

# Improving Thing Segmentation for USV Panoptic Scene Understanding with Detector-Guided Class-Specific Refinement

Anonymous MaCVi submission

Paper ID 18

## Abstract

001 *Panoptic segmentation is well suited to unmanned surface*  
002 *vehicle (USV) scene understanding because safe naviga-*  
003 *tion requires both dense parsing of background regions and*  
004 *explicit reasoning about obstacle instances. Maritime panop-*  
005 *tic benchmarks remain difficult for a familiar reason: thing*  
006 *classes are often small, sparse, and long-tailed, while stuff*  
007 *classes dominate image area. We therefore build a hybrid*  
008 *system for the LaRS maritime panoptic benchmark that com-*  
009 *bines a refined MaskDINO panoptic ensemble with an RF-*  
010 *DETR-Seg Medium thing booster and class-specific Grab-*  
011 *Cut refinement during fusion, yet can still be developed on*  
012 *a single workstation GPU. Our development experiments*  
013 *point to thing segmentation, rather than stuff parsing, as*  
014 *the main source of remaining error. This leads us to focus*  
015 *on panoptic-branch refinement, detector-guided fusion, and*  
016 *targeted boundary refinement for row boats, paddle boards,*  
017 *buoys, swimmers, and selected rare classes. On the local de-*  
018 *velopment split used for method design, the system improves*  
019 *panoptic quality from 37.50 to 45.79 and raises thing panop-*  
020 *tic quality from 17.47 to 28.49 while keeping stuff quality*  
021 *at 91.92. On the official hidden-test benchmark, the final*  
022 *anonymous submission achieves 42.6 PQ, 24.2 PQ<sub>th</sub>, 91.8*  
023 *PQ<sub>st</sub>, 51.1 RQ, 71.0 SQ, and 54.4 F1. Overall, these re-*  
024 *sults indicate that detector-guided, class-specific refinement*  
025 *can improve maritime panoptic segmentation without relying*  
026 *solely on larger end-to-end panoptic models.*

## 027 1. Introduction

028 Panoptic segmentation unifies semantic and instance seg-  
029 mentation by assigning every pixel a semantic label and, for  
030 thing regions, an instance identity [4]. For USV perception,  
031 this representation is especially attractive because navigation  
032 depends on both large background regions such as water  
033 and sky and small dynamic obstacles such as buoys, swim-  
034 mers, row boats, and floating debris. The LaRS benchmark  
035 explicitly targets this regime and is, to the best of our knowl-

edge, among the first large-scale benchmarks for maritime  
panoptic obstacle understanding [8]. 036 037

038 Despite recent progress in transformer-based segmenta-  
039 tion [3, 5], maritime panoptic segmentation remains chal-  
040 lenging. In LaRS, the three stuff classes occupy most of the  
041 scene area and are comparatively easy to predict, whereas  
042 the eight thing classes are small, rare, and visually diverse.  
043 This asymmetry suggests that, under limited compute, sim-  
044 ply scaling a unified panoptic model may not be the most  
045 effective strategy.

046 Our aim in this paper is narrower. Rather than propos-  
047 ing a new backbone or decoder, we focus on the hardest  
048 part of the task: improving thing segmentation under single-  
049 workstation-GPU constraints. The framework starts from  
050 a MaskDINO panoptic model, strengthens it through rare-  
051 thing rebalancing and thing-aware fine-tuning, and then en-  
052 sembles two complementary checkpoints at inference time.  
053 In parallel, we train an RF-DETR-Seg Medium model on  
054 thing annotations only. Its masks are refined and inserted  
055 into the panoptic prediction only when confidence, area, and  
056 overlap constraints suggest that they help without disturbing  
057 already-correct stuff regions.

058 Our contributions are threefold. First, on our develop-  
059 ment split, we show that the main gap is between thing and  
060 stuff quality rather than in overall scene parsing. Second,  
061 we present a hybrid panoptic-plus-thing pipeline that pairs a  
062 stable panoptic base with a detector-style auxiliary branch  
063 and remains workable on a single workstation GPU. Third,  
064 we provide a component-wise analysis showing that global  
065 fusion alone does not deliver the full gain; the last improve-  
066 ments in thing quality come from class-specific refinement,  
067 while stuff performance stays nearly unchanged. Taken to-  
068 gether, these results outline a simple accuracy-oriented de-  
069 sign path for this compute regime.

## 070 2. Related Work

071 **Panoptic segmentation.** Panoptic segmentation was intro-  
072 duced by Kirillov et al. [4] and has since been addressed by  
073 both multi-branch and transformer-based approaches. Early

074	systems often relied on separate semantic and instance components plus heuristic merging. More recent unified models such as Mask2Former [3] and MaskDINO [5] substantially improve performance by directly predicting masks with transformer decoders.	123
075		124
076		125
077		
078		
079	<b>Transformer detection and segmentation.</b> DETR re-framed detection as direct set prediction [2], enabling architectures that avoid many hand-crafted post-processing steps. MaskDINO extends this line to universal segmentation [5]. RF-DETR further targets real-time transformer detection [6]. In our work, we use its released segmentation variant as a practical detector-style model for instance mask prediction.	126
080		127
081		128
082		129
083		130
084		131
085		132
086	<b>Maritime visual perception.</b> Maritime perception differs from common road-scene benchmarks because targets are often tiny, sparsely distributed, and heavily affected by reflections, clutter, and large uniform water regions. SeaDronesSee [7] and MODS [1] highlight the difficulty of small obstacle perception in marine environments. LaRS extends this line by providing a diverse maritime panoptic benchmark with both thing and stuff labels [8].	133
087		134
088		
089		
090		
091		
092		
093		
094	<b>Hybrid systems.</b> Our work is closest in spirit to practical hybrid pipelines that combine a strong dense predictor with a specialized auxiliary module. We do not claim architectural novelty. Instead, we study whether, in the maritime panoptic regime, a carefully constrained fusion of a panoptic model with a thing-specialized model can provide a practical operating point under single-workstation-GPU constraints.	135
095		136
096		137
097		138
098		139
099		140
100		141
101	<b>3. Problem Setting</b>	142
102	We study the LaRS maritime panoptic benchmark [8]. Each pixel must be assigned one of 11 classes: three stuff classes (static obstacles, water, sky) and eight thing classes (boat/ship, row boat, paddle board, buoy, swimmer, animal, float, other). The output is a panoptic map in which every pixel has a semantic class, and every thing region also has an instance identifier.	143
103		144
104		145
105		146
106		147
107		148
108		149
109		150
110		151
111		152
112		153
113		154
114		155
115		156
116		157
117	<b>4. Method</b>	158
118	<b>4.1. Panoptic Base Model</b>	159
119	Our base model is MaskDINO with a ResNet-50 backbone [5]. We choose it for its strong panoptic segmentation performance and its practicality on a single workstation GPU. We train the baseline on LaRS panoptic annotations using	160
120		161
121		162
122		163
	$896 \times 896$ large-scale-jitter crops. The baseline already provides reliable stuff parsing, but small obstacle instances are often missed or under-segmented.	164
		165
	<b>4.2. Rare-Thing Rebalancing</b>	166
	The LaRS thing distribution is highly imbalanced. In the training split, boat/ship appears in 1491 images, whereas float appears in only 21. To reduce this skew, we create a rebalanced panoptic training stream by repeating images that contain rare thing categories. The repeat factors are category-dependent and largest for the rarest classes. This rebalancing increases how often the panoptic model sees small obstacle categories during fine-tuning.	167
		168
		169
		170
		171
	<b>4.3. Thing-Aware Crop Fine-Tuning</b>	172
	After rebalanced fine-tuning, we further refine the panoptic model with a thing-aware crop mapper. Instead of always sampling a random crop after resizing, the mapper selects a crop centered near a randomly chosen non-crowd thing instance with a fixed probability. The selected instance must satisfy minimum-size and maximum-area-ratio constraints so that crops preferentially expose small to medium obstacle regions rather than large easy objects. This change keeps the original training recipe largely intact while increasing the density of informative thing examples.	173
		174
		175
		176
		177
		178
		179
		180
		181
		182
		183
		184
		185
		186
		187
		188
		189
		190
		191
		192
		193
		194
		195
		196
		197
		198
		199
		200
		201
		202
		203
		204
		205
		206
		207
		208
		209
		210
		211
		212
		213
		214
		215
		216
		217
		218
		219
		220
		221
		222
		223
		224
		225
		226
		227
		228
		229
		230
		231
		232
		233
		234
		235
		236
		237
		238
		239
		240
		241
		242
		243
		244
		245
		246
		247
		248
		249
		250
		251
		252
		253
		254
		255
		256
		257
		258
		259
		260
		261
		262
		263
		264
		265
		266
		267
		268
		269
		270
		271
		272
		273
		274
		275
		276
		277
		278
		279
		280
		281
		282
		283
		284
		285
		286
		287
		288
		289
		290
		291
		292
		293
		294
		295
		296
		297
		298
		299
		300
		301
		302
		303
		304
		305
		306
		307
		308
		309
		310
		311
		312
		313
		314
		315
		316
		317
		318
		319
		320
		321
		322
		323
		324
		325
		326
		327
		328
		329
		330
		331
		332
		333
		334
		335
		336
		337
		338
		339
		340
		341
		342
		343
		344
		345
		346
		347
		348
		349
		350
		351
		352
		353
		354
		355
		356
		357
		358
		359
		360
		361
		362
		363
		364
		365
		366
		367
		368
		369
		370
		371
		372
		373
		374
		375
		376
		377
		378
		379
		380
		381
		382
		383
		384
		385
		386
		387
		388
		389
		390
		391
		392
		393
		394
		395
		396
		397
		398
		399
		400
		401
		402
		403
		404
		405
		406
		407
		408
		409
		410
		411
		412
		413
		414
		415
		416
		417
		418
		419
		420
		421
		422
		423
		424
		425
		426
		427
		428
		429
		430
		431
		432
		433
		434
		435
		436
		437
		438
		439
		440
		441
		442
		443
		444
		445
		446
		447
		448
		449
		450
		451
		452
		453
		454
		455
		456
		457
		458
		459
		460
		461
		462
		463
		464
		465
		466
		467
		468
		469
		470
		471
		472
		473
		474
		475
		476
		477
		478
		479
		480
		481
		482
		483
		484
		485
		486
		487
		488
		489
		490
		491
		492
		493
		494
		495
		496
		497
		498
		499
		500
		501
		502
		503
		504
		505
		506
		507
		508
		509
		510
		511
		512
		513
		514
		515
		516
		517
		518
		519
		520
		521
		522
		523
		524
		525
		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539
		540
		541
		542
		543
		544
		545
		546
		547
		548
		549
		550
		551
		552
		553
		554
		555
		556
		557
		558
		559
		560
		561
		562
		563
		564
		565
		566
		567
		568
		569
		570
		571
		572
		573
		574
		575
		576
		577
		578
		579
		580
		581
		582
		583
		584
		585
		586
		587
		588
		589
		590
		591
		592
		593
		594
		595
		596
		597
		598
		599
		600
		601
		602
		603
		604
		605
		606
		607
		608
		609
		610
		611
		612
		613
		614
		615
		616
		617
		618
		619
		620
		621
		622
		623
		624
		625
		626
		627
		628
		629
		630
		631
		632
		633
		634
		635
		636
		637
		638
		639
		640
		641
		642
		643
		644
		645
		646
		647
		648
		649
		650
		651
		652
		653
		654
		655
		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681</

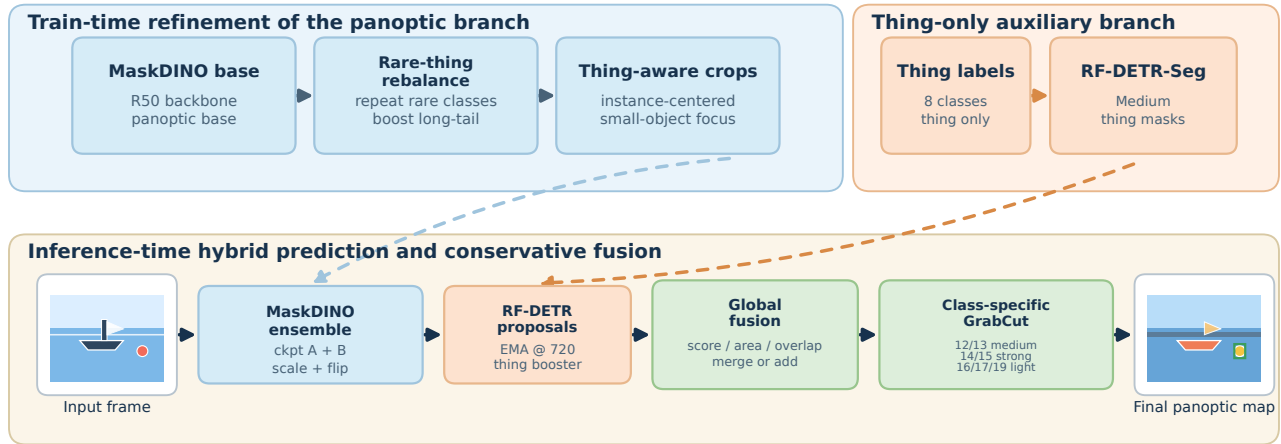


Figure 1. Overview of the proposed hybrid framework. Train-time refinement improves the MaskDINO panoptic branch through rare-thing rebalancing and thing-aware crop tuning, while a separate RF-DETR-Seg Medium model is trained only on thing classes. At inference, the refined panoptic checkpoints are ensembled, augmented with RF-DETR thing proposals, and then passed through conservative global fusion plus class-specific GrabCut refinement for the classes that benefit consistently on the development split.

172 a good trade-off between thing recall and memory usage in  
173 our development setting.

#### 174 4.6. Constrained Fusion

175 The first inference-stage fusion step inserts RF-DETR masks  
176 into the panoptic ensemble output using conservative global  
177 rules. Let  $m$  be an RF-DETR mask with confidence  $s(m)$ .  
178 We first discard low-confidence or tiny masks. If  $m$  overlaps  
179 an existing same-class panoptic instance with sufficient IoU,  
180 we merge the two. Otherwise, we insert  $m$  as a new instance  
181 only when its overlap with existing thing masks stays below  
182 a threshold. This rule is designed to add missed obstacles  
183 without degrading already-correct stuff and thing regions.

184 In the final configuration, global fusion uses a score  
185 threshold of 0.35, a minimum area of 40 pixels, a merge  
186 IoU threshold of 0.25, and a maximum overlap of 0.30 with  
187 existing panoptic instances. The key design principle is not  
188 aggressive replacement, but selective insertion of missing  
189 thing instances.

#### 190 4.7. Class-Specific GrabCut Refinement

191 Global fusion already improves thing recall, but we found  
192 that the remaining errors are often boundary-quality failures  
193 concentrated in a few thin or low-area classes. We therefore  
194 refine detector masks with a lightweight GrabCut step be-  
195 fore making the final merge decision. A global setting with  
196 margin 12, two iterations, and a minimum refined area of 48  
197 pixels provides a reasonable default configuration.

198 We then add class-specific overrides only where the de-  
199 velopment split shows consistent gains. Buoy and swimmer  
200 masks (classes 14 and 15) benefit from the strongest re-  
201 finement, with a wider context margin and more GrabCut

202 iterations. Row boat and paddle board masks (classes 12  
203 and 13) benefit from a more moderate setting. Animal, float,  
204 and other (classes 16, 17, and 19) require a more conserva-  
205 tive policy with slightly higher score thresholds and smaller  
206 valid areas. This targeted policy yields the final improvement  
207 from 45.58 PQ to 45.79 PQ on the development split while  
208 keeping  $PQ_{st}$  effectively unchanged.

## 209 5. Experimental Setup

### 210 5.1. Dataset and Evaluation

211 The local LaRS split used in our experiments contains 2605  
212 training images, 198 development images, and 1203 hidden-  
213 test images for official benchmark evaluation. We use the  
214 development split for method iteration, ablation, and fusion-  
215 threshold selection. Final submission performance is re-  
216 ported with the official hidden-test metrics: PQ,  $PQ_{th}$ ,  
217  $PQ_{st}$ , RQ, SQ, and F1.

### 218 5.2. Implementation Details

219 All experiments were developed under constrained hardware.  
220 The MaskDINO runs use a single NVIDIA RTX A4000  
221 GPU with batch size 1 and image size 896. The rare-thing  
222 fine-tuning stage runs for 16k iterations, and the thing-aware  
223 stage runs for 8k additional iterations starting from the best  
224 rare-thing checkpoint. The RF-DETR-Seg Medium model  
225 is trained at resolution 432 with batch size 1 and gradient  
226 accumulation of 8, then evaluated with the EMA checkpoint  
227 at inference resolution 720. The final fusion stack uses the  
228 global thresholds described above, followed by class-specific  
229 GrabCut overrides for classes 12–17 and 19. Here, *resource-*  
230 *efficient* refers to development feasibility on a single work-

Thing class	Train imgs	Train inst.
Boat/ship	1491	5921
Row boat	232	444
Paddle board	110	154
Buoy	734	1568
Swimmer	122	349
Animal	72	361
Float	21	23
Other	281	501

Table 1. Long-tail frequency of LaRS thing classes in the training split. The rarest classes have very limited supervision, motivating our rebalancing and thing-focused design choices.

231 station GPU rather than to a deployment-optimized real-time  
 232 system. We therefore do not report end-to-end FPS for the  
 233 final pipeline, because the submitted method combines a  
 234 two-checkpoint panoptic ensemble, horizontal-flip augmen-  
 235 tation, an auxiliary detector branch, and GrabCut-based post-  
 236 processing, all of which were tuned for benchmark accuracy  
 237 rather than throughput.

### 238 5.3. Class Difficulty Analysis

239 Table 1 summarizes the long-tail nature of thing categories.  
 240 Classes such as animal and float are extremely rare in the  
 241 training set, which supports our claim that the main difficulty  
 242 lies less in general scene parsing than in long-tail thing  
 243 segmentation.

## 244 6. Results

### 245 6.1. Main Results

246 Table 2 reports both development-set and official benchmark  
 247 results. The development rows show how the submitted  
 248 system was assembled, while the official row gives the final  
 249 blind evaluation. Relative to the base MaskDINO model, the  
 250 final system gains 8.29 PQ points on the development split,  
 251 including an 11.02-point increase in  $PQ_{th}$ . By contrast,  
 252  $PQ_{st}$  changes by less than one point across the full pipeline.  
 253 This pattern again points to thing segmentation, rather than  
 254 stuff parsing, as the main source of remaining error on the  
 255 development data.

256 The numbers also clarify where the gains come from.  
 257 The panoptic model already handles the globally dominant  
 258 stuff regions and provides a reasonable obstacle layout, so  
 259 the later stages can concentrate on underrepresented thing  
 260 instances. The biggest jumps appear in the detector-guided  
 261 stages: global RF-DETR fusion lifts PQ from 41.14 to 43.68,  
 262 and global GrabCut refinement raises it again to 45.58. Class-  
 263 specific policies contribute smaller but repeatable gains, with  
 264 the clearest increments coming from buoy and swimmer  
 265 refinement.

Split	Method	PQ	$PQ_{th}$	$PQ_{st}$	Notes
Dev	Base	37.50	17.47	90.93	single
	MaskDINO				checkpoint
Dev	+ panoptic	39.99	20.88	90.93	rebalance +
	branch re-				crop tuning
	finement				
Dev	+ panoptic en-	41.14	22.15	91.78	v6 + v31,
	semble				scale 896 +
					flip
Dev	+ global RF-	43.68	25.60	91.89	Medium,
	DETR fusion				res 720
Dev	+ global	45.58	28.20	91.92	margin 12,
	GrabCut				iters 2
	refinement				
Dev	+ class-	<b>45.79</b>	<b>28.49</b>	<b>91.92</b>	final dev
	specific				system
	refinement				
					RQ=51.1,
					SQ=71.0,
Official	Final submis-	42.6	24.2	91.8	F1=54.4
	sion				

Table 2. Main results. The development split is used for method analysis, while the official row reports the final anonymous challenge submission. The final system combines a panoptic ensemble, RF-DETR-guided fusion, and class-specific GrabCut refinement.

Variant	PQ	$PQ_{th}$	$\Delta PQ$
Panoptic ensemble only	41.14	22.15	-
+ global fusion only	43.68	25.60	+2.54
+ global GrabCut refine-	45.58	28.20	+1.90
ment			
+ class-specific on 14,15	45.71	28.38	+0.13
+ extend to 12,13	45.77	28.46	+0.06
+ conservative 16,17,19	45.79	28.49	+0.02

Table 3. Ablation of the final refinement stack on the development split. Gains are measured relative to the preceding row. The main improvements come from global fusion and global GrabCut refinement, while class-specific policies provide smaller but consistent final gains.

### 6.2. Ablation of the Refinement Policy

266 Table 3 shows that global fusion alone is not enough. The  
 267 largest increment after detector insertion comes from refining  
 268 detector boundaries before merge decisions, indicating that  
 269 recall and mask quality must be improved together. The final  
 270 class-specific steps are much smaller in magnitude, but they  
 271 remain useful because the best operating point differs for  
 272 thin classes such as row boats and paddle boards and for  
 273 buoy- and swimmer-like instances. 274



Figure 2. Qualitative comparison on validation images. The two columns show selected validation examples, while the rows show the input crop, the global-fusion baseline, the class-specific final prediction, and the ground-truth overlay. The class-specific final row is outlined in green.

### 275 6.3. Qualitative Analysis

276 Figure 2 shows two validation examples that reflect the same  
277 trend. In both cases, the globally fused baseline already  
278 captures the dominant water/sky layout, and class-specific  
279 refinement improves selected thing regions without changing  
280 the overall scene structure. This is in line with the quantita-  
281 tive results: most of the gain comes from targeted corrections  
282 on thing regions, not from large changes in stuff prediction.

283 The official benchmark result of 42.6 PQ and 24.2  $PQ_{th}$   
284 is lower than the development score, possibly due in part  
285 to blind-test distribution shift, but the overall picture is un-  
286 changed: the final system is solid on stuff and remains lim-  
287 ited mainly by rare thing segmentation.

## 7. Discussion

Our results suggest that USV panoptic segmentation is not uniformly hard across all parts of the label space. In LaRS, a strong panoptic model already handles stuff comparatively well. The harder cases are the long-tail thing categories. That is why a detector-style auxiliary model and targeted refinement help in ways that pure panoptic scaling does not fully recover. In that sense, a staged accuracy-first design remains a sensible option when memory and compute make end-to-end scaling expensive.

The study also exposes the limits of our current approach. First, the class-specific refinement policy is manually tuned on a small development split and may not transfer perfectly across domains. Second, the RF-DETR booster still struggles on the rarest categories and under severe appearance changes. Third, our system is image-based; restricted temporal experiments plateaued below the best class-specific refinement configuration and were therefore excluded from the final submission. These limitations point to clear next steps: learned refinement policies, stronger rare-class supervision, and more reliable temporal propagation.

## 8. Conclusion

We presented a hybrid framework for maritime panoptic segmentation that combines a MaskDINO panoptic base model, thing-aware fine-tuning, a panoptic ensemble, RF-DETR-guided fusion, and class-specific GrabCut refinement. The method is built around the part of LaRS that remains hardest: small and rare thing instances. On the development split, the largest gains come from detector-guided refinement, and the final anonymous submission reaches 42.6 PQ with 24.2  $PQ_{th}$  and 91.8  $PQ_{st}$  on the official hidden-test benchmark. Our main takeaway is straightforward: under single-workstation-GPU constraints, a hybrid pipeline can be a better use of limited compute than relying only on a larger unified panoptic model.

## References

- [1] Borja Bovcon, Jon Muhovič, Dušan Vranac, Dean Mozetič, Janez Perš, and Matej Kristan. MODS: A USV-oriented object detection and obstacle segmentation benchmark. *arXiv preprint arXiv:2105.02359*, 2021. 2
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020. 2
- [3] Bowen Cheng, Alexander Schwing, and Alexander Kirillov. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1290–1299, 2022. 1, 2
- [4] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *Pro-*

- 339 *ceedings of the IEEE/CVF Conference on Computer Vision*  
340 *and Pattern Recognition (CVPR)*, pages 9404–9413, 2019. 1,  
341 2
- 342 [5] Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang,  
343 Lionel M. Ni, and Heung-Yeung Shum. Mask dino: Towards a  
344 unified transformer-based framework for object detection and  
345 segmentation. In *Proceedings of the IEEE/CVF Conference*  
346 *on Computer Vision and Pattern Recognition (CVPR)*, pages  
347 3041–3050, 2023. 1, 2
- 348 [6] Isaac Robinson, Peter Robicieux, Matvei Popov, Deva Ra-  
349 manan, and Neehar Peri. RF-DETR: Neural architecture  
350 search for real-time detection transformers. *arXiv preprint*  
351 *arXiv:2511.09554*, 2025. 2
- 352 [7] Leon Amadeus Varga, Benjamin Kiefer, Martin Messmer,  
353 and Andreas Zell. Seadronessee: A maritime benchmark  
354 for detecting humans in open water. In *Proceedings of the*  
355 *IEEE/CVF Winter Conference on Applications of Computer*  
356 *Vision (WACV)*, pages 2260–2270, 2022. 2
- 357 [8] Lojze Žust, Janez Perš, and Matej Kristan. Lars: A diverse  
358 panoptic maritime obstacle detection dataset and benchmark.  
359 In *Proceedings of the IEEE/CVF International Conference on*  
360 *Computer Vision (ICCV)*, pages 20304–20314, 2023. 1, 2