# Deep Neural Networks for Object Detection

Paper by C. Szegedy, A. Toshev, D. Erhan [2013] Presentation by Joaquín Ruales

### The Problem: Object Detection

#### Identifying and locating objects in an image



### The Problem: Object Detection

• Identifying and locating objects in an image



### Previous Work in Object Detection

- Discriminative Part-based models:
  - Identifying parts of an object and their relation to identify the whole
  - Exploits domain knowledge. Uses HOG descriptors
- Some NN approaches, but used as local classifiers, or incapable of distinguishing many instances of same class of object



# Why DNN for Object Detection?

- Success of DNNs for related problem: image classification
  - <u>A. Krizhevsky, I. Sutskever, G.</u> <u>Hinton. (2012). ImageNet</u> <u>Classification with Deep</u> <u>Convolutional Neural Networks</u>
- Can take advantage of the small shift-invariance in DNN image classification
- Simpler models, easily extensible to new classes of objects



### Deep Neural Networks for Object Detection

- This paper uses DNNs to classify and precisely locate objects of 20 classes (plane, bicycle, bird, boat, etc.)
- Requires several applications of the DNNs
- Obtains state-of-the-art performance on the Pascal VOC dataset



Detection

### Detection

- For each object category X∈{plane, bicycle, bird, boat, etc.}
  - Input: Image.
  - Step 1: Generate binary masks using DNN specific to X
  - Step 2: Get bounding boxes from masks
  - Step 3: Refine bounding boxes
  - <u>Output:</u> Bounding boxes and confidence scores for all objects of type X in the image



- Same DNN structure as [<u>A. Krizhevsky, I. Sutskever, G.</u> <u>Hinton. (2012). ImageNet Classification with Deep</u> <u>Convolutional Neural Networks</u>]
  - 5 convolutional layers (3 with max pooling), 2 connected layers, ReLu nonlinearities
  - Except: replace softmax classification layer (last layer) with a regression layer that produces a binary mask



- Same DNN structure as [<u>A. Krizhevsky, I. Sutskever, G.</u> <u>Hinton. (2012). ImageNet Classification with Deep</u> <u>Convolutional Neural Networks</u>]
  - 5 convolutional layers (3 with max pooling), 2 connected layers, ReLu nonlinearities
  - Except: replace softmax classification layer (last layer) with a regression layer that produces a binary mask



- Actually, 5 DNNs trained per category
  - Full object mask, left half, bottom half, right half, top half
  - 5 masks are then merged to get the final mask
- DNN inputs are 225x225 pixels. Output masks are 24x24 pixels



- Compute these masks for many sub windows of the original image, at several scales
- (Different than sliding window approach since usually need <40 windows per image)</li>



#### Detection Step #2: Get Bounding Boxes

Find the bounding boxes with best scores for the set of  $\bullet$ 24x24px output masks Percentage of bounding box that

$$S(bb) = \sum_{h \in halves} (S(bb(h), m^h) - S(bb(\bar{h}), m^h))$$

The complement of region h

- (exhaustive search. Sped up using integral images)
- Map bounding box back to image space (note resolution loss)



#### Detection Step #3: Refine bounding boxes

- Crop original image to each bounding box
- Repeat step #1 (Generate binary masks with DNN) on the cropped image
- Repeat step #2 (Get bounding boxes) for the generated binary masks
- Discard the bounding boxes that received a low score
- Run the detected object through a classifier DNN and discard the corresponding bounding box if misclassified
- Result: Final, fine-grained bounding boxes around the object with scores



#### Precision and Recall Before and After Refinement

Based on results on VOC2007 test data



# Training

# Training

- Needs a lot of training data: Objects of different sizes at almost every location
- Use VOC2012 training and validation set (~11K images) for training
- Remember: we need to train 2 types of DNNs:
  - 1) Mask generator DNN (maps images to binary masks)
  - 2) Classifier DNN (used for final pruning of detections)

## 1) Mask Generator Training

- Krizhevsky et al. ImageNet CNN with last layer replaced by regression layer
- Minimize  $L_2$  error for predicting a ground truth mask *m* for an image *x* Regularizer in **R**<sup>+</sup>. Ground truth



# 1) Mask Generator Training

- Several thousand samples from each image (10M total)
- 60% negative examples
  - outside of bounding box of any object of interest
- 40% positive examples
  - each covers >80% of area of some ground truth bounding box of interest
- Crops sampled so that cropWidth~Uniform(minScale, imageWidth)



# 2) Classifier Training

- Krizhevsky et al. ImageNet CNN
- Several thousand samples per image (10M total)
- 60% negative examples
  - each has <0.2 Jaccard-similarity with any ground truth box
  - acts as a 21st class in the classifier
- 40% positive examples
  - each has >0.6 Jaccard-similarity with any ground truth box
  - labeled according to category of most similar bounding box



# Final Notes on Training

- CNNs, max pooling, dropout
- AdaGrad training

# Assume the gradient dx and parameter vector x
cache += dx\*\*2
x += - learning\_rate \* dx / np.sqrt(cache + 1e-8)

- A type of adaptive learning rate for SGD
- Training for localization harder than for classification, so they reuse the classification DNN weights for the localization DNN







Results

### Results

- Algorithm obtained state-of-the-art for VOC2007 (Pascal Visual Object Challenge 2007) dataset
  - Best detection for 8 of the 20 categories
  - Best detection for 5 out of 7 animal categories (bird, cat, cow, dog, sheep)
  - 5-6sec per image per class on a 12-core machine
  - More training data than others in this table. Unfair comparison?

class	aero	bicycle	bird	boat	bottle	bus	car	cat	chair	COW
DetectorNet <sup>1</sup>	.292	.352	.194	.167	.037	.532	.502	.272	.102	.348
Sliding windows <sup>1</sup>	.213	.190	.068	.120	.058	.294	.237	.101	.059	.131
3-layer model [19]	.294	.558	.094	.143	.286	.440	.513	.213	.200	.193
Felz. et al. [9]	.328	.568	.025	.168	.285	.397	.516	.213	.179	.185
Girshick et al. [11]	.324	.577	.107	.157	.253	.513	.542	.179	.210	.240
class	table	dog	horse	m-bike	person	plant	sheep	sofa	train	tv
class DetectorNet <sup>1</sup>	table .302	dog .282	horse .466	m-bike .417	person .262	plant .103	sheep .328	sofa .268	train .398	tv .470
class         DetectorNet <sup>1</sup> Sliding windows <sup>1</sup>	table .302 .110	dog .282 .134	horse .466 .220	m-bike .417 .243	person           .262           .173	plant .103 .070	sheep .328 .118	sofa .268 .166	train .398 .240	tv .470 .119
classDetectorNet1Sliding windows13-layer model [19]	table           .302           .110           .252	dog .282 .134 .125	horse .466 .220 .504	m-bike .417 .243 .384	person           .262           .173           .366	plant .103 .070 .151	sheep           .328           .118           .197	sofa .268 .166 .251	train .398 .240 .368	tv .470 .119 .393
classDetectorNet1Sliding windows13-layer model [19]Felz. et al. [9]	table           .302           .110           .252           .259	dog           .282           .134           .125           .088	horse .466 .220 .504 .492	m-bike .417 .243 .384 .412	person           .262           .173           .366           .368	plant .103 .070 .151 .146	sheep           .328           .118           .197           .162	sofa .268 .166 .251 .244	train .398 .240 .368 .392	tv .470 .119 .393 .391

Table 1: Average precision on Pascal VOC2007 test set.



Thank You