

C5: Visual Recognition

Report on Object Detection and Segmentation

Team 1: Carles Pregonas Graells, Pau Vallespí Monclús, Carlos Boned Riera, Marc Pérez Sabater

Abstract—In this report, we aim to showcase what we have learnt in the "Visual Recognition" module. We will be presenting our results, acquired knowledge, insights, and the difficulties that we encountered during this project.

I. INTRODUCTION

In the beginning of this course, we started with an exploration of PyTorch, a robust deep learning framework widely acclaimed for its versatility and performance, and will compare it with Keras.

Moving forward, we'll dive into the realm of object detection, recognition, and segmentation. Armed with tools like Detectron2 [1], equipped with cutting-edge models such as Mask R-CNN[2] and Faster R-CNN[3], we'll unravel the methodologies employed to pinpoint objects within images, discern their identities, and delineate their contours with precision. After an initial inference in this models, we will perform fine tuning on the models parameters in order to find the best configuration for our data. Alongside, we'll explore the utilization of YOLO[4], an alternative framework renowned for its efficiency and simplicity in object detection tasks.

Transitioning to the realm of image retrieval, we'll peer into the mechanisms underpinning content-based image search engines. Here, sophisticated algorithms analyze image content to retrieve relevant matches from extensive databases, facilitating seamless access to visual information.

Lastly, we'll turn our attention to cross-modal retrieval—a convergence of computer vision and natural language processing. By bridging different data modalities, such as images and text, we'll explore how meaningful cross-references and information retrieval are facilitated, paving the way for enhanced understanding and utilization of multimodal data sources.

Throughout this journey, we aim to provide a comprehensive understanding of visual recognition techniques, underpinned by the utilization of sophisticated tools and methodologies. By interpreting these diverse facets, we seek to underscore the significance of visual recognition across various domains and applications.

II. RELATED WORK

In this section, we'll delve into the narrative of object detection, exploring various perspectives such as vanilla detectors and the evolution of more sophisticated methods.

A. Vanilla Object Detectors

The transition in object detection, led by deep learning, contrasts with the pioneering era of the 1990s, where computer vision relied heavily on handcrafted features due to limited image representation. Viola-Jones Detectors in 2001 [5], [6] achieved real-time human face detection, outperforming contemporaneous algorithms through techniques like "integral image," "feature selection," and "detection cascades." Subsequent innovations, like Histogram of Oriented Gradients (HOG) [7], significantly advanced pedestrian detection. Another milestone was the introduction of Deformable Part-based Models [8], [9] in 2008, which revolutionized object detection by decomposing it into distinct parts.

In the realm of deep learning, object detectors are categorized into two main groups: *two-stage detectors* and *one-stage detectors*. The former approaches detection as a "coarse-to-fine" process, while the latter aims to complete detection in a single step.

B. Two-stage detectors

R-CNN [10] uses selective search [11]. It begins by generating a collection of 2000 object proposals (candidate boxes). These proposals are then resized to a standardized image size and passed through a pre-trained CNN model like AlexNet, trained on ImageNet, to extract features and then an SVM is used to classify them. Finally, a linear regression model is trained to generate bounding boxes for each identified object. R-CNN takes around 47 seconds for each image (inference), and the training stage is expensive and slow, as it extracts features with a CNN from 2000 regions per image.

Earlier CNN models, such as AlexNet, necessitated a fixed-size input, typically requiring images to be resized to dimensions like 224x224 pixels. A first solution to solve this problem comes from **Spatial Pyramid Pooling** [12] (**SPP**), whose goal is to get fixed-length representations for variable-size feature maps. In this case, a CNN is run just once per image to obtain a feature map, and then a (variable size) window related to the region proposals to detect the objects in the image.

This concept is further developed in **Fast R-CNN**[13], which employs RoI pooling, akin to SPP. The RoI pooling layer adjusts the region proposals to match the CNN input size. Each region is forwarded to a Fully Connected Network (FCN), where a softmax layer and a linear regression layer produce class predictions and bounding box coordinates. Fast R-CNN utilizes a combined loss function. Fast R-CNN addresses two primary challenges of R-CNN: it reduces the

Object Detection Milestones

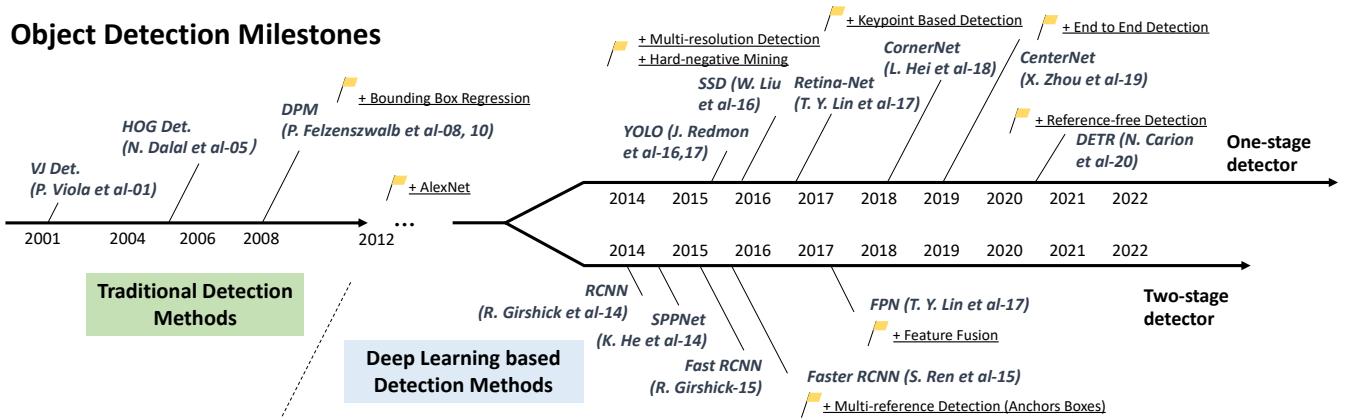


Fig. 1. Chronology of state-of-the-art object detection and instance segmentation architectures.

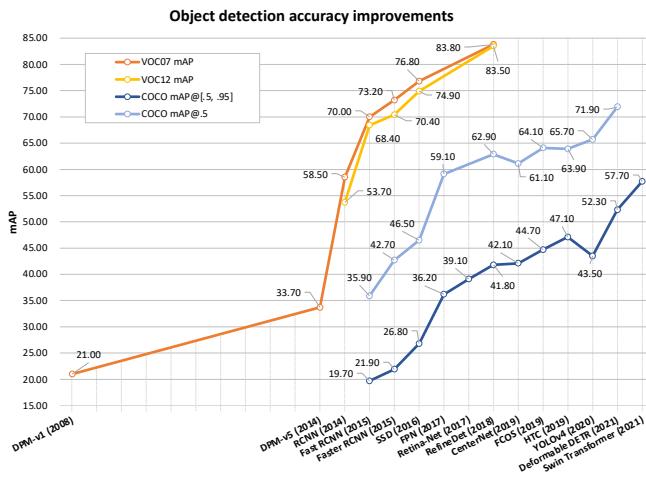


Fig. 2. Accuracy improvement in object detection across VOC07, VOC12, and MS-COCO datasets.

number of regions passed to a CNN from 2000 to one per image, and it merges feature extraction, classification, and bounding box generation into a single model. However, a new bottleneck emerges with the selective search algorithm used to locate Regions of Interest (RoIs), which is slow and time-consuming, resulting in an inference time of approximately 2 seconds per image.

To solve this problem, **Faster R-CNN** [14] introduces Region Proposal Networks (RPN): using a sliding window over the feature maps from a CNN, it generates 9 anchor boxes of different shapes and sizes at each sliding position. For each anchor, RPN predicts two things: the probability that an anchor is an object (class-agnostic), and the bounding box regressor for adjusting the anchors to better fit the object.

In Faster R-CNN, the inference time is around 0.2 seconds per image, which is much faster than R-CNN and Fast R-CNN, but is still far from real-time object detection.

C. Object instance segmentation

Moving towards segmentation, Mask R-CNN [15] builds upon Faster R-CNN by incorporating a parallel branch dedicated to predicting object masks alongside the existing

branch for bounding box detection. This model applies segmentation to the Region Proposal Network (RPN) predictions, resulting in the generation of high-quality segmentation masks for each instance.

D. One-stage detectors

Instead of using region proposals, one-stage detectors make predictions by only looking into the input image once, which can result in a faster performance. One of the state-of-the-art single pass detectors is **YOLO** [16] (and its successors [17], [18], [19]) and so on upon YOLO9, which idea is to extract a feature map using a CNN called Darknet, divide it into $S \times S$ cells (YOLO uses $S = 7$), and predict one bounding box for each cell. YOLO predicts the class of that bounding box and if it is centered at that cell, and uses Non-Maximum Suppression (NMS) to reduce the number of output bounding boxes.

Another well-known one-stage detector is **SSD** [20], which makes predictions at multiple feature maps (YOLO only does it for one) using a VGG16 network [21] as a feature extractor. The idea is that each feature map, which has different local receptive fields, specializes in objects of different sizes.

The NMS method is also used at the end of the SSD model, and Hard Negative Mining (HNM) is then used to further reduce the number of predicted negative boxes.

Finally, **RetinaNet** [22] has been formed by making two improvements over existing single stage detectors: Feature Pyramid Networks (FPN) [23] and Focal Loss. FPN replaces the feature extractor of detectors like Faster R-CNN, and focal loss is introduced to handle the class imbalance problem with one-stage object detection models.

A comparison between some one-stage and two-stage detectors on different datasets such as the COCO dataset [24] is shown in Fig. 2.

III. METHODOLOGY

In this section, we describe the methodology employed for the experimentation, including model selection, dataset preparation, training procedure, and evaluation metrics. We implemented the models using the Detectron2 framework.

It includes a model zoo of state-of-the-art object detection and instance segmentation algorithms such as Faster R-CNN, RetinaNet and Mask R-CNN. Each model configuration has four parts:

- Backbone: ResNet (R) or ResNext (X)
- Number of layers: 50 or 101
- Backbone combination: ResNet + Feature Pyramid Network (FPN), ResNet conv4 backbone with conv5 head (C4) or ResNet conv5 with dilations (DC5)
- Learning Rate (LR) Scheduler: 1x or 3x

For example, Mask R-CNN R_50_FPN_3x uses ResNet with 50 layers as backbone, a FPN and a 3x LR scheduler.

In our case we choose between Faster R-CNN & Mask R-CNN.

We fine-tuned Faster R-CNN and Mask R-CNN models for object detection and segmentation tasks. We used pre-trained weights obtained from the MSCOCO2017 dataset to initialize the models. For each task we modified the labeling of the predefined COCO classes. The models were then initialized with the pre-trained weights. In Figure 3 we can see a summary of the working flow of both models. The diagram shows schematically how an image is processed until we get the classification with bounding boxes, and the additional masks in the Mask-RCNN case.

Through iterative optimization, the entire model was fine-tuned, gradually unfreezing more layers and adjusting learning rates to learn task-specific features without forgetting pre-trained knowledge. Regularization techniques like dropout or weight decay were applied to prevent overfitting and improve generalization on unseen data.

Finally, we evaluated the trained models using the defined COCOEvaluator, measuring their performance on a validation set. Additionally, we conducted inference on the test set, generating predictions for object instances and their segmentation masks.

Out of Context Tasks

We studied the behaviour of the networks in unusual situations using **Out Of Context (OOC) datasets**. With this purpose, we use a small subset of images with OOC objects from SUN dataset[25], and we modify some specific images to perform experiments. Some example can be seen in Figures 5 and 6.

IV. EXPERIMENTAL DESIGN

In this section, we outline the experimental methodology employed in our study. We detail dataset selection, training strategies, hyperparameter tuning, and evaluation metrics.

A. Dataset

In our study, we assess the performance of pre-trained models by employing the COCO metrics provided by the Detectron2 framework [1]. To facilitate this evaluation, it was imperative to convert the KITTI-MOTS dataset into the COCO format. The KITTI-MOTS dataset encompasses two distinct annotation formats: PNG and TXT. These formats

provide access to crucial data, including detection and segmentation masks, as well as other pertinent annotations necessary for both the evaluation and training phases. Detailed descriptions of these formats are available on the dataset's official documentation. An example of an image with the annotations in the KITTI-MOTS dataset can be seen in Image 4.

The dataset comprises a collection of images depicting various traffic scenarios from the vantage point of a vehicle driver. It spans a diverse range of environments, including urban streets, inter-urban roadways, and rural routes. Within this dataset, three categories of road users are annotated, along with designated regions that are to be ignored. For the purpose of our investigation, we have chosen to focus exclusively on the classes representing cars and pedestrians.

An analysis of the dataset reveals a distribution wherein 70% of the annotated instances are vehicles, and the remaining 30% are pedestrians. This disparity in representation suggests a potential bias where models might demonstrate enhanced detection capabilities for vehicles as compared to pedestrians, due to the more abundant exemplars present in the dataset.

B. Train-Val-Test Split

We utilize the official validation set of KITTI Multiple Object Tracking and Segmentation (KITTI MOTS) as our test set. For training and validation, we employ the training sets of KITTI MOTS consisting of 12 sequences and the MOTSChallenge consisting of 4 sequences. Detailed statistics are shown in Tab. I :

	KITTI MOTS		MOTSChallenge
	train	test	train
# Sequences	12	9	4
# Frames	5,027	2,981	2,862
# Masks Pedestrian			
Total	8,073	3,347	26,894
Manually annotated	1,312	647	3,930
# Masks Car			
Total	18,831	8,068	-
Manually annotated	1,509	593	-

TABLE I
STATISTICS OF THE MODIFIED SPLIT OF KITTI MOTS AND MOTSCHALLENGE DATASETS

To assess our model's performance, we employ a 4-fold cross-validation strategy. Each fold comprises 3 sequences from KITTI MOTS and 1 sequence from MOTSChallenge. In each experiment, we train our model using 3 folds and validate it on the remaining fold, iterating through all 4 possible combinations.

C. Metrics

Our evaluation of object detection models hinges on the Average Precision (AP) metric, which is computed subsequent to establishing a cut-off point for the Intersection over Union (IoU) ratio.

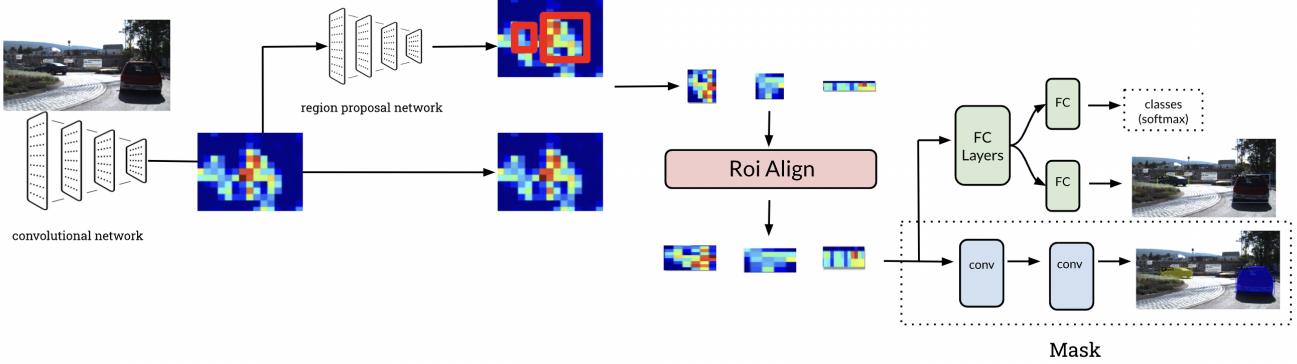


Fig. 3. Diagram for Faster R-CNN and Mask R-CNN extension



Fig. 4. Example image with annotations within the KITTI-MOTS dataset. There are some cars annotated with the confidence (in %) being each one painted with a different color.

The IoU quantifies the extent of overlap between two bounding boxes. Specifically, within the context of object detection, it assesses the degree of correspondence between the predicted bounding box and the actual, or ground truth, bounding box. To this end, the IoU is determined for every prediction, and a predetermined threshold is applied to categorize the prediction as either a true positive or a false positive.

$$IoU = \frac{\text{area of overlap}}{\text{area of union}} \quad (1)$$

For AP calculation, we construct the precision-recall graph based on a chosen IoU threshold. Although the formal definition of AP is the integral of this curve, it is often approximated for practicality. Our methodology adopts the COCO evaluator's approximation approach, which entails interpolating the precision across 101 evenly spaced recall levels after enforcing a monotonic decrease in precision to the right. AP can be computed for various IoU levels or across different scales of bounding boxes. The aggregate AP for the COCO dataset is determined by taking the mean of the APs calculated at IoU thresholds ranging from 0.5 to 0.95 in increments of 0.05.

All the experiments have been done with the default optimizer hyperparameters. The finetunne has been performed with the Adam optimizer with a batch size of 16 images per



Fig. 5. Example of an augmented image in which a car is miss-classified as a snowboard.



Fig. 6. Example of the classification of a boat which not only depends on the boat but more importantly in the context, in this case, the water.

batch and 200 epochs. The experiments has been done over a GPU RTX 3090 with 24G.

V. RESULTS

VI. CONCLUSIONS

REFERENCES

- [1] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick, "Detectron2," *arXiv preprint arXiv:1911.02549*, 2019.
- [2] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," *arXiv preprint arXiv:1703.06870*, 2017.
- [3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
- [4] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 779–788, 2016.
- [5] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proceedings of the 2001 IEEE computer society conference on computer vision and pattern recognition. CVPR 2001*, vol. 1, pp. I–I, Ieee, 2001.

- [6] P. Viola and M. J. Jones, "Robust real-time face detection," *International journal of computer vision*, vol. 57, pp. 137–154, 2004.
- [7] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, vol. 1, pp. 886–893, Ieee, 2005.
- [8] P. Felzenszwalb, D. McAllester, and D. Ramanan, "A discriminatively trained, multiscale, deformable part model," in *2008 IEEE conference on computer vision and pattern recognition*, pp. 1–8, Ieee, 2008.
- [9] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE transactions on pattern analysis and machine intelligence*, vol. 32, no. 9, pp. 1627–1645, 2009.
- [10] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 580–587, 2014.
- [11] J. Uijlings, K. Sande, T. Gevers, and A. Smeulders, "Selective search for object recognition," *International Journal of Computer Vision*, vol. 104, pp. 154–171, 09 2013.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [13] R. Girshick, "Fast r-cnn," in *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1440–1448, 2015.
- [14] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, 06 2015.
- [15] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," *CoRR*, vol. abs/1703.06870, 2017.
- [16] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 779–788, 2016.
- [17] J. Redmon and A. Farhadi, "Yolo9000: Better, faster, stronger," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6517–6525, 2017.
- [18] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," 04 2018.
- [19] A. Bochkovskiy, C.-Y. Wang, and H.-y. Liao, "Yolov4: Optimal speed and accuracy of object detection," 04 2020.
- [20] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. Berg, "Ssd: Single shot multibox detector," vol. 9905, pp. 21–37, 10 2016.
- [21] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv 1409.1556*, 09 2014.
- [22] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2999–3007, 2017.
- [23] T. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 936–944, 2017.
- [24] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Zitnick, "Microsoft coco: Common objects in context," 05 2014.
- [25] M. J. Choi, A. Torralba, and A. S. Willsky, "A tree-based context model for object recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 2, pp. 240–252, 2012.