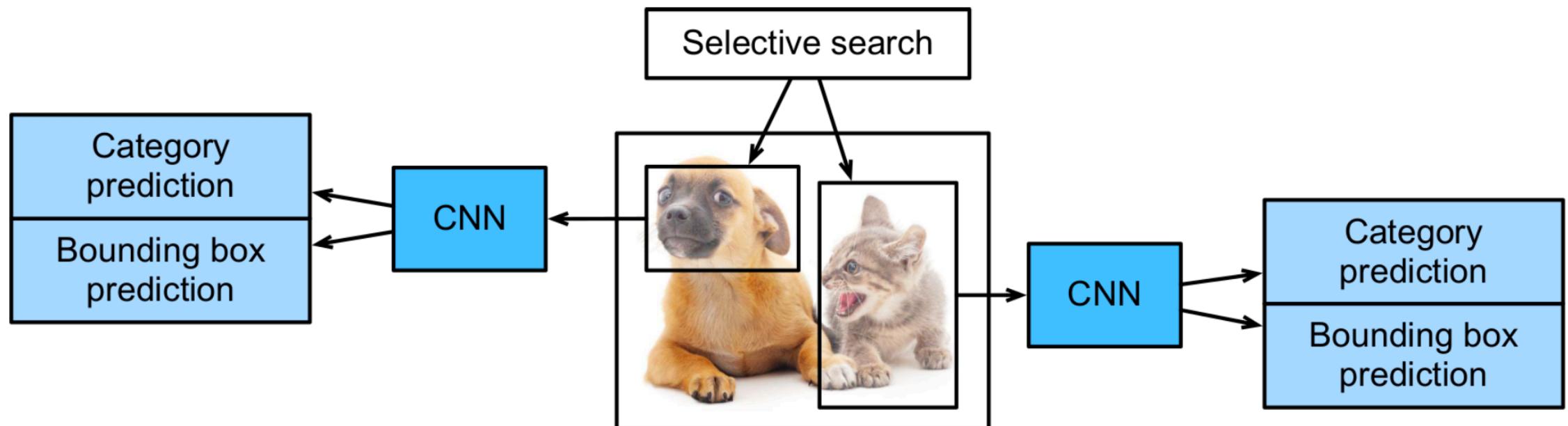
Single Shot Multibox Detection (SSD)



- The base network is used to extract features of original images.
- The closer a feature block is to the top, the better suited it is to detect larger objects.
- As the SSD generates different numbers of anchor boxes of different sizes based on the base network block and each multiscale feature block and then predicts the categories and offsets (i.e., predicted bounding boxes) of the anchor boxes in order to detect objects of different sizes, SSD is a multiscale object detection model.

Region-based CNNs (R-CNNs)

- R-CNN models first select several proposed regions from an image (for example, anchor boxes are one type of selection method) and then label their categories and bounding boxes (e.g., offsets).
- Then, they use a CNN to perform forward computation to extract features from each proposed area.
- Afterwards, we use the features of each proposed region to predict their categories and bounding boxes.

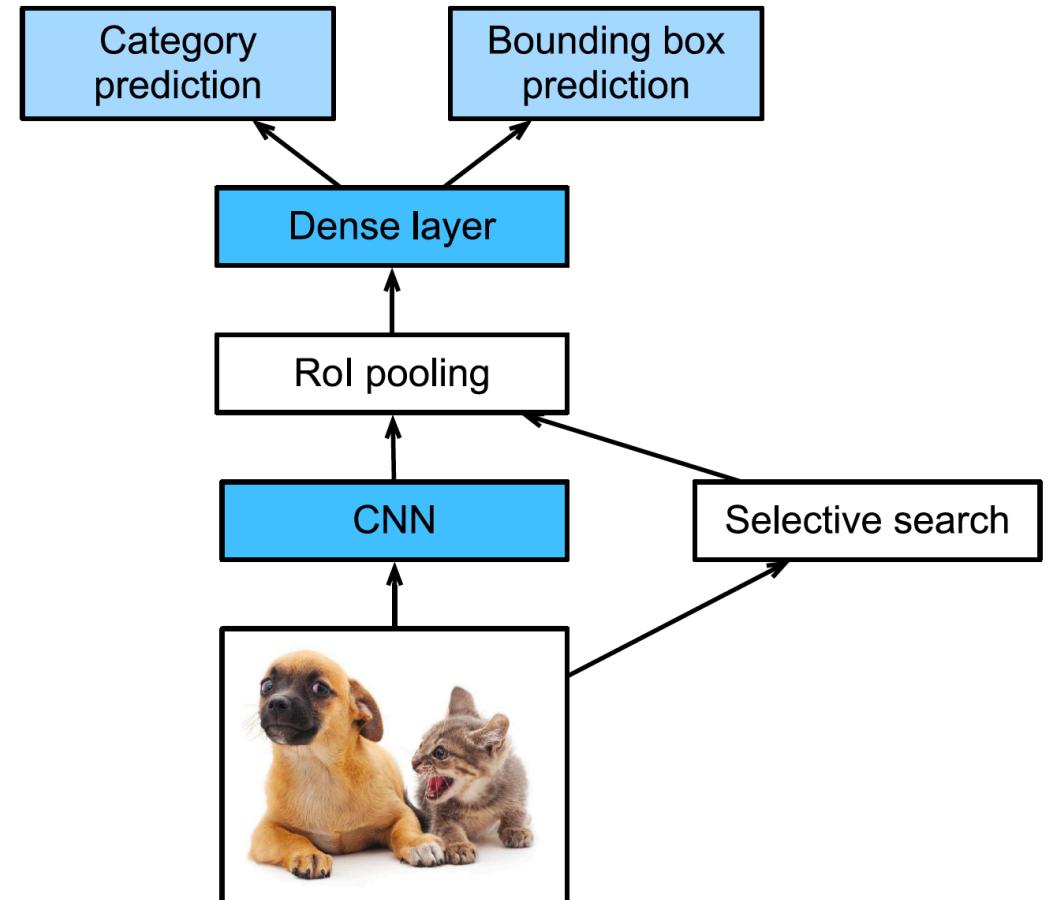


R-CNNs

- Specifically, R-CNNs are composed of three main parts:
 1. Selective search is performed on the input image to select multiple high-quality proposed regions. The category and ground-truth bounding box of each proposed region is labeled.
 2. A pre-trained CNN is selected and placed, in truncated form, before the output layer. It transforms each proposed region into the input dimensions required by the network and uses forward computation to output the features extracted from the proposed regions.
 3. The features and labeled category of each proposed region are combined as an example to train multiple support vector machines for object classification. Here, each support vector machine is used to determine whether an example belongs to a certain category.

Fast R-CNN

- The main downside of R-CNN is the slow speed.
- The main performance bottleneck of an R-CNN model is the need to independently extract features for each proposed region.
- As these regions have a high degree of overlap, independent feature extraction results in a high volume of repetitive computations.
- Fast R-CNN improves on the R-CNN by only performing CNN forward computation on the image as a whole.

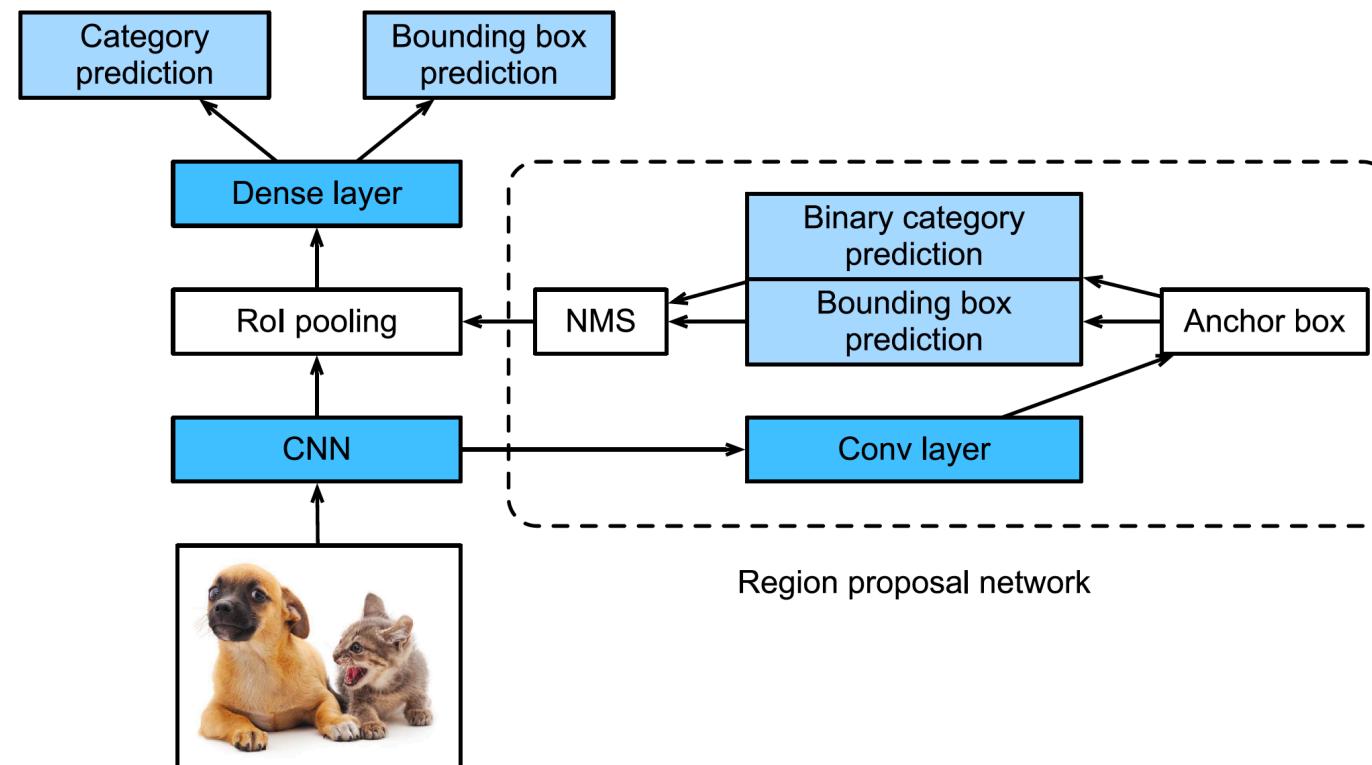


Fast R-CNN

- Fast R-CNN model. It's primary computation steps are described below:
 1. Compared to an R-CNN model, a Fast R-CNN model uses the entire image as the CNN input for feature extraction, rather than each proposed region.
 2. Assuming selective search generates n proposed regions, their different shapes indicate regions of interests (Rois) of different shapes on the CNN output.
 3. Fast R-CNN introduces RoI pooling, which uses the CNN output and Rois as input to output a concatenation of the features extracted from each proposed region.
 4. A fully connected layer is used to transform the output shape to $n \times d$, where d is determined by the model design.
 5. During category prediction, the shape of the fully connected layer output is again transformed to $n \times q$ and we use softmax regression (q is the number of categories). During bounding box prediction, the shape of the fully connected layer output is again transformed to $n \times 4$.

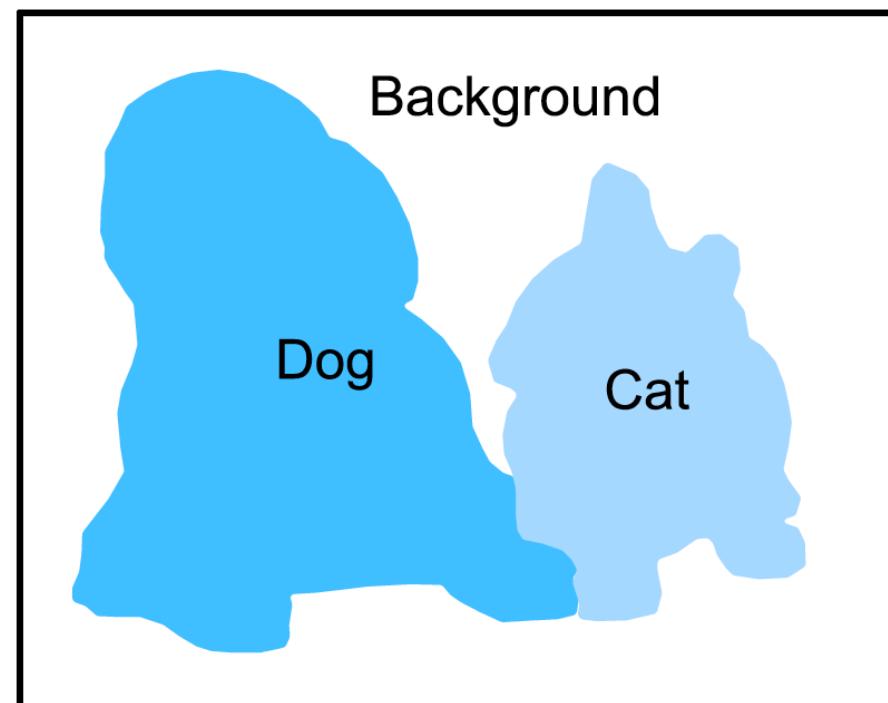
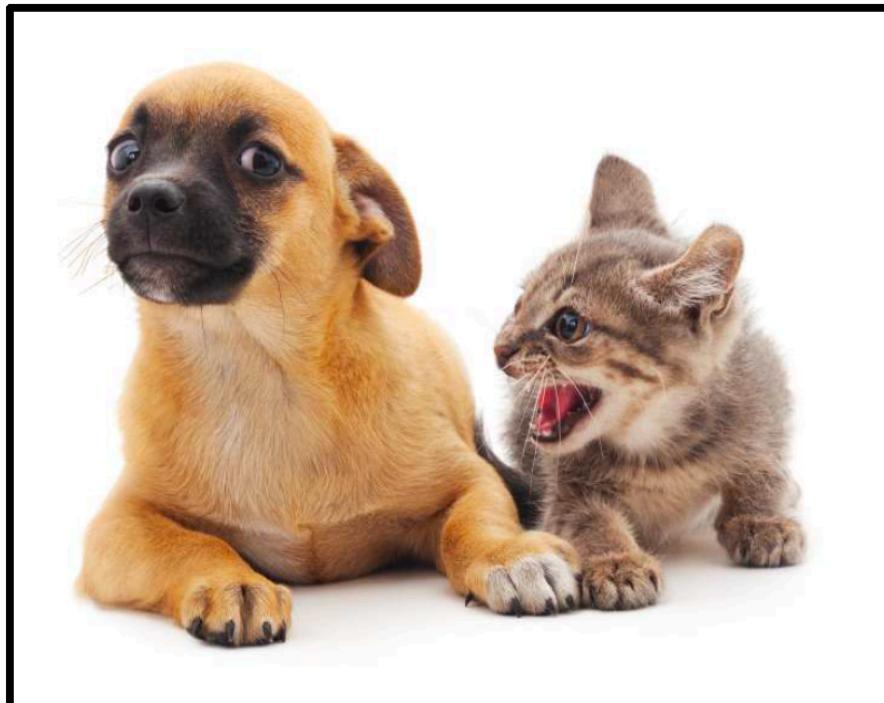
Faster R-CNN

- Fast R-CNN generally requires that many proposed regions be generated in selective search.
- Faster R-CNN replaces selective search with a region proposal network.
- This reduces the number of proposed regions generated, while ensuring precise object detection.



Semantic Segmentation and Data Sets

- So far, we only used rectangular bounding boxes to label and predict objects in images.
- Now, we will look at semantic segmentation, which attempts to segment images into regions with different semantic categories.



Transposed Convolution

- The layers we introduced so far for convolutional neural networks, including convolutional layers and pooling layers, often reduce the input width and height, or keep them unchanged.
- Applications such as semantic segmentation and generative adversarial networks however, require to predict values for each pixel and therefore needs to increase input width and height.
- Transposed convolution, also named fractionally-strided convolution or deconvolution, serves this purpose.

Basic 2D Transposed Convolution

- Let's consider a basic case that both input and output channels are 1, with 0 padding and 1 stride.
- Below figure illustrates how transposed convolution with a 2×2 kernel is computed on the 2×2 input matrix.

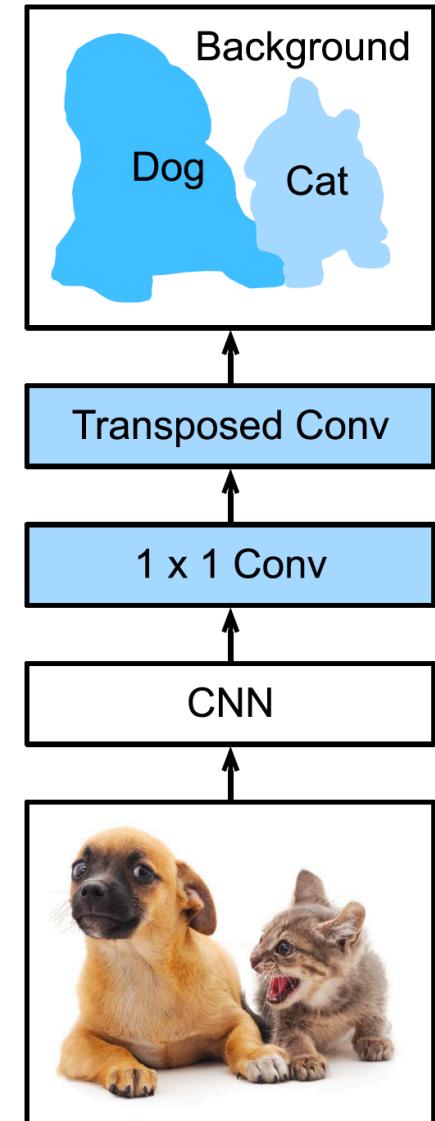
Input	Kernel	= \sum				Output																																																									
<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td></tr><tr><td>2</td><td>3</td></tr></table>	0	1	2	3	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td></tr><tr><td>2</td><td>3</td></tr></table>	0	1	2	3	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>0</td><td></td></tr><tr><td>0</td><td>0</td><td></td></tr><tr><td></td><td></td><td></td></tr></table>	0	0		0	0					$+$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td></td><td>0</td><td>1</td></tr><tr><td></td><td>2</td><td>3</td></tr><tr><td></td><td></td><td></td></tr></table>		0	1		2	3				$+$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td></td><td></td><td></td></tr><tr><td>0</td><td>2</td><td></td></tr><tr><td>4</td><td>6</td><td></td></tr></table>				0	2		4	6		$+$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td></td><td></td><td></td></tr><tr><td>0</td><td>3</td><td></td></tr><tr><td>6</td><td>9</td><td></td></tr></table>				0	3		6	9		$=$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>4</td><td>6</td></tr><tr><td>4</td><td>12</td><td>9</td></tr></table>	0	0	1	0	4	6	4	12	9
0	1																																																														
2	3																																																														
0	1																																																														
2	3																																																														
0	0																																																														
0	0																																																														
	0	1																																																													
	2	3																																																													
0	2																																																														
4	6																																																														
0	3																																																														
6	9																																																														
0	0	1																																																													
0	4	6																																																													
4	12	9																																																													

Fully Convolutional Networks (FCN)

- Semantic segmentation uses each pixel in an image for category prediction.
- An FCN uses convolutional neural network to transform image pixels to image categories.
- An FCN transforms the height and width of the intermediate layer feature map back to the size of input image through the transposed convolution layer.
- This way the predictions have a one-to-one correspondence with input image in spatial dimension (height and width).
- Given a position on the spatial dimension, the output of the channel dimension will be a category prediction of the pixel corresponding to the location.

Fully Convolutional Neural Networks

- The fully convolutional network first uses the convolutional neural network to extract image features.
- Then transforms the number of channels into the number of categories through the 1×1 convolution layer.
- And finally transforms the height and width of the feature map to the size of the input image by using the transposed convolution layer.
- The model output has the same height and width as the input image and has a one-to-one correspondence in spatial positions.
- The final output channel contains the category prediction of the pixel of the corresponding spatial position.



Generative Adversarial Networks

- Throughout most of this book, we've talked about how to make predictions.
- In some form or another, we used deep neural networks learned mappings from data points to labels.
- This kind of learning is called discriminative learning, as in, we'd like to be able to discriminate between photos of cats and photos of dogs.
- Classifiers and regressors are both examples of discriminative learning.
- But there's more to machine learning than just solving discriminative tasks.
- For example, given a large dataset, without any labels, we might want to learn a model that concisely captures the characteristics of this data.
- Given such a model, we could sample synthetic data points that resemble the distribution of the training data.
- This kind of learning is called generative modeling.

Generative Adversarial Networks

- In 2014, a breakthrough paper introduced Generative adversarial networks (GANs), a clever new way to leverage the power of discriminative models to get good generative models.
- At their heart, GANs rely on the idea that a data generator is good if we cannot tell fake data apart from real data.
- In statistics, this is called a two-sample test - a test to answer the question whether datasets X and X' were drawn from the same distribution.
- The main difference between most statistics papers and GANs is that the latter use this idea in a constructive way.
- In other words, rather than just training a model to say “hey, these two datasets don’t look like they came from the same distribution”, they use the two-sample test to provide training signal to a generative model.
- This allows us to improve the data generator until it generates something that resembles the real data.
- At the very least, it needs to fool the classifier.

Generative Adversarial Networks

- There are two pieces to GANs: the generator network and the discriminator network.
- It attempts to distinguish fake and real data from each other.
- Both networks are in competition with each other.
- The generator network attempts to fool the discriminator network.
- At that point, the discriminator network adapts to the new fake data.
- This information, in turn is used to improve the generator network, and so on.

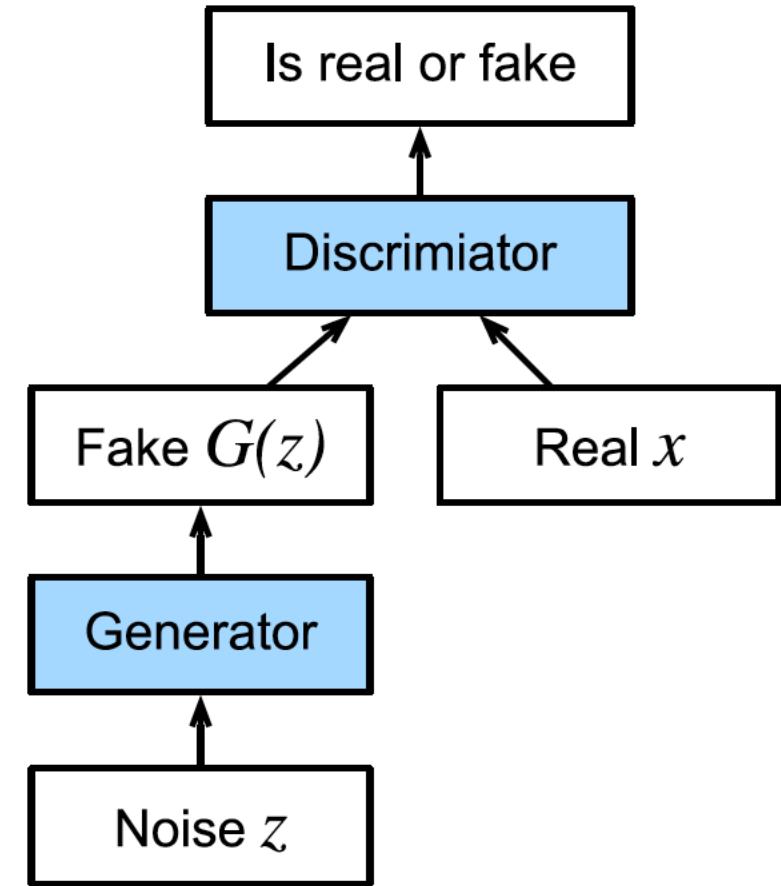


Fig. 16.1.1: Generative Adversarial Networks

Generative Adversarial Networks

- The discriminator is a binary classifier to distinguish if the input x is real (from real data) or fake (from the generator).
- Assume the label y for true data is 1 and 0 for fake data.
- We train the discriminator to minimize the cross-entropy loss

$$\min -y \log D(\mathbf{x}) - (1 - y) \log(1 - D(\mathbf{x})),$$

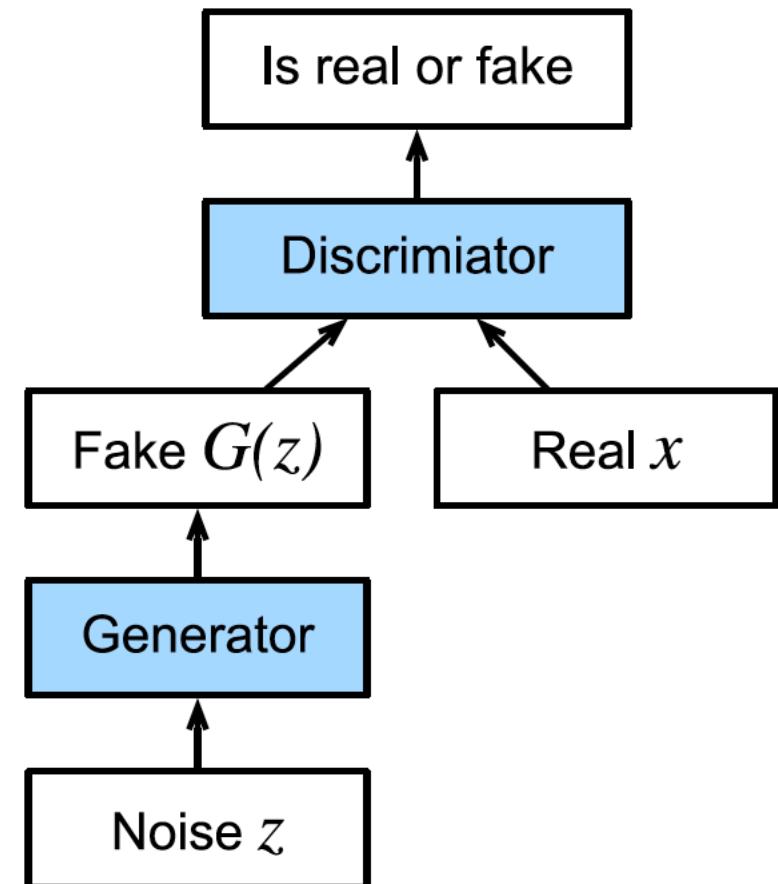


Fig. 16.1.1: Generative Adversarial Networks

Generative Adversarial Networks

- For the generator, it first draws some parameter \mathbf{z} from a source of randomness.
- We often call \mathbf{z} the latent variable.
- It then applies a function to generate $\mathbf{x}' = G(\mathbf{z})$
- The goal of the generator is to fool the discriminator to classify \mathbf{x}' as true data.
- In other words, we update the parameters of the generator to maximize the cross-entropy loss when $y = 0$.
$$\max -\log(1 - D(\mathbf{x}')).$$
- If the discriminator does a perfect job, then the above loss will be close to 0 which results the gradients are too small to make a good progress for the generator.
- So, commonly, we optimize $\max \log(D(\mathbf{x}'))$,

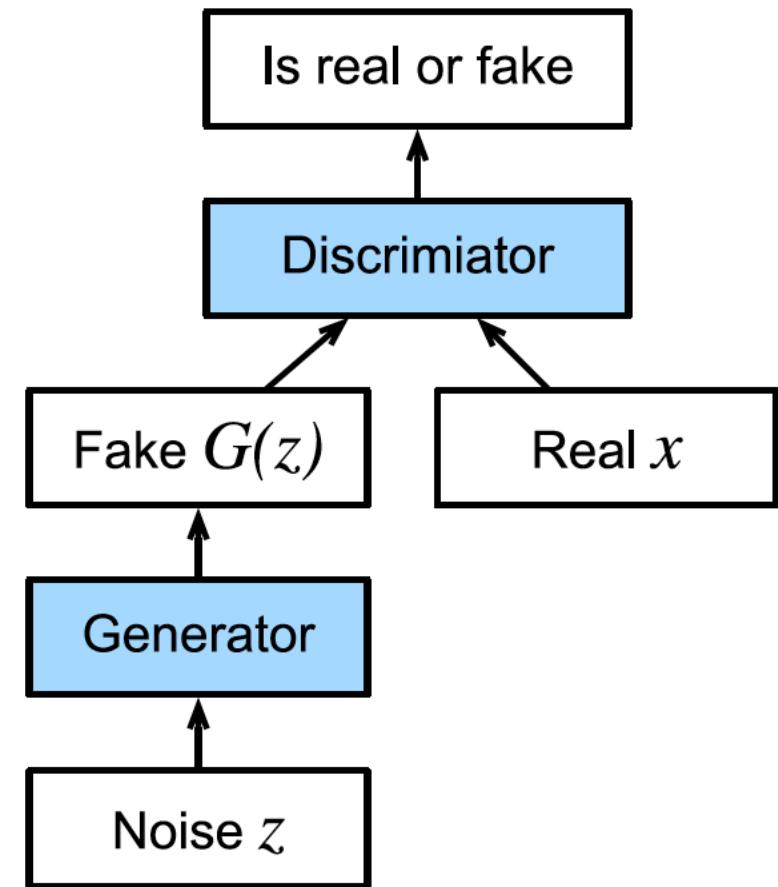


Fig. 16.1.1: Generative Adversarial Networks

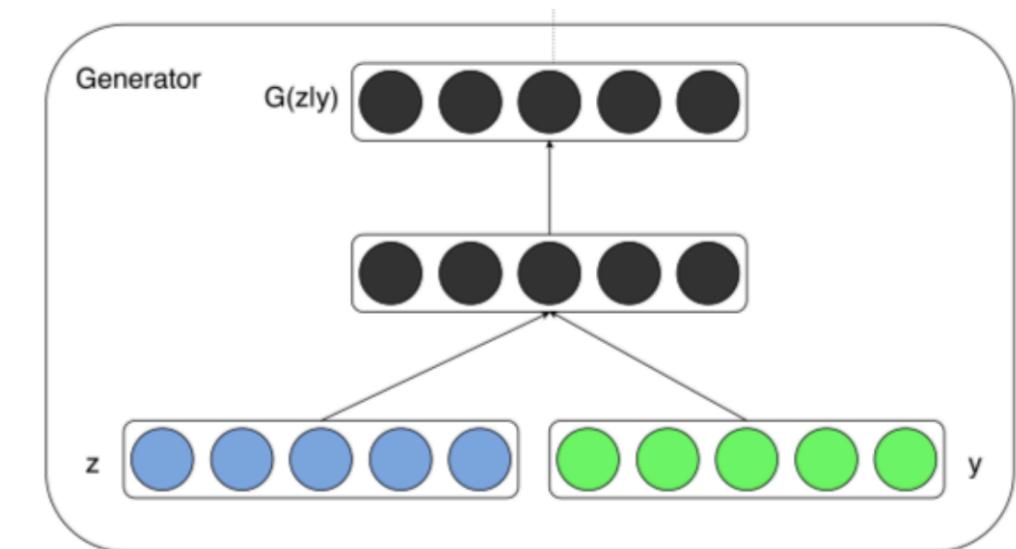
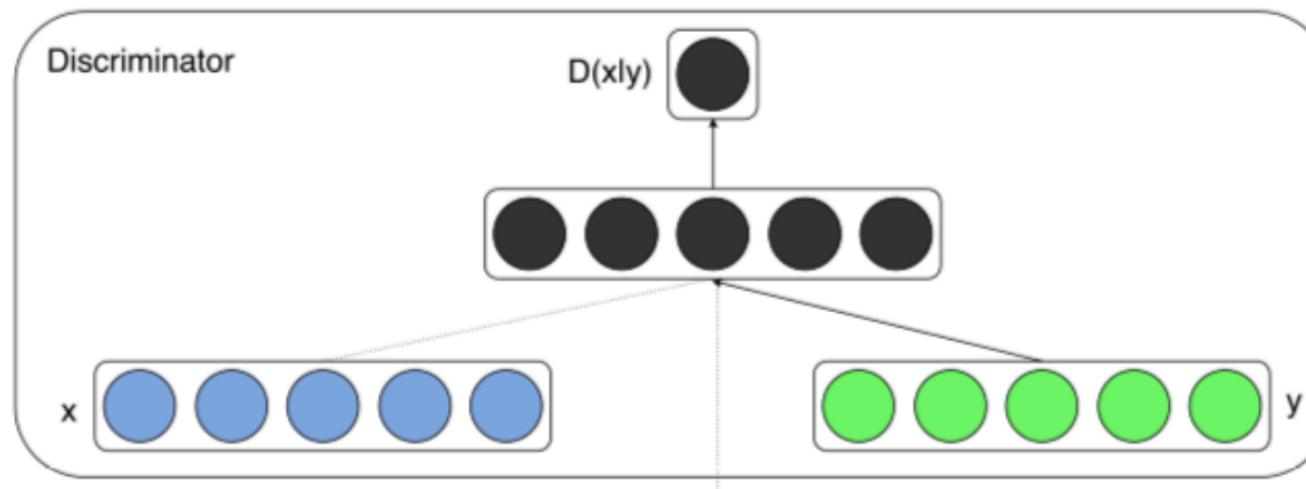
This is just the beginning of the story of GANs...

- You're currently at the second red 'X'



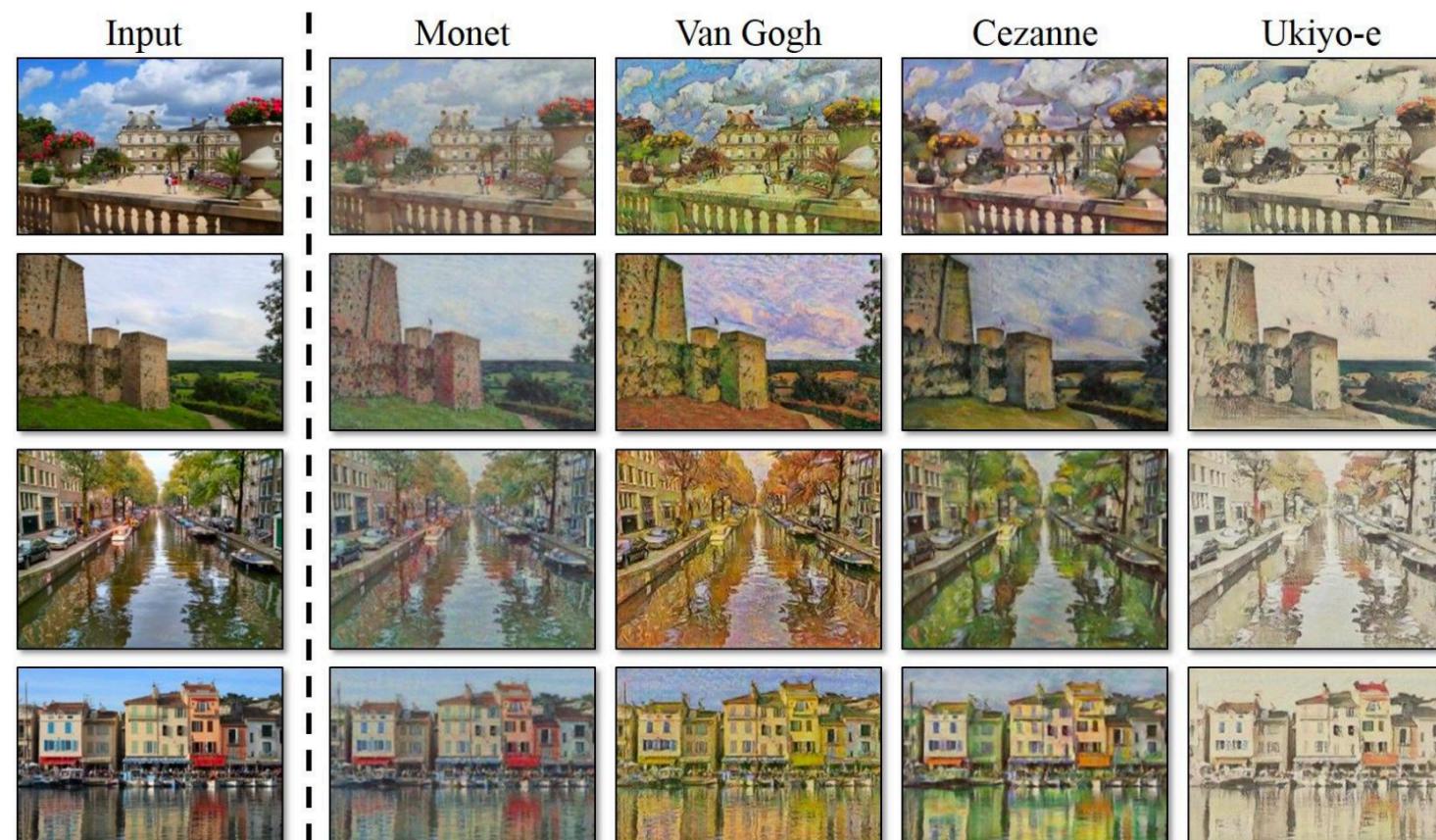
CGAN: Conditional Generative Adversarial Network

- What if your training dataset includes both cats and dogs? Your GAN will produce blurry half breeds.
- CGAN aims to solve this issue by telling the generator to generate images of only one particular class, like a cat, dog, or Nicholas Cage.
- CGAN concatenates a one-hot vector y to the random noise vector z to result in an architecture:



CycleGAN (<https://arxiv.org/abs/1703.10593v6>)

- GANs aren't used just for generating images. They can also be used to learn a mapping from one domain of images to another. So we train it on say, the set of all Monet paintings.



ProGAN: Progressive growing of Generative Adversarial Networks (<https://arxiv.org/pdf/1710.10196.pdf>)

- There are many problems with training GANs. The most important of which is the training instability.
- ProGAN is a technique that helps stabilize GAN training by incrementally increasing the resolution of the generated image.
- The intuition here is that it's easier to generate a 4x4 image than it is to generate a 1024x1024 image. Also, it's easier to map a 16x16 image to a 32x32 image than it is to map a 2x2 image to a 32x32 image.
- So ProGAN first trains a 4x4 generator and a 4x4 discriminator and adds layers that correspond to higher resolutions later in the training process.
- https://cdn-images-1.medium.com/max/1600/1*tUhgr3m54Qc80GU2BkaOiQ.gif

WGAN: Wasserstein Generative Adversarial Networks

<https://arxiv.org/abs/1701.07875v3>

- In short, WGAN (the ‘W’ stands for Wasserstein) proposes a new cost function that some nice properties
- The original GAN paper showed that when the discriminator is optimal, the generator is updated in such a way to minimize the Jensen-Shannon divergence.

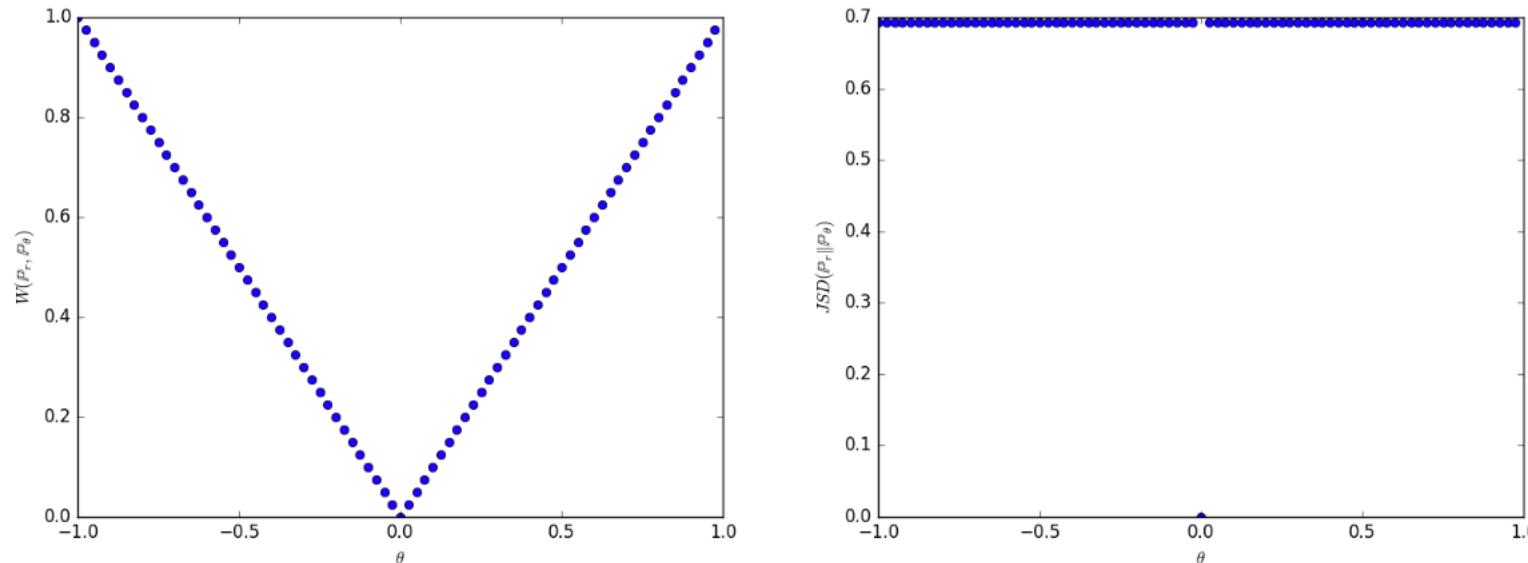


Figure 1: These plots show $\rho(\mathbb{P}_\theta, \mathbb{P}_0)$ as a function of θ when ρ is the EM distance (left plot) or the JS divergence (right plot). The EM plot is continuous and provides a usable gradient everywhere. The JS plot is not continuous and does not provide a usable gradient.