

# Explicit Shape Encoding for Real-Time Instance Segmentation

Wenqiang Xu\*

Haiyang Wang\*

Fubo Qi

Cewu Lu<sup>†‡</sup>

Department of Computer Science and Engineering  
Shanghai Jiao Tong University

{vinjohn, wanghaiyang, 727749815, lucewu}@sjtu.edu.cn

## Abstract

In this paper, we propose a novel top-down instance segmentation framework based on explicit shape encoding, named **ESE-Seg**. It largely reduces the computational consumption of the instance segmentation by explicitly decoding the multiple object shapes with tensor operations, thus performs the instance segmentation at almost the same speed as the object detection. ESE-Seg is based on a novel shape signature Inner-center Radius (IR), Chebyshev polynomial fitting and the strong modern object detectors. ESE-Seg with YOLOv3 outperforms the Mask R-CNN on Pascal VOC 2012 at mAP<sup>r</sup> @0.5 while 7 times faster.

## 1. Introduction

Instance segmentation is a fundamental task in the computer vision, which is important for many real-world applications such as autonomous driving, robot manipulation. As the task seeks to predict both the object location and the shape, the methods for the instance segmentation are generally not as efficient as the object detection frameworks. Forwarding each object instance through an upsampling network to obtain the instance shape, as mainstream instance segmentation frameworks do [12, 22, 3, 19], is quite computation-consuming, especially when compared with the object detection which only needs to regress the bounding box, *i.e.* a 4D vector for each object. Thus, if the network can also regress the object shape to a short vector, and decode the vector to the shape (see Fig. 1) in a simple way just like the bounding box, it can make the instance segmentation reach *almost* equal computational efficiency to the object detection. To achieve this goal, we propose a novel instance segmentation framework based on **Explicit**

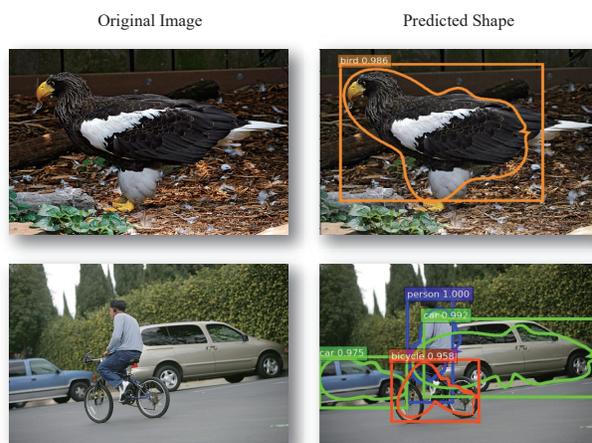


Figure 1. ESE-Seg learns to estimate the shapes of the detected objects, it can be simultaneously obtained along with the bounding boxes.

**Shape Encoding** and modern object detectors, named **ESE-Seg**.

Shape encoding is originally developed for instance retrieval [39, 17, 37], which encodes the object to a shape vector. Recently, a number of works encode the shape implicitly [9, 29, 38], which is to project the shape content to a latent vector, typically through a black-box design such as deep CNN. Thus the decoding procedure under this approach should be also put through a network, which requires several forwarding for multiple instances, and causes large computation. In pursuit of fast decoding, we employ an explicit shape encoding that involves only simple numeric transformations.

However, designing a satisfactory explicit shape encoding method is non-trivial. Concerning the CNN training, as it is known to regress with uncertainties, a preferred shape vector should be relatively **short** but contains sufficient information, **robust** to the noise, and **efficiently de-**

\*these two authors have equal contributions.

<sup>†</sup>Cewu Lu is the corresponding author.

<sup>‡</sup>Cewu Lu is a member of MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, and SJTU-SenseTime AI lab.

**codable** to reconstruct the shape. In this paper, we propose a contour-based shape signature to meet these requirements. A novel “Inner-center Radius” (IR) shape signature for instance shape representation is introduced. The IR first locates an inner-center inside the object segment, and based on this inner-center, it transforms the contour points to polar coordinates. That is, we can form a function of radius  $f(\theta)$  along the contour with respect to angle  $\theta$ . To make the shape vector even shorter and more robust, we apply the Chebyshev polynomials to conduct the function approximation on  $f(\theta)$ . As such, the IR signature is represented by a small number of coefficients with small error, and these coefficients are the shape vector to be predicted. Additionally, we also in-depth discuss about the comparison with other shape signature designs. Conventional object detector (e.g. YOLOv3 [31]) is used to regress the shape vector, along with 4D bounding box vector. To note that our shape decoding can be implemented by simple tensor operations (multiplication and addition) which are extremely fast.

The ESE-Seg itself is independent of all the bounding box-based object detection frameworks [32, 4, 8, 20, 23]. We demonstrate the generality on Faster R-CNN [32], RetinaNet [20], YOLO [30] and YOLOv3-tiny [31] and evaluate our ESE-Seg on standard public datasets, namely Pascal VOC [6] and COCO [21]. Our method achieves 69.3 mAP<sup>r</sup>, 48.7 mAP respectively with IOU threshold 0.5. The score is better than Mask R-CNN [12] on Pascal VOC 2012, and is competitive to the performance on COCO. It is decent considering it is 7 times faster than Mask R-CNN with the same backbone ResNet-50 [13]. The speed can be even faster at  $\sim 130$ fps on GTX 1080Ti when the base detector changes to YOLOv3-tiny, while the mAP<sup>r</sup>@0.5 remains 53.2% on the Pascal VOC. It is noteworthy, ESE-Seg speeds up the instance segmentation not depending on the model acceleration techniques [15, 40], but relying on a new mechanism that cut down shape prediction after object detection.

**Contributions.** We propose an explicit shape encoding based instance segmentation framework, ESE-Seg. It is a top-down approach but reconstructs the shapes for multiple instances in one pass, thus greatly reduces the computational consumption, and makes the instance segmentation reach the speed of the object detection with no model acceleration techniques involved.

## 2. Related Work

**Explicit vs Implicit Shape Representation** A previous work with similar ideology has been done by Jetley *et al.* [16]. They took the implicit shape representation path by first training an autoencoder on object binary mask. The encoded shape vector is decoded to shape mask through the decoder component. In the implementation, they adopted the YOLO [30] to regress the bounding box and the shape vector for each detected object. The YOLO structure can

thus be viewed as both detector and encoder. The encoded vector from YOLO is then decoded by the pre-trained denoising autoencoder. The major differences between our work and theirs:

- Explicit shape representation is typically based on the contour, while implicit shape representation is typically based on the mask.
- Explicit shape representation requires no additional decoder network training. Parallelizing the decoding process for all objects in the images, which is hard for network structured decoder, can be easily achieved by the explicit shape encoding. As a matter of fact, implicit decoding requires multiple passes for multiple objects, one for each, while explicit decoding can obtain all the shapes in one pass.
- The input for training autoencoder and training YOLO (viewed as an encoder) is quite different (object scales, color pattern), which may cause trouble for the decoder, since the decoder is not further optimized with YOLO training. Such an issue does not exist for explicit shape representation.

In addition to our proposed IR shape signature, there exist various methods to represent the shape, to name a few, centroid radius, complex coordinates, cumulative angle [5, 34, 39] *etc.* While such methods sample the shape related feature along the contour, only a few of them can be decoded to reconstruct the shape.

**Object detection** Object detection is a richly studied field. Object detection frameworks with CNN can be roughly divided into two categories, one-stage and multi-stage. Two-stage detection scheme is a classic multi-stage scheme, which typically learns an RPN to sample region proposals and then refine the detection with roi pooling or its variations, the representative works are Faster R-CNN [32], R-FCN [4]. Recently, some works extend the two-stage to multi-stage in a cascade form [1]. On the other hand, one-stage detectors divide the input image to size-fixed grid cells and parallelize the detection on each cell with fully convolutional operations, the representative networks are SSD [23], YOLO [30], RetinaNet [20]. Recently, point-based detections are proposed, CornerNet [18] directly detects the upper-left and bottom-right points, which is a one-stage detector. Grid R-CNN [24] regresses 9 points to construct the bounding box, which is a two-stage detector. Our method is compatible with all the bounding box-based detection networks. We experiment with Faster R-CNN, YOLO, YOLOv3, and RetinaNet to prove the generality. See Table 4. However, it is not compatible with the point-based detector, as the shape (bounding box) in this setting is not parametrized.

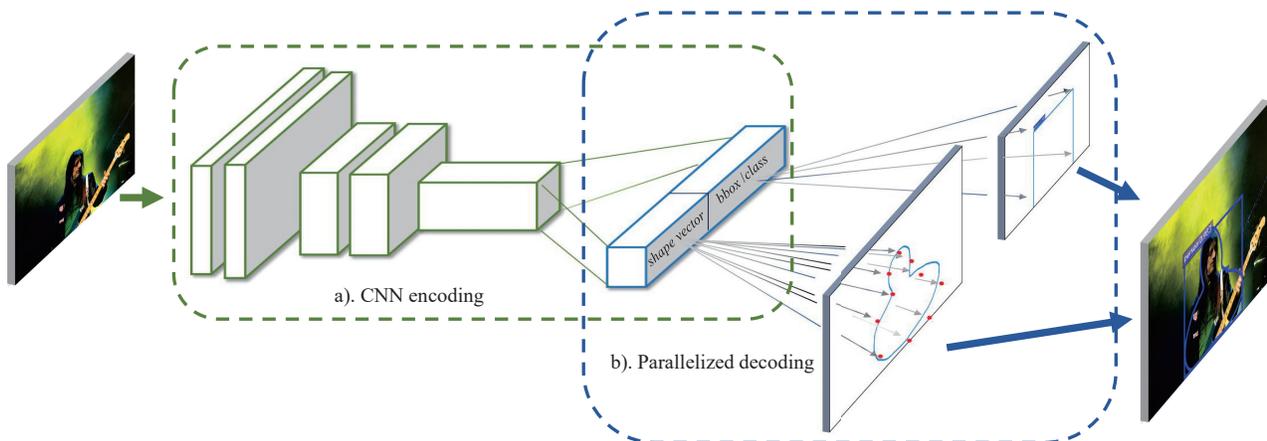


Figure 2. The pipeline of the shape detection, regression and reconstruction.

**Instance Segmentation** Instance segmentation requires not only to locate the object instance but also to delineate the shape. The mainstream methods can be roughly divided to top-down [12, 22, 3, 19, 28, 27, 2] or bottom-up [26, 35] approaches. Ours belongs to the top-down line. The top-down approaches such as MNC [3], FCIS [19], Mask R-CNN [12] are generally slowed down when the object number in an image is large, as they predict the instance mask in sequence. On the contrary, our ESE-Seg alleviates the cumbersome computation by regressing the object shapes to short vectors and decoding them simultaneously. It is also the first top-down instance segmentation framework which is not affected by the instance number in the images with respect to the inference time. Besides, the works on augmenting the performance of instance segmentation frameworks through data augmentation [7, 36], scale normalization [33] can be easily integrated to our system.

## 3. Method

### 3.1. Overview

We propose an explicit shape encoding based detection to solve the instance segmentation. It predicts all the instance segments in one forwarding pass, which can reach equal efficiency as object detection solver. Given an object instance segment, we parametrize the contour with a novel shape signature “Inner-center Radius” (IR) (Sec. 3.2.1). The Chebyshev polynomials are used to approximate the shape signature vector with a small number of coefficients (Sec. 3.2.2). Those coefficients are served as the shape descriptor, and the network will learn to regress it. (Sec. 3.3). Finally, we describe how to decode the shape descriptor under the ordinary object detection framework by simple ten-

sor operations. (Sec. 3.4). The overall pipeline is shown in Fig. 2.

**The Advantage of Explicit Shape Encoding** In object detection system (*e.g.* YOLOv3), the network regresses the bounding boxes (*i.e.* 4D vectors) and the bounding box is decoded by tensor operations, which is light to process and easy to parallelize. By contrast, conventional instance segmentation (*e.g.* Mask R-CNN) requires an add-on network structure to compute the object shape. The decoding/upsampling forwarding involves a large number of parameters, which is heavy to load in parallel for multiple instances. This is why instance segmentation is normally much slower than object detection. Therefore, if we also regress the object shape into short vectors directly, the instance shape decoding can be achieved by fast tensor operations (multiplication and addition) in a similar way. Thus the instance segmentation can reach the speed of object detection.

### 3.2. Shape Signature

#### 3.2.1 Inner-center Radius Shape Signature

In this section, we will describe the design of the “inner-center radius” shape signature and compare it to previously proposed shape signatures.

The construction of the “inner-center radius” contains two steps: First, locate an Inner center point inside the object segment as the origin point to build the polar coordinate system. Second, sampling the contour points according to the angle  $\theta$ . This signature is translation-invariant and scale-invariant after normalized.

**Inner center** The inner-center point is defined by the most far-way point from the contour, which can be obtained through distance transform [25]. To note, some commonly used center such as the center of mass, the center of the bounding box cannot guarantee to be inside the object. See Fig. 3.

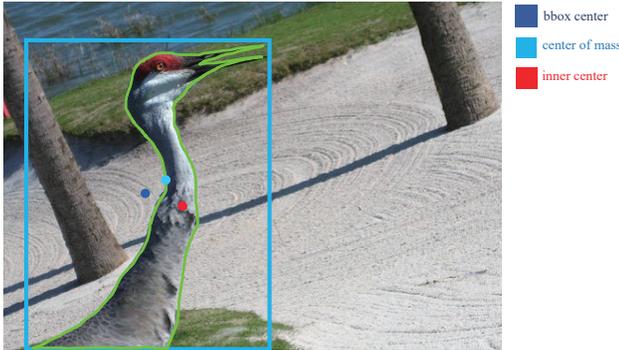


Figure 3. The center points of an object. As we can see, bounding box center and the center of mass cannot guarantee to be inside an object.

In a few cases, an object is separated into disconnected regions, resulting in multiple inner centers. To deal with such situations, we dilate the broken areas to a single one and then find the contour of the dilated shape. Of course the contour is very rough, however, it can help to reorder the contour points of the outline points. The whole process is depicted in Fig. 4. Thus inner center is computed from the completed contour.

**Dense Contour Sampling** We sample the contour points according to the angles at the interval of  $\tau$  around inner-center point, thus a contour will result in  $N = \lceil 2\pi/\tau \rceil$  points. In practice,  $\tau = \pi/180$  and thus  $N = 360$  points are sampled from an object contour. If the ray casting from the inner-center intersects more than once to the contour. We collect the point with the largest radius only. The function  $f(\theta)$  is denoted as radius at different angles  $\theta$ . To note, we are aware that the contour sampling in this way will not be perfect, however, after extensive experiments in Pascal VOC, and COCO, we find it suitable for natural objects (see Table 2). A further discussion is in the next Sec. 3.2.3.

### 3.2.2 Fitting the Signature to Coefficients

The IR makes shape representation into a vector. But, it is still too long for the network to train. Besides, the shape signature is very sensitive to the noise (see Fig. 7). Thus, we take a further step to shorten shape vector and resist noise through Chebyshev polynomial fitting.

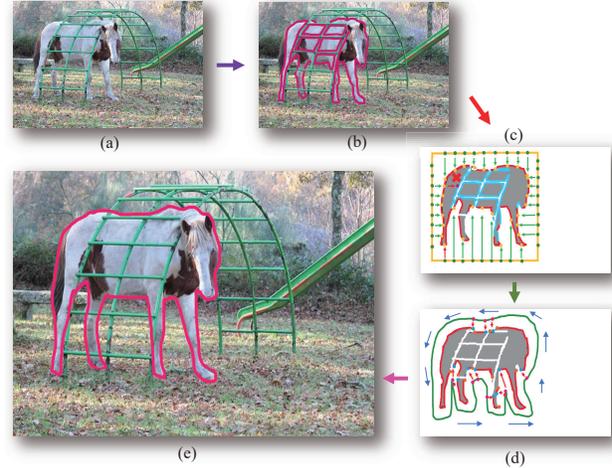


Figure 4. The process to complete the separated areas. An occluded object (a) has many separated areas (b), we split the contour points of each area into outline and inner points with the help of the bounding box (c), then we dilate the broken area into one, and reorder the outline points according to the dilated shape contour (d), finally, we complete the instance (e).

**Chebyshev polynomials** The Chebyshev polynomial is defined in recurrence:

$$T_0(x) = 1, \quad (1)$$

$$T_1(x) = x, \quad (2)$$

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x), \quad (3)$$

which is also known as **The Chebyshev polynomials of the first kind**. It can effectively minimize the problem of *Runge's phenomenon* and provides a near-optimal approximation under the maximum norm<sup>1</sup>.

Given the IR shape signature, the Chebyshev approximation is to find the coefficients in

$$f(\theta) \sim \sum_{i=0}^{\infty} c_i T_i(\theta)$$

Truncating the function with  $n$  terms, we have the approximation function  $\tilde{f}(\theta) = \sum_{i=0}^n c_i T_i(\theta)$ .  $\mathbf{k} = (c_0, \dots, c_n)$  are the shape signature vector to represent the object.

### 3.2.3 Discussion

**Comparison with Other Shape Signatures** The angle-based sampling for shape signature such as proposed IR is rarely adopted before, because it cannot perfectly fit shape segment. Actually, we compare and in-depth analyze other shape signatures and finally choose this solution. For example, a quite straight-forward design is to sample along the

<sup>1</sup>[https://en.wikipedia.org/wiki/Chebyshev\\_polynomials](https://en.wikipedia.org/wiki/Chebyshev_polynomials)

contour. The contour is represented by a set of contour polygon vertex coordinates. This method can nearly perfectly fit the object segment, especially non-convex shape. However, we find the performance of this design drops about 10 mAP and more results are reported in Table 2. The possible reason is that our angle-based sampling produces 1D sample sequence, yet, contour vertices sequence is a 2D sample sequence which is more sensitive to noise. We report the reconstruction error of these two shape signatures on Pascal VOC 2012 training in Fig. 5 (denoted as “IR” and “XY” respectively). Admittedly, the XY has less reconstruction error when sampling the same points on the contour, but when compared with the same dimension of the vector, IR is more accurate. For example, the dimension of the vector of IR at  $N = 20$  is the same as XY at  $N = 10$ , the IR has a significantly less reconstruction error. Though when the  $N$  gets larger, the difference gets smaller, a large  $N$  will make training unstable as presented in Table 2.

Other classic shape signatures such as centroid radius, cumulative angle cannot reconstruct the shape.

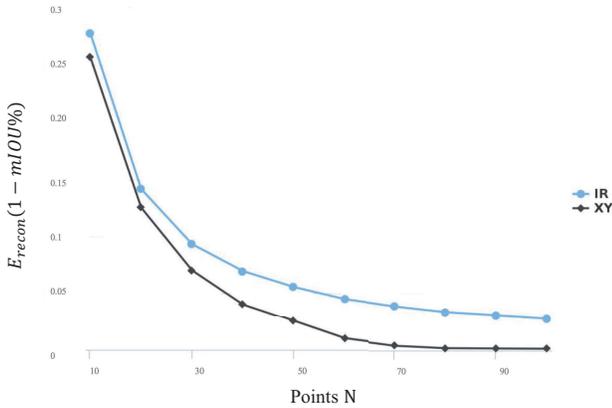


Figure 5. The reconstruction error  $E_{recon}$  of IR and XY with different sampling number points.

**Comparison with Other Fitting Methods** Other commonly used function approximation methods, namely polynomial regression and Fourier series fitting are also considered.

For polynomial regression, the goal is to fit shape vector  $\mathbf{k} = (v_0, \dots, v_n)$ , which is the coefficients of  $n$  degree polynomials,  $\mathbf{f}(\theta) = \sum_{i=0}^n v_i x^i$ . For Fourier series fitting, the shape vector is  $\mathbf{k} = (\omega, a_0, a_1, \dots, a_n, b_1, \dots, b_n)$ , the truncated  $n$  degree Fourier series is  $\mathbf{f}(\theta) = a_0/2 + \sum_{i=1}^n [a_i \cos(i\omega\theta) + b_i \sin(i\omega\theta)]$ . As the dimension of  $\mathbf{k}$  can be determined in advance, denoted as  $l$ . Thus we compare the methods from three aspects, *i.e.* the reconstruction error  $E_{recon}$ , sensitivity to the noises, and the numeric distribution of the coefficients.

The reconstruction errors  $E_{recon}$  is calculated by  $1 - mIOU$  under the same dimension  $l$  and point number  $N$  in Fig. 6. Then we set  $l = 8$  as an example to conduct the sensitivity analysis as shown in Fig. 7. For each coefficient, it is interrupted by the noise  $\varepsilon \sim N(0, \alpha \bar{k})$ ,  $\bar{k}$  is the mean of the corresponding coefficient. As we can see, the  $\omega$  of Fourier series is extremely sensitive, which may cause the Fourier fitting not suitable for the CNN training, as the CNN is known to regression with uncertainties. If we fix  $\omega = 1$ , it becomes less sensitive, but has considerably larger reconstruction error. Besides, considering the difficulty for the network to learn, we also investigate the statistic on the distribution of the fitted coefficients. See Fig. 8, Fig. 9 and Fig. 10. Chebyshev polynomials are better for shape signature fitting as it has less reconstruction error, less sensitivity to noise, better numeric distribution of coefficients.

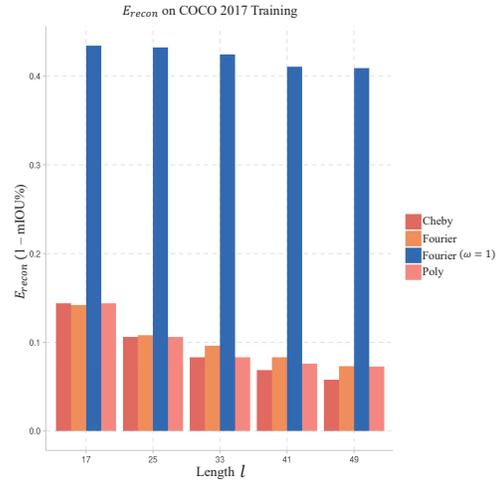


Figure 6. Comparison of  $E_{recon}$  on COCO 2017 training.

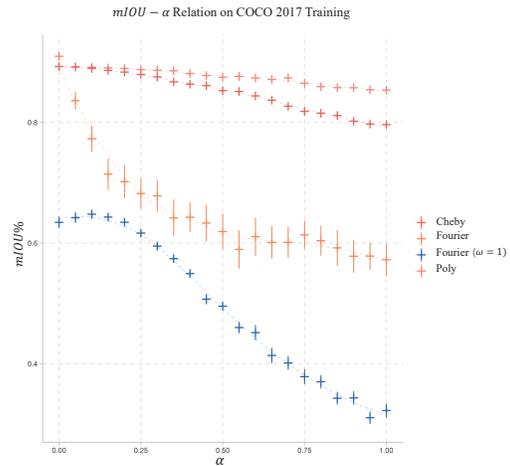


Figure 7. Comparison of the sensitivity on COCO 2017 training.

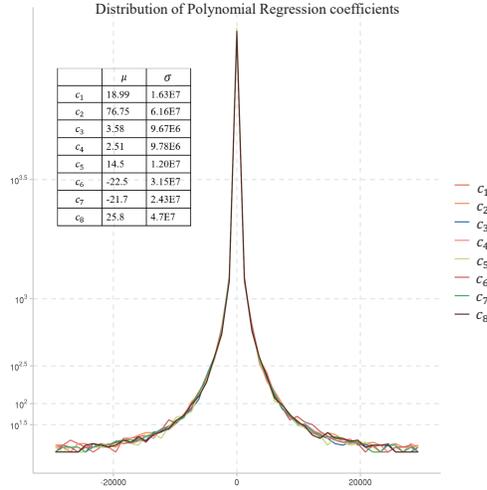


Figure 8. The overall mean of the coefficients is , and the variance is for Polynomial regression.

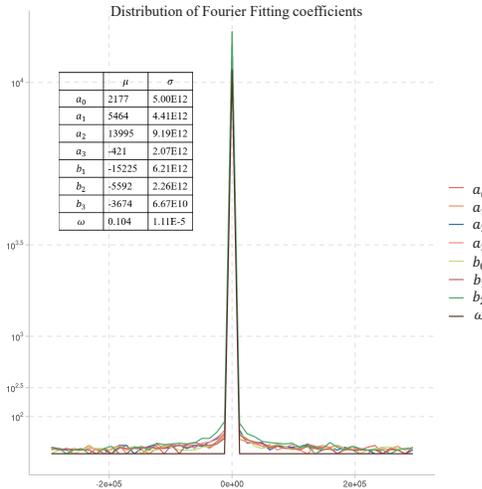


Figure 9. Coefficients distribution of Fourier series fitting on COCO training 2017.

### 3.3. Regression Under Object Detection Framework

Our network will learn to predict the inner center  $\hat{\mathbf{p}}$ , the shape vector  $\hat{\mathbf{k}}$ , along with the object bounding box. The loss function for bounding boxes regression, classification stays the same to the original object detection frameworks. For YOLOv3, the loss function for bounding box  $\mathcal{L}_{bbox}$  and classification  $\mathcal{L}_{cls}$  can be referred to [31]. As for the loss function for the shape learning:

$$\mathcal{L}_{shape} = \mathbb{1}^{obj} \|(\hat{\mathbf{p}} - \mathbf{p}) + (\hat{\mathbf{k}} - \mathbf{k})\|_2^2,$$

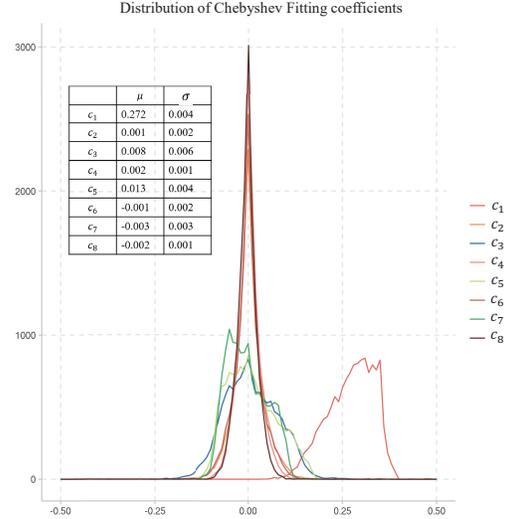


Figure 10. Coefficients distribution of Chebyshev polynomial fitting on COCO training 2017.

where  $\mathbb{1}^{obj}$  indicates the grid cells with objects for the one-stage detectors, and the proposals for the two-stage detectors. Thus the overall objective function is:

$$\mathcal{L} = \lambda_{cls} \mathcal{L}_{cls} + \lambda_{bbox} \mathcal{L}_{bbox} + \lambda_{shape} \mathcal{L}_{shape}$$

### 3.4. Decoding Shape Vector to Shape

Given the shape vector dimension  $l$ , the predicted shape vector  $\hat{\mathbf{k}} = (\hat{k}_0, \dots, \hat{k}_{l-1})^\top$ , the fitted Chebyshev polynomial is  $\hat{f}(\theta) = \sum_{i=0}^{l-1} \hat{k}_i T_i(\theta)$ . And the polar coordinate transform factor  $\mathbf{u}(\theta) = (\cos \theta, \sin \theta)$ . Thus the shape can be recovered by traversing the  $\theta \in [0, 2\pi)$

$$\hat{\mathbf{p}}_i = \hat{\mathbf{p}}_c + \hat{\mathbf{f}}(\theta) \odot \mathbf{u}(\theta).$$

$\odot$  is the Hadamard product. This calculation can be written in tensor operation form. Given the batch size  $bs$ , the corresponding tensor version are  $\Theta \in \mathbb{R}^{bs \times 1 \times N}$  for angles sampled,  $\hat{\mathbf{C}} \in \mathbb{R}^{bs \times 1 \times l}$  for the predicted shape vector,  $\hat{\mathbf{P}}_c \in \mathbb{R}^{bs \times 2 \times N}$  for the predicted inner centers and  $\hat{\mathbf{P}} \in \mathbb{R}^{bs \times 2 \times N}$  represents the decoded contour points. As expressed:

$$\hat{\mathbf{P}} = \hat{\mathbf{P}}_c + \hat{\mathbf{C}} T(\Theta) \odot \mathbf{u}^\top(\Theta).$$

In the GPU setting, the computation cost of such tensor operation is very minor. Due to this extremely fast shape decoding, our instance segmentation can achieve the same speed with object detection.

## 4. Experiment

We conduct extensive experiments to justify the descriptor choice and the efficacy of proposed methods. If not spec-

ified, the base detector is YOLOv3 implemented by GluonCV [14], the input image is  $416 \times 416$ .  $\lambda_{cls} = \lambda_{bbox} = \lambda_{shape} = 1$ . Other hyper-parameters stays the same as the YOLOv3 implementation. We trained 300 epochs and report the performance with the best evaluation results. For the model name with a bracket and a number in it, the number is the dimension of the shape vector.

#### 4.1. Explicit v.s. Implicit

We first compare the explicit shape encoding with the implicit shape encoding. As the previous work [16] provides a baseline for implicit shape representation with YOLO [30] as the base detector, to be fairly compared, we also trained the ESE-Seg with YOLO base detector, the dimension of the shape vector is also the same. We denote the model as “YOLO-Cheby (50)” and “YOLO-Cheby (20)”. The experiments are on Pascal SBD 2012 val [10].

To note, the mainstream instance segmentation based on mask, namely SDS [11], MNC [3], FCIS [19], Mask RCNN [12], can also be viewed as implicit shape encoding. We compare them with “YOLOv3-Cheby (20)” on Pascal VOC 2012 without SBD and COCO with their reported scores, which outperforms the Mask R-CNN (with ResNet50) at  $mAP^r @ 0.5$  on Pascal VOC and close to it on COCO. To note, the input image size is 800 on the shorter side for Mask R-CNN with ResNet50-FPN, which is almost 4 times to our  $416 \times 416$ . All results are reported in Table 1.

#### 4.2. On explicit descriptors

In this section, we will compare the object shape signatures and the function approximation methods quantitatively.

**On Different Shape Signatures** For object shape signatures, we compare our proposed IR with a straightforward 2D vertices representation on Pascal VOC 2012. (See Table 2) We adopt the squared boxes, *i.e.* the bounding box, as the baseline. To note, the squared boxes baseline is not the object detection scores, as the baseline computes the IoU between the bounding box and the instance mask.

For each shape signature, we compare regressing directly and regressing after Chebyshev polynomial fitting. For direct regression, we control the length of the shape signature by adjusting the  $\tau$  for each shape. We select 20 and 40 points to regress. We denote model trained on 2D vertices “XY”, the shape vector has a dimension of 40 and 80 respectively. As for the Chebyshev fitting on these signatures, we fit the  $x$  coordinates and  $y$  coordinates respectively. Denoted as “XY-Cheby (10+10)” means each fitted function has 10 coefficients.

**On Different Function Approximation Techniques** We have already compared the function approximation tech-

SBD (5732 val images)				
model \ mAP <sup>r</sup>	0.5	0.7	vol	Time (ms)
BinaryMask[16]	32.3	12.0	28.6	26.3
Radial[16]	30.0	6.5	29.0	27.1
Embedding (50) [16]	32.6	14.8	28.9	30.5
Embedding (20) [16]	34.6	<b>15.0</b>	31.5	28.0
YOLO-Cheby (50)	39.1	10.5	32.6	24.2
YOLO-Cheby (20)	<b>40.7</b>	12.1	35.3	<b>24.0</b>
Pascal VOC 2012 val				
model \ mAP <sup>r</sup>	0.5	0.7	vol	Time (ms)
SDS	49.7	25.3	41.4	48k
MNC	59.1	36.0	-	360
FCIS	65.7	<b>52.1</b>	-	160
Mask R-CNN	68.5	40.2	-	180
YOLOv3-Cheby (20)	62.6	32.4	52.0	<b>26.0</b>
+ COCO pretrained	<b>69.3</b>	36.7	54.2	<b>26.0</b>
COCO 2017 val				
model \ mAP	0.5	0.75	all	Time (ms)
FCIS	49.5	-	29.2	160
Mask R-CNN	51.2	31.5	30.3	180
YOLOv3-Cheby (20)	48.7	22.4	21.6	<b>26.0</b>

Table 1. Comparison of ESE-Seg to the previous methods on Pascal SBD 2012 val, Pascal VOC 2012 without SBD val, and COCO 2017 val.

model \ mAP <sup>r</sup>	0.5	0.7
Squared Boxes	42.3	8.6
XY (20)	46.1	10.7
XY (40)	43.5	11.2
XY-Cheby (10+10)	48.3	16.4
XY-Cheby (20+20)	53.1	20.9
IR (20)	48.8	13.5
IR (40)	52.6	19.3
IR (60)	51.7	16.4
IR-Cheby (20)	62.6	32.4

Table 2. We compare different choice of the shape signatures on Pascal VOC 2012.

niques through off-line analysis. However, it is still interesting to know performance of the neural network on the coefficients obtained by these methods.

All the function approximations are carried out on IR  $f(\theta)$ . The polynomial regression is denoted as “Poly”, while “Fourier” for Fourier series fitting and “Cheby” for Chebyshev polynomial fitting. All models have tested on Pascal VOC 2012 val. See Table 3.

Samples Selected from COCO 2017



Samples Selected from VOC 2012

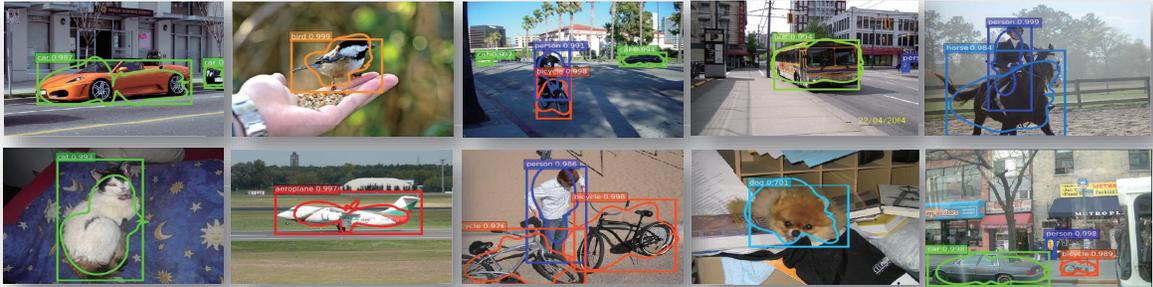


Figure 11. Qualitative results generated by our methods.

model	mAP <sup>r</sup>	
	0.5	0.7
Poly (20)	26.3	5.4
Fourier (20)	37.5	9.1
Fourier (40)	36.1	8.5
Cheby (20)	62.6	32.4
Cheby (40)	60.7	31.5

Table 3. Comparison of the performance of different shape signatures on Pascal VOC 2012 val.

model	mAP <sup>r</sup>			
	0.5	0.7	vol	Tims (ms)
YOLOv3-Cheby (20)	62.6	32.4	52.0	26
Faster-Cheby (20)	63.4	32.8	54.2	180
Retina-Cheby (20)	65.9	36.5	56.7	73
YOLOv3-tiny-Cheby (20)	53.2	15.8	42.5	<b>8</b>

Table 4. Comparison of different base object detectors with IR shape signature and Chebyshev fitting on Pascal VOC 2012 val.

### 4.3. On base object detector

To show the generality of the object shape detection, we also conduct the shape learning on Faster R-CNN (“Faster-Cheby (20)”), RetinaNet (“Retina-Cheby (20)”) and YOLOv3-tiny (“YOLOv3-tiny-Cheby (20)”). Not only the performance is stable for all these bounding box-based detectors, but the speed boost due to the detector can be enjoyed. As shown in Table 4.

### 4.4. Qualitative Results

Qualitative results are shown in Fig. 11. Obviously, the predicted shape vectors indeed capture the characteristics of the contours, not produce the random noise.

## 5. Limitations and Future Works

Our proposed ESE-Seg can achieve the instance segmentation with minor time-consumption, with a decent performance at IoU threshold 0.5. However, due to the inaccuracy  $E_{recon}$  of the shape vector, and the noise comes with the CNN regression, performance at larger IoU threshold like 0.7 drop a large margin. In the future, better ways to explicitly represent the shape, and better ways to train the CNN regression which will contribute to higher performance at high IOU threshold are of high interest.

**Acknowledgement** This work is supported in part by the National Key R&D Program of China, No. 2017YFA0700800, National Natural Science Foundation of China under Grants 61772332.

## References

- [1] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *CVPR*, 2018.
- [2] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiao-xiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4974–4983, 2019.
- [3] Jifeng Dai, Kaiming He, and Jian Sun. Instance-aware semantic segmentation via multi-task network cascades. In *CVPR*, 2016.
- [4] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016.
- [5] Edward Roy Davies. Machine vision: Theory, algorithms and practicalities. 1990, 1997.
- [6] Mark. Everingham, Luc. Van Gool, Christopher Williams, John. Winn, and Andrew. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010.
- [7] Hao-Shu Fang, Jianhua Sun, Runzhong Wang, , Minghao Gou, Yonglu Li, and Cewu Lu. Instaboost: Boosting instance segmentation via probability map guidedcopy-pasting. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2019.
- [8] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [9] Albert Gordo, Jon Almazán, Jerome Revaud, and Diane Larlus. Deep image retrieval: Learning global representations for image search. In *European conference on computer vision*, pages 241–257. Springer, 2016.
- [10] Bharath Hariharan, Pablo Arbelaez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *International Conference on Computer Vision (ICCV)*, 2011.
- [11] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *European Conference on Computer Vision*, pages 297–312. Springer, 2014.
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [14] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. *arXiv preprint arXiv:1812.01187*, 2018.
- [15] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [16] Saumya Jetley, Michael Sapienza, Stuart Golodetz, and Philip HS Torr. Straight to shapes: real-time detection of encoded shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6550–6559, 2017.
- [17] Hae-Kwang Kim and Jong-Deuk Kim. Region-based shape descriptor invariant to rotation, scale and translation. *Signal Processing: Image Communication*, 16(1):87 – 93, 2000.
- [18] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018.
- [19] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2359–2367, 2017.
- [20] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014. Springer International Publishing.
- [22] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [23] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [24] Xing Lu, Buyu Li, Yuxin Yue, Quanquan Li, and Junjie Yan. Grid r-cnn. *CoRR*, abs/1811.12030, 2018.
- [25] Calvin R Maurer, Rensheng Qi, and Vijay Raghavan. A linear time algorithm for computing exact euclidean distance transforms of binary images in arbitrary dimensions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(2):265–270, 2003.
- [26] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pages 2274–2284, USA, 2017. Curran Associates Inc.
- [27] David Novotny, Samuel Albanie, Diane Larlus, and Andrea Vedaldi. Semi-convolutional operators for instance segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 86–102, 2018.
- [28] Bo Pang, Kaiwen Zha, Hanwen Cao, Chen Shi, and Cewu Lu. Deep rnn framework for visual sequential applications. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 423–432, 2019.
- [29] Ali S Razavian, Josephine Sullivan, Stefan Carlsson, and Atsuto Maki. Visual instance retrieval with deep convolutional

- networks. *ITE Transactions on Media Technology and Applications*, 4(3):251–258, 2016.
- [30] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [31] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [32] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [33] Bharat Singh, Mahyar Najibi, and Larry S Davis. SNIPER: Efficient multi-scale training. *NIPS*, 2018.
- [34] Peter J Van Otterloo. *A contour-oriented approach to shape analysis*. Prentice Hall International (UK) Ltd., 1991.
- [35] Kilian Q. Weinberger and Lawrence K. Saul. Distance metric learning for large margin nearest neighbor classification. *J. Mach. Learn. Res.*, 10:207–244, June 2009.
- [36] Wenqiang Xu, Yonglu Li, and Cewu Lu. Srda: Generating instance segmentation annotation via scanning, reasoning and domain adaptation. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [37] Ian T. Young, Joseph E. Walker, and Jack E. Bowie. An analysis technique for biological shape. i. *Information and Control*, 25(4):357 – 370, 1974.
- [38] Joe Yue-Hei Ng, Fan Yang, and Larry S Davis. Exploiting local features from deep networks for image retrieval. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 53–61, 2015.
- [39] Dengsheng Zhang, Guojun Lu, et al. A comparative study of fourier descriptors for shape representation and retrieval. In *Proc. 5th Asian Conference on Computer Vision*, page 35. Citeseer, 2002.
- [40] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6848–6856, 2018.