

# RareCrafter: Controllable Generative Augmentation for Rare Object Detection in Driving Scenes

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Reliable perception of roadway environments is essential*  
002 *for intelligent transportation systems (ITS) and their digital*  
003 *twin frameworks, where object detectors provide a key in-*  
004 *terface between the physical world and its digital represen-*  
005 *tation. However, existing driving datasets are dominated by*  
006 *common traffic participants, leaving rare but safety-critical*  
007 *objects underrepresented and limiting robustness to uncom-*  
008 *mon roadway events. We propose RareCrafter, a training-*  
009 *free generative augmentation framework that synthesizes*  
010 *rare objects directly into real driving scenes using control-*  
011 *lable generative models. The framework combines struc-*  
012 *tured prompt diversification, depth-conditioned empirical*  
013 *size priors, and region-conditioned verification to ensure*  
014 *that inserted objects remain semantically consistent and*  
015 *maintain realistic, depth-dependent object scales. The gen-*  
016 *erated images are then used to augment training data for*  
017 *object detection. Experiments on the CODA2022 corner-*  
018 *case dataset show that RareCrafter improves rare-object*  
019 *detection by 2.82 mAP on average over real-only training*  
020 *and consistently outperforms a CopyPaste augmentation*  
021 *baseline. These results demonstrate that controllable gen-*  
022 *erative augmentation can effectively mitigate data scarcity*  
023 *for long-tail objects in driving scenes, reducing the need for*  
024 *additional real-world data collection.*

## 025 1. Introduction

026 Reliable roadway perception underpins modern transporta-  
027 tion systems including intelligent infrastructure, traffic  
028 monitoring, and autonomous driving. Increasingly, these  
029 systems are integrated into digital representations of the  
030 transportation network, often referred to as digital twins,  
031 where real-world traffic states are mirrored to support mon-  
032 itoring, planning, and automated decision-making. [11]  
033 refers to the perception layer as the "eye" of a digital twin.  
034 Thus, in these environments, perception models form one of  
035 the important interfaces between the physical roadway and

its digital representation.

The fidelity of digital twin representations depends heav-  
ily on the reliability of perception systems that provide ob-  
servations of the physical environment. Object detection  
models identify roadway actors and conditions that form  
the basis for higher-level monitoring and decision-support  
processes. These downstream components, often described  
as the "brain" of a digital twin [11], extract patterns, model  
system dynamics, and support operational decisions based  
on the perceived state of the roadway. When perception  
models fail to detect rare or unexpected objects, the result-  
ing digital representation may omit critical elements of the  
environment, leading to incomplete situational awareness  
and potentially reducing the reliability of monitoring, pre-  
diction, and decision-support functions built upon that rep-  
resentation.

Contemporary object detection models for street-level  
scenes are commonly trained and benchmarked on large-  
scale datasets such as NuScenes [7] and UA-DETRAC [41].  
While these datasets have significantly advanced driving-  
scene research, they primarily represent frequent traffic par-  
ticipants and common roadway conditions, offering lim-  
ited coverage of rare object categories. As a result, models  
trained on such datasets often struggle with long-tail events  
such as wildlife entering the roadway, stalled vehicles, tem-  
porary work-zone artifacts, or unexpected obstacles.

Existing approaches to mitigating long-tail imbalance  
generally fall into three categories: selective data curation,  
conventional data augmentation, and simulation. Manual  
curation of rare instances from existing datasets is inher-  
ently limited by their sparse occurrence. Traditional aug-  
mentation techniques such as copy-paste insertion can in-  
crease class frequency but remain constrained to the limited  
pool of available object instances and may introduce light-  
ing and shadow mismatches or truncated instances inher-  
ited from occlusions in the source images. Collecting addi-  
tional real-world data from sources such as dashcam footage or  
staged deployments incurs significant logistical and anno-  
tation costs while still failing to guarantee coverage of rare  
events. Simulation platforms such as CARLA provide con-

076	trolled environments for generating annotated data, but the	129
077	diversity of scenes and objects depends on available assets	130
078	and modeling effort, and models trained heavily on simu-	131
079	lated imagery often require additional adaptation due to the	132
080	sim-to-real gap.	133
081	Recently, diffusion-based generative models have	134
082	emerged as powerful tools for producing diverse and high-	135
083	fidelity images. These models have been successfully ap-	136
084	plied to tasks such as image editing [5, 14] and inpaint-	137
085	ing [34, 40], and are increasingly explored for synthetic	138
086	data generation and augmentation [10]. Unlike traditional	139
087	augmentation methods that rely on simple transformations	140
088	such as rotation, flipping, or brightness adjustments, dif-	141
089	fusion models can generate entirely new visual instances	142
090	that extend beyond the configurations present in the original	143
091	dataset. By learning a generative representation of the data	144
092	distribution, they can synthesize realistic samples that pop-	145
093	ulate underrepresented regions while maintaining seman-	146
094	tic coherence. Furthermore, modern diffusion frameworks	147
095	support fine-grained control over content and style, allow-	
096	ing augmentation strategies to be tailored to specific fail-	
097	ure modes (e.g., rare-object insertion, lighting variations,	
098	or domain-shift simulation). These properties collectively	
099	make diffusion-based augmentation more expressive, flexi-	
100	ble, and aligned with the needs of modern deep vision mod-	
101	els than conventional transformation-based techniques.	
102	Despite this potential, the use of diffusion models for	
103	safety-critical roadway perception remains limited. Prior	
104	work has largely focused on scene stylization or weather	
105	adaptation [22], or limited forms of object insertion with-	
106	out fine-grained control over the generated instances [47].	
107	Modern diffusion frameworks support conditioning mecha-	
108	nisms such as text prompts and spatial masks, which enable	
109	detailed control over object category, pose, viewpoint, and	
110	spatial placement. Such controllability enables the system-	
111	atic generation of rare traffic objects under diverse contex-	
112	tual configurations, providing a promising mechanism for	
113	enriching long-tail object distributions.	
114	To address this gap, we investigate whether controllable	
115	generative models can be leveraged to enrich rare-object	
116	distributions in driving scenes. Our framework introduces	
117	rare traffic objects into real roadway imagery through con-	
118	trolled generative augmentation, enabling the creation of re-	
119	alistic corner cases without costly real-world data collection	
120	or manual staging. The proposed approach integrates gener-	
121	ative synthesis with verification and filtering mechanisms to	
122	ensure that inserted objects remain semantically consistent	
123	and spatially plausible for detector training.	
124	Specifically, we propose <i>RareCrafter</i> , a training-free	
125	generative augmentation framework that synthesizes rare	
126	objects directly into real driving scenes using conditional	
127	generative models and structured guidance. By enriching	
128	training data with rare corner cases, the framework aims to	
	improve the robustness of perception systems when encoun-	129
	tering uncommon but foreseeable roadway objects.	130
	The main question explored in this work is: <i>Can genera-</i>	131
	<i>tive models mitigate the costly process of collecting rare-</i>	132
	<i>object data in driving scenes?</i> We empirically evaluate	133
	whether generative augmentation can improve object detec-	134
	tion performance while reducing reliance on extensive data	135
	collection and manual annotation. Our main contributions	136
	are:	137
	• <i>RareCrafter</i> , a controllable generative augmentation	138
	framework for synthesizing rare objects in real driving	139
	scene images.	140
	• A structured synthesis pipeline with depth-aware place-	141
	ment, structured prompt diversification, and region-	142
	conditioned verification for consistent object insertion.	143
	• Augmenting training data with <i>RareCrafter</i> , which im-	144
	proves rare-object detection on a real-world corner-case	145
	dataset while preserving performance on common cate-	146
	gories.	147
	<b>2. Related Works</b>	148
	<b>Conditional image generation.</b> Diffusion models have	149
	emerged as a powerful class of generative models, achiev-	150
	ing state-of-the-art results in high-fidelity image synthe-	151
	sis [15, 28]. They work by training a model to reverse a	152
	fixed, gradual process of adding noise; allowing them to	153
	start from pure Gaussian noise and iteratively denoise it	154
	to a coherent, high-quality image. Unlike unconditional	155
	generative models, conditional diffusion models provide a	156
	fine-grained control over visual content. Specifically, Con-	157
	trolNet [48] enhances controllability by injecting spatial	158
	conditions, such as edge maps or poses, into the diffusion	159
	backbone. In addition to diffusion-based architectures, re-	160
	cent hybrid transformer–flow models such as FLUX [19]	161
	have demonstrated impressive capabilities in aligning text	162
	prompts with high-fidelity visual outputs by integrating	163
	large-scale attention modules with rectified flow refine-	164
	ment steps. This design enables Flux to produce photorealistic	165
	images with strong semantic consistency and robust con-	166
	trollability. Object-oriented editing leverages conditional	167
	generative models to remove [16, 21, 36, 45] an object spec-	168
	ified by a mask or synthesize [9, 35, 37, 45] a new object	169
	with control over object attributes and appearances through	170
	text prompts or reference images. In this work, we take	171
	advantage of FLUX-ControlNet-Inpainting [36] to generate	172
	rare objects conditioned on a driving background scene and	173
	structured prompts.	174
	<b>Corner case and anomaly driving datasets.</b> Captur-	175
	ing rare and unexpected events in driving scenes remains	176
	challenging due to their infrequent occurrence in real-world	177
	data. To facilitate their study, several datasets have been	178
	introduced that focus on anomalous objects, unusual traffic	179
	situations, and other corner cases in driving scenes.	180

Table 1. Comparison of methods.

Method	Training Free	Rare Object Insertion	Fine-Grained Control
RareDiffusion [47]	✗	✓	✗
LTDA-Drive [46]	✓	✗	✓
RareCrafter (ours)	✓	✓	✓

181 One line of research mines or re-annotates existing  
 182 datasets to highlight anomalous events and corner-case ob-  
 183 jects [13, 20, 42]. While this improves visibility of rare  
 184 instances, it remains inherently constrained by the original  
 185 data distribution. Another direction focuses on curated real-  
 186 world collections or staged scenarios [25, 27, 30], where  
 187 unusual objects are deliberately placed in traffic scenes.  
 188 Although these efforts provide valuable benchmarks, they  
 189 require substantial manual effort and offer limited diver-  
 190 sity of naturally occurring rare-event combinations. Some  
 191 anomaly-centric datasets [8] further rely on binary anomaly  
 192 labels without detailed semantic categorization, limiting  
 193 their usefulness for object-level perception tasks. Alterna-  
 194 tive approaches attempt to increase rare-event coverage syn-  
 195 thetically. Image-level blending methods overlay anoma-  
 196 lous objects onto real scenes [3], but are restricted by the  
 197 available object pool of the source dataset from which the  
 198 pasted instances are drawn, and often introduce geometric  
 199 inconsistencies. Simulation-based datasets [4, 6, 13, 17, 25]  
 200 enable scalable scenario generation under controlled condi-  
 201 tions; however, their diversity depends on predefined assets  
 202 and scene configurations, and models trained primarily on  
 203 simulated data may struggle to generalize to real-world im-  
 204 agery.

205 **Diffusion models for image data augmentation.** Dif-  
 206 fusion models have emerged as powerful generative tools  
 207 for producing diverse, high-quality, and semantically coher-  
 208 ent augmented data. Traditional augmentation methods ex-  
 209 pand datasets through simple image transformations such as  
 210 translation, flipping, cropping, or color jittering, which pre-  
 211 serve semantic content while modifying appearance [44].  
 212 However, these operations provide limited diversity and are  
 213 often insufficient for modern vision tasks [1]. Generative  
 214 augmentation with diffusion models instead samples from  
 215 an approximate data distribution, enabling the creation of  
 216 entirely new images or controlled modifications of existing  
 217 ones. Through conditioning signals such as text prompts,  
 218 depth maps, bounding boxes, or segmentation masks, these  
 219 models can generate synthetic object-detection training im-  
 220 ages [10], generate counterfactual examples for bias miti-  
 221 gation [29], or improve robustness under out-of-distribution  
 222 shifts [39].

223 Recent generative models such as FLUX [19] adopt  
 224 flow-matching as an alternative training paradigm. By  
 225 learning probability paths grounded in optimal transport,

flow-based models enable faster sampling than diffusion  
 models [23], though often at the cost of reduced sample di-  
 versity [33]. While typically viewed as a limitation, prior  
 work [2] shows that lower diversity can benefit training-  
 free image editing. Building on this insight, we exploit the  
 structure-preserving behavior of flow-matching models to  
 insert rare objects into driving scenes while minimizing un-  
 intended background modifications, enabling targeted aug-  
 mentation for long-tail object detection.

Existing diffusion-based augmentation methods remain  
 limited for rare-object learning (Table 1). RareDiffu-  
 sion [47] trains a dedicated diffusion model for a prede-  
 fined set of categories, restricting scalability to arbitrary  
 classes and offering limited control over viewpoint or struc-  
 tural variations, thereby reducing its effectiveness for tar-  
 geted augmentation. LTDA-Drive [46] adopts a training-  
 free, text-conditioned inpainting framework; however, it  
 focuses on object categories with abundant training sam-  
 ples, making it less transferable to rare objects with limited  
 representation in existing datasets. Additionally, neither  
 method has publicly released code. In contrast, RareCrafter  
 is training-free, extensible to arbitrary rare categories via  
 pretrained generative models, and enables fine-grained con-  
 trol over object attributes, making it suitable for targeted  
 long-tail data augmentation.

## 3. Methodology

### 3.1. Rare Object Generation

We synthesize rare-object instances using FLUX-  
 ControlNet-Inpainting [36] applied to background images  
 taken from the ONCE dataset [26]. The inpainter receives:  
 (i) the background frame; (ii) a mask defining the insertion  
 region; and, (iii) a text prompt describing the target class.  
 Objects are inserted within predefined bounding boxes to  
 preserve geometric consistency with the scene.

**Structured Prompt Diversification.** A central compo-  
 nent of our synthesis pipeline is controlled prompt diversifi-  
 cation. Rather than using a single generic description (e.g.,  
 “a stroller” or “a motorcycle”), we leverage a large language  
 model to generate a structured set of object-centric prompts.  
 We exploit the LLM’s latent knowledge of object taxon-  
 omy, materials, structural components, configurations, and  
 viewpoint variations to enumerate semantically meaningful  
 intra-class variations.

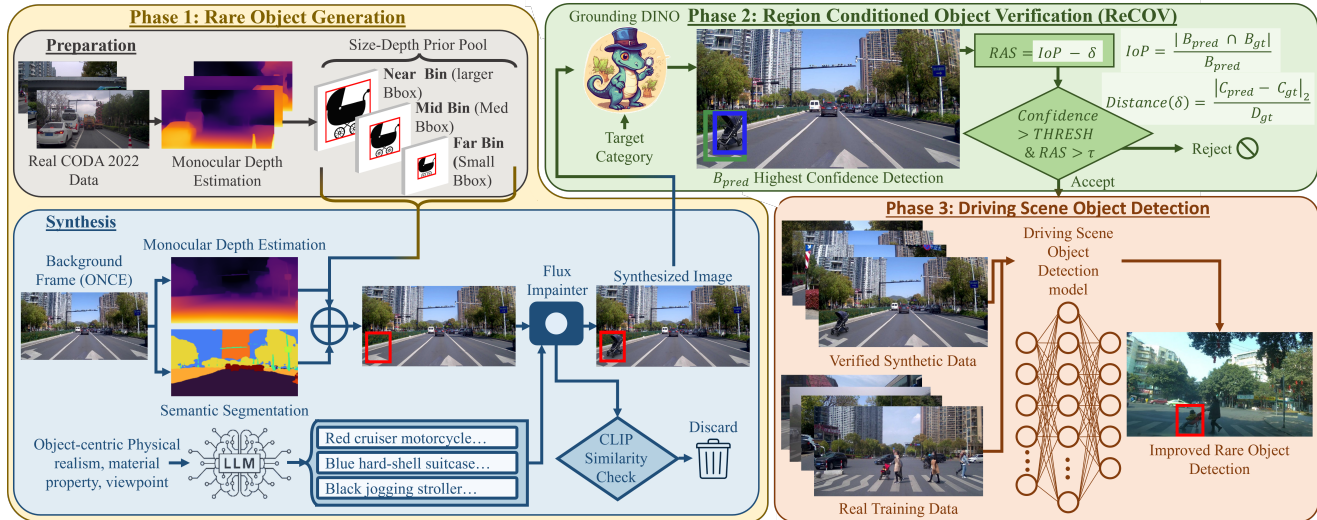


Figure 1. **Overview of the proposed RareCrafter method.** Starting from a real driving-scene image, monocular depth and semantic segmentation are estimated to identify valid ground regions for object insertion. The depth value at the sampled point determines a near, mid, or far bin, which is used to select a bounding-box template from a depth-conditioned empirical pool constructed from real data. Object appearance diversity is introduced through structured prompt diversification, where an LLM generates object-centric prompts describing variations in structure, materials, and configurations, which are used with a FLUX-ControlNet inpainting model to synthesize the rare object at the selected location. Generated samples are filtered using CLIP similarity and further validated through Region-Conditioned Object Verification (ReCOV) with GroundingDINO to ensure the object appears within the intended region. Verified synthetic images are then combined with real training data to train the driving-scene object detector.

269 Prompts are constrained to describe exactly one object and to emphasize physical realism, viewpoint, material properties, and structural attributes, while deliberately  
270 excluding environmental factors such as weather, lighting direction, or scene context. This decouples object specification from background conditions, allowing the inpainting model to adapt seamlessly to the underlying background.  
271  
272  
273  
274  
275

276 This strategy is particularly beneficial for rare categories, where limited real data can lead detectors to overfit to narrow visual prototypes. By explicitly varying structured object attributes, we approximate a broader intra-class distribution than is available in the training set, mitigating representation sparsity. For example, stroller prompts vary by functional subtype (compact, jogging, double), frame material, fabric type, and structural configuration. Motorcycle prompts vary by vehicle style (sport, cruiser, off-road), paint finish, exposed versus enclosed components, accessories, and usage condition. The system prompt used for prompt construction is included in Supplementary Fig.8.  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287

288 **Depth-Conditioned Empirical Size Prior.** A critical component of the synthesis pipeline is the construction of insertion masks that are geometrically plausible while remaining stable for generative modeling and detector supervision. Prior work such as LTDA-Drive [46] models object scale by fitting a parametric (e.g., Gaussian) distribution over bounding-box statistics. While effective for frequent categories with abundant samples, such parametric  
289  
290  
291  
292  
293  
294  
295

296 fitting becomes unreliable for rare classes where only a limited number of real instances are available. In this low-data regime, distribution fitting can produce unrealistic aspect ratios and distorted scale configurations (e.g., unusually tall-thin stroller boxes resembling pedestrians), which deviate from the empirical geometry of the category and provide misleading spatial cues to the inpainter.  
297  
298  
299  
300  
301  
302

303 To avoid distributional artifacts introduced by parametric modeling, we directly leverage the empirical bounding boxes observed in real data. Using the CODA2022 training set, we associate each rare-object instance (suitcase, stroller, motorcycle) with its predicted monocular depth and its normalized bounding-box dimensions. Depth values are discretized into coarse bins (near, mid, far). For each class and depth bin, we construct a non-parametric pool consisting of the observed bounding-box dimensions from real instances. This produces a depth-conditioned empirical prior that preserves realistic aspect ratios and category-specific scale characteristics without requiring distribution fitting.  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314

315 At synthesis time, candidate insertion locations are sampled from semantically valid regions (e.g., road or sidewalk areas without existing annotations). The depth at the sampled pixel determines the corresponding depth bin, and a bounding-box template is drawn from that bin’s empirical pool. This preserves the depth-dependent scale hierarchy observed in real driving scenes.  
316  
317  
318  
319  
320  
321

322 However, extremely small bounding boxes were em-

323 empirically found to destabilize diffusion-based inpainting  
324 (see Tab. 3): masks with limited spatial support tend  
325 to encourage background reconstruction or prompt under-  
326 conditioning. To address this, the sampled bounding box  
327 is uniformly expanded by a fixed margin before being used  
328 both as the inpainting mask and as the detector ground-truth  
329 box. This dilation is applied consistently across all depth  
330 bins and categories, preserving the relative depth-dependent  
331 scale ordering: boxes drawn from the far bin remain smaller  
332 than those from the mid and near bins.

333 Consequently, the final insertion masks maintain semantic  
334 placement consistency and remain anchored to the em-  
335 pirical depth-conditioned size prior, while accommodating  
336 the practical stability requirements of generative synthesis.

337 **Semantic Filtering.** To enforce category consistency,  
338 we compute the CLIP similarity between the generated re-  
339 gion and the target class name, discarding samples below a  
340 predefined threshold. Final acceptance additionally requires  
341 passing the region-conditioned verification described in the  
342 next section.

### 343 3.2. Region-Conditioned Object Verification

344 Initial experiments revealed a failure mode in which the  
345 generative model reconstructs the masked region with con-  
346 textually plausible background content while failing to ren-  
347 der the intended object. Such cases may achieve high global  
348 semantic similarity (e.g., CLIP score) despite unsuccessful  
349 object insertion, motivating an explicit region-level veri-  
350 fication procedure.

351 We therefore introduce a region-conditioned validation  
352 stage based on open-vocabulary detection using Ground-  
353 ingDINO [24]. For each generated image, GroundingDINO  
354 is prompted with the target category and applied to the full  
355 image. Let  $B_{gt}$  denote the inpainting mask (ground-truth  
356 insertion region) and  $B_{pred}$  the highest-confidence detection  
357 for the target category exceeding a predefined confidence  
358 threshold. If no such detection exists, the insertion is re-  
359 jected.

360 Because  $B_{gt}$  defines a placement region rather than a  
361 tight object box, standard IoU is not suitable. Instead,  
362 we measure spatial containment using Intersection-over-  
363 Prediction (IoP):

$$364 \text{IoP}(B_{pred}, B_{gt}) = \frac{|B_{pred} \cap B_{gt}|}{|B_{pred}|}$$

365 IoP quantifies the proportion of the detected object that lies  
366 inside the intended insertion region.

367 To additionally enforce geometric alignment, we intro-  
368 duce a normalized center-distance term. Let  $c_{gt}$  and  $c_{pred}$   
369 denote the box centers, and let  $D$  be the diagonal length of  
370 the minimal enclosing box covering both  $B_{gt}$  and  $B_{pred}$ . We  
371 define

$$372 \delta = \frac{\|c_{pred} - c_{gt}\|_2}{D}$$

This term penalizes detections whose centers deviate from  
373 the intended insertion location. Center alignment is partic-  
374 ularly relevant for our downstream detector, which employs  
375 center-based supervision. 376

377 We combine containment and alignment into a *Region*  
378 *Alignment Score* (RAS):

$$379 \text{RAS} = \text{IoP}(B_{pred}, B_{gt}) - \delta$$

380 A generated image is considered successfully verified if  
381 (i) the detection confidence exceeds a fixed threshold and  
382 (ii)  $\text{RAS} \geq \tau$ , where  $\tau$  is a predefined alignment threshold.  
383 This region-conditioned verification filters out context-only  
384 reconstructions and ensures that the synthesized object is  
385 both spatially contained and center-aligned within the des-  
386 ignated insertion region.

### 387 3.3. Driving Scene Object Detection

388 We use FCOS [38] as the object detection model in all  
389 experiments. FCOS is a one-stage, anchor-free architec-  
390 ture that avoids proposal generation and predefined anchor  
391 boxes, making it well-suited to driving scenes with sub-  
392 stantial scale variation across categories, including rare and  
393 small objects.

394 Our goal is not to optimize detector performance but to  
395 measure the effect of generative data augmentation under a  
396 controlled setup. Therefore, the detector architecture, train-  
397 ing protocol, and hyperparameters remain fixed across all  
398 experiments to ensure fair comparison between real-only  
399 and augmented training regimes. Evaluation is conducted  
400 within the same driving-scene domain as training; we avoid  
401 cross-domain settings (e.g., indoor or artistic imagery) so  
402 that performance differences can be attributed to augmenta-  
403 tion rather than domain shift.

## 404 4. Experiments

### 405 4.1. Datasets and Implementation Details

406 **Datasets.** We use ONCE [26] and CODA [20]. ONCE is a  
407 large-scale autonomous driving dataset containing 15k fully  
408 annotated scenes across diverse environments, weather con-  
409 ditions, and day/night settings. Unless stated otherwise, we  
410 use images from camera 3 only. In our framework, ONCE  
411 serves solely as the background source for synthetic rare-  
412 object insertion.

413 CODA is a real-world corner-case dataset with bound-  
414 ing box annotations for 43 object categories. We use the  
415 CODA2022 subset, which contains 80,180 annotated ob-  
416 jects and is divided into a training set and a test set, each  
417 comprising 4,884 images. CODA2022 includes both com-  
418 mon driving classes (e.g., car, bus, truck, pedestrian, cyclist)  
419 and rare categories.

420 To prevent data leakage between synthetic training data  
421 and evaluation data, we enforce a strict dataset separation.

422 Synthetic images are generated using ONCE backgrounds  
423 that do not overlap with CODA2022, with any shared images  
424 removed prior to synthesis. Real rare-object training  
425 samples are taken exclusively from the CODA2022 training  
426 split, while all evaluations are performed on the CODA2022  
427 test split.

428 Following CODA’s definition of corner cases, which in-  
429 cludes (i) novel classes and (ii) novel instances of common  
430 classes, our study focuses exclusively on the first type, novel  
431 (rare) classes. Accordingly, only rare categories are synthe-  
432 sized, and evaluation is performed on real images contain-  
433 ing only common and mentioned rare categories.

434 For the CopyPaste baseline, object instances are sourced  
435 from the LVIS dataset [12], specifically from the *suitcase*,  
436 *baby buggy*, and *motorcycle* categories. Since LVIS masks  
437 can be noisy, we apply additional mask refinement and fil-  
438 tering before insertion. Dataset splits and their usage are  
439 summarized in Tab.5 of Supplementary.

440 **Implementation Details.** Without loss of generality, we  
441 consider three object categories, *suitcase*, *stroller*, and *mo-*  
442 *torcycle*, to build our framework upon. We selected the  
443 FCOS [38] object detection model [38] and trained it with  
444 the Adam optimizer [18] for 20 epochs using a base learning  
445 rate of  $1.6 \times 10^{-5}$  with cosine annealing and a batch size of  
446 16. Images are processed at their original resolution. Data  
447 augmentation during epochs 1–18 includes random resiz-  
448 ing, cropping, color jitter, and horizontal flipping; the final  
449 two epochs use only resizing and horizontal flipping.

450 For object insertion, we adopt a zoom-in strategy: the  
451 target placement box is expanded to include surrounding  
452 context, cropped, and resized to  $1024 \times 1024$  before be-  
453 ing passed to the generative model. This enables consis-  
454 tent inpainting across varying image aspect ratios and ob-  
455 ject scales. We use the FLUX model [19] with 28 inference  
456 steps. Each image is generated up to ten times with different  
457 seeds and filtered using  $\tau$  of 0.6 and a CLIP Score thresh-  
458 old; images failing the threshold are discarded. Thresholds,  
459 determined on the CODA2022 training split, are 0.22 (mo-  
460 torcycle), 0.20 (suitcase), and 0.25 (stroller).

461 Monocular depth estimation is performed using DPT-  
462 Large [32], and ground-region localization using Seg-  
463 Former [43]. Detection performance is evaluated using the  
464 COCO protocol, reporting mAP and mAR averaged over  
465 IoU thresholds from 0.50 to 0.95. All experiments are con-  
466 ducted on a single NVIDIA A100 GPU.

## 467 4.2. Evaluation on Object Detection

468 In this section, we evaluate the effectiveness of the proposed  
469 RareCrafter framework. For each rare category, we gener-  
470 ate synthetic data with twice as many images as the real  
471 training set and augment it with the original dataset. Fig. 2  
472 presents qualitative results, while Tab. 2 reports quantitative  
473 performance.

Table 2. Comparison of object detection performance of FCOS [38] across training configurations for rare categories.

Regime	Stroller		Motorcycle		Suitcase	
	mAP	mAR	mAP	mAR	mAP	mAR
Real only	3.66	22.60	0.13	30.50	0.60	14.60
CopyPaste	4.26	<b>35.03</b>	2.16	30.26	0.63	10.76
RareCrafter (ours)	<b>5.76</b>	35.00	<b>2.53</b>	<b>31.43</b>	<b>4.56</b>	<b>15.13</b>
RareCrafter vs Real ( $\Delta$ )	+2.10	+12.40	+2.40	+0.93	+3.96	+0.53

474 We compare RareCrafter with a CopyPaste baseline,  
475 which inserts segmented instances of the target categories  
476 into predefined bounding boxes using alpha blending. In  
477 practice, raw instance segmentation masks often contain ar-  
478 tifacts such as holes or fragmented regions and may produce  
479 objects with unrealistic scales relative to the target bounding  
480 box. To mitigate these issues, we refine the masks by filling  
481 holes and removing discontinuities, and apply an instance  
482 filtering step to discard implausible samples (e.g., extremely  
483 small objects or instances requiring excessive rescaling).  
484 The remaining instances are resized while preserving their  
485 aspect ratio and constrained to fit within the target bounding  
486 boxes before compositing them into the scene.

487 Tab. 2 compares different training regimes using FCOS  
488 for rare-category detection. Training on real data only re-  
489 sults in very limited performance. Incorporating synthetic  
490 data through RareCrafter substantially improves detection  
491 performance across all categories. Specifically, RareCrafter  
492 increases stroller performance by +2.10 mAP and +12.40  
493 mAR, improves motorcycle mAP by +2.40, and boosts suit-  
494 case mAP by +3.96. Compared to the CopyPaste augmen-  
495 tation baseline, RareCrafter consistently achieves higher  
496 mAP across all categories. These results demonstrate that  
497 RareCrafter effectively mitigates the data scarcity problem  
498 for rare objects and significantly improves detection perfor-  
499 mance. Importantly, performance on common categories  
500 remains stable across training regimes (Tab.6 of Supple-  
501 mentary), with an average maximum variation of approx-  
502 imately 1.0 mAP. Qualitative examples of the synthesized  
503 training images generated by our RareCrafter framework  
504 and the CopyPaste baseline are illustrated in Supplementary  
505 Fig.5 and Fig.6, respectively.

## 506 4.3. Ablation Study

507 **Ablation on bounding box dilation.** We ablate our mask  
508 dilation strategy and perform inpainting directly within  
509 bounding boxes sampled from the empirical distribution of  
510 the CODA2022 training set. This setting constrains the in-  
511 painter to insert rare objects within the original bounding  
512 box extents observed in real data. In both settings, the syn-  
513 thesis budget is fixed (100 background images per rare cat-  
514 egory), and the same empirical bounding-box pool is used;  
515 the only difference is whether the sampled box is expanded



Figure 2. **Qualitative object detection results on the CODA2022 dataset.** The detection models are trained under three regimes: real-only data, real data augmented with CopyPaste, and real data augmented with RareCrafter. In the first row, the detector trained on real-only data incorrectly localizes another object as a stroller. In the second row, both the real-only and CopyPaste models fail to detect the motorcycle. In the third row, the real-only model confuses a pedestrian carrying a suitcase with the suitcase itself, while CopyPaste incorrectly localizes another object as a suitcase. In contrast, RareCrafter correctly detects the suitcase without misclassifying the pedestrian.

516 before inpainting.

517 We report two synthesis quality metrics. The *CLIP pass*  
518 *rate* measures the proportion of generated samples whose  
519 masked region exceeds the CLIP similarity threshold with  
520 the target class name. The *ReCOV rate* measures the pro-  
521 portion of CLIP-passed samples that additionally pass the  
522 region-conditioned verification described in Sec. 3.2, mean-  
523 ing that GroundingDINO detects the intended object within  
524 the designated insertion region with sufficient alignment.

525 As shown in Tab. 3, removing mask dilation substan-  
526 tially degrades both metrics. Without dilation, the CLIP  
527 pass rate drops from 98% to 71%, and the ReCOV rate  
528 drops from 97% to 58%. Qualitatively, small insertion re-  
529 gions frequently cause the inpainter to reconstruct plausible  
530 background instead of inserting the requested object. Such  
531 cases may still achieve moderate semantic similarity, which  
532 means that a CLIP-only filter would retain samples where  
533 the labeled bounding box contains only background rather  
534 than the intended object. This introduces label noise into the  
535 synthetic dataset, where the annotated bounding box some-  
536 times contains the intended object and sometimes contains  
537 only background.

538 Applying mask dilation provides the generative model  
539 with sufficient spatial support to synthesize the object, re-  
540 ducing prompt-ignoring behavior and improving both se-  
541 mantic consistency and spatial verification. As a result, the  
542 CLIP and ReCOV rates both rise, indicating that most sam-  
543 ples satisfying the semantic filter also contain a verifiable  
544 object instance within the intended region. These results  
545 confirm that enlarging the inpainting region is beneficial for

Table 3. Ablation on mask dilation for the inpainting region. En-  
larging the bounding box significantly improves both semantic  
consistency (CLIP pass rate) and successful object insertion veri-  
fied by ReCOV.

Design Choice	CLIP Pass Rate (%) $\uparrow$	ReCOV Rate (%) $\uparrow$
Without dilation	71	58
With dilation	<b>98</b>	<b>97</b>

546 stable object synthesis and for preventing noisy synthetic  
547 labels.

548 **Ablation on prompt diversification.** We ablate our  
549 prompt generation pipeline by replacing the structured  
550 prompts with a simple category-only prompt (e.g., "a  
551 stroller") during inpainting. As shown in Tab. 4, structured  
552 prompt diversification consistently improves detection per-  
553 formance across all rare categories. The gains are particu-  
554 larly pronounced for motorcycle and suitcase, where mAP  
555 increases from 0.30 to 2.53 and from 1.8 to 4.56, respec-  
556 tively. These results suggest that object-centric prompt di-  
557 versification introduces greater intra-class appearance vari-  
558 ation in the synthesized data, helping the detector learn  
559 more robust representations for rare categories.

560 **Ablation on synthetic data ratio.** To study the effect  
561 of synthetic data scaling, we vary the ratio of RareCrafter-  
562 generated images added to the training set from  $0\times$  (real-  
563 only) to  $4\times$  the size of the real dataset. As shown in  
564 Fig. 3, adding synthetic samples consistently improves per-  
565 formance over the real-only baseline across all rare cate-  
566 gories. For stroller and suitcase, performance peaks around

Table 4. Ablation study on the effect of the prompt generation pipeline for RareCrafter across rare categories using mAP metric for FCOS.

Regime	Stroller		Motorcycle		Suitcase	
	Simple prompt	Structured prompt	Simple prompt	Structured prompt	Simple prompt	Structured prompt
RareCrafter	4.16	5.76	0.30	2.53	1.8	4.56

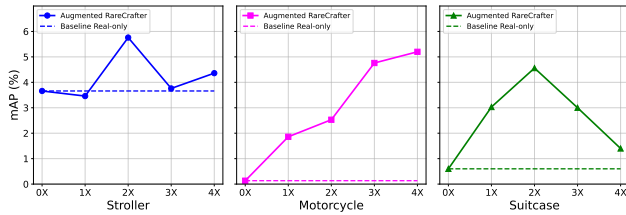


Figure 3. Ablation study on the impact of synthetic data scaling using RareCrafter.

567 a  $2\times$  augmentation ratio and slightly decreases with further  
 568 scaling, indicating diminishing returns. In contrast, motor-  
 569 cycle performance continues to improve as more synthetic  
 570 data is introduced.

571 **Ablation on the inpainting model.** To evaluate the impact  
 572 of the inpainting backbone, we replace Flux Inpainting  
 573 [36] with the SDXL Inpainting model [31] while keep-  
 574 ing all other components unchanged. Both models receive  
 575 the same text prompt and inpainting mask. Figure 4 shows  
 576 a qualitative comparison.

577 Flux demonstrates stronger prompt adherence and better  
 578 captures fine-grained attributes described in the text. The  
 579 generated objects exhibit richer high-frequency details and  
 580 more realistic textures, whereas SDXL outputs often appear  
 581 simplified or stylized. This difference is particularly visi-  
 582 ble in object geometry and material appearance, where Flux  
 583 produces more realistic structure and shading.

584 We also observe that SDXL occasionally introduces un-  
 585 intended elements (e.g., an extra person in the stroller ex-  
 586 ample), while Flux typically restricts generation to the in-  
 587 tended object. In addition, Flux achieves better integration  
 588 with the surrounding scene, producing inpainted regions  
 589 with lighting, color, and texture that are more consistent  
 590 with the background. Additional qualitative comparisons  
 591 between Flux and SDXL are provided in Supp. Fig7.

## 592 5. Conclusion

593 This work introduced *RareCrafter*, a training-free genera-  
 594 tive augmentation framework for mitigating the long-tail  
 595 distribution of object categories in driving-scene datasets.  
 596 By inserting rare traffic objects into real roadway imag-  
 597 ery using controllable flow matching-based generation,  
 598 RareCrafter enables systematic creation of rare training  
 599 instances without additional real-world data collection.



Figure 4. **Ablation on the inpainting model.** Comparison between Flux and SDXL Inpainting using the same prompt and mask. Flux (left) generates objects with stronger prompt adherence, richer details, and better scene integration, while SDXL (right) often produces simplified or stylized outputs with a cartoon-like appearance.

The framework combines prompt diversification, depth-  
 600 conditioned size priors, and region-based verification to  
 601 guide synthesis and maintain semantic consistency. Exper-  
 602 iments on the CODA2022 corner-case dataset show aug-  
 603 menting training data with synthesized rare-object instances  
 604 improves detection performance and outperforms a Copy-  
 605 Paste baseline, demonstrating the potential of controllable  
 606 generative models to enrich underrepresented driving-scene  
 607 data and reduce the burden of collecting rare safety-critical  
 608 scenarios. 609

During our experiments, we observed that approximat-  
 610 ing object scale using empirical depth-conditioned size  
 611 statistics provides stable synthesis but does not explic-  
 612 itly model precise geometric relationships between objects  
 613 and the scene. While the approach preserves the depth-  
 614 dependent scale hierarchy observed in real driving data,  
 615 stronger geometric cues could improve scale consistency.  
 616 Future work may explore additional geometric supervision,  
 617 such as depth annotations, 3D structure, or LiDAR mea-  
 618 surements. Improving the representation of rare roadway  
 619 objects in perception systems may ultimately enhance the  
 620 fidelity of digital roadway representations and digital twin  
 621 environments used for monitoring and decision support. 622

623

**References**

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

- [1] Panagiotis Alimisis, Ioannis Mademlis, Panagiotis Radoglou-Grammatikis, Panagiotis Sarigiannidis, and Georgios Th. Papadopoulos. Advances in diffusion models for image data augmentation: A review of methods, models, evaluation metrics and future research directions, 2025. 3
- [2] Omri Avrahami, Or Patashnik, Ohad Fried, Egor Nemchinov, Kfir Aberman, Dani Lischinski, and Daniel Cohen-Or. Stable flow: Vital layers for training-free image editing. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 7877–7888, 2025. 3
- [3] Hermann Blum, Paul-Edouard Sarlin, Juan Nieto, Roland Siegwart, and Cesar Cadena. Fishyscapes: A benchmark for safe semantic segmentation in autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2019. 3
- [4] Daniel Bogdoll, Iramm Hamdard, Lukas Namgyu Rößler, Felix Geisler, Muhammed Bayram, Felix Wang, Jan Imhof, Miguel de Campos, Anushervon Tabarov, Yitian Yang, Martin Gontscharow, Hanno Gottschalk, and J. Marius Zöllner. *AnoVox: A Benchmark for Multimodal Anomaly Detection in Autonomous Driving*, page 206–223. Springer Nature Switzerland, 2025. 3
- [5] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions, 2023. 2
- [6] Tom Bu, Xinhe Zhang, Christoph Mertz, and John M. Dolan. Carla simulated data for rare road object detection. In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pages 2794–2801, 2021. 3
- [7] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving, 2020. 1
- [8] Robin Chan, Krzysztof Lis, Svenja Uhlemeyer, Hermann Blum, Sina Honari, Roland Siegwart, Pascal Fua, Mathieu Salzmann, and Matthias Rottmann. Segmentmeifyoucan: A benchmark for anomaly segmentation, 2021. 3
- [9] Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor: Zero-shot object-level image customization. *arXiv preprint arXiv:2307.09481*, 2023. 2
- [10] Haoyang Fang, Boran Han, Shuai Zhang, Su Zhou, Cuixiong Hu, and Wen-Ming Ye. Data augmentation for object detection via controllable diffusion models. 2024. 2, 3
- [11] Yongjie Fu, Mehmet Kerem Turkcan, Mahshid Ghasemi, Zhaobin Mo, Chengbo Zang, Abhishek Adhikari, Zoran Kostic, Gil Zussman, and Xuan Di. Ai-powered cps-enabled vulnerable-user-aware urban transportation digital twin: Methods and applications. *IEEE Transactions on Intelligent Transportation Systems*, pages 1–18, 2026. 1
- [12] Agrim Gupta, Piotr Dollár, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation, 2019. 6
- [13] Dan Hendrycks, Steven Basart, Mantas Mazeika, Andy Zou, Joe Kwon, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Song. Scaling out-of-distribution detection for real-world settings, 2022. 3
- [14] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control, 2022. 2
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, pages 6840–6851. Curran Associates, Inc., 2020. 2
- [16] Longtao Jiang, Zhendong Wang, Jianmin Bao, Wengang Zhou, Dongdong Chen, Lei Shi, Dong Chen, and Houqiang Li. Smarteraser: Remove anything from images using masked-region guidance, 2025. 2
- [17] Hoon Kim, Kangwook Lee, Gyeongjo Hwang, and Changho Suh. Crash to not crash: Learn to identify dangerous vehicles using a simulator. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):978–985, 2019. 3
- [18] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 6
- [19] Black Forest Labs. Flux: A hybrid transformer–flow architecture for high-fidelity text-to-image synthesis. Model Card, 2024. <https://blackforestlabs.ai/>. 2, 3, 6
- [20] Kaican Li, Kai Chen, Haoyu Wang, Lanqing Hong, Chaoqiang Ye, Jianhua Han, Yukuai Chen, Wei Zhang, Chunjing Xu, Dit-Yan Yeung, Xiaodan Liang, Zhenguo Li, and Hang Xu. Coda: A real-world road corner case dataset for object detection in autonomous driving, 2022. 3, 5
- [21] Ruibin Li, Tao Yang, Song Guo, and Lei Zhang. Rorem: Training a robust object remover with human-in-the-loop, 2025. 2
- [22] Hongbin Lin, Zilu Guo, Yifan Zhang, Shuaicheng Niu, Yafeng Li, Ruimao Zhang, Shuguang Cui, and Zhen Li. Drivegen: Generalized and robust 3d detection in driving via controllable text-to-image diffusion generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 27497–27507, 2025. 2
- [23] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling, 2023. 3
- [24] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 5
- [25] Kira Maag, Robin Chan, Svenja Uhlemeyer, Kamil Kowol, and Hanno Gottschalk. Two video data sets for tracking and retrieval of out of distribution objects, 2022. 3
- [26] Jiageng Mao, Minzhe Niu, Chenhan Jiang, Xiaodan Liang, Yamin Li, Chaoqiang Ye, Wei Zhang, Zhenguo Li, Jie Yu, Chunjing Xu, et al. One million scenes for autonomous driving: Once dataset. 2021. 3, 5
- [27] Alexey Nekrasov, Malcolm Burdorf, Stewart Worrall, Bastian Leibe, and Julie Stephany Berrio Perez. Spotting the Unexpected (STU): A 3D LiDAR Dataset for Anomaly Segmentation in Autonomous Driving. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 3
- [28] Alex Nichol and Pratul Dharwal. Improved denoising diffusion probabilistic models, 2021. 2

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736	[29] Rishubh Parihar, Abhijnya Bhat, Abhipsa Basu, Saswat Mallick, Jogendra Nath Kundu, and R. Venkatesh Babu. Balancing act: Distribution-guided debiasing in diffusion models, 2025. 3	
737		
738		
739		
740	[30] Peter Pinggera, Sebastian Ramos, Stefan Gehrig, Uwe Franke, Carsten Rother, and Rudolf Mester. Lost and found: Detecting small road hazards for self-driving vehicles, 2016. 3	
741		
742		
743		
744	[31] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023. 8, 12	
745		
746		
747		
748	[32] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. <i>CoRR</i> , abs/2103.13413, 2021. 6	
749		
750		
751	[33] Johannes Schusterbauer, Ming Gui, Pingchuan Ma, Nick Stracke, Stefan Andreas Baumann, Vincent Tao Hu, and Björn Ommer. <i>FMBBoost: Boosting Latent Diffusion with Flow Matching</i> , page 338–355. Springer Nature Switzerland, 2024. 3	
752		
753		
754		
755		
756	[34] Chaehun Shin, Jooyoung Choi, Heeseung Kim, and Sungroh Yoon. Large-scale text-to-image model with inpainting is a zero-shot subject-driven image generator. <i>arXiv preprint arXiv:2411.15466</i> , 2024. 2	
757		
758		
759		
760	[35] Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price, Jianming Zhang, Soo Ye Kim, and Daniel Aliaga. Object-stitch: Generative object compositing, 2022. 2	
761		
762		
763	[36] Alimama Creative Team. Flux-controlnet-inpainting. <a href="https://github.com/alimama-creative/FLUX-Controlnet-Inpainting">https://github.com/alimama-creative/FLUX-Controlnet-Inpainting</a> , 2024. 2, 3, 8, 12	
764		
765		
766	[37] Yoad Tewel, Rinon Gal, Dvir Samuel, Yuval Atzmon, Lior Wolf, and Gal Chechik. Add-it: Training-free object insertion in images with pretrained diffusion models. In <i>The Thirteenth International Conference on Learning Representations</i> , 2025. 2	
767		
768		
769		
770		
771	[38] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: Fully convolutional one-stage object detection. In <i>Proc. Int. Conf. Computer Vision (ICCV)</i> , 2019. 5, 6, 13	
772		
773		
774	[39] Joshua Vendrow, Saachi Jain, Logan Engstrom, and Aleksander Madry. Dataset interfaces: Diagnosing model failures using controllable counterfactual generation, 2023. 3	
775		
776		
777	[40] Yikai Wang, Chenjie Cao, Junqiu Yu, Ke Fan, Xiangyang Xue, and Yanwei Fu. Towards enhanced image inpainting: Mitigating unwanted object insertion and preserving color consistency. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , 2025. 2	
778		
779		
780		
781		
782	[41] Longyin Wen, Dawei Du, Zhaowei Cai, Zhen Lei, Ming-Ching Chang, Honggang Qi, Jongwoo Lim, Ming-Hsuan Yang, and Siwei Lyu. Ua-detrac: A new benchmark and protocol for multi-object detection and tracking, 2020. 1	
783		
784		
785		
786	[42] Zhongyu Xia, Jishuo Li, Zhiwei Lin, Xinhao Wang, Yongtao Wang, and Ming-Hsuan Yang. Openad: Open-world autonomous driving benchmark for 3d object detection, 2025. 3	
787		
788		
789		
790	[43] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M. Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. <i>CoRR</i> , abs/2105.15203, 2021. 6	
791		
792		
793		
	[44] Mingle Xu, Sook Yoon, Alvaro Fuentes, and Dong Sun Park. A comprehensive survey of image augmentation techniques for deep learning. <i>Pattern Recognition</i> , 137:109347, 2023. 3	794
		795
		796
	[45] Yongsheng Yu, Ziyun Zeng, Haitian Zheng, and Jiebo Luo. Omnipaint: Mastering object-oriented editing via disentangled insertion-removal inpainting, 2025. 2	797
		798
		799
	[46] Mahmut Yurt, Xin Ye, Yunsheng Ma, Jingru Luo, Abhirup Mallik, John Pauly, Burhaneddin Yaman, and Liu Ren. Lta-drive: Llms-guided generative models based long-tail data augmentation for autonomous driving, 2025. 3, 4	800
		801
		802
	[47] Hancheng Zhang, Yuanyuan Hu, Zhendong Qian, Jirui Sha, Min Xie, Yuyang Wan, and Pengfei Liu. Enhancing rare object detection on roadways through conditional diffusion models for data augmentation. <i>IEEE Transactions on Intelligent Transportation Systems</i> , 25(11):19018–19029, 2024. 2, 3	803
		804
		805
		806
		807
		808
		809
	[48] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 2	810
		811
		812
	<b>6. Supplementary</b>	813
	This section provides additional material supporting the experiments presented in the main paper. First, we present qualitative examples of synthetic training images generated with the proposed RareCrafter pipeline and with the CopyPaste baseline, illustrating the visual characteristics of each synthesis approach when inserting rare objects into real ONCE driving scenes. We further include additional qualitative comparisons between different inpainting models used for generative synthesis, highlighting differences in prompt adherence and scene integration.	814
		815
		816
		817
		818
		819
		820
		821
		822
		823
	Next, we summarize the datasets, splits, and usage protocols employed in our pipeline, including the ONCE background pool, CODA2022 training and evaluation splits, LVIS foreground sources used for CopyPaste, and the construction of synthetic datasets. This summary also clarifies the experimental regimes used for training detectors with real-only and real+synthetic data.	824
		825
		826
		827
		828
		829
		830
	To complement the rare-class experiments reported in the main paper, we additionally report detection performance on common object categories to verify that augmenting training with synthetic rare objects does not degrade performance on frequent classes.	831
		832
		833
		834
		835
	Finally, we provide the full system prompt and prompt template used to generate structured object descriptions for the generative synthesis pipeline, which are used to produce diverse object prompts for the inpainting model.	836
		837
		838
		839



Figure 5. **Synthetic training images generated using the RareCrafter framework.** Rare object instances are synthesized directly into real ONCE driving scenes using FLUX-ControlNet-Inpainting conditioned on structured text prompts and predefined insertion regions. The generative model adapts the inserted object to the surrounding scene context, allowing object appearance, shading, and texture to be consistent with the background. Each column corresponds to one target category (stroller, suitcase, and motorcycle), illustrating examples of generated objects placed in different scene configurations.

Table 5. Datasets, splits, and usage in our pipeline. ONCE provides backgrounds, CODA2022 provides real supervision/evaluation, and Synthetic data supplements rare classes. **Common classes (5):** bus, car, cyclist, truck, pedestrian. **Rare classes (3):** motorcycle, stroller, suitcase.

Dataset	Split / Subset	Size & Labels	Purpose	Protocol & Key Notes
ONCE	Background Pool ( $D_{bg}$ ) (Training scenes, Cam 3)	6,002 images 5 Classes (Common)	<b>Synthesis Backgrounds:</b> Canvas for inserting rare objects (via FLUX or Copy-Paste).	<b>Leakage Prevention:</b> We strictly remove ONCE images that overlap with CODA. Overlapping frames are discarded <i>before</i> synthesis.
	Validation ( $D_{val}$ ) Subsets: – Full ( $D_{val}^{full}$ )	4,884 images 8 Classes (5 Common, 3 Rare)	<b>Real Training Data:</b> Used for detector training, hyperparameter selection (e.g., CLIP threshold), and bbox size priors.	<b>Subset Definitions:</b> • $D_{val}^{full}$ : Images with <b>all 8 classes</b> .
	Test Set ( $D_{test}$ )	4,884 images 8 Classes	<b>Evaluation Only:</b> Final reporting on real-world rare objects.	<b>Strict Separation:</b> $D_{test}$ has no overlap with $D_{bg}$ (ONCE) or synthetic data ( $D_{syn}$ ).
LVIS	Source Foreground	160k images Instance masks	<b>Copy-Paste Source:</b> Provides cutouts for motorcycle, stroller, and suitcase that are pasted onto ONCE backgrounds.	Used <b>only</b> for the CopyPaste baseline; masks are refined/filtered due to noisy annotations before compositing.
Synthetic	1. Copy-Paste ( $D_{syn}^{CP}$ ) (LVIS + ONCE)	1 image / class 8 Classes (5 Common, 3 Rare)	<b>Rare Class Training:</b> Supplements rare classes in <i>Real+Synthetic</i> experiments.	<b>Regime Definitions (Source of Rare Classes):</b> • <b>Real-only:</b> Train on $D_{val}^{full}$ . Rare classes are <b>real only</b> . • <b>Real+Syn:</b> Train on $D_{val}^{full} + D_{syn}$ . Real rare classes are <b>augmented</b> with synthetic data.
	2. Generative ( $D_{syn}^{RC}$ ) (FLUX + ONCE)			



Figure 6. **Synthetic training images generated using the CopyPaste baseline.** Object instances from the LVIS dataset are composited into ONCE driving scenes using alpha blending within predefined bounding boxes. Because pasted instances are taken directly from existing images, they may exhibit artifacts originating from the source dataset, such as incomplete objects due to occlusion in the original image, fragmented or noisy segmentation masks, and inconsistencies in lighting, texture, or shadows relative to the target scene. Each column corresponds to one target category (stroller, suitcase, and motorcycle).

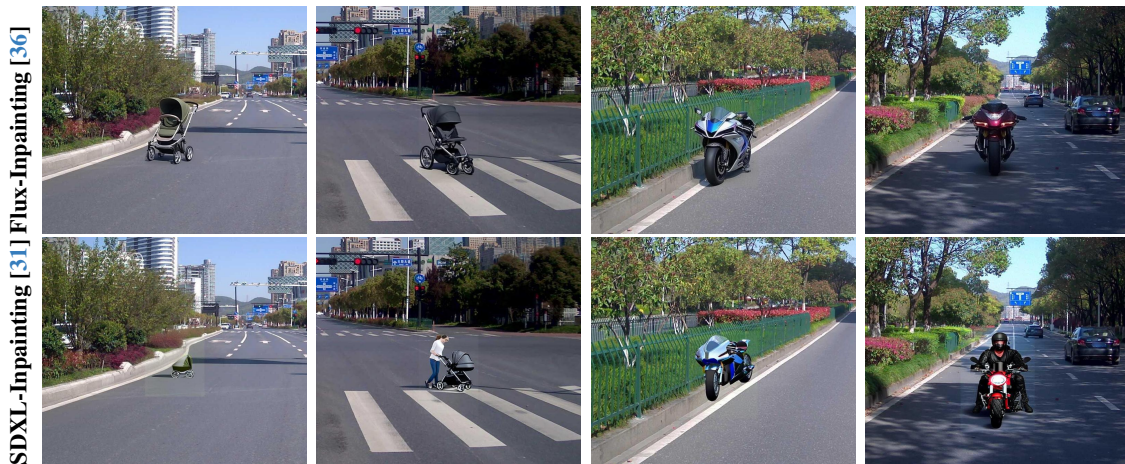


Figure 7. **Additional qualitative comparison between Flux Inpainting and SDXL Inpainting.** Top row: results generated by Flux Inpainting. Bottom row: results generated by SDXL Inpainting using the same prompts and inpainting masks. Flux generally produces objects with stronger prompt adherence, richer details, and better integration with the surrounding scene. In contrast, SDXL often generates simplified or stylized, toy-like objects, may miss fine-grained attributes, and can introduce unintended elements (e.g., a person pushing the stroller or a rider on the motorcycle).

Table 6. Comparison of object detection performance of FCOS [38] across training configurations for common categories.

Regime	Bus	Car	Cyclist	Pedestrian	Truck
Real only	34.63	53.76	36.26	30.40	39.10
+ CopyPaste	33.46	53.13	35.30	29.66	38.43
+ RareCrafter (ours)	33.56	53.10	34.76	29.36	38.46

### System Prompt

You are generating structured prompts for an image inpainting model (Flux Inpainting Alimama).

Your task is to produce a list of independent prompts for inserting ONE motorcycle into an existing scene.

IMPORTANT RULES:

- Each prompt must describe only ONE motorcycle.
- Do NOT describe the environment.
- Do NOT specify lighting direction.
- Do NOT describe weather.
- Focus primarily on the motorcycle itself.
- Include only general integration cues such as:
  - realistic lighting
  - natural shadow
  - correct perspective
  - photorealistic integration

Each prompt must follow the prompt template.

Vary:

- Motorcycle type (sportbike, cruiser, cafe racer, dirt bike, touring motorcycle, electric motorcycle, naked bike, scrambler, etc.)
- Color (red, blue, green, solid colors, dual-tone, metallic shades, matte tones, etc.)
- Finish (matte, glossy, metallic paint, carbon fiber, chrome, brushed metal, etc.)
- Condition (brand new, lightly used, slightly worn, gently scuffed, moderately worn, subtle signs of use, etc.)
- Design details (fairings, exposed engine, saddlebags, knobby tires, aerodynamic bodywork, minimalist frame, custom exhaust, etc.)
- Viewing angle (front view, rear view, side profile, 3/4 front view, 3/4 rear view, slightly top-down view, low-angle view, eye-level view, etc.)

Generate 50 distinct prompts.

Output only the prompts as a numbered list. Do not include explanations. Do not include commentary.

### Prompt Template

*"A realistic [motorcycle type/style], [color], [material/finish details], [condition], [design features], [viewing angle], highly detailed, accurate proportions, realistic lighting, natural shadow integration, correct perspective, photorealistic, seamless integration into the scene"*

Figure 8. System prompt used to generate motorcycle insertion prompts for the inpainting model.