

IntensPure: Attack Intensity-aware Secondary Domain Adaptive Diffusion for Adversarial Purification

(Appendix)

	Query	Top 10 retrieval results	SIM_{top}	$INCON_{top}$
			sum ↑	sum ↓
Benign			24.70	0.08
Metric-FGSM ($\epsilon = 4$)			20.86	3.88
Metric-FGSM ($\epsilon = 16$)			12.76	6.51
Deep Mis-Ranking ($\epsilon = 4$)			19.53	3.52
Deep Mis-Ranking ($\epsilon = 16$)			13.05	6.79
MetaAttack ($\epsilon = 4$)			21.04	2.77
MetaAttack ($\epsilon = 16$)			11.78	7.00

Figure 1: Illustration of the top 10 retrieval results for query images on Market1501 dataset. The SIM and $INCON$ are components of ID stability and attribute inconsistency, respectively. The SIM sum and $INCON$ sum represent the summation of their respective result vectors. As the attack intensity increases, there is a tendency for the SIM value to decrease, while the $INCON$ value tends to increase.

	Query	Bottom 10 retrieval results	SIM_{bottom}	$INCON_{bottom}$
			sum ↑	sum ↓
Benign			17.33	6.12
Metric-FGSM ($\epsilon = 4$)			12.03	6.98
Metric-FGSM ($\epsilon = 16$)			7.28	7.60
Deep Mis-Ranking ($\epsilon = 4$)			11.72	6.44
Deep Mis-Ranking ($\epsilon = 16$)			7.49	7.87
MetaAttack ($\epsilon = 4$)			12.26	6.69
MetaAttack ($\epsilon = 16$)			6.94	8.02

Figure 2: Illustration of the bottom 10 retrieval results for query images on Market1501 dataset. The SIM and $INCON$ are components of ID stability and attribute inconsistency, respectively. The SIM sum and $INCON$ sum represent the summation of their respective result vectors. As the attack intensity increases, there is a tendency for the SIM value to decrease, while the $INCON$ value tends to increase.

A Analysis of the Attack Intensity Estimator

A.1 Visualizing Retrieval Results

Figure 1 illustrates the inconsistency in retrieval results caused by adversarial attacks on the query set of the Market1501 dataset [Zheng *et al.*, 2015]. We covered the face area for privacy protection. In each row, the first image represents the query image, followed by the top 10 retrieved images. For benign queries, the top 10 retrieval results typically align with the query image identity (ID). Even if the retrieved image within the top 10 does not match the ID, it still returns a person who is visually very similar. However, adversarial attacks on person re-identification (Metric-FGSM [Bai *et al.*, 2020], Deep Mis-Ranking [Wang *et al.*, 2020], and MetaAttack [Yang *et al.*, 2022]), which can be abbreviated as person re-ID attacks, cause the retrieval system to return images unrelated to the query image in the top ranking list, such as retrieval results entirely different from the target person, or the wrong images, such as irrelevant objects like body parts or bicycles. Furthermore, the top 10 retrieved images not only

become irrelevant to the query but also exhibit reduced consistency within the results.

The sum values of SIM_{top} and $INCON_{top}$, which are the components of ID stability and attribute inconsistency, are also shown in Figures 1 and 2. For the benign query, the SIM_{top} value is relatively higher than the values that come from perturbed queries and tends to decrease as the attack intensity increases. Further, the $INCON_{top}$ value, which represents the error in attribute recognition, is close to 0 for benign queries, whereas the $INCON_{top}$ value tends to have large values according to the attack intensity.

Figure 2 shows that the bottom 10 ranking list exhibits consistency between retrieval results for benign queries, similar to the top 10 ranking results. However, when the query is perturbed, the bottom 10 ranking list exhibits inconsistency between retrieval results. The SIM_{bottom} and $INCON_{bottom}$ can be obtained from the bottom 10 ranking list and they show similar trends to SIM_{top} and $INCON_{top}$. By extracting features from these secondary effects of adversarial attacks, we can effectively detect attacks and estimate their intensities.

20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

ID stability (SIM)			Attribute inconsistency (INCON)			Acc \uparrow	MAE \downarrow
$Q-R$	top	bottom	$Q-R1$	top	bottom		
✓	-	-	-	-	-	92.13	4.094
✓	✓	-	-	-	-	94.89	3.508
✓	✓	✓	-	-	-	96.91	3.154
-	-	-	✓	-	-	77.14	7.887
-	-	-	✓	✓	-	89.18	5.337
-	-	-	✓	✓	✓	90.82	5.253
✓	-	-	✓	-	-	93.74	2.019
✓	✓	-	✓	✓	-	97.61	1.195
✓	✓	✓	✓	✓	✓	98.85	0.806

Table 1: Ablation study on attack intensity estimator in terms of detection accuracy (Acc) and mean absolute error (MAE) on Market1501 dataset attacked by Metric-FGSM. The hyphen - indicates without, and the check ✓ indicates with the method. The best results are in bold.

Number of rank images	1	5	10	15	20
Acc \uparrow	81.65	97.02	98.85	98.77	98.49
MAE \downarrow	3.113	1.214	0.806	0.820	0.856

Table 2: Comparison of adversarial detection accuracy (Acc) and mean absolute error (MAE) according to the number of top and bottom rank images on Market1501 dataset attacked by Metric-FGSM. The best results are in bold.

40 A.2 Ablation Study of ID Stability and Attribute 41 Inconsistency Components

42 Table 1 illustrates an ablation study to evaluate the contribution of the components used for adversarial attack intensity estimation. First, we investigated the performance of ID 43 stability and attribute inconsistency individually. Detection 44 accuracy refers to the binary classification accuracy for 45 detecting whether an image is an adversarial example and mean 46 absolute error is calculated using the estimated attack intensity 47 and ground truth attack intensity. In both individual and joint 48 experiments, the performance is more improved by additionally 49 considering features derived from the consistency between 50 retrieval results themselves than considering only the 51 features between the query and the top retrieval results ($Q-R$ 52 and $Q-R1$). Using the features derived from the top ranking 53 list, the detection accuracy is achieved at 94.89%, 89.18%, 54 and 97.61%, respectively. Additionally, using both the features 55 derived from the bottom ranking list and the top ranking list, the 56 detection performance is improved by 2.02%, 1.64%, and 1.24%, 57 respectively, and the attack intensity estimation error is also 58 reduced by 0.354, 0.084, and 0.389, respectively. This 59 suggests the effectiveness of leveraging the bottom ranking 60 list. Attribute inconsistency exhibits relatively lower 61 performance than ID stability when used alone due to the 62 inherent variability of person attributes. However, when 63 used in conjunction with ID stability, attribute inconsistency 64 significantly improves the performance of both adversarial 65 detection and attack intensity estimation.

Number of BDCT coefficients	1	3	6	10
Rank-1 accuracy (%)	49.12	60.87	66.65	65.63

Table 3: Rank-1 accuracy (%) variation with the number of BDCT coefficients used for secondary image generation of IntensPure, evaluated on Market1501 dataset attacked by Metric-FGSM with $\epsilon = 8$. The best results are in bold.

BDCT block size	Rank-1 accuracy (%) \uparrow	FLOPs (G) \downarrow	Params (M) \downarrow	Time (ms) \downarrow
4×4	66.58	154	751	98
8×8	66.65	39	751	59
16×16	60.87	10	751	34

Table 4: Rank-1 accuracy (%) and complexity variation with BDCT block size of IntensPure, evaluated on Market1501 dataset attacked by Metric-FGSM with $\epsilon = 8$. The complexity is assessed using a single diffusion time step without the estimator. The best results are in bold.

68 B Estimator Performance Comparison for the 69 Number of Retrieval Results

70 Table 2 examines the influence of rank range on the performance 71 of both adversarial detection and attack intensity estimation. Rank 72 range 1 refers to using only the top 1 and bottom 1 retrieval 73 results. In general, performance improves as the range increases 74 up to 10. We achieve the highest accuracy of 98.85% and the 75 lowest estimation error of 0.806, when the rank range is 10. 76 However, extending the rank range beyond 10 leads to a decline 77 in performance due to a weakened correlation between the retrieved 78 images.

79 C Performance Evaluation of BDCT 80 Coefficient Counts and Block Size

81 Table 3 represents the variation of rank-1 accuracy according 82 to the varying number of block discrete cosine transform (BDCT) 83 coefficients employed for generating secondary images within the 84 Market1501 dataset attacked by Metric-FGSM with attack intensity 85 $\epsilon = 8$. The rank-1 accuracy is highest when the six lowest-frequency 86 BDCT coefficients were used for generating secondary images. This 87 result indicates that there is a trade-off between using sufficient 88 coefficients to maintain benign information and using as few 89 coefficients as possible to minimize the impact of attacks. Based on 90 the experimental results, we selected the 6 BDCT coefficients 91 for secondary domain purification.

92 Table 4 shows the variation of rank-1 accuracy and complexity 93 variation according to the BDCT block sizes (4×4 , 94 8×8 , 16×16). As the block size becomes smaller, the 95 number of blocks increases, leading to higher resolution in the 96 secondary image. Therefore, when the block size is 4×4 , the 97 complexity is the largest, resulting in a slight decline in re-ID 98 performance. On the other hand, a block size of 16×16 reduces 99 complexity but marginally lowers re-ID performance, as it relies 100 on limited information. Selecting an 8×8 block 101 size maximizes the enhancement in rank-1 accuracy, as it optimally 102 utilizes information.

103 Figures 3 and 4 show the visualization of this secondary 104 domain adaptive diffusion process for the Market1501 dataset 105

Filter size	3	5	7	9
Rank-1 accuracy (%)	66.37	66.65	66.60	66.47

Table 5: Rank-1 accuracy (%) variation with the inter-block directional filter size of IntensPure, evaluated on Market1501 dataset attacked by Metric-FGSM with $\epsilon = 8$. The best results are in bold.

σ_r / σ_s	0.05 / 1	0.10 / 2	0.15 / 3	0.20 / 4
Rank-1 accuracy (%)	66.46	66.51	66.65	66.43

Table 6: Rank-1 accuracy (%) variation with the inter-block directional filter sigma values of IntensPure, evaluated on Market1501 dataset attacked by Metric-FGSM with $\epsilon = 8$. The best results are in bold.

and the DukeMTMC dataset [Ristani *et al.*, 2016], respectively. Along the diffusion process, the perturbed secondary images pass through the forward diffusion process, where random noise is injected. Subsequently, the denoising process is applied to remove both random noise and perturbations. In this process, the color change induced by the adversarial attack is erased. Finally, inter-block directional filters enhance the correlation between blocks, effectively removing remaining perturbations. After undergoing these diffusion processes, the purified secondary images are reconstructed into a spatial domain image via an inverse BDCT process and further enhanced using a perturbation constraint set. The reconstructed image is subsequently fed into the re-ID model.

106
107
108
109
110
111
112
113
114
115
116
117
118
119

D Parameters Analysis of Inter-block Filter

120 We investigate the influence of inter-block directional filter 121 size and sigma values (σ_r and σ_s) on performance. Concretely, 122 the experiments are conducted on the Market1501 123 dataset attacked by Metric-FGSM with attack intensity $\epsilon = 8$. 124

125 Table 5 presents rank-1 accuracy variations according to 126 the inter-block directional filter size (3, 5, 7, and 9). The 127 experimental results indicate that the highest rank-1 accuracy 128 is achieved with the size 5. Through these empirical results, 129 we determined 5 as the optimal filter size.

130 Table 6 evaluates effects by varying the setting of the range 131 and spatial standard deviations (σ_r and σ_s) included in inter- 132 block directional filter. In this experiment, the highest accuracy 133 of 66.65% is achieved with σ_r and σ_s at 0.15 and 3. Based 134 on these results, we experimentally choose these optimal 135 settings.

E Additional Ablation Study

136 We conduct an additional ablation study to evaluate the effectiveness of each module by removing only one module at 137 a time on the Market1501 dataset attacked by Metric-FGSM 138 with various ϵ values (0, 4, 8, 12, and 16). Table 7 shows the 139 rank-1 accuracy of IntensPure with different configurations. 140 The first row shows the accuracy of the complete model, 141 in which all modules are activated. The second row shows 142 the rank-1 accuracy when only the attack intensity estimation 143 module is deactivated and the time step is fixed at 50. 144 Additionally, the perturbation constraint set is not available 145 when the attack intensity estimator is deactivated, as it relies 146

Attack Intensity Estimation	Secondary Image Generation	Directional Diffusion	Inter-Block Directional Filter	Perturbation Constraint Set	Rank-1 Accuracy
✓	✓	✓	✓	✓	70.04
✓	✓			✓	64.01
✓	✓		✓	✓	64.20
✓	✓	✓		✓	68.64
✓	✓	✓	✓		69.26
✓	✓	✓	✓		68.31

Table 7: Ablation study of IntensPure, removing one module at a time. Model performances are evaluated on Market1501 dataset attacked by Metric-FGSM with various ϵ values (0, 4, 8, 12, and 16). The check ✓ indicates with the method.

on the estimated attack intensity. As expected, the accuracy decreases when the attack intensity estimation module is deactivated due to its fixed-step purification of all images, including benign ones, similar to previous diffusion-based purification methods. However, as presented in the main paper, Diffpure [Nie *et al.*, 2022] achieves 55.48%, GNSP [Lee and Kim, 2023] achieves 62.53%, and IntensPure still outperforms these state-of-the-art models with a rank-1 accuracy of 64.01%. This demonstrates that IntensPure can achieve notable performance gains with secondary domain adaptive diffusion alone, without relying on optimal purification strength estimation. The third row shows the performance of IntensPure when only the attack intensity estimation module for the time step optimization and perturbation constraint set are activated. This shows that IntensPure can achieve significant performance gains by adjusting purification strength. The difficulty in performing directional diffusion in the spatial domain comes from the deactivation of the secondary image generation module, because individual channels in the spatial domain is difficult to be considered that they have specific directional characteristics. Furthermore, the inter-block directional filter is not applicable to the spatial domain images because they do not need to consider the correlation between blocks. The fourth and fifth rows show the performance of IntensPure when directional diffusion and inter-block directional filter are deactivated, respectively. Directional diffusion leads to a performance improvement of 1.36%p, and inter-block directional filter leads to a performance improvement of 0.74%p. The sixth row shows that the perturbation constraint set further improves the performance by 1.79%p by suppressing over-purification and preserving the original information.

F Visual Comparison of Retrieval Results through Purification

For a qualitative comparison, we present the visual retrieval results of IntensPure and state-of-the-art models. Figures 5, 182 6, 7, 8, 9, and 10 illustrate the top 10 retrieval results on 183 Market1501 dataset, which is attacked by Metric-FGSM, Deep 184 Mis-Ranking, and MetaAttack with various attack intensities. 185 The results demonstrate that IntensPure consistently outperforms 186 other state-of-the-art models in terms of rank-1 accuracy, 187 under various attack types and intensities. Its superiority 188 is further presented via the secondary domain adaptive 189 diffusion. While other models struggle to restore original colors, 190 leading to inconsistent retrieval results, IntensPure effec- 191

136
137
138
139
140
141
142
143
144
145
146

148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179

Method	Person Re-ID Attack		Person Re-ID Attack + Attribute Recognition Attack	
	Acc \uparrow	MAE \downarrow	Acc \uparrow	MAE \downarrow
MEAAD	98.50	3.189	99.20	2.714
Ours	99.55	0.769	99.13	1.850

Table 8: Comparison of adversarial detection accuracy (Acc) and mean absolute error (MAE) according to the adaptive attack for proposed attack intensity estimator on Market1501 dataset attacked by Deep Mis-Ranking.

tively purifies color perturbations, maintaining consistency in the retrieval outcomes. Furthermore, compared to other models, IntensPure generally exhibits superior visual consistency among the top 10 retrieval results. This consistency enhances the reliability of the retrieval process.

198 G Discussion

199 G.1 Influence of Adaptive Attack on Attack 200 Intensity Estimation

201 To evaluate the robustness of the attack intensity estimator
202 against adaptive attacks, we generated images subjected to
203 both re-ID attacks on the primary task model and classi-
204 fication (attribute recognition) attacks on the auxiliary task
205 model. Subsequently, experiments are conducted on attack
206 detection and intensity estimation using these manipulated
207 images. We perturbed the query set in the Market1501 dataset
208 using Deep Mis-Ranking, which not only enables re-ID at-
209 tacks but also offers the additional option of attribute recog-
210 nition attacks. The intensity of attribute attacks and re-ID
211 attacks was set to be equal and the rest of the setting is the
212 same as the experiment in the main paper.

213 As shown in Table 8, the incorporation of attribute recog-
214 nition attacks marginally enhanced the detection performance
215 of MEAAD [Wang *et al.*, 2021], the comparative model.
216 This improvement is attributed to the fact that the perturba-
217 tion caused by attribute recognition attacks does not serve as
218 an adaptive attack against MEAAD but aids in its detection.
219 Although the proposed attack intensity estimator exhibited a
220 slight performance decline, it still demonstrated excellent ca-
221 pabilities in attack detection. The detection accuracy for ad-
222 versarial examples subjected to both types of attacks reached
223 99.13%, reflecting a marginal decrease of 0.42% p , while the
224 error in attack intensity estimation increased to 1.850, mark-
225 ing a rise of 1.081. While the estimation performance is
226 slightly lower than that of existing general attacks for adap-
227 tive attacks, it is still meaningful because it maintains superior
228 performance compared to the comparative model and exhibits
229 a relatively lower error. This is because features are extracted
230 based on the relationships between the retrieval results. How-
231 ever, there is still room for improvement in accurate attack in-
232 tensity estimation for adversarial attacks that also affect aux-
233 illary tasks.

234 G.2 Influence of Fitting Error on Optimal 235 Diffusion Time Step Selection

236 In the main paper, we derive some graphs for the optimal time
237 step t^* according to various attack intensities of each attack
238 type and fit it to a single logarithmic curve. All of the graphs

Dataset	Average fitting error (time step)	Average rank-1 accuracy error (%)
Market1501	3.6	0.05
DukeMTMC	2.9	0.08

Table 9: Rank-1 accuracy error according to the fitting error and their average values, evaluated on Market1501 and DukeMTMC dataset .

show similar and consistent logarithmic trends, so the fitted curve can approximately cover all variations. As the fitted curve is represented by a single line, discrepancies arise between t^* and the time step calculated using the fitted curve, particularly at specific outlying points. We define the error occurring between t^* and the calculated time step as fitting error, and Table 9 presents the rank-1 accuracy error based on the fitting error and their average values for Market1501 and DukeMTMC dataset. Nevertheless, the performance of the re-ID model is rarely impacted by the observed fitting error below 0.08

239 References

- [Bai *et al.*, 2020] Song Bai, Yingwei Li, Yuyin Zhou, Qizhu Li, and Philip HS Torr. Adversarial metric attack and defense for person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(6):2119–2126, 2020.
- [Lee and Kim, 2023] Minjong Lee and Dongwoo Kim. Robust evaluation of diffusion-based adversarial purification. *ICCV*, 2023.
- [Nie *et al.*, 2022] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. Diffusion models for adversarial purification. *ICML*, 2022.
- [Ristani *et al.*, 2016] 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.
- [Wang *et al.*, 2020] Hongjun Wang, Guangrun Wang, Ya Li, Dongyu Zhang, and Liang Lin. Transferable, controllable, and inconspicuous adversarial attacks on person re-identification with deep mis-ranking. In *CVPR*, 2020.
- [Wang *et al.*, 2021] Xueping Wang, Shasha Li, Min Liu, Yaonan Wang, and Amit K Roy-Chowdhury. Multi-expert adversarial attack detection in person re-identification using context inconsistency. In *ICCV*, 2021.
- [Yang *et al.*, 2022] Fengxiang Yang, Juanjuan Weng, Zhun Zhong, Hong Liu, Zheng Wang, Zhiming Luo, Donglin Cao, Shaozi Li, Shin’ichi Satoh, and Nicu Sebe. Towards robust person re-identification by defending against universal attackers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):5218–5235, 2022.
- [Zheng *et al.*, 2015] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015.

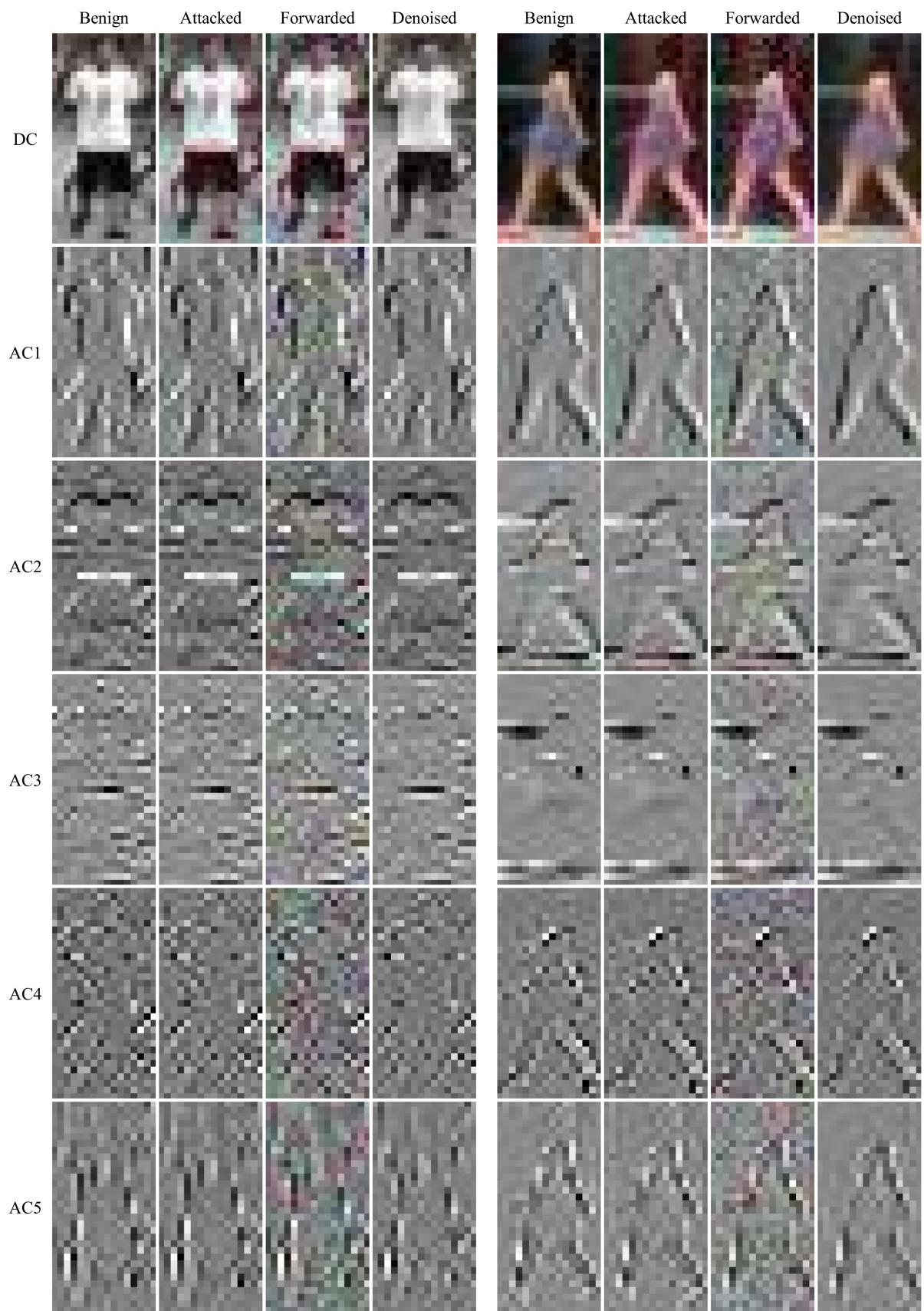


Figure 3: Illustration of the secondary diffusion process on **Market1501 dataset** attacked by Deep Mis-Ranking with $\epsilon = 16$.

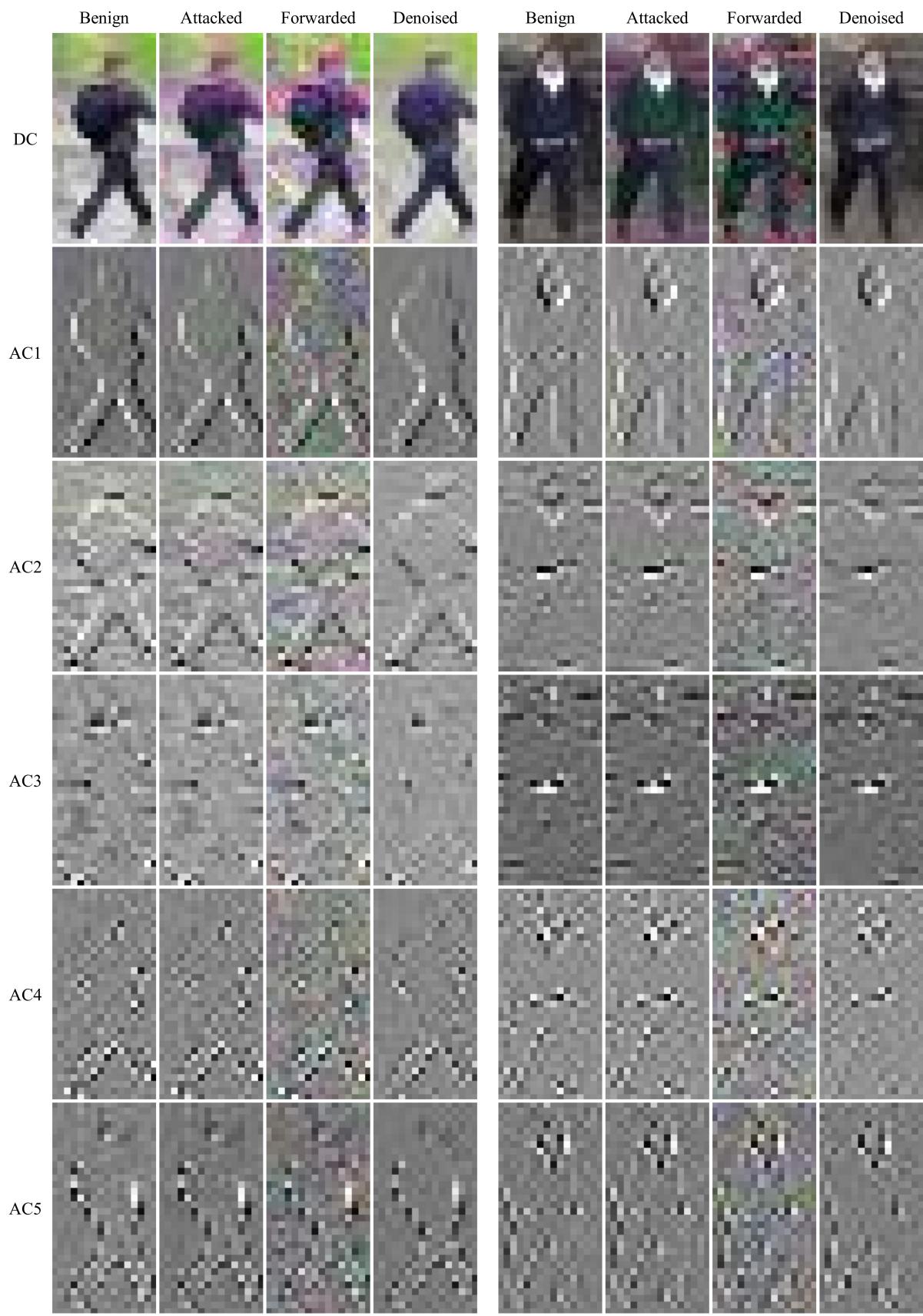


Figure 4: Illustration of the secondary diffusion process on **DukeMTMC dataset** attacked by Deep Mis-Ranking with $\epsilon = 16$.



Figure 5: Illustration of the top 10 retrieval results for images purified on Market1501 dataset attacked by Metric-FGSM with $\epsilon = 4$. **Red boxes** denote false results. **Green boxes** denote correct results.



Figure 6: Illustration of the top 10 retrieval results for images purified on Market1501 dataset attacked by Metric-FGSM with $\epsilon = 16$. **Red boxes** denote false results. **Green boxes** denote correct results.



Figure 7: Illustration of the top 10 retrieval results for images purified on Market1501 dataset attacked by Deep Mis-Ranking with $\epsilon = 4$. Red boxes denote false results. Green boxes denote correct results.



Figure 8: Illustration of the top 10 retrieval results for images purified on Market1501 dataset attacked by Deep Mis-Ranking with $\epsilon = 16$. Red boxes denote false results. Green boxes denote correct results.



Figure 9: Illustration of the top 10 retrieval results for images purified on Market1501 dataset attacked by MetaAttack with $\epsilon = 4$. **Red boxes** denote false results. **Green boxes** denote correct results.



Figure 10: Illustration of the top 10 retrieval results for images purified on Market1501 dataset attacked by MetaAttack with $\epsilon = 16$. **Red boxes** denote false results. **Green boxes** denote correct results.