001

002

003

004

005

006

007

008

009

010

011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035

036

037

038

039

040

054

#### 055 056 **Refined and Enriched Captions With Physical Scale For Dynamic Disaster Scene**<sup>057</sup> 058 059 060 Anonymous CVPR 2023 submission 061 062 Paper ID 2 063 064 065 with heavy fog and snowfall events. Moreover, unpredicted<sup>066</sup> Abstract disaster and traffic accident scenes require more pretrain-067

Vision-Language models (VLMs), i.e., image-text pairs of CLIP, have boosted image-based Deep Learning (DL). Unseen images by transferring semantic knowledge from seen classes can be dealt with the help of language models pre-trained only with texts. Two-dimensional spatial relationships and a higher semantic level have been performed. Moreover, Visual-Question-Answer (VQA) tools and openvocabulary semantic segmentation provide us with more detailed scene descriptions, i.e., qualitative texts, in captions. However, the capability of VLMs presents still far from that of human perception. This paper proposes PanopticCAP for refined and enriched qualitative and quantitative captions to make them closer to what human recognizes by combining multiple DLs and VLMs. In particular, captions with physical scales and objects' surface properties are integrated by water level, counting, depth map, visibility distance, and road conditions. Fine-tuned VLM models are also used. An iteratively refined caption model with a new physics-based contrastive loss function is used. Experimental results using images with adversarial weather conditions, i.e., rain, snow, fog, landslide, flooding, and traffic events, i.e., accidents, outperform state-of-the-art DLs and VLMs. A higher semantic level in captions for real-world scene descriptions are shown.

#### 1. Introduction

041 Segmentation has become an important task for realworld applications by Deep Learning (DL) [4, 5, 8-12, 15, 042 043 18-20, 26, 28, 33, 37, 39, 40, 50, 51, 56-58, 60, 61, 69-72, 74, 75, 79-81, 88, 93-95, 98, 106, 112, 115, 116, 120, 044 124, 126, 132–134, 152]. Segmentation has become diver-045 sified into semantic and instance segmentation. Although 046 047 many variants of segmentation models are presented, issues 048 with limits of training image datasets and robustness to illumination and noise remain unsolved. Only pretraining a 049 finite number of images could not have enhanced segmenta-050 tion accuracy in real-world scenes. Dynamic changes have 051 052 been dealt with by rain drops [98, 134], defog, and dehaze 053 [40, 61, 71, 74, 89, 132]; however, these methods fail to deal with heavy fog and snowfall events. Moreover, unpredicted<sup>066</sup> disaster and traffic accident scenes require more pretrain-<sup>067</sup> ing image datasets; however, in spite of vital events, they<sup>068</sup> are hard to collect sufficient images and videos due to rare<sup>069</sup> chances. Therefore, segmentation to such conditions and<sup>070</sup> events becomes degraded.

Recently, CV, DL, and NLP have been combined, i.e., 072 Vision Language Model (VLM). It is known that unseen 074 images that have not been pretrained have been recognized much better than only CV or DL models [3, 6, 14, 23, 25, 076 )75 34, 41–43, 48, 54, 55, 64, 67, 83, 86, 91, 96, 97, 100, 108, 077 111, 114, 117, 121, 129, 139, 140, 142, 144, 146, 147, 151] 078 VLMs can understand vision and text, allowing them to 079 perform tasks requiring multimodal understanding, i.e., Visual Question Answer (VQA), image captioning, or image retrieval. Moreover, VLMs can be pre-trained on large  $_{082}^{001}$ datasets [64, 91, 100] and fine-tuned on smaller datasets for specific tasks, allowing for efficient transfer learning [3, 6,084 14, 23, 25, 34, 41–43, 48, 54, 55, 67, 83, 86, 96, 97, 108, 085 111, 114, 117, 121, 129, 139, 140, 142, 144, 146, 147, 151].

This can be useful in various applications, such as  $object_{087}^{000}$  detection [29, 30, 32, 38, 49, 59, 77, 82, 90, 92, 102, 119,088 128, 138, 143, 149], segmentation [24, 36, 63, 76, 78, 84,089 87, 110, 127, 131, 137, 145, 150], and classification. VLMs<sub>090</sub> can save time and resources in various applications and im-<sub>091</sub> prove semantic understanding by recognizing relationships<sub>092</sub> between objects and concepts and developing a comprehen-<sub>093</sub> sive understanding of visual content.

Image captioning is an important and challenging task in<sub>095</sub> computer vision that involves generating natural language<sub>096</sub> descriptions of complex visual scenes that include objects<sub>097</sub> and their surrounding context. However, single VML is of-<sub>098</sub> ten weak for dynamic changes, i.e., disaster scenes [113].<sub>099</sub> Moreover, heavy rainfall and snowfall have been increasing,<sub>100</sub> which may cause a chain reaction of natural disasters ob-<sub>101</sub> served from the satellite images, i.e., landslides and flood-<sub>102</sub> ing [17, 44, 125].

However, camera image-based post-disaster object104 recognition for dirt, water, and rocks remains unsolved on105 the road. Since domain adaptation segmentation DL mod-106 els [46, 118] require manual selection of the optimal pre-107

113

114

115

116

117

118

119

120

121

122

131

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

169

170

108 trained model, they are not useful for dynamic changes. 109 However, few papers effectively integrate physical scale 110 into the VLM model and combine multi-VLMs. 111

For images to be best captioned, they need to depict information most similar to Human perception. Human perception can simultaneously process different visual cues, i.e., texture, shape, and depth, to identify and label objects at varying distances. Therefore, a method for refining and enriching captions with different visual cues is needed.

To this end, this paper proposes PanopticCAP: a panoptic vision-language model under adversarial visual conditions using single images. This paper proposes refined and enriched captions for scene descriptions under adversarial conditions by the proposed PanopticCAP with multiple task-oriented DLs and VLMs.

123 PanopticCAP consists of eleven modules, i.e., Deep Vi-124 sual Language Classification (Dvlc), Deep Visual Language 125 Segmentation (Dvls), Visual-Query-Answer (VQA), Con-126 trastive Language Physical scale Pretraining (CLPP), Deep 127 Visibility estimation (Dvis), Deep Road conditions (Droad), 128 Deep Depth (Ddepth), Deep anomaly (Danomal), Deep 129 water-level (Dwater), Deep snowfall (Dsnow), and Deep 130 Scene (DScene). The branched architecture allows us to efficiently maintain and upgrade each of the eleven modules. 132 Contributions of this paper are fourfold:

> 1. Multiple vision language and Transformer-based Deep Learning (DL) models with branched structures for efficiency in light of memory, training, and maintenance. Danomal excludes difficult images, i.e., lenz reflection, to stabilize the overall system. Due to enormous datasets of VLMs, Dvls, and Dvlc are fine-tuned VLMs from SOTA models for segmentation and classification, respectively.

2. It is the first time to contain dynamic changes with physical scales, i.e., depth by Ddepth, fog visibility distance by Dvis, weather conditions by Dsnow, water level by Dwater, and road conditions by Droad. Unseen images like adversarial weather and disaster conditions can be dealt with. Moreover, more detailed scene descriptions of traffic accident events are shown. Captions with 3D-related adverbs, i.e., behind, rear, in front of, and far, enable to generate as SOTAs have used 2D-related adverbs, i.e., left and right.

3. More refined and enriched captions are generated 153 based on fixed queries at VQA, CLPP, and the above 154 155 multiple modules. A new contrastive loss function is 156 proposed in CLPP to refine and enrich captions with the object's physical scale, i.e., size and position, un-157 der an iterative refinement process. API tools, i.e., 158 Visual ChatGPT [123], may be hard to generate dy-159 160 namic scene changes with physical scales as this paper 161 presents.

4. Many experimental results show the superiority of the<sup>162</sup> proposed PanopticCAP over SOTA DLs and VLMs.<sup>163</sup> The proposed PanopticCAP will help notify detailed<sup>164</sup> scene descriptions, i.e., more quantitative texts, to<sup>165</sup> drivers, auto-driving, and rescue workers from camera<sup>166</sup> 167 images. 168

#### 2. Related Work

This section briefly describes Computer Vision (CV),171 Deep Learning (DL), and Vision Language Model (VLM)172 concerning methods and issues in scene understanding of173 camera images under various conditions. Visibility lev-174 els are one of the most important visual factors to esti-175 mate for monitoring and auto-driving. To estimate visibil-176 ity, segmentation-based DL models have been reported and 177 used by Dvis [2] and Droad [1]. 178

An all-in-one image restoration network for unknown179 corruption has been proposed [62]; however, this method 180 can be degraded heavy fog and snowfall, as shown in Dvis181 [2] and Droad [1]. Dehazing in [33] is limited to closer182 views of daytime lighter foggy scenes, i.e., indoor and gar-183 den, unlike our proposed method [2] for distant scenes with184 heavy fog at night, i.e., highway. A unified framework for185 depth-aware panoptic segmentation has been reported [57]186 under clear weather conditions. 187

Although Cityscapes with 3000 images [21], Foggy188 Cityscape DBF with 500 synthetic foggy images [103], and 189 Foggy Zurich with 3800 real light foggy images [104] are 190 publicly available, they are almost all daytime and lighter191 fog data. Moreover, image datasets for road conditions have192 not been built, unlike Droad [1]. 193

In recent years, the Vision Language model (VLM)194 field has undergone significant progress [141, 147, 148].195 But most of them are pre-trained with large-scale training196 datasets and fine-tuned with task-specific annotated train-197 ing data. The pre-training of VLMs has been explored using198 three main approaches: contrastive objectives [22, 73, 99],199 generative objectives [47, 122][108, 109], and alignment200 objectives [35, 52, 85, 101, 130, 136]. VLMs are trans-201 ferred by Text-Prompt Tuning [3, 6, 23, 42, 43, 55, 83, 86,202 108, 111, 129, 139, 142, 146, 147, 151]. 203

Besides finetuning, knowledge distillation is a method204 to improve VMLs for downstream tasks, including object205 detection [29, 30, 32, 38, 49, 59, 77, 82, 90, 92, 102, 119,206 128, 138, 143, 149] and semantic segmentation [24, 36, 63, 207 76, 78, 84, 87, 110, 127, 131, 137, 145, 150]. 208

Unseen images that have not been pre-trained have be-209 come recognized by VLM frameworks [13, 31]. More di-210 verse and out-of-distribution data for pre-training and eval-211 uation are used [45]. Prompt learning to adapt VLMs to212 new tasks without fine-tuning is also shown [53]. Contents213 of captions have been enhanced for better descriptions of214 real-world objects [31]. 215

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

281

282

Geometric reasoning or depth estimation to infer 3-D information from 2-D images [135, 141] is shown using 3D point-cloud data and indoor scenes. Pretraining VLMs require over 100 million image-text paired datasets for high accuracy, more than DL models require. Therefore, many efficient models have been introduced [13, 31, 65, 68, 100, 105, 109, 153]. However, laborious and time-consuming tasks remain unsolved in pretraining VLMs.

Visual ChatGPT API tool has become famous as the image-text captioning tool. The advantage of Visual Chat-GPT [123] is that it can produce acceptable results on the general scene and unseen classes. However, since Visual ChatGPT [123] is trained on the limited data of the year 2021, it generates captions under older datasets. So far, Visual ChatGPT [123] is weak at generating dynamic scene descriptions like weather and road conditions. Moreover, the physical size of objects, fog visibility distance, and water level are contained.

Therefore, as aforementioned above, no SOTA VLM papers and API tools have challenged images with the physical scale in natural phenomena, i.e., fog visibility in 3D depth and water level. More refined and enriched captions will be provided by the proposed PanopticCAP in the following sections.

#### 3. Proposed Method

This section describes the proposed PanopticCAP method/system for refinement and enrichment of captioning and classes from a single image input. Instead of using only vision models or a single vision-language model, this paper proposes a new architecture that integrates multiple Deep Learning and vision-language modules.



Figure 1: Overview of the proposed PanopticCAP model.

Figure 1 shows an overview of the proposed Panoptic-CAP. Since this paper deals with many challenging scenes with disasters and car accidents, adversarial conditions are taken into account. Further detailed explanations of the multiple modules will be given in Sub-sections 3.1 to 3.4.

#### 3.1. Proposed Danomal

To identify and reject adversarial images, as shown in Figure 2, an algorithm to reject such images is proposed to avoid the degradation of the cascaded other recognition

modules. Three major adversarial image patterns have been<sup>270</sup> selected: a) lens reflection, b) strong headlight, and c) rain-<sup>271</sup> drops. These adversarial images were collected from over<sup>272</sup> 2500 images and used to train by Swin Transformer [37]<sup>273</sup> into 3 classes.<sup>275</sup>

A COLORADO		and the second	CALLS.	213
60	De la	1 × 1		276
	100			277
	(a)	(b)	(c)	278

Figure 2: Examples of rejected images: (a) Lens reflection. (b) Strong279 headlight. (c) Raindrops. 280

#### **3.2. Proposed Dvlc and Dvls**

Dvlc is a vision-language model trained on image and<sup>283</sup> text pairs that can predict the most relevant text given an<sup>284</sup> image. It does not need to be directly optimized for this<sup>285</sup> task and can perform "zero-shot" learning like GPT-3 and<sup>286</sup> -4. Dvlc matches the performance of the original ResNet50<sup>287</sup> on ImageNet "zero-shot" without using any of the original<sup>288</sup> 1.28M labeled examples, which is a significant accomplish-<sup>289</sup> ment in Computer Vision. 290

Dvlc utilizes the input texts of five distinct disaster cate-<sup>291</sup> gories: *car crashes, flooding, fog, landslide, and rain.* Tai-<sup>292</sup> lored textual input descriptions are employed for each dis-<sup>293</sup> aster category to enhance natural language processing tech-<sup>294</sup> niques in analyzing disaster-related data. These scenes are<sup>295</sup> associated with domain-specific terms to improve the accu-<sup>296</sup> racy of automated disaster detection and classification. <sup>297</sup>

Dvls is proposed to obtain semantic segmentation of<sup>298</sup> these scenes. Dvls is finetuned from OvSeg [76] by adding<sup>299</sup> a new physical constraint to the loss function. To obtain de-<sup>300</sup> scriptions of disasters for the Dvlc, a classification task is<sup>301</sup> performed using keywords corresponding to each disaster<sup>302</sup> scene. These texts are used to generate text descriptions of<sup>303</sup> the disasters that are fixed for each type of scene. <sup>304</sup>

Therefore, since Dvlc and Dvls recognize texts and seg-<sup>305</sup> mentation from a single image, this paper proposes to com-<sup>306</sup> bine respective outputs. 307

#### 3.3. Proposed CLPP

309 310

The proposed Contrastive Language Physical-Scale Pre-311 training (CLPP) is a VLM with inputs from a depth map, 312 object location from image and text description pairs, and 313 a modified contrastive loss function. Unlike SOTA VLMs 314 with no physical models in contrastive loss functions, this 315 paper proposes CLPP with additional physical constraints, 316 as shown in Figure 3. The original contrastive loss function 317 of CLIP [100] is defined by

$$L = \frac{1}{2}(1-Y).D^2 + \frac{1}{2}Y.max(0,m-D)^2$$
(1)<sup>319</sup>  
320

where Y is the binary label indicating whether the text and 321 image are similar or dissimilar, D is the distance between 322 the learned embeddings of the text and image, and m is the 323

335

336

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

385

386

387

388

389

390

391

392

393

margin hyperparameter, i.e., 0.2. In order to incorporate the physical scale, including the size and location of objects, i.e., meters, a similarity metric is added. The modified contrastive loss is then defined as

$$L = \frac{1}{2}(1-Y).D^{2}.(1-sim) + \frac{1}{2}Y.max(0,m-D)^{2}.sim$$
(2)

where sim is the physical similarity between the text description and the image with object location. sim is computed as the Euclidean distance between the location of objects in the image and its description in the text. sim is defined by

$$sim = w_s \cdot E(S_T, S_I) + w_l \cdot E(R_T, R_I)$$
 (3)

where:  $w_s$  is the weight of an object physical size, and  $w_l$ 337 is the weight of object's physical location, normally,  $w_s$  and 338  $w_l$  are both set equal to 0.5.  $E(S_T, S_I)$  is the Euclidean dis-339 tance between the physical size in image  $S_I$  and in the text 340 description  $S_T$ .  $E(R_T, R_I)$  is RMSE between the physical 341 object location in the image  $R_I$  and in text description  $R_T$ . 342 The physical size of the object is determined based on the 343 ratio between the object size in pixels and the object size in 344 meters as labeled in the dataset. 345

When using cosine similarity as the distance metric Din the contrastive loss function, which ranges from -1 to 1, the margin hyperparameter is typically set to a small value, i.e., 0.2 to 0.5.



Figure 3: Proposed contrastive language for pre-training in physical scale.

#### 3.4. Proposed Droad, Dsnow, and Dvis

This section discusses the proposed Droad, Dsnow, and Dvis as explained in Section 3. Unlike SOTA papers in DLs and VLMs, this paper aims to generate dynamic scene changes with the weather conditions, i.e., rain, snow, and fog, and road conditions, i.e., dry, wet, and snow. Dvis [2] and Droad [1] are applied for further detailed classes of segmented objects. Dscene [1, 2] is also applied to ensure snow conditions.

In Dvis [2], Droad [1], and Dscene [1, 2], Swinformer [126] is trained from over 7500 winter road images. It is noted that since publicly available annotation datasets are insufficient, various weather and road scenes from different countries under adversarial conditions have been collected and used to train.

#### 3.5. Dwater

375Dwater is a simple combination of Dscene and a376transformer-based classifier, specifically the ViT [27] clas-377sifier. Dwater estimates the water level based on reference

objects, such as cars, buses, humans, trees, poles, and traf-<sup>378</sup> fic signs, through two steps: In Step 1, objects are extracted <sup>379</sup> from the image using Dscene. In Step 2, the extracted ob-<sup>380</sup> jects are classified to their respective water levels, which<sup>381</sup> are pre-defined for each object. Table 1 shows the physical <sup>382</sup> height of the reference object for water levels. <sup>383</sup>

Level/Objects	human (m)	car (m)	pole (m)
Lv1	0.3	0.3	0.5
Lv2	0.6	0.6	1
Lv3	0.9	0.9	2
Lv4	1.3	1.2	3
Lv5	1.7	1.5	4

Table 1: Physical height of reference objects.

#### 4. Experiments and Discussion

#### 4.1. Danomal for adversarial conditions

This section evaluates the performance of the proposed <sup>394</sup> Danomal. The dataset comprises over 2500 images with <sup>395</sup> 1 normal condition and 3 different adversarial conditions: <sup>396</sup> lens reflection, strong light, and raindrops. In a compara-<sup>397</sup> tive study, Droad, Dsnow, Dvis, and Dvlc are applied with <sup>398</sup> and without Danomal. Evaluation is conducted using three <sup>399</sup> classes of road condition, three classes of snowfall, four <sup>400</sup> classes of visibility, and five classes of scene types. <sup>401</sup>

Table 2 shows the results that Danomal can effectively <sup>402</sup> reject images with adversarial conditions, where the accuracy using Danomal becomes better for each DL model <sup>404</sup> without Danomal. No thresholding setting is required. <sup>405</sup> Therefore, the proposed Danomal has been proven useful <sup>406</sup> in rejecting such three adversarial factors in road images. <sup>407</sup> <sup>408</sup>

 Table 2: Accuracy comparison with and without DAnomal under adversar-409
 ial conditions for Droad, Dsnow, Dvis, and Dvlc.
 410

	Without Danomal (%)	With Danomal (%)	41
Droad	81.31	86.07	
Dsnow	72.90	78.22	41
Dvis	75.58	80.86	41
Dvlc	91.78	92.65	/11

#### 4.2. Scene Recognition Capability of CLPP and 415 416 417 417

To evaluate the performance of scene analysis of the pro-418 posed system, this section demonstrates the qualitative and419 quantitative accuracies of CLPP and Dvlc. The test dataset420 consists of six selected categories: a car crash in snow con-421 ditions, flooding with rain, low visibility with fog, land-422 slide, wet road with rain, and traffic flow. Figure 4 (a)-(f)423 shows such images. In the test dataset, ground truth cap-424 tions have been manually annotated and contain the follow-425 ing key criteria: the number of vehicles, scene category, and426 road conditions. Totally, 3500 images are collected. The ex-427 periments have been conducted to evaluate the performance428 of scene understanding using VLMs on dynamic scenes as429 well as the level of detail of the description based on the430 generated captions. BLIP [66] model has been chosen to431

433

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

485

Table 3: Comp	parison of captions	among grou	nd truth, propose	d CLPP, and
BLIP [66]				
	Ground truth	CLPP	BLIP [66]	

	Ground truth	CLPP	BLIP [00]
(a)	7 cars under heavy snowfall in a crashed scene	7 cars on the frozen road, 1 severely damaged car	A car is stuck in the snow
(b)	4 vehicles under flooding scene	Cars and motorcycle on the flooded highway	A man is crossing the street in the rain
(c)	Empty highway under heavy foggy scene	Highway under heavy foggy at daytime	Foggy road in the mountains
(d)	One car on the road, in a landslide scene	One car on the damaged road occluded by the rock	A tractor is parked on the side of a road next to a pile of rocks
(e)	One car on the rainy road	A street under a light rain at night	The image is of a street intersection at night

compare since the achievement of high capability of transfer flexibly for vision language understanding and captioning tasks.

Table 3 shows the captions containing ground truth and two captions generated by the proposed CLPP (m = 4 in eq. 2) and BLIP [66] on six images in Figure 4. CLPP and BLIP [66] can recognize the scene's context and key objects, i.e., traffic conditions and several objects. However, BLIP [66] is insufficient to capture detailed scene features compared with CLPP. In (a), (b), and (f), the counting number of vehicles by CLPP has been improved from one to many. In (b), no road conditions by BLIP [66] have been recognized, but CLPP recognizes the flooding highway. In (e) and (f), CLPP recognizes rain and sunny, respectively.

However, BLIP [66] shows no weather conditions. In (c), it is shown that the degree of fog is recognized as heavy fog by CLPP with time. In (d), CLPP recognizes that the road is occluded by the rock, but BLIP [66] shows a pile of rocks. Therefore, the proposed CLPP has demonstrated more refined and enriched captions over BLIP [66].



Figure 4: Results of captions in images: (a) Car crash. (b) Flooding. (c) Fog. (d) Landslide. (e) Rain. (f) Traffic flow.

471 Table 3 shows the captions containing Ground truth and two captions generated by CLPP and BLIP [66] on five im-472 473 ages in Figure 4. CLPP and BLIP [66] are able to recognize and count on things of the scenes; however, the captions 474 generated by CLPP are semantically closer to the class of 475 ground truths. Moreover, CLPP shows a higher level of de-476 tail in the caption since the reflection of the condition of the 477 road surface is due to the effects of weather and disaster. In 478 479 order to evaluate the classification capability of CLPP for 480 dynamic scenes quantitatively, the similarity between the captions generated by CLPP and the annotated class of the 481 images is evaluated based on the adjective or noun sets, i.e., 482 frozen, flooded, damaged, and light as well as synonyms 483 words occurring in the caption. 484

Table 4 demonstrates a comparison of the classification

 Table 4: Accuracy comparison among ViT, Resnet, Vgg18, and Dvlc. The

 bold font indicates the best score.

 487

Classes/model	Dvlc (%)	ViT (%)	Resnet101(%)	Vgg19(%)
Car crashes	93.37	92.41	91.12	87.67
Flooding	90.69	89.23	87.83	86.54
Fog	92.98	91.19	86.77	85.23
Landslide	89.52	87.63	87.19	84.89
Rain	92.33	87.58	88.92	83.11
Average	91.78	89.61	88.37	85.49

results among the proposed Dvlc, CLPP, and SOTA clas-494 sifiers, i.e., ViT [27], VGG19, and Resnet101. Dvlc has495 classified best in the accuracy of 91.78% accuracy without496 retraining. Therefore, the effectiveness of Dvlc as prompt497 engineering facilitating to Dvls has been proven. 498

#### 4.3. Refined Semantic Segmentation by Prompt En-500 gineering 501

This section denotes the proposed Dvls and how to obtain the final refined captions using prompt engineering. The prompt for each scene is pre-defined as a list of words, i.e., (1) car crashes: ["pedestrian", "car", "car crash", "road", "bike", "tree"]; (2) flooding: ["water", "car", "person", "tree", "sky"]; (3) fog: ["foggy", "mountain", "road", "car", "wet"]; (4) landslide: ["landslide", "debris flow", "rocks", "road", "dirt"]; (5)rain: ["water", "rain", "umbrella", "road", "person"]. Prompts are selected respectively by classification results from Dvlc.

Figure 5 illustrates the effectiveness of our approach on 513 images with foggy and traffic accident scenes. (a) shows 514 the input images, while (c) displays the segmentation results generated by the transformer-based SOTA segmentation model, i.e., Mask2former [19], which shows generic 517 classes, i.e., "sky-other-merged", and "car". (b) presents 518 improved segmentation results, and achieved prompt engineering, which provides more detailed semantic segmentation results, i.e., more detail from "sky-other-merged", 521 to "foggy" for the foggy scene and from "car" to "car crash" 522 for the traffic accident scene.

It has been demonstrated that prompt tuning for  $Dvlc_{524}^{524}$  is helpful for detailing segmentation results under dynamic 525 conditions. 526

527

528

#### 4.4. Image Captioning by CLPP

This section challenges image captioning with obsta-529 cles on the road using the proposed Contrastive Language530 Physical-Scale Pretraining (CLPP) (m = 4 in eq. 2) in a 3D531 manner, unlike GraphVQA [1] with the 2D graph of on, left,532 or right. Figure 6 (a) shows post-disaster images suffered533 from the enormous typhoon, where many obstacles like dirt534 and rocks piled up on the road and other regions. In or-535 der to recognize whether obstacles are present on the road536 or not, the occluded road surfaces must first be identified.537 For this, DeepDepth (b) with vertical (Oy), horizontal (Ox),538 and depth (Oz) coordinates are used to recognize nearly flat539

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

608

645



Figure 5: Results of segmentation by SOTA and proposed PanopticCAP: (a) Original image. (b) Proposed refined semantic segmentation. (c) Mask2Former [19]

road surfaces that are assumed to be the normal road (xOz) below the obstacles.

On the other hand, Dscene (c) can provide segmented objects. Depth maps (b) and objects from (c) are fed into the image encoder of a VLM to make contrastive with description texts, respectively. The presentation of the depth map in 3D (d) shows the relative location of objects. The location of objects like tree and mountain classes are included in the caption for the image when encoding depth map with objects mask respectively with relative location by adverbs in texts, i.e., "on", "before", "above", "bottom", and "side" for 2D and adverbs for 3D are "behind," "in front of", and "rear".

The refined caption and segmentation (e) show that the proposed method can enrich the caption by adding detailed locations of objects and road conditions. SOTA single VLMs, i.e., BLIP [66], could not estimate the depth and road condition. It has been proven that CLPP can estimate physical factors on images, and CLPP caption result is more detailed than single VLM.



Figure 6: Proposed CLPP applied to post-disaster scenes to identify road regions with various obstacles (dirt, water, rock): (a) Input image. (b) Depth map. (c) Panoptic segmentation. (d) Refined road surface. (e) The left image is the input image with the BLIP caption result. The right image is the refined caption: "Rock on the dry road".

bie 5:	Com	parison o	of the refined	captions	with E	SLIP Ca	iptior	results.	55
									50
	1	1	D I I I	1	DI	ID [CC]			23

	Proposed method	BLIP [66]	
(1)	Rocks lay on the flooding road, within 938m in visibility	A flooded road in the rain	596
(2)	Rock debris lay on the wet road, within clear visibility	A road in the rain with rocks and debris on the side	597 598
(3)	15 vehicles on the wet highway, under heavy snowfall and within 637m in visibility	A snowstorm on a highway	599
(4)	A truck on the wet highway, snow on the side of the highway, under heavy snowfall and within 512m in visibility	A snow plow clears a road in the snow	600 601
(5)	12 people stand on a flooded road, within 812m in visibility, and 0.5m water level (Lv2)	A group of people on flooded road	602
	The highway under light snowfall	Snow-covered road	603
(6)	with the snow on the side of the road, within 748m in visibility	with a fence and a street light	604

4.5. Dynamic Captions with Weather and Road<sup>605</sup> Conditions by Proposed Dvis, Dsnow, Droad,<sup>606</sup> and Dwater 607

To provide further complicated captions, this section<sub>609</sub> conducts experiments on various weather conditions with traffic and disaster scenes. The proposed Dsnow, Droad, 611 and Dvis are used by comparing a SOTA VL captioning model, BLIP [66]. Figure 7 shows six scenes. As a result, 613 road conditions by Droad (1)-(6) are wet in blue and snow 614in yellow. Dsnow's indicators (3)-(6) present light to heavy snowfall. Dvlc recognizes overall scene objects like moun-616 tains, rivers, rocks, sky, and trees. Proposed Dvis [2] from<sub>617</sub> images (1)-(6) can estimate a physical scale from weather<sub>618</sub> phenomenon, i.e., foggy visibility distances in meters: 938,619 clear, 637, 512, 812, and 748 m, respectively. Therefore, the road condition and visibility distance have been included in  $\frac{1}{621}$ the captions of Dvlc. Table 4 shows a comparison of the re- $\frac{1}{622}$ fined captions and a SOTA BLIP [66] result using six scenes of Figure 7. The comparison results show that a refined cap- $_{624}$ tion is detailed about the scene by adding road conditions, 625 snowfall status, location of objects, and exact visibility in meters. Besides, the caption from BLIP lacks a description. The result has proven that the proposed method integrating Droad, Dsnow, and Dvis outperforms single VML, i.e., 629 BLIP [66] 630



Figure 7: Results of proposed Dvls with refined and enriched captions in 641 dynamic scenes: (1) Flooding road. (2) Landslide on the road. (3), (4) Heavy snowfall on the highway. (5) Flooded scene with water level, Lv2 (6) Light snowfall on the highway. 643 **5** A blottion study. 644

#### 5. Ablation study

To justify the proposed PanopticCAP, many additional646 experiments are ablated below. 647

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699 700

701

# 5.1. Caption Refinement by Dvls

To show the usefulness of refined Dvls, many unseen 650 disaster scenes that have not been pre-trained are used to 651 segment with classes. As shown in Figure 8 (a), images 652 present disaster events. Two SOTAs of (c) MaskDINO [7]. 653 654 (d) OVSeg [76] are compared. As a result, Table 6 summarizes classes of (b) proposed Dvls and (c), (d) two SOTAs. 655 In (1), a track (c) or boat (d) has been annotated, whereas 656 the proposed Dvls has refined to "car crash" over water (b). 657 In (2)-(5), snow to water, landslide to rocks, pavement to 658 rain, and tree to strong wind have been annotated by (b) the 659 proposed Dvls, respectively. 660

Therefore, refined texts from SOTAs' texts could en-661 hance original to higher semantic texts. In particular, (5) 662 tree (c) is normal segmentation, but strong wind (b), (d) 663 stands for intuitive weather conditions as humans may an-664 nounce. When depth maps are added, combined with lo-665 cation prompts, Dvls can label segmented objects more se-666 mantically. Therefore, it has been proven that the proposed 667 Dvls with texts will play an important role in messaging 668 heavy disaster events more clearly than SOTAs' texts. 669



Figure 8: Comparison of the proposed method, MaskDINO [7], and OVSeg [76] (a) Input image. (b) Proposed Dvls. (c) MaskDINO [7] (d) OVSeg [76]

Table 6: Comparison of classes by SOTAs and proposed Dvls.

Image	SOTA	Proposed
(1)	boat, truck	car crash
(2)	snow, rain	water
(3)	rock-merged, rain	landslide
(4)	pavement-merged, rain	rain
(5)	tree-merged, typhoon	strong wind

### 5.2. Caption Enrichment Using Refined Segmentation and Query Loop

This section describes a combination of refined segmentation from the proposed Dvls with a query loop from VQA to further refine the caption for an image. A loop is necessary because the resulting segmentation of Dvls depends on the newly added class. By selecting the most appropriate word, the best segmentation result can be achieved during the looping. Figure 9 depicts a looping scheme between VQA and Dvls. The list of objects/events consists of words with similar meanings that describe an object/event. The<sup>702</sup> process of selecting the best answer involves evaluating the<sup>703</sup> cosine similarity between the answer and the image. The<sup>704</sup> answer with the highest cosine similarity score is consid-<sup>705</sup> ered the best. To ensure the most accurate selection, all<sup>706</sup> words with similar meanings are looped through until the<sup>707</sup> best answer has been chosen. The effectiveness of these in-<sup>708</sup> tegrated models is compared to an online image captioning<sup>709</sup> API, i.e., visual ChatGPT. The queries are designed to cor-<sup>710</sup> respond to different types of events, such as "car crashes,"<sup>711</sup> "flooding," "fog," "landslide," "rain," and "snowfall." Pre-<sup>712</sup> defined prompts for these scenes are used, including cars,<sup>713</sup> water, snow, rocks, and debris flow. Figure 10 presents a<sup>714</sup>



comparison between two approaches: (a) visual ChatGPT722 [123] caption results and (b) finally refined segmentation723 using a query loop with consideration of the relative size724 and location of objects. The comparison reveals that vi-725 sual ChatGPT [123] only provides an overview of scene de-726 scriptions with no physical scales. On the other hand, the727 proposed model presents more detailed physical scales of728 sizes and locations for accident events. Thus, the proposed729 method has proven capable of handling dynamic captions.730 More refined and enriched captions by the proposed model731 have been generated for such traffic accident scenes than732 visual ChatGPT [123]. 733



Figure 10: Comparison of vision-language models between visual Chat-743 GPT and the proposed method: (a) Image captioning using visual Chat-744 GPT. (b) Refined segmentation and captions with physical scales by the745 proposed CaptionCAP. (c) Captions from visual ChatGPT and the pro-746 posed CaptionCAP.

### 5.3. Compare the Proposed PanopticCAP with 748 SOTA VMLs 749

To justify the performance of the proposed Panoptic-750 CAP, more complicated scenes under adversarial weather751 and disaster conditions are selected to provide captions,752 as shown in Figure 11 (a). VLM-based SOTAs, i.e., (d)753 ZegFormer [24], (e) OvSeg [76], and (f) ClipSeg [84], are754 compared. Dwater is used for water level estimation us-755

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779 780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799 800

801

802 803

804

805

806

807 808

809

ing pre-registered known standard objects' sizes, i.e., buildings, cars, poles, and human heights. Results by (f) ClipSeg [84] show the worst segmentation, particularly flooding regions. On the other hand, results by (d) ZegFormer [24] segment are better than results by (e) OvSeg [76]. However, all SOTAs are limited to generating captions with physical scales. Results by (b) the proposed PanopticCAP present segmented objects with bikes, cars, and water. In captions (c) from (b), water levels from 0.3 m to 1.56 m and several cars have been contained. These quantitative captions can describe more details than only qualitative texts, i.e., flooding. Therefore, it has been demonstrated that the proposed PanopticCAP provides dynamic physical scales that are helpful in understanding how cars and bikes are damaged. It is the first time to generate dynamic captions from such disaster scenes. When videos are available, spatiotemporal changes will also be captioned to monitor for rescue purposes.



Figure 11: Results of SOTA VL-based segmentation: (a) Input image. (b) Proposed PanopticCAP. (c) Refined caption by the proposed PanopticCAP with scene descriptions and physical scale in water level. (d) ZegFormer [24] (e) OvSeg [76] (f) ClipSeg [84]

#### 5.4. CLPP with Different Loss Function Parameters

This section presents an experimental comparison of various parameters for the contrastive loss function (L) used in the proposed CLPP. CLPP is applied in Sections 4.2 and 4.4 with m = 0.4 in equations 2 and 3. In this experiment, m is used from array list values [0.2, 0.3, 0.4, 0.5]. A comparison of equation (1) and the proposed equations (2) and (3) is carried out. Table 7 shows a comparison of the modified loss function and the original one for the physical scale-generated caption. The result shows that the performance of CLPP is the lowest in RMSE when m is set to 0.4. Therefore, it has been reconfirmed that the selected m is optimal.

Table 7: RMSE of different values m with/without sim.

m	Modified	Original
0.2	0.1985	0.2214
0.3	0.2043	0.2375
0.4	0.1894	0.2018
0.5	0.1964	0.2145
-	A D	

#### 5.5. Overall Evaluation PanopticCAP

This section presents an experiment that evaluates the final output of all eleven modules. The experiment measures Table 8: Performance evaluation of proposed panopticCAP on public image-text datasets

			010
Datas of Moth of	DemonstiaCAD	Visual	81
Dataset/Methoa	ranopucCAr	ChatGPT	
COCO Caption	0.3854	0.4415	812
Conceptual Caption	0.3659	0.4235	813

 Table 9: Performance evaluation of proposed panopticCAP on collected814 datasets.
 815

	Dataset/Method	PanopticCAP	Visual ChatGPT	816
	Disaster	0.4521	0.3124	01/
	Traffic accident	t 0.4315	0.3254	818
Tabl	e 10: Computation	al cost and memory	usage compari	sons. 819
	Daufaum (Madal	Computational cost	Memory	820
	Perioriii/Wodei	(second)	usage (Mb)	821
	Proposed method	9.423	11231	
	Visual ChatGPT	8.123	6132	822
	BLIP[66]	1.432	3214	823

823 824

performance using the BLEU score and is conducted on825 two datasets. The first dataset is publicly available and in-826 cludes the COCO Caption dataset [16] and the Conceptual827 Captions dataset [107], both of which contain image-text828 pairs. The second dataset includes two images with accom-829 panying text descriptions describing snowfall status, water830 level, and physical scale. These collections are the Disas-831 ter dataset (1850 image-text pairs) and the Traffic accident832 dataset (2130 image-text pairs). 833 According to the results in Table 8, PanopticCAP does not834

perform as well as Visual ChatGPT. This could be due to835 the fact that the text descriptions in the public image set do836 not include information about road conditions, water levels,837 snow conditions, or visibility, whereas PanopticCAP is ca-838 pable of generating captions with these details. However,839 Table 9 presents contradictory results, where PanopticCAP840 outperforms Visual ChatGPT on datasets featuring disaster841 or traffic accident conditions. It has been proven that Panop-842 ticCAP can provide detailed semantics about the physical843 aspects of scenes. These can be highly useful for tasks such844 as traffic coordination and rescue operations. 845

In addition, we compared the computational cost and mem-846 ory usage of the proposed system with that of SOTA meth-847 ods on the same hardware device. Table 10 presents a com-848 parison of the computational cost and memory usage for849 these methods. 850

#### 6. Conclusion

851 852

853

This paper has proposed PanopticCAP with multiple DL854 and VLM models, which consist of branched structures for855 efficiency in light of memory, training, and maintenance. It856 is the first time to contain dynamic changes in captions with857 physical scales, i.e., depth, fog visibility distance, weather858 conditions, water level, and road conditions. A physics-859 based loss function generates more refined and enriched860 captions at a contrastive loss. PanopticCAP will help no-861 tify detailed scene descriptions to drivers, auto-driving, and862 rescue workers from camera images. 863

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

# <sup>864</sup> References

- [1] Anonymous. Panopticroad: Integrated panoptic road segmentation under adversarial conditions. *in CVPR Workshop*, 2023.
- [2] Anonymous. Panopticvis: Integrated panoptic segmentation for visibility estimation at twilight and night. *in CVPR Workshop*, 2023.
- [3] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Exploring visual prompts for adapting large-scale models. *arXiv preprint arXiv:2203.17274*, 1(3):4, 2022.
- [4] Gedas Bertasius and Lorenzo Torresani. Classifying, segmenting, and tracking object instances in video with mask propagation. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, pages 9736–9745. Computer Vision Foundation / IEEE, 2020.
- [5] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. YOLACT: real-time instance segmentation. In 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019, pages 9156–9165. IEEE, 2019.
- [6] Adrian Bulat and Georgios Tzimiropoulos. Language-aware soft prompting for vision & language foundation models. *CoRR*, abs/2210.01115, 2022.
- [7] Christoph Busch and Eric Debes. Wavelet transform for analyzing fog visibility. *IEEE Intell. Syst.*, 13(6):66–71, 1998.
- [8] Anh-Quan Cao and Raoul de Charette. Monoscene: Monocular 3d semantic scene completion. In IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pages 3981–3991. IEEE, 2022.
- [9] Jiale Cao, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shahbaz Khan, Yanwei Pang, and Ling Shao. Sipmask: Spatial information preservation for fast image and video instance segmentation. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XIV, volume 12359 of Lecture Notes in Computer Science, pages 1–18. Springer, 2020.
- [10] Nicolas Carion, Francisco Massa, Gabriel Synnaeve,
  Nicolas Usunier, Alexander Kirillov, and Sergey
  Zagoruyko. End-to-end object detection with transformers. In Andrea Vedaldi, Horst Bischof, Thomas
  Brox, and Jan-Michael Frahm, editors, *Computer*

Vision - ECCV 2020 - 16th European Conference, 918 Glasgow, UK, August 23-28, 2020, Proceedings, 919 Part I, volume 12346 of Lecture Notes in Computer 920 Science, pages 213–229. Springer, 2020. 921

- [11] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, <sup>922</sup> Nicolas Usunier, Alexander Kirillov, and Sergey<sup>923</sup> Zagoruyko. End-to-end object detection with trans-<sup>924</sup> formers. In Andrea Vedaldi, Horst Bischof, Thomas<sup>925</sup> Brox, and Jan-Michael Frahm, editors, *Computer*<sup>926</sup> *Vision - ECCV 2020 - 16th European Conference*, <sup>927</sup> *Glasgow, UK, August 23-28, 2020, Proceedings*, <sup>928</sup> *Part I*, volume 12346 of *Lecture Notes in Computer*<sup>929</sup> *Science*, pages 213–229. Springer, 2020. <sup>930</sup>
- [12] João Carreira, Viorica Patraucean, Laurent Mazaré,<sup>931</sup> Andrew Zisserman, and Simon Osindero. Massively<sup>932</sup> parallel video networks. In Vittorio Ferrari, Martial<sup>933</sup> Hebert, Cristian Sminchisescu, and Yair Weiss, ed-<sup>934</sup> itors, Computer Vision - ECCV 2018 - 15th Euro-<sup>935</sup> pean Conference, Munich, Germany, September 8-<sup>936</sup> 14, 2018, Proceedings, Part IV, volume 11208 of<sup>937</sup> Lecture Notes in Computer Science, pages 680–697.<sup>938</sup> Springer, 2018.
- [13] Feilong Chen, Duzhen Zhang, Minglun Han, Xiu-Yi<sup>940</sup> Chen, Jing Shi, Shuang Xu, and Bo Xu. VLP: A<sup>941</sup> survey on vision-language pre-training. *Int. J. Autom.*<sup>942</sup> *Comput.*, 20(1):38–56, 2023.
- [14] Guangyi Chen, Weiran Yao, Xiangchen Song,944
   Xinyue Li, Yongming Rao, and Kun Zhang. Prompt945
   learning with optimal transport for vision-language946
   models. *CoRR*, abs/2210.01253, 2022. 947
- [15] Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen,948 Yongming Huang, and Youliang Yan. Blendmask:949 Top-down meets bottom-up for instance segmenta-950 tion. In 2020 IEEE/CVF Conference on Computer951 Vision and Pattern Recognition, CVPR 2020, Seat-952 tle, WA, USA, June 13-19, 2020, pages 8570–8578.953 Computer Vision Foundation / IEEE, 2020. 954
- [16] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakr-955 ishna Vedantam, Saurabh Gupta, Piotr Dollár, and 956 C Lawrence Zitnick. Microsoft coco captions: Data 957 collection and evaluation server. arXiv preprint 958 arXiv:1504.00325, 2015. 959
- Y. H. Chen, P.J. Lee, and T.A. Buil. Multi-scales fea-960 ture extraction model for water segmentation in the961 satellite image. *In Proc. IEEE Int. Conf. Consumer*962 *Electornics (ICCE)*, 2023. 963
- [18] Bowen Cheng, Maxwell D. Collins, Yukun Zhu, Ting<sub>964</sub> Liu, Thomas S. Huang, Hartwig Adam, and Liang-<u>965</u> Chieh Chen. Panoptic-deeplab: A simple, strong,<u>966</u> and fast baseline for bottom-up panoptic segmenta-<u>967</u> tion. In 2020 IEEE/CVF Conference on Computer<u>968</u> Vision and Pattern Recognition, CVPR 2020, Seat-<u>969</u> tle, WA, USA, June 13-19, 2020, pages 12472–12482.<u>970</u> Computer Vision Foundation / IEEE, 2020. <u>971</u>

- [19] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Maskedattention mask transformer for universal image segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 1280– 1289. IEEE, 2022.
  [20] Bowen Cheng, Alexander G. Schwing, and Alexan
- [20] Bowen Cheng, Alexander G. Schwing, and Alexan-980 der Kirillov. Per-pixel classification is not all you 981 need for semantic segmentation. In Marc'Aurelio 982 Ranzato, Alina Beygelzimer, Yann N. Dauphin, 983 Percy Liang, and Jennifer Wortman Vaughan, ed-984 itors, Advances in Neural Information Processing 985 Systems 34: Annual Conference on Neural Informa-986 tion Processing Systems 2021, NeurIPS 2021, De-987 cember 6-14, 2021, virtual, pages 17864-17875, 988 2021.
- 989 [21] Marius Cordts, Mohamed Omran, Sebastian Ramos, 990 Timo Rehfeld, Markus Enzweiler, Rodrigo Benen-991 son, Uwe Franke, Stefan Roth, and Bernt Schiele. 992 The cityscapes dataset for semantic urban scene un-993 derstanding. In 2016 IEEE Conference on Computer 994 Vision and Pattern Recognition, CVPR 2016, Las Ve-995 gas, NV, USA, June 27-30, 2016, pages 3213-3223. 996 IEEE Computer Society, 2016.
- 997 [22] Quan Cui, Boyan Zhou, Yu Guo, Weidong Yin, Hao 998 Wu, Osamu Yoshie, and Yubo Chen. Contrastive 999 vision-language pre-training with limited resources. 1000 In Shai Avidan, Gabriel J. Brostow, Moustapha 1001 Cissé, Giovanni Maria Farinella, and Tal Hassner, 1002 editors, Computer Vision - ECCV 2022 - 17th Eu-1003 ropean Conference, Tel Aviv, Israel, October 23-27, 1004 2022, Proceedings, Part XXXVI, volume 13696 of 1005 Lecture Notes in Computer Science, pages 236–253. 1006 Springer, 2022.
- 1007 [23] Mohammad Mahdi Derakhshani, Enrique Sanchez,
  1008 Adrian Bulat, Victor Guilherme Turrisi da Costa,
  1009 Cees G. M. Snoek, Georgios Tzimiropoulos, and
  1010 Brais Martínez. Variational prompt tuning improves
  1011 generalization of vision-language models. *CoRR*,
  1012 abs/2210.02390, 2022.
- 1013 [24] Jian Ding, Nan Xue, Gui-Song Xia, and Dengxin
  1014 Dai. Decoupling zero-shot semantic segmentation.
  1015 In *IEEE/CVF Conference on Computer Vision and*1016 *Pattern Recognition, CVPR 2022, New Orleans, LA,*1017 USA, June 18-24, 2022, pages 11573–11582. IEEE,
  1018 2022.
- [25] Kun Ding, Ying Wang, Pengzhang Liu, Qiang Yu, Haojian Zhang, Shiming Xiang, and Chunhong Pan.
  Prompt tuning with soft context sharing for visionlanguage models. *CoRR*, abs/2208.13474, 2022.
- 1023 [26] Alexey Dosovitskiy, Lucas Beyer, Alexander
  1024 Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
  1025 Thomas Unterthiner, Mostafa Dehghani, Matthias

Minderer, Georg Heigold, Sylvain Gelly, Jakob<sup>1026</sup> Uszkoreit, and Neil Houlsby. An image is worth<sup>1027</sup> 16x16 words: Transformers for image recognition at<sup>1028</sup> scale. In 9th International Conference on Learning<sup>1029</sup> Representations, ICLR 2021, Virtual Event, Austria,<sup>1030</sup> May 3-7, 2021. OpenReview.net, 2021.

- [27] Alexey Dosovitskiy, Lucas Beyer, Alexander<sup>1032</sup> Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,<sup>1033</sup> Thomas Unterthiner, Mostafa Dehghani, Matthias<sup>1034</sup> Minderer, Georg Heigold, Sylvain Gelly, Jakob<sup>1035</sup> Uszkoreit, and Neil Houlsby. An image is worth<sup>1036</sup> 16x16 words: Transformers for image recognition at<sup>1037</sup> scale. In 9th International Conference on Learning<sup>1038</sup> Representations, ICLR 2021, Virtual Event, Austria,<sup>1039</sup> May 3-7, 2021. OpenReview.net, 2021.
- [28] Dawei Du, Yuankai Qi, Hongyang Yu, Yi-Fan Yang, 1041 Kaiwen Duan, Guorong Li, Weigang Zhang, Qing-1042 ming Huang, and Qi Tian. The unmanned aerial 1043 vehicle benchmark: Object detection and tracking. 1044 In Vittorio Ferrari, Martial Hebert, Cristian Smin-1045 chisescu, and Yair Weiss, editors, Computer Vision 1046 - ECCV 2018 - 15th European Conference, Munich, 1047 Germany, September 8-14, 2018, Proceedings, Part 1048 X, volume 11214 of Lecture Notes in Computer Sci-1049 ence, pages 375–391. Springer, 2018.
- [29] Yu Du, Fangyun Wei, Zihe Zhang, Miaojing Shi, Yue<sup>1051</sup>Gao, and Guoqi Li. Learning to prompt for open-1052 vocabulary object detection with vision-language<sup>1053</sup>model. In *IEEE/CVF Conference on Computer Vi*-1054 sion and Pattern Recognition, CVPR 2022, New Or-1055 leans, LA, USA, June 18-24, 2022, pages 14064–1056 14073. IEEE, 2022.
- [30] Chengjian Feng, Yujie Zhong, Zequn Jie, Xiangx-1058 iang Chu, Haibing Ren, Xiaolin Wei, Weidi Xie,1059 and Lin Ma. Promptdet: Towards open-vocabulary1060 detection using uncurated images. In Shai Avi-1061 dan, Gabriel J. Brostow, Moustapha Cissé, Gio-1062 vanni Maria Farinella, and Tal Hassner, editors,1063 *Computer Vision - ECCV 2022 - 17th European Con-*1064 *ference, Tel Aviv, Israel, October 23-27, 2022, Pro-*1065 *ceedings, Part IX*, volume 13669 of *Lecture Notes in*1066 *Computer Science*, pages 701–717. Springer, 2022. 1067
- [31] Jonathan Francis, Nariaki Kitamura, Felix Labelle,1068
   Xiaopeng Lu, Ingrid Navarro, and Jean Oh. Core1069
   challenges in embodied vision-language planning.1070
   *CoRR*, abs/2106.13948, 2021.
- [32] Mingfei Gao, Chen Xing, Juan Carlos Niebles,1072 Junnan Li, Ran Xu, Wenhao Liu, and Caim-1073 ing Xiong. Open vocabulary object detection1074 with pseudo bounding-box labels. In Shai Avi-1075 dan, Gabriel J. Brostow, Moustapha Cissé, Gio-1076 vanni Maria Farinella, and Tal Hassner, editors,1077 Computer Vision - ECCV 2022 - 17th European Con-1078 ference, Tel Aviv, Israel, October 23-27, 2022, Pro-1079

*ceedings, Part X,* volume 13670 of *Lecture Notes in Computer Science*, pages 266–282. Springer, 2022.

[33] Naiyu Gao, Fei He, Jian Jia, Yanhu Shan, Haoyang Zhang, Xin Zhao, and Kaiqi Huang. Panopticdepth: A unified framework for depth-aware panoptic segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 1622–1632. IEEE, 2022.

**CVPR** 

#2

1080

1081

- [34] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *CoRR*, abs/2110.04544, 2021.
- 1094 [35] Yuting Gao, Jinfeng Liu, Zihan Xu, Jun Zhang, Ke
  1095 Li, and Chunhua Shen. Pyramidclip: Hierarchi1096 cal feature alignment for vision-language model pre1097 training. *CoRR*, abs/2204.14095, 2022.
- 1098 [36] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-1099 Yi Lin. Scaling open-vocabulary image segmen-1100 tation with image-level labels. In Shai Avi-1101 dan, Gabriel J. Brostow, Moustapha Cissé, Gio-1102 vanni Maria Farinella, and Tal Hassner, editors, 1103 Computer Vision - ECCV 2022 - 17th European 1104 Conference, Tel Aviv, Israel, October 23-27, 2022, 1105 Proceedings, Part XXXVI, volume 13696 of Lec-1106 ture Notes in Computer Science, pages 540-557. 1107 Springer, 2022.
- [37] Jiaqi Gu, Hyoukjun Kwon, Dilin Wang, Wei Ye, 1108 Meng Li, Yu-Hsin Chen, Liangzhen Lai, Vikas 1109 1110 Chandra, and David Z. Pan. Multi-scale highresolution vision transformer for semantic segmenta-1111 1112 tion. In IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, 1113 LA, USA, June 18-24, 2022, pages 12084–12093. 1114 IEEE, 2022. 1115
- 1116[38] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin1117Cui. Open-vocabulary object detection via vision1118and language knowledge distillation. In The Tenth1119International Conference on Learning Representa-1120tions, ICLR 2022, Virtual Event, April 25-29, 2022.1121OpenReview.net, 2022.
- [39] Chunle Guo, Chongyi Li, Jichang Guo, Chen Change 1122 Loy, Junhui Hou, Sam Kwong, and Runmin Cong. 1123 Zero-reference deep curve estimation for low-light 1124 image enhancement. In 2020 IEEE/CVF Conference 1125 on Computer Vision and Pattern Recognition, CVPR 1126 2020, Seattle, WA, USA, June 13-19, 2020, pages 1127 1777-1786. Computer Vision Foundation / IEEE, 1128 2020. 1129
- [40] Chunle Guo, Qixin Yan, Saeed Anwar, Runmin Cong, Wenqi Ren, and Chongyi Li. Image dehazing transformer with transmission-aware 3d position embedding. In 2022 IEEE/CVF Conference on Com-

*puter Vision and Pattern Recognition (CVPR)*, pages 1134 5802–5810, 2022. CVPR #2

- [41] Zixian Guo, Bowen Dong, Zhilong Ji, Jinfeng Bai, <sup>1136</sup> Yiwen Guo, and Wangmeng Zuo. Texts as images<sup>1137</sup> in prompt tuning for multi-label image recognition. <sup>1138</sup> *CoRR*, abs/2211.12739, 2022. <sup>1139</sup>
- [42] Ziyu Guo, Renrui Zhang, Longtian Qiu, Xi-1140 anzheng Ma, Xupeng Miao, Xuming He, and Bin<sup>1141</sup> Cui. CALIP: zero-shot enhancement of CLIP with<sup>1142</sup> parameter-free attention. *CoRR*, abs/2209.14169,<sup>1143</sup> 2022.
- [43] Xuehai He, Diji Yang, Weixi Feng, Tsu-Jui Fu, Ar-1145 jun R. Akula, Varun Jampani, Pradyumna Narayana, 1146 Sugato Basu, William Yang Wang, and Xin Wang. 1147 CPL: counterfactual prompt learning for vision and 1148 language models. In Yoav Goldberg, Zornitsa 1149 Kozareva, and Yue Zhang, editors, *Proceedings of* 1150 *the 2022 Conference on Empirical Methods in Natu-*1151 *ral Language Processing, EMNLP 2022, Abu Dhabi*, 1152 *United Arab Emirates, December 7-11, 2022*, pages 1153 3407–3418. Association for Computational Linguis-1154 tics, 2022.
- [44] Daniel Hernández, José M. Cecilia, Juan-Carlos1156 Cano, and Carlos T. Calafate. Flood detection using1157 real-time image segmentation from unmanned aerial1158 vehicles on edge-computing platform. *Remote. Sens.*,1159 14(1):223, 2022.
- [45] Jeremy Howard and Sebastian Ruder. Universal lan-1161 guage model fine-tuning for text classification. In1162 Iryna Gurevych and Yusuke Miyao, editors, *Pro*-1163 *ceedings of the 56th Annual Meeting of the Associa*-1164 *tion for Computational Linguistics, ACL 2018, Mel*-1165 *bourne, Australia, July 15-20, 2018, Volume 1: Long*1166 *Papers*, pages 328–339. Association for Computa-1167 tional Linguistics, 2018. 1168
- [46] Lukas Hoyer, Dengxin Dai, Haoran Wang, and 169 Luc Van Gool. MIC: masked image consistency 1170 for context-enhanced domain adaptation. *CoRR*, 1171 abs/2212.01322, 2022.
- [47] Runhui Huang, Yanxin Long, Jianhua Han, Hang<sub>1173</sub> Xu, Xiwen Liang, Chunjing Xu, and Xiaodan<sub>1174</sub> Liang. NLIP: noise-robust language-image pre-1175 training. *CoRR*, abs/2212.07086, 2022.
- [48] Tony Huang, Jack Chu, and Fangyun Wei. Unsu-1177 pervised prompt learning for vision-language mod-1178 els. *CoRR*, abs/2204.03649, 2022.
- [49] Dat Huynh, Jason Kuen, Zhe Lin, Jiuxiang Gu, and 1180 Ehsan Elhamifar. Open-vocabulary instance segmen-1181 tation via robust cross-modal pseudo-labeling. In 1182 *IEEE/CVF Conference on Computer Vision and Pat*-1183 *tern Recognition, CVPR 2022, New Orleans, LA*, 1184 USA, June 18-24, 2022, pages 7010–7021. IEEE, 1185 2022.
- [50] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and

1189

1190

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

Alexei A Efros. Image-to-image translation with conditional adversarial networks. pages 1125–1134, 2017.

- [51] Jitesh Jain, Jiachen Li, MangTik Chiu, Ali Hassani,
  Nikita Orlov, and Humphrey Shi. Oneformer: One
  transformer to rule universal image segmentation. *CoRR*, abs/2211.06220, 2022.
- 1195 [52] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana 1196 Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, 1197 Zhen Li, and Tom Duerig. Scaling up visual and 1198 vision-language representation learning with noisy 1199 text supervision. In Marina Meila and Tong Zhang, 1200 editors, Proceedings of the 38th International Con-1201 ference on Machine Learning, ICML 2021, 18-24 1202 July 2021, Virtual Event, volume 139 of Proceedings 1203 of Machine Learning Research, pages 4904–4916. 1204 PMLR, 2021.
- [53] Jingjing Jiang, Ziyi Liu, and Nanning Zheng. Correlation information bottleneck: Towards adapting pre-trained multimodal models for robust visual question answering, 2023.
- [54] Jonathan Kahana, Niv Cohen, and Yedid Hoshen.
  Improving zero-shot models with label distribution priors. *CoRR*, abs/2212.00784, 2022.
- [55] Muhammad Uzair Khattak, Hanoona Abdul
  Rasheed, Muhammad Maaz, Salman Khan, and
  Fahad Shahbaz Khan. Maple: Multi-modal prompt
  learning. *CoRR*, abs/2210.03117, 2022.
- 1216 [56] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and
  1217 In So Kweon. Video panoptic segmentation. In
  1218 2020 IEEE/CVF Conference on Computer Vision and
  1219 Pattern Recognition, CVPR 2020, Seattle, WA, USA,
  1220 June 13-19, 2020, pages 9856–9865. Computer Vision Foundation / IEEE, 2020.
  - [57] Alexander Kirillov, Kaiming He, Ross B. Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. *CoRR*, abs/1801.00868, 2018.
  - [58] K. J. Kucharik and J. A. Norman. Effect of light and shade on visibility. *J. the Illuminating Engineering Society*, 13(3):117–125, 04 1984.
  - [59] Weicheng Kuo, Yin Cui, Xiuye Gu, A. J. Piergiovanni, and Anelia Angelova. F-VLM: openvocabulary object detection upon frozen vision and language models. *CoRR*, abs/2209.15639, 2022.
    - [60] Hyunmin Lee and Jaesik Park. Instance-wise occlusion and depth orders in natural scenes. pages 21178–21189, 2022.
- [61] Sohyun Lee, Taeyoung Son, and Suha Kwak. FIFO: learning fog-invariant features for foggy scene segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 18889– 18899. IEEE, 2022.
- [62] Boyun Li, Xiao Liu, Peng Hu, Zhongqin Wu,

Jiancheng Lv, and Xi Peng. All-in-one image restoration for unknown corruption. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, *CVPR 2022, New Orleans, LA, USA, June 18-24*, 2022, pages 17431–17441. IEEE, 2022.

- [63] Boyi Li, Kilian Q. Weinberger, Serge J. Belongie, <sup>1247</sup>
   Vladlen Koltun, and René Ranftl. Language-<sup>1248</sup>
   driven semantic segmentation. In *The Tenth Inter*-<sup>1249</sup>
   *national Conference on Learning Representations*, <sup>1250</sup>
   *ICLR 2022, Virtual Event, April 25-29, 2022.* Open-<sup>1251</sup>
   Review.net, 2022. 1252
- [64] Chunyuan Li, Haotian Liu, Liunian Harold Li,<sup>1253</sup> Pengchuan Zhang, Jyoti Aneja, Jianwei Yang, Ping<sup>1254</sup> Jin, Yong Jae Lee, Houdong Hu, Zicheng Liu, and<sup>1255</sup> Jianfeng Gao. ELEVATER: A benchmark and toolkit<sup>1256</sup> for evaluating language-augmented visual models.<sup>1257</sup> *CoRR*, abs/2204.08790, 2022.
- [65] Feng Li, Hao Zhang, Yi-Fan Zhang, Shilong Liu, 1259 Jian Guo, Lionel M. Ni, PengChuan Zhang, and 1260 Lei Zhang. Vision-language intelligence: Tasks, 1261 representation learning, and large models. *CoRR*, 1262 abs/2203.01922, 2022.
- [66] Junnan Li, Dongxu Li, Caiming Xiong, and Steven1264
  C. H. Hoi. BLIP: bootstrapping language-image1265 pre-training for unified vision-language understand-1266 ing and generation. In Kamalika Chaudhuri, Ste-1267 fanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu,1268 and Sivan Sabato, editors, *International Conference*1269 on Machine Learning, ICML 2022, 17-23 July 2022,1270 Baltimore, Maryland, USA, volume 162 of Proceed-1271 ings of Machine Learning Research, pages 12888–1272 12900. PMLR, 2022.
- [67] Junnan Li, Silvio Savarese, and Steven C. H. Hoi.1274 Masked unsupervised self-training for zero-shot im-1275 age classification. *CoRR*, abs/2206.02967, 2022. 1276
- [68] Manling Li, Ruochen Xu, Shuohang Wang, Luowei1277 Zhou, Xudong Lin, Chenguang Zhu, Michael Zeng,1278 Heng Ji, and Shih-Fu Chang. Clip-event: Connecting1279 text and images with event structures. In *Proceed*-1280 *ings of the IEEE/CVF Conference on Computer Vi*-1281 *sion and Pattern Recognition (CVPR)*, pages 16420–1282 16429, June 2022.
- [69] Qin Li, Yi Li, and Bin Xie. Single image-based scene<sub>1284</sub> visibility estimation. *IEEE Access*, 7:24430–24439,1285 2019.
- [70] Ruoteng Li, Robby T. Tan, and Loong-Fah Cheong.<sub>1287</sub> All in one bad weather removal using architectural<sub>1288</sub> search. In 2020 IEEE/CVF Conference on Computer<sub>1289</sub> Vision and Pattern Recognition, CVPR 2020, Seat-<sub>1290</sub> tle, WA, USA, June 13-19, 2020, pages 3172–3182.<sub>1291</sub> Computer Vision Foundation / IEEE, 2020. 1292
- [71] Yi Li, Yi Chang, Yan Gao, Changfeng Yu, and 1293 Luxin Yan. Physically disentangled intra- and inter-1294 domain adaptation for varicolored haze removal. In 1295

1296IEEE/CVF Conference on Computer Vision and Pat-1297tern Recognition, CVPR 2022, New Orleans, LA,1298USA, June 18-24, 2022, pages 5831–5840. IEEE,12992022.

- [72] Yanwei Li, Xinze Chen, Zheng Zhu, Lingxi Xie, Guan Huang, Dalong Du, and Xingang Wang.
  Attention-guided unified network for panoptic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 7026–7035. Computer Vision Foundation / IEEE, 2019.
- 1307 [73] Yangguang Li, Feng Liang, Lichen Zhao, Yufeng 1308 Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, 1309 and Junjie Yan. Supervision exists everywhere: 1310 A data efficient contrastive language-image pre-1311 training paradigm. In The Tenth International Con-1312 ference on Learning Representations, ICLR 2022, 1313 Virtual Event, April 25-29, 2022. OpenReview.net, 1314 2022.
- 1315 [74] Yu Li, Shaodi You, Michael S. Brown, and Robby T.
  1316 Tan. Haze visibility enhancement: A survey and quantitative benchmarking. *Computer Vision and Im-age Understanding*, 165:1–16, 2017.
- [75] Yanwei Li, Hengshuang Zhao, Xiaojuan Qi, Liwei
  Wang, Zeming Li, Jian Sun, and Jiaya Jia. Fully convolutional networks for panoptic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021,*pages 214–223. Computer Vision Foundation / IEEE,
  2021.
- [76] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li,
  Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted CLIP. *CoRR*,
  abs/2210.04150, 2022.
- 1331[77]Chuang Lin, Peize Sun, Yi Jiang, Ping Luo, Lizhen1332Qu, Gholamreza Haffari, Zehuan Yuan, and Jianfei1333Cai. Learning object-language alignments for open-1334vocabulary object detection. *CoRR*, abs/2211.14843,13352022.
- [78] Yuqi Lin, Minghao Chen, Wenxiao Wang, Boxi Wu,
  Ke Li, Binbin Lin, Haifeng Liu, and Xiaofei He.
  CLIP is also an efficient segmenter: A text-driven approach for weakly supervised semantic segmentation. *CoRR*, abs/2212.09506, 2022.
- [79] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In 2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021, pages 9992–10002. IEEE, 2021.
- [80] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei,
  Zheng Zhang, Stephen Lin, and Baining Guo. Swin

transformer: Hierarchical vision transformer using <sup>1350</sup> shifted windows. *CoRR*, abs/2103.14030, 2021.

- [81] Jonathan Long, Evan Shelhamer, and Trevor Darrell. 1352
   Fully convolutional networks for semantic segmen-1353
   tation. In *IEEE Conference on Computer Vision and* 1354
   *Pattern Recognition, CVPR 2015, Boston, MA, USA*, 1355
   *June 7-12, 2015*, pages 3431–3440. IEEE Computer 1356
   Society, 2015.
- [82] Yanxin Long, Jianhua Han, Runhui Huang, Xu<sup>1358</sup> Hang, Yi Zhu, Chunjing Xu, and Xiaodan Liang.<sup>1359</sup> P<sup>3</sup>ovd: Fine-grained visual-text prompt-driven self-<sup>1360</sup> training for open-vocabulary object detection. *CoRR*,<sup>1361</sup> abs/2211.00849, 2022.
- [83] Yuning Lu, Jianzhuang Liu, Yonggang Zhang, Ya-1363 jing Liu, and Xinmei Tian. Prompt distribution learn-1364 ing. In *IEEE/CVF Conference on Computer Vision*1365 and Pattern Recognition, CVPR 2022, New Orleans, 1366 LA, USA, June 18-24, 2022, pages 5196–5205. IEEE, 1367 2022.
- [84] Timo Lüddecke and Alexander S. Ecker. Image<sup>1369</sup> segmentation using text and image prompts. In<sup>1370</sup> *IEEE/CVF Conference on Computer Vision and Pat-*<sup>1371</sup> *tern Recognition, CVPR 2022, New Orleans, LA*, 1372 USA, June 18-24, 2022, pages 7076–7086. IEEE, 1373 2022.
- [85] Huaishao Luo, Junwei Bao, Youzheng Wu, Xi-1375 aodong He, and Tianrui Li. Segclip: Patch aggrega-1376 tion with learnable centers for open-vocabulary se-1377 mantic segmentation. *CoRR*, abs/2211.14813, 2022. 1378
- [86] Chengcheng Ma, Yang Liu, Jiankang Deng, Lingxi1379 Xie, Weiming Dong, and Changsheng Xu. Under-1380 standing and mitigating overfitting in prompt tuning1381 for vision-language models. *CoRR*, abs/2211.02219,1382 2022.
- [87] Chaofan Ma, Yuhuan Yang, Yan-Feng Wang, Ya1384 Zhang, and Weidi Xie. Open-vocabulary seman-1385 tic segmentation with frozen vision-language mod-1386 els. In 33rd British Machine Vision Conference 2022, 1387 BMVC 2022, London, UK, November 21-24, 2022,1388 page 45. BMVA Press, 2022.
- [88] Long Ma, Tengyu Ma, Risheng Liu, Xin Fan, and 1390 Zhongxuan Luo. Toward fast, flexible, and robust 1391 low-light image enhancement. In 2022 IEEE/CVF 1392 Conference on Computer Vision and Pattern Recog-1393 nition (CVPR), pages 5627–5636, 2022. 1394
- [89] Xianzheng Ma, Zhixiang Wang, Yacheng Zhan, Yin-1395 qiang Zheng, Zheng Wang, Dengxin Dai, and Chia-1396 Wen Lin. Both style and fog matter: Cumulative1397 domain adaptation for semantic foggy scene under-1398 standing. In *IEEE/CVF Conference on Computer*1399 Vision and Pattern Recognition, CVPR 2022, New1400 Orleans, LA, USA, June 18-24, 2022, pages 18900–1401 18909. IEEE, 2022.
- [90] Zongyang Ma, Guan Luo, Jin Gao, Liang Li, Yuxin<sub>1403</sub>

1404Chen, Shaoru Wang, Congxuan Zhang, and Weim-1405ing Hu. Open-vocabulary one-stage detection with1406hierarchical visual-language knowledge distillation.1407In IEEE/CVF Conference on Computer Vision and1408Pattern Recognition, CVPR 2022, New Orleans, LA,1409USA, June 18-24, 2022, pages 14054–14063. IEEE,14102022.

- [91] Antoine Miech, Jean-Baptiste Alayrac, Ivan Laptev, Josef Sivic, and Andrew Zisserman. Rareact: A video dataset of unusual interactions. *CoRR*, abs/2008.01018, 2020.
- [92] Matthias Minderer, Alexey A. Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil Houlsby. Simple open-vocabulary object detection with vision transformers. *CoRR*, abs/2205.06230, 2022.
- Yongguang Mo, Jianjun Huang, and Gongbin Qian.
  Deep learning approach to UAV detection and classification by using compressively sensed RF signal. *Sensors*, 22(8):3072, 2022.
- 1426 [94] Mihai Negru and Sergiu Nedevschi. Image based fog detection and visibility estimation for driving assistance systems. In 2013 IEEE 9th International Conference on Intelligent Computer Communication and Processing (ICCP), pages 163–168, 2013.
- [95] Lucas Prado Osco, José Marcato Junior, Ana Paula Marques Ramos, Lúcio André de Castro Jorge, Sarah Narges Fatholahi, Jonathan de Andrade Silva, Edson Takashi Matsubara, Hemerson Pistori, Wesley Nunes Gonçalves, and Jonathan Li. A review on deep learning in UAV remote sensing. *Int. J. Appl. Earth Obs. Geoinformation*, 102:102456, 2021.
- [96] Omiros Pantazis, Gabriel J. Brostow, Kate E.
  Jones, and Oisin Mac Aodha. Svl-adapter: Self-supervised adapter for vision-language pretrained models. In *33rd British Machine Vision Conference* 2022, *BMVC 2022, London, UK, November 21-24,* 2022, page 580. BMVA Press, 2022.
- [97] Fang Peng, Xiaoshan Yang, and Changsheng Xu.
  Sgva-clip: Semantic-guided visual adapting of
  vision-language models for few-shot image classification. *CoRR*, abs/2211.16191, 2022.
- [98] Ruijie Quan, Xin Yu, Yuanzhi Liang, and Yi Yang. Removing raindrops and rain streaks in one go. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021,* pages 9147–9156. Computer Vision Foundation / IEEE, 2021.
- [99] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models

from natural language supervision. In *International*<sup>1458</sup> *conference on machine learning*, pages 8748–8763.<sup>1459</sup> PMLR, 2021.<sup>1460</sup>

- [100] Alec Radford, Jong Wook Kim, Chris Hallacy,<sup>1461</sup> Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,<sup>1462</sup> Girish Sastry, Amanda Askell, Pamela Mishkin, Jack<sup>1463</sup> Clark, Gretchen Krueger, and Ilya Sutskever. Learn-<sup>1464</sup> ing transferable visual models from natural language<sup>1465</sup> supervision. In Marina Meila and Tong Zhang, ed-<sup>1466</sup> itors, *Proceedings of the 38th International Con*-<sup>1467</sup> *ference on Machine Learning, ICML 2021, 18-24*<sup>1468</sup> *July 2021, Virtual Event*, volume 139 of *Proceedings*<sup>1469</sup> *of Machine Learning Research*, pages 8748–8763.<sup>1470</sup> PMLR, 2021.
- [101] Kanchana Ranasinghe, Brandon McKinzie, Sachin<sup>1472</sup> Ravi, Yinfei Yang, Alexander Toshev, and Jonathon<sup>1473</sup> Shlens. Perceptual grouping in vision-language mod-<sup>1474</sup> els. *CoRR*, abs/2210.09996, 2022.
- [102] Hanoona Abdul Rasheed, Muhammad Maaz,1476 Muhammad Uzair Khattak, Salman H. Khan, and1477 Fahad Shahbaz Khan. Bridging the gap between1478 object and image-level representations for open-1479 vocabulary detection. *CoRR*, abs/2207.03482,1480 2022.
- [103] Christos Sakaridis, Dengxin Dai, and Luc Van Gool.1482 Semantic foggy scene understanding with synthetic1483 data. *Int. J. Comput. Vis.*, 126(9):973–992, 2018. 1484
- [104] Christos Sakaridis, Dengxin Dai, Simon Hecker, and 1485 Luc Van Gool. Model adaptation with synthetic and 1486 real data for semantic dense foggy scene understand-1487 ing. In Vittorio Ferrari, Martial Hebert, Cristian 1488 Sminchisescu, and Yair Weiss, editors, Computer Vi-1489 sion - ECCV 2018 - 15th European Conference, Mu-1490 nich, Germany, September 8-14, 2018, Proceedings, 1491 Part XIII, volume 11217 of Lecture Notes in Com-1492 puter Science, pages 707–724. Springer, 2018. 1493
- [105] Aditya Sanghi, Hang Chu, Joseph G. Lambourne,1494 Ye Wang, Chin-Yi Cheng, Marco Fumero, and Ka-1495 mal Rahimi Malekshan. Clip-forge: Towards zero-1496 shot text-to-shape generation. In *Proceedings of the*1497 *IEEE/CVF Conference on Computer Vision and Pat-*1498 *tern Recognition (CVPR)*, pages 18603–18613, June1499 2022. 1500
- [106] Aashish Sharma, Loong-Fah Cheong, Lionel Heng, 1501 and Robby T. Tan. Nighttime stereo depth estima-1502 tion using joint translation-stereo learning: Light ef-1503 fects and uninformative regions. In 2020 Interna-1504 tional Conference on 3D Vision (3DV), pages 23–31, 1505 2020.
- [107] Piyush Sharma, Nan Ding, Sebastian Goodman, and<sub>1507</sub> Radu Soricut. Conceptual captions: A cleaned, hy-1508 pernymed, image alt-text dataset for automatic image<sub>1509</sub> captioning. In Iryna Gurevych and Yusuke Miyao,<sub>1510</sub> editors, *Proceedings of the 56th Annual Meeting of*<sub>1511</sub>

1512the Association for Computational Linguistics, ACL15132018, Melbourne, Australia, July 15-20, 2018, Vol-1514ume 1: Long Papers, pages 2556–2565. Association1515for Computational Linguistics, 2018.

- [108] Sheng Shen, Shijia Yang, Tianjun Zhang, Bohan
  [108] Sheng Shen, Shijia Yang, Tianjun Zhang, Bohan
  Zhai, Joseph E. Gonzalez, Kurt Keutzer, and Trevor
  Darrell. Multitask vision-language prompt tuning. *CoRR*, abs/2211.11720, 2022.
- [109] Hengcan Shi, Munawar Hayat, Yicheng Wu, and Jianfei Cai. Proposalclip: Unsupervised opencategory object proposal generation via exploiting clip cues. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), pages 9611–9620, June 2022.
- [110] Gyungin Shin, Weidi Xie, and Samuel Albanie.
  Reco: Retrieve and co-segment for zero-shot transfer. *CoRR*, abs/2206.07045, 2022.
- [111] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu,
  Tom Goldstein, Anima Anandkumar, and Chaowei
  Xiao. Test-time prompt tuning for zero-shot generalization in vision-language models. *CoRR*,
  abs/2209.07511, 2022.
- [112] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9,* 2015, Conference Track Proceedings, 2015.
- [113] S. Sreelakshmi and S.S. V. Chandra. Machine learning for disaster management: insights from past research and future implications. *In Proc. IEEE Int. Conf. Comput. Communication, Security, and Intelligent Systems (IC3SIS)*, 2022.
- 1545[114] Ximeng Sun, Ping Hu, and Kate Saenko. Dualcoop:1546Fast adaptation to multi-label recognition with lim-1547ited annotations. CoRR, abs/2206.09541, 2022.
- [115] A. Tafreshi and M. Shahraeeni. Effect of atmospheric haze on visibility distance. *J. Environmental Science and Health*, 44(1):44–48, 04 2009.
- [116] Mingxing Tan and Quoc V. Le. Efficientnet: Re-1551 thinking model scaling for convolutional neural net-1552 works. In Kamalika Chaudhuri and Ruslan Salakhut-1553 dinov, editors, Proceedings of the 36th International 1554 Conference on Machine Learning, ICML 2019, 9-1555 15 June 2019, Long Beach, California, USA, vol-1556 ume 97 of Proceedings of Machine Learning Re-1557 search, pages 6105-6114. PMLR, 2019. 1558
- [117] Vishaal Udandarao, Ankush Gupta, and Samuel Albanie. Sus-x: Training-free name-only transfer of vision-language models. *CoRR*, abs/2211.16198, 2022.
  [118] Oin Wang, Dengrin Dei, Lalaga Haran, Lug Yang, Sungaran, Sungaran,
- [118] Qin Wang, Dengxin Dai, Lukas Hoyer, Luc Van
   [164] Gool, and Olga Fink. Domain adaptive semantic segmentation with self-supervised depth estimation. In

2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, 1567 October 10-17, 2021, pages 8495–8505. IEEE, 2021.

- [119] Tao Wang and Nan Li. Learning to detect and seg-<sup>1569</sup> ment for open vocabulary object detection. *CoRR*, <sup>1570</sup> abs/2212.12130, 2022.
- [120] Sanghyun Woo, Dahun Kim, Joon-Young Lee, and<sup>1572</sup> In So Kweon. Learning to associate every segment<sup>1573</sup> for video panoptic segmentation. In *IEEE Confer-*1574 *ence on Computer Vision and Pattern Recognition*,<sup>1575</sup> *CVPR 2021, virtual, June 19-25, 2021*, pages 2705–1576 2714. Computer Vision Foundation / IEEE, 2021. 1577
- [121] Mitchell Wortsman, Gabriel Ilharco, Jong Wook1578 Kim, Mike Li, Simon Kornblith, Rebecca Roelofs,1579 Raphael Gontijo Lopes, Hannaneh Hajishirzi,1580 Ali Farhadi, Hongseok Namkoong, and Ludwig1581 Schmidt. Robust fine-tuning of zero-shot models.1582 In *IEEE/CVF Conference on Computer Vision and*1583 *Pattern Recognition, CVPR 2022, New Orleans, LA*,1584 USA, June 18-24, 2022, pages 7949–7961. IEEE,1585 2022.
- Bichen Wu, Ruizhe Cheng, Peizhao Zhang, Tianren1587 Gao, Joseph E. Gonzalez, and Peter Vajda. Data effi-1588 cient language-supervised zero-shot recognition with1589 optimal transport distillation. In *The Tenth Inter*-1590 *national Conference on Learning Representations*,1591 *ICLR 2022, Virtual Event, April 25-29, 2022.* Open-1592 Review.net, 2022.
- [123] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong1594
   Wang, Zecheng Tang, and Nan Duan. Visual chatgpt:1595
   Talking, drawing and editing with visual foundation1596
   models. *CoRR*, abs/2303.04671, 2023.
- [124] Xinyi Wu, Zhenyao Wu, Hao Guo, Lili Ju, and Song1598
   Wang. Dannet: A one-stage domain adaptation net-1599
   work for unsupervised nighttime semantic segmenta-1600
   tion. *CoRR*, abs/2104.10834, 2021.
- [125] Deqiang Xiang, Xin Zhang, Wei Wu, and Hongbin<sub>1602</sub> Liu. Denseppmunet-a: A robust deep learning net-1603 work for segmenting water bodies from aerial im-1604 ages. *IEEE Trans. Geosci. Remote. Sens.*, 61:1–11,1605 2023.
- [126] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anand-1607 kumar, Jose M. Alvarez, and Ping Luo. Segformer: 1608 Simple and efficient design for semantic segmenta-1609 tion with transformers. In Marc'Aurelio Ranzato, 1610 Alina Beygelzimer, Yann N. Dauphin, Percy Liang, 1611 and Jennifer Wortman Vaughan, editors, Advances 1612 in Neural Information Processing Systems 34: An-1613 nual Conference on Neural Information Processing 1614 Systems 2021, NeurIPS 2021, December 6-14, 2021, 1615 virtual, pages 12077–12090, 2021.
- [127] Jinheng Xie, Xianxu Hou, Kai Ye, and Linlin Shen.
   [127] CLIMS: cross language image matching for weakly
   [1618] supervised semantic segmentation. In *IEEE/CVF*

1620Conference on Computer Vision and Pattern Recog-<br/>nition, CVPR 2022, New Orleans, LA, USA, June 18-<br/>24, 2022, pages 4473–4482. IEEE, 2022.

- [128] Johnathan Xie and Shuai Zheng. ZSD-YOLO: zero-shot YOLO detection using vision-language knowl-edgedistillation. *CoRR*, abs/2109.12066, 2021.
- [129] Yinghui Xing, Qirui Wu, De Cheng, Shizhou Zhang, Guoqiang Liang, and Yanning Zhang. Class-aware visual prompt tuning for vision-language pre-trained model. *arXiv preprint arXiv:2208.08340*, 2022.
- [130] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin
  Byeon, Thomas M. Breuel, Jan Kautz, and Xiaolong
  Wang. Groupvit: Semantic segmentation emerges
  from text supervision. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*2022, New Orleans, LA, USA, June 18-24, 2022,
  pages 18113–18123. IEEE, 2022.
- 1637 [131] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong 1638 Lin, Yue Cao, Han Hu, and Xiang Bai. A sim-1639 ple baseline for open-vocabulary semantic segmen-1640 tation with pre-trained vision-language model. In 1641 Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, 1642 Giovanni Maria Farinella, and Tal Hassner, editors, Computer Vision - ECCV 2022 - 17th Euro-1643 1644 pean Conference, Tel Aviv, Israel, October 23-27, 1645 2022, Proceedings, Part XXIX, volume 13689 of 1646 Lecture Notes in Computer Science, pages 736–753. 1647 Springer, 2022.
- [132] Wending Yan, Aashish Sharma, and Robby T. Tan.
  Optical flow in dense foggy scenes using semisupervised learning. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition
  (CVPR), pages 13256–13265, 2020.
- [133] Li Yang, Radu Muresan, Arafat Al-Dweik, and Leontios J. Hadjileontiadis. Image-based visibility estimation algorithm for intelligent transportation systems. *IEEE Access*, 6:76728–76740, 2018.
- [134] Wenhan Yang, Robby T. Tan, Shiqi Wang, Yuming Fang, and Jiaying Liu. Single image deraining: From model-based to data-driven and beyond. *CoRR*, abs/1912.07150, 2019.
- [135] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19313–19322, June 2022.
- [136] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, Ce Liu, Mengchen Liu, Zicheng Liu, Yumao Lu, Yu Shi, Lijuan Wang, Jianfeng Wang, Bin Xiao, Zhen Xiao, Jianwei Yang, Michael Zeng, Luowei Zhou, and Pengchuan Zhang. Florence: A new foundation

model for computer vision. *CoRR*, abs/2111.11432, 1674 2021. 1675

- [137] Nir Zabari and Yedid Hoshen. Semantic segmenta-<sup>1676</sup> tion in-the-wild without seeing any segmentation ex-<sup>1677</sup> amples. *CoRR*, abs/2112.03185, 2021.
- [138] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen<sup>1679</sup> Huang, and Chen Change Loy. Open-vocabulary<sup>1680</sup> DETR with conditional matching. In Shai Avi-<sup>1681</sup> dan, Gabriel J. Brostow, Moustapha Cissé, Gio-<sup>1682</sup> vanni Maria Farinella, and Tal Hassner, editors,<sup>1683</sup> *Computer Vision - ECCV 2022 - 17th European Con*-<sup>1684</sup> *ference, Tel Aviv, Israel, October 23-27, 2022, Pro*-<sup>1685</sup> *ceedings, Part IX,* volume 13669 of *Lecture Notes in*<sup>1686</sup> *Computer Science*, pages 106–122. Springer, 2022. <sup>1687</sup>
- [139] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang,<sup>1688</sup>
   and Chen Change Loy. Unified vision and language<sup>1689</sup>
   prompt learning. *CoRR*, abs/2210.07225, 2022.
- [140] Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, 1691
   Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng 1692
   Li. Tip-adapter: Training-free clip-adapter for better 1693
   vision-language modeling. *CoRR*, abs/2111.03930, 1694
   2021. 1695
- [141] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, 1696
   Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and 1697
   Hongsheng Li. Pointclip: Point cloud understand-1698
   ing by clip. In *Proceedings of the IEEE/CVF Con*-1699
   *ference on Computer Vision and Pattern Recognition*1700
   (CVPR), pages 8552–8562, June 2022.
- [142] Renrui Zhang, Longtian Qiu, Wei Zhang, and Ziyao1702 Zeng. VT-CLIP: enhancing vision-language mod-1703 els with visual-guided texts. *CoRR*, abs/2112.02399,1704 2021.
- [143] Shiyu Zhao, Zhixing Zhang, Samuel Schulter, Long<sub>1706</sub> Zhao, B. G. Vijay Kumar, Anastasis Stathopoulos,<sub>1707</sub> Manmohan Chandraker, and Dimitris N. Metaxas.<sub>1708</sub> Exploiting unlabeled data with vision and lan-<sub>1709</sub> guage models for object detection. In Shai Avi-<sub>1710</sub> dan, Gabriel J. Brostow, Moustapha Cissé, Gio-<sub>1711</sub> vanni Maria Farinella, and Tal Hassner, editors,<sub>1712</sub> *Computer Vision - ECCV 2022 - 17th European Con*-<sub>1713</sub> *ference, Tel Aviv, Israel, October 23-27, 2022, Pro*-<sub>1714</sub> *ceedings, Part IX*, volume 13669 of *Lecture Notes in*<sub>1715</sub> *Computer Science*, pages 159–175. Springer, 2022. 1716
- [144] Chong Zhou, Chen Change Loy, and Bo Dai. Ex-1717 tract free dense labels from CLIP. In Shai Avi-1718 dan, Gabriel J. Brostow, Moustapha Cissé, Gio-1719 vanni Maria Farinella, and Tal Hassner, editors, 1720 Computer Vision - ECCV 2022 - 17th European1721 Conference, Tel Aviv, Israel, October 23-27, 2022, 1722 Proceedings, Part XXVIII, volume 13688 of Lec-1723 ture Notes in Computer Science, pages 696–712.1724 Springer, 2022.
- [145] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from CLIP. In Shai Avi-

dan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXVIII, volume 13688 of Lecture Notes in Computer Science, pages 696–712. Springer, 2022.*

- [146] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *IEEE/CVF Conference* on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pages 16795–16804. IEEE, 2022.
- [147] Kaiyang Zhou, Jingkang Yang, Chen Change Loy,
  and Ziwei Liu. Learning to prompt for visionlanguage models. *Int. J. Comput. Vis.*, 130(9):2337–
  2348, 2022.
- [148] Wangchunshu Zhou, Yan Zeng, Shizhe Diao, and
  Xinsong Zhang. VLUE: A multi-task benchmark
  for evaluating vision-language models. *CoRR*,
  abs/2205.15237, 2022.
- [149] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, Computer Vision - ECCV 2022 - 17th European Con-ference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part IX, volume 13669 of Lecture Notes in Computer Science, pages 350–368. Springer, 2022.
- [150] Ziqin Zhou, Bowen Zhang, Yinjie Lei, Lingqiao
  Liu, and Yifan Liu. Zegclip: Towards adapting
  CLIP for zero-shot semantic segmentation. *CoRR*,
  abs/2212.03588, 2022.
- [151] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu,
  and Hanwang Zhang. Prompt-aligned gradient for
  prompt tuning. *CoRR*, abs/2205.14865, 2022.
- [152] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks.
  In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017,* pages 2242–2251. IEEE Computer Society, 2017.
- [153] Xiangyang Zhu, Renrui Zhang, Bowei He, Ziyao
  Zeng, Shanghang Zhang, and Peng Gao. Pointclip V2: adapting CLIP for powerful 3d open-world learning. *CoRR*, abs/2211.11682, 2022.