

CMD: Self-supervised 3D Action Representation Learning with Cross-modal Mutual Distillation

Yun Yao Mao¹, Wengang Zhou^{1,2,*}, Zhenbo Lu²,
Jiajun Deng¹, and Houqiang Li^{1,2,*}

¹ CAS Key Laboratory of Technology in GIPAS, EEIS Department,
University of Science and Technology of China

² Institute of Artificial Intelligence, Hefei Comprehensive National Science Center
myy2016@mail.ustc.edu.cn, zhgw@ustc.edu.cn, luzhenbo@iaai.ustc.edu.cn,
dengjj@ustc.edu.cn, lihq@ustc.edu.cn

Abstract. In 3D action recognition, there exists rich complementary information between skeleton modalities. Nevertheless, how to model and utilize this information remains a challenging problem for self-supervised 3D action representation learning. In this work, we formulate the cross-modal interaction as a bidirectional knowledge distillation problem. Different from classic distillation solutions that transfer the knowledge of a fixed and pre-trained teacher to the student, in this work, the knowledge is continuously updated and bidirectionally distilled between modalities. To this end, we propose a new **Cross-modal Mutual Distillation (CMD)** framework with the following designs. On the one hand, the neighboring similarity distribution is introduced to model the knowledge learned in each modality, where the relational information is naturally suitable for the contrastive frameworks. On the other hand, asymmetrical configurations are used for teacher and student to stabilize the distillation process and to transfer high-confidence information between modalities. By derivation, we find that the cross-modal positive mining in previous works can be regarded as a degenerated version of our CMD. We perform extensive experiments on NTU RGB+D 60, NTU RGB+D 120, and PKU-MMD II datasets. Our approach outperforms existing self-supervised methods and sets a series of new records. The code is available at: <https://github.com/maoyun Yao/CMD>

Keywords: Self-supervised 3D action recognition, contrastive learning

1 Introduction

Human action recognition, one of the fundamental problems in computer vision, has a wide range of applications in many downstream tasks, such as behavior analysis, human-machine interaction, virtual reality, *etc.* Recently, with the advancement of human pose estimation algorithms [3,14,59], skeleton-based 3D human action recognition has attracted increasing attention for its light-weight and

* Corresponding authors: Wengang Zhou and Houqiang Li

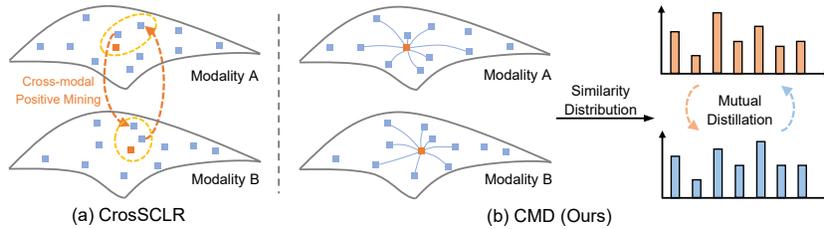


Fig. 1. CrossSCLR [23] vs. CMD (Ours). Given handful negative samples, CrossSCLR performs cross-modal positive mining according to the cosine similarity between embeddings. The nearest neighbor of the positive query in modality A will serve as an additional positive sample in modality B, and vice versa. In our approach, we reformulate cross-modal interaction as a bidirectional knowledge distillation problem, with similarity distribution that models the modality-specific knowledge.

background-robust characteristics. However, fully-supervised 3D action recognition [7,13,20,24,25,30,44,46,47,62,63,64] requires large amounts of well-annotated skeleton data for training, which is rather labor-intensive to acquire. In this paper, we focus on the self-supervised settings, aiming to avoid the laborious workload of manual annotation for 3D action representation learning.

To learn robust and discriminative representation, many celebrated pre-texts like motion prediction, jigsaw puzzle recognition, and masked reconstruction have been extensively studied in early works [27,32,33,34,50,65]. Recently, the contrastive learning frameworks [4,17,35] have been introduced to the self-supervised 3D action recognition community [27,41]. It achieves great success thanks to the capability of learning discriminative high-level semantic features. However, there still exist unsolved problems when applying contrastive learning on skeletons. On the one hand, the success of contrastive learning heavily relies on performing data augmentation [4], but the skeletons from different videos are unanimously considered as negative samples. Given the limited action categories, it would be unreasonable to just ignore potential similar instances, since they may belong to the same category as the positive one. On the other hand, cross-modal interactive learning is largely overlooked in early contrastive learning-based attempts [27,41], yet integrating multimodal information [8,11,12,26,45,56] is the key to improving the performance of 3D action recognition.

To tackle these problems, CrossSCLR [23] turns to cross-modal positive mining (see Figure 1 (a)) and sample reweighting. Though effective, it suffers the following limitations. Firstly, the positive sample mining requires reliable preliminary knowledge, thus the representation in each modality needs to be optimized independently in advance, leading to a sophisticated two-stage training process. Secondly, the contrastive context, defined as the similarity between the positive query and negative embeddings, is treated as individual weights of samples in complementary modalities to participate in the optimization process. Such implicit knowledge exchange lacks a holistic grasp of the rich contextual information. Besides, the cross-modal consistency is also not explicitly guaranteed.

In this work, we go beyond heuristic positive sample mining and reformulate cross-modal interaction as a general bidirectional knowledge distillation [18] problem. As shown in Figure 1 (b), in the proposed Cross-modal Mutual Distillation (CMD) framework, the neighboring similarity distribution is first extracted in each modality. It describes the relationship of the sample embedding with respect to its nearest neighbors in the customized feature space. Compared with individual features [18] or logits [42], such relational information is naturally suitable for modeling the knowledge learned with contrastive frameworks. Based on the relational information, bidirectional knowledge distillation between each two modalities is performed via explicit cross-modal consistency constraints. Since the representation in each skeleton modality is trained from scratch and there is no intuitive teacher-student relationship between modalities, embeddings from the momentum updated key encoder along with a smaller temperature are used for knowledge modeling on the teacher side, so as to stabilize the distillation process and highlight the high-confidence information in each modality.

Compared to previous works, the advantages of our approach are three-fold: **i)** Instead of heuristically reweighting training samples, the contextual information in contrastive learning is treated as a whole to model the modality-specific knowledge, explicitly ensuring the cross-modal consistency during distillation. **ii)** Unlike cross-modal positive sample mining, our approach does not heavily rely on the initial representation, thus is free of the sophisticated two-stage training. This largely benefits from the probabilistic knowledge modeling strategy. Moreover, the positive mining is also mathematically proved to be a special case of the proposed cross-modal distillation mechanism under extreme settings. **iii)** The proposed CMD is carefully designed to be well integrated into the existing contrastive framework with almost no extra computational overhead introduced.

We perform extensive experiments on three prevalent benchmark datasets: NTU RGB+D 60 [43], NTU RGB+D 120 [28], and PKU-MMD II [9]. Our approach achieves state-of-the-art results on all of them under all evaluation protocols. It’s worth noting that the proposed cross-modal mutual distillation is easily implemented in a few lines of code. We hope this simple yet effective approach will serve as a strong baseline for future research.

2 Related Work

Self-supervised Representation Learning: Self-supervised learning methods can be roughly divided into two categories: generative and contrastive [29]. Generative methods [2,16,36] try to reconstruct the original input to learn meaningful latent representation. Contrastive learning [4,17,35] aims to learn feature representation via instance discrimination. It pulls positive pairs closer and pushes negative pairs away. Since no labels are available during self-supervised contrastive learning, two different augmented versions of the same sample are treated as a positive pair, and samples from different instances are considered to be negative. In MoCo [17] and MoCo v2 [5], the negative samples are taken from previous batches and stored in a queue-based memory bank. In contrast,

SimCLR [4] and MoCo v3 [6] rely on a larger batch size to provide sufficient negative samples. Similar to the contrastive context in [23], the neighboring similarity in this paper is defined as the normalized product between positive embedding and its neighboring anchors. Our goal is to transfer such modality-specific information between skeleton modalities to facilitate better contrastive 3D action representation learning.

Self-supervised 3D Action Recognition: Many previous works have been proposed to perform self-supervised 3D action representation learning. In LongT GAN [65], an autoencoder-based model along with an additional adversarial training strategy are proposed. Following the generative paradigm, it learns latent representation via sequential reconstruction. Similarly, P&C [50] trains an encoder-decoder network to both predict and cluster skeleton sequences. To learn features that are more robust and separable, the authors also propose strategies to weaken the decoder, laying more burdens on the encoder. Different from previously mentioned methods that merely adopt a single reconstruction task, MS²L [27] integrates multiple pretext tasks to learn better representation. In recent attempts [23,41,53,54], momentum encoder-based contrastive learning is introduced and better performance is achieved. Among them, CrosSCLR [23] is the first to perform cross-modal knowledge mining. It finds potential positives and re-weights training samples with the contrastive contexts from different skeleton modalities. However, the positive mining performed in CrosSCLR requires reliable initial representation, two-stage training is indispensable. Differently, in this paper, a more general knowledge distillation mechanism is introduced to perform cross-modal information interaction. Besides, the positive mining performed in CrosSCLR can be regarded as a special case of our approach.

Similarity-based Knowledge Distillation: Pairwise similarity has been shown to be useful information in relational knowledge distillation [38,40,55]. In PKT [39], CompRes [1], and SEED [15], similarities of each sample with respect to a set of anchors are converted into a probability distribution, which models the structural information of the data. After that, knowledge distillation is performed by training the student to mimic the probability distribution of the teacher. Recently, contextual similarity information has also shown great potential in image retrieval [37,58] and representation learning [22,51]. Our approach is partially inspired by these works. Differently, the cross-modal mutual distillation in our approach is designed to answer the question of *how to transfer the biased knowledge between complementary modalities during 3D action pre-training*.

3 Method

3.1 Framework Overview

By consolidating the idea of leveraging complementary information from cross-modal inputs to improve 3D action representation learning, we design the Cross-modal Mutual Distillation (CMD) framework. As shown in Figure 2, the proposed CMD consists of two key components: Single-modal Contrastive Learning (SCL) and Cross-modal Mutual Distillation (CMD). Given multiple skeleton

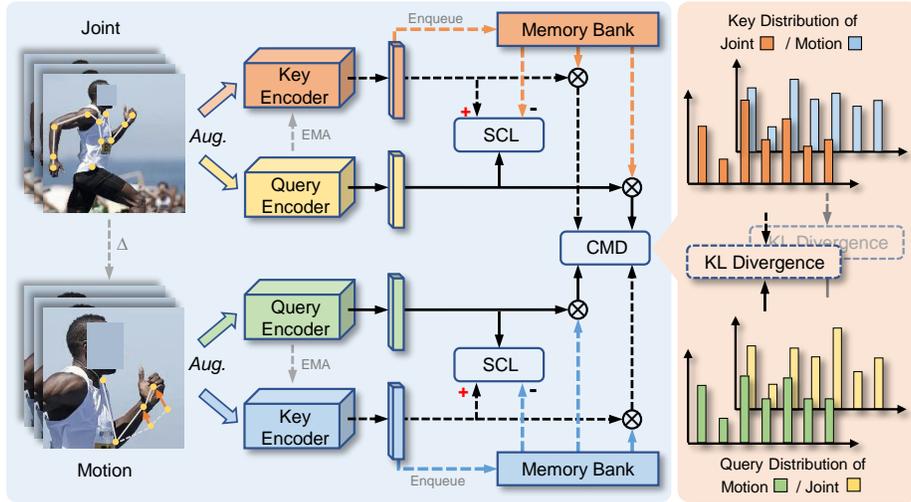


Fig. 2. The overall pipeline of the proposed framework. It contains two modules, Single-modal Contrastive Learning (SCL) and Cross-modal Mutual Distillation (CMD). Given multiple skeleton modalities (*e.g.* joint and motion) as input, the SCL module performs self-supervised contrastive learning in each modality and the CMD module simultaneously transfers the learned knowledge between modalities. SCL and CMD work collaboratively so that each modality learns more comprehensive representation.

modalities (*e.g.* joint, motion, and bone) as input, SCL is applied to each of them to learn customized 3D action representation. Meanwhile, in CMD, the knowledge learned by SCL is modeled by the neighboring similarity distributions, which describe the relationship between the sample embedding and its nearest neighbors. Cross-modal knowledge distillation is then performed by bidirectionally minimizing the KL divergence between the distributions corresponding to each modality. SCL and CMD run synchronously and cooperatively so that each modality learns more comprehensive representation.

3.2 Single-modal Contrastive Learning

In this section, we revisit the single-modal contrastive learning as the preliminary of our approach, which has been widely adopted in many tasks like image/video recognition [10,21,49] and correspondence learning [57]. In self-supervised 3D action recognition, previous works like AS-CAL [41], CrosSCLR [23], ISC [53], and AimCLR [54] also take the contrastive method MoCo v2 [5] as their baseline.

Given a single-modal skeleton sequence x , we first perform data augmentation to obtain two different views x_q and x_k (query and key). Then, two encoders are adopted to map the positive pair x_q and x_k into feature embeddings $z_q = E_q(x_q, \theta_q)$ and $z_k = E_k(x_k, \theta_k)$, where E_q and E_k denote query encoder and key encoder, respectively. θ_q and θ_k are the learnable parameters of the two encoders.

Note that in MoCo v2, the key encoder is not trained by gradient descent but the momentum updated version of the query encoder: $\theta_k \leftarrow \alpha\theta_k + (1 - \alpha)\theta_q$, where α is a momentum coefficient that controls the updating speed. During self-supervised pre-training, the noise contrastive estimation loss InfoNCE [35] is used to perform instance discrimination, which is computed as follows:

$$\mathcal{L}_{\text{SCL}} = -\log \frac{\exp(z_q^\top z_k / \tau_c)}{\exp(z_q^\top z_k / \tau_c) + \sum_{i=1}^N \exp(z_q^\top m_i / \tau_c)}, \quad (1)$$

where τ_c is a temperature hyper-parameter [18] that scales the distribution of instances and m_i is the key embedding of negative sample. N is the size of a queue-based memory bank \mathbf{M} where all the negative key embeddings are stored. After the training of the current mini-batch, z_k is enqueued as a new negative key embedding and the oldest embeddings in the memory bank are dequeued.

Under the supervision of the InfoNCE loss, the encoder is forced to learn representation that is invariant to data augmentations, thereby focusing on semantic information shared between positive pairs. Nevertheless, the learned representation is often modally biased, making it difficult to account for all data characteristics. Though it can be alleviated by test-time ensembling, several times the running overhead will be introduced. Moreover, the inherent limitations of the learned representation in each modality still exist. Therefore, during self-supervised pre-training, cross-modal interaction is essential.

3.3 Cross-modal Mutual Distillation

While SCL is performed within each skeleton modality, the proposed CMD models the learned knowledge and transfers it between modalities. This enables each modality to receive knowledge from other perspectives, thereby alleviating the modal bias of the learned representation. Based on MoCo v2, CMD can be easily implemented in a few lines of code, as shown in Alg. 1.

Knowledge Modeling: To perform knowledge distillation between modalities, we first need to model the knowledge learned in each modality in a proper way. It needs to take advantage of the existing contrastive learning framework to avoid introducing excessive computational overhead. Moreover, since the distillation is performed cross-modally for self-supervised learned knowledge, conventional methods that rely on individual features/logits are no longer applicable.

Inspired by recent relational knowledge distillation works [38,40,55], we utilize the pairwise relationship between samples for modality-specific knowledge modeling. Given an embedding z and a set of anchors $\{n_i\}_{i=1,2,\dots,K}$, we compute the similarities between them as $\text{sim}(z, n_i) = z^\top n_i, i = 1, 2, \dots, K$.

In the MoCo v2 [5] framework, there are a handful of negative embeddings stored in the memory bank. We can easily obtain the required anchors without additional model inference. Note that if all the negative embeddings are used as anchors, the set $\{z^\top m_i\}_{i=1,2,\dots,N}$ is exactly the contrastive context defined in [23]. In our approach, we select the top K nearest neighbors of z as the

anchors. The resulting pairwise similarities are further converted into probability distributions with a temperature hyper-parameter τ :

$$p_i(z, \tau) = \frac{\exp(z^\top n_i / \tau)}{\sum_{j=1}^K \exp(z^\top n_j / \tau)}, i = 1, 2, \dots, K. \quad (2)$$

The obtained $\mathbf{p}(z, \tau) = \{p_i(z, \tau)\}_{i=1,2,\dots,K}$ describes the distribution characteristic around the embedding z in the customized feature space of each modality.

Knowledge Distillation: Based on the aforementioned probability distributions, an intuitive way to perform knowledge distillation would be to directly establish consistency constraints between skeleton modalities. Different from previous knowledge distillation approaches that transfer the knowledge of a fixed and well-trained teacher model to the student, in our approach, the knowledge is continuously updated during self-supervised pre-training and each modality acts as both student and teacher.

To this end, based on the contrastive framework, we make two customized designs in the proposed approach: **i)** Different embeddings are used for teacher and student. As shown in Figure 2, in MoCo v2 [5], two augmented views of the same sample are encoded into query z_q and key z_k , respectively. In our approach, the key distribution obtained in one modality is used to guide the learning of query distribution in other modalities, so that knowledge is transferred accordingly. Specifically, for the key embedding z_k^a from modality A and the query embedding z_q^b from modality B, we select the top K nearest neighbors of z_k^a as anchors and compute the similarity distributions as $\mathbf{p}(z_q^b, \tau)$ and $\mathbf{p}(z_k^a, \tau)$ according to Eq. 2. Knowledge distillation from modality A to modality B is performed by minimizing the following KL divergence:

$$\text{KL}(\mathbf{p}(z_k^a, \tau) \parallel \mathbf{p}(z_q^b, \tau)) = \sum_{i=1}^K p_i(z_k^a, \tau) \cdot \log \frac{p_i(z_k^a, \tau)}{p_i(z_q^b, \tau)}. \quad (3)$$

Since the key encoder is not trained with gradient, the teacher is not affected during unidirectional knowledge distillation. Moreover, the momentum updated key encoder provides more stable knowledge for the student to learn. **ii)** Asymmetric temperatures τ_t and τ_s are employed for teacher and student, respectively. Considering that there is no intuitive teacher-student relationship between modalities, a smaller temperature is applied for the teacher in CMD to emphasize the high-confidence information, as discussed in [52].

Since the knowledge distillation works bidirectionally, given two modalities A and B, the loss function for CMD is formulated as follows:

$$\mathcal{L}_{\text{CMD}} = \text{KL}(\mathbf{p}(z_k^a, \tau_t) \parallel \mathbf{p}(z_q^b, \tau_s)) + \text{KL}(\mathbf{p}(z_k^b, \tau_t) \parallel \mathbf{p}(z_q^a, \tau_s)). \quad (4)$$

Note that Eq. 4 can be easily extended if more modalities are involved. The final loss function in our approach is the combination of \mathcal{L}_{SCL} and \mathcal{L}_{CMD} :

$$\mathcal{L} = \mathcal{L}_{\text{SCL}}^a + \mathcal{L}_{\text{SCL}}^b + \mathcal{L}_{\text{CMD}}, \quad (5)$$

where the superscripts a and b denote modality A and B, respectively.

Algorithm 1 Pseudocode of the CMD module in a PyTorch-like style.

```

1 # z_q_a, z_q_b, z_k_a, z_k_b: query/key embeddings in modality A/B (BxC)
2 # queue_a, queue_b: queue of N keys in modality A/B (CxN)
3 # tau_s, tau_t: temperatures for student/teacher (scalars)
4
5 l_a, lk_a = torch.mm(z_q_a, queue_a), torch.mm(z_k_a, queue_a) # compute similarities
6 l_b, lk_b = torch.mm(z_q_b, queue_b), torch.mm(z_k_b, queue_b)
7
8 lk_a_topk, idx_a = torch.topk(lk_a, K, dim=-1) # select top K nearest neighbors
9 lk_b_topk, idx_b = torch.topk(lk_b, K, dim=-1)
10
11 loss_cmd = loss_kld(torch.gather(l_b, -1, idx_a) / tau_s, lk_a_topk / tau_t) # A to B
12 + loss_kld(torch.gather(l_a, -1, idx_b) / tau_s, lk_b_topk / tau_t) # B to A
13
14 def loss_kld(inputs, targets):
15     inputs, targets = F.log_softmax(inputs, dim=1), F.softmax(targets, dim=1)
16     return F.kl_div(inputs, targets, reduction='batchmean')
```

3.4 Relationship with Positive Mining

Cross-modal Positive Mining: Cross-modal positive mining is the most important component in CrosSCLR [23], where the most similar negative sample is selected to boost the positive sets for contrastive learning in complementary modalities. The contrastive loss for modality B is reformulated as:

$$\begin{aligned}
\mathcal{L}_{\text{CPM}}^b &= -\log \frac{\exp(z_q^{b\top} z_k^b / \tau_c)}{\exp(z_q^{b\top} z_k^b / \tau_c) + \sum_{i=1}^N \exp(z_q^{b\top} m_i^b / \tau_c)} \\
&\quad - \log \frac{\exp(z_q^{b\top} m_u^b / \tau_c)}{\exp(z_q^{b\top} z_k^b / \tau_c) + \sum_{i=1}^N \exp(z_q^{b\top} m_i^b / \tau_c)} \\
&= \mathcal{L}_{\text{SCL}}^b - \log \frac{\exp(z_q^{b\top} m_u^b / \tau_c)}{\exp(z_q^{b\top} z_k^b / \tau_c) + \sum_{i=1}^N \exp(z_q^{b\top} m_i^b / \tau_c)},
\end{aligned} \tag{6}$$

where u is the index of most similar negative sample in modality A.

CMD with $\tau_t = 0$ and $K = N$: Setting temperature $\tau_t = 0$ and $K = N$, the key distribution $\mathbf{p}(z_k^a, \tau_t)$ in Eq. 4 will be an one-hot vector with the only 1 at index u , and thus the loss works on modality B will be like:

$$\begin{aligned}
\mathcal{L}^b &= \mathcal{L}_{\text{SCL}}^b + \mathcal{L}_{\text{CMD}}^b \\
&= \mathcal{L}_{\text{SCL}}^b + \sum_{i=1}^N p_i(z_k^a, 0) \cdot \log \frac{p_i(z_k^a, 0)}{p_i(z_q^b, \tau_s)} \\
&= \mathcal{L}_{\text{SCL}}^b + 1 \cdot \log \frac{1}{p_u(z_q^b, \tau_s)} \\
&= \mathcal{L}_{\text{SCL}}^b - \log \frac{\exp(z_q^{b\top} m_u^b / \tau_s)}{\sum_{j=1}^N \exp(z_q^{b\top} m_j^b / \tau_s)}.
\end{aligned} \tag{7}$$

We can find that the loss $\mathcal{L}_{\text{CMD}}^b$ is essentially doing contrastive learning in modality B with the positive sample mined by modality A. Compared with Eq. 6, the

only difference is that when the mined m_u^b is taken as the positive sample, the key embedding z_k^b is excluded from the denominator. The same result holds for modality A. Thus we draw a conclusion that the cross-modal positive mining performed in CrosSCLR [23] can be regarded as a special case of our approach with the temperature of teacher $\tau_t = 0$ and the number of neighbors $K = N$.

4 Experiments

4.1 Implementation Details

Network Architecture: In our approach, we adopt a 3-layer Bidirectional GRU (BiGRU) as the base-encoder, which has a hidden dimension of 1024. Before the encoder, we additionally add a Batch Normalization [19] layer to stabilize the training process. Each skeleton sequence is represented in a two-actor manner, where the second actor is set to zeros if only one actor exists. The sequences are further resized to a temporal length of 64 frames.

Self-supervised Pre-training: During pre-training, we adopt MoCo v2 [5] to perform single-modal contrastive learning. The temperature hyper-parameter in the InfoNCE [35] loss is 0.07. In cross-modal mutual distillation, the temperatures for teacher and student are set to 0.05 and 0.1, respectively. The number of neighbors K is set to 8192. The SGD optimizer is employed with a momentum of 0.9 and a weight decay of 0.0001. The batch size is set to 64 and the initial learning rate is 0.01. For NTU RGB+D 60 [43] and NTU RGB+D 120 [28] datasets, the model is trained for 450 epochs, the learning rate is reduced to 0.001 after 350 epochs, and the size of the memory bank N is 16384. For PKU-MMD II [9] dataset, the total epochs are increased to 1000, and the learning rate drops at epoch 800. We adopt the same skeleton augmentations as ISC [53].

4.2 Datasets and Metrics

NTU RGB+D 60 [43]: NTU-RGB+D 60 (NTU-60) is a large-scale multi-modality action recognition dataset which is captured by three Kinect v2 cameras. It contains 60 action categories and 56,880 sequences. The actions are performed by 40 different subjects (actors). In this paper, we adopt its 3D skeleton data for experiments. Specifically, each human skeleton contains 25 body joints, and each joint is represented as 3D coordinates. Two evaluation protocols are recommended by the authors: cross-subject (x-sub) and cross-view (x-view). For x-sub, action sequences performed by half of the 40 subjects are used as training samples and the rest as test samples. For x-view, the training samples are captured by camera 2 and 3 and the test samples are from camera 1.

NTU RGB+D 120 [28]: Compared with NTU-60, NTU-RGB+D 120 (NTU-120) extends the action categories from 60 to 120, with 114,480 skeleton sequences in total. The number of subjects is also increased from 40 to 106. Moreover, a new evaluation protocol named cross-setup (x-set) is proposed as a substitute for x-view. Specifically, the sequences are divided into 32 different setups

Table 1. Performance comparison on NTU-60, NTU-120, and PKU-II in terms of the linear evaluation protocol. Our approach achieves state-of-the-art performance on all of them, both when taking single skeleton modality as input and when ensembling multiple modalities during evaluation. The prefix “3s-” denotes multi-modal ensembling.

Method	Modality	NTU-60		NTU-120		PKU-II
		x-sub	x-view	x-sub	x-set	x-sub
LongT GAN [65]	Joint only	39.1	48.1	-	-	26.0
MS ² L [27]	Joint only	52.6	-	-	-	27.6
P&C [50]	Joint only	50.7	76.3	42.7	41.7	25.5
AS-CAL [41]	Joint only	58.5	64.8	48.6	49.2	-
SeBiReNet [33]	Joint only	-	79.7	-	-	-
AimCLR [54]	Joint only	74.3	79.7	-	-	-
ISC [53]	Joint only	76.3	85.2	67.1	67.9	36.0
CrosSCLR-B	Joint only	77.3	85.1	67.1	68.6	41.9
CMD (Ours)	Joint only	79.8	86.9	70.3	71.5	43.0
3s-CrosSCLR [23]	Joint+Motion+Bone	77.8	83.4	67.9	66.7	21.2
3s-AimCLR [54]	Joint+Motion+Bone	78.9	83.8	68.2	68.8	39.5
3s-CrosSCLR-B	Joint+Motion+Bone	82.1	89.2	71.6	73.4	51.0
3s-CMD (Ours)	Joint+Motion+Bone	84.1	90.9	74.7	76.1	52.6

according to the camera distances and background, with half of the 32 setups (even-numbered) used for training and the rest for testing.

PKU-MMD [9]: PKU-MMD is a new benchmark for multi-modality 3D human action detection. It can also be used for action recognition tasks [27]. PKU-MMD has two phases, where Phase II is extremely challenging since more noise is introduced by large view variation. In this work, we evaluate the proposed method on Phase II (PKU-II) under the widely used cross-subject evaluation protocol, with 5,332 skeleton sequences for training and 1,613 for testing.

Evaluation Metrics: We report the top-1 accuracy for all datasets.

4.3 Comparison with State-of-the-art Methods

In the section, the learned representation is utilized for 3D action classification under a variety of evaluation protocols. We compare the results with previous state-of-the-art methods. Note that during evaluation, we only take single skeleton modality (joint) as input by default, which is consistent with previous arts [27,50,53]. Integrating multiple skeleton modalities for evaluation can significantly improve the performance, but it will also incur more time overhead.

Linear Evaluation Protocol: For linear evaluation protocol, we freeze the pre-trained encoder and add a learnable linear classifier after it. The classifier is trained on the corresponding training set for 80 epochs with a learning rate of 0.1 (reduced to 0.01 and 0.001 at epoch 50 and 70, respectively). We evaluate the proposed method on the NTU-60, NTU-120, and PKU-II datasets. As shown in Table 1, we include the recently proposed CrosSCLR [23], ISC [53], and AimCLR

Table 2. Performance comparison on NTU-60 and NTU-120 in terms of the KNN evaluation protocol. The learned representation exhibits the best performance on both datasets. Surpassing previous state-of-the-art methods by a considerable margin.

Method	NTU-60		NTU-120	
	x-sub	x-view	x-sub	x-set
LongT GAN [65]	39.1	48.1	31.5	35.5
P&C [50]	50.7	76.3	39.5	41.8
ISC [53]	62.5	82.6	50.6	52.3
CrosSCLR-B	66.1	81.3	52.5	54.9
CMD (Ours)	70.6	85.4	58.3	60.9

Table 3. Performance comparison on PKU-II in terms of the transfer learning evaluation protocol. The source datasets are NTU-60 and NTU-120. The representation learned by our approach shows the best transferability.

Method	To PKU-II	
	NTU-60	NTU-120
LongT GAN [65]	44.8	-
MS ² L [27]	45.8	-
ISC [53]	51.1	52.3
CrosSCLR-B	54.0	52.8
CMD (Ours)	56.0	57.0

[54] for comparison. Our approach outperforms previous state-of-the-art methods by a considerable margin on all the three benchmarks. Note that ISC and the proposed CMD share the same BiGRU encoder, which is different from the ST-GCN [60] encoder in CrosSCLR. For a fair comparison, we additionally train a variation of CrossSCLR with BiGRU as its base-encoder (denoted as CrosSCLR-B). We can find that our method still outperforms it on all the three datasets, which shows the superiority of the proposed cross-modal mutual distillation.

KNN Evaluation Protocol: An alternative way to use the pre-trained encoder for action classification is to directly apply a K-Nearest Neighbor (KNN) classifier to the learned features of the training samples. Following [50], we assign each test sample to the most similar class where its nearest neighbor is in (*i.e.* KNN with $k=1$). As shown in Table 2, we perform experiments on the NTU-60 and NTU-120 benchmarks and compare the results with previous works. For both datasets, our approach exhibits the best performance, surpassing CrosSCLR-B [23] by 4.5%~6% in the more challenging cross-subject and cross-setup protocols.

Transfer Learning Evaluation Protocol: In transfer learning evaluation protocol, we examine the transferability of the learned representation. Specifically, we first utilize the proposed framework to pre-train the encoder on the source dataset. Then the pre-trained encoder along with a linear classifier are finetuned on the target dataset for 80 epochs with a learning rate of 0.01 (reduced to 0.001 at epoch 50). We select NTU-60 and NTU-120 as source datasets, and PKU-II as the target dataset. We compare the proposed approach with previous methods LongT GAN [65], MS²L [27], and ISC [53] under the cross-subject protocol. As shown in Table 3, our approach exhibits superior performance on the PKU-II dataset after large-scale pre-training, outperforming previous methods by a considerable margin. This indicates that the representation learned by our approach is more transferable.

Semi-supervised Evaluation Protocol: In semi-supervised classification, both labeled and unlabeled data are included during training. Its goal is to train a classifier with better performance than the one trained with only labeled

Table 4. Performance comparison on NTU-60 in terms of the semi-supervised evaluation protocol. We randomly select a portion of the labeled data to fine-tune the pre-trained encoder, and the average of five runs is reported as the final performance. Our approach exhibits the state-of-the-art results compared with previous methods.

Method	NTU-60							
	x-view				x-sub			
	(1%)	(5%)	(10%)	(20%)	(1%)	(5%)	(10%)	(20%)
LongT GAN [65]	-	-	-	-	35.2	-	62.0	-
MS ² L [27]	-	-	-	-	33.1	-	65.1	-
ASSL [48]	-	63.6	69.8	74.7	-	57.3	64.3	68.0
ISC [53]	38.1	65.7	72.5	78.2	35.7	59.6	65.9	70.8
CrosSCLR-B [23]	49.8	70.6	77.0	81.9	48.6	67.7	72.4	76.1
CMD (Ours)	53.0	75.3	80.2	84.3	50.6	71.0	75.4	78.7
3s-CrosSCLR [23]	50.0	-	77.8	-	51.1	-	74.4	-
3s-Colorization [61]	52.5	-	78.9	-	48.3	-	71.7	-
3s-AimCLR [54]	54.3	-	81.6	-	54.8	-	78.2	-
3s-CMD (Ours)	55.5	77.2	82.4	86.6	55.6	74.3	79.0	81.8

samples. For a fair comparison, we adopt the same strategy as ISC [53]. The pre-trained encoder is fine-tuned together with the post-attached linear classifier on a portion of the corresponding training set. We conduct experiments on the NTU-60 dataset. As shown in Table 4, we report the evaluation results when the proportion of supervised data is set to 1%, 5%, 10%, and 20%, respectively. Compared with previous methods LongT GAN [65], MS²L [27], ASSL [48], and ISC [53], our algorithm exhibits superior performance. For example, with the same baseline, the proposed approach outperforms ISC and CrosSCLR-B by a large margin. We also take 3s-CrosSCLR [23], 3s-Colorization [61], and recently proposed 3s-AimCLR [54] into comparison. In these methods, test-time multi-modal ensembling is performed and the results of using 1% and 10% labeled data are reported. We can find that our 3s-CMD still outperforms all of these methods after ensembling multiple skeleton modalities.

4.4 Ablation study

To justify the effectiveness of the proposed cross-modal mutual distillation framework, we conduct several ablative experiments on the NTU-60 dataset according to the cross-subject protocol. More results can be found in the supplementary.

Number of neighbors: The number of nearest neighbors controls the abundance of contextual information used in the proposed cross-modal mutual distillation module. We test the performance of the learned representation with respect to different numbers of nearest neighbors K under the linear evaluation protocol. As shown in Figure 3, on the downstream classification task, the performance of the pre-trained encoder improves as K increases. When K is large

Table 5. Ablative experiments of modality selection and bidirectional distillation. The performance is evaluated on the NTU-60 dataset according to the cross-subject protocol. J, M, and B denote joint, motion, and bone modality respectively. The horizontal arrows indicate the direction of distillation.

Modality & Direction	Linear Evaluation				KNN Evaluation			
	Bone	Motion	Joint	Δ	Bone	Motion	Joint	Δ
Baseline	74.4	73.1	76.1		62.0	56.8	63.4	
J \leftarrow B	74.4	-	76.5		62.0	-	64.3	
J \rightleftharpoons B	76.6	-	77.7	\uparrow 1.2	65.9	-	66.5	\uparrow 2.2
J \leftarrow M	-	73.1	78.9		-	56.8	64.8	
J \rightleftharpoons M	-	77.5	79.8	\uparrow 0.9	-	67.0	68.7	\uparrow 3.9
J \leftarrow M, J \leftarrow B	74.4	73.1	78.8		62.0	56.8	66.5	
J \rightleftharpoons M, J \rightleftharpoons B, M \rightleftharpoons B	77.8	77.1	79.4	\uparrow 0.6	69.5	68.7	70.6	\uparrow 4.1

enough ($K \geq 8192$), continuing to increase its value hardly contributes to the performance. This is because the newly added neighbors are far away and contain little reference value for describing the distribution around the query sample. In addition, we can also find that when the value of K varies from 64 to 16384, the performance of our approach is consistently higher than that of CrosSCLR-B [23] and our baseline. This demonstrates the superiority and robustness of the proposed approach.

Modality Selection: In our approach, we consider three kinds of skeleton modalities for self-supervised pre-training as in [23]. They are joint, motion, and bone, respectively. Our approach is capable of performing knowledge distillation between any two of the above modalities. As shown in Table 5, we report the performance of the representation obtained by pre-training with different combinations of skeleton modalities. Note that the joint modality is always preserved since it is used for evaluation. There are several observations as follows: **i)** Cross-modal knowledge distillation helps to improve the performance of the representation in student modalities. **ii)** Under the linear evaluation protocol, knowledge distillation between joint and motion achieves the optimal performance, exceeding the baseline by 3.7%. **iii)** Under the KNN evaluation protocol, the learned representation shows the best results when all the three modalities are involved in knowledge distillation, which outperforms the baseline with an absolute improvement of 7.2%.

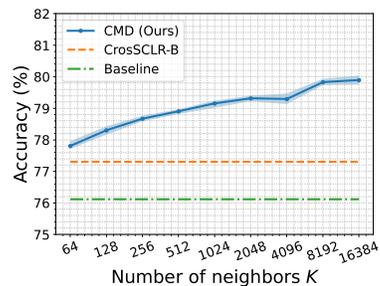


Fig. 3. Ablative study of the number of neighbors K in the cross-modal mutual distillation module. The performance is evaluated on the cross-subject protocol of the NTU-60 dataset.

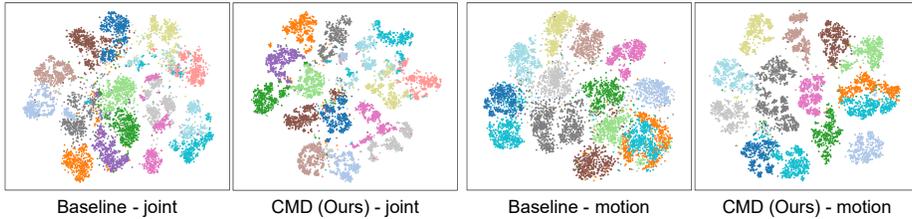


Fig. 4. t-SNE [31] visualization of feature embeddings. We sample 15 action classes from the NTU-60 dataset and visualize the features extracted by the proposed CMD and its baseline respectively. Compared with the baseline, CMD learns more compact and more discriminative representation in both joint and motion modalities.

Bidirectional Distillation: In addition to modality selection, we also verify the effectiveness of bidirectional distillation. It enables the modalities involved in the distillation to interact with each other and progress together, forming a virtuous circle. In Table 5, the last column of each evaluation protocol reports the performance gain of bidirectional mutual distillation over unidirectional distillation in the joint modality. Results show that regardless of which skeleton modalities are used during pre-training, bidirectional mutual distillation further boosts the performance, especially under the KNN evaluation protocol.

Qualitative Results: We visualize the learned representation of the proposed approach and compare it with that of the baseline. The t-SNE [31] algorithm is adopted to reduce the dimensionality of the representation. To obtain clearer results, we select only 1/4 of the categories in the NTU-60 dataset for visualization. The final results are illustrated in Figure 4. For both joint and motion modalities, the representation learned by our approach is more compactly clustered than those learned by the baseline in the feature space. This brings a stronger discrimination capability to the representation, explaining the stunning performance of our approach in Table 2.

5 Conclusion

In this work, we presented a novel approach for self-supervised 3D action representation learning. It reformulates cross-modal reinforcement as a bidirectional knowledge distillation problem, where the pairwise similarities between embeddings are utilized to model the modality-specific knowledge. The carefully designed cross-modal mutual distillation module can be well integrated into the existing contrastive learning framework, thus avoiding additional computational overhead. We evaluate the learned representation on three 3D action recognition benchmarks with four widely adopted evaluation protocols. The proposed approach sets a series of new state-of-the-art records on all of them, demonstrating the effectiveness of the cross-modal mutual distillation.

Acknowledgement: This work was supported by the National Natural Sci-

ence Foundation of China under Contract U20A20183 and 62021001. It was also supported by the GPU cluster built by MCC Lab of Information Science and Technology Institution, USTC.

References

1. Abbasi Koohpayegani, S., Tejankar, A., Pirsiavash, H.: Compress: Self-supervised learning by compressing representations. In: Proceedings of the Advances in Neural Information Processing Systems (NeurIPS). vol. 33, pp. 12980–12992 (2020) [4](#)
2. Ballard, D.H.: Modular learning in neural networks. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). vol. 647, pp. 279–284 (1987) [3](#)
3. Cao, Z., Hidalgo, G., Simon, T., Wei, S., Sheikh, Y.: Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **43**(01), 172–186 (2021) [1](#)
4. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: Proceedings of the International Conference on Machine Learning (ICML). pp. 1597–1607 (2020) [2](#), [3](#), [4](#)
5. Chen, X., Fan, H., Girshick, R., He, K.: Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297* (2020) [3](#), [5](#), [6](#), [7](#), [9](#)
6. Chen, X., Xie, S., He, K.: An empirical study of training self-supervised vision transformers. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 9640–9649 (2021) [4](#)
7. Chen, Y., Zhang, Z., Yuan, C., Li, B., Deng, Y., Hu, W.: Channel-wise topology refinement graph convolution for skeleton-based action recognition. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 13359–13368 (2021) [2](#)
8. Cheng, K., Zhang, Y., He, X., Chen, W., Cheng, J., Lu, H.: Skeleton-based action recognition with shift graph convolutional network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 183–192 (2020) [2](#)
9. Chunhui, L., Yueyu, H., Yanghao, L., Sijie, S., Jiaying, L.: Pku-mmd: A large scale benchmark for continuous multi-modal human action understanding. *arXiv preprint arXiv:1703.07475* (2017) [3](#), [9](#), [10](#)
10. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 248–255 (2009) [5](#)
11. Deng, J., Pan, Y., Yao, T., Zhou, W., Li, H., Mei, T.: Single shot video object detector. *IEEE Transactions on Multimedia* **23**, 846–858 (2021) [2](#)
12. Deng, J., Yang, Z., Liu, D., Chen, T., Zhou, W., Zhang, Y., Li, H., Ouyang, W.: Transvg++: End-to-end visual grounding with language conditioned vision transformer. *arXiv preprint arXiv:2206.06619* (2022) [2](#)
13. Du, Y., Wang, W., Wang, L.: Hierarchical recurrent neural network for skeleton based action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1110–1118 (2015) [2](#)
14. Fang, H.S., Xie, S., Tai, Y.W., Lu, C.: Rmpe: Regional multi-person pose estimation. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 2334–2343 (2017) [1](#)
15. Fang, Z., Wang, J., Wang, L., Zhang, L., Yang, Y., Liu, Z.: Seed: Self-supervised distillation for visual representation. In: Proceedings of the International Conference on Learning Representations (ICLR) (2021) [4](#)

16. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 16000–16009 (2022) [3](#)
17. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9729–9738 (2020) [2](#), [3](#)
18. Hinton, G., Vinyals, O., Dean, J., et al.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015) [3](#), [6](#)
19. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: Proceedings of the International Conference on Machine Learning (ICML). pp. 448–456 (2015) [9](#)
20. Ke, Q., Bennamoun, M., An, S., Sohel, F., Boussaid, F.: A new representation of skeleton sequences for 3d action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3288–3297 (2017) [2](#)
21. Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T.: Hmdb: a large video database for human motion recognition. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 2556–2563 (2011) [5](#)
22. Li, J., Selvaraju, R.R., Gotmare, A.D., Joty, S., Xiong, C., Hoi, S.: Align before fuse: Vision and language representation learning with momentum distillation. In: Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) (2021) [4](#)
23. Li, L., Wang, M., Ni, B., Wang, H., Yang, J., Zhang, W.: 3d human action representation learning via cross-view consistency pursuit. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4741–4750 (2021) [2](#), [4](#), [5](#), [6](#), [8](#), [9](#), [10](#), [11](#), [12](#), [13](#)
24. Li, M., Chen, S., Chen, X., Zhang, Y., Wang, Y., Tian, Q.: Actional-structural graph convolutional networks for skeleton-based action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3595–3603 (2019) [2](#)
25. Li, T., Ke, Q., Rahmani, H., Ho, R.E., Ding, H., Liu, J.: Else-net: Elastic semantic network for continual action recognition from skeleton data. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 13434–13443 (2021) [2](#)
26. Liang, D., Fan, G., Lin, G., Chen, W., Pan, X., Zhu, H.: Three-stream convolutional neural network with multi-task and ensemble learning for 3d action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 934–940 (2019) [2](#)
27. Lin, L., Song, S., Yang, W., Liu, J.: Ms2l: Multi-task self-supervised learning for skeleton based action recognition. In: Proceedings of the 28th ACM International Conference on Multimedia (ACM MM). pp. 2490–2498 (2020) [2](#), [4](#), [10](#), [11](#), [12](#)
28. Liu, J., Shahroudy, A., Perez, M., Wang, G., Duan, L.Y., Kot, A.C.: Ntu rgb+d 120: A large-scale benchmark for 3d human activity understanding. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **42**(10), 2684–2701 (2020) [3](#), [9](#)
29. Liu, X., Zhang, F., Hou, Z., Mian, L., Wang, Z., Zhang, J., Tang, J.: Self-supervised learning: Generative or contrastive. IEEE Transactions on Knowledge and Data Engineering (TKDE) (2021) [3](#)
30. Liu, Z., Zhang, H., Chen, Z., Wang, Z., Ouyang, W.: Disentangling and unifying graph convolutions for skeleton-based action recognition. In: Proceedings of the

- IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 143–152 (2020) [2](#)
31. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *Journal of Machine Learning Research (JMLR)* **9**(11), 2579–2605 (2008) [14](#)
 32. Misra, I., Zitnick, C.L., Hebert, M.: Shuffle and learn: unsupervised learning using temporal order verification. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 527–544 (2016) [2](#)
 33. Nie, Q., Liu, Z., Liu, Y.: Unsupervised 3d human pose representation with viewpoint and pose disentanglement. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 102–118 (2020) [2](#), [10](#)
 34. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 69–84 (2016) [2](#)
 35. Van den Oord, A., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1809.03327* (2018) [2](#), [3](#), [6](#), [9](#)
 36. van den Oord, A., Vinyals, O., Kavukcuoglu, K.: Neural discrete representation learning. In: *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)* (2017) [3](#)
 37. Ouyang, J., Wu, H., Wang, M., Zhou, W., Li, H.: Contextual similarity aggregation with self-attention for visual re-ranking. In: *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)* (2021) [4](#)
 38. Park, W., Kim, D., Lu, Y., Cho, M.: Relational knowledge distillation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 3967–3976 (2019) [4](#), [6](#)
 39. Passalis, N., Tefas, A.: Learning deep representations with probabilistic knowledge transfer. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 268–284 (2018) [4](#)
 40. Peng, B., Jin, X., Liu, J., Li, D., Wu, Y., Liu, Y., Zhou, S., Zhang, Z.: Correlation congruence for knowledge distillation. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. pp. 5007–5016 (2019) [4](#), [6](#)
 41. Rao, H., Xu, S., Hu, X., Cheng, J., Hu, B.: Augmented skeleton based contrastive action learning with momentum lstm for unsupervised action recognition. *Information Sciences* **569**, 90–109 (2021) [2](#), [4](#), [5](#), [10](#)
 42. Romero, A., Ballas, N., Kahou, S.E., Chassang, A., Gatta, C., Bengio, Y.: Fitnets: Hints for thin deep nets. In: *Proceedings of the International Conference on Learning Representations (ICLR)* (2015) [3](#)
 43. Shahroudy, A., Liu, J., Ng, T.T., Wang, G.: Ntu rgb+d: A large scale dataset for 3d human activity analysis. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 1010–1019 (2016) [3](#), [9](#)
 44. Shi, L., Zhang, Y., Cheng, J., Lu, H.: Skeleton-based action recognition with directed graph neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 7912–7921 (2019) [2](#)
 45. Shi, L., Zhang, Y., Cheng, J., Lu, H.: Two-stream adaptive graph convolutional networks for skeleton-based action recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 12026–12035 (2019) [2](#)
 46. Shi, L., Zhang, Y., Cheng, J., Lu, H.: Adasgn: Adapting joint number and model size for efficient skeleton-based action recognition. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. pp. 13413–13422 (2021) [2](#)

47. Si, C., Chen, W., Wang, W., Wang, L., Tan, T.: An attention enhanced graph convolutional lstm network for skeleton-based action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1227–1236 (2019) [2](#)
48. Si, C., Nie, X., Wang, W., Wang, L., Tan, T., Feng, J.: Adversarial self-supervised learning for semi-supervised 3d action recognition. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 35–51 (2020) [12](#)
49. Soomro, K., Zamir, A.R., Shah, M.: Ucf101: A dataset of 101 human actions classes from videos in the wild. arXiv preprint arXiv:1212.0402 (2012) [5](#)
50. Su, K., Liu, X., Shlizerman, E.: Predict & cluster: Unsupervised skeleton based action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9631–9640 (2020) [2](#), [4](#), [10](#), [11](#)
51. Tejankar, A., Koohpayegani, S.A., Pillai, V., Favaro, P., Pirsiavash, H.: Isd: Self-supervised learning by iterative similarity distillation. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 9609–9618 (2021) [4](#)
52. Tejankar, A., Koohpayegani, S.A., Pillai, V., Favaro, P., Pirsiavash, H.: Isd: Self-supervised learning by iterative similarity distillation. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 9609–9618 (2021) [7](#)
53. Thoker, F.M., Doughty, H., Snoek, C.G.: Skeleton-contrastive 3d action representation learning. In: Proceedings of the 29th ACM International Conference on Multimedia (ACM MM). pp. 1655–1663 (2021) [4](#), [5](#), [9](#), [10](#), [11](#), [12](#)
54. Tianyu, G., Hong, L., Zhan, C., Mengyuan, L., Tao, W., Runwei, D.: Contrastive learning from extremely augmented skeleton sequences for self-supervised action recognition. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI) (2022) [4](#), [5](#), [10](#), [11](#), [12](#)
55. Tung, F., Mori, G.: Similarity-preserving knowledge distillation. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 1365–1374 (2019) [4](#), [6](#)
56. Wang, M., Ni, B., Yang, X.: Learning multi-view interactional skeleton graph for action recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) (2020) [2](#)
57. Wang, N., Zhou, W., Li, H.: Contrastive transformation for self-supervised correspondence learning. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). pp. 10174–10182 (2021) [5](#)
58. Wu, H., Wang, M., Zhou, W., Li, H., Tian, Q.: Contextual similarity distillation for asymmetric image retrieval. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9489–9498 (2022) [4](#)
59. Xu, J., Yu, Z., Ni, B., Yang, J., Yang, X., Zhang, W.: Deep kinematics analysis for monocular 3d human pose estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 899–908 (2020) [1](#)
60. Yan, S., Xiong, Y., Lin, D.: Spatial temporal graph convolutional networks for skeleton-based action recognition. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). pp. 7444–7452 (2018) [11](#)
61. Yang, S., Liu, J., Lu, S., Er, M.H., Kot, A.C.: Skeleton cloud colorization for unsupervised 3d action representation learning. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 13423–13433 (2021) [12](#)
62. Zhang, P., Lan, C., Xing, J., Zeng, W., Xue, J., Zheng, N.: View adaptive neural networks for high performance skeleton-based human action recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **41**(8), 1963–1978 (2019) [2](#)

63. Zhang, P., Lan, C., Zeng, W., Xing, J., Xue, J., Zheng, N.: Semantics-guided neural networks for efficient skeleton-based human action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1112–1121 (2020) [2](#)
64. Zhang, X., Xu, C., Tao, D.: Context aware graph convolution for skeleton-based action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 14333–14342 (2020) [2](#)
65. Zheng, N., Wen, J., Liu, R., Long, L., Dai, J., Gong, Z.: Unsupervised representation learning with long-term dynamics for skeleton based action recognition. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). pp. 2644–2651 (2018) [2](#), [4](#), [10](#), [11](#), [12](#)