
iShape: A First Step Towards Irregular Shape Instance Segmentation

Lei Yang **Ziwei Yan** **Wei Sun** **Yisheng He**
yanglei@megvii.com yanzw@buaa.edu.cn sunwei@megvii.com yhebk@connect.ust.hk

Zhenhang Huang **Haibin Huang** **Haoqiang Fan**
zhhuang@buct.edu.cn jackiehuanghaibin@gmail.com fhq@megvii.com

Abstract

1 In this paper, we introduce a brand new dataset to promote the study of instance
2 segmentation for objects with irregular shapes. Our key observation is that though
3 irregularly shaped objects widely exist in daily life and industrial scenarios, they
4 received little attention in the instance segmentation field due to the lack of corre-
5 sponding datasets. To fill this gap, we propose iShape, an irregular shape dataset
6 for instance segmentation. Unlike most existing instance segmentation datasets of
7 regular objects, iShape has many characteristics that challenge existing instance
8 segmentation algorithms, such as large overlaps between bounding boxes of in-
9 stances, extreme aspect ratios, and large numbers of connected components per
10 instance. We benchmark popular instance segmentation methods on iShape and
11 find their performance drop dramatically. Hence, we propose an affinity-based
12 instance segmentation algorithm, called ASIS, as a stronger baseline. ASIS ex-
13 plicitly combines perception and reasoning to solve Arbitrary Shape Instance
14 Segmentation including irregular objects. Experimental results show that ASIS
15 outperforms the state-of-the-art on iShape. Dataset and code are available at
16 <http://ishape.github.io>

17 1 Introduction

18 Instance segmentation aims to predict the
19 semantic and instance labels of each im-
20 age pixel. Compared to object detection
21 [1, 2, 3, 4, 5, 6, 7, 8] and semantic segmen-
22 tation [9, 10, 11], instance segmentation
23 provides more fine-grained information but
24 is more challenging and attracts more and
25 more research interests of the community.
26 Many methods [12, 13, 14, 15] and datasets
27 [16, 17, 18] continue to emerge in this field.
28 However, most of them focus on regularly
29 shaped objects and only a few [19, 18]
30 study irregular ones, which are thin, curved,
31 or having complex boundary and can not
32 be well-represented by regularly rectangu-
33 lar boxes. We think the insufficient explo-

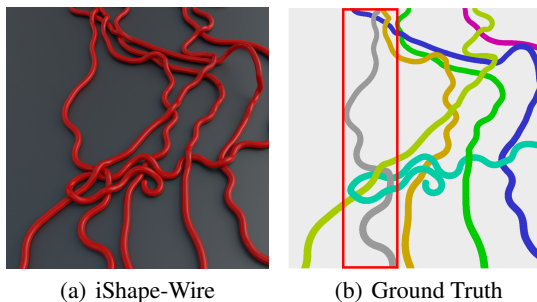


Figure 1: A typical scene of objects with irregular shape and similar appearance. It has many characteristics that challenge instance segmentation algorithms, including the large overlaps between bounding boxes of objects, extreme aspect ratios (bounding box of the grey mask), and large numbers of connected components in one instance (green and blue masks).

34 ration of this direction is caused by the lack
35 of corresponding datasets.

36 In this work, we present iShape, a new dataset designed for irregular **Shape** instance segmentation.
37 Our dataset consists of six sub-datasets, namely iShape-Antenna, iShape-Branch, iShape-Fence,
38 iShape-Log, iShape-Hanger, and iShape-Wire. As shown in Figure 2, each sub-dataset represents
39 scenes of a typical irregular shape, for example, strip shape, hollow shape, and mesh shape. iShape
40 has many characteristics that reflect the difficulty of instance segmentation for irregularly shaped
41 objects. The most prominent one is the large overlaps between bounding boxes of objects, which is
42 hard for proposal-based methods[12, 14] due to feature ambiguity and non-maximum suppression
43 (NMS [20]). Meanwhile, overlapped objects that share the same center point challenge center-based
44 methods[21, 22, 23]. Another characteristic of iShape is a large number of objects with similar
45 appearances, which makes embedding-based methods[24, 25] hard to learn discriminative embedding.
46 Besides, each sub-dataset has some unique characteristics. For example, iShape-Fence has about 53
47 connected components per instance, and iShape-Log has a large object scale variation due to various
48 camera locations and perspective transformations. We hope that iShape can serve as a complement of
49 existing datasets to promote the study of instance segmentation for irregular shape as well as arbitrary
50 shape objects.

51 We also benchmark existing instance segmentation algorithms on iShape and find their performance
52 degrades significantly. To this end, we introduce a stronger baseline considering irregular shape in
53 this paper, which explicitly combines perception and reasoning. Our key insight is to simulate how a
54 person identifies an irregular object. Taking the wire shown in Figure 1 for example, one natural way
55 is to start from a local point and gradually expand by following the wire contour and figure out the
56 entire object. The behavior of such “following the contour” procedure is a process of **continuous**
57 **iterative reasoning based on local clues**, which is similar to the recent affinity-based approaches
58 [26, 27]. Under such observation, we propose a novel affinity-based instance segmentation baseline,
59 called ASIS, which includes principles of generating effective and efficient affinity kernel based on
60 dataset property to solve Arbitrary Shape Instance Segmentation. Experimental results show that the
61 proposed baseline outperforms existing state-of-the-art methods by a large margin on iShape.

62 Our contribution is summarized as follows:

- 63 • We propose a brand new dataset, named iShape, which focuses on irregular shape instance
64 segmentation and has many characteristics that challenge existing methods. In particular,
65 we analyzed the advantages of iShape over other instance segmentation datasets.
- 66 • We benchmark popular instance segmentation algorithms on iShape to reveal the drawbacks
67 of existing algorithms on irregularly shaped objects.
- 68 • Inspired by human’s behavior on instance segmentation, we propose ASIS as a stronger
69 baseline on iShape, which explicitly combines perception and reasoning to solve Arbitrary
70 Shape Instance Segmentation.

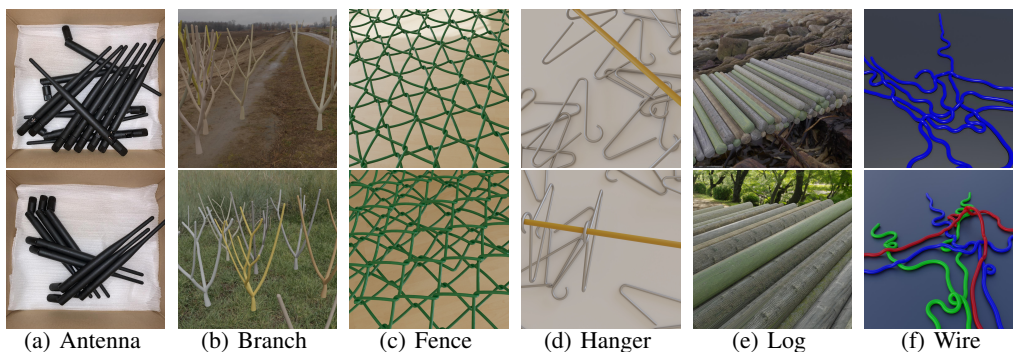


Figure 2: The six sub-datasets in iShape.

71 2 Related Work

72 2.1 Existing Datasets

73 There are several benchmark datasets collected to promote the exploration of instance segmentation.
74 The COCO [16] and the Cityscapes [17] are the most popular ones among them. However, the
75 shapes of target objects in these datasets are too regular. The connected components per instance
76 (CCPI) and average MaxIoU are low in the datasets and state-of-the-art algorithms selected from
77 them can not generalize to more challenging scenarios. Instead, in the scenario of human detection
78 and segmentation, the OC human [19] and the Crowd Human [28] introduce datasets with larger
79 MaxIoU. Nevertheless, the OC human dataset only provides a small number of images for testing,
80 and the number of instances per image is too small to challenge instance segmentation algorithms.
81 While the crowd human dataset only provides annotations of object bounding boxes, limiting their
82 application to the instance segmentation field. In the area of photogrammetry, the iSAID [18] dataset
83 is proposed to lead algorithms to tackle objects with multi scales. However, shapes of objects in this
84 dataset are common, most of which are rectangular, and the lack of instance overlapping reduces
85 its challenge to instance segmentation algorithms as well. Under the observation that these existing
86 regular datasets are not enough to challenge algorithms for more general scenarios, we propose
87 iShape, which contains irregularly shaped objects with large overlaps between bounding boxes of
88 objects, extreme aspect ratios, and large numbers of CCPI to promote the capabilities of instance
89 segmentation algorithms.

90 2.2 Instance Segmentation Algorithms

91 Existing instance segmentation algorithms can be divided into two classes, proposal-based and
92 proposal-free.

93 **Proposal-based approaches** One line of these approaches [12, 14, 29] solve instance segmentation
94 within a two-stage manner, by first propose regions of interests (RoIs) and then regress the semantic
95 labels of pixels within them. The drawback of these approaches comes from the loss of objects by
96 NMS due to large IoU. Instead, works like [15] tackle the problem within a single-stage manner. For
97 example, PolarMask [15] models the contours based on the polar coordinate system and then obtain
98 instance segmentation by center classification and dense distance regression. But the convex hull
99 setting limits its accuracy.

100 **Proposal-free approaches** To shake off the rely on proposals and avoid the drawback caused by
101 them, many bottom-up approaches like [22, 23, 24, 25] are introduced. These works are in various
102 frameworks. The recent affinity-based methods obtain instance segmentation via affinity derivation
103 [26] and graph partition[30]. This formulation is more similar to the perception and reasoning
104 procedure of we human beings and can handle more challenging scenarios. GMIS [26] utilizes both
105 region proposals and pixel affinities to segment images and SSAP [27] outputs the affinity pyramid
106 and then performs cascaded graph partition. However, The affinity kernels of GMIS and SSAP are
107 sparse in angle and distance, leading to missing components of some instances due to loss of affinity
108 connection. To this end, we propose ASIS which includes principles of generating effective and
109 efficient affinity kernel based on dataset property to solve Arbitrary Shape Instance Segmentation and
110 achieve great improvement on iShape.

111 3 iShape Dataset

112 3.1 Dataset Creation

113 iShape consists of six sub-datasets. One of them, iShape-Antenna, is collected from real scenes, which
114 are used for antenna counting and grasping in automatic production lines. The other five sub-datasets
115 are synthetic datasets that try to simulate five typical irregular shape instance segmentation scenes.

116 **iShape-Antenna Creation.** For the creation of iShape-Antenna, we first prepare a carton with a
117 white cushion at the bottom, then randomly and elaborately place antennas in it to generate various
118 scenes. Above the box, there is a camera with a light that points to the inside of the box to capture the
119 scene images. We collect 370 pictures and annotate 3,036 instance masks then split them equally
120 for training and testing. The labeling is done by our supplier. We have checked all the annotations

121 ourselves, and corrected the wrong annotations. Although iShape-Antenna only contains 370 images,
 122 the number of instances reaches 3,036 which is more than most categories in Cityscapes [17] and
 123 PASCAL VOC [31].

124 **Synthetic Sub-datasets Creation** There are lots of typical irregular shape instance segmentation
 125 scenes. Consequently, it is impractical to collect a natural dataset for each typical scene. Since it
 126 is traditional to study computer vision problems using synthetic data [32, 33], we synthesize five
 127 sub-datasets of iShape which include iShape-Branch, iShape-Fence, iShape-Log, iShape-Hanger, and
 128 iShape-Wire, by using CG software Blender. In particular, We build corresponding 3D models and
 129 placement they appropriate in Blender with optional random background and lighting environment,
 130 optional physic engine, and random camera position. The creation configs of synthesis sub-datasets
 131 are listed in the appendix. After setting up the scene, we use a ray tracing render engine to render the
 132 RGB image. Besides, We build and open source a blender module, bpycv [34], to generate instance
 133 annotation. We generate 2500 images for each sub-dataset, 2000 for training, 500 for testing.

134 3.2 Dataset Characteristics

135 In this sub-section, we analyze the characteristics of iShape and compare it with other instance
 136 segmentation datasets. Since each sub-dataset represents irregularly shaped objects in different
 137 scenes, we present the statistical results of each sub-dataset separately.

138 **Dataset basic information.** As summarized in Table 1, iShape contains 12,870 images with 175,840
 139 instances. All images are 1024×1024 pixels and annotated with pixel-level ground truth instance
 140 masks. Since iShape focus on evaluating the performance of algorithms on the irregular shape, each
 141 scene consists of multiple instances of one class, which is also common cases in industrial scenarios.

142 **Instance count per image.** A larger instance count is more challenging. Despite iSAID getting
 143 the highest instance count per image, it is unfair for extremely high-resolution images and normal-
 144 resolution images to be compared on the indicator. Among iShape, the instance count per image of
 145 iShape-Log reaches 28.86 that significantly higher than other normal-resolution datasets.

146 **The large overlap between objects.** We introduce a new indicator, Overlap of Sum (OoS), which
 147 aims to measure the degree of occlusion and crowding in a scene, defined as follows:

$$Overlap\ of\ Sum = \begin{cases} 1 - \frac{|\bigcup_{i=1}^n C_i|}{\sum_{i=1}^n |C_i|}, & n > 0 \\ 0, & n = 0 \end{cases} \quad (1)$$

148 where C means bounding boxes(bbox) or convex hulls(convex) of all instances in the image, n means
 149 number of instances, \bigcup means union operation, and $|C_i|$ means to get the area of C_i . The statistics of
 150 average OoS for bounding box and convex hull are presented in Table. 1. For bounding box OoS,
 151 All iShape sub-datasets are higher than other datasets, which reflects the large overlap characteristic
 152 of iShape. Thanks to the large-area hollow structure, iShape-Fence gets the highest average convex
 153 hull OoS 0.63. Moreover, The Average MaxIoU [19] of all images also reflects the large overlap
 154 characteristic of iShape.

Table 1: Comparison of statistics with different datasets.

Dataset	Images	Ins.	Ins./image	OoS		AvgMIoU	Aspect ratio	CCPI
				bbox	convex			
Cityscapes	2,975	52,139	17.52	0.14	0.07	0.394	2.29	1.34
COCO	123,287	895,795	7.26	0.15	0.09	0.210	2.59	1.41
CrowdHuman	15,000	339,565	22.64	-	-	-	-	-
OC Human	4,731	8,110	1.71	0.25	0.20	0.424	2.28	3.11
iSAID	2,806	655,451	233.58	-	-	-	2.40	-
Antenna	370	3,036	8.20	0.62	0.23	0.655	9.86	2.45
Branch	2,500	26,046	10.14	0.62	0.52	0.750	2.47	10.88
Fence	2,500	7,870	3.15	0.65	0.63	0.983	1.05	53.65
Hanger	2,500	49,275	19.71	0.53	0.34	0.685	3.28	4.94
Log	2,500	72,144	28.86	0.73	0.06	0.843	34.14	2.64
Wire	2,500	17,469	6.99	0.74	0.60	0.795	3.32	4.76
iShape	12,870	175,840	13.66	0.65	0.42	0.806	15.84	6.99

155 **The similar appearance between object instances.** Instances from the same object class in iShape
 156 share similar appearance, which is challenging to embedding-based algorithms. In particular, any two
 157 object instance in iShape-Antenna, iShape-Fence and iShape-Hanger are indistinguishable according
 158 to their appearance. They are generated from either industrial standard antennas or copies of the same
 159 mesh models. Meanwhile, the appearance of objects in iShape-Branch, iShape-Log, and iShape-Wire
 160 are slightly changeable to add some variances, but appearances of different instances are still much
 161 more similar than those from other existing datasets in Table 1.

162 **Aspect ratio.** Table 1 presents statistics on the average aspect ratio of the object’s minimum bounding
 163 rectangle for each dataset. Among them, iShape-Log’s aspect ratio reaches 34.14, which is more
 164 than 10 times of other regularly shaped datasets. Such a gap is caused by two following reasons:
 165 Firstly, the shape of logs has a large aspect ratio. Secondly, partially occluding logs leads to a higher
 166 aspect ratio. iShape-Antenna also has a high aspect ratio, 9.86, which exceeds other regularly shaped
 167 datasets.

168 **Connected Components Per Instance (CCPI).** Larger CCPI poses a larger challenge to instance
 169 segmentation algorithms. Due to the characteristics of irregular shaped objects and the occlusion of
 170 scenes, the instance appearance under the mesh shape tends to be divided into many pieces, leading
 171 to large CCPI of iShape-Fence. As is shown in Table 1, the result on CCPI of iShape-Fence is 53.65,
 172 about 5 times higher than the second place. iShape-Branch, iShape-Hanger, and iShape-Wire also
 173 have a large CCPI that exceeds other regularly shaped datasets.

174 4 Baseline Approach

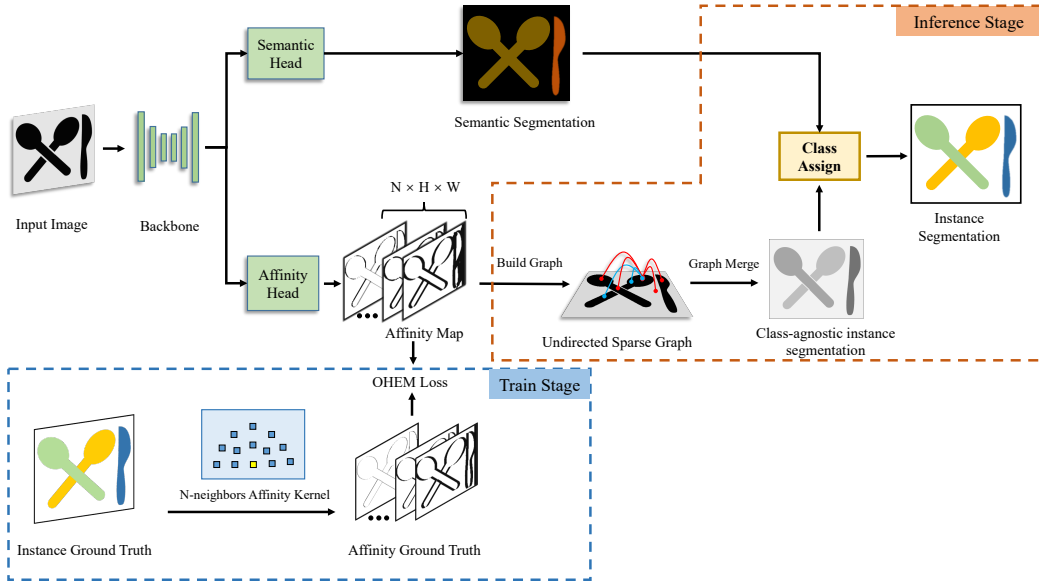


Figure 3: **Overview of ASIS.** In the training stage, the network learns to predict the semantic segmentation as well as the affinity map where the ground truth of affinity can be generated by affinity kernel and instance ground truth. In the inference stage, the predicted affinity map will be used to construct a sparse and undirected graph, with pixel as node and affinity map as edge. The final instance label then can be generated by applying a class assign module on top of the constructed graph and semantic segmentation map.

175 Inspired by how a person identifies a wire shown in Figure 1, We propose an affinity-based instance
 176 segmentation baseline, called ASIS, to solve Arbitrary Shape Instance Segmentation by explicitly
 177 combining perception and reasoning. Besides, ASIS includes principles of generating effective and
 178 efficient affinity kernel based on dataset property. In this section, an overview of the pipeline is
 179 firstly described in Subsection 4.1, then design principles of the ASIS affinity kernel are explained in
 180 Subsection 4.2.

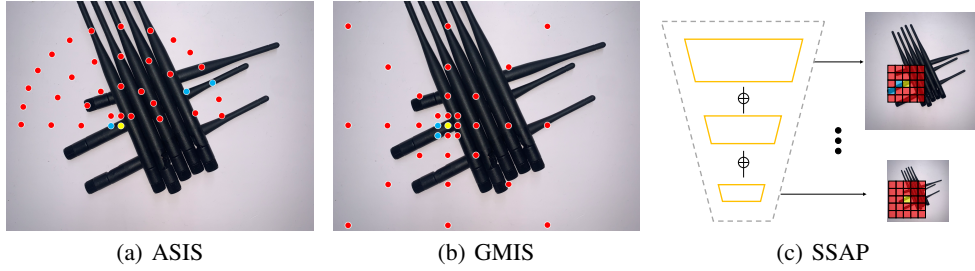


Figure 4: Illustration for affinity kernels. (a) ASIS affinity kernel could connect these two segments with two neighbors (blue points). (b) GMIS affinity kernel cannot reach the right segment. (c) Examples of failure case for SSAP affinity kernel. For higher resolutions (top), 5×5 affinity window cannot reach the segment on the right. For lower resolutions (bottom), **the view of thin antennas are lost in the resized feature maps.**

181 4.1 Overview of ASIS

182 As shown in Figure 3, we firstly employ the PSPNet [11] as the backbone and remove its last softmax
 183 activation function to extract features. The semantic head, which combines a single convolution
 184 layer and a softmax activation function, will input those features and output a $C \times H \times W$ semantic
 185 segmentation probability map where C means the total categories number. The affinity head that
 186 consists of a single convolution layer and a sigmoid activation function will output a $N \times H \times W$
 187 affinity map, where N is the neighbor number of affinity kernel. **Affinity kernel [26] defines a set of**
 188 **neighboring pixels that needs to generate affinity information. Examples of affinity kernels can found**
 189 **in Figure 4. Each channel of the affinity map represents a probability of whether the neighbor pixel**
 190 **and the current one belong to the same instance.**

191 During the training stage, we apply the affinity kernel on the instance segmentation ground truth
 192 to generate the affinity ground truth. Since affinity ground truth is extremely imbalanced, an
 193 OHEM [35] loss is calculated between the predicted affinity map and the affinity ground truth to
 194 effectively alleviate the problem. For affinity map with input size $S = N \times H \times W$, we define
 195 $A = \{a_1, a_2, \dots, a_S\}$ and $Y = \{y_1, y_2, \dots, y_S\}$ the sets of each pixel of the predicted affinity map
 196 and the corresponding ground truth. The loss of the i_{th} pixel L_i is defined as:

$$L_i = -y_i \log(a_i) - (1 - y_i) \log(1 - a_i). \quad (2)$$

197 Assume that the set L' is the *Topk* value in $L = \{L_1, L_2, \dots, L_S\}$. K takes the top ten percent. The
 198 OHEM loss is as follows:

$$\mathcal{L}_{aff} = \frac{1}{|L'|} \sum_{l' \in L'} l', \quad (3)$$

199 Affinities that connect segments of fragmented instances are important but hard to learn. Thanks
 200 to the difficulty of learning these affinities, the OHEM loss pays more attention to these important
 201 affinities. Besides, a standard cross-entropy loss for pixels \mathcal{L}_{sem} is applied to semantic segmentation
 202 output. The final training loss \mathcal{L} is defined as:

$$\mathcal{L} = \lambda \mathcal{L}_{aff} + (1 - \lambda) \mathcal{L}_{sem} \quad (4)$$

203 For the inference stage, we firstly take pixels as nodes and affinity map as edges to build an undirected
 204 sparse graph. The undirected sparse graph in Figure 3 shows an example of how a pixel node on
 205 the spoon should connect the other pixel nodes. **Then, we apply the graph merge algorithm [26] on**
 206 **the undirected sparse graph. The algorithm will merge nodes that have a positive affinity to each**
 207 **other into one supernode, by contrast, keep nodes independent if their affinity is negative. Pixels that**
 208 **merged to the same supernodes are regarded as belonging to the same instance. In this way, we obtain**
 209 **a class-agnostic instance map.** A class assign module [26] will take the class-agnostic instance map
 210 and the semantic segmentation result as input, then assign a class label with a confidence value to
 211 each instance.

Table 2: Qualitative results on iShape. We report the mmAP of six sub-datasets and the average of mmAP.

Method	Backbone	Antenna	Branch	Fence	Hanger	Log	Wire	Avg
SOLOv2 [21]	ResNet-50	6.6	27.5	0.0	28.8	22.2	0.0	14.07
PolarMask [15]	ResNet-50	0.0	0.0	0.0	0.0	18.6	0.0	3.10
SE [22]	-	38.3	0.0	0.0	49.8	20.9	0.0	18.17
Mask RCNN [12]	ResNet-50	16.9	4.2	0.0	22.1	32.6	0.0	12.63
DETR [38]	ResNet-50	2.1	2.6	0.0	32.2	46.2	0.0	13.85
ASIS(ours)	ResNet-50	77.5	25.1	37.1	53.1	69.3	64.9	54.50

212 4.2 ASIS Affinity Kernel

213 Since instances could be divided into many segments, it is important to design an appropriate affinity
 214 kernel to connect those segments that belong to the same instance. As shown in Figure 4(b) and
 215 Figure 4(c), The yellow point is the current pixel. Red points belong to different instances and blue
 216 points belong to the same instance of the current pixel. The antenna that the current pixel (yellow
 217 point) belongs to has two segments that need to be connected by affinity neighbor. The previous
 218 affinity-based approaches [26, 27] don't take into account such problems and cause some failures.
 219 Hence, we propose principles of generating effective and efficient affinity kernel based on dataset
 220 property to solve Arbitrary Shape Instance Segmentation. Our affinity kernel is shown in 4(a).

221 Affinity kernels of GMIS and SSAP are centered symmetric, unfortunately, that will cause redundant
 222 outputs. For example, the affinity of pixel (1, 1) with its right side pixel and the affinity of pixel (1, 2)
 223 with its left side pixel both mean the probability of these two pixels belonging to one instance. A
 224 detailed description of redundant affinity can be found in the appendix. To reduce the network's
 225 outputs, redundant affinity neighbors are discarded in the ASIS affinity kernel. As shown in 4(a),
 226 affinity neighbors of ASIS are distributed in an asymmetric semicircle structure. Besides, the area
 227 covered by asymmetric semicircle affinity kernel is reduced by half, in other words, the demand for
 228 receptive fields is reduced, which further reduces the difficulty of CNN learning affinity.

229 Two main parameters determine the shape of the ASIS affinity kernel. Kernel radius r_k controls the
 230 radius of the kernel and determines how far the farthest of two segments can be reached. Affinity
 231 neighbor gap g represents the distance between any two nearly affinity neighbors, thus, g controls
 232 the sparseness of the affinity neighbor. Since each dataset has its optimal affinity kernel, we propose
 233 another algorithm that could adaptively generate appropriate r_k and g based on the dataset property.
 234 Detailed descriptions of these two algorithms can be found in the appendix.

235 5 Experiments

236 In this section, we choose representative instance segmentation methods in various paradigms and
 237 benchmark them on iShape to reveal the drawbacks of existing methods on irregularly shaped objects.
 238 All the existing methods are trained and tested on six iShape sub-datasets with their defaults setting.
 239 And we further study the effect of our baseline method, ASIS.

240 **Evaluation Metrics** The evaluation metric is mainly Average Precision (AP), which is calculated by
 241 averaging the precision under mask IoU (Intersection over Union) thresholds from 0.50 to 0.95 at the
 242 step of 0.05.

243 **Implementation Details** The input image resolution of our framework is 512×512 . The image data
 244 augmentation is flipped horizontally or vertically with a probability of 0.5. We use the ResNet-50
 245 [36] as our backbone network and the weight is initialized with ImageNet [37] pretrained model. All
 246 experiments are trained in 4 2080Ti GPUs and batch size is set to 8. The stochastic gradient descent
 247 (SGD) solver is adopted in 50K iterations. The momentum is set to 0.9 and weight decay is set to
 248 0.0005. The learning rate is initially set to 0.01 and decreases linearly.

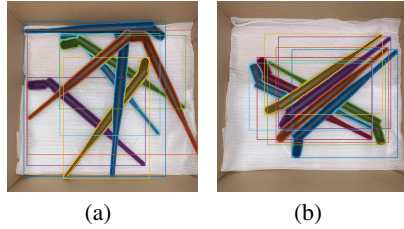


Figure 5: Two example false cases of ASIS on iShape-Antenna. (a) Two antennas merged into one (blue and orange). (b) ASIS fails to connect the right parts of an object (red and sky blue).

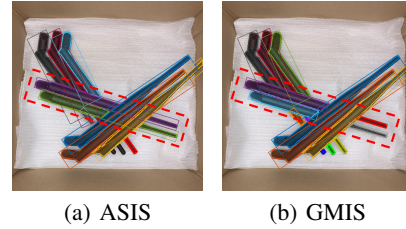


Figure 6: Results compared with GMIS kernel. As shown in (b), GMIS fail to connect segments that belong to one instance.

249 5.1 Experiment Results

250 We evaluate the proposed ASIS and other popular approaches on iShape. The quantitative results are
 251 shown in Table 2 and some qualitative results are reported in Figure 7.

252 As is shown in Table 2, the performance of Mask R-CNN [12] is far from satisfactory on iShape. We
 253 think the drop in performance mainly comes from three drawbacks of the design. Firstly, the feature
 254 maps suffer from ambiguity when the IoU is large, which is a common characteristic of crowded
 255 scenes of irregular shape objects. Also, Mask R-CNN depends on the proposals of RoI, which may
 256 be abandoned by the NMS algorithm due to large IoU and lead to missing of some target objects.
 257 Moreover, many thin objects can not be segmented by Mask R-CNN because of its RoI pooling,
 258 which resizes the feature maps and lost the view of thin objects. The recent proposed end-to-end
 259 object detection approach, DETR [38], shake of the reliance of NMS and can better deal with objects
 260 with large IoU and achieve better performance, as shown in the table. However, DETR still suffers
 261 from the RoI pooling problems and performs badly on thin objects, as shown in Figure 7.

262 We also report some qualitative results of SE [22] in Figure 7. As is shown in the figure, one common
 263 failure case of SE is that when the length of irregular objects is longer than a threshold, the object
 264 will be split into multi instances, for example, the wire in Figure 7. We think that’s because SE
 265 will regress a circle of the target instance and then calculate its IoU with the mask for supervision.
 266 However, for long and thin irregular objects, the radius of the center circle can not reach the length
 267 of the target object, leading to a multi-split of a long instance. Also, instances that share the same
 268 center may cause ambiguity to SE, such as hanger and fence in Figure 7. Moreover, many centers of
 269 irregular objects lie outside the mask, making it hard to match them to the objects themselves.

270 We evaluate SOLO v2 [21] on the proposed iShape and find that it failed to segment instances that
 271 share the same center, for example, fences in Figure 7. Also, since SOLO V2 depends on the center
 272 point as SE, it also suffers from performance drop caused by object centers that lie outside the mask.

273 In Table 2, we report the performance of PolarMask [15] on our dataset. As is shown in the table,
 274 PolarMask can not solve the instance segmentation of irregular objects. That is because PolarMask
 275 can only represent a thirty-six-side mask due to its limited number of rays. Hence, it can not handle
 276 objects with hollow, for example, the fences. Also, they distinguish different instances according to
 277 center regression, which, however, can not handle instances that share the same center. We also find
 278 that PolarMask can only tackle some cases of logs in iShape, which looks like circles on the side and
 279 fit its convex hull mask setting.

280 Thanks to the perception and reasoning mechanism as well as the well-designed affinity kernels of
 281 our ASIS, it obtained the best performance on iShape. In Table 2, ASIS advances other approaches
 282 by 36% on the mAP metric. However, there are still some drawbacks to the design of ASIS and
 283 some failure cases caused by them. For example, in Figure 5(a), two instances are merged into one.
 284 We think that’s because the graph merge algorithm is a kind of greedy algorithm, while the greedy
 285 algorithm makes optimal decisions locally instead of looking for a global optimum. Hence, ASIS is
 286 not robust to false-positive (FP) with high confidence. Also, ASIS fails to connect the two parts of an
 287 object if they are far away from each other, for example, the antenna on Figure 5(b). We think that’s
 288 because CNN is not good at learning long-range affinity.

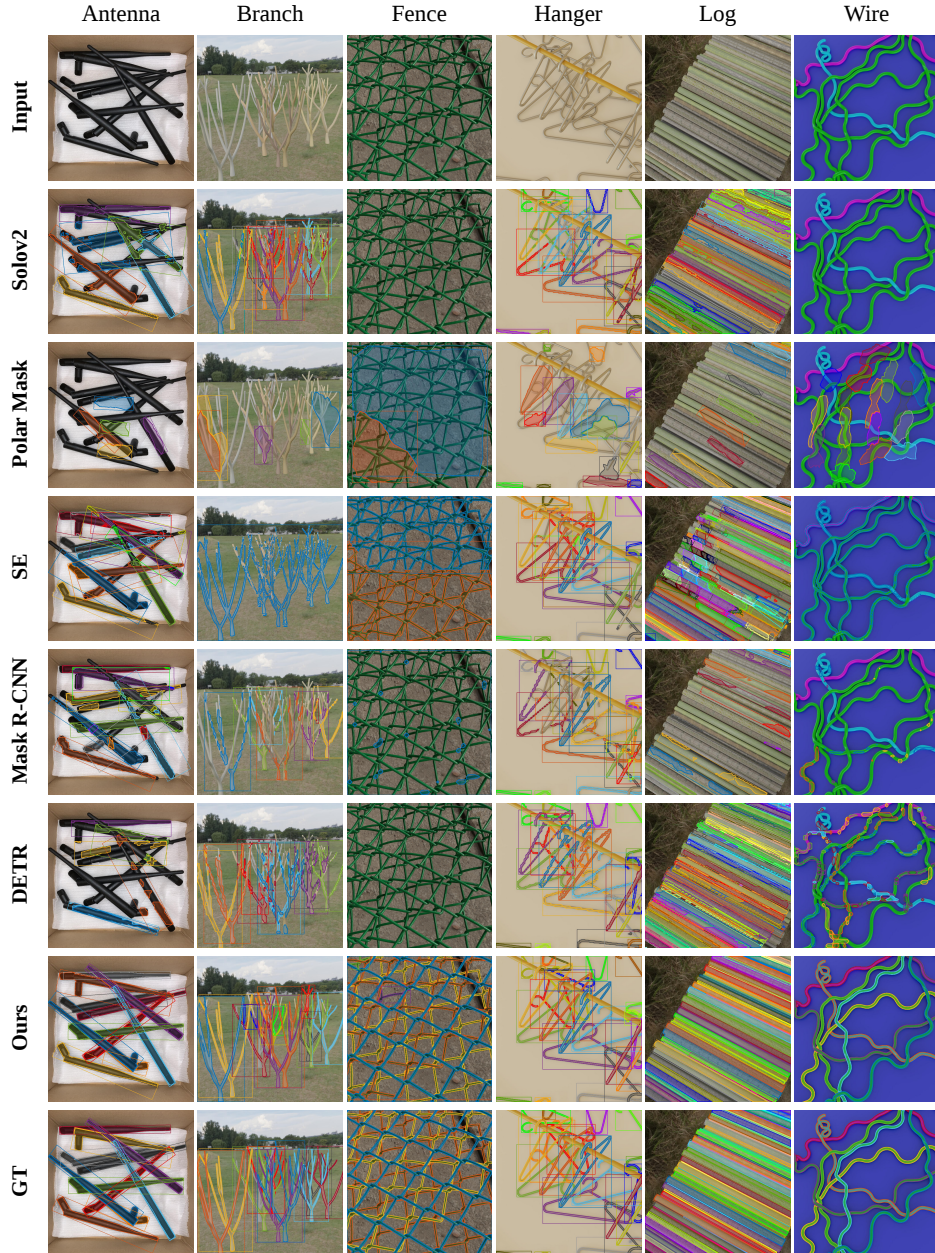


Figure 7: Qualitative results of different instance segmentation approaches on iShape.

289 **5.2 Ablation Study**

290 **Effect of ASIS.** We study the effect of ASIS in Table 3, where ASIS advances GMIS by 10.9% on
 291 iShape-Antenna. We think that is because our well designed affinity kernels based on dataset property
 292 can better discover the connectivity of different parts of an object. While GMIS suffers from its
 293 sparsity in distance and angle, results are shown in Figure 6(b). We also use ground truth affinity map
 294 to explore the upper bound of ASIS, where a 98.5% mAP is achieved, showing its great potential.
 295 Moreover, we find our non-centrosymmetric design of affinity kernels outperform centrosymmetric
 296 ones in the table. We think such a design cut off the output and calculation redundancy and reduce
 297 the requirement of large receptive field from CNN, simplifying representation learning.

298 **Effect of OHEM.** Table 3 shows that OHEM boosts the performance of GMIS and ASIS by a large
 299 margin. We think that is because OHEM can ease problems caused by imbalance distribution of
 positive and negative affinity.

Table 3: Comparison result of GMIS and ASIS. “SY” and “ASY” indicate a centrosymmetric or
 asymmetric affinity kernel respectively. \checkmark denotes equipped with and \circ not.

Affinity Kernel	Neighbors	Affinity GT	OHEM	mAP
GMIS [26]	56 (SY)	\circ	\circ	44.5
		\circ	\checkmark	69.9
		\checkmark	-	90.2
	28 (ASY)	\circ	\checkmark	72.7
ASIS(ours)	53 (ASY)	\circ	\circ	58.4
		\circ	\checkmark	77.5
		\checkmark	-	98.5

300

301 6 Conclusion

302 In this work, we introduce a new irregular shape instance segmentation dataset (iShape). iShape has
 303 many characteristics that challenge existing instance segmentation methods, such as large overlaps,
 304 extreme aspect ratios, and similar appearance between objects. We evaluate popular algorithms
 305 on iShape to establish the benchmark and analyze their drawbacks to reveal possible improving
 306 directions. Meanwhile, we propose a stronger baseline, ASIS, to better solve Arbitrary Shape Instance
 307 Segmentation. Thanks to the combination of perception and reasoning as well as the well-designed
 308 affinity kernels, ASIS outperforms the state-of-the-art methods on iShape. We believe that iShape and
 309 ASIS can serve as a complement to existing datasets and methods to promote the study of instance
 310 segmentation for irregular shape as well as arbitrary shape objects.

References

- 311 [1] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger, 2016.
- 312 [2] Ross Girshick. Fast r-cnn, 2015.
- 313 [3] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate
314 object detection and semantic segmentation, 2014.
- 315 [4] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection
316 with region proposal networks, 2016.
- 317 [5] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and
318 Alexander C. Berg. Ssd: Single shot multibox detector. *Lecture Notes in Computer Science*, page 21–37,
319 2016.
- 320 [6] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional
321 networks, 2016.
- 322 [7] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection, 2017.
- 323 [8] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time
324 object detection, 2016.
- 325 [9] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmenta-
326 tion, 2015.
- 327 [10] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Deeplab:
328 Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,
329 2017.
- 330 [11] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing
331 network, 2017.
- 332 [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE
333 international conference on computer vision*, pages 2961–2969, 2017.
- 334 [13] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation, 2019.
- 335 [14] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: High quality object detection and instance segmenta-
336 tion, 2019.
- 337 [15] Enze Xie, Peize Sun, Xiaoge Song, Wenhai Wang, Ding Liang, Chunhua Shen, and Ping Luo. Polarmask:
338 Single shot instance segmentation with polar representation, 2020.
- 339 [16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
340 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on
341 computer vision*, pages 740–755. Springer, 2014.
- 342 [17] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson,
343 Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding,
344 2016.
- 345 [18] Syed Waqas Zamir, Aditya Arora, Akshita Gupta, Salman Khan, Guolei Sun, Fahad Shahbaz Khan, Fan
346 Zhu, Ling Shao, Gui-Song Xia, and Xiang Bai. isaid: A large-scale dataset for instance segmentation in
347 aerial images, 2019.
- 348 [19] Song-Hai Zhang, Ruilong Li, Xin Dong, Paul L. Rosin, Zixi Cai, Han Xi, Dingcheng Yang, Hao-Zhi
349 Huang, and Shi-Min Hu. Pose2seg: Detection free human instance segmentation, 2019.
- 350 [20] A. Neubeck and L. Van Gool. Efficient non-maximum suppression. In *18th International Conference on
351 Pattern Recognition (ICPR'06)*, volume 3, pages 850–855, 2006.
- 352 [21] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance
353 segmentation, 2020.
- 354 [22] Davy Neven, Bert De Brabandere, Marc Proesmans, and Luc Van Gool. Instance segmentation by jointly
355 optimizing spatial embeddings and clustering bandwidth. In *Proceedings of the IEEE Conference on
356 Computer Vision and Pattern Recognition*, pages 8837–8845, 2019.
- 357

- 358 [23] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-
359 Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation.
360 In *CVPR*, 2020.
- 361 [24] Bert De Brabandere, Davy Neven, and Luc Van Gool. Semantic instance segmentation with a discriminative
362 loss function. *arXiv preprint arXiv:1708.02551*, 2017.
- 363 [25] Shu Kong and Charless C Fowlkes. Recurrent pixel embedding for instance grouping. In *Proceedings of*
364 *the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9018–9028, 2018.
- 365 [26] Yiding Liu, Siyu Yang, Bin Li, Wengang Zhou, Jizheng Xu, Houqiang Li, and Yan Lu. Affinity derivation
366 and graph merge for instance segmentation, 2018.
- 367 [27] Naiyu Gao, Yanhu Shan, Yupei Wang, Xin Zhao, Yinan Yu, Ming Yang, and Kaiqi Huang. Ssap: Single-
368 shot instance segmentation with affinity pyramid. In *Proceedings of the IEEE International Conference on*
369 *Computer Vision*, pages 642–651, 2019.
- 370 [28] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. Crowdhuman: A
371 benchmark for detecting human in a crowd, 2018.
- 372 [29] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmen-
373 tation, 2018.
- 374 [30] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern*
375 *analysis and machine intelligence*, 22(8):888–905, 2000.
- 376 [31] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal
377 visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136,
378 January 2015.
- 379 [32] Sergey I Nikolenko. Synthetic data for deep learning. *arXiv preprint arXiv:1909.11512*, 2019.
- 380 [33] Stephan R. Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *IEEE International*
381 *Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2232–2241, 2017.
- 382 [34] Lei Yang. bpycv. <https://github.com/DIYer22/bpycv>.
- 383 [35] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with
384 online hard example mining, 2016.
- 385 [36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition,
386 2015.
- 387 [37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
388 Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large
389 scale visual recognition challenge, 2015.
- 390 [38] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey
391 Zagoruyko. End-to-end object detection with transformers, 2020.

392 **Checklist**

- 393 1. For all authors...
- 394 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
395 contributions and scope? [Yes]
- 396 (b) Did you describe the limitations of your work? [Yes]
- 397 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 398 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
399 them? [Yes]
- 400 2. If you are including theoretical results...
- 401 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 402 (b) Did you include complete proofs of all theoretical results? [N/A]
- 403 3. If you ran experiments (e.g. for benchmarks)...
- 404 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
405 mental results (either in the supplemental material or as a URL)? [Yes] All code and
406 data are available at <https://ishape.github.io/>
- 407 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
408 were chosen)? [Yes] See the beginning of the section Experiments.
- 409 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
410 ments multiple times)? [No]
- 411 (d) Did you include the total amount of compute and the type of resources used (e.g., type
412 of GPUs, internal cluster, or cloud provider)? [Yes] For the proposed ASIS algorithm,
413 we use 4 2080Ti GPUs, 64-core CPUs, and it takes about 1 day to train a sub-dataset.
414 Other benchmark methods use 8 2080Ti GPUs for training each sub-dataset.
- 415 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 416 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 417 (b) Did you mention the license of the assets? [Yes] iShape dataset will be released under
418 CC0 license
- 419 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 420 (d) Did you discuss whether and how consent was obtained from people whose data you're
421 using/curating? [N/A]
- 422 (e) Did you discuss whether the data you are using/curating contains personally identifiable
423 information or offensive content? [N/A]
- 424 5. If you used crowdsourcing or conducted research with human subjects...
- 425 (a) Did you include the full text of instructions given to participants and screenshots, if
426 applicable? [No] In this paper, the iShape-Antenna dataset are collected and annotated
427 by ourselves, and it only took a few days to complete. The remaining dataset is
428 synthesized using open source software Blender.
- 429 (b) Did you describe any potential participant risks, with links to Institutional Review
430 Board (IRB) approvals, if applicable? [N/A]
- 431 (c) Did you include the estimated hourly wage paid to participants and the total amount
432 spent on participant compensation? [N/A]