

# Object Detection, Segmentation & Counting Using Deep Learning

Nandini N<sup>1</sup>, Nandini C S<sup>2</sup>, Dr. K R Nataraj<sup>3</sup>

<sup>1</sup>Mtech Digital Communication & Networking, SJBIT, Bangalore

<sup>2</sup>Jounier Data Scientist, Curl Analytics, Bangalore,

<sup>3</sup>Head of the Department, ECE, SJB Institute of Technology, Bangalore

\*\*\*

**Abstract** - Object counting is a not easy task in image processing. It is usually agreed in different areas of industries, research institutes, laboratories, agriculture industries among others. Accurately counting objects instances in a given picture or video frame is a difficult problem to solve in machine learning. Identifying the number of objects present in the image can be helpful for extra investigation in a spacious set of applications. In this project we propose a simple method for automatically detect the object, segmenting by using pixel wise mask and determining the number of objects in an image. This approach competently detects physical object in associate with the image at the same time creating the good quality partition mask for each occurrence of object and count the number of object detected. We use the Mask R-CNN method to detect and generate instance mask for each object. This method is straightforward for training and gives a little transparency of Faster R-CNN, operated by 5 fps.

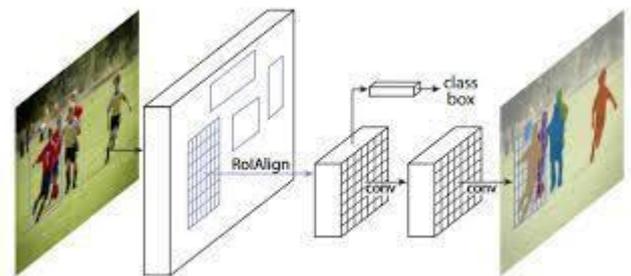
**Key Words:** Mask R-CNN, Object detection, Instance Segmentation, Object counting

## 1. INTRODUCTION

The vision group of people has fast enhanced object recognition and semantic division fallout all over a small duration of the moment. In huge section of these in advance has determined aside strong touchstone methods, for example the Faster RCNN [34] and Fully Convolutional Network (FCN) [29] structures for target recognition and also for semantic division, correspondingly. These types are theoretically sensitive, propose whippiness and strength, mutually with the quick preparation and inference point. The idea is to build up the comparably enabling structure for instance division.

Instance division are tough since this technology needs accurate perception of all the objects present in a picture while also accurately segments every occurrence. For that reason merges essentials by the traditional computer vision responsibilities of physical object sensing, the purpose is to knowing the categorise each physical entity's and localise to each one by means of bounding box. Then next is to use the semantic sectionalization to categorise every picture element into a assured equipoise collection and also performs lack off differentiating physical object occurrence. Certain this might be a difficult process to reach good outcomes. But, we illustrate that an unexpectedly straightforward, flexible, and high-speed system that can exceed previous state-of-the-art occurrence segmentation outcome.

The technique, described as Mask R-CNN, which is broaden of Faster R-CNN [34] by addition of branch for estimating segmentation masks on every Region of Interest (RoI), in corresponding with the offered branch for classification and bounding box regression with Count of the object(Fig- 1). The mask subdivision ca obtained by a tiny FCN useful to every RoI, generating a partition cover in a pixel-to-pixel method. Mask RCNN is straightforward to execute and train than the Faster R-CNN structure, which help the extensive range of flexible structural designs. Furthermore, the cover subdivision only adds a little computational overhead, allows a high-speed structure and quick testing.



**Fig -1:** The Mask RCNN model for performs instance division

In standard Mask RCNN was a discerning porch of Faster RCNN, so far making a mask division correctly is important for great outcomes. Almost, mask RCNN was intended for pixel to pixel association in between system inputs & yield [41]. This is just about how plain RoI-Pool is [18, 12], the factual set of process to the presence of instances, performs crude attribute division for characteristic dections. To identify the misplacement, we suggest the easy, quantization free level, named RoI Align, which authentically conserve accurate spatial positions. In spite of being a apparently slight alteration, RoIAlign has huge wallop, it develops mask precision by proportional 10 to 50 percentage, viewing better profit beneath stricter localization poetics. Subsequent, we establish it important to separate mask & class calculation: we prognosticate a twofold mask to every classes separately, not including struggle to among all the classes, & trust the connected network's RoI categorization branch for expecting collection. In distinguish, FCNs generally make per pixel multi class classification, which connect sectionalization and categorization, and founded on this research works badly for occurrence sectionalization.

Lacking doorbell and whistles, Mask RCNN exceed the all earlier progressive one-on-one model outcome on the COCO occurrence sectionalization job [28], as well as a lot of engineered access from the year 2016 contest achiever. As our technique besides stand out on COCO physical object perception job. In ablation research, we estimate numerous necessary instantiations which permit to show the robustness to examine the effects of core aspects. This model can excite at speed of 200ms/frame using a GPU system.

At last, we proposals the generalization this program via mission of human-like judgment on the COCO fundamental component present in dataset [28]. By presentating each key position as one hot twofold mask with least variation Mask RCNN that can be useful to find instance peculiar poses. Mask RCNN better & the champion of 2016 COCO contest and by the aforesaid moment runs at five frame/sec. Mask RCNN has a flexible structure for instance level identification & might be eagerly prolonged to the most difficult work.

## 2. RELATED WORK

**RCNN:** The Region based CNN (RCNN) uses [13] a bounding box object categorisation is to focus to controllable count of applicant object areas [42, 20] & calculate convolutional system [25, 24] separately to every RoI. RCNN was comprehensive [18, 12] & agree to attending RoIs into characteristic maps by use of RoIPool, important to the fast rate & better precision. Faster RCNN is superior to this flowed by knowing the significance method called Region Proposal Network. Faster RCNN is lithe & forceful to upgrading (e.g., [38, 27, 21]), & the present most important framework in numerous standard.

**Instance Sectionization:** Motivated by the success of RCNN, different ways to deal with instance division depend on portion plot. Previous technique [13, 15, 16, 9] resorted to base up segments [42, 2]. DeepMask [33] & subsequent works [34, 8] find out to offer segment candidates, which as categorized by Fast RCNN. In those technique segmentation leads identification, & it is slow and fewer precise. Similarly, Dai et al. [10] projected a multifaceted multiple step flow that can determins segment scheme by bounding box proposals, after that categorization. As an alternative, our process is depends on the equivalent identification by masks and category tag, that are frank & more bendable.

Mainly, Li et al. [26] pooled the section proposal system [8] & object recognition system [11] for "fully convolutional instance segmentation". The frequent plan in [8, 11 & 26] is to forecast set of location receptive output channels by means of fully convolutionally methode. These channels concurrently deal with the object category, boxes, & masks, creating a system with high-speed. But FCIS shows methodical errors in infringement instances & makes false edges, viewing that is braved by essential difficulties in segmenting instances.

One more relatives solutions [23, 4, 3, 29] to instance segmentation are determined by the achievement of semantic segmentation. Initial from per pixel categorization fallout, these schemes try to cut pixels of smilar group within different instances. In distinguish the segmentation the first approach is Mask RCNN which is derived from an instance first approach. We look ahead to deeper integration of both strategies is to be considered in the future.

## 3. MASK RCNN

Mask RCNN is theoretically straightforward. Faster RCNN consist of two outcomes to every aspirant object, a category symbol & a bounding box compensates, to do this we include third division that outputs an object mask. Mask RCNN has a usual and perceptive tips. But the supplementary mask yeild is diverse by class & box yeilds, involving withdrawal of superior spatial formation to an object. Next we initiate the type of Mask RCNN, as well as pixel-pixel line-up, that is the major mislaid piece to Fast/Faster R-CNN.

**Faster RCNN:** We start with the in brief reviewing of the Faster RCNN detector [36]. Faster RCNN include the two phases. The first phase, named a Region Proposal Network, offers aspirant object/target bounding boxes. Second phase is Fast RCNN [12], takes out features via RoIPool by use of every candidate box & carry out classification & bounding box regression. The features worn by both phases can be mutual for faster implication. We submit person who reads to [21] for newest, complete evaluations to Faster RCNN & other frameworks.

**Mask RCNN:** Mask RCNN acknowledges a similar two-stage process, with an indistinct first stage which is RPN. In the second stage, in comparing to anticipating the class and box balance, Mask RCNN additionally yields a twofold cover for every rous. This is in recognizing to most current strategies, where orders rely upon cover expectations (e.g. [33, 10, 26]). Our strategy tracks the quality of Fast R-CNN that relates jumping enclose game plan and relapse parallel.

Officially, throughout training, we uses multi task loss to every illustration ROI as  $L = L_{cls} + L_{box} + L_{mask}$ .

categorization loss  $L_{cls}$  & bounding box loss  $L_{box}$  are the same as those distinct in [12]. The mask bough has a  $Km^2$

- dimensional yeild for every ROI, to encodes  $K$  twofold masks of pledge  $m \times m$ , for every  $K$  classes. To do this we concern a per pixel sigmoid, & define  $L_{mask}$  as the standard binary crossentropy loss. For a ROI related with groundtruth class  $k$ ,  $L_{mask}$  is simply defined on  $k$ -th mask.

The  $L_{mask}$  is characterized as to interface the framework to make the covers for each class without rivalry among classes; we depend on the committed grouping branch to anticipate the class mark used to choose the yeild veil. This decouples cover and class forecast. This is not quite the same

as basic practice while applying FCNs [30] to semantic division, which normally utilizes a for every pixel softmax and a multinomial cross-entropy misfortune. All things considered, veils crosswise over classes contend; for our situation, with a for every pixel sigmoid and a twofold misfortune, they don't. We appear by tests that this definition is key for good case division comes about.

**Mask Representation:** A cover/mask encodes an information protest's/object spatial design. Consequently, dissimilar to class names or box counterbalances that are unavoidably fell into short yield vectors by completely associated (fc) layers, removing the spatial structure of covers can be tended to normally by the pixel-to-pixel correspondence gave by convolutions. In particular, we foresee a  $m \times m$  cover from every rous utilizing a FCN [30]. This permits each layer in the veil branch to keep up the unequivocal  $m \times m$  protest spatial format without crumbling it into a vector portrayal that needs spatial measurements. Dissimilar to past techniques that fall back on fc layers for cover forecast [33, 34, 10], our completely convolutional portrayal requires less parameters, and is more precise as showed by tests. This pixel-to-pixel conduct requires our RoI highlights, which they are have little element maps must be fit to sensibly secure the exact per-pixel spatial association. This urged us to develop the consequent RoIAlign store layer that performance center a contribution to veil count.

**RoIAlign:** RoIPool [12] is normal function used for pull out a petite feature map (e.g.,  $7 \times 7$ ) starting to every RoI. RoIPool at first quantizes a coasting number RoI to the unmistakable granularity of the component graph, this quantizing RoI is a while later keep on sub-partitioned the spatial canisters which are then quantized, and finally highlight esteems encased by everybody case are amassed. Then the quantization was performed, on a uninterrupted match of  $x$  by calculating  $\lfloor x/16 \rfloor$  where sixteen is the feature chart step and is rounding; the same way, quantization is carry out when isolating into bins (e.g.,  $7 \times 7$ ). These quantization's begin with the misalignments in the RoI and the pull out the features. At the same time this could not crash categorization, which is tough to petite translations, it have a bulky negative cause on estimating pixel-accurate cover. To deal with this, we suggest the RoIAlign deposit layer that eliminate the insensitive quantization of RoIPool and accurately arrange in a line to the pulled out features throughout the input. Our projected changes are simple: we pass up any quantization of the RoI limitations or bins. We make use of bilinear interpolation [22] to figure out the accurate values of the incoming features at four frequently sampled positions in every RoI case, and combined the effect, we make a note of the outcome are not responsive to the accurate sampling positions, or how many positions be sampled, provided that no quantization is achieved.

RoIAlign show the way to big developments. We furthermore evaluate to the RoIWarp process projected in Distinct to RoIAlign, RoIWarp ignored the position problem and was realized in [10] seeing that quantizing RoI at

present similar to RoIPool. Subsequently even though RoIWarp also implements bilinear resampling provoked by [22], it carry out on average with RoIPool as exposed by experimentsts, representative the critical role of arrangement.

**Network Architecture:** To express, overview of our approach is that we represent the Mask R-CNN by various architectures. The intelligibility, we distinguish among: (i) the convolutional networks style used for feature taking out more an whole picture, and (ii) the system beginning for bounding-box detection (grouping and regression) and mask calculation i.e. applied independently to every ROI. We signify that the backbone structural design gives the classification system-depth-features.

[19] We estimate ResNet [19] as well as ResNeXt [45] set-up of strength of 50 or 101 layers. As the unusual accomplishment of the Faster R-CNN by way of ResNets extort the features starting from the last convolutional layer of the 4th phase, as we describe C4. This backbone by means of ResNet-50, such as, is denoting by the ResNet-50-C4. This was a frequent option utilizes within [19, 10, 21, 39].

[20] We furthermore investigate one great efficient anchor recently projected at Lin et al. [27], named as Characteristic Pyramid System (FPN). FPN utilizes the up and down structural design by means of tangential relations to make in an system attribute pyramid starting with one input. Faster R-CNN among an FPN anchor pull out the ROI features from unusual levels of the feature pyramid reported to their range, other than the remains of the coming be comparable toward vanilla ResNet. By means of a ResNet-FPN is backbone/system intended for feature action by means of Mask RCNN to gives outstanding gains in both accuracy and momentum.

Intended which is the network head we intimately go behind the architectures existing in prior to job then we have to insert a fully conventional pretense calculation branch. Purposely, After pull out the fast/Faster R-CNN envelope heads before, first starting with the ResNet [19] and FPN [27] documents. Information is as shown in Fig- The top on the ResNet-C4 backbone contains the fifth step of, which is computed intensive. Designed FPN, is the support previously consists of res5 and those permits for a additional most efficient skull that requires smaller number of devices/filters.

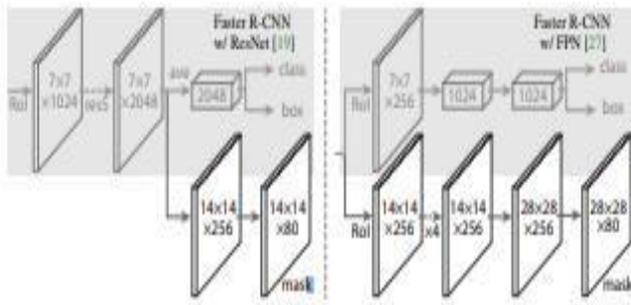


Fig -2: Head Architecture

We collected our mask branch contains the clear-cut constitution. More difficult develops include the possible toward improve the performance except that not the focal point of this effort.

4. IMPLEMENTATION DETAILS

We ready frenzied constraints subsequent to active Fast or Faster R-CNN job [12, 36, 27]. Even though these results be done intended for objective finding in primary credentials [12, 36, 27], we grow our instance cell division/segmentation methods are very well strong to them.

**Training details:** Seeing that in Fast R-CNN, a RoI is well thought-out constructive and it has an IoU through a land-truths boxes of the smallest amount of 0.5 or else antagonistic or else. If the positive RoIs defines the mask loss Lmask. The mask object is the meeting point that connecting to a RoI and its related argument precision mask.

We taken the picture-centric were training [12] in that Images to be resized with the intention of their size is eight hundred pixels [27]. Every minimum-batch has two pictures per GPU and everyone have N number of sampled ROI's, by the ratio is 1:3 of constructive to antagonistic

[12]. N be 64 designed of the C4 anchor (i.e., in [12, 36]) and 512 used for FPN ([27]). We prepare training on 8

GPU's intended for 160k repeat, with a acquisition speed of 0.02 which is reduced of ten at 120k iterations. We are make use of a weight decompose of 0.0001 and force of 0.9. With ResNeXt [45], we trained by 1 picture/image for each GPU and as similar amount of ingeminate, with a preliminary acquisition/learning speed is 0.01.

5. MAIN RESULTS

We evaluate Mask R-CNN to the progressive technique inside instance cell division/segmentation and also gives the number of item present in an image. Each and every observe of our representation break standard modifications of earlier progressive methods. That contains MNC [10] and FCIS [26],

the champions of the COCO 2015 as well as 2016 cell deviation/ segmentation computations, correspondingly. Mask R-CNN among ResNet-101-FPN backbone column exceed FCIS+++ [26], which contains multi-scale train/trial, the crosswise flip test and the online awkward example mining (OHem)

[38]. Mask R-CNN get good quality of outcome even below challenging situations.

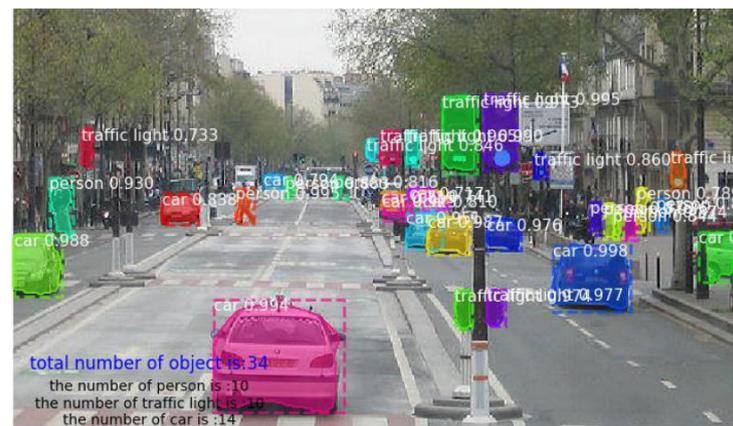


Fig -2: Outcomes of the project

6. CONCLUSIONS

This project formulated, developed and evaluated two versions of an instance segmentation algorithm. The first version used the output from an object detector CNN and a mask proposal subsystem to generate the instances, whereas the second one, to which three deferent possible architectures were devised, used also the output from a semantic segmenter to obtain initial guesses for the masks of each instance. Then finally calculate the count of each detected object class.

REFERENCES

[1] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2D human pose estimation: New benchmark and state of the art analysis. In CVPR, 2014. 8

- [2] P. Arbelaez, J. Pont-Tuset, J. T. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping. In CVPR, 2014. 2
- [3] A. Arnab and P. H. Torr. Pixelwise instance segmentation with a dynamically instantiated network. In CVPR, 2017. 3, 9
- [4] M. Bai and R. Urtasun. Deep watershed transform for instance segmentation. In CVPR, 2017. 3, 9
- [5] S. Bell, C. L. Zitnick, K. Bala, and R. Girshick. Insideoutside net: Detecting objects in context with skip pooling and recurrent neural networks. In CVPR, 2016. 5
- [6] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multiperson 2d pose estimation using part affinity fields. In CVPR, 2017. 7, 8
- [7] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The Cityscapes dataset for semantic urban scene understanding. In CVPR, 2016. 9
- [8] J. Dai, K. He, Y. Li, S. Ren, and J. Sun. Instance-sensitive fully convolutional networks. In ECCV, 2016. 2
- [9] J. Dai, K. He, and J. Sun. Convolutional feature masking for joint object and stuff segmentation. In CVPR, 2015. 2
- [10] J. Dai, K. He, and J. Sun. Instance-aware semantic segmentation via multi-task network cascades. In CVPR, 2016. 2, 3, 4, 5, 6
- [11] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: Object detection via region-based fully convolutional networks. In NIPS, 2016. 2
- [12] R. Girshick. Fast R-CNN. In ICCV, 2015. 1, 2, 3, 4, 6
- [13] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In CVPR, 2014. 2, 3
- [14] R. Girshick, F. Iandola, T. Darrell, and J. Malik. Deformable part models are convolutional neural networks. In CVPR, 2015. 4
- [15] B. Hariharan, P. Arbelaez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In ECCV, 2014. 2
- [16] B. Hariharan, P. Arbelaez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In CVPR, 2015. 2
- [17] Z. Hayder, X. He, and M. Salzmann. Shape-aware instance segmentation. In CVPR, 2017. 9
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In ECCV, 2014. 1, 2
- [19] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In CVPR, 2016. 2, 4, 7, 10
- [20] J. Hosang, R. Benenson, P. Dollar, and B. Schiele. What makes for effective detection proposals? PAMI, 2015. 2
- [21] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, et al. Speed/accuracy trade-offs for modern convolutional object detectors. In CVPR, 2017. 2, 3, 4, 6, 7
- [22] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial transformer networks. In NIPS, 2015. 4
- [23] A. Kirillov, E. Levinkov, B. Andres, B. Savchynskyy, and C. Rother. Instancecut: from edges to instances with multicut. In CVPR, 2017. 3, 9
- [24] A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. In NIPS, 2012. 2
- [25] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. Neural computation, 1989. 2
- [26] Y. Li, H. Qi, J. Dai, X. Ji, and Y. Wei. Fully convolutional instance-aware semantic segmentation. In CVPR, 2017. 2, 3, 5, 6
- [27] T.-Y. Lin, P. Dollar, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In CVPR, 2017. 2, 4, 5, 7
- [28] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and C. L. Zitnick. Microsoft COCO: Common objects in context. In ECCV, 2014. 2, 5
- [29] S. Liu, J. Jia, S. Fidler, and R. Urtasun. SGN: Sequential grouping networks for instance segmentation. In ICCV, 2017. 3, 9
- [30] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In CVPR, 2015. 1, 3, 6
- [31] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In ICML, 2010. 4

- [32] G. Papandreou, T. Zhu, N. Kanazawa, A. Toshev, J. Tompson, C. Bregler, and K. Murphy. Towards accurate multiperson pose estimation in the wild. In CVPR, 2017. 8
- [33] P. O. Pinheiro, R. Collobert, and P. Dollar. Learning to segment object candidates. In NIPS, 2015. 2, 3
- [34] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollar. Learning to refine object segments. In ECCV, 2016. 2, 3
- [35] I. Radosavovic, P. Dollar, R. Girshick, G. Gkioxari, and K. He. Data distillation: Towards omni-supervised learning. arXiv:1712.04440, 2017. 10
- [36] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In NIPS, 2015. 1, 2, 3, 4, 7
- [37] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In TPAMI, 2017. 10
- [38] A. Shrivastava, A. Gupta, and R. Girshick. Training regionbased object detectors with online hard example mining. In CVPR, 2016. 2, 5
- [39] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta. Beyond skip connections: Top-down modulation for object detection. arXiv:1612.06851, 2016. 4, 7
- [40] C. Sun, A. Shrivastava, S. Singh, and A. Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In ICCV, 2017. 10