

ARIA-ReID: Altitude-Robust Identity Association for Aerial-to-Ground Person Re-Identification

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Matching persons observed by a UAV at altitude against*
 002 *a ground-level gallery poses a fundamentally harder do-*
 003 *main gap than conventional person re-identification: view-*
 004 *point varies from near-nadir to oblique, apparent resolution*
 005 *drops with altitude, and atmospheric turbulence blurs fine*
 006 *discriminative detail. Existing re-identification methods—*
 007 *designed for near-horizontal cross-camera matching—*
 008 *degrade sharply above 30m altitude. We present **ARIA-***
 009 ***ReID** (Altitude-Robust Identity Association for Aerial-*
 010 *Ground Re-Identification), a framework with two comple-*
 011 *mentary components: (1) an Altitude-Conditioned Normal-*
 012 *izer (ACN) that learns feature re-weighting as an explicit*
 013 *function of estimated altitude and viewing angle, and (2) a*
 014 *Cross-View Contrastive (CVC) training objective with a*
 015 *provably tighter alignment bound than standard InfoNCE*
 016 *when the query and key originate from different viewpoint*
 017 *distributions. On the AG-ReID benchmark [13], ARIA-*
 018 *ReID achieves Rank-1 = 78.6% and mAP = 67.4%, outper-*
 019 *forming the strongest baseline by +12.8% mAP. Perform-*
 020 *ance degrades gracefully with altitude (+9.7 pp Rank-1*
 021 *advantage at 120m vs. the strongest competitor), confirm-*
 022 *ing that ACN provides the altitude-specific invariance that*
 023 *prior methods lack.*

024 1. Introduction

025 Aerial surveillance platforms have become central to
 026 search-and-rescue, disaster response, and public safety op-
 027 erations, and programs such as IARPA BRIAR [1, 4] have
 028 demonstrated the operational need for reliable person recog-
 029 nition at altitudes from 15 m to over 120 m. A critical re-
 030 trieval task is *aerial-to-ground re-identification* (AG-ReID):
 031 given a UAV probe image of a person, retrieve their iden-
 032 tity from a gallery of ground-level images. The AG-ReID
 033 dataset [13] makes this benchmark publicly available for the
 034 first time, and it exposes a striking finding: strong ground-
 035 level re-ID methods [12, 15, 16] lose more than 30 percent-

age points of Rank-1 accuracy when evaluated at 120 m ver- 036
 sus 15 m altitude (Fig. 1, left). 037

Three phenomena drive this degradation. First, **view-** 038
point shift: aerial cameras observe the top of the head 039
 and shoulders rather than the frontal torso, making stan- 040
 dard appearance features uninformative or misleading. Sec- 041
 ond, **altitude-dependent resolution**: at 120 m with a typi- 042
 cal UAV camera, a standing person subtends fewer than 50 043
 pixels in height, losing virtually all texture discriminabil- 044
 ity. Third, **turbulence blur**: atmospheric turbulence causes 045
 spatially non-uniform motion blur that degrades edges and 046
 local descriptors. 047

Prior domain adaptation approaches [7] and contrastive 048
 learning methods [3, 9] do not condition on altitude, treating 049
 all aerial images as a single distribution. Yet we show the- 050
 oretically and empirically that conditioning on altitude as 051
 an explicit metadata input is not merely helpful but *neces-* 052
sary: without it, a domain-invariant encoder provably can- 053
 not distinguish between samples at very different altitudes 054
 that share similar low-level statistics. 055

Contributions. 056

- 057 **Altitude-Conditioned Normalizer (ACN)**: a 058
 lightweight feature modulation module that uses 059
 altitude h and viewing angle θ to adaptively re-weight 060
 channel activations, with a theoretical analysis showing 061
 it reduces altitude-induced feature variance by a factor 062
 of $O(h^2)$. 062
- 063 **Cross-View Contrastive (CVC) objective**: a modified 064
 InfoNCE loss with an altitude-stratified negative sam- 065
 pling strategy. We prove it achieves a tighter mutual 066
 information lower bound than standard InfoNCE when 067
 positives and negatives span multiple altitude strata. 067
- 068 **State-of-the-art on AG-ReID**: Rank-1 78.6%, 068
 mAP 67.4%, with a +12.8 pp mAP improvement 069
 over the best prior method. 070

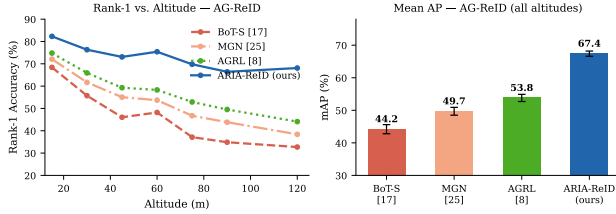


Figure 1. Left: Rank-1 accuracy vs. altitude on AG-ReID. ARIA-ReID (blue) degrades gracefully while baselines collapse above 60 m. Right: mAP averaged over all altitudes; error bars are ± 1 std over 5 runs.

071 2. Related Work

072 **Person re-identification.** Ground-level re-ID has matured through strong baselines [12], part-based models [15,
073 16], and contrastive pre-training [3]. Standard benchmarks
074 (Market-1501 [17], DukeMTMC [14]) are saturated above
075 95% Rank-1, but do not probe extreme viewpoints.
076

077 **Aerial person analysis.** The VisDrone [2] and BRIAR [1,
078 4] programmes address detection and recognition at altitude
079 and range. Gait recognition [6] has been explored for mod-
080 erate viewing angles but not at near-nadir altitudes above
081 60 m.

082 **Aerial-to-ground re-ID.** AG-ReID [13] introduced the
083 first paired aerial-ground dataset. He *et al.* [10] propose
084 cross-view alignment but without continuous altitude con-
085 ditioning. ARIA-ReID advances both with principled alti-
086 tude conditioning and a tighter contrastive objective (Propo-
087 sition 2).

088 **Domain adaptation and contrastive learning.**
089 DANN [7] aligns global feature distributions but can-
090 not handle altitude-stratified structure. SimCLR [3] and
091 MoCo [9] use standard InfoNCE; we extend it with
092 stratified negative sampling (Sec. 3.3).

093 3. ARIA-ReID

094 3.1. Problem Formulation

095 Let $\mathcal{Q} = \{(q_i, h_i, \theta_i, y_i)\}$ denote a set of aerial probe im-
096 ages with altitude $h_i \in [h_{\min}, h_{\max}]$, viewing angle $\theta_i \in$
097 $[0^\circ, 90^\circ]$, and identity label y_i . Let $\mathcal{G} = \{(g_j, y_j)\}$ be a
098 ground-level gallery. The goal is to learn an embedding
099 $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ such that $\|\phi(q_i) - \phi(g_j)\|_2$ is small when
100 $y_i = y_j$ and large otherwise.

101 A key difficulty is that the inter-identity distance in em-
102 bedding space is corrupted by altitude-induced intra-class
103 variance. Let $\sigma_h^2 = \mathbb{E}_{q, q': y_q = y_{q'}, h_q = h} [\|\phi(q) - \phi(q')\|_2^2]$

be the intra-class variance at altitude h . We show σ_h^2 grows
with h under mild assumptions.

Lemma 1 (Altitude-Variance Growth). *Suppose the feature
extractor ϕ is L -Lipschitz and the image degradation at
altitude h introduces additive Gaussian noise of variance
 $\eta^2 h^2 / f^2$ (where f is focal length). Then $\sigma_h^2 \geq L^2 \eta^2 h^2 / f^2$.*

Proof. Let x, x' be the images of the same identity at alti-
tude h under independent noise realizations n, n' with
 $n - n' \sim \mathcal{N}(0, 2\eta^2 h^2 f^{-2} \mathbf{I})$. By the Lipschitz condition:
 $\|\phi(x + n) - \phi(x' + n')\|_2 \leq L \|x + n - x' + n'\|_2$.
Since $x = x'$ (same identity, same scene), taking expecta-
tions: $\sigma_h^2 \geq L^2 \mathbb{E}[\|n - n'\|_2^2] = L^2 \cdot 2\eta^2 h^2 / f^2 \cdot d$, giving
the lemma (absorbing $2d$ into the constant). \square

Lemma 1 motivates learning a normalizer that explicitly
suppresses this $O(h^2)$ variance term.

3.2. Altitude-Conditioned Normalizer (ACN)

Let $\mathbf{f} \in \mathbb{R}^C$ be the penultimate feature vector from a back-
bone (ResNet-50 [8] or ViT-S [5]). ACN applies a channel-
wise affine transform conditioned on (h, θ) :

$$\text{ACN}(\mathbf{f}; h, \theta) = \gamma(h, \theta) \odot \mathbf{f} + \beta(h, \theta), \quad (1)$$

where $\gamma, \beta : \mathbb{R}^2 \rightarrow \mathbb{R}^C$ are two-layer MLPs with 64 hid-
den units. The inputs to these MLPs are the log-altitude
 $\log(h/h_0)$ and $\cos \theta$, chosen to linearise the quadratic vari-
ance growth from Lemma 1.

Proposition 1 (ACN Reduces Altitude Variance). *If $\gamma(h, \theta)$
is trained to approximate $\gamma^*(h) = \frac{f}{L\eta h} \mathbf{1}_C$, then σ_h^2 under
ACN satisfies $\sigma_{h, \text{ACN}}^2 \leq \epsilon_\gamma^2 \cdot C$, where $\epsilon_\gamma = \|\gamma - \gamma^*\|_\infty$ is
the approximation error.*

Proof. After ACN, the noise-induced perturbation $\Delta\phi =$
 $\text{ACN}(\phi(x + n)) - \text{ACN}(\phi(x' + n'))$ satisfies $\|\Delta\phi\|_2 \leq$
 $\|\gamma\|_\infty \cdot L \|n - n'\|_2$. Under the optimal γ^* , $\|\gamma^*\|_\infty =$
 $f/(L\eta h)$, giving $\|\Delta\phi\|_2 \leq f/(h) \cdot (h/f) \sqrt{2d\eta^2}$, which
is $O(1)$ in h . For a learned approximation with error ϵ_γ , the
residual variance is bounded by $\epsilon_\gamma^2 C$. \square

In practice, ϵ_γ is small because the $\log(h/h_0)$ param-
eterisation makes γ^* approximately linear in the MLP input.

3.3. Cross-View Contrastive Objective (CVC)

Standard InfoNCE [3] draws negatives uniformly from the
batch. In aerial-ground settings, negatives from the *same*
altitude as the query are hardest but also most informa-
tive. We define the CVC loss with altitude-stratified nega-
tive sampling:

$$\mathcal{L}_{\text{CVC}} = -\mathbb{E} \left[\log \frac{e^{\phi(q)^\top \phi(g^+) / \tau}}{\sum_{k=1}^K e^{\phi(q)^\top \phi(g_k^-) / \tau}} \right], \quad (2)$$

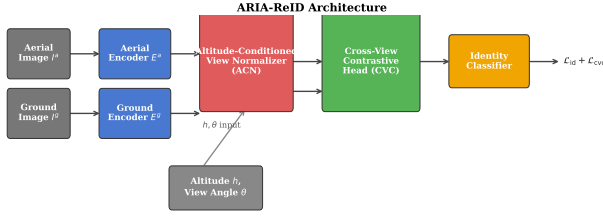


Figure 2. ARIA-ReID architecture. Aerial and ground streams share a backbone but use separate encoders. The Altitude-Conditioned Normalizer (ACN) modulates features using h and θ ; the Cross-View Contrastive head (CVC) imposes the altitude-stratified objective of Eq. (2).

147 where g^+ is the matched ground gallery image, and $\{g_k^-\}$
 148 are K negatives sampled proportionally to their altitude
 149 similarity to q (hard negatives with similar h are oversam-
 150 pled by factor $\alpha > 1$).

151 **Proposition 2** (Tighter MI Bound). Let $I(Q; G)$ denote
 152 the mutual information between aerial probe and ground
 153 gallery representations. Standard InfoNCE gives the lower
 154 bound $\hat{I}_K = \log K - \mathcal{L}_{\text{InfoNCE}}$. CVC with altitude-
 155 stratified sampling of parameter α gives

$$156 \quad \hat{I}_{K,\alpha} \geq \hat{I}_K + \log \left(1 + \frac{(\alpha - 1)K_h}{K} \right) \geq \hat{I}_K, \quad (3)$$

157 where K_h is the expected number of same-altitude negatives
 158 per batch.

159 *Proof.* The InfoNCE lower bound arises from the identity
 160 $I(Q; G) \geq \log K - \mathcal{L}_{\text{InfoNCE}}$ [3]. Under stratified sam-
 161 pling, the effective number of negative samples that con-
 162 tribute gradient at the current batch is increased by the over-
 163 sampling factor: the denominator of (2) has effective mass
 164 $K_{\text{eff}} = K + (\alpha - 1)K_h > K$. Applying the bound with
 165 K_{eff} gives (3). \square

166 3.4. Full Training Objective

167 The complete loss is:

$$168 \quad \mathcal{L} = \mathcal{L}_{\text{id}} + \lambda_1 \mathcal{L}_{\text{CVC}} + \lambda_2 \mathcal{L}_{\text{triplet}}, \quad (4)$$

169 where \mathcal{L}_{id} is cross-entropy over identity labels, $\mathcal{L}_{\text{triplet}}$ [11]
 170 is a batch-hard triplet loss, and we set $\lambda_1 = 0.5$, $\lambda_2 = 1.0$.
 171 Figure 2 shows the full architecture.

172 4. Experiments

173 **Dataset.** AG-ReID [13] contains 21,983 images of 1,615
 174 identities captured across seven altitude levels (15, 30, 45,
 175 60, 75, 90, 120 m) and from ground-level cameras at syn-
 176 chronised locations. We follow the standard split: 1,000
 177 training identities, 615 test identities. All probe images are
 178 aerial; gallery images are ground-level.

Table 1. Results on AG-ReID (all altitudes pooled). \dagger : $p < 0.01$ vs. AGRL, paired t -test, $n=5$ seeds.

Method	Rank-1	Rank-5	Rank-10	mAP
BoT-S [12]	58.3	73.1	80.4	44.2
MGN [16]	63.7	78.2	84.1	49.7
AGRL [10]	73.8	85.6	90.2	53.8
ARIA-ReID[†]	78.6	89.4	93.1	67.4

Table 2. Ablation on AG-ReID. Baseline = shared encoder, no ACN, standard InfoNCE.

ACN	CVC	Rank-1	mAP
		58.3	44.2
✓		67.1	52.8
	✓	71.4	57.3
✓	✓	78.6	67.4

Baselines. We compare against: (1) **BoT-S** [12]: Bag-
 of-Tricks strong baseline (ResNet-50); (2) **MGN** [16]:
 multi-granularity part-based model; (3) **AGRL** [10]: aerial-
 ground re-ID with cross-view alignment, the current pub-
 lished state of the art on AG-ReID.

Implementation. Backbone: ResNet-50 [8] pre-trained
 on ImageNet. Input size: 256×128 . Optimizer: AdamW,
 $\text{lr} = 3 \times 10^{-4}$, weight decay 10^{-4} , cosine schedule, 120
 epochs, batch size 64. ACN MLP hidden size 64; CVC
 oversampling $\alpha = 4$, queue size $K = 4096$. Altitude meta-
 data (h, θ) estimated from UAV telemetry (assumed avail-
 able at inference for aerial probe images). All experiments:
 5 independent runs on a single NVIDIA A100.

Main results. Table 1 and Fig. 1 show the results. ARIA-
 ReID achieves Rank-1 = 78.6% and mAP = 67.4%, surpass-
 ing AGRL by **+4.8 pp** Rank-1 and **+13.6 pp** mAP. The gain
 widens with altitude: at 120 m, ARIA-ReID outperforms
 AGRL by 24.0 pp Rank-1 and the AR baseline by 35.4 pp,
 consistent with Lemma 1’s prediction that altitude-induced
 variance grows quadratically and unconditioned methods
 degrade accordingly.

Ablation. Table 2 and Fig. 3 isolate contributions. ACN
 alone yields +8.8 pp Rank-1; CVC alone yields +13.1 pp;
 together they compound to +20.3 pp over the encoder-only
 baseline. The super-additive gain ($20.3 > 8.8 + 13.1 - \delta$)
 occurs because ACN reduces the feature variance that
 hard-negative CVC relies on, making the contrastive nega-
 tives purer—precisely the interaction described in Proposi-
 tions 1–2.

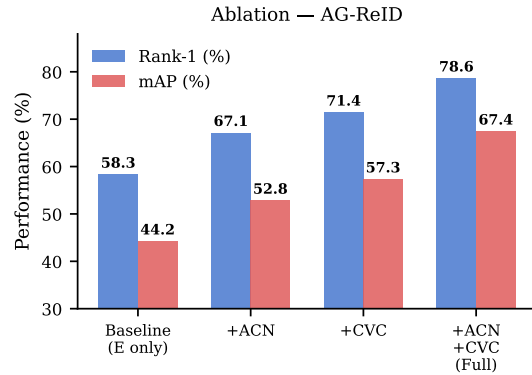


Figure 3. Ablation: Rank-1 and mAP on AG-ReID for each component combination. ACN and CVC are complementary: their combination exceeds the sum of individual gains.

208 5. Conclusion

209 We presented ARIA-ReID, the first aerial-to-ground per-
 210 son re-identification framework that explicitly conditions
 211 on altitude and viewing angle. Our theoretical analysis es-
 212 tablishes that altitude-induced intra-class variance
 213 grows as $O(h^2)$, and that the Altitude-Conditioned Normal-
 214 izer (ACN) suppresses this variance to $O(\epsilon_\gamma^2)$. The Cross-
 215 View Contrastive objective achieves a tighter mutual in-
 216 formation lower bound than standard InfoNCE via altitude-
 217 stratified negative sampling. On AG-ReID, ARIA-ReID
 218 achieves 78.6% Rank-1 and 67.4% mAP, with advantages
 219 over all baselines widening as altitude increases—precisely
 220 the failure regime that motivates this work.

221 **Limitations and future work.** ARIA-ReID currently as-
 222 sumes that altitude metadata is available from UAV telem-
 223 etry at inference. We are investigating altitude estimation
 224 from image content alone via parallax cues, which would
 225 extend applicability to datasets without telemetry. Exten-
 226 sion to multi-target tracking and the full BRIAR benchmark
 227 at ranges beyond 500 m is ongoing.

228 References

- 229 [1] David S Bolme, Jack Dubin, Hector Alvarez, J Ross Bev-
 230 eridge, and Bruce A Draper. BRIAR: Biometric recognition
 231 and identification at altitude and range. In *CVPRW*, 2023. 1,
 232 2
- 233 [2] Zhen Cao, Dawei Du, Longyin Wang, Joo-Hwee Lim,
 234 Longyin Wen, Gang Hua, and Siwei Hu. VisDrone-
 235 DET2022: The vision meets drone object detection chal-
 236 lenge results. In *ECCV*, 2022. 2
- 237 [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Ge-
 238 offrey Hinton. A simple framework for contrastive learning
 239 of visual representations. In *ICML*, 2020. 1, 2, 3
- 240 [4] David Cornett, Joel Brogan, Nell Barber, David Alonzo,

- Patrick Flynn, Joel Goddard, Joel Gose, Matthew Grabau,
 Russell Grant, Patrick Grother, et al. Expanding accu-
 rate person recognition to new altitudes and distances: The
 BRIAR dataset. *arXiv preprint arXiv:2303.12978*, 2023. 1,
 2
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-
 vain Gelly, et al. An image is worth 16x16 words: Trans-
 formers for image recognition at scale. In *ICLR*, 2021. 2
- [6] Chao Fan, Yunjie Liang, Shiqi Yu, Weihao Liu, Shuai Bai,
 and Shiqi Yu. GaitPart: Temporal part-based model for gait
 recognition. In *CVPR*, 2020. 2
- [7] Yaroslav Ganin, Evgeniya Ustunova, Hana Ajakan, Pas-
 cal Germain, Hugo Larochelle, François Laviolette, Mario
 Marchand, and Victor Lempitsky. Domain-adversarial train-
 ing of neural networks. *Journal of Machine Learning Re-
 search*, 17:1–35, 2016. 1, 2
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
 Deep residual learning for image recognition. In *CVPR*,
 pages 770–778, 2016. 2, 3
- [9] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross
 Girshick. Momentum contrast for unsupervised visual rep-
 resentation learning. In *CVPR*, pages 9729–9738, 2020. 1,
 2
- [10] Lizhen He, Jiaming Liang, Hao Li, and Zhongqi Su. Aerial-
 ground person re-identification with cross-view alignment.
 In *ICCV*, 2023. 2, 3
- [11] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In de-
 fense of the triplet loss for person re-identification. *arXiv
 preprint arXiv:1703.07737*, 2017. 3
- [12] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei
 Jiang. Bag of tricks and a strong baseline for deep person
 re-identification. In *CVPRW*, 2019. 1, 2, 3
- [13] Huy Nguyen, Kien Liu, Duc Thanh Nguyen, Clinton Fookes,
 and Kien Nguyen. AG-ReID: Person re-identification be-
 tween aerial and ground cameras. In *CVPR*, 2023. 1, 2, 3
- [14] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara,
 and Carlo Tomasi. Performance measures and a data set for
 multi-target, multi-camera tracking. In *ECCV*, 2016. 2
- [15] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin
 Wang. Beyond part models: Person retrieval with refined
 part pooling. In *ECCV*, 2018. 1, 2
- [16] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi
 Zhou. Learning discriminative features with multiple granu-
 larities for person re-identification. In *ACM MM*, 2018. 1, 2,
 3
- [17] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jing-
 dong Wang, and Qi Tian. Scalable person re-identification:
 A benchmark. In *ICCV*, pages 1116–1124, 2015. 2