# TOWARDS ROBUSTNESS OF PERSON SEARCH AGAINST CORRUPTIONS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Person search aims to simultaneously detect and re-identify a query person within an entire scene, involving detection and re-identification as a multi-task problem. While existing studies have made significant progress in achieving superior performance on clean datasets, the challenge of robustness under various corruptions remains largely unexplored. To address this gap, we propose two benchmarks, CUHK-SYSU-C and PRW-C, designed to assess the robustness of person search models across diverse corruption scenarios. Previous studies on corruption have been conducted independently for single tasks such as re-identification and detection. However, recent advancements in person search adopt an end-to-end multitask learning framework that processes the entire scene as input, unlike the combination of single tasks. This raises the question of whether independent achievements can ensure corruption robustness for person search. Our findings reveal that merely combining independent, robust detection and re-identification models is not sufficient for achieving robust person search. We further investigate the vulnerability of the detection and representation stages to corruption and explore its impact on both foreground and background areas. Based on these insights, we propose a foreground-aware augmentation and regularization method to enhance the robustness of person search models. Supported by our comprehensive robustness analysis and evaluation framework our benchmarks provide, our proposed technique substantially improves the robustness of existing person search models. Code will be made publicly available.

031 032

033

003 004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

# 1 INTRODUCTION

Person search is a task that involves detecting individuals in complex scenes and subsequently reidentifying the same individuals from a gallery of scene images. Significant advancements in this task have been made through the discriminative capabilities of deep neural networks (DNNs). Since person search requires both pedestrian detection and person re-identification simultaneously, earlier studies have adopted a two-step framework (Lan et al., 2018; Chen et al., 2018; Han et al., 2019; Wang et al., 2020), where independent detection and re-identification models are applied sequentially. More recent models have adopted a one-step approach (Li & Miao, 2021; Lee et al., 2022; Cao et al., 2022; Yu et al., 2022), which processes the entire scene in a single pass and leverages joint end-to-end multi-task learning, leading to significant improvements.

043 Despite these advancements, DNNs have shown vulnerability to common corruptions such as noise, 044 blur, and compression artifacts, often resulting in significant performance degradation (Liu et al., 2024a; Chen et al., 2024; Kong et al., 2024; He et al., 2023; Schiappa et al., 2022; Yi et al., 2021). This vulnerability leads to the necessity for models to maintain robustness under such challenging 046 conditions, driving numerous studies focused on enhancing corruption robustness (Mintun et al., 047 2021; Kar et al., 2022; Dooley et al., 2022; Dong et al., 2023). This line of research has been 048 explored in the fields of detection (Michaelis et al., 2019; Mao et al., 2023; Lee et al., 2024) and re-identification (Chen et al., 2021). However, to our knowledge, the impact of data corruption on person search models remains largely unexplored, highlighting the need for further investigation in 051 this area. 052

To tackle this issue, we introduce two new benchmarks – CUHK-SYSU-C and PRW-C – that extend existing popular person search datasets (Xiao et al., 2017; Zheng et al., 2017) to address their lack

055

070

071

073

074 075

076

077

078

084







Figure 2: Performance evaluation of SoTA person search models on the proposed corruption
benchmarks: CUHK-SYSU-C and PRW-C. We employ six state-of-the-art person search models:
OIMNet (Xiao et al., 2017), NAE (Chen et al., 2020), OIMNet++ (Lee et al., 2022), SeqNet (Li &
Miao, 2021), COAT (Yu et al., 2022), PSTR (Cao et al., 2022). Models are trained on CUHK-SYSU
and PRW, then evaluated on the proposed benchmarks CUHK-SYSU-C and PRW-C, respectively.

of robustness evaluation under corrupted conditions. These benchmarks incorporate 18 types of cor-085 ruption with five severity levels, enabling a detailed evaluation under diverse corruption scenarios, as shown in Figure 1. Given these proposed benchmarks, we extensively evaluate the robustness 087 of corruption in seminal state-of-the-art person search models. Moreover, we explore a straightfor-088 ward solution that naturally arises: integrating a robust detection model (Lee et al., 2024) with a re-identification model (Chen et al., 2021), both of which are designed for corruption robustness. 090 However, our experimental results reveal that existing person search models remain highly vulnera-091 ble to corruption, and a simple integration approach is insufficient to address this issue (Section 3.1). 092 For instance, as shown in Figure 2, we observe a notable performance drop in existing person search models, with relative mAP declines of up to 80%.

094 To investigate the underlying reasons for this phenomenon, we conduct further analysis of corrup-095 tion vulnerability, considering the unique aspects of this task. Since person search models typically 096 localize pedestrian candidates through a detection head and then extract person representations from the detected regions, we explore the sensitivity of the detection and re-identification stages to corrup-098 tion (Section 3.2). Additionally, given the complex scene inputs and large receptive fields in person 099 search models, we assess how corrupted regions in the input images, including the foreground (containing the person) and the background, affect performance (Section 3.3). Our experimental findings 100 indicate that the re-identification stage and the foreground regions are particularly vulnerable to cor-101 ruption, leading to significant performance drops. 102

To this end, we propose a simple yet effective method to enhance corruption robustness foreground-aware augmentation and regularization for robust person representation, which can be easily applied to various existing person search models. Specifically, we apply selective augmentations to the foreground in the input scene images and compute a regularization between the person representations of the original and foreground-augmented images to improve robustness against corruption. Thanks to its simple and easy-to-implement design, extensive experimental results on  CUHK-SYSU-C and PRW-C demonstrate that our method consistently improves the corruption robustness of five state-of-the-art person search models.

In summary, the contributions of this paper are as follows:

- We propose new benchmarks, CUHK-SYSU-C and PRW-C, to evaluate corruption robustness and reveal the significant vulnerability of state-of-the-art person search models to corruption.
- We analyze the sensitivity of both the detection and re-identification stages to corruption and explore the impact of corrupted foreground and background regions.
- Based on this analysis, we propose a foreground-aware augmentation method and a regularization specifically designed to enhance robust person representation.
- Extensive experiments on CUHK-SYSU-C and PRW-C demonstrate that our method significantly improves the corruption robustness of five state-of-the-art person search models.
- 120 121 122

112

113

114

115

116

117

118

119

# 2 RELATED WORKS

123 **Person Search.** Existing person search works are typically categorized into two approaches: the 124 two-step and the one-step methods. Two-step approaches involve two independent detection and 125 reID models, where a detection model first detects people appearing in the scene, and then a reID 126 model recognizes the detected individuals. The primary goal of this two-step approach (Han et al., 127 2019; Dong et al., 2020; Wang et al., 2020) is to specialize these models, which were initially designed for independent tasks, for the person search task by employing techniques like adapting the 128 prediction by a detector to facilitate the reID process. On the other hand, the one-step approach (Xiao 129 et al., 2017; Munjal et al., 2019; Kim et al., 2021; Zhang et al., 2021a; Han et al., 2021; Li & Miao, 130 2021; Yu et al., 2022; Lee et al., 2022; Cao et al., 2022; Li et al., 2022; Han et al., 2022; Yan et al., 131 2023), which performs both sub-tasks with a single model, has recently gained popularity due to its 132 simplicity and efficiency, reducing the number of parameters by nearly half. 133

- One-step approaches take the entire scene as input and jointly train the detection and reID heads in 134 an end-to-end manner. This is commonly achieved by employing a Faster R-CNN (Ren et al., 2016) 135 detection head and reID head with the OIM (Xiao et al., 2017) loss. The one-step approach focuses 136 on addressing challenges that arise from the joint learning of the two tasks. Chen et al. (2020) 137 address conflicting learning objectives, and SeqNet (Li & Miao, 2021) and OIMNet++ (Lee et al., 138 2022) incorporate the detection quality on reID head training and search performance. PSTR (Cao 139 et al., 2022) and COAT (Yu et al., 2022) enable recent one-step person search frameworks to leverage 140 the advantages of transformers. Another challenge unique to person search, unlike conventional 141 reID, is the influence of background information on extracting person representations from larger 142 receptive fields. GLCNet (Qin et al., 2023) mitigates this by extracting features from the entire 143 global scene, including both the background and foreground, to improve the final person embedding. 144 While these studies have made significant advancements in the person search field, they have not 145 been sufficiently investigated for corruption robustness.
- 146 Benchmarking Robustness to Corruption. Deep neural networks are often considered general-147 izable, but they are not as robust to corruption as humans are (Leveque et al., 2022). The pioneering 148 study by Hendrycks & Dietterich (2019) introduced a benchmarking paradigm, where synthetic 149 corruptions are used to evaluate model robustness. This benchmark, ImageNet-C, employs algorith-150 mically generated corruptions, revealing the unexpected vulnerability of networks to even simple 151 perturbations. Following this paradigm, robustness research has expanded into broader areas of vi-152 sion tasks, including object detection (Michaelis et al., 2019), pose estimation (Wang et al., 2021), semantic segmentation (Kamann & Rother, 2020), person re-identification (Chen et al., 2021), and 153 depth estimation (Kong et al., 2024). 154

While significant progress has been made in improving the robustness of specific vision tasks, existing research primarily focuses on individual tasks (Wang et al., 2021; Kamann & Rother, 2020; Kong et al., 2024). Person search, however, is a multi-task problem involving both detection and re-identification, each requiring different input types and resulting in conflicting optimization (Lin et al., 2021). Although previous studies (Michaelis et al., 2019; Chen et al., 2021) have contributed substantially to our understanding of how models for individual tasks perform under various corruption scenarios, the robustness of the multi-task framework in person search against corrupted environments has yet to be sufficiently explored. Table 1: **Performance of state-of-the-art person search models under corruption.** The column 'CUHK-SYSU' ('PRW') denotes the performance measured on the clean images, and 'CUHK-SYSU-C' ('PRW-C') indicates the performance measured on the corrupted images.

		CUHK	-SYSU	SYSU CUHK-SYSU-C								
Method	Sea	urch	Detec	tion		Sear	ch			Dete	ction	
	R@1	mAP	Recall	AP	rR@1	rmAP	R@1	mAP	rRecall	rAP	Recall	AP
OIMNet	87.7	86.2	87.5	81.1	37.7	36.7	33.0	31.7	78.8	77.4	68.9	62.7
NAE	92.3	91.4	92.3	86.9	44.3	42.4	40.9	38.8	72.7	71.8	67.1	62.3
OIMNet++	94.0	93.2	92.4	88.9	38.6	37.0	36.3	34.5	77.3	76.3	71.4	67.8
SeqNet	94.5	93.8	92.0	89.2	46.2	44.4	43.6	41.6	74.5	74.1	68.6	66.1
COAT	94.7	94.2	91.3	88.1	52.6	50.4	49.8	47.5	76.4	75.9	69.8	66.9
PSTR	94.9	93.6	89.5	66.9	49.5	46.3	47.0	43.3	71.6	74.1	64.0	49.6
OADG+CIL	94.0	92.9	94.7	90.9	48.5	46.0	45.6	42.8	89.8	88.4	85.0	80.4

			PF	RW					PRV	W-C			
	Method	Sea	urch	Detec	tion		Sea	rch			Deteo	ction	
		R@1	mAP	Recall	AP	rR@1	rmAP	R@1	mAP	rRecall	rAP	Recall	AP
	OIMNet	76.7	37.3	93.7	85.0	45.8	23.6	35.2	8.8	79.4	78.8	74.4	66.9
	NAE	80.6	42.8	93.3	88.7	42.6	20.9	34.3	8.9	67.3	64.9	62.8	57.6
C	IMNet++	83.2	47.3	96.3	93.2	43.9	19.9	36.5	9.4	75.2	74.2	72.4	69.2
	SeqNet	83.4	46.7	96.3	93.9	46.7	22.6	38.9	10.5	74.3	73.6	71.5	69.1
	COAT	87.4	53.3	94.9	92.6	49.9	23.7	43.6	12.6	67.1	66.6	63.6	61.7
	PSTR	88.1	50.0	90.4	77.7	51.0	23.7	44.9	11.8	69.9	61.7	63.2	47.9
0.	ADG+CIL	85.6	41.6	91.7	89.4	52.0	27.5	44.5	11.5	94.6	93.8	86.7	83.9

179 180 181

162

163

164

165 166 167

### 181 182 183

185

# 3 ROBUSTNESS ANALYSIS ON PERSON SEARCH

## 3.1 INVESTIGATION OF CORRUPTION ROBUSTNESS ON PERSON SEARCH

186 Benchmark Dataset Design. To evaluate the robustness of person search models, we propose two 187 benchmarks: CUHK-SYSU-C and PRW-C. These benchmarks are built upon the widely adopted 188 CUHK-SYSU (Xiao et al., 2017) and PRW (Zheng et al., 2017) datasets, which have been exten-189 sively used in numerous person search studies (Han et al., 2019; Chen et al., 2018; Li & Miao, 2021; 190 Kim et al., 2021). CUHK-SYSU dataset offers a diverse range of backgrounds from various urban scenarios, CUHK-SYSU-C retains this diversity while introducing corrupted scenes captured in 191 those environments. The PRW dataset includes images captured from six distinct viewpoints at each 192 location, and PRW-C preserves this multi-view feature, adding corruption across various angles. For 193 evaluation using our benchmarks, we utilize the same test splits from these datasets. 194

195 We extend the corruptions used in the pioneering study by Hendrycks & Dietterich (2019) to the 196 person search task. Considering the characteristics of urban outdoor environments in person search, we incorporate rain (Chen et al., 2021) and dark (Kong et al., 2024) as corruption types of our bench-197 marks. The corruption types of our benchmark include: 'gaussian noise', 'speckle noise', 'defocus blur', 'glass blur', 'motion blur', 'gaussian blur', 'snow', 'frost', 'fog', 'brightness', 'spatter', 'rain', 199 'dark', 'contrast', 'elastic', 'pixelate', 'jpeg compression', and 'saturate'. Each corruption is applied 200 at five distinct severity levels, with higher severity indicating greater image degradation (*i.e.*, more 201 severe corruption). The severity levels are established based on the traits of corruption scenarios ob-202 served in real-world. As an example, the rain corruption depicted in Figure 1 has various attributes 203 including slope of rain droplets, color of rain droplets, drop length & width, overall blurriness and 204 brightness to simulate different severity. Details for all corruptions are available in Appendix. 205

In person search, the goal is to match a given query image to its corresponding gallery image in a database; corruption can be applied to either the query or the gallery image. Query and gallery images may not always exhibit the same type or severity level of corruption. We construct the benchmarks by randomly applying different corruption types and severity levels to query and gallery images, similar to prior work in instance retrieval (Chen et al., 2021). We repeat the evaluation process five times and report the mean performance for experimental results. For benchmark statistics and further discussion of this design choice, see Appendix C and B, respectively.

212

Evaluation Models. To explore the corruption robustness of person search models, we employ six
seminal state-of-the-art models including both CNN and transformer architectures: OIMNet (Xiao et al., 2017), NAE (Chen et al., 2020), OIMNet++ (Lee et al., 2022), SeqNet (Li & Miao, 2021), COAT (Yu et al., 2022), and PSTR (Cao et al., 2022). More details of these models are provided



Figure 3: **Person search model performance across various corruption severities**, evaluated on the CUHK-SYSU-C benchmark.

in the appendix. For NAE, OIMNet++, SeqNet, PSTR, and COAT, we utilize the official checkpoints and configurations provided by the authors. Considering the evolution of training techniques
since OIMNet was published, we re-implement and report the results. We first train the model on
the training splits of CUHK-SYSU and PRW, then evaluate it on the CUHK-SYSU-C and PRW-C
benchmarks. Unless otherwise specified, SeqNet is used for most experiments, as it serves as a
baseline in several preceding works (Li et al., 2022; Jaffe & Zakhor, 2023; Li et al., 2023), thanks
to its simple yet effective architecture.

234 **Evaluation Metrics.** Given that person search inherently involves both detection and reID, we 235 report the performance metrics in two categories: detection and search performance. For search 236 performance, we adopt the widely used metrics of mean Average Precision (mAP) and Cumulative 237 Matching Characteristic at Rank-1 (R@1). For detection performance, we utilize Recall and Aver-238 age Precision (AP). Following existing corruption studies (Michaelis et al., 2019; Wang et al., 2021; 239 Schiappa et al., 2022), we also evaluate the relative performance drop caused by corruption, comparing the performance on the corrupted set to the clean set. This includes *relative mAP* (rmAP), 240 relative Rank-1 (rR@1), relative Recall (rRecall), and relative AP (rAP). For instance, relative mAP 241 (rmAP) is calculated by dividing the mAP from the corrupted set by the mAP from the clean set 242 (*i.e.*, rmAP = 'mAP from corrupted set'/'mAP from clean set'). 243

244 **Evaluation Results.** We evaluate the performance of state-of-the-art person search models under 245 corruption. Table 1 presents the performance of each model on both CUHK-SYSU (PRW) and 246 CUHK-SYSU-C (PRW-C). The results show a significant drop in performance under corruption, 247 with up to an 80% decline in mAP on PRW-C compared to the clean set. This reveals the vulnerability of current person search models to corruption and highlights the necessity of developing 248 corruption-robust person search models. We also evaluate these models across five different severity 249 levels of corruption. As shown in Figure 3, both the search and detection capabilities of person 250 search models decrease as the severity level increases. Specifically, we observe a 20% to 30% de-251 cline in search performance with each increase in severity, while detection performance degrades at 252 a more gradual pace. Performance begins to degrade from the first severity level, and when severity 253 reaches level 5, overall search performances drop to approximately 10% of mAP. 254

255 **Evaluation of Combining Existing Robustness Methods.** In the last row of the Table 1, we 256 further investigate whether combining existing corruption-robust detection and reID models can 257 achieve robust person search. Specifically, we combine two independent models: OADG (Lee et al., 258 2024), a robust detection model designed for corrupted environments, and CIL (Chen et al., 2021), 259 a reID model built for re-identification under corruption scenarios. For a fair comparison, we train 260 both models on the person search dataset. Specifically, OADG and CIL are initialized with ImageNet (Deng et al., 2009) pre-trained ResNet-50 (He et al., 2016) and trained on the clean CUHK-261 SYSU (PRW) dataset. For evaluation, OADG is first used to detect pedestrians on CUHK-SYSU-C 262 (PRW-C), and the detected individuals are then input into CIL to extract person representations, 263 following the standard protocol of the two-step approach. The results in Table 1 demonstrate that 264 this simple integrating approach is not sufficient to achieve robustness to corruption. Although it 265 shows the best rR@1 and rmAP on PRW-C, its overall performance on the corrupted benchmarks, 266 especially in R@1 and mAP, remains less effective compared to other person search models. 267

In summary, these results indicate that, despite the corruption-robust design of these methods, the simple integration lacks robustness for the multi-task nature of person search. This suggests the need for methods that enhance corruption robustness while considering the unique aspects of this

284

286

287 288

289

Table 2: Individual evaluation of detection and reID stage of end-to-end person search framework
on CUHK-SYSU-C. 'Representation against Corruption' refers to the search performance when extracting person representation in corrupted images, assuming that detection was performed on clean
images. 'Detection against Corruption' denotes the search performance when detection is conducted
on corrupted images, coupled with the representations that are extracted from clean images.

	Repres	sentatior	n against C	Corruption	Detection against Corruption				
Severity	Search		Detection		Search		Detection		
	R@1	mAP	Recall	Recall AP R		mAP	Recall	AP	
Oracle	95.2	94.5	100.0	100.0	95.2	94.5	100.0	100.0	
Clean	94.6	93.7	92.1	89.2	94.6	93.7	92.1	89.2	
Level 1	79.1	77.6	92.1	89.2	93.5	92.2	86.2	83.6	
Level 2	60.5	59.5	92.1	89.2	91.7	89.1	79.2	76.8	
Level 3	42.6	42.5	92.1	89.2	88.6	84.8	71.2	68.7	
Level 4	28.8	28.1	92.1	89.2	82.1	74.4	59.6	57.1	
Level 5	19.3	18.6	92.1	89.2	69.8	58.9	46.0	44.0	

task. In the following sections, we explore the underlying reasons for corruption vulnerabilities in the person search for developing a method tailored to its unique challenges.

## 3.2 SENSITIVITY TO CORRUPTION IN DETECTION AND REPRESENTATION STAGES

The experiments in the previous section demonstrate that existing person search models are vulnerable to corruption. In a typical person search framework, the *detection* head first identifies person candidates, and then the reID head extracts *representations* from the detected regions. Therefore, we assess the individual sensitivity of both the detection and representation stages to corruption.

To evaluate detection sensitivity to corruption, we first apply the detection head to corrupted images to obtain predicted bounding boxes. We then extract features from the corresponding regions in clean images using these box coordinates and perform the search process. Conversely, to evaluate the representation sensitivity to corruption, we perform detection on clean images to obtain predicted boxes, then extract features from the corresponding regions in corrupted images, followed by the search process. We evaluate the model across five severity levels, providing a separate benchmark for each to observe how the model responds to increasing corruption severity.

301 The results of these individual sensitivity studies are in Table 2. The 'Oracle' row represents per-302 formance using ground truth bounding boxes, while the 'Clean' row refers to the results where both 303 detection and search are performed on clean images. The five remaining rows (Level 1 through Level 304 5) show the impact of increasing the severity level of corruption. The results reveal that performance 305 degrades more significantly when person representations are extracted from corrupted images. At level 5 of 'Detection against Corruption', although detection performance declines, search perfor-306 mance remains relatively stable, with 70% of R@1, indicating that the model retains the search 307 capability despite noisy detection results. However, when person representations are corrupted (*i.e.*, 308 'Representation against Corruption'), search performance drops significantly, even with relatively 309 strong detection performance. These findings reveal that both detection and representation stages are 310 sensitive to corruption, further highlighting that the *representation stage is particularly vulnerable*. 311

312

#### 312 3.3 INFLUENCE OF FOREGROUND & BACKGROUND FOR ROBUST PERSON REPRESENTATION 313

Unlike person re-identification, which relies on cropped images to extract person representations, person search processes entire scenes and thus uses larger receptive fields to extract representations. As a result, the background can significantly impact person search performance. To investigate the individual influence of corruption on the foreground and background, we selectively apply corruption to either region, as illustrated in Figure 4. We then measure search performance based on the representations obtained from the corrupted images across five severity levels.

Figure 4 presents search and detection performance under three corruption scenarios: corruption
 in the foreground ('fg corrupted'), corruption in the background ('bg corrupted'), and corruption
 on the entire scene ('both corrupted'). The results show that search performance remains largely
 unaffected when corruption is limited to the background. However, when corruption is applied
 to the foreground – directly impacting the person's appearance – there is a notable drop in both



Figure 4: **Performance evaluation on images from CUHK-SYSU where corruption is applied to either the background or the foreground.** 'bg-corrupted' indicates performance on images where corruption is applied only to the background, while 'fg-corrupted' refers to performance on images where corruption is applied only to the foreground (person) regions. 'both' represents performance on the CUHK-SYSU-C dataset, where corruption is applied to the entire image.

search and detection performance, comparable to the degradation observed when the entire scene is corrupted. This suggests that the integrity of the background plays a minor role in robust person representations, while *corruption on the foreground is much more detrimental*. Based on these observations, in the following section, we introduce a novel foreground augmentation approach with tailored regularization to achieve more robust person representations under corruption.

342 343

334

335

336

337 338

339

340

341

344 345

346

347

# 4 PROPOSED METHOD FOR ROBUST PERSON SEARCH TO CORRUPTION

### 4.1 FOREGROUND-AWARE AUGMENTATION

348 Data augmentation has proven to be a valuable approach for boosting model robustness in computer 349 vision (Shorten & Khoshgoftaar, 2019; Rebuffi et al., 2021; Liu et al., 2024b). However, as shown 350 in Figure 5 (b), applying data augmentation to the entire scene results in severe semantic corruption and unreliable bounding boxes (e.g., overlapping pedestrians). Our analysis highlights that search 351 performance remains largely unaffected when corruption is restricted to the background. To this end, 352 we introduce a foreground-aware augmentation to generate appropriate augmented counterparts for 353 given input images. We define the areas containing people as foreground and determine these using 354 ground truth bounding boxes during training. To minimize the role of ground truth bounding boxes 355 in generating shortcuts during training, we apply a translation to the box coordinates of a person 356 before applying augmentation to the corresponding region. 357

In our training process, we incorporate both clean scenes and their augmented counterparts. Let the 358 input image be denoted by  $x^c$ , and the set of people appearing in the scene be  $P_{x^c} = \{p_1^c, \ldots, p_{n_x}^c\}$ . 359 Here,  $n_x$  represents the total number of people appearing in the scene. We create  $x^a$  by applying  $\mathcal{T}$ 360 to each person, where  $\mathcal{T}$  represents the transformation by augmentation functions. The set of people 361 appearing in  $x^a$  is  $P_{x^a} = \{p_1^a, p_2^a, \dots, p_{n_{x^a}}^a\}$  and  $p_j^a = \mathcal{T}(p_j^c)$ . We employ Augmix (Hendrycks 362 et al., 2020) and random erasing (Zhong et al., 2020) for our augmentation functions. For a rigorous robustness evaluation, we exclude augmentations that operate on similar principles to the corruptions 364 used in creating CUHK-SYSU-C and PRW-C. With this in mind, the augmentations we use in our 365 method are rotation, shearing, translation, solarization, autocontrast, equalization, and posterization. 366 We use  $x^c$  and  $x^a$  to train the model with the OIM (Xiao et al., 2017) loss. Let  $z_i \in \{z | z \in$  $f(x^c;\theta) \cup f(x^a;\theta)$  be the normalized representation of a person  $p_i$  from a scene, where  $f(\cdot;\theta)$  is 367 the function for extracting a set of person candidates on the given scene. Let  $v_l$  be a vector in the 368 lookup table where  $l \in \{1, ..., L\}$ , and  $u_q$  be the q-th vector in the queue where  $q \in \{1, ..., Q\}$ . 369 This lookup table serves as a memory bank containing representations for labeled persons, while the 370 queue is a memory bank that stores representations for unlabeled persons. L and Q are the size of 371 the lookup table and queue. The OIM loss is adopted as follows: 372

373 374

$$\mathcal{L}_{OIM} = \mathbb{E}_z[-\log h_i],\tag{1}$$

$$b_{i} = \frac{\exp(v_i^\top z_i/\tau)}{\sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_$$

$$h_{i} = \frac{\exp(v_{i}^{\top} z_{i}/\tau)}{\sum_{l=1}^{L} \exp(v_{l}^{\top} z_{i}/\tau) + \sum_{q=1}^{Q} \exp(u_{q}^{\top} z_{i}/\tau)},$$
(2)





(a) Augmentation applied to foreground

(b) Augmentation applied to entire image

Figure 5: **Illustration of augmentation on foreground (a) and entire image (b)**. Unlike (a), naïve augmentation on the entire image can generate severe semantic perturbations. We define the fore-ground using ground-truth bounding boxes, and augmentation is applied to each area individually.

where  $\tau$  is a temperature parameter. Each representation in the lookup table is updated with the momentum parameter  $\eta$ :

$$v_l \leftarrow \eta z_i + (1 - \eta) v_l. \tag{3}$$

## 4.2 REGULARIZATION FOR ROBUST PERSON REPRESENTATION

In Section 3.2, we observed that the representation stage was proven to be more susceptible to cor-ruption than the detection stage. In this context, we propose a regularization method for person search to achieve robust person representation. Our goal is to achieve robust person representation through regularization by learning invariance between clean input images and their augmented ver-sions. In person search frameworks, models typically generate multiple detected results for each person. Our regularization approach exploits this feature by utilizing these multiple detected results. To limit the excessive contribution of low-quality results, we utilize detected results where the In-tersection over Union (IoU) with the ground truth bounding box is larger than  $\alpha$ . Let  $h_i^c$  and  $h_i^a$ be normalized representations obtained from the *j*-th person of original scene and its augmented counterpart. Using  $h_j^c$ , we define a mixture representation  $M_j$  for the *j*-th person. 

$$M_j = \sum_m s_{j,m} h_{j,m}^c, \tag{4}$$

For  $h^a$  corresponding to detection results satisfying the constraint, we apply the Kullback-Leibler divergence. The regularization loss  $L_{reg}$  is formulated as follows:

$$\mathcal{L}_{reg} = \frac{1}{n} \sum_{j} \sum_{r} s_{j,r} \mathrm{KL}[h_{j,r}^a || M_j],$$
(5)

415 where

$$s_{j,*} = \frac{\exp(s_{j,*}/\tau_{iou})}{\sum_{k}^{n} \exp(s_{j,k}/\tau_{iou})},$$
(6)

 $s_{j,*}$  denotes the IoU score of  $h_j^c$  or  $h_j^a$  and n indicates the number of detected results. Here, we detach  $M_j$ , aiming for invariant representation towards augmentation, and utilize the IoU scores  $s_j$  as a weight, considering each target's quality.

ł

Our proposed method can be seamlessly integrated into existing person search models. In the following section, we apply and validate our method to several state-of-the-art person search models.

## 5 EXPERIMENTS

Implementation Details. We use an ImageNet (Deng et al., 2009) pre-trained ResNet50 (He et al., 2016) as the backbone for all methods in our experiments. The queue size is set to 5000 and 500 for CUHK-SYSU and PRW, respectively. We use a gallery size of 100 when evaluating CUHK-SYSU-C. For PRW-C, the gallery size matches the total number of images in the test set. An initial learning rate of 0.003 is used. We employ SGD with a momentum of 0.9 and a weight decay of

Method	CUHK-SYSU CUHK-SYSU-C			PRW PRW-C								
	R@1	mAP	rR@1	rmAP	R@1	mAP	R@1	mAP	rR@1	rmAP	R@1	mAP
OIMNet (Xiao et al., 2017)	87.7	86.2	37.7	36.7	33.0	31.7	76.7	37.3	45.8	23.6	35.2	8.8
+Ours	90.2	89.3	54.0	53.0	48.7	47.4	78.0	40.7	54.3	32.0	40.8	13.0
NAE (Chen et al., 2020)	92.3	91.4	44.3	42.4	40.9	38.8	80.6	42.8	42.6	20.9	34.3	8.9
+Ours	94.3	93.7	65.9	63.7	62.2	59.7	81.1	44.0	56.6	35.2	45.9	15.5
OIMNet++ (Lee et al., 2022)	94.0	93.2	38.6	37.0	36.3	34.5	83.2	47.3	43.9	19.9	36.5	9.4
+Ours	94.2	93.8	62.5	61.2	58.8	57.4	84.5	48.1	57.4	33.8	47.8	16.3
SeqNet (Li & Miao, 2021)	94.5	93.8	46.2	44.4	43.6	41.6	83.4	46.7	46.7	22.6	38.9	10.5
+Ours	94.9	94.3	70.3	68.9	66.7	65.0	84.0	47.0	57.0	33.8	47.8	15.9
COAT (Yu et al., 2022)	94.7	94.2	52.6	50.4	49.8	47.5	87.4	53.3	49.9	23.7	43.6	12.6
+Ours	92.8	92.1	68.9	67.4	63.9	62.0	86.9	53.3	55.4	40.1	55.4	21.4

Table 3: **Evaluation of the proposed method on various person search models** across clean and corrupted benchmarks.

Table 4: Ablation study of our method on CUHK-
SYSU and CUHK-SYSU-C. Baseline refers to Se-
qNet (Li & Miao, 2021).

Method	CUHK	S-SYSU	CUHK-SYSU-C		
	R@1	mAP	R@1	mAP	
Baseline	94.5	93.8	43.6	41.6	
+ Regularization	94.2	93.4	55.6	53.7	
+ Foreground-aware	94.5	94.0	65.0	63.1	
+ IoU Score	94.9	94.3	66.7	65.0	

Table 5:	Evaluation	under	real	cor-
ruption s	scenarios.			

Method	Da	ark	Ra	Rain		
	R@1	mAP	R@1	mAP		
Baseline	33.3	34.7	31.6	35.9		
+ Ours	45.1	47.0	41.9	44.1		
Mathad	B	lur	F	og		
Method	B R@1	lur mAP	Fe R@1	og mAP		
Method Baseline	B R@1 64.0	lur mAP 64.0	F@1 56.8	og mAP 47.8		

0.0005. For the hyper-parameters used in the training, the threshold value  $\alpha$  for the IoU is set to 0.6. The momentum parameter  $\eta$  is set to 0.5, and a temperature value  $\tau$  is set to 0.2. The temperature parameter for the IoU score  $\tau_{iou}$  is set to 0.6. We use an NVIDIA RTX 3090 GPU for the experiments.

**Results.** The experimental results are shown in Table 3. Our method shows promising performance improvements under corruption for all tested models. Specifically, the R@1 for NAE in-creases by 52% on CUHK-SYSU-C and 34% on PRW-C with our approach. It should be noted that the impact on clean dataset performance appears to be limited, with variations within  $\pm 5\%$ . On the CUHK-SYSU dataset, our approach enhances the performance of all the original models except for COAT. COAT's method incorporates token mixup, which already offers inherent advantages in augmentation. These results suggest that our approach improves the robustness of person search models in corrupted environments without compromising their performance in clean conditions. 

Note that our method performs better than the combination of well-known robustness methods in
 detection and re-identification (OADG+CIL) shown in Table 1. This suggests that our method is
 highly competitive with existing approaches to handling corruption.

Ablation Study. We conduct an ablation study on CUHK-SYSU and CUHK-SYSU-C to evaluate the effectiveness of each technique in Table 4, using SeqNet as our baseline. The second row, labeled 'Regularization', represents the method discussed in Section 4.2 that excludes the IoU score. The third row, labeled 'Foreground-aware', details changes in augmentation application from the entire scene to just the foreground area, while the fourth row, labeled 'IoU Score', shows results incorporating the IoU score.

478
479 Our findings indicate that all techniques—'Regularization', 'Foreground-aware' and 'IoU Score'—enhance search performance in corrupted environments, highlighting the effectiveness of our proposed method in tackling the challenges of person search under corrupted conditions.

Validation of Proposed Method towards Real Corruptions. To assess our method's effective ness in real-world corruption cases, we collect images from BDD100K (Yu et al., 2020) and use them
 for our experiments. Since BDD100K is a comprehensive dataset acquired from tens of thousands
 of drivers across various locations, weather conditions, and time frames, it is commonly adopted to
 evaluate model robustness in corruption and associated research areas (Kim & Shin, 2024; Cygert

486 & Czyżewski, 2021; Kim & Shin, 2023). Although BDD100K is mainly used for detection and 487 segmentation tasks, as it includes a diverse array of weather and temporal conditions, we gather 488 scenes corresponding to several corruption scenarios from it and perform manual annotation. From 489 this dataset, we gather scenes representing dark, rain, blur, and fog scenarios, then manually an-490 notate the bounding boxes and identity labels of individuals present in these images. We adhere to the labeling protocol used in the creation of CUHK-SYSU. For evaluation, we use 100 images 491 from each of the dark, rain, blur, and fog scenarios (400 total), following the default gallery size 492 of the CUHK-SYSU evaluation protocol. Following the PRW evaluation protocol, we employ all 493 images except the one containing the query person as the gallery and use all labeled persons as a 494 query once. We train the models on the uncorrupted data (CUHK-SYSU) and evaluate their perfor-495 mances in each real corruption scenario. We make our annotations publicly available (annotation 496 only), which include the information on the samples and label data we use (download link in the 497 appendix). 498

Table 5 presents the model's performance evaluated across four real-world corruption scenarios.
Our proposed method demonstrates performance improvements across all real corruption scenarios.
Notably, in the dark corruption scenario, our method achieves a 12.3% mAP performance gain, showing the effectiveness of the proposed method leveraging our analyses. These results indicate that the proposed method can work in real-world corruption scenarios. Our method, which performs effectively under our proposed corruption, shows robustness in real-world corruption scenarios as well. See the appendix for a qualitative comparison of the proposed and real-world corruptions.

506 Frequency Sensitivity Analysis. To provide theoretical depth to our experiments, 507 we analyze our proposed method from the 508 perspective of frequency sensitivity. The 509 25x25 heatmap in Figure 6 represents the er-510 ror rate (mAP) when evaluating the model us-511 ing data corrupted with different frequency 512 bases. The edges of the heatmap represent 513 evaluations using data corrupted with high-514 frequency bases, while the center represents 515 experiments performed with low-frequency 516 bases. Red colors indicate higher error rates, 517 while blue colors indicate lower error rates. Compared to the left heatmap (Baseline), the 518 right heatmap (Ours) shows overall perfor-519 mance improvement across low, middle, and 520 high frequencies. Yin et al. (2019) analyzed 521



Figure 6: **Qualitative analysis of frequency sensitivity** between baseline (SeqNet) and our method.

the common corruptions in the frequency domain, showing that fog and contrast have relatively low-frequency components, noise-related corruptions have high-frequency components, and blur and pixelate have middle-frequency components. We think that our model's performance improvement on corruption benchmarks stems from its resilience to various frequency perturbations.

526

527

# 6 CONCLUSION

In this work, we present two benchmarks for evaluating the robustness of person search, CUHK-SYSU-C and PRW-C, to assess and analyze the robustness of person search models. We explore how various features of person search influence robustness with the expectation where our findings will be valuable lessons to the research community in related fields. From the findings obtained through our experiments, we propose a solution for robust person search that not only achieves competitive performance on clean datasets, but also demonstrates effective robustness enhancement in corrupted environments.

Potential Limitations. While our method performs well in simulated environments, we have not
 extensively tested its performance under real-world corruption scenarios. To mitigate this gap, we
 collect and evaluate our method on real corruption images, which shows the validity of our proposed
 approach in the real world. We have not considered scenarios with multiple simultaneous corrup tions. Additionally, while we categorize corruption severity into five levels, this discretization may
 not fully capture the continuous nature of real-world corruption.

# Ethics Statement

-This section is not included in the page limit.-

#### 7 SOCIAL IMPACT

546 The development and deployment of person search systems require careful consideration due to their 547 potential impact on individual privacy and societal norms. While these systems offer substantial ben-548 efits for security and surveillance by enhancing the ability to locate and identify individuals across 549 diverse environments, they simultaneously pose a profound risk to personal privacy. The ability to 550 track and identify individuals without their consent can lead to a range of privacy violations, from 551 the unwarranted monitoring of public movements to the potential for misuse in stalking and harass-552 ment. It is, therefore, imperative that developers, implementers, and policymakers involved in the creation and use of person search systems consider these ethical implications from the outset. The 553 development of these systems should be guided by ethical principles that prioritize the well-being 554 and privacy of individuals, incorporating mechanisms for accountability and oversight to prevent 555 misuse. It is our collective responsibility to balance innovation with the imperative to safeguard 556 human dignity and privacy.

558 By doing the above, we can harness the benefits of person search systems while mitigating the risks, 559 ensuring these technologies enhance societal welfare without compromising individual freedoms. For instance, person search technology can expedite victim identification and rescue efforts in natu-560 ral disasters or accidents, improving the effectiveness of emergency responses and potentially saving 561 lives. Disaster environments often involve challenges like heavy rainfall or snowfall, while accident 562 scenes may present issues such as varying lighting conditions or spatter-covered camera lenses. 563 Our study aims to optimize person search techniques to function effectively under these varied and adverse conditions.

565 566 567

540

541 542

543 544

#### FURTHER DISCUSSIONS ABOUT DATASETS 8

568 569 570

Privacy. Testing person search models in diverse situations and environments requires collecting new data for various situations with people in them, which can raise ethical issues. Since we provide 18 different scenarios for evaluating the person search models, our benchmarks serve as a good proxy 571 for evaluating models across diverse environments without collecting new data. The parent datasets 572 CUHK-SYSU (Xiao et al., 2017), PRW (Zheng et al., 2017) expose people's faces, and our datasets 573 (CUHK-SYSU-C, PRW-C) inevitably inherit this issue. To mitigate the issue, we mask faces that 574 are captured prominently and distinctly visible for the figures used in the paper to ensure their 575 identities are not recognizable. Future researchers using the proposed CUHK-SYSU-C and PRW-C 576 must be aware of this and take precautionary measures. When using them in research papers, we 577 strongly recommend to de-identify individuals when their faces are clearly recognizable. Regarding 578 the subset collected from BDD100K for Section 5, we ensure that the samples we utilize are from a data source that has already been publicly available, carefully adhering to its licensing terms. 579

580

License. We will release them in the form of code to generate our benchmarks from the parent datasets(CUHK-SYSU, PRW) rather than as raw images. The authors of the parent datasets provide 582 guidelines for dataset usage on their respective websites, which are similar to the terms of CC BY-583 NC: For CUHK-SYSU, users are permitted to use the data only for non-commercial research and 584 educational purposes. Users are not permitted to distribute the data. For PRW, they emphasize the 585 purpose of non-commercial research applications. The dataset requires citation. For these reasons, 586 we will also release CUHK-SYSU-C and PRW-C under CC BY-NC.

- 587 588

- 592

#### 594 REFERENCES 595

601

618

632

640

- Jiale Cao, Yanwei Pang, Rao Muhammad Anwer, Hisham Cholakkal, Jin Xie, Mubarak Shah, and 596 Fahad Shahbaz Khan. Pstr: End-to-end one-step person search with transformers. In *Proceedings* 597 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9458–9467, 2022. 598
- Di Chen, Shanshan Zhang, Wanli Ouyang, Jian Yang, and Ying Tai. Person search via a mask-600 guided two-stream cnn model. In Proceedings of the european conference on computer vision (ECCV), pp. 734–750, 2018. 602
- Di Chen, Shanshan Zhang, Jian Yang, and Bernt Schiele. Norm-aware embedding for efficient per-603 son search. In Proceedings of the IEEE/CVF conference on computer vision and pattern recogni-604 tion, pp. 12615–12624, 2020. 605
- 606 Minghui Chen, Zhiqiang Wang, and Feng Zheng. Benchmarks for corruption invariant person re-607 identification. Advances in Neural Information Processing Systems, 2021. 608
- Shuo Chen, Jindong Gu, Zhen Han, Yunpu Ma, Philip Torr, and Volker Tresp. Benchmarking 609 robustness of adaptation methods on pre-trained vision-language models. Advances in Neural 610 Information Processing Systems, 36, 2024. 611
- 612 Sebastian Cygert and Andrzej Czyżewski. Robustness in compressed neural networks for object 613 detection. In 2021 International Joint Conference on Neural Networks, pp. 1–8. IEEE, 2021. 614
- 615 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-616 erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, 617 pp. 248–255. Ieee, 2009.
- Wenkai Dong, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Instance guided proposal network 619 for person search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 620 Recognition, pp. 2585–2594, 2020. 621
- 622 Yinpeng Dong, Caixin Kang, Jinlai Zhang, Zijian Zhu, Yikai Wang, Xiao Yang, Hang Su, Xingxing 623 Wei, and Jun Zhu. Benchmarking robustness of 3d object detection to common corruptions. 624 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1022-1032, 2023. 625
- 626 Samuel Dooley, George Z Wei, Tom Goldstein, and John Dickerson. Robustness disparities in face 627 detection. Advances in Neural Information Processing Systems, 35:38245–38259, 2022. 628
- 629 Byeong-Ju Han, Kuhyeun Ko, and Jae-Young Sim. End-to-end trainable trident person search net-630 work using adaptive gradient propagation. In Proceedings of the IEEE/CVF International Con-631 ference on Computer Vision, pp. 925–933, 2021.
- Chuchu Han, Jiacheng Ye, Yunshan Zhong, Xin Tan, Chi Zhang, Changxin Gao, and Nong Sang. 633 Re-id driven localization refinement for person search. In Proceedings of the IEEE/CVF Interna-634 tional Conference on Computer Vision, pp. 9814–9823, 2019. 635
- 636 Chuchu Han, Zhedong Zheng, Kai Su, Dongdong Yu, Zehuan Yuan, Changxin Gao, Nong Sang, 637 and Yi Yang. Dmrnet++: Learning discriminative features with decoupled networks and enriched 638 pairs for one-step person search. IEEE Transactions on Pattern Analysis and Machine Intelli-639 gence, 45(6):7319-7337, 2022.
- Haodong He, Jian Ding, and Gui-Song Xia. On the robustness of object detection models in aerial 641 images. arXiv preprint arXiv:2308.15378, 2023. 642
- 643 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-644 nition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 645 770-778, 2016.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common cor-647 ruptions and perturbations. In International Conference on Learning Representations, 2019.

648 649 650	Dan Hendrycks, Norman Mu, Ekin Dogus Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi- narayanan. Augmix: A simple method to improve robustness and uncertainty under data shift. In <i>International conference on learning representations</i> , volume 1, pp. 5, 2020.
651 652 653	Lucas Jaffe and Avideh Zakhor. Gallery filter network for person search. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pp. 1684–1693, 2023.
654 655 656 657	Christoph Kamann and Carsten Rother. Benchmarking the robustness of semantic segmentation models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 8828–8838, 2020.
658 659 660	Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 18963–18974, June 2022.
661 662 663 664	Hanjae Kim, Sunghun Joung, Ig-Jae Kim, and Kwanghoon Sohn. Prototype-guided saliency feature learning for person search. In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4863–4872, 2021.
665 666 667	Youngjun Kim and Jitae Shin. Robust object detection against multi-type corruption without catastrophic forgetting during adversarial training under harsh autonomous-driving environments. <i>IEEE Access</i> , 11:26862–26876, 2023.
668 669 670 671	Youngjun Kim and Jitae Shin. Efficient and robust object detection against multi-type corruption using complete edge based on lightweight parameter isolation. <i>IEEE Transactions on Intelligent Vehicles</i> , 2024.
672 673 674	Lingdong Kong, Shaoyuan Xie, Hanjiang Hu, Lai Xing Ng, Benoit Cottereau, and Wei Tsang Ooi. Robodepth: Robust out-of-distribution depth estimation under corruptions. <i>Advances in Neural</i> <i>Information Processing Systems</i> , 36, 2024.
675 676 677	Xu Lan, Xiatian Zhu, and Shaogang Gong. Person search by multi-scale matching. In <i>Proceedings</i> of the European conference on computer vision (ECCV), pp. 536–552, 2018.
678 679 680	Sanghoon Lee, Youngmin Oh, Donghyeon Baek, Junghyup Lee, and Bumsub Ham. Oimnet++: Prototypical normalization and localization-aware learning for person search. In <i>European Conference on Computer Vision</i> , pp. 621–637. Springer, 2022.
681 682 683 684	Wooju Lee, Dasol Hong, Hyungtae Lim, and Hyun Myung. Object-aware domain generalization for object detection. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 2947–2955, 2024.
685 686 687	Lucie Leveque, François Villoteau, Emmanuel V. B. Sampaio, Matthieu Perreira Da Silva, and Patrick Le Callet. Comparing the robustness of humans and deep neural networks on facial expression recognition. <i>Electronics</i> , 11(23), 2022.
688 689 690	Junjie Li, Yichao Yan, Guanshuo Wang, Fufu Yu, Qiong Jia, and Shouhong Ding. Domain adaptive person search. In <i>European Conference on Computer Vision</i> , pp. 302–318. Springer, 2022.
691 692 693 694	Junjie Li, Guanshuo Wang, Yichao Yan, Fufu Yu, Qiong Jia, Jie Qin, Shouhong Ding, and Xiaokang Yang. Generalizable person search on open-world user-generated video content. <i>arXiv preprint arXiv:2310.10068</i> , 2023.
695 696	Zhengjia Li and Duoqian Miao. Sequential end-to-end network for efficient person search. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 2011–2019, 2021.
697 698 699 700	Xiangtan Lin, Pengzhen Ren, Yun Xiao, Xiaojun Chang, and Alex Hauptmann. Person search challenges and solutions: A survey. <i>Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21</i> , 2021.
701	Xiaoqiong Liu, Yunhe Feng, Shu Hu, Xiaohui Yuan, and Heng Fan. Benchmarking the robustness

702 703 704 705	Zhendong Liu, Jie Zhang, Qiangqiang He, and Chongjun Wang. Understanding data augmentation from a robustness perspective. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pp. 6760–6764. IEEE, 2024b.
705 706 707 708	Xiaofeng Mao, Yuefeng Chen, Yao Zhu, Da Chen, Hang Su, Rong Zhang, and Hui Xue. Coco-o: A benchmark for object detectors under natural distribution shifts. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 6339–6350, 2023.
709 710 711	Claudio Michaelis, Benjamin Mitzkus, Robert Geirhos, Evgenia Rusak, Oliver Bringmann, Alexan- der S Ecker, Matthias Bethge, and Wieland Brendel. Benchmarking robustness in object detec- tion: Autonomous driving when winter is coming. <i>arXiv preprint arXiv:1907.07484</i> , 2019.
712 713 714 715	Eric Mintun, Alexander Kirillov, and Saining Xie. On interaction between augmentations and corruptions in natural corruption robustness. <i>Advances in Neural Information Processing Systems</i> , 34:3571–3583, 2021.
716 717 718	Bharti Munjal, Sikandar Amin, Federico Tombari, and Fabio Galasso. Query-guided end-to-end person search. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 811–820, 2019.
719 720 721	Jie Qin, Peng Zheng, Yichao Yan, Rong Quan, Xiaogang Cheng, and Bingbing Ni. Movienet- ps: a large-scale person search dataset in the wild. In <i>ICASSP 2023-2023 IEEE International</i> <i>Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pp. 1–5. IEEE, 2023.
722 723 724 725	Sylvestre-Alvise Rebuffi, Sven Gowal, Dan Andrei Calian, Florian Stimberg, Olivia Wiles, and Timothy A Mann. Data augmentation can improve robustness. <i>Advances in Neural Information Processing Systems</i> , 34:29935–29948, 2021.
726 727 728	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 39(6):1137–1149, 2016.
729 730 731	Ujjwal Saxena. Automold-road augmentation library. https://github.com/ UjjwalSaxena/AutomoldRoad-Augmentation-Library, 2023.
732 733 734	Madeline Schiappa, Shruti Vyas, Hamid Palangi, Yogesh Rawat, and Vibhav Vineet. Robustness analysis of video-language models against visual and language perturbations. <i>Advances in Neural Information Processing Systems</i> , 35:34405–34420, 2022.
735 736 737 738	Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local- ization. In <i>Proceedings of the IEEE International Conference on Computer Vision (ICCV)</i> , Oct 2017.
739 740 741	Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. Journal of big data, 6(1):1–48, 2019.
742	A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
743 744 745 746	Cheng Wang, Bingpeng Ma, Hong Chang, Shiguang Shan, and Xilin Chen. Tcts: A task-consistent two-stage framework for person search. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 11952–11961, 2020.
747 748 749	Jiahang Wang, Sheng Jin, Wentao Liu, Weizhong Liu, Chen Qian, and Ping Luo. When human pose estimation meets robustness: Adversarial algorithms and benchmarks. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 11855–11864, 2021.
750 751 752 753	Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. Joint detection and identification feature learning for person search. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 3415–3424, 2017.
754 755	Yuanlu Xu, Bingpeng Ma, Rui Huang, and Liang Lin. Person search in a scene by jointly model- ing people commonness and person uniqueness. In <i>Proceedings of the 22nd ACM international</i> <i>conference on Multimedia</i> , pp. 937–940, 2014.

756 757 758	Yichao Yan, Jinpeng Li, Jie Qin, Peng Zheng, Shengcai Liao, and Xiaokang Yang. Efficient person search: An anchor-free approach. <i>International Journal of Computer Vision</i> , 131(7):1642–1661, 2023.
759 760 761	Chenyu Yi, Siyuan Yang, Haoliang Li, Yap-peng Tan, and Alex Kot. Benchmarking the robustness of spatial-temporal models against corruptions. <i>arXiv preprint arXiv:2110.06513</i> , 2021.
762 763 764	Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective on model robustness in computer vision. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
765 766 767 768 769	Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madha- van, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learn- ing. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 2636–2645, 2020.
770 771 772	Rui Yu, Dawei Du, Rodney LaLonde, Daniel Davila, Christopher Funk, Anthony Hoogs, and Brian Clipp. Cascade transformers for end-to-end person search. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 7267–7276, 2022.
773 774 775 776	Xinyu Zhang, Xinlong Wang, Jia-Wang Bian, Chunhua Shen, and Mingyu You. Diverse knowledge distillation for end-to-end person search. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 3412–3420, 2021a.
777 778 779	Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. Fairmot: On the fairness of detection and re-identification in multiple object tracking. <i>International journal of computer vision</i> , 129:3069–3087, 2021b.
780 781 782	Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian. Person re-identification in the wild. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1367–1376, 2017.
783 784 785 786 787	Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmen- tation. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pp. 13001– 13008, 2020.
788 789	
790 791 792	
793 794	
795 796 707	
797 798 799	
800 801	
803	

# Appendix

#### FURTHER DISCUSSIONS А

813 814 815

816

826 827

828

829

836 837

841

847

850

851

852

810

811 812

#### A.1 FURTHER DISCUSSION OF BENCHMARK DATASET DESIGN IN SECTION 3.1

Table 6 presents the case study that highlights a potential bias in person search models when the 817 same type of corruption is applied to both query and gallery images. The table compares two sce-818 narios: the case where corruption is applied only to the gallery images ('Gallery corrupted'), and the 819 other case where the same corruption is applied to both query and gallery images ('Both corrupted'). 820 Four different types of corruptions are examined: Brightness, Contrast, Saturate, and Spatter. In-821 terestingly, for four corruption types, we observe higher performance when the same corruption is 822 applied to both query and gallery images. This pattern suggests that when both query and gallery 823 images are subjected to the same type of corruption, the model's ability to match them could im-824 prove. We consider this potential bias when designing our benchmark dataset, which is to randomly 825 apply corruption types and severity levels to query and gallery images.

Table 6: Case studies when the same corruption is applied to both query and gallery images, it could lead to another bias, such that the similarity of the two images increases. CUHK-SYSU-C, severity level 5, SeqNet are used for the case studies, R@1 is used as an evaluation metric.

Corruptions	Brightness	Contrast	Saturate	Spatter
Gallery corrupted	71.5	11.7	55.9	60.9
Both corrupted	77.8	12.8	62.8	63.6

## A.2 MODEL EXPLANATION IN SECTION 3.1

In this section, we describe person search models that we examine in Section 3.1. These seminal 838 state-of-the-art works have contributed to various aspects of the person search field. OIMNet (Xiao 839 et al., 2017) firstly proposes an end-to-end person search framework by jointly training the detection 840 and reID head. NAE (Chen et al., 2020) addresses the issue of conflicting learning objectives, a common challenge in person search and related fields (Xu et al., 2014; Zhang et al., 2021b; Lin et al., 2021), that arise during the joint learning of the detector and reID head. SeqNet (Li & Miao, 842 2021) improves the quality of detection results through a stronger detection head by considering 843 that the detection result influences the training of the reID head. This simple yet effective concept 844 has prompted subsequent studies to adopt its design (Li et al., 2022; Jaffe & Zakhor, 2023; Li 845 et al., 2023). OIMNet++ (Lee et al., 2022) improves the widely used OIM (Xiao et al., 2017) loss 846 and considers the quality of detection results in the training of the reID head. PSTR (Cao et al., 2022) and COAT (Yu et al., 2022) enable recent one-step person search frameworks to leverage the 848 advantages of the Transformer (Vaswani, 2017). 849

#### A.3 FURTHER DISCUSSION AND VALIDATION FOR EXPERIMENTS IN SECTION 3.2 AND SECTION 3.3

853 In Section 3.2 and 3.3 of the main paper, we analyzed which stage (representation or detection) is 854 more susceptible to corruption and examined the influence of foreground and background regions 855 on robust person representation. To validate these analyses and provide further insights, we conduct experiments that simultaneously examine both aspects. We analyze four scenarios: 'bg-corrupted & 856 Representation against Corruption', 'bg-corrupted & Detection against Corruption', 'fg-corrupted 857 & Representation against Corruption', and 'fg-corrupted & Detection against Corruption'. This 858 experimental design allows us to examine the relative susceptibility of detection and representation 859 stages under both background-only and foreground-only corruptions. 860

861 Table 7 presents our experimental results, revealing several key findings: In the case of 'fg-corrupted & Representation against Corruption', search performances align with the trends observed in both 862 Figure 4's 'fg-corrupt' and Table 2's 'Representation against Corruption' cases; The performance 863 from 'fg-corrupted & Detection against Corruption' shows similar the patterns seen in Figure 4's

866

867

868

877 878 879

883

884

885 886

887

Table 7: Analysis of corruption effects on detection and representation stages under foreground 865 and background corruptions. The terms 'bg-corrupted', 'fg-corrupted', 'Representation against Corruption', and 'Detection against Corruption' have the same meanings as those used in Table 2 and Figure 4.

	bg-corrupted								fg-corrupted							
-	a	Repres gainst C	entation Corruptic	on	Detection against Corruption			Representation against Corruption			Detection against Corruption					
	Sea	rch	Dete	ction	Sea	ırch	Deteo	ction	Sea	irch	Deteo	ction	Sea	rch	Dete	ction
severity	R@1	mAP	recall	AP	R@1	mAP	recall	AP	R@1	mAP	recall	AP	R@1	mAP	recall	AP
oracle clean	95.2 94.6	94.5 93.7	100.0 92.1	100.0 89.2	95.2 94.6	94.5 93.7	100.0 92.1	100.0 89.2	95.2 94.6	94.5 93.7	100.0 92.1	100.0 89.2	95.2 94.6	94.5 93.7	100.0 92.1	100.0 89.2
level1 level2 level3 level4	93.8 93.7 93.7 92.6 02.3	93.1 92.8 92.6 92.0	92.1 92.1 92.1 92.1 92.1	89.2 89.2 89.2 89.2	94.6 94.6 94.5 94.5	93.7 93.7 93.7 93.6 03.6	91.5 90.9 90.3 89.6	89.2 89.1 88.9 88.3	73.3 54.9 42.5 29.6	72.5 52.5 40.8 28.2	92.1 92.1 92.1 92.1 92.1	89.2 89.2 89.2 89.2	93.3 90.5 86.0 80.0 70.5	91.8 87.7 80.0 71.2	85.2 77.4 65.8 55.2	81.6 73.0 60.8 49.4 27.8

'fg-corrupt' and Table 2's 'Detection against Corrupt scenarios'; Both search and detection performances under background corruption maintain tendencies similar to those obtained with clean set. These results are consistent with previous findings, validating our observations.

#### В FURTHER ANALYSIS

HYPERPARAMETER ANALYSIS **B**.1

889 890

We conduct experiments to analyze four hyperparameters: the temperature for IoU score  $\tau_{iou}$ , types 891 of augmentation used in  $\mathcal{T}, \tau$ , and the threshold  $\alpha$ . Table 8 shows the results of our analysis on 892 the effect of  $\tau_{iou}$ . In this analysis, 'reverse' refers to the effect of our  $\tau_{iou}$  applied in reverse. 893 We implement this by applying  $1 - \tau_{iou}$ , which inverts the shape of the IoU distribution, before 894 applying  $\tau_{iou}$ . The 'uniform' case represents an uniform IoU distribution, equivalent to not using 895 the IoU score at all. Our results show that performance improves when using the IoU score (with 896  $\tau_{iou}$  values of 0.4, 0.6, and 0.8) compared to the 'uniform' case where IoU is not used. Moreover, 897 the 'uniform' case outperforms the 'reverse' case, which aligns with our intuition. The performance 898 remains relatively stable across different hyperparameter values when using the IoU score.

899 Table 9 presents our analysis of the three hyperparameters: types of augmentation used in  $\mathcal{T}, \tau$ , 900 and  $\alpha$ . For the analysis of augmentation  $\mathcal{T}$ , 'color' represents the evaluation result of a model 901 trained using solarization, autocontrast, equalization, and posterization, while 'geometric' refers to 902 the result of a model trained using rotation, shearing, and translation. The results indicate that the 903 'both' setting, which combines color and geometric augmentations, yields the best performance. 904 This suggests that both types of augmentation contribute positively, with geometric augmentations 905 showing more effectiveness. Regarding  $\tau$  and  $\alpha$ , we set these parameters as values of 0.2 and 0.6, respectively, which shows the best performance. 906

907 908

909

010

Table 8: Analysis for  $\tau_{iou}$  impact. The number indicates the value set for  $\tau_{iou}$ , uniform denotes not to apply the concept of IoU score, *reverse* indicates applying it in reverse.

310			
911 912	$ au_{iou}$	CUHK- R@1	-SYSU-C mAP
913	reverse	64.2	62.4
014	uniform	65.0	63.1
914	0.4	66.2	64.2
915	0.6	66.7	65.0
916	0.8	66.1	64.4
917			

Table 9: Analysis	for im	pacts of $\mathcal{T}_{i}$	, $ au$ , and $lpha$ .
-------------------	--------	----------------------------	------------------------

Type	Settings	CUHK-SYSU-C		Type	Settings	CUHK-SYSU-C		Type	Settings	CUHK	-SYSU-C
Type	bettings	R@1	mAP	Type	Settings	R@1	mAP	Type	Settings	R@1	mAP
	color	54.5	52.3		0.1	63.6	61.6		0.5	64.8	63.2
$\mathcal{T}$	geometric	57.5	55.3	au	0.2	66.7	65.0	$\alpha$	0.6	66.7	65.0
	both	66.7	65.0		0.3	63.2	62.0		0.7	64.2	62.0



Figure 7: Comparison of proposed corruptions with real-world corruption scenarios for dark, blur, fog, and rain conditions.

## B.2 QUALITATIVE RESULTS

We present the qualitative results of evaluating different models on corruption scenarios. The results on PRW-C are shown in Figure 9. Each row represents the results of different models, and each column shows different query examples and the search results of each model accordingly. Blue indicates the location of the query person, red represents incorrect detection results by the model, and green represents correct detection results. By looking at examples 1 to 3, where our model successfully retrieves a person while other models fail, we can see the efficacy of our method in extracting robust representations of the person when corruption is applied. We also provide the result of the real-world corruption data (introduced in Section 5) in Figure 10.

949 950 951

918 919

927 928

938 939 940

941 942

943

944

945

946

947

948

## B.3 QUALITATIVE COMPARISON BETWEEN PROPOSED AND REAL-WORLD CORRUPTIONS

952 In this section, we compare our proposed corruptions with real-world corruptions gathered as de-953 scribed in Section 5. Figure 7 illustrates four corruption scenarios: dark, blur, rain, and fog. Our 954 proposed corruption can capture distinct features in real-world corruption scenes. For example, 955 the proposed dark reflects the reduced overall brightness observed in real-world blur scenes, the 956 proposed blur presents decreased sharpness, and motion traces observed in real-world blur scenes. 957 While our corruption benchmark assumes the presence of only one type of corruption in a scene, 958 corruption in real-world images can be more complex. As shown in Figure 7, the real-world fog image contains some raindrops as well as a fog effect, while the real-world rain image exhibits blur-959 riness. We aim to investigate these multiple corruption scenarios in our future research endeavors. 960 The annotations we use for real-world corruption samples can be accessed in this link<sup>1</sup>. 961

962 963

## B.4 PERFORMANCE COMPARISON ACROSS VARIOUS SEVERITY LEVELS

We conduct the experiment to evaluate the robustness of our proposed method across various severity levels when applied to five different person search models. The graph in Figure 8 illustrates the mAP performance of five person search models enhanced with our approach, compared against the OADG+CIL combination, across five levels of corruption severity. The graph shows a general downward trend in performance as severity increases. SeqNet+ours and COAT+ours show strong resilience, maintaining high mAP scores even at the most severe corruption levels. We observe

<sup>&</sup>lt;sup>1</sup>Url:https://drive.google.com/drive/folders/13z7nn9gesSTzXHKKSNk131zxvYMRTBoL? usp=drive\_link



Figure 8: (a) Evaluation of our method across various severities. We apply our method to existing five different person search models. OADG+CIL refers to the combination of existing works mentioned previously. (b) Qualitative analysis of foreground-aware augmentation. We compare the contents in bounding box regions after augmentation applied based on two regional criteria: entire image and foreground region. The last two columns denote the Grad-CAM analysis of two types of augmentations on two corruption scenarios (spatter, fog).

that even the basic OIMNet model, when enhanced with our method, consistently surpasses the performance of the OADG+CIL combination. This underscores the effectiveness of our approach, 995 even when applied to simpler baseline models. 996

## **B.5** FURTHER ANALYSIS ABOUT FOREGROUND-AWARE AUGMENTATION

999 To provide a deeper insight into our proposed foreground-aware augmentation, we present addi-1000 tional qualitative analysis. Figure 8 presents the results with and without foreground-aware strategy, along with corresponding Grad-CAM (Selvaraju et al., 2017) visualization results. The 'Entire Im-1001 age Augmented' represents cases without foreground-aware strategy, while the 'Foreground-aware 1002 Augmented' shows cases with its application. The third, fourth, and fifth columns in each row show 1003 corrupted input images and their corresponding Grad-CAM results under two different corruptions 1004 (spatter, fog). We observe the bounding box regions cropped from the full images. In the first col-1005 umn, we can see that the naïve use could lead to the problem, such as severe semantic perturbation 1006 and unreliable bounding boxes. In contrast, the second column shows that our foreground-aware 1007 augmentation successfully applies augmentation while ensuring discriminative parts of the person 1008 remain within the bounding box. Accordingly, the results obtained from our strategy in the fifth 1009 column show that the model better captures the person's discriminative information. 1010

1011 С **BENCHMARK DETAILS** 1012

987

988

989

990

991

992 993

994

997

998

1020

1013 C.1 **BENCHMARK STATISTICS** 1014

1015 Table 5 presents the statistics of CUHK-SYSU-C and PRW-C. The statistics regarding images, iden-1016 tities, and pedestrians are the same as those for the test split of CUHK-SYSU and PRW. As men-1017 tioned in the main paper, each image in CUHK-SYSU-C and PRW-C is randomly assigned to one of the 18 corruptions and randomly assigned to one of the five severity levels. We repeat this process 1018 several times (5) and report the average performance. 1019

	# images	# identities	# pedestrians	# images per corruption	# images per se
CUHK-SYSU-C	6,978	2,900	40,871	388.6	1,395.6
PRW-C	6.112	544	25,062	339.5	1.222.4

1026 C.2 BIAS

1028 CUHK-SYSU was collected from various locations in Chinese urban cities and movie scenes, while
 1029 PRW was filmed at a Chinese university. Both datasets contain a sufficiently large number of indi 1030 viduals representing both genders. While PRW features a relatively younger demographic, CUHK 1031 SYSU includes people of various ages. Both datasets have an ethnicity bias, predominantly featuring
 1032 Asian individuals.

- 1033
- 1034 C.3 CORRUPTION IMPLEMENTATION DETAILS

In creating the benchmarks for the person search, we consider the typical attributes of scenes used
in person search, where both prominent subjects and multiple small background figures coexist in
the same scene. We conduct the data quality check to ensure that the persons in the images are
still detectable and re-identifiable by humans, even after corruption is applied. The implementation
details of the corruption in CUHK-SYSU-C and PRW-C are as follows:

1040

1041 **Snow.** We use the method proposed in Hendrycks & Dietterich (2019) to implement the snow 1042 corruption. To represent the diverse features of the snowy scene in the real world, we express 1043 various attributes of snowy scenes and adjust their intensity for different severity levels. Flake Size: 1044 Determines the average size or thickness of the snowflakes. As this value increases, the size of the 1045 snowflakes increases. We set the parameters (0.1, 0.2, 0.55, 0.55, 0.55) to adjust its intensity. Size 1046 Variation: Represents the standard deviation of the size distribution of snowflakes. A larger value results in greater variation in the sizes of snowflakes. We set 0.3 as a parameter for the intensity. 1047 *Snowfall Intensity*: Indicates the degree of snowfall intensity applied to the image. Higher values 1048 simulate heavier snowfall. We set the parameters (3, 2, 4, 4.5, 2.5) to adjust its intensity. Snow 1049 Coverage Threshold: Sets the minimum value for snow generation. Snow will not be generated 1050 below this threshold, simulating areas with less snow coverage. We set the parameters (10, 12, 12, 1051 12, 12) to adjust its intensity. Wind Effect: Determines the radius of motion blur. This simulates 1052 the direction and speed of falling snow, affected by wind. We set the parameters (0.5, 0.5, 0.9, 0.85, 1053 0.85) to adjust its intensity. *Blurriness*: Determines the intensity of motion blur. Higher values 1054 make the snowflakes appear more blurred, simulating faster-falling snow or stronger wind. We set 1055 the parameters (4, 4, 8, 8, 12) to adjust its intensity. Snow Opacity: Determines the mixing ratio 1056 between the original image and the snow effect layer. Values closer to 1 show more of the original 1057 image, while values closer to 0 intensify the snow effect, simulating denser snowfall. We set the parameters (0.8, 0.7, 0.7, 0.65, 0.65) to adjust its intensity. 1058

1059

Frost. We use the method proposed in Hendrycks & Dietterich (2019) to implement the frost corruption. This simulates the effect of frost or ice forming on the surface of the image, giving the appearance of a cold and frosty environment. *Frost Intensity:* This determines the strength of the frost effect applied to the image. Higher values result in more pronounced frost, simulating thick ice or first on the surface. We set the parameters (1, 0.8, 0.7, 0.65, 0.6) to adjust its intensity. *Blending Ratio:* Controls how much the frost image is blended with the original image. A lower value results in more frost coverage, while a higher value reveals more of the original image beneath the frost layer. We set the parameters (0.4, 0.6, 0.7, 0.7, 0.75) to adjust the blending ratio.

1067 1068

Fog. We use the method proposed in Hendrycks & Dietterich (2019) to implement the fog corruption. This simulates the effect of fog or mist, reducing visibility and softening the details in the image. *Fog Density:* This determines the density of the fog applied to the image. Higher values result in denser fog, which obscures more of the image. We set the parameters (1.5, 2.0, 2.5, 2.5, 3.0) to adjust the density. *Fog Smoothness:* Controls the rate of decay for the fractal noise used to generate the fog. Lower values create fog with sharper, more defined transitions, while higher values produce smoother fog with more gradual transitions. We set the parameters (2, 2, 1.7, 1.5, 1.4) to adjust its smoothness.

1076

Rain. We use the method proposed in Saxena (2023) to implement the rain corruption. To simulate the rain effect on images, we introduce various attributes related to rain and adjust their intensity at different severity levels. *Slope of Rain Droplets:* Determines the inclination of rain droplets. This simulates how much the rain is tilted by the wind. At lower intensities, droplets are lighter and more

1080 easily tilted by wind, forming steeper slopes. At higher intensities, droplets are heavier, causing rain 1081 to fall more vertically and be less affected by wind. We randomly select the slope from within the 1082 range (-20, 20) to adjust its intensity. Color of Rain Droplets: Determines how the rain droplets 1083 reflect light. As the severity increases, the color of the droplets changes from light gray to dark gray, simulating the visual effect of increasing rain density. This change reflects the reduction of light 1084 reflection and transmission through the air due to heavier rain. The default color is set to light gray (200, 200, 200), while for heavy rain, it is defined as medium gray (150, 150, 150), and for torrential 1086 rain, it is represented as dark gray (80, 80, 80). Blur Value: Indicates the degree of blur effect ap-1087 plied to the image. Higher values make the image blurry, simulating reduced visibility during heavy 1088 rainfall. We set the parameters (2, 3, 4, 5, 6) to adjust its intensity. **Drop Length:** Sets the length of 1089 rain droplets. At higher intensities, the length of droplets increases, creating a more dramatic rain 1090 effect. We set the parameters (10, 20, 30, 40, 50) to adjust its intensity. Drop Width: Determines the 1091 width of rain droplets. At higher intensities, the width of droplets increases, enhancing the more dis-1092 tinct rain effect. We set the parameters (1, 2, 3, 4, 5) to adjust its intensity. Brightness Adjustment: 1093 Adjusts the overall brightness of the image to simulate the dark environment of a rainy day. Lower 1094 values decrease the overall image brightness, creating a gloomy atmosphere. We set the parameters 1095 (0.7, 0.6, 0.5, 0.4, 0.3) to adjust its intensity.

1096

1097 Dark. We use the method proposed in Kong et al. (2024) to implement the dark corruption. To simulate darkness, we reduce the overall brightness of the image, making it appear darker. The lower the value, the darker the image becomes. We set the parameters (0.60, 0.54, 0.48, 0.42, 0.36) to adjust its intensity.

1101

Contrast. We use the method proposed in Hendrycks & Dietterich (2019) to implement the contrast corruption. To simulate contrast, we adjust the difference in brightness between pixels based on the average brightness of the image. Lower values result in a reduction of contrast, causing the image to appear blurrier and colors to become more uniform. We set the parameters (0.4, 0.33, 0.26, 0.18, 0.1) to adjust its intensity.

1107

Gaussian Noise. We use the method proposed in Hendrycks & Dietterich (2019) to implement the gaussian noise corruption. To simulate the gaussian noise, we adjust the noise intensity, which controls the standard deviation of the gaussian noise distribution applied to the image. The higher the value, the more pronounced the noise becomes. We set the parameters (0.05, 0.07, 0.09, 0.12, 0.15) to adjust its intensity.

1113

Speckle Noise. We use the method proposed in Hendrycks & Dietterich (2019) to implement the speckle noise corruption. To simulate the speckle noise, we adjust the noise intensity that is multiplied by the pixel values themselves. The higher the value, the more the image is disrupted and appears to have a grainy mixture. We set the parameters (0.1, 0.2, 0.3, 0.4, 0.5) to adjust its intensity.

**Gaussian Blur.** We use the method proposed in Hendrycks & Dietterich (2019) to implement the gaussian blur corruption. To simulate the gaussian blur, we adjust the standard deviation of the Gaussian filter applied to the image. The higher the value, the more blurred the image appears. We set the parameters (1, 2, 3, 4, 5) to adjust its intensity.

1123

1124 **Motion Blur.** We use the method proposed in Hendrycks & Dietterich (2019) to implement the 1125 motion blur corruption. This simulates the blur caused by the movement of the camera or objects 1126 in the scene during exposure, producing a streaking effect. *Motion Trace Length:* This controls the 1127 length of the motion blur, representing how far objects have moved during the exposure. A higher 1128 radius results in a more pronounced blur effect, simulating faster movement. We set the parameters 1129 (10, 15, 15, 15, 20) to adjust its intensity. Blur Sharpness: Determines the sharpness of the motion 1130 blur. Higher values result in a smoother blur, while lower values retain more definition along the 1131 motion streak. We set the parameters (3, 5, 8, 12, 15) to adjust its sharpness. Angle Direction: Controls the angle at which the motion blur is applied, simulating motion in different directions. 1132 The angle is randomized within a certain range to simulate natural motion blur effects caused by 1133 different movements.

1134 **Defocus Blur.** We use the method proposed in Hendrycks & Dietterich (2019) to implement the 1135 defocus blur corruption. To represent the diverse features of defocus blur in real-world photography, 1136 we express two attributes of defocus and adjust their intensity for different severity levels. Blur 1137 Kernel Size: Determines the size of the blur kernel. As the radius increases, the blurring spreads 1138 over a larger area, causing details to fade into the background. We set the parameters (3, 4, 5, 7, 9) to adjust its intensity. Blur Smoothness: Controls the smoothness or sharpness of the blur effect. 1139 Higher values produce a smoother, more gradual blur, while lower values retain sharper transitions 1140 at the edges of the blur. This affects the overall softness of the defocus effect. We set the parameters 1141 (0.1, 0.5, 0.5, 0.5, 0.5) to adjust its intensity. 1142

1143

Glass Blur. We use the method proposed in Hendrycks & Dietterich (2019) to implement the glass 1144 blur corruption. To represent the diverse features of the glass blur scene in the real world, we express 1145 various attributes of glass blur and adjust their intensity for different severity levels. *Blur Strength:* 1146 Determines the strength of the glass blur applied to the image, simulating distortion as seen through 1147 the glass. Higher values result in a more blurred and diffused image, where details become softer 1148 and less defined. As sigma increases, the overall smoothness of the blur effect intensifies. We set the 1149 parameters (0.7, 0.9, 1, 1.1, 1.5) to adjust its intensity. *Glass Distortion Magnitude:* Represents the 1150 degree of distortion caused by imperfections in the glass. As this value increases, the image appears 1151 more warped, simulating the effect of viewing through glass with varying thickness or composition. 1152 We set the parameters (1, 2, 2, 3, 4) to adjust its intensity. *Distortion Repetitions (Iterations):* 1153 Determines how many times the displacement effect is applied, creating multiple layers of distortion. More iterations result in a more pronounced and complex glass-like effect. We set the parameters 1154 (2, 1, 3, 2, 2) to adjust its intensity. 1155

1156

Elastic Transform. We use the method proposed in Hendrycks & Dietterich (2019) to implement the elastic transform corruption. To simulate elastic distortion, we apply random, small-scale deformations to the image pixels, adding a flexible, rubber-like warping effect. The higher the value, the more pronounced the distortions become, making the image appear more heavily warped. We set the parameters (12.5, 16.25, 21.25, 25, 30) to adjust its intensity.

1162

**Spatter.** We use the method proposed in Hendrycks & Dietterich (2019) to implement the spatter 1163 corruption. This method simulates liquid splashes or mud spatters on the image, which can occur 1164 in outdoor or dirty environments, distorting visibility and adding a natural effect of environmental 1165 interference. Liquid Amount: This determines the average amount of liquid spattered on the image. 1166 Higher values simulate heavier splashes or more liquid, resulting in larger and more widespread 1167 spatter areas. We set the parameters (0.65, 0.65, 0.65, 0.65, 0.67) to adjust its intensity. Size Vari-1168 ation: Represents the standard deviation of the size of spatter drops. A larger value results in more 1169 variation in the size of splatter particles, simulating irregular drops. We set the parameters (0.3, 0.3, 1170 0.3, 0.3, 0.4) to adjust its intensity. **Blur Radius:** Determines the size of the gaussian blur applied to the liquid layer. Higher values create more diffused spatter, simulating less defined edges and softer 1171 splashes. We set the parameters (4, 3, 2, 1, 1) to adjust the blurriness. *Coverage Threshold:* This 1172 sets the minimum intensity threshold for spatter formation. Spatter below this threshold is not visi-1173 ble, simulating splatters that did not adhere to the surface. We set the parameters (0.69, 0.68, 0.68, 1174 0.65, 0.65) to adjust its coverage. **Opacity:** This controls the opacity of the spatter layer. Higher 1175 values create more opaque spatter, making it more prominent on the image. Lower values create 1176 more transparent splatter effects, simulating thinner layers of liquid. We set the parameters (0.6, 1177 0.6, 0.5, 1.5, 1.5) to adjust its visibility. *Mud vs. Water Effect:* This determines whether the spatter 1178 simulates water (0) or mud (1). When set to 1, the spatter is brown and more opaque, simulating 1179 thick mud. When set to 0, the spatter simulates pale water splashes. We set the parameters (0, 0, 0, 0, 0)1180 1, 1) to adjust its effect.

1181

Saturate. We use the method proposed in Hendrycks & Dietterich (2019) to implement the saturate corruption. To represent varying levels of color saturation in real-world scenarios, we express different attributes that influence color intensity and adjust them for different severity levels. *Saturation Scale:* This determines the strength of the color saturation applied to the image. Higher values result in more vivid and intense colors, while lower values make the image appear more desaturated or washed out. Both high or low saturation scales can lead to distortion and degradation of the clean image. We set the parameters (2, 0.2, 0.1, 0.3, 5) to adjust its intensity. *Offset:* Represents a con-

stant value added to the saturation scale. It controls the base level of saturation applied uniformly across the image, ensuring that even low saturation images retain some color vibrance. We set the parameter (0, 0, 0, 0, 0.1) to adjust the base intensity.

**Pixelate.** We use the method proposed in Hendrycks & Dietterich (2019) to implement the pix-elate corruption. To simulate pixelation, we reduce the resolution of the image, converting it into larger blocks of pixels, then resizing it back to the original resolution, removing details. The lower the value, the stronger the pixelation effect, making the image appear more blocky. We set the parameters (0.6, 0.5, 0.4, 0.3, 0.25) to adjust its intensity.

JPEG Compression. We use the method proposed in Hendrycks & Dietterich (2019) to imple-ment the JPEG compression corruption. To simulate JPEG compression, we adjust the compression quality level to observe its effects on image quality. The lower the value, the more the image quality degrades, leading to more artifacts and damage. We set the parameters (25, 18, 15, 10, 7) to adjust its intensity. 

Brightness. We use the method proposed in Hendrycks & Dietterich (2019) to implement the brightness corruption. To simulate brightness, we adjust the overall brightness of the image by adding a constant value to the pixel values. The higher the value, the brighter the image becomes, and the lower the value, the darker the image appears. We set the parameters (0.1, 0.2, 0.3, 0.4, 0.5)to adjust its intensity.

1210	Table 11. Fertormance by unrefert types of corruptions.										
1211	Туре	Snow	Frost	Fog	Rain	Dark	Contrast				
1212	R@1/mAP	59.9 / 55.5	64.6 / 63.1	84.8 / 83.3	72.5 / 70.0	83.5 / 82.3	62.9 / 58.5				
1213	Туре	Gaussian Noise	Speckle Noise	Gaussian Blur	Motion Blur	Defocus Blur	Glass Blur				
1214	R@1/mAP	32.3 / 27.8	71.0 / 68.9	79.2 / 75.6	79.1 / 74.5	79.3 / 76.5	81.5 / 79.0				
1215	Туре	Elastic Transform	Spatter	Saturate	Pixelate	JPEG Compression	Brightness				
1216	R@1/mAP	92.6/91.7	79.5 / 77.7	54.2 / 45.0	78.8 / 76.7	78.3 / 73.2	90.3 / 89.3				

Table 11: Performance by different types of corruptions



Figure 9: **Qualitative results of person search** on the PRW-C Dataset. The first row displays query images, while the second, third, fourth, and fifth rows show the results from OADG+CIL (Lee et al., 2024; Chen et al., 2021), SeqNet (Li & Miao, 2021), OIMNet++ (Lee et al., 2022), COAT (Yu et al., 2022), and Ours with SeqNet, respectively. Each column indicates the different query and the corresponding retrieval results of various models. The blue color denotes the box for a query, the red color indicates the box for failure cases, and the green color represents the box for success cases. A total of 18 types of corruption and 5 levels of severity are involved in establishing the PRW-C dataset.

1294



Figure 10: **Qualitative results of baseline and our method on real-world corruptions.** Baseline refers to SeqNet (Li & Miao, 2021).