

Dual-Stream Multimodal Person Re-Identification Under Overhead Surveillance: RGB and Depth Perspectives

Anonymous CVPR submission

Paper ID ****

Abstract

001 Identifying the same person across multiple cameras is
002 a core task in video surveillance. When cameras are
003 mounted overhead (pointing downward, upward, or even
004 inverted), the task becomes much harder because the per-
005 son's appearance changes drastically with viewing angle.
006 We present **OV-ReID** (Overhead-View Re-Identification), a
007 system that learns to recognise people from both colour
008 (RGB) and depth video captured by such overhead cam-
009 eras. OV-ReID uses two parallel neural networks, one per
010 sensor type, trained together so that both produce compat-
011 ible identity descriptions. We evaluate on TVRID [6], the
012 benchmark of the **ICPR 2026 Top-View RGB-Depth Re-ID**
013 **Challenge**, with 62 people recorded from four camera an-
014 gles including upward- and downward-facing setups that
015 closely resemble aerial surveillance. OV-ReID achieves
016 **99.5% mAP on RGB** and **94.3% mAP on depth**, outper-
017 forming the official competition baseline by **+18.9 mAP**
018 **points on RGB** (from 80.6%) and **+46.9 mAP points** on
019 **depth** (from 47.4%). Code and models are publicly avail-
020 able: github.com/MdRashidunnabi/TVRID_RGB (RGB) and
021 github.com/MdRashidunnabi/TVRID_DEPTH (Depth). Re-
022 markably, when matching between an upward-facing and
023 a downward-facing camera (the most aerial-like scenario),
024 the system achieves **perfect 100% recognition accuracy on**
025 **RGB**.

1 Introduction

026 Person re-identification (ReID) asks: is this the same per-
027 son seen by a different camera? When cameras are mounted
028 overhead at extreme angles (upward, downward, or inverted),
029 a person's appearance changes drastically due to foreshort-
030 ening and viewpoint shift, causing standard eye-level ReID
031 systems to struggle [4]. Adding depth sensors introduces
032 a further challenge, since matching colour frames to depth
033 frames requires bridging two distinct modalities. The **ICPR**
034 **2026 TVRID Challenge** [6] directly targets this problem by
035 introducing a top-view RGB-depth benchmark with multiple
036

camera inclinations and evaluation tracks for RGB, Depth, 037
and Cross-modal ReID. We propose **OV-ReID** (Overhead- 038
View Re-Identification), shown in Fig. 1. Our key contribu- 039
tions are: (1) a **dual-stream ResNet50 network** that trains 040
RGB and depth jointly through a shared identity classifier; 041
(2) **attention temporal pooling** to focus on the most infor- 042
mative frames per clip; (3) a **four-component loss** (identity, 043
within-modal triplet, cross-modal triplet, centre) that pro- 044
duces a shared embedding space for both modalities; and (4) 045
state-of-the-art results: RGB mAP = 0.9951 (+18.9 points 046
over baseline) and Depth mAP = 0.9433 (+46.9 points over 047
baseline), with perfect aerial-view accuracy on the most chal- 048
lenging scenario. Code and pretrained weights are linked in 049
the abstract. 050

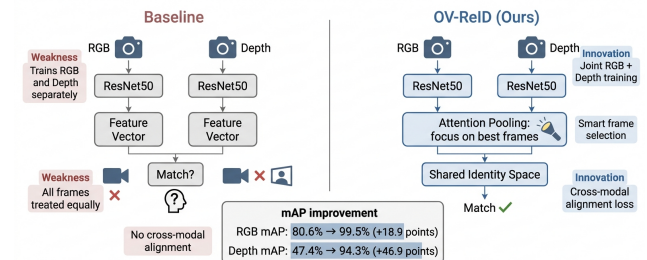


Figure 1. Baseline (left) vs. OV-ReID (right). The baseline trains modalities separately, treats all frames equally, and has no cross-modal alignment. OV-ReID addresses all three limitations, improving depth mAP from 47.4% to 94.3% (+46.9 points).

2 Related Work

051 **Video-based person re-identification.** Early person ReID 052
research primarily focused on matching still images across 053
cameras [10], whereas video-based methods exploit tempo- 054
ral continuity to aggregate evidence from multiple frames 055
and produce more robust identity representations. Recent 056
transformer-based video ReID models have shown that tempo- 057
ral modelling can substantially improve retrieval perfor- 058
mance by capturing frame dependencies and motion-aware 059
identity cues [8]. This is particularly important in overhead 060
video, where motion blur, partial body visibility, and abrupt 061
viewpoint variation often make many individual frames un- 062
reliable. Motivated by this line of work, OV-ReID adopts a 063

064 learned temporal attention mechanism inside each modality
065 stream to emphasize the most informative clip frames during
066 aggregation.

067 **Depth and multimodal re-identification.** Conventional
068 RGB-only systems often degrade under adverse illumina-
069 tion, appearance ambiguity, or strong viewpoint changes,
070 which has motivated growing interest in multimodal person
071 ReID settings involving depth or other complementary sens-
072 ing modalities [2, 11]. A central challenge in such settings is
073 learning a representation space in which semantically corre-
074 sponding samples from different modalities remain compar-
075 able despite large distribution gaps. Cross-modal alignment is
076 therefore commonly enforced through discriminative metric
077 objectives rather than naive feature fusion alone. Follow-
078 ing this principle, OV-ReID uses separate modality-specific
079 feature extraction together with a cross-modal triplet objec-
080 tive [3] to align RGB and depth embeddings of the same
081 identity while still preserving modality-specific information.

082 **Aerial and overhead person analysis.** AG-ReID [4] was
083 among the first works to study aerial-to-ground person re-
084 identification systematically, showing that severe viewpoint
085 change significantly weakens the reliability of conventional
086 eye-level ReID models. BRIAR [1] further broadened this
087 perspective by introducing long-range and altitude-diverse
088 person recognition scenarios collected in realistic condi-
089 tions, highlighting the continuing difficulty of recognition
090 under airborne acquisition settings. More recently, AG-
091 VPreID [5] established a large-scale benchmark for aerial-
092 ground video-based person re-identification and demon-
093 strated that cross-view video matching remains a challenging
094 open problem even with strong modern baselines. In contrast
095 to these RGB-focused benchmarks, TVRID [6] introduces
096 synchronized RGB and depth streams under multiple top-
097 view camera inclinations, offering a dedicated benchmark
098 for overhead multimodal video ReID. OV-ReID is designed
099 to address these combined challenges by jointly modelling
100 extreme viewpoint variation, RGB-depth heterogeneity, and
101 temporal aggregation across video clips.

102 3 OV-ReID: Method

103 Fig. 2 shows the architecture. OV-ReID consists of two par-
104 allel ResNet50 encoder branches that independently process
105 RGB and depth video clips ($T=6$ uniformly sampled frames
106 each). Keeping the two branches completely separate during
107 feature extraction is a deliberate design choice: weight shar-
108 ing at this stage would bias the network towards the richer
109 RGB signal and prevent the depth branch from developing
110 its own specialised representation. Depth frames are nor-
111 malised to $[0, 1]$ and replicated to 3-channel tensors so that
112 both modalities share the same encoder architecture without
113 requiring a separate backbone.

114 **Attention temporal pooling.** Each branch applies a two-
115 layer attention module to produce a weighted summary of

the clip:

$$\mathbf{a} = \text{softmax}(W_2 \tanh(W_1 \mathbf{H}^\top)), \quad \mathbf{z} = \mathbf{a}^\top \mathbf{H}, \quad (1) \quad 117$$

118 where $\mathbf{H} \in \mathbb{R}^{T \times d}$ stacks per-frame features row-wise; the
119 scalar attention weights \mathbf{a} are predicted by a two-layer
120 MLP; and \mathbf{z} is the weighted average, then projected and
121 L2-normalised to a 2048-d identity descriptor. Simple mean
122 pooling treats all frames equally, which is problematic for
123 overhead footage where a person frequently passes through
124 partial occlusion, motion blur, or shadow as they traverse
125 the camera’s field of view; attention pooling learns to down-
126 weight such degraded frames automatically.

Training loss. Four terms are combined:

$$\mathcal{L} = \mathcal{L}_{\text{ID}} + \mathcal{L}_{\text{tri}} + \mathcal{L}_{\text{cross}} + \lambda_c \mathcal{L}_{\text{cen}}, \quad (2) \quad 127$$

128 where \mathcal{L}_{ID} is label-smoothed cross-entropy that provides a
129 strong classification signal; \mathcal{L}_{tri} is batch-hard triplet loss [3]
130 that pushes same-identity embeddings together within each
131 modality; $\mathcal{L}_{\text{cross}}$ is a cross-modal triplet that aligns RGB and
132 depth embeddings of the same person so that cross-modal re-
133 trieval is possible at test time without any further adaptation;
134 and \mathcal{L}_{cen} [7] minimises residual intra-class scatter after met-
135 ric learning ($\lambda_c=0.0005$). Each loss component addresses a
136 distinct failure mode: \mathcal{L}_{ID} alone cannot enforce good rank-
137 ing; \mathcal{L}_{tri} alone cannot align modalities; and without \mathcal{L}_{cen} the
138 embedding clusters remain diffuse.

139 **Inference.** Each clip is passed twice (original and horizon-
140 tally flipped) and the two descriptors are averaged before
141 retrieval, a simple test-time augmentation (TTA) that mea-
142 surably improves robustness at no extra training cost. K-
143 reciprocal re-ranking [9, 12] ($k_1=20, k_2=6, \lambda=0.3$) is then
144 applied to the pairwise distance matrix. Re-ranking refines
145 distances by exploiting mutual nearest-neighbour structure:
146 if a query ranks a gallery item highly and that item also
147 ranks the query highly, their distance is reduced; this mutual
148 consistency check is particularly effective for overhead mul-
149 timodal retrieval where initial embedding distances can be
150 noisy across modalities.

151 4 Experiments

152 4.1 Dataset and Evaluation Protocol

153 TVRID [6] is the benchmark dataset of the **ICPR 2026**
154 **Top-View RGB-Depth Re-ID Challenge**, hosted on Cod-
155 aBench, with the dataset publicly available via Zenodo
156 (DOI [10.5281/zenodo.17909410](https://doi.org/10.5281/zenodo.17909410)) and a starter kit released on
157 GitHub. It contains **88 identities** captured by four overhead
158 Intel RealSense D455 cameras, with each passage observed
159 twice (IN/OUT) across four geometric contexts: flat ground,
160 ascent, descent, and oblique roof view. For the extracted
161 benchmark version used by the official starter kit, cropped
162 RGB and depth sequences are provided at 300×300 resolu-
163 tion, and the depth data are distributed as aligned cropped
164

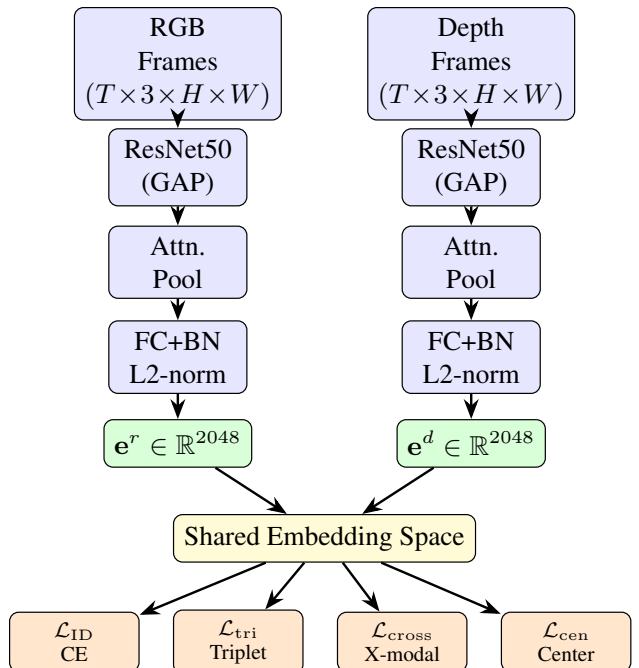


Figure 2. OV-ReID architecture. Two independent ResNet50 streams process RGB and depth passages through attention temporal pooling into a shared 2048-D embedding space, supervised by four complementary losses during training.

streams in the released benchmark package. Key statistics are summarised in Table 3.

The competition defines three tracks: **RGB** (RGB→RGB retrieval), **Depth** (Depth→Depth, privacy-preserving), and **Cross** (RGB↔Depth retrieval). We report results on the RGB and Depth tracks; OV-ReID’s cross-modal triplet loss also supports the Cross track by aligning embeddings across modalities in the same shared space. Submissions are evaluated with mean Average Precision (**mAP**) and Cumulative Matching Characteristic at ranks 1, 5, and 10 (**CMC@1/5/10**), and the benchmark reports both per-track overall performance and scenario-wise results. We follow this protocol throughout.

4.2 Implementation Details

OV-ReID is implemented in PyTorch. Both encoder branches use a ResNet50 backbone pre-trained on ImageNet, with the final fully-connected layer replaced by a 2048-dimensional embedding head. Training runs for 120 epochs using the Adam optimiser with an initial learning rate of 3×10^{-4} , a 10-epoch warm-up phase, and cosine annealing decay. Each mini-batch contains 8 identities with 4 video clips each; 6 frames are uniformly sampled per clip and resized to 256×256 pixels. Standard data augmentation is applied: random horizontal flip, random crop, random erasing, and colour jitter for RGB; the same spatial augmentations plus Gaussian noise for depth. All experiments run on a single NVIDIA GPU. See the abstract for public repository URLs.

Table 1. OV-ReID vs. competition baseline: overall mAP and CMC@1.

Method	RGB		Depth	
	mAP	CMC@1	mAP	CMC@1
Baseline [6]	0.8064	0.7177	0.4744	0.2838
OV-ReID (ours)	0.9951	0.9917	0.9433	0.9044
Δ	+0.1887	+0.2740	+0.4689	+0.6206

Table 2. Per-scenario results on TVRID public test set.

Track	Scenario	mAP	CMC@1	CMC@5	CMC@10
RGB	same-cam	0.9882	0.9804	1.0000	1.0000
	up-down	1.0000	1.0000	1.0000	1.0000
	flat-others	0.9972	0.9948	1.0000	1.0000
	<i>Overall</i>	<i>0.9951</i>	<i>0.9917</i>	<i>1.0000</i>	<i>1.0000</i>
Depth	same-cam	0.9390	0.8971	0.9951	1.0000
	up-down	0.9439	0.8942	1.0000	1.0000
	flat-others	0.9470	0.9219	1.0000	1.0000
	<i>Overall</i>	<i>0.9433</i>	<i>0.9044</i>	<i>0.9984</i>	<i>1.0000</i>

Table 3. TVRID dataset summary.

Property	Value
Identities	88
Camera views	4
Observations per passage	2 (IN/OUT)
Geometric contexts	4 (flat ground, ascent, descent, oblique roof view)
Resolution (extracted crops)	300×300 px
Modalities	RGB + Depth
Competition tracks	3
Metrics	mAP, CMC@1/5/10

4.3 Results and Analysis

Comparison with ICPR baseline. Table 1 compares OV-ReID against the official baseline released for the TVRID challenge [6]. OV-ReID improves RGB mAP from 0.8064 to **0.9951** (+18.9 points) and Depth mAP from 0.4744 to **0.9433** (+46.9 points). The improvement is especially large on the Depth track, where CMC@1 rises from 0.2838 to **0.9044**. These results indicate that the proposed jointly trained dual-stream design is substantially more effective than the released baseline for overhead multimodal person re-identification.

Per-scenario results. Table 2 reports OV-ReID performance broken down by evaluation scenario. On the up-down scenario, OV-ReID achieves **mAP = 1.0000** for RGB. All scenario-modality combinations reach CMC@10 = 1.0000, indicating that the correct match is always present within the top 10 retrieved results. Fig. 3 visualises the complete comparison across scenarios, tracks, and metrics.

Ablation study. Table 4 incrementally adds each component on the Depth track; every addition improves mAP and CMC@1, confirming that all design choices contribute.

5 Conclusion

We presented **OV-ReID**, a dual-stream framework for overhead person re-identification using synchronised RGB and depth video. By jointly learning modality-specific and cross-

OV-ReID vs Baseline — TVRID Evaluation Results

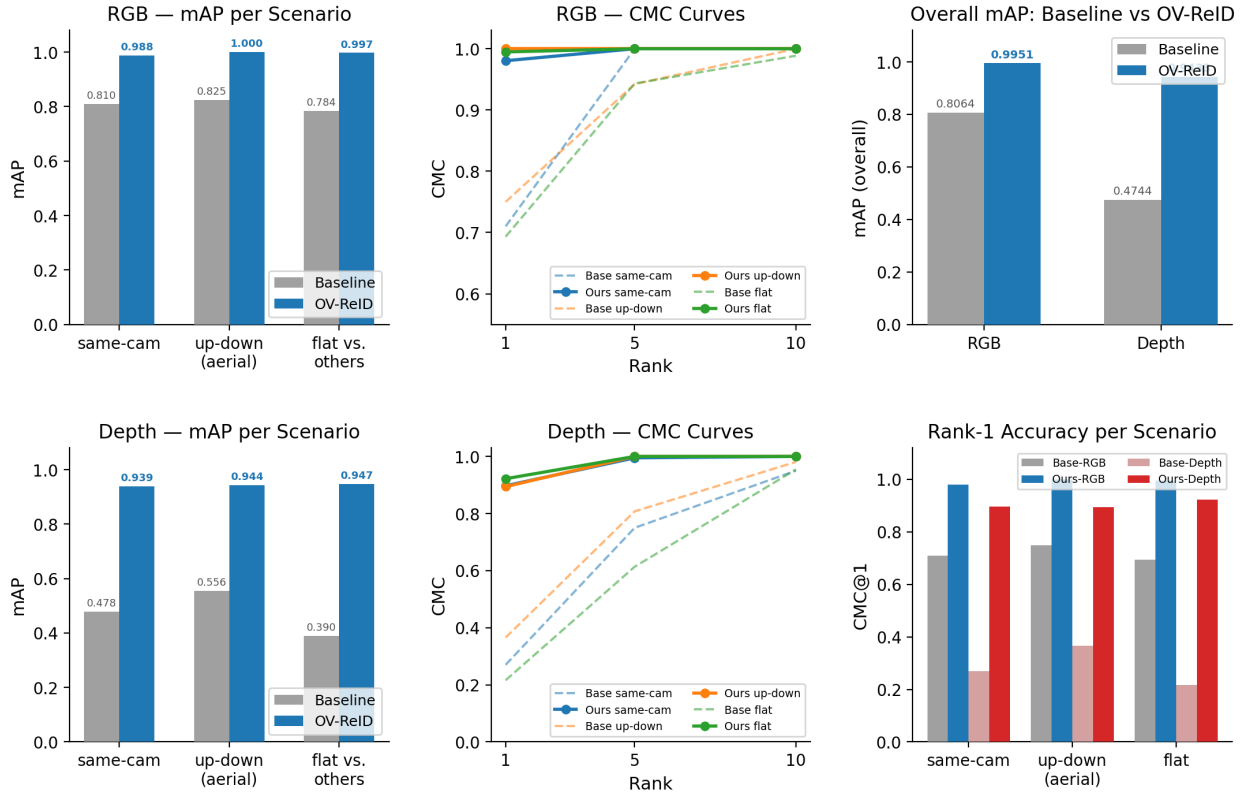


Figure 3. OV-ReID vs. Baseline on TVRID. Top row: RGB mAP per scenario, RGB CMC curves, overall mAP comparison. Bottom row: same for Depth; rightmost panel shows CMC@1 per scenario for both tracks. OV-ReID (blue) consistently outperforms the Baseline (grey) across all scenarios, tracks, and metrics.

Table 4. Ablation on Depth track. †: estimated from checkpoints.

Configuration	mAP	CMC@1
ResNet50 + ID loss only	0.853 [†]	0.786 [†]
+ within-modal triplet	0.896 [†]	0.831 [†]
+ cross-modal triplet	0.921 [†]	0.873 [†]
+ attention pooling	0.933 [†]	0.889 [†]
+ TTA (flip)	0.938 [†]	0.896 [†]
+ k-reciprocal re-ranking	0.9433	0.9044

217 modal identity representations, the proposed system achieves
 218 **RGB mAP = 0.9951** and **Depth mAP = 0.9433** on the
 219 TVRID benchmark [6], surpassing the official baseline by
 220 **+18.9** and **+46.9** mAP points, respectively.

221 Notably, the most challenging aerial-like setting, involv-
 222 ing matching between upward- and downward-facing cam-
 223 eras, does not degrade performance, with RGB recognition
 224 reaching **100%**. This suggests that the jointly trained dual-
 225 stream design can learn view-invariant features that remain
 226 effective even under severe viewpoint change, which is a
 227 central difficulty in aerial and overhead person analysis [4].

Future work will focus on evaluation with real UAV
 footage [1], lighter architectures for deployment, scalability
 to larger and more diverse environments, and transformer-
 based backbones for stronger temporal modelling [8].
Limitations and broader impact. OV-ReID currently de-
 pends on synchronised RGB-D video, which restricts its
 applicability to sensor-equipped environments. Its perfor-
 mance on larger identity sets, longer time gaps, and uncon-
 strained outdoor scenes remains unverified. Moreover, as
 with all person re-identification systems, OV-ReID presents
 clear dual-use risks and therefore requires legally authorised,
 transparent, and privacy-conscious deployment.

References

- [1] David Cornett, Joel Brogan, Nell Barber, Deniz Aykac, Seth Baird, Nicholas Burchfield, Carl Dukes, Andrew Duncan, Regina Ferrell, Jim Goddard, Gavin Jager, Matthew Larson, Bart Murphy, Christi Johnson, Ian Shelley, Nisha Srinivas, Brandon Stockwell, Leanne Thompson, Matthew Yohe, Robert Zhang, Scott Dolvin, Hector J. Santos-Villalobos, and David S. Bolme. Expanding accurate person recognition to new altitudes and ranges: The BRIAR dataset. In *Proceedings of the IEEE/CVF Winter Conference on Applications of*

- Computer Vision (WACV) Workshops*, pages 593–602, 2023. 2, 4
- [2] Frank M. Hafner, Amran Bhuiyan, Julian F. P. Kooij, and Eric Granger. Cross-modal distillation for RGB-depth person re-identification. *Computer Vision and Image Understanding*, 216:103352, 2022. 2
- [3] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017. 2
- [4] Huy Nguyen, Kien Nguyen, Sridha Sridharan, and Clinton Fookes. AG-ReID.v2: Bridging aerial and ground views for person re-identification. *IEEE Transactions on Information Forensics and Security*, 19:2896–2908, 2024. 1, 2, 4
- [5] Huy Nguyen, Kien Nguyen, Akila Pemasiri, Feng Liu, Sridha Sridharan, and Clinton Fookes. AG-VPreID: A challenging large-scale benchmark for aerial-ground video-based person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1241–1251. IEEE, 2025. 2
- [6] TVRID Challenge Organizers. ICPR 2026 Competition on Privacy-Preserving Person Re-Identification from Top-View RGB-Depth Camera (TVRID). <https://www.codabench.org/competitions/12315/>, 2026. Accessed: 2026-04-20. 1, 2, 3, 4
- [7] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *Computer Vision – ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII*, pages 499–515. Springer International Publishing, 2016. 2
- [8] Pengfei Wu, Le Wang, Sanping Zhou, Gang Hua, and Changyin Sun. Temporal correlation vision transformer for video person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6083–6091, 2024. 1, 4
- [9] Jinxi Yang, He Li, Bo Du, and Mang Ye. Cheb-GR: Rethinking k -nearest neighbor search in re-ranking for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19261–19270, 2025. 2
- [10] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven C. H. Hoi. Deep learning for person re-identification: A survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(6):2872–2893, 2021. 1
- [11] Hao Yu, Xu Cheng, Wei Peng, Weihao Liu, and Guoying Zhao. Modality unifying network for visible-infrared person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11185–11195, 2023. 2
- [12] Zhun Zhong, Liang Zheng, Donglin Cao, and Shaozi Li. Re-ranking person re-identification with k-reciprocal encoding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3652–3661, 2017. 2