

Adversarial Robustness in RGB-Skeleton Action Recognition: Leveraging Attention Modality Reweigher

Chao Liu¹ Xin Liu^{2*} Zitong Yu^{3*} Yonghong Hou¹ Huanjing Yue¹ Jingyu Yang¹
¹Tianjin University ²Lappeenranta-Lahti University of Technology ³Great Bay University
xin.liu@lut.fi, yuzitong@gbu.edu.cn

Abstract

Deep neural networks (DNNs) have been applied in many computer vision tasks and achieved state-of-the-art (SOTA) performance. However, misclassification will occur when DNNs predict adversarial examples which are created by adding human-imperceptible adversarial noise to natural examples. This limits the application of DNN in security-critical fields. In order to enhance the robustness of models, previous research has primarily focused on the unimodal domain, such as image recognition and video understanding. Although multi-modal learning has achieved advanced performance in various tasks, such as action recognition, research on the robustness of RGB-skeleton action recognition models is scarce. In this paper, we systematically investigate how to improve the robustness of RGB-skeleton action recognition models. We initially conducted empirical analysis on the robustness of different modalities and observed that the skeleton modality is more robust than the RGB modality. Motivated by this observation, we propose the **Attention-based Modality Reweigher (AMR)**, which utilizes an attention layer to re-weight the two modalities, enabling the model to learn more robust features. Our AMR is plug-and-play, allowing easy integration with multimodal models. To demonstrate the effectiveness of AMR, we conducted extensive experiments on various datasets. For example, compared to the SOTA methods, AMR exhibits a 43.77% improvement against PGD20 attacks on the NTU-RGB+D 60 dataset. Furthermore, it effectively balances the differences in robustness between different modalities.

1. Introduction

Deep Neural Networks (DNNs) have achieved SOTA performance in various vision tasks, like segmentation [31], detection [32], and super-resolution [12]. However, DNNs are vulnerable to imperceptible adversarial perturbations [13]. A finely crafted adversarial perturba-

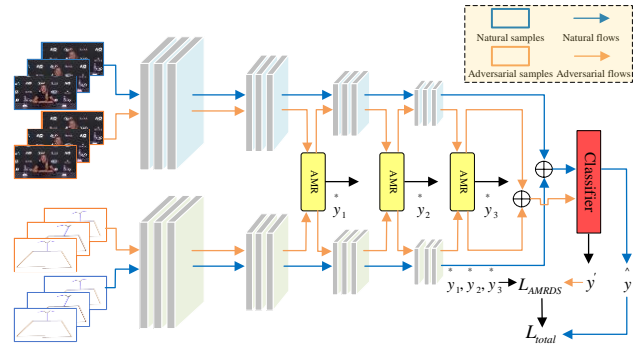


Figure 1. The framework of the proposed method AMR. During the training, we feed-forward both natural and adversarial samples from two modalities into the multimodal network in parallel. Features extracted from adversarial samples are input into AMR, and their weighted counterparts enter respective subsequent layers of the network. These features are simultaneously used for classification to generate auxiliary predictions denoted as \hat{y}^* , contributing to the loss function. The natural samples go through regular forward propagation to yield prediction results.

tion can easily fool the neural network. Adversarial attacks covered many deep learning tasks including video action recognition [41], person re-identification [38], and natural language processing [21]. Given the DNN's widespread applications, especially in safety-critical domains like autonomous driving [5, 10, 25, 49] and medical diagnostics [2, 18, 26, 23, 24], enhancing neural network robustness is crucial.

To address DNN vulnerability, efforts are categorized into adversarial attacks and defenses. Various methods are proposed to generate diverse adversarial samples [35, 3, 28, 30, 34]. Many works aim to defend against these attacks, with adversarial training showing the most reliable robustness. Adversarial training inputs adversarial examples into each training step. Natural examples are augmented with worst-case perturbations within a small l_p -norm ball, smoothing the loss landscape and improving classification performance in boundary regions.

As DNNs advances, multimodal technology has become

*Xin Liu and Zitong Yu are corresponding authors

a widely studied field with successful applications in vision-language navigation [39], image-text matching[11], and medical image diagnosis [16, 29]. Multimodal technology offers advantages over unimodal technology. Previous research on model robustness has focused mainly on unimodal domains, like video comprehension [17], with limited attention to multimodal models. Thus, research on enhancing the robustness of multimodal models is urgently needed. This paper systematically analyzes the robustness of action recognition models based on RGB-skeleton data and proposes methods to enhance their robustness.

To understand the robustness of various modalities, we analyzed the RGB and skeleton modalities under different attacks and intensities. Figure 2 shows our experiments on the NTU-RGB+D [33] and iMiGUE [22] datasets, evaluating their response to three types of attacks. As attack intensity increases, the skeletal modality’s accuracy declines slower and smoother than the RGB and multimodal modalities, which drop sharply. This indicates that slight perturbations cause significant errors in the RGB modality, highlighting the skeletal modality’s higher robustness. Both datasets show the same trend. This aligns with findings by Yan *et al.* [44], who noted that bone and joint trajectories are robust to lighting and scene variations. Skeleton data lacks intricate information like lighting and background details, making it easier to capture human actions. Multimodal fusion results in the lowest accuracy, as found in [36], suggesting it does not enhance and may even diminish robustness.

Based on this observation, we propose an Attention-based Modality Reweighter (AMR) to enhance the robustness of RGB-skeleton action recognition models, shown in Figure 1. Our backbone consists of two branch networks, I3D [4] and HCN [20], extracting features from RGB videos and skeleton, respectively. Although I3D and HCN are not SOTA, their simple design suits our approach. Our goal is to enhance model robustness, not just achieve the highest action recognition rate. Using a simple model eliminates unnecessary factors. In the final blocks, we incorporate AMR to connect a two-stream network. AMR has two attention layers assigning weights to features from both modalities, helping the model learn robust features, especially against adversarial attacks. After reweighted features pass through global pooling layers (GAP if necessary), two classification vectors are derived, concatenated, and fed into a fully connected layer (FC) for auxiliary classification results. AMR generates three outputs: reweighted features of both modalities and auxiliary classification predictions. The adjusted features continue forward propagation, while auxiliary classification results are used in the loss function. Further details are in Section 3.

In summary, our contributions are:

- We conducted empirical analysis comparing the ro-

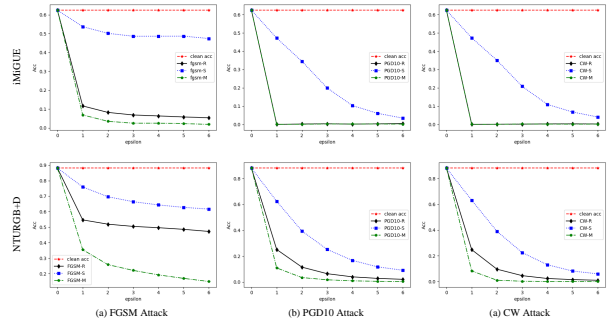


Figure 2. Accuracy against three types of adversarial attacks based on the iMiGUE dataset [22] and the NTU-RGB+D dataset [33], namely adversarial robustness. The x-axis indicates the attack strength $\epsilon (\times \frac{2}{255})$. FGSM-S, and FGSM-M respectively denote the FGSM attack on the RGB modality, skeleton modality, and multiple modalities simultaneously, and so on. ‘Clean acc’ represents the accuracy after multimodal fusion on clean data. The findings manifest that, as the intensity of attacks increases, the robustness of the skeleton modality exhibits a comparatively gradual and smooth decline in contrast to the sharp decrease observed in both the RGB modality and multimodal fusion. Therefore, we confirm the assertion that the skeleton modality shows higher robustness than the RGB modality.

bustness of different modalities, showing the superior robustness of the skeleton modality over RGB.

- We propose a novel module, Attention-based Modality Reweighter (AMR), with a new loss function. It supports both adversarial and standard training, reweighting modalities to enhance robustness.
- We conducted comprehensive experiments on popular benchmark datasets, evaluating performance against classical adversarial attacks. The results show that our method enhances the robustness of multimodal action recognition models, addresses robustness discrepancies across modalities, and achieves new SOTA performance without additional data.

2. Related Work

2.1. RGB-Skeleton Action Recognition

Conventional approaches to action recognition primarily rely on RGB video as the main input modality. Yan *et al.* [44] first introduced deep learning into skeleton-based action recognition. RGB videos provide comprehensive visual information about actions, including human attributes, attire, colors, spatial relationships, and contextual background details. Skeletal data focuses on pose information, offering insights into human posture and alignment through skeletal key-points. RGB and skeletal modalities provide synchronized and complementary information. Their inte-

gration enhances action recognition by offering more comprehensive and accurate information. Numerous studies in RGB-skeleton action recognition have emerged [8, 45, 37]. While clean data fusion aids action recognition, the impact on robustness when RGB-skeleton modalities face attacks is unclear. This paper uses RGB-skeleton action recognition to investigate the robustness of multimodal learning against adversarial attacks.

2.2. Adversarial Attack

Since [13] finds the adversarial perturbations in neural networks, generating adversarial images to attack deep networks have attracted great interests. Lots of attack methods have been proposed to craft adversarial examples which can fool deep neural networks. The most famous attack is Fast Gradient Sign Method (FGSM) [13], which only needs one step to create adversarial examples. Further, Madry *et al.* [28] proposed the Project Gradient Descent (PGD) attack, which is the most used adversarial attack in the adversarial research field. Boundary-based attacks such as deepfool [30] and CW [3] also make the model more challenging. Duan *et al.* [9] proposed a multi-sample generation model for black-box model attacks, called MsGM. AutoAttack (AA) is an ensemble attack, which consists of three white-box attacks (APGD-CE [7], APGD-DLR [7], and FAB [6]) and one black-box attack (Square attack [1]) for a total of four attack methods. Xu *et al.* [43] investigate the undiscovered adversarial attacks in the unsupervised clustering domain, a field of equal significance, and propose the concept of the adversarial set and adversarial set attack for clustering.

2.3. Adversarial Defense

With the development of adversarial attacks, there is also a lot of work to improve the adversarial robustness of deep neural networks [14]. Adversarial training [28] is the most effective way to defend against adversarial examples and serves as our fundamental training approach. TRADES [47] analysis the trade-off between clean accuracy and adversarial accuracy, and uses the Kullback-Leibler divergence [19] to balance them. Misclassification Aware adversarial Training (MART) [40] emphasizes the misclassified examples. Adversarial logit pairing (ALP) [15] and its variant Adaptive Adversarial Logits Pairing (AALP) [42] define pairing loss that pulls adversarial logit and natural logit together. Adversarial Neural Pruning (ANP) [27] uses Bayesian methods to prune vulnerable features. Zhang *et al.* [48] uses sample reweighting techniques to improve the adversarial robustness against PGD attack.

All the above works are in the fields of image and video processing. Tian *et al.* [36] is most relevant to our research, investigating audio-visual model robustness under multimodal attacks. Their experiments show that integrat-

ing audio and visual components under multimodal attacks does not enhance but decreases model robustness. This indicates that audio-visual models are vulnerable to multimodal adversarial attacks and that integration may decrease robustness. To mitigate these attacks, they propose an multimodal defense approach based on an audio-visual dissimilarity constraint (MinSim) and External Feature Memory banks (ExFMem), representing the SOTA in multimodal defense. However, ExFMem’s storage of external features requires additional disk space and does not consider variations in robustness across different modalities. Our proposed method addresses these challenges.

3. Methodology

In this section, we first introduce the Standard Adversarial Training (SAT) [28]. Then, we will describe the Attention-based Modality Reweigher (AMR) in detail.

3.1. Standard Adversarial Training

Given a standard training data set $D = \{(x_i, y_i)\}_{i=1}^n$ with n examples and C classes, where $x_i \in \mathbb{R}^d$ is the natural example and $y_i \in \{0, 1, \dots, C - 1\}$ is corresponding ground-truth label. Standard training (ST) achieves good classification performance by minimizing the empirical risk on the training dataset, formulated as Equation (1).

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(f_{\theta}(x_i), y_i), \quad (1)$$

where l is the cross-entropy loss widely used in classification tasks, $f: \mathbb{R}^d \rightarrow \mathbb{R}^C$ is the neural network parameterized by θ and Θ is the parameter space of θ . The neural network trained in this manner, however, does not exhibit satisfactory prediction accuracy when faced with adversarial examples. To solve this problem, Madry *et al.* [28] used adversarial data to train the neural network, which can be formulated as a min-max optimization problem:

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n I(f_{\theta}(x'_i), y_i), \quad (2)$$

where

$$x'_i = \arg \max_{x' \in \mathcal{B}_{\varepsilon}[x_i]} l(f_{\theta}(x'), y_i). \quad (3)$$

$\mathcal{B}_{\varepsilon}[x_i]$ is the sampling space of the adversarial example x'_i , which is bounded in the L_p -norm neighbor space of the natural example x_i , *i.e.*, $\mathcal{B}_{\varepsilon}[x_i] = \{x' \in \mathbb{R}^d \mid \|x' - x_i\|_p \leq \varepsilon\}$, and is the closed ball of radius $\varepsilon > 0$ centered at x_i . $\delta = x'_i - x_i$ is the adversarial perturbation. In this paper, only $p = \infty$ is considered. It implies the optimization of adversarial robust network, with one step maximizing loss to find adversarial data and one step minimizing loss on the adversarial data *w.r.t.* the network parameters θ . To generate

adversarial data, SAT [28] uses PGD to approximately solve the inner maximization of Equation (3). Namely, given a starting point $x'_0 \in \mathbb{R}^d$ and step size $\alpha > 0$, PGD works as follows:

$$x'_{t+1} = \prod_{\mathcal{B}_\epsilon[x'_0]} (x'_t + \alpha \text{sign}(\nabla_{x'_t} l(f_\theta(x'_t), y))), \forall t \geq 0 \quad (4)$$

until a certain stopping criterion is satisfied. For example, the criterion can be a fixed number of iterations K , namely the PGD-K algorithm [28]. In Equation (4), l is the loss function in Equation (3); x'_0 refers to natural data or natural data corrupted by a small Gaussian or uniform random noise; y is the corresponding label for natural data; x'_t is adversarial data at step t ; and $\prod_{\mathcal{B}_\epsilon[x'_0]} (\cdot)$ is the projection function that projects the adversarial data back into the ϵ -ball centered at x'_0 if necessary.

3.2. Multimodal Adversarial Training

Let $x_{R,i}$ be an input RGB video frame, $x_{s,i}$ be an input skeleton data, and y_i be the corresponding groundtruth label for the multisensory input: $\{x_{R,i}, x_{s,i}\}$. Since there are multiple inputs, we can divide our multimodal attack into two categories: unimodality attacks that only generate RGB adversarial example $x'_{R,i}$ or skeleton adversarial example $x'_{s,i}$, and RGB-skeleton multimodal attacks that generate both RGB and skeleton adversarial examples: $\{x'_{R,i}, x'_{s,i}\}$. We can extend the min-max formulation Equation (2) and Equation (3) to multimodal data [36], formulated as follows:

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(f_\theta(x'_{R,i}, x'_{s,i}), y_i), \quad (5)$$

where

$$\{x'_{R,i}, x'_{s,i}\} = \arg \max_{\substack{x'_R \in \mathcal{B}_{\delta_R}[x_{R,i}], \\ x'_s \in \mathcal{B}_{\delta_s}[x_{s,i}]}} l(f(x'_R, x'_s), y_i), \quad (6)$$

where $\mathcal{B}_{\delta_R}[x_{R,i}] = \{x'_R \in \mathbb{R}^{d_R} \mid \|x'_R - x_{R,i}\|_p \leq \delta_R\}$, $\mathcal{B}_{\delta_s}[x_{s,i}] = \{x'_s \in \mathbb{R}^{d_s} \mid \|x'_s - x_{s,i}\|_p \leq \delta_s\}$. $\delta_R = x'_{R,i} - x_{R,i}$ and $\delta_s = x'_{s,i} - x_{s,i}$ are video adversarial perturbation and skeleton adversarial perturbation, respectively. The symbols used in the formulas presented in this section retain the same meanings as those discussed in Section 3.1. With the adversarial objective, the attacker will maximize the loss function by seeking small perturbations within allowed budgets, and try to push the trained model to make incorrect predictions. For unimodality attacks, either δ_R or δ_s is 0. Likewise, Equation (4) can be extended to generate

multimodal adversarial samples:

$$\begin{cases} x'_{R,t+1} = \prod_{\mathcal{B}_{\epsilon_R}[x_{R,0}']} (x'_{R,t} + \alpha \text{sign}(\nabla_{x'_{R,t}} l(f_\theta(x'_{R,t}, x'_{s,t}), y))), \\ x'_{s,t+1} = \prod_{\mathcal{B}_{\epsilon_s}[x_{s,0}']} (x'_{s,t} + \alpha \text{sign}(\nabla_{x'_{s,t}} l(f_\theta(x'_{R,t}, x'_{s,t}), y))), \end{cases} \quad (7)$$

where, $\forall t \geq 0$.

3.3. Attention-based Modality Reweigher

Reviewing Figure 2, numerous experimental results show that the robustness of the skeleton and RGB modalities differs: the skeleton modality is significantly more robust. As attack intensity increases, the skeleton modality's accuracy declines gradually and smoothly, while the RGB and multimodal approaches drop sharply. These results align with the conclusions in [44, 36, 46]. Motivated by this finding, we aim to enable the multimodal model to assign appropriate weights to the features of the two modalities, making the model learn more robust features against adversarial attacks. Specifically, more robust features should receive greater weight, while less robust features should get smaller weights.

Our proposed AMR is illustrated in Figure 3. Specifically, given a mini-batch of multi-modal natural sample $\{x_R, x_s\}$ (B in total), we can obtain corresponding adversarial $\{x'_R, x'_s\}$ samples by using algorithms such as PGD, as mentioned in Equation (6). Then, we extract latent features for both the video and skeleton adversarial samples from the hidden layers of the two branch networks, denoted as $X'_R \in \mathbb{R}^{B \times C_R \times T_R \times H \times W}$ and $X'_s \in \mathbb{R}^{B \times C_s \times T_s \times V}$ respectively, where B is the input batch size, C_R and C_s are respectively the number of output channels for the video and skeleton, T_R and T_s are the number of frames of the video sequence and the skeleton sequence, respectively. H and W are the spatial dimensions of the video. V is the number of vertices in each frame of skeleton data. Subsequently, we create two attention weight matrices, $W_R \in \mathbb{R}^{1 \times C_R \times T_R \times H \times W}$ and $W_s \in \mathbb{R}^{1 \times C_s \times T_s \times V}$, based on the size of the features, noting that these matrices will be updated alongside the other parameters of the network. After replicating the two attention weight matrices B times along the batch size dimension, they are respectively multiplied with the corresponding extracted latent features, formulated as follows:

$$\begin{cases} \tilde{X}'_R = X'_R \otimes W_R, \\ \tilde{X}'_s = X'_s \otimes W_s, \end{cases} \quad (8)$$

where \otimes represents element-wise multiplication. At present, $\tilde{X}'_R \in \mathbb{R}^{B \times C_R \times T_R \times H \times W}$ and $\tilde{X}'_s \in \mathbb{R}^{B \times C_s \times T_s \times V}$ embody features adjusted by the attention weight matrices. Consecutively, they are separately fed into

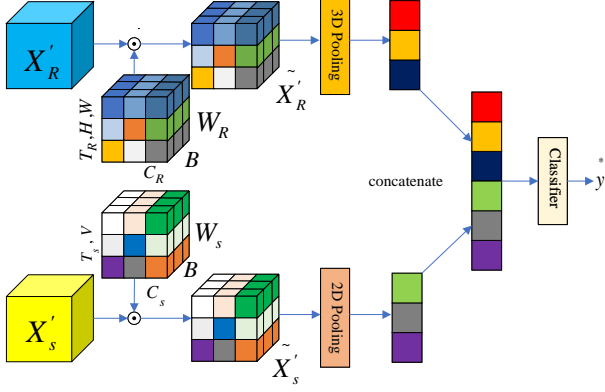


Figure 3. Architecture of AMR for two modalities. X'_R and X'_s , that represent the features at a given layer of two unimodal network, are the inputs to the module. For better visualization, we represent the spatiotemporal dimensions on a single axis.

a 3D pooling layer and a 2D pooling layer, and after flattening, the results are as follows:

$$\begin{cases} Z_R = \frac{1}{T_R \times H \times W} \sum_{t=1}^{T_R} \sum_{h=1}^H \sum_{w=1}^W \tilde{X}'_R, \\ Z_s = \frac{1}{T_s \times V} \sum_{t=1}^{T_s} \sum_{v=1}^V \tilde{X}'_s, \end{cases} \quad (9)$$

where

$$\begin{cases} Z_R \in \mathbb{R}^{B \times C_R}, \\ Z_s \in \mathbb{R}^{B \times C_s}, \end{cases} \quad (10)$$

Afterward, concatenate Z_R and Z_s along their respective channel dimensions to obtain $Z_{mul} \in \mathbb{R}^{B \times (C_R + C_s)}$, which is a vector that can be leveraged for classification. Finally, we can feed Z_{mul} into a fully connected layer (FC) to generate auxiliary classification prediction $\hat{y} \in \mathbb{R}^{B \times N}$. The above is a detailed description of the structure of AMR. In this module, three outputs will be provided: the features $\{\tilde{X}'_R, \tilde{X}'_s\}$ of the two modalities after attention weight adjustment, as well as the auxiliary classification prediction \hat{y} . The reweighted features continue to undergo forward propagation, with \hat{y} being utilized for the computation of the loss function.

We can integrate S AMRs into distinct hidden layers of a multimodal network. AMR can be deemed as an auxiliary component of the network, amenable to standard training or adversarial training in conjunction with the backbone network. We devise a novel loss function. Taking S instances of AMR inserted into the network as an example, the AMR loss function can be defined as follows:

$$L_{AMR} = l(y', y) + \frac{1}{S} \sum_{i=1}^S l(\hat{y}_i, y), \quad (11)$$

where $y' = f_{\theta}(x'_R, x'_s)$ signifies the ultimate predictive output for adversarial samples. y is the corresponding groundtruth label. $l(\cdot)$ denotes the cross-entropy loss function. Equation (11) encompasses the auxiliary classification results from multiple hidden layers, serving the purpose of deep supervision. It compels the network to learn more robust features at intermediate layers. To simultaneously consider both clean accuracy and robust accuracy, we also incorporate the outputs of natural samples into the loss function. The overall objective function for adversarial training of our AMR strategy is as follows:

$$L_{total} = l(\hat{y}, y) + \lambda L_{AMR}, \quad (12)$$

where $\hat{y} = f_{\theta}(x_R, x_s)$ denotes the prediction output for nature data, and λ functions as a hyperparameter to modulate the balance between the two loss components, allowing for customization to accommodate diverse datasets.

4. Experiment

4.1. Dataset

NTU-RGB+D [33] is a well-known large-scale multimodal human action recognition dataset. It contains 56,880 samples captured from 40 subjects performing 60 classes of activities at 80 view-points. Each action clip includes up to two people on the RGB video as well as 25 body joints on 3D coordinate space. The dataset is split in two ways: Cross-subject (X-Sub) and Cross-view (X-View), for which action subjects, camera views are different in training and validation. Unless otherwise specified, we conduct experiments on the X-sub splits for NTU-RGB+D.

iMiGUE [22] is a new dataset for the emotional artificial intelligence research: identity-free video dataset for Micro-Gesture Understanding and Emotion analysis (iMiGUE). iMiGUE comprises 36 female and 36 male players whose ages are between 17 and 38. After processing, a total of 18,419 video clips were included, divided into 32 action categories. The corresponding human skeleton data is also available.

4.2. Experimental Setup

Training setup: The overall architecture utilized in the experiments is illustrated in Figure 1. For iMiGUE and NTU-RGB+D datasets, we respectively insert 1 and 3 AMRs forward from the last fully connected (FC) layer of both the I3D and HCN branch networks. We employ SGD (momentum 0.9, batch size 16) to train the model for 40 epochs on the iMiGUE dataset and 80 epochs on the NTU-RGB+D dataset with weight decay $5e-4$ and initial learning rate 0.01. To adjust the learning rate adaptively, we utilize the cosine annealing learning rate decay strategy, where T_{max} corresponds to the respective number of epochs. For

Dataset	Attack	✓RS	✗R	✗S	✗RS	Avg	Unimodal ✓R	Unimodal ✓S
NTURGB+D	FGSM		50.48	66.42	22.12	46.34	81.46	81.46
	PGD10	88.09	6.39	25.18	1.66	11.08		
	CW		4.55	22.46	0.21	9.03		
iMiGUE	FGSM		6.86	48.61	2.52	19.03	59.48	46.86
	PGD10	62.37	0.35	19.94	0.24	6.84		
	CW		0.20	20.82	0.26	7.09		

Table 1. Various accuracy(%) on NTU-RGB+D and iMiGUE datasets under different attack methods. ✗R ✗S and ✗RS denote that only RGB, only skeleton, and both RGB video and skeleton inputs are attacked, respectively. We set ϵ_R and ϵ_s as $6/255$ respectively for ✗R ✗S and ✗RS . The symbol: ✓ means that inputs are clean. The baselines: Unimodal ✓R and Unimodal ✓S models are their respective unimodality branch models.

the internal maximization process, we use PGD10 attack to simultaneously obtain adversarial samples for two modalities, with a random start, step size $2.0/255$, and perturbation size $\epsilon_R, \epsilon_s = 8.0/255$. We use $\lambda = 1$ in all experiments.

Evaluation setup: We report the clean accuracy on natural examples and the adversarial accuracy on adversarial examples. We follow the widely used protocols in the adversarial research field. We consider three popular attack methods: FGSM [13], PGD20 [28], and CW [3] (optimized by PGD30).

4.3. Model Robustness under Multimodal Attacks

In Figure 2, we qualitatively observed the variation trends of model robustness under different attack scenarios. Here, we quantitatively analyze the diverse facets of model robustness under multimodal adversarial attacks. Section 4.3 presents various aspects of the model’s robustness on the NTU-RGB+D and iMiGUE datasets under different attack scenarios. To better interpret the multimodal robustness, we also include results from two baselines: Unimodal R and Unimodal S, which are two unimodality models and only use video and skeleton modalities, respectively. Clearly, all of the three attack methods can significantly decrease recognition results. The results demonstrate that at the same level of attack intensity, the accuracy of the skeleton modality is notably higher than that of the RGB modality and the multimodal, indicating the stronger resilience of the skeleton modality to interference. When we leverage different attack strategies to perform unimodality attacks on the NTU-RGB+D and iMiGUE datasets, RGB-skeleton models: ✗R and ✗S are always inferior to Unimodal ✓S and Unimodal ✓R, respectively. Note that ✗R and ✗S have clean skeleton and RGB modalities, respectively. This discovery suggests that when one modality input is subjected to an attack, the multimodal integration weakens the recognition performance, consistent with the conclusion in [36].

4.4. Robustness Analysis and Evaluation

To validate the effectiveness of the proposed AMR, we compared it with several baselines and SOTA methods. For

Defence(NTURGB+D)	✓RS	✗R	✗S	✗RS	Avg	RI
None	88.09	2.09	11.84	0.15	4.69	0
AT[28] with ✗R	75.59	75.59	4.75	4.74	28.36	11.17
AT[28] with ✗S	78.44	2.09	78.31	2.06	27.49	13.15
AT[28] with ✗RS	61.42	61.41	40.41	40.32	47.38	16.02
MinSim[36]	80.85	2.93	80.85	3.03	28.85	16.92
ExFMem[36]	84.19	0.00	60.04	0.00	20.01	11.42
MinSim+ExFMem[36]	82.81	2.02	75.77	2.14	26.64	16.67
AMR(Ours)	81.65	68.54	80.40	66.52	71.82	60.69
Defence(iMiGUE)	✓RS	✗R	✗S	✗RS	Avg	RI
None	62.37	0.11	6.97	0.07	2.38	0
AT[28] with ✗R	47.32	47.32	0.00	0.00	15.77	-1.66
AT[28] with ✗S	60.20	0.00	59.87	0.00	19.96	15.41
AT[28] with ✗RS	39.01	39.01	31.89	31.95	34.28	8.54
MinSim[36]	60.60	0.00	53.28	0.00	17.76	13.61
ExFMem[36]	60.53	0.00	43.00	0.00	14.33	10.11
MinSim+ExFMem[36]	60.62	0.00	55.05	0.00	18.35	14.22
AMR(Ours)	43.89	43.89	32.85	32.94	36.56	15.70

Table 2. The various accuracy(%) metrics on the NTU-RGB+D and iMiGUE datasets are presented with different defense methods. Here, we utilize PGD20 (perturbation budgets $\epsilon_R, \epsilon_s = 8/255$, step size $\alpha = 2/255$) to generate adversarial samples for assessing robustness accuracy. The best results among the metrics are denoted in bold.

the sake of fairness, the SOTA methods we compared to are also multimodal defense strategies, rather than unimodal defense approaches.

Baselines and SOTAs: 1) None: RGB-skeleton multimodal network without any defense; 2) AT [28] with ✗R : adversarial training using only adversarial samples from the RGB video modality; 3) AT [28] with ✗S : adversarial training using only adversarial samples from the skeleton modality; 4) AT [28] with ✗RS : adversarial training using adversarial samples from both the RGB video and the skeleton modality simultaneously; 5) MinSim [36]: a recent SOTA multimodal defense mechanism, namely the dissimilarity constraint to encourage multimodal dispersion and unimodal compactness; 6) ExFMem [36]: a recent SOTA multimodal defense mechanism, namely adopting external feature memory banks to denoise attacked video and skeleton examples at a feature level; 7) MinSim+ExFMem [36]: the combination of methods MinSim and ExFMem; 8) AMR: our proposed defence method (certainly based on AT with ✗RS)

Evaluation Metrics: To evaluate the performance of different defense methods and facilitate a fair comparison, we follow the protocol proposed in [36]. We use recognition accuracy as the metric. Results from both the nature samples: ✓RS and adversarial samples: ✗R , ✗S , and ✗RS are computed. Since there are multiple defense results under different multimodal attacks for a single method, we also use the averaged accuracy:

$$Avg = \frac{1}{3}(\text{✗R} + \text{✗S} + \text{✗RS}), \quad (13)$$

as an overall metric to evaluate robustness of different de-

iMiGUE	AMR 1	AMR 2	AMR 3
Mean value of W_R	0.4666	0.4663	0.4665
Mean value of W_s	0.4693	0.4939	0.4905
NTURGB+D	AMR 1	AMR 2	AMR 3
Mean value of W_R	0.3067	0.3065	0.3067
Mean value of W_s	0.3187	0.3275	0.3240

Table 3. The mean weights of different AMRs trained on the iMiGUE and NTURGB+D datasets.

fenses. To comprehensively assess the model’s clean accuracy and robust accuracy, Tian *et al.* [36] propose a relative improvement (RI) metric:

$$RI = (\check{RS}_m + Avg_m) - (\check{RS}_n + Avg_n), \quad (14)$$

where results from both clean samples and adversarial samples are considered, and the m refers to a defense method and n refers to a base model, which is the baseline: None in our experiments. If a defense method decreases clean data performance, the RI will penalize it accordingly.

Result Analysis: Table 2 and display various accuracy metrics against PGD20 attack on the NTU-RGB+D and iMiGUE datasets with different defense methods. See more results in Table 6 and Table 7 in the Supplementary. The \check{RS} , representing the accuracy on clean data, is highest with the None, *i.e.*, the baseline. This aligns with adversarial training and robustness theory, as incorporating adversarial examples during training is bound to compromise clean accuracy. AT enhances robustness that None lacks. It can be observed that training with adversarial samples from different modalities significantly enhances the robustness of the corresponding modality. AT with $\mathcal{X}R$ and AT with $\mathcal{X}S$ achieve the highest values on $\mathcal{X}R$ and $\mathcal{X}S$, respectively. However, they exhibit almost no robustness on the modalities without the use of adversarial samples. AT with $\mathcal{X}RS$ strikes a balance between the robustness of various modalities. The most relevant method to our work, MinSim/ExFMem, achieves high performance in \check{RS} and $\mathcal{X}S$, but its accuracy is very poor in $\mathcal{X}R$ and $\mathcal{X}RS$. This is because this method does not consider the differences in the inherent robustness of different modalities and does not specially handle the modalities that are more susceptible to attacks. From the results, it is evident that our AMR outperforms all compared methods in terms of $\mathcal{X}RS$ and Avg (representing robust accuracy) as well as the comprehensive performance indicator RI. AMR visibly boosts the robustness of the vulnerable RGB modality by taking into account the characteristics of each modality.

4.5. Weight Distribution Visualization

As mentioned in Section 3.3, our AMR can assign appropriate weights to the features of two modalities. Specifically, more robust features will receive higher weights,

Defence(NTU-RGB+D)	\check{RS}	$\mathcal{X}R$	$\mathcal{X}S$	$\mathcal{X}RS$	Avg	RI
None	88.09	0.98	9.00	0.00	3.33	0
AT with $\mathcal{X}RS$ (0 AMR)	61.42	61.41	38.74	38.82	46.32	16.32
1 AMR	76.44	76.44	37.89	37.83	50.72	37.74
2 AMRs	81.61	67.93	80.53	66.00	71.49	61.68
3 AMRs	81.65	68.62	80.35	66.61	71.86	62.09
Defence(iMiGUE)	\check{RS}	$\mathcal{X}R$	$\mathcal{X}S$	$\mathcal{X}RS$	Avg	RI
None	62.37	0.13	6.67	0.02	2.27	0
AT with $\mathcal{X}RS$ (0 AMR)	39.01	39.01	30.94	30.97	33.64	8.01
1 AMR	43.89	43.89	32.46	32.50	36.28	15.53
2 AMRs	42.93	42.93	31.27	31.16	35.12	13.41
3 AMRs	42.38	42.38	33.33	33.42	36.38	14.12

Table 4. The impact of the number of AMRs on accuracy against CW (perturbation budgets $\epsilon_R, \epsilon_s = 8/255$, step size $\alpha = 2/255$, optimized by PGD30) attack on both datasets.

λ	\check{RS}	$\mathcal{X}R$	$\mathcal{X}S$	$\mathcal{X}RS$	Avg	RI
None	62.37	0.11	6.97	0.07	2.38	0
0.1	60.25	36.23	59.15	32.02	42.47	37.97
0.5	45.08	45.08	27.39	27.24	33.24	13.57
1	43.89	43.89	32.85	32.94	36.56	15.70
2	42.76	42.76	34.82	34.80	37.46	15.47
5	41.64	41.64	37.50	37.56	38.9	15.79

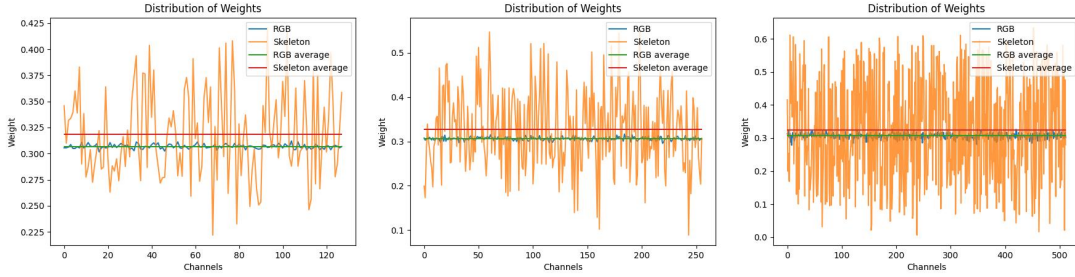
Table 5. Robust comparison of different λ on iMiGUE with 1 AMR.

while less robust ones will receive smaller weights. To verify if our results align with the proposed idea, we visualize the attention weight distribution of the AMR in the pretrained network, denoted as W_R and W_s in the Equation (8), as shown in the Figure 4. Further, we also calculate the overall mean of weight matrices within each AMR. The results in Table 3 indicate that in each AMR, the average weight values corresponding to the skeletal modality are greater than the average weight values for the RGB modality. The experimental results align with our expectations. Observing Figure 4, we notice a more substantial fluctuation in the weight distribution of the skeleton modality compared to the RGB modality. This phenomenon can be attributed to the greater impact of adversarial attacks on the RGB modality. Additionally, the skeleton encompasses a broader range of motion features, leading the model to incline towards adjusting the weights of the skeleton modality during the training process to enhance overall robustness.

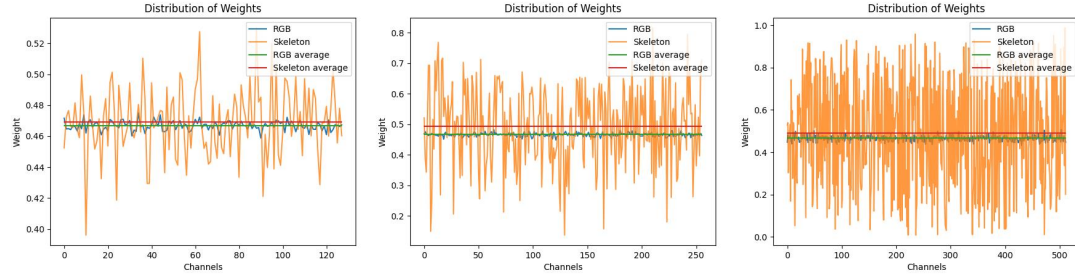
4.6. Ablation Study

4.6.1 The Impact of the number of AMR

Our proposed AMR is an AT-based method. In order to evaluate the effectiveness of AMR and investigate the impact of the number of AMRs on overall performance, we conduct ablation experiments using None and AT with $\mathcal{X}RS$ (which can be considered as having no AMRs, *i.e.*, 0 AMR) as baselines. We start from the last FC layer of the branch network and insert 1 to 3 AMRs in a forward manner.



(a) Distribution of Wights in different AMR trained on iMiGUE



(b) Distribution of Wights in different AMR trained on NTURGB+D

Figure 4. The distribution of weight matrices in different AMRs trained on two datasets. The x -axis represents the channels, and the y -axis represents the average weight values corresponding to each channel.

The other parameters remain the same as described in Section 4.2. The results in Table 4 demonstrate that, with an increasing number of AMRs, the RI continuously improves on the NTU-RGB+D dataset. On the iMiGUE dataset, the RI reaches its peak with one AMR, and further increasing the number of AMRs leads to a decrease in RI. This is because NTU-RGB+D is a large-scale dataset containing 57k samples, while iMiGUE has only 18k samples. We can conclude that for larger datasets, we require more AMRs to achieve better performance. However, regardless of the number of AMRs, the use of AMR consistently outperforms not using AMR by a significant margin.

4.6.2 The Impact of the hyperparameter λ

In this part, we evaluate the impact of λ in Equation (12). λ acts as a hyperparameter, enabling us to fine-tune the trade-off between the two loss component. We train our models with one AMR on the iMiGUE dataset with 5 different values $\lambda = [0.1, 0.5, 1, 2, 5]$, and evaluate their accuracy under PGD20 attacks. The natural accuracy decrease with the increase of λ . This is because the increase of λ reduces the proportion of natural samples in the loss function. The highest RI is achieved when $\lambda = 0.1$.

5. Conclusion

In this paper, we investigate the differences in robustness among various data modalities and find that the skeleton modality is more robust than the RGB. Based on this observation, we creatively propose the Attention-based Modality Reweigher (AMR). It autonomously reweights different modalities, enhancing overall robustness. To accommodate AMR, we introduce a new loss function that encompasses multiple auxiliary classification results, serving as deep supervision while considering both clean accuracy and robust accuracy. We conduct extensive experiments on popular benchmark datasets and evaluate performance against SOTA attack methods. The results demonstrate that our approach significantly improves the robustness of multimodal action recognition models compared to other methods. Importantly, it balances the discrepancy in robustness between different modalities and achieves new SOAT performance without any additional data.

6. Acknowledement

This work was supported by National Natural Science Foundation of China under Grant 62171309 and 62306061, and Guangdong Basic and Applied Basic Research Foundation (Grant No. 2023A1515140037).

References

- [1] M. Andriushchenko, F. Croce, N. Flammarion, and M. Hein. Square attack: A query-efficient black-box adversarial attack via random search. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII*, page 484–501, Berlin, Heidelberg, 2020. Springer-Verlag.
- [2] V. H. Buch, I. Ahmed, and M. Maruthappu. Artificial intelligence in medicine: current trends and future possibilities. *British Journal of General Practice*, 68(668):143–144, 2018.
- [3] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57, 2017.
- [4] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4733, 2017.
- [5] J. Chen, B. Yuan, and M. Tomizuka. Model-free deep reinforcement learning for urban autonomous driving. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pages 2765–2771, 2019.
- [6] F. Croce and M. Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In *Proceedings of the 37th International Conference on Machine Learning, ICML’20*. JMLR.org, 2020.
- [7] F. Croce and M. Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. *ICML’20*, pages 2206–2216. JMLR.org, 2020.
- [8] H. Duan, Y. Zhao, K. Chen, D. Lin, and B. Dai. Revisiting skeleton-based action recognition. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2959–2968, 2022.
- [9] M. Duan, K. Li, J. Deng, B. Xiao, and Q. Tian. A novel multi-sample generation method for adversarial attacks. *ACM Trans. Multimedia Comput. Commun. Appl.*, 18(4), mar 2022.
- [10] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song. Robust physical-world attacks on deep learning visual classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1625–1634, 2018.
- [11] Z. Fu, Z. Mao, Y. Song, and Y. Zhang. Learning semantic relationship among instances for image-text matching. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15159–15168, 2023.
- [12] D. Glasner, S. Bagon, and M. Irani. Super-resolution from a single image. In *2009 IEEE 12th international conference on computer vision*, pages 349–356. IEEE, 2009.
- [13] I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *CoRR*, abs/1412.6572, 2014.
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [15] H. Kannan, A. Kurakin, and I. Goodfellow. Adversarial logit pairing. *ArXiv*, abs/1803.06373, 2018.
- [16] F. Khader, G. Mueller-Franzes, T. Wang, T. Han, S. T. Arasteh, C. Haarburger, J. Stegmaier, K. Bresslem, C. Kuhl, S. Nebelung, J. N. Kather, and D. Truhn. Medical diagnosis with large scale multimodal transformers: Leveraging diverse data for more accurate diagnosis. *ArXiv*, abs/2212.09162, 2022.
- [17] K. A. Kinfu and R. Vidal. Analysis and extensions of adversarial training for video classification. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 3415–3424, 2022.
- [18] B. Kong, X. Wang, Z. Li, Q. Song, and S. Zhang. Cancer metastasis detection via spatially structured deep network. In M. Niethammer, M. Styner, S. Aylward, H. Zhu, I. Oguz, P.-T. Yap, and D. Shen, editors, *Information Processing in Medical Imaging*, pages 236–248, Cham, 2017. Springer International Publishing.
- [19] S. Kullback and R. A. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 22:79–86, 1951.
- [20] C. Li, Q. Zhong, D. Xie, and S. Pu. Co-occurrence feature learning from skeleton data for action recognition and detection with hierarchical aggregation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18*, page 786–792. AAAI Press, 2018.
- [21] L. Li, R. Ma, Q. Guo, X. Xue, and X. Qiu. Bert-attack: Adversarial attack against bert using bert. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202. Association for Computational Linguistics, 2020.
- [22] X. Liu, H. Shi, H. Chen, Z. Yu, X. Li, and G. Zhao. imigue: An identity-free video dataset for micro-gesture understanding and emotion analysis. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10626–10637, 2021.
- [23] X. Liu, K. Yuan, X. Niu, J. Shi, Z. Yu, H. Yue, and J. Yang. Multi-scale promoted self-adjusting correlation learning for facial action unit detection. *arXiv preprint arXiv:2308.07770*, 2023.
- [24] X. Liu, Y. Zhang, Z. Yu, H. Lu, H. Yue, and J. Yang. rppg-mae: Self-supervised pretraining with masked autoencoders for remote physiological measurements. *IEEE Transactions on Multimedia*, 2024.
- [25] H. Lu, X. Niu, J. Wang, Y. Wang, Q. Hu, J. Tang, Y. Zhang, K. Yuan, B. Huang, Z. Yu, et al. Gpt as psychologist? preliminary evaluations for gpt-4v on visual affective computing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 322–331, 2024.
- [26] X. Ma, Y. Niu, L. Gu, Y. Wang, Y. Zhao, J. Bailey, and F. Lu. Understanding adversarial attacks on deep learning based medical image analysis systems. *Pattern Recognition*, 110:107332, 2021.
- [27] D. Madaan, J. Shin, and S. J. Hwang. Adversarial neural pruning with latent vulnerability suppression. In *Proceedings of the 37th International Conference on Machine Learning, ICML’20*, pages 6575–6585. JMLR.org, 2020.
- [28] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning*

- Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [29] J. H. Moon, H. Lee, W. Shin, Y.-H. Kim, and E. Choi. Multi-modal understanding and generation for medical images and text via vision-language pre-training. *IEEE Journal of Biomedical and Health Informatics*, 26(12):6070–6080, 2022.
- [30] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2574–2582, 2016.
- [31] N. R. Pal and S. K. Pal. A review on image segmentation techniques. *Pattern recognition*, 26(9):1277–1294, 1993.
- [32] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [33] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1010–1019, 2016.
- [34] J. Su, D. V. Vargas, and K. Sakurai. One pixel attack for fooling deep neural networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841, 2019.
- [35] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [36] Y. Tian and C. Xu. Can audio-visual integration strengthen robustness under multimodal attacks? In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5597–5607, 2021.
- [37] H. R. Vaezi Joze, A. Shaban, M. L. Iuzzolino, and K. Koishida. Mmtm: Multimodal transfer module for cnn fusion. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13286–13296, 2020.
- [38] H. Wang, G. Wang, Y. Li, D. Zhang, and L. Lin. Transferable, controllable, and inconspicuous adversarial attacks on person re-identification with deep mis-ranking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 342–351, 2020.
- [39] L. Wang, Z. He, J. Tang, R. Dang, N. Wang, C. Liu, and Q. Chen. A dual semantic-aware recurrent global-adaptive network for vision-and-language navigation. In E. Elkind, editor, *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pages 1479–1487. International Joint Conferences on Artificial Intelligence Organization, 8 2023. Main Track.
- [40] Y. Wang, D. Zou, J. Yi, J. Bailey, X. Ma, and Q. Gu. Improving adversarial robustness requires revisiting misclassified examples. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [41] X. Wei, J. Zhu, S. Yuan, and H. Su. Sparse adversarial perturbations for videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8973–8980, 2019.
- [42] S. Wu, J. Sang, K. Xu, G. Zheng, and C. Xu. Adaptive adversarial logits pairing. *ACM Transactions on Multimedia Computing, Communications and Applications*, 20(2), oct 2023.
- [43] Y. Xu, X. Wei, P. Dai, and X. Cao. A2sc: Adversarial attacks on subspace clustering. *ACM Trans. Multimedia Comput. Commun. Appl.*, 19(6), jul 2023.
- [44] S. Yan, Y. Xiong, and D. Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence, AAAI’18/IAAI’18/EAAI’18*. AAAI Press, 2018.
- [45] B. X. Yu, Y. Liu, X. Zhang, S.-h. Zhong, and K. C. Chan. Mmnet: A model-based multimodal network for human action recognition in rgb-d videos. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3):3522–3538, 2023.
- [46] K. Yuan, Z. Yu, X. Liu, W. Xie, H. Yue, and J. Yang. Auformer: Vision transformers are parameter-efficient facial action unit detectors. *arXiv preprint arXiv:2403.04697*, 2024.
- [47] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. E. Ghaoui, and M. Jordan. Theoretically principled trade-off between robustness and accuracy. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7472–7482. PMLR, 09–15 Jun 2019.
- [48] J. Zhang, J. Zhu, G. Niu, B. Han, M. Sugiyama, and M. S. Kankanhalli. Geometry-aware instance-reweighted adversarial training. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [49] Y. Zhang, H. Lu, X. Liu, Y. Chen, and K. Wu. Advancing generalizable remote physiological measurement through the integration of explicit and implicit prior knowledge. *arXiv preprint arXiv:2403.06947*, 2024.

Adversarial Robustness in RGB-Skeleton Action Recognition: Leveraging Attention Modality Reweigher

Supplementary Material

Defence (NTURGB+D)	✓RS	✗R	✗S	✗RS	Avg	RI
None	88.09	49.08	64.05	17.36	43.50	0
AT[28] with ✗R	75.59	75.59	22.18	22.20	39.99	-16.01
AT[28] with ✗S	78.44	8.36	78.35	8.33	31.68	-21.47
AT[28] with ✗RS	61.42	61.41	47.04	47.00	51.82	-18.35
MinSim[36]	80.85	13.80	80.85	13.80	36.15	-14.59
ExFMem[36]	84.19	20.96	75.27	10.59	35.61	-11.79
MinSim+ExFMem[36]	82.81	17.39	80.27	16.25	37.97	-10.81
AMR(Ours)	81.65	71.99	80.77	70.63	74.46	24.97
Defence (iMiGUE)	✓RS	✗R	✗S	✗RS	Avg	RI
None	62.37	5.98	48.85	2.26	19.03	0
AT[28] with ✗R	47.32	47.32	10.54	9.77	22.54	-11.54
AT[28] with ✗S	60.20	2.78	60.11	3.11	22.0	0.8
AT[28] with ✗RS	39.01	38.98	34.43	34.36	35.92	-6.47
MinSim[36]	60.60	5.19	57.59	3.62	22.13	1.33
ExFMem[36]	60.53	6.14	55.64	2.98	21.59	0.72
MinSim+ExFMem[36]	60.62	4.30	58.36	3.29	21.98	1.2
AMR(Ours)	43.89	43.89	36.18	36.45	38.84	1.17

Table 6. The various accuracy(%) metrics on the NTU-RGB+D and iMiGUE datasets are presented with different defense methods. Here, we utilize FGSM (perturbation budgets $\epsilon_R, \epsilon_s = 8/255$, step size $\alpha = 2/255$) to generate adversarial samples for assessing robustness accuracy. The best results among the metrics are denoted in bold.

Defence(NTURGB+D)	✓RS	✗R	✗S	✗RS	Avg	RI
None	88.09	0.98	9.00	0.00	3.33	0
AT[28] with ✗R	75.59	75.59	2.80	2.89	27.09	11.26
AT[28] with ✗S	78.44	1.41	78.31	1.41	27.04	14.06
AT[28] with ✗RS	61.42	61.41	38.74	38.82	46.32	16.32
MinSim[36]	80.85	0.50	80.84	0.48	27.27	16.70
ExFMem[36]	84.19	0.00	60.98	0.00	20.33	13.10
MinSim+ExFMem[36]	82.81	0.00	75.60	0.00	25.20	16.59
AMR(Ours)	81.65	68.62	80.35	66.61	71.86	62.09
Defence(iMiGUE)	✓RS	✗R	✗S	✗RS	Avg	RI
None	62.37	0.13	6.67	0.02	2.27	0
AT[28] with ✗R	47.32	47.32	0.00	0.00	15.77	-1.55
AT[28] with ✗S	60.20	0.00	59.79	0.00	19.93	15.49
AT[28] with ✗RS	39.01	39.01	30.94	30.97	33.64	8.01
MinSim[36]	60.60	0.00	53.74	0.02	17.92	13.88
ExFMem[36]	60.53	0.00	45.19	0.00	15.06	10.95
MinSim+ExFMem[36]	60.62	0.00	54.72	0.00	18.24	14.20
AMR(Ours)	43.89	43.89	32.46	32.50	36.28	15.53

Table 7. The various accuracy(%) metrics on the NTU-RGB+D and iMiGUE datasets are presented with different defense methods. Here, we utilize CW attack (perturbation budgets $\epsilon_R, \epsilon_s = 8/255$, step size $\alpha = 2/255$, optimized by PGD30) to generate adversarial samples for assessing robustness accuracy. The best results among the metrics are denoted in bold.

Defence(NTURGB+D)	✓RS	✗R	✗S	✗RS	Avg	RI
None	88.09	2.09	11.84	0.15	4.69	0
AT with ✗RS(0 AMR)	61.42	61.41	38.74	38.82	47.38	16.02
1 AMR	76.44	76.44	38.29	38.32	51.02	34.68
2 AMRs	81.61	67.80	80.58	65.88	71.42	60.25
3 AMRs	81.65	68.54	80.40	66.52	71.82	60.69
Defence(iMiGUE)	✓RS	✗R	✗S	✗RS	Avg	RI
None	62.37	0.13	6.67	0.02	2.27	0
AT with ✗RS(0 AMR)	39.01	39.01	31.89	31.95	34.28	8.54
1 AMR	43.89	43.89	32.85	32.94	36.56	15.70
2 AMRs	42.93	42.93	32.35	32.43	35.92	14.08
3 AMRs	42.38	42.38	34.34	34.32	37.01	14.64

Table 8. The impact of the number of AMRs on accuracy against PGD20 (perturbation budgets $\epsilon_R, \epsilon_s = 8/255$, step size $\alpha = 2/255$) attack on both datasets.