

A APPENDIX

A.1 IMPLEMENTATION DETAILS

This section supplements the pipeline description in Section 5 with additional implementation details that could not be included earlier due to space constraints. *Italicized* text denotes tunable hyperparameters, with sensitivity analyses provided in Sections A.2 and A.4.

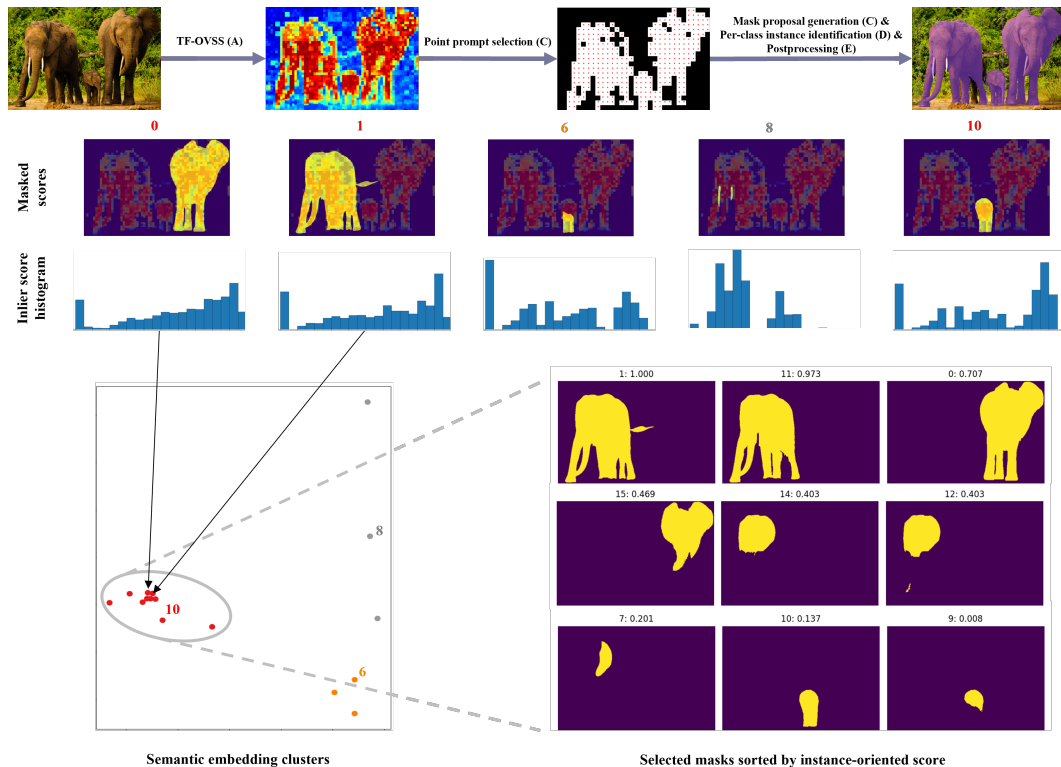


Figure 7: Sample outputs from each step of TIDES. **Top:** Overall pipeline — input image, patch-text semantic alignment scores, foreground patches & selected point prompts for mask generation, and final output. **Middle:** Sample predictions after PSM postprocessing with relaxed criteria and the corresponding histograms. **Bottom:** Semantic embedding clusters (colors denote different clusters). The masks on the right are from the largest cluster, sorted by instance-oriented score. Final postprocessing filters these masks, yielding the outputs shown in the top row.

Semantic feature extraction (A). Following standard dual encoder usage, alignment scores between text and image embeddings are computed via cosine similarity, yielding values in $[-1, 1]$. The input image is first resized to have a longer side of 512, then we adopt a sliding-window approach with multiple *window scales* and *sizes*, using a stride equal to half the *window size*. Patch-text scores from all windows are upsampled to the original resolution, averaged across scales to form pixel-text scores, and converted back to patch-text scores using the greatest common divisor of patch sizes. This allows TIDES to use a patch size different from that of the dual encoder and improves its ability to capture fine-grained detail.

In-context class filtering (B). Sliding-window inference yields multiple $[\text{CLS}]$ tokens per image. For each window, we compute the similarity between the $[\text{CLS}]$ token and target class tokens, softmaxed against the *background* token. A class is marked present if any window produces a score greater than that of *background*.

Point prompt selection and mask proposal generation (C). Foreground patches are identified by taking the patch-wise product of per-class min-max normalized scores and softmax-normalized scores

against background, followed by thresholding at the *foreground region threshold*. All foreground patches across classes are then collected as a batch and passed to the PSM with a single point prompt at the center of each patch. Using every raw PSM prediction significantly increases computation without improving results, so we relax PSM’s postprocessing by lowering the *stability score threshold* and raising the *box NMS threshold* to retain diverse masks at lower latency (The exact trade-offs are discussed in Section A.4). Among the generated masks, those whose IoU with the corresponding foreground region falls below the *foreground-mask IoU threshold* are discarded.

Per-class instance identification (D). For each mask, patch–text alignment scores from its foreground patches are collected, binned into a histogram (*bin size*), and normalized by count to produce fixed-length semantic embeddings. KDE-based clustering is then applied using *bandwidth* and *density threshold* as hyperparameters. Only masks from the largest cluster are kept, filtering out partial or incorrect predictions. The local density of each selected mask is used as its IO score.

Postprocessing (E). Overlapping masks are removed by applying NMS with the IO score and an IoU threshold (*NMS threshold*). Each mask is then assigned a semantic score by averaging its min–max normalized patch–text alignment scores, and those below the *semantic score filtering threshold* are discarded to suppress noise.

A.2 HYPERPARAMETER SENSITIVITY ANALYSIS

Table 4 reports the range of values explored for each hyperparameter along with the corresponding sensitivity of the pipeline, summarized using mean, median, minimum, and maximum AP. The **bolded** value indicates the setting that achieved the highest performance for each hyperparameter.

The search and analysis are conducted using a pipeline built with CS (ViT-B/16) (Li et al., 2025) and EfficientViTSAM (XL1) (Sun et al., 2024) on a subset of the MS COCO 2017 validation set (Lin et al., 2014), selected to include at least three instances of each of the 80 categories. Empirically, the best-performing configurations on this subset are reasonably consistent with those observed on the full validation set of 5,000 samples.

Table 4: Impact of hyperparameter choices on performance. [a:b:c] denotes a range from a to c (inclusive) with step size b. The **bolded** choices indicates the one that led to the best performance.

Hyperparameter	Stage	Choices (selected)	AP			
			Minimum	Median	Maximum	Mean
<i>window scales</i>	A	1, 2, 3				
<i>window sizes</i>	A	36, 112 , 168, 224, 336, 448	14.8	25.9	30.2	25.3
<i>PSM box IOU threshold</i>	C	Section A.4 (0.80)	26.6	28.2	30.1	28.0
<i>PSM stability score threshold</i>	C	Section A.4 (0.90)				
<i>foreground region threshold</i>	C	[0.10:0.10:0.90] (0.10)	22.8	30.2	30.4	29.3
<i>foreground-mask IoU threshold</i>	C	[0.10:0.10:0.90] (0.20)	25.3	29.4	30.7	29.0
<i>histogram bin size</i>	D	0.025 , 0.050, 0.100	10.7	21.9	31.3	21.6
<i>clustering bandwidth</i>	D	[0.50:0.20:1.50] (0.50)	10.6	22.1	30.3	21.9
<i>clustering density threshold</i>	D	0.75, 0.80, 0.85, 0.90	11.6	22.4	31.6	22.1
<i>NMS threshold</i>	E	[0.10:0.10:0.90] (0.70)	27.5	29.6	30.4	29.5
<i>semantic score filtering threshold</i>	E	[0.10:0.10:0.90] (0.10)	1.6	29.1	30.3	23.3

A.3 ADDITIONAL STUDIES ON PSM BEHAVIOR AND THE IMPACT OF IO SCORING

Per-Prompt Prediction Distributions of Different PSMs. Figure 8 presents the same subplots as Figure 3—*All Predictions*, *Best IoU*, and *Best Score*—but includes two additional PSMs: MobileSAM (Tiny) (Zhang & Jiao, 2023) and SAM2.1 (Large) (Ravi et al., 2024). The distributions from these additional models are consistent with the original three, confirming that the findings in Section 3.2 generalize across other PSMs.

Out of the 80 COCO classes, the following 18 are considered *non-uniform*: person, bicycle, car, motorcycle, bus, airplane, train, truck, boat, hot dog, chair, couch, potted plant, bed, dining table, bench, skis, and book.

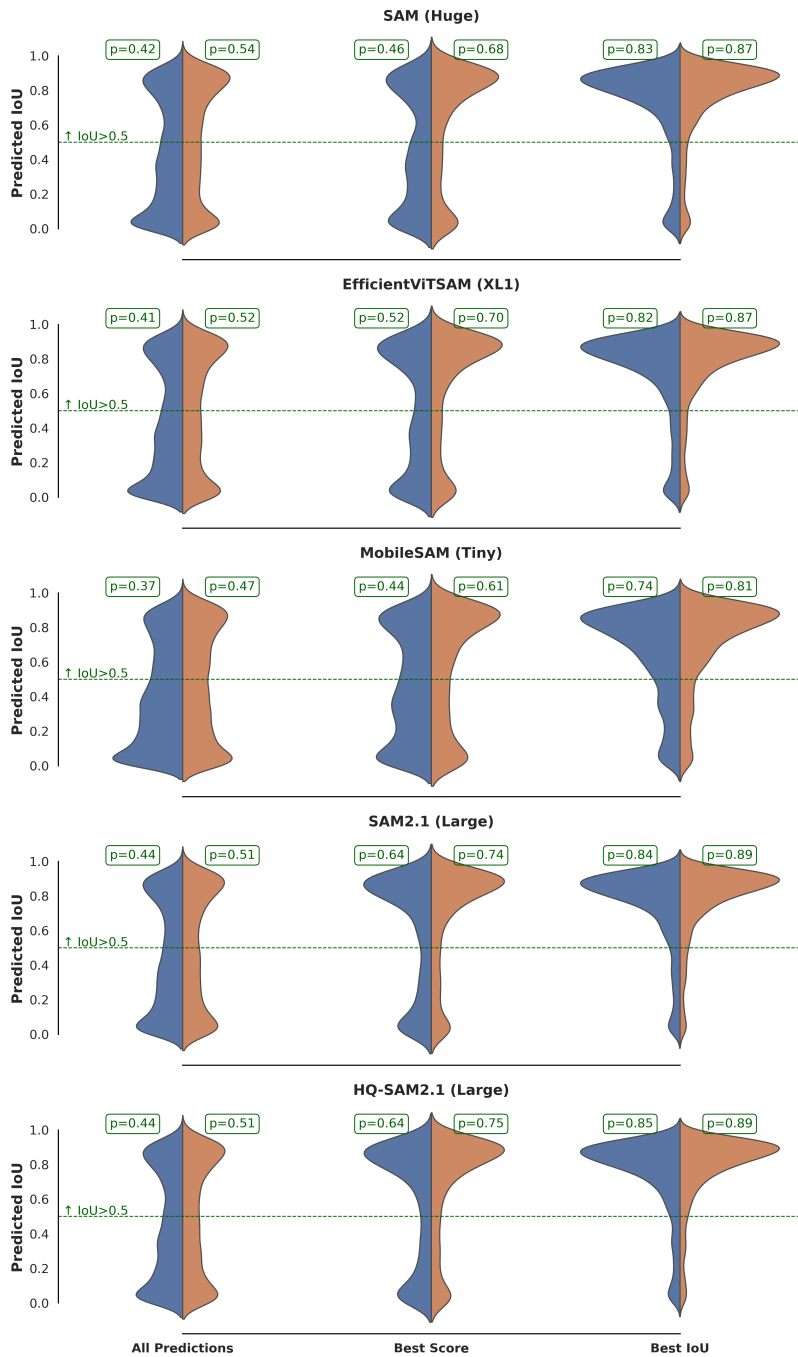


Figure 8: Per-point prediction distributions of different PSMs across IoU values for *uniform* and *non-uniform* objects (p indicates percentage of predictions with IoU > 0.5).

Effect of Different Scoring Configurations on Per-Instance Distributions. To illustrate the effect of each assumption described in Section 4, we presents *Best Score* distributions under different point sampling and scoring configurations using the five PSMs.

1. *Random & Best Raw Score*: n random point prompts within GT masks & original scores.
2. *Random & Best Per-Instance IO Score*: n random point prompts within GT masks & IO scores from instance-level clustering ($n \times m$ embeddings).

3. *TF-OVSS & Best Per-Instance IO Score*: n high-score prompts selected from foreground regions identified with CS (Li et al., 2025) & IO scores from instance-level clustering ($n \times m$ embeddings).
4. *TF-OVSS & Best Per-Image IO Score*: n high-score prompts selected from foreground regions identified with CS (Li et al., 2025) & IO scores from image-level clustering ($n \times m$ embeddings).

Please note that, since the primary focus of this study is to demonstrate the implicit bias of each PSM and the impact of the new scoring method, we consider only a subset of the MS COCO 2017 validation set (Lin et al., 2014), consisting of 500 samples. Among these, we include only the instances for which CS (Li et al., 2025) successfully identified $n = 5$ point prompts. In other words, the PSM tuned for TIDES pipeline may exhibit a slightly different distribution compared to this optimal setting, due to the inclusion of PSM’s postprocessing logic under a relaxed configuration (Section A.4).

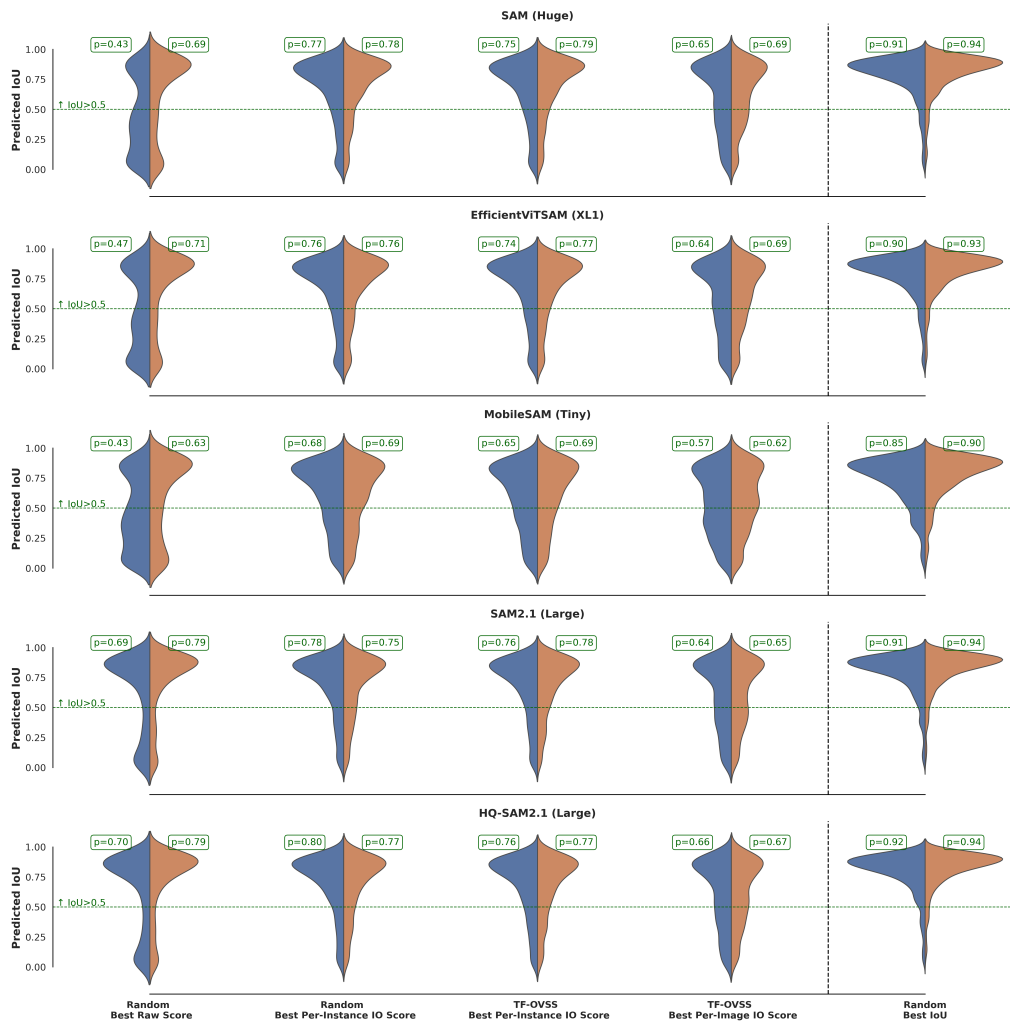


Figure 9: Effect of different point sampling and scoring configurations on the distribution of *Best Score* predictions across PSMs for *uniform* and *non-uniform* objects (p indicates percentage of predictions with IoU > 0.5).

Figure 9 reveals that the changes in distributions resulting from different sampling and scoring configurations exhibit similar patterns across PSMs. Therefore, we describe the identified pattern based on the SAM (Huge) (Kirillov et al., 2023) graph.

From *Random & Best Raw Score* to *Random & Best Per-Instance IO Score*, we observe a 9% increase for *uniform* objects and a substantial 33% increase for *non-uniform* objects, resulting in 78% and 77% of predictions above $\text{IoU} > 0.5$ threshold, respectively. Results from *TF-OVSS & Best Per-Instance IO Score* remain comparable, suggesting that point prompts derived from high patch-text alignment regions do not introduce semantic bias.

When predictions are aggregated across multiple instances (*TF-OVSS & Best Per-Image IO Score*), we observe a 10% decrease for both *uniform* and *non-uniform* objects. Nonetheless, the resulting distributions remain skewed toward higher IoU values (69% for *uniform* and 65% for *non-uniform*), indicating that the misalignment observed in *Random & Best Raw Score* for *non-uniform* objects has been substantially mitigated.

Figure 9 also includes the distributions of *Random & Best IoU* derived from random point prompts with raw scores. The high values of 94% (*uniform*) and 91% (*non-uniform*) represent empirical upper bounds, highlighting the potential for further gains through improved scoring.

A.4 INFERENCE TIME ANALYSIS

In this section, we analyze TIDES’ behavior with respect to PSM postprocessing configurations, aiming to identify the setup that yields the best TIDES performance in practical scenarios.

- *PSM default*: the default postprocessing configuration.
 - *PSM stability score threshold* = 0.95, *PSM box NMS threshold* = 0.70
- *LooseNMS*: only the NMS filtering configuration relaxed.
 - *PSM stability score threshold* = 0.95, *PSM box NMS threshold* = 0.90
- *Balanced*: both the NMS and stability score filtering configurations slightly relaxed.
 - *PSM stability score threshold* = 0.80, *PSM box NMS threshold* = 0.90
- *MinFilter*: both the NMS and stability score filtering configurations relaxed to the extreme.
 - *PSM stability score threshold* = 0.60, *PSM box NMS threshold* = 1.00

Inference latency and corresponding instance segmentation performance are measured by running the pipeline on the same subsampled images used for hyperparameter tuning in Section A.2, but with only the ground-truth classes as the vocabulary. All experiments are performed on a machine equipped with an Intel Core i9-9940X CPU @ 3.30GHz, 125 GB RAM, and a single NVIDIA Quadro RTX 6000 GPU with 24 GB of memory. We report the mean and 95% confidence interval in seconds for the latency, and the AP value for the segmentation performance.

Table 5: Inference latency (in seconds) and segmentation performance (AP) under different PSM configurations. The *Balanced* setting shows a notable performance gain with a moderate increase in latency, and is therefore selected as the default for TIDES.

PSM	Params.	Inference latency (s) ↓ / segmentation performance (AP) ↑			
		<i>PSM default</i>	<i>LooseNMS</i>	<i>Balanced</i>	<i>MinFilter</i>
SAM (Base)	93.7M	2.05 ± 0.31 / 21.1	2.03 ± 0.33 / 21.4	3.51 ± 0.63 / 22.6	5.59 ± 0.99 / 24.7
SAM (Huge)	636.2M	3.01 ± 0.40 / 26.9	3.15 ± 0.42 / 26.8	4.48 ± 0.68 / 26.7	6.20 ± 1.04 / 27.3
E-SAM (L0)	34.8M	1.21 ± 0.06 / 10.5	1.22 ± 0.07 / 10.0	1.28 ± 0.08 / 13.4	1.30 ± 0.08 / 13.5
E-SAM (XL1)	203.3M	1.94 ± 0.14 / 26.7	2.03 ± 0.15 / 26.6	2.72 ± 0.25 / 30.1	3.54 ± 0.37 / 29.2
M-SAM (Tiny)	10.1M	1.98 ± 0.23 / 21.9	2.07 ± 0.26 / 22.1	3.54 ± 0.55 / 25.0	4.96 ± 0.92 / 24.5
SAM2.1 (Base+)	80.8M	1.09 ± 0.05 / 4.8	1.10 ± 0.04 / 5.2	1.11 ± 0.04 / 7.9	1.12 ± 0.04 / 8.8
SAM2.1 (Large)	224.4M	1.25 ± 0.05 / 10.6	1.23 ± 0.05 / 10.5	1.36 ± 0.06 / 20.6	1.55 ± 0.08 / 19.7
HQ-SAM2.1 (Large)	224.7M	1.74 ± 0.06 / 17.8	1.74 ± 0.06 / 17.8	1.79 ± 0.07 / 20.4	1.97 ± 0.09 / 20.5
Average	-	1.78 / 17.5	1.82 / 17.6	2.47 / 20.8	3.28 / 21.0

As shown in Table 5, TIDES inference time closely correlates with both the size and architectural design of the underlying PSM. Models optimized for speed, such as E-SAM (L0) and M-SAM (Tiny), achieve lower latency with fewer parameters compared to their respective baselines, SAM (Base) and SAM (Huge). Despite having more parameters, SAM2.1 also exhibits improved latency over SAM, owing to its redesigned decoder, optimized prompt encoding, and reduced redundant computation.

As expected, inference latency increases as the filtering conditions are relaxed, due to the higher number of masks processed. Segmentation performance also follows the same pattern, increasing as the filtering conditions are relaxed. However, the performance gain from *Balanced* to *MinFilter* is minimal, showing that a *PSM stability score threshold* of 0.80 and a *PSM box NMS threshold* of 0.90 provide the most return in performance on the investment in inference latency. Therefore, we select this setting as the default for our TIDES pipeline.

A.5 RESOURCE CONSUMPTION ANALYSIS

To better understand which factors drive resource consumption, we conduct two studies on computational efficiency and memory usage. For all experiments, we evaluate TIDES configured with CS (ViT-B/16) and SAM (Base) on a subset of the LVIS dataset (Gupta et al., 2019). We report per-sample peak GPU memory usage and latency as mean \pm standard deviation.

Sensitivity to the Number of Classes. Table 6 summarizes how peak GPU memory and latency scale with the number of classes. Both metrics grow steadily as the number of classes increases, showing that TIDES is sensitive to this factor.

Table 6: Sensitivity of TIDES to the number of classes.

# of Classes	1	25	50	75	100
Peak GPU Memory (GB)	3.3 ± 1.9	4.1 ± 0.7	5.5 ± 0.9	7.5 ± 1.9	8.7 ± 2.2
Latency (sec)	3.0 ± 0.1	43.8 ± 2.1	85.8 ± 4.1	128.7 ± 4.6	166.6 ± 12.2

Sensitivity to the Number of Instances. Table 7 shows how resource usage changes with respect to the average number of instances per image. Here we use a 50-class subset of LVIS. While peak GPU memory grows slightly, latency remains largely unaffected.

Table 7: Sensitivity of TIDES to the number of instances per image.

Avg. Instances/Image	2	4	6	8
Peak GPU Memory (GB)	5.5 ± 0.7	5.3 ± 0.5	5.9 ± 0.6	6.3 ± 0.9
Latency (sec)	84.5 ± 4.2	80.2 ± 4.1	84.9 ± 3.5	87.3 ± 8.4

Improving scalability with respect to both class count and instance density remains an important area to work on for making TIDES more practical and broadly deployable.

A.6 ADDITIONAL QUANTITATIVE ANALYSIS

To further back up our findings in Section 6, we evaluate the performance of TIDES designed with MobileSAM (Zhang & Jiao, 2023) and SAM2 (Ravi et al., 2024). In addition, we evaluate TIDES linked with smaller version of SAM (Kirillov et al., 2023) (Base) and SAM2.1 (Base+). Our complete evaluation, summarized in Table 8, further confirms that our IO scoring consistently improves performance over raw scores. The evaluation setting is identical to that described in Section 6.

Table 8: Zero-shot TF-OVIS performance of TIDES with different TF-OVSS and PSM combinations on MS COCO 2017. Each includes two variants: Raw and IO scores. ‘‘SS perf.’’ denotes the average TF-OVSS performance on standard benchmarks. ‘‘SOTA’’ refers to Zip (Shi & Yang, 2024).

TF-OVSS (SS Perf.)	PSM	Score	AR ₁₀	AR ₁₀₀	AP _s	AP _m	AP _l	AP ₇₅	AP ₅₀	AP _↑	vs. Naive	vs. SOTA (11.8)
CS (37.2)	SAM (Base)	Raw	34.9	36.1	4.6	22.4	34.4	12.3	21.7	13.0		
		IO	34.1	35.0	4.6	21.5	34.1	14.2	24.8	15.2	+2.2	+3.4
	SAM (Huge)	Raw	40.0	41.7	5.7	25.9	37.4	14.1	22.9	14.6		
		IO	39.8	40.9	6.7	27.6	39.2	18.7	29.4	19.2	+4.6	+7.4
	E-SAM (L0)	Raw	35.6	36.3	5.3	24.1	36.5	14.6	25.2	15.6		
		IO	35.5	36.1	5.4	24.4	37.5	17.5	29.3	18.5	+2.9	+6.7
	E-SAM (XL1)	Raw	39.4	40.3	6.4	27.9	40.3	17.8	27.3	18.2		
		IO	38.5	39.3	6.9	27.6	40.8	20.5	31.2	21.0	+2.8	+9.2
	M-SAM (Tiny)	Raw	34.9	35.9	5.2	20.2	31.4	10.8	21.8	12.4		
		IO	34.3	35.0	5.8	21.9	32.8	13.6	26.5	15.5	+3.1	+3.7
	SAM2.1 (Base+)	Raw	5.8	5.8	0.0	0.6	13.3	5.5	6.5	5.4		
		IO	5.8	5.8	0.0	0.6	13.2	5.4	6.5	5.4	0.0	-6.4
	SAM2.1 (Large)	Raw	19.7	19.7	1.5	13.2	30.9	14.0	20.5	14.1		
		IO	20.2	20.2	1.5	13.9	32.1	14.5	21.2	14.6	+0.5	+2.8
HQ-SAM2.1 (Large)	Raw	25.6	25.7	3.4	15.4	31.8	12.3	25.4	14.6			
	IO	25.9	25.9	3.5	15.8	31.8	12.5	26.3	15.0	+0.4	+3.2	
SCLIP (39.1)	SAM (Base)	Raw	28.7	30.9	1.5	12.9	34.5	8.9	16.1	9.6		
		IO	29.9	32.0	1.7	12.5	33.3	10.8	20.8	12.4	+2.8	+0.6
	SAM (Huge)	Raw	36.4	39.2	2.6	20.9	40.8	13.2	20.9	13.7		
		IO	35.8	37.9	3.0	21.6	41.5	17.6	26.7	18.2	+4.5	+6.4
	E-SAM (L0)	Raw	29.5	30.6	2.0	15.1	37.3	12.0	21.0	13.0		
		IO	29.4	30.2	2.3	15.3	36.1	15.2	25.2	16.3	+3.3	+4.5
	E-SAM (XL1)	Raw	35.2	36.8	3.1	22.7	40.9	15.2	23.4	15.7		
		IO	34.3	35.5	3.4	23.2	41.8	18.9	28.1	19.4	+3.7	+7.6
	M-SAM (Tiny)	Raw	31.8	33.4	2.2	15.9	33.9	9.8	19.8	11.5		
		IO	31.1	32.4	2.6	16.8	34.7	12.9	24.6	14.7	+3.2	+2.9
	SAM2.1 (Base+)	Raw	6.4	6.4	0.0	0.7	13.9	5.8	7.1	5.9		
		IO	6.4	6.4	0.0	0.7	13.9	5.9	7.1	5.9	0.0	-5.9
	SAM2.1 (Large)	Raw	19.0	19.0	0.8	10.2	32.1	13.6	20.0	13.6		
		IO	18.8	18.8	0.8	10.3	32.0	13.7	20.1	13.9	+0.3	+2.1
HQ-SAM2.1 (Large)	Raw	22.2	22.2	1.6	12.1	30.4	11.8	22.4	13.3			
	IO	21.8	21.8	1.7	12.2	30.5	12.4	23.1	13.9	+0.6	+2.1	
SC-CLIP (43.9)	SAM (Base)	Raw	33.2	35.1	3.3	20.2	35.7	11.7	21.2	12.6		
		IO	32.0	33.3	3.0	18.9	35.0	13.9	24.2	15.0	+2.4	+3.2
	SAM (Huge)	Raw	40.0	42.5	4.3	25.1	40.4	14.6	23.7	15.3		
		IO	38.2	40.0	4.7	25.1	39.6	18.5	28.8	19.1	+3.8	+7.3
	E-SAM (L0)	Raw	34.0	35.0	4.0	22.1	38.4	14.8	25.2	15.7		
		IO	31.7	32.4	4.0	20.4	36.0	16.1	27.0	17.2	+1.5	+5.4
	E-SAM (XL1)	Raw	38.6	40.1	5.1	26.8	40.6	16.8	26.0	17.4		
		IO	37.1	38.2	5.5	26.5	40.0	20.1	30.2	20.6	+3.2	+8.8
	M-SAM (Tiny)	Raw	34.9	36.4	4.0	19.5	33.7	11.2	22.7	13.0		
		IO	33.4	34.5	4.2	19.8	33.2	13.6	26.5	15.6	+2.6	+3.8
	SAM2.1 (Base+)	Raw	5.8	5.8	0.0	0.6	13.3	5.5	6.5	5.4		
		IO	5.8	5.8	0.0	0.6	13.2	5.4	6.5	5.4	0.0	-6.4
	SAM2.1 (Large)	Raw	19.5	19.5	1.1	12.6	31.7	14.3	20.5	14.2		
		IO	19.4	19.4	1.1	12.5	32.1	14.5	20.7	14.4	+0.2	+2.6
HQ-SAM2.1 (Large)	Raw	23.3	23.3	2.1	14.4	30.8	12.2	23.5	14.0			
	IO	23.0	23.0	2.2	14.9	31.5	13.2	24.5	14.8	+0.8	+3.0	

1026 A.7 ADDITIAONL QUALITATIVE ANALYSIS

1027 A.7.1 TIDES WITH OPTIMAL SUBCOMPONENT CONFIGURATIONS

1028 To complement the ablation study in Section 7, we present qualitative results of TIDES (CS (ViT-B/16)
1029 & E-SAM (XL1)) configured with the optimal alternative for each subcomponent (Figures 10, 11,
1030 and 12). The figures illustrate how improvements to individual subcomponents contribute to overall
1031 performance gains, while also providing insights into common failure cases:
1032
1033

- 1034 1. Instance label incompleteness in ground truth (false negatives):
1035 Not all instances in the COCO images are labeled, leading to potential underestimation of
1036 TIDES’ performance in quantitative metrics. TIDES often detects more instances than are
1037 annotated, particularly in cluttered scenes (e.g., samples 181666, 238866). While these
1038 predictions are semantically valid, they reduce measured precision. For example, a shelf with
1039 multiple items labeled as one object may be segmented by TIDES into multiple instances.
1040
- 1041 2. Non-uniform objects and over-segmentation:
1042 Objects with distinct subcomponents (e.g., sofas with separate cushions or a person with
1043 accessories) may be over-segmented due to their multi-component appearance. For instance,
1044 in sample 219578, each sofa cushion is predicted as a separate instance.
1045
- 1046 3. Class filtering and confusion:
1047 In the absence of ground-truth semantic supervision, in-context class filtering can produce
1048 false positives or misclassifications. For example, in sample 184791, an unlabeled large jug
1049 is detected as a bowl.
1050
- 1051 4. Small object detection:
1052 While not a primary bottleneck, TIDES shows lower accuracy on small objects, consistent
1053 with trends observed in other segmentation models (Table 1).

1054 Throughout our experiments with Zip, we observed that its performance on multi-instance or complex
1055 scenes falls significantly below a usable level. High false positives caused by confusion between
1056 foreground and background (e.g., samples 238866, 181666) indicate that the main issues stem from
1057 an inability to distinguish foreground from background and to correctly identify in-context classes.
1058 This highlights the stability and advancements introduced by TIDES.

1059 It is also worth noting that these failures arise from fundamental limitations in Zip’s design. Zip
1060 focuses on detecting boundary pixels to group non-boundary pixels into instances, which results in
1061 aggressive discarding of pixels near instance boundaries. Recovery of these pixels relies on SAM,
1062 introducing additional instability in complex scenarios.

1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Sample ID	181666	289393	463730
1134			
1135			
1136			
1137			
1138			
1139	GT		
1140			
1141			
1142			
1143	Zip		
1144			
1145	GT vocab. ✓		
1146	GT alignments ✗		
1147	GT mask ✗		
1148			
1149			
1150	TIDES		
1151			
1152	GT vocab. ✗		
1153	GT alignments ✗		
1154	GT mask ✗		
1155			
1156			
1157	TIDES		
1158			
1159	GT vocab. ✓		
1160	GT alignments ✗		
1161	GT mask ✗		
1162			
1163			
1164	TIDES		
1165			
1166	GT vocab. ✗		
1167	GT alignments ✓		
1168	GT mask ✗		
1169			
1170			
1171	TIDES		
1172			
1173	GT vocab. ✗		
1174	GT alignments ✗		
1175	GT mask ✓		
1176			
1177			
1178	TIDES		
1179			
1180	GT vocab. ✓		
1181	GT alignments ✓		
1182	GT mask ✓		
1183			
1184			
1185			

Figure 11: Qualitative results of TIDES, its variants with subcomponents replaced by optimal alternatives, and Zip (Shi & Yang, 2024).

Sample ID	475779	511321	219578
1188			
1189			
1190			
1191			
1192	GT		
1193			
1194			
1195			
1196			
1197	Zip		
1198	GT vocab. ✓		
1199	GT alignments ✗		
1200	GT mask ✗		
1201			
1202			
1203			
1204	TIDES		
1205	GT vocab. ✗		
1206	GT alignments ✗		
1207	GT mask ✗		
1208			
1209			
1210			
1211	TIDES		
1212	GT vocab. ✓		
1213	GT alignments ✗		
1214	GT mask ✗		
1215			
1216			
1217			
1218	TIDES		
1219	GT vocab. ✗		
1220	GT alignments ✓		
1221	GT mask ✗		
1222			
1223			
1224			
1225	TIDES		
1226	GT vocab. ✗		
1227	GT alignments ✗		
1228	GT mask ✓		
1229			
1230			
1231			
1232	TIDES		
1233	GT vocab. ✓		
1234	GT alignments ✓		
1235	GT mask ✓		
1236			
1237			
1238			
1239			

Figure 12: Qualitative results of TIDES, its variants with subcomponents replaced by optimal alternatives, and Zip (Shi & Yang, 2024).

1242 A.7.2 EFFECT OF PSM PERFORMANCE ON TIDES

1243
1244 The first set of figures (Figure 13, 14, and 15) illustrates how different PSMs affect the final segmen-
1245 tation outputs of TIDES, underscoring the system’s reliance on PSM for robust performance. To
1246 isolate the impact of the PSM, we use GT class names as prompts and apply the same TF-OVSS
1247 method (CS (ViT-B/16)) to derive the patch-text alignment scores.

1248 Based on our comprehensive quantitative and qualitative analysis, EfficientViT-SAM is identified
1249 as the most suitable PSM for TIDES. Given HQ-SAM2.1’s ability to capture fine-grained details,
1250 this difference needs careful examination. Therefore, we summarize the top three masks produced
1251 by each PSM in Figure 16, 17, 18, and 19. Our observations reveal that each PSM has distinct
1252 characteristics that affect the types of masks retained after its postprocessing logic and TIDES’
1253 instance identification. For example, SAM2.1 tends to generate many masks that span multiple
1254 objects. This behavior negatively impacts TIDES’ clustering process, often causing incorrect masks
1255 to be ranked higher than the correct ones. Developing a pipeline that explicitly accounts for the
1256 implicit characteristics of each PSM presents a promising direction for future work.

1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Sample ID	184791	238866	286994
1296			
1297			
1298			
1299			
1300	GT		
1301			
1302			
1303			
1304			
1305	TIDES (SAM Huge)		
1306	GT vocab. ✓		
1307	GT alignments ✗		
1308	GT mask ✗		
1309			
1310			
1311			
1312	TIDES (E-SAM XL1)		
1313	GT vocab. ✓		
1314	GT alignments ✗		
1315	GT mask ✗		
1316			
1317			
1318			
1319	TIDES (M-SAM Tiny)		
1320	GT vocab. ✓		
1321	GT alignments ✗		
1322	GT mask ✗		
1323			
1324			
1325			
1326	TIDES (SAM2 Large)		
1327	GT vocab. ✓		
1328	GT alignments ✗		
1329	GT mask ✗		
1330			
1331			
1332			
1333	TIDES (HQ-SAM2 Large)		
1334	GT vocab. ✓		
1335	GT alignments ✗		
1336	GT mask ✗		
1337			
1338			
1339			
1340	Zip (SAM Base)		
1341	GT vocab. ✓		
1342	GT alignments ✗		
1343	GT mask ✗		
1344			
1345			
1346			
1347			

Figure 13: Difference in TIDES performance across underlying PSMs, along with Zip (Shi & Yang, 2024).

Sample ID	181666	289393	463730
1350			
1351			
1352			
1353			
1354	GT		
1355			
1356			
1357			
1358			
1359	TIDES (SAM Huge)		
1360			
1361	GT vocab. ✓		
1362	GT alignments ✗		
1363	GT mask ✗		
1364			
1365			
1366	TIDES (E-SAM XL1)		
1367			
1368	GT vocab. ✓		
1369	GT alignments ✗		
1370	GT mask ✗		
1371			
1372			
1373	TIDES (M-SAM Tiny)		
1374			
1375	GT vocab. ✓		
1376	GT alignments ✗		
1377	GT mask ✗		
1378			
1379			
1380	TIDES (SAM2 Large)		
1381			
1382	GT vocab. ✓		
1383	GT alignments ✗		
1384	GT mask ✗		
1385			
1386			
1387	TIDES (HQ-SAM2 Large)		
1388			
1389	GT vocab. ✓		
1390	GT alignments ✗		
1391	GT mask ✗		
1392			
1393			
1394	Zip (SAM Base)		
1395			
1396	GT vocab. ✓		
1397	GT alignments ✗		
1398	GT mask ✗		
1399			
1400			
1401			

Figure 14: Difference in TIDES performance across underlying PSMs, along with Zip (Shi & Yang, 2024).

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

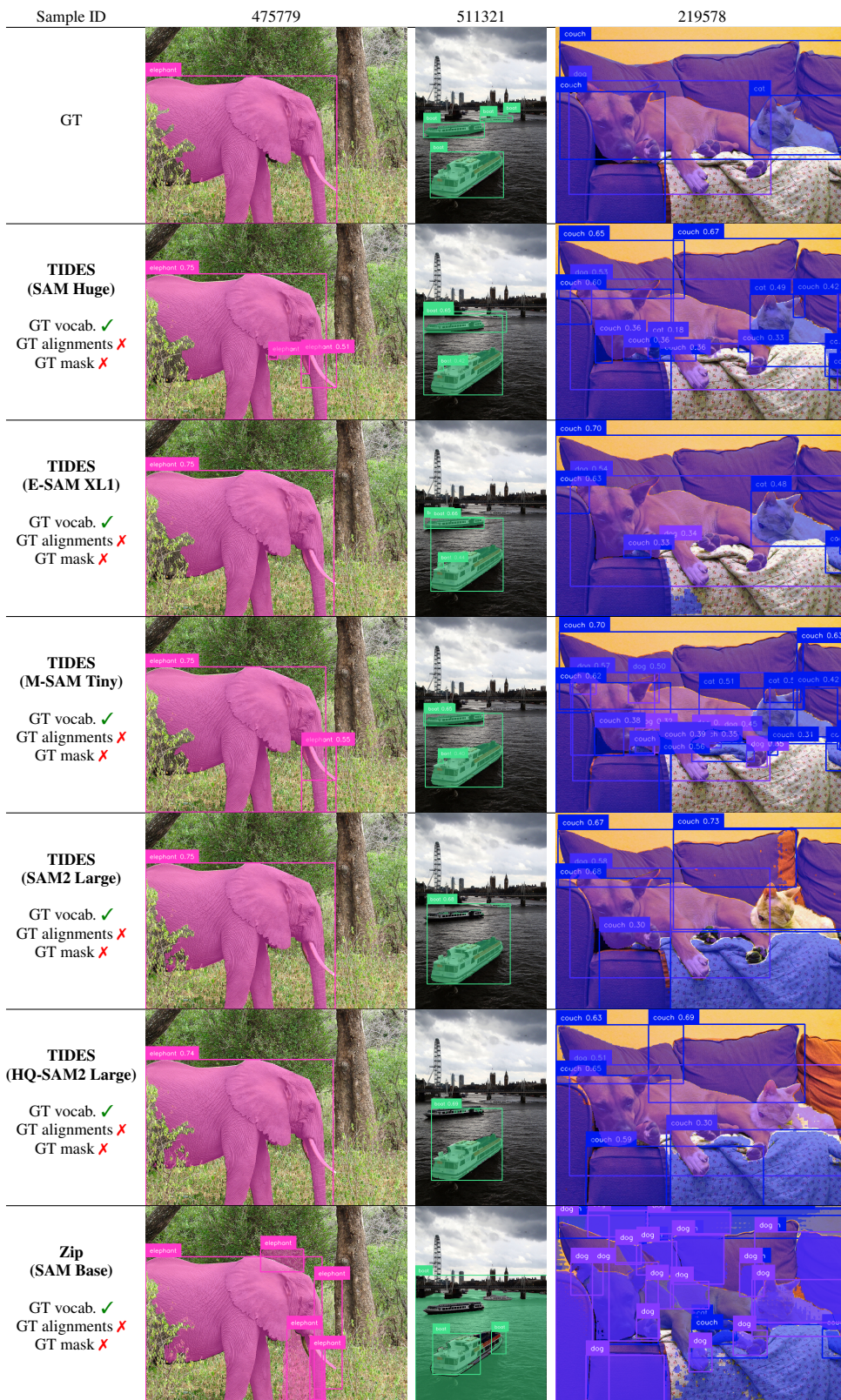


Figure 15: Difference in TIDES performance across underlying PSMs, along with Zip (Shi & Yang, 2024).

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

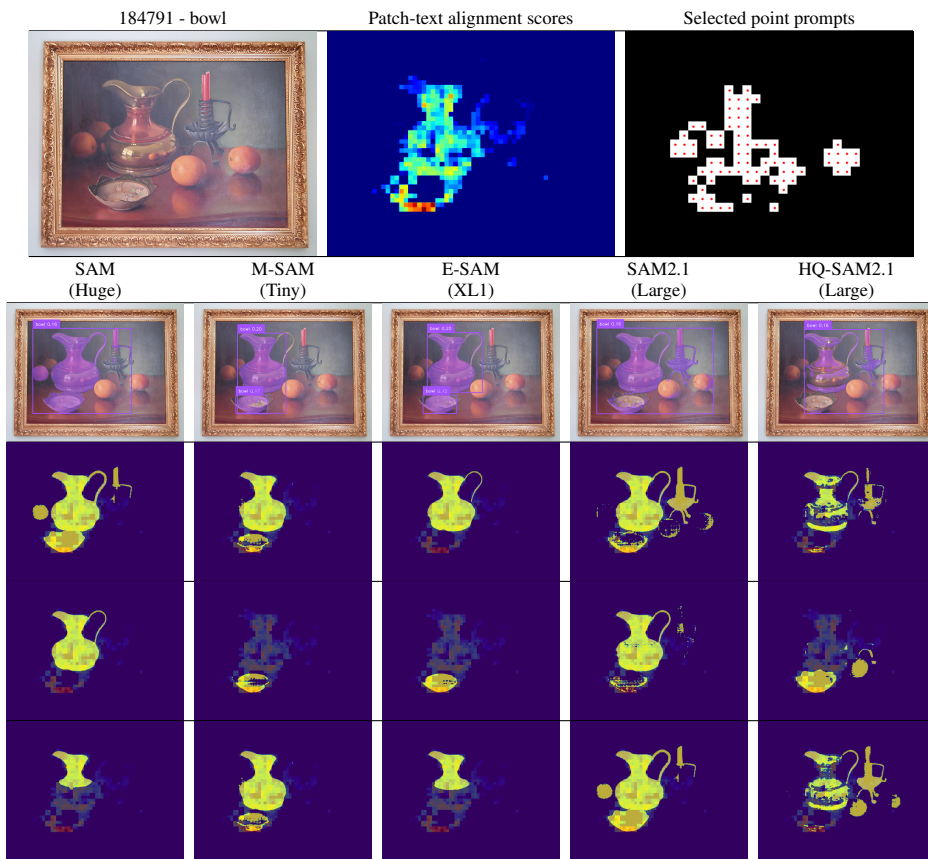


Figure 16: Patch-text alignment scores derived by CS (Li et al., 2025), selected point prompts, and PSM predictions corresponding to the top three TIDES predictions for different PSM.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

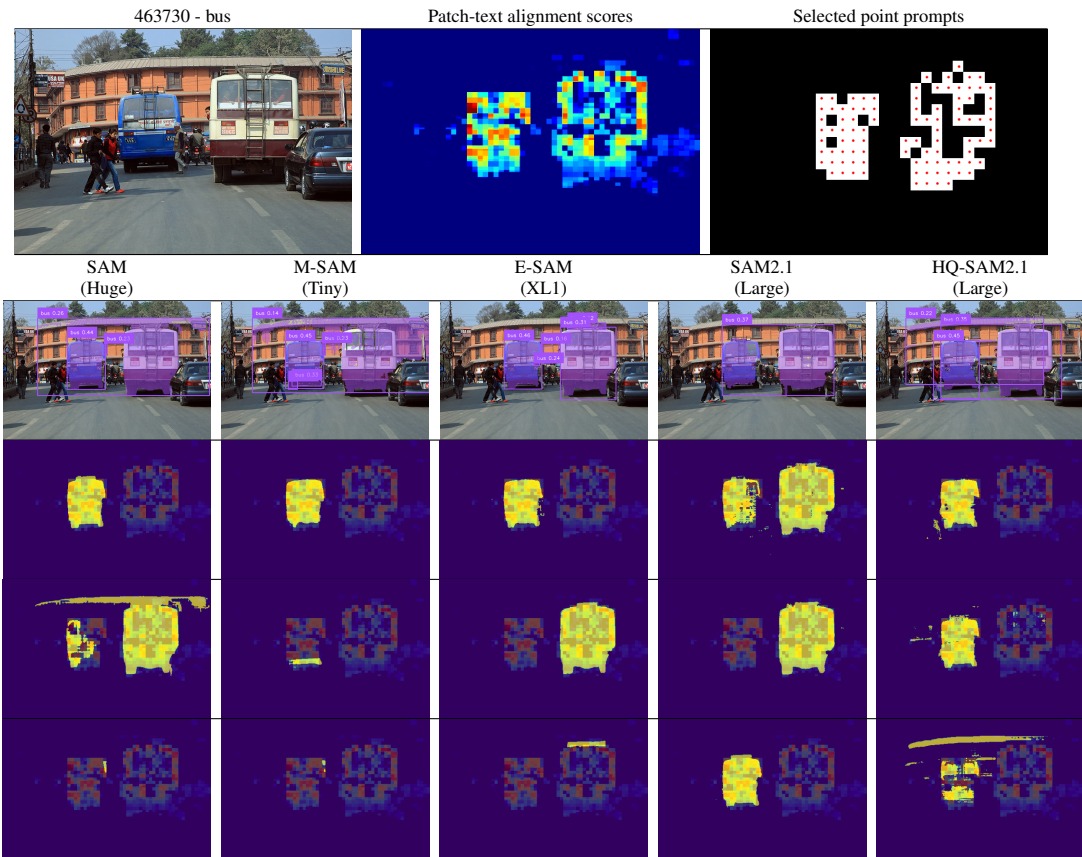


Figure 17: Patch-text alignment scores derived by CS (Li et al., 2025), selected point prompts, and PSM predictions corresponding to the top three TIDES predictions for different PSM.

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

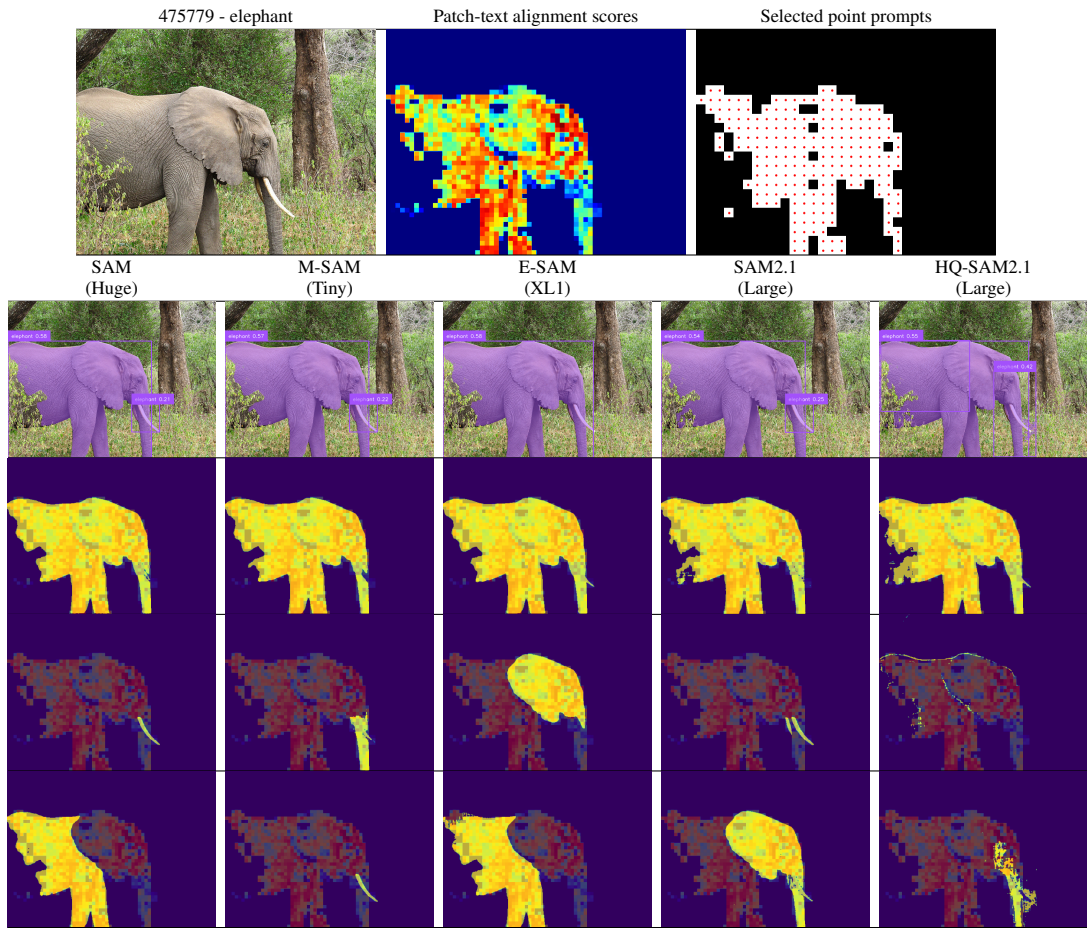


Figure 18: Patch-text alignment scores derived by CS (Li et al., 2025), selected point prompts, and PSM predictions corresponding to the top three TIDES predictions for different PSM.

1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

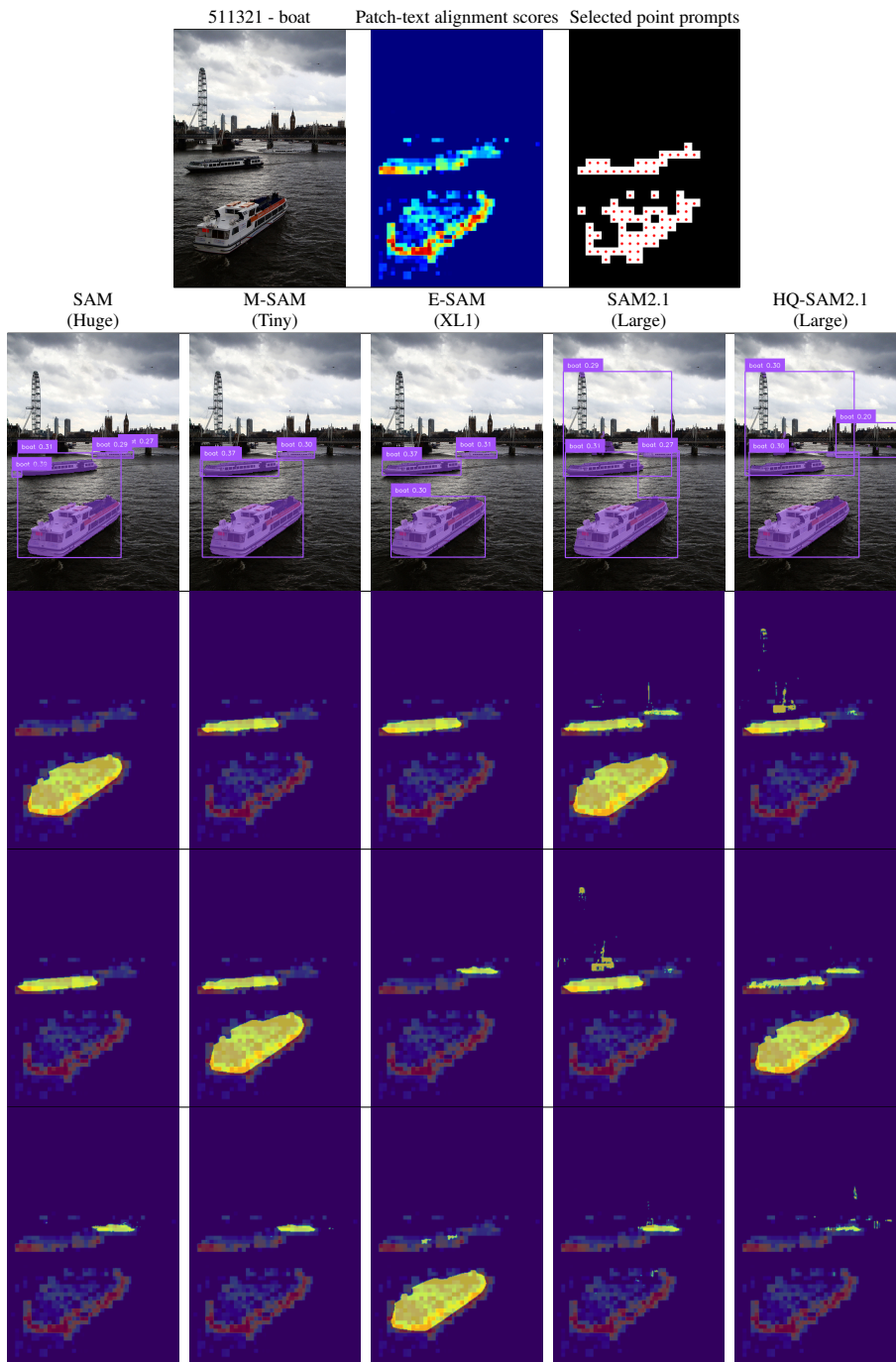


Figure 19: Patch-text alignment scores derived by CS (Li et al., 2025), selected point prompts, and PSM predictions corresponding to the top three TIDES predictions for different PSM.