

Scale-Aware Vision-Language Adaptation for Extreme Far-Distance Video Person Re-identification

Anonymous CVPR submission

Paper ID 3

Abstract

001 *Extreme far-distance video person re-identification (ReID)*
002 *is particularly challenging due to scale compression, res-*
003 *olution degradation, motion blur, and aerial-ground view-*
004 *point mismatch. As camera altitude and subject distance*
005 *increase, models trained on close-range imagery degrade*
006 *significantly. In this work, we investigate how large-*
007 *scale vision-language models can be adapted to operate*
008 *reliably under these conditions. Starting from a CLIP-*
009 *based baseline, we upgrade the visual backbone from ViT-*
010 *B/16 to ViT-L/14 and introduce backbone-aware selec-*
011 *tive fine-tuning to stabilize adaptation of the larger trans-*
012 *former. To address noisy and low-resolution tracklets,*
013 *we incorporate a lightweight temporal attention pooling*
014 *mechanism that suppresses degraded frames and empha-*
015 *sizes informative observations. We retain adapter-based*
016 *and prompt-conditioned cross-view learning to mitigate*
017 *aerial-ground domain shifts, and further refine retrieval*
018 *using improved optimization and k-reciprocal re-ranking.*
019 *Experiments on the DetReIDX stress-test benchmark show that*
020 *our approach achieves an average mAP of 35.73 across*
021 *aerial-ground (A2G), ground-aerial (G2A), and aerial-*
022 *aerial (A2A) protocols, improving substantially over the ex-*
023 *isting baseline. These results show that large-scale vision-*
024 *language backbones, when combined with stability-focused*
025 *adaptation, significantly enhance robustness in extreme far-*
026 *distance video person ReID.*

027 1. Introduction

028 Person re-identification (ReID) aims to match individ-
029 uals across non-overlapping camera views [1] and has
030 progressed rapidly with vision transformers, CLIP-based
031 text encoders, and large-scale pretraining [2]. These ad-
032 vances have significantly improved performance on stan-
033 dard ground-level benchmarks and enabled applications in
034 surveillance, crowd analysis, and public safety. However,
035 extreme far-distance aerial scenarios operate under funda-



Figure 1. Illustration of extreme far-distance conditions in the DetReIDX dataset. As UAV altitude and horizontal distance increase, pedestrians undergo severe scale compression and resolution degradation. The red bounding boxes highlight the same group of individuals as their pixel footprint shrinks dramatically, showing the visual challenges of aerial-ground video person ReID.

036 mentally different visual constraints. As the altitude of
037 unmanned aerial vehicles (UAVs) and subject distance in-
038 crease, pedestrians often occupy only a few pixels and are
039 affected by severe scale compression, resolution degrada-
040 tion, motion blur, and strong viewpoint differences between
041 aerial and ground cameras, as shown in Figure 1. Under
042 such conditions, fine-grained appearance cues become un-
043 reliable, and models developed for close-range imagery de-
044 grade sharply. The assumptions underlying conventional
045 ReID, sufficient spatial detail, and consistent viewpoints no
046 longer apply.

047 Recent stress-test datasets and benchmarks for UAV-
048 based person analysis highlight this limitation explicitly.
049 In particular, the DetReIDX was introduced to evaluate ro-
050 bustness under long-range viewpoints, aerial-ground cross-
051 domain matching, and session-level appearance changes.

052 The benchmark reveals that strong detection and ReID base-
053 lines can collapse under extreme distance and viewpoint
054 variation. Improving performance, therefore, requires adap-
055 tation strategies that explicitly account for scale degradation
056 and noisy video tracklets, rather than direct transfer from
057 conventional settings.

058 In this work, we explore how large-scale vision-language
059 backbones can be effectively adapted for extreme far-
060 distance video ReID. Starting from the official CLIP ViT-
061 B/16 baseline for DetReIDX, we investigate backbone scal-
062 ing and stability-oriented fine-tuning. We upgrade the vi-
063 sual encoder to ViT-L/14 and introduce backbone-aware
064 selective fine-tuning to preserve pretrained representations
065 while enabling high-level domain adaptation. To im-
066 prove robustness to frame-level degradation, we incorpo-
067 rate a lightweight temporal attention pooling mechanism
068 that emphasizes informative frames within noisy track-
069 lets. We further retain prompt-based cross-view condi-
070 tioning and adapter tuning to mitigate aerial-ground do-
071 main shifts, and refine retrieval performance through op-
072 timized training strategies and k-reciprocal re-ranking [3].
073 On the DetReIDX benchmark, evaluated across aerial-to-
074 aerial (A2A), aerial-to-ground (A2G), and ground-to-aerial
075 (G2A) protocols [4], our approach improves the average
076 mAP from 28.11 to 35.73, surpassing the recent work
077 results (32.89). These findings demonstrate that large-
078 scale CLIP backbones, when carefully adapted, substan-
079 tially improve robustness in extreme far-distance video per-
080 son ReID. We summarize the contribution of this work be-
081 low:

- 082 1. We present a systematic study of CLIP backbone scaling
083 (ViT-B/16 to ViT-L/14) for extreme far-distance aerial-
084 ground video person ReID, demonstrating that increased
085 model capacity improves robustness under severe scale
086 compression.
- 087 2. We propose a stability-focused adaptation strategy com-
088 bining backbone-aware selective fine-tuning, temporal
089 attention pooling, and prompt/adapter tuning to handle
090 noisy, low-resolution video tracklets.
- 091 3. We introduce practical optimization refinements, includ-
092 ing cosine scheduling and enhanced data augmentation,
093 and enable k-reciprocal re-ranking at inference.
- 094 4. Our method achieves 35.73 mAP on the benchmark, im-
095 proving substantially from the baseline (28.11) and the
096 highest publicly reported score (32.89).

097 The remainder of this paper is organized as follows. Sec-
098 tion 2 reviews related work in ground-level and aerial video
099 person ReID, vision-language adaptation, and retrieval re-
100 ranking. Section 3 describes the proposed scale-aware
101 adaptation framework, including backbone scaling, selec-
102 tive fine-tuning, temporal attention pooling, and optimiza-
103 tion refinements. Section 4 presents experimental results on
104 the DetReIDX benchmark, including implementation de-

tails and quantitative comparisons. Section 5 discusses the
results, limitations and future improvements. Finally, Sec-
tion 6 concludes the paper.

2. Related Work

Person ReID has been extensively studied in both image-
based and video-based settings. With the adoption of
transformer backbones and large-scale pretraining, signif-
icant progress has been achieved on conventional ground-
level benchmarks. However, robustness under extreme far-
distance aerial conditions remains largely underexplored.
We review related work in ground-level ReID, video-based
temporal modeling, aerial and cross-view benchmarks,
vision-language adaptation, and retrieval re-ranking.

2.1. Ground-level Person ReID

Early ReID research focused primarily on ground-level
camera networks, where pedestrians are captured at moder-
ate distances with sufficient spatial resolution. Benchmarks
such as Market-1501 [5] and DukeMTMC4ReID [6, 7] es-
tablished standardized evaluation protocols and drove rapid
development of discriminative feature learning methods.
More recently, transformer-based architectures, including
TransReID [8] and related global-context modeling frame-
works, have demonstrated strong performance by leverag-
ing self-attention and large-scale pretraining.

Despite these advances, most methods are evaluated un-
der relatively stable imaging conditions. When applied
to extreme far-distance aerial scenarios, where person in-
stances are severely scale-compressed and contain limited
visual detail, their performance degrades substantially. This
limitation highlights the challenge of directly transferring
ground-level transformer models to aerial settings without
explicit adaptation for scale degradation and viewpoint dis-
crepancy.

2.2. Video-based ReID and Temporal Modeling

Video-based ReID extends image-based matching by ex-
ploiting temporal cues across tracklets. The MARS bench-
mark introduced large-scale video evaluation and moti-
vated sequence aggregation strategies beyond simple frame
averaging [9]. Early approaches employed recurrent ar-
chitectures or heuristic pooling, while more recent works
adopt attention-based mechanisms to emphasize informa-
tive frames and suppress noise. Methods such as spatio-
temporal attention (STA) [10] and temporal complementary
learning networks (TCLNet) [11] demonstrate the effective-
ness of adaptive aggregation for handling occlusion, pose
variation, and motion blur. In extreme far-distance aerial
footage, tracklets often contain heavily degraded frames
due to motion instability and severe resolution loss. Un-
der such conditions, uniform aggregation can amplify noise,
making adaptive temporal weighting particularly important.

Table 1. Comparison of aerial-ground benchmark datasets for person detection, ReID, tracking, and action recognition.

Dataset	Camera	Format	PIDs	BBoxes	Height (m)	Distance (m)
AG-ReID.v2	UAV+CCTV	Still	1615	100.6K	15–45	–
G2APS-ReID	UAV+CCTV	Still	2788	200.8K	20–60	–
DetReIDX	DSLR+UAV	Video+Still	334	13M	5–120	10–120

155 2.3. Aerial and Cross-View ReID Benchmarks

156 UAV-based person analysis introduces challenges that are
 157 not captured by conventional ground-only datasets. UAV-
 158 Human [12] and P-DESTRE [13] extend ReID to aerial
 159 platforms, incorporating detection, tracking, and both short
 160 and long-term matching tasks. For explicit aerial-ground
 161 cross-view matching, AG-ReID [14] and related datasets re-
 162 veal the substantial appearance gap between top-down UAV
 163 imagery and horizontal ground cameras.

164 More recently, DetReIDX (see Table 1) was proposed as
 165 a stress-test benchmark [4] targeting extreme far-distance
 166 conditions. It explicitly models altitude variation, scale
 167 compression, cross-view mismatch, and session-level ap-
 168 pearance drift. Experimental findings show that strong de-
 169 tection and ReID baselines can collapse under such condi-
 170 tions, underscoring the need for scale-aware and stability-
 171 oriented adaptation strategies. These characteristics make
 172 DetReIDX an appropriate benchmark for evaluating robust-
 173 ness in extreme far-distance aerial-ground ReID.

174 2.4. Vision-Language Pretraining and Efficient 175 Adaptation

176 Vision-language models (VLMs), particularly CLIP, have
 177 demonstrated strong transferability through large-scale
 178 image-text pretraining [2]. CLIP-based features have been
 179 incorporated into ReID frameworks to improve generaliza-
 180 tion and mitigate the absence of semantic class labels. Ap-
 181 proaches such as CLIP-ReID leverage prompt learning to
 182 bridge identity supervision and semantic embedding spaces.

183 Beyond prompt optimization, parameter-efficient adap-
 184 tation methods, including adapters, low-rank updates, and
 185 conditional prompts (e.g., CoOp and CoCoOp) [15], enable
 186 stable fine-tuning of large pretrained backbones while re-
 187 ducing overfitting risk. Such strategies are particularly re-
 188 levant in extreme far-distance aerial-ground ReID, where se-
 189 vere degradation and limited effective visual detail increase
 190 the risk of representation drift during adaptation.

191 2.5. Re-Ranking for Retrieval-based ReID

192 Post-processing remains an important component of
 193 retrieval-based ReID systems. In particular, k-reciprocal
 194 re-ranking refines similarity relationships by exploiting re-
 195 ciprocal nearest neighbors in the embedding space, lead-
 196 ing to consistent improvements in mAP and ranking accu-
 197 racy. Careful parameter tuning is often required to max-

imize gains for a given backbone and feature distribution 198
[3]. 199

200 In contrast to prior work, which primarily focuses on
 201 ground-level settings or moderate aerial conditions, we
 202 target the extreme far-distance regime evaluated by De-
 203 tReIDX. We investigate how large-scale vision-language
 204 backbones can be systematically scaled and stably adapted
 205 to address severe scale compression and noisy video track-
 206 lets. Section 3 describes our proposed scale-aware adapta-
 207 tion framework in detail.

208 3. Scale-Aware Adaptation Framework

209 We build upon the official CLIP-based [2] video ReID base-
 210 line for the DetReIDX benchmark and introduce a series of
 211 scale-aware modifications tailored for extreme far-distance
 212 (XFD) aerial-ground scenarios. Our approach focuses on
 213 three key aspects: (1) increasing backbone capacity through
 214 model scaling, (2) stabilizing adaptation of the larger trans-
 215 former, and (3) improving robustness to degraded video
 216 tracklets via adaptive temporal aggregation and refined op-
 217 timization. We first describe the baseline framework and
 218 then present the proposed modifications.

219 3.1. Baseline Framework

220 We adopt the official CLIP-based video ReID baseline re-
 221 leased for DetReIDX. The framework performs tracklet-
 222 level retrieval under aerial-ground cross-view conditions
 223 and includes a vision transformer backbone, prompt-based
 224 cross-view conditioning, and parameter-efficient adaptation
 225 modules.

226 **Backbone and frame-level feature extraction.** The base-
 227 line employs a Vision Transformer backbone with patch
 228 size 16×16 (ViT-B/16). All frames are resized to 256×128
 229 during both training and testing [16].

230 **Tracklet construction and sampling.** Video ReID is per-
 231 formed at the tracklet level. For each training instance, a
 232 fixed-length sequence of $L = 16$ frames is sampled and
 233 encoded independently by the backbone. A softmax triplet
 234 sampling strategy is adopted to support both identity clas-
 235 sification and metric learning [17]. The batch size is set to
 236 16 tracklets. During inference, the same sequence length

Table 2. Baseline training configuration.

Parameter	Stage 1	Stage 2
Optimizer	Adam	Adam
Base Learning Rate	3.5e-4	1.0e-4
Max Epochs	120	120
Weight Decay	1e-4	2.5e-4
Weight Decay (Bias)	1e-4	1e-4
Images Per Batch	16	16
ID Loss Weight (λ_{id})	0.25	
Triplet Loss Weight (λ_{tri})	1.0	
I2T / T2I Weights	1.0	

237 is used, and retrieval is conducted using cosine similarity
238 between ℓ_2 -normalized features.

239 **Prompt-based cross-view conditioning (PBP).** To reduce
240 aerial-ground domain mismatch, the baseline incorporates
241 prompt-based cross-view conditioning. A prompt of length
242 1 is inserted across 9 transformer layers. Metadata signals,
243 including altitude, horizontal distance, and viewing angle,
244 are discretized into bins (18 altitude, 18 distance, and 3 angle
245 bins) and encoded as conditioning inputs. This design
246 injects physical context related to viewpoint and scale variation
247 into the representation learning process.

248 **Training strategy.** The baseline is trained using a multi-
249 term objective combining identity classification, triplet met-
250 ric learning, and CLIP-style cross-modal alignment. The
251 overall loss is defined as:

$$252 \quad \mathcal{L} = \lambda_{id}\mathcal{L}_{id} + \lambda_{tri}\mathcal{L}_{tri} + \lambda_{i2t}\mathcal{L}_{i2t} + \lambda_{t2i}\mathcal{L}_{t2i}, \quad (1)$$

253 where the loss weights and optimization parameters are
254 summarized in Table 2.

255 Stage 1 of training focuses on stable representation
256 learning with cosine-style decay, with some warm-up.
257 Stage 2 training performs further fine-tuning using a re-
258 duced learning rate and multi-step scheduling to refine re-
259 trieval performance and improve ranking stability.

260 **Inference settings.** During evaluation, tracklet features
261 are extracted using fixed-length sequences of 16 frames and
262 ℓ_2 normalized before retrieval. Cosine similarity is used
263 for distance computation. Re-ranking is disabled in the de-
264 fault baseline configuration, considering it may degrade the
265 learned representation. The full inference configuration is
266 summarized in Table 3.

267 **Parameter-efficient adaptation via adapters.** The base-
268 line further uses lightweight adapter modules [18] as a

Table 3. Baseline inference configuration.

Parameter	Setting
Batch Size	16
Sequence Length	16
Feature Normalization	ℓ_2 normalization
Distance Metric	Cosine
Re-Ranking	Disabled

parameter-efficient mechanism for adapting the pretrained
transformer to the target domain. Other optional compo-
nents provided by the baseline, such as VCAH and QATW,
are disabled initially. Temporal Attention pooling is also
disabled in the baseline, resulting in the use of a standard
mean pooling.

3.2. Scale-Aware Adaptation

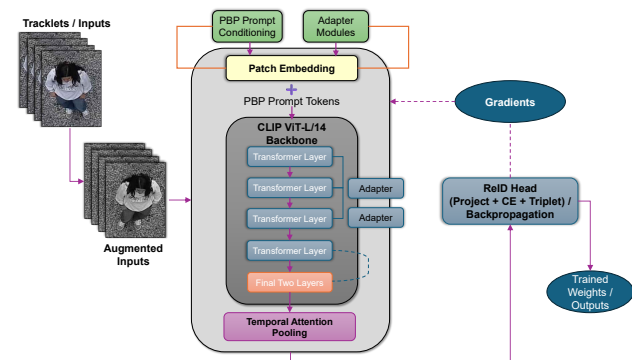


Figure 2. Overview of the proposed CLIP ViT-L/14-based video ReID framework. Augmented tracklets are encoded by a CLIP backbone enhanced with PBP prompt conditioning and adapter modules, with only the final two transformer blocks unfrozen during fine-tuning. Frame-level features are aggregated via temporal attention pooling and optimized using identity and triplet losses.

While the baseline provides a strong foundation for extreme far-distance video ReID, performance remains limited under severe scale compression and heavy frame degradation. We therefore introduce several targeted modifications focused on backbone capacity, stable large-model adaptation, robust temporal aggregation, re-ranking and optimized training strategies. An overview of the proposed pipeline is shown in Figure 2.

3.2.1. Backbone Scaling: ViT-B/16 \rightarrow ViT-L/14

We modify the baseline configuration by switching the CLIP visual encoder from ViT-B/16 to the larger ViT-L/14 backbone, thereby increasing representational capacity under severe resolution degradation and scale compression. In extreme far-distance settings, pedestrian regions often oc-

copy only a small number of pixels, and fine-grained appearance cues such as texture and local structure become unreliable. A larger transformer backbone provides improved modeling capacity to extract robust global and mid-level features from such degraded inputs [2, 16].

ViT-B/16 produces 768-dimensional embeddings, whereas ViT-L/14 outputs 1024-dimensional embeddings. Consequently, we adjust the associated projection layers, classification heads, prompt embeddings, and adapter modules to ensure dimensional consistency throughout the architecture. These adjustments are strictly architectural and do not alter the training objective or data sampling strategy. Apart from the backbone upgrade and the necessary dimensional alignment, the overall training pipeline is preserved. This design allows us to isolate the impact of backbone scaling while maintaining comparability with the original baseline.

3.2.2. Backbone-Aware Selective Fine-Tuning

Fine-tuning all parameters of a large transformer can degrade the strong feature representations learned during large-scale pretraining and lead to unstable optimization and overfitting under extreme far-distance conditions. Hence, instead of training the entire architecture on this dataset, the baseline had most backbone parameters frozen, focusing on adapter and prompt training, but we unfreeze the high-level blocks (i.e., blocks 22 and 23) of the architecture for adaptation. The final two blocks were unfrozen since the earlier layers encode general visual features. This strategy preserves the pretrained model’s generalization ability while enabling high-level representations to adapt to the target domain [18, 19]. We further use differential learning rates by assigning a smaller learning rate to the unfrozen backbone blocks and a larger learning rate to lightweight modules and heads (e.g., prompts/adapters and the identity classifier). Specifically, the unfrozen backbone blocks are trained using a learning rate scaled to $0.1\times$ of the base learning rate. This improves training stability when adapting the larger ViT-L/14 model and helps avoid catastrophic drift from pretrained weights.

3.2.3. Temporal Attention Pooling

The baseline uses a standard mean pooling because attention pooling is disabled. However, in extreme far-distance aerial settings, individual frames often suffer from severe blur, motion artifacts, occlusion, or extreme scale compression. Standard mean pooling can, therefore, allow degraded frames to negatively influence the final representation. To address this limitation, we enable and introduce a lightweight temporal attention mechanism that adaptively weights frames based on their feature responses [10, 11]. Let $\{f_t\}_{t=1}^T$ denote the sequence of frame-level embeddings extracted from a tracklet, where $T = 16$ in our setting and $f_t \in \mathbb{R}^C$.

Table 4. Optimization refinements compared to the baseline configuration.

Parameter	Baseline (ViT-B/16)	Ours (ViT-L/14)
Backbone	ViT-B/16	ViT-L/14
Stride Size	16	14
PBP Prompt Length	1	4
QATW Module	Disabled	Enabled
Instance Norm	Disabled	Enabled
Stage 1 Training		
Optimizer	Adam	Adam
Images Per Batch	16	48
Base Learning Rate	3.5e-4	2.0e-4
Max Epochs	120	50
LR Schedule	Cosine Annealing	Cosine Annealing
Weight Decay	1e-4	1e-4
Stage 2 Fine-Tuning		
Optimizer	Adam	Adam
Images Per Batch	16	24
Base Learning Rate	1.0e-4	1.0e-4
Max Epochs	120	40
LR Schedule	Multi-step (60, 90)	Cosine Annealing

We learn a scalar attention score for each frame using a linear projection:

$$s_t = w^\top f_t, \quad (2)$$

where $w \in \mathbb{R}^C$ is a learnable parameter vector implemented as a fully connected layer. The attention weights are obtained via softmax normalization across the temporal dimension:

$$\alpha_t = \frac{\exp(s_t)}{\sum_{k=1}^T \exp(s_k)}. \quad (3)$$

The final tracklet representation is computed as a weighted sum:

$$z = \sum_{t=1}^T \alpha_t f_t. \quad (4)$$

This mechanism enables the model to suppress degraded or irrelevant frames while emphasizing temporally consistent and discriminative observations. In our implementation, temporal attention is applied consistently to the intermediate and projected feature representations before classification and retrieval. Compared to the baseline, this adaptive pooling strategy improves robustness under extreme resolution degradation and unstable aerial motion.

3.2.4. Optimization Strategy

In addition to backbone scaling and selective fine-tuning, we refine the optimization strategy to better accommodate the larger ViT-L/14 backbone and improve training stability under extreme far-distance conditions. Compared to the

366 baseline configuration, we adjust batch size, learning rates,
 367 training duration, and scheduling strategy to balance con-
 368 vergence speed and generalization. Table 4 summarizes the
 369 key differences between the baseline and our refined opti-
 370 mization setup. Specifically, we increase the Stage 1 batch
 371 size from 16 to 48 to expose the larger backbone to more
 372 identity diversity per iteration. The base learning rate for
 373 Stage 1 is slightly reduced (from 3.5×10^{-4} to 2.0×10^{-4})
 374 to stabilize optimization when training the deeper ViT-L/14
 375 architecture. The total number of training epochs is also re-
 376 duced in both stages (Stage 1: $120 \rightarrow 50$, Stage 2: $120 \rightarrow$
 377 40), reflecting faster convergence with the scaled backbone.

378 In Stage 1, the selectively unfrozen backbone blocks are
 379 trained with a reduced learning rate scaled to $0.1 \times$ the base
 380 learning rate, as described in Section 3.2.2, to prevent desta-
 381 bilization of pretrained representations. Unlike the baseline,
 382 which uses multi-step decay during Stage 2, our refined
 383 configuration adopts cosine annealing [20] for both training
 384 stages. This smooth decay schedule avoids abrupt learning
 385 rate drops and provides more stable adaptation of the selec-
 386 tively unfrozen transformer blocks. Additionally, instance
 387 normalization at the neck is enabled to further improve ro-
 388 bustness under extreme scale degradation. These refine-
 389 ments collectively stabilize large-model adaptation and im-
 390 prove convergence behavior under severe resolution loss.

391 3.2.5. Data Augmentation

392 To improve robustness under diverse aerial capture condi-
 393 tions, we extend the baseline augmentation pipeline by in-
 394 corporating color jitter during training. Specifically, we ap-
 395 ply controlled perturbations in brightness, contrast, satura-
 396 tion, and hue to simulate illumination variation and camera-
 397 induced color shifts commonly observed in extreme far-
 398 distance imagery. In addition to random horizontal flipping
 399 and random erasing used in the baseline, these color aug-
 400 ments encourage the model to rely less on fragile color
 401 cues and more on structural and identity-consistent features.
 402 During inference, we further apply horizontal flip augmen-
 403 tation and average the features extracted from the original
 404 and flipped sequences to improve orientation robustness.

405 3.2.6. k-Reciprocal Re-Ranking

406 To further refine retrieval results, we enable k-reciprocal re-
 407 ranking as a post-processing step during inference. After
 408 computing cosine distances between ℓ_2 -normalized query
 409 and gallery embeddings, the initial distance matrix is re-
 410 fined using k-reciprocal encoding to improve neighborhood
 411 consistency in the embedding space. Specifically, we ap-
 412 ply re-ranking with parameters ($k_1 = 28$, $k_2 = 6$, $\lambda =$
 413 0.28). This procedure re-evaluates similarity relationships
 414 by considering reciprocal nearest neighbors, reducing the
 415 impact of local feature noise, and improving retrieval ro-
 416 bustness. Re-ranking is applied only during evaluation and
 417 does not affect model training. This inference refinement

yields additional improvements in the mean Average Preci- 418
 sion (mAP) score. 419

420 4. Experiment & Results

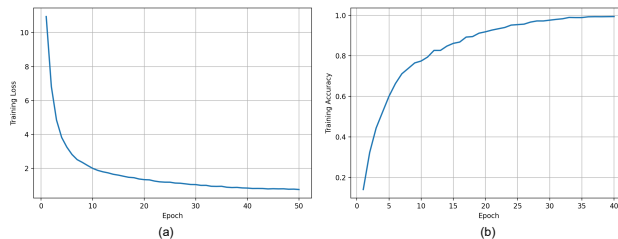


Figure 3. Training across the two-stage optimization process. (a) Stage 1 shows smooth loss convergence under cosine annealing scheduling, while (b) Stage 2 illustrates accuracy improvement during fine-tuning.

The experiments were conducted on a high-performance 421
 computing (HPC) cluster using NVIDIA A100 GPUs. Each 422
 training run utilized 8 CPU cores and 48 GiB of system 423
 memory. This setup was sufficient to train the scaled ViT- 424
 L/14 backbone under the proposed adaptation strategy. 425

426 4.1. DetReIDX Dataset

We train and evaluate our model on the DetReIDX stress- 427
 test benchmark [4], which is designed to evaluate aerial 428
 surveillance person ReID under extreme far-distance degra- 429
 dation and aerial-ground cross-view mismatch. The final 430
 benchmark score is reported as the average mAP across the 431
 three domains (A2G, G2A and A2A). We use mean Average 432
 Precision (mAP) as the primary evaluation metric, follow- 433
 ing standard ReID protocols. The final score is computed 434
 as the average mAP across the three evaluation cases. 435

$$\text{mAP}_{\text{avg}} = \frac{1}{3} (\text{mAP}_{\text{A2G}} + \text{mAP}_{\text{G2A}} + \text{mAP}_{\text{A2A}}). \quad (5) \quad 436$$

437 4.2. Implementation Details

Backbone and Training Setup. The visual encoder is ini- 438
 tialized from the CLIP ViT-L/14 checkpoint pretrained on 439
 large-scale image-text data. Training is conducted in two 440
 stages. In Stage 1, the model is trained for 50 epochs 441
 with a batch size of 48 tracklets and a base learning rate of 442
 2×10^{-4} . Stage 2 performs fine-tuning for 40 epochs using a 443
 batch size of 24 and a base learning rate of 1×10^{-4} . Selec- 444
 tive fine-tuning is applied by unfreezing only the final two 445
 transformer blocks (resblocks.22 and resblocks.23), while 446
 earlier layers remain frozen. The unfrozen blocks are op- 447
 timized with a learning rate scaled to $0.1 \times$ the base rate to 448
 preserve pretrained representations and stabilize adaptation. 449
 Adam is used as the optimizer in both training stages. Fig- 450
 ure 3 illustrates the smooth convergence behavior of Stage 1 451
 training loss and Stage 2 training accuracy. 452

Table 5. Main results on DetReIDX using the official averaged mAP across A2G/G2A/A2A (Eq. 5).

Method	mAP
CLIP ViT-B/16 (Official baseline)	28.11
DetReIDXv1	32.89
Scale-aware adaptation w/ CLIP ViT-L/14 (Ours)	35.73

453 **Prompt and Conditioning Configuration.** Prompt-based
454 cross-view conditioning is enabled in the final model with a
455 prompt length of 4 and deep prompt insertion across 9 trans-
456 former layers. Camera and metadata conditioning-based on
457 discretized altitude, distance, and viewing angle bins, are
458 activated as described in Section 3.1. Adapter modules, the
459 QATW module, and instance normalization are also enabled
460 in the final configuration.

461 **Learning Rate Scheduling and Regularization.** Cosine
462 annealing is used in both training stages to provide smooth
463 learning rate decay and stable convergence. Weight decay
464 is set to 1×10^{-4} in Stage 1 and 2.5×10^{-4} in Stage 2.

465 **Data Processing and Augmentation.** All frames are re-
466 sized to 256×128 and normalized. Each tracklet consists
467 of 16 sampled frames. During training, we apply random
468 horizontal flipping ($p = 0.5$), random erasing ($p = 0.5$),
469 color jitter with brightness, contrast, saturation, and hue pa-
470 rameters of $(0.1, 0.1, 0.1, 0.05)$, and padding of 10 pixels.
471 These augmentations improve robustness to noise and en-
472 hance generalization under aerial degradation conditions.

473 **Temporal Aggregation and Inference.** Temporal at-
474 tention pooling aggregates frame-level embeddings into
475 tracklet-level representations. A lightweight linear attention
476 module computes adaptive frame weights through softmax
477 normalization along the temporal dimension, enabling the
478 model to emphasize informative frames while suppressing
479 degraded ones. During inference, tracklet embeddings are
480 ℓ_2 -normalized, and cosine similarity is used for retrieval.

481 **K-reciprocal re-ranking.** To further improve retrieval
482 consistency, we apply k-reciprocal re-ranking as a post-
483 processing step during inference. After computing cosine
484 distances between ℓ_2 -normalized query and gallery em-
485 beddings, the distance matrix is refined using k-reciprocal
486 encoding. Unless otherwise specified, we set $k_1 = 28$,
487 $k_2 = 6$, and $\lambda = 0.28$. Re-ranking is applied only during
488 evaluation and does not affect model training.

Table 6. Ablation study on DetReIDX using averaged mAP across A2G/G2A/A2A (Eq. 5).

Variant	Backbone	mAP (%)	Δ
CLIP ViT-B/16		28.11	+0.00
+ Optimization		28.18	+0.07
+ k-reciprocal re-ranking	ViT-B/16	29.40	+1.29
+ Re-ranking tuned		29.47	+1.36
+ Temporal attention pooling		29.71	+1.60
+ Re-ranking tuned (final)		29.84	+1.73
+ Backbone scaling + Adaptation	ViT-L/14	35.73	+7.62

4.3. Results 489

We report results on DetReIDX using the official evaluation 490
metric defined in Eq. (5), namely the averaged mAP across 491
the three protocols (A2G, G2A, and A2A). Since the bench- 492
mark ranking is determined by this aggregated score, we fo- 493
cus on mAP_{avg} for all comparisons. Our full model achieves 494
an averaged mAP of 35.73, representing a +7.62 improve- 495
ment over the official CLIP ViT-B/16 baseline (28.11). 496
Compared to the strongest publicly reported result (32.89), 497
our method yields a further +2.84 gain. These results indi- 498
cate that backbone scaling combined with stability-focused 499
adaptation substantially enhances robustness under extreme 500
far-distance degradation and aerial-ground cross-view mis- 501
match. 502

4.4. Ablation Study 503

We perform ablation experiments to evaluate the contribu- 504
tion of each component in the proposed scale-aware adap- 505
tation framework. All results are reported using the of- 506
ficial DetReIDX metric, mAP_{avg} , averaged across A2G, 507
G2A, and A2A (Eq. 5). The ablation results indicate 508
that k-reciprocal re-ranking consistently improves perfor- 509
mance in the ViT-B/16 setting. Incorporating temporal at- 510
tention pooling yields additional gains by adaptively down- 511
weighting degraded frames within noisy tracklets. The most 512
substantial improvement comes from scaling the backbone 513
to ViT-L/14, demonstrating that increased representational 514
capacity, when paired with stability-oriented selective fine- 515
tuning, plays a critical role in handling extreme far-distance 516
scale compression. 517

5. Discussion 518

Our experiments suggest that extreme far-distance (XFD) 519
video person ReID benefits primarily from increased rep- 520
resentational capacity, provided that adaptation remains 521
stable. Scaling the CLIP visual encoder from ViT-B/16 522
to ViT-L/14 yields the largest improvement in the offi- 523
cial DetReIDX averaged score, indicating that additional 524
transformer depth and width help recover discriminative 525
cues even when pedestrian regions are severely scale- 526
compressed and fine-grained appearance details are unreli- 527

528 able. Importantly, these gains are realized only when back-
529 bone scaling is coupled with stability-oriented fine-tuning.
530 Selectively unfreezing the final transformer blocks and ap-
531 plying a reduced learning rate to the backbone preserves the
532 pretrained feature space while enabling high-level domain
533 adaptation.

534 Beyond backbone scaling, the ablation study high-
535 lights the complementary role of temporal attention and
536 re-ranking. Temporal attention pooling improves robust-
537 ness by down-weighting degraded or blurred frames within
538 tracklets, mitigating the impact of motion instability and
539 resolution loss during feature aggregation. In contrast, k-
540 reciprocal re-ranking enhances neighborhood consistency
541 in the embedding space at inference time, improving re-
542 trieval reliability. Although each component contributes a
543 smaller absolute gain than backbone scaling, together they
544 establish a stronger and more stable foundation for per-
545 formance improvements under extreme far-distance condi-
546 tions.

547 **Limitations.** This study is evaluated exclusively on the
548 DetReIDX stress-test benchmark using the official averaged
549 metric across the A2G, G2A, and A2A protocols. While
550 this aggregated score enables standardized comparison, it
551 limits detailed analysis of protocol-specific behavior and
552 makes it difficult to determine which domain transfer set-
553 ting benefits most from individual components. Moreover,
554 the final system employs a larger ViT-L/14 backbone and
555 optional re-ranking during inference, which increases com-
556 putational cost and latency compared to the baseline con-
557 figuration.

558 **Future work.** Several direction's may further advance ex-
559 treme far-distance ReID. First, multi-granularity temporal
560 modeling such as hierarchical aggregation across short-
561 and long-term temporal windows, could better handle long
562 tracklets with intermittent visibility. Second, more ex-
563 pressive scale-aware prompting or metadata-conditioned
564 adapters may further improve cross-view alignment under
565 substantial altitude and distance variation. Third, evalu-
566 ating cross-dataset generalization across additional aerial
567 benchmarks would provide stronger evidence of robustness
568 beyond DetReIDX. Finally, given the privacy implications
569 of aerial surveillance, future research should consider re-
570 sponsible deployment practices and privacy-aware evalua-
571 tion protocols when developing aerial person recognition
572 systems.

573 6. Conclusion

574 We presented a scale-aware adaptation framework for
575 CLIP-based video person re-identification in extreme far-
576 distance aerial-ground scenarios. By scaling the back-

bone from ViT-B/16 to ViT-L/14 and introducing stability-
577 oriented selective fine-tuning, temporal attention pooling,
578 and optimized training strategies, we significantly improve
579 robustness under severe scale compression and frame-level
580 degradation. On the DetReIDX stress-test benchmark,
581 our approach achieves an averaged mAP of 35.73 across
582 the three official protocols, substantially outperforming
583 the baseline. These findings demonstrate that large-scale
584 vision-language backbones, when paired with stable adap-
585 tation and adaptive temporal aggregation, provide an effec-
586 tive foundation for extreme far-distance video person ReID.
587

588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644**References**

- [1] Lin Wu, Yang Wang, Jianhuang Gao, and Xuelong Li. Deep learning for person re-identification: A survey and outlook. *arXiv preprint arXiv:2001.04193*, 2020. 1
- [2] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 1, 3, 5
- [3] 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)*, 2017. 2, 3
- [4] Khushal Hambarde, Emmanuel Mbongo, Nirav Menghani, Ajay Ramesh, Rogerio Schmidt Feris, and Hugo Proença. DetReIDX: A stress-test dataset for real-world UAV-based person recognition. *arXiv preprint arXiv:2505.04793*, 2025. 2, 3, 6
- [5] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015. 2
- [6] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *European Conference on Computer Vision Workshops (ECCV Workshops)*, 2016. 2
- [7] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by GAN improve the person re-identification baseline in vitro. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017. 2
- [8] Shuting He, Haoxi Wu, Peng Wang, Mengdan Zhang, Zhongzhan Huang, and Yonghong Tian. Transreid: Transformer-based object re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2
- [9] Liang Zheng, Zhi Bie, Yifan Sun, Jingdong Wang, Chi Su, Shengjin Wang, and Qi Tian. MARS: A video benchmark for large-scale person re-identification. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [10] Yang Fu, Xiaoyang Wang, Yunchao Wei, and Thomas Huang. Spatial-temporal attention model for video-based person re-identification. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2019. 2, 5
- [11] Ruibing Hou, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. TCLNet: Temporal complementary learning for video person re-identification. In *European Conference on Computer Vision (ECCV)*, 2020. 2, 5
- [12] Tianjiao Li, Jun Liu, Wei Zhang, Yun Ni, Wenqian Wang, and Zhiheng Li. UAV-Human: A large benchmark for human behavior understanding with unmanned aerial vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [13] S. V. Aruna Kumar, Ehsan Yaghoubi, Abhijit Das, B. S. Harish, and Hugo Proença. The P-DESTRE: A fully annotated dataset for pedestrian detection, tracking, and short/long-term re-identification from aerial devices. *IEEE Transactions on Information Forensics and Security (TIFS)*, 16:1696–1708, 2021. 3
- [14] Huy Nguyen, Kien Nguyen, Sridha Sridharan, and Clinton Fookes. Aerial-ground person re-identification. In *IEEE International Conference on Multimedia and Expo (ICME)*, pages 2585–2590, 2023. 3
- [15] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. *arXiv preprint arXiv:2203.05557*, 2022. 3
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. 3, 5
- [17] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017. 3
- [18] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mohammad Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning (ICML)*, 2019. 4, 5
- [19] Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations (ICLR)*, 2022. 5
- [20] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR)*, 2017. 6