

# Supplementary Material for A<sup>2</sup>RNet: Adversarial Attack Resilient Network for Robust Infrared and Visible Image Fusion

Anonymous submission

## Abstract

Here, we provide supplementary materials for the main text. First, we detail the specific process of pseudo-label generation. Then, additional experimental details and results are provided to facilitate a better understanding. All in all, thank you for taking the time to read this document.

## Pseudo-label Generation

As shown in Fig. 1, we use a common CNN to generate pseudo-labels.  $l_1$ , SSIM, and gradient loss are used as the loss functions during training. We utilize these pseudo-labels as the necessary references for adversarial attacks and training, which introduces a new paradigm for research on adversarial robustness in the IVIF task. Moreover, these pseudo-labels are more “moderate” and less likely to lead to overfitting.

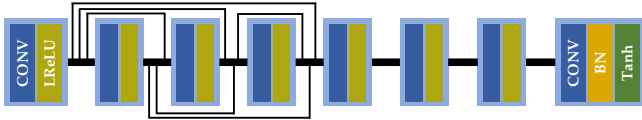


Figure 1: Pipeline of the pseudo-label generation.

## More Details of Experiments

### Quantitative Metrics

In the quantitative comparison, we choose Entropy (EN), Standard Deviation (SD), Peak Signal-toNoise Ratio (PSNR), Correlation Coefficient (CC) (Shah, Merchant, and Desai 2013) and the Sum of the Correlations of Differences (SCD) (Aslantas and Bendes 2015). Specifically, EN represents the measure of information contained in fused images. SD quantifies the dispersion of pixel values around the mean in fusion results. PSNR reflects the level of distortion by comparing the peak signal power to the noise power. CC represents the degree of linear correlation between source and fused images. SCD denotes the characterization of differences between source and fused images. Clearly, EN and SD are metrics calculated based on the fused results themselves, while PSNR, CC and SCD evaluate the relationship between source images and fused images. For all metrics, higher values indicate better image quality.

## Configurations of Downstream Tasks

We select the M<sup>3</sup>FD (Liu et al. 2022) and MFNet (Ha et al. 2017) datasets, which provide annotated detection and segmentation labels. Note that we use the fusion results as inputs to retrain the downstream task models. For the detection task, we employ YOLOv5 (Redmon et al. 2016) as the detector. The optimizer, learning rate, epochs, and batch size are set to SGD optimizer, 1e-2, 300, and 8, respectively. 4200 image pairs from the M<sup>3</sup>FD dataset are divided into training, validation, and test sets in an 8:1:1 ratio. In the segmentation task, DeepLabV3+ (Chen et al. 2018) is introduced with 300 epochs and a batch size of 8, keeping the other parameters consistent with the original model. A total of 1,083 fusion results from the MFNet dataset are used as the training set, while the remaining 361 are used as the test set.

## Fusion Results under Different Perturbations

For a more comprehensive comparison, we have included an additional set of experiments with 8/255 attacks in the supplementary materials. Note that adversarial examples during the adversarial training are produced using the same perturbations. As shown in Fig. 2, two sets from the M<sup>3</sup>FD and MFNet datasets are presented. Clearly, our method does not exhibit undesirable artifacts when subjected to stronger perturbations. However, other methods exhibited more pronounced noise or distortions. It can demonstrate that our method maintains robustness even under stronger perturbations.

## Quantitative Comparisons in Downstream tasks

In this supplementary material, we provide detailed quantitative results for each category in the detection and segmentation tasks, where the subscript indicates the degree of change under adversarial conditions. As shown in Table. 1 and 2, our method achieves the best performance across all categories, demonstrating that fused images generated by A<sup>2</sup>RNet remain robust in downstream task performance. Combined with the qualitative results in the main text, the downstream task performance of our method remains a leading position under attacks.

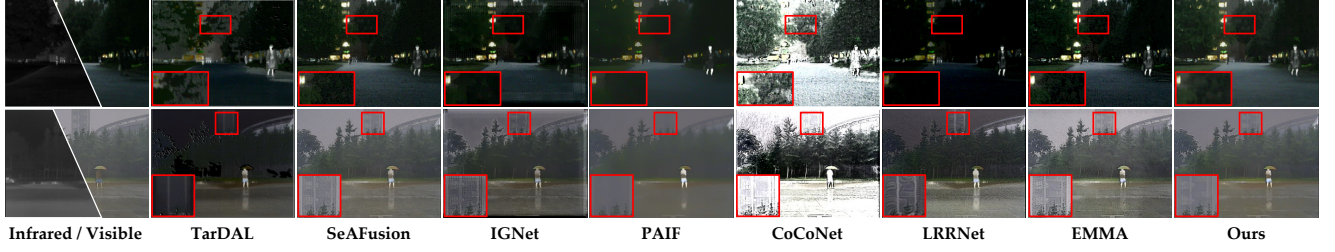


Figure 2: Fusion comparisons with SOTA methods in MFNet and M<sup>3</sup>FD datasets. We apply PGD to clean samples and add stronger perturbations with  $\epsilon = 8/255$  to generate AEs.

Category	AP@.5							
	TarDAL	SeAFusion	IGNet	PAIF	CoCoNet	LRRNet	EMMA	Ours
People	0.530 <sub>↓0.171</sub>	<b>0.704</b> <sub>↓0.041</sub>	0.559 <sub>↓0.134</sub>	0.703 <sub>↓0.034</sub>	0.209 <sub>↓0.133</sub>	0.680 <sub>↓0.038</sub>	0.676 <sub>↓0.032</sub>	<b>0.764</b> <sub>↓0.011</sub>
Car	0.597 <sub>↓0.187</sub>	<b>0.860</b> <sub>↓0.040</sub>	0.825 <sub>↓0.042</sub>	0.859 <sub>↓0.029</sub>	0.680 <sub>↓0.063</sub>	0.841 <sub>↓0.024</sub>	0.801 <sub>↓0.039</sub>	<b>0.893</b> <sub>↓0.001</sub>
Bus	0.491 <sub>↓0.284</sub>	0.776 <sub>↓0.073</sub>	0.582 <sub>↓0.200</sub>	0.790 <sub>↓0.049</sub>	0.456 <sub>↓0.000</sub>	<b>0.779</b> <sub>↓0.045</sub>	0.751 <sub>↓0.052</sub>	<b>0.866</b> <sub>↓0.012</sub>
Motorcycle	0.348 <sub>↓0.093</sub>	0.584 <sub>↓0.018</sub>	0.411 <sub>↓0.158</sub>	0.554 <sub>↓0.015</sub>	0.203 <sub>↓0.025</sub>	<b>0.620</b> <sub>↑0.014</sub>	0.600 <sub>↓0.007</sub>	<b>0.647</b> <sub>↓0.021</sub>
Truck	0.518 <sub>↓0.123</sub>	0.771 <sub>↓0.050</sub>	0.529 <sub>↓0.205</sub>	0.751 <sub>↓0.080</sub>	0.238 <sub>↓0.200</sub>	<b>0.785</b> <sub>↓0.028</sub>	0.711 <sub>↓0.054</sub>	<b>0.809</b> <sub>↓0.021</sub>
Lamp	0.289 <sub>↓0.106</sub>	<b>0.699</b> <sub>↓0.027</sub>	0.436 <sub>↓0.186</sub>	0.631 <sub>↓0.069</sub>	0.077 <sub>↓0.138</sub>	0.518 <sub>↓0.073</sub>	0.635 <sub>↓0.046</sub>	<b>0.712</b> <sub>↓0.038</sub>
<b>mAP@.5</b>	0.462 <sub>↓0.161</sub>	<b>0.732</b> <sub>↓0.042</sub>	0.557 <sub>↓0.154</sub>	0.715 <sub>↓0.046</sub>	0.311 <sub>↓0.093</sub>	0.704 <sub>↓0.032</sub>	0.696 <sub>↓0.038</sub>	<b>0.781</b> <sub>↓0.024</sub>

Table 1: Quantitative results of detection. **Red** and **blue** denote the optimal and suboptimal results, respectively. The subscripts indicate the change compared adversarial conditions with the clean.

Category	IoU							
	TarDAL	SeAFusion	IGNet	PAIF	CoCoNet	LRRNet	EMMA	Ours
Background	0.963 <sub>↓0.005</sub>	0.971 <sub>↓0.004</sub>	0.967 <sub>↓0.014</sub>	0.974 <sub>↓0.004</sub>	0.962 <sub>↓0.014</sub>	<b>0.978</b> <sub>↓0.000</sub>	0.973 <sub>↓0.010</sub>	<b>0.979</b> <sub>↓0.005</sub>
Car	0.685 <sub>↓0.058</sub>	0.802 <sub>↓0.024</sub>	0.749 <sub>↓0.127</sub>	<b>0.831</b> <sub>↓0.056</sub>	0.704 <sub>↓0.041</sub>	0.792 <sub>↓0.057</sub>	0.808 <sub>↓0.086</sub>	<b>0.845</b> <sub>↓0.055</sub>
Person	<b>0.663</b> <sub>↓0.040</sub>	0.634 <sub>↓0.064</sub>	0.469 <sub>↓0.213</sub>	0.643 <sub>↓0.078</sub>	0.383 <sub>↓0.113</sub>	0.615 <sub>↓0.099</sub>	0.625 <sub>↓0.106</sub>	<b>0.686</b> <sub>↓0.040</sub>
Bike	0.340 <sub>↓0.011</sub>	0.570 <sub>↓0.083</sub>	0.511 <sub>↓0.169</sub>	0.547 <sub>↓0.132</sub>	0.434 <sub>↓0.096</sub>	<b>0.597</b> <sub>↓0.009</sub>	0.591 <sub>↓0.094</sub>	<b>0.645</b> <sub>↓0.053</sub>
Curve	0.136 <sub>↓0.145</sub>	0.449 <sub>↓0.089</sub>	0.281 <sub>↓0.249</sub>	0.306 <sub>↓0.201</sub>	0.166 <sub>↓0.065</sub>	0.427 <sub>↓0.012</sub>	<b>0.488</b> <sub>↓0.092</sub>	<b>0.508</b> <sub>↓0.091</sub>
Car Stop	0.252 <sub>↓0.094</sub>	0.600 <sub>↓0.042</sub>	0.462 <sub>↓0.198</sub>	0.492 <sub>↓0.157</sub>	0.294 <sub>↓0.092</sub>	0.594 <sub>↓0.009</sub>	<b>0.676</b> <sub>↓0.021</sub>	<b>0.685</b> <sub>↓0.031</sub>
Guardrail	0.153 <sub>↓0.656</sub>	0.626 <sub>↓0.066</sub>	0.594 <sub>↓0.137</sub>	0.602 <sub>↓0.129</sub>	0.295 <sub>↓0.049</sub>	<b>0.663</b> <sub>↓0.015</sub>	0.645 <sub>↓0.119</sub>	<b>0.682</b> <sub>↓0.087</sub>
Color Cone	0.427 <sub>↓0.084</sub>	0.509 <sub>↓0.053</sub>	0.418 <sub>↓0.210</sub>	0.506 <sub>↓0.100</sub>	0.257 <sub>↓0.125</sub>	0.520 <sub>↓0.007</sub>	<b>0.581</b> <sub>↓0.034</sub>	<b>0.604</b> <sub>↓0.041</sub>
Bump	0.256 <sub>↓0.116</sub>	0.405 <sub>↓0.270</sub>	0.176 <sub>↓0.335</sub>	0.302 <sub>↓0.342</sub>	0.192 <sub>↓0.135</sub>	0.422 <sub>↓0.044</sub>	<b>0.453</b> <sub>↓0.210</sub>	<b>0.461</b> <sub>↓0.205</sub>
<b>mIoU</b>	0.415 <sub>↓0.069</sub>	0.618 <sub>↓0.078</sub>	0.514 <sub>↓0.184</sub>	0.578 <sub>↓0.134</sub>	0.409 <sub>↓0.080</sub>	0.623 <sub>↓0.027</sub>	<b>0.649</b> <sub>↓0.086</sub>	<b>0.677</b> <sub>↓0.068</sub>

Table 2: Quantitative results of segmentation. **Red** and **blue** denote the optimal and suboptimal results, respectively. The subscripts indicate the change compared adversarial conditions with the clean.

## References

- Aslantas, V.; and Bendes, E. 2015. A new image quality metric for image fusion: The sum of the correlations of differences. *Aeu-international Journal of electronics and communications*, 69(12): 1890–1896.
- Chen, L.-C.; Zhu, Y.; Papandreou, G.; Schroff, F.; and Adam, H. 2018. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, 801–818.
- Ha, Q.; Watanabe, K.; Karasawa, T.; Ushiku, Y.; and Harada, T. 2017. MFNet: Towards real-time semantic segmentation for autonomous vehicles with multi-spectral scenes. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 5108–5115. IEEE.
- Liu, J.; Fan, X.; Huang, Z.; Wu, G.; Liu, R.; Zhong, W.; and Luo, Z. 2022. Target-aware dual adversarial learning and a multi-scenario multi-modality benchmark to fuse infrared and visible for object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5802–5811.
- Redmon, J.; Divvala, S.; Girshick, R.; and Farhadi, A. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 779–788.
- Shah, P.; Merchant, S. N.; and Desai, U. B. 2013. Multi-focus and multispectral image fusion based on pixel significance using multiresolution decomposition. *Signal, Image and Video Processing*, 7: 95–109.