

# PartSTAD: 2D-to-3D Part Segmentation Task Adaptation — Supplementary Material

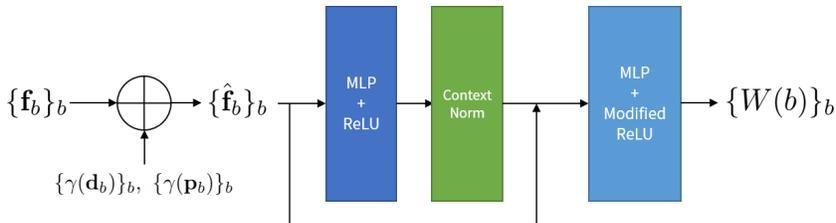
Hyunjin Kim<sup>1†</sup> and Minhyuk Sung<sup>2</sup>

<sup>1</sup> KRAFTON Inc., South Korea

<sup>2</sup> KAIST, South Korea

rlaguswls98@krafton.com, mhsung@kaist.ac.kr

In this supplementary material, we present additional implementation details (Sec. S.1), additional results on the ablation studies (Sec. S.2, Sec. S.3, Sec. S.4, Sec. S.5, and Sec. S.6), how bounding box weights change compared to GLIP confidence scores (Sec. S.7), the reason for not replacing GLIP with SAM (Sec. S.8), more analysis on results (Sec. S.9), outcomes with scanned data (Sec. S.10), and comprehensive results covering all PartNet-Mobility categories (Sec. S.11, Sec. S.12, and Sec. S.13).



**Figure S1:** Architecture of the weight prediction network used in PARTSTAD.  $\mathbf{f}_b$ ,  $\mathbf{d}_b$ , and  $\mathbf{p}_b$  represent the bounding box feature, view direction, and position in the 2D image of the bounding box  $b$ , respectively.  $\gamma$  denotes the positional encoding function. This network takes all bounding box features of a single object as input and outputs bounding box weights.

## S.1 Additional Implementation Details

**Network Architecture.** Fig. S1 shows the detailed network architecture of the weight prediction network of PARTSTAD. As mentioned in Sec. 4.2 of the main paper, the network consists of small shared two-layer MLP with the context normalization [10] layer between them to embed context information. This light architecture design is inspired by LoRA [1] and PartSLIP [5] which add small learnable parameters while keeping the original pretrained model parameters.

We add positional encoded vectors to each bounding box feature to embed the positional information. For given bounding box feature  $\mathbf{f}_b$ , we do not directly

<sup>†</sup> This work was conducted when the author was at KAIST.

feed  $\mathbf{f}_b$  to the network but feed  $\hat{\mathbf{f}}_b$  which concatenates  $\mathbf{f}_b$ , positional encoding of view direction  $\mathbf{d}_b$ , and positional encoding of 2D bounding box position  $\mathbf{p}_b$ . Positional encoding  $\gamma$  is defined as below:

$$\gamma(x_1, x_2, \dots, x_n) = \bigoplus_{i=1}^n (x_i, \sin(2^0 \pi x_i), \cos(2^0 \pi x_i), \dots, \sin(2^{L-1} \pi x_i), \cos(2^{L-1} \pi x_i)), \quad (1)$$

where  $\oplus$  indicates the concatenation operation and  $L$  is set to 10 in our experiments. Thus the input  $\hat{\mathbf{f}}_b$  of the weight prediction network can be written as below:

$$\hat{\mathbf{f}}_b = \mathbf{f}_b \oplus \gamma(\mathbf{d}_b) \oplus \gamma(\mathbf{p}_b). \quad (2)$$

We initialize all network parameters  $\theta_i \sim \mathcal{N}(0, \epsilon)$ , where  $\epsilon$  is a very small number ( $\epsilon$  is set to 0.0001 in our experiments) so that the initial MLP output becomes 0. Since the last layer is modified ReLU layer  $\phi$  (Eq. 7 of the main paper), the initial weight is set to  $\tau$ . This is inspired by zero convolution of ControlNet [11], and this makes the training more stable by preventing drastic weight changes.

Instead of using an attention-based network to consider the relations between bounding boxes, we opt to add context normalization [10] between two MLP networks in the weight prediction network. This allowed us to keep the network lightweight while still considering the relations between bounding boxes.

**Design of Modified ReLU.** The modified ReLU function (Eq. 7 of the main paper) is designed to set the initial value of the bbox weight, which is the output of the weight prediction network, to a value  $\tau > 0$ .

**Training Details.** We generate 2D images of size  $800 \times 800$  from 10 fixed viewpoints for each object that is normalized to fit in a unit sphere using the Pytorch3D point cloud renderer with the fixed camera distance 2.2 following the same procedure as described in PartSLIP [5]. After rendering, we obtain bounding boxes for each image using the GLIP [4]. Subsequently, all bounding box features corresponding to bounding boxes from a single object are simultaneously fed into the weight prediction network to calculate the weights. Training is conducted using a single RTX 3090 GPU.

## S.2 Results with Varying Parameters

**Number of Views.** Tab. S1 presents the results of the ablation study on the number of views. It is observed that as the number of views increases, mIoU also increases, with the most significant difference observed when the view changes from 5 to 10. This highlights that having too few samples of bounding boxes used in training can lead to suboptimal results.

**Table S1:** Ablation study on the number of views.

# of Views	5	10	15	20
mIoU	59.4	65.0	66.3	<b>67.2</b>

**Table S2:** Ablation study on the number of training data. The experiment is conducted on the StorageFurniture category as it only has more than 128 shapes (346 in total).

# of Training Data	8	16	32	64	128
mIoU	57.0	56.7	58.0	57.8	<b>60.1</b>

**Table S3:** Ablation study on the hyperparameter  $\tau$ . The mIoU is measured for five object categories: Chair, Table, StorageFurniture, Faucet, and TrashCan, which are the five categories with the most test data.

Initial Weight ( $\tau$ )	1	5	10	15	20
mIoU (5 categories)	54.8	55.6	<b>55.8</b>	55.6	55.6

**Number of Training Data.** Tab. S2 presents the results of the ablation study based on the number of training data. Only the StorageFurniture category has more than 128 data, so the experiment is conducted only for this category. There is a tendency for mIoU to increase as the number of training data increases, but the difference is not significant. This demonstrates that even using only 8 data points can yield sufficiently good results.

**Hyperparameter  $\tau$ .** Tab. S3 shows the ablation results for the hyperparameter  $\tau$ . The results are best when  $\tau$  is 10, but they also demonstrate that the  $\tau$  value does not significantly impact the results when it is greater than zero.

### S.3 Random Viewpoints vs. Fixed Viewpoints

**Table S4:** Comparison between cases of rendering with random viewpoints and fixed viewpoints.

Category	Chair	Kettle	Lamp	Suitcase	Scissors
Random	<b>85.5</b>	82.5	66.8	<b>71.3</b>	64.8
Fixed (Ours)	85.3	<b>84.2</b>	<b>68.4</b>	68.3	<b>68.5</b>

Tab. S4 shows the results with random viewpoints and shows that the outcomes are not sensitive to the choice of viewpoints.

## S.4 Cross-Entropy Loss vs. mRIoU Loss

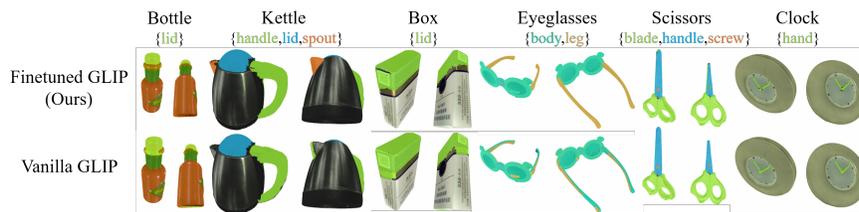
**Table S5:** Comparison with the cases of using cross-entropy loss and mRIoU loss (Ours).

Method	mIoU (%)
PartSLIP [5] + <i>SAM Mask Integration</i> (Baseline)	61.9
PARTSTAD + Cross-Entropy - mRIoU	63.5
PARTSTAD + Cross-Entropy	64.5
PARTSTAD (Ours)	<b>65.0</b>

As mentioned in Sec. 4.1 of the main paper, the use of mRIoU loss is crucial for achieving significant improvement in our task adaptation. To demonstrate the effectiveness of our mRIoU loss, we conduct an experiment comparing it with the alternative, cross-entropy loss.

Tab. S5 shows the ablation results for different loss types. Baseline at the 1st row in Tab. S5 represents the result which only applies *SAM* [3] *mask integration* to PartSLIP [5] (same as our method *without weight prediction*). When using the commonly used cross-entropy loss for segmentation tasks, the mIoU decreases by 1.5%p compared to using the mRIoU loss. Even when both losses are used together, the mIoU decreases by 0.5%p. This indicates that the mRIoU loss is more effective for 3D segmentation task adaptation, and it shows its highest effectiveness when used alone.

## S.5 Vanilla GLIP vs. Finetuned GLIP



**Figure S2:** Qualitative comparison between vanilla GLIP and finetuned GLIP on OmniObject3D [8], a real-scanned dataset.

Tab. S6 compares results using vanilla GLIP and finetuned GLIP. Vanilla GLIP yields significantly worse results, emphasizing the substantial impact of

**Table S6:** Quantitative comparison between vanilla GLIP and finetuned GLIP on PartNet-Mobility [9] dataset.

Method	Vanilla GLIP	Finetuned GLIP (Ours)
PartSLIP [5]	27.2	58.0
PARTSTAD (Ours)	48.9 (+ 21.7)	<b>65.0 (+ 7.0)</b>

bounding box prediction on the final outcome. At the same time, when we use the vanilla GLIP, our weight prediction network significantly improves the mIoU from 27.2 to 48.9. This indicates that our weight prediction network is more effective when the 2D prediction is inaccurate.

In our experiments with scanned objects (Sec. 5.4 of the main paper), we used a finetuned GLIP instead of the vanilla GLIP, as it exhibited better performance in detecting the parts, even for real images, due to its finetuning for the specific parts. Fig. S2 illustrates the superior performance of the finetuned GLIP compared to the vanilla GLIP for the OmniObject3D dataset.

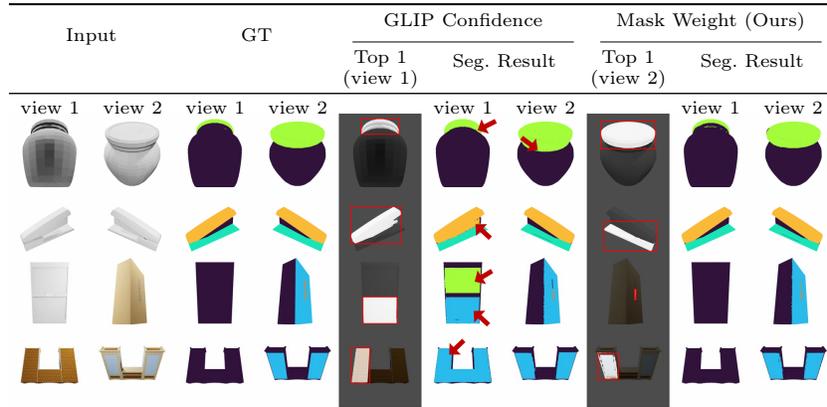
## S.6 GLIP Confidence Score vs. Mask Weight

**Table S7:** Comparison with the cases of using GLIP confidence score as weight and predicted mask weight (Ours).

Method	mIoU (%)
PartSLIP [5] + <i>SAM Mask Integration</i> (Baseline)	61.9
PARTSTAD + GLIP Conf.	53.3
PARTSTAD + Normalized GLIP Conf.	62.3
PARTSTAD (Ours)	<b>65.0</b>

It is worth noting that the GLIP [4] model also outputs a confidence score for each predicted bounding box. This implies that we can consider using the GLIP confidence scores as weights in the voting scheme ( $W(b)$  in Eq. 5 of the main paper). Thus, we compare the results when using GLIP confidence scores and the predicted mask weights from our method.

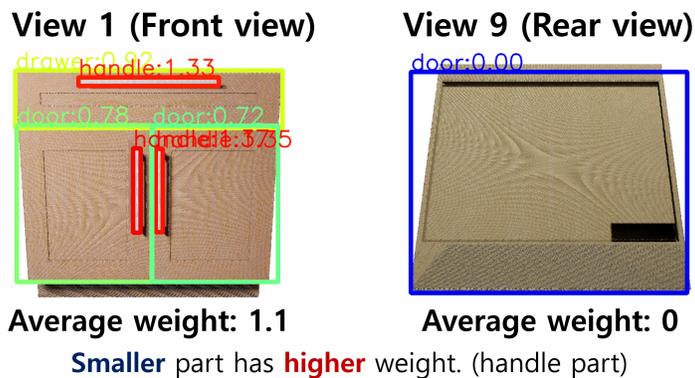
Tab. S7 presents the comparison results between using GLIP confidence as weights and using weights predicted from the network. We compare those methods with the baseline (1st row in Tab. S7) which only applies *SAM* [3] *mask integration* to PartSLIP [5] (same as our method *without weight prediction*). As the confidence scores from GLIP range in  $[0, 1]$ , utilizing them directly as weights causes an overall score reduction, resulting in segments are not generated. Consequently, the outcome is notably poor, with a 53.3 mIoU. To ensure a



**Figure S3:** Comparison of results using GLIP confidence scores and Mask Weights. The 5th and the 8th columns depict masks with the highest scores (weights), where red rectangles represent bounding boxes from GLIP, and the white regions are segmentation masks after integrating SAM. From the top row to the bottom, each corresponds to the Bottle, Stapler, and two StorageFurniture categories. When using GLIP confidence, the highest score mask for Bottle (1st row) and Stapler (2nd row) includes an incorrect region, leading to inaccurate segmentation (denoted as red arrows). In contrast, our method assigns the highest score to the correct mask, indicating that the incorrect mask has a lower score. Additionally, when using the GLIP confidence score, the highest score mask for the 3rd and the 4th rows each indicates a completely wrong part (the backside of StorageFurniture). However, our method assigns the highest score to the handle and the correct door part at the front side for the 3rd row and the 4th row, respectively.

fair comparison, we normalize the weights to maintain the same sum as before. With these normalized weights, as presented in the second row of Tab. S7, the result becomes 62.3 mIoU, a slight increase of 0.4%p compared to the baseline. However, this is still 2.7%p lower than utilizing predicted mask weights from a network trained with 3D mRIoU loss. In conclusion, our method provides results more optimized for 3D segmentation than GLIP confidence scores, demonstrating the effectiveness of our method.

As shown in Fig S3, GLIP confidence scores occasionally assign the highest score to incorrect bounding boxes, leading to suboptimal segmentation results. In contrast, the weights predicted by our method consistently assign the highest weight to the correct regions. For instance, using the GLIP confidence score as the weight results in the highest score masks for Bottle (1st row) and Stapler (2nd row) including incorrect regions, leading to inaccuracies in segmentation (indicated by red arrows). In contrast, our method assigns the highest score to the correct mask, indicating that the incorrect mask has a lower score. Additionally, with the GLIP confidence score, the highest score masks for the 3rd and the 4th rows each indicate completely wrong parts (the backside of StorageFurniture). However, our method assigns the highest score to the handle and the correct door part at the front side for the 3rd row and the 4th row, respec-



**Figure S4:** The smaller part (handle) has a higher weight compared to the bigger part (door), and the rear view (fewer GT parts) has a lower weight compared to the front view (more GT parts).

tively. Those results demonstrate that utilizing the mask weight predicted from the network trained with 3D mRIoU loss produces more accurate predictions compared to using the GLIP confidence score as the weight.

## S.7 Learned Bounding Box Weights

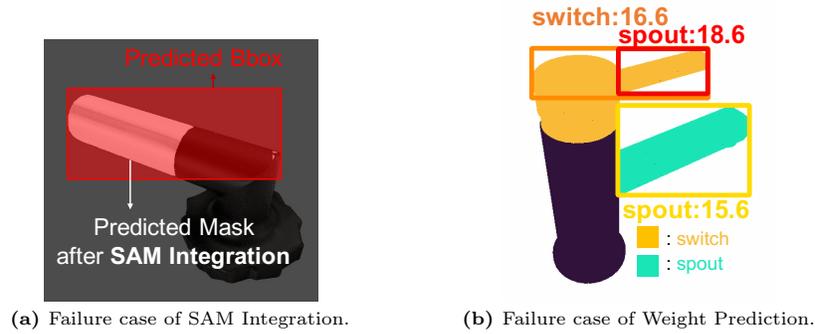
As seen in the example of Fig. S4 for the storage furniture category, smaller parts like handles tend to have relatively higher weights compared to larger parts like doors. Additionally, in rear views where there are no ground truth parts, there is a tendency for the average weight to be lower compared to front views with many ground truth parts. This indicates that learned weight is influenced by both view direction and part labels, unlike the GLIP confidence score, which has a uniform average without distinct trends regarding view and parts.

## S.8 Reasons for Not Replacing GLIP [4] with SAM [3]

SAM [3] allows text prompts as inputs, which could enable direct replacement of GLIP with SAM. However, since the pretrained model supporting text prompts has not been released, we resort to an alternative approach. We serialize GLIP and SAM by using a bounding box predicted by GLIP as an input prompt for SAM.

Note that Grounded-SAM [7] and other recent 2D segmentation methods based on text prompts (e.g., SAM-HQ [2]) also involve the serialization of a bounding box prediction network (like GLIP) and SAM. Grounded-SAM specifically uses GroundingDINO [6] instead of GLIP. We believe that the selection of the bounding box prediction network should not impact our contributions.

## S.9 More Analysis on Results



**Figure S5:** Failure cases of PARTSTAD.

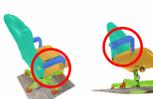
The qualitative and quantitative results in the main paper demonstrate that our PARTSTAD provides more specialized 2D predictions tailored to 3D segmentation compared to PartSLIP [5]. However, in some categories, there are cases where removing specific components from PARTSTAD leads to better results or even where PartSLIP outperforms PARTSTAD (e.g., Faucet category). Fig. S5 illustrates cases where each component performs worse predictions than the baseline.

When the initially given bounding box contains few wrong points but includes many correct points, there are cases where the new mask obtained through SAM does not include the previously contained correct points (Fig. S5a). In such cases, the performance may deteriorate when *SAM mask integration* is applied. In the Faucet category, both the switch part and the spout part protrude prominently, resulting in the initial bounding box containing few irrelevant points. Therefore, it appears that the performance deteriorates when correct points are excluded rather than irrelevant points through SAM mask integration.

Secondly, in visually similar parts, weights might be inaccurately predicted. Fig. S5b illustrates the predicted bounding boxes and weights of the Faucet object, showing that the switch part is predicted as the spout with the highest bounding box weight. In such cases where parts are not visually distinguishable, the weight prediction may not be properly learned. Additionally, adding weight prediction in these cases may lead to a decrease in performance.

Note that such cases are rare and do not significantly impact the overall improvement, as shown in Tab. 1 and Tab. 2 in the main paper.

Additionally, for some parts such as door and drawer sometimes have extremely low IoU. This is mainly caused by GLIP, as it fails to detect the parts due to a lack of data. For example, in the Table class, there is no training data that includes the door part.

Category	Cart	Chair	Dispenser	Display	Faucet
Text Prompt	wheel	arm,back,leg, seat,wheel	head,lid	base,screen, support	spout,switch
Input					
PartSLIP [5]					
PARTSTAD (Ours)					

Category	Kettle	KitchenPot	Storage Furniture	Suitcase	TrashCan
Text Prompt	lid,handle, spout	lid,handle	door,drawer, handle	handle, wheel	footpedal, lid,door
Input					
PartSLIP [5]					
PARTSTAD (Ours)					

**Figure S6:** Qualitative results on real-world scan data. In the highlighted red circle, it is evident that our method achieves more accurate segmentation than PartSLIP [5].

## S.10 Results on Real-World Scanned Data

Fig. S6 illustrates the results of semantic segmentation on the real-world scan data used in PartSLIP [5], which is captured by smartphone. As seen in the figure, our method demonstrates its robustness by successfully predicting not only with higher-quality real-world scans like OmniObject3D [8] as illustrated in Fig. 5 of the main paper but also with lower-quality scan data. Also, our method provides more accurate segmentation than PartSLIP [5]. In the case of the Chair, our method accurately segments the arm part, while PartSLIP fails to do so. For the Kettle, our method better identifies the spout compared to PartSLIP. Additionally, in the cases of StorageFurniture and TrashCan, PartSLIP segments

parts that should not be segmented (the backside of StorageFurniture and the lid of TrashCan). On the other hand, for the KitchenPot, while PartSLIP finds the lid part that our method misses, its boundary is still not perfect. We demonstrate that our approach identifies more accurate parts while simultaneously predicting more precise boundaries.

## S.11 Complete Quantitative Results of PARTSTAD

Tab. S8, Tab. S9, and Tab. S10 show the full table of quantitative results for semantic segmentation, part-aware instance segmentation, and part-agnostic instance segmentation, respectively. Overall, our method demonstrates the best results across whole categories and parts. Please refer to the complete table on the subsequent page for comprehensive information. Moreover, after quantitative result tables, additional qualitative results are illustrated in Sec. S.12 and Sec. S.13.



**Table S9:** Full table of part-aware instance segmentation results on the PartNet-Mobility [9] dataset.

		Baselines		Ablations					Baselines		Ablations		
Category	Part	PartSLIP	<i>w/o</i> Weight Pred.	<i>w/o</i> SAM Integ.	Ours		Category	Part	PartSLIP	<i>w/o</i> Weight Pred.	<i>w/o</i> SAM Integ.	Ours	
Bottle	lid	73.4	74.9	<b>79.5</b>	77.4								
Box	lid	56.3	<b>72.2</b>	60.7	62.9		Microwave	display	33.7	33.7	33.7	<b>38.1</b>	
Bucket	handle	10.2	<b>78.2</b>	39.3	78.2	door		39.0	23.2	<b>42.7</b>	5.0	5.0	5.0
						handle		50.5	50.5	<b>60.4</b>	50.5	50.5	50.5
						button	12.0	12.0	12.9	<b>13.6</b>			
Camera	button	37.6	36.7	<b>39.7</b>	38.7		Mouse	button	<b>5.0</b>	3.0	5.0	5.0	
	lens	<b>35.8</b>	32.7	29.6	29.8	cord		<b>66.3</b>	66.3	66.3	66.3	66.3	
						wheel	<b>50.5</b>	50.5	50.5	50.5	50.5		
Cart	wheel	69.5	<b>74.0</b>	68.3	71.6		Oven	door	27.9	38.1	32.0	<b>42.5</b>	
						knob		66.0	69.4	<b>72.2</b>	71.7		
Chair	arm	<b>59.1</b>	50.2	58.9	50.2		Pen	cap	51.0	30.3	<b>54.6</b>	25.9	
	back	94.2	92.3	<b>94.8</b>	93.2	button		48.0	48.5	<b>52.4</b>	52.1		
	leg	72.8	79.8	72.0	<b>81.8</b>		Phone	lid	<b>28.8</b>	28.1	16.5	23.0	
	seat	90.1	<b>95.0</b>	89.1	91.1	button		34.2	34.8	35.0	<b>35.7</b>		
	wheel	94.8	95.3	<b>96.1</b>	96.1								
Clock	hand	18.7	<b>32.3</b>	18.6	18.5		Pliers	leg	3.2	4.1	<b>31.1</b>	31.1	
Coffee Machine	button	1.3	1.3	1.3	<b>1.4</b>		Printer	button	1.1	1.1	<b>1.7</b>	1.6	
	container	23.6	23.2	<b>24.2</b>	20.6		Refrige- rator	door	<b>27.2</b>	25.6	20.4	19.6	
	knob	<b>13.4</b>	12.5	12.5	11.0	handle		36.5	<b>39.3</b>	21.3	30.0		
	lid	24.8	<b>24.9</b>	24.0	22.3		Remote	button	19.9	20.3	29.6	<b>33.7</b>	
Dishwasher	door	49.0	49.0	48.8	<b>53.5</b>		Safe	door	70.7	<b>76.2</b>	65.4	71.1	
	handle	31.7	34.4	40.6	<b>47.1</b>	switch		19.3	19.3	20.2	<b>21.8</b>		
Dispenser	head	36.1	39.6	<b>49.8</b>	44.3		button	<b>1.0</b>	1.0	0.0	1.0		
	lid	76.6	81.8	85.0	<b>86.6</b>								
Display	base	96.1	94.0	<b>98.6</b>	94.1		Scissors	blade	13.2	14.8	13.0	<b>15.5</b>	
	screen	49.9	59.0	68.0	<b>70.7</b>	handle		52.9	<b>55.7</b>	47.1	54.0		
	support	<b>60.9</b>	52.0	45.6	51.9	screw	4.8	6.0	9.2	<b>9.9</b>			
Door	frame	5.1	5.5	8.5	<b>11.2</b>		Stapler	body	86.1	86.6	<b>87.0</b>	87.0	
	door	13.9	15.1	16.5	<b>20.2</b>	lid		73.4	83.2	80.0	<b>94.1</b>		
	handle	23.9	23.8	27.7	<b>28.1</b>								
Eyeglasses	body	57.2	51.8	55.7	<b>60.6</b>		Storage Furniture	door	16.2	20.6	21.8	<b>25.6</b>	
	leg	<b>82.5</b>	<b>85.2</b>	84.9	83.5	drawer		8.9	9.2	7.7	<b>9.9</b>		
Faucet	spout	<b>54.4</b>	50.8	46.4	43.6		handle	62.1	<b>71.5</b>	58.7	71.1		
	switch	<b>31.5</b>	28.9	30.5	27.0								
Folding Chair	seat	86.4	89.6	<b>100.0</b>	91.3		Suitcase	handle	74.3	72.0	<b>74.9</b>	73.1	
						wheel	39.4	33.7	30.1	<b>43.0</b>			
Globe	sphere	92.1	<b>100.0</b>	98.8	84.3		Switch	switch	21.2	<b>22.2</b>	19.7	22.1	
Kettle	lid	67.7	<b>88.9</b>	88.9	82.0		Table	door	<b>0.0</b>	0.0	0.0	0.0	
	handle	66.0	<b>74.9</b>	74.9	74.9	drawer		8.9	10.9	<b>11.6</b>	10.1		
	spout	<b>68.6</b>	68.6	57.3	66.3	leg		40.2	<b>42.4</b>	40.4	40.4		
						tabletop		61.4	63.1	67.7	<b>68.3</b>		
Keyboard	cord	78.6	80.2	<b>95.2</b>	91.6		wheel	<b>73.0</b>	54.3	61.9	66.6		
	key	34.4	31.2	47.4	<b>49.0</b>	handle	12.1	<b>14.3</b>	12.5	13.9			
KitchenPot	lid	<b>95.1</b>	95.1	69.8	76.4		Toaster	button	36.8	30.1	<b>37.1</b>	31.6	
	handle	39.9	53.8	49.6	<b>57.4</b>	slider	29.6	32.1	<b>49.0</b>	45.6			
Knife	blade	<b>44.5</b>	41.3	39.2	32.3		Toilet	lid	49.7	54.9	47.5	<b>57.2</b>	
						seat		2.1	<b>6.1</b>	1.4	5.0		
Lamp	base	<b>84.4</b>	81.5	71.6	76.5		button	56.9	<b>61.3</b>	59.2	59.5		
	body	85.2	85.2	70.4	<b>85.5</b>		TrashCan	footpedal	<b>0.0</b>	0.0	0.0	0.0	
	bulb	<b>15.8</b>	15.8	7.6	7.6	lid		32.8	34.6	<b>37.0</b>	32.0		
	shade	89.7	<b>90.2</b>	86.6	90.2	door		2.0	2.0	1.0	<b>2.9</b>		
Laptop	keyboard	54.5	67.6	62.4	<b>70.1</b>		USB	cap	17.1	15.3	<b>25.3</b>	19.6	
	screen	24.0	42.2	42.2	<b>60.0</b>	rotation		24.7	19.5	24.2	<b>33.4</b>		
	shaft	2.0	2.0	3.2	<b>3.5</b>		Washing Machine	door	35.3	27.4	<b>35.6</b>	22.9	
	touchpad	7.4	9.0	9.4	<b>10.0</b>	button		12.4	12.4	<b>17.3</b>	17.3		
	camera	<b>1.0</b>	1.0	1.0	1.0		Window	window	23.6	19.3	<b>26.1</b>	20.0	
Lighter	lid	<b>38.9</b>	38.9	17.2	29.8								
	wheel	34.9	<b>70.1</b>	33.3	56.1		<b>Mean</b>		41.6	44.7	44.2	<b>45.6</b>	
	button	28.6	<b>35.8</b>	31.3	30.1								

**Table S10:** Full table of part-agnostic instance segmentation results on the PartNet-Mobility [9] dataset.

Category	Baselines		Ablations		Ours	Category	Baselines		Ablations		Ours
	SAM3D	PartSLIP	<i>w/o</i> Weight Pred.	<i>w/o</i> SAM Integ.			SAM3D	PartSLIP	<i>w/o</i> Weight Pred.	<i>w/o</i> SAM Integ.	
Bottle	15.3	73.4	74.9	<b>79.5</b>	76.2	Microwave	1.4	15.0	12.9	<b>18.7</b>	16.2
Box	18.2	56.3	<b>72.2</b>	60.7	69.6	Mouse	2.5	12.9	11.2	20.5	<b>21.0</b>
Bucket	12.0	10.2	<b>78.2</b>	39.3	78.2	Oven	1.0	46.8	55.6	56.4	<b>62.0</b>
Camera	2.2	37.1	35.5	<b>37.8</b>	37.0	Pen	8.2	40.8	35.8	<b>47.2</b>	34.5
Cart	8.5	69.5	<b>74.0</b>	68.3	71.2	Phone	1.1	31.0	33.3	31.6	<b>34.7</b>
Chair	5.7	80.7	81.1	81.2	<b>83.3</b>	Pliers	22.2	3.2	4.1	<b>31.1</b>	31.1
Clock	1.0	18.7	<b>32.3</b>	18.6	18.5	Printer	1.0	1.1	1.1	<b>1.7</b>	1.5
Coffee Machine	3.6	<b>9.8</b>	9.5	9.4	9.4	Refrigerator	1.8	27.7	<b>30.5</b>	19.4	22.0
Dishwasher	6.7	39.0	41.9	47.1	<b>49.7</b>	Remote	1.0	19.9	20.3	29.6	<b>33.7</b>
Dispenser	20.9	59.0	59.9	<b>69.0</b>	65.5	Safe	2.0	15.5	<b>15.6</b>	14.1	15.2
Display	34.1	61.5	66.4	66.0	<b>68.4</b>	Scissors	7.4	27.0	<b>29.4</b>	24.9	26.2
Door	<b>17.2</b>	11.5	11.6	13.9	16.0	Stapler	45.1	76.0	80.4	<b>85.0</b>	82.7
Eyeglasses	27.6	71.8	72.1	72.6	<b>73.9</b>	Storage Furniture	1.2	29.7	35.9	34.8	<b>41.7</b>
Faucet	23.7	<b>42.8</b>	37.1	37.9	35.6	Suitcase	1.8	<b>47.3</b>	46.2	47.1	43.8
Folding Chair	27.2	86.4	89.6	<b>100.0</b>	91.3	Switch	5.0	21.2	22.2	19.7	<b>22.4</b>
Globe	4.2	92.1	<b>100.0</b>	98.8	84.3	Table	10.8	28.7	<b>29.2</b>	27.6	28.2
Kettle	41.8	69.5	75.3	<b>76.3</b>	75.9	Toaster	1.8	34.8	35.9	38.6	<b>39.2</b>
Keyboard	1.0	34.1	31.0	47.3	<b>48.9</b>	Toilet	3.6	35.4	40.0	36.6	<b>41.0</b>
KitchenPot	31.9	55.4	<b>66.2</b>	52.3	65.4	TrashCan	7.6	22.6	<b>24.5</b>	24.3	22.8
Knife	2.2	<b>44.5</b>	41.3	39.2	32.1	USB	23.4	20.5	15.2	<b>27.3</b>	23.5
Lamp	48.4	77.0	<b>77.7</b>	66.5	75.9	Washing Machine	2.8	17.5	14.6	<b>21.4</b>	19.6
Laptop	1.9	19.5	26.5	27.4	<b>34.1</b>	Window	18.0	23.6	19.3	<b>26.1</b>	19.5
Lighter	8.9	32.7	<b>49.9</b>	26.6	39.8						
						<b>Mean</b>	12.1	38.9	42.6	42.6	<b>44.1</b>

### S.12 Additional Qualitative Results for Semantic Segmentation

We present additional part segmentation results below, encompassing semantic segmentation. Each row shows the same object with different views.

view 1						view 2					
Input	GT	SATR	SATR +SP	PartSLIP	Ours	Input	GT	SATR	SATR +SP	PartSLIP	Ours

view 1						view 2					
Input	GT	SATR	SATR +SP	PartSLIP	Ours	Input	GT	SATR	SATR +SP	PartSLIP	Ours

view 1						view 2					
Input	GT	SATR	SATR +SP	PartSLIP	Ours	Input	GT	SATR	SATR +SP	PartSLIP	Ours

view 1						view 2					
Input	GT	SATR	SATR +SP	PartSLIP	Ours	Input	GT	SATR	SATR +SP	PartSLIP	Ours

view 1						view 2					
Input	GT	SATR	SATR +SP	PartSLIP	Ours	Input	GT	SATR	SATR +SP	PartSLIP	Ours

view 1						view 2					
Input	GT	SATR	SATR +SP	PartSLIP	Ours	Input	GT	SATR	SATR +SP	PartSLIP	Ours

view 1						view 2					
Input	GT	SATR	SATR +SP	PartSLIP	Ours	Input	GT	SATR	SATR +SP	PartSLIP	Ours

### S.13 Additional Qualitative Results for Instance Segmentation

We present additional part segmentation results below, encompassing instance segmentation. Each row shows the same object with different views.

Input	view 1				Input	view 2			
	GT	SAM3D	PartSLIP	Ours		GT	SAM3D	PartSLIP	Ours

Input	view 1				Input	view 2			
	GT	SAM3D	PartSLIP	Ours		GT	SAM3D	PartSLIP	Ours

Input	view 1				Input	view 2			
	GT	SAM3D	PartSLIP	Ours		GT	SAM3D	PartSLIP	Ours

Input	view 1				Input	view 2			
	GT	SAM3D	PartSLIP	Ours		GT	SAM3D	PartSLIP	Ours
									
									
									
									

## References

1. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: LoRA: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021)
2. Ke, L., Ye, M., Danelljan, M., Liu, Y., Tai, Y.W., Tang, C.K., Yu, F.: Segment anything in high quality. In: NeurIPS (2023)
3. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. arXiv preprint arXiv:2304.02643 (2023)
4. Li, L.H., Zhang, P., Zhang, H., Yang, J., Li, C., Zhong, Y., Wang, L., Yuan, L., Zhang, L., Hwang, J.N., et al.: Grounded language-image pre-training. In: CVPR (2022)
5. Liu, M., Zhu, Y., Cai, H., Han, S., Ling, Z., Porikli, F., Su, H.: PartSLIP: Low-shot part segmentation for 3d point clouds via pretrained image-language models. In: CVPR (2023)
6. Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al.: Grounding DINO: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499 (2023)
7. Ren, T., Liu, S., Zeng, A., Lin, J., Li, K., Cao, H., Chen, J., Huang, X., Chen, Y., Yan, F., Zeng, Z., Zhang, H., Li, F., Yang, J., Li, H., Jiang, Q., Zhang, L.: Grounded SAM: Assembling open-world models for diverse visual tasks. arXiv preprint arXiv:2401.14159 (2024)
8. Wu, T., Zhang, J., Fu, X., Wang, Y., Ren, J., Pan, L., Wu, W., Yang, L., Wang, J., Qian, C., Lin, D., Liu, Z.: OmniObject3D: Large-vocabulary 3D object dataset for realistic perception, reconstruction and generation. In: CVPR (2023)
9. Xiang, F., Qin, Y., Mo, K., Xia, Y., Zhu, H., Liu, F., Liu, M., Jiang, H., Yuan, Y., Wang, H., Yi, L., Chang, A.X., Guibas, L.J., Su, H.: SAPIEN: A simulated part-based interactive environment. In: CVPR (2020)
10. Yi, K.M., Trulls, E., Ono, Y., Lepetit, V., Salzmann, M., Fua, P.: Learning to find good correspondences. In: CVPR (2018)
11. Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: ICCV (2023)