

6. Appendix

In this appendix, we first provide more results of integrating BSR with DIM, TIM, Admix and PAM, the adversaries are crafted on three other models *i.e.* Inc-v4, IncRes-v2 and Res-101. Then we provide more comparison on the heatmap of various images using different methods.

6.1. More Evaluations on Combined Input Transformation

We provide the attack results of adversarial examples generated by integrating BSR with DIM, TIM, *Admix*, and PAM on the other three models, namely Inc-v4, IncRes-v2 and Res-101. As shown in Tab. 5-7, BSR can remarkably boost the transferability of various input transformation based attacks when generating the adversarial examples on the three models. Zhang et al. [55] have demonstrated that the combination of PAM-TI-DIM achieves state-of-the-art transferability. However, its performance on adversarially trained models remains relatively weak. Our proposed BSR-SI-DI-TIM consistently outperforms PAM-TI-DIM with a margin of at least 6.4%. Such excellent performance improvement further validates the high effectiveness of the proposed BSR for boosting the adversarial transferability.

6.2. Visualization on Attention Heatmaps

To further support our motivation, we visualize more attention heatmaps of sampled benign images on both source model Inc-v4 and target model Inc-v3 in Fig. 5. We can observe that the attention heatmaps of adversarial examples crafted by BSR on source model are more similar with heatmap on target model than the results of other attacks. These results validate our motivation that our proposed method can eliminate the variance of attention heatmaps among various models, which is of benefit to generate more transferable adversarial examples.

Attack	Inc-v3	IncRes-v2	Res-101	Inc-v3 _{ens3}	Inc-v3 _{ens4}	IncRes-v2 _{ens}
BSR-DIM	96.8 ^{↑24.8}	94.1 ^{↑30.3}	87.8 ^{↑30.4}	59.7 ^{↑37.1}	53.6 ^{↑32.5}	37.6 ^{↑11.7}
BSR-TIM	93.8 ^{↑34.7}	91.7 ^{↑42.7}	84.6 ^{↑42.7}	69.4 ^{↑42.6}	63.4 ^{↑40.5}	51.1 ^{↑34.5}
BSR-SIM	99.0 ^{↑18.6}	97.9 ^{↑24.5}	95.2 ^{↑25.8}	86.7 ^{↑38.1}	85.2 ^{↑40.0}	70.2 ^{↑40.6}
BSR-Admix	99.1 ^{↑10.1}	98.5 ^{↑13.2}	97.0 ^{↑18.0}	88.0 ^{↑32.5}	86.0 ^{↑34.3}	74.0 ^{↑41.7}
BSR-PAM	99.2 ^{↑12.5}	98.2 ^{↑16.6}	96.2 ^{↑20.3}	87.2 ^{↑31.8}	81.4 ^{↑30.9}	62.8 ^{↑29.6}
Admix-TI-DIM	91.0	88.6	83.2	76.0	74.7	64.3
PAM-TI-DIM	91.5	88.5	83.5	77.1	73.2	62.1
BSR-TI-DIM	91.2	88.5	83.5	77.1	73.2	62.1
BSR-SI-TI-DIM	98.1	96.2	94.7	90.1	89.0	80.0

Table 5. Attack success rates (%) on seven models under single model setting with various input transformations combined with BSR. The adversaries are crafted on Inc-v4. [↑] indicates the increase of attack success rate when combined with BSR.

Attack	Inc-v3	Inc-v4	Res-101	Inc-v3 _{ens3}	Inc-v3 _{ens4}	IncRes-v2 _{ens}
BSR-DIM	95.1 ^{↑24.8}	95.0 ^{↑30.3}	91.5 ^{↑33.5}	71.7 ^{↑41.3}	65.4 ^{↑41.9}	54.8 ^{↑37.9}
BSR-TIM	93.4 ^{↑31.2}	93.2 ^{↑37.6}	89.1 ^{↑38.8}	80.3 ^{↑47.9}	75.9 ^{↑48.4}	70.0 ^{↑47.4}
BSR-SIM	98.5 ^{↑12.6}	98.4 ^{↑18.4}	97.4 ^{↑21.3}	92.5 ^{↑36.3}	89.5 ^{↑40.4}	81.9 ^{↑39.4}
BSR-Admix	99.1 ^{↑8.3}	99.3 ^{↑13.0}	98.1 ^{↑15.9}	93.7 ^{↑30.1}	92.1 ^{↑35.5}	87.3 ^{↑37.9}
BSR-PAM	99.3 ^{↑10.7}	99.2 ^{↑12.9}	98.5 ^{↑16.9}	94.2 ^{↑28.2}	91.2 ^{↑32.9}	83.3 ^{↑32.3}
Admix-TI-DIM	90.9	89.5	86.5	81.7	77.7	76.4
PAM-TI-DIM	92.2	90.5	86.9	84.2	80.9	78.4
BSR-TI-DIM	94.5	93.6	90.6	81.0	77.5	72.9
BSR-SI-TI-DIM	98.6	98.1	96.3	95.4	94.0	91.7

Table 6. Attack success rates (%) on seven models under single model setting with various input transformations combined with BSR. The adversaries are crafted on IncRes-v2. [↑] indicates the increase of attack success rate when combined with BSR.

Attack	Inc-v3	Inc-v4	IncRes-v2	Inc-v3 _{ens3}	Inc-v3 _{ens4}	IncRes-v2 _{ens}
BSR-DIM	97.7 ^{↑21.7}	96.6 ^{↑28.2}	97.1 ^{↑26.8}	82.7 ^{↑48.0}	75.8 ^{↑44.0}	59.7 ^{↑40.1}
BSR-TIM	96.7 ^{↑36.8}	95.6 ^{↑43.4}	94.9 ^{↑43.0}	86.6 ^{↑52.2}	83.4 ^{↑52.2}	75.3 ^{↑51.6}
BSR-SIM	98.8 ^{↑20.3}	98.3 ^{↑25.8}	98.1 ^{↑26.9}	91.5 ^{↑46.7}	89.1 ^{↑49.0}	75.8 ^{↑48.5}
BSR-Admix	99.6 ^{↑15.1}	99.2 ^{↑19.0}	99.4 ^{↑18.7}	94.2 ^{↑42.6}	92.5 ^{↑47.8}	82.4 ^{↑52.5}
BSR-PAM	98.7 ^{↑21.3}	97.2 ^{↑23.3}	97.7 ^{↑22.0}	90.0 ^{↑38.8}	86.8 ^{↑40.3}	70.3 ^{↑38.1}
Admix-TI-DIM	89.5	85.6	87.5	80.4	75.2	67.5
PAM-TI-DIM	85.8	83.2	84.9	81.1	76.1	66.9
BSR-TI-DIM	97.3	95.1	95.6	88.3	85.4	76.1
BSR-SI-TI-DIM	97.4	97.0	96.6	94.8	93.2	89.1

Table 7. Attack success rates (%) on seven models under single model setting with various input transformations combined with BSR. The adversaries are crafted on Res-101. [↑] indicates the increase of attack success rate when combined with BSR.

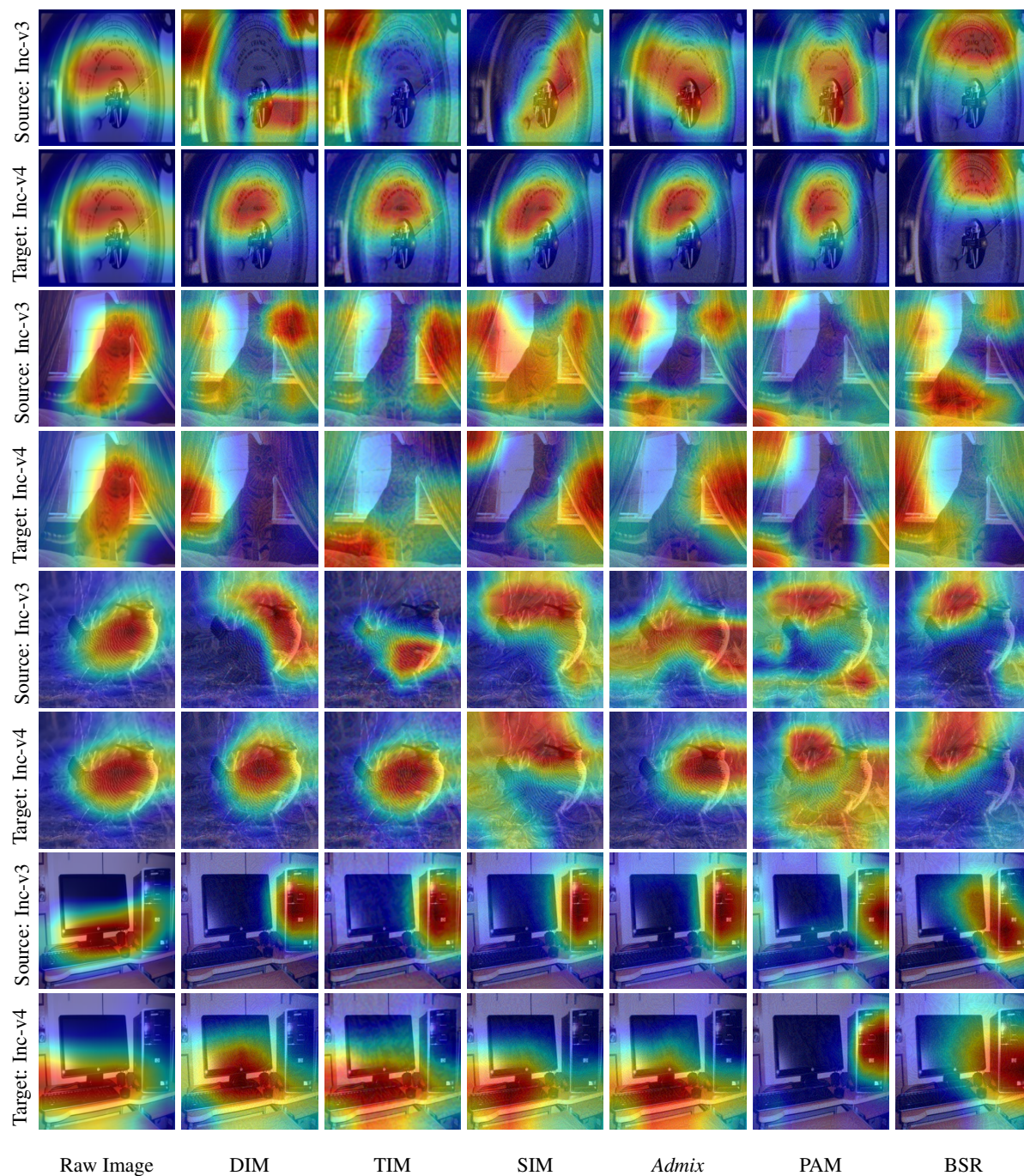


Figure 5. The visualization of attention heatmaps [28] of adversaries crafted by various input transformation methods on source model Inc-v3 [36], and target model Inc-v4 [35]