

# Adversarial Attack on Skeleton-based Human Action Recognition

Jian Liu, Naveed Akhtar, and Ajmal Mian,

**Abstract**—Deep learning models achieve impressive performance for skeleton-based human action recognition. However, the robustness of these models to adversarial attacks remains largely unexplored due to their complex spatio-temporal nature that must represent sparse and discrete skeleton joints. This work presents the first adversarial attack on skeleton-based action recognition with graph convolutional networks. The proposed targeted attack, termed Constrained Iterative Attack for Skeleton Actions (CIASA), perturbs joint locations in an action sequence such that the resulting adversarial sequence preserves the temporal coherence, spatial integrity, and the anthropomorphic plausibility of the skeletons. CIASA achieves this feat by satisfying multiple physical constraints, and employing spatial skeleton realignments for the perturbed skeletons along with regularization of the adversarial skeletons with Generative networks. We also explore the possibility of semantically imperceptible localized attacks with CIASA, and succeed in fooling the state-of-the-art skeleton action recognition models with high confidence. CIASA perturbations show high transferability for black-box attacks. We also show that the perturbed skeleton sequences are able to induce adversarial behavior in the RGB videos created with computer graphics. A comprehensive evaluation with NTU and Kinetics datasets ascertains the effectiveness of CIASA for graph-based skeleton action recognition and reveals the imminent threat to the spatio-temporal deep learning tasks in general.

**Index Terms**—Adversarial attack, Adversarial examples, Action recognition, Skeleton actions, Adversarial perturbations, Spatio-temporal.

## I. INTRODUCTION

Skeleton representation provides the advantage of capturing accurate human pose information while being invariant to action-irrelevant details such as scene background, clothing patterns and illumination conditions. This makes skeleton-based action recognition an appealing approach [1]–[6]. The problem is also interesting for multiple application domains, including security, surveillance, animation and human-computer interactions etc. Recent contributions in this direction predominantly exploit deep models to encode spatio-temporal dependencies of the skeleton sequences [7]–[10], and achieve remarkable recognition accuracy on benchmark action datasets [11]–[14].

Although deep learning has been successfully applied to many complex problems, it is now known that deep models are vulnerable to adversarial attacks [15], [16]. These attacks can alter model predictions at will by adding imperceptible perturbations to the input. After the discovery of this intriguing weakness of deep learning [15], many adversarial attacks have surfaced for a variety of vision tasks [17]–[20]. Developing and investigating these attacks not only enhances our understanding of the inner workings of the neural networks [21], but

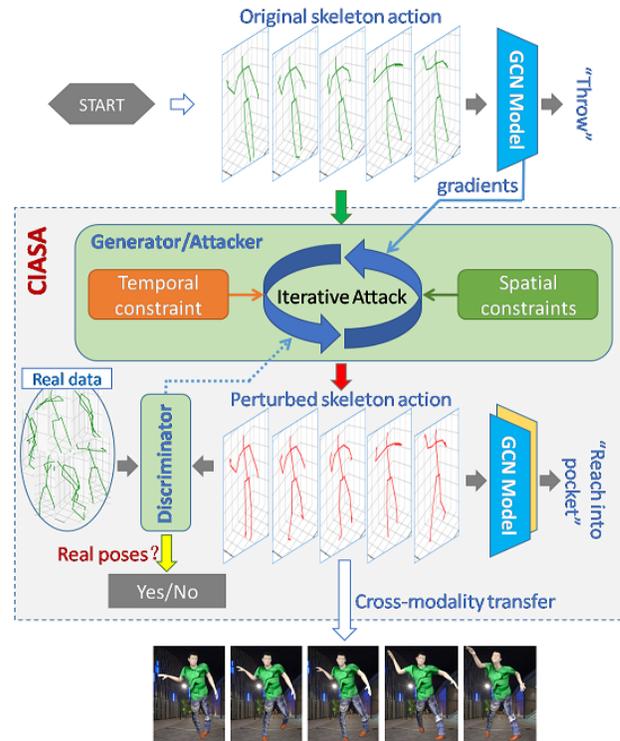


Fig. 1. Constrained Iterative Attack for Skeleton Actions (CIASA) schematics. Model gradients are computed for input action sequence to iteratively minimise the model’s loss for a target label in small step, while accounting for the relevant spatio-temporal constraints. A generator-discriminator framework further ensures anthropomorphic plausibility of the skeletons. Besides cross-model transferability, the attack can also affect RGB videos generated with computer graphics using the skeletons perturbed by CIASA.

also provides valuable insights for improving the robustness of deep learning in practical adversarial settings.

Deep models for skeleton-based action recognition may also be vulnerable to adversarial attacks. However, adversarial attacks on these models are yet to be explored. A major challenge in this regard is that the skeleton data representation differs significantly from image representation, for which the existing attacks are primarily designed. Human skeleton data is sparse and discrete that evolves over time in rigid spatial configurations. This prevents an attacker from freely modifying the skeletons without raising obvious attack suspicions. Skeleton actions also allow only subtle perturbations along the temporal dimension to preserve the natural action dynamics. In summary, adversarial attacks on skeleton data must carefully account for the skeleton’s spatial integrity, temporal coherence and anthropomorphic plausibility. Otherwise, the attack

may be easily detectable. These challenges have so far kept skeleton-based action recognition models away from being scrutinized for adversarial robustness.

In this work, we present the *first* adversarial attack on deep skeleton action recognition. In particular, we attack the (most) promising branch of graph convolutional networks [22] for skeleton-based action recognition [8]. These models represent actions as spatio-temporal graphs that encode intra-body and inter-frame connections as edges, and body joints as nodes. Graph convolution operations are leveraged to model the spatio-temporal dependencies within the skeleton sequences. The physical significance of nodes and edges in these models imposes unique constraints over the potential attacks. For instance, the graph nodes for a skeleton sequence can not be added or removed because the number of joints in the skeleton must remain fixed. Similarly, the lengths of intra-body edges in the graph can not be altered arbitrarily as they represent bones. Moreover, inter-frame edges must always connect the same joints along the temporal dimension. Rooted in the skeleton data, such constraints thoroughly distinguish the adversarial attacks on skeleton-based action recognition models from the attacks developed for other kinds of graph networks [23].

We develop an iterative scheme called Constrained Iterative Attack for Skeleton Actions (CIASA), to generate the desired adversarial skeleton sequences, see Fig. 1. For a given action, CIASA iteratively perturbs its skeleton sequence in small steps to minimize the model prediction loss for a pre-selected target class while satisfying multiple physical constraints to keep the resulting adversarial sequence natural. In particular, it accounts for spatio-temporal constraints that preserve intra-skeleton joint connections, inter-frame joint connections, and the skeleton bone lengths using a mechanism termed ‘spatial skeleton realignment’. For perturbation imperceptibility, it restricts the  $\ell_\infty$ -norm of the added noise. Additionally, it imposes external temporal dynamics constraints for imperceptible evolution of the adversarial patterns in the skeleton sequence. To further ensure anthropomorphic plausibility of the adversarial skeleton sequence, it exploits the Generative Adversarial Network (GAN) framework [24]. The used GAN configuration reduces the difference between the distribution of adversarial samples generated by our iterative scheme and the clean ground truth samples.

We analyze the proposed attack by allowing different modes in which CIASA can be used by an attacker. Analogous to standard image based attacks, we allow perturbation of all skeleton joints in the *basic* mode. In a *localized* mode, we provide the flexibility of perturbing only localized regions, e.g. legs of skeleton. This type of attack is particularly suitable to skeleton actions where an attacker may independently alter motion of the least relevant joints for an action to change the prediction. We also introduce an *advanced* attack mode that further allows a hierarchical magnitude variation in joint perturbations based on the graph structure of the joints.

The notion of localized perturbation also leads to *semantically* imperceptible perturbations under CIASA where significant perturbation still remains hard to perceive because it is applied to the least significant joints for the original action semantics. Besides demonstrating high fooling rates for the

state-of-the-art graph skeleton action recognition model ST-GCN [8] on NTU [11] and Kinetics [25] datasets, we also show high cross-model transferability of the proposed attack. Additionally, we show that videos generated (using computer graphics) from the adversarial skeletons (CIASA’s advanced mode) result in lower action recognition accuracy implying that the attack can be launched in the real world. To the best of our knowledge, this is the first of its kind demonstration of transferability of adversarial attacks beyond a single data modality.

The rest of this article is organized as follows. We review the related literature in Section II. The relevant concepts of graph skeleton action recognition are revisited in Section III along with the problem formulation. In Section IV, we give the implementation details of the proposed attack scheme. Experimental results are provided in Section V. The article concludes in Section VI.

## II. RELATED WORK

### A. Skeleton-based Action Recognition

The use of skeleton data in action recognition becomes popular as reliable skeleton data can be obtained from modern RGB-D sensors (e.g. Microsoft Kinect), or extracted from images taken from a single RGB camera [26]. A skeleton action is represented as a sequence of human skeletons, which encode rich spatio-temporal information regarding human motions. Early research in skeleton-based action recognition formulated skeleton joints and their temporal variations as trajectories [2]. Huang *et al.* [27] incorporated the Lie group structure into the task, and transformed the high-dimensional Lie group trajectory into temporally aligned Lie group features for skeleton-based action recognition.

To leverage the power of convolutional neural network, Du *et al.* [3] represented a skeleton sequence as a matrix by concatenating the joint coordinates. The matrix is arranged as an image which can be fed into CNN for recognition. Similarly, Ke *et al.* [5] transformed a skeleton sequence into three clips of gray-scale images that encode spatial dependencies between the joints by inserting reference joints. To fit the target neural networks, these methods re-size the transformed images. Liu *et al.* [28] proposed a universal unit ‘skepxel’ to create images of arbitrary dimensions for CNN processing. In addition to CNNs, Recurrent Neural Networks are also employed to model temporal dependencies in skeleton based human action analysis [29]–[31].

To directly process the sparse skeleton data with neural networks, graph convolutional network (GCN) [22] is used for action recognition. Since GCN is particularly relevant to this work, we review its relevant literature and application to action recognition in more detail.

### B. Graph Convolution Networks

The topology of human skeleton joints is a typical graph structure, where the joints and bones are respectively interpreted as graph nodes and edges. Consequently, there have been several recent attempts in modeling human skeleton

actions using graph representation and exploiting the spatio-temporal dependencies in skeleton sequences with the help of graph-based convolutional network (GCN).

Yan *et al.* [8] used graph convolutional networks as a spatial-temporal model (ST-GCN) that aims to capture embedded patterns in the spatial configuration of skeleton joints and their temporal dynamics simultaneously. Along the skeleton sequence, they defined a graph convolution operation, where the input is the joint coordinate vectors on the graph nodes. The convolution kernel samples the neighboring joints within the skeleton frame as well as the temporally connected joints at a defined temporal range.

Tang *et al.* [32] incorporated deep reinforcement learning with graph neural network to recognize skeleton-based actions. Their model distills the most informative skeleton frames and discards the ambiguous ones. As opposed to previous works where joints dependency is limited in the real physical connection (intrinsic dependency), they proposed extrinsic joint dependency, which exploits the relationship between joints that have physical disconnection. Since, graph representation of skeleton is crucial to graph convolution, Gao *et al.* [33] formulated the skeleton graph representation as an optimization problem, and proposed graph regression to statistically learn the underlying graph from multiple observations. The learned sparse graph pattern links both physical and non-physical edges of skeleton joints, along with the spatio-temporal dimension of the skeleton action sequences.

To justify the importance of bones' motions in skeleton action recognition, Zhang *et al.* [34] focused on skeleton bones and extended the graph convolution from graph nodes to graph edges. Their proposed graph edge convolution defines a receptive field of edges, which consists of a center edge and its spatio-temporal neighbours. By combining the graph edge and node convolutions, they proposed a two-stream graph neural network, which achieved remarkable performances on benchmark datasets. Similarly, Shi *et al.* [35] also proposed a two-stream framework to model joints and bones information simultaneously.

### C. Adversarial Attacks on Graph Data

Adversarial attacks [15] have recently attracted significant research attention [36], resulting in few attacks on graph data as well. However, compared to the adversarial attacks for image data [16], [37]–[39], several new challenges appear in attacking graph data [40]. First, the graph structure and features of graph nodes are in discrete domain with certain pre-defined structures, which leaves a lower degree of freedom for creating adversarial perturbations. Second, the imperceptibility of adversarial perturbations in graph data is neither easy to define nor straightforward to achieve, as the discrete graph data inherently prevents infinitesimal small changes [23].

Dai *et al.* [41] focused on attacking structural information, i.e. adding/deleting graph edges, to launch adversarial attacks on graph structured data. Given the gradient information of target classifier, one of their proposed attacks modifies the graph edges that are most likely to change the objective. In addition to modifying graph edges, Zügner *et al.* [23] adopted

an attack strategy to modify the graph node features as well as graph edge structure. To ensure the imperceptibility of adversarial perturbations, they designed constraints based on power-law [42] to preserve the degree distribution of graph structures and feature statistics.

Being atypical graph data, human skeletons have several unique properties. In a human skeleton, the graph edges represent rigid human bones, which connect finite number of human joints to form a standard spatial configuration. Unlike graph data with mutable graph structure (e.g. social network graph [43]), the human bones are fixed in terms of both joint connections and bone lengths. This property implies that attacking human skeletons by adding or deleting bones will be detected easily by observers. The hierarchical nature of human skeleton data is also different from normal graph data, as in human skeleton the motion of children joints/bones are affected by their parents' behaviours. This chain-like motion kinetics of human skeletons must be considered when launching adversarial attacks on skeleton actions. Hence, despite the existence of adversarial attacks on graph data, robustness of skeleton based human action recognition against adversarial attacks remains largely unexplored.

In this work, we specifically focus on adversarial attacks on human skeleton sequences to fool skeleton-based action recognition models. To design effective and meaningful attacks, we take the spatial and temporal attributes of skeleton data into account while creating the adversarial perturbations. Due to its wide-spread use in graph convolution network based action recognition, we select ST-GCN [8] as our target model, and launch our attack against it. However, our attack is generic for similar graph based model. In the section to follow, we formulate our problem in the context of skeleton based human action recognition.

## III. PROBLEM FORMULATION

To formulate the problem, we first briefly revisit the spatio-temporal graph convolutional network ST-GCN [8] for skeleton-based action recognition. Using this prerequisite knowledge, we subsequently formalize our problem of adversarial attacks on skeleton action recognition.

### A. Revisiting ST-GCN

An action in skeleton domain is represented as a sequence of  $T$  skeleton frames, where every skeleton consists of  $N$  body joints. Given such  $N \times T$  volumes of joints, an undirected spatio-temporal graph  $G = (V, E)$  can be constructed, where  $V$  denotes the node set of graph and  $E$  is the edge set. Here,  $V = \{v_{ti} | t = 1, \dots, T, i = 1, \dots, N\}$  encodes the skeleton joints. An element ' $v$ ' of this set can also be considered to encode a joint's Cartesian coordinates. Two kinds of graph edges  $E$  are defined for joints, namely; intra-body edge  $E^S$  and inter-frame edge  $E^F$ . Specifically,  $E^S$  is represented as an  $N \times N$  adjacency matrix of graph nodes, where the matrix element  $E_{ij}^S = 1 | i \neq j$  identifies that a physical bone connection exists between the body joint  $v_i$  and  $v_j$ . The inter-frame edges  $E^F$  denotes the connections of the same joints

between consecutive frames, which can also be treated as temporal trajectories of the skeleton joints.

Given the spatio-temporal skeleton graph  $G$ , a graph convolution operation is defined by extending the conventional image-based convolution. Along the spatial dimension, graph convolution is conducted on a graph node  $v_i$  around its neighboring nodes  $v_j \in B(v_i)$ :

$$f_{out}(v_i) = \sum_{v_j \in B(v_i)} \frac{1}{Z_i(v_j)} f_{in}(v_j) \cdot w(l_i(v_j)), \quad (1)$$

where  $B$  is the sampling function to define a neighboring node set for the joint  $v_i$ ,  $f_{in}$  is the input feature map,  $w$  is the weight function which indexes convolution weight vectors based on the labels of neighboring nodes  $v_j$ , and  $Z_i(v_j)$  is the number of neighboring nodes to normalize the inner product. The labels of neighboring nodes are assigned with a labelling function  $l : B(v_i) \rightarrow \{0, \dots, K-1\}$ , where  $K$  defines the spatial kernel size. ST-GCN employs different skeleton partitioning strategies for the labelling purpose. To conduct graph convolution in spatio-temporal dimensions, the sampling function  $B(v)$  and the labelling function  $l(v)$  are extended to cover a pre-defined temporal range  $\Gamma$ , which decides the temporal kernel size.

ST-GCN [8] adopts the implementation of graph convolution network in [22] to create a 9-layer neural network with temporal kernel size  $\Gamma = 9$  for each layer. Starting from 64, the number of channels is doubled for every 3 layers. The resulting tensor is pooled at the last layer to produce a feature vector  $f_{final} \in \mathbb{R}^{256}$ , which is fed to a Softmax classifier for predicting the action label. The network mapping function is compactly represented as:

$$Z_{G,c} = \mathcal{F}_\theta(V, E) = \arg \max < \text{softmax}(f_{final}) >, \quad (2)$$

where  $\theta$  denotes the network parameters. We use  $Z_{G,c}$  to denote the probability of assigning spatio-temporal skeleton graph  $G$  to class  $c \in C = \{1, 2, \dots, c_k\}$ . After training, the network parameters are fine-tuned to minimize the cross entropy loss between the predicted class  $c$  and the ground truth  $c_{gt}$  that maximizes the probability  $Z_{G,c}|c = c_{gt}$  for the dataset under consideration.

### B. Adversarial Attack on Skeleton Action Recognition

Given an original spatio-temporal skeleton graph  $G^0 = (V^0, E^0)$ , and a trained ST-GCN model  $\mathcal{F}_\theta$ , our goal is to apply adversarial perturbation to the graph  $G^0$ , resulting in a perturbed graph  $G' = (V', E')$  that satisfies the following broad constraint:

$$Z_{G',c} = \mathcal{F}_\theta(V', E'), \text{ s.t. } c \neq c_{gt} \quad (3)$$

Below, we examine this objective from various aspects to compute effective adversarial perturbations for the skeleton action recognition.

1) *Feature and structure perturbations* : As explained in Section III-A,  $V$  denotes the skeleton joints whose elements can be represented as the Cartesian coordinates of joints, e.g.  $v_{ti} : \{x_{ti}, y_{ti}, z_{ti}\}$ . For a particular node  $v_{ti}$  in the skeleton graph  $G$ , an adversarial attack can change its original location

such that  $v'_{ti} = v_{ti}^0 + \rho_{ti}$ , where  $\rho_{ti} \in \mathbb{R}^3$  is the adversarial perturbation for the node  $v_{ti}$ . We refer to this type of perturbation as *feature perturbation*. Alternatively, one can define *structure perturbation* that aims at changing the adjacency relationship in a graph such that  $E'_{ij} \neq E_{ij}^0 | i, j \in \mathcal{V}$ , where  $\mathcal{V}$  denotes the set of affected graph nodes.

In a spatio-temporal skeleton graph  $G$ , perturbing edges have strong physical implications. Recall that intra-body connections of joints define the rigid bones within a skeleton, and inter-frame connections define the temporal movements of the joints. Changes to these connections can lead to skeleton sequences that cannot be interpreted as any meaningful human action. Hence, the objective in Eq. 3 must further be constrained to preserve the graph structure while computing the perturbation. To account for that, we must modify the overall constraint to:

$$Z_{G',c} = \mathcal{F}_\theta(V', E^0), \text{ s.t. } c \neq c_{gt}. \quad (4)$$

2) *Perturbation imperceptibility*: Imperceptibility is an important attribute of adversarial attacks, as adversaries are likely to fool deep models in unnoticeable ways. Here, we explore perturbation imperceptibility in the context of skeleton actions. This leads to further constraints that must be satisfied when launching adversarial attacks on a skeleton graph  $G$ .

For the conventional image data, imperceptibility of perturbations is typically achieved by restricting  $\|\rho\|_p < \xi$ , where  $\|\cdot\|_p$  denotes the  $\ell_p$ -norm of a vector with  $p \in [0, \infty)$ , and  $\xi$  is a pre-defined constant [44]. For the skeleton graph data, however, the graph structure is discrete and graph nodes are dependent on each other, which makes it more challenging to keep a valid perturbation fully imperceptible. We tackle the challenge of perceptibility for skeleton perturbations from multiple point of views that results in multiple constraints for the overall objective, as explained in the following paragraphs.

### Joints variation constraint:

Focusing on *feature perturbation* on skeleton graph, the location of a target skeleton joint is changed such that  $v'_{ti} = v_{ti}^0 + \rho_{ti}$ . It is intuitive to constrain  $\rho$  of every target joint in a small range to avoid breaking the spatial integrity of the skeleton. Hence, we employ the following constraint:

$$\|\rho_{ti}\|_\infty \leq \epsilon_i \quad | \quad t \in [1, \dots, T]; i \in [1, \dots, N], \quad (5)$$

where  $\|\cdot\|_\infty$  denotes  $\ell_\infty$ -norm, and  $\epsilon_i$  is a pre-fixed constant. By restricting the joint variations to a small  $\ell_\infty$ -ball, we encourage perturbation imperceptibility. From the implementation view point, when the ball radius  $\epsilon$  is constant for all joints, we call it *global clipping* of the perturbed joints, and when the value of  $\epsilon_i$  is joint-dependent, we call it *hierarchical clipping*.

### Bone length constraint:

In a skeleton graph  $G$ , the intra-body graph connections  $E^S$  represent rigid human bones, hence their lengths must be preserved despite the perturbations. In the case of  $E_{ij}^S = 1 | i \neq j$ , the length of the bone between joint  $i$  and  $j$  at frame  $t$  can be calculated as  $B_{ij,t} = \|v_{ti} - v_{tj}\|_2$ . After applying perturbations to the graph, the new bone length  $B'_{ij,t} = \|v'_{ti} - v'_{tj}\|_2$  should

satisfy the following:

$$B_{ij,t} = B'_{ij,t} \mid t \in [1, \dots, T] \text{ s.t. } E_{ij}^S = 1. \quad (6)$$

**Temporal dynamics constraint:** Due to the spatio-temporal nature of skeleton action graphs, we disentangle the restrictions over perturbations into spatial and temporal constraints. Previous paragraphs mainly focused on the spatial constraints. Here, we analyze the problem from a temporal perspective.

A skeleton action is a sequence of skeleton frames that transit smoothly along the temporal dimension. A skeleton perturbation may lead to random jitters in the temporal trajectories of the joints and compromise the smooth temporal dynamics of the target skeleton action. To address this problem, we impose an explicit temporal constraint over the perturbations. Inspired by [26], we penalize acceleration of the perturbed joints to enforce temporal stability. Given consecutive perturbed skeleton frames  $f'_{t-1}, f'_t$ , and  $f'_{t+1}$ , the acceleration is calculated as  $\ddot{f}'_t = f'_{t+1} + f'_{t-1} - 2f'_t$ . Note that,  $f'_t = \{v'_{ti} \mid i = 1, \dots, N\}$ , where  $N$  is the number of perturbed skeleton joints. The calculation of acceleration is conducted on individual joints. We optimize our attacker,  $\mathcal{A}$  (discussed further below) by including the following temporal smoothness loss in the overall objective:

$$\mathcal{L}_{smooth}(\mathcal{A}) = \frac{1}{T-1} \sum_{t=2}^T \ddot{f}'_t = \frac{1}{T-1} \sum_{t=2}^T \sum_{i=1}^N v''_{ti}, \quad (7)$$

where  $T$  denotes the number of time steps considered. In the text to follow, we use  $\ddot{f}'_t$  to denote the joint acceleration for notational simplification.

3) *Anthropomorphic plausibility:* After adversarial perturbation is applied to a skeleton, the resulting skeleton can become anthropomorphically implausible. For instance, the perturbed arms and legs may bend unnaturally, or significant self-intersections may occur within the perturbed armature. Such unnatural behaviour can easily raise attack suspicions. Therefore, this potential behavior needs to be regularized while computing the perturbations.

Let  $\mathcal{P}$  define the distribution of natural skeleton graphs. A sample graph  $G^0$  is drawn from this distribution with probability  $\mathcal{P}(G^0)$ . We can treat an adversarial skeleton's graph  $G'$  to be a sample of another similar distribution  $\mathcal{P}'$ . The latter distribution should closely resemble the former under the restriction of minimal perturbation of joints and anthropomorphic plausibility of the skeletons. Hence, to obtain effective adversarial skeletons we aim at reducing the distribution gap between  $\mathcal{P}$  and  $\mathcal{P}'$ . To that end, we employ a Generative Adversarial Network (GAN) [24] to learn appropriate distribution in a data-driven manner.

Specifically, we model a skeleton action ‘attacker’ as a function  $\mathcal{A}$  such that  $G' = \mathcal{A}(G^0)$ . In the common GAN setup, the attacker can be interpreted as a generator of perturbed skeletons (see Fig. 1). We set up a binary classification network as the discriminator  $\mathcal{D}$ . The discriminator accepts either the natural graph  $\tilde{G}$  or the perturbed graph  $G'$  as its input, and predicts the probability that the input graph came from  $\mathcal{P}$ . The  $\tilde{G}$  and  $G'$  are kept ‘unpaired’, implying  $\tilde{G}$  and  $G^0$  are different graphs sampled from the distribution  $\mathcal{P}$ . To formulate

the adversarial learning process, we leverage the least squares objective [45] to train the attacker  $\mathcal{A}$  and the discriminator  $\mathcal{D}$  using the following loss functions:

$$\mathcal{L}_{adv}(\mathcal{A}) = \mathbb{E}_{G' \sim \mathcal{P}'}[(\mathcal{D}(G') - 1)^2], \quad (8)$$

$$\mathcal{L}_{adv}(\mathcal{D}) = \mathbb{E}_{\tilde{G} \sim \mathcal{P}}[(\mathcal{D}(\tilde{G}) - 1)^2] + \mathbb{E}_{G' \sim \mathcal{P}'}[\mathcal{D}(G')^2]. \quad (9)$$

During training,  $\mathcal{A}$  and  $\mathcal{D}$  are optimized jointly. We discuss the related implementation details in Section IV.

4) *Localized joint perturbation:* Unlike the pixel space of images, a skeleton action graph has highly discrete structure along both spatial and temporal dimensions. This discreteness poses unconventional challenges for adversarial attacks in this domain. Nevertheless, it also gives rise to interesting investigation directions. For instance, it is intriguing to devise a *localized* adversarial attack which fools the model by perturbing only a particular part of the skeleton graph. If we closely observe a skeleton action, it is clear that different body joints contribute differently to our perception of actions. Additionally, most of the human actions are recognizable by the motion patterns associated with the dominant body parts, e.g. arms and legs. Such observations make localized perturbations particularly relevant to the skeleton data.

Localized joint perturbations allow for less variations in the overall skeleton, which is beneficial for imperceptibility. They also provide a controlled injection of regional modification to the target skeleton action. To allow that, we define a subset of joints within a skeleton as the attack region. Only the joints in that region are modified for localized perturbations. Consequently, all the constraints in Section III-B2 still hold for the attack.

## IV. ATTACKER IMPLEMENTATION

### A. One-Step Attack

First, we adopt the Fast Gradient Sign Method (FGSM) [16] as a primitive attack to create skeleton perturbation  $V'$  in a single step. This adoption allows us to put our attack in a better context for the active community in the direction of adversarial attacks. For the FGSM based attack in our setup, the perturbation computation can be expressed as:

$$V' = V^0 + \epsilon \text{sign}(\nabla_{V^0} \mathcal{L}(\mathcal{F}_\theta(V^0, E^0), c_{gt})) \quad (10)$$

where  $\mathcal{F}_\theta$  denotes trained ST-GCN [8] model,  $\mathcal{L}$  is the cross-entropy loss for action recognition, and  $\nabla_{V^0}$  is the derivative operation that computes the gradient of ST-GCN loss w.r.t.  $V^0$ , given the current model parameters  $\theta$  and the ground truth action label  $c_{gt}$ . The sign of gradient is scaled with a parameter  $\epsilon$ , and added to the original graph  $V^0$ . The FGSM-based attack is computationally efficient as it takes only a single step in the direction of increasing the recognition loss of the target model.

The basic FGSM attack does not specify the label for the misclassified action, and therefore is a ‘non-targeted’ attack. If we specify a particular label for  $c_{gt}$  in Eq. 10, and subtract the gradient’s sign from the original graph  $V^0$  (instead of adding it, as in Eq. 10) the resulting attack becomes a targeted attack [46] that is likely to change the predicted label of the considered action to a pre-specified label.

## B. Iterative Attack

The FGSM attack takes a single step over the model cost surface to increase the loss for the given input. An intuitive extension of this notion is to iteratively take multiple steps while adjusting the step direction [47]. For the iterative attack, we also adopt the same technique for the skeleton graph input. However, here we focus on targeted attacks. This is because (a) targeted attacks are more interesting for the real-world applications, and (b) non-targeted attacks can essentially be considered a degenerate case of the targeted attack, where the target label is chosen at random. Hence, an effective targeted attack already ensures non-targeted model fooling. To implement, we specify the desired target class and take multiple steps in the direction of decreasing the prediction loss of the model for the target class.

We implement the iterative targeted attack while enforcing the constraints discussed in Section III-B2. The resulting algorithm is termed as Constrained Iterative Attack for Skeleton Actions (CIASA). At the core of CIASA is the following iterative process:

$$V'_0 = V^0; V'_{N+1} = \mathcal{C}(V'_N - \alpha (\nabla_{V'_N} \mathcal{L}_{\text{CIASA}}(V'_N, c_{\text{target}}))), \quad (11)$$

At each iteration,  $V'_N$  is adjusted towards the direction of minimizing the overall CIASA loss  $\mathcal{L}_{\text{CIASA}}$  using a step size  $\alpha$ . This is equivalent to a gradient descent iteration with  $\alpha$  as the learning rate, where the skeleton graph  $V'_N$  is treated as the model parameter. Hence, we directly exploit the Adam Optimizer [48] in the PyTorch library<sup>1</sup> for this computation. The operation  $\mathcal{C}(\cdot)$  in Eq. 11 truncates and realigns the values in its argument with pre-set conditions, explained below.

In Eq. 11, the overall CIASA loss  $\mathcal{L}_{\text{CIASA}}$  consists of the following components:

$$\mathcal{L}_{\text{CIASA}} = \mathcal{L}_{\text{pred}} + \lambda(\mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{adv}}) \quad (12)$$

where  $\mathcal{L}_{\text{pred}}$  is the cross-entropy loss of the model prediction on  $V'$  for the desired target class  $c_{\text{target}}$ .  $\mathcal{L}_{\text{smooth}}$  is the temporal smoothness loss calculated according to Eq. 7. GAN regularization loss  $\mathcal{L}_{\text{adv}}$  is a combination of  $\mathcal{L}_{\text{adv}}(\mathcal{A})$  and  $\mathcal{L}_{\text{adv}}(\mathcal{D})$  given in Eq. 8 and Eq. 9.  $\lambda$  is a weighting hyperparameter to balance the individual loss components.

Implementing the process identified by Eq. 11 produces the perturbed skeleton  $V'$  that fools the model into misclassifying the original action as  $c_{\text{target}}$ , while complying to the spatio-temporal constraints derived in the previous Sections. The pseudo-code of implementing the process of Eq. 11 as CIASA is presented in Algorithm 1. The algorithm starts with a forward-pass of  $V'$  through the target model  $\mathcal{F}_\theta(\cdot)$ , i.e. ST-GCN. The respective losses are then computed to form the overall CIASA loss  $\mathcal{L}_{\text{CIASA}}$ . At line 6,  $\mathcal{L}_{\text{adv}}$  is computed as the accumulation of the losses defined in Eq. 8 and Eq. 9. Here, we replace  $G$  with  $V$  based on the algorithm context.  $\mathcal{D}_\omega$  denotes the discriminator network which is parameterized by  $\omega$ . Note that, the real data  $\tilde{V}$  and the perturbed data  $V'$  are unpaired, as discussed in Sect. III-B3. On line 8, the gradient information is obtained through the back propagation

---

**Algorithm 1** Constrained iterative attacker  $\mathcal{A}$  to fool skeleton-base action recognition.

---

**Input:** Original graph nodes  $V^0 \in \mathbb{R}^{3 \times N \times T}$ , trained ST-GCN model  $\mathcal{F}_\theta(\cdot)$ , desired target class  $c_{\text{target}}$ , perturbation clipping factor  $\epsilon$ , max\_iter= $M$ , learning rate  $\alpha$

**Output:** Perturbed graph nodes  $V' \in \mathbb{R}^{3 \times N \times T}$ .

```

1: set initial  $V' = V^0$ 
2: while  $i < M$  do
3:   feed forward  $Z = \mathcal{F}_\theta(V')$ 
4:    $\mathcal{L}_{\text{pred}} = \text{CrossEntropyLoss}(Z, c_{\text{target}})$ 
5:    $\mathcal{L}_{\text{smooth}} = \frac{1}{T-1} \sum_{t=2}^T \ddot{f}'_t$ 
6:    $\mathcal{L}_{\text{adv}} = (\mathcal{D}_\omega(V') - 1)^2 + (\mathcal{D}_\omega(\tilde{V}) - 1)^2 + \mathcal{D}_\omega(V')^2$ 
7:    $\mathcal{L}_{\text{CIASA}} = \mathcal{L}_{\text{pred}} + \lambda(\mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{adv}})$ 
8:    $(\mathcal{L}_{\text{CIASA}}).\text{Backward}() \Rightarrow \text{gradients}$ 
9:    $V', \omega = \text{AdamOptimizer}([V', \omega], \text{gradients})$ 
10:  if  $|V' - V^0| > \epsilon$  then
11:     $V' = \text{Clip}(V') \sim [V^0 - \epsilon, V^0 + \epsilon]$ 
12:  end if
13:  Skeleton realignment  $V' = \text{SSR}(V')$ 
14:   $i = i + 1$ 
15: end while
16: return  $V'$ 

```

---

operation denoted as ‘.Backward()’. We employ the Adam Optimizer [48] to update the skeleton joints  $V'$  and the discriminator parameters  $\omega$ . Clipping operation is then applied to truncate  $V'$  to pre-set ranges. In our case, the scaling factor  $\epsilon$  restricts the  $\ell_\infty$ -norm of the perturbation at graph nodes. For global clipping,  $\epsilon \in \mathbb{R}$  is a scalar value that results in equal clipping on all joints. For the hierarchical clipping,  $\epsilon \in \mathbb{R}^N$  defines different clipping strengths for different joint. The clipping imposes the *joint variation constraint* over the perturbations. To impose the *bone length constraint*, Spatial Skeleton Realignment (SSR) is proposed to realign the skeleton bones within the clipped  $V'$  according to the original bone lengths. Note that the operations of clipping and realignment constitute the function  $\mathcal{C}(\cdot)$  shown in Eq. 11. We empirically set the weight factor  $\lambda$  as 10, and the base learning rate for the Adam Optimizer  $\alpha$  as 0.01. Below we discuss the implementation of SSR and discriminator network  $\mathcal{D}$ .

1) *Spatial Skeleton Realignment:* We propose Spatial Skeleton Realignment (SSR) to preserve the *bone length constraint* as we perturb the skeleton graph. SSR is executed at each iteration after  $V'$  is updated and clipped in order to realign every perturbed skeleton frame based on the original bone lengths. Specifically, for every updated skeleton joint  $v'_j$ , we find its parent joint  $v'_i$  along the intra-body edge  $E^S$ . The bone between joints  $i$  and  $j$  is defined as a vector  $b'_{ij} = v'_j - v'_i$ . Then, we modify the joint  $v'_j$  along the vector direction  $b'_{ij}$  to meet the constraint in Eq 6. The modification applied to  $v'_j$  is also applied to all of its children/grandchildren joints. To complete the SSR, the above process starts from the root joint and propagates through the whole skeleton.

<sup>1</sup><https://pytorch.org/>

2) *GAN Regularization*: To enforce the anthropomorphic plausibility of the perturbed skeleton action, the adversarial regularization term  $\mathcal{L}_{adv}$  is optimized jointly with the other attack objectives. Taking per-frame skeleton feature map, say  $X$  as the input, a discriminator network  $\mathcal{D}$  is trained to classify the skeleton as *fake* or *real* (i.e. perturbed v.s original), while the attacker  $\mathcal{A}$  is competing with  $\mathcal{D}$  to increase the probability of the perturbed skeleton being classified as *real*.

We leverage the angles between skeleton bones to construct the feature map  $X$ . For a pair of bones  $b_{ij}$  and  $b_{uv}$ , the corresponding element in the feature map is defined as the cosine distance between the bones as:

$$x_{ij-uv} = \frac{b_{ij} \cdot b_{uv}}{\|b_{ij}\| \|b_{uv}\|} \quad (13)$$

We select a group of major bones to construct the feature map  $X$ , while insignificant bones of fingers and toes are excluded to avoid unnecessary noise. The resulting feature map has dimension  $X \in \mathbb{R}^{C,H,W}$ , where  $C = 1$ , and  $H = W$  equals to the number of selected bones. We model  $\mathcal{D}$  as a binary classification network that consists of two convolution layers and one fully-connected layer. The convolution kernel size is 3, and the number of channels produced by the convolution is 32.  $\mathcal{D}$  outputs values in the range  $[0, 1]$ , signifying the probability that  $X$  is a real sample.

## V. EXPERIMENTS

Below we evaluate the effectiveness of the proposed attack for skeleton-based action recognition. We examine different attack modes on standard skeleton action datasets. We also demonstrate the transferability of attack and explore generalization of the computed adversarial perturbations beyond the skeleton data modality. Lastly, an ablation study is provided to highlight the contributions of various constraints to the overall fooling rate achieved by the proposed attack.

### A. Dataset and Evaluation Metric

**NTU RGB+D**: NTU RGB+D Human Activity Dataset is collected with Kinect v2 camera and includes 56,880 action samples. Each action has RGB, depth, skeleton and infra-red data associated with it. However, we are only concerned with the skeleton data in this work. For the skeleton-based action recognition with ST-GCN, we follow the standard protocols defined in [11], i.e. cross-subject and cross-view recognition. Accordingly, two different ST-GCN models are used in our experiments, one for each protocol. We denote these models as  $\text{NTU}_{XS}$  and  $\text{NTU}_{XV}$  for cross-subject and cross-view recognition. While the original dataset is split into training and testing sets, we only manipulate the testing set, as no separate training data is required for the attack.

**Kinetics**: Kinetics dataset [25] is a large unconstrained action dataset with 400 action classes. For skeleton-based action recognition using this data, the original ST-GCN [8] first uses OpenPose [49], [50] to estimate 2D skeletons with 18 body joints. Then, the estimation confidence ‘c’ for every joint is concatenated to its 2D coordinates (x, y) to form a tuple (x,y,c). The tuples for all joints in a skeleton are collectively

considered as an input sample by the ST-GCN model. For the adversarial attack, we mask the channel of confidence values and only perturb the (x,y) components for the Kinetics dataset.

**Evaluation metric**: The evaluation metric used to evaluate the success of adversarial attacks is known as *fooling rate* [36]. It indicates the percentage of data samples over which the model changes its predicted label after the samples have been adversarially perturbed. In the adversarial attacks literature, this is the most commonly used metric to evaluate an attack’s performance [36]. In the case of targeted attacks, it determines the percentage of the samples successfully misclassified as the target label after the attack.

### B. Non-targeted Attack

Since this is the first work in the direction of attacking skeleton-based action recognition, it is important to put our attacking technique into perspective. Hence, we first conduct a simpler non-targeted attack on NTU and Kinetics datasets using the one-step attack discussed in Section IV-A, Eq. (10). We compute the fooling rates for both datasets under different values of the perturbation scaling factor  $\epsilon$ . Both cross-view and cross-subject protocols were considered in this experiment for the NTU dataset. The fooling rates achieved with the one-step method for various  $\epsilon$  values are summarized in Fig. 2. As can be seen, the non-targeted fooling is reasonably successful under the proposed formulation of the problem for skeleton-based action recognition. The fooling rates for all protocols remain higher than 90% once the  $\epsilon$  value reaches 0.02. This is still a reasonably small perturbation value that is equivalent to one twentieth of the average skeleton height.

To visualize perturbed skeletons, Fig. 3(a) shows a successful attack on NTU dataset for cross-view fooling. The original and perturbed skeletons are plotted with green and red colors respectively. Note that, in this illustration and the examples to follow, we provide a positional offset between different skeletons for better visualization. For the shown sequence of skeleton frames, the original label is ‘Brush hair’, that is predicted as ‘Wipe face’ after the attack is performed. The temporal dimension evolves from left to right. Ignoring the positional offset, it is easy to see that the perturbation generally remains hard to perceive in the skeleton.

### C. Targeted Attack

We use the proposed CIASA attacker explained in Section IV-B, Alg. 1 to conduct targeted attacks on both NTU and Kinetics datasets. We specify the least-likely action prediction of the ST-GCN models as the target label  $c_{target}$  as described in Eq. 11, implying that the most challenging misclassification target is chosen to launch attacks. CIASA is configured to launch attacks in three modes; namely, *basic* mode, *localized* mode, and *advanced* mode. Below we discuss these modes along the experimental results.

Figure 3(b) shows an example of CIASA attack in the *basic* mode. We apply the global clipping discussed in Section III-B2 in this attack mode, where all the skeleton joints are perturbed with the same scaling factor  $\epsilon = 0.02$ . With this setting, the original action of ‘Cheer up’ in Figure 3(b) is misinterpreted

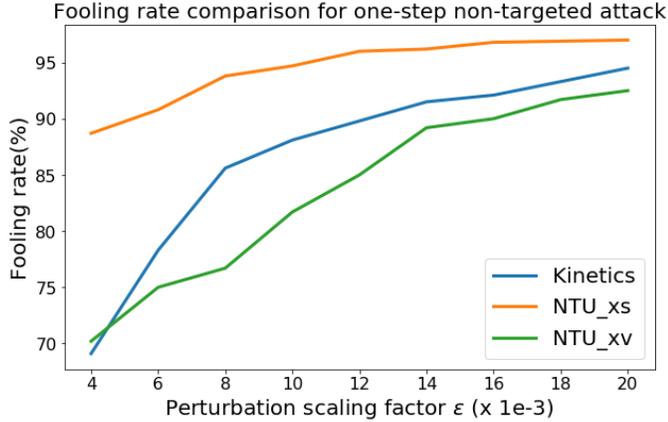


Fig. 2. Fooling rates (%) achieved by one-step non-targeted attack with different perturbation scaling factors for NTU and Kinetics datasets. Both cross-subject  $NTU_{XS}$  and cross-view  $NTU_{XV}$  protocols are considered for the NTU dataset.

TABLE I  
FOOLING RATES (%) ACHIEVED BY CIASA TARGETED ATTACK (BASIC MODE) WITH DIFFERENT GLOBAL CLIPPING STRENGTH  $\epsilon$  FOR NTU AND KINETICS DATASETS. BOTH CROSS-SUBJECT  $NTU_{XS}$  AND CROSS-VIEW  $NTU_{XV}$  PROTOCOLS ARE CONSIDERED FOR THE NTU DATASET.

$\epsilon$ ( $\times 1e-3$ )	4	6	8	10	12
Kinetics	82.5	92.5	96.5	97.5	99.3
$NTU_{XS}$	89.4	96.6	98.7	99.2	99.8
$NTU_{XV}$	78.2	85.5	93.3	98.9	99.6

as ‘Kicking’ with confidence score 99.4%. In the basic mode, the comparison of fooling rates with different  $\epsilon$  values for the two benchmark datasets are summarized in Table I. Firstly, the results demonstrate successful fooling even for very low  $\epsilon$  values. Secondly, it is noteworthy that for similar  $\epsilon$  values, higher fooling rates are generally achieved by CIASA for targeted fooling as compared to the non-targeted fooling of the one step method in Fig. 2. This demonstrates the strength of CIASA as a targeted attack. In our experiments, we observed that the least-likely label of ST-GCN model remains similar for multiple actions. Whereas the presented results do not diversify the target labels of such actions to strictly follow the evaluation protocol, it is possible to manually do so. Loosening the evaluation criterion on these lines will further improve the fooling rate of CIASA.

In Fig. 3(c), we show an example of CIASA attack in the *localized mode*, where the localized joint perturbation discussed in Section III-B4 is applied. In this example, two legs of skeleton are set to be the attack regions, which allow 8 active joints for perturbations. The remaining joints are unaffected by the computed perturbations. Compared to the basic mode, fewer joints contribute to the overall perturbation in the localized mode. To compensate for the reduced number of active joints, we loose the perturbation scaling factor and set  $\epsilon$  to 0.08 for this experiment. For the shown example, CIASA achieves fooling with 93.2% confidence for this mode, which is still competitive to the 99.4% confidence in the basic mode.

TABLE II  
FOOLING RATE(%) ACHIEVED BY CIASA TARGETED ATTACK (LOCALIZED MODE) WITH DIFFERENT ATTACK REGIONS ON NTU DATASET. BOTH CROSS-SUBJECT AND CROSS-VIEW PROTOCOLS ARE EVALUATED. GLOBAL CLIPPING STRENGTH IS SET TO  $\epsilon = 0.04$ .

Attack region	set-1	set-2	set-3	set-4
$NTU_{XS}$	90.8	93.3	61.3	83.3
$NTU_{XV}$	85.2	91.7	60.0	81.7

To further evaluate the localized mode of CIASA with different attack regions, we split the skeleton joints into 4 sets, as illustrated in Fig. 4. Then, we conduct CIASA localized attack on NTU dataset for the 4 sets separately. Global clipping is applied for these experiments with the scaling factor  $\epsilon = 0.04$ . The chosen value of  $\epsilon$  is intentionally kept lower than that in Fig. 3(c) because we focus on analysing the fooling prowess of different attack regions instead of simply achieving high fooling rates for all the regions. The results of our experiments are summarized in Table II. It is clear that the CIASA localized attack achieves impressive fooling rates by perturbing only a small set of joints within the skeleton. In addition, different sets of active joints affect the fooling performance differently. In Table II, set-1 and set-2 achieve higher fooling rates than the other two sets. This can be explained by the observation that many dominant movements in the NTU dataset occur at the upper part of human body.

We also extend the localized mode of CIASA to an *advanced mode* by replacing the global clipping by hierarchical clipping discussed in Section III-B2. In that case, the scalar clipping value  $\epsilon$  is replaced by  $\epsilon \in \mathbb{R}^N$ , where N is the number of active joints to be perturbed. Here, we allow various active joints to change in pre-defined ranges by using differentiated clipping values. One strategy to differentiate the clipping strength is applying incremental  $\epsilon$  variables from parent joints to their children joints, based on the observation that children joints normally move in larger ranges than their parents. Figure 3(d) illustrates an example of successful advanced attack on NTU dataset with two legs activated for the attack. The  $\epsilon$  variables are set to 0.01, 0.05, 0.15, 0.25 for the joint *hips*, *knees*, *ankles*, and *feet*, respectively. Note that we intentionally amplify the perturbation ranges at certain joints such as *ankles* and *feet*, which results in noticeable perturbations at the attack region. We will justify the intuition behind this differential clipping in the paragraphs to follow. For now, notice that in Fig. 3(d), the original label of ‘Cheer up’ is misclassified as ‘Kicking’ with a confidence score 96.1% with the advanced attack.

Although the CIASA attack in *advanced mode* apparently sacrifices the visual imperceptibility of the perturbation, it is able to maintain the “semantic imperceptibility” for the perturbed skeleton. We corroborate this claim with the following observations. First, in Fig. 3(d), the dominant body movements for ‘Cheer up’ action mainly occur in the upper part of the skeleton, while the fooling is conducted by perturbing the lower body to which less attention is paid for this action. Consequently, the attack does not incur significant

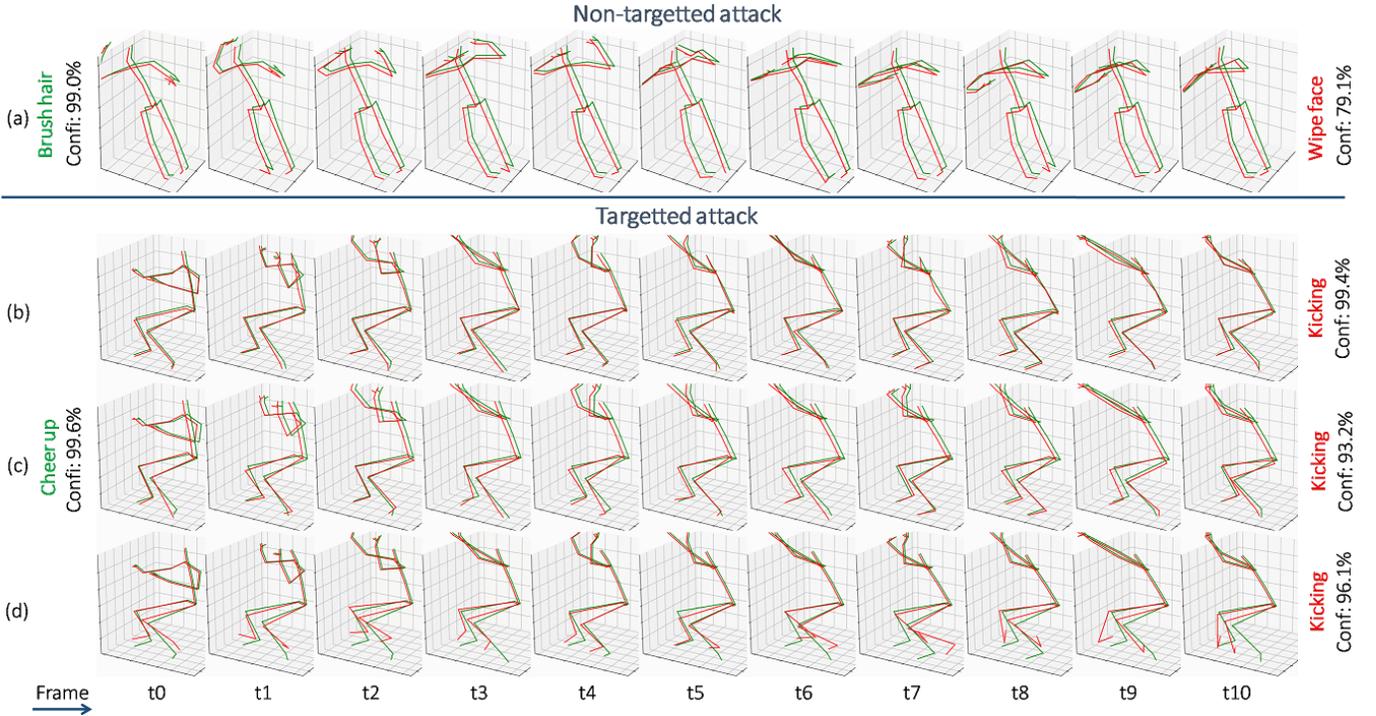


Fig. 3. TOP: (row a) One-step attack with  $\epsilon = 0.02$  is shown where “brush hair” action is misclassified as “wipe face”. BOTTOM: CIASA targeted attack in different modes are shown. (row b) The *basic* mode that perturbs all joints with  $\epsilon = 0.01$ . (row c) The *localized* mode with only two legs allowed to be perturbed. Global clipping is applied with  $\epsilon = 0.08$ . (d) The *advanced* mode where the same two legs are perturbed with hierarchical clipping. The attacks in all modes successfully fool the recognition model with confidences higher than 90%. The temporal dimension evolves from left to right.

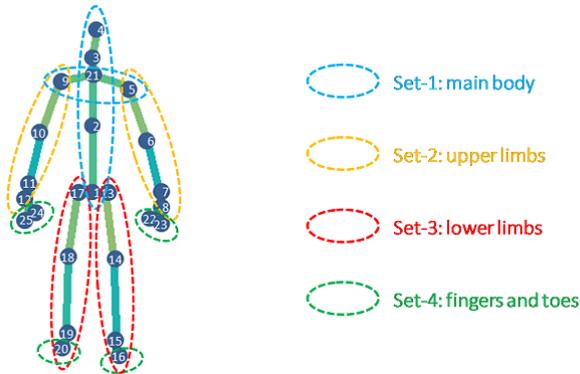


Fig. 4. The skeleton of NTU dataset is spitted into 4 attack regions, each of which is activated to apply CIASA localized attacks. Every attack region consists of roughly the same number of joints.

perceptual attention in the first place. Furthermore, due to the spatio-temporal constraints with CIASA attacks, the injected perturbation patterns remain smooth and natural. This further reduces the attack suspicions, as compared to any small but unnatural perturbations, e.g. shakiness around the joints.

Further to the above discussion, the perturbations generated in the advanced mode can not only fool the recognition model in skeleton spaces, but can also be imitated and reproduced in the Physical world. Imagine an end-to-end skeleton-based action recognition system using a monocular camera as its input sensor. For that, RGB images taken from the Physical world are first converted to skeleton frames, which are then

passed through the skeleton-based action recognition model. For this typical pipeline, it may be inconvenient to interfere with the intermediate skeleton data for the attacking purpose. However, the adversarial perturbations can be injected into the input RGB data by performing an action in front of the camera while imitating the perturbation patterns with selective body parts. The advanced mode of CIASA allows the discovery of perturbation patterns for such attacks. This is elaborated further in Section V-D2 with relevant context.

#### D. Transferability of Attack

We examine the transferability of the proposed CIASA attack from two perspectives. First, we evaluate the cross-model transferability of the generated perturbations. Concretely, we attack a skeleton action recognition model A to generate perturbed skeletons. Then, we predict the label of the perturbed skeletons using model B and examine the fooling rate for model B. We respectively chose ST-GCN and 2s-AGCN [10] as model A and B in our experiments.

Second, we analyze the cross-modality transferability of CIASA attack. i.e. we generate perturbations for one data modality and test their fooling capability in another data modality. We formulate this task as transferring perturbations from skeleton data to RGB data, as RGB cameras are widely used as input sensors for the real world systems. For the cross-modality test, we generate perturbed skeletons by attacking the ST-GCN. Then, those skeletons are converted to RGB actions using a graphics rendering pipeline. To examine whether the

TABLE III  
COMPARISON OF CROSS-MODEL RECOGNITION ACCURACY (%) AND FOOLING RATE (%) ON THREE CONFIGURATIONS OF 2s-AGCN FOR CROSS-VIEW NTU PROTOCOL. ‘ORIGINAL ACCURACY’ IS ON CLEAN DATA. ‘ATTACKED ACCURACY’ IS ON PERTURBED DATA.

Model	Js-AGCN	Bs-AGCN	2s-AGCN
Original Accuracy	93.7	93.2	95.1
Attacked Accuracy	13.5	6.8	11.8
Fooling rate (%)	86.1	93.1	88.4

adversarial information can be preserved during the conversion, we predict the label of RGB actions under the usual skeleton-based action recognition pipeline for the ST-GCN.

1) *Cross-Model Transferability*: The 2s-AGCN [10] is a two-stream adaptive graph convolutional network for skeleton-based action recognition. This network is significantly different from the ST-GCN [8] as it models a learnable topology of the skeleton graph. In addition to the joint locations, 2s-AGCN also models the bone directions, which results in a two-stream network structure.

We first generate perturbed skeleton actions based on ST-GCN model. The *basic mode* of CIASA with global clipping is employed, where the perturbation scaling factor  $\epsilon$  is empirically set to 0.012. The cross-view protocol of NTU dataset is adopted to create perturbed skeletons, which are then evaluated by 2s-AGCN models. We compare the change of recognition accuracy before and after the attack, and record the fooling rates for three different configurations of the 2s-AGCN, i.e. joint only (Js-AGCN), bone only (Bs-AGCN) and ensemble (2s-AGCN). The results in Table III show that the perturbations generated with ST-GCN significantly degrades the recognition performance of 2s-AGCN. This demonstrate that the proposed CIASA attacker is able to generalize well on ‘unseen’ action recognition models.

2) *Cross-Modality Transferability*: To transfer the perturbations from skeleton to RGB space, we adopt a human pose synthesis technique [51] to create RGB actions based on the perturbed skeleton sequences generated with the advanced mode of CIASA. The adopted synthesis pipeline can produce realistic RGB actions with diversified human models, backgrounds, cloth textures and illuminations. Moreover, the temporal dynamics of the underlying action is also reproducible in the synthesized RGB video. We demonstrate successful cross-modality transferability in Fig. 5. The rows (a) and (d) are the original and perturbed skeleton sequences respectively. (b) and (e) show the RGB actions generated using [51] with (a) and (d) used as the inputs skeleton sequences.

Firstly, the successful generation of realistic RGB videos in (b) and (e) affirms that the skeleton perturbations generated by CIASA are useful in producing action perturbations in the Physical world beyond the skeleton space. Secondly, we observe that the adversarial information remains largely preserved during the cross-modality transfer. In Fig. 5, we use VNect [26] as a 3D pose extractor to recover 3D skeletons directly from the synthesized RGB actions. The recovered skeleton sequences are then fed to the trained ST-GCN model

for action recognition, mimicking the typical pipeline for the skeleton-based action recognition for RGB sensors.

The VNect-recovered 3D skeletons from clean and perturbed RGB data are respectively shown in rows (c) and (f) of the figure. As can be seen, the recovered skeletons generally follow the motion patterns encoded in the respective source skeletons. For the clean data, the recovered skeletons in (c) and the source skeletons in (a) are both correctly recognized as ‘Throw’ action. For the perturbed data, the recovered skeleton sequence in (f) has fooled the ST-GCN into misclassifying the action as ‘Back pain’. Although the fooling is not in the exact least likely class as in row (d), misclassification due to CIASA attack for this very challenging scenario is still intriguing. We note that the attack here is naturally degenerating into an untargeted attack.

To further scale up the cross-modality experiment, we randomly select 240 skeleton actions for the cross-view protocol of the NTU dataset. Then, we conduct the cross-modality transfer for all those sequences. We only use a subset of the NTU dataset because of the unreasonable computational time required to render videos for the complete dataset. Subsequently, we predict action labels with ST-GCN on the VNect-recovered skeleton sequences for both clean and perturbed data. With this setting, the recognition accuracy is recorded as 53.3% for the clean data, and 38.9% for the perturbed data. Compared to the original NTU cross-view accuracy of 88.3% [8], lower performance is observed on the clean data due to inaccurate 3D pose extraction by VNect. Nevertheless, the proposed attack is still able to preserve its adversarial characteristics to further cause a significant accuracy drop in this challenging scenario.

### E. Ablation Study

For the CIASA attack, we have proposed a set of spatio-temporal constraints to achieve high-quality adversarial perturbations in terms of both temporal coherence and spatial integrity of the perturbed skeletons. Here, we provide an ablation study to compare the contributions of these constraints in the overall results.

To enforce temporal smoothness in the perturbed skeleton sequences, we penalize the joint accelerations between the consecutive skeleton frames. Figure 6 compares the perturbed skeletons with and without this temporal constrains in the *basic mode* of CIASA, where the original and perturbed skeletons are highlighted with green and red color respectively. It is apparent that the perturbed skeletons in (b) move more smoothly than those in (a) along the temporal dimension. This ascertains the effectiveness of temporal smoothing in our attack. Both perturbations in (a) and (b) successfully fool the recognition model to misclassify “Drink action” as the “Jump up” action.

To enforce spatial integrity and anthropomorphic plausibility, we use spatial skeleton realignment (SSR) and GAN regularization. Such spatial constrains are particularly important for the CIASA localized attacks, where only a given subset of the joints is permitted to be changed. Figure 7 compares the perturbation results with and without the spatial constraints

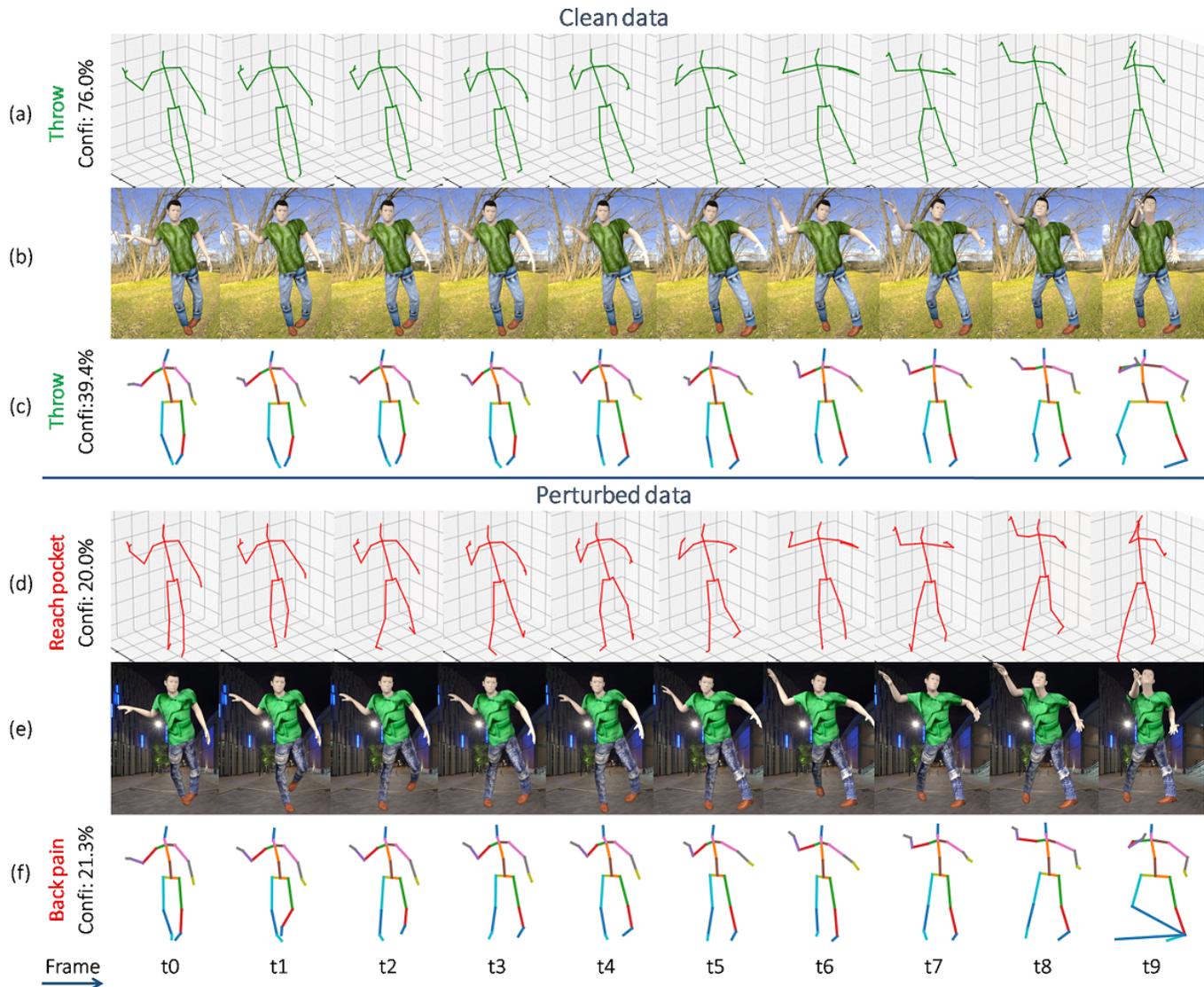


Fig. 5. TOP: Clean data of different modalities. (a) Original skeleton sequences. (b) RGB video rendered from the original sequence. (c) Recovered 3D pose sequence extracted from (b) using VNect [26]. BOTTOM: Perturbed data of different modalities. (d) Perturbed skeleton sequences created with the advanced mode of CIASA. (e) RGB video rendered from (d). (f) 3D poses extracted from (e) using VNect [26].

for a localized attack on skeleton legs. Without any spatial constrains, the perturbed skeletons in (a) shows unrealistic pose configurations and arbitrary lengths of bones. With only SSR enabled in (b), lengths of the perturbed bones are more consistent with their original values, however, the resulting poses are still not realistic in terms of plausibility. By adding the GAN regularization, the skeletons in (c) are more realistic. The skeleton sequences in the figure clearly demonstrates the effectiveness of SSR and GAN regularization in our attack. All sequences in Fig. 7 (a), (b) and (c) successfully fool the recognition model in predicting the label “Drink water” as “Jump up”.

## VI. CONCLUSION

We present the first systematic adversarial attack on skeleton-based action recognition. Unlike the existing attacks

that target non-sequential tasks, e.g. image classification, semantic segmentation and pose estimation, we attack deep sequential models from a spatio-temporal perspective. With skeleton-based action recognition model ST-GCN [8] as the target, we demonstrate its successful fooling by mainly perturbing the joint positions. The proposed attack algorithm CIASA imposes spatio-temporal constraints on the adversarial perturbations to produce perturbed skeleton sequences with temporal smoothness, spatial integrity, and anthropomorphic plausibility. The proposed algorithm works in different modes based on the needs of the attack. With the *localized* mode of CIASA, we are able to perturb only a particular set of the body joints to launch localized attack. Such attacks can be used to inject regional perturbations to pre-specified parts of the body, without interfering with the dominant action patterns that are performed by the other joints. Compared to the *basic* mode that perturbs all the joints with global

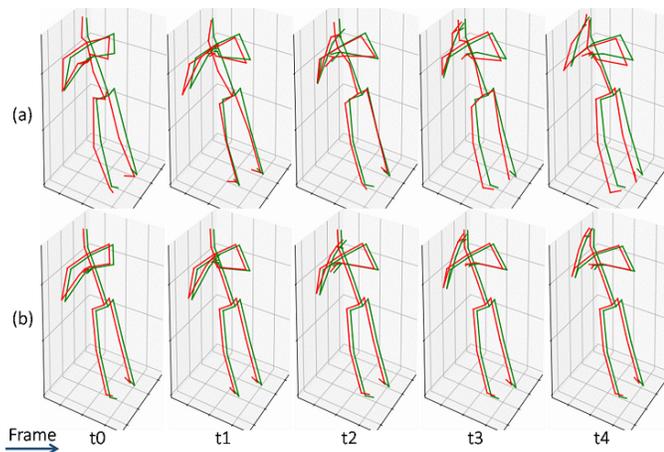


Fig. 6. Temporal smoothness in CIASA. (a) Perturbed skeleton sequence without temporal smoothness constrains. (b) Perturbed sequence with temporal smoothness constrains. The original and perturbed skeletons are shown in Green and Red colors respectively.

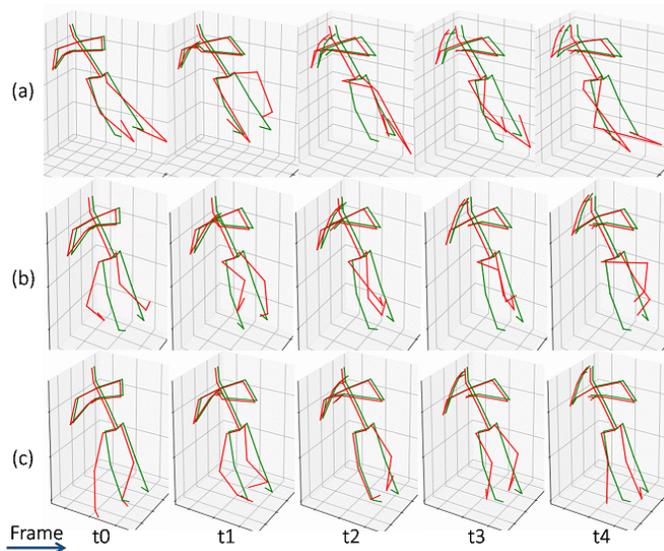


Fig. 7. Effectiveness of spatial constraints in CIASA. A localized attack is launched on two legs of the skeleton. (a) No spatial constraints: Pose configuration and bone lengths change randomly. (b) Spatial Skeleton Realignment (SSR): Constrained consistent bone lengths, but unnatural poses. (c) GAN regularization: Realistic poses that can correspond to the real-world skeleton motions.

clipping, an *advanced* mode utilizes localized attacks with hierarchical joint variations to disguises the attack intentions with realistic motion patterns. Our experiments show that the proposed CIASA perturbations generalize well across different recognition models. Moreover, they also have the ability to transfer to RGB video modality under graphics rendering pipeline. This indicates that CIASA generated perturbations can allow attackers to mimic semantically imperceptible adversarial patterns in the real world to fool skeleton based action recognition systems.

#### ACKNOWLEDGMENT

This research is supported by the Australian Research Council (ARC) grant DP190102443. The Tesla K-40 GPU

used for this research is donated by the NVIDIA Corporation.

#### REFERENCES

- [1] M. Zanfir, M. Leordeanu, and C. Sminchisescu, "The moving pose: An efficient 3d kinematics descriptor for low-latency action recognition and detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 2752–2759.
- [2] M. Devanne, H. Wannous, S. Berretti, P. Pala, M. Daoudi, and A. Del Bimbo, "3-d human action recognition by shape analysis of motion trajectories on riemannian manifold," *IEEE transactions on cybernetics*, vol. 45, no. 7, pp. 1340–1352, 2015.
- [3] Y. Du, Y. Fu, and L. Wang, "Skeleton based action recognition with convolutional neural network," in *Pattern Recognition (ACPR), 2015 3rd IAPR Asian Conference on*. IEEE, 2015, pp. 579–583.
- [4] R. Vemulapalli and R. Chellapa, "Rolling rotations for recognizing human actions from 3d skeletal data," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4471–4479.
- [5] Q. Ke, M. Bennamoun, S. An, F. Sohel, and F. Boussaid, "A new representation of skeleton sequences for 3d action recognition," *arXiv preprint arXiv:1703.03492*, 2017.
- [6] Y. Yang, C. Deng, D. Tao, S. Zhang, W. Liu, and X. Gao, "Latent max-margin multitask learning with skeletons for 3-d action recognition." *IEEE Trans. Cybernetics*, vol. 47, no. 2, pp. 439–448, 2017.
- [7] J. Liu, A. Shahroudy, D. Xu, A. C. Kot, and G. Wang, "Skeleton-based action recognition using spatio-temporal lstm network with trust gates," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 12, pp. 3007–3021, 2017.
- [8] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [9] C. Cao, Y. Zhang, C. Zhang, and H. Lu, "Body joint guided 3-d deep convolutional descriptors for action recognition," *IEEE transactions on cybernetics*, vol. 48, no. 3, pp. 1095–1108, 2018.
- [10] L. Shi, Y. Zhang, J. Cheng, and H. Lu, "Two-stream adaptive graph convolutional networks for skeleton-based action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12026–12035.
- [11] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "Ntu rgb+d: A large scale dataset for 3d human activity analysis," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1010–1019.
- [12] K. Yun, J. Honorio, D. Chattopadhyay, T. L. Berg, and D. Samaras, "Two-person interaction detection using body-pose features and multiple instance learning," in *2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 2012, pp. 28–35.
- [13] F. Offi, R. Chaudhry, G. Kurillo, R. Vidal, and R. Bajcsy, "Berkeley mhad: A comprehensive multimodal human action database," in *2013 IEEE Workshop on Applications of Computer Vision (WACV)*. IEEE, 2013, pp. 53–60.
- [14] J.-F. Hu, W.-S. Zheng, J. Lai, and J. Zhang, "Jointly learning heterogeneous features for rgb-d activity recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5344–5352.
- [15] 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.
- [16] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [17] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2016, pp. 372–387.
- [18] J. Su, D. V. Vargas, and K. Sakurai, "One pixel attack for fooling deep neural networks," *IEEE Transactions on Evolutionary Computation*, 2019.
- [19] V. Mirjalili and A. Ross, "Soft biometric privacy: Retaining biometric utility of face images while perturbing gender," in *2017 IEEE International joint conference on biometrics (IJCB)*. IEEE, 2017, pp. 564–573.
- [20] Y.-C. Lin, Z.-W. Hong, Y.-H. Liao, M.-L. Shih, M.-Y. Liu, and M. Sun, "Tactics of adversarial attack on deep reinforcement learning agents," *arXiv preprint arXiv:1703.06748*, 2017.
- [21] N. Akhtar, A. Jalwana, M. Bennamoun, and A. Mian, "Label universal targeted attack," *arXiv preprint arXiv:1905.11544*, 2019.
- [22] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.

- [23] D. Zügner, A. Akbarnejad, and S. Gunnemann, “Adversarial attacks on neural networks for graph data,” in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 2018, pp. 2847–2856.
- [24] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [25] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, “The kinetics human action video dataset,” *arXiv preprint arXiv:1705.06950*, 2017.
- [26] D. Mehta, S. Sridhar, O. Sotnychenko, H. Rhodin, M. Shafiei, H.-P. Seidel, W. Xu, D. Casas, and C. Theobalt, “Vnect: Real-time 3d human pose estimation with a single rgb camera,” *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, p. 44, 2017.
- [27] Z. Huang, C. Wan, T. Probst, and L. Van Gool, “Deep learning on lie groups for skeleton-based action recognition,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [28] J. Liu, N. Akhtar, and A. Mian, “Skepxels: Spatio-temporal image representation of human skeleton joints for action recognition,” *arXiv preprint arXiv:1711.05941*, 2017.
- [29] Y. Du, W. Wang, and L. Wang, “Hierarchical recurrent neural network for skeleton based action recognition,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1110–1118.
- [30] V. Veeriah, N. Zhuang, and G.-J. Qi, “Differential recurrent neural networks for action recognition,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4041–4049.
- [31] A. Shahroudy, T.-T. Ng, Q. Yang, and G. Wang, “Multimodal multipart learning for action recognition in depth videos,” *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 10, pp. 2123–2129, 2016.
- [32] Y. Tang, Y. Tian, J. Lu, P. Li, and J. Zhou, “Deep progressive reinforcement learning for skeleton-based action recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5323–5332.
- [33] X. Gao, W. Hu, J. Tang, P. Pan, J. Liu, and Z. Guo, “Generalized graph convolutional networks for skeleton-based action recognition,” *arXiv preprint arXiv:1811.12013*, 2018.
- [34] X. Zhang, C. Xu, X. Tian, and D. Tao, “Graph edge convolutional neural networks for skeleton based action recognition,” *arXiv preprint arXiv:1805.06184*, 2018.
- [35] L. Shi, Y. Zhang, J. Cheng, and H. Lu, “Adaptive spectral graph convolutional networks for skeleton-based action recognition,” *arXiv preprint arXiv:1805.07694*, 2018.
- [36] N. Akhtar and A. Mian, “Threat of adversarial attacks on deep learning in computer vision: A survey,” *IEEE Access*, vol. 6, pp. 14410–14430, 2018.
- [37] Y. Liu, X. Chen, C. Liu, and D. Song, “Delving into transferable adversarial examples and black-box attacks,” *arXiv preprint arXiv:1611.02770*, 2016.
- [38] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, “Practical black-box attacks against machine learning,” in *Proceedings of the 2017 ACM on Asia conference on computer and communications security*. ACM, 2017, pp. 506–519.
- [39] C. Xiao, J.-Y. Zhu, B. Li, W. He, M. Liu, and D. Song, “Spatially transformed adversarial examples,” *arXiv preprint arXiv:1801.02612*, 2018.
- [40] L. Sun, J. Wang, P. S. Yu, and B. Li, “Adversarial attack and defense on graph data: A survey,” *arXiv preprint arXiv:1812.10528*, 2018.
- [41] H. Dai, H. Li, T. Tian, X. Huang, L. Wang, J. Zhu, and L. Song, “Adversarial attack on graph structured data,” *arXiv preprint arXiv:1806.02371*, 2018.
- [42] A. Bessi, “Two samples test for discrete power-law distributions,” *arXiv preprint arXiv:1503.00643*, 2015.
- [43] M. E. Newman, D. J. Watts, and S. H. Strogatz, “Random graph models of social networks,” *Proceedings of the National Academy of Sciences*, vol. 99, no. suppl 1, pp. 2566–2572, 2002.
- [44] N. Akhtar, J. Liu, and A. Mian, “Defense against universal adversarial perturbations,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3389–3398.
- [45] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, “Least squares generative adversarial networks,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2794–2802.
- [46] A. Kurakin, I. Goodfellow, and S. Bengio, “Adversarial machine learning at scale,” *arXiv preprint arXiv:1611.01236*, 2016.
- [47] —, “Adversarial examples in the physical world,” *arXiv preprint arXiv:1607.02533*, 2016.
- [48] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [49] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, “Realtime multi-person 2d pose estimation using part affinity fields,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7291–7299.
- [50] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, “Openpose: realtime multi-person 2d pose estimation using part affinity fields,” *arXiv preprint arXiv:1812.08008*, 2018.
- [51] J. Liu, N. Akhtar, and A. Mian, “Learning human pose models from synthesized data for robust rgb-d action recognition,” *arXiv preprint arXiv:1707.00823*, 2017.