# Amodal Segmentation for Laparoscopic Surgery Video Instruments

Ruohua Shi[1,2*], Zhaochen Liu[1,2,3*], Lingyu Duan[1,2], Tingting Jiang[1,2,4✉]

[1] National Engineering Research Center of Visual Technology, School of Computer Science, Peking University, Beijing, China
[2] State Key Laboratory of Multimedia Information Processing, School of Computer Science, Peking University
[3] AI Innovation Center, School of Computer Science, Peking University
[4] National Biomedical Imaging Center, Peking University
{shiruohua,dreamerliu,lingyu,ttjiang}@pku.edu.cn

**Abstract.** Segmentation of surgical instruments is crucial for enhancing surgeon performance and ensuring patient safety. Conventional techniques such as binary, semantic, and instance segmentation share a common drawback: they do not accommodate the parts of instruments obscured by tissues or other instruments. Precisely predicting the full extent of these occluded instruments can significantly improve laparoscopic surgeries by providing critical guidance during operations and assisting in the analysis of potential surgical errors, as well as serving educational purposes. In this paper, we introduce Amodal Segmentation to the realm of surgical instruments in the medical field. This technique identifies both the visible and occluded parts of an object. To achieve this, we introduce a new Amoal Instruments Segmentation (AIS) dataset, which was developed by reannotating each instrument with its complete mask, utilizing the 2017 MICCAI EndoVis Robotic Instrument Segmentation Challenge dataset. Additionally, we evaluate several leading amodal segmentation methods to establish a benchmark for this new dataset.

## 1 Introduction

As minimally invasive surgical robots advance, computer vision and machine learning-based assistive systems have become increasingly crucial in boosting surgeon performance and patient safety. However, numerous challenges arise in analyzing the data captured by surgical cameras. Surgical instrument segmentation serves as a critical and indispensable element in a range of computer-assisted interventions.

Recent research has made strides in addressing these challenges, with many researchers exploring solutions [36,25,10,8,23,33]. Notably, a U-Net-based model [25] clinched the top spot at the 2017 MICCAI EndoVis Robotic Instrument Segmentation Challenge [1]. Subsequent studies, such as [23], have a generative-adversarial approach for unsupervised surgical tool segmentation of optical-flow

---
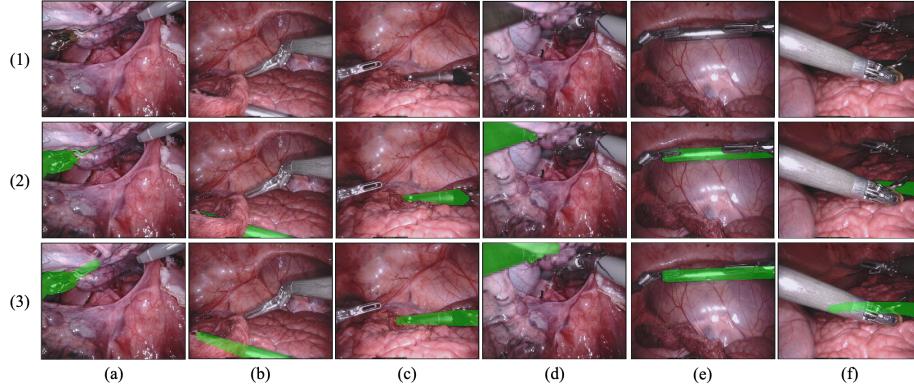* These authors contributed equally.

**Fig. 1.** Segmentation examples. (1) Frames from laparoscopic surgery; (2) Instance segmentation ground truth; (3) Segmentation masks with occluded parts.

images. Yet, these advancements predominantly focus on the visible portions of surgical instruments. Occlusions, whether between instruments and tissues or among the instruments themselves, can obscure critical data during clinical procedures, complicating the surgical process. For illustration, Fig. 1 shows the selected frames from laparoscopic surgery videos. The second and third rows distinctly display instrument masks with and without occlusions. It is evident that standard instance segmentation masks fail to retain the original shape of the instruments, thereby complicating the identification of instrument types from a single frame.

Therefore, accurately predicting occluded instruments can be beneficial in surgeries. Specifically, there are three primary applications: (a) During surgery, predicting the obscured parts of surgical instruments can assist surgeons by providing crucial visual cues about the position and orientation of the tools, ensuring the accuracy of their maneuvers. For example, if the occluded head of an instrument is accurately predicted, surgeons can better judge the alignment of the instrument relative to the target area. (b) Post-operatively, accurate reconstruction of the surgical process through video analysis allows medical professionals to assess whether the instruments were used correctly and if the procedures met the required standards. This includes determining the types of surgical instruments employed and evaluating the precision of the interactions with the patient, thereby assessing the potential impact on patient safety, including any adverse events. (c) For educational purposes, accurately tracking the key points and trajectories of surgical instruments allows instructors to dissect and elucidate optimal surgical techniques. It also highlights crucial visual indicators during procedural training.

Segmenting both the visible and occluded regions of each object instance simultaneously is referred to as **Amodal Segmentation**. This concept is a recent development that builds upon instance-aware segmentation. Numerous method-

ologies have been introduced to tackle this issue [4,11,28,31]. To the best of our knowledge, the recent amodal algorithms that achieved state-of-the-art performance is SAM-based [18], exploiting the powerful feature extraction capability provided by the large-scale foundation model [12]. However, all these amodal segmentation approaches have predominantly been applied to natural images and have not been adapted for the medical field.

In this work, for the first time, we introduce the amodal segmentation to the area of laparoscopic surgery, and propose the first amodal dataset for surgical instruments, named **AIS**, which is based on the 2017 Robotic Instrument Segmentation Challenge dataset [1]. Unlike traditional datasets, AIS includes labels for the full mask of each instrument, covering both visible and occluded regions. Additionally, we evaluated several leading amodal segmentation methods to establish a benchmark for this innovative dataset.

In summary, our paper contains the following contributions:

- For the first time, amodal segmentation is applied to surgical instruments within the medical field.
- A novel medical amodal segmentation dataset is developed specifically tailored for this task.
- A benchmark is introduced to evaluate the accuracy of predicting instruments, including their occluded parts, using our new dataset.

## 2   Related Work

### 2.1   Surgery Video Instrument Segmentation

Various advancements have been made in semantic segmentation of surgical surgery, incorporating various techniques and methodologies, mainly including combining segmentation networks with attention mechanisms [20,24] and exploring some synthetic data for semi-supervised training [23,6]. Ni et al. propose an attention-guided lightweight network, utilizing depth-wise separable convolution as the basic unit to reduce computational costs, thereby performing surgical instrument segmentation in real-time [20]. Shen et al. employ a lightweight encoder and branch aggregation attention mechanism to remove noise caused by reflection, and water mist to improve segmentation accuracy and achieve a lightweight model [24]. Sestini et al. designed a fully unsupervised method for the segmentation of binary surgical instruments relying only on implicit motion information and a priori knowledge of the shape of the instrument [23]. Luis et al. facilitate this task by generating large amounts of trainable data by synthesizing surgical instruments with real surgical backgrounds [6]. However, previous approaches mainly focus on real-time or semi-supervised, unsupervised learning, and there has been no research on interactive segmentation. Moreover, their methods cannot distinguish the occluded instruments. Therefore, we explored the segmentation of both visible and occluded instruments in surgery.

## 2.2   Amodal Segmentation

Amodal segmentation task was proposed in 2016 to predict the complete shape of target object including both the visible parts and the occluded parts [15]. Possessing great significance to our seeking visual intelligence, amodal segmentation arouses increasing attention in the academic community. Besides direct methods [15,36,22], a series of elaborate approaches have been designed with diverse concepts involved to achieve better performance, such as depth relationship [35], region correlation [4,11,28], shape priors [32,16,5], and compositional models [31]. Recently, a SAM-based approach [18] achieved state-of-the-art performance, well exploiting the mighty feature extraction capability provided by the large-scale foundation model [12]. Noticing the labor-intensive and error-prone challenges in the annotation of amodal masks, researchers also propose many weakly supervised approaches for amodal segmentation using only box-level supervision or self supervision [34,19,13,14,26,17].

Based on the advances of amodal segmentation algorithms, researchers develop various applications. For example, utilizing amodal segmentation methods, intelligence systems like autonomous driving and robotic grasping can enhance the safety and the reliability [22,2,29,30,9], while many new implementation paths emerge in the fields of diminished reality (DR) and novel view synthesis [7,16,21]. Though numerous related work is delivered, the potential of amodal segmentation in the medical field has never been explored. With our newly released dataset, we select several typical segmentation methods for experiments, revealing the notable value in medical applications.

## 3   Dataset

We introduce the AIS dataset by re-annotating the 2017 Robotic Instrument Segmentation Challenge dataset [1]. The dataset consists of 10 videos, each with 300 frames. It covers different abdominal porcine procedures recorded using the da Vinci Xi systems. In this paper, we relabeled the 3000 frames at a resolution of 1024×1280. The annotation includes semantic instance-level amodal segmentation and each object is manually assigned a class label. There are a total of 7084 objects. The dataset is split into training and test sets based on the splitting rule of the 2017 Robotic Instrument Segmentation Challenge: take the first 225 frames of 8 sequences as training data and keep the last 75 frames of those 8 sequences as test data. 2 of the full 300 frame sequences were kept as test sequences.

### 3.1   Dataset Acquisition

To annotate the dataset, we employed a widely-used public annotation tool named *labelme* [27]. This tool allows for the interactive labeling of surgical instruments on a frame-by-frame basis, enabling manual annotations directly within the software interface.
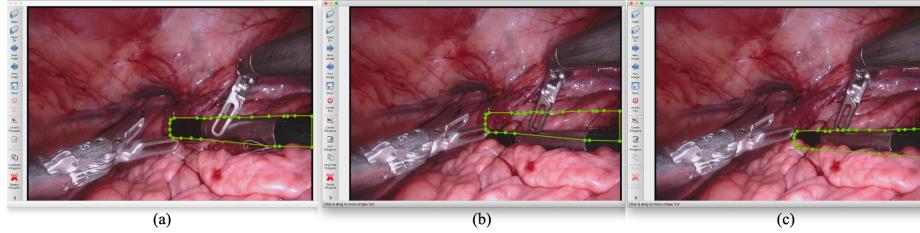
**Fig. 2.** Illustration of annotation process. (a) displays a polygon annotation formed by various key points and line segments. Once the labeling of this frame is complete, the annotation can be carried forward to the next frame for further labeling, as depicted in (b). By adjusting the positions of the key points within the polygon, the final annotation of the instrument in (b) is presented in (c).

Consider a practical example of our annotation methodology as illustrated in Fig.2(a), the middle instrument is annotated with a polygon that includes several key points and line segments. In this interface, annotations can be refined by adjusting the positions of key points (illustrated as green points). After completing the annotation for one frame, it can be carried over to the next frame, as depicted in Fig. 2(b). Given that the movement of surgical instruments between consecutive frames is typically minimal, the annotation for the current frame can often be achieved by simply repositioning the key points from the previous frame to align with the instrument's shape. This process completes the annotation for the middle instrument in the current frame, with the final result displayed in Fig. 2(c). By employing this annotation strategy, we attempt to reconstruct the original form of the instruments as accurately as possible, even for occluded parts. Our guiding principle is to ensure that the prediction of the occluded sections adheres closely to the instrument's actual shape and follows its natural motion trajectory.

**Table 1.** Amodal segmentation dataset statistics.

| Dataset | COCOA | COCOA-cls | D2S | Ours |
|---|---|---|---|---|
| Image/Video | Image | Image | Image | Video |
| Resolution | 275K pix | 275K pix | 3M pix | 1M pix |
| | - | - | 1440×1920 | 1024×1280 |
| # of images | 5073 | 3499 | 3499 | 3000 |
| # of instances | 46314 | 10562 | 28720 | 7084 |
| # of occluded instances | 28106 | 5175 | 16337 | 1455 |
| Avg. occlusion rate | 18.8% | 10.7% | 15.0% | 20.54% |

## 3.2   Dataset Analysis

We compare the proposed AIS dataset with other prominent amodal segmentation datasets from real-world scenes, such as COCOA[37], COCOA-cls[37,4], and D2S [4,3]. Specifically, we examine the differences in resolution, the number of instances included, and the rate of occlusion across these datasets. As shown in Table 1, our dataset exhibits a higher average occlusion rate and represents the first amodal segmentation dataset specifically developed for the medical field.

# 4   Benchmarks

## 4.1   Selected Methods

Our newly released dataset brings novel medical application scenarios for amodal segmentation. To better reveal the capabilities of existing approaches, we select several typical methods and benchmark their performance on this dataset, which are SAM [12], AISFormer [28], C2F-Seg [5], and PLUG [18].

SAM, the segment anything model, is the most influential foundation model for general segmentation task, which is trained on billions of object masks thus possessing strong abilities on diverse downstream tasks. The other three methods are specifically designed for amodal segmentation. PLUG is the state-of-the-art approach now, which is based on SAM and introduces the parallel LoRA structure and the uncertainty guidance. AISFormer and C2F-Seg are other two most recent methods. AISFormer first injects the transformer backbone into the amodal segmentation task, while C2F-Seg introduces the vector-quantization model and the coarse-to-fine framework.

## 4.2   Evaluation Metrics

As related work, we choose intersection-over-union (IoU) as the evaluation metric. The IoU of a predicted mask equals the ratio of the intersecting area with corresponding ground truth mask to the area of their union. In order to better reflect the performance, we calculate the average IoU on each sub-testset (from 1 to 10) and the mean IoU of all sub-testsets.

## 4.3   Experimental Setup

For SAM, we adopt the largest ViT-H version pretrained model to achieve its best performance. For the other three methods, we conduct training on our proposed dataset. The settings and parameters adopted during the training process are consistent with the original papers.

It is worth explaining that the input mode of bounding boxes. Since the models require clear regions of interest, we input bounding boxes as guidance in all these methods. SAM, AISFormer, and PLUG directly receive input bounding boxes, while C2F-Seg uses the predictions of SAM as input thus indirectly
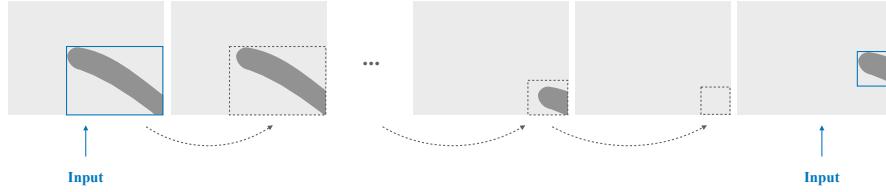
**Fig. 3. The input mode of bounding boxes.** In order to reduce the number of manual inputs, we use the bounding box of the predicted mask in the previous frame as the input bounding box of the next frame when possible.

receiving the bounding boxes (C2F-Seg generates amodal mask based on the input visible mask. For fairness, we adopt the output of SAM as the input visible mask of C2F-Seg.). Considering practical needs, we hope to reduce the number of manual inputs as much as possible for the convenience of users. Actually, we notice that the difference between two adjacent frames of a video is very small, and we can approximately assume that the bounding boxes in the two frames are almost the same. So as shown in Fig. 3, we first input the bounding box of the first frame, and then use the bounding box of the predicted mask of the previous frame for each subsequent frame until a certain frame without the corresponding bounding box in the previous frame occurs, we then manually input the bounding box again.

## 5    Results

### 5.1    Performance comparison

As shown in Table 2, all these segmentation methods can basically predict the amodal masks of surgery instruments quite well. The pretrained SAM without any fine-tuning can already reach practical performance, achieving 85.94 on mean IoU. Specifically designed, the other three methods can provide more accurate predictions. AISFormer and C2F-Seg attain 86.65/88.17 on mean IoU, which are 0.71/2.23 higher than SAM. The state-of-the-art method PLUG maintains the leading position and reaches 89.25 on the mean IoU, which beats SAM by 3.31. It can be observed that the performance is close on some relatively simple sub-testsets, such as testset-4, while on some complex sub-testsets like testset-1, the advantage of tailored amodal segmentation methods especially PLUG is evident.

To intuitively display the performance of these methods, we choose one frame from each sub-testset, totaling ten frames. The ten sets of qualitative results are shown in Fig. 5. Relatively speaking, though the predictions of other methods are acceptable, PLUG can determine more precise and complete shapes. For the rod-shaped instrument in the 1st row, the prediction of PLUG is clear to identify showing smooth boundary and covering more area. For the semi-inserted instrument in the 8th row, C2F-Seg and PLUG can segment the needle head, while the prediction of PLUG is thinner and closer to reality.

**Table 2.** The performance comparison of selected segmentation methods on the surgery dataset. Numbers 1-10 represent ten sub-testsets.

| Methods | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | Mean |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SAM [12] | 75.1 | 80.5 | 91.5 | 93.1 | 81.1 | 85.8 | 83.7 | 91.2 | 85.9 | 91.5 | 85.94 |
| AISFormer [28] | 74.6 | 82.6 | 91.0 | 93.6 | 83.0 | 83.9 | 88.6 | 93.0 | 85.2 | 91.0 | 86.65 |
| C2F-Seg [5] | 75.8 | 85.6 | 92.6 | 93.5 | 85.0 | 86.1 | 89.1 | 94.2 | 87.0 | 92.8 | 88.17 |
| PLUG [18] | 79.8 | 86.2 | 93.2 | 93.6 | 86.1 | 87.1 | 90.6 | 94.4 | 88.4 | 93.1 | 89.25 |

### 5.2   Hard Case

In the experiment, we find that some excessively occluded instruments (as shown in Fig. 4) in a fraction of frames are a common challenge for these methods. These instruments are almost completely occluded, making it difficult to obtain effective visible mask or other information, thus bringing confusion to methods inspired by extracted clues within the frame. In the example shown in Fig. 4, even the state-of-the-art PLUG approach predicts a quite tortuous boundary. To tackle the hard case, a pertinent method designed for video input that integrates the context of the previous and subsequent frames may be needed. How to implement such a specific method? We leave this issue to future work.
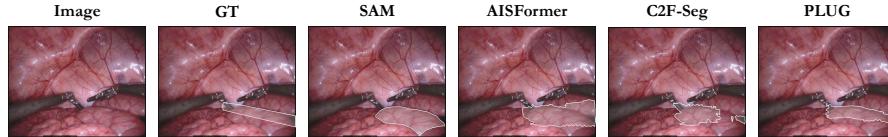


**Fig. 4. An example of the excessively occluded case.** The instrument is almost completely occluded in this frame.
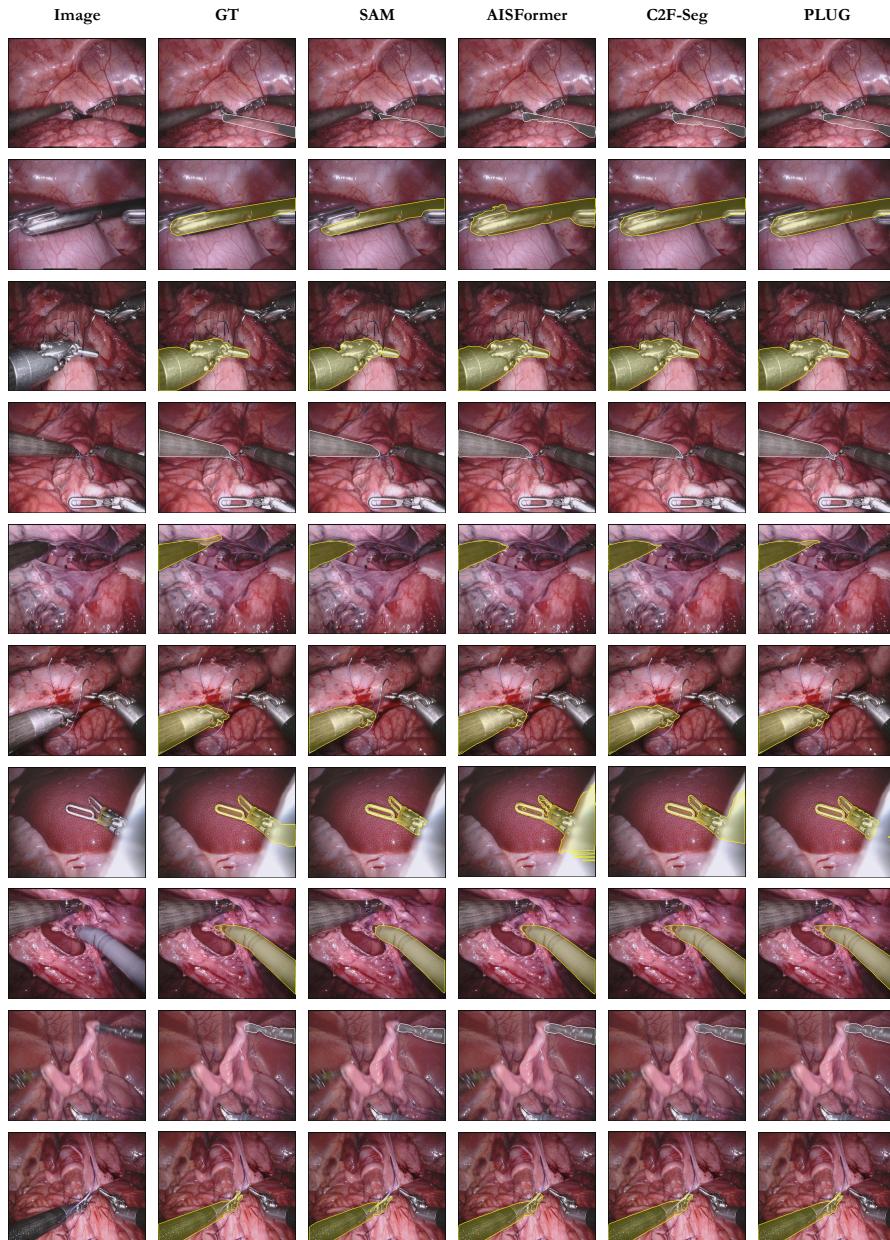
**Fig. 5. Qualitative results.** The qualitative comparison of predicted amodal masks from SAM, AISFormer, C2F-Seg and PLUG. These ten rows are chosen from the first to tenth sub-testset in order from top to bottom.

# References

1. Allan, M., Shvets, A., Kurmann, T., Zhang, Z., Duggal, R., Su, Y.H., Rieke, N., Laina, I., Kalavakonda, N., Bodenstedt, S., et al.: 2017 robotic instrument segmentation challenge. arXiv preprint arXiv:1902.06426 (2019)
2. Breitenstein, J., Fingscheidt, T.: Amodal cityscapes: a new dataset, its generation, and an amodal semantic segmentation challenge baseline. In: IEEE Intelligent Vehicles Symposium. pp. 1018–1025 (2022)
3. Follmann, P., Bottger, T., Hartinger, P., Konig, R., Ulrich, M.: MVTec D2S: densely segmented supermarket dataset. In: Proceedings of the European Conference on Computer Vision. pp. 569–585 (2018)
4. Follmann, P., König, R., Härtinger, P., Klostermann, M., Böttger, T.: Learning to see the invisible: End-to-end trainable amodal instance segmentation. In: Proceedings of the IEEE Winter Conference on Applications of Computer Vision. pp. 1328–1336 (2019)
5. Gao, J., Qian, X., Wang, Y., Xiao, T., He, T., Zhang, Z., Fu, Y.: Coarse-to-fine amodal segmentation with shape prior. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1262–1271 (2023)
6. Garcia-Peraza-Herrera, L.C., Fidon, L., D'Ettorre, C., Stoyanov, D., Vercauteren, T., Ourselin, S.: Image compositing for segmentation of surgical tools without manual annotations. IEEE Transactions on Medical Imaging 40(5), 1450–1460 (2021)
7. Gkitsas, V., Sterzentsenko, V., Zioulis, N., Albanis, G., Zarpalas, D.: Panodr: Spherical panorama diminished reality for indoor scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3716–3726 (2021)
8. González, C., Bravo-Sánchez, L., Arbelaez, P.: Isinet: an instance-based approach for surgical instrument segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 595–605 (2020)
9. Inagaki, Y., Araki, R., Yamashita, T., Fujiyoshi, H.: Detecting layered structures of partially occluded objects for bin picking. In: Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems. pp. 5786–5791 (2019)
10. Jin, Y., Cheng, K., Dou, Q., Heng, P.A.: Incorporating temporal prior from motion flow for instrument segmentation in minimally invasive surgery video. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 440–448 (2019)
11. Ke, L., Tai, Y.W., Tang, C.K.: Deep occlusion-aware instance segmentation with overlapping bilayers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4019–4028 (2021)
12. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., Dollar, P., Girshick, R.: Segment anything. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4015–4026 (2023)
13. Kortylewski, A., He, J., Liu, Q., Yuille, A.L.: Compositional convolutional neural networks: A deep architecture with innate robustness to partial occlusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8940–8949 (2020)
14. Kortylewski, A., Liu, Q., Wang, A., Sun, Y., Yuille, A.: Compositional convolutional neural networks: A robust and interpretable model for object recognition under occlusion. International Journal of Computer Vision 129, 736–760 (2021)

15. Li, K., Malik, J.: Amodal instance segmentation. In: Proceedings of the European Conference on Computer Vision. pp. 677–693 (2016)
16. Li, Z., Ye, W., Jiang, T., Huang, T.: 2D amodal instance segmentation guided by 3D shape prior. In: Proceedings of the European Conference on Computer Vision. pp. 165–181 (2022)
17. Liu, Z., Li, Z., Jiang, T.: BLADE: Box-level supervised amodal segmentation through directed expansion. In: Proceedings of the AAAI Conference on Artificial Intelligence (2024)
18. Liu, Z., Qiao, L., Chu, X., Jiang, T.: PLUG: Revisiting amodal segmentation with foundation model and hierarchical focus. arXiv preprint arXiv:2405.16094 (2024)
19. Nguyen, K., Todorovic, S.: A weakly supervised amodal segmenter with boundary uncertainty estimation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7396–7405 (2021)
20. Ni, Z.L., Bian, G.B., Hou, Z.G., Zhou, X.H., Xie, X.L., Li, Z.: Attention-guided lightweight network for real-time segmentation of robotic surgical instruments. In: 2020 IEEE International Conference on Robotics and Automation. pp. 9939–9945 (2020)
21. Pintore, G., Agus, M., Almansa, E., Gobbetti, E.: Instant automatic emptying of panoramic indoor scenes. IEEE Transactions on Visualization and Computer Graphics 28(11), 3629–3639 (2022)
22. Qi, L., Jiang, L., Liu, S., Shen, X., Jia, J.: Amodal instance segmentation with kins dataset. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3014–3023 (2019)
23. Sestini, L., Rosa, B., De Momi, E., Ferrigno, G., Padoy, N.: Fun-sis: A fully unsupervised approach for surgical instrument segmentation. Medical Image Analysis 85, 102751 (2023)
24. Shen, W., Wang, Y., Liu, M., Wang, J., Ding, R., Zhang, Z., Meijering, E.: Branch aggregation attention network for robotic surgical instrument segmentation. IEEE Transactions on Medical Imaging (2023)
25. Shvets, A.A., Rakhlin, A., Kalinin, A.A., Iglovikov, V.I.: Automatic instrument segmentation in robot-assisted surgery using deep learning. In: IEEE International Conference on Machine Learning and Applications. pp. 624–628 (2018)
26. Sun, Y., Kortylewski, A., Yuille, A.: Amodal segmentation through out-of-task and out-of-distribution generalization with a bayesian model. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1215–1224 (2022)
27. Torralba, A., Russell, B.C., Yuen, J.: Labelme: Online image annotation and applications. Proceedings of the IEEE 98(8), 1467–1484 (2010)
28. Tran, M., Vo, K., Yamazaki, K., Fernandes, A., Kidd, M., Le, N.: Aisformer: Amodal instance segmentation with transformer. In: Proceedings of the British Machine Vision Conference (2022)
29. Wada, K., Kitagawa, S., Okada, K., Inaba, M.: Instance segmentation of visible and occluded regions for finding and picking target from a pile of objects. In: Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems. pp. 2048–2055 (2018)
30. Wada, K., Okada, K., Inaba, M.: Joint learning of instance and semantic segmentation for robotic pick-and-place with heavy occlusions in clutter. In: Proceedings of the International Conference on Robotics and Automation. pp. 9558–9564 (2019)
31. Wang, A., Sun, Y., Kortylewski, A., Yuille, A.L.: Robust object detection under occlusion with context-aware compositionalnets. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12645–12654 (2020)

32. Xiao, Y., Xu, Y., Zhong, Z., Luo, W., Li, J., Gao, S.: Amodal segmentation based on visible region segmentation and shape prior. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 2995–3003 (2021)
33. Yue, W., Zhang, J., Hu, K., Xia, Y., Luo, J., Wang, Z.: Surgicalsam: Efficient class promptable surgical instrument segmentation. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 38, pp. 6890–6898 (2024)
34. Zhan, X., Pan, X., Dai, B., Liu, Z., Lin, D., Loy, C.C.: Self-supervised scene de-occlusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3784–3792 (2020)
35. Zhang, Z., Chen, A., Xie, L., Yu, J., Gao, S.: Learning semantics-aware distance map with semantics layering network for amodal instance segmentation. In: Proceedings of the 27th ACM International Conference on Multimedia. pp. 2124–2132 (2019)
36. Zhu, Y., Tian, Y., Metaxas, D., Dollár, P.: Semantic amodal segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1464–1472 (2017)
37. Zhu, Y., Tian, Y., Metaxas, D., Dollár, P.: Semantic amodal segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1464–1472 (2017)